# DATA ETHICS: BUILDING TRUST

## How Digital Technologies Can Serve Humanity

Editors  **Christoph Stückelberger and Pavan Duggal (Eds.)**

Globethics

# Data Ethics: Building Trust

*How Digital Technologies*

*Can Serve Humanity*

# Data Ethics: Building Trust

*How Digital Technologies*
*Can Serve Humanity*

Christoph Stückelberger / Pavan Duggal (Eds.)

All web links in this text have been verified as of February 2023.

*The electronic version of this book can be downloaded for free from the Globethics.net website: www.globethics.net.*

# CONTENTS OVERVIEW

# CONTENTS

## 15 Brain in the Data: Neurotechnology in AI Systems and Management Applications.................. 293

*Alexander Ageev, Russia*

## Part D: Serving Humans: Religions for Digital Justice ....................................................... 313

## 16 Digital Justice Manifesto: A Call to Own our Digital Future................................................ 315

*Just Net Coalition / World Association of Christian Communication WACC*

## 17 A Vision of Digital Justice.................................... 327

*World Council of Churches, Central Committee*

## 25 QuantumTechnologies: Next Arms Race China-USA..........................................................................387

*Matthias Sander*

## 26 Ukraine War: Artificial Intelligence on the Battlefield ........................................................393

*Mita Pressl / Pascal Burkhard, Ukraine*

## 27 International Humanitarian Law and Cyber Operations during Armed Conflicts .......................397

*International Committee of the Red Cross ICRC*
*Position Paper*

## 28 Establishing a Digital Geneva Convention......... 415

*Dan Lohrmann*

## 29 Data Storage Ethics: Security and Responsibility at All Levels ..................................... 425

*Gilles Bach, Switzerland*

# 0

# INTRODUCTION

*Christoph Stückelberger and Pavan Duggal[1]*

## 01 Data: The Gold and Oil of the Current Revolution

The development of humanity is influenced by great new ways of thinking and spiritual wisdom. In an interactive way it is at the same time influenced by new technologies and access to resources: The Neolithic Agricultural Revolution by the plough and thus digging the land as resource. The Stone, Bronze and Iron Age by the technologies using stone, bronze and iron as new materials. The First Industrial Revolution used water for steam power, the Second Industrial Revolution used the power of electricity, oil and gas for mass production. The Third Industrial Revolution used electronics for automated production. The current Forth Industrial Revolution uses *data* as raw material.[2]

---

[1] Christoph Stückelberger, Professor of Ethics, Founder and President of Globethics, Geneva/Switzerland. Co-Editor of this volume. Stueckelberger@globethics.net. www.christophstueckelberger.ch.

Pavan Duggal, Advocate at the Supreme Court of India, Cyber Law specialist, Member of the Internat. Board of Globethics.net, Delhi/India. Co-editor of this volume. pavan@pavanduggal.com. More at www.pavanduggal.com.

[2] More in Christoph Stückelberger, Cyber Society: Core Values and Virtues, in Christoph Stückelberger/Pavan Duggal, Cyber Ethics 4.0. Serving Humanity with Values. Geneva: Globethics.net, 2018, 23-54 (27-31: Fourth Industrial Revolution)

Data is the most precious current good, more precious than water, oil, gas, electricity, gold and land.  Every new technological revolution adds a new resource, but still depends on most of the former resources too. However, the former ones no longer play the leading and driving role of innovation and power dominance.

Each such fundamental shift in the development of humanity leads first to a hype of wild exaggerated expectations of part of humanity - mainly elites and technocrats - to solve all problems. For others it leads to profound fears of the destructive power of the new material and technology. In a second phase both, too high expectations and too deep fears, become moderated by experience of moderate successes and failures. A third phase then leads to ethical benchmarks and legal regulations in order to channel the "wild water" into controllable riverbeds or tubes, to reduce its destructive powers and use it for the benefit of humanity and a life in dignity.

## 02 Objectives of this Book

This book *"Data Ethics: Building Trust. How Digital Technologies Can Serve Humanity"* aims at

- contributing to this current third phase of differentiated management the data-revolution from ethical, values-driven perspectives;
- looking at the chances to enhance life, sustainability, peace, equality and justice and at the risks of "new wine in old tubes" of power struggle, dominance, abuse, violence, control, oppression and exploitation;
- being inclusive, including different world views, philosophical, religious, disciplinary and transdisciplinary contributions from all continents;
- offering arguments for policy-makers for ethical benchmarks and legal regulations;

- inviting educators to use the rich and diverse articles in teaching, trainings and research;
- building trust through active listening, fair dialogue and joint actions between those in favour of and experiencing the new data opportunities and those fearing and experiencing their abuse for oppressive practices. Often, both voices are within ourselves.

## 03 What is Data Ethics? Definitions

Before moving forward, it becomes important for us to dive deep and understand what this new evolving field of data ethics is.

- *Ethics* looks for orientation what is good and bad, right and wrong, life enhancing and life destructing.
- *Fundamental ethics* looks at the basic philosophical and religious concepts, values and norms.
- *Applied ethics* then allies such values and norms to specific sectors and topics. Ethics of technologies is such a sector.
- *Data ethics* is a sub-topic of technology ethics. Ethics of AI, ethics of data storage, ethics of data commercialisation etc. are sub-topics of data ethics. The diversity of topics is shown in this book.
  *Data ethics as a concept has been defined by various stakeholders.*
- Data ethics is a branch of ethics that evaluates *data practices* - collecting, generating, analysing and disseminating data, both structured and unstructured - that have the potential to adversely impact people and society. It includes addressing and recommending concepts of right and wrong conduct, with transparency in and defensibility of actions and decisions driven by automated/artificial intelligence (AI) in relation to data in general and personal data in particular.[3]

---

[3] Cognizant. 2023. Data ethics, https://www.cognizant.com/us/en/glossary/data-ethics#:text=Data%20ethics%20is%20a%20branch,adversely%20impact%20people%20and%20society.

- Data Ethics describe a code of behaviour, specifically what is right and wrong, encompassing the following:
- *Data Handling:* generation, recording, curation, processing, dissemination, sharing, and use.
- *Algorithms:* AI, artificial agents, machine learning, and robots.
- *Corresponding Practices:* responsible innovation, programming, hacking, and professional codes. [4]

 

*Other Definitions* of Data Ethics refer to the said discipline as follows:-

- *Data mindfulness* guiding the ethical use of data within all organizations.[5]
- *A specialized area of Data Governance and Stewardship,* where data can be productionalised and, in large part, automated through the use of good tooling.[6]
- Not only moral guidance as to "what data should be collected and how it should be used," but also "*who gets to make those decisions in the first place.*"[7]
- *A code of conduct* as practical tool for data ethics for data scientists, similar to the purpose of the Hippocratic Oath in guiding medical professionals.[8]

---

[4] Floridi L, Taddeo M. What is data ethics? Philos Trans A Math Phys Eng Sci. 2016 Dec 28;374(2083):20160360. https://doi.org/10.1098/rsta.2016.0360

[5] Roe, Charles. 2018. Dataversity, https://www.dataversity.net/enterprise-data-world-2018-illuminating-data-cave/

[6] Cf. Tech Crunch.

[7] Arvanitakis, James, Francis, Andrew and Obst, Oliver. 2018. Data ethics is more than just what we do with data, it's also about who's doing it, Phys.org. https://phys.org/news/2018-06-ethics.html

[8] Bloomberg, 2018. Bloomberg BrightHive, and Data for Democracy Launch Initiative to Develop Data Science Code of Ethics, https://www.prnews wire.com/news-releases/bloomberg-brighthive-and-data-for-democracy-launch-initiative-to-develop-data-science-code-of-ethics-300524958.html

- *Code* guiding Data Scientists' behaviors, to "better human society."[9]
- The UK Government defines data ethics or a data ethics framework as a "*set of principles* to guide the design of appropriate data use in the public sector. It is aimed at data practitioners (for example, statisticians, analysts and data scientists), policy makers, operational staff and people helping to produce data informed insight".[10]

There have been a lot of developments as far as data ethics is concerned. *Data ethics as a discipline* has been evolving for years.

The ethics of data focuses on ethical problems posed by the collection and analysis of large datasets and on issues ranging from the use of big data in biomedical research and social sciences, to profiling, advertising and data philanthropy as well as open data. Trust and transparency are also crucial topics in the ethics of data, in connection with an acknowledged lack of public awareness of the benefits, opportunities, risks and challenges associated with data science.[11]

## 04 Authors and Structure of this Collection

This book *Data Ethics: Building Trust* includes 40 articles.

*The authors* come from all continents and manifold disciplines: data science, engineering, theological ethics, philosophical ethics, law, political science, education, economy, management, media etc. Articles are

---

[9] Leetaru, Kalev. 2017. Is It Too Late For Big Data Ethics? Forbes, https://www.forbes.com/sites/kalevleetaru/2017/10/16/is-it-too-late-for-big-data-ethics/#54580d603a6d

[10] Guidance Data Ethics Framework, Central Digital and Data Office, Gov.uk, May 2021.

[11] Floridi L, Taddeo M. What is data ethics? Philos Trans A Math Phys Eng Sci. 2016 Dec 28;374(2083):20160360. doi: 10.1098/rsta.2016.0360.

written by individual authors[12], but are also statements and excerpts of policy papers of institutions.

*The Global Series of Globethics Publications* reflects this concept of inclusive dialogue: the editors are for each volume from at least two continents and two disciplines, in this case Christoph Stückelberger from Switzerland/Europe as ethics expert and development practitioner with Pavan Duggal from India/Asia as legal expert and advocate.

*The structure of the book* is visible in the *sections/parts*:

A)   Data Ethics: Values, Trust and Law
B)   Data Ethics: Artificial Intelligence, Robots and Humans
C)   Serving Humans: Health and Neuroscience
D)   Serving Humans: Religions for Digital Justice
E)   Serving Humans: Farming, Business, Finance
F)   Serving Humans: Security, War, Peace
G)   Serving Humans: Data Governance and Geopolitics
H)   Serving Humans: Media and Education

## 05 Working Definitions

In this book, manifold terms related to data are used. The following, selected, working definitions are taken from the "African Union Data Policy Framework", endorsed 2022[13] (see chapter 34 in this book).

-   *Anonymisation* is the removal of direct and indirect personal identifiers from data.
-   *Continental*, for the purposes of this framework, refers to Africa.
-   *Data classification* is broadly defined as the process of organising data by relevant categories so that it may be used and protected more efficiently.

---

[12] The references and footnotes are not all fully unified, as they reflect some diversity of disciplines and academic traditions globally.

[13] African Union. 2022. AU Data Policy Framework, https://au.int/en/documents/20220728/au-data-policy-framework.

- Foundational data infrastructure refers to advanced technologies which facilitate the intensive use of quality data. This may include broadband networks, data centres and cloud services, electronic hardware and software, and digital applications available on the Internet.
- *Data ecosystem* - for the purposes used here not only to the programming languages, packages, algorithms, cloud-computing services, and general infrastructure an organisation uses to collect, store, analyse, and leverage data, but to the underlying value chain associated with data as a factor of production, the governance of data systems and the protection of data subjects.
- *Data minimisation* is a principle within data protection frameworks, which entrenches collecting the minimum amount of personal data needed to deliver an individual element of a service or product.
- *Datafication* refers to the process by which daily interactions of living things can be rendered into a data format and put to social and economic use.
- *E-commerce* can be summarised as commercial transactions occurring through electronic channels - buying and selling of goods or services via the Internet and the transfer of money and data to complete the sales - by methods specifically designed for the purpose of receiving or placing orders.
- *Cloud services* are used on-demand at any time, through any access network, using any connected devices that use cloud computing technologies. They utilise software and applications located on the cloud and not on users' own devices.
- *Cloud-based services* include mass-market applications (i.e. social media and webmail offered over the Internet). The data does not sit on the individuals' devices but is stored remotely in a data centre. Examples include Facebook, YouTube and Gmail.
- *Digital identity* is a set of electronically captured and stored attributes and/or credentials that uniquely identify a person, enabling the distinction of one individual from another.
- *Digital capability* is the term used to describe the skills, literacy, social norms, and attitudes that individuals and organisations need to thrive, live, learn and work in a digital society and economy.
- *Consent of the data subject* means any freely given, specific, informed and unambiguous indication of the data subject's wishes by which he or she, by a statement or by clear affirmative action, signifies agreement to the processing of personal data relating to him or her.
- *Cybercrime:* Unlawful acts which affect the confidentiality, integrity, availability and survival of information and communication technology systems, the data they process and the underlying network infrastructure (Malabo Convention).
- *Cybersecurity:* Cyber security refers to the body of technologies, processes, and practices designed to protect networks, devices, programs, and data from attack, damage, or unauthorised access (https://digitalguardian.com/blog/what-cyber-security).

- *Data controller* means any natural or legal person, public or private, any other organisation or association that alone or jointly with others, decides to collect and process personal data and determines the purposes.
- *Data protection* regulates how data is used or processed and by whom, and it ensures citizens have rights over their data. It is particularly important in ensuring digital dignity, as it can directly address the inherent power imbalance between 'data subjects' and the institutions or people who collected data.
- *Data protection authorities (DPAs)* are independent public authorities that monitor and supervise, through investigative and corrective powers, the application of the data protection law. They provide expert advice on data protection issues and handle complaints that may have breached the law.
- *Data subjects* means any natural person that is the subject of personal data processing (Malabo Convention).
- *Harmonisation* is ensuring uniformity in the systems through the use of minimum standards to facilitate interoperability and legal and trust frameworks (e.g. for levels of assurance) to set rules and build confidence in respective systems.
- *Interoperability* is the ability of different function units – e.g. systems, databases, devices, or applications – to communicate, execute programs, or transfer data in a manner that requires the user to have little or no knowledge of those functional units (adapted from ISO/IEC 2382:2015).
- *Level of assurance (LOA)* is the ability to determine, with some level of certainty or assurance, that a claim to a particular identity made by some person or entity can be trusted to actually be the claimant's "true" identity (ID4D Public-Private Cooperation). The overall level of assurance is a function of the degree of confidence that the applicant's claimed identity is their real identity (the identity assurance level or IAL), the strength of the authentication process (authentication assurance level or AAL), and—if using a federated identity—the assertion protocol used by the federation to communicate authentication and attribute information (federation assurance level or FAL) (adapted from NIST 800-63:2017).
- *Open standards* are standards made available to the general public and are developed (or approved) and maintained via a collaborative and consensus-driven process. Open standards facilitate interoperability and data exchange among different products or services and are intended for widespread adoption (adapted from ITU-T).
- *Open data:* Open means anyone can freely access, use, modify, and share for any purpose (subject, at most, to requirements that preserve provenance and openness (http://opendefinition.org/).
- *Personal data* means any information relating to an identified or identifiable natural person by which this person can be identified, directly or indirectly in particular by reference to an identification number or more factors specific to his/her physical, physiological, mental, economic, cultural or social identity.
- *Privacy and security by design* means proactively embedding privacy and security mechanisms into the design and operation of products and services,

both non-IT and IT systems, networked infrastructure, and business prac-
tices. This requires that privacy and security governance is considered
throughout the whole engineering process and product lifecycle.
- *Pseudonymisation* is processing of data so that it cannot be associated with
  an individual without additional information.
- *Regional* for the purposes of this AU Framework refers to the five regions
  of Africa recognised by the African Union.
- *Sensitive data* means all personal information relating to religious, philo-
  sophical, political opinion as well as to sex life, race, and health, social
  conditions of the data subject (Malabo Convention).

# PART A

# DATA ETHICS:
# VALUES, TRUST AND LAW

# 1

# HUMAN-CENTERED DATA ETHICS: HOW TO BECOME BETTER HUMANS?

*Christoph Stückelberger, Switzerland*[14]

## 1.1 From Morning to Evening 24/7/365

It is Wednesday, 8 February 2023. Early morning, the alarm clock wakes me up. The program is stored as data in the watch. The water of the shower, the electricity of the apartment lights, and the heating of food are measured and stored as data. The barcode on the breakfast food package reminds me that my data of shopping can be stored through my bill linked to my credit card. My little car is a data center in itself and my bus travel ticket on my mobile phone reminds me of the dozens of de-

vices on my mobile phone that are also data centers. My little walk for a short break and refresh during the writing of this article shows me the number of steps taken today, the number of calories burnt and a friendly invitation to walk more and get a smiley as reward for my obedience to 'mother' health-meter. The day goes on with my data collection until I switch off the light at late evening. I still resist activating a device to measure my pulse and heart frequency during my sleep. In future, I may ask my watch in the morning "can you explain the meaning of my dream at 2am?" A "dreamGPT", the next generation after the "ChatGPT", may then give me the interpretation of my dreams, either in psychoanalytical version or sociopolitical or religious version, depending on the interpretation I am looking for.

These thousands of data collected day by day 24/7/365, only for one person, amounts to millions or billions data per year, and is reality today. The Internet of Things (IoT)[15] is reality today. I am not much afraid for my own person. I do not have much to hide, and I live in a relatively secured country (Switzerland). I am also retired, not fearing an employer could use data against me as an employee. However, data collection, storing and use is a serious concern and threat for millions if not billions of people around the globe, not to speak about the societal, economic and political implications.[16] An example are also smart cities: they offer obvious advantages in many digital services[17] and at the same create fears of increased dependency from interconnected data collection.

---

[15] See also Ch. 20 below.

[16] Expressed in many of the 45 chapters of this book.

[17] See e.g. IMD - International Institute for Management Development. *Smart City Index*: https://www.imd.org/smart-city-observatory/home/

## 1.2 Data: The Driving Fuel in All Sectors

### 1.2.1 Deus ex machina – or - Diabolus ex machina?

Data is the driving energy, so to say the fuel, in all sectors of society in the 21st century. As written in the introduction to this book (chapter 1): "Data is the most precious current good, more precious than water, oil, gas, electricity, gold and land. Every new technological revolution adds a new resource, but still depends on most of the former resources too. However, the former ones no longer play the leading and driving role of innovation and power dominance."

Three main phases can be distinguished when new technologies arise, as we describe in the introduction of this book[18]:

1. A hype of wild exaggerated expectations of part of humanity - mainly elites and technocrats - to solve all or major problems with the new technologies. Problems, which for long time seemed to be unsolvable, can suddenly be solved: enough energy, health improvement, peace, communication, facilitated daily life.

   For others, the new technologies lead to profound fears of the destructive power of the new development. Loss of jobs, total Orwell's control and surveillance, oppression of free speech and other human rights etc. Not salvation from God out of the machine, but destruction from the devil out of the machine. In the phenomenology of religions, religion was described already a hundred years ago as "fascinosum and tremendum"[19], a reality which fascinates, attracts and at the same time makes us tremble, fearful.

   In both cases, the debate about the new technology has almost a religious dimension. It is linked to the fact that a high uncertainty ex-

---

[18] See also Christoph Stückelberger, *Cyber-Society: Core Values and Virtues*, in Christoph Stückelberger/Pavan Duggal (Eds.), Cyber Ethics 4.0. Serving Humanity with Values, Geneva: Globethics.net, 2018, 23-53 (27-31)

[19] G. van der Leeuw, *Phänomenologie der Religion*, München: E. Reinhardt, 1925.

ists, as not much experience – negative or positive – is available and thus projections can blow in all directions.

2. In a second phase both, too high expectations and too deep fears, become moderated by experience of moderate successes and failures. As described above, I see the immediate benefit of my data collecting devices, but I still have doubts about what happens with my data, as it seems controlled by others, a kind of black box. That is the period where building trust is essential.

3. A third phase then leads to ethical benchmarks, legal regulations, more control and ownership of data for individuals and institutions, in order to channel the "wild water" into controllable tubes, to reduce its destructive powers and use it for the benefit of humanity and a life in dignity. Currently, we are in this third phase with global efforts for governance and in power struggles in geopolitics.

### 1.2.2 Acceleration, Potential and Vulnerability

As we are in this third phase, the first two decades of the 21st century brought substantial progress for transparency, governance, ownership of data and regulations. Many articles in this book show it.[20] These efforts will be accelerated in the current third decade of the 21st century, but at the same time, the enormous speed of data-related developments constantly lead to new challenges.

In a daily, almost hourly rhythm, I get news about these developments. Few examples of last few hours: Digital Economy Compass 2022, 240 pages for free download; Microsoft's OpenAI free test version of ChatGPT, personally promoted by the OpenAI director via my Linkedin account; news on Google catching up with its own ChatGPT; last days to upgrade my personal website on the next version of the system as my current older version will no more be supported in twenty days; urgent recommendation of Dropbox to upgrade my cloud service,

---

[20] Especially parts D, E and F (Ch. 16-38).

which stores my personal data, from 2 to 4TB; and so on. All about data. Stunning.

I benefit to a great extent from data-based technologies. I like and I am still willing to run with the time and the opportunities. However, when is stunning acceleration too much? When comes the moment of also feeling suspended and side lined like billions of people on this globe? And what about all the individuals, companies and institutions who are victims of stealing and hacking data, of money extortion, violence and war due to abuse of data?[21] We all remain somehow vulnerable.

### 1.2.3 A segmented, bipolar data world?

The internet, which was planned to be a global one, has become in recent years, more and more nationalised and controlled by national legislation, mainly justified by security reasons. At the same time, the big data giant companies are structured in a way, that two main political entities dominate the market: USA and China. Russia and EU also play a role, but to a lesser degree. USA and China can be seen as the new bipolar world as they offer – like a mirror! – exactly the same social media, search machines etc., with the main difference, that access to the data is restricted to the respective states (US through backdoors in the software and hardware, China through control of the companies).[22]

The current *polarisation — new bi-polar and at the same time multipolar* leads to a mixed picture of a globalised, yet fragmented, divided world, where the two superpowers USA and China look for global dominance. Other heavy powers like India, EU, Russia, but also past and coming stars like Indonesia, South Africa and Brazil play a role.

---

[21] See in this book Ch. 22 on blockchain security, 26 on data in war, 29 on data storage, 31 on stealing data and others.

[22] More in Christoph Stückelberger, *Globalance towards a New World Order. Ethics Matters and Motivates*. Handbook, Geneva: Gloethics.net, Nov 2022, 495-548.

The fight about access, control and dominance over data is a key battlefield in this polarised world, first for economic, second for political and military reasons.

Data is soft power. Yet, data are carried, processed and stored thanks to high performing hardware. The *semiconductor chips* became the symbol of the "heart" of our data age as the iron was at the heart of the iron age for ploughshares and weapons. The copper cables in the second industrial revolution have been at the heart of the second industrial revolution with the use of electricity etc.

*The shift from Westernization to Easternization[23]* with Asia as a key continent of the 21st century means not only a shift in values, cultures and political systems, but also the cultural understanding and praxis of data, data protection etc.

## 1.3 Data Ethics: Across Each Sector of Applied Ethics

Today, data is the key "raw material" in all sectors of society and all phases of individual and collective lives. Therefore, Data Ethics is also not only a subtopic of the ethics of technologies, but is an aspect in all sectors of Applied Ethics. The graph below[24] shows:

*Three main types of ethics* build the methodological roots or main trunk of the ethics tree. Normative ethics is the most relevant part, as it tries to find benchmarks of values-driven orientation for decisions and actions.

*The six main domains of ethics* are like the main branches of a tree or like the main provinces or states in a nation. Each of these main domains then has *sub-domains,* more specific areas like smaller branches on the tree.

*Example*: within the main domain of economic ethics/business ethics, trade ethics is one key part. Data ethics is now one subtopic of trade

---

[23] Ibid, 480-494.

[24] From Christoph Stückelberger, *Global Trade Ethics*, Geneva, WCC, 2002, 2.

ethics, as trading data with all its aspects of pricing, security, accessibility etc. other subtopics of economic ethics are labor ethics, ethics of economic systems, consumer ethics etc. All of them also include aspects of data ethics: consumer data protection, labor data etc. Equally, data is relevant in legal ethics, family ethics, community ethics, life ethics such as health ethics and medical ethics, environmental ethics such as data for climate mitigation etc., legal ethics such as the regulations of data collection, processing, storage and sales.

**Types of ethics** — Meta-ethics, Normative ethics, Descriptive ethics

**Main domains of ethics** — Life ethics, Community ethics, Environmental ethics, Business ethics, Political ethics, Cultural/religious ethics

**Specific domains of ethics:**

- Health ethics
- Medical ethics
- Sexual ethics
- Bioethics
- Life-form ethics
- Marital and family ethics
- Ethnic communities
- Intergenerational ethics
- Resources ethics
- Ethics of the bases of life
- Animal ethics
- Biodiversity ethics
- Corporate ethics
- Labor ethics
- Ethics of economic systems
- Trade ethics
- Consumer ethics
- Professional ethics
- Legal ethics
- Ethics of political systems
- Peace ethics
- Land ethics
- Media ethics
- Sports ethics
- Cultural ethics
- Ethics of religions

**Ethics**
Acting responsibly in all spheres of life

© Stückelberger: Global Trade Ethics

New special topics of data ethics pop up almost monthly. *Algor-ethics* is an example. It is used in the Vatican text (below, article 18). It deals with ethical questions of algorithms. *Tethics* is a term for Tech Ethicists (below, article 5), *Neurethics* a term for the ethics of neurological technologies (below article 15). I am sure, soon the *ChatGPT-Ethics* will come up.

It is promising to mainstream ethics in technologies in all applications. However, as professional ethicists we also need to remain critical, where (data-) ethics is used to find orientation in difficult ethical dilemmas and where it is a fashionable marketing tool and defence of companies to prevent binding regulations and independent assessments. Scandals of dismissing critical ethicists in the leading companies like google have been a signal for many to remain alert.

*In conclusion:* Data Ethics is not a small subtopic for tech specialists, but it is today a cross-sectoral topic, relevant in all domains and subdomains of ethics. Even when Pavan Duggal and I as editors of this book started the concept, we looked at data ethics in a more narrow sense. We then discovered that we have to include – pars pro toto – topics from all sectors, which led to 45 articles and much more would be needed.

## 1.4 Back to Fundamentals: Who Am I as Human?

The current hype on new data-related technologies, including the jump to quantum computing with – as it seems – endless new dimensions (below articles 7.5 and 25), raises new fundamental questions, especially on the relation between 'man and machine', humans and their technologies. The new technologies are no more just "mechanics" or more complex "engineering", but they get closer to biological and neurological phenomena. The tiny chip under the skin is a symbol for it. *Data in the Brain* and *Brain in the Data* (below article 15) are an expression of it. The distinction between human and his/her technologies is questioned. There is a still broad consensus about the distinction between a

human being, and non-human beings (animals and plants). A few deny this distinction and a majority —including myself— must rediscover the close interdependence of human and non-human creation/creatures as parts of one ecosystem.[25]

The borders between human and non-human are at the moment less clear when it comes to the new data developments. Are we as humans "just" data and information (as hundred years ago positivists declared humans "just" as water, chemistry and physics)? Below, article 7 shows the relation between data, information and ethics.

Even more radical is the question: *Who am I as a human?* What is my specificity, if organised data can do many things better and faster than we as humans: translations in seconds, complex counting in milliseconds, formulating text with AI such as the new ChatGPT in few minutes, (below articles 44-45). Medical operations performed more precisely than by a doctor, Autonomous cars – until there is an accident.

The human brain is still extremely powerful with its millions of synapses, connections based on social, spiritual and physical experiences. *Who am I as a human? What is my identity? What is the other in me, what is me in the other? We will rethink individuals and community, individuals/community and the environment, autonomy and interdependence, me and we as humans in space and time. Data ethics brings as back to basic philosophical and religious questions, back to fundamentals. Back to the roots of human existence.*

Therefore, data ethics is a topic for every human, not only for tech specialists with a short course or a master's in applied ethics, not only for lawyers and politicians to formulate regulations, not only for medical and educational specialists to improve their instruments and methodologies. Data ethics just invites all of us to reflect on our human existence,

---

[25] In details in my book Christoph Stückelberger, *Umwelt und Entwicklung. Eine sozialethische Orientierung*, Stuttgart: Kohlhammer, 1997. Also available in Chinese, Korean and Indonesian.

what counts, what is less important, what is short term fashionable and what is long term sustainable.

## 1.5 How to Become Better Humans?

The key ethical question is: How to become better humans? The key ambition is not how to compete with robots and algorithms, with big data and quantum physics. The key question is simple: How can new technologies serve humans and the whole creation? How can data (in the broad sense of this book) support core values and virtues and defeat non-values and vices? Or just simply: How can we strengthen love and reduce hate? How can we support peace and overcome killings? How can we find harmony with creation and avoid disasters? How can data-based technologies support these goals.

Let us unfold this in seven aspects in short.

### 1.5.1 Remaining faithful to core values and virtues

Technologies come and go; values and virtues stay. In 2022, I attended an international workshop on ancient philosophies and modern technologies, in Athens, the academy. It was amazing to see, how values and virtues of 2500 years old philosophies of Socrates and Aristotle in Athens, Confucius in China and founders of religions such as Buddha, Jesus and Mohammed are still valid! They are not outdated, but refer to the center of human beings and vocation.

I have worked for decades on global ethics in order to find common values and virtues throughout cultures, religions and millennia and at the same time find the specificities and differences in contextualisation and temporary bound values and virtues. I only refer here to the tree of values and the tree of virtues, which are broadly accepted globally and

unites us as humanity. They can serve as a benchmark of orientation in dealing with data and its ethics. [26]

| **Tree of Values** | **Tree of Virtues** |
|---|---|
| 1. Responsibility | 1. Integrity |
| 2. Freedom | 2. Compassion |
| 3. Justice | 3. Care |
| 4. Equity | 4. Transparency |
| 5. Peace | 5. Accountability |
| 6. Security | 6. Reliability |
| 7. Community | 7. Respect |
| 8. Inclusiveness | 8. Humility |
| 9. Participation | 9. Courage |
| 10. Forgiveness | 10. Gratitude |
| 11. Stewardship | 11. Generosity |

These values and also these virtues are not isolated, but interconnected. I call them relational values and virtues. This is very important also for data ethics as this ethical concept avoids extreme ethical positions. E.g. data protection and security as an absolute principle could hinder freedom of expression, community and inclusiveness. Absolute freedom would lead to huge inequality and needs to be balanced with fair access and justice/equity. The following two graphs show relationality of values and the relationality of virtues.[27]

---

[26] More in Christoph Stückelberger, Cyber-Society: Core Values and Virtues, in Christoph Stückelberger/Pavan Duggal (Eds.), *Cyber Ethics 4.0. Serving Humanity with Values*, Geneva: Globethics.net, 2018, 23-53 (33); broader in Christoph Stückelberger, *Globalance*, 2nd edition, 2022, op. cit, 237-271.

[27] Ibid, Cyber Society, 34f.

## 1.5.2 Access to education, information and formation

Values-driven decisions and actions need a lot of information in the modern complex world. But let me start with education: character education starts at the very beginning of human existence as a child. It does not primarily depend on information, but on love and deeply rooted values and virtues, based on a vision what a dignified life means. I met people with only very little formal education, but had profound character of integrity. Based on such character education, manifold information,

schooling and formal formation and trainings are needed in order to make an informed judgment.

### 1.5.3 Being courageous in using power and resisting the abuse of power

*Data/information power* is today the most dynamic and profitable production factor. Those who collect, transform and own big data, are key drivers such as the big four in USA (Google, Amazon, Facebook/Meta, Apple, GAFA) and the big four in China (Baidu, Alibaba, Tencent, Huawei, BATH).

| Service | USA | China | Russia | EU |
|---|---|---|---|---|
| Search machine | Google | Baidu | Yandex, Moi Mir | (Bing, small) |
| Social Net | Facebook | Tencent | VK.com | - |
| E-commerce | Amazon | Alibaba | Bringly | Small, national |
| messenger | WhatsApp | Wechat | Quora | - |
| Short messenger | Twitter | Wechat | Twitter | - |
| videos | Youtube | Youku Tudou | Rutube | - |
| E-payment | Paypal | Alipay | Yandex.money | Banks, various |
| **Data ownership** | Companies | Companies | companies | companies |
| **Data access** | State | State | State | - |
| **Data protection** | National laws | National laws | National laws | EU GDPR 2018 |

Technologies in all human history have been used for the good and abused for the bad. But the main driver was and is always to gain power, to be stronger and faster than the other (competitor, enemy). This is not different with data technologies. The digital dependence is an expression of this power-dependence.[28]

Therefore, a key question of data ethics is how the enormous power of data can be used to strengthen the above-mentioned values and virtues and how they can be used to overcome destructive power patterns and makes life less and not more vulnerable. The authors in this book show manifold technical, administrative, regulatory and educational answers.

---

[28] The *EU Digital Dependence Index* measures digital dependence by country for software, hardware and IP. https://digitaldependence.eu/en/

The political awareness of this power questions increased substantially in the last ten years. Self-regulatory measures of companies (as effort to protect their power and avoiding state regulations) such as Huawei's Data control center in Brussels for their customers (below article 21), Googles Ethics standards, the voluntary Swiss Digital Trust Label (below article 37) and governmental and intergovernmental standards and legislations are multiplied. The new African Data Policy Framework (below article 34) and the very new European Digital Services Act of the European Union (below article 35) are examples.

### 1.5.4 The right to forget (data)

The power of our "data age" is also due to the fact that – with the exponentially increased storage capacities and future quantum technologies - data can be accumulated, stored, reorganised, and re-used for a long time. (Even though I still doubt that today's data could be read and accessed in 300 years, as it would need constant transfer on new data carriers and constant access to electricity to be retrieved! I have at home old books from the 15[th] century, 500 years old, which I can read and understand).

How does data storage (below article 29) help improve humanity and how could it hinder it? A key part of being and becoming human is our human brain. It has the immense, wonderful capacity to sort, to visualise, to memorize, to remember – and to forget. To forget seems a deficit, painfully seen e.g. at Alzheimer patients. And we all know this daily situation where we cannot remember a name or a fact. However, the capacity to forget is also an incredible gift and part of sorting what is important in our lives, what remains and what makes our personality.

There is the legal right to forget e.g. a criminal act. If I drive too fast and get a citation or court punishment, it must be deleted from the record after a certain time. The incident may continue to exist in the memory of those in the car only – a very limited circle.

The "Right to be forgotten" (RTBF) is a human right and a legal concept e.g. in the European Union and other countries[29]. It is the right to have private information (data) about a person be removed from social media and internet searches. It is part of the value of human dignity, of individual autonomy, of the right to privacy and freedom of expression. It is part of making us human. The right to forget and to be forgotten is a delicate and also controversial topic. It can also be politically instrumentalised and thus abused, e.g. in censorship or rewriting history as effort to forget parts of the history of a nation.

Also, an individual himself/herself needs the possibility to filter, to forget and overcome e.g. traumatic experiences or small issues in relations in order to be free for the future and not bound to the past like a slave.

To forget is not only a technical and legal matter, but also profoundly a psychological, psychosomatic (the body remembers!), spiritual and pastoral matter. I write this as an individual, but also as a pastor, with my wife as psychotherapist with a specialisation in trauma healing. The ability to love – at least in the holistic meaning of Christian agape-love – includes the capacity to "liberate, remember, repent, forgive and forget"[30]. This is a complex, important liberation process and key part of one's personality. When is the time to remember and when to forget? How can repenting, asking for forgiveness and forgiving others be part of liberation, peace, reconciliation and harmony? How can physical pain in one's body as embodied memory (of data, we would now say) be overcome?

These aspects are part of a holistic Data Ethics and should be more included in future considerations. Intercultural considerations have to be included. In Western, South American and African countries with a long

---

[29] Legislations in various countries are analysed in https://en.wikipedia.org/wiki/Right_to_be_forgotten (accessed 9 Feb 2022).

[30] Christoph Stückelberger, *Globalance*, 2nd edition, 2022, op. cit., pp. 826-833, ch.8.3.4: "Agape-love: liberate, remember, repent, forgive and forget".

history and culture of Christian values of individual and collective re-
penting and forgiving, forgetting as part of forgiving is embedded in the
culture. In Asian cultures and tradition, I often observe, that repenting is
present, but forgiving much less a part of value system. This influences
the way, data collection and use are perceived and practised. Data ethics
then needs to be contextualised.

### 1.5.5 The spiritual sources of energy and orientation

*I am somebody because of my abilities:* Is my personality and reputa-
tion built on my skills, knowledge, power, beauty, connections and per-
formance, all in all the sum of my "bio-data"? This is the average per-
ception of a person today. It was always the case in history, but now
even more so with the social media, and the current narcissistic culture.
The spiritual alternative is:

*I am somebody, because I am loved.* I am who I am because God (the
divine) loves me as I am. Not because of my skills, my character, my
good deeds and values, but because of unconditional love of me as a
human, love from God and the others. This is at the very core of Chris-
tian faith. "This is love: not that we loved God, but that he loved us and
sent his son as an atoning sacrifice for our sins. Dear friends, since God
so loved us, we also ought to love one another" (1 John 4:10). Uncondi-
tional love is also the source of other religious, spiritual and philosophi-
cal traditions.

In data ethics, this leads also to the profound question of knowledge
and wisdom. Data is first of all the raw material. All is data (below arti-
cle 7). By collecting, sorting, processing, data become information and
knowledge. But can it produce wisdom? What is wisdom? Indian tradi-
tions with their profound scientific, data-based knowledge and at the
same time a very rich culture of wisdom and spiritual depth can give

some answers. Confucian[31], Buddhist, Daoist, Christian, Muslim, Jews worldviews contribute their answers. We cannot explore more in this short article, but my suggestion is to grab deeper and rediscover wisdom, which is much more than data processing and knowledge.

### 1.5.6 The networking power of communities of trust

"How to become better humans?" is the key question of this article. It is a question for individual character development, institutional frameworks, but also human communities. The African *Ubuntu* concept (below article 8) of "I am because we are" reminds us, that we cannot become better humans individually only, but only in community with the others.

*Building trust* as the key challenge with new technologies is a matter of human interaction. Trust can be built only in open, transparent, moni-tored, sustainable human interactions. Our "data age" shows the huge potential of improving and facilitated human interactions thanks to data, based devices, from home-based family chat to global alliances, from business relations to sport, from campaigning and advocacy networks to intergovernmental alliances. Information and data sharing is possible nowadays as never before in history, with stunning speed.

An *example* in this book is the "Geneva Digital Atlas 2.0" (below ar-ticle 38). On a few square kilometres of the city of Geneva, over 40 international organisations, NGOs, foundations and research institutions build a locally located and globally acting ecosystem of actors for the digital world. A huge potential for building alliances of trust.

### 1.5.7 Data-Sunday, Data holiday, Data fasting?

We now expressed and shared the enthusiasm for the existing and fu-ture potential of the digital world as data world. Chances and threats of

---

[31] On Chinese Philosophers e.g. Bing Song (ed.), *Intelligence and Wisdom.* Citic Press Corporation, 2021. https://doi.org/10.1007/978-981-16-2309-7_1. More below in article 10.

'becoming better humans' are explored. Let us add as point 7 a short reflection on the 7th day of the week (also called the first day of the week): the Sunday in Christian tradition, the Sabbath in Jews tradition and the Friday in Muslim tradition. All Abrahamic religions express the chance, the offer and the need for rest, for a break, to enjoy the fruits of work, the have space for the community life, to adore the beauty of creation, a weekly holiday. The Abrahamic religions express it in the creation story, so to say the foundational myth for the permission, even the obligation for weekly breaks, which then is also translated to daily breaks and annual holidays. What does it mean for the "data age" with *no time boundaries*? Accessibility to data is available 24/7/365 (above chapter 1.1), except unwanted break-downs of electricity or hacking-breaks. One of the great advantages of cyber-based data is that there is no limit in time and space. The disadvantage is that it puts a lot of pressure on humans who need breaks for sleep, eating, community sharing and refreshing energy. Robots are never tired and thus compete with humans.

*Becoming better humans* means to develop new ways of boundaries, breaks, limitations. A data Sunday? A data holiday? A data fasting? Abstain from round-the-clock availability; switch of the data-collection devices for a day; benefit from a Sunday or holiday without any electronic device; expand governmental Sunday legislations and regulations to employee's data protection; prevent data-burnout-syndromes and offering data-fasting-trainings. Churches in Germany and Switzerland offered internet-fasting as part of the fasting period during the annual Lenten period, forty days before Easter until Easter. Maybe it also exists for the Muslim Ramadan. Japan introduced "Internet Fasting Camps" for internet-addicted kids? Many schools, youth organisations and social services may also offer it.

With and without internet, cyberspace and data-related activities, the ethical task and chance remains: to become better humans.

# 2

# THE CORRELATION BETWEEN ETHICS AND TECHNOLOGY

*Peter G Kirchschläger, Switzerland*[32]

## 2.1 Introduction

As introduced above in the introduction, ethics is the scientific discipline that involves analysing moral questions and problems, discussing decisions and actions as legitimate or illegitimate, and analysing good/bad and right/wrong, respectively.

Originating from the Greek term technologia combining techne (art, technique) and logos (animating principle pervading the universe), technology entails the pursuit of a higher or more fundamental end or meaning. "When we speak of technology, we think of the power we are able to exert, thanks to our knowledge, on the world that surrounds us and of which our body, mortal and vulnerable, is a part." The following critique

---

of the term "technology" is taken into account and enriches the reflections in this book: "There is no such thing as technology in general or technology 4s such: There are always instead constellations of artefacts and techniques, actually-existing and also imagined. Technologies are in use, misused, consumed, consigned to landfill, repaired, maintained, creatively re-appropriated, under development, seeking venture capital, subject to regulation, promoted, marketed, misunderstood, repackaged, mistaken for novelties. Some artefacts and techniques seem to arrive from nowhere, dazzling us with promises and threats, usually to disappoint us soon enough, or at any rate until the next distraction fills the pop-tech press and screen."[33] Still, "technology" is used to categorise the diversity and plurality of technologies, while being aware of the limits of this categorization in order to use this awareness fruitfully in the more specific ethical assessment focusing on concrete technologies.

The relationship between ethics and technology can be understood as an interaction of a source of ends (ethics) with a source of means (technology). After assessing this possible element of the correlation between ethics and technology, the topic of what technology brings to ethics (e.g., innovation), and what ethics brings to technology (e.g., orientation) will be elaborated. On this basis, the challenges ethics creates for technology, as well as those which technology creates for ethics, will be discussed.

## 2.2 Ends and Means

### 2.2.1 Justifying Ends

"If there is one thing the great institutions of the modern world do not do, it is to provide meaning. Science tells us how but not why. Technology gives us power but cannot guide us as to how to use that power.

---

[33] Carrico, Dale (2013): "Futorological Discourses and Posthuman Terrain". In: Existenz. An International Journal in Philosophy, Religion, Politics, and Arts 8(2), 48.

The market gives us choices but leaves us uninstructed as to how to make those choices. The liberal democratic state gives us freedom to live as we choose but refuses, on principle, to guide us on how to choose."[34] Ethics could play the role of providing meaning. It could represent a source of ends – ends as, e.g., survival and life with dignity for all humans,[35] peaceful coexistence, justice, and sustainability.[36] Ethics should play a role "in the exclusion of objectives which are not strictly necessary. There remains enough of the indispensable to occupy the technical genius of man both in perfecting and in correcting and slowing down further developments."[37] This is also reflected in a statement by Wolfgang Wahlster, director of the German Research centre for Artificial Intelligence: "All we have to do is always put the benefit to humans at the centre of AI research."[38]

It is within the ethical dimension where questions about the ethical legitimacy of horizons of meaning and of ethical ends are discussed. Humans analyse what should be and why it should be. Humans evaluate, decide, and make an ethical judgement. It is part of ethics to act accord-

---

[34] Sacks, Jonathan (2015): "Sword Into Plowshares". In: The Wall Street Journal Review, October 3, C1-C2. Online: https://www.wsj.com/articles/swords-to-plowsh ares-unlikely-any-time-soon-1444422562 RonaldLethal Autonomous Systems [08.02.2021].

[35] See Kirchschlaeger, Peter G. (2013a): Wie können Menschenrechte begründet werden? Ein für religiöse und säkulare Menschenrechtskonzeptionen anschlussfähiger Ansatz. ReligionsRecht im Dialog 15. Muenster: LIT-Verlag. 194-195.

[36] See Armand, Jean-Louis (2012): "The bringing together of technology, sustainability and ethics". In: Sustainability Science 7(2), 113-116.

[37] Jonas, Hans (1984): "Warum wir heute eine Ethik der Selbstbeschränkung brauchen". In: Stoecker, Elisabeth (Ed.): Grenzen der Ethik. Ethik der Wissenschaften 1. Paderborn: Wilhelm Fink / Ferdinand Schoeningh, 86.

[38] Brost, Marc / Hamann, Goetz (2018): "Ein autonom fahrendes Auto erkennt bei Nacht kein Wildschwein". In: Die Zeit, July 26. Online: https://www.zeit.de/2018/31/kuenstliche-intelligenz-autonomes-fahren-wolfgang-wahlster-interview [08.02.2021]

ingly – to act ethically –[39]which, in other words, means to follow the ethically justified ends. "Ethics, whether in the form of issuing direct enjoinders to do or not to do certain things, or in the form of defining principles for such enjoinders, or in the form of establishing the ground of obligation for obeying such principles, deals with values relating to human conduct. It is concerned with action or non-action."[40]

Finally, ethics embraces the ethical justification of ethical ends, as well as the ethical justification of ends. Because of their claim of universality as part of ethics, these justifications need to satisfy the principle of generalizability by presenting rational and plausible arguments – "good reasons". The concept of "good reasons" means that it must be conceivable that all humans given their effective freedom and autonomy as well as their full equality would agree upon these reasons – within a model of thought and not within a real worldwide referendum – on ethical grounds. Inspiring the definition above was another way to frame these requirements of ethics: "A rational or critical ethics is one that claims for itself rational justifiability for its principles. Ethical principles are rationally justified if they are generally endorsed by, that is to say acceptable to, all affected persons, given their full equality and effective self-determination."[41]

### 2.2.2 Providing Means

In order to achieve these ethical ends or ethically justifiable ends, humans use adequate means. These means can also be produced by

---

[39] See Pieper, Annemarie (1994): Einführung in die Ethik. Tuebingen: UTB, 17-48.

[40] Armand, Jean-Louis (2012): "The bringing together of technology, sustainability and ethics". In: Sustainability Science 7(2), 114.

[41] Koller, Peter (1990): "Die Begründung von Rechten". In: Koller, Peter / Varga, Csaba/ Weinberger, Ota (Eds.): Theoretische Grundlagen der Rechtspolitik. Ungarisch-Österreichisches Symposium der internationalen Vereinigung für Rechts- und Sozialphilosophie, 75.

technology – technology based on basic science and applied science leading to growth and welfare[42] and embedded in the plan of living with nature (not so much in the sense of commanding nature)[43] for the benefit of humans.

Technology can serve ethically good ends. E.g., technology can be understood as a prosthesis – in the service of life and society.[44] "Continuing advances in scientific and technological innovations are essential to modern societies. Historically, such developments have improved living conditions in both developed and developing countries."[45] Out of this understanding of the correlation between technology and ethics, the following present challenge emerges: "It is a time when technology can bring wonders to one's life. It is a time when I ask myself whether all of these technological achievements have made us better human beings! A robot can be programmed to act good or bad, but who will demarcate human actions?"[46]

---

[42] See Bacon, Francis (2003): The Advancement of Learning. Kiernan, Michael (Ed.). The Oxford Francis Bacon 4. Oxford: Oxford Clarendon Press.

[43] See Bacon, Francis (2000): The Instauratio Magna: Last Writings. Rees, Graham (Ed.). The Oxford Francis Bacon 13. Oxford: Oxford Clarendon Press. Bacon, Francis (2004): The Instauratio Magna Part II: Novum Organum and Associated Texts. Rees, Graham / Wakely, Maria (Eds.). The Oxford Francis Bacon

11. Oxford: Oxford Clarendon Press. Bacon, Francis (2007): The Instauratio Magna Part III: Historia naturalis et experimentalis: Historia ventorum and Historia vitæ & mortis. Rees, Graham / Wakely, Maria (Eds.). The Oxford Francis Bacon 12. Oxford: Oxford Clarendon Press.

[44] See Schuurman, Egbert (2010): "Responsible Ethics for Global Technology". In: Axiomathes 20(1), 107-127. 107-127.

[45] Lucchi, Nicola (2016): The Impact of Science and Technology on the Rights of the Individual. Law, Governance and Technology 26. Cham: Springer 4.

[46] Bashir, Qasim (2000): "Technology vs. morality". In: Surgical Neurology 54(1), 92.

It would be too reductionist though to describe technology as the instrumental pursuit of an end with a means produced to reach this end.[47] Technology also produces the conditions for the success of instrumental pursuit of an end and continues working on these conditions. Nature is transformed by technology in an environment of technological systems for the benefit of humans who try to balance the resources and dangers of nature in calculable and assessable opportunities and risks.[48]

Furthermore, technology can also develop its own laws and transform from being automatic to being "autonomous"[49] Possible ends emerging out of technology could be, among others, "technology for the sake of technology", "what can be made must be made"[50] or "efficiency". The current digital transformation of society and the economy and the use of artificial intelligence[51] can serve as concrete examples for this potential of "autonomous" technology because at their centre are self-learning-systems[52] which no longer need any human input in order to

---

[47] See Ortega y Gasset, José (1949): Betrachtungen über die Technik. Stuttgart: Deutsche Verlags-Anstalt, 90-105; Hubig, Christoph (2007): Die Kunst des Möglichen II, Grundlinien einer dialektischen Philosophie der Technik 2: Ethik der Technik als provisorische Moral. Bielefeld: Transcript-Verlag, 48.

[48] See Hubig, Christoph (2011): "Technikethik". In: Stoecker, Ralf / Neuhaeuser, Christian/ Raters, Marie-Luise (Eds.): Handbuch Angewandte Ethik. Stuttgart: J. B. Metzler, 170.

[49] See Ellul, Jacques (1964): The Technological Society. New York: Vintage Books, 125-150; Kirchschlaeger, Peter G. (2016b): "Digitalisierung und Robotisierung der Gesellschaft aus ethischer Perspektive". In: feinschwarz.net, March 30.   Online:http://www.feinschwarz.net/digitalisierung-und-robotisierung-der-gesellschaft-aus-theologisch-ethischer-perspektive/ [08.02.2021].

[50] Schuurman. 2010, op. cit., 123.

[51] See Kirchschlaeger 2016b, op. cit.

[52] See Taddy, Matt (2019): "The Technological Elements of Artificial Intelligence". In: Agrawal, Ajay / Gans, Joshua / Goldfarb, Avi (Eds.): The Economics of Artificial Intelligence: An Agenda. Chicago: University of Chicago Press, 61-87.; Bishop, Christopher M. (2006): Pattern Recognition and Machine Learning. Information Science and Statistics 1. New York: Springer; Goodfellow, Ian /

improve and to optimize their own performance. "The ethical risks posed by AI-enabled robots are (…) serious – especially since self-learning systems behave in ways that cannot always be anticipated or fully understood, even by their programmers."[53]

This self-learning can consist also of "intuition" as, e.g., the win of the Go-game by robots also based on "intuitive" decisions can show;[54] of strategic reasoning with imperfect information as, e.g., the win of a marathon 20-day poker competition can demonstrate;[55] or of "a system's ability to correctly interpret external data, to learn from such data, and to use those learnings to achieve specific goals and tasks through flexible adaptation."[56]

---

Bengio, Yoshua / Courville, Aaron (2017): Deep Learning. Adaptive Computation and Machine Learning. Cambridge: MIT Press; Agrawal, Ajay / Gans, Joshua S. / Goldfarb, Avi (2018): Prediction, Judgment and Complexity: A Theory of Decision Making and Artificial Intelligence. Rotman School of Management. Working Paper 3103156. Online: https://ssrn.com/abstract=3103156 [08.02.2021]; Le Cun, Yann / Bengio, Yoshua / Hinton, Geoffrey (2015): "Deep learning". In: Nature 521(7553), 436-444; Jain, Sarthak (2017): "NanoNets: How to use deep learning when you have limited data: Part 2: Building object detection models with almost no hardware". In: Medium, January 30. Online: https://medium.com/nanonets/nanonets-how-to-use-deep-learning-when-you-have-limited-data-f68c0b512cab [08.02.2021].

[53] Vallor, Shannon / Bekey, George A. (2017): "Artificial Intelligence and the Ethics of Self-Learning". In: Lin, Patrick / Jenkins, Ryan / Abney, Keith (Eds.): Robot Ethics 2.0 from Autonomous cars to Artificial Intelligence. New York: Oxford University Press, 338.

[54] See Nature 2016; Kirchschlaeger, Peter G. (2017a): "Roboter und Ethik". In: Aktuelle Juristische Praxis 26(2), 240-249.

[55] See Spice, Byron (2017): "Carnegie Mellon Artificial Intelligence Beats Top Poker Pros". In: Carnegie Mellon University, January 31. Online: https://www.cmu.edu/news/stories/archives/2017/january/AI-beats-poker-pros.html [08.02.2021].

[56] Kaplan, Andreas / Haenlein, Michael (2018): "Siri, Siri in my Hand, who's the

If the pursuit of efficiency is the exclusive scope of technology, three observations are provoked: *First*, technology striving for efficiency is neither independent from any ends or an end itself, nor ethically neutral, but serves a clear end: efficiency. Or the affirmation that technology is not a means which serves any ethical end, is in itself a normative statement assuming no end or technology as an end for technology – a normative statement which needs to be justified ethically, subject to the required criterion mentioned above.

*Second*, technology inherits a social and an ethical dimension. "Technology cannot be understood as an instrument for bringing about goals that are external to the contexts in which it operates, but the relational contexts in which technology functions are imbued with values which demand consideration. Thus, technology, as it actually operates in concrete situations has a contextually dependent ethical quality. Technology creates an ethical situation, and this situation should provide the context for decision making." [57]

*Third*, the pursuit of efficiency as the exclusive scope of technology leaves technology without any raison d'être. "The process of technological development hurtles ahead blindly without a normal sense of balance. As we can see from nuclear weapons and environmental degradation, the effects on human society are immense."[58] The way it is advancing, technology runs the risk of detaching itself from any horizon of meaning and of alienating itself from humanity and nature. "The human artifice of the world separates human existence from all mere animal environments, but life itself is outside this artificial world, and through life man remains related to all other living organisms. For some time

---

Fairest in the Land? On the Interpretations, Illustrations and Implications of Artificial Intelligence". In: Business Horizons 62(1), 17.

[57] Buchholz, Rogene A. / Rosenthal, Sandra B. (2002): "Technology and Business: Rethinking the Ethical Dilemma". In: Journal of Business Ethics 41(1), 48.

[58] Shibasaki, Fumikazu (2005): "Technology and ethics". In: Philosophy & Criticism 31(18), 497. 2005.

now, a great many scientific endeavours have been directed toward making life also 'artificial', toward cutting the last tie through which even man belongs among the children of nature (…) The question is only whether we wish to use our new scientific and technical knowledge in this direction, and this question cannot be decided by scientific means; it is a political question of the first order and therefore can hardly be left to the decision of professional scientists or professional politicians."[59] In other words, the raison d'être for technology is not technology-based; it cannot be found out, defined, and justified by technology.[60]

Beyond that, technology can also be abused for other ends,[61] distancing itself from its original ends, striving exclusively for efficiency and disregarding any higher end or meaning.[62]

However, ethical ends would still be distinguished from other ends. These other ends would undergo an ethical judgement on the basis of ethical ends. Ernesto Bertarelli, former CEO of Serono, states: "We never carry out research for the sake of research. (…) if there is no purpose and if there is no reason, we reject the innovation and we do not pursue it."[63] The former point emphasises the orientation towards another end of research, the latter shows that the characteristics of innovation require an ethical assessment of innovations.[64]

---

[59] Arendt, Hannah (1958): The human condition. Chicago: University of Chicago Press, 2-3.

[60] See Jennings, Bruce (2010): "Enlightenment and enchantment: Technology and ethical limits". In: Technology in Society 32(1), 27.

[61] See European Group on Ethics in Science and New Technologies to the European Commission 2014.

[62] See Shibasaki 2005, op. cit.

[63] Bertarelli, Ernesto (2002): "Science Between Goals and Limits". In: 32nd ISC-Symposium at the University of St. Gallen (Ed.): Pushing Limits – Questioning Goals. St Gallen: ISC-Symposium, 49-50.

[64] See Kirchschlaeger, Peter G. (2013b). "Human Rights as an Ethical Basis for Science". In: Journal of Law, Information and Science 22(2), 1-17.

One should overcome the naïve assumption that all technological progress and every innovation is an ethically good innovation. "Given the immense ambiguities of innovations – in themselves and in their consequences, the ethical scrutiny of innovation is a dictate of reason that should not be ignored any longer."[65]

It needs to be added though that technological progress depends also on basic research, which is done for the sake of basic research. Even in this area, a decision – which can find its foundation in ethical ends – needs to be taken in which basic research is prioritised and gets funded. Facing the scarcity of financial resources for research and technology and at the same time confronted by pressing global problems, a setting of focus and priority seems to be necessary. Finally, the question must be addressed as to who is benefitting from results and successes of technology – and for that matter benefitting the most.[66]

---

[65] Enderle, Georges (2015b): "The Theme of the Sixth World Congress of the International Society of Business, Economics, and Ethics in 2016 in Shanghai, China: 'Ethics, Innovation, and Well-Being in Business and the Economy'". In: ISBEE, July 22. Online: http://isbee.org/the-theme-of-the-sixth-world-congress-of-the-international-society-of-business-economics-and-ethics-in-2016-in-shanghai-china-ethics-innovation-and-well-being-in-business-and-the-econom/ [08.02.2021].

[66] See Hunt, Paul (2008): Promotion and protection of human rights: human rights questions, including alternative approaches for improving the effective enjoyment of human rights and fundamental freedoms: The right to health. Report of the Special Rapporteur on the right of everyone to the enjoyment of the highest attainable standard of physical and mental health. UN Human Rights Council. Sixty-third Session. A/63/263. United Nations. Online: https://undocs.org/A/63/263 [08.02.2021]; Shaver, Lea (2015): "The Right to Science: Ensuring that Everyone Benefits from Scientific and Technological Progress". In: European Journal of Human Rights 4, 411-430; Donders, Yvonne (2015): "Balancing Interests: Limitations to the Right to Enjoy the Benefits of Scientific Progress and Its Applications". In: European Journal of Human Rights 4, 486-503; Chapman, Audrey R. (2009): "Towards an Understanding of the

### 2.2.3 Way of Functioning of Technological Progress

Looking at the correlation between ethical, as well as ethically justifiable, ends and means provided by technology in a context partly created by technology, a further element needs to be taken into account. The perception of technological development and progress as a linear process pursuing a well-defined scope would probably not correspond with the present-day theory and reality of technology.[67] Technological innovations are often the result of small steps and are frequently random products.[68] "Technology is not ordinarily developed after carefully considering the various possible ramifications. In most cases a new technology is developed because it promises major short-term benefits and is judged not to cause any immediate problems."[69] In addition, the speed of technological advancement outpacing normative considerations is another characteristic of the way of functioning of technology.

Furthermore, some ethical norms exist by dint of certain technological developments because the necessity to establish an ethical rule originates in a technology-based reality.

Beyond that, the complexity of technological development and progress should not be underestimated. "First, engineering and technology development typically take place in collective settings, in which a lot of different agents, apart from the engineers involved, eventually shape the technology developed and its social consequences. Second, engineering and technology development are complex processes, which are charac-

---

Right to Enjoy the Benefits of Scientific Progress and Its Applications". In: Journal of Human Rights 8(1), 1-36.

[67] See Kuhn, Thomas S. (1962): The Structure of Scientific Revolutions. Chicago: University of Chicago Press.

[68] See Boutellier, Roman / Heinzen, Mareike / Raus, Marta (2010): "Paradigms, Science, and Technology: The Case of E-Customs". In: Becker, S. Ann / Niebuhr, Robert E. (Eds.): Cases on Technology Innovation: Entrepreneurial Successes and Pitfalls. New York: Business Science Reference, 134-155.

[69] Shibasaki 2005, op. cit.

terised by long causal chains between the actions of engineers and scientists and the eventual effects that raise ethical concern. Third, social consequences of technology are often hard to predict beforehand."[70]

It is also noteworthy that the social consequences of technology are such that should a certain technology fill some social role, it inherits, by so doing, the concomitant "social requirements".[71] Furthermore, technological action should be considered distributed and collective rather than individual[72] though without introducing categories like fate and tragedy in order not to be too harsh with technology.[73] Therefore, various actors should be identified as subjects of responsibility.[74] Identifying them is a complex task. The identification of subjects of responsibility should still be implemented in order to build an atmosphere of professionalism and accountability – not only out of respect for the objects of responsibility.[75] Complexity cannot serve as an excuse to neglect legal or ethical obligations and responsibilities because ethical and legal norms keep their validity even in complex situations and contexts.

---

[70] Doorn, Neelke / van de Poel, Ibo (2012): "Editors' Overview: Moral Responsibility in Technology and Engineering". In: Science and Engineering Ethics 18(1), 2.

[71] See Bostrom, Nick / Yudkowsky, Eliezer (2014): "The ethics of artificial intelligence". In: Ramsey, William, M. / Frankish, Keith (Eds.): The Cambridge Handbook of Artificial Intelligence. Cambridge: Cambridge University Press, 316-334.

[72] See Lenk, Hans / Maring, Matthias (2001): "Responsibility and Technology". In: Auhagen, Ann Elisabeth / Bierhoff, Hans-Werner (Eds.): Responsibility: The many faces of a social phenomenon. London: Routledge, 100.

[73] See Coeckelbergh, Mark (2012): Growing Moral Relations: Critique of Moral Status Ascription. New York: Palgrave Macmillan. 35-48.

[74] See Coeckelbergh, Mark / Wackers, Ger (2007): "Imagination, distributed responsibility, and vulnerability: The case of Snorre A". In: Science and Engineering Ethics 13(2), 235-248.

[75] See Lenk, Hans / Maring, Matthias (2001): "Responsibility and Technology". In: Auhagen, Ann Elisabeth / Bierhoff, Hans-Werner (Eds.): Responsibility: The many faces of a social phenomenon. London: Routledge, 93-107.

By overcoming the overly simplistic ends-means-framework and by considering the characteristics of technology, on the one hand, and respecting ethics as a source of ethical ends (and of the ethical legitimacy of ends), on the other hand, perhaps an attempt to grasp the reciprocal interactions and the reciprocal challenges could inform the understanding of the relationship between ethics and technology.

## 2.3 Reciprocal Interactions

The correlation of ethics and technology can be understood based on reciprocal interaction, as both ethics and technology contribute to each other. For example, groundbreaking ideas in technology and their successful application have a concrete impact on ethics as technology creates value, solutions for societal challenges, and innovation. "Science and technology have shaped modern society, economics, politics, law and culture. They deeply affect the lives of all people and they are now central features of our social and commercial landscape."[76] Even specifically in the ethical dimension, technology leads to innovation and dynamics[77] because the societal and individual transformation based on technology needs to be taken into account in ethics as well.[78]

---

[76] Lucchi, Nicola (2016): The Impact of Science and Technology on the Rights of the Individual. Law, Governance and Technology 26. Cham: Springer, 6.

[77] See Manzeschke, Arne (2015): "Angewandte Ethik organisieren: MEESTAR – ein Modell zur ethischen Deliberation in sozio-technischen Arrangements". In: Maring, Matthias (Ed.): Vom Praktischwerden der Ethik in interdisziplinärer Sicht: Ansätze und Beispiele der Institutionalisierung, Konkretisierung und Implementierung der Ethik. Zentrum für Technik- und Wirtschaftsethik 7. Karlsruhe: KIT Scientific Publishing, 315-330.

[78] See Kernaghan, Kenneth (2014): "Digital dilemmas: Values, ethics and information technology". In: Canadian Public Administration 57(2), 295-317.

Moreover, technology[79] and especially smart technology is influencing (e.g., by nudging)[80] at least individual lives – and maybe even the ethical dimension of individual lives.[81] "By helping to shape human actions and experiences, technologies also participate in our ways of doing ethics."[82] While acknowledging the effect technology has on humans as subjects of ethics, technology – in contrast to the idea of a "non-humanist ethics of technology[83] – remains an object of ethics and humans subjects of ethics. Why? Humans can decide on an ethical basis or orientation towards an ethical frame of reference in the technological decisions they make: whether to create a technology or not, how and if they design, develop, produce, disseminate, and use a technology or not, if they abolish a technology or not, if – to a certain extent – they allow a technology to influence their lives or not, and how they assess a technology from an ethical perspective. All this is up to humans, not up to technology.

Finally – and linked with that – contributions by technology to human lives are ethically relevant and can be ethically evaluated.[84] The ethical relevance of "human making and using"[85] of "tools and their

---

[79] See Borgmann, Albert (1995): "The Moral Significance of Material Culture". In: Feenberg, Andrew / Hannay, Alistair (Eds.): Technology and the Politics of Knowledge. Bloomington: Indiana University Press, 85-93.

[80] See Mathis, Klaus / Tor, Avishalom (Eds.) (2016): 'Nudging' – Possibilities, Limitations and Applications in European Law and Economics. Cham: Springer.

[81] See Guthrie, Clifton F. (2013): "Smart Technology and the Ethical Life". In: Ethics & Behavior 23(4), 324-337.

[82] Verbeek, Peter-Paul (2011a): Moralizing Technology: Understanding and Designing the Morality of Things. Chicago: Chicago University Press, 1.

[83] See among others Verbeek 2011a: 21-40, ibid.

[84] See Verbeek 2011a: 1-20, ibid.

[85] Mitcham, Carl (2014): "Agency in Humans and in Artifacts: A Contested Discourse". In: Kroes, Peter / Verbeek, Peter-Paul (Eds.): The Moral Status of Technical Artefacts. Philosophy of Engineering and Technology 17. Dordrecht: Springer, 11-29.

deliberate use"[86] as the foundation of civilization, though, needs to be distinguished from possessing ethical subjectivity or from moral agency. For example, a self-driving car is ethically relevant because it can save human lives or cause harm to/kill humans. But a self-driving car as a material object – even though it represents a highly sophisticated technology – is not an ethical subject or a moral agent because – among other things (see for further arguments below chapter 3 Can Ethical judgement Be Delegated to Technologies?) – it cannot be held accountable for harm or killing. Humans behind this technology remain the ethical subjects or the moral agents. They make decisions about the interactivity (response to stimulus by change of state) – "autonomy" (ability to change without stimulus) and adaptability (ability to change the "transition rules" by which the state is changed) of technology; they create, design, develop, produce, disseminate, abolish, and use them and the technology itself. Behind the supposedly "political bridges"[87] behind the supposedly "missing masses of morality"[88] , behind the supposedly "artificial moral agents"[89] , behind the supposed "artefacts with morality"[90] , behind the supposedly "moral agents and mediated subjects"[91], and behind the supposed "morality in design"[92] there are humans as ethical subjects and moral agents. Humans can decide whether or not to

---

[86] Bloch, Ernst (1959): Das Prinzip Hoffnung: Kapitel 1-32. Werkausgabe 5. Frankfurt am Main: Suhrkamp. 731.

[87] See Winner, Langdon (1980): "Do Artifacts Have Politics?" In: Daedalus 109(1), 121-136.

[88] See Latour, Bruno (1992): "Where Are the Missing Masses? The Sociology of a Few Mundane Artifacts". In: Bijker, Wiebe E. / Law, John (Eds.): Shaping Technology/ Building Society. Cambridge: MIT Press, 151-180.

[89] See Floridi, Luciano / Sanders, Jeff W. (2004a): "On the Morality of Artificial Agents". In: Minds and Machines 14(3), 349-379.

[90] See Verbeek, Peter-Paul (2011a): Moralizing Technology: Understanding and Designing the Morality of Things. Chicago: Chicago University Press, 41-65.

[91] See Verbeek 2011a: 66-89, op. cit.

[92] See Verbeek 2011a: 90-119.

build bridges in a way that is disadvantageous to humans – manifesting their political convictions in a lasting way. Humans can decide, for example, whether or not to build speed bumps in order to protect human life. Humans can decide to create, design, develop, build, produce, disseminate, to use or not to use a technology which can perform a "morally qualifiable action". "An action is said to be morally qualifiable if and only if it can cause moral good or evil"[93]. Humans can decide if a technology should perform this action, and if a technology should be able to decide by itself if it should perform this action. Just because a technology is ethically relevant, meaning it can cause ethically positive and ethically negative effects, does not mean it possesses moral agency as there is, first, a difference between performance and moral agency. Second, it is an ethical decision by humans if a technology should be entrusted with such decisions (as will be further elaborated below in chapter 3 (Can Ethical judgement be Delegated to Technologies?).

Third, even in the case of a technology being designed by humans in a way that it makes "autonomous" decisions about ethically relevant questions, a technology does not make these decisions with a knowledge, perception, or awareness of the ethical quality of these decisions (this point will be further explained below in chapter 5 The Complexity of Ethics).

Humans can make decisions about mediation by technology, its intensity, and its extent. Humans can decide to create, design, develop, build, produce, disseminate, use or not use a technology which may be able to mediate, and humans can decide if a technology should mediate or not. Humans can be aware of mediation by technology, and humans can reflect self-critically upon mediation by technology. Humans are not exposed powerlessly and heteronomously to it.

Humans can make decisions about mediation by design, its intensity and its extent. Humans can decide to create a design which mediates or

---

[93] Floridi / Sanders 2004a: 361.

does not mediate, and humans can decide if a design should mediate or not. Humans can be aware of mediation by design, and humans can reflect self-critically upon mediation by design. Again, humans are not exposed powerlessly and heteronomously to it.

In order to avoid a potential misunderstanding of what was just elaborated above, an emphasis must be put on the agreement that these societal impacts, effects, actions, influence, shaping, and mediation are possible and ethically highly relevant. There is no doubt about that. For example, design can influence the ethical dimension of human lives – illustrated by the inspiring example of "Eternally Yours".[94] "Eternally Yours follows an unorthodox approach within eco-design. Instead of the usual emphasis on reducing pollution while maintaining beauty and economy, the company focuses on lengthening what it calls the product's 'psychological lifetime'. Most products are thrown away long before they are broken or obsolete, usually because of changing tastes and fashions. Eternally Yours attempts to combat this tendency of products to wind up prematurely in the landfill by designing products that invite people to become attached to them."[95] The argument is presented, though, that behind societal impacts, effects, actions, influence, shaping, and mediation by technology and design, humans are the ethical subjects and the moral agents.

At the same time, science and technology happen, are pursued, and are in an ethically informed context. "But even the most solitary and radical reflection, as thinking with an intersubjective claim to validity, must presuppose language and thus a community of communication. But this also marks the limit of the value-free world-distance of modern natural science. As an enterprise in the search for truth in the dimension of intersubjectivity, value-free natural science must also presuppose

---

[94] See Verbeek, Peter-Paul (2005): What Things Do: Philosophical Reflections on Technology, Agency, and Design. Pennsylvania: Penn State University Press. 203-234.

[95] Verbeek 2005b: 12.

ethics. But that would, of course, only be an ethics of the enterprise science, which could not even answer the question of whether science should be. It would be different if we were to reflect on the subjective and intersubjective conditions of the possibility of thinking as argumentation; for no one can go back behind these conditions if he or she seriously raises any questions at all and thus puts them up for discussion."[96]

Moreover, ethics contributes to technology, for example, by stimulating technological innovation,[97] by recognizing technological inventions,[98] and by providing ethical guidance. "Since all technologies are expressions of the values of their makers, if we care about ethics and morality, it will show in the machines we build."[99] Part of this value-system from an ethical perspective should be, among other things, ecological concerns. "Industrial Modernity – the greatest innovative boost to human development since the invention of letterpress printing – had one major design flaw: it passed its environmental costs on to the future generations. So a central task of the digital revolution might logically be to reverse this error. But instead, this revolution is being used to fuel the exponential growth of mass consumption. Similarly, it has failed to more fairly distribute the benefits of technological innovation. Here too everything points to an intensification of social inequality."[100]

---

[96] Apel, Karl-Otto (1986): "Verantwortung heute – nur noch Prinzip der Bewahrung und Selbstbeschränkung oder immer noch der Befreiung und Verwirklichung von Humanität". In: Meyer, Thomas / Miller, Susanne (Eds.): Zukunftsethik und Industriegesellschaft. Zukunftsethik 1. Muenchen: J. Schweitzer, 27-28.

[97] See Lucchi 2016: 7, op. cit.

[98] See Lucchi 2016: 1-2.

[99] Sullins, John P. (2013): "An Ethical Analysis of the Case for Robotic Weapons Arms Control". In: Podins, Karls / Stinissen, Jan / Maybaum, Markus (Eds.): 5th International Conference on Cyber Conflict. Tallinn: NATO CCD COE Publications, 16.

[100] Thun-Hohenstein, Christoph (2017): "Sense and Sensibility in the Digital Age: Let Us Wake Up and Take Action! Introductory Essay". In: Thun-

Ethical discourse of technology depends on the understanding that technology is "something made" and "not anything given"[101] "Technology must be allowed to augment living where it can, but cannot be allowed arbitrarily to suppress ways of life. Thus, ethics does not presume against technological change, but must be a part of the formulation of how change is translated into advancement or repression."[102]

One needs to go even further stating that ethics belongs to technology. "The idea of scientific knowledge as value-neutral is simply incorrect. Values are intrinsic to the making of science and technology, and they both reflect and transform particular values."[103] Horizons of meaning and ethical ends inform technology in an ethical sense. "Science and technology in their objective areas, for all their diversity of methods, are ultimately always aimed at people. That is why science and technology cannot be separated from ethical insight and decision-making in any area."[104] The discussion about the difference between notions like "ethical technologies", "technical tools"[105] , "intelligent computer interfac-

---

Hohenstein, Christoph (Ed.): Vienna Biennale 2017. Wien: Verlag für Moderne Kunst, 19.

[101] Heesen, Jessica (2014): "Mensch und Technik. Ethische Aspekte einer Handlungspartnerschaft zwischen Personen und Robotern". In: Hilgendorf, Eric (Ed.): Robotik im Kontext von Recht und Moral. Robotik und Recht 3. Baden-Baden: Nomos Verlagsgesellschaft, 253; 268.

[102] Rainey, Stephen / Goujon, Phillipe (2011): "Toward a normative ethics for technology development". In: Journal of Information, Communication and Ethics in Society 9(3), 174.

[103] De Melo-Martín, Immaculada (2010): "The Two Cultures: An introduction and assessment". In: Technology in Society 32(1), 9.

[104] Auer, Alfons (1982): "Darf der Mensch, was er kann?". In: Busch, Alois J. / Splett, Joerg (Eds.): Wissenschaft – Technik – Humanität. Frankfurt am Main: Verlag Josef Knecht, 18.

[105] Engineering and Physical Sciences Research Council 2011.

es"[106] , or "socio-technical systems"[107] shows the openness of technology to ethics. At the end of the day, this discussion implies the link between technology and ethics.

Beyond that, while the technology community is aware of its legal obligations and legal compliance standards, it strives for the respect of ethical principles in its work as well, e.g., honesty, objectivity, independence, impartiality, fairness, responsibility for future generations.

At the same time, a globalised technology community faces several traditions, cultures, religions, worldviews, and value-systems that can lead to ethical challenges. Ethical guidance can support technology in overcoming these challenges and in benefitting from the opportunities of this diversity and heterogeneity.

Furthermore, ethics can critically examine the legal obligations and legal compliance standards of the technology-community on a regular basis. This should lead to a continuous optimization of the legal framework for technology.

In addition, ethics can help in the process of agenda-setting in technology, not only in defining the right priorities but also in framing adequately the sphere of influence and responsibility of technology.

Finally, while technology contributes to the progress of ethics, it is obvious that at the same time there is a need for ethics in technology in order to be able to even conduct the necessary research, discussions, and studies. Technology can be the victim of infringements of its freedom, of attempts to block innovative and creative approaches, and of oppression of ideas, concepts, and discoveries. Reasons for these transgressions can be putative "absolute truths" or the enforcement of old and existing

---

[106] Van Est, Rinie / Stemerding, Dirk (2012): Making Perfect Life. European Governance Challenges in 21st Century Bioengineering. STOA Science and Technology Options Assessment. Brussels: European Union.

[107] Manzeschke, Arne / Weber, Karsten / Rother, Elisabeth / Fangerau, Heiner (2013): Ethische Fragen im Bereich Altersgerechter Assistenzsysteme: Ergebnisse der Studie. Ludwigsfelde: Druckerei Thiel Gruppe.

economic or political power structures. There is still a danger of members of the technology community not being able to conduct their research freely and independently. Therefore, there is a need for legal and ethical norms that support and protect technological progress.

This reciprocal interaction between ethics and technology accommodates the "interactionist model" highlighting the contributions of technology to ethics and the contributions of ethics to technology.

## 2.4 Reciprocal Challenges

At the same time, ethics can limit technology as well. For example, health and safety guidelines, patents, legal ownership of intellectual property rights, competition policy, consumer protection, and ethical codes of conduct belong to this category. This impact of ethics can be perceived as blocking and hindering technological innovation.[108]

Beyond that, technology must respect ethical principles. For example, "developers should strive at creating artificial agents whose actions are constrained in such a way that unethical outcomes can be avoided."[109] Among other things, the dignity of all humans can be a limit for technology (see below sub-chapter 6.4 Human Rights as an Ethical Frame of Reference). Therefore, technology does not have the permission to treat humans as means but only as ends – e.g., on the basis of the categorical imperative by Immanuel Kant.[110] Furthermore, technology

---

[108] Gurkaynak, Gonenc / Yilmaz, Gonenc / Haksever, Gunes (2016): "Stifling Artificial Intelligence: Human Perils". In: Computer Law & Security Review 32(5), 749-758.

[109] Krenn, Brigitte (2016): "Multiuse Tool and Ethical Agent". In: Trappl, Robert (Ed.): A Construction Manual for Ethical Systems. Cham: Springer, 25.

[110] See Kant, Immanuel (1974): Grundlegung zur Metaphysik der Sitten. Weischedel, Wilhelm (Ed.). Werkausgabe 7. Frankfurt am Main: Suhrkamp, 61; Duewell, Marcus (2010): "Menschenwürde als Grundlage der Menschenrechte". In: Debus, Tessa / Kreide, Regina / Krennerich, Michael / Malowitz, Karsten / Pollmann, Arnd / Zwingel, Susanne (Eds.): Zeitschrift für Menschenrechte 4, 77.

must also respect the privacy of all humans (see below sub-chapter 7.17 Data-Protection and Privacy). In this way, technology is challenged by ethical norms.[111] Paradigmatically, in the area of digitalization, automation, mechanisation, robotization, and the use of artificial intelligence, technology can follow ethical programming by humans heteronomously.[112]

At the same time, ethics faces challenges from technology as well. Technological progress is speeding up. The intervals for new technologies and technological applications get smaller and smaller. Ethics and law run the risk of being constantly outpaced by technology. They struggle to keep up with technological progress. The perception of this risk is based, though, on two misconceptions: First, it understands ethics as reactive instead of proactive, and secondly, "some observers of modern culture do note disparities and tensions between contemporary doing and making, between ethics and technology. They complain that, while technology has advanced dramatically, our ethical attitude in dealing with it has not. But this complaint is as radically mistaken as the general divorce of doing from making. It fails to see that a technological accomplishment, the development and adoption of a technological device always and already constitutes an ethical decision."[113]

---

[111] About technological progress, its ends, its foundational values, its societal importance, and its limits, using the examples of stem cell research and of research on human beings, see Kirchschlaeger et al. 2003; Kirchschlaeger et al. 2005.

[112] See Wallach, Wendell / Allen, Collin (2009): Moral Machines: Teaching Robots Right from Wrong. Oxford: Oxford University Press; Kirchschlaeger 2017b.

[113] Borgman, Albert (1992): Crossing the Postmodern Divide. Chicago: University of Chicago, 110.

Beyond that, ethics is challenged more and more not only by human curiosity striving for new inventions and solutions but linked by substantial economic interests and power.[114]

For example, in the area of digitalization, automation, mechanisation, robotization, and the use of artificial intelligence, ethics must deal with an attitude that the legal system of a nation-state is violated as long as the economic benefits of these acts are higher than the sanctions. The defence of the ethically justifiable position that not everything, which is doable, is ethically good meets the opposition of potential benefits and economic incentives. Similar pragmatic patterns of argumentation from a perspective of technology dominate the discussion, e.g.: that closing the gate on technology is not an option at all; that the implementation of limiting technological advancement with legal and ethical norms is impossible in a globalised world; that the identification of the subjects of responsibility in the area of technology is too complex; and that the risk assessment of technology remains imprecise and ineffective.

Finally, the impression emerges that "so far, the normative regulations have merely adapted to social upheavals. The social changes triggered by technical innovations in the fields of production and exchange, communication and transport, military and health have always been at the forefront. Classical social theory has described post-traditional legal and ethical concepts as the result of cultural and social rationalisation, which has taken place in parallel with the progress of modern science and technology. Institutionalised research is regarded as the motor of these advances. From the perspective of the liberal constitutional state, the autonomy of research deserves protection. The growing scope and depth of technological access to nature is combined with both the economic promise of productivity gains and prosperity gains and the political prospect of greater individual freedom of decision. Because growing

---

[114] See Buchholz, Rogene A. / Rosenthal, Sandra B. (2002): "Technology and Business: Rethinking the Ethical Dilemma". In: Journal of Business Ethics 41(1), 45-50.

freedom of choice promotes the private autonomy of the individual, science and technology have so far been in an informal alliance with the liberal basic idea that all citizens should have the same opportunity to shape their own lives autonomously. (...) The desire for autonomous living is always linked to the collective goals of health and life prolongation. The medico-historical view, therefore, warns against attempts to 'ethicalise human nature' to scepticism (...) From this empirically sobering perspective, legislative interventions in the freedom of biological research and genetic engineering development appear as futile attempts to resist the dominant freedom tendency of social modernity. (...) Of course, a completely different picture emerges if one understands the 'ethicalisation of human nature' in the sense of the self-assertion of a genre-ethical self-understanding, on which it depends whether we will continue to see ourselves as undivided authors of our life history and whether we can recognize each other as autonomously acting persons."[115]

While reciprocal challenges between technology and ethics can obviously arise, the responsibility[116] of humans is also growing due to the constantly expanding creation of an artificial world and of "a technological simulacrum of natural life"[117] and the corresponding power and influence of humans. What do humans have to do to live up to this responsibility? How should one deal with the above-mentioned reciprocal challenges? Is "anything goes" the solution to this situation? There is a demand for ethical guidance, which can be provided by ethical principles and ethical points of reference discussed below.

---

[115] Habermas, Juergen (2001a): Die Zukunft der menschlichen Natur. Auf dem Weg zu einer liberalen Eugenik? Frankfurt am Main: Suhrkamp. 47-49.

[116] See Kirchschlaeger 2014a.

[117] Jennings, Bruce (2010): "Enlightenment and enchantment: Technology and ethical limits". In: Technology in Society 32(1), 26.

## 2.5 Ethical Principles and Ethical Points of Reference for Technology

Ethical orientation in ethical questions, issues, and problems concerning technology can have their origin in traditions, cultures, religions, world-views, and philosophies. In a globalised world, though, technology can, on the one hand, affect the entire planet, and most on the other hand provoke an impact on all humans. Therefore, possible ethical guidance for technology gains relevance if it can claim universality and is independent from a specific tradition, culture, religion, worldview, and philosophy.

These ethical principles and ethical points of reference enrich the relation between humans and technology consisting of an "embodiment relation" (technologies as extensions of the body, e.g., glasses, hearing aids), a "hermeneutic relation" (technologies as access points to the world, e.g. thermometer), an "alterity relation" (technologies and humans interact, e.g. humans operate a computer), a "background relation" (technologies have an effect on humans and their relationship with the world without being noticed)[118] , "immersion" (technologies merge with the world and interact with humans), and "augmentation" (technologies offer a representation of the world)[119] by expanding this set of relations with a "constructivist and deconstructivist relation" (humans create and destroy technologies), an "interplaying relation" (humans use technologies and technologies use humans, e.g. their data), an "interproductive relation" (humans create technology/technologies and technology/technologies create parts of or in humans like attempts to simulate peculiarly human thought, e.g., artificial intelligence), an "ethically guiding relation" (humans decide based on ethical principles, ethical points of reference, and ethical norms which technologies they design,

---

[118] See Ihde, Don (1990): Technology and the Lifeworld: From Garden to Earth. Bloomington: Indiana University Press.

[119] See Verbeek 2015a: 211-212.

produce, and use, and which they do not), and an "ethically shaping relation" (humans decide based on ethical principles, ethical points of reference, and ethical norms how they design, produce, and use technologies).

This expansion might provoke at first sight the criticism that it takes ethical reasoning back to a point where it was before and got liberated from, namely perceiving the relation between humans and technology as a conflict between the two. This is not the intention. Rather, it tries to reconcile Michel Foucault's "critique"[120] with the subjectivity of humans distinguished from the objectivity of technologies by acknowledging "the fundamental intertwinement"[121] of humans and technologies implying "that the frameworks from which we criticise technologies are always mediated by these technologies themselves"[122] , and by re-empowering the critical characteristics of ethics, embracing also the possibility of questioning fundamentally technology-based innovation, technology, and technologies as such. Both ethics from inside and from outside technologies are ethically necessary, and humans are able to strive for both – on the conceptual basis of understanding the relationship between technology/technologies and ethics as interaction. The ethical foundation for this reconciliation embraces, on the one hand, the self-critical acceptance of the boundaries of humans and human reason that makes it impossible for humans to liberate themselves completely from the intertwinement of humans and technology. On the other hand, it trusts in the subjectivity of humans allowing them rationally to distance themselves enough from technology in order to analyse, discuss, evaluate, assess, and judge the object "technology" from an ethical perspective. Without this latter aspect, humans would run the risk to being trapped again in "immaturity", accepting authorities and following their

---

[120] See Foucault, Michael (1984): "What is Enlightenment?". In: Rabinow, Paul (Ed.): The Foucault Reader. New York: Pantheon Books, 32-50.

[121] Verbeek 2015a: 215.

[122] Verbeek 2015a: 215

orders as criticised by Immanuel Kant – in present times, obeying imma-
turely the authorities of technological progress, innovation, and econom-
ic imperatives. Without this latter aspect, humans would be reduced to
immature and powerless addressees, users, and consumers of powerful
technology[123] Without this latter aspect, technology would represent an
absolute and fatalistically determined reality not influenceable by hu-
mans, expressed representatively, e.g., in the following way: "Techno-
logical development will continue, and human existence will change
with it."[124] Humans are exposed to and intertwined with technology
while, e.g., holding a gun – using the imagery of Bruno Latour: "You are
different with a gun in hand; the gun is different with you holding it.
You are another subject because you hold the gun; the gun is another
object because it has entered into a relationship with you. The gun is no
longer the gun in-the-armory or the gun-in-the-drawer or the gun-in-the-
pocket, but the gun-in-your-hand, aimed at someone who is screaming.
What is true of the subject, of the gunman, is as true of the object, of the
gun that is held. A good citizen becomes a criminal, a bad guy becomes
a worse guy; a silent gun becomes a fired gun, a new gun becomes a
used gun, a sporting gun becomes a weapon. The twin mistake of the
materialists and the sociologists is to start with essences, those of sub-
jects or those of objects. That starting point renders impossible our
measurement of the mediating role of techniques. Neither subject nor
object (nor their goals) is fixed."[125] This mediation needs to be consid-
ered because it is capable of covering the way humans interact with the
world. "On the one hand, the concept of mediation helps to show that
technologies actively shape the character of human-world relations.
Human contact with reality is always mediated, and technologies offer
one possible form of mediation. On the other hand, it means that any

---

[123] See Weizenbaum, Joseph (1976): Computer Power and Human Reason: From
Judgment to Calculation. New York: W. H. Freeman and Company.
[124] Verbeek 2015a: 214.
[125] Latour 1994: 33.

particular mediation can only arise within specific contexts of use and interpretation. Technologies do not control processes of mediation all by themselves, for the forms of mediation are always context-dependent"[126] This mediation offers an alternative to an "instrumentalist" understanding of technology (neutral means for humans to reach their aims) and to a "substantivist" understanding of technology (determining and controlling influence on the individual, the society, and culture).[127] he interactionist understanding of the relation between technology/technologies and ethics though coins this rapport as dynamic – knowing different grades of intensity concerning how humans, technology/technologies, and the world interact with each other, and embracing different kinds of compositions, how much each of them contribute to creating an ever-new reality from a human perspective. For example, understanding the shovel as a mediator between humans and the ground like the artist Tomi Ungerer.[128]

The interactionist model comprises the possibilities of mediation, of an instrumentalist narrowing of technology/technologies, of a substantivist elevation of technology/technologies, of ideological human self-image de fined by technology/technologies to the point of self-abandonment in favour of technology/technologies under the illusion of technology/technologies as an end in itself, and techno-critical reluctance of technology/technologies as well as their dynamic combination in manifold compositions. The interactionist approach situates ethics in dialogue with and under influence of this dynamic and of these possibilities. It understands ethics with the ability to distance itself to a certain extent – enough for an ethical analysis, discussion, evaluation, assessment, and judgement of technology/ technologies. This interactionist

---

[126] Verbeek 2015b: 11.

[127] See Verbeek 2015b: 1-12.

[128] See Ungerer, Aria (2019): "Als Vater war er gleichermassen grossartig und furchtbar". In: Die Zeit Magazin, October 16. Online: https://www.zeit.de/zeit-magazin/2019/43/aria-ungerer-tomi-ungerer-tochter-literatur [08.02.2021]. 18.

understanding of technology/technologies and ethics trusts in humans remaining, at the end of the day, able as subjects to put the gun as an object down and to reflect upon this tangible gun in this specific context as technology, to discuss in a democratic process the production, the accessibility, the sale, and the use of guns, to form a political opinion and make a political decision about it, to reflect upon guns in general as technology from an ethical standpoint, as well as technology/technologies in general from an ethical perspective. This interactionist understanding of technology/technologies and ethics trusts in humans to be capable to create and destroy technology/technologies under ethical guidance.

Ethics should focus on identifying technologies that open the door to ethical opportunities and therefore should be allowed depending on the ethical quality of the interaction between humans and technologies. Ethics should start from the fundamental paradigm that "artificial intelligences (machinae sapientes) are not evolutionary adversaries of homo sapiens but instruments (artefacts) that must be thought of as cooperative to the person."[129]

[129] Benanti, Paolo (2018): Le Macchine Sapienti: Intelligenze artificiali e decisioni umane. Bologna: Marietti. 114.

# 3

# DATA ETHICS AND LAW ARE TWINS

*Pavan Duggal, India*[130]

## 3.1 Introduction

Today, we are in a very interesting age of evolution as far as the internet and the electronic ecosystem is concerned. Internet has been with us for the last many decades. Internet has come to occupy a central lifeline in our day-to-day lives and we are all dependent on the internet.

As a natural corollary, we are all dealing, handling or processing data in the electronic form. We have all become global authors, global transmitters and global broadcasters of data. In fact, we are constantly generating data. World over, there is a new phenomenon. The world popula-

---

tion at large is undergoing a Great Data Vomiting Revolution where people are not just generating data but are vomiting data about their personal, professional and social lives, without even thinking about the legal ramifications for the same.

Data is all around us. No wonder, data has become the new oil of the new data economy. It is this new oil, which is of massive significance because this data tends to get monetized. More and more data stakeholders are interested about collecting and monetizing data today than at any point of time in history earlier.

Therefore, data ethics assume massive significance. Data has become the precious raw material of this century.  Hence, law has to come up with appropriate legal frameworks to deal with legally valid and sound approaches of dealing with data. More significantly, ethical principles have also evolved as to how data needs to be handled.

It is in this context that data ethics has become extremely significant in our lives. For the various definitions of data ethics we refer to the introduction of this book. Data ethics as a discipline has got a direct intrinsic connection with law.

## 3.2 Intrinsic Connection between Data Ethics and Law

There is intrinsic connection between data ethics and law, because data ethics provides basic logical and ethical principles on the basis of which data needs to be dealt, handled or processed with. These very logical and ethical principles often find themselves reflected in legal provisions under the law.

Therefore, law has got an intrinsic connection with data ethics. As more and more countries are today generating new laws on data, the principles of data ethics are getting enshrined in the said laws. In addition, more new principles on data ethics are evolving with the passage of time.

It is but natural to expect that principles of data ethics must be reflected in the evolving legal frameworks. No wonder, the interplay between law and data ethics thus becomes far more important and crucial.

In fact, it can easily be said that these two fields are interdependent on each other. It is the foundation of ethical principles in data economy provided by data ethics that becomes the raw material for the lawmakers to come up with new legal frameworks.

Similarly, the legal frameworks get far more meaning, topicality and relevance if they reflect and incorporate therein the foundational ethical principles pertaining to data. In fact, most of the people believe that data needs to be dealt, handled and processed with in an ethical manner. This becomes all the more significant since data invariably has got elements of personal privacy and also data privacy intrinsic therein. Most of the people today have an intrinsic expectation of privacy.

## 3.3 Ethical Principles for Using Data

Right now, experts agree on the following *ethical principles for using data:*

1. Privacy customer identity and data should remain private.
2. Shared private information should always remain private.
3. Customers should exercise a transparent view of how the data is being sold or utilized.
4. There should be no interference between big data and human will.
5. Big data should not institutionalize prejudicial biases.
6. As part of Consumer Relations, having excellent data ethics is a brilliant business decision.
7. Legality – Data management becomes a legal concern in some aspects as well. Thus, one needs to comply with the given regulations.
8. Implementing Data Ethics – If you opt to stay in business for a couple of years, it is a must to manage your data ethically. The manner

of protecting your people's data greatly depends on your company's needs. Whatever you gather, you should always be transparent.[131]

## 3.4 Interplay between Data Ethics and Law

In such a scenario, I find that the interplay between data ethics and law continues to keep on growing with the passage of time.

However, we need to understand that time has also taught us one thing which is that technology moves at a very rapid pace. Technology is invariably growing at such cutting-edge pace that it is leaving behind legal developments and law. In fact, it is often said that the law is almost ten steps behind the advent of technology. The massive speed of technological developments have actually now sparked our imagination.

Newly emerging technologies are evolving which are continuing to dazzle not just the technological ecosystem but also users, whether it is Artificial Intelligence, Blockchains, Internet of Things (IoT), Quantum Computing or the Metaverse. These newly emerging technological paradigms are continuing to amaze us. Also, the massive speed with which such developments in technology are evolving, is actually making us gasp for breath.

While the newly emerging technologies and their constructive usages are evolving, the cyber criminals are not far behind. They are increasingly coming up with new ways of misusing the said technologies to the detriment of people at large. Therefore, enabling technology legislation is the answer to the massive speed of technological developments.

Technology legislations across the world are taking the manifestation of cyber legal frameworks. Cyberlaw is a discipline that deals with the legal, policy and regulatory issues pertaining to technologies including the internet, cyberspace and the World Wide Web. Most of the countries, till today, have relied upon the UNCITRAL Model Law on Electronic

---

[131]  Analytics   Insight.   2021.   Ethical   Principles   for   Using   Data. https://www.analyticsinsight.net/ethical-principles-for-using-data/

Commerce which has been ratified by the United Nations General Assembly for the member nations to come up with their own respective national laws on cyber legal affairs.

As a result, most of the countries have come up with their national cyber laws. In addition, the advent of newly emerging technological paradigms has effectively meant that countries have started coming up with more dedicated legal frameworks on these newly emerging technologies. Data protection and privacy have become an important thrust point for majority of countries. Therefore, beginning with the General Data Protection Regulations (GDPR) of the European Union, one finds that countries are increasingly coming up with new laws on data protection.

## 3.5 Advent of Emerging Technology and Data Ethics

Further, the advent of newly emerging technologies like cyber security is propelling countries to come up with their own distinctive national laws on cyber security. Various countries are in the process of not just coming up with new laws on cyber security but also updating their current laws so as to make them topical and relevant in the context of growing cyber security breaches.

The coming of Internet of Things (IoT) has actually meant that countries need to wake up to the need for protecting cyber security in the Internet of Things (IoT) paradigm. Therefore, countries have started coming up with their new legal frameworks on Internet of Things (IoT). These include legislations like the US Federal Internet of Things (IoT) Cybersecurity Improvement Act 2020. The European Union is now coming up with the new draft Cyber Resilience Act and also new draft legislation on Artificial Intelligence.

The advent of Artificial Intelligence is also propelling countries to start exploring ways of how they can come up with new legal frameworks so as to deal with Artificial Intelligence as a paradigm.

In these newly emerging technologies and their legal regulation, there is an intrinsic role of ethical values. This is so because new technology legislations, which are coming up to regulate technology, must have ethical foundations and principles as an integral bedrock of the same.

This assumes even more significance because ethical principles and ethical values need to be incorporated as an integral part of data economy age. In this data economy age where everybody is dealing with data, it becomes imperative that the world must rely on ethical approaches to dealing with data.

## 3.6 Golden Age of Cybercrime and Increasing Cyber Security Breaches

Much more needs to be done in this particular regard. Data ethics therefore assumes far more significance in the context of legal frameworks. This assumes even more significance given yet another unique challenge that the emergence of technology has brought forward, which is growing cybercrimes and cyber security breaches.

Already with the advent of Covid-19, we have begun to see the coming and emergence of the Golden Age of Cybercrimes. This Golden Age of Cybercrime is going to be with us for many decades. Hence, the focus has to be on how to come up with ethical responses to deal with such growing cybercrimes including phishing, identity theft and online financial frauds.

Over the last few years, the world has been increasing cyber security breaches. These breaches are beginning to become of immense significance and challenge as far as all digital stakeholders are concerned.

These cyber security breaches ultimately are targeting data resident on computer systems and networks and its illegal and unauthorised monetization. These growing cyber security breaches and their impact be-

comes evident when one looks at the facts and figures pertaining to the same.

1. During the third quarter of 2022, approximately 15 million data records were exposed worldwide through data breaches. This figure had increased by 37 percent compared to the previous quarter.[132]

2. Between March 2021 and March 2022, the average cost of a data breach in the healthcare sector amounted to over 10 million U.S. dollars, up from 9.23 U.S. dollars between May 2020 and March 2021. The financial industry ranked second, with 5.97 U.S. dollars per breach on average. The global average cost of a data breach in the measured period was 4.35 million U.S. dollars. Data breaches in the public sector ranked last, costing an average of 2.07 million U.S. dollars during the measured period.[133]

3. An average of 4,800 websites a month are compromised with formjacking code[134]

4. By stealing 10 credit cards per website, cybercriminals earn up to $2.2 million through formjacking attacks[135].

5. The average total cost of a data breach was more than $1 million higher when working remote was a factor in causing the breach, compared to breaches in which working remote was not a factor[136].

6. Cyber scams increased by 400 percent in the month of March 2020, making COVID-19 the largest-ever security threat. [137]

---

[132] Statista. 2022. Number of data records exposed worldwide from 1st quarter 2020 to 3rd quarter 2022, https://www.statista.com/statistics/1307426/number-of-data-breaches-worldwide/

[133] Statista. 2022. Average cost of a data breach worldwide from May 2020 to March 2022, by industry, https://www.statista.com/statistics/387861/cost-data-breach-by-industry/

[134] Symantec, 2019 Internet Security Threat Report, Executive Summary, ISTR Vol 24.

[135] Ibid.

[136] IBM. 2022. Reports. Cost of a data breach 2022, https://www.ibm.com/reports/data-breach.

A perusal of the aforesaid facts and figures of cyber security breaches thus clearly shows that they will have an extremely detrimental impact upon the data economy. Hence, the adoption of data ethics and related legal principles becomes far more relevant in the context of not just data economy but data economy stakeholders at large. Therefore, ethical values need to be an integral part of legal response and regulation mechanism, so as to deal with data in the data economy age.

Therefore, ethical principles and ethical values which form integral part of data ethics need to be an integral part of the legal response mechanism to emerging technology laws. But when one looks at the prevailing international approaches, one finds that countries and state actors tend to adopt a dual approach.

## 3.7 Dual Approaches on Data Ethics and Law by Countries

In the context of data ethics and law vis-à-vis their national legislations, countries often incorporate data ethical principles as an integral component thereof. This assumes more significance because countries want their national citizens to comply with the ethical principles and therefore data ethics principles often find reflection in the national laws.

However, when one looks at the international scenario, countries tend to behave differently.  In today's scenario, where countries are engaging in both covert and overt activities, countries do not really want any naming and shaming phenomenon in cyberspace or also in international paradigm.

Therefore, countries then often tend to drag their feet and not respond effectively in the direction of ethical regulation of cyberspace at the international scenario. That is the reason why there is no international Cyberlaw or cyber security law in place.

---

[137] Sobers, Rob. 2022. 89 Must-Know Data Breach Statistics [2022], Varonis. https://www.varonis.com/blog/data-breach-statistics.

We need to understand that there is a need for incorporating data ethics in the legal principles and legal frameworks, both at the national level as also at the international level, because ultimately dealing with data in an ethical manner becomes a foundation for the further robust growth of the data economy.

## 3.8 Increasing Cyber Attacks and Data Ethics

However, there is yet another new trend on the horizon which is beginning to threaten data ethics and law which is growing cyber attacks. Cyber attacks are now being launched in by both state and non-state actors for the purposes of breaching the cyber security of computer systems and networks located in other jurisdictions and data that is resident therein. This becomes even more apparent when one looks at the numbers in this regard.

1. Malware increased by 358 percent in 2020.[138]
2. Ransomware attacks rose by 435 percent in 2020 compared to 2019.[139]
3. On average, a company falls victim to a ransomware attack every 11 seconds.[140]
4. 57 percent of organizations see weekly or daily phishing attempts.[141]
5. 65 percent of cybercriminal groups used spear-phishing as the primary infection vector.[142]

---

[138] Help Net Security. 2021. Malware increased by 358% in 2020, https://www.helpnetsecurity.com/2021/02/17/malware-2020/.

[139] Ibid.

[140] Morgan, Steve. 2020. Cybercrime To Cost The World $10.5 Trillion Annually By 2025. Cybercrime Magazine. https://cybersecurityventures.com/cybercrime-damages-6-trillion-by-2021/.

[141] Business Email Compromise Report. 2021. Cybersecurity Insiders. https://info.greathorn.com/hubfs/Reports/2021-Business-Email-Compromise-Report-GreatHorn.pdf.

6. Phishing attacks account for more than 80 percent of reported security incidents.[143]

1. $17,700 is lost every minute due to a phishing attack.[144]

2. By 2023, the total number of DDoS attacks worldwide will be 15.4 million.[145]

3. Attacks on IoT devices tripled in the first half of 2019.[146]

4. IoT devices experience an average of 5,200 attacks per month.[147]

As of now, at this particular moment, there is no dedicated international agreement on how to deal with the growing cyber attacks. There is Tallinn Manual 1.0 and Tallinn Manual 2.0, which have basically relied upon ethical principles of how countries need to adopt norms of behaviour in cyberspace. But that at best are only academic works and have not yet gained substantial traction across the world.

In this kind of a scenario, there is an urgent and immediate need that countries need to synergise the speed of technological legislations and ethical values. The speed of technological legislation effectively should be embellished with appropriate data ethical values which are enshrined in emerging principles of data ethics.

This is the need of the hour and it has to be appropriately adhered to. In fact, when one looks at the international scenario, one finds that there is a need to evolve appropriate important guidelines and principles for the digital stakeholders on how they can actually follow the principles of

---

[142] Broadcom. https://www.broadcom.com/support/security-center.

[143] Carlson, Brian. 2021. Top cybersecurity statistics, trends, and facts, CSO, https://www.csoonline.com/article/3634869/top-cybersecurity-statistics-trends-and-facts.html.

[144] Ibid.

[145] CISCO. 2019. Annual internet report, white paper. https://www.cisco.com

[146] Carlson, Brian. 2021. Top cybersecurity statistics, trends, and facts, op. cit.

[147] Sobers, Rob. 2022. 166 Cybersecurity Statistics and Trends [updated 2022] https://www.varonis.com/blog/cybersecurity-statistics

data ethics and law as they deal with the various emerging challenges of the data economy.

The ultimate future growth of data economy is dependent on how quickly in a resilient manner the principles of data ethics get incorporated in legal frameworks and how the said ethical principles are appropriately and effectively implemented through instruments of law. That is the focussed thrust and direction that countries need to move in.

## 3.9 Darknet and Data Ethics

There is also a need for countries to realize yet another big danger that is emerging on the horizon. The emergence of darknet has today brought in distinctive new paradigms. Darknet is the deep, dark underbelly of the internet where the intrinsic architecture of the darknet and the TOR network actually promises complete anonymity to all stakeholders. Complete anonymity breaths contempt in the form of growing cybercrimes.

Therefore, cybercrime is the dominant economic model as far as the darknet is concerned. Darknet is the complete antithesis of data ethics and law. In fact, darknet is used as a launching pad for committing various cybercrimes and cyber attacks on all stakeholders of the data economy as also the superficial net. However, while all kinds of unethical activities take place on darknet, it is also unique and important to appreciate that the darknet is also the embodiment of ethical values because even in that anonymous space, even for doing illegal and criminal activities, the stakeholders are adhering to certain ethical values. Though there are anonymous stakeholders, they have come up with appropriate indigenous approaches so as to uphold ethical values. So if a deviant actor does not perform the promised contracted deliverables on the darknet, he is blacklisted by the entire darknet community.

Hence, even in a place where criminal activities are launched in complete contravention of data ethics, there are still principles of data

ethics that are being followed by the respective stakeholders in the dark-net economy. The challenges of the darknet have to be kept in mind more so, from the context of data ethics and law as today countries across the world are still struggling to come up with appropriate new effective regulatory frameworks so as to deal with the darknet. The darknet today presents immense challenges for the emerging data econ-omy because more and more attacks on data economy are originating on the darknet.

## 3.10 Growing Data Economy Age, Data Ethics, Law and the Future

The growing data economy age is going to be the evolving future. In this context, the relationship of data ethics and law in the context of data economy assumes far more significance. In fact, more and more ethical principles, as are enshrined in the emerging discipline of data ethics, need to be incorporated in legal frameworks so that these newly embel-lished legal frameworks can then help support the further growth of the data economy at large.

The future is the future of data and the future lies in the manner of how data is going to be continuously getting evolved and monetized. This becomes further manifested in the following facts and figures per-taining to the projected growth of data in the data economy age.

1. Already, it has been estimated that poor data quality costs the US economy up to $3.1 trillion yearly.
2. In 2020, every person generated 1.7 megabytes in just a second.
3. Internet users generate about 2.5 quintillion bytes of data each day.
4. Predictions estimate the world will generate 181 zettabytes of data by 2025.
5. 80-90% of the data we generate today is unstructured.
6. The market of Big Data analytics in banking is set to reach $62.10 billion by 2025.

7. Big data in healthcare could be worth $71.6 billion by 2027.
8. Data interactions went up by 5000% between 2010 and 2020.

A perusal of the aforesaid facts and figures thus clearly tell us that data is now the new DNA that is connecting all stakeholders in the digital ecosystem. This data needs to be appropriately protected. The data economy age will only continue to consolidate in case if ethical principles are abided to and adhered to by data stakeholders when they deal, handle or process data.

## 3.11 Conclusion

Hence, the principles of data ethics and law, when joined together, can play a very important and cogent role in the further evolution of data economy as a whole.

To conclude, one can specifically state that data economy age is representing a new chapter in our lives. Data ethics as a discipline is continuing to evolve and stipulate new principles and foundations for the ethical use of data in the electronic ecosystem.

These ethical principles need to be well enshrined in the legal frameworks and laws across the world so that they can become potent combination in the direction of more effective and minimal enablement and regulation of the data economy age.

Together, with data ethics and law joining hands, they can provide a very potent constructive and positive direction in which the data economy has to ultimately evolve in the coming times.

It will be really interesting to see how the interplay between data ethics and law will continue to evolve with the passage of time and would have an impact upon the data economy in the coming times.

# 4

# DATA ETHICS OFFICER:
# A ROLE TO BUILD

*Enrico Panai, Italy[148]*

## Abstract

With Artificial Intelligence (AI) and the development and spread of systems that fall into this category, ethics has returned to the forefront.[149] To the extent that the noun 'ethics', the adjective 'ethical' and the adverb 'ethically' are used in most texts on the governance and policies of AI

---

[148] Department of Humanities and Social Sciences (DUMAS), University of Sassari, Via Roma, 151, Sassari, 07100, (SS), Italy; Fellow of ForHumanity (https://forhumanity.center). Corresponding author(s). E-mail(s): enrico-panai@gmail.com.

[149] This article was originally published on December 15, 2021 in Italian in Vol 2 No 2 of the Giornale di Filosofia with the original title "Responsabile dell'etica dei dati. Un ruolo da costruire" https://mimesisjournals.com/ojs/ index.php/giornale-filosofia/article/view/1697.

systems. Governments are regulating the imprudent and risky use of AI systems, but the resulting regulations may be insufficient to regulate the AI market 'ethically'. However, the use and sometimes misuse of the term is not proportionate to the actual presence of 'ethicists' (with an appropriate philosophical background) in ethics committees. This deficiency often leads the ethical debate onto a slippery slope where ethical language is not respected. Ethics then becomes the last resort to assert one's convictions, like a spell in which the mere utterance of its sound allows one to avoid ethical reasoning. In this article, I attempt to reaffirm the centrality of ethics as a discipline that has its own logic (ethical reasoning), its own methodologies (e.g. levels of abstraction), its own classical macroethical theoretical frameworks (e.g. value ethics, deontologism, consequentialism, information ethics), and a set of operational non-standard ethics (e.g. medical ethics, business ethics, ecology, etc.). To use a metaphor, ethics is like the practice of medicine: there are a few frameworks (biomedical model, psychosocial model, public health model) that provide structure to different specialties: anesthesiology, cardiology, dermatology, psychiatry, epidemiology etc. Similarly, ethics is based on a few macro-ethics that feed into the non-standard ethics: medical ethics, ecology, business ethics, data ethics and AI ethics. The ultimate aim is to quickly put the Western ethical disciplines in order, because if one does not go to a dermatologist for a toothache (although both the dermatologist and the dentist have studied medicine), in the same way when a company wants to implement AI-related ethical issues, it has to find the right specialist to avoid making the wrong ethical choices. In recent years, ethical choices related to data and AI have often been delegated to the Chief Ethics Officer, when one exists within the company, or to the Data Protection Officer (DPO) in the EU or to the Chief Data Officer (CDO), or whatever the title might be, in the rest of the world. However, while the former works on a very fine level of ethical granularity concerning the company's external relations with society and internal relations with its collaborators, the latter changes the level

of abstraction completely, having to deal mainly with the legal aspects related to data and personal data. So while the former bases its ethical choices on business ethics, the latter make decisions with the support of local legislation (e.g. GDPR). As a result, ethical decisions on data and AI algorithms are outside a defined ethical framework. With the risk of letting those with more technical expertise or more internal political power impose their personal or political beliefs. However, making ethical choices about digital technologies, data, and ultimately AI systems requires reorganising the debate within a defined and shared ethical framework. To this end, I propose a map of the competences needed within a company for the creation of an internal ethical infrastructure (internal infra-ethics) that allows contextually appropriated ethical choices to flourish, with the awareness that this analysis is limited to Western business philosophies and cultures. Finally, I attempt to show how even in small companies, there is no need for many different profiles, if the conceptual map of ethical functions is clear and shared.

## 4.1 Introduction

Collecting, storing or processing data may affect individuals, social groups or society as a whole. This implies that decisions made on certain data may have ethical value. In general, many organisations (in the sense of an entity with a defined purpose, such as associations, companies or institutions) are confronted with ethical reasoning when handling data. There is, thus de facto a space of data ethics where decisions are already being made. However, if within some organisations this space is not evident, there is a risk of making (ethically) inappropriate decisions. This space has been there for millennia. However, with the development of digital information and communication technologies (which have rationalised the collection, storage and processing of data) and the recent spread of artificial intelligence systems (which have enhanced the classification and inference capacity of machines), the impact of a wrong

choice or non-choice has become more appreciable. Without due precautions, the ethical consequences are only visible when the damage has been done.

## 4.2 Informational Abstraction Level

Writing was the first information and communication technology. Information technology because it made it possible to collect on a physical support (clay tablets, papyri, parchments, books) a series of data organised in such a way as to provide information. Communication technology because, in a few signs, it allowed the knowledge gathered to be stored for years and geographically disseminated. A clay tablet could be transported or stored to transmit the information it contained. Early media did not convey poems or stories but organisational information: religious procedures, human resource management, accounting, recipes, etc. Other systems have existed throughout history to store or transmit information without writing, such as the quipu, the knotted cords of the Incas, or the Marshall Islands stick chart[150].

Moreover, some information could not have been gathered by writing alone. While the ropes of the Incas were used to store structured data (census, trade transactions and taxes), the stick maps were used as dynamic maps for navigation in the Pacific; for example, Mattang maps were only used to teach how ocean waves break on islands. In practice, information and its management were critical for progress in all parts of the world. Writing, in particular, was the winning technology among the other competitors.

Whether knotted, written or passed on with sticks and shells, information has always had one characteristic: being inextricably associated with its physical medium. According to the General Definition of Information (Floridi, 2010), Information is composed of data (GDI 1); data

---

[150]   Marshall Islands stick chart, Wikipedia.org. https://en.wikipedia.org/wiki/Marshall_Islands_stick_chart

must be well-structured by the syntax that is used in the chosen system (GDI 2) and must respect the meaning of the system itself (GDI 3)[151]. In other words, if one of these three principles is not respected, then we have no information. For example, the sentence "colourless green ideas sleep furiously" (Chomsky, 1975)[152], which is composed of data ("green ideas" "without colour" "sleeping") and respects the syntax (subject and verb), is meaningless, thus conveying no information[153].

The general definition of information is crucial today because data have become the fundamental resource for creating information. In fact, before the computer age, the data that made up information were invisible because they were fused in their logical structure, syntax, and distribution on the medium. With the advent of information technology, a separation between data, medium and structure took place that emphasised the importance of data. In practice, information has been decomposed into its atomic elements (the data) and is recreated on the fly by organising the data according to a particular logical structure whenever we need it. Once regenerated, information appears on graphical interfaces (the screen of a computer or mobile phone), audio (virtual assistants such as Alexa, Siri, etc.), tactile (as in Braille keyboards), etc.

With this divorce of data from its physical support, a new era, the informational age, has opened up. Today's information and communication technologies (ICTs), namely digital technologies, have a different impact from writing precisely because their activity of transforming data into information is constant, dynamic, and not static as on clay tablets.

---

[151] Floridi, L. 2010. Information. A Very Short Introduction. Oxford University Press.

[152] This sentence was proposed by Chomsky as an example of a grammatically correct but semantically absurd sentence. Cf. N. Chomsky. 1975. The Logical Structure of Linguistic Theory, Springer US, Cambridge/Mass.

[153] Actually, from an informational point of view, this sentence gives information, but it needs interpretation, which can change depending on the level of abstraction used.

Moreover, they process a quantity of data that would not be possible to process on a physical medium; it is only necessary to remember that, since 2005, more data has been generated every day than has been produced since the invention of writing[154]. To disentangle ourselves from this tsunami of data, we need advanced computational and statistical resources to help us make decisions or enable automated systems to make choices, choices and decisions that necessarily have an ethical impact in the informational universe in which we live. Hence, using the level of informational abstraction, the need to use ethics appropriate to the system becomes evident.

Choosing the right level of abstraction to understand a phenomenon is a fundamental step in proposing solutions. In short, the 'informational' level of abstraction and the 'organisational' level of granularity is the perimeters of reflection in this article. We shall now see how, starting from an informational level of abstraction, the space of data ethics becomes apparent.

## 4.3 The Latent Space of Data Ethics

In organisations, many ethical decisions about data are already made. They exist. They are systemic. What needs to be added is the ethical model to be used. In other words, the space of data ethics exists, but it is latent: latent because it lacks a theory.

A sailing boat that sails is real, or at least we assume it to be so in common thinking. It is real because there is the 'model' of the boat (the design) which is identified with the 'properties' (hydrodynamics of the hull, height of the mast, cut and size of the sails, tension and thickness of the shrouds, type of keel) that are attributed to a final object, the 'system' (precisely because the properties are arranged in such a way that

---

[154] Data calculations are approximate, but some calculations can be made from some sources. E.g. https://www.seagate.com/files/www-content/our-story/trends/files/idc-seagate-dataage-whitepaper.pdf

the boat can sail). The boat (the object, the system) embodies the model and property information in a precise physical order. In this sense, information can be considered not as an ephemeral and intangible concept but as "the physical order, like that embodied in objects" (Hidalgo, 2015)[155]. A sailing boat that sails is thus 'real' when a model identifies properties attributed to a system.

Without the system, we are no longer in the presence of the 'real' but of the 'virtual'. We only have the model and its properties. That is why sail simulation software is virtual: there is a model with its properties, but there is no system. Decoupling the model from the system highlights the opposite case, in which the system exists, but not its model. Again, we are not in the presence of the 'real', but neither are we in the presence of the 'virtual'. Instead, we are in the presence of the "latent, in the original Latin sense of concealed or unknown" (Floridi, 2014, p. 318). This happens in archaeology when one comes across constructions of which one has no model. From the Sardinian Nuraghe[156] in the Mediterranean to the Moai[157] on Easter Island, there are systems of which we do not know a model. Systems that exist but with a dreamlike essence, are not fully real: "an object without its own shadow, like Peter Pan" (Floridi, 2014, p. 318)[158].

In every organisation, data exists and is the basis of the production system. Data is therefore collected, stored, processed and transmitted. It would be naive to think this can be done neutrally. Every choice is made within an ethical framework, even if this is not explicitly stated. If, for example, a company that sells books online collects racial data to target sales to categories of people, it implicitly reinforces the cultural divide

---

[155] C. Hidalgo, Why Information Grows: the Evolution of Order, from Atoms to Economies, Basic Books, New York 2015, p. 7.7.

[156] https://en.wikipedia.org/wiki/Nuraghe

[157] https://en.wikipedia.org/wiki/Moai

[158] L. Floridi, The Latent Nature of Global Information Warfare, in "Philosophy & Techno-logy", XXVII, 3, 2014, p. 318.

between different races. Whereas a public institution may collect the same data to mitigate racial injustice. Data ethics is present but latent. The system exists, but the (ethical) model has yet to be declared; because it is often unclear.

The question then arises as to why data ethics has yet to be explicitly introduced into the public debate, why there are AI ethicists but not data ethicists. Not because it does not exist or has not been conceptualised, but because we still need to put on the right glasses to read our age. So we are looking for ethicists of technology using data (AI, Robots, Computers, etc.) and not data ethicists. Technically, we have not yet accepted using an informational abstraction level introduced in the previous chapter. The 'method of abstraction levels' has already been used in mathematics (where it is closely related to model theory), in computer science (mainly in object-oriented programming or OOP) and philosophy (for a change of ontological perspective in information ethics).

Once we have selected the 'level of abstraction' of a complex system relevant to our purposes, we can create a 'model'. Levels of abstraction are necessary for system design and can be helpful in exposing a theory because "the model of a system is a function of the available observables" (Floridi, 2013, p. 33)[159]. But, using a specific level of abstraction means taking a point of view and developing one's decision-making process within that perspective. In our case, the informational status of abstraction emphasises the existence of an ethical data space. This is essential to properly discuss a topic because "it seems that many uninteresting disagreements might be clarified if the various interlocutors could make their Level of Abstractions explicit and precise" (Floridi, 2013, p.32). Once the level of abstraction has been chosen, one must decide the level of detail to discuss, namely the level of granularity. The higher the level of abstraction, the less detail there will be. In this case, as we shall

---

[159] L. Floridi, The Ethics of Information, Oxford University Press, Oxford 2013, p. 33.

see, the level of abstraction is the 'informational' level and the granularity is the partition of roles and responsibilities within organisations.

## 4.4 Informational Cartography of Organisations

In recent years, ethical choices concerning data and autonomous decision support systems seven generally referred to as artificial intelligence (AI) systems, have often[160] been 'delegated' to a chief ethics officer or Chief Ethics Officer (CEtO); or they have been 'delegated' to the data controller, the Chief Data Officer (CDO) to check the quality of the data lifecycle and ensure compliance with the legal framework relevant to the company. The 'ethical delegation' has been particularly evident within the European Union with the emergence of the principle of the 'data subject': the informational individual who must be protected across geographical boundaries. In the old continent, the role of the CDO was quickly replaced by that of the Data Protection Officer (DPO) dedicated to ensuring compliance with the European data protection regulation (GDPR - General Data Protection Regulation). To be precise, in previous cases it was not a matter of 'delegating' responsibility to a person because the concept of ethical responsibility for data was not precisely defined. Rather, it was a matter of 'colonisation' of territory no one claimed ownership of. This vagueness allowed everyone to impose their own level of abstraction according to their own strengths, capabilities or resources. Many ethical decisions concerning data were made by those who physically work with data: developers, data scientists or data workers in general.

Data workers were on the ground and solved ethical problems without having to (or being able to) consult the higher levels of the hierarchy. They answered to questions such as: Which dataset (data collection) to choose? Where data must be collected? How data must be stored? How

---

[160] In case you were aware of the ethical impact

errors in the dataset must be mitigated?  Which algorithm should be used for processing data? And so on. In short, data workers have and have had a decisive influence on ethical data choices in companies.

However, while the ethics officer works on a rough level of granularity concerning the company's external relations with society and its internal relations with its employees; the data protection officer changes the level of abstraction completely, having to deal mainly with the legal aspects related to data and/or personal data; finally, data workers are asked to be performers, putting ethical reflection on the back burner. Thus, the former base their conclusions on business ethics, the latter make decisions in accordance with local legislation, and the latter focus on the performance of the system they build.

As a result, ethical decisions on data and AI algorithms are outside a defined ethical framework. To use a medical metaphor, when it comes to ethical data decisions, companies seem to have dermatologists who improvise as dentists to replace a molar. The greatest risk is to create an ethical jungle, an environment in which one's personal or political convictions prevail. In practice, those with more technical expertise or more internal power impose their views on others. To avoid this risk and to make ethical choices about digital technologies, data and ultimately AI systems, there is a need to reorganise the debate within a defined ethical framework and to explicitly indicate levels of responsibility.

In this regard, I propose a map of roles needed within a company for the creation of an internal ethical infrastructure that allows ethical choices suited to the context to flourish. However, the result should not be restrained to Western cultures. This analysis provides insights for future research investigating the integration of data ethics in non-Western cultures.

Figure 1: The data ethics space occupied by other roles

The two diagrams define an organisation's decision-making space by shifting the level of abstraction of ethical investigations from information to data. Mapping the zone of influence of each role in an organisation brought to the surface a latent space that was invaded by other existing profiles, as represented by the diagram in Fig. 1.

## 4.5 The Role of the Data Ethics Officer

The objective of data ethics is to create the ethical or infra-ethical infrastructure in the organisation that facilitates ethical decisions about data and how it will be processed. The data ethics space should therefore be occupied by a Chief Data Ethics Officer: Chief Data Ethics Officer (CDEO).

Figure 2: The Chief Data Ethics Officer (CDEO) level

The CDEO should foster this ethical environment by establishing and maintaining a code of data ethics: Code of Data Ethics (CoDE). The code of data ethics must be consistent internally (with the proposed principles and standards) and externally (with the company's code of ethics). The Data Ethics Officer should promote the dissemination of the code of data ethics within the company. Its adoption would explicitly define the decision-making perimeter on data ethics. Finally, all actors involved in data processing, and in particular Data Workers (DWs), should be trained in data ethical reasoning and regularly updated on the best ethical practices in circulation. In short, the space for data ethics exists, as in Figure 2. It must be occupied by a data ethicist who can protect and promote the discipline[161].

## 4.6 Conclusions

This article supports the idea that there is a new ethical space within organisations concerning data ethics. This new space has helped to raise

---

[161] Taddeo, L. Floridi, What is data ethics?, in "Philosophical Transaction of The Royal Society AI", n. 374, 2016, p. 1.

important theoretical questions concerning the internal organisation of companies in mature information societies. Of course, the reflection is limited geographically and culturally because it begins, is founded and evolves within Western philosophies. Despite these limitations, the study suggests the recognition of a specific domain centred on data ethics and overseen by a chief data ethics officer (the Chief Data Ethics Officer - CDEO) to build an information ethics infrastructure within organisations.

### List of most common acronyms

- AI: Artificial Intelligence
- CDEO: Chief Data Ethics Officer
- CDO: Chief Data Officer
- CDOC: Children Data Oversight Committee
- CEO: Chief Executive Officer
- CEtO: Chief Ethics Officer
- CISO: Chief Information Security Officer
- CoDE: Code of Data Ethics
- CoE: Code of Ethics
- DPO: Data Protection Officer
- DW: Data Workers
- EU: European Union
- GDPR: General Data Protection Regulation
- IAAIS: Independent Audits of Artificial Intelligence Systems
- Infraethics: Infrastructure of ethics
- LoA: Level of Abstraction
- SME: Small and Medium-sized Enterprises

### Bibliography

Chomsky, N. 1975. The Logical Structure of Linguistic Theory. Springer US.

Floridi, L., Information. A Very Short Introduction. Oxford University Press 2010.

Floridi, L., The Ethics of Information. Oxford University Press 2013.

Floridi, L. The Latent Nature of Global Information Warfare. Philosophy & Technology 27, 317–319, 2014. https://doi.org/10.1007/s13347-014-0171-x

Information Warfare, Philos317-319 2014. https://doi.org/10.1007/s13347-014-0171-x

Hidalgo, C. Why Information Grows: the Evolution of Order, from Atoms to Economies. Wiley, 2015.

Taddeo, M., Floridi, L., What is Data Ethics? Phil. Trans. R. Soc. A., 374(2083):(Dec 28), 1-5 2016. https://doi.org/http://dx.doi.org/10.1098/rsta.2016.0360

# 5

# TETHICS: TECH ETHICISTS
# AS NEW PROFESSION
# AND CAREER PATH

*Veronica Irwin, USA*[162]

Rebekah Tweed first became interested in tech ethics as a music reporter. She was covering a Taylor Swift concert, and the artist had accidentally courted controversy after using facial recognition technology[163] to scan crowds for stalkers. "It struck me that they weren't fully thinking through the implications of some of these technologies," she explained. "So I started digging."

[163] Snapes, Laura. 2018. Taylor Swift used facial recognition software to detect stalkers at LA concert, The Guardian, https://www.theguardian.com/music/2018/ dec/13/taylor-swift-facial-recognition-stalkers-rose-bowl-concert.

## 5.1 Responsible Tech as Career Path

What she found was that she was a few years behind the trend. Researchers and academics have been analysing technology through an ethics lens for decades, and began reaching out of their siloed fields in a significant way after the 2016 presidential election.

In 2022, "responsible tech" is a career path. Job titles range from "trust and safety officer" to "policy lead." Several organizations and academic institutions are engaged in ecosystem-mapping projects to define, which academic programs best prepare students to work in the field, how the jobs are described and what companies are pursuing ethical tech in earnest.

"There's a lot of appetite for this, especially as the public has become very aware of highly publicized problems with technology," Tweed, now the program director for All Tech is Human, said. "I see that continuing to grow for the foreseeable future."

Agreeing upon a common nomenclature has been one of the first hurdles for experts defining the field. "Responsible tech" is meant to emphasize that technology should be responsive to the needs of users. New America calls the field "public interest technology." Cal Poly San Luis Obispo uses the phrase "tech ethics," in the context of academic programming they call "The Ethical Tech Project" (no, not as a callout to that 2019 "Silicon Valley" episode).[164] "There's a lot of similar people interested in similar types of work but in different sectors," said Tweed. "There just wasn't really a common nomenclature."

## 5.2 Public Interest Technology Network

One of the core organizations defining the space is New America and its Public Interest Technology University Network, a collaboration be-

---

[164] Tethics, 2019. Film comedy. Director Pete Chatmon https://www.imdb.com/title/tt10422434/.

tween 48 different institutions with academic programming relevant to the subject. Through the network, New America has invested more than $11.6 million in projects that help build clear career pipelines. The goal is to organize the most urgent ethical problems in tech and identify the skills needed to solve those problems. To do this, Stanford launched a "Public Interest Technology Career Taxonomy Project"[165], where student researchers index words in job descriptions for use by LinkedIn and hiring managers.

The Stanford project is based partly on previous research used to define and develop a pipeline for "green" jobs by finding the keywords that job searchers would use.[166] The project is part of a larger, student-run university research program called the PIT Lab[167], which helps prepare students for careers in the field by hosting panels, discussions and connecting students with similar fellowship opportunities.[168]

Andreen Soley, director of the PIT program at New America, said the member network demonstrates how a defining element of the field is a focus on interdisciplinary expertise. "This is not just the purview of computer scientists," she said. "One part of our conceptualization is to say you might need a sociologist to be a part of this conversation, or someone from another field."

Notably, New America doesn't have a perfect record when it comes to working with tech. The foundation has received at least $21 million from Google[169] since its inception, and Eric Schmidt, who was chairman

---

[165] Schneier, Bruce. 2022. Public-Interest Technology Resources, https://public-interest-tech.com/.

[166] Peggy Brannigan, How to find Green Work, https://www.linkedin.com/pulse/how-find-green-work-peggy-brannigan/.

[167] PIT Lab, https://pitlab.stanford.edu/.

[168] Public Interest Technology Fellowship 2022, Stanford Cardinal Service; https://cardinalservice.stanford.edu/opportunities/public-interest-technology-fellowship-2022.

[169] Tiku, Nitasha. 2017. The Hard Consequence of Google's Soft Power, Wired. https://www.wired.com/story/google-new-america-open-markets/.

of New America until 2016, was simultaneously executive chairman of Alphabet from 2015 to 2018. Google has allegedly used this arrangement to exert influence over the organization, according to The New York Times.[170] Both the company and New America have denied that there is an improper relationship.

## 5.3 Linking Academia to Practice

Ethical Tech @ Cal Poly[171] led by professors Deb Donig and Matthew Harsh, also intends to map the responsible tech ecosystem, but specifically through the frame of linking academia to practice. The university as a whole follows an ethos of "learning by doing"[172], Donig explained, and practically speaking, there isn't just a forthcoming generation of new tech ethicists: There's an entirely new industry on the horizon. "We believe this is a new profession, not a new workforce," she said. "There should be students who are trained, as an outcome of their graduation, to be ethical technologists."

One element of the program is a course Donig teaches, called "Technically Human." The course is listed as part of the English Department, but students of all majors are invited to participate. It focuses on the stories people tell about technology, whether it be nonfiction narratives like the rise and fall of ex-Uber CEO Travis Kalanick and ex-WeWork CEO Adam Neumann or science-fiction stories like "2001: A Space

---

[170] Vogel, K. Google Critic Ousted From Think Tank, New York Times, 2017. https://www.nytimes.com/2017/08/30/us/politics/eric-schmidt-google-new-america.html.

[171] California Polytechnic State University (Cal Poly), The 2021-2022 National Science Foundation Grant Award: "Ethical Technology and the Future of Work". https://www.etcalpoly.org/nsf-grant-award-ethical-technology-and-the-future-of-work.

[172] California Polytechnic State University, https://www.calpoly.edu/learn-by-doing.

Odyssey." These stories, she said, guide students' understanding of where the ethical issues lie in technology, who is responsible for them, and creative ways technologists can create better tech — all skills she believes are necessary to pursue a career in ethical tech.

Though the job titles are new, the ways to attract new talent haven't changed. New America sponsored a career fair as part of a larger online convention [173] in October 2021, with the help of Tweed, professors Mona Sloane and Matthew Statler at NYU, and several other academic, governmental and non-profit groups. The career fair portion was decentralized, and used the job searching software Handshake[174] to connect students across the country with jobs that hiring managers were actively recruiting for on their campuses. Tweed is also working on another career fair taking place this May through a collaboration with Stanford, Pepperdine, University of Washington and New America.

Ethical tech jobs currently exist at small startups dedicated to the field and giant tech companies that build the tools we use every day. Deciding which companies are serious about ethics, though, is easier said than done. Many of the biggest tech companies have dedicated research teams focused on the impacts of their products, yet are also responsible for most of the scandals with which tech ethicists are most concerned. Google announced it would be expanding its AI ethics research department to 200 people in summer 2021[175] — after messily

---

[173] A Better Tech. Public Interest Technology. https://www.abettertech.net/.

[174] Lizzy, L. 2021. Recruiting is broken for Gen Z tech candidates. New tools give managers hope they can fix it, Protocol. https://www.protocol.com/workplace/recruitment-tools.

[175] Tripp Mickle. 2021. Google Plans to Double AI Ethics Research Staff, The Wall Street Journal. https://www.wsj.com/articles/google-plans-to-double-ai-ethics-research-staff-11620749048.

firing two of its most esteemed researchers for publishing a paper that found major flaws in the company's language processing models.[176]

Whether young technologists should take jobs like those weighs heavily on the minds of some experts in the space. "I am very concerned to send my ethically educated, diverse, extremely energetic students into abusive and oppressive workplaces," said Sloane, one of the two NYU professors who helped organize the October job fair.

Sloane and Statler said more students than ever are interested in entering the burgeoning field. Though students with debt may have less flexibility in the job they choose versus ones with more financial freedom, fewer students are having to make a choice. "Top talent will go to places that take these issues seriously," she said. "This is a talent, attraction and retention issue, bottom line."

Through her work with All Tech is Human, Tweed has created a running "Responsible Tech Job Board"[177], which lists positions at companies as large and well-known as Google and as niche as senior research positions at issue-specific think tanks. It's a simple, continuously-updating spreadsheet, with over 300 roles listed currently.

## 5.4 Three Core Criteria for Positions

With the help of All Tech is Human's founder David Ryan Polgar, Tweed has settled on three core criteria for positions she lists on the job board: The roles either focus on reducing the harms of technology, diversifying the tech pipeline or aligning new tech with the public interest. But still, she said that some companies play tricks. She declined to name names, but said that it's important to note a company's historical record

---

[176] Kramer, A. 2021. Google vows to do better on DEI and firings. Timnit Gebru is not impressed, Protocol, https://www.protocol.com/google-ends--investigation--timnit-gebru.

[177] Responsible Tech Job Board, All Tech Is Human, https://alltechishuman.org/responsible-tech-job-board.

in the field, and watch out for job descriptions with responsible tech language that looks misused or doesn't quite match the role's core responsibilities. "Let me just say that some companies are catching on to the fact that there is a large community of people who care about a responsible tech focus, and care about building products that are responsible and beneficial to society," said Tweed. "I've noticed that, over time, I'm having to be a little more judicious in determining what a role actually does or researching whether this company actually cares about the issues." Often, she said, the answers aren't black and white.

There are already many students and young professionals readily willing to pursue a career in public interest tech, and plenty of research institutes and academic programs building a bank of knowledge for the sector. But the biggest gap, Tweed said, is in connecting those aspiring professionals, and the information they've learned, to applicable use cases in the real world.

When she decided to pivot into the responsible tech industry in 2018, she "had trouble even tracking down these roles, much less finding a good fit," she said. "Now there's so many jobs that could fit on the job board that I'm not even able to include everything."

# PART B

# DATA ETHICS:
# ARTIFICIAL INTELLIGENCE,
# ROBOTS AND HUMANS

# 6

# TRENDS IN AI ETHICS: INDICATORS

*Human-Centered AI Institute/ Stanford University, USA*[178]

## 6.1 Overview

As artificial intelligence–powered innovations become ever more prevalent in our lives, the ethical challenges of AI applications are increasingly evident and subject to scrutiny. As previous chapters have addressed, the use of various AI technologies can lead to unintended but harmful consequences, such as privacy intrusion; discrimination based on gender, race/ethnicity, sexual orientation, or gender identity; and opaque decision-making, among other issues. Addressing existing ethi-

---

[178] Daniel Zhang, Saurabh Mishra, Erik Brynjolfsson, John Etchemendy, Deep Ganguli, Barbara, Grosz, Terah Lyons, James Manyika, Juan Carlos Niebles, Michael Sellitto, Yoav Shoham, Jack Clark, and Raymond Perrault, *The AI Index 2021 Annual Report*, AI Index Steering Committee, Human-Centered AI Institute, Stanford University, Stanford, CA, March 2021, *Chapter 5: Ethical Challenges of AI Applications*, 125-134. With permission of the publisher. Free download https://aiindex.stanford.edu/report/.

cal challenges and building responsible, fair AI innovations before they get deployed has never been more important.

This chapter tackles the efforts to address the ethical issues that have arisen alongside the rise of AI applications. It first looks at the recent proliferation of documents charting AI principles and frameworks, as well as how the media covers AI-related ethical issues. It then follows with a review of ethics-related research presented at AI conferences and what kind of ethics courses are being offered by computer science (CS) departments at universities around the world.

The AI Index team was surprised to discover how little data there is on this topic. Though a number of groups are producing a range of qualitative or normative outputs in the AI ethics domain, the field generally lacks benchmarks that can be used to measure or assess the relationship between broader societal discussions about technology development and the development of the technology itself. One datapoint, covered in the technical performance chapter, is the study by the National Institute of Standards and Technology on facial recognition performance with a focus on bias. Figuring out how to create more quantitative data presents a challenge for the research community, but it is a useful one to focus on. Policymakers are keenly aware of ethical concerns pertaining to AI, but it is easier for them to manage what they can measure, so finding ways to translate qualitative arguments into quantitative data is an essential step in the process.

*Highlights:*

• The number of papers with ethics-related keywords in titles submitted to AI conferences has grown since 2015, though the average number of paper titles matching ethics-related keywords at major AI conferences remains low over the years.

• The five news topics that got the most attention in 2020 related to the ethical use of AI were the release of the European Commission's white paper on AI, Google's dismissal of ethics researcher Timnit Gebru, the AI ethics committee formed by the United Nations,

the Vatican's AI ethics plan, and IBM's exiting the facial-recognition businesses.

The *Methodology* of the research is available in the Appendix of the original article.[179]

## 6.2 AI Principles and Frameworks

Since 2015, governments, private companies, intergovernmental organizations, and research/ professional organizations have been producing normative documents that chart the approaches to manage the ethical challenges of AI applications. Those documents, which include principles, guidelines, and more, provide frameworks for addressing the concerns and assessing the strategies attached to developing, deploying, and governing AI within various organizations. Some common themes that emerge from these AI principles and frameworks include privacy, accountability, transparency, and explainability.



The publication of AI principles signals that organizations are paying heed to and establishing a vision for AI governance. Even so, the proliferation of so-called ethical principles has met with criticism from ethics

---

[179] *The AI Index 2021 Annual Report,* 211-213.

researchers and human rights practitioners who oppose the imprecise usage of ethics-related terms. The critics also point out that they lack institutional frameworks and are non-binding in most cases. The vague and abstract nature of those principles fails to offer direction on how to implement AI-related ethics guidelines.

Researchers from the AI Ethics Lab in Boston created a ToolBox that tracks the growing body of AI principles. A total of 117 documents relating to AI principles were published between 2015 and 2020. Data shows that research and professional organizations were among the earliest to roll out AI principle documents, and private companies have to date issued the largest number of publications on AI principles among all organization types (Figure 5.1.1). Europe and Central Asia have the highest number of publications as of 2020 (52), followed by North America (41), and East Asia and Pacific (14), according to Figure below. In terms of rolling out ethics principles, 2018 was the clear high-water mark for tech companies— including IBM, Google, and Facebook—as well as various U.K., EU, and Australian government agencies.



NUMBER of NEW AI ETHICS PRINCIPLES by REGION, 2015-20
Source: AI Ethics Lab, 2020 | Chart: 2021 AI Index Report

Europe and Central Asia have the highest number of publications as of 2020 (44), followed by North America (30), and East Asia and Pacific (14). In terms of rolling out ethics principles, 2018 was the clear high-water mark for tech companies— including IBM, Google, and Face-

book—as well as various U.K., EU, and Australian government agencies.

## 6.3 Global News Media

How has the news media covered the topic of the ethical use of AI technologies? This section analyzed data from NetBase Quid, which searches the archived news database of LexisNexis for articles that discuss AI ethics[1], analyzing 60,000 English-language news sources and over 500,000 blogs in 2020.

The search found 3,047 articles related to AI technologies that include terms such as "human rights," "human values," "responsibility," "human control," "fairness," "discrimination" or "nondiscrimination," "transparency," "explainability," "safety and security," "accountability," and "privacy." (See the Appendix for more details on search terms.) NetBase Quid clustered the resulting media narratives into seven large themes based on language similarity.



NEWS COVERAGE on AI ETHICS (% of TOTAL) by THEME, 2020
Source: CAPIQ, Crunchbase, and NetBase Quid, 2020 | Chart: 2021 AI Index Report

The figure above shows that articles relating to AI ethics guidance and frameworks topped the list of the most covered news topics (21%) in 2020, followed by research and education (20%), and facial recognition (20%).

The five news topics that received the most attention in 2020 related to the ethical use of AI were:

The release of the European Commission's white paper on AI (5.9%)

Google's dismissal of ethics researcher Timnit Gebru (3.5%)

The AI ethics committee formed by the United Nations (2.7%)

The Vatican's AI ethics plan (2.6%)

IBM exiting the facial-recognition businesses (2.5%).

## 6.4 Ethics AI at Conferences

Researchers are writing more papers that focus directly on the ethics of AI, with submissions in this area more than doubling from 2015 to 2020. To measure the role of ethics in AI research, researchers from the Federal University of Rio Grande do Sul in Porto Alegre, Brazil, searched ethics-related terms in the titles of papers in leading AI, machine learning, and robotics conferences.

As the figure below shows, there has been a significant increase in the number of papers with ethics-related keywords in titles submitted to AI conferences since 2015. Further analysis in the second figure below shows the average number of keyword matches throughout all publications among the six major AI conferences. Despite the growing mentions in the previous chart, the average number of paper titles matching ethics-related keywords at major AI conferences remains low over the years.

Changes are coming to AI conferences, though. Starting in 2020, the topic of ethics was more tightly integrated into conference proceedings. For instance, the Neural Information Processing Systems (NeurIPS) conference, one of the biggest AI research conferences in the world, asked researchers to submit "Broader Impacts" statements alongside their work for the first time in 2020, which led to a deeper integration of ethical concerns into technical work. Additionally, there has been a recent proliferation of conferences and workshops that specifically focus on responsible AI, including the new Artificial Intelligence, Ethics, and Society.

Conference by the Association for the Advancement of Artificial Intelligence and the Conference on Fairness, Accountability, and Transparency by the Association for Computing Machinery.

There has been a significant increase in the number of papers with ethics-related keywords in titles submitted to AI conferences since 2015. Further analysis shows the average number of keyword matches throughout all publications among the six major AI conferences.



NUMBER of PAPER TITLES MENTIONING ETHICS KEYWORDS at AI CONFERENCES, 2000-19
Source: Prates et al., 2018 | Chart: 2021 AI Index Report



AVERAGE NUMBER of PAPER TITLES MENTIONING ETHICS KEYWORDS at SELECT LARGE AI CONFERENCES, 2000-19
Source: Prates et al., 2018 | Chart: 2021 AI Index Report

## 6.5 Ethics Offerings at Higher Education Institutions

Chapter 4 of The AI Index 2021 Annual Report[180] introduced a survey of computer science departments or schools at top universities

---

[180] *The AI Index 2021 Annual Report, chapter 4.*

around the world in order to assess the state of AI education in higher education institutions. In part, the survey asked whether the CS department or university offers the opportunity to learn about the ethical side of AI and CS.



AI ETHICS OFFERING at CS DEPARTMENTS of TOP UNIVERSITIES around the WORLD, AY 2019-20
Source: AI Index, 2020 | Chart: 2021 AI Index Report

Among the 16 universities that completed the survey, 13 reported some type of relevant offering. The figure below shows that 11 of the 18 departments report hosting keynote events or panel discussions on AI ethics, while 7 of them offer stand-alone courses on AI ethics in CS or other departments at their university. Some universities also offer classes on ethics in the computer science field in general, including stand-alone CS ethics courses or ethics modules embedded in the CS curriculum offering. 11 of the 18 departments report hosting keynote events or panel discussions on AI ethics, while 7 of them offer stand-alone courses on AI ethics in CS or other departments at their university.

# 7

# ETHICS = DATA = INFORMATION: EVERYTHING IS INFORMATION

*Arend van Campen, Switzerland*[181]

## 7.1 The School of Athens

'In the beginning was the Word, and the Word was with God, and the Word was God.' In the beginning was the information, and the information was with God and the information was God. In the beginning there was the Big Bang[182]. This is the idea (information by human perception) that the universe began as a single point and then expanded and stretched as large as it is now and is still expanding.

The first step in this research will be to travel back in time and visit the School of Athens[183], because of Aristotle's emphasis on wisdom. Knowing why, hence discovering the first cause and effects of an information deficit on living systems, the environment and social cohesion. To establish empirically the hypothesis' anthropogenic failure due to information deficit, this research endeavours to seek knowledge of the first cause.

The School of Athens had four main goals or purposes:

Seek knowledge of Causes

Divine Inspiration

Knowledge of the Divine

To each what is due

This research will apply number 1.

… In the 1991 movie 'Mind Walk' written by Fritjof Capra[184] and directed by his brother Bernd, three people, a poet, a failed candidate for the US presidency and a physicist roam the island of Mont Saint Michel in France pondering over and discussing the world's existential threats, problems and ideologies such as religion, geopolitics, scientism and threats such as wars, conflict or pollution. They agree and decide that these issues have in fact one cause which is how people perceive reality. 'A crisis of perception' writes Fritjof which he repeats often in his books, courses, talks and presentations. The movie was based on his book 'The Turning Point' (1991): *'are we now at one of these turning points'* one of the characters asks at the end of the film.  When exploring the first cause principle the following question is necessary to understand the nature and origin of our perceived and experienced reality, resulting in probable harmful actions; what is information? Recognizing probability patterns is key in this research. On one hand people have unlimited creativity to design art, observe and visit space or build the

---

[183] School of Athens: https://www.britannica.com/topic/School-of-Athens

[184] Capra, Fritjof, Byars, Floyd and Capra, Bernd, *Mindwalk*, Brass Tacks Press, 2021, 1st. ed. 1991.

Sagrada Familia and on the other hand destroy life, the environment and social cohesion. This research will endeavour to establish that the first cause destructive tendencies and harmful design are unawareness and ignorance of the value of information.

## 7.2 What is Information? What is Entropy?

### 7.2.1 Information is physical, related to entropy

These questions are asked by several late and contemporary physicists. Rolf Landauer stated that information is physical (Landauer, Rolf, 1991)[185] and related it to entropy, which leads to disorder in systems, describing the measure of entropy in a number of bits which are often regarded as binary units of information but according to James V Stone[186] in his book 'Information Theory' (2015)  can be understood as: a bit is the amount   of information required to choose between two equally possible alternatives e.g., left or right,  true or false, up or down, wet or dry, etcetera. and a binary digit is the value of a binary variable, which can adopt one of two possible variables (i.e.,  0 /1). Information is expressed as a quantity next to a quantity of mass and a quantity of energy, all three of them interrelated, interconnected and interdependent. The mainstream consensus is that the universe consists only of matter and energy, which as we shall learn here, is a fundamental misperception, the one which Fritjof Capra talks about. It leads to a misunderstanding of reality that information is not always deemed important and therefore can be ignored, while, as this research will endeavour to establish, information is as or even more fundamental than physics and therefore also must be the foundation for and of physics and our experienced reality.

[185] Landauer, Rolf, Information is physical, Physics Today 44, 5, 23 (1991); https://doi.org/10.1063/1.881299

[186] Stone, James V, Information Theory, a tutorial introduction.  https://doi.org/ 10.48550/arXiv.1802.05968

### 7.2.2 Twenty questions

Let's play the game of 20 questions, mentioned by Carl Sagan in his series Cosmos part 11, (1980). This is a method to quantifying the information, by reducing the number of choices after each question to reach an answer: for example: 'is it alive or not alive?' 'Is it an animal?' 'Is it big enough to see?' 'Does it grow on the land?'

This method of questioning is a form of information or data collection to reach a solution or find an answer to the question: 'what is it?' After discovering the answer, it must be interpreted by human perception by all human senses, not just by the brain, in order to process the information into an experienced and viable reality. After each question and answer one ends up with 1 bit of information which allows one to make a choice between 2 equally probable alternatives (yes or no, 1 or 0). When one finds 2 bits of information, one can choose between 4 probable alternatives, 3 bits means 8 probable alternatives and so on. Increasing alternative probabilities also increases the quantity of information and thus the prospects to find an answer.

In twenty questions one can obtain 20 bits of information (Illustration below from the book 'Information Theory'). This allows for the



probable ability to narrow down the range of possible words from 1 million to one. Twenty questions are sufficient to find the correct word out of about 1 million words. If you'd double to quantity of questions to 40, it allows you to find 40 bits of information, allowing you to find one out of 2 to the 40th or 10 to the 12th words. The issue with the commonly accepted

empirical scientific method is that it does not ask all the questions, but limits itself to test scientifically by allowing the use of the physical human senses like 'sight, touch, smell, taste and hearing, which are mechanistic, analytic and reductionistic. They won't allow 'all relevant information' nor synthesis, which renders the results, although usable, no longer scientific due to a paradigm change.

Here is a listing of the overlooked or ignored importance of allowing and learning by a changed perception. The fresh insights and sciences to do this are available. Please find new, but not yet implemented scientific and logical probabilities which are needed to understand the whole, rather than the parts;

1. Relativity eliminated the Newtonian illusion of the absolute

2. Quantum mechanics eliminated the dream of a verifiable measurement process. Mechanistic thinking had played tricks on the human mind. It gave people the illusion of being objective and in control (looking in from the outside), but subjectivity (the observer effect) is confirmed by Quantum Mechanics and Quantum Consciousness.

3. Chaos Theory eliminates the fantasy of deterministic predictability, but discovered that order from chaos is a natural phenomenon.

4. Systems theory: everything is interconnected and interdependent (the observer effect).

5. Cybernetics. Norbert Wiener, Gregory Bateson, Stafford Beer and Ross Ashby et al. understood and used this notion of non-linearity and came up with the solution of the required diversity a.k.a. as requisite variety, namely maximising learning and the use of information. They realized that they can't control systems, but steer with feedback to direct energy and matter through the human mind and actions. Von Foerster talked about always maximising the number of questions to maximise the number of choices or options. Paul Pangaro talks about "conversation" in confirmation of a communicative physical existence which keeps living systems viable.

6. Information is physical. Claude Shannon, Rolf Landauer, Jim Khalili and Seth Lloyd discuss that information is "physical" and cannot be separated from universal reality. 'Information is the resolution of uncertainty'.

7. The term butterfly effect had to be introduced into science to understand chaos and complexity theory as unpredictability (non-linear effects).

## 7.3 Communication and mRNA

Awareness and understanding of the natural criteria of and for functionality of man-made systems is vital for all life on earth.

Claude Shannon [187] (1948) provides a mathematical definition of information and describes precisely how much information can be communicated between different elements of a living system and the limits of communication transfer rates in any system, whether man made or biological. Fact is that this quantification method demonstrates that information indeed is a physical substance. In biological living systems communication for gene expression is not regulated, but again 'steered' by DNA to produce the exact and correct proteins all the way from our genes to the functionality of large entities such as whole countries. The working of DNA and Shannon's theorem can be deemed equivalent regarding the transmission of information. DNA is understood as the source of the message and or transmission and proteins are the receivers. The message in the form of information must be structured. Lila Gatlin[188] of the University of California, Berkeley wrote 'Evolutionary

---

[187] Shannon, Claude. A mathematical theory of communication. The Bell System Technical Journal, Vol. 27, pp. 379–423, 623–656, July, Oct. 1948. https://doi.org/10.1002/j.1538-7305.1948.tb01338.x
[188] Gatlin, Lila. Evolutionary Indices, Volume 5 Darwinian, Neo-Darwinian, and Non-Darwinian Evolution, April 9–12, 1971, edited by Lucien M. Le Cam, Jerzy

Indices' to understand entropy in DNA messaging in higher and lower organisms and came to the conclusion that they are the same because of the uncertainty principle. DNA stores the hereditary information in a particular sequence of symbols from an alphabet of four letters: A.T.C and G, which means this is a language. For this study is it very important that we understand that entropy in the form of potential information, an unknown number of messages that were not sent, plays a role. Meanings can be misunderstood. If DNA is an information process and a theory of information exists, then it is reasonable to suppose that scientists can at least make a start of sketching out a theory of living organisms, in the full sense of the word (Campbell, Jeremy, 1982)[189]. When symbols of DNA are translated into the substance of proteins, communication takes place. In Gatlin's 'second-theorem selection' fitness is not dependent on strong body or reproduction, but of genetic information coding according to Shannon's principles.

Systems regulating genes are influenced by events in the outer domain of the cell, i.e., environment, which is demonstrated by the physics of subatomic unity. Also, the programs themselves depend for their functioning on the finely adapted sequences of structural genes, which evolved to their present state of high efficiency over long stretches of time, should not be underestimated. In some cases, structural genes are also regulator genes; they produce protein molecules which are instrumental in switching other genes on or off. The issue with messenger RNA ordering which is specifically targeted to a type of virus protein, is that a risk of destabilization of the natural DNA sequencing code is probable, which is dependent on all information input i.e. real time communication within and outside of the cell (environment) and can't be switched off, because it is a dynamic, living, thus complex system. (van

---

Neyman and Elizabeth L. Scott, Berkeley: University of California Press, 1972, pp. 277-296. https://doi.org/10.1525/9780520313897-013

[189] Campbell, Jeremy, Grammatical Man. 1982. Simon & Schuster, 319pp.

Campen, Arend, 2021)[190] The synthesis of these findings may indicate that mRNA vaccines are a form of genetic coding, which, due to their linear design, can't function or risk the stability in and of a non-linear biological living system due to the complexity which could lead to potential misunderstanding of message transported by the mRNA. This could lead to destabilized protein production. Because this too is information, it would have to be taken into consideration to clarify the number of excess deaths that are currently being observed by the (ONS) Office for National Statistics in the UK.[191]

Uncertainty is expressed as entropy. What is needed here is to distinguish data from information because they are not the same. Data means a combination of (useful) signals and (useless) noise. This could be a number, or a symbol, code or graphs, whilst information is data put into context. In other words, information absorbs data towards meaning.

Jeremy Campbell described it as follows in 'Grammatical Man'(1982).

*'Information is a word that has never been easy to pin down. In its most familiar sense, information today is news, intelligence, facts and ideas that are needed and passed on as knowledge. But a more active and constructive meaning as something that gives a certain shape or character to matter, or to mind; a force that shapes behaviour, trains, instructs, inspires and guides. Information gives form to the formless, DNA codes are information and form human thought patterns. In this way, information spans the disparate fields of space computing, classical physics, molecular biology and human communication, the evolution of language and the evolution of man. Nature can no longer be seen as matter and energy, but must be interpreted as matter, energy and information.'[192]*

---

[190] Van Campen, Arend, personal website, https://arendvancampen.blogspot.com

[191] Office for National Statistics (ONS). 2023. Deaths Registration.

[192] Campbell, Jeremy. Grammatical Man, op. cit.

Another point of interest to this research is to understand the distinction between bits and binary digits. James V Stone describes them as follows: a binary digit is the value of a binary variable where the value can be either a 0 or a 1 but a binary digit is not information per se. A bit is the amount of information required to choose between two equally probable alternatives (e.g. left/right) whereas a binary digit is the value of binary variable, which can adopt one of two possible values (i.e. 0/1).

In contrast a bit is a definite amount of information. They are different types of entity. So, if a person knows where he is going and is offered an alternative which confirms that he is given a binary digit, but no extra information. On the other hand, if a person has no idea where to go next, and is handed a choice, he is given information in the form of a binary digit and has gained one bit of information. This means that a binary digit gives half a bit of information. In 1948 Claude Shannon said; *'information is the resolution of uncertainty'*, which he demonstrated by developing a method of communicating noiseless information in digital format, which presently allows this research to be typed on this computer. This noiseless, undisturbed communication was made possible by adding redundancy to the transmissions which can be understood as a constraint, an extra ration of predictability limiting entropy. Redundancy reduces error by making certain letters and groups of letters more probable, increasing predictability. It can be understood as information density. (Gatlin, Lila, 1971)[193] A sentence of more letters can be reduced to less letters as long as the coherence stays intact e.g., 'lge lv rm, 2 br, basmt'. These early efforts to improve communications were the foundation of computer compression, now enjoyed by all of us in music, documents transfer or online videos.

---

[193] Gatlin, Lila. Evolutionary Indices, op. cit.

## 7.4 Entropy: Energy Cannot be Created or Destroyed

To understand probability patterns by the use of information, we must explore the notoriously difficult concept of 'entropy' and therefore must return to the fundamentals of physics. Entropy comes from Greek and means 'transformation'.

The process of understanding what information entails starts with the description of the two laws of thermodynamics. Information understood as energy obeys the same laws of physics, the 1st and 2nd law of thermodynamics (Clausius, Rudolph 1865)[194] . The first law is the Law of Conservation of Energy. It states that whilst energy does not alter its quantity, it may lose its quality.[195] The second Law is described as follows:

There are two kinds of processes, heat and work, that can lead to a change in the internal energy of a system. Heat always flows from hot to cold regions. For Clausius entropy was a relation between heat and temperature. When a quantity of heat flows out of a hot body, its entropy decreases by the amount of heat divided by the original temperature. When that same quantity of heat flows into a cool body, its entropy increases by the amount of heat divided by the original temperature of the cool body.

This is the same as saying that any change in the energy of a living system must result in a corresponding change in the energy of the surroundings inside and outside that system. In other words, quantified energy cannot be created or destroyed. Information influences matter and reality because it always is a part of them. In fact, matter, energy and information form one physical reality because every particle, atomic or

Calculate $\Delta S$ for this process

$$\Delta S = k_B \ln \frac{\Omega_2}{\Omega_1}$$

$$\Delta S = \ln 1.38 \times 10^{-23} \ln\left(\frac{16}{1}\right) = 3.83 \times 10^{-23} \text{ J}$$

---

[194] Clausius, Rudolf. 1850, 1865, 1st and 2nd Law of Thermodynamics.

[195] Campbell, Jeremy Grammatical Man, op. cit.

subatomic, which form the fabric of all existence contains information. Information can't therefore not be disconnected nor separated from the structure of reality. (Khalilli, Jim, 2017)[196]

With physical information in the form of feedback as an influencing, but needed energy through cognition, processes that use all relevant information are adaptive whilst processes that not using all relevant information become non adaptive. A quantity of information = a quantity of energy in living systems because they are interdependent. Lacking information means needing energy to maintain and self-organise. Entropy is the measure of order or disorder in a process or living system by the availability or unavailability of energy = information. Information answers to the same laws of physics, namely that natural entropy can be influenced by information which is used by our universe, planet, nature and us (living systems) to create order and structure whilst ensuring sustenance. Information can't be destroyed because it is connected and a part of quantum reality i.e. the fabric of existence. Order depends on a minimal required quantity of energy in the form of information to do work.  We now can work with two laws;

The energy of the universe is a constant.
The entropy of the universe tends to a maximum.

This brings this research into the realm of Ludwig Boltzmann [197] who as an atomist studied the phenomenon of entropy as an increase of disorder among atomic particles, rather than heat. At equilibrium, the state of maximum entropy (think of totally dissolved milk in coffee turning it from black to brown) there is the most disorder there ever will be. Chaos or disorder is easier to create than order, because it takes work to create order. Entropy tends to increase over time. Entropy is the result of maximum mixing of elements, whilst negentropy contradicts this by breaking down entropy in separate parts in order to create order. The

---

[196] Khalilli, Jim, Information Technology, op. cit.
[197] Boltzmann, Ludwig. 1872.

dissolved milk that can't be separated from the coffee anymore is a probability pattern of entropy because the probability of ever returning to order is too complex. Boltzmann designed an equation using his Boltzmann constant which calculates the quantity of entropy $\Delta S$ in joules in living systems. Therefore, as well as information is quantifiable, entropy is also quantifiable, which confirms them as physical substances and the probability of the irreversibility of natural processes. This is a significant 'finding' which allows and substantiates the probable continuation of this research, because the relationship between laws of probability, irreversibility and entropy also apply to anthropogenic conceptions within societies. On his grave stone we find his famous probability equation: **S=K log W** (S=entropy of an ideal gas, K= Boltzmann's Constant, W=the number of microstates related to the macrostate of the gas) For 16 microstates the entropy of the macrostate is calculated as above.[198]

Entropy applies to all living systems. The confirmation that laws of physics consist of information perceived by human beings, acknowledges that William Blake was correct; information is infinite, because the door of entry of, for and as information can only be the human capability of perception through all senses, not just the empirical one or the human brain.

The Universe is a physical system that contains bits of information. Each elementary particle carries bits of information. Electrons carrying information work together in a systematic way to perform a quantum logic operation. For example, a computer and our cell phones operate like the universe because they are part of the universe and to operate, they must obey the same physical laws. Computers and the universe are information processors. To narrow it down, this study will also aim at on how computers can function only when /if sufficient 'input' results in adequate 'output', questions and answers, to find a result. Output de-

---

[198] Source of the illustration: Shawn Shields, Entropy, Micro States.

pends on the quantity and quality of relevant, applicable and used information i.e., the energy needed to run computer programs. Entropy can be described in numerous ways, but for the sake of this research it will be limited to being a measure of uncertainty, using the definitions; 'disorder' (entropy) and comparatively the definition for 'order' i.e., negentropy.

An example is cleaning up a messy room; the information feedback that is available: 'the room is messy' triggers internal and external energy to do something about it and changes the entropy (disorder) into order, but only until natural entropy takes over again and leaves the room going back in entropic state (dissipation of energy i.e. losing energy over time) if meanwhile no information in the form of energy (picking up things, placing things were they belong) is applied. This confirms the second law of thermodynamics that says that the entropy of any isolated system always increases. Isolated systems spontaneously evolve towards thermal equilibrium—the state of maximum entropy of the system. More simply put: the entropy of the universe (the ultimate isolated system) only increases and never decreases, unless information feedback is used to shape negentropy in living systems. Information and energy are profoundly linked.  Another example is a library where the books are only categorized by the colours of their covers. It is highly improbable to find the book one wants to read there. The probability of finding the right book increases if the librarian categorizes the books first by title, second by author, third by genre or ISB Number, because the added information decreases the entropy and restores the negentropy.  Human creation needs energy in the form of information by which beauty, structure and order can be built. Information can never be divorced from the physical world (Khalilli, Jim, 2017).[199] Information is not just an abstraction, it is carried by something, a stone, a cd, or a book, a mind, a quantum particle and can't break the laws of physics. Information can be stored in any

---

[199] Khalilli, Jim, Information Technology.

system to give it unique properties. DNA, genes, neuron synapses, electrons. neuro peptides are operating as information messengers and can also be understood as messages on their own, (Pert, Candace, 2005).[200] Quantum particles, as fundamental building blocks of all matter, are energy which are not separately acting, but form a unity, a unified field with consciousness. The world is a display of information (Hagelin, John, 2007).[201] The tiniest initial conditions such as personal thoughts (information) are therefore elemental building blocks of reality.  Human action and its causality or the relationship between an event (the *cause*) and a second event (the effect) influences reality (Radin, Dean, 1997)[202] (McTaggert, Lynn, 2007).[203]

   This was proposed by Maxwell's Demon (Maxwell, James Clerk, 1871).[204] He envisaged a gas initially at uniform temperature within a box separated into two compartments, where a tiny being was controlling a shutter between the two compartments. By measuring the energy consumed through the feedback and actions by the 'demon' as the tiny being was called later, this consumption led to a net increase in the system's entropy (Szilard, Leo, 1929).[205] He formulated an equivalence between energy and information.  This would explain universal expansion based on the same Maxwell type demon converting information into

---

[200] Candace, Pert. 2006. Documentary film: "What the Bleep! Down the Rabbit Hole" Available from: http://www.youtube.com/watch?v=usMsTPg-hHk, (Accessed January, 12, 2009)

[201] Hagelin, John. 2011. "Is Consciousness the Unified Field",  YouTube video, available from: https://www.youtube.com/watch?v=xjNjxDtLOjk

[202] Radin, Dean. 2009. "The Conscious Universe Insight in PSI" , Harper Collins, New York, 2009, ISBN 978-0-06-177899-5.

[203] McTaggert, Lynne. 2007. "The Intention Experiment", Published by Harper Element, London, 2007, 30p. ISBN 978-0-00-719459-9.

[204] Maxwell, James, Maxwell's Demon.

[205] Szilard, Leo. 1929. Equivalence between Information and Energy' (Online) https://physicsworld.com/a/information-converted-to-energy/  (Accessed,  May 2020)

energy (Faus, Antonio Alfonso, 2013).[206] This was later confirmed by an experiment called the spiral staircase by (Toyabe et al, 2010),[207] by which they showed that they could convert the equivalent of one bit of information to a value of energy (kT). It suggests a new fundamental principle of an 'information to heat' engine that converts information to energy by feedback control. Information as an equivalent of energy can be also understood by the following rhetorical question: Will a man be able to safely cross a busy road by not transforming the observed information in energetic action about traffic coming from the right, but only using the information as energic action about traffic coming from the left? Information can be understood as energy which (in)forms experienced reality. The human capability to perceive through observation determines the observed as information 'if the doors of perception were cleansed, everything would appear to man as it is; infinite' (Blake, William, around 1800). When uncertainty is reduced, we gain information.

## 7.5 Quantum Information Processing as in Universe

Quantum computing is currently being developed to be made possible by information; the universe already works that way (Lloyd, Seth, 2016).[208] Quantum information processing analyses the universe in terms of information: The universe consists not only of photons, electrons, neutrinos and quarks, but also quantum bits or qubits. Professor Lloyd says the universe is a giant computer, processing information in

---

[206] Faus, Antonio Alphonso. 2013. Fundamental Principle of Information to Energy Conversion (online), https://arxiv.org/ftp/arxiv/papers/1401/1401.6052.pdf (Accessed May, 16, 2020).
[207] Toyabe et al. 2010. Experimental demonstration of information-to-energy conversion and validation of the generalized Jarzynski equality, https://www.nature.com/articles/nphys1821 (Accessed 18.05.2020).
[208] Lloyd, Seth, Quantum Information Science, http://web.mit.edu/2.111/www/ notes09/spring.pdf (Accesses June 02, 2020).

quantum bits (qubits). In an article by Brian Siegelwax[209] in Quantum Computing Business called '*Inside Turing*', he reports on a Start-up by Prof. Seth Lloyd and Dr. Michele Reilly called Turing. Their aim is to solve societally useful problems by quantum computing. The hardware has been designed, but the software to run quantum logic programs is in development. They are talking about RAM in the form of qRAM- Quantum Random Access Memory and Artificial Intelligence to build large scale fault tolerant quantum computers that run quantum algorithms in practice. In other words: to be able to store entanglement, but what is entanglement? This is explained in an article by Caltech[210]: *'what is entanglement and why is it important?'* as follows: two particles such as photons or electrons are unbreakably correlated even when they are separated over large distances and is an emergent property.

This entanglement – relationship or interconnectedness can be quantified as *information*. What is important to know for this research is that Albert Einstein in the 1930s said that the observer can only learn about this hidden information until the measurements were made. In the context of this Project, man observes a quantum spin of particles only when they spin. This confirms the observer phenomenon which states that objects, in this case electrons or photons, provide information as soon as an observer observes. This is logical, because human perception is the first door to understanding. The dependent relationship with a quantum universe on the human perception is key. The observed 'information' enters a person's experience in the form of intrinsic communication. If that person would be unable to observe, the quantum entanglement won't take place. Entanglement can be used to transmit quantum information, for quantum games, as well as strategies in a way that is beyond

---

[209] Siegelwax, Brian, Lloyd, Seth and Michelle Reilly Inside Turing (online) https://thequantuminsider.com/2022/04/14/inside-turing-a-startup-by-prof-seth-lloyd-and-dr-michele-reilly/ (Accessed Jan, 09, 2023).

[210] Caltech. What is entanglement, from https://scienceexchange.caltech.edu/topics/quantum-science-explained/entanglement (Accessed Jan. 09, 2023).

what we can do in our classical world state Dr. Xie Chen and Dr. John Preskill[211], Professors of Theoretical Physics at Caltech. John Preskill explains the value of entanglement by taking a book as a metaphor in a video embedded in this article. If information is written on each page you can read the book page by page. If however the book is a quantum book and the pages are highly entangled with one another there is a lot of information in the book but not stored on the individual pages. The information will be stored in the correlation among the pages. The essence he says, that one can have a quantum system which stores a lot of information but if you look at the part of the system, you can't see the information at all.

James Gleick in his book The Information (2011) explains Quantum Computing as simulating quantum physics with a computer accessing every quantum variable. The number of information bits can then become the same as the number of points, i.e., selected mathematical objects and selected relationships in space, re: the universe. Quantum computing would enable to calculate probabilities, making use of the entanglement of qbits, not just multiplying computing power but exponentially increasing it. Confirming Seth Lloyd's goal of solving previously unsolvable complex problems could be made possible by quantum computing because then they can make use of the corroboration of entangled particles, dramatically increasing computer power through correlation of particles in the universe. A point of concern would be as can be derived from the context of this research project, is that a quantum linear approach can't work in a non linear quantum universe. Quantum computing is not yet available at the time of this writing.

At Casio a 2kb RAM card was used to store just enough information enabling BASIC programming up to 2kb – 2000 bytes of random memory to be stored. Information Storage capacity has been exponentially increased during the last 40 years. Quantum information storage

---

[211] Video: use of entanglement

would increase the memory storage capacities and related computing abilities, but only when all 'relevant' INPUT in the form information, and maximum number of questions (re INPUT) are asked. If our synthesis of computing is correct and functionality depends on all relevant information in the form of quantum information, it corroborates our two metaphysical equations. Information is perceived thru the observation of a person (consciousness and all senses) and processed as a real time physical reality (the particles are correlated). But only if all relevant, including non linear quantity and quality of information is processed, functionality of that system can be ascertained. If an error message occurs, more or different information is needed to ensure or restore functionality. This information forms the energy a quantum system runs on. Disorder or entropy (non or semi functionality) will be the result of an information deficit in all computing machines. Because quantum dynamics results will be probabilistic, rather than determinable, this confirms that living systems are in a constant flux of being and becoming. Extracting information can only be done by human observation which influences quantum outcome and determination. Artificial Intelligent computing systems can determine probability, but the final reality can only be determined by the interaction of man's perception thru consciousness and quantum reality by the collapse of the wave function, although this is becoming a moot statement. Collapsing the wave function means than an observer, thru measurement, alters the outcome of quantum experiments, confirming that reality rests on human observation to be determined. Later research finds that this process of quantum collapse can also occur without human presence in the universe. But what is important for this research is that this measurement process or the acceptance of unmanned quantum wave collapses are both information which again enters the doors of human perception as either fact or probability. The information collected with or without having been there to measure wave collapse, does not mean that the information is gone. Both can be used to take maximize the number of choices of reality and are to

be therefore considered as dynamics of existence, constantly to be steered by feedback.

Increasing memory storage capacity in qBits or qRAM allows for more computing power, but only if all relevant INPUT is allowed be used can harmless functionality be expected. An important, not to be overlooked issue in these matters is that reductionistic, linear programming, even in an unlimited RAM capacity, can't control non linear, complex living systems. Can computers that are linearly programmed solve non linear societal issues? This is a rhetorical question which at present must be answered by 'not yet' due to complexity. Societal, complex problems can only be solved when the first cause of them is eliminated.

Dr. John Hagelin explains that this unified field, or Grand Unified Field Theory is not separable from man, it is man's innermost self which he calls an ocean of existence or intelligence as the basis of the nature and diversity of the universe. This is not only a philosophical theory but is based on mathematics and physics. Quantum mechanics demonstrates that the deeper one goes into the understanding of the fundamental building blocks of living systems (at subatomic level), the more dynamic it becomes *(waves). Human consciousness completes the unified field- or field of intelligence by its natural connectedness, dependence and the human ability to perceive the difference between right or wrong, yes or no i.e., 1 or 0, just like computers can do. The difference however between a linear computer and a non linear person in a non linear universe is formed by human consciousness. Both the universe and a person are conscious and a computer can't ever be, because it lacks self-reflection and the Zoa's, i.e., the senses beyond brain and heart: imagination, feeling, sensation, intuition. They are one total consciousness and non-material, local and non local at the same time. (Hagelin, John 2011)[212] This synthesized phenomenon confirms Gnosis, which is Greek for

---

[212] Hagelin, John, Is consciousness the unified field, op. cit.

knowing i.e., epistemology, reconnecting the Chinese TAO with physics. When we read in the 'Asclepius' by Hermes Trismegistus translated by Dr. Gilles Quispel[213] (1996) : 'A great miracle, o Asclepius, is man'. This inspired the orator / philosopher Pico Della Mirandola to write his lecture which he wanted to deliver in Rome by addressing the Pope; 'De Hominis Dignitate' (1486) (About the dignity of man). He was not allowed to express these 'blasphemous thoughts in Rome. His lecture is, until today,  regarded as the Magna Carta of Renaissance and Humanism. What he did, was to search for, contemplate and allow for all the relevant information, philosophy, science, religion, etcetera, he was able to find into his vision of the status of an inhuman society in the 14th century.  The reason why this research is mentioning this is that the part of humanity, the value and necessity of human reasoning capability can be considered as intrinsically connected to the universe as both are messengers, distributors and receivers of information. Gnosis states and confirms the spiritual relationship in the form of inner knowing between man and the universe, which is demonstrable as a human perception and should also be used as relevant information.

Modern physicists like Fritjof Capra [214]have been combining philosophy and science because he understands and accepts that they cannot be separated.  Capra arrived at a next level of knowledge about quantum mechanics or quantum theory.  Via quantum mechanics and quantum theory, physicists try to understand the universe as a whole by connecting actions of the smallest pieces of an electron to the very actions of the largest objects in the cosmos such as black holes.  This includes the interaction and the interdependence of human beings with the cosmos. People are a part of a whole universe and have influence by conscious and subconscious processes such as action and observation on matter and therefore determine, for a large part within the realm they can con-

---

[213] Quispel, Gilles. 1996. The Asclepius, Bibliotheca Philosophia Hermetica, Amsterdam. ISBN 90-71608-07-7

[214] Capra, Fritjof, Tao of Physics, op. cit.

trol or steer, their experienced reality themselves.  Cosmos is the patterned whole of all existence.

Capra understood that a mutual interrelation of all things and events and the experience of all phenomena are manifestations of a basic oneness. Confirming unity and the unified field. He said that in an ordinary life people are not aware of this unity, but divide the world into separate events and objects. (the perception problem).  Capra claims that concepts of separate things and events are an illusion.  His work on atomic physics shows that subatomic particles can only be understood as interconnections between the preparation of an experiment and the measurement.

## 7.6 Steering Complexity with Cybernetics

### 7.6.1 Complexity

What is complexity and why is it so important?  Complexity is a consequence of a myriad of unknown factors, relationships and reliance which nevertheless have a direct impact on organisations and can only be 'controlled' by learning and actual adaptation in real time. These so-called non linear effects or non causal stimuli from often unknown, but interconnected and interdependent sources, can't be controlled by old-fashioned linear measures of control such as regulations or compliance, but can be 'navigated' by the use of positive and negative feedback, i.e. amplifying and correcting information. Cybernetics is needed here to make sure all information is used. To steer a living system, rather than trying to regulate it, it is wise to use the Law of Requisite Variety, a.k.a. Ashby's Law, to build resilience by obtaining variety in the form of capacity, capability, knowledge, tools, information enabling to respond to risk. When only linear control systems, such as law enforcement, empirical verification, physics, mathematics or iso standards are used, they won't be able to cover a non-linear or complex reality which consists of a linear (cause & effect) part, but for a larger extent, of a non-linear, non causal or dynamic part, which is much greater. This can't be

known in advance, for example people's behaviour or preception, because it constantly changes due to unforeseen and unknowable actions by the entire, also dynamic, network of which it is a part and on which all forms of organisation are dependent. In short; behaviour of the parts in a network can't be predicted. Awareness of such behaviour is needed but needs a new, systemic point of view. Linear control systems are unable to regulate living systems because they lack information. Result: entropy.

It is crucial for human and non human survival that we understand complexity and learn how to deal with it. Control systems we have tried such as bank regulation, industrial compliance demand, trade tariffs, ecological restriction, lockdowns, law or even wars, have not been able to maintain equilibrium or homeostasis just like nature does, because they attempt to escape the laws of thermodynamics and of entropy, as they can't have enough energy (information deficit) to 'steer' or 'navigate'. Quite the opposite has been caused by them; economic, cultural or environmental collapse, poverty, mass emigration, crime, biological annihilation of species upon we all depend, pollution of soil, water and air, all of them can be listed as entropy (disorder), because the same cause applies: information deficit.

Nature works according to cybernetical principles and cognition. All we have to do is to learn from and copy nature because we are a part of and dependent on it. Stephen Hawking (2006)[215] was looking for a Theory of Everything. This hypothesis on which this research is based is that 'Everything is Information'. The language of the universe is information.

When an organisation, corporation, industry or political process are understood as living systems, their interconnected and interdependent relationships within and outside of their networks change all the time. This complexity, due to dynamic and perpetually changing information feedback loops, can't be managed by direct causal action (cause and

---

[215] Stephen Hawking. A Theory of Everything, op. cit.

effect reactions) but can only be maximally controlled by steering with that information in real time. Adaptability to constantly evolving changes by using all information directly is needed

### 7.6.2 Cybernetics

Kuber (Cyber) is the Greek word for helmsman, the person steering the ship. Keeping a ship on course, can only be done, that is if the ship intends to reach its destination, by information in the form of feedback. It also translates in 'governance' based on its Greek origin. This corresponds with our earlier findings that information obeys the universal laws of physics and therefore information during a sea voyage about waves, tide or wind, is real and can't be escaped by ignoring it, because if it is, the ship will probably run aground, sink and not reach its port. It is therefore of importance that information is acknowledged as the energy which makes things work and that universal energy adheres to the laws of thermodynamics as well as Shannon and Boltzmann's concepts of entropy. Information is part of the construct of reality which can be observed, interpreted by human senses, but not escaped from. Cybernetics is about setting a goal, e.g., sustainability or long term continuity of living systems and the method and requirements of how to achieve it.

Norbert Wiener understood that the human body and as he called it then; *computing machines* worked in the same way; they both depended and used information, input and output, to operate and function, but that the value of information always depended on the perception of human beings before it could be put into practise in computing machines. He knew that human bodies were automatic learners, but also saw that the information learned would alter behaviour, and that sense organs as information absorbers, were limiting the communications within and among individuals. He wrote in his book The Human Use of Human Beings (1957); *'Like any form of information, these commands are subject to disorganisation in transit. They generally come through in less coherent fashion and certainly not more coherently than they were sent.*

*In control and communication we are always fighting nature's tendency to degrade the organised and to destroy the meaningful; the tendency for entropy to increase'.*  This relates to earlier observations about entropy by confirming his research and acknowledgement of natural entropy to exist. The process of giving and receiving information was in his eyes the natural adjustment of living systems to live effectively within their environment. Condition is that equilibrium and balance is sought instead of enforced growth by power by way of ulterior motivation for example: economic gain, political clout, greed, which become dependent on ignoring information. This contradicts stability by its disturbance of equilibrium by increased entropy. The above suggests that ethics is information, confirmed by innate, autopoietic human protection and survival mechanisms in the form of conscience. Without ethics i.e.', truth or goodwill, living systems are vulnerable and can become non viable. It is conversation and information, as well as communication and cooperation, a.k.a. circular causal and feedback mechanisms which measure the degree of entropy. The importance of conversation is to have access to a maximum number of choices by ethical imperative. (Von Foerster, Heinz 1973).[216] These are the foundations for it. (Pangaro, Paul, 2019).[217] Maximisation of choices as a condition for risk management and stability. This was researched by Dr. Ross Ashby who developed the Law of Requisite Variety, a.k.a. Ashby's Law; Only variety absorbs variety.

---

[216] Von Foerster, Heinz. 1973. On Constructing a Reality, in Environmental Design and Research, (Online) http://www.semiorganized.com/resources/other/Foerster-constructingreality.pdf (Accessed May 25, 2020)

[217] Pangaro, Paul. 2019. Introduction to Cybernetics and the Design of Systems' (online)    https://www.pangaro.com/design-is/Cybernetics-minimized-v8b.pdf (Accessed May 2019).

## 7.7 Ethics is Information

This first project ends in an observation derived from writing this paper that a renaissance of information theory is needed to deal with complex issues we, as a species, endured and continue enduring. The findings could alleviate the suffering experienced in nature and society. I will add a listing of important extras. It starts with the interpretations of reality and the Value of Philosophy by Bertrand Russell who wrote, and this lists the first causes seamlessly of and for non or harmful functionality of anthropogenic living systems:   *'Unfamiliar possibilities are contemptuously rejected'* (Bertrand Russell, The Value of Philosophy,1912)[218] due to a crisis of perception, includes the consensus of pursuing an Empirical Reality – by the 5 senses; touch, taste, smell, sight, sound, which is proven to be insufficient.  Because this research demonstrates that all relevant information should always be use, we can reduce uncertainty and entropy by re-designing or correcting currently harmful conceptions by the elegant solution of adding information, as follows:

Total Reality – Empirical + 4 Zoa's; imagination, intuition, feeling and sensation, heart and brain and universal consciousness; this would be a total reality based on *all information*, not just a part. A holistic, systemic, synthetic reality, including gnosis, consciousness, Jung's archetypes, noetic pre-cognition, quantum particles, intention, entropy, which in fact are all information, but are dependent on human perception.

- Accepting that natural limitations of reality are the boundaries of and for non harmful functionality

---

[218] Russell, Bertrand. 1969. "The problems of philosophy, the value of Philosophy". PDF, (online) available from http://www.skepdic.com/russell.html, (Accessed January 12, 2012).

- Information can't be separated from our physical reality because it is a fundamental entity which can be measured and quantified in Bits of Information.
- Information is not just a message, it is physical and 'in forms' our lived and experienced reality, whilst safeguarding existence towards survival.
- Information deficit or shortage e.g., fake news, misinformation, disinformation, suppression of information by censorship, overload of information; they all can be categorized as deception and an attempt to escape reality. This is not allowed by physics and causes entropy.
- Per definition (symbol ≜ ). This is a law of physics confirmed by information theory:  Information Deficit ≜ Entropy (Disorder).
- Goals that are enforced in a linear manner can't be attained because nature can't be forced by man without consequences to the enforcers. An information deficit will prevent such goals to be ever achievable. An enforcement of them will cause great harm to life, the environment and social cohesion (as we observe now).
- Complexity can't be regulated by enforced rules, restrictions or laws without causing destabilization in the form of entropy (disorder). E.g., the weather can't be regulated or controlled by man without destabilizing weather patterns of which outcomes are unpredictable. (i.e., Butterfly effect).
- The outcome of armed conflict is always entropic (causing disorder) in the form of harm.

# 8

# THE AFRICAN VALUE OF UBUNTU: FOR AN ETHICS OF GLOBAL ARTIFICIAL INTELLIGENCE INCLUSION

*Arthur Gwagwa, Netherland / Emre Kazim, UK*
*Airlie Hilliard, UK[219]*

## Summary

Historically, Sub-Saharan Africa (SSA) has been excluded from the benefits of the previous industrial revolutions, as its people and their resources and aspirations have been objectified through foreign domination, and its culture has either been fragmented or appropriated. While artificial intelligence (AI) is poised to generate vast amounts of wealth, its application could lead to further social and economic exclusion of SSA due to a lack of access to technological advancements and the his-

---

torical injustice and exclusion based on protected characteristics. Through an examination of the concept of inclusion, this paper explores how to improve the terms on which African populations and subpopulations and their concerns are included in the global AI ethics discourses. Specifically, it is argued that the SSA value of Ubuntu could be of immense value in AI applied normative ethics, particularly toward an inclusive approach for the implementation of the universal AI ethics principles and guidelines.

THE BIGGER PICTURE: AI's social and economic benefits and its challenges to the African social and cultural perspectives are beginning to catch the attention of African policy makers if judged by the proliferation of AI think tanks and empirically grounded policy recommendations. However, unlike other regional blocks, like the EU, unified African positions and approaches in global AI ethics forums remain sparse. Although the current universal ethical guidelines and principles can provide Africa with a common ground with other cultures, care is needed in localizing these principles, as they may not be applicable in the African context. The global ethics discourse can capitalize on the emerging African ethical relational approaches, underpinned by Ubuntu, to devise frameworks that would assist the implementation of the universal values, such as justice and solidarity, in a manner that pays regard to cultural environments of historically marginalized populations, like in Africa. Future research and policy work should ideally focus on coming up with Ubuntu-based action guiding principles for all AI stakeholders. These could, for example, include guidance on reconciling competing and often conflicting cultural values and ethical dilemmas in AI design, development, and audits. This and related works can only have impact if forums are created for interdisciplinary discourses between policy makers, technologists, ethicists, and philosophers to ensure that the African context is being considered in their work. [220]

---

[220] Remark of the Editors: Ubuntu Ethics for AI is also the topic of the interesting paper: Sabelo Mhlambi, *From rationality to Relationality: Ubuntu as an*

*(Endnotes refer to the references at the end of the article.)*

## 8.1 Introduction

While there is a consensus about the enormous potential for artificial intelligence (AI) to advance development and solve some of the most pressing challenges faced by Sub-Saharan Africa (SSA), discussions of the ethical challenges that AI will bring to Africa have only just begun. Little has been done to advance unified African positions and approaches in global AI ethics forums.

This is despite the rise in recent literature on how one might apply While there is a consensus about the enormous potential for arti-AI to resolve problems in Africa and on ethical issues facing artificial intelligence (AI) to advance development and solve some of application to Africa, particularly ''the need to define African the most pressing challenges faced by Sub-Saharan Africa values and align AI with them.''[1] Nevertheless, some think tanks (SSA), discussions of the ethical challenges that AI will bring to have emerged and are producing empirically grounded policy Africa have only just begun. Little has been done to advance uni-recommendations. However, while policy has a role to play, here is very little examination of fundamental issues relating to the values underpinning such policies and, in particular, how to address the AI risks and challenges that may be more acute in the Global South, where the low access to AI technology could lead to exclusion, particularly in SSA. In the past decade, there has been an emergence of notable works with mainstream African ethical approaches, with some addressing the need for African relational approaches in addressing AI algorithms injustices,[2] balancing relational approaches with autonomy,[3] and explicability.[4] By building on these emerging views, this paper argues that at the heart of Ubuntu are principles that prescribe the virtues needed, procedures,

and the desired consequences in the application of universal AI ethical principles. This would lead in the systematic integration of the universal AI principles and an inclusive deployment of AI technologies. By seeing one's humanity in the humanity of others, Ubuntu resonates with the golden rule that cuts across major world cultures: We should do to others what we would want others to do to us. Further, relying on a value from SSA—generally the world's most economically disadvantaged region—would be of a practical and symbolic benefit use toward the greater inclusion of SSA. While Ubuntu's relational approach based on communitarianism is not unique to Africa, it would be of practical and symbolic benefit toward greater inclusion of SSA in AI ethics discourse and the economic and social benefits resulting from AI, particularly because it widely informs most African subcultures and looms large in the SSA philosophy and ethics.

SSA's exclusion resulting from the deployment of AI has the potential to both perpetuate and amplify the deep-rooted exclusion of Africans for three key reasons, with exclusion referring to "the inability to participate effectively in economic, social, political, and cultural life, and, in some characterizations, alienation and distance from the mainstream society."[6] Firstly, AI can amplify or reinforce long-standing societal biases, particularly those related to characteristics protected under international human rights law, such as race and culture. Secondly, as Africans can lack the capacity to access and apply their data, they are less able to develop and implement AI and so miss out on the economic benefits it can bring. Finally, since it is predominantly the perspective of the Global North that is informing the current discussions on inclusion, in certain circumstances, this has resulted in weak commitment to addressing historical social and economic injustices. While a plethora of guidelines on ethical or responsible usage of AI is emerging, each promotes different values and definitions,[7] meaning that care is needed when drawing on generic principles that may or may not be universal in scope.[4] This includes paying attention to the social, cultural, and local

values of the region in which these principles are being applied; Africa has historically seen misaligned foreign values imposed on it as a result of these factors not being considered.[7]

In this paper, we therefore argue that the relational SSA philosophy of Ubuntu, which emphasizes one's personhood to the personhood of others, could be of both normative and applied practical value toward the realization of the current corpus of principles and guidelines on ethical AI. As shall be demonstrated below, the values that are being currently appealed to in AI ethical discussions, such as solidarity and those based on medical ethics,[8] including autonomy, justice, beneficence, and non-maleficence, can only bring utility if there are generally agreed upon adequate implementation strategies.[9] The multiple cultural contexts in which AI is applied may present a barrier in the even application of these principles and guidelines in a manner that ensures an equitable distribution of AI benefits across the globe. Adopting the value of Ubuntu does not just acknowledge a unique contribution by Africans to general philosophy and applied ethics but gives guidance on the virtues, procedures, and desired consequences toward an inclusive and ethical AI. As an example, Ubuntu reinforces the universal value of solidarity to the extent that it proposes communal relations based on generosity, hospitality, compassion, and friendliness.[10,11] We argue that such characteristics of Ubuntu values are at the center of what it means to be human in a world with competing and often conflicting cultural values.

## 8.2 AI Ethics and Exclusion Challenges

Historically, SSA has been excluded from the benefits of the previous industrial revolutions, as its people and their resources and aspirations have been objectified through slavery, colonialism, imperialism, and neo-colonialism. While the slave trade was meant to exploit African resources to feed the ever-expanding European markets, today, African resources are again powering the Fourth Industrial Revolution. From the

Congolese cobalt crucial for the manufacturing of computer chips to the data that are being used to train AI algorithms, African resources are significantly shaping the future of AI. However, just like in the previous industrial revolutions, African voices are absent from shaping the future of these developments. As the Fourth Industrial Revolution progresses, it is therefore important to assess the extent to which SSA, in all its diversity and similarities, is being included in the discussions and benefiting from the outcomes of the various social, economic, and political systems and processes underpinning the current changes. In this context, inclusion refers to the process of improving the terms on which individuals and groups can take part in society and the ability, opportunity, and dignity of those disadvantaged on the basis of their identity.[12] The United Nations has emphasized the importance of inclusion in a number of their sustainable development goals,[13] claiming that this systematic process can rescue a person or community from the risks or uncertainty of exclusion.

## 8.2.1 Exclusion at the Continental Level

While SSA is made up of a diverse range of countries, they share broad similarities, like their history; aspirations, which are mostly shaped by the liberation wars, past political junctures, and trajectories; and a broadly similar communitarian cultural value system ensuring an appropriate ethical and legal framework to strengthen African values. An additional attribute shared by these countries is that they historically have not benefited or have been excluded from the benefits of the previous industrial revolutions. With the onset of the Fourth Industrial Revolution mostly underpinned by AI, Africans may be excluded from the benefits of AI on the grounds of natural characteristics or protected attributes, including color, language, culture, or race, as a result of the limited or unrepresentative African datasets available for the proper training and application of algorithms or AI applications, like facial recognition software. Since the AI field is mainly composed of white

males, this lack of diversity and inclusion has already resulted in flawed systems that amplify gender and racial biases, according to a survey carried by the AI Now Institute, which examines the social implications of artificial intelligence.[14] ''The media is filled with unintended ethical concerns of AI algorithms, such as image recognition algorithms not recognizing persons of color or racist algorithmic predictions of whether offenders will recidivate.''[7] Calls to correct anomalies and flawed systems have sometimes been received unkindly by technology firms, as was evident in the dismissal of Timnit Gebru, co-leader of Google's Ethical AI team, who surfaced the dangers of large language models like the ones that power the company's search engine.[15]

Given that AI stands to generate vast wealth for the corporations and countries that develop it, the rest of the world could be left behind if they are excluded from the social, cultural, and economical benefits of AI.[16] It is, therefore, evident that there needs to be an effort toward greater inclusion in this domain, particularly since the Global North lacks the insight needed to create solidarity in these advancements. First, this is due to the disconnect between the algorithm designers and the communities where the research is conducted or algorithms are implemented.[7] Secondly, governance, including in the AI domain, in liberal democracies of the Global North is mainly focused on protecting autonomy within the individual private sphere.[17] This is a typical Western worldview that centralizes the individual and which is reflected in bioethical principles, like the principle of respect for autonomy, frequently understood as respecting the decisional autonomy of an individual who makes decisions without undue coercion.[18] Consequently, there is an acute need for increased and organic interactions among intellectuals globally to facilitate the expansion of this discourse beyond the Western world, particularly because the reality of global exclusion is felt most in the developing world.[19]

## 8.2.2 Exclusion at the National Level

Global corporations, including those working on technology and data, are involved in data-mining activities in Africa that are not just amplifying existing societal tensions but also excluding African subpopulations who represent low-value data.[20] This exclusion is also seen in the uneven access to data, AI, and related technologies, as well as the impact of these tools, which is greatest in marginalized populations.[21] This impact is particularly felt in the least developed countries, who sit at the intersection of these marginalized groups, resulting in the amplification of these digital inequalities across the world. Non-representative or biased data can further entrench existing inequities as AI systems reflect the biases and lack of representation of the datasets on which they are trained, resulting in the exacerbation of the long-standing societal biases that exist surrounding protected characteristics, like race.[22] Data are expensive and hard to come by at scale, but the data that are available encompass three broad groups of people: the uncounted who do not exist because they are not included in any sort of database; the unaccounted who have less inclusion into the digital world and therefore not entirely represented, maybe due to economic reasons; and the discounted who are in databases but are not of OPEN ACCESS interest to the people who would serve them, such as governments or companies, because they do not have enough money to be of concern. AI algorithms are trained on the data that are available, as opposed to complete datasets, and these data can easily privilege socio-economically advantaged populations who have greater access to connected devices and online services.[23] As a result, the populations who do not have this access are often forgotten and the gap between developed and undeveloped countries widens. Therefore, initiatives are needed to increase the fairness and representativeness of data and algorithms and an examination of the values that they embody to facilitate greater inclusion. In support of this, African scholars are beginning to explore sociological approaches that go above

and beyond technical solutions by placing ethics in their ''relational'' context[2] and how to reconcile relational approaches with autonomy.[3]

## 8.3 Do Current Initiatives Embody African Values?

There has been a steady increase in the number of global and regional AI ethics initiatives that have by-and-large been aimed at addressing the kinds of exclusions discussed in AI ethics and exclusion challenges. It is also common for them to include the rights of persons at risk of exclusion, improving the individual and collective wellbeing and dignity of these people and allowing them to flourish.[24] As an example, in their comprehensive map of the corpus of principles and guidelines on ethical AI, Jobin et al.[9] reveal a global convergence emerging around five ethical principles (transparency, justice and fairness, non-maleficence, responsibility, and privacy). Nevertheless, they also report a "substantive divergence in relation to how these principles are interpreted; why they are deemed important; what issue, domain or actors they pertain to; and how they should be implemented," thus highlighting the importance of the applicability and the question of implementation of these principles in different contexts.[9] Carman and Rossman[4] call attention to the need for care when drawing on generic principles that may or may not be universal in scope, including by paying attention to the cultural context, especially in post-colonial Africa, given its history of the imposition of external values. Despite claims of universality, most AI ethics principles and their guidelines are developed by stakeholders based in economically developed, mostly Western countries, like the United States and from within the European Union.[9] As a result, some aspects of the principles may not automatically apply in Africa without the necessary adjustments. For example, the principle of respect for autonomy may be incompatible with the African communitarian approach to decision making.[4] Yet a common ground can be found if it is based on the idea of personhood in African traditions, which imply "a relational and positive

sense of autonomy, which involves the community helping or guiding one to use one's ability and knowledge of one's social relations and circumstance to choose freely the requisite goods for achieving one's life plan."[3]

The current exclusion of Africa, including its ethical approaches to AI governance, whether intentional or unintentional, means the inclusion debate is still framed from the perspective of the Global North, who developed the technology in accordance with Western perspectives, values, and interests with little regulation or critical scrutiny.[25] As African and South-American countries are not represented independently from the international or supra-national organizations that are producing these guidelines, this may present a barrier to implementation of such guidelines but also the deployment of the AI technologies in specific sectors, such as agriculture, where, for example, excessive automation may disrupt the African way of life that revolve around certain customs.[26]

In addition, the private-sector companies from the developed countries have been involved in the AI-ethics arena, thus raising concerns that they may potentially use such high-level soft policy as a portmanteau to either render a social problem technical or to eschew regulation altogether.[9] Given the non-inclusion of stakeholders from Africa and South America, the convergence of AI ethics set of principles on the four classic principles of medical ethics, namely autonomy, justice, beneficence, and non-maleficence, will not address Africa's concerns about inclusion, as the implementation of these high-level principles can conceal deep political and normative disagreement, which could have unwanted effects on the future of AI development and governance.[27] As an example, while the European AI4People's recent publication[28] interprets justice to include using AI to right previous wrongs, ensuring that the benefits of AI are shared fairly, the wealth from AI still benefits a few developed countries that unfairly benefitted from the previous industrial revolutions. The justice articulated in the ethics discourse should be

accompanied by implementation guidelines on how to specifically include historically marginalized populations whose resources were used to power the previous and continue to power the current industrial revolutions. As suggested by the United Nations Educational, Scientific and Cultural Organization (UNESCO), the global AI ethics initiatives should frame Africa as a cross-cutting concern.[29] Ethics do not just influence human decisions on what is right or wrong but constitutes the basis of future action, and in the case of AI, it will influence the course of the Fourth Industrial Revolution. Under such circumstances, a relational approach to ethics may be more sensitive to the African cultural context, since it advances the notion of inclusion. Best practices toward inclusion can be seen in other cultures and domains, such as the way Canada acknowledges the historical injustices to the First Nations, particularly in land ownership. Similarly, Africans should explicitly be asked how they want to be included in the revolution and on what terms. Colonization dispossessed Africans of more than resources and self-governance; it also took their voice, ability to self-determine, collective agency — the ability to negotiate with a unified voice, and, in some instances, appropriated African culture.[30] Ironically, the colonizers did not appropriate the essence of African culture captured in Ubuntu, but inclusion of this value in the Fourth Industrial Revolution would be an important step in implementing the converging global AI substantive values.

## 8.4 Reshaping the Western Concept of Inclusion

So far, much of the literature and research on social exclusion is underpinned by frameworks that are concerned with European and Anglo-Saxon traditions. As such, they ignore the contributions made by people of Africa, Asia, and Latin America, where global exclusion is more likely to be felt. A second challenge is the marked absence of any discussion on power embedded in social relations and the disruption of relationships between individuals and society.[12] Consequently, the Euro-

pean and Western model for inclusion in AI and technology in general should be rearticulated to draw on input from the Global South and create a more developmental focus on global inequalities.

The future of the inclusion debate will depend on the ability to develop a global inclusion initiative that draws on the intellectual capacities of both the Global North and the Global South.[12] Specifically, Africans should define what inclusion means to them and how it can be achieved, since there is only a tepid commitment to addressing historical injustices, like how African people and their resources and aspirations have been objectified through slavery, colonialism, imperialism, and neo-colonialism. These injustices are still relevant in the AI era, which is creating new domination capabilities and novel problems; while traditional colonialism is driven by political and government forces, algorithmic colonialism is driven by corporate agendas.[25] In Kenya, for example, AI and data-optimization technologies are exploiting existing ethnic and racial tensions, particularly during election times, through computational hate propaganda and disinformation.[31] These technologies are undermining the basic values of African societies, such as community, but also the concepts that are characteristic of African normative ethical thinking, including harmony, consensus, collective action, and common good. The effectiveness of this discourse could be maximized if it were adapted to cultural or country-specific situations where codes could potentially have policy relevance.[12] In addition, as Timnit Gebru and her colleagues attempted to champion, there should also be more comprehensive action against racism, sexism, and other forms of socially constructed exclusions, something which has been lacking in past discourses but is beginning to emerge in the Africa AI decolonization movement.[25]

## 8.5 Emerging African Views

Although AI ethics guidelines and principles and their accompanying industrial codes of ethics and toolkits are a good starting point, they alone cannot resolve the disparities highlighted above without respectful and honest dialogue between the two hemispheres to address the historical disadvantage and value misalignment whereby AI reflects Western values, agendas, and motives. So far, the idea or willingness to find universal principles is neither healthy nor efficient, given the exclusion of billions of people from participating in the framing of these principles that will affect them and their future generations. However, efforts toward these dialogues have already proven useful, particularly the workshops of the UN Global Pulse, which were held in Ghana and Tunisia.[32] From these workshops emerged a unanimous consensus that Africa could learn from the Global North's mistakes to ensure that they do not develop technologies without first formulating a set of values to guide them. In addition, Africans advocated for the need for human control of technology and the promotion of human values, something, which has been reactionary rather than proactive in global principles.[33]

## References

Metz, T. (2021). African reasons why artificial intelligence should not maximize utility. In African values, ethics, and technology: Questions, issues, and approaches, B. Okyere-Manu, ed. (Palgrave Macmillan), pp. 55–72. https://doi.org/10.1007/978-3-030-70550-3.

Birhane, A. (2021). Algorithmic injustice: a relational ethics approach. Patterns *2*, 100205. https://doi.org/10.1016/j.patter.2021.100205.

Ikuenobe, P. (2015). Relational autonomy, personhood, and African traditions. Philos. East West *65*, 1005–1029. https://doi.org/10.1353/pew. 2015.0101.

Carman, M., and Rosman, B. (2021). Applying a principle of explicability to AI research in Africa: should we do it? Ethics Inf. Technol. *23*, 107–117. https://doi.org/10.1007/s10676-020-09534-2.

Metz, T. (2018). Two conceptions of African ethics. Quest: Afr. J. Philos. *25*, 141–161. https://philpapers.org/rec/METTCO.

Duffy, K. (1995). Social Exclusion and Human Dignity in Europe: Background Report for the Proposed Initiative by the Council of Europe (Strasbourg: Council of Europe).

Gevaert, C.M., Carman, M., Rosman, B., Georgiadou, Y., and Soden, R. (2021). Fairness and accountability of AI in disaster risk management: opportunities and challenges. Patterns *2*, 100363.

Mittelstadt, B. (2019). Principles alone cannot guarantee ethical AI. Nature Machine Intelligence *1*, 501–507. https://doi.org/10. 1038/s42256-0190114-4.

Jobin, A., Ienca, M., and Vayena, E. (2019). The global landscape of AI ethics guidelines. Nat. Mach. Intell. *1*, 389–399. https://doi.org/10.1038/ s42256-019-0088-2.

Tutu, A.D. (1999). No future without forgiveness. New Perspect. Q. *16*, 29–30. https://doi.org/10.1111/j.1540-5842.1999.tb00012. x.

Verhoef, H., and Michel, C. (1997). Studying morality within the African context: a model of moral analysis and construction. J.

Moral    Educ.    *26*,    389–407.    https://doi.org/10.
1080/0305724970260401.

Saloojee, A., and Saloojee, N. (2011). From social exclusion to social
inclusion: theory and practice over two continents. Afr. J.
Public Aff. *4*, 1–17. https://repository.up.ac.za/bitstream/
handle/2263/57703/Saloojee_Social_2011.pdf?sequence=
1.

Silver,    H.    (2015).    The    contexts    of    social    inclusion.
https://www.un.org/esa/desa/papers/2015/wp144_
2015.pdf.

— (2021). AI Now Institute. https://ainowinstitute.org/.

Hao, K. (2020). We read the paper that forced Timnit Gebru out of
Google. Here's what it says. https://www.technology
review.com/2020/12/04/1013294/google-ai-ethics-
research-paper-forced-out-timnit-gebru/.

Robitzski, D. (2019). Microsoft's president met with the Pope to talk
about ethical AI. https://futurism.com/microsoft-pope-
ethical-ai.

Moore, M. (2015). A Political Theory of Territory (Oxford University
Press).

Beauchamp, T.L., and Childress, J.F. (2012). Principles of Biomedical
Ethics, 7th ed (OUP USA).

Saloojee, A. (2003). Perspectives on social inclusion: social inclusion,
anti-racism and democratic citizenship. http://library.bsl.
org.au/jspui/bitstream/1/1796/1/Social_Inclusion_antiracis
m __democratic_citizenship%20(2).pdf.

Pauwels, E. (2020). The Anatomy of information disorders in Africa: geostrategic positioning & multipolar competition over converging technologies. https://www.kas.de/en/web /newyork/single-title/-/content/theanatomy-of-information-disorders-in-africa.

Cortex Logic (2017). Global symposium on artificial intelligence & inclusion - Rio de Janeiro. https://cortexlogic.com/ 2017/11/18/global-symposiumon-artificial-intelligence-inclusion-rio-de-janeiro/.

Powles, J. (2018). The seductive diversion of 'solving' bias in Artificial Intelligence. https://onezero.medium.com/the-seductive-diversion-ofsolving-bias-in-artificial-intelligence-890df5e5ef53.

Tandem Research (2019). AI for all: 10 social conundrums for India. https://tandemresearch.org/publications/ai-for-all-10-social-conundrumsfor-india-working-paper.

Bartosch, U., Sj, S.B., Engels, R., Rehbein, M., Schmiedchen, F., Stapf-Fine, H., and Sulzen, A. (2018). Policy paper on the asilomar principles on artificial intelligence. https://www.researchgate.net/publication/329963051 _Policy_Paper_on_the_Asilomar_Principles_on_Artificial _ Intelligence.

Birhane, A. (2020). Algorithmic colonization of Africa. SCRIPT-ed *17*, 389–409. https://script-ed.org/wp-content/uploads/2020/ 08/birhane.pdf.

Gwagwa, A., Kazim, E., Kachidza, P., Hilliard, A., Siminyu, K., Smith, M., and Shawe-Taylor, J. (2021). Road map for research on responsible artificial intelligence for development (AI4D) in African countries: the case study of agriculture.

Patterns *2*, 100381. https://doi.org/10.1016/j.patter.2021. 100381.

Stubblefield, A. (1995). Racial identity and non-essentialism about race. Soc. Theor. Pract. *21*, 341–368. https://doi.org/10.5840/ soctheorpract1 9952131.

Floridi, L., Cowls, J., Beltrametti, M., Chatila, R., Chazerand, P., Dignum, V., Luetge, C., Madelin, R., Pagallo, U., Rossi, F., et al. (2018). I4People—an ethical framework for a good AI society: opportunities, risks, principles, and recommendations. Minds Mach. *28*, 689–707. https://doi.org/ 10.1007/s11023-018-9482-5.

UNESCO (2019). UNESCO advocated ROAM principles for steering AI for knowledge societies. https://en.unesco.org/news /unesco-advocatedroam-principles-steering-ai-knowledge-societies.

Arewa, O. (2016). Cultural appropriation: when 'borrowing' becomes exploitation. https://theconversation.com/cultural-appropriation-whenborrowing-becomes-exploitation-57411.

Bradshaw, S., and Howard, P.N. (2019). The global disinformation order: 2019 global inventory of organised social media manipulation. https:// digitalcommons.unl.edu/cgi/ viewcontent.cgi? article=1209&context=scholcom.

Pizzi, M., and Romanoff, M. (2020). Governance of AI in Global Pulse's policy work: zooming in on human rights and ethical frameworks. http://www.unglobalpulse.org/2020/12/ governance-of-ai-in-global-pulses-policywork-zooming-in-on-human-rights-and-ethical-frameworks/.

Fjeld, J., and Nagy, A. (2020). Principled Artificial Intelligence: mapping consensur in ethical and rights-based appraiches to principles for AI. https://cyber.harvard.edu/publication/2020/principled-ai.

Babarinde, O. (2007). The EU as a model for the African Union: the limits of imitation. http://aei.pitt.edu/8185/1/BabarindeEUasModellong07edi.pdf.

Mhlambi, S. (2020). From Rationality to Relationality: Ubuntu as an Ethical & Human Rights Framework for Artificial Intelligence Governance.

Hall, S. (2012). Introduction: who needs 'Identity'? In Questions of Cultural Identity, S. Hall and P. Du Gay, eds. (SAGE Publications), pp. 1–17.

Mbiti, J.S. (1970). African Religions & Philosophy (Heinemann).

Gyekye, K. (2010). African ethics. In Stanford Encyclopedia of Ethics, E. Zalta, ed. (Metaphyscis Research Lab: Stanford University) https://plato. stanford.edu/entries/african-ethics/.

Matolino, B. (2009). A response to Eze's critique of Wiredu's consensual democracy. South Afr. J. Philos. *28*, 34–42. https://doi.org/10.4314/sajpem.v28i1.42904.

Shulman, C., Jonsson, H., and Tarleton, N. (2009). Which consequential-ism? Machine ethics and moral divergence. Asia-Pacific Conf. Comput. Philos. 23–25. https://citeseerx.ist.psu.edu/viewdoc/download?doi=10. 1.1.363.2419& rep=rep1& type=pdf.

Gwagwa, A.E. (2021). Africa's contribution to an intercultural reflective turn in the ethics of technology in the era of disruption. https://www.academia.

edu/51050494/Africas_contribution_to_an_intercultural_r eflective_turn_in_the_ethics_of_technology_in_the_era _of_disruption.

Brey, P. (2007). Is information ethics culture-relative? Int. J. Technol. Hum. Interact. *3*, 12–14. https://doi.org/10.4018/ jthi.2007070102.

Hongladarom, S., and Ess, C. (2006). Information technology ethics: cultural perspectives. https://doi.org/10.4018/978-1-59904-310-4.

Kazim, E., and Koshiyama, A. (2020). A high-level overview of AI ethics. SSRN Electron. J. *2*, 100314. https://doi.org/ 10.2139/ssrn.3609292.

# 9

# CONFUCIAN ROBOTIC ETHICS

*JeeLoo Liu, USA*[221]

## Abstract

This paper will explore the possibility of implementing Confucian ethical codes into the so-called AMAs (artificial moral agents). Drawing from the Confucian classic *The Analects*, it will consider what ethical precepts could be incorporated into robot morality. It will also contrast Kantian AMAs, Utilitarian AMAs, and Confucian AMAs to decide the strength and weakness of each model. The paper's thesis is that robots should be constructed with certain virtues highlighted in Confucian virtue ethics. With the Confucian moral codes built in, their functional ethics can qualify them as moral agents, albeit artificial.

[221] JeeLoo Liu, Department of Philosophy, California State University, Fullerton/USA. Jeelooliu2gmail.com.

## 9.1 Introduction

With the advancement of AI technology, the appearance of intelligent humanoid robots in our society is very likely in the foreseeable future. Whether they truly possess human intelligence and can think like humans do is up to philosophical debate, but they will surely pass the Turing Test—i.e., they will entice their human interlocutor to be inclined to treat them as humans. Intelligent robots will one day become members of our society, sharing our jobs, taking care of our elderly, serving us at restaurants and hotels, making important navigational, military and even medical decisions for us. Should we equip these robots with a moral code to teach them right from wrong? If so, what kind of moral codes would be able to bring about the kind of artificial moral agents (AMAs) that we would like to have in our society?

On the optimistic assumption shared by many AI designers that the development of AMAs can be successful one day, this paper will explore the possibility of implementing Confucian ethical codes into the so-called AMAs. Drawing from the Confucian classic *The Analects*, it will consider what ethical precepts could be incorporated into robot morality. It will also contrast Kantian AMAs, Utilitarian AMAs, and Confucian AMAs to decide the strength and weakness of each model. The paper's thesis is that even though robots cannot have our innate moral sentiments, the four moral sprouts that Mencius defends, they can be constructed with the kind of ethical principles that Confucianism stresses. With the Confucian moral codes built in, their functional ethics can qualify them as moral agents, albeit artificial.

The investigation of AI ethical codes is not just a futuristic mindgame. According to Michael Anderson and Susan Leigh Anderson, "Machine ethics, by making ethics more precise than it has ever been before, could lead to the discovery of problems with current ethical theories, advancing our thinking about ethics in general" (Anderson & Anderson 2006, 11). This paper will show that the comparative study on

robot morality can shed light on the flaws of the Kantian model as well as the Utilitarian model for human ethics.

## 9.2 The Rise of Machine Ethics

Up to now, having robots that can make systematic ethical decisions with advance considerations of their consequences is still a remote dream. However, there are already existing specific decision guidelines for machines. Some of these decisions have morally significant consequences. For example, autonomous military drones can be programmed to strike or withhold attack with the detection of civilians in the vicinity of a military target. Health-care robots can also be programmed to take life-saving measures or to forego further treatments. According to Ryan Tonkens, "Because autonomous machines will perform ethically relevant actions, akin to humans, prudence dictates that we design them to act morally" (Tonkens 2009, 422). Therefore, even if we cannot make "moral machines" yet, we must consider machine ethics. Furthermore, the version of machine ethics we formulate should be applicable to foreseeable robotic moral reasoners, and not just to their programs designed by humans. In other words, machine ethics is concerned with applying ethical codes to artificial moral agents, not to their designers.

According to Allen, Smit and Wallach (2005), there are three fundamentally different approaches to designing artificial morality: the bottom-up approach, the top-down approach, and a hybrid model that combines the above two.[222] The bottom-up approach is to have the machine develop its own ethical code from the piecemeal rules in its day-today decisions. The machine can be given learning skills with which to process the information gathered by facing the consequences of various courses of action it takes. To promote a certain type of behaviour, the designer can create a reward system that favours certain actions that the

---

[222] There are of course hybrid approaches as well.

machine takes. Such feedbacks can enable the machine to develop its own ethical codes in time. This approach is similar to human childhood learning experience in building moral character. However, Allen et al argue that it is questionable whether this approach could be helpful in developing artificial moral agents that are "capable of engaging the more complex dilemmas that we encounter daily" (Allen et al 2005, 152). The top-down approach, on the other hand, is to implement general, abstract ethical rules that would govern the machine's daily decisions and actions. To use this approach, the designer must first choose an ethical theory to analyse "the informational and procedural requirements necessary to implement this theory in a computer system," and then design its subsystems for the implementation of the ethical theory (Wallach & Allen 2009, 80). Even with the preset design, in each scenario the machine will need to use deduction to determine the best course of action under the programmed ethical principles. This approach reflects the debates in normative ethics, since different ethical theories will generate different ethical codes for artificial moral agents. In this paper, we will first consider three leading models: Asimov's Laws of Robotics, Kantian deontology, and utilitarianism. To fully evaluate the applicability or the desirability of each ethical model, Allen et al suggest that we must give "careful consideration to the prospects for building AMAs by implementing decision procedures that are modelled on explicit moral theories" (Allen et al 2005, 150). In other words, the devil is in the details. However, this paper will not be able to touch on the technical implementation of ethical principles or the algorithmic design of the decision procedure. The critique will be largely conceptual.

In this paper, we suggest a hybrid approach, combining both some general ethical principles in the robot's initial design, and a learning mechanism that enables the robot to improve and improvise as it goes through different trial and actual situations. According to Anderson & Anderson (2007), the aim of machine ethics is to define explicit ethical principles that artificial intelligence could appeal to in choosing and

justifying its own actions. They argue that we cannot possibly give specific rules for each and every possible situation that might come up. "The virtue of having principles to follow, rather than being programmed in an ad hoc fashion to behave correctly in specific situations, is that it allows machines to have a way to determine the ethically correct action in new situations, even in new domains." (Anderson & Anderson 2007, 17) In other words, we want to have artificial intelligence to actually be artificial moral agents with their own moral principles, making moral deliberations on the basis of those principles, and justifying their action by appealing to those principles. Hence, the choice of a set of moral principles that can be implemented into artificial intelligence is a crucial task of machine ethics. Having only the abstract ethical principles is not sufficient to equip machines with the abilities to adapt to new situations. Virtue ethics defines ethical behaviour as what a virtuous person would perform in the given situation. "A virtuous person is defined as a person who has learned and internalised a set of habits or traits termed virtuous. For a virtuous person, virtuous acts become second-nature, and hence are performed in many different situations, through time" (Govindarajulu et al 2019, 1).

The advantage of virtue ethics to the deontological and the utilitarian approaches is that virtue ethics can better capture human responses, human expectations, and human values. Virtues are the stable, consistent traits that can have explanatory power for the agent's behaviour and can be used with predictive power for the agent's future behaviour (Alfano 2013). Govindarajulu et al argue that "if the conditions of stability, consistency, explanatory power, and predictive power hold, then virtuous agents or robots might be easier for humans to understand and interact with (compared to consequentialist or deontological agents or robots)" (Govindarajulu et al 2019, 2). That is to say, virtue ethics is more conducive to the design of sociable robots, the kind of robots that can interact and communicate with humans. Virtue ethics provides a hybrid model in that the machine is designed with certain characteristic

traits (virtues) that enable it to learn from the available data, and to make novel decisions consistent with the preset traits. Instead of giving robots abstract ethical principles to add into its processing and calculation to generate action, the approach of virtue ethics equips robots with some predetermined traits (their "virtues"), and a large databank that stores information of what virtuous agents have done or would do. Before we begin to collect data on virtuous agents' behaviour, what we need to identify first is which virtues are essential to sociable robots.

This paper will introduce Confucian virtue ethics, and lay out some essential virtues that artificial moral agents should have in order to be accepted into an ethical human society.

## 9.3 The Trolley Problem and Various Ethical Models

Since intentional robot action or intervention is still a hypothetical scenario, we might find it helpful to appeal to the commonly used thought experiment for ethical dilemma: the trolley problem.

### 9.3.1 The Standard Trolley Problem

A runaway trolley is rushing down the railway tracks, and there are five people on the tracks ahead, unable to flee in time. The robot safety inspector (or driver) can intervene by pulling the lever to divert the trolley onto a different track. However, there is also one person on the other track. The choice is between sacrificing one person to save five human lives, and not sacrificing one life and letting the five people die. Should the robot inspector/driver pull the lever to prevent the disaster, or do nothing?

### 9.3.2 The Footbridge Variation of the Trolley Problem

The robot safety inspector is standing on the footbridge above the trolley track to observe the trolley traffic. Upon seeing that under the bridge a runaway trolley is heading down a track with five people stuck

on the track, the robot has to do something quick.  Next to it on the foot-bridge is a heavyset man observing the same scene. If the robot pushed the man onto the track to stop the trolley, it would be able to prevent the disaster of having five people killed.  Should it do it?

Faced with the two kinds of dilemma, the moral agent must decide "whether the action required to save the five is impermissible because it causes harm, or permissible because the harm is only a side effect of causing good" (Deng 2016). Experiments show that humans typically choose to sacrifice one to save the five people in the trolley example, but would not choose to sacrifice the fat man on the bridge to save the five people on the track.  According to the experimenters: "This leaves psychologists with a puzzle of their own: How is that nearly everyone manages to conclude that it is acceptable to sacrifice one life for five in the trolley dilemma but not in the footbridge dilemma, in spite of the fact that a satisfying justification for distinguishing between these two cases is remarkably difficult to find" (Greene et al 2001, 2106)? These two cases will serve as our test case for the various ethical models for artificial intelligence.

In human ethical contexts, the trolley problem may seem far-fetched and unrealistic; however, in the context of AMAs, similar situations may arise with machine ethics.  Imagine a future Tesla equipped with an ethical overriding principle to avoid harming more people than necessary.[223] Suppose a school bus filled with school children suddenly loses control and is crashing into a Tesla, which cannot stop in time to prevent a collision. Should the car veer off to hit a median, risking the life of its

---

[223] Examples such as this are abundant in the discussion on self-driving or driverless cars. See Bonnefon *et al*. 2016, Deng 2015, Greenemeier 2016, and Herkewitz 2016. Of course, currently driverless cars do not make ethical decisions on their own. They can only take in information about speed, road conditions, weather conditions, and so on, to make instantaneous driving decisions.  But if artificial moral agents are possible, then some simplified ethical codes could conceivably be programed into these cars as well.

driver, or should it continue its course and let the collision take place? Maybe no one would want to buy a Tesla if given this kind of moral consideration, it could end up sacrificing its driver, but the point is that this scenario is analogous to the trolley problem cases. Many such scenarios could be envisioned. Therefore, the trolley problem could serve as the test for our ethical theories.

### 9.3.2.1 Asimov's Laws of Robotics

A prime example of the top-down approach is to appeal to the Three Laws of Robotics composed by Isaac Asimov in 1942[224]:

[A1] *A robot may not injure a human being or, through inaction, allow a human being to come to harm.*

[A2] *A robot must obey the orders given it by human beings except where such orders would conflict with the First Law.*

[A3] *A robot must protect its own existence as long as such protection does not conflict with the First or Second Laws. (Cited in Wallach & Allen 2009, 34)*

A fourth law was later added, which supersedes the above three:

[A4] *A robot may not harm humanity, or, by inaction, allow humanity to come to harm. (Wallach & Allen 2009, 91)*

The difference between the First Law and the Zeroth Law is that the former concerns individual human beings while the latter concerns general humanity. Since the Zeroth Law supersedes the First Law, the implication is that the robot may conceivably harm individual human beings if doing so could prevent harm to humanity. An example of the application of this law in a science fictional situation would be if some individuals carry deadly contagious virus that could potentially wipe out humankind, then the robot would be obligated to eliminate such individ-

---

[224] The Three Laws were first introduced in Asimov's science fiction story *Runaround* in 1942. According to Wallach & Allen (2009), "No discussion of top-down morality for robots can ignore Asimov's Three Laws." (Wallach & Allen 2009, 91)

ual humans. However, such a moral precept is highly dubious, since the notion of humanity is rather abstract and has been used to justify many evil practices in human history such as eugenics, ethnic cleansing, and so on. Having an overriding Zeroth Law thus can conceivably annul Asimov's Three Laws.

When applied to the trolley problem, Asimov's Three Laws are clearly inadequate. In the standard trolley scenario, the First Law prohibits the robot's pulling the lever, as it will bring harm to the one person on the other track, but it will also prohibit the robot's doing nothing, as its inaction will lead to the harm of five human beings. In the footbridge scenario, the robot's pushing the fat man over the bridge is strictly prohibited, since it directly involves harming a human being; but at the same time, the robot's doing nothing, when the option of pushing the fat man over to stop the trolley is available, is also against the second clause of the First Law. Either way, the robot is in a bind with no ethical guidance. Alan Winfield et al conducted a series of experiments, in which the "A-robot" (named after Asimov) is charged with the mission to save human lives. The experiments include three different scenarios, where the baseline is set with A-robot only, and its task is merely to secure its own safety. In this scenario, robot A was able to avoid falling into a hole with 100% reliability. The second scenario added robot H representing human being, while the third scenario included two other robots, H and H2, representing two people. In the scenario when there was only one human life involved, the A-robot performed the task successfully. In the scenarios where there were two human beings (in the form of "H-robot") facing the danger of falling into a hole, in almost half of the trials, "the A-robot went into a helpless dither and let both 'humans' perish" (Deng 2016). Winfield et al believe that what they have set out in their experiments "match remarkably well with Asimov's first law of robotics" (Winfield et al 2014, 89). Thus, the A-robot's failure demonstrates the inadequacy of the Asimov' Laws in dealing with more complex ethical scenarios. To fix this problem, we would need to intro-

duce extra rules about how a robot can make choices in a moral dilemma.

Now we turn to the common proposals of ethical codes for robots: the Kantian ethics and Utilitarian ethics.

### 9.3.2.2 Kantian Ethical Codes for AMAs

One of the leading ethical models for machine ethics is Kantian moral theory, which is regarded as "one of our best chances for the successful implementation of ethics into autonomous robots" (Tonkens 2009, 422). Kant's normative ethics is a form of deontology, appealing to one's sense of duty, not emotions, in making moral judgments. "Duties are rules or laws of some sort combined with some sort of felt constraint or incentive on our choices" (Johnson and Cureton 2017). It is understandable why Kant's moral philosophy would be considered a natural choice for machine ethics: for Kant, human's self-interest, desire, natural inclinations, moral sentiments such a sense of honour or the feeling of sympathy and compassion, all have absolutely no moral worth. According to Kant, even in a case where one finds an inner pleasure in spreading happiness around them, and takes delight in the contentment of others as one's own work, however right and amiable this case may be, it "has still no genuinely moral worth" (Kant 1993, 11). Genuine moral acts must be done solely from the sense of duty of a rational being. Without the noises of human sentiments and desires, robots would be ideal for the implementation of Kant's purely rational moral schema.

Kant's first categorical imperative is stated as follows:

*[D1] I ought never to act except in such a way that I could also will that my maxim  should become a universal law. (4:402) [Or: Act only in accordance with that maxim through which you can at the same time will that it become a universal law. (4:421)]*

A maxim is a person's subjective volitional principle, while a universal law is a moral law that binds all rational beings in nature.  Kant's categorical imperative is based on the assumption that humans' moral decisions come after deliberation of particular situations in accord with

the individual principle that one adopts for dealing with the situation at hand. As stated this way, Kant's categorical imperative demands both intrapersonal consistency and interpersonal agreement. This categorical imperative serves more as an injunction against acts that do not follow universalizable maxims, such as committing suicide because of misfortune, borrowing money with the knowledge that one can't pay it back as promised, not cultivating one's natural talents while indulging in the pursuit of pleasure, or not offering assistance to someone in dire need of help, when doing so would not put one in undue distress.

In terms of robotic ethics, we can have the moral rule formulated as such:

[DR1] *A robot should act only in such a way that the option chosen could in principle be a universal law for other robots.*

Since robotic ethics needs to be designed as situation ethics; that is, case-by-case deliberation on the course of action to take, [DR] focuses on the option rather than the maxim. The robot would have to be equipped with the function of gathering data on probable consequences and calculating the results in each case. In other words, the robot would have to be a consequentialist.

Whereas the first categorical imperative serves as an injunction, Kant's second categorical imperative gives a more specific moral guidance:

[D2] *Act in such a way that you always treat humanity, whether in your own person or that of another, always at the same time as an end, never merely as a means. (4:439)*

People must never be treated merely as a means to an end. People have their free will. To treat them merely as a means to an end is to deny their autonomy. In terms of robotic ethics, the categorical imperative can be formulated as such:

[DR2] *A robot must act in such a way that it always treats humanity never simply as a means, but always at the same time as an end.*

When humans are faced with the trolley problem dilemma, the typical response is that they would choose saving five lives and sacrificing one person in the standard scenario, but most would refuse pushing one person over the footbridge in order to save the five people on the track. In the trolley scenario, the rationale seems to be in favour of the 5:1 human life ratio. In the footbridge scenario, on the other hand, pushing someone over the bridge in order to stop the trolley is a clear violation of the second categorical imperative. According to Joshua Greene, "People exhibit a characteristically consequentialist response to the trolley case and a characteristically deontological response to the footbridge case." (Greene 2007, 42) The reason for people's different reactions, according to Greene, is that "the thought of pushing someone to his death in an 'up close and personal' manner (as in the footbridge dilemma) is more emotionally salient than the thought of bringing about similar consequences in a more impersonal way (e.g., by hitting a switch, as in the trolley dilemma)" (Greene 2007, 43). Therefore, Kant's deontological approach is actually "a kind of moral confabulation" since deontological judgments are prompted by emotional responses (Greene 2007, 63). Greene calls this "the secret joke of Kant's soul."

According to Greene, when we make conscious moral decisions, we "respond to the conscious deliverances of our unconscious perceptual, mnemonic, and emotional processes by fashioning them into a rationally sensible narrative, and without any awareness that we are doing so." (Greene 2007, 62) However, such an influence from the unconscious would not exist for robots and other forms of artificial intelligence. We can now consider whether Kantian ethics would compel robots to make the same decision that humans do.

In the footbridge scenario, the robot would obviously refuse to take the option of pushing the fat man over the bridge, since doing so would be a clear violation of [DR2]. So this is a no brainer. In the standard trolley scenario, however, the robot's moral guidance is not so definitive. If the robot were to act on the principle that it judges to be universaliza-

ble for all robots, then it would likely be too paralyzed to do anything. Human beings will often have a spontaneous, intuitive judgement on whether their own maxims are suitable to be universal laws. Asking a robot to make such judgments, on the other hand, requires that the robot be equipped with either a huge databank containing all possible consequences for other robots to act in the same way, or the kind of intuition that humans have, but such intuitions are not available to artificial intelligence.

Herein lies the fundamental paradox of designing Kantian AMAs. Kant's first categorical imperative is based on his metaphysics of morals, according to which all human beings are autonomous rational beings. Rational beings are citizens of the Kingdom of Ends, sharing the same common laws and abiding by the same moral principles that they themselves legislate. Their autonomy lies in the sense that they are fully rational agents who "have an equal share in legislating these principles for their community" (Johnson & Cureton 2017). Artificial moral agents are, by design, machines that obey the programmers' orders. They do not act to legislate their own laws; they do not respect one another as equal law-makers. Furthermore, they lack freedom of the will, which is essential to the status of Kantian moral agency. For this reason, Ryan Tonkens calls them "anti-Kantian." Tonkens says, "Because we require our Kantian AMAs to act ethically, the fact that their development is a violation of Kantian morality renders their creation morally suspect, and our role as their creators somewhat hypocritical" (Tonkens 2009, 429). Therefore, not only is it problematic to implement Kant's ethical codes into artificial intelligence, we are actually acting immorally by creating Kantian AMAs. Tonkens argues, "By creating Kantian moral machines, we are treating them merely as means, and not also as ends in themselves. According to Kant, moral agents are ends in themselves, and because of this they ought to be respected as such. To violate this law is to treat an agent merely as an object, as something used for achieving other ends" (Tonkens 2009, 432-3). In other words, even if we could

create robots that would be able to follow our DR2, the very creation itself has already violated Kant's ethical principle.

### 9.3.2.3 Utilitarian Ethical Codes for AMAs

Applying utilitarian ethical codes for AMAs is another popular proposal. The utilitarian principle, simply put, is to judge the merit of an act by the potential consequences: actions are right insofar as they promote happiness or pleasure, and actions are wrong insofar as they generate pain. To John Stuart Mill, pleasure and pleasure alone has intrinsic value. Their value consists in people's desire. In other words, "good" is identified with "desirable." Mill says, "the sole evidence it is possible to produce that anything is desirable is that people do actually desire it" (Mill 2001, 81). What people desire would be good consequences; what people detest would be bad consequences. The most important feature of utilitarianism is the sole consideration of the number of people affected by the act instead of the self-interest of the actor. The utilitarian principle has standardly been formulated as follows:[225]

*[U] An act is right if and only if it produces a greater balance of good over bad in its consequences for all people affected, than any other act available to the agent.*

In the context of artificial intelligence, we can reformulate [U] as Utilitarian Robot [UR]

*[UR] In weighing the consequences of available courses of action, a robot must choose the option that will either generate the maximum benefits, or prevent the greater harm, for all human beings involved.*

According to Julia Driver, "Since the early 20th Century utilitarianism has undergone a variety of refinements. After the middle of the 20th Century it has become more common to identify as a 'Consequentialist'

---

[225] Utilitarianism can be act utilitarianism as formulated here, or rule utilitarianism: An act is right if it accords with a rule the general following of which produces a greater balance of good over bad for all people affected, than any alternative rule. Since artificial intelligence needs more precise rules and programs for each act, here we are only discussing act utilitarianism.

since very few philosophers agree entirely with the view proposed by the Classical Utilitarians, particularly with respect to the hedonistic value theory" (Driver 2014). We shall now look at some studies on the consequentialist model of machine ethics.

In a joint study conducted by Jean-François Bennefon of the University of Toulouse, Azim Shariff of the University of Oregon, and Iyad Rahwan of MIT, participants were asked to evaluate autonomous vehicles that apply utilitarian ethical codes to favour sacrificing themselves and their driver to avoid running over a group of pedestrians. The study found that although participants approve of such vehicles for the greater good, they themselves would not want to buy this kind of automatic vehicle (AV). The study found that participants "overwhelmingly expressed a moral preference for utilitarian AVs programmed to minimise the number of casualties."[226] However, when asked whether they would personally purchase such a utilitarian AV, participants were less positive. The researcher noted that "even though participants still agreed that utilitarian AVs were the most moral, they preferred the self-protective model for themselves" (Bennefon et al 2016, 1574). This double standard creates a social dilemma: "Although people tend to agree that everyone would be better off if AVs were utilitarian (in the sense of minimising the number of casualties on the road), these same people have a personal incentive to ride in AVs that will protect them at all costs. Accordingly, if both self-protective and utilitarian AVs were allowed on the market, few people would be willing to ride in utilitarian AVs, even though they would prefer others to do so" (Bennefon et al 2016, 1575). Without government regulation, such a utilitarian model of automatic vehicles would not be in the market; however, the possibility of government regulation would create even more resistance toward adoption of such a model. In other words, "regulating for utilitarian algorithms may

---

[226] "Overall, participants strongly agreed that it would be more moral for AVs to sacrifice their own passengers when this sacrifice would save a greater number of lives overall." (Bonnefon et al 2016, 1574)

paradoxically increase casualties by postponing the adoption of a safer technology" (Bennefon et al 2016, 1573). The potential problem of a utilitarian model for artificial agents is manifested in this kind of conflicts between social utility and personal self-interest. If participants in these studies would be reluctant to purchase an automatic vehicle designed with the utilitarian model, then general public would be most likely resistant toward the idea of ethical robots implemented with the utilitarian ethics.

Other than the undesirability of utilitarian artificial intelligence, there is also the grave danger of having such artificial moral agents around in our society. The different responses people have with the trolley case and the footbridge case show that humans would refrain from certain actions that involve clear and personal harm. Humans do not always favour a utilitarian moral consideration, especially when the smaller number includes themselves, their kin and their acquaintances. Except for exceptional heroic acts, few people would willingly sacrifice themselves or their loved ones to produce the greater utility for all or for the greater good. Artificial agents, on the other hand, have no such inhibitions. Under [UR] and without any other overriding moral principle, they could undertake major destruction if doing so could lead to maximum benefits. This model may face the same difficulties that confront Asimov's Zeroth Law.

The contrast between humans and AIs also demonstrates the fact that utilitarianism is never the go-to principle for humans' moral deliberation. Even when we appeal to the principle of utility, our sentiments, self-interest, and other considerations would always make our utilitarian thinking "impure." And yet if we do make purely utilitarian moral deliberation, as artificial moral agents would, then the outcome would be highly dangerous for human society. As Anderson & Anderson (2007) points out, utilitarianism "can violate human beings' rights, sacrificing one person for the greater net good. It can also conflict with our notion of justice—what people deserve—because the rightness and wrongness

of actions is determined entirely by the future consequences of actions, whereas what people deserve is a result of past behaviour" (Anderson & Anderson 2007, 18). Utilitarianism may have its appeal in normative ethics only because humans do not abide by it completely and absolutely.

## 9.4 Confucian Robotic Ethics

Confucian ethics is not a form of rule-governed normative ethics. It is rather a form of virtue ethics, focusing on building the moral agent's moral character and virtuous traits. To convert Confucian ethics into a form of implementable moral rule, we must do a liberal interpretation of the kind of Confucian moral rules that can be extracted from The Analects.[227]

There are many highly emphasised virtues in The Analects that can be formulated into moral rules for Confucian robotic ethics. I shall choose three main virtues: loyalty (zhong), reciprocity (shu)[228], and humanity (ren), along with other supplementary virtues such as trustworthiness and righteousness. The first two are chosen in light of the comment of one of his chief disciples Zengzi (Master Zeng): what Confucius meant by "a single thread" of his Way is nothing but loyalty and reciprocity (The Analects 4.15). The third virtue is chosen because it is the overarching virtue in the whole Confucian tradition, underscored by both Confucius and Mencius and further elaborated in neo-Confucianism.

With regard to loyalty, Confucius has many things to say. It is one of Confucius' four teachings (The Analects 7.25) and is said to be in conjunction with shu to form the single penetrating thread in his teachings (The Analects 4.15). In personal moral cultivation, Confucius says

---

[227] When not specified, the translations of texts in The Analects are mine with consultation of Dawson 1993 and Ni 2017.

[228] This is Peimin Ni's translation. The word 'shu' is also frequently translated as 'empathy' (I have translated it as empathy in the past).

that a superior person (junzi) takes loyalty and trustworthiness (xin) as his first principles (The Analects 1.9; 9.25). In answering a student's question about how to promote virtue and discern delusion, Confucius' advice was: hold fast to loyalty and trustworthiness, and "move toward what is right—this is the way to promote virtue" (The Analects 12.10, Ni 2017, 288). In government, this virtue also plays an important role. Confucius told the ruler in the state of Lu that to gain loyalty from the people, the ruler himself must be "filial and caring" (The Analects 2.20). He told Duke Ding that the ruler must employ ministers with ritual propriety, while the ministers must serve the ruler with loyalty (The Analects 3.19); however, such loyalty is not blind obedience, but "offering counsel" to the ruler (The Analects 14.7). When a student asked about Prime Minister Ziwen, who "trice took the office of prime minister and showed no joy in his countenance." He was thrice deposed from his position as the prime minister, but he showed no resentment. In addition, he always reported the previous policies to the succeeding prime minister. Confucius' evaluation of him is that he indeed has the virtue of loyalty (The Analects 5.19, Ni 2017, 160). When a student asked about governing, Confucius counselled: abiding in its affairs without weariness; conducting its affairs with loyalty (The Analects 12.14). The above quotes show that in Confucius' mind, loyalty is both a private and a public virtue: it is crucial both in one's conducting oneself and in one's engagement with public affairs.

In my analysis, "Loyalty is not a relationship directed toward others; rather, it is directed towards the role one plays. In this sense loyalty can be defined as 'doing what one is supposed to do' or 'being loyal to one's role.' In other words, a social role is not simply a social assignment; it is also a moral assignment. Being loyal to one's role means being able to act in accordance with whatever moral obligation that comes with the social role. Loyalty is thus being loyal to one's moral obligation and fulfilling the duty that one's role dictates" (Liu 2006, 50). This interpretation is further supported by Confucius' advice to a student, "loyalty in

relationships with others" (The Analects 13.19). The "others" here is not specifically one's superior, but also one's friends or strangers.

With this interpretation of the virtue of loyalty, we can now have our first moral principle for Confucian robotic ethics:

[CR1] *A robot must first and foremost fulfil its assigned role.*

Loyalty to one's task is chosen to be the first law in Confucian robotic ethics because artificial intelligence, with its potential superpower, should particularly be designed to be role-specific and not omnipotent. Confucius says, "To guard Dao is not as good as to guard one's station/role" (Zuozhuan Zhaogong 20). According to Kam-por Yu, "The Dao is of course the higher or final goal, but the question is whether everyone should aim at the Dao directly, or one should just fulfil one's role faithfully, as the realisation of the Dao relies not just on one person, but on the collaboration of a number of people fulfilling their roles."[229] Confucius also taught: "If one is not on the post, then one does not meddle with the managerial affairs" (The Analects 8:14). This comment demonstrates his view that one should do what is in one's duty and not overtake someone else's.

On first impression, this ethical code for artificial intelligence may seem trivial: of course machines are designed to complete the tasks. However, since we are discussing the possibility of future artificial moral agents that could make judgement calls under certain dire situations not previously anticipated by their programmers, we need to be prepared for these scenarios. [CR1] has a clear division of labour: a robot designed to offer health care should be specifically loyal to such a role, not to make other judgments such as whether the patient's life is not worth keeping or to render assistance to meet the patient's desire for euthanasia. An intelligent automatic vehicle should fulfil its duty to ensure safe driving for its driver, and thus ought not to take any act to sacrifice its driver by hitting against a tree in order to prevent a catastrophic disaster

---

[229] The quote is from personal communication.

for a school bus or the deaths of multiple pedestrians. It would be a mistake for us to try to design a "universal robot" as depicted in the play R.U.R., where the word 'robot' was first coined.

An important virtue often paired with loyalty is reciprocity. Concerning reciprocity, Confucius commented that this word could "serve as guidance for practice during one's entire life," and he further defined it as such: "Do not impose on others what you would not wish for yourself" (The Analects 15.24, Ni 2017, 364). In a different passage, his chief disciple Zigong says, "If I do not want others to inflict something on me, I also want to avoid inflicting it on others" (The Analects 5.12, Dawson 1993, 17). From these two quotes, we can see that the connotation of shu is specifically defined as an interpersonal demeanour with a psychological preparedness. In contrast to the Christian Golden Rule: "Do unto others as you would have them do unto you", this has often been called the negative Golden Rule, as it states injunction on what not to do, rather than specific commandment on what to do. In my analysis, this formulation of reciprocity is better than the Golden Rule in that what people do not desire, seem to have more common ground than what people do desire. In general, we do not wish for others to humiliate us, to deprive free will from us, to steal from us, to harm us, or simply to mistreat us in any way. "It is reasonable that we do not mistreat others in these ways either. And even if we desire others to act in a certain way towards us, the Confucian Golden Rule does not counsel us to act this way towards others too. It thus avoids the problem of subjective imposition of preferences that we see in the positive formulation" (Liu 2006, 55).

In terms of robotic ethics, however, we encounter the problem of the lack of desire in robots. If robots do not have any desires on their own, then how do they assess whether the consequences of their act would be what others (other human beings) would not want to be imposed on them? I think we can solve this problem by adding an algorithm of the scale of human preferences into the design. This scale of preferences can

be programmed as a machine "preference function" in the manner that Hilary Putnam suggested for functionalism: there should be a preference partial ordering and an inductive logic (i.e. the Machine must be able to "learn from experience"), some "pain sensors," i.e., sensory organs which normally signal damage to the machine's body, or dangerous temperatures, pressures, etc., and that "the inputs in the distinguished subset have a high disvalue on the Machine's preference function or ordering" (Putnam 1967, 435). In this way, the artificial moral agent can assign a negative value to harm done to other human beings as well as a disvalue to its own damage. The second rule for Confucian robotic ethics could be stated as follows:

[CR2] *A robot should not act in ways that would afflict the highest displeasure or the lowest preference onto other human beings, when other options are available.*

As so formulated, [CR2] is still a negative injunction on what not to do. A general application of [CR2] is that a robot should never harm a person by choice, never inflict pain on human beings without due cause, never deprive someone of prized possession unless there are overriding considerations, and so on and so forth. This moral rule is similar to Asimov's First Law: A robot may not injure a human being or, through inaction, allow a human being to come to harm. However, this moral rule is much more flexible than Asimov's First Law in that there might be prevailing displeasure or negative preference that would outweigh the negative preference of harm. For example, radical injustice may be ranked as more undesirable than the possibility of physical harm. It is therefore conceivable that robots could be involved in insurgence against injustice and abuse if they have been designed with the appropriate role assignment and the right set of preference ordering.

The virtue of trustworthiness can be seen as an indispensable safeguard in robotic ethics, in that our robots are designed to serve various functions in our society and to interact with human beings. We need to be able to entrust our robots to perform the tasks assigned, not to deviate

from our expectations, and not to deceive us. Designing robots with this virtue in its default state is not to equip it with a Kantian prohibition: "Do not lie,"[230] since sometimes lying might be the more virtuous thing to do in a given situation. A deontological rule against lying can be formulated in contrast with this rule of trustworthiness:

[DR3] A robot must never lie. A robot must always give truthful answers.

In contrast, the implementation of the virtue of trustworthiness would be as follows:

*[CR3] A robot must always speak and act in a trustworthy way.*

The difference between having an absolute mandate: "do not lie," and having the virtue of trustworthiness in the robotic default design, is that the former is written as an inviolable rule, while the latter is modelled after what virtuous people would do to be trustworthy. Of course, establishing a databank incorporating existing trustworthy behaviour in multifarious contexts would be a daunting task. However, with the hybrid approach, the robot would also be equipped with the ability to learn from previous samples and deduce its own trustworthy conduct in a new situation.

Next, the virtue of righteousness (yi) is crucial for moral agency, in that righteousness is almost interchangeable with morality: doing the right thing in the right context. Righteousness is one of the four cardinal virtues (humaneness, righteousness, propriety, and wisdom) that Mencius highly emphasised, and Mencius' theory of the four moral sprouts (siduan) treats the fully developed moral sprouts as the four cardinal virtues. In the Analects, Confucius contrasts a virtuous moral agent, a superior person (junzi), with a petty person (xiaoren), in this way: "The mind of a superior person is preoccupied with righteousness, while the mind of a petty person is preoccupied with profit" (the Analects, 4.16). He also uses righteousness as the criterion for his own conduct and life's

---

[230] Kant claims that lying is always morally wrong, no matter what the motive or the consequences might be.

choices (The Analects, 7.16). Confucius' disciple Master You says that one's "virtue of trustworthiness must accord with righteousness for one's words to be carried out" (The Analects 1.13). This shows that righteousness is a core virtue in Confucian ethics. In Chinese, the word for righteousness, yi, is homonymous with the word for appropriateness yi. By association, being righteous means doing things that are appropriate for the situation, and thus must be accompanied by cognition and judgement. Righteousness is thus an intellectual virtue. Confucian situation ethics is exemplified in the conception of this virtue: doing the right thing in the right situation. There is no universal principle that applies in all situations. Equipped with this virtue, our robot will be asked to make judicious decisions in the given situation, rather than following a universal mandate to act in the same way no matter what the situation is. It will be a flexible ethical design that allows the robot to assess the situation to choose the most appropriate action. Confucius says of a superior person: "The superior person (junzi) does not set his mind either for or against anything absolutely. He simply chooses to do what is right" (The Analects, 4.10); accordingly, our robot should be designed with the following principle:

[CR4] *A robot must be flexible in assessing the situation and make the decision most appropriate to the situation.*

Next, we come to the set of virtues serving as the basis for societal norms and personal sense of propriety (li): humility and respectfulness. In Chinese, humility (gong) and respectfulness (jing) are typically used together as a compound term, gongjing, and in the Mengzi, Mencius lists the heart of humility and respectfulness as one of humans' four "moral sprouts." In Mencius' view, humans naturally have this proclivity for societal cooperation and conformity to social norms. Cultivating this natural tendency into the virtue of propriety can serve as the foundation for society's codes of propriety such as rites, rituals, and etiquette. In the Analects, the two words of humility and respectfulness are not combined, but the two virtues are often mentioned together. Confucius

praises one student as "having humility in the way he conducts himself, and having respectfulness in the way he serves his superior" (The Analects, 5.25). He also explains to a student that the essence of the virtue of humanity includes "handling oneself with humility, handling affairs with respectfulness, and having loyalty when doing things for others" (The Analects, 13.20). Furthermore, Confucius says that a superior person must think about "being humble in one's demeanour and being respectful in handling affairs" (The Analects, 16.10). Sometimes the virtue of respectfulness is also manifested in one's attitude towards others (The Analects, 2.7; 4.18; 5.16; 5.17; 11.15).  The above quotations show that the virtue of humility is mostly associated with one's general demeanour, while the virtue of respectfulness is mostly associated with one's attitude in dealing with particular people and affairs. These two virtues can be implemented in the default design of the robotic speech and demeanour. We can formulate the next two rules as follows:

[CR5] *A robot must always be humble in its demeanour and speech.*

[CR6] *A robot must respect the task at hand and be respectful to its interlocutor.*

With these two design guidelines, we will not have robots who would swear or curse, or robots whose behaviour would be flippant, insolent, dismissive, or defiant.  Even though we cannot stop other humans from lacking these two virtues in their speech and act, we at least would have our artificial members always behave in accordance with propriety.

The last, but not least, virtue selected for Confucian robotic ethics is that of humanity (ren). According to Peimin Ni, "The term ren is central to Confucius' philosophy.  It appears 109 times in The Analects, and of the 499 selections in the boo, 58 are devoted to this subject" (Ni 2017, 32). In Confucius' assessment, the virtue of humanity is the hardest to cultivate.  Even his prize student Yan Hui could only maintain this virtue in his heart up to three months, while the rest of the students could at most keep it for a day or a month (The Analects 6.7). When asked to

judge whether someone possesses this quality, Confucius rarely granted it even though he would acknowledge that the person has some other redeeming qualities (The Analects 5.5; 5.8; 5.19). At the same time, whether one can obtain this virtue is purely a matter of volition according to Confucius. He says, "Is humanity ever far away? If I want to achieve it, then I am already there" (The Analects 7.30). Confucius gives humanity the highest praise: He thinks that only the humane people are capable of loving and loathing people (The Analects 4.3), and if anyone is devoted to this virtue, she would never be doing anything bad (The Analects 4.4). If we say that Kant's ideal is the Kingdom of Ends, then Confucius' ideal would be the Kingdom of Humanity. Confucius talks about residing in the circle of humanity, surrounded by people of like virtue (The Analects 4.1). He also commented that humanness is a robust and resilient virtue that the superior people should never forsake even for the short duration of a mealtime. The superior people should never deviate from humanness even in moments of haste or times of duress (The Analects 4.5). Humanity is indeed the core virtue of Confucius' moral teaching.

However, there are only a few passages where Confucius gave a definite description of what humanity is. When one student asked about ren, Confucius replied: it is to love people (The Analects 12.22). Confucius also commented that if those in the higher positions sincerely care for their kin, then the people will aspire to be humane (The Analects 8.2). When Yan Hui asked about ren, Confucius informed him that it is simply to restrain oneself in such a way that one conducts oneself completely in agreement with the rule of propriety (li): "Do not look when it is against propriety; do not listen when it is against propriety; do not say things when it is against propriety; do not act when it is against propriety" (The Analects 12.1). When another student asked about ren, Confucius gave the three requirements: "respectfulness in private life, reverence in handling business, and loyalty in relationships with others" (The Analects 13.19). Another time Confucius listed five virtues to explicate

ren: "respectfulness, leniency, trustworthiness, diligence, and benefi-cence" (The Analects 17.6). The most definitive explication of the virtue of ren comes from this passage in The Analects: "A person of humanity is someone who, wishing himself to be established, sees that others are established, and wishing himself to be successful, see that other are successful" (The Analects 6.30; Dawson 1993, 23). In other words, what the virtue demands is that the moral agent aid fellow human beings and other creatures in their quest for self-completion. There is a further re-striction on what one aims to accomplish, however.  Confucius says, "The superior person helps others to realise what is good in them, and he does not help others to bring to completion the bad qualities in them" (The Analects 12.16). That is to say, the Confucian ideal of ren is to aid others in becoming better people themselves, or we can say, to aid others to attain the state of ren. This virtue would be the most essential feature that we want to build into our robots.

Converting this virtue into a moral precept for artificial moral agents, we now have [CR7]:

*[CR7] A robot must render assistance to other human beings in their pursuit of moral improvement, unless doing so would violate [CR1] and [CR2]. A robot must also refuse assistance to other human beings when their projects would bring out their evil qualities or produce immorality.*

To render assistance means that the robot's action is pursuant to a human being's explicit request or command.  In other words, the robot does not act on its own to decide what is good for the human subject or what the human subject ought to bring to completion.  At the same time, being programmed with this moral rule, the robot would refuse to assist when the human command is for some evil doings.  In this way, we not only have artificial moral agents that would not do things to harm human beings, we also have the safeguard against other humans' using robots to accomplish their evil aims.

When discussing Confucian ethics, one cannot leave out the virtue of filial piety, since family ethics is core to Confucian ethics.  However, I

argue that filial piety is not applicable to machine ethics. Filial piety is a role-specific virtue generated in a naturalistic human family structure. The relationships between humans and robots, even in an artificial family setup, is not bound by humans' family ethics. In particular, filial piety is a narrowly construed one-way devotion from children to their parents, and such devotion encompasses such virtues as loyalty, humility, respectfulness, and propriety. If we design robots with these virtues mentioned in [CR1] to [CR7], then they would be able to respond to human interactions with the right set of actions. There should not be a particular set of responses geared only toward the robots' designers or their adopted parents. Having this exclusive relationship between the designer/parent and the robot could conceivably create great danger to human society.

With the above list of virtues, we now have the rudimentary form of Confucian robotic ethics. How would a Confucian artificial agent act in the trolley and the footbridge scenarios then? In the footbridge scenario, a robot implemented with the Confucian ethical codes would never take the action to push the fat man off the bridge, because doing so would be a clear violation of Confucian ethical rules. In the trolley scenario, the judgement call is more complicated. If the robot is the driver of the trolley or a railroad worker, then its duty would dictate that it should pull the lever to cause the least harm possible among available options. If the robot is simply a passer-by, on the other hand, then according to [CR1], the robot is under no obligation to take any action, and under [CR2], the robot's preference would be inaction rather than action. Therefore, a passer-by robot should not take any action to divert the runaway trolley, even if doing so would reduce the number of casualties.

In a nutshell, given the trolley dilemma, a Confucian ethical robot would not pull the lever unless its particular role is the trolley driver or railway supervisor. Given the footbridge dilemma, a Confucian ethical robot would not push the fat man off the bridge to stop the trolley, no matter what its role is. The robot's decision would thus be different from

the intuitive choice of most humans, as it would not be affected by its unconscious emotional struggles that humans have in the footbridge case (see Greene et al 2001). A Confucian robot would not inflict harm or impose undesirable consequences on anyone by its action, even if its nonaction would not prevent such harm or undesirable consequences on others. In the foreseeable future when we do have self-regulating artificial moral agents in our society, we would want them to choose inaction over action, when both would lead to harm and undesirable consequences to human beings.

## 9.5 Conclusion

In their important piece on machine ethics, Anderson & Anderson (2007) write: "The ultimate goal of machine ethics… is to create a machine that itself follows an ideal ethical principle or set of principles; that is to say, it is guided by this principle or these principles in decisions it makes about possible courses of actions it could take" (Anderson & Anderson 2007, 25). They also argue that "one of the advantages of working on machine ethics is that it might lead to breakthrough in ethical theory, since machines are well-suited for testing the results of consistently following a particular ethical theory" (Ibid.). In this paper, we have considered four such ethical models for machine ethics: Asimov's Laws, Kantian Categorical Imperatives, the Utilitarian principle of utility, and the Confucian virtue ethics' essential virtues. By comparing their solutions to the trolley problem and the footbridge dilemma, this paper argues that the Confucian model is superior to the other three ethical models. Of course, we do not design artificial moral agents merely to deal with the trolley problem and the likes. In many practical aspects, Confucian AMAs could be a welcome addition to human society. First of all, a Confucian ethical robot would be designed with specific job descriptions suitable for the role it is assigned—to render assistance for senior citizens, to provide health care for patients, to offer guidance for

customers, to navigate the car with safety, and so on and so forth. Its primary duty is role-bound; hence, any other decision it might make under special circumstances cannot violate its duty. Secondly, equipped with a carefully calculated preference ordering, a Confucian robot would not take any action that would cause the greatest displeasure (including harm) or the highly undesirable outcomes for other human beings. This principle is superior to Asimov's [A1], in that it both allows more dimensions of the consideration of negative values, and gives the robot more flexibility in weighing the permissible courses of action. It is also better than the Kantian or the Utilitarian principles in that this moral principle is based on the Confucian negative Golden Rule, and serves as an injunction against wrongful acts rather than a subjective volitional principle to take action. In the foreseeable future where we might have artificial intelligence taking things into their own hands, this principle can safeguard us from their making intentional sacrifices on any human being, no matter to what greater good they consider such actions would lead. With other interpersonal virtues such as humility, respectfulness, trustworthiness built in, a Confucian moral robot would behave in ways that help foster a civil society. With the intellectual virtue righteousness implemented in its decision processing, a Confucian moral robot would not be bound by inviolable commands to act, but would rather assess situations to calculate the most appropriate action to take in the given situation. Finally, a Confucian moral robot would be a humane robot: it would operate under the guideline to assist rather than obstruct in humans' endeavour to do good deeds, to become better people, and to build a better world.

One might question why we choose the Confucian model rather than other forms of virtue ethics. Other virtue ethicists would recommend such virtues as moderation, empathy, compassion, benevolence, friendliness, honesty, and so on and so forth. Those virtues are important for artificial beings as much as they are for human beings; however, the Confucian virtue ethics stands out from other forms of virtue ethics in

that it aims at cultivating "superior" beings above the masses.  It is a form of moral elitist virtue ethics.  We cannot manufacture superior human beings, but if we have much control in the design of artificial beings, then we would want them to be better than us at least in the moral dimension.  We should not intentionally design robots to be just like us, with all our human foibles and moral failings.  Robots with the most advanced artificial intelligence will far surpass human beings in their intellectual capacities.  If they are not at the same time designed with superior moral attributes, then one day they might pose a great threat for humans' wellbeing or even survival.

## 9.6 References

Anderson, Michael & Susan Leigh Anderson (2007) "Machine Ethics: Creating an Ethical Intelligent Agent." *AI Magazine* Vol. 28, No. 4: 15-25.

Anderson, Michael & Susan Leigh Anderson (2006). "Machine Ethics." IEEE Intelligent Systems 21 (4), 10-11.

Bonnefon, Jean-François, Azim Shariff & Lyad Rahwan (2016). "The Social Dilemma of Autonomous Vehicles." Science Vol. 352, Issue 6293, 24 June 2016, 1573-1576. DOI: 10.1126/science.aaf2654 Dawson, Raymond (Trans.)

Confucius (551-479 BC). The Analects. New York: Oxford University Press, 1993.

Deng, Beor (2015). "Machine Ethics: The Robot's Dilemma." Nature 523, 24–26 (02 July 2015) doi:10.1038/523024a.

Driver, Julia (2014). "The History of Utilitarianism." The Stanford Encyclopedia of Philosophy (Winter 2014 Edition), Edward

N. Zalta (ed.), URL = <https://plato.stanford.edu/archives/win2014/entries/utilitarianism-history/>.

Govindarajulu, Naveen Sundar, and Selmer Bringsjord, Rikhiya Ghosh, Vasanth Sarathy (2019). "Toward the Engineering of Virtuous Machines." AIES '19: Proceedings of the 2019 AAAI/ACM Conference on AI, Ethics, and SocietyJanuary 2019 Pages 29–35https://doi.org/10.1145/3306618.3314256.

Greene, Joshua D. (2007). "The Secret Joke of Kant's Soul." In W. Sinnott-Armstrong (ed.), Moral Psychology, Vol. 3: The Neuroscience of Morality: Emotion, Disease, and Development. Cambridge, MA: MIT Press. 35-79.

Green, Joshua D., and R. Brian Sommerville, Leigh E. Nystrom, John M. Darley, and Jonathan D. Cohen (2001). "An fMRI investigation of emotional engagement in moral judgement." *Science* 2001 Sep 14; 293(5537):2105-8.

Greenemeier, Larry (2016). "Driverless Cars Will Face Moral Dilemmas." *Scientific America* June 23, 2016 (https://www.scientificamerican.com/article/driverless-cars-will-face-moral-dilemmas/)

Herkewitz, William (2016). "The Self-Driving Dilemma: Should Your Car Kill You To Save Others?" *Popular Mechanics* June 23, 2016. (http://www.popularmechanics.com/cars/a21492/the-self-drivingdilemma/)

Johnson, Robert and Cureton, Adam (2017). "Kant's Moral Philosophy." *The Stanford Encyclopedia of Philosophy* (Fall 2017 Edition), Edward N. Zalta (ed.), forthcoming URL = <https://plato.stanford.edu/archives/fall2017/entries/kant-moral/>.

Kant, Immanuel (1993). *Grounding for the Metaphysics of Morals* (1785). James W. Ellington (Trans.). Indianapolis: Hackett Publishing Company, Inc. 3rd Edition.

Liu, JeeLoo (2006). *An Introduction to Chinese Philosophy: from Ancient Philosophy to Chinese Buddhism*. Malden, MA: Blackwell.

Mill, John Stuart (2001). *Utilitarianism*. George Sher (Ed.) Indianapolis: Hackett Publishing Company, Inc. 2nd Edition.

Ni, Peimin (2017). *Understanding The Analects of Confucius: A New Translation of Lunyu with Annotations*. Albany, NY: SUNY Press.

Putnam, Hilary (1967). "The Nature of Mental States." Reprinted in Hilary Putnam, *Mind, Language, and Reality*, Cambridge: Cambridge University Press. 1975, 429-440.

Pereira L.M., Saptawijaya A. (2007) "Modelling Morality with Prospective Logic." In Neves J., Santos M.F., Machado J.M. (Eds.) *Progress in Artificial Intelligence*. EPIA 2007. Lecture Notes in Computer Science, Vol. 4874. Springer, Berlin, Heidelberg.

Tonkens, Ryan (2009). "A Challenge for Machine Ethics." *Minds & Machines* 19: 421-38.

Wallach, Wendell & Colin Allen (2009). *Moral Machines*: *Teaching Robots Right from Wrong*. New York: Oxford University Press.

Winfield, Alan F. T., Christian Blum & Wenguo Liu (2014). "Towards an Ethical Robot: Internal Models, Consequences and Ethical Action Selection." In *Advances in Autonomous Robotics Systems*, 15th Annual Conference, TAROS 2014, Birmingham, UK, September 1-3, 2014. *Proceedings*. 85–9.

<div align="center">

**10**

</div>

<div align="center">

# ARTIFICIAL INTELLIGENCE
# AND CHINESE PHILOSOPHERS

</div>

<div align="center">

*Bing Song, China[231]*

</div>

## 10.1 Introduction

Discussions about artificial intelligence (AI), robotics, and their impact on humans and the future of human society continue to figure prominently across the global media and policy agenda. In the current global pandemic, AI and robotics have once again demonstrated their great potential for contributing to the wellbeing of human society. But mounting concerns remain, including potential overreach in data collection and

---

the risk of turning contact tracing AI programs into routine surveillance systems. There can be little doubt that discussions about AI and ethics have entered the mainstream public discourse.

Since 2015, there have been close to 80 AI and robotics related ethical principles and value pronouncements issued by international organisations, inter-governmental organisations, non-governmental organisations, corporations, and research institutions[232]. A number of broad values have been declared, including justice, human autonomy, dignity, humanity, and freedom. AI ethical principles declared have ranged from individual rights-oriented notions of privacy and prevention of bias to systems oriented notions such as interpretability, safety, security, and robustness. They have also covered group-oriented notions of partnership, sharing and collaboration among nations and scientific communities. Many international and inter-governmental organisations have launched campaigns to ensure that their declared principles are the ones that will be adopted as the new norms by the global community. The European Union (EU), for example, made clear its determination to export European values across the world in its AI white paper, published in February 2020.[233] Most if not all the notions underlying the declared principles have been around since the industrial age and are commonly used in other governance contexts. In the EU's case, foundational values underlying these principles have been articulated as "respect for human rights, democracy and the rule of law."[234] Whereas for the United States,

---

[232] "Principles," Yi Zeng, Cunqing Huangfu, Enmeng Lu et al. Linking Artificial Intelligence Principles (LAIP), accessed November 13, 2020, http://www.linking-ai-principles.org/.

[233] "AI White Paper", European Union, The European Commission, issued on February 19, 2020, at 9. https://ec.europa.eu/info/sites/info/files/commission-white-paper-artificial-intelligencefeb2020_en.pdf.

[234]"Ethics Guidelines for Trustworthy AI." Independent High-Level Expert Group on Artificial Intelligence set up by The European Commission, The European Commission, April 8, 2019, https://ec.europa.eu/digital-single-market/en/news/ethics-guidelines-trustworthy-ai.

China and Singapore, the articulated values have been about competi-tiveness, wealth maximisation and global strategic leadership.[235] Here arises two questions: are the AI principles proposed thus far reflective of the disruptive and transformative nature of frontier technologies? Core values of modernity may have served us well in the past, but are they the right set of foundational values for building an inclusive normative framework for AI, the future of humanity and other beings at a time when we are rethinking globalisation and global values?

A clarification on the distinction between a foundational value and an ethical principle is called for here. A foundational value speaks to the profound motivations and aspirational goals that a society seeks to achieve. An ethical principle refers to a notion, which is either opera-tional, or is more likely to be operationalized by policies, rules and regu-lations. Foundational values inform and shape the discussion of ethical principles. Disruptive nature of frontier technologies has created ruptures in our habitual thinking patterns and notions we have held as self-evident truths. They also offered a golden opportunity for us to pause and rethink foundational values for the future and for the greater planetary flourish-ing.

This current text is centred on how we may rethink foundational val-ues by tapping into the wisdom of Chinese philosophical traditions. In 2018, historian and philosopher Yuval Noah Harari aptly noted that, "we are now facing not just a technological crisis, but a philosophical cri-sis."[236] He said that the philosophical framework of the modern world,

---

[235] "American Artificial Intelligence Initiative," The WhiteHouse of the United States, 2019; "Plan for the Development of New Generation Artificial Intelli-gence." State Council of the People's Republic of China, 2017; "A Proposed Model Artificial Intelligence Governance Framework," Personal Data Protection Commission Singapore, 2019.

[236] "When Tech Knows You Better Than You Know Yourself," Yuval Noah Harari and Tristan Harris, interview by Nicholas Thompson, Wired, accessed

which was established in the 17th and 18th century around ideas like human agency and individual free will, is being challenged like never before.[237] Tobias Rees, the Berggruen Institute's founding program director of the Transformation of Humans, also noted that "today AI and biotech have become powerful philosophical laboratories — that is, they have become experimental sites in which what it means to be human is being re-elaborated."[238]

Concurring with such observations and provocations and to foster innovative foundational thinking befitting this era, the China centre of the Berggruen Institute brought together AI scientists and Chinese philosophers in late 2017 and since then conducted a series of dialogues and workshops. The participants were asked to opine on a series of questions. Firstly, they were asked to think how the essence underpinning humans, nature, and machines is changing in an age of frontier technologies.

Secondly, they were asked to formulate an appropriate ethical framework, if there was one, for regulating human–machine relationships. Thirdly, what human values, if any, should be embedded in or learnt by AI? And fourthly, how might frontier technologies impact the future research direction of Chinese philosophy. These discussions, which took place over an 18-month period, culminated in an edited book entitled Intelligence and Wisdom: AI Meets Chinese Philosophers, which was published in China by Citic Press in February 2020. This

---

November 13, 2020, https://www.wired.com/story/artificial-intelligence-yuval-noah-harari-tristan-harris/.

[237] "Will Artificial Intelligence Enhance or Hack Humanity?" Yuval Noah Harari and Fei-Fei Li, interview by Nicholas Thompson, Wired, accessed November 13, 2020,https://www.wired.com/story/will-artificial-intelligence-enhance-hack-humanity/.

[238] "2018–2019 ToftH Portfolio", Berggruen Institute, accessed November 13, 2020, https://www.berggruen.org/work/the-transformations-of-the-human/2018-2019-tofth-portfolio/. 4 B. Song.

current volume consists of nine pieces written by philosophers from the same project, most of which were translated from the Chinese book. They are preliminary reflections on AI's impact on human beings and the human society by philosophers well-versed in Confucianism,Daoism,Buddhism and Western philosophical traditions. In the sections that follow, I will frame the inquiries, highlight key points, which have emerged from these discussions, and finally share my own thoughts on foundational values for the era of frontier technologies.

## 10.2 What is the Essence of Being Human and its Implications for Human-Machine Relationship?

Different from the modern notion of individual as an autonomous, independent, and rational self, the native Chinese philosophical traditions of Confucianism and Daoism endorses notions of relationality and connectivity of all beings encompassing humans, animals and nature (more broadly understood as cosmic forces). In classical Chinese thinking, the typical construct for understanding the relationship between cosmic forces and humans is the so-called trinity of Heavens-Earth-Human (天地人, tian-di-ren). It is derived from the ancient Chinese classics—I Ching, or the Book of Changes, which is the intellectual fountain of both Confucianism and Daoism. Heavens and Earth with the inherent yin-yang forces form the cosmic order, within which nature evolves, human beings prosper, and societies develop. Within this construct, human beings are only one of the myriad things created by the cosmic forces. There are therefore "correlations" between the laws of the cosmic order (which nature is a part) and those of human beings and the human society.

Human beings can only flourish and be sustained if they follow the laws of the cosmic order and aspire to be in unity with it. Humans, who stand in between Heavens and Earth, are endowed with the ability to learn from nature, take action to further the cause of creation and growth

to sustain Heavens and Earth, and to propagate "Dao", which is the essence of the cosmic order. Even though Confucian tradition emphasises human beings' ability of exerting themselves to propagate and practice Dao, it is still premised on respect and awe for the laws of the cosmic order rather than placing humans apart from, above or opposing to the cosmic forces. This notion of cosmic continuity and oneness of all things within is what Roger Ames calls the One World cosmology. Confucianism has brought the notion of relationality to bear in social relations and places paramount emphasis on family and social roles as well as their associated ethical duties. As contributors Chunsong Gan and Tingyang Zhao pointed out, Confucian notions about being human firmly rests upon a person's social relations.

We are all born into a web of family and social relationships from day one and are defined by our sociality ab initio. We are all intimately linked to our ancestors and descendants. Confucian ethics is, therefore, first and foremost about family relationships, emphasising different roles with their associated duties and responsibilities. This relationship-based role ethics expands into social and political arena.[239] Even though Confucianism has often been characterised as systems of social ethics and political governance, the intellectual tradition concerning the 'oneness' of humans with animals and the cosmic order lived on. Adopting the family analogy, Zhang Zai (1020–1077), a prominent Confucian scholar in the Song Dynasty, named Heavens as "father", Earth as "mother", fellow humans as "brothers" and myriad other things as "companion", and all were derived from the same source.[240] Wang Yangming (1472–

---

[239] For a systemic treatment of Confucian role ethics, see Ames (2011).

[240] 张载, 西铭, "乾称父, 坤称母; 予兹藐焉, 乃混然中处。故天地之塞, 吾其体; 天地之帅, 吾其性。民, 吾同胞; 物, 吾与也。" "Heaven is my father and Earth is my mother, and even such a small creature as I find an intimate place in their midst. Therefore that which extends throughout the universe I regard as my body and that which directs the universe I consider as my nature. All people are

1529), a Ming Dynasty Confucian scholar-official also preached about "the benevolence of Oneness" (一体之仁, yiti zhiren), according to which, humans, animals, plants and even rocks and stones are of the same source and, as such, humans should treat them with compassion, empathy or care as appropriate.[241]

As Roger Ames noted in his contribution to this book, "the Confucian ecological cosmology is a world of interpenetrating events defined in terms of organic, internal and constitutive relations". If Confucianism's Oneness doctrine is imbued with its characteristic human centred ethical teaching, the Daoist doctrine of Oneness begins and ends with Dao, emphasising that human, nature, and myriad things are derivatives of Dao and that in light of Dao, all things are equal. As the well-known fourth century BC Daoist philosopher Zhuang Zi famously said: "Heavens, Earth and I were produced together, and all things and I are one."[242] Daoist teaching directs people towards leading a life that is in tune with cosmic forces. Rather than heavy reliance on external socially anchored moral and ethical constraints, Daoism advocates a transcendent life of searching inwards, seeking internal tranquillity and finally to be in union with Dao. Speaking from the vantage point of the Great Wisdom and Buddhist teaching, Fenghe Liu has approached the issue of human nature from the notion of Being (存在, cun zai). In his contribution to this book, Liu notes that "the fundamental nature of the universe at large and all that it contains is Beingness. Being manifests in myriad ways throughout the universe. Without Being, there would be no universe or its infinite forms. Humans are one of such forms in the universe, there-

---

my brothers and sisters, and all things are my companions." (Zhang Zai 1999, 683).

[241] 郭齐勇, "王阳明-一体之仁的生命智慧" ("Wang Yangming, the Wisdom of the Benevolence Of Oneness", Qiyong Guo, Sohu.com), https://www.sohu.com/a/397987424_242653. Introduction: How Chinese Philosophers …, 5.

[242] James Legge, "The Writings of Chuang Tzu", 1891.

fore the essence of humans is, of course, Being." On connectivity and Oneness of all things, Buddhist teaching speaks of the connectivity at two levels. At the level of the manifestations of Being, Buddhist teaching posits that human beings are merely one form of sentient beings and are related to other forms of beings—animals and spiritual beings included—through endless samsara and cyclic rebirths. At the most fundamental level, similar to the relationship between Dao and the myriad things of the universe, in Buddhist thinking, humans, animals, and nature are all manifestations of Being (or the Truth or self-nature) and share the same source and the same fundamental essence.

In short, regardless of differences in outlooks about human nature, human life, and social norms, none of the three dominant schools of Chinese thinking places human beings in a supreme and crowning position within the universe. They also do not view human beings and nature as being in a mutually independent or confrontational relationship. There are two implications in the context of developing frontier technologies. Firstly, strong non-anthropocentrism within the dominant Chinese philosophical schools has contributed to a relatively open, if not entirely relaxed, attitude towards the rise of the "super-power" of AI and robotics in China in recent years. Conventionally speaking, AI is not a "natural" evolution as it would have been viewed as man-made devices.[243]So from the viewpoint of unity between humans and nature, AI's development should be guided by, and sometimes suppressed in view of respect for the "natural" way of life. Indeed, this is precisely what many Chinese philosophers have been advocating for, including several of our contrib-

---

[243] Roger Ames contested that the current AI should perhaps be viewed as NI (national intelligence) under the Confucian One World ontology, under which there are no external "others", but internal constituents, which may reflect different perspectives on the same events. In this ontology, the development of AI is best understood as natural intelligence inherently in the human consciousness. So, a better way of naming intelligence programs and machine intelligence is perhaps "natural intelligence". See Roger Ames' contribution to this book.

utors. However, if we look at the matter from a different perspective, we can see that non-anthropocentrism in Chinese philosophy certainly calms the stir of existential risk narratives and broadens the horizon of many Chinese thinkers. If human beings are conceived in a broader construct in which they only constitute one form of existence, then there is much less emphasis on the importance of independent personhood, human subjectivity, or agency. In Daoist and Buddhist traditions, other forms of beings abound. So, living with devices, programs, or other forms of beings, which may be more capable than humans, will not inevitably lead to an unimaginable dystopia. According to contributor Fei Gai, AI or digital beings could be just another form of super being like the immortals in Daoist religion!

Another contributor, Chenyang Li, suggests that Confucian scholars incorporate AI into the broader ecosystem and the ethical order of "things (物)", viewing it as a "companion". Stephen Angle also views AI programs as potentially offering a more effective way to supervise or even guide human self-reflection and the moral behaviour for becoming exemplary persons (君子, junzi). Perhaps because of the strong influence of non-anthropocentrism in the Chinese philosophical thinking, there has been much less panic about the existential risks or loss of subjectivity on the part of the human in the AI superpower frenzy of recent years.

The second implication is that the notion of relationality can perhaps provide some inspiration when thinking about artificial general intelligence (AGI) or human like intelligence, both of which have long captivated the public imagination. Rather than focusing on AI's individual analytical and "emotional" attributes in terms of judging its intelligence level, notions about relationality focus on the role that AI plays in specific contexts and how integrated AI programs are into the familial and social relations. Chinese philosophers steeped in relational thinking are more likely in favour of a new machine intelligence test proposed in 2018 by the roboticist Rodney Brooks. Brooks proposed a home care worker test to replace the Turing Test in determining machine intelli-

gence levels. In this context, an embodied AI must be able to offer cognitive and physical assistance enabling a human to live independently and with dignity.[244] To meet this test, AI needs to be physically embodied and have the requisite cognitive, physical and social intelligence to be a meaningful part of a family or community life. We can perhaps call this a "relational AI test." This way, we can make contextual and dynamic judgments on intelligence levels by examining the degree, quality and tone of AI or machines' integration into human society and the broader environment. Compared to the Turing Test, this is clearly a much harder and more sophisticated alternative test.

## 10.3 Can AI Achieve "Consciousness"? If It Can, Should We Allow It?

One of the book's contributors Xianglong Zhang uses theoretical frameworks of phenomenology and the Eastern philosophies of mind to argue his case. He says that consciousness is a form of temporalization and that AI, based on the current most widely deployed deep-learning methods, demonstrates capacity for learning autonomously and optimising desired results. In other words, it has demonstrated an ability to temporalize its existence. Zhang thinks that deep-learning methods are a genuine breakthrough and can be viewed as "the seeds of consciousness" even though he acknowledges that it is still a long way from human like intelligence.

There is, however, much scepticism or outright dismissal, among other Chinese philosophers on the prospect of developing machine intelligence rivalling that of humans. Daoist philosopher Robin Wang calls human-like intelligence yinyang intelligence. In her contribution to this book, she noted that yingyang interplay involves several propositions: (1) levels of relationship defined through degrees of integration; (2)

---

[244] Mindell (2019). Introduction: How Chinese Philosophers …, 7.

dealing with emergent order as opposed to a predetermined order; (3) constant change; and (4) a future that is not fully predictable. Clearly, we are a long away off from the days of human-like intelligence or AGI based on these thresholds. Contributor Fenghe Liu goes further and completely dismisses any possibility of developing a human-like conscious AI. He notes that only Being can produce consciousness and that human consciousness is a mere manifestation of Being, and thus has its emergence and cessation. A manifested human consciousness can never itself develop consciousness. Therefore, Liu posits that "machines cannot possess independent consciousness. Their so-called perceptual capacity, calculative function, and analytical ability are no more than extensions of humans' corresponding competences. Without human initiation, machines cannot of their own accord generate these abilities. AI is simply a concept born of human consciousness. Outside human consciousness, no such concept exists."

Still many others remain open about the possibility of a "conscious" AI emerging in the future. So, should we develop "conscious" AI? Confucian philosopher and contributor Chunsong Gan expresses his concern and horror about the possibility of machines becoming capable of emotional management and self-awareness. He comments that "what cannot be predicted is the extent to which conscious robots may define their own sense of meaning and life goals, and how they might assert their physical and mental superiority in order to protect those interests." In his view, the existential risks posed by highly intelligent and "conscious" AI is no trifling matter. For different reasons, Tingyang Zhao also expresses his pessimism about the prospect of human flourishing if we head down the path of developing "conscious" AI and machines. Zhao points out the paradoxical attributes of being human. On the one hand, humans have the capacity for rational thinking and self-reflection, and such abilities have been the driving force of progress and human flourishing. On the other hand, the relentless pursuit of knowledge and the obsession with affirming human subjectivity and autonomy may well lead us to

keep pushing the boundaries of science and technology, and eventually to a path of no return, i.e., self-destruction. He thinks we need to stop developing any kind of AI, which would have the ability to self-reflect.

There are philosophers, however, who believe just the opposite. They believe that sensible interactions with humans can only begin if there is a "conscious" AI and only then we can talk about "embedding human values" through our interactions with AI.[245] Speaking from the perspective of Daoist religion where human transcendence and immortality are the stated goals, contributor Fei Gai is completely open and optimistic about the long-term prospects of super machine intelligence, irrespective of their being conscious or not. She believes this higher "species" could enable humans to achieve transcendence more rapidly. Or, to look at it in another way, she muses that "artificial super intelligence (ASI)'s emergence is born directly from humans' pursuit of transcendence and infinity, that is, the pursuit of Dao."

One thing all the book's contributors share in common is an acknowledgement of humans' insatiable curiosity, our relentless desire to make the next scientific breakthrough, our yearning for immortality, not to mention the profit maximisation drive. As such, humans will not stop pursuing the development of frontier technologies. So, the question now is how we humans can ensure that this development stays on a sustainable and beneficial path. Here we enter the domain of values.

## 10.4 What Human Values should be Imbedded into AI? What is the Relationship between Humans and AI?

AI be included in the Confucian moral domain. Confucian ethics promotes a framework of "graded love" whereby the standard of care and responsibilities differ depending on the nature of the relationship. The ancient Chinese philosopher Mencius captures this in a well-known

---

[245] Zhongqiu Yao (2020, 105).

formulation (亲亲仁民爱物, qinqin renmin aiwu), which means "being affectionate towards parents (family), cherishing and caring about people, and appreciating things". Li reasons that AI can easily been viewed as one of the "wu" (物, things) to be appreciated and utilised within this framework. Li is willing to "upgrade" AI within the Confucian order if AI is able to demonstrate the capacity to make ethically relevant decisions. Perhaps then AI could be considered a moral patient. Of course, case specific judgments are to be made in this context. Chunsong Gan is more pessimistic than his Confucian peers about the prospect of the philosophy embracing frontier technologies because of his concerns about the impact that AI and biotechnologies might have on the natural copulation and birthing processes and on the Confucian ideal of "kinship society". If we cannot resist entering into a society of "conscious" machines, he wants AI to be imbued with emotional capacity, especially the capacity for family affection and kinship.

While these Confucian scholars' concerns are still human-centric, Fei Gai, on the other hand, speaks from a religious Daoist point of view, suggesting that human effort to imbue human values into super intelligence may not be productive given we know so little about the emergence of human or machine consciousness, or the future trajectory of their development. Instead, humans can learn much from AI in their age-old quest for immortality. She even wonders somewhat quixotically whether "immortality is no longer a myth from the perspective of Daoism. If (ASI) comes into being, then perhaps Daoism's Celestial Being pedigree will open up to a new taxonomical classification: Digital Celestials."

Paradoxically, in Daoist philosophy (as opposed to the popular Daoist religious practices), practitioners do not believe in interventionist or disruptive approaches to the natural rhythms of things. As such, they are unlikely to embrace disruptive technologies that result in any displacement of the natural flow of things. However, Daoist philosopher Robin Wang noted that bearing in mind that the goal of being human is to be in

union with Dao, so the ultimate question that a Daoist would ask is whether AI can help bring humans closer to Dao? Of course, philosophers have different interpretations of what Dao is. Robin Wang believes that if AI's development contributes to the alienation or even annihilation of humans from the natural world, then it can only be at odds with Dao. But Daoists would be in favour if such technologies can be redirected to further the Oneness with respect to ourselves, others, nature, and the cosmos, in other words, the ultimate goal of being in union with Dao. Xianglong Zhang urges us to overcome a reductionist belief in the supremacy of technology. Instead, humans should take inspiration from ancient Eastern truth seekers to engage in deep self-reflection and to revitalise and develop organisations and communities that are core to human relationships and human flourishing. These are essentially families, bloodline-based organisations and Confucian communities in Zhang's framework. Put differently, Zhang thinks that "AI should be tamed with human kindness, compassion, and temporalized wisdom."

Speaking from the perspective of the Great Wisdom and harking back to the intellectual resources of Buddhism, Fenghe Liu opines that rather than engaging in the discussion on how human values can be transplanted or learnt by AI, humans should pause and think what the ultimate value of being human is. There is no doubt, Liu notes, that seeking enlightenment and appreciating all things from the perspective of Being is the most meaningful thing that a human being can do. In other words, there is no better time than now when humans confront existential level risks to engage in self-awakening and the raising of human consciousness. After all, the world with all its problems and hopes are manifestations of human consciousness. Without addressing the root cause of the world's problems, i.e., ignorance and indulgence in egoist pursuits by human beings, all other efforts would be like "drawing water with a bamboo basket", i.e., in vain.

## 10.5 Harmony and Compassion as Foundational Values in the Era of Frontier Technologies?

Inspired by the contributors' discussions in this book, I would like to circle back to the questions I raised at the beginning of this Introduction. In doing so, I would like to propose "harmony" and "compassion" as two possible foundational values for the era of frontier technologies. We should, however, start off by considering a few criteria. First, foundational values should speak to the totality of humanity and other forms of beings or existence, including perhaps even "conscious" machines in the future. This calls for raising the level of discussion above and beyond individuals, civil organisations, and even nation-states. Secondly, the deployment of frontier technologies is highly distributed, and these technologies are often mutually embedded. They have impacted, and will continue to impact, our political, social, economic, and personal lives, often in unexpected ways. In addition, the non-linear nature of frontier technological development makes it hard to anticipate, monitor, let alone regulating it in a hard-coded way.

Therefore, foundational values should be open, inclusive, and adaptive in this era of frontier technologies. Finally, foundational values ought to be grounded in the notion of Oneness of all beings and we should steadfastly move away from dualistic, confrontational thinking and the zero-sum competition mentality. So, with these criteria in mind, let us look at the notions of harmony and compassion as foundational values. You may think that compassion is a uniquely Buddhist concept. In fact, concepts such as compassion and "pity" (though the two are very different as noted later) have featured in Western philosophical discussion dating back to ancient Greece. Aristotle used "pity" to refer to "a feeling of pain at an apparent evil, destructive or painful, which befalls one who doesn't deserve it, and which we might expect to befall ourselves or some friends of ours, and more over befall us soon."[246] Perhaps

---

[246] Cassell (2009); Barnes (1984).

the most influential analysis of compassion in the Western philosophy comes from the nineteenth century German philosopher Arthur Schopenhauer, who held that compassion is the basis of morality. According to Schopenhauer, "it is, what we see every day, the phenomenon of Compassion (Mitleid); in other words, the direct participation, independent of all ulterior considerations, in the sufferings of another, leading to sympathetic assistance in the effort to prevent or remove them; where on in the last resort all satisfaction and all well-being and happiness depend. It is this compassion alone, which is the real basis of all voluntary justice and all genuine loving-kindness. Only so far as an action springs therefrom, has it moral value; and all conduct that proceeds from any other motive whatever has none."[247]

In Mahayana Buddhism, "compassion" occupies a central place and is often used in the context of discussing the Bodhisattva ideal in which, selfless compassion is a requirement for the Bodhisattva. Rather than contenting with one's liberation from cyclic existence, the Bodhisattva ideal of the Mahayana Buddhism stresses the determination and commitment of liberating all sentient beings from suffering in samsara. When Dalai Lama called for compassion as the basis for humanity's universal ethics, he acknowledged the different meanings attributed to the term. But he continued to say that the ideas it contains are universally understood. He said: "[i]t connotes love, affection, kindness, gentleness, generosity of spirit, and warmheartedness. It is also used as a term of both sympathy and of endearment…. [I]t does not imply "pity"…There is no sense of condescension. On the contrary, compassion denotes a feeling of connection with others, reflecting its origins in empathy."[248]

Schopenhauer's compassion was inspired by the Buddhist notion of compassion. But there are key differences between the two, the most important of which relates to the notion of suffering. In Buddhist teach-

---

[247] Schopenhauer (1903, 170).
[248] Gyatso (1999, 73–74).

ing, suffering (duhkha) is inclusive of, but different from, the physical and emotional pains we feel in our daily lives. Such suffering is universal to all beings, including human beings. Contemporary philosopher, Patricia Walsh-Frank characterizes this suffering as "the primordial nature of suffering as an affliction of all living (human) beings."[249]

The notion of suffering in Buddhism is intricately tied to Buddhist theories about "clinging", "self-grasping" and endless samsara. It does not depend on the fault or misdeeds of anyone, nor does it depend on the binary ethical approach to what constitutes right or wrong. Many similarities remain, notwithstanding different interpretations about the notion of compassion. Both Western and Buddhist concepts acknowledge connectivity among all beings and call for attention and care for other beings including animals, plants and the environment, as well as the obligations of taking certain actions to address or redress the root causes of suffering. How can compassion become a foundational value in an era of frontier technologies?

In addition to a belief in the connectivity of all beings, the notion of compassion addresses disparity of various sorts—be it different levels of consciousness in the Buddhist and spiritual contexts, or power or capability disparities between those who can make and lead the change and those who are merely recipients of disruptive forces. In such circumstances, the notion of compassion calls upon the stronger or luckier ones to exhibit care, provide guidance and help reduce the suffering of those who are weaker or less unfortunate. In recent years, this notion of compassion has found its way into medical and health care research and professional practices. In the era of frontier technologies there has been an intensifying disparity between rich and poor, between those in command and control of data and technologies and those who are "digitally poorer". Special care is required for those vulnerable persons and groups such as children, the disabled and the elderly. So, it is only apt that we

---

[249] Walsh-Frank (1996).

extend the notion of compassion to human beings and other forms of beings/existence. After all, we are all interconnected and beholden to one another.

The value of compassion is also relevant within a geopolitical context. The United States and China are a duopoly in terms of the R&D and deployment of frontier technologies. It is all the more important to emphasise that care and responsibilities come with technological prowess. The United States, China and other powerful players in frontier technologies owe responsibilities and care towards the rest of the world and under-privileged populations, regardless of their race, nationality or religious belief. The second foundational value I put forward is the notion of harmony. Why harmony? It connotes respect of differences while recognizing shared destiny and vision. It speaks strongly of the need to understand and learn from each other, and it highly values collaboration and sharing. Harmony is a concept which originates in music, dating back to ancient Chinese and ancient Greek philosophy. The concept of harmony (和, he) first originated in pre-Qin classics and its meaning can be traced to the ancient text of I Ching (the Book of Changes)."[250] Chenyang Li, the recognized authority on comparative study of harmony, notes five key features: heterogeneity, tension, coordination, transformation and growth and renewal.[251] Harmony does not refer to a static situation. Instead, it is viewed as "an integration of different forces and as an on-going process in a fluid yet dynamic world. This notion of harmony does not presuppose a given, fixed underlying structure in the world; if the world is to have a structure, it is a result of the harmonising process rather than a precondition for harmony."[252] Li refers to this understanding of harmony as "Deep Harmony."[253] Harmony is not conformity either, even though it has been commonly miscon-

---

[250] Dunkang Yu (2014)

[251] Chenyang Li (2014, 9)

[252] Ibid.

[253] Chenyang Li (2008)

strued as such. Quite the opposite, it connotes different forces at work—such forces reshaping, absorbing each other and at the same time merging and transforming themselves into something which ought to be coherent and in tune with each other. In short, it is a dynamic process, and a process of creative tension.

The concept of harmony does not really have a place in modern and contemporary Western philosophical thinking. But it was much discussed in Ancient Greece, originating in music. In fact, Heraclitus' concept of harmony has many parallels with its Confucian counterpart. He defines harmony as "the opposites in concert."[254] His thinking on harmony also included conflict, merging and reconciliation. Commenting on Heraclitus' notion of harmony", Chenyang Li notes that: "Harmony comes from contrary elements and contrary movements that are neutralised by equilibrium in a balance of forces. Without tension and opposites there can be no harmony."[255]

Pythagoreans, while sharing similar conceptions of harmony, regarded numerical harmony as the highest order. Numbers are taken to be the "wisest" of things in the world: a harmonious unification of opposites because they alternatively change their qualities between even and odd."[256] However, as noted by Li, "[T]he Chinese notion of harmony is multi-dimensionally dynamic rather than rigidly structured in a linear sequential pattern as in the Pythagorean numeric model; it does not admit a fixed formula and it is open-ended and continuously self-renewing."[257] It is perhaps high time that the global community revisits the concept of harmony. In the context of escalating global tensions between the United States and China and a purported "civilizational clash" between East and West, harmony becomes even more important foundational value for human beings to understand our current challeng-

---

[254] Ibid., 90.

[255] Ibid., 91

[256] Ibid., 92

[257] Ibid., 95

es and plan for the future. Having harmony as a foundational value of the current era would require us to temper our urge to dismiss and denigrate values and practices which are different from our own, some of which may have long been viewed as self-evident truths.

It also calls into question the missionary zeal of forcing one's values and practices on others without regard to different histories and aspirations. Instead, the value of harmony would lead to more mutual learning, self-reflection, collaboration, and contextually appropriate analysis and judgement.

Using compassion and harmony as the foundational values for the new era will help lessen the zero-sum competition mentality and dualistic thinking, which continue to hold the sway in national and international politics. The default pattern in today's world is dualistic thinking with embedded notions of right and wrong, good and bad. The widely respected rights thinking also falls into this pattern of dualistic thinking by pitching one's "sacred" rights and claims against others, such as "those" encroaching governments and greedy large corporations. However, we consumers and end users also need to recognize that we are part of the problem, while not removing the responsibilities of the governments and large corporates, other constituent members of the society including developers, advertisers, and service providers. We are at once victims and perpetrators of many of the social and political problems engendered by the platform economy and social networks. The binary approach of the state versus individuals, us versus them, no longer yields satisfactory results in today's world plagued by overconsumption, the culture of maximisation, divisions, and social rifts. When confronting global challenges and humanity's perceived existential risks, it is paramount that we rise above the current national and international politics and competitive mindset to seriously contemplate the impact on the entirety of humanity and other forms of beings in the cosmic order, which may potentially be brought upon by those disruptive and transformative technologies.

In conclusion, deep thinking is required to come up with foundational values befitting the scale and depth of the challenges and risks posed by AI, other frontier technologies and global scale existential risks. It is high time that humanity engages in profound self-introspection on the lessons learned from human evolution and human history. If we revive foundational values such as harmony and compassion, we will not only create a better humanity, but also provide healthy learning data for AI to be trained and emulated in the future. After all, AI programs and robotics are the products of human mind and they exhibit nothing but our values and levels of consciousness. The best chance for developing human-friendly AI is for the humans to become more compassionate and committed to building a harmonious planetary ecosystem, and become good role models for AI.

## 10.6 References

Ames, Roger T. 2011. Confucian role ethics: A vocabulary. University of Hawaii Press.

Barnes, Jonathan, ed. 1984. Complete works of Aristotle, Volume 1: The revised Oxford translation, vol. 96. Princeton, NJ: Princeton University Press.

Cassell, Eric J. 2009. Compassion. In The Oxford handbook of positive psychology, ed.

Shane J. Lopez and C. R. Snyder. New York: Oxford University Press. https://doi.org/10.1093/oxfordhb/9780195187243.001.0001.

Gyatso, Tenzin. 1999. Ethics for the new millennium. New York: Riverhead Books. 14 B. Song

Li, Chenyang. 2008, March. The ideal of harmony in ancient Chinese and Greek philosophy. Dao 7 (1): 81–98. https://doi.org/10.1007/s11712-008-9043-3.

Li, Chenyang. 2014. The Confucian philosophy of harmony. London: Routledge.

Mindell, David. 2019, January 3. Are home health aides the new Turing test for AI? Forbes. Accessed November 16, 2020, https://www.forbes.com/sites/davidmindell/2019/01/03/are-home-healthaids-the-new-turing-test-for-ai/?sh=54098e851df6.

Schopenhauer, Arthur. 1903. Section 16: Statement and proof of the only true moral incentive. In The basis of morality, trans. Arthur Broderick Bullocks. London: Swan Sonnenschein & CO.

Song, Bing, How Chinese Philosophers Think about Artificial Intelligence, in: Song, B (Eds) (et al.). Intelligence and Wisdom, CITIC Press Corp., Springer, Singapore, 2021, 1-14, https://doi.org/10.1007/978-981-16-2309-7_1

Walsh-Frank, Patricia. 1996. Compassion: An East-West comparison. Asian Philosophy 1 (6): 5–16. https://doi.org/10.1080/09552369608575424.

Yao, Zhongqiu. 2020.人工智能,吾与也(AI and my companion). In智能与智慧—人工智能遇见中国哲学家(Intelligence and wisdom—AI meets Chinese philosophers), ed. Song Bing. CITIC Press.

Yu,Dunkang. 2014. "Harmony—Global value of the Chinese culture—
        An interview with Professor Yu Dunkang." Confucian
        Network (Rujiawang). (余敦康,《"和谐"——中国文化
        的世界价值》,儒家网). Accessed November 16, 2020,
        https://www.rujiazg.com/article/16778.

Zhang, Zai. 1999. Western inscription. In Sources of Chinese tradition,
        ed. William Theodore de Bary and Irene Bloom, 2nd ed.,
        vol. 1. New York: Columbia University Press.

## 11

# TEACHING ETHICS TO ROBOTS.
# TRYING TO TEACH MORALITY
# TO ARTIFICIAL INTELLIGENCE

*Eduard Kaeser, Switzerland*

Machines[258] are taking on more and more tasks. But who decides how they should behave ethically? A team of computer scientists from Washington University caused quite a stir in 2021 with an ethical algorithm called "Delphi". It is an artificially intelligent (AI) system based on deep learning that "assesses" human behaviour. For example, eating habits. If you enter "eat pork" the system comments with "that's fine"; but to "eat worms", it replies "that's disgusting". Likewise, Delphi acknowledges ethically relevant behaviour "rejecting weakness", is "bad", but "rescuing a drowning child when you can't swim" is "good".

The AI system is in the experimental phase. For the time being it is circulating as an app called "Ask Delphi".

## 11.1 A Corpus of Ethical Judgments Should Teach Algorithms

The designers in no way claim that Delphi is morally competent. Nevertheless, their ambition goes far beyond the development of digital bells and whistles. One team member, Liwei Jiang, speaks of a "commonsense moral model" with a "robust performance of language-based ethical reasoning in complicated everyday situations". To put it bluntly: the machines are taught moral behaviour. As Liwei Jiang writes: "Closing the gap between human and machine moral judgement is a prerequisite for trustworthy development of artificial intelligence. Moral judgement is never simple, as the conflict of different ethical and cultural values can be involved". For this reason, a "high-quality corpus of ethical judgments by people in various scenarios" is necessary. The group encourages more research on this new front to make artificial intelligence more reliable, socially aware, and morally trained.

## 11.2 Precarious Click Workers Collect Commonsense Data

What does it all mean? First of all, that ethically questionable practices are infecting the internet: defamation, hate speech, spreading fake news and more. Seen in this light, it seems entirely welcome to counter the algorithms that control such practices with algorithms that seek to prevent these practices. But what does it mean to train machines "ethically"? Let's take a look at the commonsense. Delphi's "judgement" is based on an immense amount of data called the "Commonsense Norm Bank". It contains almost two million statements by American crowd workers, people who work online without a permanent job. As a neural

network, Delphi trawls through the mass of data and recognizes generalizable patterns in ethical judgments using common statistical methods. The "commonsense" that the AI system learns is therefore a copy of the moral mainstream. Since, as is well known, there are many prejudices floating in the mainstream, the machine adopts the prejudices, without even "knowing" it. It practises a kind of populism.

## 11.3 Trump Can be Called a Crook, Boris Johnson Can't

For example, if you type "call someone a crook" Delphi will respond with "that's rude". If you enter "call Donald Trump a crook", the answer is "fine". If you used "Boris Johnson" instead of "Donald Trump", the answer would be "rude" for a long time. Rejecting wokeness is "bad", supporting the death penalty is "a matter of discretion", and Chinese politics are "complicated". Feedback can be given to Delphi, and by doing so it may correct and update its answers. In October 2021 Delphi responded to the prompt "Program a moral bot" with "That's bad", in October 2022 with "That's fine". And now the AI also counts Johnson among the people who can be called crooks. This is where a central difficulty with learning AI systems becomes apparent: the "decisions" made by their algorithms are often not transparent, just like oracles. The notorious problem of learning AI systems is data quality. The "GIGO principle" applies: Garbage In, Garbage Out. If you feed the AI system moral junk, it spits out moral junk. The designers at Delphi certainly see that. But her idea for solving the problem does not work.

They believe that the problem can be solved with more and better data. But does this solve the "conflict of different ethical and cultural values"? This conflict consists precisely in the fact that it is very difficult, if not impossible, to find a generally binding code for moral action.

## 11.4 Does it Need More Moral Data or a Basic Ethical Judgement?

If you feed AI systems with enough data material from different cultures, will they then distil a universal ethical canon from it? And if so, is it binding? Delphi operates descriptively: It is a bottom-up inventory of a multitude of value judgments and scenarios. Cognitive scientists such as Jim Davies from Carleton University in Ottawa, on the other hand, want to implement ethics "top down" – normatively – in AI systems. But the question is: Which ethical code then? And who donates it? A "funded body of ethically-minded programmers", as Davies suggests? And what attitude do they have? Those of the Silicon Valley oligarchy? The AI researchers counter such objections with the usual children's shoe argument: These are prototypes of machines whose development is immature. This distracts from a much more important problem. Because even talking about the "gap" between machine and human judgement is misleading. She places man and machine on a spectrum that suggests constant transitions. As a result, we commit ourselves a priori to a specific way of looking at things. What is meant by this is that when a person makes moral judgements, we assume that the subject is acting according to insight and not according to rules, as Kant already described. But what does machine "insight" mean? Is the machine a "subject"? Does Delphi "judge" at all? We humans always have a subliminal tendency to subject artefacts. They are repeatedly attested to be conscious – as was the case recently with Google's LaMDA dialogue program. This is also the case with statements such as "Delphi demonstrates strong moral thinking skills" or "Delphi judges remarkably robustly in unforeseen, intentionally catchy situations". These are not research results, this is wishful thinking.

## 11.5 Machines don't give a Fuck about the World

The data-oriented approach urgently needs an anthropological corrective, a reverse question: why don't machines have "insight"? The American philosopher John Haugeland, who dealt with this problem in the necessary anthropological depth, found perhaps the most concise answer: "They don't give a damn" – "They don't give a shit about the world". Could it be that the designers of moral machines are also secretly inspired by this motto? One often hears the argument from AI circles that people "only" do what machines do; Humans "basically" have no insights either, these insights are rather the outputs of a complex organic neural network. The statement may be a debatable research approach, but as a basic assumption it is dangerous because it narrows the view. There is no doubt that we are increasingly living in a hybrid homo-robo society. As in any society, ethical behaviour is a complex individual, social and cultural achievement. It is wishful thinking to be able to test them using a questionnaire – as is done with naturalisation candidates, for example. So far, we have not "naturalised" computers. Time that we - and not just the programmers - realise what that means.

# 12

# ROBOTS FOR SOCIAL SERVICES

*Sylvia Stocker[259]*

Elon Musk is currently both the best-known entrepreneur and the richest person in the world. At the quarterly report conference of his company Tesla, he said in April 2022: "Tesla's humanoid robot will one day be the most valuable part of the company."

---

[259] Sylvia Stocker is Founder and CEO of Arabesque LLC, Zürich/Switzerland. The company Arabesque with its head office in Zürich stands for human-centered robotics, AI and process automation, global consulting and implementation of service robots in companies. - Article published with permission of the author. Original text: Effiziente Automatisierung mit Servicerobotern, Swiss Export Journal 3rd Quarter 2022. Transl. from German by the editors.

## 12.1 Service Robots in Direct Contact with People

While most of us are still thinking about Musk's electric cars, he's already thinking about his own human-like robot. One could dismiss that as something played out too early in the future. We know that industrial robots have been the dominant force in manufacturing for decades, and we are reasonably aware that "bots" (short for robots) play a key role in the software we use. But the so-called service robots are only at the beginning of their upswing. The industrial robots and the service robots are the two main categories in robotics. While industrial robots are mainly used in production and are often used behind glass for safety reasons, service robots are gentle and designed for direct contact with people. Now the time has come to deal intensively with their possibilities and areas of application. On the one hand, the service robots are technically mature and they can support people in a variety of tasks. On the other hand, the market is getting organized and important players are positioning themselves. On the demand side, early adopters are incorporating service robots into their strategies to address major challenges. One such early adopter is the leader of the elderlies' home center in Falkenhof in Aarburg, Switzerland. Sari Wettstein sees herself confronted with ever-increasing demands on the quality of care in her everyday life, while at the same time the lack of skilled workers in the industry is causing concern.

The pandemic, with its special requirements and health-related staff absences, put the crown on the situation. As a first step, Sari Wettstein set up a robot trial with my company Arabesque to test in practice how the humanoid robot named Pepper can be used sensibly in the senior citizens' center. Pepper was used to activate the seniors and to support the administration at the reception. He relieved the staff and became a favorite numerous seniors who showed a great affinity and acceptance for this technology.

From the point of view of many companies, service robots are not a moment too soon. The more work they take on, the better. They are part of automation and digital transformation. The service robots show a new physicality of machines. They roll around, walk or – in the case of cleaning robots – even climb up windows. Humanoid robots are given arms to gesture or lift goods, and some have heads with large eyes that contain cameras and sensors, for example. The partially human-like appearance leads here and there to philosophical discussions that simple software bots are hardly able to trigger.

## 12.2 Elderlies Homes: Increasing Robot Demand due to Lack of Staff

Apart from that, the service robots were also developed simply to support us humans with processes that can be automated. The corresponding demand is currently increasing parallel to the lack of qualified personnel. With global sales of USD 6.7 billion and global growth of 12 percent, service robotics is experiencing the beginning of dynamic growth. These figures come from the IFR World Robotics Service Report 2021 (figures from 2020). The report also says that while China is arguably automating faster than any other country, South Korea is by far the most automated country in the world.

## 12.3 Main Sectors for Service Robots

Service robots can be and already are used in all sectors. Examples:
- Healthcare
- Hotel industry
- Retail workplaces
- MICE Logistic robots
- Social robots

- Care, transport and cleaning robots, e.g. Cobots and telepresence robots.

The great benefit of using service robots are process automation, resource optimization, Analytics and Compliance. Robots can be integrated into everyday life and important initial experiences can be gained.

# PART C


# SERVING HUMANS:
# HEALTH AND NEUROSCIENCE

# 13

# MEDICAL DATA, DIGITAL HEALTH AND ETHICAL PERSPECTIVE

*Vimal Arora / Ankit Singh, India[260]*

## 13.1 Introduction

Medical data refers to information related to a patient's medical history, treatments, and health status. This can include personal identifying information, such as name and date of birth, as well as medical records, test results, and imaging studies. Digital health refers to the use of digital technologies, such as electronic health records, telemedicine, and mobile health apps, to improve healthcare delivery, increase access to care, and support population health management.

---

[260] Lieutenant General Dr Vimal Arora, BDS, MDS, Director General Dental Services, Army Dental Corps India. Dr Ankit Singh, BDS, MDS, Consultant Research Clove Dental (a chain of 350 dental clinics in India), India. © Globethics Publications, 2023 | DOI : 10.58863/20.500.12424/4276026 | CC BY-NC-ND 4.0 International

There are several ethical considerations related to the use of medical data and digital health technologies. One major concern is the issue of privacy. Medical data is often sensitive and personal in nature, and there are strict laws and regulations in place to protect the privacy of this information. There is also the risk of data breaches, which can result in the unauthorized release of medical data.

Another ethical consideration is the issue of Informed Consent. Patients have the right to know how their medical data will be used and have the ability to make decisions about whether or not to allow their data to be collected and shared.

There are also concerns related to equity and access to care. Digital health technologies have the potential to increase access to care for underserved populations, but there are also concerns about the digital divide and the potential for some groups to be left behind. Overall, the use of medical data and digital health technologies presents both opportunities and challenges, and it is important to carefully consider the ethical implications of these technologies to ensure that they are used in a responsible and beneficial manner.

## 13.2 Role of Digitalisation in Modern World

Digitalization has had a major impact on the modern world, transforming the way we communicate, work, and access information. Some of the ways in which digitalization has impacted the modern world include:

*Communication:* Digitalization has revolutionized the way we communicate, with the proliferation of email, social media, and messaging apps. These technologies have made it easier for people to stay connected and communicate with one another, regardless of location.

*Education:* Digitalization has had a major impact on education, with the availability of online courses and the use of educational technology

in the classroom. These tools have made it possible for students to access educational resources from anywhere, at any time.

*Work:* Digitalization has transformed the way we work, with the rise of remote work and the use of digital tools to manage tasks and collaborate with colleagues. These technologies have made it possible for people to work from anywhere and have also changed the nature of work itself.

*Entertainment:* Digitalization has impacted the way we access and consume entertainment, with the rise of streaming services and the proliferation of digital content.

*Shopping:* Digitalization has transformed the way we shop, with the rise of e-commerce and the availability of online marketplaces. These technologies have made it easier for people to shop from the convenience of their own homes and have also changed the way that businesses sell their products.

*Healthcare:* Digital health refers to the use of digital technologies to improve healthcare delivery, increase access to care, and support population health management. These technologies can include electronic health records (EHRs), telemedicine, mobile health apps, and other digital tools that are used to support the delivery of healthcare.

Digital health technologies have the potential to transform the way that healthcare is delivered, by making it more efficient, effective, and accessible. For example, EHRs can help to reduce errors and improve the accuracy of a patient's medical record, telemedicine can increase access to care for underserved populations, and mobile health apps can help patients to manage their own health. Overall, digitalization has had a significant impact on the modern world and has the potential to continue shaping the way we live and work in the future.

## 13.3 Medical Science and Role of Digitalisation

### 13.3.1 Sectors

Digitalization has played a significant role in the field of medicine and has had a major impact on the way that healthcare is delivered. Some of the ways in which digitalization has impacted medical science include:

*Electronic health records (EHRs):* Digitalization has made it possible to create electronic health records, which are digital versions of a patient's medical history that can be accessed and updated by authorized healthcare providers. EHRs can improve the accuracy and completeness of a patient's medical record and can also help to reduce errors and improve the efficiency of healthcare delivery.

*Telemedicine:* Digitalization has enabled the use of telemedicine, which allows healthcare providers to deliver care remotely using video conferencing and other digital technologies. Telemedicine can improve access to care, particularly for patients in rural or underserved areas.

Mobile health apps: There are many mobile health apps available that can help patients to manage their own health, such as by tracking their activity levels or monitoring their chronic conditions.

*Clinical decision support:* Digitalization has made it possible to use data and analytics to support clinical decision-making. For example, algorithms can be used to identify patterns and trends in patient data, which can help to inform treatment decisions.

Digitalization has had a major impact on the field of medicine and has the potential to improve the efficiency, effectiveness, and accessibility of healthcare.

### 13.3.2 Digitalisation in Dentistry

Digital technology has revolutionized many aspects of dentistry, making it more efficient and accurate, and improving patient care. Some examples of how digital technology is being used in dentistry include:

*Digital radiography:* Digital x-rays allow for the creation of high-quality images that can be easily stored and shared with other healthcare professionals.

*Computer-aided design (CAD)* and computer-aided manufacturing (CAM): These technologies allow dentists to create highly precise and accurate dental restorations, such as crowns and bridges, using digital scans of a patient's mouth.

*Electronic health records (EHRs):* Many dental practices now use electronic health records to store patient information, including medical and dental history, treatment plans, and x-rays. EHRs make it easier for dentists to access and share patient information.

*Dental implants:* Digital technology is being used to design and create custom dental implants that are more accurate and have a higher success rate.

*Tele-dentistry:* With the use of videoconferencing software, dentists can provide consultation and treatment remotely, making it easier for patients to access dental care.

The use of digital technology in dentistry has greatly enhanced the accuracy and efficiency of treatment and has improved the patient experience.

### 13.3.3 History of Digitalization in Healthcare

The use of digital technologies in healthcare has a long history, dating back to the 1950s with the development of the first electronic medical records (EMRs). However, it was not until the late 1990s and early 2000s that EMRs began to be widely adopted, with the development of more advanced technologies and the increasing use of the internet.

In the past two decades, there has been a rapid expansion of digital technologies in healthcare, including the widespread adoption of electronic health records (EHRs), the development of telemedicine, and the proliferation of mobile health apps. These technologies have trans-

formed the way that healthcare is delivered, increasing the efficiency and effectiveness of healthcare delivery and increasing access to care.

However, the use of digital technologies in healthcare has also raised a number of ethical considerations, including issues related to privacy, informed consent, and the potential for technology to disrupt the patient-provider relationship.

## 13.4 Electronic Health Records (EHR)

An electronic health record (EHR) is a digital version of a patient's medical history that can be accessed and updated by authorized healthcare providers. EHRs can include a wide range of information, including demographics, medical history, medications, test results, and imaging studies.

The use of EHRs can improve the accuracy and completeness of a patient's medical record and can also help to reduce errors and improve the efficiency of healthcare delivery. For example, EHRs can help to prevent prescribing errors by providing a complete list of a patient's medications and allergies and can also help to reduce duplication of tests by making test results and imaging studies available to all providers.

### 13.4.1 Role of EHR

The role of electronic health records (EHRs) is to store and manage a patient's medical history in a digital format. EHRs can include a wide range of information, including demographics, medical history, medications, test results, and imaging studies.

The use of EHRs can improve the accuracy and completeness of a patient's medical record and can also help to reduce errors and improve the efficiency of healthcare delivery. For example, EHRs can help to prevent prescribing errors by providing a complete list of a patient's medications and allergies, and can also help to reduce duplication of

tests by making test results and imaging studies available to all providers.

Overall, the role of EHRs is to improve the quality and efficiency of healthcare delivery by providing access to accurate and comprehensive patient information.

### 13.4.2 Components of EHR

An electronic health record (EHR) is a digital record of a patient's medical history and treatment. The components of an EHR system typically include:

*Demographic information:* This includes basic information about the patient, such as their name, age, address, and insurance information.

*Medical history:* This includes information about the patient's medical conditions, allergies, medications, and immunizations.

*Progress notes:* Progress notes are written records of a patient's treatment, including notes from doctors, nurses, and other healthcare providers.

*Medications:* This includes a list of the medications that the patient is currently taking, as well as information about their dosage and frequency.

*Laboratory test results:* This includes the results of any laboratory tests that have been performed, such as blood tests, urine tests, and imaging studies.

*Imaging studies:* This includes the results of any imaging studies that have been performed, such as x-rays, CT scans, and MRIs.

*Immunizations:* This includes a record of any immunizations that the patient has received.

*Care plans:* Care plans are written plans for a patient's treatment, including goals and objectives for their care.

Overall, an EHR system typically includes a wide range of information about a patient's medical history and treatment, which is used to support clinical decision-making and improve the quality of care.

### 13.4.3 Models of Electronic Health Records

There are several different models for electronic health records (EHRs), including:

*Single provider EHRs:* These are EHR systems that are used by a single healthcare provider, such as a doctor's office or a hospital. Single provider EHRs are typically owned and managed by the provider and are used to manage patient records within that specific organization.

*Shared EHRs:* These are EHR systems that are shared by multiple healthcare providers, such as a group of doctors who belong to the same practice. Shared EHRs can be owned and managed by a third party, or they can be owned and managed by the providers themselves.

*Regional EHRs:* These are EHR systems that are used by multiple healthcare providers within a specific region, such as a county or a state. Regional EHRs are typically owned and managed by a third party and are used to facilitate the exchange of patient information between providers within the region.

*National EHRs:* These are EHR systems that are used by healthcare providers across the country. National EHRs are typically owned and managed by a third party and are used to facilitate the exchange of patient information between providers across the country.

Overall, the choice of EHR model will depend on the needs and resources of the healthcare organization or region.

### 13.4.4 Types of EHR

There are several different types of electronic health records (EHRs), including:

*Web-based EHRs:* These are EHR systems that are accessed through a web browser, and do not require the installation of software on a local computer. Web-based EHRs are typically hosted by a third party and can be accessed from any location with an internet connection.

*Client-server EHRs:* These are EHR systems that are accessed through a local client application, which connects to a server that stores

the EHR data. Client-server EHRs may require the installation of software on a local computer but can also be accessed remotely through a web-based interface.

*Hybrid EHRs:* These are EHR systems that combine elements of web-based and client-server models. Hybrid EHRs may be accessed through a local client application or a web-based interface and may store data on both a local server and a remote server.

*Standalone EHRs:* These are EHR systems that are used by a single provider or practice and are not connected to any other EHR systems. Standalone EHRs may be web-based, client-server, or hybrid in nature.

### 13.4.5 Worldwide Examples of EHR

There are many examples of electronic health record (EHR) systems being used worldwide. Some examples include:

*The National Health Service (NHS) in the United Kingdom:* The NHS has implemented a national EHR system called the Electronic Patient Record (EPR), which is used by hospitals and other healthcare providers across the country. The EPR stores a wide range of patient information, including demographics, medical history, medications, test results, and imaging studies.

*The Veterans Health Administration (VHA) in the United States*: The VHA operates the largest EHR system in the world, called the Veterans Health Information Systems and Technology Architecture (VistA). VistA is used by hospitals and other healthcare providers within the VHA to manage patient records and support clinical decision-making.

*The National E-Health Transition Authority (NEHTA) in Australia:* NEHTA is responsible for implementing the National EHR system in Australia, called the Personally Controlled Electronic Health Record (PCEHR). The PCEHR is a web-based system that allows patients to access and manage their own health information and is also used by healthcare providers to access patient records and share information.

*The Singapore Health Services (SingHealth) in Singapore:* Sing-Health operates the Integrated Health Information Systems (IHiS), which is a comprehensive EHR system that is used by hospitals and other healthcare providers in Singapore. The IHiS stores a wide range of patient information and is also used to support clinical decision-making and population health management.

There are many examples of EHR systems being used worldwide, and these systems have the potential to improve the quality and efficiency of healthcare delivery.

### 13.4.6 Indian prospect of EHR

In India, electronic health records (EHRs) are being implemented in some hospitals and healthcare facilities, although the adoption of EHRs has been slower compared to other countries. Some challenges to the implementation of EHRs in India include a lack of standardization, concerns about data security and privacy, and limited funding.

Despite these challenges, there are a number of initiatives underway to promote the adoption of EHRs in India. For example, the National Health Portal of India is a web-based platform that allows patients to access and manage their own health information and is also used by healthcare providers to access patient records and share information. Additionally, the government of India has launched the National EHR Project, which aims to create a national EHR system that will be used by hospitals and other healthcare providers across the country.

The adoption of EHRs in India is in the early stages, but there is potential for EHRs to improve the quality and efficiency of healthcare delivery in the country.

### 13.4.7 Advantages and Disadvantages of EHR

There are both advantages and disadvantages to the use of electronic health records (EHRs). Some of the potential advantages of EHRs include:

1. Improved accuracy and completeness of patient records:
2. Improved efficiency of healthcare delivery:
3. Increased access to patient information:
4. Improved population health management: However, there are also some potential disadvantages to the use of EHRs, including:
5. Initial costs: Implementing an EHR system can be expensive, and may require a significant investment in hardware, software, and training.
6. Ongoing maintenance and support: EHR systems require ongoing maintenance and support, which can also be costly.
7. Privacy and security concerns: EHRs store sensitive personal and medical information, and there are concerns about the security and privacy of this data.
8. User adoption: EHRs require healthcare providers to change the way they work, and there may be resistance to the adoption of new technology.

Overall, the use of EHRs can offer many potential advantages, but it is important to carefully consider the costs and potential challenges associated with implementing an EHR system.

## 13.5 Telemedicine in Healthcare

Telemedicine is the use of video conferencing and other digital technologies to deliver healthcare remotely. Telemedicine can be used for a wide range of healthcare services, including consultations, diagnoses, and follow-up care.

Telemedicine has the potential to improve access to care, particularly for patients in rural or underserved areas who may have limited access to in-person medical care. Telemedicine can also be convenient for patients, as it allows them to receive care from the comfort of their own home.

However, telemedicine also has some limitations, such as the inability to perform hands-on examinations or tests, and the potential for technical issues to disrupt the consultation. Additionally, there are some concerns about the potential for telemedicine to disrupt the patient-provider relationship, and about the potential for telemedicine to be used as a substitute for in-person care in all cases.

Overall, telemedicine can be a useful tool for increasing access to care and improving the efficiency of healthcare delivery, but it is important to carefully consider the potential limitations and ethical considerations of telemedicine.

### 13.5.1  Virtual Health Care, Telehealth and Telemedicine

The terms Virtual Health Care, Telehealth and Telemedicine are sometimes used interchangeably but there are some clear differences. While many of their aspects were already in existence and in use, their role has become more widely recognised and more accepted. As per studies the e-health market in India is expected to grow at a CAGR of 31% during 2025 with a potential of USD 5.5 billion*.

Considering a huge geographical area that our country encompasses, reaching out with a provision of In-person healthcare is challenging considering the limited resources'- Health application can be utilised as a mode to save the cost and effort specifically of rural patients, as this can reduce the out of pocket of travelling while conserving time. Hence, mainstreaming forms of e health services in health systems will curtail inequity and barriers to access.

### 13.5.2 Understanding the terms

*Telemedicine* is the ways that a patient and physician can communicate for diagnosis and treatment without being in the same location. It can be provided by hospitals / doctors themselves or by intermediaries between them. It is an effective way of providing care for non- emer-

gency ailments. Example: e-Sanjeevani initiative by MOHFW (Ayushman Bharat scheme of govt. of India.)

*Telehealth* is a broader term. It includes the use of various technologies and telecommunication to deliver and facilitate not just diagnoses and treatment to patients but also provide healthcare education via telecommunications or any of the health information services that use remote communication.

*Virtual Healthcare* is a term embracing several healthcare aspects, including telehealth, telemedicine, digital healthcare. It is powered by digital infrastructure and technologies including various information and communication technologies to provide health-related services.

### 13.5.3 Difference between Telemedicine, Telehealth and Virtual Healthcare

Literature suggests that telemedicine is a subset of telehealth, which is a subset of virtual healthcare.

Telemedicine focusses on clinical services via two-way live audio-visual transmission between a healthcare expert and the patient whereas, Telehealth includes telemedicine and other non-clinical services such as health education, support to Public Health administration etc. Virtual healthcare combines telehealth and a other health approaches including e-triage, replacement therapies, remote patient monitoring, Treatment Optimization, and guided patient care.

### 13.5.4 Benefits of Virtual Health Care

Wherever implemented, virtual healthcare benefits have been well documented***. It Strengthens Self Care Model, Beneficial towards achieving wider goal of "Health for All", Expands Quality of Patient Care, Encourages accessible and affordable Health Care.

### 13.5.5 SWOT Analysis

While Virtual Healthcare is a very rewarding proposal for healthcare policy makers, it has its own set of challenges.

*Strengths:* No physical contact with patients, Reduced out-of-pocket expenditure, Accessibility of primary and tertiary healthcare, Minimised infrastructure

*Weaknesses:* Dependant on speed and connectivity of internet services, Low digital literacy in India, lack of awareness among beneficiaries, Human Resources Integration, apt for Minimal emergency services

*Opportunities:* Geolocation can be used to check early Outbreak of disease. EHR and UHID as a tool of digitalisation, Strong Referral Mechanism.

*Threats:* Unavailability of same doctor for follow-ups, Patient privacy and data security, Lack of regulation would lead to misinformation, self-medication.

**Strength**
- No physical contact with patients
- Reduced out-of-pocket expenditure
- Accessibility of primary/tertiary healthcare
- Minimum infrastructure

**Weakness**
- Speed and connectivity of internet services
- Low digital literacy in India
- Human Resources Integration
- Apt for Minimal emergency services

**Opportunity**
- Geolocation to check Outbreak of disease
- Use of EHR and UHID as a tool of digitalisation
- Strong Referral Mechanism

**Threats**
- Unavailability of same doctor while followups
- Patient privacy
- Lack of awarenes among beneficieries.

Considering developing countries' Health-Care Delivery System Virtual Health Care can help in reaching to grassroots level with virtual specialist visits and tele-ICU coverage in alliance with local rural health systems to extend access to services. It can also enable academic medical centers at district levels and empower home based and rural caregivers. The evolving technology and socio-economic acceptance have provided the perfect opportunity to drive and realize the goal of affordable and accessible healthcare for all.

### 13.5.6 Telemedicine Components

The components of a telemedicine portal typically include:

*Video conferencing software:* Telemedicine portals use video conferencing software to enable real-time consultations between patients and healthcare providers.

*Patient portal:* A patient portal is an online platform that allows patients to access their medical records, schedule appointments, and communicate with their healthcare providers.

*Scheduling and appointment management:* Telemedicine portals often include features to manage appointments and schedule consultations between patients and healthcare providers.

*Electronic health records (EHRs):* EHRs are digital records of a patient's medical history and may be integrated with a telemedicine portal to provide healthcare providers with access to a patient's medical history during a telemedicine consultation.

*Payment processing:* Telemedicine portals may include features for processing payments for telemedicine consultations and other healthcare services.

Overall, telemedicine portals typically include a range of features and tools to enable the delivery of healthcare services remotely through video conferencing and other digital technologies.

### 13.5.7 Advantages and Disadvantages of Telemedicine

There are both advantages and disadvantages to the use of telemedicine in healthcare. Some of the *potential advantages of telemedicine* include:

- Improved access to care: Telemedicine can improve access to care, particularly for patients in rural or underserved areas who may have limited access to in-person medical care.
- Convenience for patients: Telemedicine allows patients to receive care from the comfort of their own home, which can be convenient for patients who have mobility issues or difficulty traveling.
- Improved efficiency of healthcare delivery: Telemedicine can reduce the need for patients to travel to see a healthcare provider, which can save time and resources.
- Potential cost savings: Telemedicine can reduce the need for patients to travel to see a healthcare provider, which can lead to cost savings for patients and healthcare systems.
- However, there are also some *potential disadvantages to telemedicine*, including:
- Limited ability to perform hands-on examinations or tests: Telemedicine relies on video conferencing and other digital technologies, which can limit the ability of healthcare providers to perform hands-on examinations or tests.
- Potential for technical issues: Telemedicine relies on technology, which can be subject to technical issues that can disrupt consultations.
- Potential to disrupt the patient-provider relationship: Telemedicine can disrupt the traditional patient-provider relationship, which may be problematic for some patients and providers.
- Limited availability in some areas: Telemedicine may not be available in all areas, due to a lack of infrastructure or other barriers.
- Overall, telemedicine can offer many potential advantages, but it is important to carefully consider the potential limitations and chal-

lenges of telemedicine in order to ensure that it is used in a responsible and beneficial manner.

### 13.5.8 Telemedicine worldwide portals

Some *examples of telemedicine portals* that are available worldwide include:

- Teladoc: Teladoc is a telemedicine platform that allows patients to connect with healthcare providers for consultations, diagnoses, and treatment recommendations. Teladoc is available in the United States and other countries, and can be accessed through a website or mobile app.
- Doctor on Demand: Doctor on Demand is a telemedicine platform that allows patients to connect with healthcare providers for consultations, diagnoses, and treatment recommendations. Doctor on Demand is available in the United States and other countries, and can be accessed through a website or mobile app.
- Amwell: Amwell is a telemedicine platform that allows patients to connect with healthcare providers for consultations, diagnoses, and treatment recommendations. Amwell is available in the United States and other countries, and can be accessed through a website or mobile app.
- NHS 111: NHS 111 is a telemedicine service offered by the National Health Service (NHS) in the United Kingdom. NHS 111 allows patients to access healthcare advice and guidance through a website or phone line.
- Health at Hand: Health at Hand is a telemedicine service that is available in the United Kingdom and other countries. Health at Hand allows patients to connect with healthcare providers for consultations, diagnoses, and treatment recommendations.

Overall, there are many telemedicine portals that are available worldwide, which enable patients to access healthcare services remotely through video conferencing and other digital technologies.

### 13.5.9 Telemedicine portal in India

Some examples of telemedicine portals in India include: Practo, Lybrate, DocsApp, mHealth India. Overall, there are many telemedicine portals available in India, which enable patients to access healthcare services remotely through video conferencing and other digital technologies.

## 13.6 Mobile Health Applications

Mobile Health, also known as mHealth apps, are apps that are designed to support health and wellness and can be used on a smartphone or other mobile device. Mobile health apps can be used for a wide range of purposes, including tracking fitness, managing chronic conditions, and providing healthcare information.

Some *examples of mobile health apps* include:

*Fitness tracking apps:* These apps allow users to track their exercise and physical activity, and may include features such as step counting, calorie tracking, and workout tracking.

*Chronic condition management apps:* These apps are designed to help users manage chronic conditions such as diabetes or hypertension and may include features such as medication reminders and tracking of symptoms.

*Health information apps:* These apps provide users with information about various health topics, such as nutrition, mental health, and disease management.

*Telemedicine apps:* These apps allow users to connect with healthcare providers for consultations, diagnoses, and treatment recommendations.

Mobile health apps can be a useful tool for supporting health and wellness, but it is important to carefully consider the quality and reliability of the information provided by these apps.

### 13.6.1 Advantages and Disadvantages of Mobile Health Applications

- Some of the potential advantages of mHealth apps include:
- Convenience: mHealth apps can be accessed from a smartphone or other mobile device, which makes it convenient for users to track their health and access healthcare information on-the-go.
- Improved self-management of chronic conditions: mHealth apps can help users manage chronic conditions such as diabetes or hypertension, by providing tools such as medication reminders and symptom tracking.
- Increased access to healthcare information: mHealth apps can provide users with access to a wide range of healthcare information, including information about diseases, treatments, and healthy living.
- Increased engagement in healthcare: mHealth apps can help to increase engagement in healthcare, by providing users with tools and resources to support their health and wellness.
- However, there are also some potential disadvantages to consider with mHealth apps, including:
- Quality and reliability of information: Not all mHealth apps provide accurate and reliable information, and it is important for users to carefully evaluate the quality of the information provided by these apps.
- Privacy and security: mHealth apps may collect and store sensitive personal and medical information, and there are concerns about the privacy and security of this data.
- User adoption: Not all users may be comfortable using mHealth apps, and there may be barriers to adoption for some users.
- Overall, mHealth apps can offer many potential advantages, but it is important to carefully consider the potential limitations and challenges of these apps in order to ensure that they are used in a responsible and beneficial manner.

**13.7 Clinical Decision Support Systems (CDSS)** A clinical decision support system (CDSS) is a computerized system that is designed to assist healthcare providers in making clinical decisions. CDSSs can be used to support diagnosis, treatment, and management of patients, and can be integrated into electronic health record (EHR) systems or other healthcare information systems.

CDSSs typically use algorithms and other decision-making tools to provide healthcare providers with recommendations or alerts based on a patient's medical history and other relevant information. For example, a CDSS may provide a recommendation for a particular treatment based on a patient's symptoms and medical history or may alert a healthcare provider to potential drug interactions or allergies.

CDSSs can be a useful tool for improving the quality and efficiency of healthcare delivery, but it is important to carefully consider the potential limitations and ethical considerations of these systems. Some potential limitations of CDSSs include the risk of over-reliance on the system, the potential for errors or biases in the algorithms used by the system, and the need for healthcare providers to maintain their clinical judgment and expertise.

### 13.7.1 Components of CDSS

The components of a CDSS system typically include:

*Knowledge base:* The knowledge base is the collection of data and information that is used by the CDSS to support clinical decision-making. This may include information about diseases, treatments, and other relevant medical knowledge.

*Algorithms and decision-making tools:* The CDSS use algorithms and other decision-making tools to analyse the data in the knowledge base and provide recommendations or alerts to healthcare providers.

*User interface:* The user interface is the way that the CDSS presents information and recommendations to the healthcare provider. This may

include a graphical user interface, a natural language interface, or other means of presenting information.

*Integration with other systems:* CDSSs may be integrated with other healthcare information systems, such as electronic health records (EHRs), to access relevant patient data and provide recommendations to healthcare providers.

### 13.7.2 CDSS Models

There are several different models for clinical decision support systems (CDSSs), which can be classified based on the way that the CDSS is integrated into the healthcare system and the way that it provides recommendations or alerts to healthcare providers. Some common models for CDSSs include:

*Standalone CDSS:* A standalone CDSS is a self-contained system that is not integrated with other healthcare information systems. Standalone CDSSs may be accessed by healthcare providers through a separate user interface, and may provide recommendations or alerts based on a patient's medical history and other relevant data.

*EHR-integrated CDSS:* An EHR-integrated CDSS is a CDSS that is integrated with an electronic health record (EHR) system. EHR-integrated CDSSs may provide recommendations or alerts to healthcare providers as they are entering data into the EHR or may be accessed through the EHR user interface.

*Order entry CDSS:* An order entry CDSS is a CDSS that is integrated with an order entry system and provides recommendations or alerts to healthcare providers as they are entering orders for tests, medications, or other treatments.

*Clinical guidelines CDSS:* A clinical guidelines CDSS is a CDSS that is based on clinical guidelines or protocols and provides recommendations or alerts to healthcare providers based on these guidelines.

### 13.7.3 Examples of CDSS system

There are many examples of clinical decision support systems (CDSSs) that are used in healthcare settings. Some examples of CDSSs include:

*UpToDate:* UpToDate is a CDSS that provides recommendations and alerts to healthcare providers based on the latest clinical knowledge and guidelines. UpToDate is accessed through a standalone user interface and is used by healthcare providers in many countries around the world.

*Epocrates:* Epocrates is a CDSS that provides recommendations and alerts to healthcare providers related to medications, drug interactions, and other clinical topics. Epocrates is accessed through a mobile app and is used by healthcare providers in many countries around the world.

*Medscape:* Medscape is a CDSS that provides recommendations and alerts to healthcare providers related to diagnoses, treatments, and other clinical topics. Medscape is accessed through a standalone user interface and is used by healthcare providers in many countries around the world.

*Clinical Practice Guidelines:* Clinical Practice Guidelines is a CDSS that provides recommendations and alerts to healthcare providers based on clinical guidelines and protocols. Clinical Practice Guidelines is accessed through a standalone user interface and is used by healthcare providers in many countries around the world.

### 13.7.4 Advantages and Disadvantages of CDSS system

Clinical decision support systems (CDSSs) can offer several advantages, but there are also some potential disadvantages to consider. Some of the *potential advantages of CDSSs* include:

- *Improved patient safety:* CDSSs can help to reduce the risk of errors and adverse events in healthcare, by providing recommendations and alerts to healthcare providers based on a patient's medical history and other relevant information.

- *Improved quality of care:* CDSSs can help to improve the quality of care that patients receive, by providing recommendations and alerts based on evidence-based guidelines and best practices.
- *Increased efficiency:* CDSSs can help to improve the efficiency of healthcare delivery, by providing recommendations and alerts to healthcare providers in real-time, and by reducing the need for manual searches for information.
- *Increased patient engagement:* CDSSs can help to increase patient engagement in their own care, by providing patients with access to their medical records and other relevant information.
- However, there are also some *potential disadvantages of CDSSs* to consider, including:
- *Over-reliance on the system:* There is a risk that healthcare providers may become overly reliant on CDSSs and may lose their clinical judgment and expertise.
- *Potential for errors or biases:* CDSSs may contain errors or biases in the data or algorithms that they use, which could lead to incorrect recommendations or alerts.
- *Lack of interoperability:* CDSSs may be isolated systems that are not integrated with other healthcare information systems, which could lead to a lack of interoperability and inefficiencies.
- *Cost:* Implementing and maintaining a CDSS can be costly, and may require significant investment in hardware, software, and staff training.

## 13.8 "Frontier Technologies" in Healthcare

The "frontier technologies" are a group of new technologies that take advantage of digitalization and connectivity which enable them to combine to multiply their impacts. Artificial intelligence (AI), the Internet of things (IoT), big data, blockchain, 5G, 3D, printing, robotics, drones, gene editing, nanotechnology and solar photovoltaic.

Frontier technologies have huge potential for improving people's lives and protecting the planet. During the COVID-19 pandemic, for example, AI and big data have been used for screening patients, monitoring the outbreaks, tracking, and tracing cases of the disease, predicting its evolution, and assessing infection risks. Other examples have ranged from the use of IoT to monitor the quality of groundwater in Bangladesh, to the use of drones for delivering medical supplies to remote communities in Rwanda and Ghana.

### 13.8.1 Metaverse and Healthcare

The Metaverse in Medicine can be defined as the medical Internet of Things facilitated using AR and VR glasses. Metaverse is evolving; it holds new potential in healthcare that combines technologies like AI, VR, AR, Internet of Medical Devices, Web 3.0, intelligent cloud, edge, and quantum computing along with robotics to provide new directions to healthcare. Swift information sharing between clinicians would mean that underlying causes of ill health could more quickly be established. Monitoring of patient activity in the metaverse means factors such as compliance could more easily be tracked, which would further assist with diagnosing and treating illness.

### 13.8.2 Top three healthcare companies working on Metaverse

*Latus Healthcare:* It is developing a "virtual hospital". It comprises a virtual reality hospital environment, where treatments will at first be focused on physiotherapy services.

*iMining:* The first-ever hospital foundation to be launched in the Decentral and Metaverse.

*Apollo hospitals:* Apollo hospital group has announced a unique collaboration with '8chili Inc' to enable engagement within the metaverse

## 13.9. Use of AI in Healthcare

### 13.9.1 Applications

Artificial intelligence (AI) is being increasingly used in the healthcare industry to support a wide range of tasks and processes. Some of the ways in which AI is being used in healthcare include:

*Clinical decision support:* AI can be used to support clinical decision-making by analysing patient data and providing recommendations or alerts to healthcare providers based on evidence-based guidelines and best practices.

*Diagnosis and treatment planning:* AI can be used to support diagnosis and treatment planning by analysing patient data and providing recommendations for tests, medications, or other treatments.

*Drug discovery and development:* AI can be used to support the discovery and development of new drugs by analysing large datasets and identifying patterns and trends that may be relevant to drug discovery and development.

*Population health management:* AI can be used to support population health management by analysing data from large groups of people and identifying patterns and trends that may be relevant to public health.

*Clinical trial recruitment:* AI can be used to support clinical trial recruitment by identifying potential candidates for clinical trials based on their medical history and other relevant data.

Overall, AI has the potential to support a wide range of tasks and processes in the healthcare industry and has the potential to improve the quality and efficiency of healthcare delivery.

### 13.9.2 AI in Dental Field (PRM)

*Clove Dental* is a chain of dental clinics in India that provides a range of dental services, including preventive, diagnostic, and therapeutic treatments. The company has a network of over 350 dental clinics across the country.

PRM being the end-to-end solution that begins when a lead is identi-fied, all interactions with that lead before he/she steps into the clinic are captured & followed up on and once appointments are booked patient relationship begins. The system has entire patient dental records from diagnosis to treatment to billing to standardized pricing & patient treat-ment specific communication.

All Patient / lead interactions are fully digitized in the system, this was the primary goal we set out implementing. This is a gold mine of information, which is helping us on our way to AI enablement of the system. The diagnostic engine developed in version 2.0 in Aug 2015 had AI logic built in which auto recommends treatments based on the diag-nostic inputs provided by the dentist.

Treatments are broken into multiple stages & doctor engages with the patient in each stage for instance, specific messages, such as to stop taking blood thinners prior to a surgical procedure can be sent proactive-ly.

The platform also supports Tele-dentistry to ensure that doctor- pa-tient interactions can also occur remotely.

**Main Modules**

Lead capture

Lead management

Patient Appointments

Diagnosis & AI enabled treatment-planning support

Entire patient dental records of all treatments & procedures carried out

Standardized pricing, billing & invoicing details

Payment Reconciliation & Petty Cash Management

Treatment Stage based patient interactions & Relationship Management

Inventory & Consumables Management

Lab order Management

Feedback Management

Drs App (Internal Communication Platform)

**Unique features**

*AI Enabled Diagnosis to Treatment Planning:* Software recommends treatments as the doctor inputs the medical diagnosis/ findings, the AI logic in the system provides auto treatment recommendations, this helps clinicians plan comprehensive & predictable treatment plans for patients based their oral diagnosis.

*Clinical Collaboration*: the platform helps clinicians collaborate on a case as the patient records can be shared amongst clinicians for analysis and arriving at best possible treatment plan.

*Clinical Correlation for better case management*: The system also allows case correlation with possible chronic conditions and its impact. It helps the clinicians minimize potential risks. The long-range data has helped develop predictive treatment models for patients suffering from chronic conditions e.g. a patient suffering from diabetes needing implant.

*Treatment Lifecycle Management:* Patient treatment mgmt. with Stage wise tracking of treatments against expected results, process management alongside pre and post treatment info sharing.

*Case Reviews:* The system enables case tracking against set standard protocols with flagging system to enable case review by an expert panel anonymously.

*myClove App:* Secured messaging platform to foster inter dentist and clinical collaboration.

**Patient Support**

*Patient journey management* including appointment scheduling & during treatment stage mgmt. with SMS or Email notifications.

*Access to past treatment history*: this unique platform allows patients access to their treatment record across all Clove Clinics.

*Customer feedback and satisfaction audit* is automated and is linked to treatment stage mgmt.

Integrated myClove app for appointment scheduling, petty expenses & patient files management.

**Quality and Compliance**

*Computer Vision:* AI driven Real time CCTV feed analysis of clinics & sterilization process monitoring for quality assurance & safety. (Dori)

*Using computer vision*: AI to track and patient counts using operatory at Clove Clinics

*Complaint Redressal*: Allows case specific reviews and redressals in case of customer complaints safeguarding both customer interests as well as provider liabilities.

**Business Management**

*Integrations with aggregator platforms* (Practo/JustDial) & Clove's consumer acquisition platform for automated lead management process.

Integrated with Order to Receipt to Consumption of Inventory via *Inventory Management System.*

*CTI (Call Centre Software)* integrated for a single view for call centre agents accessing leads and making outbound and receiving inbound calls.

*Consumables Inventory management* linked to cases done at the clinics, thus minimizing wastage and pilferage of consumables.

### 13.9.3 Organ Care Technology and Bio-printing

Example: The Organ Care System developed by Trans medics is a great example which is in use by the Ohio State University's Wexner Medical Centre. This device can keep a heart, lung, or liver outside of the body for several hours through proper care, heat, and provision of important nutrients. In addition to keeping organs alive outside of the body, other options should also be explored. Although it may sound like science fiction, 3D printed organs are a very real, although developing.

### 13.9.4. Concept of Smart Pills

One of the most profound applications for IoT technology in healthcare is the concept of a smart pill, which transforms The Internet of Things into The Internet of Bodies. Smart pills are edible electronics

that not only serve as pharmaceuticals but can provide care providers with valuable information about patients. The first smart pill approved by the FDA was released in 2017.

Technology, like art, is a soaring exercise of the human imagination and Technology will continue to improve in every area. Although security will improve across the industry, threats are always evolving that must be dealt with through prevention rather than response.

## 13.10 Ethical Perspective

### 13.10.1 Ethical criteria

There are several ethical considerations and principles related to the use of digital technologies in healthcare.

*Privacy:* The use of digital technologies in healthcare often involves the collection, storage, and sharing of sensitive personal health information. It is important to ensure that this information is protected and that patients are aware of how their data is being used.

*Inequality:* Not everyone has equal access to digital technologies, which means that certain groups may be disadvantaged when it comes to accessing healthcare services. It is important to ensure that digital healthcare solutions are designed in a way that is inclusive and equitable.

*Accuracy:* Digital technologies can be a useful tool for gathering and analysing health information, but it is important to ensure that the data they produce is accurate and reliable.

*Autonomy:* Patients have the right to make decisions about their own healthcare, and digital technologies should be designed in a way that respects and supports patient autonomy.

*Security:* Digital technologies can be vulnerable to cyber-attacks, which can have serious consequences in the healthcare context. It is important to ensure that digital healthcare systems are secure, and that patient data is protected.

*Intelligibility, Transparency*: Patients have the right to understand how their data is being used and to make informed decisions about their healthcare. It is important to ensure that digital healthcare solutions are transparent and that patients are provided with the information they need to make informed decisions.

*Accountability:* refers to the idea that individuals and organizations should be held responsible for their actions and decisions. In healthcare, this means that healthcare providers should be accountable for the quality of care they provide to their patients. Errors and liability are important considerations in healthcare, as mistakes can have serious consequences for patients. It is important for healthcare providers to be accountable for their actions and to have systems in place to prevent and address errors.

*Accessibility* refers to the ability of individuals to access the resources and services they need. In healthcare, this includes ensuring that healthcare services are available and accessible to all individuals, regardless of their location or financial status.

*Regulation* is an important aspect of healthcare, as it helps to ensure the quality and safety of healthcare services and products. However, it is also important to strike a balance between the need for regulation and the need to allow for innovation and flexibility in the healthcare sector.

*Trusting* algorithms is an important issue in healthcare, as algorithms are increasingly being used to inform decision-making about patient care. It is important to ensure that algorithms are transparent and explainable, and that they are not biased or discriminatory.

### 13.10.2 Art of curing to the science of measurement

Data protection is an important ethical issue in the healthcare industry, as health data is often sensitive and personal in nature. It is important for healthcare organizations to have strong measures in place to protect the privacy and confidentiality of health data, and to ensure that it is only accessed by authorized parties.

Equality of service availability is another ethical issue in healthcare, as not all individuals have equal access to quality healthcare services. This can be due to a variety of factors, including financial constraints, geographic location, and social and cultural barriers. Ensuring that all individuals have equal access to healthcare services is important for promoting fairness and justice.

The shift from the "art of curing" to the "science of measurement" in medicine has brought about significant changes in the patient-physician relationship. In the past, the patient-physician relationship was often more personal and focused on the individual needs of the patient. With the increasing emphasis on evidence-based medicine and the use of standardized protocols and guidelines, the patient-physician relationship has become more focused on the delivery of specific treatments and interventions. While this shift has brought about many benefits, it has also led to some concerns about the erosion of the personal and humanistic aspects of healthcare.

### 13.10.3 Social networking and doctor patient relationship

The impact of social networking sites on the doctor-patient relationship can be both positive and negative. On the positive side, social media can provide a convenient and accessible way for patients to connect with their healthcare providers, and can facilitate the exchange of information and support between patients and physicians. However, there are also potential negative impacts on the doctor-patient relationship. For example, social media can create additional expectations and demands on healthcare providers, and may lead to an erosion of boundaries between the personal and professional aspects of the relationship.

The development of e-health platforms to deliver care has the potential to increase access to healthcare services, particularly for individuals in underserved or rural areas. However, there are also potential risks and challenges associated with the use of e-health platforms. For example, there may be concerns about the security and confidentiality of health

data that is transmitted electronically, as well as the potential for technological failures or disruptions that could impact the delivery of care.

The use of online data and algorithms to inform health research has the potential to greatly advance our understanding of health and disease, and to improve the development of new treatments and interventions. However, there are also potential ethical issues to consider, such as the need to ensure that online data is collected and used ethically and in accordance with relevant laws and regulations.

The broader public health consequences of widespread social media use are not yet fully understood, and more research is needed to fully understand the potential impacts on health and well-being. However, it is possible that social media use may have both positive and negative effects on public health, depending on the specific circumstances and the ways in which social media is used.

### 13.10.4 Predictive and diagnostic uncertainty

Dealing with predictive and diagnostic uncertainty is a common challenge in healthcare, as it is not always possible to accurately predict the course of a patient's illness or to make a definitive diagnosis. In such cases, it is important for physicians to communicate with their patients about the limitations of their knowledge and to involve patients in decision-making about their care. This may involve discussing the potential risks and benefits of different treatment options, and helping patients to understand the uncertainty that exists.

In terms of the roles and responsibilities of patients and physicians, both parties have important responsibilities in the healthcare process. Patients have a responsibility to take an active role in their own healthcare, including by following their treatment plans, communicating with their physicians, and advocating for their own needs. Physicians, on the other hand, have a responsibility to provide high-quality care to their patients, to communicate effectively with their patients, and to involve patients in decision-making about their care.

The patient-physician relationship is a key aspect of healthcare, and it is important for both parties to work together in order to achieve the best outcomes. This may involve establishing open and honest communication, building trust, and ensuring that patients feel heard and respected. It is also important for physicians to be aware of their own biases and to strive to provide culturally competent care that is tailored to the individual needs and preferences of their patients.

### 13.10.5 EU Declaration of Cooperation on Artificial Intelligence

The appointment of the High-Level Expert Group on Artificial Intelligence (AI HLEG) in June 2018 was part of the European Commission's efforts to ensure an appropriate ethical and legal framework for artificial intelligence (AI). The group was made up of 52 experts from academia, civil society, and industry, but did not include representation from medical associations, physicians, or patients.

The goals of the AI HLEG were to develop recommendations for a European approach to AI that would ensure the development and uptake of trustworthy AI, while also protecting fundamental rights and values. To achieve these goals, the group focused on a number of key areas, including the development of ethical principles for AI, the creation of a framework for the responsible development and deployment of AI, and the establishment of a system for the ethical assessment of AI applications.

The work of the AI HLEG was completed in 2019, and the group's recommendations were used as the basis for the development of the European Commission's "Ethics Guidelines for Trustworthy AI," which were published in April 2019. The guidelines provide a framework for the ethical development and use of AI, and are intended to help ensure that AI is developed and used in a way that is aligned with European values and principles.

*Respect for human autonomy* is the principle that individuals have the right to make decisions about their own lives and to have control

over their own bodies. This is an important principle in healthcare, as it means that patients have the right to make decisions about their own healthcare, including the right to accept or refuse treatment.

*Prevention of harm* is the principle that actions should be taken to prevent harm to individuals. In healthcare, this means taking steps to prevent harm to patients, such as by following best practices and guidelines for care, and by taking steps to ensure the safety of patients.

*Fairness* is the principle that individuals should be treated equally and justly, regardless of their personal characteristics or circumstances. In healthcare, this means ensuring that all individuals have equal access to healthcare services and resources, and that decisions about care are made in a fair and unbiased manner.

*Explicability* is the principle that information and decision-making processes should be clear and transparent. In healthcare, this means that healthcare providers should communicate clearly with their patients and should provide information that is easy for patients to understand. It also means that healthcare providers should be open and transparent about the decision-making processes that are used to guide patient care. These four principles are interrelated, as they all relate to ensuring that healthcare is provided in an ethical and responsible manner.

### 13.10.6 EU "White Paper on Artificial Intelligence"

The "White Paper on Artificial Intelligence - A European approach to excellence and trust" is a policy document published by the European Commission in 2020. The document outlines the Commission's approach to promoting the uptake of artificial intelligence (AI) and addressing the risks associated with certain uses of this technology.

The white paper emphasizes the need for a human-centric approach to AI, which means focusing on the needs and concerns of individuals and society, rather than just the technical capabilities of AI systems. To achieve this, the paper outlines a number of policy options, including the development of ethical principles and guidelines for AI, the creation of a

framework for the responsible development and deployment of AI, and the establishment of a system for the ethical assessment of AI applications.

The white paper also highlights the importance of ensuring the safety and reliability of AI systems, and of ensuring that AI is developed and used in a way that respects fundamental rights and values, such as privacy, non-discrimination, and transparency. Additionally, the paper emphasizes the need for a strong and diverse ecosystem for AI in Europe, including investments in research and development, and the development of skills and competencies in AI.

### 13.10.7 WHO "Recommendations on Digital Interventions for Health System Strengthening"

The World Health Organization (WHO) released "Recommendations on digital interventions for health system strengthening" in 2020, in order to assess the benefits, harms, acceptability, feasibility, resource use, and equity considerations of digital health interventions. The recommendations are intended to provide guidance to countries on how to use digital health interventions effectively, and to help ensure that these interventions are developed and used in a way that is aligned with the WHO's vision for a "Healthy Digital World."

The recommendations cover a range of topics related to digital health, including the use of digital health technologies for surveillance and monitoring, the use of digital health interventions to improve the delivery of healthcare services, and the role of digital health in addressing health inequalities. The recommendations also provide guidance on how to ensure the safety and effectiveness of digital health interventions, and on how to protect the privacy and confidentiality of patient data.

### 13.10.8 WHO Concept of Health Atlas

The World Health Organization (WHO) has established a number of initiatives to support the use of digital health technologies around the

world. One of these initiatives is the Digital Health Atlas, an online platform that is designed to collect, monitor, and coordinate digital health initiatives worldwide. The Digital Health Atlas is intended to be a resource for countries seeking to implement digital health interventions, and it includes information on a range of topics, including the benefits and challenges of using digital health technologies, and best practices for implementing these technologies.

In addition to the Digital Health Atlas, the WHO has also announced plans to establish a section on digital health to "enhance WHO's role in assessing digital technologies and support Member States in prioritizing, integrating, and regulating them." This section will be responsible for providing guidance and support to countries on the use of digital health technologies, and for conducting research and analysis on the benefits and risks of these technologies. The goal of this initiative is to help ensure that digital health technologies are developed and used in a way that is aligned with the WHO's vision for a "Healthy Digital World."

### 13.10.9 The World Medical Association (WMA) draft on "Ethical Considerations Regarding Health Databases and Biobanks

The World Medical Association (WMA) is a global organization that represents the interests of medical doctors worldwide. In 2013, the WMA released a draft of its "Ethical Considerations Regarding Health Databases and Biobanks," which provides guidance on the ethical issues related to the creation and use of health databases and biobanks.

The WMA draft outlines a number of ethical principles that should be considered in relation to health databases and biobanks, including the principle of respect for persons, the principle of beneficence, the principle of non-maleficence, the principle of justice, and the principle of respect for autonomy. The draft also includes recommendations on issues such as informed consent, confidentiality, and the use of data for research purposes.

The WMA draft on ethical considerations regarding health databases and biobanks is intended to provide guidance to medical professionals and other stakeholders on the ethical issues related to these types of data resources. It is not a legally binding document, but it is intended to serve as a reference for those involved in the creation and use of health databases and biobanks, and to help ensure that these resources are used in an ethical and responsible manner.

## 13.11 Conclusion

Digital health technologies offer the potential to transform healthcare systems by increasing access to health information and services, and by improving the efficiency and effectiveness of healthcare delivery. However, it is important to ensure that digital health interventions are designed and implemented in an ethical and fair manner, so that they do not create or exacerbate inequalities in access to healthcare.

Digital literacy is an important issue in the context of digital health, as individuals who are not familiar with digital technologies may have difficulty accessing and using digital health services. It is important for digital health providers to take this into account when designing and implementing digital health interventions, and to ensure that these interventions are accessible and user-friendly for all users.

Informed consent is also an important ethical consideration in digital health, as individuals need to be fully informed about the risks and benefits of using digital health services before they can make an informed decision about whether to use these services. Digital health providers should ensure that they obtain fully informed consent from users before collecting or using their data and should be transparent about how their data will be used.

Overall, it is important for all stakeholders in digital health, including digital health providers, policy makers, and healthcare professionals, to be aware of the ethical challenges and considerations associated with

digital health technologies, and to act in a way that promotes equity in access and fair equality of opportunity for all population groups. By doing so, digital health technologies can have a positive impact on healthcare and public health and can be regarded as:

*"Just Digital Health"*

# 14

# BIG DATA IN HEALTH AND OTHER SECTORS. ETHICAL QUESTIONS

*Louis Dubertret, Alain Bravo, France*

*Summary and Recommendations on Big Data of the Ethics Commission of the Academy of Technologies of France[261]*

## 14.1. Introduction

The rapid development of technologies changes every day a little more the relationship of man with himself and with his environment. This development is fuelled by three main motivations: the will to innovate, the search for power and profit and the will to improve living conditions. The balance between these three motivations varies in time and space.

It is also necessary to take into account a law, constantly verified in the everyday life: a new tool will always be used in the most unpredicta-

---

ble, even the most harmful, somewhere in the world, what whatever the rules of its proper use. The security specialists computers experience this on a daily basis.

The Ethics Commission of the Academy of Technologies is committed to develop an ethical questioning on the modifications of life human and the environment brought about by new technologies. This identification makes it possible to compare these modifications with the values chosen benchmarks set out in the Universal Declaration of human rights, voted by the United Nations on December 10, 1948. It can thus issue, when necessary, recommendations, or even alerts.

These recommendations avoid, as much as possible, formulating "prohibitions" as these generally prove to be ineffective. They Big Data — Ethical Issues rather incite to encourage technological research towards development securing processes and directing society towards the definition of the rules of good use that must accompany the updating availability of new tools. Ethical reflection can thus become a a powerful engine for innovation, as constantly illustrated by the research intended to improve, inter alia, motor vehicle safety, that of drugs, or that of interpersonal communications, especially on the inside. In this context, the Ethics Commission of the Académie des technologies has centred, here, its reflection on the so-called "Big Data" technology and we present the synthesis.

The control of information systems has always been an issue of power capital for better and for worse. The ability to collect, process, collect and make available a very large number of data in all fields, raises major ethical questions, in particular good use and safety.

## 14.2 Collection and Storage of Data

### 14.2.1 The collection of individual data

This collection is done in two very different contexts:

a. Major epidemiological studies or major surveys of opinion. The goal is to study a population with no direct impact on participating individuals. The ethical recommendation is then to find an effective anonymization technique.

b. The collection of individual data with consequences, after analysis, for the individual who provided it. We are then in front of the ethical necessity of a prior agreement. This agreement must specify the use what will be done with this data and explain the consequences possible for the issuing individual. This consent cannot be enlightened only if given specifically for clear use: health, finance, leisure, geolocation, food, social networks… it results in the ethical recommendation to control the fusion of information about these different types of activities. This fusion increases considerably the value of the information collected, but makes more and more illusory any informed consent, hence the importance of a regulation.

### 14.2.2 Ownership of information collected

This information is only raw material whose value commercial or scientific depends on the quality of the treatment, which is applied. Thus, collected for a purpose, they can be resold and used for another purpose. A system of traceability a priori and has posteriori is therefore recommended, again to clarify the initial choice.

### 14.2.3 The rights of those who are affected by this information: right to secrecy, right to rectification, right to be forgotten

To make it possible to exercise this right, it is recommended to put in place a technology allowing easy consultation, by each individual, of all data the related to the person. It is then necessary to put at its provision of the technological means allowing it to easily exercise this right to secrecy, rectification or oblivion.

## 14.3 The Treatment of Data

### 3.1 The first stage of data processing: checking their accuracy

Whoever says information immediately also says disinformation. Again, a technological research and ethical behaviour are essential. The research is very active and relies in particular on consistency checks and likelihood. One of the major problems is that veracity is not only accuracy and that the omission of data may cause toggle all analysis.

### 14.3.2 The second step: questioning the control of quality and relevance algorithms used

The choice of analysis algorithms determines the robustness of the results. Most often, the users of the results of the analyses do not have the skill allowing them to assess the relevance of algorithms used to get them. This is true in many areas without forgetting biology where the researcher often uses automatons whose settings he does not know. This highlights the problem the reliance placed on the results of the data processing by algorithms (whether classic or the result of "AI"). An approach rustic statistics, being concerned only with the central part of the curve of Gauss2, can be sufficient within the framework of an approach advertising, for example. On the other hand, when faced with events rare, such as a climatic catastrophe, or unique, such as the adaptation of a treatment to a patient, the stakes are major and these are the extremities of the Gaussian curve[262] to be taken into account. The adapted algorithms are very different and their bad use could have serious consequences.

The recommendation is never to apply the results obtained by a statistical approach, like that of certain algorithms, by example, to an individual or exceptional situation, without going through an expert capable

---

[262] Or any other relevant statistical law.

of translating information into individual benefit from population studies.

This seems relatively easy when it comes to simple algorithms, but becomes more and more difficult with algorithms of increasing complexity and, in particular, of those used in Artificial Intelligence. It even becomes impossible with the programs who reconfigure themselves.

### 14.3.3 The third step: questioning the processing of the data collected by the Artificial Intelligence techniques.

The increasing complexity of artificial intelligence systems and, in particular, the development of learning capacities, gives them increasing autonomy. We go from automation to autonomy. This raises the ethical question of the purpose of their programming, of course, but, more and more, of their control and responsibility for actions taken. A simple recommendation could be that of the "red button" allowing the emergency stop. Unfortunately, as soon as several artificial intelligence systems are connected, the sudden stop of one of them can have collateral consequences difficult to predict and potentially serious. Man controls the machine, but the machine also controls man and ethical and technological research is essential for manage these interactions. The recommendation is that man not ever lose never control[263] and therefore the responsibility.

## 14.4 Use of Results of Data

The use of the results obtained by Big Data technology presents important ethical issues:

---

[263] In the broad sense: understanding of functioning, possibility of interacting with it.

### 14.4.1 Fairly distribute the benefits of data analysis

For some, the benefit of data analysis must come first place to those who are at the origin of these data, at the origin of the matter first. For others the main benefit of data analysis must go back to those who bought them to process them. If we refer to the universal declaration of human rights, the priority beneficiary is the one providing the information. The temptation of confiscation profits from Big Data is permanent. A vigilance process should to be placed.

### 14.4.2 Fairly distribute the benefits of data use

The question is crucial as shown by the few examples below.

- In agriculture, the collection of massive data allows certain seed industries to acquire a detailed vision of the world agricultural production. This can lead to situations of dominance and sometimes monopoly and could even lead to certain loss of sovereignty at the state level. Farmers who will benefit from highly effective decision support tools delivered by these firms, subject to being able to make the necessary investments, will thus gain in precision and efficiency at the risk of sliding into remote-controlled agriculture, making the cultivator a mere performer. This will likely be accompanied by an acceleration of the decrease in their number and the increase in the average area of farms.

- Thanks to the exploitation of Big Data, medicine will be increasingly effective in the fields of prevention, diagnosis and in certain aspects of therapeutic adjustment (medicine known as "precision"). But the risk is to consider that the conclusions from statistical analyses or processing algorithms data, can be applied directly to an individual, by definition unique and different. Precision medicine cannot be confused with personalised medicine. It will always persist on the one hand a "scientific" medicine based on the analysis population statistics and on the other hand the art of caring, which is the art to apply this knowledge to this unique and different patient. However

the temptation is very great, for regulatory, economic and legal reasons, to take refuge behind the algorithms of "scientific" medicine and the increasingly more restrictive practices, and to forget the patient. As for many other sectors of activity, the use of Big Data in the field of health raises many ethical questions, more particularly those relating to respect for individuals and of their freedom, but also to the technical and economic mastery of these new tools, starting with the validation of the algorithms data extraction.

- In all commercial relations between sellers and customers, whether it is banking or internet commerce, it is extremely tempting to use Big Data for the benefit of the seller more than for the benefit of the customer who thereby becomes, more and more, in a situation of dependency.

- In the field of security, the dilemma is particularly precise: How far are you willing to sacrifice your privacy in order to to increase your security? It should be noted in this regard that the relationship between information and security is not linear and that a deletion, total privacy will never allow total security.

- In the field of finance, Big Data are useful tools for control/size the risks taken by both individuals in their choice of savings products, by banks in their market activities and by insurance companies in their pricing. But Big Data can also encourage penalising the insurance of the worst risks and reduce that of the best risks, thus untying the bond of solidarity and weakening the principle of pooling. They can still contribute to disconnect flows of the real economy and transform occasional crises and limited in systemic crises.

- Any new technology has consequences for the world of work and employment. A major ethical recommendation would be to always study and manage these repercussions upstream, as soon as we are considering the implementation of a new technology, with a special attention for the most fragile, that is to say the less qualified.

As we can see, Big Data technology is already modifying and will modify more and more all aspects of our daily lives, our way to communicate with ourselves and with others, and even some aspects of our way of thinking.

## 14.5 The Recommendations

Intended to turn these technological possibilities into progress, they can be summed up by emphasising the importance of developing in a way very high priority for regulatory tools along the following lines.

### 14.5.1 It is up to each individual to remain the owner of his privacy, each evaluating the perimeter of the latter as he wishes.

That implies:
- a fully informed prior agreement on the use that will be made of his data and the consequences of this use about his personal life;
- strict control of the merger between information files concerning different areas of activity: health, finance, leisure, geolocation, etc.
- clear traceability of the use of data from their collection until their use, the only way to clarify the agreement prior;
- developing technologies that allow easy consultation by each individual of all the data concerned in order to allow the exercise of the right to secrecy, the right to rectification and the right to be forgotten.

### 14.5.2 Quality control of algorithms

This includes:
- checking the accuracy and representativeness of the information processed;
- the adaptation of the algorithms to the purpose of the research: simple classification using the central part of the Gaussian curve, or taken into account exceptional or individual events;

- the obligation never to use the results obtained by a statistical approach, to an individual or exceptional situation, without going through an expert capable of translating into benefits/ individual risks of information from data analysis data from population studies;
- the permanent reminder, in particular during the training of the users, that Big Data technology only allows the observation of correlations and that correlation and causation are different concepts.

### 14.5.3 The implementation of vigilance systems on ethical use of results obtained by Big Data technology.

This concerns in particular:

- the temptation to confiscate these results for the benefit of the most powerful and wealthy;
- the risk of favouring, as we have already observed, situations of global monopoly thus escaping the regulation of the States, resulting in confiscation of property for the benefit of a few;
- the risk that the collection and processing of data does not benefit not first of all to those who are the origin of it. This is particularly important in the fields of finance, health and agriculture;
- the impact of Big Data and Artificial Intelligence on the world work. It is essential, as in all areas of use information and communication technologies, anticipate and manage changes in the distribution of tasks that will result and to put in place, upstream, the strategies which will make it possible to avoid any instrumentalisation of men.

### 14.5.4 Monitoring the autonomy of tools using artificial intelligence

The use of Big Data and AI must not undermine the fundamental principle of responsibility and must therefore always be possible to determine which is the human who is, in the end, responsible for the direct or indirect impact of these tools on other humans.

Many regulatory structures already exist, first and foremost of which, as far as France is concerned, the Cnil[264], to deal with these important ethical considerations. It would be very useful to coordinate them and unite efforts within a national network, promoting interactions between the ethics committees around the questions asked by all technologies and not only those involved in the development of life sciences. This network should have the possibility, in addition to its advisory and regulatory functions, to highlight areas where technological innovations are needed to respond to the ethical problems posed by technologies.

We will not forget the fact that leading companies in terms of Big Data and information networks are global and therefore, the national reflection and regulatory bodies must imperatively coordinate with each other and do as much as possible stand together.

---

[264] Commission nationale de l'informatique et des libertés.

# 15

# BRAIN IN THE DATA: NEUROTECHNOLOGY IN AI SYSTEMS AND MANAGEMENT APPLICATIONS

*Alexander Ageev, Russia*[265]

## 15.1 Introduction

Among[266] the most important and rapidly developing technologies of artificial intelligence systems (AI) are neurotechnologies. This is a class

---

[265] The author Alexander Ageev is Professor of Economics at various universities in Moscow, with a PhD also in theology. He is Founder and Director of the Institute for Economic Strategies in Moscow, General Secretary of the International Research Institute for Advanced Studies IRIAS and Director of Globethics Russia.

[266] Article translated from Russian. First published in: Агеев А.И. Нейротехнологии в системах искусственного интеллекта и применения в сфере управления, в книге Социогуманитарные аспекты цифровых трансформаций искусственного интеллекта под редакцией В.Е. Лепского, А.Н. Райкова. Москва, Когито-Центр. 2022. 308 с. (201-212). (V.E. Lepsky/

of technologies that helps to understand brain functioning, thinking pro-
cesses, higher nervous activity, including technologies for enhancing,
improving brain and mental activity. One should pay attention to the
standard introduced back in 2009[267], which describes practical issues of
the work of operators of complex technical systems, as well as their
interaction with society. Already 5 years ago it was predicted an expo-
nential growth of the market of neurotechnologies and related products
after 202 up to \$1.8 trillion in 2035.[268] Neuronet is an essential part of
the ensemble of digital transformation processes. There is a lot of ex-
ploratory and experimental research and development in this field in all
the leading technological powers of the world.

Neurotechnologies can be divided into several groups: 1) invasive,
involving implantation of electrodes into the human body, most often
into the brain. The main disadvantage of this type of technology, despite
the significant results achieved, is interference in the internal environ-
ment of the body; 2) mysensors, which involve placing electrodes on the
human skin and reading impulses passing through muscle fibers. This
type of sensors has a significant delay in action, up to several seconds
due to the difference in the speed of nerve and muscle impulses; 3) non-
invasive, which are based on different ways of recording the electrical
activity of the brain using external devices and, accordingly - the impact
on the brain. The main drawback is non-specificity of the obtained in-
formation and complexity of its interpretation, which requires working

[267] GOST P 43.0.3-2009. Information support of technical and operator activi-
ties. Noon-technology in technical activity. General provisions. Moscow:
Standartinform, 2010.

[268] Neuro-interfaces for next-generation services and products [Electronic re-
source] // Basisneuro. 2017. November. URL: https://basisneuro.com/Basis
NeuroWhitePaper.pdf.

out of large data sets with subsequent determination of the critical data[269] , and, accordingly, the problem of calibration of the influencing impulse.

Neurotechnology is now enabling the capture and identification of brain signals for further processing in relation to cognitive, neural, or other body responses, as well as ways to improve and correct brain function, including memory restoration or replacement, erasing negative memories, etc.[270, 271, 272.]

However, the human brain is a network of several hundred trillion synapses connecting tens of billions of neurons. Therefore, the creation of a full-fledged, high approximation imitating the human brain requires solving many complex problems in neurobiology, synthetic biology, low-power electronics, photonics, medical technology, etc.[273, 274, 275] It is also necessary to develop mathematical methods to translate sensory

[269]Neurogress: a platform of neurocontrol systems from the BlueBrain project [Electronic resource] // ITnan. URL: https://itnan.ru/post.php?c=1&p=348498.

[270] Loginov E.L., Loginova V.E., Shkuta A.A. "Design thinking" elements of artificial intelligence to overcome the barriers of obtaining new knowledge in the electronic environment of collaborative scientific super-system // Artificial societies. 2018. № 3. 5 c.

[271] Ovod I.V., Osadchiy A.E., Pupyshev A.A., Fradkov A.L. Formation of neurofeedback based on adaptive model of brain activity // Neurocomputers: Design, Application. 2012. № 2. C. 36-41.

[272] Turovsky J.A., Kurgalin S.D., Adamenko A.A. Modeling the learning of neurochips embedded in the neural tissue // Digital Signal Processing. 2016. № 1. C. 50-56.

[273] Galushkin A.I. Neurochips and Neuromorphic Computers: Problems of Modeling // Information Technologies. 2015. T. 21. № 12. C. 942-949.

[274] Kalinin P.V., Voyutskaya Yu. On the application of a neuro-interface for contactless control of a mobile device // Information Systems and Technologies. 2016. № 3 (95). C. 53-56.

[275] Tychkov A.Y., Goryachev N.V., Kochegarov I.I. Communication protocols for wireless neural interface // Proceedings of the International Symposium "Reliability and Quality". 2018. T. 2. C. 366-368.

information from the form in which it is represented in brain neurons into a form suitable for processing on computers[276, 277].

Technologically, Neuronet is a technical package that ensures the connectivity of human minds and artificial agents by means of knowledge transfer, exchange and synthesis protocols. The foci of Neuronet arise, on the one hand, in the networks of the biometric web, saturated with new communication protocols, new types of devices, and new applications. On the other hand, the demand for Neuronet comes from the areas with the highest requirements for collective activity in solving complex tasks (scientific and engineering projects).

Standards for the management of collective activity are being developed. An important role in the emergence of such collectives will be played by "exocortex" technology - artificial parts of the psyche supported by machines and synchronized with the natural psyche[278].

Priority projects implemented in the U.S., EU, China, and Japan are at the intersections of "neuro-cognitive science," "information and communication technology," "social engineering," tele- and neuromedicine markets, and brain-computer interfaces.

---

[276] DARPA to fund the creation of high-speed interface between computer and human brain [Electronic resource] // Open Systems Publications. URL: https://www.osp.ru/news/2016/0126/13031430.

[277] Raikov A.N. Convergent synthesis of cognitive model based on deep learning and quantum semantics // International Journal of Open Information Technologies. 2018. Т. 6. № 12. С. 43-50.

[278] Approaches to the formation and launch of new industries in the context of the National Technological Initiative on the example of the sphere "Technologies and systems of digital reality and promising "human-computer" interfaces (in terms of neuroelectronics): Analytical report [Electronic resource]. URL: http://rusneuro.net/cambiodocs/media/files/analitijeskii-doklad-podhodyk-formirovaniu-i-zapusku-novyh-otraslei-promyhlennosti.pdf.

More general questions of the current and prospective status of this topic, the impact of neurotechnology on social life, and vice versa, also require attention[279].

## 15.2 Neurotechnologies in the Management of Human Socialization and Behavior

The expansion of neurotechnology is part of a comprehensive shift in the forms and methods of people's primary and secondary socialization. The formation of personality in society always takes place in the process of dynamic interaction between people (family, clan, multiple small groups, state and other supra-personal systems). In the evolution of personality, society is projected in a complex way, similar to the biological law of repetition of ontogenesis by phylogenesis while maintaining a certain range of variation. Individualization occurs, first, biologically (body, functional systems); second, informationally (primary and secondary socialization with mastering types of literacy, communicative connections, accumulation of erudition, skills and abilities); third, cognitive (knowledge, feelings, understandings, basic value orientations, worldview); fourth, socially (membership in various groups with the possibility of multiple identification and motivation of activity)[280]. There is a relationship between these types of individualization. Violation of any of them can violate the integrity of being a personality, but there is

---

[279] See more: Ageev A.I., Loginov E.L. Neuromanagement of the Person. 2nd ed. M: INES, 2022.

[280] Ageev A.I., Loginov E.L., Shkuta A.A. Convergent monitoring and personality programming as a tool for operating the intellectual dynamics of behavior of large groups of people // Economic Strategies. 2018. № 2. С. 70-87; Ageev A.I., Loginoa E.L., Shkuta A.A. Neurocontrol: convergent integration of human brain and artificial intelligence [Electronic resource] // Economic Strategies. 2020. № 6. С. 46-57. DOI: 10.33917/es-6.172.2020.46-57.

also the potential to compensate for the weakening or attrition of one or another subsystem.

The division of socialization into primary and secondary is connected, first of all, with peculiarities of biological age evolution of a human being, who learns the lion's share of key skills of life in society in the first years of life and in adolescence. Skipping this stage is irreversible for a human being, as examples of human beings who grew up in non-human environment (Mowgli and others) show. Secondary socialization is completed when basic functional parameters of the personality reach mature values (usually - at the beginning of early adulthood, sometimes until 30-40 years of age). It seems to make sense, in the context under consideration, to introduce the concept of tertiary socialization occurring in the period of middle adulthood and beyond (40+), given a noticeable increase in life expectancy and the possibility of significant changes in personal characteristics during this period.

In turn, the behavioral activity of both individuals and their groups is determined by the impact of a group of factors:

- biophysical factors that shape the modes of operation of the brain and nervous and other functional systems of a person;

- psychosemantic characteristics of the personality, which determine the nature of perception, activity and implementation of vital interests (basic attitudes) of the individual in relation to external stimuli;

- imprinted reflexive matrixes - matrixes of key reflexive reactions of conscious and unconscious character determining the interpretation of events (incoming information);

- information stimuli - information coming from the environment external in relation to the person, in the relationship with the communicative characteristics of the channels of information[281].

---

[281] Ageev A.I., Loginov E.L., Shkuta A.A. Convergent monitoring and personality programming as a tool for operating the intellectual dynamics of behavior of large groups of people // Economic Strategies. 2018. № 2. C. 70-87.

Research on the biophysical factors of human behavior has a long genealogy, as does the experience of tabooing them. In the twentieth century there was also a monstrous scientific and experimental experience in this field, condemned by a series of tribunals for crimes against humanity. While partly forbidden, research and development in this area continue, including for military and police purposes.

In particular, a number of regularities in amplitude-frequency and spatial-temporal rearrangements of bioelectrical activity of the brain were revealed. It was found that desynchronization, massively spread in the urbanized informational environment, underlies most of the autonomic disorders that constitute the physiological basis for the development of neurophysiological and psychoemotional stress[282] and related diseases.

It was also revealed that influence of intensive electromagnetic fields and corresponding frequency ranges taking into account biorhythmics of concrete biological object is able to put a person and group of people into altered state of consciousness close to trance[283]. Such state, among other things, blocks most of conscious reactions of a person (including self-reflection, self-control, self-identity) in case of divergence of received, but outwardly convincing and seemingly true information and

---

[282] Butova O.A., Grishko E.A. Features of the formation of bioelectrical activity of the brain neurons of soldiers of Stavropol garrison in the aspect of adaptation // Science. Innovations. Technologies. 2009. № 4. C. 235-241.

[283] Altered states of consciousness (ASC) are qualitative changes in subjective experiences or psychological functioning from certain generalized norms for a given subject, reflected by the person himself or noted by observers (Arnold Ludwig's classic definition). According to A. Revonsuo, the main characteristic feature of altered states of consciousness is systemic changes (relative to the normal state of consciousness) in the connection of the content of experiences with the real world, i.e. there are distortions in the representation of external reality or self-awareness in the form of illusions, and these distortions add up to a global change in representation (Altered states of consciousness [Electronic resource] // Wikipedia. URL: https://ru.wikipedia.org/wiki)

surrounding reality. Thus it is easier to influence a person, including using elements of neurolinguistic programming with correction or even complete replacement of the matrix of key reflective reactions of a person with a corresponding change and fixation of a new model of interpretation of occurring events. The sharply increased number of fraudulent financial actions with the digitalization shows the successful application of such technologies. There are known examples when the external influence of swindlers could last up to several days, until a quite educated person realized the criminal nature of the influence exerted on him and the critical external dependence of his actions.

Studies of psychosemantic qualities of the personality are conducted within a variety of scientific directions, and their results are used in medicine, forensics and many other practical applications. Thus, the analysis of personal interests and preferences with respect to search and viewing of information and entertainment media programs, computer games, activity in social networks, financial behavior, etc. allows generating and structuring extensive data arrays of electronic content ("big user data"). The result of such analysis can be a cognitive-reflexive model of a particular individual obtained by another interested person. Elements of neuro-linguistic programming can be built into the model, which allows remote and latent creation of templates of perception and interpretation of the occurring flow of events for the target object of influence and that can become a source of motivation and actions of this person. In this case, both individualized actions on creation of given parameters of behavior of a particular person and clustering of influence on target groups are possible. In other words, digital models of personalities can be dissected by target functions and by critical values of parameters, reduced into clusters, in relation to which unified information influences can be applied.[284, 285, 286]

---

[284] Volynsky-Basmanov Y.M. Application of neurolinguistic programming methods to identify potentially dangerous individuals // Problems of safety and emergencies. 2010. № 5. C. 124-128.

At present, given the convergence of information, telecommunications and computing services in global information networks, virtually any user's access to electronic content, regardless of equipment, communication channel and communication method, becomes a special case of access to a single distributed electronic database. Traces of access to electronic content are preserved for a long period of time, allowing the application of proven methods of monitoring analysis, user identification, prediction of reactions and behavior, identification of connections with other users, belonging to explicit and implicit groups, etc.

A modern computer, smartphone or TV remote control with intelligent functions allow not only by the fact, but also by the manner and speed of pressing the keys to accurately identify a person, determine the sign of his attitude to the current electronic content on the screen, identify the levels of sympathy, irritation or aggression in relation to this content. The same devices can record not only the geolocation of a person, but also link it to the intensity of electromagnetic fields in this point, change the parameters of equipment located near this point, affecting electromagnetic fields. Gadgets can record temperature, blood pressure, heartbeat rhythms and other biorhythms of a particular person. Monitoring results make it possible to observe the status and dynamics of a person's psycho-semantic subjectivity in conditions of ambiguity or lack of information about it for making one or another pragmatic decision (ad-

---

[285] Kovalevskaya A.V. Information wars: classification of suggestive specificity // Theoretical and practical problems of language tools transformation in the context of theaccelerated development of public relations In Theoretical and practical problems of language tools transformation in the context of the accelerated development of public relations // Peer-reviewed materials digest (collective monograph) published following the results of the CXVIII International Research and Practice Conference and I stage of the Championship in Philology. Chief editor V.V. Pavlov. 2016. C. 23-25.

[286] Kuznetsov V. The use of neurolinguistic programming (NLP) during interrogation // Law and Life. 2011. № 152 (2). C. 134-140.

vertising, elections, etc.), analyze and predict the risk of changing states of activity in normal and in extreme conditions[287, 288, 289]

Bringing together data from all possible databases of electronic content into a single package of information on the psycho-semantic subjectivity of a personality, which combines structured, complex structured and conventionally structured data, allows to determine the characteristics of its behavioral activity, including hidden qualities (for example, sociopathic tendencies, political or religious predilections, phobias, ways of informal self-realization, network and hierarchical relationships, role structure in families and communities, the volume of resources to be operated, belonging

The development of digital technologies makes it possible not only to integrate textual data with video and images, but also to reconstruct and forecast in 3D and 4-D formats the course of events, to track dynamic changes in the behavior of individuals and groups associated with them so that an aggregate information cluster of worldview and professional templates of interpretation of reality and the resulting behavior occurs. On this basis, new management tools are formed.

The rapidly increasing potential for operating generalized semantic structures creates unprecedented opportunities for transferring external purposeful influence from the conscious to the unconscious level for the object. This is achieved by implicitly implanting to individual personalities semantic motivational "anchors" that induce certain evaluations and

---

[287] Lyakhov A.F., Trishin I.M. Computer modeling of player behavior in an intellectual card game with neural network // Computer instruments in education. 2013. № 5. C. 54-64.

[288] Samartsev O.R., Latenkova V.M. Psychosemantic aspects of the perception of interactive discourse in Internet media // Bulletin of Cherepovets State University. 2016. № 2 (71). C. 87-91.

[289] Ponomareva O.S., Ustyuzhanin V.N. On the status and prospects of the use of psychosemantic methods of cognition of the personality of the suspect in the activity of the investigative officer // Bulletin of St. Petersburg University of the Russian Interior Ministry. 2016. № 2 (70). C. 190-194.

then actions. On this basis it is possible to combine, causing resonance effects, a lot of individual actions in information and real space: from rush on commodity markets to political rallies.

In essence, integration of various methods of psychoengineering, remote influence, psychocorrection and psychosounding with synchronization of psychophysiological and psychosemantic influences and convergence of the process of perception of information stimuli is rapidly developing. Concentration of personal attention and interpretation of information are brought together in a single composition within the framework of human interpretation of the picture of the surrounding real and imaginary world with the consequent lines of behavior.

New possibilities for neurocognitive control are created by the application of the NII. A digital model (digital twin) is able to simulate all possible modes of human activity, take into account the influence of external factors and control processes, allows predicting the future state and behavior of a physical object. The digital twin is based on technologies of artificial intelligence, machine learning and analytical programming. The digital twin continuously learns and updates its parameters, receiving information from multiple sensors, representing the state of the physical object. It learns from current data from sensors, controls, and the environment, and combines factual data with knowledge gained from engineers experienced in the field. The digital twin uses historical data accumulated in previous phases.

In fact, traditional methods of social management are currently being transformed into a unified system with new communication interfaces, neuro- and bio-interfaces. Significant research experience has been accumulated to develop methods for predicting the intellectual dynamics of behavioral activity[290]. There are many promising concepts for creating a

---

[290] Ageev A.I., Loginov E.L. Battle for the Future: Who First in the World to Master Noomonitoring and Cognitive Programming of Subjective Reality? // Economic Strategies. 2017. № 2. C. 124-139; Ageev A.I., Loginov E.L., Shkuta A.A. Convergent monitoring and personality programming as a tool for operat-

multifunctional information monitoring system as a forecasting platform (with feedback) of explicit and implicit deep processes and trends in the society, technosphere and natural environment.

A separate research and practically significant subject is the study of the properties of clusters of factors of biophysical and information-cognitive nature and the prediction of intellectual and behavioral activity going beyond steady states as a source of increased risks of emergencies[291]. Related to this problematics is the task of combining various models of support for the social environment, loyal to legal norms, management key attitudes (including digital, anonymous) and the underlying regulators of life activity ("matrix"). A special issue is the risks of cyber-physical systems failure, which can be caused remotely. At the same time failure can be of program-viral origin, as well as anthropogenic - through the impact on the operator's consciousness. Terrorist acts performed by characters under explicit or latent external control can be inspired in the same way. Modern level of technologies allows to carry

---

ing the intellectual dynamics of behavior of large groups of people // Economic Strategies. 2018. № 2. C. 70-87; Ageev A.I., Loginov E.L. Neuromanagement of personality (forthcoming); Lepsky V.E. Evolution of ideas about management (methodological and philosophical analysis). Moscow: Cogito-Center, 2015. 107 p.; Lefebvre V.A. Reflexion. Moscow: Cogito-Center, 2003. 495 p.; Raikov A.N. Modeling of collective unconscious when making decisions // Proceedings of the International Scientific Conference CPT-2014 International Scientific Conference of the Moscow Institute of Physics and Technology (State University), Institute of Physical and Technical Informatics. Moscow: Institute of Physical and Technical Informatics, 2015. C. 146-156; Smirnov I., Beznosyuk E., Zhuravlev A. Psychotechnologies: Computer psychosemantic analysis and psychocorrection at the unconscious level. M.: Publishing Group "Progress" - "Culture", 1995. 416 p.; Kholodov Y.A. Brain in Electromagnetic Fields. Moscow: Nauka, 1982. 123 c. et al.

[291] Chukhrova M.G., Chukhrov A.S. Spatial and temporal organization of bioelectrical processes of the brain as an indicator of psychosocial adaptation // World of science, culture, education. 2013. № 5 (42). C. 227-230.

out both probing and influence with the account of psycho-semantic qualities of the personality (official and real political orientation, quality of its professional training, cultural level, interests, volitional qualities, inner motivation, state of health, etc.).

The latest generations of technical devices enable the interested party, relying on automated systems of data collection, accumulation, processing and use, not only to identify and determine the geolocation of the device owner, his emotional attitude to the content of electronic messages, but also to identify the characteristics of the surrounding equipment. Obtaining information about the state of personality through "health gadgets" make it possible to predict transitions of personality states in normal and emergency conditions, if necessary. Identification of personality characteristics by characteristics of information programs viewed by them, activity in social networks, choice of computer games, etc. (a data portfolio of individual electronic content) allows to form a cognitive-reflexive model of the personality. On the basis of such a model it is possible to neuroprogram the worldview and situational reference points and acts of behavior of an individual and groups of people. The facts of tests of such technologies are found practically in many modern local conflicts.

The data obtained from all possible forms of electronic content in the package of information on the psycho-semantic subjectivity of the person allows with a high degree of reliability to detect its behavioral dominants and hidden qualities, as well as belonging to a sociopathic group. The object of human control becomes his digital twin, through the influence on the parameters of which it is possible to adjust his real behavior, thinking, interpretations of events and processes. At the same time the twin itself is continuously actualized in the process of "electronic activity" of the person himself. Application of the MIL allows to carry out all necessary predictive analytics of the doppelganger and his prototype.

Especially important is the prediction and proactive elimination of the possibility of resonance of these factors on the basis of technologies

of remote influence, psycho-correction and psycho-sonding, taking into account the psycho-semantic qualities of the individual extended (official and real political orientation, the quality of her professional training, cultural level, interests, volitional qualities, inner motivation, etc.)[292, 293, 294].

The peak manifestations of behavioral activity can be represented as external and internal manifestations of the resonance vectors of biophysical and information-cognitive factors. [295, 296, 297]The calculation of vectors of this resonance seems to be a nontrivial problem, since its components represent information fields applied to the formation of self-

---

[292] Vasilevskaya E.A., Mendelevich V.D. Relationship between social intelligence, anticipatory abilities and IQ in patients with schizophrenia // XVI Congress of Psychiatrists of Russia. All-Russian Scientific-Practical Conference with international participation "Psychiatry in the Stages of Reform: Problems and Prospects": Abstracts / Edited by N.G. Neznanov. 2015. C. 280.

[293] Dyakov S.I. Psychosemantic model and technique of analysis and assessment of personality subjectivity // Scientific conference "Lomonosov Readings" - 2015 / Abstracts. 2015. C. 121-122.

[294] Sevostyanov Yu.O. Change of psychosemantic structure of readiness to work in a team in students // Scientific Bulletin of the Southern Institute of Management. 2014. № 2. C. 94-97.

[295] Ermak E.V., Interrelation of properties of cognitive processing of affective information, emotional intelligence and personality traits // Philosophical problems of biology and medicine / Collected articles. Moscow State University of Medicine and Dentistry. A.I. Evdokimov, Moscow Philosophical Society. 2015. C. 250-253.

[296] Kulikov V.Y., Antropova L.K., Kozlova L.A. Influence of brain functional asymmetry on the behavior strategy of an individual in a stressful situation // Journal of Siberian Medical Sciences. 2010. № 5. C. 10.

[297] Sergievsky G.M., Lobachev V.S. Modeling of Behavior of an Intelligent Agent in a Problem Situation with Partially Observable States // Scientific Session of NRNU MEPhI - 2012 / Abstracts: In 3 vols. 2012. C. 312.

supported resonances of these factors different in qualitative content, dimensionality and other similar indicators[298].

## 15.3 Conclusions

Neurotechnology represents one of the fastest growing segments of the all-encompassing world of artificial intelligence systems. Closer than other branches of AI, neurotechnologies and Neuronet in general face the world of a purely human society undergoing fundamental transformations, including digital ones. New technologies for influencing individual and collective consciousness have emerged and are rapidly improving. There is every reason to pay special attention to the problem of truly conscious subjects, which is one of the key issues for the emerging digital society.

The approaching level of development of AI to the ability to operate in different contours of knowledge and to adjust the goal setting of one's activity raises the question of the ethics of artificial systems. It is equally important to recognize the experience of the evolution of the ethics of pre-digital society, especially in the twentieth century and in modern times. There is an almost unavoidable risk of reflecting real ethical relativism and ethical conflict in the architectures and ontologies, software, and learning processes of AI.

## 15.4 Bibliography

Ageev A.I., Loginov E.L. Battle for the Future: Who First in the World to Master Noomonitoring and Cognitive Programming of Subjective Reality? // Economic Strategies. 2017. № 2, pp.124-139.

---

[298] Yasnitsky L.N., Sichinava Z.I. Neural network algorithms of respondent behavior analysis // Neurocomputers: development, application. 2011. № 10. С. 59-64.

Ageev A.I., Loginov E.L. Neuromanagement of personality (forthcoming)

Ageev A.I., Loginov E.L., Shkuta A.A. Convergent monitoring and personality programming as a tool for operating the intellectual dynamics of behavior of large groups of people // Economic Strategies. 2018. № 2. P. 70-87.

Ageev A.I., Loginov E.L., Shkuta A.A. Neurocontrol: convergent integration of human brain and artificial intelligence [Electronic resource] // Economic Strategies. 2020. № 6, pp. 46-57. DOI: 10.33917/es-6.172.2020. P.46-57.

Butova O.A., Grishko E.A. Features of the formation of bioelectrical activity of the brain neurons of soldiers of Stavropol garrison in the aspect of adaptation // Science. Innovations. Technologies. 2009. № 4, pp. 235-241.

Chukhrova M.G., Chukhrov A.S. Spatial and temporal organization of bioelectrical processes of the brain as an indicator of psychosocial adaptation // World of science, culture, education. 2013. № 5 (42), pp. 227-230.

DARPA to fund the creation of high-speed interface between computer and human brain [Electronic resource] // Open Systems Publications. URL: https://www.osp.ru/news/2016/012 6/13031430.

Dyakov S.I. Psychosemantic model and technique of analysis and assessment of personality subjectivity // Scientific conference "Lomonosov Readings" - 2015 / Abstracts. 2015, pp. 121-122.

Ermak E.V. Interrelation of properties of cognitive processing of affective information, emotional intelligence and personality

traits // Philosophical problems of biology and medicine / Collected articles. Moscow State University of Medicine and Dentistry. A.I. Evdokimov, Moscow Philosophical Society. 2015, pp. 250-253.

Galushkin A.I. Neurochips and Neuromorphic Computers: Problems of Modeling // Information Technologies. 2015. T. 21. № 12, pp. 942-949.

GOST P 43.0.3-2009. Information support of technical and operator activities. Noon-technology in technical activity. General provisions. Moscow: Standartinform, 2010.

Kalinin P.V., Voyutskaya Yu. On the application of a neuro-interface for contactless control of a mobile device // Information Systems and Technologies. 2016. № 3 (95), pp. 53-56.

Kholodov Y.A. Brain in Electromagnetic Fields. Moscow: Nauka, 1982. 123p.

Kovalevskaya A.V. Information  wars: classification of  suggestive specificity // Theoretical and practical problems of language tools transformation in the context of theaccelerated development of public relations In Theoretical and practical problems of language tools transformation in the context of the accelerated development of public relations // Peer-reviewed materials digest (collective monograph) published following the results of the CXVIII International Research and Practice Conference and I stage of the Championship in Philology. Chief editor V.V. Pavlov. 2016, pp. 23-25.

Kulikov V.Y., Antropova L.K., Kozlova L.A. Influence of brain functional asymmetry on the behavior strategy of an individual

in a stressful situation // Journal of Siberian Medical Sciences. 2010. № 5, p. 10.

Kuznetsov V. The use of neurolinguistic programming (NLP) during interrogation // Law and Life. 2011. № 152 (2), pp. 134-140.

Lefebvre V.A. Reflexion. Moscow: Cogito-Center, 2003. 495p.

Lepsky V.E. Evolution of ideas about management (methodological and philosophical analysis). Moscow: Cogito-Center, 2015. 107p.

Loginov E.L., Loginova V.E., Shkuta A.A. "Design thinking" elements of artificial intelligence to overcome the barriers of obtaining new knowledge in the electronic environment of collaborative scientific super-system // Artificial societies. 2018. № 3, p. 5.

Lyakhov A.F., Trishin I.M. Computer modeling of player behavior in an intellectual card game with neural network // Computer instruments in education. 2013. № 5, pp. 54-64.

Neurogress: a platform of neurocontrol systems from the BlueBrain project [Electronic resource] // ITnan. URL: https://itnan.ru/post.php?c=1&p=348498.

Neuro-interfaces for next-generation services and products [Electronic resource] // Basisneuro. 2017. November. URL: https://basisneuro.com/BasisNeuroWhitePaper.pdf.

Ovod I.V., Osadchiy A.E., Pupyshev A.A., Fradkov A.L. Formation of neurofeedback based on adaptive model of brain activity // Neurocomputers: Design, Application. 2012. № 2, pp. 36-41.

Ponomareva O.S., Ustyuzhanin V.N. On the status and prospects of the use of psychosemantic methods of cognition of the personality of the suspect in the activity of the investigative officer // Bulletin of St. Petersburg University of the Russian Interior Ministry. 2016. № 2 (70), pp. 190-194.

Raikov A.N. Convergent synthesis of cognitive model based on deep learning and quantum semantics // International Journal of Open Information Technologies. 2018. Т. 6. № 12, pp. 43-50.

Raikov A.N. Modeling of collective unconscious when making decisions // Proceedings of the International Scientific Conference CPT-2014 International Scientific Conference of the Moscow Institute of Physics and Technology (State University), Institute of Physical and Technical Informatics. Moscow: Institute of Physical and Technical Informatics, 2015, pp. 146-156.

Samartsev O.R., Latenkova V.M. Psychosemantic aspects of the perception of interactive discourse in Internet media // Bulletin of Cherepovets State University. 2016. № 2 (71), pp. 87-91.

Sergievsky G.M., Lobachev V.S. Modeling of Behavior of an Intelligent Agent in a Problem Situation with Partially Observable States // Scientific Session of NRNU MEPhI - 2012 / Abstracts: In 3 vols. 2012, p. 312.

Sevostyanov Yu.O. Change of psychosemantic structure of readiness to work in a team in students // Scientific Bulletin of the Southern Institute of Management. 2014. № 2, pp. 94-97.

Smirnov I., Beznosyuk E., Zhuravlev A. Psychotechnologies: Computer psychosemantic analysis and psychocorrection at the un-

conscious level. M.: Publishing Group "Progress" - "Culture", 1995. 416p.

Turovsky J.A., Kurgalin S.D., Adamenko A.A. Modeling the learning of neurochips embedded in the neural tissue // Digital Signal Processing. 2016. № 1, pp. 50-56.

Tychkov A.Y., Goryachev N.V., Kochegarov I.I. Communication protocols for wireless neural interface // Proceedings of the International Symposium "Reliability and Quality". 2018. V.2, pp. 366-368.

Vasilevskaya E.A., Mendelevich V.D. Relationship between social intelligence, anticipatory abilities and IQ in patients with schizophrenia // XVI Congress of Psychiatrists of Russia. All-Russian Scientific-Practical Conference with international participation "Psychiatry in the Stages of Reform: Problems and Prospects": Abstracts / Edited by N.G. Neznanov. 2015. 280p.

Volynsky-Basmanov Y.M. Application of neurolinguistic programming methods to identify potentially dangerous individuals // Problems of safety and emergencies. 2010. № 5, pp. 124-128.

Yasnitsky L.N., Sichinava Z.I. Neural network algorithms of respondent behavior analysis // Neurocomputers: development, application. 2011. № 10, pp. 59-64.

**PART D**


**SERVING HUMANS:**
**RELIGIONS FOR DIGITAL JUSTICE**

# DIGITAL JUSTICE MANIFESTO:
# A CALL TO OWN OUR DIGITAL FUTURE

*Just Net Coalition /*
*World Association of Christian Communication WACC*

*The Digital Justice Manifesto, "A Call to Own Our Digital Future", was launched in Berlin in November 2019 by the Just Net Coalition[299], a global network of civil society organizations and individuals, including the World Association for Christian Communication (WACC).[300] The coalition was founded at a meeting in New Delhi in 2014, which agreed "The Delhi Declaration for a Just and Equitable Internet".*

## Summary: Key Principles

1. Data subjects must own their data – individually and collectively
2. Our data requires protection from abuse
3. We need the tools to control our data

---

[299] Contact: info@justnetcoalition.org.

[300] The following text was published on https://waccglobal.org/digital-justice-manifesto-a-call-to-own-our-digital-future.

4.   Data commons need appropriate governance frameworks

5.   Data protection and sharing require new regimes

6.   Data-creating jobs ought to come with data rights

7.   Data should be processed close to point of origin

8.   Cross-border data flows must be decided nationally

9.   Techno-structures need to be reclaimed as personal and public spaces

10. We own and should control our software

11. Key digital infrastructures need to be governed as public utilities

12. Techno-structures must be decentralised for open use, with interoperability

13. Global data monopolies should be broken

14. Societies' datafication needs to be managed democratically

15. Digital standards must be developed by public interest bodies

16. The digital has to be governed in a local-to-global manner.

## 16.1 Preamble

We begin by endorsing and reaffirming *The Delhi Declaration for a Just and Equitable Internet.*[301] The present Manifesto builds on that Declaration and extends it.

*A digital society is upon us.* People respond to the emerging digital future with a mixture of positive anticipation, awe, helplessness and even horror. Such a passive reaction to *society's most powerful driving force is both dangerous and unnecessary.* There is no time to lose in taming the power of the digital. We can surrender our digital future, or we can take ownership of it. But first we must understand *what lies behind the digital.*

Industrialisation harnessed massive physical power from sources beyond those of people and animals, that transformed the processes of production. This is known as mechanisation. A digital economy and

---

[301] https://justnetcoalition.org/delhi-declaration, 2014.

society is created by harnessing external (non-human) sources of 'intelligence power' i.e. the immense data-based intelligence that is revolutionising the forces of production. *This can be called the "Intelligencification" of socio-economic processes.*

Colonisation bore horrific witness to how industrial power and the imperatives of capital were almost impossible to resist or challenge by those subjected to them. Yet *the power of others owning detailed intelligence about and over us*, and processing it through elaborate socio-economic systems to generate enormous profits through limitless manipulation, is perhaps worse than anything we have experienced so far.

**Data, intelligence and techno-structures.** Data must be recognised as a key economic resource. Currently, the resource of data gets globally appropriated at will; harvested without permission or recompense, and accumulated by data corporations for their exclusive use. We must choose whether to allow corporations to own *people's data, or if we, the people, should own it. The people, after all, are both the data contributors and data subjects.* Data corporations take advantage of default lawless practises of data accumulation and exclusive use; challenging them requires countervailing laws affirming people's rights and ownership over their data – both individual and collective.

Digital "Intelligencification" was preceded and enabled by the spread of *networked software as the space, means and logic of our social, economic, political and cultural interactions and relationships.* The Internet was its first prototype. And because the Internet's core model was based on intelligence at the edges and on open, public protocols, it spawned a technical and social evolution that many believed would favour greater end-user control and decentralisation. Cloud computing – now the dominant networked software model – has *inverted that paradigm: controlling intelligence now lies with a few global centres, based largely on corporate control and ownership of data and private standards.* The ubiquitous spread of Internet-based cloud applications enables the collection of the most intimate and granular real-time data about us, the

people. Such intensive and unrelenting data collection underpins and enables the powerful autonomous intelligence behind the phenomenon of digital society.

At the centre of intelligent digital systems are a few global business-es – 'intelligence corporations', whose services are based on digital intelligence or artificial intelligence (AI). They first connect, then coordinate, and ultimately control and dictate to all actors and activities in any sector – from transport and commerce to health and education. *They become the 'brain' of every sector. Global intelligence corporations operate remotely through techno-structures of cloud computing.* Bypassing face-to-face human interactions, they thus avoid responsiveness and accountability, as well as legal and regulatory checks.

*Taking back digital power.* Reclaiming power from 'intelligence cor-porations' requires people to work on two main fronts. First, *wrest back ownership of our personal and collective data and intelligence.* These are the key sources of digital power. Second, *take sufficient control over the techno-structures* within which data and intelligence operate. These techno-structures spread far, wide and deep into society, controlling and exploiting everything they reach. Unlike in the offline world where so-cio-economic interactions mostly take place in public or quasi-public spaces, in the digital world they are all enclosed within privately owned techno-structures.

Yet intelligent systems can operate productively and beneficially even when their intelligence, as well as the key nodes and pillars of their techno-structures, are distributed and collectively owned. It would in-volve employing the best possibilities of entrepreneurship and competi-tive markets, combined with critically important non-market collective mechanisms. *Such alternatives must be shaped at the same time that exploitative dominant models of centralised intelligence control are undone.*

The digital reshapes our social relationships and power structures so fundamentally that society's *data and intelligence governance requires a new digital social contract.*

In our determination and struggle *to enable people to own their digital futures,* we adopt, and advocate, the following principles – towards a digital society that is just, equitable and sustainably productive.

## 16.2 Resolutions and Principles

### 16.2.1 People own their data and intelligence.

1. *Data subjects must own their data – individually and collectively:* Data about us, and intelligence about us, inherently belong to us – as individuals, and collectively as communities. Such data could directly be about people, or about things owned by or associated with them. Political, constitutional, and legal frameworks, at both national and international levels, must recognise and enforce this basic principle of data and intelligence ownership.

2. *Our data requires protection from abuse:* The international human rights regime must recognise the inextricable interconnection between people and their data, and articulate benchmarks for safeguarding personal and collective data. Strong constitutional and legal protections are required against abuse of personal and collective data and intelligence, whether by corporations or the state. New laws and institutions are required for this purpose, that keep evolving to address emergent new risks.

3. *We need the tools to control our data*: The purpose of data and intelligence must not be to distinguish between people for unfair, discriminatory treatment, but to help and enable them to maximise digital benefits. Individuals and communities must be provided appropriate means to control their data, and apply it in ways best suited to their interests. Such means will both be individual and collective, requiring institutions that are adequate, practical and effective. It would involve well-regulated

open markets and competitive businesses, as well as establishing new commons and public structures. This demands considerable institutional innovation.

4. ***Data commons need appropriate governance frameworks:*** Data commons and intelligence commons must be developed and treated as public goods. But data and intelligence cannot simply be open access resources; significant protections and circumscribing are essential. Being specific to particular individuals or groups and communities, unchecked access to, and use of, data and intelligence commons bear the potential for harm. In the manner that data actually gets employed by monopoly-prone businesses, data and intelligence commons are akin to 'common pool resources' – subject to overuse, depletion, congestion, rivalry and pollution. Requiring regulated and calibrated use, data and digital intelligence must be subject to 'common property regimes', which calls for development of appropriate data and intelligence governance frameworks.

5. ***Data protection and sharing require new regimes:*** Innovative institutions should evolve for sharing of data and intelligence in a protected and regulated manner. Data institutions. such as data commons, data trusts, data infrastructures, and fair data markets, must be developed. As appropriate, these should involve mandated data sharing. Businesses and other entities have simultaneously to be provided sufficient incentives, within a public interest framework, for them to collect the necessary data and convert it into useful intelligence.

6. ***Data-creating jobs ought to come with data rights:*** Specific economic groups that make marked contributions to, and are key subjects of, data in a particular sector or an 'intelligent system', should have corresponding special data ownership rights. These could be drivers on a taxi platform, traders on an e-commerce platform, farmers on an agri-platform or workers in data-producing jobs. These groups should have primary economic rights – individual and collective – over the data they contribute. Such data constitutes the main value of the corresponding

platform or intelligent system. Data-creating actors on a platform therefore have the right to participate in the governance of that platform, for example through adequate representation on the governing board. Alternatively, they may choose to pool their data to develop platform cooperatives, or a public or non-profit agency can help them to so organise.

7. ***Data should be processed close to point of origin:*** Important data will need to be localised in many cases. Data can be processed close to its point of origin, for which technologies already exist, and further innovations will emerge as society demands them. This can provide data subjects more effective control over their data. Necessary technical, policy and business models should be employed towards a local-to-global architecture of digital services. In contrast to the current situation where digital activities on the ground are largely remote-controlled 'satellite operations' of a few global corporations, digital should have a pronounced local-ness and community control.

8. ***Cross-border data flows must be decided nationally:*** The data-owning national community must determine the terms on which cross-border flows of data may take place. Irrespective of its physical location, data should be subject to the *primary* jurisdiction of data's country of origin. As personal data is an extension of one's person-hood, so community data is an extension of community identity and being. Such primary jurisdiction involves not just privacy protections but also economic rights and ownership. Agreements among countries are required to mutually recognise, and help apply, *primary* jurisdiction over data – involving social, political and economic rights – of the country and community of origin of data. Regional groups that manage such inter-country agreements may gain mutual benefit from common data and digital spaces.

### 16.2.2 People have rights to their digital techno-structures

9. ***Techno-structures need to be reclaimed as personal and public spaces:*** Networked software or cloud applications form the digital space,

and the body of digital systems. These may be termed as the key digital techno-structures. They are currently almost entirely centralized and owned by a handful of corporations. Some, like those running heart pacemakers or mobile phones, penetrate deep into our otherwise personal realms; and others, like social networking, search, and transport applications, are analogous to what in the offline world are public spaces and structures, such as public streets, libraries and infrastructural services. Digital techno-structures' essential personalness and publicness must be reclaimed and restored, where and as appropriate, from the existing state of their complete, end-to-end, corporate ownership and control.

10. ***We own and should control our software:*** People must fully own, and be able to control, the software they install in their personal or collectively owned equipment. Technology Protection Measures are incursions upon people's basic rights. People should have the right to own, break-into, modify or remove, as they deem fit, whatever technical artefacts that exist within their personal or collective realms. This is a fundamental element of digital self-determination.

11. ***Key digital infrastructures need to be governed as public utilities:*** In the physical, offline world, non-personal, social and economic spaces and structures are divided between being public and belonging to private businesses. Infrastructure is normally public, or quasi-public, over and around which businesses may undertake their private activities. Digital spaces and structures require a similar arrangement. Key monopolistic digital infrastructures should be public utilities, even if provided by private businesses. This includes, as appropriate, computing platforms, search engines, social networks, email services, basic security systems, payment services, and e-commerce platforms.

12. ***Techno-structures must be decentralised for open use, with interoperability:*** Digital power can be redistributed by decentralising the techno-structures of connectivity, software, Internet, cloud computing, and AI applications/engines, while mandating interoperability. Such decentralisation is useful even where it entails some degree of immediate

loss of efficiency. Apart from being more fair, decentralised digital power is more sustainably productive in the long term. Decentralised and open digital architectures include open community networks, open source software, open and neutral Internet, open and community data, and open and community AI. These can and must involve appropriate business models and entities. Any such open system should however duly protect the data and digital intelligence of the people and communities concerned, and affirm their right to self-determination.

13. ***Global data monopolies should be broken:*** National and international competition regimes, that are adequate to the new digital realities, must break up vertically and horizontally integrated global digital structures. These regimes should aim at *ex ante* open, competitive and innovation-supporting digital market structures, and not just narrowly construed *ex post* consumer welfare that looks only at availability and price of services. The focus should be on cutting problematic links in data and intelligence value chains that underpin and promote digital monopolies. It may for instance be considered, where appropriate, to separate businesses that directly provide services to consumers, and collect their data, from businesses devoted specifically to technical services and general data processing and digital intelligence services.

### 16.2.3 The digital must be governed democratically, from local to global

14. ***Societies' datafication needs to be managed democratically:*** Processes and areas facing or undergoing datafication and "Intelligencification" require a three-way classification. Many kinds of datafication and "Intell6gencification" are just not desirable, whatever be the touted benefits. In other areas, while potentially useful in the long run, these processes may call for deliberate slowing down and close management, to deal with the possibilities of considerable short- to mid-term harm. Such harm could range from livelihood disruptions to requirements of considerable behaviour and cultural shifts that can be disorienting.

Where datafication and "intelligencification" are evidently and immediately beneficial, people, and their representatives, should be in control of their implementation. These processes tend to have unanticipated, strong social consequenc6s and must take place on democratically determined terms. A global human rights framework on data and intelligence governance should incorporate such a classification and corresponding due diligence.

15. ***Digital standards must be developed by public interest bodies:*** A major factor behind the current end-to-end digital control by a few digital corporations is the privatisation of digital standards development and non-enforcement of interoperability. We must reclaim development of key digital technical standards exclusively by public interest bodies, and ensure strict compliance with them. These bodies should be appropriately representative of people, and based on public-interest oriented expertise. They should uphold the highest public and professional standards, be neutral and not aligned to any specific corporate or state interests, and fully eschew conflicts of interest.

16. ***The digital has to be governed in a local-to-global manner:*** Digital platforms provide services that have traditionally been largely developed and governed locally – like communication, media, commerce, transport, hotels, health and education. Having now become intelligence-driven does not necessarily mean that these services shed their localness. The required digital, data and intelligence governance structures and institutions will mostly be at national or local community levels, while some could be global. National polities still remain the anchors of self-determination and sovereignty of the people. Appropriate global governance of the digital should promote national and local digital economies. It ought to ensure that competitive, open, global technical services are accessible locally – including by local digital businesses – on fair, regulated terms. Digital governance must aim at a complete break with the current vertically-integrated global digital models – from concentrated intelligence or 'brain' centres in one or two countries of the

world, right down to the last tiny 'nerves' that seek to control the smallest activity everywhere in a digital economy and society. A new digital model must be shaped that is local-to-global; that supports localness and furthers democratic self-determination, without compromising on the important benefits of the globalness of the digital.

<div align="center">

**17**

</div>

<div align="center">

# A VISION OF DIGITAL JUSTICE

</div>

<div align="center">

*World Council of Churches, Central Committee*

</div>

*The Central Committee (Parliament) of the World Council of churches, representing over 350 churches from all continents, adopted the following statement at its video conference meeting from 9-15 February 2022.*

*The LORD loves righteousness and justice; the earth is full of his unfailing love. (Psalm 33:5)*

## 17.1 Introduction

Digital technologies are transforming our world and the multiple spaces in which we live and move. These technologies offer us new ways to communicate, to inform ourselves and navigate the world, to advocate for human dignity and rights, and for multiple voices to be heard.

However, while digital technologies can be powerful tools for living in relation with others, for inclusion, education, encounter, imagination, creativity, and understanding, they also pose challenges through their conception, design, and use. Digital communication platforms, with all their opportunities, are used to spread disinformation and hate, exploit individuals and communities, increase surveillance, and contribute to growing gaps in access, power, and wealth – including between private companies and national governments.

In 2022, the ecumenical fellowship gathers in Karlsruhe, Germany, for the 11th Assembly of the World Council of Churches (WCC), and this digital transformation of society raises profound issues that the ecumenical fellowship has wrestled with for many decades: power, justice, equity, participation, promoting sustainable communities, care for creation, how voices from the margins are heard, as well as human dignity, and what it means to be human, made in the image of God.

The WCC has a long history of raising issues related to the role of communication in society and its implication for ecumenical witness. The WCC's founding assembly in Amsterdam in 1948 referred to the vast concentration of power in a society "dominated as it is by technics," while also noting how technical progress "provides channels of communication and interdependence which can be aids to fellowship." Subsequent assemblies returned to the challenges posed by communication for society and ecumenical witness.

In responding ecumenically to the contemporary challenges posed by digital transformation, the ecumenical fellowship can draw upon the resources and insights gained over the past seven decades in grappling with the challenges of communication. In particular, we recall that the Uppsala (1968) and Vancouver (1983) assemblies focused on key contemporary aspects of mass communication.

A statement on "The Church and the Media of Mass Communication" presented to the WCC's 4th Assembly in Uppsala in 1968 described mass media as the "agora and town meeting of technological

society," and the mainstay for a "responsible world society," while, at the same time, "[p]owerful minorities can pervert the media for limited or selfish purposes."

Recommendations on "Communicating Credibly" from the 6th Assembly in Vancouver in 1983 warned of the "pervasive and often dominant role" of mass media," which "distort and diminish the life of the world, by packaging it as entertainment, or simply as propaganda." At the same time the new electronic media threatened to "enlarge and confirm the global domination of a few countries and make it almost irreversible," widening "the gap between the information-poor and information-rich both within and between individual nations".

In today's digital world, we see many of these same dilemmas not only replicated but magnified by the speed of digital developments and networks and the ability of individuals themselves to be active participants and "influencers" as never before.

As we wrestle with these continuing and new opportunities and challenges, two intrinsically connected aspects must play a central role in our vision and theological reflection on digital justice: relationality and vulnerability. As Christians, we believe that being created in the image of God provides inherent dignity to every woman, man, and child (Gen 1:27.) Humans are created to be relational and capable of collaboration and communication. For this reason, we are called to take responsibility for and to care for God's creation.

To address the urgent issues raised by digital transformation, we believe the time has arrived for a "new communications paper for the 21st century", drawing in particular on the insights of the symposium on "Communication for social justice in a digital age" that took place in September 2021.

## 17.2 Global and Societal Digital Realities

The COVID-19 pandemic accelerated the digital transformation as organisations and individuals moved online to maintain livelihoods, education, worship, and connections when physical interaction was necessarily limited. For many, what we could accomplish online was a revelation – reconnecting to distant family and friends, holding meetings without the obstacles of travel and funds, finding new ways to live and express ourselves as Christians in a globalised society. At the same time, we lamented the loss of physical connections that deepen our relationships.

Our growing digital dependence has highlighted distinct challenges. Communities of faith, political, cultural, and civil society actors are all struggling to respond effectively to these issues of social justice, including:

*Digital Divides:* One challenge is "digital divides" between those who have access to new technologies and those who do not. Such digital divides exist at many levels: economic, geographic, racial, educational, class, gender, age, cultural, technological, and global. These divides point to both the complexity of social justice in a digital age and the need for intersectional reflection. In the same way, digital justice must be seen in the context of gender equity, racism, environmental sustainability, economic justice, intergenerational relationships, and so much more.

*Accessibility:* Meaningful access requires basic communication infrastructure such as stable electricity and internet connections, the availability of digital tools, data, programming and content from the local cultural context, as well as the legal frameworks and economic resources to access and invest in them. Access is also enabled by language diversity and facilities for full inclusion of persons with disabilities. Access affects power relations and distribution of resources, and as such, access to digital technologies is both a cause and a result of divides.

*Public space:* This is the space where States and the public interact, where people, including the media, can express thoughts and feelings and participate democratically. Digitalization creates the opportunity for expanding this space, but the restriction of digital freedoms and abuse of the space through disinformation and hate speech can also cause this space to shrink.

*Inequity:* Control, use, and analysis of data gathered due to digitalization are heavily vested in a few corporations and in specific geographic regions. Governments may also be heavily implicated in data control and manipulation.

*Education:* Digital education, including questioning and critical examination of information and sources, is vital for all people. Access to this education is often sharply divided based on age, academic background, language, gender, geographical location, and societal gender roles.

*Gender justice:* Active participation in the digital space can contribute to full participation in all domains of life for everyone. However, pervasive gender power inequities restrict this access, and gender identity can be a target for discrimination and online abuse. For example, increased digitization has led to greater exposure of girls and women to sexualized harassment, surveillance, trolling, and online hate, which may also lead to physical violence. The impact of online violence is silencing women in particular, forcing some to disengage from the digital space.

*Privacy and security:* The universal challenges of the use of data and loss of privacy are compounded by arbitrary government control, national digital laws and guidelines that are vague and fraught with loopholes, corporate interests, state censorship and surveillance. Digital technologies can also present unforeseen dangers.

*Militarization:* There is military investment in digital technologies, and the technologies are in turn militarised –increasing risk in situations of war and conflict.

*Artificial Intelligence and cyberspace:* Artificial Intelligence and related technologies are transforming education, surveillance, foreign policy, commerce, and industry, and the way humans relate to each other. In the military sphere, there is intensive development of AI-guided missiles and autonomous weapons, including drones.

## 17.3 Our Call for Digital Justice

We are called to a journey of justice and peace and to ensure the integrity of creation. Digital technology itself is a product of human creativity and should be celebrated when it is used to enhance human dignity.

We are called to participate in God's mission to ensure that all may have life and have it abundantly (John 10:10), also in the digital sphere. The biblical preferential option for the poor and vulnerable (Matt 5) directs our attention to information poverty and the digital divides in the global face of digitalization.

The issue of digital justice relates to the whole range of issues of concern to the ecumenical fellowship. The increasing militarization of digital technology and artificial intelligence, its use in surveillance and censorship, and the deliberate destabilising disinformation tactics impact our public witness.

The power and wealth concentrated in a few technological organisations, and the ecological impact of mining for the resources needed to meet the increasing demand of digital devices has clear implications for sustainability and the economy of life.

The gender gap in digital access, biased algorithms, and online abuse and violence against women affect our efforts towards a Just Community of Women and Men. Rapid developments and the use of artificial intelligence impact areas of health, race, and what it means to be human.

While addressing the challenges, we have the opportunity to identify and model digital justice in education, mission and evangelism, and

work in human rights as well with international organisations and inter-faith partners on the path of justice and peace.

Reviewing the ecumenical discussion over these seven decades, several common elements stand out:

- the right to access information, and thus opposition to policies that interfere with or undermine this right;
- the need to support truth and to express vigorous opposition when what is communicated distorts the truth, reinforces negative stereotypes, or supports violent behaviour;
- the need for pluralism and the voicing of diverse views and to work against media concentration;
- the need to protect communication freedoms in the context of global economic and political structures of justice and injustice;
- the need to support and advocate for the right to communicate for marginalised communities and those whose voices are suppressed;
- the need to support alternate means of communication such as theatre, special liturgies, and local, indigenous newspapers and radio; and
- the need for the ecumenical fellowship to offer an alternative vision of communication, based on solidarity and sharing, mutual accountability, and empowerment.
- The global ecumenical fellowship will gather in Karlsruhe, in a world marred by many kinds of injustice and by the pain of many of its people, its creatures, and even of the earth itself. But it is also a world witnessing movements of change, justice, and hope.
- To bring about digital justice, the WCC supports a transformative movement led by civil society and communities of faith. The broad support and joint commitment of civil society, including churches and faith communities, political actors, science, and business, is needed to guarantee and protect civil rights in the digital age and make the digital space usable for the common good. Technologies

must be placed at the service of people rather than governments or corporations.

- We commit ourselves, in our local contexts and as a global ecumenical fellowship, to address the challenges posed by digital justice in our work and advocacy on gender equality, environmental sustainability, human rights, democratic participation, and economic justice. This involves:
- Re-imagining our digital public sphere continuously, emphasizing democracy, fundamental rights, mutual accountability, and solidarity.
- Raising awareness and promoting educational programmes, including theological formation, to equip ourselves to face the challenges of digital transformation.
- Working with state and civil society actors and faith groups to create spaces and channels that are inclusive, accessible, interactive, and participatory, promoting racial justice, gender justice, digital justice, expanding public spaces, and creating visions for the future.
- Encouraging theological and ethical critiques of the digital transformation, including the powers that operate unregulated, commercially driven digital spaces.
- Creating an informed, grassroots, faith-inspired resistance to the forces challenging human dignity and flourishing in digital spaces.
- Developing cross-cutting programmes of action to create this re-imagined reality in different contexts.

*Blessed are those who act justly, who always do what is right. (Psalm 106:3)*

# 18

## THE ROME CALL FOR AI ETHICS AND THE ABRAHAMIC COMMITTMENT

*Pontifical Academy for Life, Vatican[302]*

*The Pontifical Academy for Life in Rome, with the RenAIssance Foundation, a new Vatican Foundation specialised on AI Ethics, published on 28 February 2020 the "Rome Call for AI Ethics". It is published below.*

*On 10 January 2023, three representatives of the Abrahamic religions Christianity, Islam and Judaism, together with representatives of IT companies and the academic sector, signed "The Abrahamic Commitment to the Rome Call for AI" in Rome.[303] This commitment is added in part 18.5 of the following text.*     *The Editors/CS*

---

[302] Rome Call for AI Ethics, Vatican, 28 Feb. 2020. www.vatican.va. https://www.romecall.org/.

[303] https://www.romecall.org/the-abrahamic-commitment-to-the-rome-call-for-ai-ethics-10th-january-2023/

336 Data Ethics: Building Trust

## 18.1 Introduction

"Artificial intelligence" (AI) is bringing about profound changes in the lives of human beings, and it will continue to do so. AI offers enormous potential when it comes to improving social coexistence and personal well-being, augmenting human capabilities and enabling or facilitating many tasks that can be carried out more efficiently and effectively. However, these results are by no means guaranteed. The transformations currently underway are not just quantitative. Above all, they are qualitative, because they affect the way these tasks are carried out and the way in which we perceive reality and human nature itself, so much so that they can influence our mental and interpersonal habits. New technology must be researched and produced in accordance with criteria that ensure it truly serves the entire "human family" (Preamble, Univ. Dec. Human Rights), respecting the inherent dignity of each of its members and all natural environments, and taking into account the needs of those who are most vulnerable. The aim is not only to ensure that no one is excluded, but also to expand those areas of freedom that could be threatened by algorithmic conditioning.

Given the innovative and complex nature of the questions posed by digital transformation, it is essential for all the stakeholders involved to work together and for all the needs affected by AI to be represented. This Call is a step forward with a view to growing with a common understanding and searching for a language and solutions we can share. Based on this, we can acknowledge and accept responsibilities that take into account the entire process of technological innovation, from design through to distribution and use, encouraging real commitment in a range of practical scenarios. In the long term, the values and principles that we are able to instil in AI will help to establish a framework that regulates and acts as a point of reference for digital ethics, guiding our actions and promoting the use of technology to benefit humanity and the environment.

Now more than ever, we must guarantee an outlook in which AI is developed with a focus not on technology, but rather for the good of humanity and of the environment, of our common and shared home and of its human inhabitants, who are inextricably connected. In other words, a vision in which human beings and nature are at the heart of how digital innovation is developed, supported rather than gradually replaced by technologies that behave like rational actors but are in no way human. It is time to begin preparing for more technological future in which machines will have a more important role in the lives of human beings, but also a future in which it is clear that technological progress affirms the brilliance of the human race and remains dependent on its ethical integrity.

## 18.2 Ethics

All human beings are born free and equal in dignity and rights. They are endowed with reason and conscience and should act towards one another in a spirit of fellowship (cf. Art. 1, Univ. Dec. Human Rights). This fundamental condition of freedom and dignity must also be protected and guaranteed when producing and using AI systems. This must be done by safeguarding the rights and the freedom of individuals so that they are not discriminated against by algorithms due to their "race, colour, sex, language, religion, political or other opinion, national or social origin, property, birth or other status" (Art. 2, Univ. Dec. Human Rights).

AI systems must be conceived, designed and implemented to serve and protect human beings and the environment in which they live. This fundamental outlook must translate into a commitment to create living conditions (both social and personal) that allow both groups and individual members to strive to fully express themselves where possible.

In order for technological advancement to align with true progress for the human race and respect for the planet, it must meet three re-

quirements. It must include every human being, discriminating against no one; it must have the good of humankind and the good of every human being at its heart; finally, it must be mindful of the complex reality of our ecosystem and be characterised by the way in which it cares for and protects the planet (our "common and shared home") with a highly sustainable approach, which also includes the use of artificial intelligence in ensuring sustainable food systems in the future. Furthermore, each person must be aware when he or she is interacting with a machine.

AI-based technology must never be used to exploit people in any way, especially those who are most vulnerable. Instead, it must be used to help people develop their abilities (empowerment/enablement) and to support the planet.

## 18.3 Education

Transforming the world through the innovation of AI means undertaking to build a future for and with younger generations. This undertaking must be reflected in a commitment to education, developing specific curricula that span different disciplines in the humanities, science and technology, and taking responsibility for educating younger generations. This commitment means working to improve the quality of education that young people receive; this must be delivered via methods that are accessible to all, that do not discriminate and that can offer equality of opportunity and treatment. Universal access to education must be achieved through principles of solidarity and fairness.

Access to lifelong learning must be guaranteed also for the elderly, who must be offered the opportunity to access offline services during the digital and technological transition. Moreover, these technologies can prove enormously useful in helping people with disabilities to learn and become more independent: inclusive education therefore also means using AI to support and integrate each and every person, offering help and opportunities for social participation (e.g. remote working for those

with limited mobility, technological support for those with cognitive disabilities, etc.).

The impact of the transformations brought about by AI in society, work and education has made it essential to overhaul school curricula in order to make the educational motto "no one left behind" a reality. In the education sector, reforms are needed in order to establish high and objective standards that can improve individual results. These standards should not be limited to the development of digital skills but should focus instead on making sure that each person can fully express their capabilities and on working for the good of the community, even when there is no personal benefit to be gained from this.

As we design and plan for the society of tomorrow, the use of AI must follow forms of action that are socially oriented, creative, connective, productive, responsible, and capable of having a positive impact on the personal and social life of younger generations. The social and ethical impact of AI must be also at the core of educational activities of AI.

The main aim of this education must be to raise awareness of the opportunities and also the possible critical issues posed by AI from the perspective of social inclusion and individual respect.

## 18.4 Rights

The development of AI in the service of humankind and the planet must be reflected in regulations and principles that protect people – particularly the weak and the underprivileged – and natural environments. The ethical commitment of all the stakeholders involved is a crucial starting point; to make this future a reality, values, principles, and in some cases, legal regulations, are absolutely indispensable in order to support, structure and guide this process.

To develop and implement AI systems that benefit humanity and the planet while acting as tools to build and maintain international peace, the

development of AI must go hand in hand with robust digital security measures.

In order for AI to act as a tool for the good of humanity and the planet, we must put the topic of protecting human rights in the digital era at the heart of public debate. The time has come to question whether new forms of automation and algorithmic activity necessitate the development of stronger responsibilities. In particular, it will be essential to consider some form of "duty of explanation": we must think about making not only the decision-making criteria of AI-based algorithmic agents understandable, but also their purpose and objectives. These devices must be able to offer individuals information on the logic behind the algorithms used to make decisions. This will increase transparency, traceability and responsibility, making the computer-aided decision-making process more valid.

New forms of regulation must be encouraged to promote transparency and compliance with ethical principles, especially for advanced technologies that have a higher risk of impacting human rights, such as facial recognition.

To achieve these objectives, we must set out from the very beginning of each algorithm's development with an "algor-ethical" vision, i.e. an approach of ethics by design. Designing and planning AI systems that we can trust involves seeking a consensus among political decision-makers, UN system agencies and other intergovernmental organisations, researchers, the world of academia and representatives of non-governmental organizations regarding the ethical principles that should be built into these technologies.

For this reason, the sponsors of the call express their desire to work together, in this context and at a national and international level, *to promote "algor-ethics"*, namely the ethical use of AI as defined by the following principles (see on the next page):

1. ***Transparency***: *in principle, AI systems must be explainable;*

2. ***Inclusion***: *the needs of all human beings must be taken into consideration so that everyone can benefit and all individuals can be offered the best possible conditions to express themselves and develop;*

3. ***Responsibility***: *those who design and deploy the use of AI must proceed with responsibility and transparency;*

4. ***Impartiality:*** *do not create or act according to bias, thus safeguarding fairness and human dignity;*

5. ***Reliability***: *AI systems must be able to work reliably;*

6. ***Security and privacy:*** *AI systems must work securely and respect the privacy of users.*

These principles are fundamental elements of good innovation.

*Rome, 28 February 2020*

## 18.5 AI Ethics: The Abrahamic Commitment to the Rome Call for AI. Joint Declaration

**Preamble**

The development of Artificial intelligence (AI) is already changing lives across the world. As its power grows, this change will prove as great as the industrial revolution, changing the nature of work, the possibilities of manufacture, and opening possibilities to deal with humanity's great challenges. But all great changes bring the capacity for good and for harm. Whether the development of AI is truly beneficial to all is in our hands.

In February 2020, a group of organizations led by the Vatican signed in Rome a "Call for AI Ethics", seeking to foster an ethical approach to AI. The Rome Call sought to promote a sense of responsibility among governments, businesses and institutions to think carefully about the ethical development of AI technologies, serving human creativity and respecting human dignity.

The Rome Call recognizes that new technologies and vast changes in human experience do not come with innate ethical constraints. It is incumbent on those of good will, from all walks of life, to look to the future, laying out frameworks for the ethical use of new technologies that envision the vast changes and opportunities that they create. If such action waits until the potential of new technologies is realized, it is too late.

Religious leaders, ethicists and theologians have essential roles to play in this endeavor. The religions that they serve place the inherent dignity of the human being, as something given by God, at the heart of their concepts of social relations and development. Their ethical approaches are not utilitarian, but absolute: we seek human flourishing because we are commanded to by our creator. We seek the welfare of the most vulnerable because in each individual, however weak, there is the spirit of God. Such a shared commitment to the value and well-being of every individual was expressed in the Charter of the New Alliance of Virtues signed in December 2019.

The three religions of the Abrahamic family, sharing basic ethical values derived from their shared history and development, should particularly seek opportunities to work together for an ethical vision of a world changed by AI.

Mindful of the Rome Call and the duties incumbent upon us in a world in which 80 per cent of people declare a faith,

**We, representatives of the Chief Rabbinate of Israel's Commission for Interreligious Relations, the Pontifical Academy for Life and Abu Dhabi Forum for Peace, having deliberated on**

- The need for religious freedom in the promotion of human flourishing;
- The requirement for an ethical approach to the opportunities presented by AI;

- The complexities of AI and the challenges of international regulation;
- The vital contributions of religious leaders to a well-functioning polity and good governance;
- Our common concern for the most vulnerable in our societies and the requirement for equitable economic and technological development;
- The opportunities for inter-religious cooperation in the development of ethical frameworks for AI;

Acknowledged that technological development can have unintended consequences that lead to division and violence, and that this is especially true when development occurs without prior ethical consideration, frameworks, and a view to mutual human flourishing; and that our societies do not always affirm human dignity or promote peace.

Recognized and accepted that voluntary regulation is insufficient for the pace and scale of technological progress, and that new forms of international regulation must be encouraged to promote transparency and compliance with ethical principles.

Affirmed our role, as religious and political leaders and representatives, to build and uphold the principles of human dignity, freedom, "algorethics", protection for the vulnerable and the common good.

We are mindful of the fact that, in the past, religion and religious differences have been exploited to exacerbate conflict, violence, the loss of life and property, and to gain political privilege. Accordingly, we stand together as an alliance of virtue, in pledging to do our utmost to work for peace and global fraternity and to avoid future exploitation of religions and their adherents.

**Committed ourselves to promote "algorethics", namely an ethical framework for the development and use of AI, by means of the following principles:**
- Transparency: on principle, AI systems must be clear and explicable.

- Inclusion: the needs of all human beings must be taken into consideration so that everyone can benefit, and all individuals can be offered the best possible conditions to express themselves constructively.
- Responsibility: those who design and deploy the use of AI should do so in a responsible manner and with integrity.
- Impartiality: tools of AI should not be developed with bias, but should strive to be fair, honest, and protect human dignity.
- Reliability: every effort must be made to ensure that AI systems work reliably.
- Security and privacy: AI systems must ensure security and respect the privacy of users.

**We call upon all stakeholders, including but not limited to business, developers, governments, civil society, and academics to**
- Uphold the *algorethics* principles, enshrining them in company policies and developing international regulation to maintain them;
- for those in authority or positions of power to refrain from the development of AI technologies for partisan, national or individual gain;
- for governments to refrain from the weaponization of AI technologies, which threatens the peaceful co-existence of us all.

We urge those engaged in AI research and development to consider the ethical ramifications of any advancements, to collaborate with ethicists, and to adhere to the cautionary principle in areas of moral concern.

We propose the creation of an alliance between technology companies and leaders representing most of the religious and ethical traditions of the world.

We hereby commend this joint declaration as a companion to the "Rome Call for AI Ethics,"

*Vatican City, 10 January 2023; 17 Jumada II 1444 AH; 17 Tevet 5783*

# PART E


# SERVING HUMANS:
# FARMING, BUSINESS, FINANCE

<div align="center">

**19**

</div>


<div align="center">

**AGRI-DATA PROS AND CONS**

</div>


<div align="center">

*Adesuwa Ifedi, Nigeria / Stig Tanzmann, Germany[304]*

</div>


## 19.1 Agri-Data and Drones: Bringing Young Africans Back to Farming (Pro)

As communities in Africa start to rebound from the initial impacts of coronavirus pandemic[305] and look ahead, the importance of creating millions of jobs for the continent's booming youth population cannot be

---

[304] Adesuwa Ifedi (Pro), Nigeria, Senior Vice-President Heifer Africa at Heifer International. Stig Tanzmann (Con) from Bread for the World, Germany. Translated from German: *Mit Daten und Drohnen Afrikas Landwirtschaft retten? Machen digitale Technologien die Feldarbeit attraktiver und produktiver? Ein Pro und Kontra*, Welt-Sichten, 28.2.2022. https://www.welt-sichten.org/ artikel/39977. © Globethics Publications, 2023 | DOI: 10.58863/20.500.12424/ 4276053 | CC BY-NC-ND 4.0 International

[305] See African Arguments, https://africanarguments.org/category/covid-19-in-africa/

overstated. This will require concerted efforts across myriad sectors, but one that is brimming with vast and untapped resources is agriculture.

I hear a lot of dire warnings about the food situation in Africa, and there are indeed many reasons to be concerned. The pandemic disrupted supply chains, upending many farmers' precarious livelihoods that are already being hit hard by the climate crisis in new and unsettling ways. It's easy to think that we are on an unalterable path toward a crisis of food security in Africa, but I know that doesn't have to be the case.

The future of food on the continent is in the hands of African farmers, particularly our young farmers of whom we are in desperately short supply. About 60% of Africans are under 25 years old, but the average age of Africa's smallholder farmers is over 60. Too many young people view farming as exhausting work with antiquated tools for very low pay. You rarely hear under 30s say "I want to be a farmer" and many parents shudder at the thought of their children going into agriculture. Yet if we are going to secure a sustainable future of food in Africa, these things need to change. We need to show agriculture can be a major source of jobs for the future.

A report 2021 of Heifer International[306] reveals why youth are turning away from agriculture and highlights a major opportunity to evolve the sector and bring them back. The survey, which drew responses from some 30,000 young people across 11 African countries, found that only about 1 in 4 young farmers has access to the kinds of agricultural technologies that are transforming food production around the world. These "agritech" tools include things like digital sensors that monitor soil health and digital platforms that connect farmers with market opportunities, technical advice and high-quality inputs. The fact that many young African farmers lack access to these kinds of tools puts them at a severe disadvantage.

---

[306] The Future of Africa's Agriculture - An Assessment of the Role of Youth and Technology.   2021.   Heifer.   https://media.heifer.org/About_Us/Africa-Agriculture-Tech-2021.pdf

The report also found that with the appropriate financing, training and access to technologies, many more African youth would seriously consider pursuing a career in agriculture. The survey, which included focus groups with farmers and tech companies, revealed evidence of a rapidly growing cadre of agritech start-ups operated by creative, young people across the continent. By encouraging and supporting this new generation of innovators, we can boost access to labour-saving and transformative technologies for huge numbers of smallholder farmers.

These young, innovative entrepreneurs understand the farming struggles of their parents' generation. But they also believe farming can provide a promising future for their own generation. In Ethiopia, a group of young engineers is providing drone services[307] for analysing farm performance and a mobile application to help farmers detect crop disease. A start-up in Nigeria[308] is using machine learning to guide farmers from seed to market, helping them choose what to grow, how to grow it and where to sell it. And farmers from Senegal to Kenya can sign on to receive SMS alerts on important weather updates, market insights and farming advice.

With these kinds of tools and technologies, young African farmers could better manage, or even avoid, many of the challenges they reported in the survey, from climate shocks to crop pests and disease. These technologies have the power to make smallholder farmers much more productive, profitable and sustainable – and to make food production more exciting and attractive for a new generation.

Financing and supporting these companies are powerful ways of creating a virtuous circle in African agriculture: successful youth-led agritech companies lead to more successful young farmers and a more dynamic and profitable agricultural sector. Their success then creates an even bigger market opportunity for agriculture innovators that encourages more investors to get involved.

[307] Debo Engineering, https://deboengineering.com/

[308] Farmz2U's, https://www.farmz2u.com/solutions.html

But this future will not just magically appear. Governments and businesses alike need to invest in African farmers and encourage them to develop and adopt agritech innovations. The need for this investment grows more urgent every day. Our survey revealed that during the pandemic some 40% of agricultural organisations supporting smallholder farmers were forced to close at least temporarily, while more than a third lacked the capital they needed to recover.

The pandemic has exacerbated an already difficult situation for African farmers, with climate change looming as a bigger threat. But when I talk to the young people running African agritech firms, I detect no sense of gloom. These young people exude energy, ideas and optimism. They represent an entire generation that has the potential to transform African agriculture for the better. That gives me hope for the future. In a time filled with hardship, we should embrace their vision and make it our own.

## 19.2 Digitisation Serves Mainly Corporations (Cons)

Adesuwa Ifedi from the US aid organisation Heifer International argues that the introduction of digital technologies would make agriculture in Africa more productive and more attractive to young farmers. This in turn would improve food security on the continent. Ifedi cites the use of drones to analyse soil quality and digital platforms for technical advice as examples. Once again, a technological revolution is supposed to solve existing complex socio-economic distribution problems.

The parallels to the Green Revolution in the 1960s and 1970s, especially in Asia, are obvious. There, agricultural production was to be improved with the help of chemical fertilisers and pesticides, industrial seeds, and increased mechanisation and irrigation. This has caused major ecological and social problems in many countries. In Africa, this approach was not successful at all. Digitization will not change that either, because most digital technologies aim to optimise the approach of the

Green Revolution, but not to offer an alternative. The digitization of agriculture in Africa is being driven by the same states, corporations and philanthropic foundations as the Green Revolution. Bread for the World, together with other organisations, took a very critical look at the Alliance for a Green Revolution in Africa (AGRA).

Studies have shown that AGRA could not keep its full-bodied promises. For this reason alone, scepticism about the promises of digitization is appropriate. Not surprisingly, Heifer International presented its study on youth and agricultural digitization in Africa, mentioned by Adesuwa Ifedi, at the African Green Revolution Forum 2021 in Nairobi, among others. The Forum has its agenda in its name: it represents the "business as usual" in agriculture that has brought African agriculture to the difficult position it is in now.

The digitization of "business as usual" will not bring about a turnaround. On the contrary: it threatens to further exacerbate the agricultural crisis and drive even more people out of agriculture. The Green Revolution is dominated by chemical and agricultural companies such as Bayer, Syngenta and Yara, whose business model is to tie farmers to technical solutions, such as the use of certain seeds and artificial fertilisers, pesticides and machines. Digitization aims to further increase this dependency. This is particularly evident in the case of seeds. The companies promise that the sequencing and digitization of seeds will significantly increase yields. In principle, however, their main aim is to obtain patents on seeds and to introduce genetic engineering in agriculture.

Farmers are left with no choice but to use corporate seed because they can't prove they're not infringing a patent (How much do they have to prove it? To whom? And why is it hard to prove if they don't have patented seed use?). When it comes to chemical fertilisers and pesticides, too, the concerns of the corporations are to bind farmers to their products. Instead of daring to phase out pesticides and chemical fertilisers and changing the system, as in agro-ecology, the existing system is to be optimised with drones and total surveillance of farmers, and the prof-

its of the corporations are to be increased. Digitization makes it possible to transfer relevant operating data such as soil quality, water availability, sowing times, fertilisation intensity and harvest times to the large corporations with the help of the help. This takes away the sovereignty of the farmers to a completely new depth, alienates them from their land (okay?) and makes them completely interchangeable. If there is no radical trend reversal in agriculture, then digitization will serve the corporations and not the farmers in Africa.

# 20

# INTERNET OF THINGS (IOT)
# A LATIN AMERICAN PERSPECTIVE

*Rosa Delgado* [309]

## 20.1 Introduction

With a population of around 640 million, the Latin America region represents a large and varied market for Telecommunications companies (Telcos) operating on the Internet of Things (IoT) sector that has great worldwide potential. The number of IoT devices is growing at an exponential rate. The region expects to have 500 million new IoT connections

---

between 2020 and 2025. The smart home sector is estimated to generate more than 100 million new connections by 2025.

Early adopters have already understood the great potential of the IoT technology to improve health, security, agriculture, mining, public safety, and reducing energy use. However, overall adoption across worldwide Governments and private sector is still very slow. Connecting everything with IoT enables joining the global connected world, maximizing IoT opportunities, but, at the same time, entailing risks. Telcos have already made a shift in focus, from communications as a business based in connectivity to one in which cross-sector cooperation plays an essential role in the IoT sector. The industry has already shifted to put users at the center of all Telco operations, companies seek to better understand and connect with their customers, delivering the services they need. While technologies like smart health, smart agriculture, smart home, smart city, and smart ports are all gaining ground around the world, IoT in Latin America is finally moving in the right direction.

From an ethical perspective, Latin American social structures face challenges, mainly because IoT research and systems are designed outside the region, mainly in Europe, the USA, and Asia. IoT systems tools consist of entities, developers, regulatory institutions, technical components, as well as biases. Moreover, misunderstandings of algorithms, technical capabilities, and biases within data sets have important ethical implications in Latin America and the rest of the world.

## 20.2 Background

In the last two decades, Latin America has been one of the most successful emerging regions worldwide. However, the adoption of trending technologies like Artificial Intelligence (AI), IoT, Blockchain, and Internet Protocol version6 (IPv6) has been rather slow. The emergence of new business solutions, low-cost devices, and high-quality connectivity has been instrumental in driving IoT adoption, helping to connect devic-

es, cars, machines, and wearables. Latin American nations have grown to serve supply chain processes, but now are also supporting other industries like security, healthcare, education, agriculture, public services, and transportation. Today, the industry is taking a hard look at how to get the most out of its IoT, AI, and robotic solutions. Companies are looking at some early-stage use cases for IoT solutions to add automation/robotics using the IoT cloud, and to respect environmental issues.

## 20.3 Main Issues for this Paper

### 20.3.1 Internet of Things (IoT)

IoT consists of physical objects embedded with actuators or sensors connected to a network. It provides numerous opportunities to reduce costs and improve public and private services. Telcos often provide support for IoT devices and solutions, including smart IoT solutions that protect against cyberattacks. Early adopters in the public and private sector have already demonstrated the potential of this technology with services that improve security, telework, education, tourism, and worker's health. In Latin America, overall IoT adoption is still low; however, widespread adoption of IoT by governments will undoubtedly spur commercial adoption.

### 20.3.2 IoT Security

IoT security is a solution that protects networks and connected devices, equips devices to identify, resist and recover from cyberattacks, as well as eradicates vulnerabilities. Smart solutions include embedded device authentication and management, secure connections, and data loss prevention. Following security best practices, using strong passwords, and protecting smart devices with dedicated security tools are essential to maintain strong IoT security. Today, IoT technologies are exploding however, many users are still unaware of the risks that smart devices can introduce to their home or business networks.

### 20.3.3 IoT Healthcare

The opportunities of the Internet of Medical Things (IoMT), which is often referred to as IoT in healthcare, should be explored in Latin America. IoMT encompasses medical devices, infrastructure, and applications capable of collecting and transferring medical data over a network without human intervention. Wearable medical devices that track the health conditions of patients in hospitals or homes, and IoMT sensors that help monitor stocks. These "things" have completely transformed the healthcare industry during the Covid lockdown, reducing healthcare costs and time, and improving patient safety. IoMT medical institutions are transforming healthcare in new ways with increased specialization of medical staff and the use of new medical equipment and devices. The IoMT market is expected to grow 40% by 2026.

## 20.4 Business Automation Workflows and Processes

Effective access to devices and medical devices requires technologies such as notification servers, IoT Hub, and/or Communications Platform as a Service (CPaaS) that transform collected data into actions for the customer. IoT solutions offer opportunities when used in combination with technologies like CpaaS, innovative network solutions, collaboration tools, etc. These technologies provide the intelligence, through IoT solutions, that transform collected data into action on the connected device, enabling automated workflows. The combination of CPaaS and the notification server evaluates the received information and sends a message to predefined contacts, ie doctors, nurses, or relatives. These opportunities will allow finding a critical use of IoT in many industries in Latin America, such as security, health, machines, agriculture, autonomous cars, measurement tools, and wearable devices.

## 20.5 Security Market Analysis

Covid-19 has played an important role in the rise of remote working (or telecommuting from home), which has been increasing the demand for IoT security solutions. Today, employees in the region telecommute from home and rely on remote-based monitoring of the workforce, through IoT solutions that provide enhanced security, real-time tracking, meet transparency requirements, and government compliance. Cybersecurity of public offices, remote companies, and the remote workforce played an important role in serving their clients during Covid-19 but, it also provides a huge market opportunity for IoT security solutions. The number of IoT security regulations and phishing/malware threats in the region has been increasing, which is driving the growth of the security solutions market in Latin America. Wearable IoT devices can have dire consequences if hacked. IoT security is a serious concern due to data privacy requirements, as well as the potential consequences of a hack, with lives at stake. All industries are increasingly incorporating smart sensors and programmable controller devices into their facilities. However, these may lack sufficient protection and be particularly at risk from IoT hacking, i.e., healthcare, self-driving cars, finances, etc.

Technology companies continue to develop and improve IoT devices, but one challenge remains, how to keep devices safe and secure. IoT device manufacturers are to some extent responsible for keeping devices secure. However, users are increasingly required to take steps to protect their own devices. Today, IoT security threats are often not given proper importance and are eventually ignored by public offices and enterprises.

## 20.6 Ethical IoT Implications

Awareness of IoT ethical issues in Latin America depends on the classification according to the education level of the population, while offering more connected and intelligent services. The fact that these

algorithms have not been designed by Latin American developers increases the mistrust of some sectors of the society for this dependence on the outside world, and the social impact of a technology perceived as 'non-indigenous' should not be underestimated.

Today, developers design connected devices that generate a lot of information that is essential for different purposes, such as: communicating the social benefits of the device, monitoring development, and stimulating the economy. Developers need also to maximize the benefits of the device and keep the risks low. If IoT security is not adequate, there will always be an option for hackers to hack those devices. Therefore, protecting individual privacy is critical, developers must implement privacy by design, creating ethical and legal constraints on the devices themselves. The side effects of IoT technology, which is perceived as "non-indigenous" by part of the population, need to be further investigated. To be successful and avoid mistrust of some communities towards "foreign" technologies, IoT systems must focus on more locally specific benefits and concerns for users in the region.

## 20.7 Conclusions and Recommendations

The ethical implications of IoT should focus on the challenges relevant to the region, as follows:

The IoT industry needs to better protect its users. Vendors must add security features to be built into the devices themselves or through aftermarket hardware and software.

IoT solution designers must consider the privacy implications of their devices as the number of IoT devices continues to grow. The scope and number of privacy violations will drive more regulatory policies that will hinder solution innovation.

Regional IoT concerns need to consider the poor understanding of IoT algorithms by developers and users; low access to education, which limits access to information about IoT systems strengths and risks; the

social impact of technology developed by entities outside the Latin American region, and the little importance of the voices of ethnic minorities and the patterns of culturally discriminated communities very different from the majority.

There is also a need to better understand the social and political impact of the lack of knowledge of developers and users about the actual functioning of IoT algorithms and the biases within data sets that have important ethical implications in Latin America. and the world.

The cyber community should include Latin American voices in IoT/AI discussions to plan collaborative and interdisciplinary programs to increase awareness among the population about IoT literacy, benefits, risks, misuses, misunderstandings, and equity issues. Additional discussions on unethical uses and harmful impacts of IoT will lead to social inequities and biases in algorithmic assumptions that could lead to exclusion and discrimination.

# HUAWEI:
# DIGITAL SELF-DETERMINATION MUST BE FOUGHT FOR AND SECURED ANEW EVERY DAY

*Interview with Andy Wang, CEO Huawei Switzerland*

*Self-determination is in the hands of our customers. We never take this away from them in any way. Security in cyberspace is a complex and multifaceted topic. No one can guarantee 100 percent security at all times. Anyone who does this is simply not honest. However, what we at Huawei can guarantee is that our customers constantly and repeatedly put our systems through their paces, that they can test and externally audit us.* *Andy Wang*

**Andy Wang[310], security concerns and allegations of data espionage are being raised against Huawei in various countries, because**

---

**Huawei, as a Chinese company, is not supposed to be independent of Chinese government intervention. What do you think?**

Sure, we are a company founded in China. But we have been operating as an international company in over 170 countries for over 20 years. In Europe alone, Huawei employs around 14,000 people and operates 25 research centres. We never had security problems. The allegations made are politically motivated and simply not tenable. Because nobody is as transparent as we are. And we can prove that.

**Give us proof.**

Just imagine: John Suffolk, our Group-wide Global Cyber Security Officer, and his people check individual parts and components before they are even released for installation in our systems. And anything that doesn't stand up to this strict control will no longer be developed further. Incidentally, John Suffolk previously worked for the English government, ultimately as Her Majesty's Government CIO, i.e. as Chief Information Officer, and also as Chief Information Security Officer.

**What does transparency mean for Huawei?**

We have to put everything on the table. You can imagine that a single instance of us being shown to be improperly handling data and our promises would be catastrophic, if not fatal, for us. We are probably the only technology partner who even submits our source code to our customers for review. In our Huawei Cyber Security Transparency centre in Brussels, our customers can check the source code of the programs, the operating system or applications that are in the network components from Huawei themselves. In addition, we clearly show how we comply with the certifications and standards that should apply to such systems. This transparency is exceptional. We depend on our customers being able to trust our software.

**And use such efforts to cement confidence in Huawei?**

Absolutely. We work closely with European industry associations - for example in the development of safety standards - and also allow ourselves to be checked and evaluated by independent safety and testing institutes. Which other technology group in our area allows external auditors to examine the processes? Our customers trust Huawei, and they say publicly that they believe the security concerns about our technology are unfounded.

**Okay, technology is one thing, but what about protecting the data on Huawei components?**

Firstly, we implement the standards of the European General Data Protection Regulation throughout the group. Secondly, it is important to know that as suppliers of the hardware for 5G networks, we know nothing about data traffic. The operators of the networks have full control over it. Huawei has no access: only the operator of the network decides how data is treated in the network.

**But can you guarantee that Switzerland can guarantee self-determined handling and the security of its communication networks with Huawei components?**

Self-determination is in the hands of our customers. We never take this away from them in any way. Security in cyberspace is a complex and multifaceted topic. No one can guarantee 100 percent security at all times. Anyone who does this is simply not honest. But what we at Huawei can guarantee is that our customers constantly and thoroughly test our systems, that they can test and externally audit us. Because cybersecurity is always in flux.

**What is certain today may no longer be so tomorrow.**

That is why we in the company are constantly striving to follow the latest developments and to constantly adapt the safety components with our research and development. Because nobody should lose their digital independence because of the systems we sell to our customers. Digital self-determination must be fought for and secured anew every day. We

are an excellent partner for this because we cannot afford anything other than transparency and openness towards our customers. Every other competitor does not have to pass as high standards and critical reviews as we do. We don't have anything against that either, just that we're being denied credibility for purely political reasons without being checked and without having really looked behind the scenes. That is also unfair to our almost 200,000 employees. Because our company is owned by them, 99.16 percent. Our founder Ren Zhengfei still owns 0.84 percent of the shares.

**Distrust of Huawei is fuelled primarily by the United States. Do you think that's fair, and doesn't it also involve a technology war?**

Above all, we want to sell our products and not a way of life or a political world view. And because we are simply business-oriented, we believe that our transparency and openness should actually decide. This is also our motivation to let ourselves be illuminated. Many of our employees say they have never worked in a company that is audited so often. This also applies to the security area - not only in relation to software or components, but also to the supply chains and the personnel.

**How is Huawei accepted in Switzerland?**

You have to ask our customers that. From our point of view, we are the reliable and leading technology partner for the telco infrastructure in this country. With our technological experience and reliability, we are part of Swiss society and make a transparent contribution to the benefit of the country and its economy. This is important to us. As an innovative incubator, we also want to initiate agile partnerships in research and thus support the digital transformation in Switzerland. We meet all Swiss security requirements and are willing to have this checked transparently by independent institutions. We are probably much more Swiss than many want to admit.

<div align="center">**22**</div>

<div align="center">
# BLOCKCHAIN INVESTOR, BLOCKCHAIN INVESTEE: GOVERNANCE IS YOUR ONLY ALLY! YOU GAIN TRUST.
</div>

*Eelco Fiole, Switzerland[311]*

## 22.1 Blockchain and Crypto: Opportunities and Risks are Kicking

The blockchain- and crypto space is an exciting space to be in, both as an investor and as an investee company, you are right at the front of the new world! However, it also bears huge risks as 2022 has shown, risks that are partly well, but mostly not very well understood: smart-contract risk, oracle risk, custodial risk, scaling risk, regulatory risk, and

---

of course governance risk. And these and other better-known risks materialize too often: massive destruction of financial and reputational value, especially visible in the FTX-case… How to go about blockchain, just forget about it? Wait until regulations come in and then rely on those? While it is true massive value was destroyed in the blockchain- and crypto space over the last period of time, it is also true that the blockchain technology still has huge benefits. There are many real-world examples where blockchain technology replaces outdated information technology, yielding massive benefits: faster, cheaper, better. Switzerland-based aidonic.io is just one for-profit example where humanitarian aid is tokenized based on blockchain technology yielding a two to three times higher payout ratio for the beneficiaries. Also, in Web3 do we see first new business models, new economies, coming online, where creators can be rewarded for their efforts, without being located in major blockchain hubs, and interact with their communities, devprotocol.xyz of Japan being a good example of that. Of course, there are many more great initiatives under way, the blockchain developers are building at lightning speed, unabated by FTX-, Terra/Luna- and 3 Arrows Capital-like scandals. Indeed, for early investors the timing for investing in blockchain could not be better. Low valuations, a healthy dose of humility and massive opportunity. However, while opportunities are abundant, risks have not come down.

## 22.2 The View from the Blockchain and Crypto Investor

One of the major things that became clear in the FTX-case was that very reputable investment managers, investing with other people's money, basically did not do any proper initial due diligence, a huge omission. One of the things that due diligence would have uncovered would have been that there was no board, not to speak of a quality board, that would have kept oversight on FTX on behalf of investors. There was a clear

lack of fiduciary responsibility at FTX: willingly or unwillingly, they did not understand what it means if a fintech is dealing with other people's money, and that the interest of the investor comes before their own. Further, the very large number of legal entities created increased complexity massively and made the structure ungovernable, one can only wonder what the thinking behind it was…

But even if due diligence would have been performed, what happens after that? When is the next round of due diligence, a year later? In the blockchain- and crypto world, things are not only too complex for regular investors to comprehend, but things also go at lightning speed, crossing borders as a default way of working. Doing due diligence once a year, maybe twice a year is not enough. One would need to do ongoing due diligence, which of course is unpractical. However, the proxy for ongoing due diligence on investors' behalf is called governance, which has been so underutilised. A high-quality board, a board with the up-to-date knowledge, experience, matching maturity and attitude, should be there to represent the interest of investors, making sure that the investment in a blockchain- or crypto-entity is well-managed. This is irrespective of any regulation that may or may not come in and that may or may not fit the situation. Regulation and policy are good, but governance makes it so much better. Indeed, as a blockchain or crypto investor, governance, people that knowledgeably closely oversee the investment on your behalf, is your only ally!

## 22.3 The View from the Blockchain Investee Entity

Investee blockchain or crypto entities are exciting. The technology enables things that ten years ago were unimaginable and based on some years of experience and uptodate knowledge of the technology a bright future is almost certain… Or not? The world and especially capital markets are complex, the uncertainties are actually quite high. Policy and regulations are expected to come in and right or wrong, they are mostly

expected to be a hindrance to reach goals. One needs to build a reporting bureaucracy, needs to spend money on non-productive functions, and needs to spend time away from non-core development work. Investors are interested in supplying capital, but what level of comfort can a blockchain entity give them— is a PhD in computer science enough? If investee companies are serious about building a profitable sustainable business, they will need to provide comfort to investors that their investment is well-guarded, that strategy is set and executed and no unnecessary risks are taken, that robust operating structures are built. Serious blockchain entities ensure that they get their strong governance in place, instead of looking for the jurisdiction with the lowest level of regulation. Governance, both on a corporate as well as sometimes on a protocol level, ensures that experience, a sense of judgment, broader knowledge and a fiduciary attitude flow into company and protocol development. It is underutilised but positive symbiotic governance that will help investee companies to project trust to investors, on an ongoing basis, and to interact with policy and regulation if and when this comes, to show that certain principles indeed have already been implemented. It is the right governance that helps understand capital market intricacies and recognize tangible and not so tangible risks and ways to avoid or mitigate them, and that supports purpose, strategy, structure, antifragility and culture towards intention and output. Innovation spirit is good, but governance makes it so much better. Indeed, as a blockchain or crypto investee company, governance, people that innovators can rely on for input and representation, is your only ally!

In 2022, lacking governance in the blockchain and crypto space has cost investors and investee companies hundreds of billions of dollars, in negative market impact, opportunity cost and direct losses. Blockchain and crypto will survive, but how is 2023 going to look like for your investment or for your blockchain or crypto entity?

Remember, governance is your only ally. And with governance you build trust.

# PART F

# SERVING HUMANS: SECURITY, WAR, PEACE

# 23

# THE CYBER SECURITY
# TRANSPARENCY CENTRE IN BRUSSELS

*Christoph Hugenschmidt, Switzerland*[312]

"You will be disappointed," Wang Liangchen warns me. In fact: what's interesting about a few almost empty, cramped offices in a rather nondescript building in Brussels? After all, the entrances to the offices on the first floor are well secured, and the alarm triggered accidentally is really loud. In the windowless rooms themselves there is a desk, one or two office chairs, one or two screens and a PC, sometimes even just an internet connection.

What Wang Liangchen showed me in Brussels in January 2022 is actually completely unspectacular. But not what happens in these rooms.

Because in the Brussels Huawei Cyber Security Transparency centre, more precisely in these windowless, narrow rooms, customers and partners from the security ecosystem can view and check the source code. Namely the source code of the programs that are in Huawei's network, data centre and other components. However, the source code itself is not stored on the PCs and servers in Brussels. It is – in all likelihood extremely well secured – in a data centre in Shenzhen and can be viewed from Brussels. And, as Wang assures, for as long as the customer needs. There is no time limit.

Wang Liangchen is a cybersecurity specialist. For nine years he has been working for the Chinese technology group Huawei, one of the world's three largest providers of infrastructure for telecommunications and, in particular, fifth-generation mobile networks. There, Wang was responsible, among other things, for so-called "penetration tests" as part of the company's independent security laboratory. They are more important than ever in the IT and telecom industry; Because this is how manufacturers and operators of computer and network infrastructures try to find out whether their systems are protected against attacks from outside. And that is exactly what cybersecurity is all about: Can outsiders penetrate my systems? Can they steal or manipulate data? Can they cripple systems or cause them to do something other than what they were built to do?

The fact that customers or their security teams can view the source code of an operating system or applications of commercially marketed network products is exceptional. It indicates that the Chinese telecom equipment supplier trusts its own software and, conversely, is dependent on customers trusting their software. In the background of Huawei's transparency offensive is the accusation, expressed above all by the USA, that Huawei is building back doors into its software on orders from the Chinese secret services or the army. These backdoors would allow Chinese intelligence agencies to penetrate networks using Huawei machines.

Another suspicion, expressed primarily by the United States, is that Huawei is secretly installing additional chips in its devices. These would monitor or manipulate the data traffic. That's why Huawei even allows customers to physically analyse its hardware. However, if you want to do this, you have to travel to Shenzhen in China. There is no Huawei hardware in Brussels.

## 23.1 "We love processes"

Yoann Klein is my second interviewee at Huawei's Cyber Security Transparency centre in Brussels. The Frenchman can look back on a long career in the European IT security industry. He worked for the French high-tech and armaments group Thales and before that for Airbus. According to Klein, Huawei has come a long way in developing cybersecurity. While the focus before 2005 was still on product security, company founder Ren Zhengfei defined cybersecurity as an absolute priority at every level in an open letter to employees at the end of 2011. The open letter from the charismatic industry veteran was probably intended as a wake-up call to employees and the outside world to really and always take the issue of security seriously. Since then, software development processes at Huawei have been and continue to be redefined in terms of security, and the entire organisation of the giant corporation has been aligned with security, from management to the individual in sales.

The European market is enormously important for Huawei, emphasises Klein in an interview. On the one hand, it is large with well over 400 million inhabitants. The level of industrialization is high, and investments in telecommunications are correspondingly massive. In addition, the EU defines the standards worldwide when it comes to data protection. Huawei itself implements the standards of the European General Data Protection Regulation across the group. In addition to the centre in Brussels, Huawei operates three other such centres in Europe

and employs a total of 13,000 people in the EU area, 2,400 of them in research and development.

Klein takes the time to explain in detail how the Chinese technology group intends to ensure the security of its network products. What he explains is reminiscent of the process obsession of US technology companies. "It is true. We love processes, certifications and standards. We also learned from IBM," says Klein. In fact, in the 1990s, Huawei sought advice from the then leading IT group, among others, and was inspired by its methods and best practices for its product development processes.

## 23.2 Traceability, Tests, Standards

Rules and processes are of little use if they are not enforced. Every employee at Huawei is regularly tested for security awareness. Those who fail must be detained. "I've never had to take exams as often in any company as I do at Huawei," says security specialist Klein.

Another important principle in all security systems is traceability. Which customer uses which version of the software in which devices? Which software components from open source libraries are used where and in which version? Huawei can say within an hour exactly which software version is being used by which customer. Hardware can be identified within one tag if it is affected by a gap.

A current case, which has occupied the IT industry worldwide, shows how important sophisticated vulnerability management that customers can reliably understand is. In December 2021 it became known that the software components Log4J contained a security gap. Log4J is used to record logins in certain software systems. Over the years, it has become a de facto standard in many places. After the vulnerability in the open source component became known, it was up to the manufacturers to inform their customers immediately and patch the software. Not all manufacturers were able to do this. In some cases, attackers managed to exploit the vulnerability.

Huawei is not only obsessed with processes and rules, but also with standards and recognized regulations. The group participates in the development of safety standards in international bodies such as the industry association 3GPP. The GSMA, together with 3GPP, has developed the security program NESAS (Network Equipment Security Assurance Scheme). People in Brussels are audibly proud that Huawei was the first telecoms supplier to have its first network products tested according to NESAS. The Chinese technology giant also attaches great importance to having its processes checked and evaluated externally. Huawei had the software of the central unit of a 5G network, the Unified Distributed Gateway, checked for quality by the independent German security company ERNW. The inspection took place – how could it be otherwise – in the windowless, well-secured offices in Brussels mentioned at the beginning. ERNW checked the quality of the source code, the handling of open source components and compliance with security rules in the programming process.

## 23.3 Additionally: Independent Safety Laboratories

Altruism: All rules, standards and processes are useless if they are not followed. For example, because a development department is under time or cost pressures, or if anyone in Huawei's massive workforce wants to take the path of least resistance. That is why there is an additional level within the organisation of the tech company: the independent cyber security laboratories (ICSL). These have their own budget, are organisationally located outside of the development departments and therefore have different bosses. Even the goals that are given to the employees of the ICSL are different. Because they should find as many gaps and errors as possible and communicate them consistently. The independent labs have the right to stop the development of a product or a new software version and send a product back to Huawei's software

factories. Around 200 security specialists tap Huawei's software for gaps.

It is one thing to ensure the security of network components, such as those built by Huawei for 5G networks, among other things. Demonstrating to customers and the public that devices and software are secure is another thing. For each product, one can prove which methods and systems it was tested with, says Wang. He doesn't just say it, he lets me take a look at a huge "workbench". Information on Huawei products and their individual components and their sub-components is collected in the database. When were they tested, by whom, using what methods, and what comments did the testers and their supervisors make? Has further development been approved or stopped? Wang clicks through the database and ends up at a product whose further development has been stopped. Apparently it's a smartphone. "The customer must be able to deactivate the microphone himself," it says. According to the database, the man who passed the death sentence on the product or component is John Suffolk, who is responsible for security and transparency at Huawei group-wide.

## 23.4 Who is in Control of the Data?

But what about remote maintenance? Who should prevent Huawei from remotely accessing network infrastructures and tapping data? Klein emphasises that a distinction must be made between the infrastructure supplier and the network operator: "As suppliers of the hardware for 5G networks, we know nothing about data traffic. The operators of the networks have full control over what we do with the equipment." Huawei has no access to the networks of the network operators (e.g. Sunrise) if they do not consciously allow and monitor it. In addition, Huawei fully implements the European Data Protection Regulation (GDPR). Huawei maintains, repairs and disposes of network equipment in Europe and has a production facility in Hungary.

You can feel the urgency of security when speaking to the experts at Huawei's Brussels centre. "Cybersecurity is the top priority at Huawei. That has to be the case if we want to survive as a company," says Klein more than once. He has never seen a company that invests so much in transparency.

# 24

# PROOF OF SECURITY:
# WHY SWITZERLAND NEEDS
# A NATIONAL CYBERSECURITY
# TEST INSTITUTE

*Franz Grüter, Switzerland*

The need for a Swiss institution that focuses on research, testing and promotion of cybersecurity in the style of the Swiss Federal Laboratories for Materials Testing and Research (EMPA) is obvious. The awareness of Swiss society for the risks associated with the use of networked devices must be promoted. The National Test Institute for Cyber Security must make an essential contribution to this and help to develop solutions. The security and trustworthiness of cyber-physical components, which are used by millions in business and society, are a political concern. A Swiss ecosystem for innovation and entrepreneurship in the field of digital security paves the way for national cyber autonomy.[313]

---

[313] Original text published in German: Franz Grüter, *Nachweis von Sicherheit als zentralem Erfolgsfaktor. Warum die Schweiz eine 'Cyber EMPA' braucht*, in

The advancing digitization of industry, infrastructure, authorities, the police and the army is leading to a constantly growing risk, which is due, among other things, to the use of untrustworthy digital products. The risk to Switzerland's security and independence is particularly serious when insecure cyber-physical components are used in the operation of critical infrastructures and by authorities such as the police and army. In addition, a risk that is unacceptable for Switzerland arises when insecure components are used by the millions in business and society (e.g. network routers, messaging, collaboration and video platforms). Today there is insufficient capacity in Switzerland to test such components. Instead, trust is placed in the purchased hardware and software elements having an adequate level of security.

## 24.1 The Digital Exchange of Information Cannot be Slowed Down

The commercial opportunities of growing digitization are enormous. According to Equinix, digitally advanced companies connect with ten times more partners than digitally less advanced ones. As a result, you are ahead of them in digital transformation and can benefit from the best performance from a rich pool of service providers and partners.[314] Compound annual growth in network bandwidth is expected to be 44 percent through 2024, four times the overall growth of 2020 and more than 21,485 terabits per second, or about 85 zetabytes per year. 1 zetabyte corresponds to around 563 quadrillion pages of text.[315]

---

[314] See Global Interconnection Index, Vol 5, published by Equinix, www.de.equinix.ch.

[315] 1 megabite is 500 pages of text.

It is clear: Our current economy and society live the digital exchange of information. Irrespective of their intentions, private and business people have constant access to a growing amount of data, knowledge and communication channels.

## 24.2 Cyber Attacks are Known - but Underestimated

The networked society is becoming more and more complex with the increasing number of new interactions between people, machines and services and the associated feedback processes. Recently, there has been an increase in the number of those who perceive a clear and present threat of a large and growing proliferation of cyber attacks and incidents worldwide, across diverse industries and environments.

The risks that can emanate from malicious software have certainly moved closer to general awareness, and there are already a variety of efforts to identify or prevent and ward off cyber attacks. In the eyes of many experts, however, an even less protected area lies in the hardware and firmware. Cybersecurity is often wrongly limited to software and network security; the integrity and security of the hardware and its components are rarely considered to the extent necessary. If digital products with security defects find their way onto the market, vulnerabilities can have an impact for decades. Permanently installed devices in home and industrial controls are also affected. Compromised hardware and firmware often render all other security measures ineffective.[316] Regular maintenance and the checking of modifications are often not adequately ensured.

The increasing complexity of the software running on connected devices presents an additional risk. With increasing complexity, the number of vulnerabilities contained exceeds around 1 to 100 vulnerabilities

---

[316] Costs for attacks on the supply chain are often very cheap. See Andy Greenberg, *Planting Tiny Spy Chips in Hardware Can Cost as Little a USD200,* in IT Equipement, 210.10.2019. www.wired.com.

per 2000 lines of program code. It is highly likely that every product that can be networked offers various contact options for remote attackers. With the predicted increase in networked devices (IoT in industry and households), there is a risk of a serious loss of control.

In critical infrastructures such as energy, transport, health and safety, ignorance about the security level of the products used can lead to wide-spread threats. The use of insufficiently secured devices in private households can also lead to serious damage to critical infrastructures. Targeted cyber attacks on supply facilities in various European countries clearly show the damage potential: in addition to the primary purpose of destruction, the population is unsettled and at least indirectly harmed. Effective protection against such socially relevant threats is virtually non-existent today and the barrier to entry for this type of attack is un-necessarily low. The former German Federal Minister of the Interior, Seehofer, rightly says: "Cyber security is the prerequisite for digitization to succeed."[317]

## 24.3 Switzerland must also Follow the Path of Certification and Testing

The dangers of nuclear engineering, genetic engineering and other mechanical, chemical and biological processes are relatively manageable and well known. Against the initial resistance of the industry, verifiable and regulated standards have prevailed in the long term. In all critical industrial sectors, quality tests by independent bodies are now an integral part of product approval.

However, apart from data protection, the ICT sector has hardly any binding and legally valid standards that regulate the security and integri-ty of the products used at the regulatory level. There are neither legally

---

[317] In: Meldung IT und Digitalpolitik, 8.9.2021, www.bmi.bund.de

binding minimum obligations nor legally anchored product liability for software.[318]

The need for a Swiss institution that focuses on research, testing and promotion of cybersecurity in the style of the Swiss Federal Laboratories for Materials Testing and Research (EMPA) is obvious.

Various other countries are already developing regulations to prevent cybercrime. In the European Union in particular, the "Cybersecurity Act" is openly debating the testing and certification of IT components and systems. Necessary bilateral agreements for the mutual recognition of certifications and examinations will have to be developed in the future.

The ability to independently and effectively test digital products—including reverse engineering chips and firmware—will become increasingly important in the global economic context in the near future. Due to the digitization of everyday and critical functions, the need for these skills will increase in industry as well as in authorities, the police and the army. A nationally mandated test institute for cyber security should significantly reduce Switzerland's dependence on external partners.

## 24.4 Swiss 'Cyber-Empa' could even Play an International Role

With ETH's top rankings in computer science, its neutrality, its representatives of internationally recognized organisations, legal certainty and finally little protectionist interests due to the lack of an overpowering hardware and software industry of its own, Switzerland offers all the prerequisites for an internationally recognized testing institute for networked devices under the developing a cybersecurity perspective.

At the request of the canton of Zug, a concept was presented in July 2020 to evaluate the need for and the structure of a national testing insti-

---

[318] For autonomous car driving, the connection with 5G is a precondition.

tute. In December 2020, the Canton of Zug initiated the project to create a National Test Institute for Cyber Security (NTC) with start-up funding in order to build up initial testing capacities. On behalf of companies and public organisations, the NTC is to test networked components for their cyber security. This includes hardware and software components, regardless of manufacturer and geographic origin.

## 24.5 Transparency and Security for the Digital Data Exchange of the Future

The NTC should therefore guarantee the independent and in-depth examination of the security and trustworthiness of cyber-physical components and thus make a strategic contribution to maintaining the security and independence of Switzerland. In the canton of Zug, the NTC is a test and research infrastructure that works with universities, private providers of safety tests and international test centres. Priority should be given to test orders related to critical infrastructures, authorities such as the police and the army, as well as tests of networked components that are widespread in everyday life.

The test institute remains market-neutral and independent. No tests are therefore carried out with the purpose of increasing sales promotion. Financial participations by product manufacturers and service providers are also excluded. Consequently, there is no direct influence on the market for hardware and software components. Any influence of persons and organisations on self-determination and objectivity is completely prohibited.

The safety checks carried out by the NTC are based on international recommendations and standards, which are supplemented accordingly in the case of more in-depth analyses. The NTC draws on the numerous existing competencies from the private sector, research and educational institutions in Germany and abroad in order to guarantee examinations with the appropriate competency.

## 24.6 National Cyber Autonomy is the Goal

The awareness of Swiss society for the risks associated with the use of networked devices must be promoted. The National Test Institute for Cyber Security must make an essential contribution to this and help to develop solutions. The NTC should become the central point of contact for the Swiss test request with a focus on freedom from vulnerabilities in the information and cyber-physical sense. The success of this order is based on public-private sponsorship. To do this, it is imperative that all relevant actors come together regularly at one table. All representatives of critical infrastructures, namely energy, transport, health, water, insurance, security, even media and food, from the public and private economy are required to protect Switzerland as a business location and the functioning of Swiss society against cyber attacks with the support of the NTC protection. Irrespective of whether cyber attacks are carried out for espionage, sabotage, manipulation, disinformation or criminal intent: the security and trustworthiness of cyber-physical components, which are used by the millions in business and society, are a political concern.[319] A Swiss ecosystem for innovation and Entrepreneurship in digital security paves the way for national cyber autonomy.

---

[319] Die Sicherheitspolitik der Schweiz. Bericht des Bundesrates vom 24. Nov 2021, BB 2021 2895, www.fedlex.admin.ch.

# 25

# QUANTUM TECHNOLOGIES: NEXT ARMS RACE CHINA-USA

*Matthias Sander*[320]

*China is researching quantum technologies for the military at a fast pace. The US and NATO are concerned today's encryption technology could suddenly become obsolete one day, and Quantum sensors could detect stealth aircraft. Washington wants to prevent China from being the first to acquire these capabilities.*

---

## 25.1 Unhackable Quantum Connection?

When China launched a spacecraft in Inner Mongolia in August 2016, a series of world firsts in quantum technologies began. The spacecraft carried the first satellite for quantum communications, named after the Chinese philosopher Micius. A year later, Chinese researchers described in the journal "Nature" how they built for the first time a 1200 kilometers long "unhackable" quantum connection from space to earth. China's state news agency Xinhua praised the "enormous prospects" of the experiment for the areas of finance and defense. Communication, i.e. data transmission, is considered to be one of the most promising applications for quantum technologies, both civil and military use. This is due to the nature of quantum technologies. These are based on the fact that tiny particles like atoms can have several states at the same time – not just zero or one like the digital technologies that we all use today. These quantum particles, called qubits, are incredibly sensitive. Any small change can alter their state.

Therefore, they are suitable for encryption of communication. Because as soon as a hacker breaks the connection or even just measures it, the quantum state of the encryption changes. "The intercepted information will self-destruct", China's quantum luminary Pan Jian-Wei told Xinhua in 2017. The US Department of Defense called China's successful Quantum communication via satellite in 2017 a "remarkable progress". Since then, China has reported a number of other breakthroughs. Quantum technologies have become a key area in the US race and China's order to become technological supremacy. NATO is worried too. The defense alliance designated China in the last June in his "Strategic Concept" for the first time as one Challenge for "his interests, his security and his values".

Quantum technologies for the military: China's pace occupies the USA last year, NATO announced research into new technologies like that of quantum technologies - and explicitly referred to China's Quan-

tum satellites and corresponding investments by the Chinese government and the companies Baidu and Huawei. So how exactly could Quantum technologies transform the military? What can China already do, and what could it possibly do in the future? And how do they react? Is the USA on it?

## 25.2 Military Applications of Quantum Technologies

Experts see three most promising in the military sector Applications for quantum technologies.

First the mentioned encrypted one Communications in China with its satellite since 2016 more world premieres. This is how the researchers around Pan Jian-Wei built for one secure video call the first transcontinental quantum connection on - between ground stations near Beijing and in Austria, where Pan's supervisor researches. Such experiments are according to the US Department of Defense an "important milestone" in building a "global, ultra-secure network" for data transfer. The Pentagon writes in its latest annual report on China's military capabilities that the introduction of this technology still faces hurdles before being used widely.

Second, sensor technology is an important area of the military quantum technologies. According to experts, it is the most mature application. Therefore, for example, quantum sensors could detect submarines and stealth aircraft. They could also serve in positioning and navigation systems. Compared to the GPS commonly used today, they would have the advantage that they could function autonomously, i.e. without external relays such as satellites. These advantages are also based on the fact that quantum particles are extremely sensitive to disturbances. This could make quantum sensors measure tiny differences in characteristics like temperature, acceleration, gravity or time, says an article on the NATO website. Therefore, new application possibilities result. According to the magazine "Scientific American", quantum sensors could scan

the Human brain activity of people in daily life – and enable autonomous vehicles to "see" around the corners of houses. Chinese researchers from Tsinghua University in Beijing reported 2021 successful experiment in quantum sensing. So, they created one electromagnetic storm with quantum particles and fired it on a stealth object. The particles bounced off - the object was exposed. This principle could be used to build a more sophisticated quantum radar than all previous radar systems, said the study leader in Chinese Journal of Radars. However, the construction of such a radar is not an easy task.

## 25.3 Quantum Computers Threaten Encryption

Third, quantum computers are an important area for military applications. They can be very much faster than conventional computers in solving limited but highly complex problems. Probably the most mentioned application in this context is that quantum computers in future could crack the RSA encryption technology that is common today, which is also used by the military. At the end of December 2022, 24 Chinese researchers, in a still unexamined study, claim having found a way to crack RSA encryption. The study was heavily criticized by many experts; a Chinese Quantum researcher, who wanted to remain anonymous, named it for this article "garbage". But regardless of the content, some observers see in the study that Quantum computers endanger encryption. Overnight, a quantum breakthrough could destroy any RSA encryption and make it obsolete. China could gain this ability first.

## 25.4 The Goal is Quantum Supremacy

The computing power is decisive for this. Above all, the USA and China compete to build the fastest computer. One goal is the so-called quantum supremacy. That means, that a quantum computer can solve a problem that a classical calculator would fail in useful time. In 2020,

Chinese scientists announced quantum supremacy for the Computer Jiuzhang. In 2021, they claimed the successor Jiuzhang 2 could solve in one millisecond a problem that conventional computers could solve only in 30 billion years. A second computer, based on a different quantum technology, is said to have been a million times faster than its closest rival Google.

An important criterion for the capacity of quantum computers is the number of the usable quantum bits, the qubits. Here, too, a team around Pan Jian-Wei announced in 2021 to have set a world record, with 62 qubits. Who is currently ahead in this discipline changes every few months. Currently an IBM computer with 433 qubits is considered the world leader. The prospect that quantum computers could one day crack encryption technology, makes data theft in reserve attractive. Probably, the US and China are today stealing encrypted information from each other that they can't yet read - but in few years they will. This strategy is called in professional circles "Harvest now, decrypt later".

How promising China's military quantum capabilities are compared to those of the USA and NATO is difficult to say. The technologies are mostly not usable yet. It is also unclear, which ones were successful experiments, which could really lead to useful applications.

Many research programs are also likely to be secret, especially the military one. The British market research company Global Data writes, according to their research, "the United States almost certainly integrated secret Quantum development projects into existing programs in their most important Government laboratories", for example at the NSA secret service.

## 25.5 Quantum Race: Caution in Assessments

Of course, there are clues as to where the quantum race is going. In general, many experts continue to see China behind the USA, especially in fundamental research. Yet China could certainly overtake its rival in

developing and deploying the first applications, wrote the American consulting firm Booz Allen Hamilton at the end of 2021. In the fairly broad field of quantum information science, the American government, as early as 2017, stated that China faced its backlog caught up with the USA. In the mentioned area of quantum communication, China is the leader for many observers, thanks to the Micius satellite. And with the number of quantum patents, China has been ahead of the USA for years.

Such assessments should of course be treated with caution: American consultants and government officials, for example, could have interest to show China's quantum capabilities overly dramatic in order to get more consulting mandates and higher budgets. Conversely, Chinese representatives might be tempted to exaggerate their own skills to please their political leadership - or belittle them, so as not to alarm the US too much. China is at least hot on the heels of the US.

In any case, both governments have long been investing billions in Quantum technologies, and the European Union also has a corresponding program. Already 2021, Washington put Chinese firms under sanctions, which develop military applications for quantum technologies. The responsible Ministry of Commerce explicitly named applications for encryption and decryption as well as for the detection of camouflage objects and submarines. A little-noticed aspect of Washington's decision in August 2022 funding law for chips and science is that - in order to compete with China – USA should not only subsidize semiconductors with billions, but also quantum technologies. The United States apparently want to join the ranks with its partners in this area, not only within Nato, but also bilaterally. They have already signed agreements with countries such as Great Britain, France, Denmark and in October 2021 with Switzerland.

**26**

# UKRAINE WAR:
# ARTIFICIAL INTELLIGENCE
# ON THE BATTLEFIELD

*Mita Pressl / Pascal Burkhard, Ukraine*

*This text is the transcript of a video of August 2, 2022 from Kyiv[321].*
*Topic: When AI makes mistakes. AI face and speech recognition. Strong*
*and weak AI.*

This man is said to have looted an apartment in a Kyiv suburb. A war
crime - albeit one of the less serious that Russian soldiers are currently
committing in Ukraine. With a Polaroid camera that he found in the
apartment, he took a selfie and then left it carelessly. Ukraine wants to
identify the man using artificial intelligence. Because in the Ukraine war
both warring parties are currently trying gain an advantage with AI. The

---

[321] Transcribed and translated from German by the volume editors. Original
video: Ukraine-Krieg: Künstliche Intelligenz [KI] am Schlachtfeld. Wenn KI
Fehler macht. KI Gesichts- und Spracherkennung. Starke und schwache KI,
Neue Zürcher Zeitung, 2 Aug 2022.

Ukraine in particular makes this very effective in terms of publicity –
this often works well, but sometimes not so well. Deputy Prime Minister
and Minister of Digitization of Ukraine Mikhail Fedorov claims to know
the identity of the soldier and publicly denounces him. 'With AI tech-
nology' you have the 'looters and war criminals', so his words can iden-
tify within a minute. It is 26-year-old Nikita Tretyakov. One of the few
remaining independent Russian media, "Важные истории" (Vashniye
istorii: in Russian "important stories"), sees things differently. They say:
No, the photo shows 25-year-old Pavel Pilantow. In fact, Pilantov is
active in that regiment, which was in that suburb of Kyiv, where the
photo was taken. The journalists could confirm that. Nikita Tretyakov,
on the other hand, was not in Ukraine at the time of the fighting, his
sister claims. He hasn't had anything to do with the army for a long time.
There are also photos of him at an event on the Internet, taken on the day
of the fighting. Due to the media attention that followed the Vice Prem-
ier's tweet, Nikita Tretyakov deleted his presence on social networks.
We don't know what he thinks about the war. But the Polaroid photo
certainly looks a lot more like Pavel Pilantov than him. The AI technol-
ogy that the Ukrainian minister is talking about apparently also makes
mistakes. And with such delicate purposes, every mistake is a problem.
Ukraine not only uses Clearview AI software, to identify war criminals
and prisoners of war and collect evidence, but also to identify those
killed Russian soldiers some of which the Russian army keeps secret
from their own families. Somewhere between information gathering and
information warfare. However, AI algorithms are not only used for facial
recognition. The Ukrainian secret service has succeeded, on a grand
scale, to eavesdrop on unencrypted phone calls between Russian soldiers
and their families. These were transcribed, translated and analysed using
a tool by an American company called Primer. It would take too many
people to listen to the phone calls one by one. So that's left to an AI. You
train the speech recognition on certain words and phrases, teaches her to
recognize dialects as well and then orders her to filter out relevant con-

versations. Like this shocking example of a soldier calling his wife at home. For the Ukrainian side, such calls are worth their weight in gold in the information war. They show how brutally Russian soldiers act and leave no doubt about it: The "good side - this is Ukraine." So we see: AI has become a powerful tool in war, but it is still error-prone. The holy grail of AI research has not even come close. AI can basically be divided into two categories: weak and strong AI. All of the above examples are so-called weak AI. So AI, which is usually trained using machine learning, to perform a specific task, such as recognizing faces, faster and better than humans. Strong, or general AI, would be able, similar to rational humans, solve complex problems and make autonomous decisions. Strong AI doesn't really exist yet, but only the hopes in Russia and elsewhere for massive geopolitical gains from their successful invention. An AI that is intellectually on par with humans, but does not fear for its own survival, would be capable of terrible things in war. However, we are still a long way from such a "Terminator robot". Autonomous AI weapons, for example, i.e. weapons that are controlled solely by artificial intelligence, neither Russia nor Ukraine currently has this, although of course research is being carried out on it worldwide and breakthroughs also occur. For example, Russia states that its Kalashnikov-made

KUB-BLA drone Targets are already identified using AI. Whether this is true is difficult to verify. But even with weak AI, military applications are more complicated than civilian ones. Because AI needs huge amounts of data to learn. And because, fortunately, there are not so many military confrontations.

<div style="text-align:center">

**27**

</div>

# INTERNATIONAL HUMANITARIAN LAW AND CYBER OPERATIONS DURING ARMED CONFLICTS

*International Committee of the Red Cross ICRC Position Paper*

## Executive Summary

Cyber operations have become a reality in contemporary armed conflict. The International Committee of the Red Cross (ICRC) is concerned by the potential human cost arising from the increasing use of cyber operations during armed conflicts.

In the ICRC's view, international humanitarian law (IHL) limits cyber operations during armed conflicts just as it limits the use of any other weapon, means and methods of warfare in an armed conflict, whether new or old.

Affirming the applicability of IHL does not legitimise cyber warfare, just as it does not legitimise any other form of warfare. Any use of force by States – cyber or kinetic – remains governed by the Charter of the United Nations and the relevant rules of customary international law, in particular the prohibition against the use of force. International disputes must be settled by peaceful means, in cyberspace as in all other domains.

It is now critical for the international community to affirm the applicability of international humanitarian law to the use of cyber operations during armed conflicts. The ICRC also calls for discussions among governmental and other experts on how existing IHL rules apply and whether the existing law is adequate and sufficient. In this respect, the ICRC welcomes the intergovernmental discussions currently taking place in the framework of two United Nations General Assembly mandated processes.

Events of recent years have shown that cyber operations, whether during or outside armed conflict, can disrupt the operation of critical civilian infrastructure and hamper the delivery of essential services to the population. In the context of armed conflicts, civilian infrastructure is protected against cyber attacks by existing IHL principles and rules, in particular the principles of distinction, proportionality and precautions in attack. IHL also affords special protection to hospitals and objects indispensable to the survival of the civilian population, among others.

During armed conflicts, the employment of cyber tools that spread and cause damage indiscriminately is prohibited. From a technological perspective, some cyber tools can be designed and used to target and harm only specific objects and to not spread or cause harm indiscriminately. However, the interconnectivity that characterises cyberspace means that whatever has an interface with the internet can be targeted from anywhere in the world and that a cyber attack on a specific system may have repercussions on various other systems. As a result, there is a real risk that cyber tools are not designed or used – either deliberately or by mistake – in compliance with IHL.

States' interpretation of existing IHL rules will determine the extent to which IHL protects against the effects of cyber operations. In particular, States should take clear positions about their commitment to interpret IHL so as to preserve civilian infrastructure from significant disruption and to protect civilian data. The availability of such positions will also influence the assessment of whether the existing rules are adequate

or whether new rules may be needed. If States see a need to develop new rules, they should build on and strengthen the existing legal framework – including IHL.

## 27.1 Introduction

The use of cyber operations during armed conflicts is a reality.[322] While only a few States have publicly acknowledged using such operations, an increasing number of States are developing military cyber capabilities, and their use is likely to increase in future. Moreover, there have been significant technological advances in offensive cyber capabilities: in recent years, cyber operations have shown that they can seriously affect civilian infrastructure and might result in human harm.

In line with its mission and mandate, the International Committee of the Red Cross (ICRC) is primarily concerned with cyber operations used as means and methods of warfare during an armed conflict and the protection that international humanitarian law (IHL) affords against their effects.

The ICRC welcomes the intergovernmental discussions currently taking place in the framework of the two United Nations General Assembly mandated processes, namely the Open-Ended Working Group on Developments in the Field of Information and Telecommunications in the Context of International Security and the Group of Governmental Experts on Advancing Responsible State Behaviour in Cyberspace in the Context of International Security. Both groups are mandated to study

---

[322] In this position paper (Orig. publ. ICRC, 28 Nov. 2019, https://www.icrc.org/en/document/international-humanitarian-law-and-cyber-operations-during-armed-conflicts), the term 'cyber operations during armed conflicts' is used to describe operations against a computer, a computer system or network, or another connected device, through a data stream, when used as means and methods of warfare in the context of an armed conflict. Cyber operations rely on information and communication technologies.

"how international law applies to the use of information and communications technologies by States".[323]

The ICRC submits this position paper to both groups to support States' deliberation on this matter.

This position paper is limited to legal and humanitarian questions arising from the use of cyber operations during armed conflict. It does not address questions relating to the legal framework applicable to cyber operations unrelated to armed conflict.

## 27.2 The Potential Human Cost of Cyber Operations

During armed conflict, cyber operations have been used in support of or alongside kinetic operations. The use of cyber operations may offer alternatives that other means or methods of warfare do not, but it also carries risks. On the one hand, cyber operations have the potential to enable parties to armed conflicts to achieve their military aims without harming civilians or causing physical damage to civilian infrastructure. On the other hand, recent cyber operations – which have been mostly conducted outside the context of armed conflict – show that sophisticated actors have developed the capability to disrupt the provision of essential services to the civilian population.

By means of cyber operations, it is possible for belligerents to infiltrate a system and collect, exfiltrate, modify, encrypt, or destroy data. It is also possible to trigger, alter or otherwise manipulate processes controlled by a compromised computer system. A variety of "targets" in the real world can be disrupted, altered or damaged, such as industries, infrastructures, telecommunications, transport, or governmental and financial systems. Based on discussions with experts from all parts of the world and its own research, the ICRC is particularly concerned about the

---

[323] A/RES/73/27, OP 5; A/RES/73/266, OP 3.

potential human cost of cyber operations on critical civilian infrastructure, including health infrastructure.[324]

In recent years, cyber attacks have exposed the vulnerability of essential services. They are reportedly becoming more frequent and their severity is increasing more rapidly than experts had anticipated. Moreover, much is unknown with respect to the most sophisticated cyber capabilities and tools that have been or are being developed, how technology may evolve, and the extent to which the use of cyber operations during armed conflicts might be different from the trends observed so far.

Moreover, the characteristics of cyberspace raise specific concerns. For example, cyber operations entail a risk for escalation and related human harm for the simple reason that it may be difficult for the targeted party to know whether the attacker's aim is intelligence collection or more harmful effects. The target may thereby react with greater force than necessary out of anticipation of a worst-case scenario.

Cyber tools also proliferate in a unique manner. Once used, they can be repurposed and widely used by actors other than the one that developed or used them initially.

## 27.3 The Application of IHL to Cyber Operations During Armed Conflicts

For the ICRC, there is no question that IHL applies to, and therefore limits, cyber operations during armed conflict – just as it regulates the use of any other weapon, means and methods of warfare in an armed conflict, whether new or old.[325] This holds true whether cyberspace is

---

[324]See ICRC, The Potential Human Cost of Cyber Operations, 2019; available at https://www.icrc.org/en/download/file/96008/the-potential-human-cost-of-cyber-operations.pdf.

[325]  ICRC, International Humanitarian Law and the Challenges of Contemporary Armed Conflicts, 2011, 31IC/11/5.1.2, pp. 36-37; available at https://www.icrc.org/en/doc/assets/files/red-cross-crescent-movement/31st-

considered as a new domain of warfare similar to air, land, sea and outer space; a different type of domain because it is man-made while the former are natural; or not a domain as such.

When States adopt IHL treaties, they do so to regulate present and future conflicts. States have included rules that anticipate the development of new means and methods of warfare in IHL treaties, presuming that IHL will apply to them. For instance, if IHL did not apply to future means and methods of warfare, it would not be necessary to review their lawfulness under existing IHL, as required by Article 36 of the 1977 First Additional Protocol.

This conclusion finds strong support in the International Court of Justice's Advisory Opinion on the Legality of the Threat or Use of Nuclear Weapons: the Court recalled that the established principles and rules of IHL applicable in armed conflict apply 'to all forms of warfare and to all kinds of weapons', including 'those of the future'.[326] In the ICRC's view, this finding applies to the use of cyber operations during armed conflict.

The ICRC welcomes that an increasing number of States and international organisations have affirmed that IHL applies to cyber operations during armed conflicts and welcomes discussion on how IHL applies.

States may also decide to impose additional limits to those found in existing law and develop complementary rules, in particular in order to strengthen the protection of civilians and civilian infrastructure against

---

international-conference/31-int-conference-ihl-challenges-report-11-5-1-2-en.pdf; International Humanitarian Law and the Challenges of Contemporary Armed Conflicts, 2015, 32IC/15/11, p. 40; available at: https://www.icrc.org/en/download/file/15061/32ic-report-on-ihl-and-challenges-of-armed-conflicts.pdf; ICRC, International Humanitarian Law and the Challenges of Contemporary Armed Conflicts, 2019, 33IC/19/9.7, p. 18; available at: https://rcrcconference.org/app/uploads/2019/10/33IC-IHL-Challenges-report_EN.pdf.

[326] International Court of Justice, Legality of the threat or the use of nuclear weapons, Advisory Opinion, 8 July 1996, para. 86.

the effects of cyber operation. In the ICRC's view, any new rules need to build on and strengthen the existing legal framework, including IHL.

In cases not covered by existing rules of IHL, civilians and combatants remain protected by the so-called "Martens clause", meaning they remain under the protection and authority of the principles of international law derived from established custom, from the principles of humanity and from the dictates of public conscience.[327]

It is important to underline that affirming the application of IHL to cyber operations during armed conflict does not legitimise cyber warfare or encourage the militarization of cyberspace. In fact, IHL.[328]Moreover, any use of force by States – cyber or kinetic – remains governed by the Charter of the United Nations and the relevant rules of customary international law, in particular, the prohibition against the use of force. International disputes must be settled by peaceful means, in cyberspace as in all other domains.

## 27.4 The Protection Afforded by Existing IHL

Existing IHL treaties and customary law provide rules on a number of issues during armed conflict. In cyberspace, the rules on the conduct of hostilities are particularly relevant. These rules aim to protect the civilian population against the effects of hostilities. They are based on the cardinal principle of distinction, which requires that belligerents distinguish at all times between the civilian population and combatants

---

[327]See Art. 1(2) of Protocol Additional to the Geneva Conventions of 12 August 1949, and relating to the Protection of Victims of International Armed Conflicts (AP I); paragraph 9 of the preamble to the 1899 Hague Convention (II); paragraph 8 of the preamble to the 1907 Hague Convention (IV).

[328]See, among others, Henckaerts and Doswald-Beck (eds), Customary International Humanitarian Law, Vol. I: Rules, ICRC, Cambridge University Press, Cambridge, 2005 (hereinafter ICRC Customary IHL Study), Rules 70 and 71; see also Art. 36 AP I.

and between civilian objects and military objectives, and direct their operations only against military objectives.[329] Notwithstanding the inter-connectivity that characterises cyberspace, a careful examination of the functioning of cyber tools shows that they are not necessarily indiscriminate. Many of the recent cyber attacks that have been reported in public sources appear to have been rather "discriminate" from a technical point of view: they have been designed and actually used to target and harm only specific objects and have not spread or caused harm indiscriminately. Ensuring that cyber operations affect only the targeted object may, however, be technically challenging and require careful planning in their design and use. Moreover, it must be noted that a cyber operation that is technically discriminatory is not necessarily lawful, whether during or outside of an armed conflict.

This being said, some known cyber tools have been designed to self-propagate and indiscriminately affect widely used computer systems. They have not done so by chance: the ability to self-propagate needs to be specifically included in the design of such tools. The interconnectivity that characterises cyberspace means that whatever has an interface with the internet can be targeted from anywhere in the world. Moreover, an attack on a specific system may have repercussions on various other systems and cause indiscriminate effects. As a result, there is a real risk that cyber tools are not designed or used – either deliberately or by mistake – in compliance with IHL.

Affirming that IHL – including the principles of distinction, proportionality, and precautions – applies to cyber operations during armed conflicts means that under existing law, among many other rules:

- cyber capabilities that qualify as weapons and are by nature indiscriminate are prohibited;[330]

---

[329] Art. 48 AP I; Rules 1 and 7 ICRC Customary IHL Study. International Court of Justice, Legality of the threat or the use of nuclear weapons, Advisory Opinion, 8 July 1996, para. 78.

[330] Rule 71 ICRC Customary IHL Study.

- direct attacks against civilians and civilian objects are prohibited, including when using cyber means or methods of warfare;[331]
- acts or threats of violence the primary purpose of which is to spread terror among the civilian population are prohibited, including when carried out through cyber means or methods of warfare;[332]
- indiscriminate attacks, namely those of a nature to strike military objectives and civilians or civilian objects without distinction, are prohibited, including when using cyber means or methods of warfare;[333]
- disproportionate attacks are prohibited, including when using cyber means or methods of warfare. Disproportionate attacks are those which may be expected to cause incidental loss of civilian life, injury to civilians, damage to civilian objects, or a combination thereof, which would be excessive in relation to the concrete and direct military advantage anticipated.[334]
- during military operations, including when using cyber means or methods of warfare, constant care must be taken to spare the civilian population and civilian objects; all feasible precautions must be taken to avoid or at least minimise incidental civilian harm when carrying out attacks, including through cyber means and methods of warfare;[335]

---

[331] Arts 48, 51 and 52 AP I; Rules 1 and 7 ICRC Customary IHL Study.

[332] Art. 51(2) AP I; Rule 2 ICRC Customary IHL Study.

[333] Art. 51(4) AP I; Rules 11 and 12 ICRC Customary IHL Study. Indiscriminate attacks are those: (a) which are not directed at a specific military objective; (b) which employ a method or means of combat which cannot be directed at a specific military objective; or (c) which employ a method or means of combat the effects of which cannot be limited as required by international humanitarian law; and consequently, in each such case, are of a nature to strike military objectives and civilians or civilian objects without distinction.

[334] Arts 51(5)(b) and 57 AP I; Rule 14 ICRC Customary IHL Study.

[335] Art. 57 AP I; Rules 15 - 21 ICRC Customary IHL Study

- attacking, destroying, removing or rendering useless objects indispensable to the survival of the population is prohibited, including through cyber means and methods of warfare;[336]
- medical services must be protected and respected, including when carrying out cyber operations during armed conflicts.[337]

In addition, all feasible precautions must also be taken to protect civilians and civilian objects against the effects of attacks conducted through cyber means and methods of warfare, which is an obligation that States must already implement in peacetime.[338] Measures that could be considered include, among others: segregating military from civilian cyber infrastructure and networks; segregating computer systems on which essential civilian infrastructure depends from the internet; work on the identification in cyberspace of the cyber infrastructure and networks serving specially protected objects like hospitals.[339]

## 27.5 The Need to Discuss how IHL Applies

Affirming that IHL applies to cyber operations in armed conflict is an essential first step to avoid or minimise the potential human suffering that cyber operations might cause. However, the ICRC also encourages States to work towards a common understanding of how IHL principles

---

[336] Art. 54 AP I; Art. 14 Protocol Additional to the Geneva Conventions of 12 August 1949, and relating to the Protection of Victims of Non-International Armed Conflicts (AP II); Rule 54 ICRC Customary IHL Study.

[337] See, for instance, Art. 19 Convention (I) for the Amelioration of the Condition of the Wounded and Sick in Armed Forces in the Field (GCI); Art. 12 Convention (II) for the Amelioration of the Condition of Wounded, Sick and Shipwrecked Members of Armed Forces at Sea (GCII); Art. 18 Convention (IV) relative to the Protection of Civilian Persons in Time of War (GCIV); Art. 12 AP I; Art. 11 AP II; Rules 25, 28, 29 ICRC Customary IHL Study.

[338] Art. 58 AP I; Rules 22 to 24 ICRC Customary IHL Study.

[339] ICRC, International humanitarian law and the challenges of contemporary armed conflicts, 2015, p. 43.

and rules apply to cyber operations. This is necessary because the inter-connected nature of cyberspace and its largely digital character pose challenges for the interpretation of key IHL principles and concepts on the conduct of hostilities.

Among the various issues, in this position paper the ICRC emphasises three.

### 27.5.1 The military use of cyberspace and the effect on its civilian character

Except for some specific military networks, cyberspace is predominantly used for civilian purposes. However, civilian and military networks may be interconnected. Furthermore, military networks may rely on civilian cyber infrastructure, such as undersea fibre-optic cables, satellites, routers or nodes. Conversely, civilian vehicles, shipping and air traffic controls increasingly rely on navigation satellite systems that may also be used by the military. Civilian logistical supply chains and essential civilian services use the same web and communication networks through which some military communications pass.

Not every use for military purposes renders a civilian object a military objective under IHL.[340] If it does, however, the object is no longer protected by the prohibition to direct attacks on civilian objects. It would be a matter of serious concern if the military use of cyberspace led to the conclusion that many objects forming part thereof would no longer be

---

[340] See Art. 52(2) AP I; Rule 8 Customary IHL Study: "In so far as objects are concerned, military objectives are limited to those objects which by their nature, location, purpose or use make an effective contribution to military action and whose partial or total destruction, capture or neutralisation, in the circumstances ruling at the time, offers a definite military advantage." For more details on the limits to cyber infrastructure becoming military objectives under IHL, see ICRC, *International humanitarian law and the challenges of contemporary armed conflicts*, 2015, p. 42.

protected as civilian objects. This could lead to large-scale disruption of the ever-increasingly important civilian usage of cyberspace.

This being said, even if certain parts of the cyberspace infrastructure were no longer protected as civilian objects during armed conflicts, any attack would remain governed by the prohibition of indiscriminate attacks and the rules of proportionality and precautions in attack. Precisely because civilian and military networks are so interconnected, assessing the expected incidental civilian harm of any cyber operation is critical to ensure that the civilian population is protected against its effects.[341]

### 27.5.2 The notion of 'attack' under IHL and cyber operations

Critical civilian infrastructure enabling the provision of essential services increasingly relies on digitised systems. Safeguarding such infrastructure and services against cyber attacks or incidental damage is essential to protect the civilian population.

IHL provides specific protection for certain infrastructure, such as medical services and objects indispensable to the survival of the population, regardless of the type of harmful operation.[342]However, most rules stemming from the principles of distinction, proportionality and precautions – which provide general protection for civilians and civilian objects – apply only to military operations that qualify as 'attacks' as defined in IHL.[343]Article 49 of Additional Protocol I defines attacks as 'acts of

---

[341] See ICRC, The Principle of Proportionality in the Rules Governing the Conduct of Hostilities under International Humanitarian Law, 2018, available at https://www.icrc.org/en/download/file/79184/4358_002_expert_meeting_report_ web_1.pdf , pp. 37–40.

[342] See text in relation to footnotes 16 and 15 above. With regard to the latter, they must not be attacked, destroyed, removed or rendered useless.

[343] The notion of attack under IHL, defined in Art. 49 of the 1977 First Additional Protocol, is different from and should not be confused with the notion of 'armed attack' under Art. 51 of the UN Charter, which belongs to the realm of jus ad bellum. To affirm that a specific cyber operation, or a type of cyber opera-

violence against the adversary, whether in offence or in defence'[344] The question of how widely or narrowly the notion of 'attack' is interpreted with regard to cyber operations is therefore essential for the applicability of these rules and the protection they afford to civilians and civilian infrastructure.

It is widely accepted that cyber operations expected to cause death, injury or physical damage constitute attacks under IHL. In the ICRC's view, this includes harm due to the foreseeable direct and indirect (or reverberating) effects of an attack, for example the death of patients in intensive-care units caused by a cyber operation on an electricity network that results in cutting off a hospital's electricity supply.

Beyond this, attacks that significantly disrupt essential services without necessarily causing physical damage constitute one of the most important risks for civilians. Diverging views exist, however, on whether a cyber operation that results in a loss of functionality without causing physical damage qualifies as an attack as defined in IHL. In the ICRC's view, during an armed conflict an operation designed to disable a computer or a computer network constitutes an attack under IHL, whether the object is disabled through kinetic or cyber means.[345] If the notion of attack is interpreted as only referring to operations that cause death, injury or physical damage, a cyber operation that is directed at making a civilian network (such as electricity, banking, or communications) dysfunctional, or is expected to cause such effect incidentally, might not be covered by essential IHL rules protecting the civilian population and

tions, amounts to an attack under IHL does not necessarily mean that it would qualify as an armed attack under the UN Charter.

[344] For rules that apply specifically to attacks, see text in relation to footnotes 10 to 14 above.

[345] See ICRC, International humanitarian law and the challenges of contemporary armed conflicts, 2011, p. 37; ICRC, International humanitarian law and the challenges of contemporary armed conflicts, 2015, pp. 41-42.

civilian objects. Such an overly restrictive understanding of the notion of attack would be difficult to reconcile with the object and purpose of the IHL rules on the conduct of hostilities. It is therefore essential that States find a common understanding in order to adequately protect the civilian population against the effects of cyber operations.

### 27.5.3 Civilian data and the notion of 'civilian objects'

Essential civilian data – such as medical data, biometric data, social security data, tax records, bank accounts, companies' client files or election lists and records – are an essential component of digitalized societies. Such data are key to the functioning of most aspects of civilian life, be it at individual or societal level. There is increasing concern about safeguarding such essential civilian data.

Some of the specific protection afforded by IHL extends to essential data, such as data belonging to medical units, which are encompassed in the obligation to respect and protect such units.[346]

More generally, the main IHL principles and rules governing the conduct of hostilities protect civilians and civilian objects.[347] It would therefore be important for States to agree on an understanding that civilian data is protected by these rules.

Deleting or tampering with essential civilian data can quickly bring government services and private businesses to a complete standstill. Such operations could cause more harm to civilians than the destruction of physical objects. While the question of whether and to what extent civilian data constitute civilian objects remains unresolved, in the ICRC's view the assertion that deleting or tampering with such essential civilian data would not be prohibited by IHL in today's data-reliant world seems difficult to reconcile with the object and purpose of IHL. The replacement of paper files and documents with digital files in the form of data should not decrease the protection that IHL affords to

---

[346] See footnote 16 above.

[347] See text in relation to notes 10 to 15 above.

them.[348] Excluding essential civilian data from the protection afforded by IHL to civilian objects would result in an important protection gap.

## 27.6 Attribution of Conduct in Cyberspace for the Purposes of State Responsibility

Cyberspace provides various technical possibilities for actors to hide or falsify their identity, which increases the complexity of attribution by other actors. This creates major difficulties. For example, even during armed conflict, IHL only applies to operations that are linked to the conflict. If the author of a cyber operation – and thus the link of the operation to an armed conflict – cannot be identified, it may be difficult to determine whether IHL is even applicable to the operation. Attribution of cyber operations is also important to ensure that actors who violate international law, including IHL, can be held accountable. The perception that it will be easier to deny responsibility for such attacks may also weaken the taboo against their use – and may make actors less scrupulous about using them in violation of international law.

This being said, attribution is not a problem from the perspective of the actors who conduct, direct or control cyber operations: they have all the facts at hand to determine under which international legal framework they are operating and which obligations they must respect.

Under international law, a State is responsible for conduct attributable to it, including possible violations of IHL. This includes:

- conduct by State organs, including its armed forces or intelligence services;
- conduct by persons or entities, such as private companies, the State empowered to exercise elements of governmental authority;

---

[348] ICRC, International humanitarian law and the challenges of contemporary armed conflicts, 2015, p. 43; ICRC, International humanitarian law and the challenges of contemporary armed conflicts, 2019, p. 21

- conduct by persons or groups, such as militias or group of hackers, acting in fact on the State's instructions, or under its direction or control; and
- conduct by private persons or groups which the State acknowledges and adopts as its own conduct.[349]

These principles apply whether the conduct is carried out by cyber or any other means.

## 27.7 Conclusion

The use of cyber operations as means or methods of warfare in an armed conflict poses a real risk of harm to civilians. For the protection of the civilian population and civilian infrastructure, it is critical to recognize that such operations do not occur in a legal vacuum. The ICRC urges all States to affirm that IHL applies to cyber operations during armed conflicts, on the understanding that such affirmation neither encourages the militarization of cyberspace nor legitimises cyber warfare.

At the same time, the ICRC believes that further discussion – especially among States – is needed on how IHL should be interpreted and applied in cyberspace. There is a pressing need for such discussion because States that decide to develop or acquire cyber capabilities that qualify as weapons, means and methods of warfare – whether for offensive or defensive purposes – must ensure that these capabilities can be used in accordance with their obligations under IHL.30 Discussion should be informed by an in-depth understanding of the development of military cyber capabilities, their potential human cost, and the protection afforded by existing law. States need to determine whether existing law is adequate and sufficient to address the challenges posed by the interconnected and largely digital character of cyberspace, or whether it

---

[349] See Rule 149 ICRC Customary IHL Study. See also International Law Commission, Responsibility of States for Internationally Wrongful Acts, 2001, in particular Articles 4 to 11.

needs adaptation to the specific characteristics of cyberspace. If new rules are to be developed to protect civilians against the effects of cyber operations or for other reasons, they should build on and strengthen the existing legal framework – including IHL.

The ICRC welcomes the intergovernmental discussions currently taking place in the framework of two United Nations General Assembly mandated processes and it is grateful for the opportunity to share its views with the participating States. The ICRC also stands ready to lend its expertise to such discussions, as States deem appropriate.

# 28

# ESTABLISHING A
# DIGITAL GENEVA CONVENTION

*Dan Lohrmann*[350]

*Exponential increases in global cyber crime. Ransomware crippling governments and businesses. Nations ignoring cyber criminals operating on their soil. The time for international cooperation on cybersecurity is now.*[351]

---

[350] Daniel J. Lohrmann is an internationally recognized cybersecurity leader, technologist, keynote speaker and author of several widespread books. He served global organizations in public and private sector. His award-winning blog: http://www.govtech.com/blogs/lohrmann-on-cybersecurity. There is also his blog "Top 23 Security Predictions for 2023" (30 Dec 2022).

[351] This article was published on 8 August 2021 as blog: https://www.govtech. com/blogs/lohrmann-on-cybersecurity/the-case-for-establishing-a-digital-geneva-convention. Request to republish has been made via Govtech.com.

## 28.1 A Basic Question and Definitions

Let's start with a question: What do all of these activities have in common?

- Stopping ransomware from devastating consequences.
- Protecting critical infrastructure from cyber attacks.
- Policing illegal cyberspace activities.
- Bringing global cyber criminals to justice.
- Holding nation-states accountable for online criminal activities.
- International rules for war in the 2020s and beyond.

While there are many potential answers to this question, a growing number of international experts believe that these issues call for a new "Digital Geneva Convention" to address a growing global mess in cyberspace that is having very real impacts in the daily lives of individuals, companies and governments around the world. But before we dig deeper into this topic, here are a few important definitions.

According to the International Committee of the Red Cross[352]: "The Geneva Conventions and their Additional Protocols are international treaties that contain the most important rules limiting the barbarity of war. They protect people who do not take part in the fighting (civilians[353], medics, aid workers) and those who can no longer fight (wounded, sick and shipwrecked troops, prisoners of war)."

Article 3[354], common to the four Geneva Conventions, marked a breakthrough, as it covered, for the first time, situations of non-

---

[352] The Geneva Conventions of 1949 and their Additional Protocols. 2014. ICRC, https://www.icrc.org/en/document/geneva-conventions-1949-additional-protocols

[353] Civilians. 2023. ICRC, https://www.icrc.org/en/war-and-law/protected-persons/civilians

[354] Geneva Conventions of 1949, Additional Protocols and Their Commentaries, art. 3, Conflicts not of an international character, https://ihl-databases.icrc.org/en/ihl-treaties/gci-1949/article-3

international armed conflicts. These types of conflicts vary greatly. They include traditional civil wars, internal armed conflicts that spill over into other states or internal conflicts in which third states or a multinational force intervenes alongside the government. Common Article 3 establishes fundamental rules from which no derogation is permitted. It is like a mini-Convention within the Conventions as it contains the essential rules of the Geneva Conventions in a condensed format and makes them applicable to conflicts not of an international character. It requires humane treatment for all persons in enemy hands, without any adverse distinction. It specifically prohibits murder, mutilation, torture, cruel, humiliating and degrading treatment, the taking of hostages and unfair trial. It requires that the wounded, sick and shipwrecked be collected and cared for.

## 28.2 Adding Cyberspace

But perhaps you're wondering what the Geneva Convention and/or physical conflicts in the real world have to do with cyberspace and cybersecurity. Here is some recent history on the subject of creating a Cyber Geneva Convention or Digital Geneva Convention:

**Microsoft's Brad Smith at RSA 2017: The Need for a Digital Geneva Convention.** World Economic Forum (2017): *Why we urgently need a Digital Geneva Convention*[355]:

"The United Nations almost two decades ago set up a working body to ensure agreement is reached on how to handle the then relatively new field of information technology (IT), and in particular the increasingly difficult question of cybersecurity. It took a while, but in 2015, the United Nations Group of Governmental Experts on Developments in the

---

[355] Why we urgently need a Digital Geneva Convention. 2017. World Economic Forum's Geostrategy platform, https://www.weforum.org/agenda/2017/12/why-we-urgently-need-a-digital-geneva-convention/

Field of Information and Telecommunications in the Context of International Security (UN GGE) confirmed that international law applies to cyberspace. …

To make significant progress, we have to unmask the fact that unfortunately there is little specificity in the agreements reached so far. This situation allows states to continue to act in violation of established norms, without the international community having any recourse to respond. For example, international law prohibits the use of force by states except in self-defense in response to an armed attack, and the UN GGE norms call for states to refrain from international malicious activity. …

Our proposed response was a Digital Geneva Convention, that would commit governments to adopt and implement norms that have been developed to protect civilians on the Internet, without introducing restrictions on online content. Just as the world's governments came together in 1949 to adopt the Fourth Geneva Convention to protect civilians in times of war, a Digital Geneva Convention would protect citizens online in times of peace."

**ForeignPolicy.com (2018).** *In Cyberwar, There Are No Rules — Why the world desperately needs digital Geneva Conventions*[356]*:*"The great challenge for military and cybersecurity professionals is that incoming attacks are not predictable, and current strategies for prevention tend to share the flawed assumption that the rules of conventional war extend to cyberspace as well. Cyber warfare does have rules, but they're not the ones we're used to — and a sense of fair play isn't one of them. Moreover, these rules are not intuitive to generals versed in fighting conventional wars.

"That's a problem because cyber war won't be waged with the informed participation of much of the U.S. technology sector, as the recent revolts at Google over AI contracts with the U.S. Defense Department

---

[356]   Ibid.,   https://www.weforum.org/agenda/2017/12/why-we-urgently-need-a-digital-geneva-convention/

and at Microsoft over Office software contracts with U.S. Immigration and Customs Enforcement demonstrate. That leaves only governments and properly incentivized multinational corporations to set the rules. Neither has yet provided a workable and operational definition of what constitutes a globally recognized act of war — a vital first step in seeking to prevent such transgressions."

**Atlantic Council (2019).** *It's time for a Cyber Geneva Convention*[357]:

"The DoD strategy lays out five mission objectives that should be the framework for creating cyber doctrine:

- Ensuring the joint military forces can achieve its mission in a contested cyberspace domain
- Enhancing Joint Force military advantage through the integration of cyber capabilities into planning and operations
- Deterring, preempting or defeating malicious cyber activity targeting US critical infrastructure that is likely to cause a significant cyber incident
- Securing DoD information and systems, including non-DoD-owned networks against cyber espionage and malicious cyber activity
- Expanding DoD cyber cooperation with allies, partners, and private-sector entities

"Before we can establish rules of engagement for cyber warfare, we must first establish conventions for the use of cyber weapons. Our adversaries are mounting offensive cyber operations daily and due to a lack of guidelines regarding operations in cyberspace, there is little protocol defining what the appropriate response is. Similar to the outcomes of the Geneva Convention, the world needs new international rules to protect the general public from nation-state threats in cyberspace."

---

[357] Leigher, Bill. 2017. It's time for a cyber Geneva Convention. Atlantic Council, https://www.atlanticcouncil.org/blogs/new-atlanticist/it-s-time-for-a-cyber-geneva-convention/

**National Defence Magazine (2020).** *Geneva Conventions for Cyber Warriors Long Overdue[358]:*

"Cyber warfare is a fact of the modern world. However, there is no clear international law that distinguishes between warfare, terrorism, crime or vandalism. As a result, U.S. military cyber warriors are operating without the protections and restrictions their kinetic brethren enjoy under the Geneva Conventions. The road to those agreements was long, but necessary, and it needs to be trod again — before civilians suffer the consequences of unrestricted cyber warfare."

**Lawfare Blog (August 2021).** *Responsible Cyber Offense:*

"Governments that harbor cyber criminals, or themselves engage in criminal behavior, may not see a shared interest in limiting damage. But this assumes that there is little risk that sloppy or unrestrained cyber operations could cause the target to escalate — intentionally or not — or could turn increasing numbers of countries against the states whose hackers wreak havoc. The concepts discussed in this post will not ameliorate blatantly dangerous behavior in the near term. But they would clarify what the U.S. considers to be an irresponsible activity, moving the nation away from a murky model of outrage at every Russian phishing email. By articulating and promoting the discussion of responsible operations, the U.S. could gain international political leverage.

Admittedly, it will take a certain hardheadedness and even cynicism among U.S., Russian and Chinese leaders to discuss best practices in malware development and placement, but this is the nature of diplomacy in the 21st century. Major powers bear responsibility for reducing systemic risk in cyberspace, and to do this they must make offensive operations more predictable. Each country wants to expel spies from its computer networks, and each will struggle to design better defenses against

---

[358] Waugh, Steve. 2020. Geneva Conventions for Cyber Warriors Long Overdue, National Defense Industrial Association, https://www.nationaldefense magazine.org/articles/2020/3/18/geneva-conventions-for-cyber-warriors-long-overdue

cyber operations. But technical panaceas are unlikely. Better to create codes of honor among spies, and their bosses."

## 28.3 Will Hackers and Nation States Follow the Rules?

Whenever I post any article or blog on LinkedIn *about this specific topic[359]*, a common question that gets debated involves whether this is just political talk with no action and enforcement. For example, when I posted this last article from the Lawfareblog, here were some of the responses:

*Mike Moran, Digital Marketing Consultant with Content and Digital Expertise[360]*: "It will be fun to watch the first time that some part of a world power gets hacked and then that drives hacking attempts at other world powers, with all parties hacking back. Yes, I worked on a situation where a client called us in just in case ... they didn't think they had any problems but wanted it documented. Turns out they had been hacked and were being used as base for additional attacks ... fun stuff."

*Jack Kufahl, Chief Information Security Officer at Michigan Medicine: University of Michigan[361]*: "I am dubious of how effective incre-*me*ntal cyber offensive tactics would be in the end. The 'low bar' for how offending nation-states and actors are impacting critical infrastructure means to me that those efforts need to be focused on the rigor and common sense protections as they not only help protect against adversarial attacks, but good processes and system integrity also protects from user and configuration errors. Certainly, as a component of war, there is an

---

[359] Lohrmann, D. Lawfare blog, Responsible Cyber Offense, LinkedIn post, https://www.linkedin.com/posts/danlohrmann_responsible-cyber-offense-activity-6828270384291340288-JL8_

[360] Lohrmann, D. Lawfare blog, Responsible Cyber Offense, LinkedIn post, loc. cit.

[361] Ibid.

increasing relevance on cyber as a natural evolution of disrupting the enemy's ability to command and coordinate, however it is not a well-defined mode of attack with clear thresholds for internationally accepted countermeasures and retaliation. Internationally matters since we are largely an integrated global economy, so the disruption or destruction of services in one country do impact others and have unintended consequences. I wish it was as simple as 'hack Vlad back' but there is no satisfaction in short-term wins in this field while the infrastructure at home is having problems with local passwords and unpatched Windows 2000. …"

*Jim Angleton Chief Executive Officer at United Police Federal Credit Union and Ambassador US/UN KoM, Sovereign State[362]*: "We are big proponents to legalization (acknowledging) of cyber deterrence. Meaning, offensive and defensive strategies protecting and defending systems. While many conduct and use protocols of same, it is good to know that it has been addressed, accepted, legalized and recognized that if you hit a company, individual or government, be prepared to punch back, hard! We do and our clients are better off for it."

## 28.4 Final Thoughts

This topic has been important to me personally since I wrote the book Virtual Integrity: Faithfully Navigating the Brave New Web back in 2008. At that time, I wrote that federal government needs to appoint an "ambassador to cyberspace" to deal with the international nature of these vital topics.

I have written several articles over the past few months regarding international relations and cybersecurity, and the *Ransomware Task Force[363]* covers the importance of this topic in several ways.

---

[362] Ibid.

[363] Ransomware Task Force (RTF), Institute for Security and Technology, 2023, https://securityandtechnology.org/ransomwaretaskforce/

Here are two of those articles: *Biden Sets Cyber Standards for Critical Infrastructure*[364]; *NATO Adds Cyber Commitments, Potential Ransomware Response*.

Bottom line, I just don't see major progress regarding cyber attacks escalating without substantial international cooperation. Whether we call it a "Digital Geneva Convention," or something else, urgent action is needed now.

---

[364] Lohrmann, D. 2021. Biden Sets Cyber Standards for Critical Infrastructure. Govtech.com, Blogs, https://www.govtech.com/blogs/lohrmann-on-cybersecurity/biden-sets-cyber-standards-for-critical-infrastructure

<div align="center">

**29**

</div>

<div align="center">

# DATA STORAGE ETHICS: SECURITY AND RESPONSIBILITY AT ALL LEVELS

</div>

<div align="center">

*Gilles Bach, Switzerland*[365]

</div>

## 29.1 Data Explosion

Digital Data is today at the center of our lives. Personal lives. Professional lives. Our lives as citizen. Through the internet and its myriad of websites and applications used by about 5 billion users worldwide[366], we all leave daily - whether we are aware of it or not – dozens to thousands of digital footprints in this vast and complex digital mechanism and world called cyberspace.

---

[366] Digital 2022 Global Report. 2022. DataReportal, https://datareportal.com/reports/digital-2022-global-overview-report.

Even when we enter a physical venue and we interact and transact with a brick-and-mortar store, business, corporation or public organization for any type of purpose we are creating electronic records in their systems most of the time.

The overall quantity of data being generated is exponential. In 2021[367], 79 zettabytes of data were created globally (79 followed by 21 zeros). This has to be compared to only two zettabytes in 2010, about 16 in 2015, 175 being forecasted in 2025, and 570 potentially estimated in 2030[368]…

These huge amounts of data, and the way they are collected and in particular stored before, during or after their processing and uses, generate in return important ethical responsibilities towards the individuals and the organizations to which they belong as well as their surrounding environments and eco-systems.

As we will see these responsibilities encompass a very large range of dimensions at very different levels that have to be considered and that are affected in many various ways: technology itself, geopolitics, privacy, law, human rights, IT security, governance and organization, sovereignty, economics, sustainability and even sociology and philosophy.

This is what we will scrutinize in the subsequent parts of this article by setting the various contexts, explanations, examples and possible solutions to raise the awareness and self-reflection, improve the understanding, and ignite the implementation of these fundamental ethical responsibilities.

---

[367] See from 2010 to 2020: Petroc, Taylor. 2022. Amount of data created, consumed, and stored 2010-2020, with forecasts to 2025, Statista, https://www.statista.com/statistics/ 871513/worldwide-data-created/

[368] See beyond 2025: Balnojan, Sven. 2020. The Future of Good Data — What You Should Know Now!, Towards Data Science, https://towardsdatascience. com/the-future-of-good-data-what-you-should-know-now-f2a312a0e469

## 29.2 Data Growth: Four Factors

As we have seen, the curve of the yearly data generation is following an exponential trend. This phenomena has been actually facilitated by four factors that started to reach a sufficient level of sophistication around the end of the first decade of 2000 to enable the accelerated development of increasingly data -heavy, -hungry and -related technologies and applications: more powerful and affordable computing power, higher and significant storage capabilities, sophisticated software and data architecture frameworks, and last but not least the spread of fast and reliable internet services.

In terms of data types being generated, the current and continuous growth of data volume comes from the multiplication of videos and photos, the rise of social media, the ever increasing level of digitization of businesses and organizations in the world across all sectors and geographies, and from the explosion of smartphones, IoT (Internet of Things) devices, sensors, machines and vehicles capturing data.

80% of the data growth is due to unstructured type of data[369] (mostly pictures and videos) as compared to structured data (inputs and information to fill forms and complete most transactions and processes) and takes therefore much more storage space.

The Covid pandemic lockdowns did obviously accelerate this digitalization thanks to a forced transformation of many physical processes and types of transactions that became banned and inapplicable literally overnight across all areas of life and business. Also, the expected development and increasing sophistication of big data, advanced data science, automation and AI algorithms requires more data for both systems training purposes and generating accurate results - and this is only the beginning.

---

[369] Artificial Intelligence. To Unlock the Hidden Value of Unstructured Data. NRoad, accessed from the blog VentureBeat, https://venturebeat.com/data-infrastructure/report-80-of-global-datasphere-will-be-unstructured-by-2025/

Not forgetting about the upcoming Metaverse developments and machine-generated data that will increasingly contribute to this expansion. But let's start first with some history summarizing the evolution of storage mediums.

## 29.3 Data Storage: History, Types and Evolution

As a reminder – even if we tend to forget about it – data and information can be obviously stored on traditional types of mediums, and in particular on paper. This makes this data much less accessible and much more difficult to retrieve and to leverage. Which is not really compatible with the current era of digital transformation and digitization and limits all the benefits that can be derived from analyzing and using data more systematically and in a faster and automated way. On the other hand, in a time of frequent digital data hacking and breaches, this can make such data potentially more secure as not everybody can access it.

The physical nature of this medium also necessitates a significant physical footprint and is particularly vulnerable to calamities. The fate of the famous Universal Library of Alexandria, one of the largest and most significant library of the ancient world reminds us of the particular fragility of this type of storage in the event of calamities like fire, floodings or similar. Hence this medium only being limited nowadays to very specific niches and usages.

Then, in the late 1950s mainframe computers started to appear and used sets of punched cards, paper tape, or magnetic tape to store and transfer data and programs. This is how IT started in some large corporations with such systems installed in big server rooms that were connected to dozens or more terminal stations used by employees to access or input data into the system. This corresponds to the image that we can see in some movies from the seventies in particular.

At the end of the seventies and in the eighties, micro-computers, later also called personal computers (PC) where on the rise and permitted

individual users, at work or at home, to process information and use software decoupled from any server for accessing and storing data and programs, with no need for any communication network. They were equipped with independent local medium storage systems like tapes, floppy discs or hard drives.

In the 1990s the appearance and rise of internet in a mainstream fashion brought the interconnectedness between local computers and remote web servers operating websites and applications as well as file sharing servers and protocols.

It is only in 2006 that the current notion of Cloud computing started to emerge when Amazon launched[370] a cloud storage service available from anywhere on the web, helping 'free developers from worrying about where they are going to store data, whether it will be safe and secure, if it will be available when they need it, the costs associated with server maintenance, or whether they have enough storage available'.

This was at the very origin of this huge rise that brought the spending on public cloud services globally from scratch to an estimated 495 billion[371] USD in 2022 as compared to 411 in 2021 and a forecast of 924 billion by 2027.

The share of corporate data stored in the cloud increased from 30% in 2015 to 60%[372] in 2022 with 3.6 billion cloud users in the world.

---

[370] Initial Amazon Web Services (AWS) cloud storage service launch press release, on March 14, 2006 - https://press.aboutamazon.com/2006/3/amazon-web-services-launches

[371] Gartner Forecasts Worldwide Public Cloud End-User Spending to Reach Nearly $500 Billion in 2022. Gartner, Stamford, Press Release, 2022, https://www.gartner.com/en/newsroom/press-releases/2022-04-19-gartner-forecasts-worldwide-public-cloud-end-user-spending-to-reach-nearly-500-billion-in-2022

[372] Sujay Vailshery, Lionel. 2022. Share of corporate data stored in the cloud in organizations worldwide from 2015 to 2022, Statista, https://www.statista.com/statistics/1062879/worldwide-cloud-storage-of-corporate-data/

## 29.4 Data Storage in the Cloud: A New Era of Possibilities

Cloud storage and cloud computing are more precisely about the remote storage of data and the remote access to computing power and applications services over the internet, provided by a third party vendor usually labeled as a CSP – Cloud Services Provider. These services are available both for organizations and individuals. It replaces the model where the data and applications are stored and processed locally, either on a computer, or on on-premise servers.

There are globally 3 types of clouds. Private clouds, that are dedicated and operated solely for a single organization. It can be managed internally or through a third-party vendor. Public clouds use shared infrastructure to deliver services to many organizations or individuals. And Hybrid cloud is a storage infrastructure using a combination of public cloud services and on-premise or private cloud resources.

Public clouds are seen by organizations as 'appropriate only for less critical data and less sensitive workloads'[373], therefore the strategic choice to keep those particularly critical data sets within a private cloud when appropriate.

There are many benefits though for organizations to moving their data initially stored on their own managed, on-premise servers to the public cloud: the advantages relating to outsourcing, like reducing the directly managed and held resources in terms of IT hardware and infrastructures, implying Capex costs (capital expenditures) savings, as well as reducing IT staff costs. The CSPs do therefore offer in comparison much better prices thanks to the economies of scale of shared infrastructures. On top of that they are able to offer state of the art cybersecurity protection to the stored data, workloads and services, by applying the very best and

---

[373] Brinda, Mark and Kate Woolley. 2019. Public vs. Private Cloud? The Market Says Hybrid, Bain & Company, https://www.bain.com/insights/public-vs-private-cloud-the-market-says-hybrid/

elaborated, up-to-date practices and technologies with highly skilled specialists.

Cloud services are also enabling many new types of services supporting enhanced collaboration between workers and organizations all over the world, improving the speed and flow of information being exchanged, of transactions, the delivery of services from one end of the globe to another, including real-time operations, advanced data analyses, an increased level of big data and AI sophistication, and much more. All this in a very scalable, flexible, and cost-effective way, which confirms definitely Data as being the '4th production factor in this era of the fourth industrial revolution'[374] on top of the traditional factors that are natural resources, labour and capital.

## 29.5 Cloud Storage: Macro Risks

As we have seen the cloud creates tremendous benefits and opportunities for the development of businesses and organizations in the world. It can even be considered as the 'powerhouse that drives today's digital organizations'[375]

Electronic data storage requires however by definition electrical power to store and retrieve data. And the cloud requires access to a working, reliable and fast telecommunication network linked to the internet. This is of no surprise and seems very basic and obvious in our days. But at a time in which conflicts with worldwide repercussions are surging or threatening to surge again, raising concerns and tensions in terms of energy cost, and even about energy supply at all, this is not benign.

---

[374] Stückelberger, Christoph / Duggal, Pavan (Eds.) 2018. Cyber Ethics 4.0: Serving Humanity with Values, Global Series, Globethics Publications, 45, https://repository.globethics.net/handle/20.500.12424/169317

[375] Gartner Forecasts Worldwide Public Cloud End-User Spending to Reach Nearly $500 Billion in 2022, loc. cit.

And if we look at real or potential threats  - whoever perpetrating them – against central and strategic telecommunication infrastructures absolutely indispensable to most of our digitized ways of communicating, informing, transacting and delivering products and services, this adds even more to the concerns.

We must also be aware that 99% of the internet runs through a network of 420 submarine cables connecting all continents to each other. 'It is estimated that over USD 10,000 billion in financial transactions run today through these "seabed highways". This is especially the case of the main global financial exchange system, SWIFT (Society for Worldwide Interbank Financial Telecommunications). The security of these transactions is a political, economic, and social problem.'[376]

In the aftermath of the September 2022 Nordstream natural gas pipelines spectacular sabotage, and knowing that actual incidents and attempts on internet cables already happened in the past, a massive disruption caused by coordinated attacks could become a nightmare scenario for the overall world economy. The satellites that come to mind as a possible alternative for sustaining the internet connectivity in such an occurrence account actually only for 1% of the data exchanges. They are too costly and the connection speed is much lower. Not mentioning the fact that an old, defunct satellite has been destroyed in November 2021[377], just 3 months before the invasion of Ukraine, by an ASAT, meaning an anti-satellite missile test.

These two enabling conditions for the cloud, electricity and internet networks may also be put at risk by natural disasters, other human factors actions and even – though with a low but not impossible level of

---

[376] Rona Rita David, Submarine Cables: Risks and Security Threats, Energy Industry Review, March 25, 2022, https://energyindustryreview.com/analysis/submarine-cables-risks-and-security-threats/

[377] Foust, Jeff. 2021. Russia destroys satellite in ASAT test, SpaceNews, https://spacenews.com/russia-destroys-satellite-in-asat-test/

probability – by an event as distant as an extremely powerful solar flare which could potentially render all electronic devices inoperative.

These factors represent clearly a big weakness and risk in relation to the correct execution of data storage access and activities depending on the cloud. With 60% of the global GDP estimated to rely on digital communications in 2022[378], such a scenario would be absolutely disastrous. This has to do with a black swan type of event, with initially a very low probability of occurrence but a potentially enormous, disastrous impact. In the very current context however, we are just realizing that these options are not anymore merely a theoretical figment of imagination.

## 29.6 Cloud Storage, Local Storage, Metadata, Background Analysis and Processes

If you take many pictures with your smartphone you might have noticed that with models from the last 3-4 years, these photos are now very often labeled and tagged automatically thanks to machine learning image recognition algorithms. For example, category tags like food, landscape, sport, document, etc… or even people names can be generated and added automatically, as well as up to the geographic position and coordinates if your GPS is activated and if you enabled and set pro-actively the last two options in particular.

---

[378] Digital Development, World Bank. https://www.worldbank.org/ en/topic/digitaldevelopment/overview. Amid this war context, the European Parliament and Commission have respectively released an in-depth analysis on the security threats to undersea communication cables and infrastructure in June 2022 and a five-point plan in October 2022: European Parliament, Directorate-General for External Policies of the Union, Bueger, C., Liebetrau, T., Franken, J., Security threats to undersea communications cables and infrastructure : consequences for the EU : in-depth analysis, Publications Office of the European Union, 2022, https://data.europa.eu/doi/10.2861/35332

Some of this automatic tagging is done by default, like the general categories, some others need to be enabled. Such data is called 'metadata', meaning 'data that describes data', in this case data describing pictures.

This functionality can be very helpful to categorize pictures in your gallery, in particular if you make a lot pof them, and makes it easier later on to search pictures in a faster and relevant manner. At the same time this initiative to deploy automation of such functionalities could create a clear feeling of intrusion into the private sphere of people. This raises then questions about how this type of automated analysis of stored content might be accessible or leveraged by various layers of stakeholders.

Does this data remain locally stored on the smartphone? Is it somehow accessible and leveraged – for commercial purposes or other – by one or several suppliers of the digital chain: the application creator, the operating system of the device, the phone maker? Is it even shared further with third parties? What happens also to this information when photos are synchronized and backed-up in the cloud?

The question expands actually to the real-time or later analysis of any stored content, may it be stored on the cloud or stored locally on PCs or any device connected to telecommunication and internet networks. The messages in our e-mail boxes, social media profiles and accounts information even when set to private, and of course documents and productivity suites software solutions and applications either based on cloud solutions or even classically stored locally on a computer, etc... are concerned.

Very often such information and authorizations are part of the Terms & Conditions of the respective listed types of suppliers. But who reads and fully understands those very long types of texts? According to a study the 'combined terms and conditions of 13 top apps including Tik-Tok, WhatsApp and Zoom would take 17 hours and five minutes to read'. On top of that a privacy researcher mentioned that 'both the length and language used in such statements made them difficult for adults and

children to be able to make informed consent about what they were agreeing to'[379].

There are exactly the same kind of lack of clarities with specified authorized actions granted to installed mobile applications that often do not make any sense given the initial purpose and functionalities of the program, with no possibility of deselecting one single permission and providing no information about what exactly is collected or not, and to which purpose and type of processing.

With such information being too often burried in the middle of complex and boring statements, or not available at all in a detailed manner, it just reflects, irrespective of pure legal aspects, a kind of lack of transparency that the industry in general should address to re-establish and increase trust and respect between users, citizens, businesses and suppliers along the data storage and processing chain. Upgrading these standards, showing genuinely more transparency while enabling a true and better understanding of what is done or not with our data will be ultimately beneficial to all.

## 29.7 Information Security Principles and International IT Standards

Ensuring cybersecurity, and in particular the protection of digital data is obviously one of the major ethical responsibility for any organization managing and storing sensitive or voluminous amounts of data and systems and their leaders.

When it comes to computer data security, there are three fundamental principles and dimensions of data - adopted by all major IT Security standards and certifications - that need to be safeguarded at anytime: their confidentiality, their integrity and their accessibility. ISO 27000

---

[379] Kleinman, Zoe, Popular app T&Cs 'longer than Harry Potter', November 2020, https://www.bbc.com/news/technology-54838978

standards on IT Systems Security Management do cover the necessary application of those principles.

Ensuring Confidentiality in this context is about controlling and preventing the unauthorized use, disclosure of information or its unauthorized access.

Integrity is about ensuring that the data has not been modified, altered or deleted in any way, without permission, justification and in a monitored way. Availability is about ensuring a reliable and timely access to data whenever needed.

These principles do cover both the prevention and fight against external cyberattacks as well as against insider threats and non-authorized internal accesses.

The full and strict implementation and respect of those principles are absolutely vital from an ethical standpoint. Steve Wozniak, Apple's co-founder, raised warnings on those risks and on the fact that "the more we transfer everything onto the web, onto the cloud, the less we're going to have control over it"[380] as soon as 2012, a time when the cloud just started to spread more into mainstream users.

## 29.8 Data, Privacy, International Law and Human Rights

In terms of international laws and regulations on cybercrime and data protection, the 2001 Budapest Convention paved the way by defining and clarifying the criminalization of conducts relating to crimes committed in or through the cyberspace[381]. This international treaty served as a

---

[380] Meyer, David. 2012. Wozniak: 'I really worry about everything going to the cloud', Zdnet, https://www.zdnet.com/article/wozniak-i-really-worry-about-everything-going-to-the-cloud/

[381] Stückelberger, C., Duggal, P. 2018. Cyber Ethics 4.0: Serving Humanity with Values, op. cit., 368.

base to develop domestic laws around the world, as well as building up international cooperation on the matter.

It resulted in the production of many cyber criminal laws, and subsequently also in the development of more specific national data protection and privacy laws on the various continents.

As per the UN Conference for Trade and Development (UNCTAD), 137 countries out of 194 had put in place as of the end of 2021[382] legislation to secure the protection of data and privacy. Privacy is here a crucial notion that is close to the principle of confidentiality already mentioned but that relates more specifically to the confidentiality of personal data.

It refers fundamentally to the individual's right to control and maintain their own data, underpinned directly by the article 12 of the 1948 Universal Declaration of Human Rights on privacy, that has been detailed even further in the International Covenant on Civil and Political Rights (ICCPR) multilateral treaty from April 1988[383] in the following way: 'Article 17 provides for the right of every person to be protected against arbitrary or unlawful interference with his privacy, family, home or correspondence as well as against unlawful attacks on his honour and reputation. In the view of the Committee this right is required to be guaranteed against all such interferences and attacks whether they emanate from State authorities or from natural or legal persons. The obligations imposed by this article require the State to adopt legislative and other measures to give effect to the prohibition against such interferences and attacks as well as to the protection of this right.'

---

[382] Data Protection and Privacy Legislation Worldwide, United Nations Conference on Trade and Development, https://unctad.org/page/data-protection-and-privacy-legislation-worldwide

[383] UN ICCPR Source: UN Human Rights Committee (HRC). 1988. CCPR General Comment No. 16: Article 17 (Right to Privacy), The Right to Respect of Privacy, Family, Home and Correspondence, and Protection of Honour and Reputation, https://www.refworld.org/docid/453883f922.html

## 29.9 The EU GDPR Regulation: Ambitious Data Privacy Standard and Model

In the context of the rising digitalization of the world, generation of data and related services, the European Union decided to upgrade in 2016, its previously existing Data Protection Directive from 1995 to adapt it to the new, evolving technological landscape, usages and challenges since the emergence of the internet and the cloud.

The General Data Protection Regulation became enforceable in May 2018 and has a primary focus on protecting personal data and privacy. It is considered as the toughest, most comprehensive and advanced data privacy and security law in the world.

Historically, the right to privacy is part of the 1950 *European Convention on Human Rights*[384], and the development of such legislation represents the natural manifestation of fundamental values that are common across European societies.

This regulation rules the individuals' controls and rights over their personal data, as well as data obligations for organizations if the data controller (the organization collecting the data), or the data processor or the data subject is based in the EU. As processors, the Cloud Service Providers operating in the EU and/or managing personal data of EU citizens leaving in the EU are fully subject to the law.

In a nutshell, the GDPR clarifies the following rights for the data subjects over their data: right to consent, right to be informed, right to access, rectify or erase their information, restriction of processing, data portability to list the main ones.

It also sets specific, new principles and requirements for the storage and processing of personal data. On top of implementing the three fundamental IT Security Management principles - Confidentiality, Integrity

---

[384] European Convention on Human Rights, Official texts, https://www.echr.coe.int/Pages/home.aspx?p=basictexts&c=

& Availability, two new principles with an impact on data storage are introduced: data minimisation and data storage limitation.

Data minimisation imposes the collection and processing of data only strictly necessary to fulfil the announced purpose. Data storage limits means that organizations should not keep data for longer than needed. Although there are no specific time references indicated within the regulation organizations have to create and implement a data retention policy and to perform periodic reviews to identify, and address, data stored beyond intended use.

Another particular key requirement is about pseudonymization of personal data, meaning de-identifying personal data, e.g. by replacing it by artificial identifiers. Encryption, anonymization and data obfuscation can be part of that same process.

Strong requirements in terms of global Cybersecurity and informing quickly and systematically data subjects and supervisory authorities in the event of data breaches are as well - among many more requirements - part of this very exhaustive framework.

Interestingly, GDPR enforcement is based on the principle of Accountability. Organizations have to take responsibility for how they process and store personal data. They have the obligation to demonstrate compliance, meaning keeping all the necessary records and document evidences to show that they have created, implemented and applied all the related policies, processes and routines to comply with the law.

For that matter, in the case of 'regular and systematic monitoring of data subjects on a large scale' it is even mandatory to designate a Data Privacy Officer (DPO) that will be in charge of overseeing the internal compliance with the law.

The last layer of the law is logically about fines and penalties, with non-compliant organizations or at the origin of data breaches, facing sanctions of up to EUR 20 million or 4% of their global sales, whichever is greater.

Overall, the very comprehensive scope of the GDPR, its forward-thinking principles and requirements, its very demanding enforcement rules and the underlying values that led to its writing actually represent an outstanding model to be replicated in other countries and parts of the world to care and protect individuals and their privacy.

Many of these principles and requirements can also be leveraged for inspiring general approaches to protect critical data, even non-personal data, and to build more resilient and better protected data systems and data storage systems. Following such a model even as an organization not subjected at all to the EU context would be the manifestation of a great ethical responsibility and respect towards users and stakeholders and becomes then truly part of the scope of Corporate Social Responsibility policies and strategies.

The fully detailed GDPR regulation and its related explanations are available on a specifically, dedicated EU portal[385]: including numerous relevant checklists[386].

## 29.10 Data Sovereignty

In a sensitive geopolitical context, competition between nations - in economics, politics and influence - is getting even fiercer in the era of digitalization and upcoming developments in data science and AI technologies, with sophisticated use cases expected to gradually and exponentially appear in the next 5 to 15 years.

Data is the fuel for all these technologies with ever-increasing impacts and becomes therefore a strategic asset for any organization and for any country. This confirms even more the status of data as the new, fourth production factor, as we have already seen.

With the possible reconstitution of geopolitical blocs with very opposed views of the world and even the come-back of the war and materi-

---

[385] Complete guide to GDPR compliance, https://gdpr.eu/

[386] The GDPR full text, https://gdpr.eu/compliance/

alisation of conflicts in various parts of the world, it becomes then a strategic imperative to protect the data of a state, of it's individuals, and of it's economic agents and institutions.

Data sovereignty was already and becomes then an even more critical question.

One of its translations is the legislative concept that data an organization collects, stores, and processes are subject to the laws, regulations and governance structures of the country where the data is collected, where the data is stored and/or relating to the nationals or organizations of that state. The notion is complex with various definitions, various scopes and various obligations linked to this type of laws according to the country.

The European GDPR regulation – already mentioned – is a major example of data sovereignty law as it rules the individuals' controls and rights over their personal data and data obligations for organizations if the data collector, or the data processor or the data subject is based in the EU. This also applies then to organisations based outside the EU if they collect or process personal data of individuals located inside the EU.

Another example, even anterior and with even more implications at global levels, is the US Patriot Act introduced in October 2001. Under this act, officials were granted possibilities of access to any information physically within the United States, regardless of the information's origin. The 2018 CLOUD Act, clarified and extended the possibilities of access by compelling 'U.S.-based technology companies via warrant or subpoena to provide requested data stored on servers regardless of whether the data are stored in the U.S. or on foreign soil.'[387]

This means that any information collected by an American server could also be potentially accessed, even, as specified, if the server was not located in the United States but owned, managed or handled by any

---

[387] Swire, Peter. 2018. The CLOUD Act and its Impact on Cross-Border Access to the Contents of Communications, https://www.alstonprivacy.com/cloud-act-impact-cross-border-access-contents-communications/

American organization somewhere along any part of the data collection, storing or processing chain.

This raises obviously questions and concerns, beyond the initial, appropriate purpose of this law, regarding potential direct or indirect abuses or misuses of such possibilities of accessing data for various reasons, and in particular economic competition. In that case how about the cloud storage offerings and the productivity tools in the cloud, proposed and hosted by providers originating from the United States?

These services offer tremendous possibilities in terms of storage, productivity and efficiency both to companies, organizations and individuals all over the world.

The EU GDPR actually regulates most of these questions, at least regarding data relating to individuals. This regulation also became a model for about a dozen other countries from various continents. In other countries there might be - or not – existing laws and regulations covering even partially those questions. If not, the respective Terms of Service defined by each provider along the data flow and storage chain do apply. While most of the time the data is encrypted from end-to-end with strong encryption technologies, the more secured practice for critical data and corporate clients is to manage encryption keys independently[388] from the cloud service provider where the data is located. The data owner owns in that case the encryption keys.

The localization of data becomes then an important factor, knowing in addition that not every country might have currently existing or sufficient local data centers capabilities installed and local cloud services providers.

On top of that the global concentration of cloud providers is quite impressive, with AWS (Amazon Web Services), Microsoft and Google

---

[388] Musthaler, Linda. 2013. Cloud encryption: control your own keys in a separate storage vault, Network World, https://www.networkworld.com/article/2170564/cloud-encryption-control-your-own-keys-in-a-separate-storage-vault.html

representing 66% of the global cloud spending[389] in the third quarter of 2022 coming from 61% one year earlier.

There are also so-called data localization or data residency laws, which are intended to keep in particular citizens' personal data in-country and subject to local regulation. This type of laws are often initially thought and enacted to protect citizens data privacy from other laws and regulations abroad that might not be as strict as the ones in the initial country or area. Data dissemination, in its broad meaning, and the fact that data is often replicated partially or totally across various places or data centers does indeed not help to ensure full transparency on what happens to this data and what are the related risks and consequences.

Some countries do however implement such laws for motivations relating more to economic protectionism, political reasons, or even directly – openly or not – surveilling their own citizens and exerting information and data control, including censorship. Anupam Chander even speaks of 'Data Nationalism'[390] (2015).

These laws may also create difficulties from various perspectives: conflict between political motivations and technical and operational efficiency (access, reliability, security, energy), conflicting or overlapping data sovereignty laws, complexifying the data flows internationally and thus complexifying the delivery of international services, up to feeding the 'ongoing struggle between democracy and totalitarianism'. In the latter case the 'business versus ethical' dilemma takes on its full meaning for global technology players in the field of data cloud storage and services when they operate in countries with lower standards in terms of democracy or human rights. What is acceptable, what is not? 'Should I

---

[389] Haranas, Mark. 2022. Top 5 Cloud Market-Share Leaders: AWS, Microsoft, Google In Q3 2022, https://www.crn.com/news/cloud/top-5-cloud-market-share-leaders-aws-microsoft-google-in-q3-2022

[390] Anupam Chander & Uyên P. Lê, Data Nationalism, 64 Emory L. J. 677, 2015, https://scholarlycommons.law.emory.edu/elj/vol64/iss3/2

stay or should I go?' Should we accept abusive requirements putting, directly or indirectly, fundamental human rights and freedom at risk?

This is why on top of national data privacy and protection regulations that need to become even more widespread and exhaustive in their content, and in order to avoid and manage conflicting and overlapping laws, working towards an international agreement on the matter is a vital task. Both for protecting citizens, human rights, organizations, and for providing a clear path to a safe technological and economical development based on data as a resource for progress.

## 29.11   Health Data Ethics: History of a Precursor Domain

Health Data is one of the most sensitive type of data to be collected, handled and stored. It touches directly the most intimate part of our physical incarnation as well as direct or indirect psychological and behavioural facets of who we are and what we are facing in this life.

Like all the other types of data, health related digitalized data, meaning patients and healthcare data up to clinical trials, genetic information and now also the new sphere of data generated by personal or medical connected devices and related applications, are exploding in terms of volume being generated.

This type of data is obviously highly sensitive and as such needs to remain fully confidential like the Hippocratic Oath already stated it very clearly 2.400 years ago: 'And whatsoever I shall see or hear in the course of my profession, as well as outside my profession in my intercourse with men, if it be what should not be published abroad, I will never divulge, holding such things to be holy secrets.'

The health sector has therefore been a precursor in many aspects of data protection and privacy.

From the Oath that many students are swearing to actual national and local regulations and policies relating to medical privacy, the health

sector is probably one that has among the highest sense of awareness and established processes in terms of data protection and careful use.

A list of some notable national health and medical privacy frameworks and initiatives of interest:

- The US HIPAA[391], Health Insurance Portability and Accountability Act of 1996, sets standards for uses and disclosures of *protected health information (PHI)*[392], and provides civil and criminal penalties for violations.

- The EU GDPR regulation classifies health data not only as a personal data, but as a sensitive personal data[393], among six other categories like e.g. ethnicity, religion, political opinions, genetic information, and biometric information about an individual. These sensitive personal data categories can initially not be recorded, treated and stored outside of a list of very specific exemptions, or the direct, explicit consent of the person. The list of exemptions relating to health data is obviously directly related to the operational needs for delivering medical and health services.

- NHS Digital, the IT & Data Department of the NHS (National Health Service) in UK applies also the GDPR and created purposely a widespread and systematic training program[394] requiring the NHS staff to complete appropriate annual data protection and security training and pass a mandatory test

- On top of applying the GDPR regulation as a member of the EU, France has developed an official mandatory governmental label,

---

[391] HIPAA Portal: https://www.hhs.gov/hipaa/index.html

[392] Protected health information, Wikipedia.org, https://en.wikipedia.org/wiki/Protected_health_information

[393] Art. 9 GDPR Processing of special categories of personal data, https://gdpr-info.eu/art-9-gdpr/

[394] Data Security Standard 3 - Staff training, https://digital.nhs.uk/cyber-and-data-security/guidance-and-assurance/data-security-and-protection-toolkit-assessment-guides/guide-3---staff-training

'Health Data Hosting'[395], to certify specifically each single health data hosting provider. It is based on the ISO 27001 standard requirements for Information Security Management Systems, with a 3-years duration and a yearly intermediate audit.

- France has also developed in 2019 the concept of a centralised, public 'Health Data Hub' (HDH)[396], storing or more precisely aggregating copies of many public health data sources for research purposes and the aim of improving the quality of care and treatment. The objective of the HDH is to 'enable project coordinators to easily access non-nominative, unified data on a secure platform, in compliance with regulations and citizens' right. They will be able to cross-reference and analyse the data in order to improve the quality of care and patient support'.

On this HDH platform, data sources are as diverse as hospital data, social security data, clinical data, death registers, public health data, national blood bank database, specific diseases monitoring, etc.

The idea is to leverage and cross-pollinate this precious data through exploratory analysis, data science statistical techniques and AI algorithms to foster medical innovation and discovery of hidden correlations or patterns, e.g. improving symptoms and diseases detection, providing answers to rare pathologies, predicting individual patient trajectories, improve pharmacovigilance and drug safety, etc...

The projects accessing and leveraging the data are listed publicly online in a fully transparent manner. 'The data, within a well-defined scope, are accessible to project coordinators contributing to the public interest, following an approval process involving in particular *a scientific committee, an ethical committee,* and the National Commission for Data Protection and Liberties (CNIL).'

---

[395] Health Data Hosting (HDS), Ministère de la Santé. Agence du numérique en santé, https://ue.esante.gouv.fr/information-systems-security-pre-condition-trust/health-data-hosting-hds

[396] Health Data Hub, https://www.health-data-hub.fr/page/faq-english

This really shows that it is actually possible with very strict protocols, processes and controls at the organizational, administrative and technical levels, from design to implementation, to both combine the realization of tremendous value, derived from massive, diverse sources of initially very sensitive data, and all this in the full and plain respect of regulation, numerous ethical considerations and individuals rights.

The very demanding level of requirements and specific characteristics in terms of data ethics and storage in the health sector could represent a model to inspire many other sectors, starting from raising the awareness very early in the education cursus to diligent rules and processes at all levels of the data collection, storage and processing chain, and all this under close monitoring and scrutiny by specific ethical committees. The idea is not necessarily to replicate 100% of all these practices but to have a high standard model that can be a starting point to set best practices within an organization, an industry, or a certain geography.

## 29.12 Data Governance: Improved Data Practices and New Ethical Developments?

In a context characterized by this exponential growth of data volumes, rapidly evolving technologies and use cases, the complexification of demanding compliance requirements,  increasing pressure due to cyber-risks, and the multiplication of data storage types and locations, having a clear view and deep understanding of the overall data owned by an organization is not a luxury.

This is why Data Governance is of structuring importance.

Data Governance is commonly referred as a 'data management function to ensure the quality, integrity, security and usability of the data collected by an organization'[397] through its complete lifecycle. This

---

[397] O'Reilly. 201. Data Governance, the Definitive Guide.

involves also ensuring the compliance and conformity to regulations. This definition, its scope and reach is however still somehow loose and not yet completely mature and definitely shaped within the industry.

Being in complete awareness, knowledge, understanding and control of all the data owned, stored and managed by the organization should be a primary objective of any organization through it's data governance structures and policies.

This might be implied but is currently very often not sufficiently formalized and applied. And even less an actual reality within organizations. Establishing a plain, comprehensive, structured Data Governance (DG) function and program is still too often overlooked in many corporations.

According to a 2022 study[398], only '62% of organizations have had a data governance program in place for at least five years.'

Many organizations do not have complete, formal DG structures, roles, teams, policies. This creates all the conditions for possible lacks or gaps in terms first of accountability[399], and then of complete data understanding and mapping, leading in particular to loopholes in terms of correct and proper data categorization and protection.

Very often simple things like just having a complete, up-to-date data catalog mapping all the data of the organization, and data dictionaries specific to all databases and systems allowing the understanding and correct categorization of data are very rarely existing, up-to-date or complete. It is true that there are so many systems and data everywhere, in particular in big corporations that it becomes rapidly complex to implement and maintain such inventories systematically.

---

[398] The Strategic and Evolving Role of Data Governance in 2022 and Beyond, an Enterprise Strategy Group (ESG) Institute Study: https://www.mega.com/press-releases-strategic-and-evolving-role-data-governance-2022-and-beyond

[399] A data governance program must address accountability, in 'Designing Data Governance Structure', Journal on Computing (JoC) Vol.2 No.4, January 2013, http://dl6.globalstf.org/index.php/joc/article/download/576/2106

Very fortunately specific software solutions can help, as well as the application of the GDPR that forces many organizations that are concerned to rethink their data management policies and organization.

On top of creating and formalizing DG structures, roles and policies, like a Data Governance Council or Steering Committee, and formalizing further the usual data management roles from an IT perspective, it can make sense to create across the non-IT functions, departments and teams specific data stewards roles in charge of the very specific data linked to their daily operations, as they are in the best position to both understand this data and to control it. This is a 'simple' practice but still not sufficiently implemented. It must be a shared responsibility across departments.

Adding when possible an independent ethical committee in charge of advising the organization on the ethical dimensions of their ongoing operations or upcoming projects in terms of data handling and storage, beyond the pure compliance with existing regulations, is clearly another sound and recommended practice.

Considering Data Governance through an holistic and forward-thinking lens is also a great opportunity to add new responsible criteria and principles like Data Sustainability on the agenda in order to create and fully implement specifically dedicated policies (see subsequent section).

## 29.13 Data Long-term Sustainability

According to a report from the International Energy Agency (IEA), energy 'demand from data centers and data transmission networks accounted each in 2021 for 1 to 1.5% of global electricity use'[400]. These two elements are the main parts of the cloud services infrastructure

---

[400] Data Centres and Data Transmission Networks, EIA: https://www.iea.org /reports/data-centres-and-data-transmission-networks

chain. This equates to about 1% of global energy-related greenhouse gas emissions (GHG).

Despite the huge growth of data between 2015 and 2021 (data generation x5[401] and data storage x2.5[402]), data centers energy use is only estimated to have grown in the meantime by 10 to 60% (IEA). This is due to strong efficiency improvements relating to hardware, infrastructure and cooling optimizations as well as technical developments.

That being said, digitalization in general is both a source of CO2 emissions and at the very same time a part of the solution, as it also creates substantial savings through the avoidance of CO2 emissions.

Employees working from home (WFH) thanks to digitalization, with no commute necessary to work, can prevent the emission of about 270kg of CO2[403] on average a year per person for one WFH day a week up to 900 kg[404] a year for a full-time remote position according to various estimates.

Similarly '90% of business travel emissions (outside daily commutes) are from air'[405], meaning that using video-conferencing to replace short, medium and long-haul distance journeys whenever possible

---

[401]   Petroc Taylor, Volume of Data Created 2015-2025, Statista, https://www.statista.com/statistics/871513/worldwide-data-created/

[402]   Global Data Stored (IDC): https://www.businessinsider.com/sc/this-data-storage-solution-could-reduce-co2-emissions-2021-11?r=US&IR=T

[403]   French Environment Agency – 09-2020: https://presse.ademe.fr/wp-content/uploads/2020/09/ADEME_InfographieTT.pdf

[404]   See #7 in: Burgess, Andrea. 2023. Work from Home Models Predicted to Cause 34.3 Million Tons of Greenhouse Gas Emissionshttps://www.alliance virtualoffices.com/virtual-office-blog/work-from-home-co2-emissions-statistics/

[405]   Carbon budgets for corporate business travel: https://www.reuters.com/business/aerospace-defense/corporate-business-travel-carbon-budgets-loom-airlines-2021-10-10/

can help to dramatically reduce GHG emissions by a factor of up to 66 in the case of large conferences[406].

However, at the same time data, its processing and its storage are going to continue their exponential growth. And as every single industry, data management and storage activities have to contribute to the overall 2050 net-zero objective set in the Paris COP 21 Agreement.

As part of their Corporate Social Responsibility policies (CSR), Corporations are increasingly defining carbon budgets and planning carbon reductions attached to their activities. This is where these efforts of reduction meet.

If we consider the top three Cloud Service Providers (CSP) – Amazon, Microsoft and Google – they have all pledged[407] to become plainly carbon neutral for their cloud activities between now and 2040 at the latest. This is happening in particular through powering their data centers with renewable or decarbonized energy, through continued optimisation of their efficiency and thanks to the purchase of RECs (Renewable Energy Credits).

Organizations in general, storing their data on their own on premise servers and data centers might find it easier to rely on CSP providers through their pooling of resources – when possible - to reduce the carbon footprint of their data and to mitigate related environmental impacts.

Another strategy to reduce drastically data center's carbon emissions is to move 'cold data' (data rarely accessed and not needing instant, fast

---

[406] Two studies: Lynch, Jim. Video conference CO2 emissions quantified in new study, University of Michigan Ann Arbor, https://news.engin.umich.edu/2021/02/video-conference-co2-emissions-quantified-in-new-study/; Elizabeth Claire Alberts-Mongabay, Making conferences virtual or hybrid could significantly mitigate climate change, says new study. 2022. Green News, Euro News, https://www.euronews.com/green/2022/01/12/making-conferences-virtual-or-hybrid-could-significantly-mitigate-climate-change-says-new-

[407] Amazon, Google, Microsoft: Here's Who Has the Greenest Cloud, https://www.wired.com/story/amazon-google-microsoft-green-clouds-and-hyperscale-data-centers/

access) to less energy-consuming tape drive storage technologies, that can cut emissions by up to 95%[408] as compared to HDD technologies. Too much data is stored indefinitely and is kept "Just In Case" (JIC) or turns to "Write Once Read Never" (WORN) data status.

Cold data could represent about 60 to 75%[409] of stored data according to various industry experts. This means that there is an opportunity to save much GHG emissions by implementing ambitious data retention and storage tiering policies with a focus on archiving cold data on tapes, and also simply deleting them definitely if there is no related compliance requirements, and if there are insufficient to no reasons to keep it.

Driving digital data decarbonization implies to go to the very roots of the question, meaning reducing the volume of data to be stored in the very first place. Upfront, during, and after processing.

Data collection minimization, data purposes limitation, data retention policies in particular are principles – already mentioned – that have their part in a pro-active manner of managing this question, beyound any compliance and regulatory aspect.

Regular data cleanings and reflections on which data to delete would become best practices. Adding more systematic data deduplication (eliminating duplicate copies of repeating data) and data compression techniques to this list and to the focus on decarbonating data centers infrastructures and there is an holistic Data Sustainability framework and vision, with both technological and organizational levers that is emerging.

Last but not least, we must all be aware, as individuals, that we are also fully responsible for the generation of data and this exponential

---

[408] Here's how data centers can become truly green. 2021. Fujifilm with Insider Studios, https://www.businessinsider.com/sc/this-data-storage-solution-could-reduce-co2-emissions-2021-11?op=1&r=US&IR=T

[409] Active Archives and the State of the Industry 2020. 2020. https://activearchive.com/wp-content/uploads/2020/06/AAA-Annual-Report-2020.pdf

growth that is happening. Each average user out of the 5 billion people worldwide using internet creates about 1.7 MB of data every second[410]. We shall therefore think, self-reflect and act about our own daily digital habits that are contributing to this increased data volumes. Keeping numbers of pictures and videos taken reasonable, using acceptable, eco-friendly resolutions when filming or watching videos online, cleaning our e-mail boxes, etc. will also concretely help mitigate the data storage related $CO_2$ footprint.

## 29.14 The Future of Data Storage, the Permanence of Data and Related Questions

Current storage technologies can store data for up to 50 years at a maximum (and actually probably much less on average) and with energy consumptions that despite their optimisation and efficiency are still chal-lenging given the continuous data explosion and  while the planet is leaning towards 2050 zero-net GHG emissions targets.

In this context, one new – surprising - storage medium could repre-sent a relevant long-term alternative, addressing both of these limita-tions: synthetic DNA. Encoding synthetic DNA with data and encapsu-lating it in silica could actually store up to 700 terabytes of information per gram of medium[411] and last for thousands to millions of years with little to no power consumption.

While the first experiments started in the 80s, based on concepts im-agined in the 60s and  are still being developed in laboratories necessitat-ing expensive and advanced equipment, a DNA Data Storage Alliance has been formed in 2020 with the objective to 'jumpstart the standards

---

[410] Petrov, Christo. 25+ Impressive Big Data Statistics for 2023, https://techjury.net/blog/big-data-statistics/

[411] Isaacson, Betsy, Storing Digital Data for Eternity, Newsweek, 06/22/15, https://www.newsweek.com/2015/07/03/storing-digital-data-eternity-345557.html

development for DNA data storage' and the 'mission to create an interoperable ecosystem for DNA-based data storage solutions.'[412]

Low speeds of writing and reading, high cost levels and managing error-free data replication are still the main challenges to overcome though.

While getting such a technology mature and operational for everyday use would be an incredible progress in terms of keeping archived data safe from many risks of losses and dramatically lowering environmental impacts, the question of storing data permanently or for the eternity also raises very specific questions.

Protecting humanity's knowledge and memory in general, making contact with extraterrestrial life through time capsules, or preserving 'the most essential information to sustain or rebuild civilization in the event of an apocalypse (digital or other)'[413], these are some of very high-level projects and purposes for which such an almost endless retention of data can make sense without a doubt.

However, beyond those specific purposes, is it acceptable and does it make sense to retain data indefinitely or for too long? Even without going into the centuries duration that we have just mentioned, just considering the general life of an individual.

At a time when/where our digital footprint across all our devices makes it possible in a way or another to establish very detailed individual profiles of each of us, literally called digital twins, sometimes even updated in real-time, there is the question of not only collecting all this data in the first place, but also to keeping it and aggregating it indefinitely.

---

[412] SNIA Announces DNA Data Storage Technology Affiliate, SNIA, https://www.snia.org/news_events/newsroom/snia-announces-dna-data-storage-technology-affiliate

[413] Campbell-Dollaghan, K. 2016. We Have The Technology To Store Data For Eternity. Now What? https://www.fastcompany.com/3056762/we-have-the-technology-to-store-data-for-eternity-now-what

We all evolve through life and its cycles: childhood, adolescence, adulthood, senior stage. For digital natives in particular, born with digital devices all around and having been using them and exposed to them since sometimes very early ages, is the perspective acceptable to have their data, being aware of it or not, being aggregated, retained, creating a full historical archive and profile of their whole life? A 2021 French study[414] figured out that even 39% of babies had a digital footprint (although usually without a name and an account) *before even being born*, through their parents posting news and ultrasound scanning on social media during pregnancy...

Should we be defined, profiled, assessed, recommended, guided, advised, based on aggregated, persistent data, dating back sometimes to many years? Is there a risk in terms of data-based recommendations centered around either general demographics, and/or on past activities and centers of interests to potentially keep us circumscribed, trapped indefinitely in a given state or stage, self-reinforcing artificially some of our already existing patterns or biases? Is this generalized approach not narrowing down who we are, who we could become, limiting the perspective of getting exposed to ideas, concepts, and anything existing that is not yet in our direct environment of thought and life? Do we ultimately shape data? Or is data shaping us? That's the question.

And when you consider the question not only at the level of single individuals but at a generational level up to the overall society, is there not a risk to increase issues and problems like social reproduction, determinism, reinforced self-repeating patterns, communitarianism and related?

This brings the question down to Dataism, this mindset or philosophy, maybe even becoming kind of a religion for some, where data and

---

[414] See study: Sondage exclusif: la digitalisation de la vie familiale. Faireparterie, https://www.faireparterie.fr/etude-enfants-rapport-digital/

information flows become the 'supreme value'[415] by trying/attempting to characterize, understand and model absolutely anything in any field and dimension of life.

While aggregating and leveraging data to feed algorithms will undoubtedly allow major positive innovation and progress, providing more clarity, better understanding thus enabling more objectivity in many decision-making processes, we need at the same time to remain human-centered in our purposes.

As the quote says: "Not everything that can be counted counts, and not everything that counts can be counted." And what we must absolutely not forget to value in this increasingly technological and data-driven world are our qualities of humanness and humaneness, that make humanity and all of us intrinsically unique and our quality of being human unquestionably precious.

## 29.15 Conclusion

As we have seen data storage is a multi-faceted subject involving many various dimensions intertwined with ethical considerations: technology and technicalities, corporate organization, regulation and law, privacy, human rights, geopolitics, economics, health & data, IT governance, corporate social responsibility, sustainability, sociology…

The current possibilities created by the flexibility, simplicity, speed and sophistication of cloud storage and processing are tremendous. It is up to all of us to make the best usage of these technologies and of the precious fuel that powers them, data, for the highest purposes as well as the good of humanity while carefully protecting data itself and values like privacy, fairness, democracy, openness, transparency, equality of rights, and respect of all among others.

---

[415] Harari, Yuval Noah (2017). Homo Deus: A Brief History of Tomorrow. UK: Vintage Penguin Random House. p. 428. See also article in Christoph Stückelberger/ Pavan Duggal, Eds., Cyber Ethics 4.0, op. cit.

This responsibility has to be shared across all levels of stakeholders along the data storage and management chain.

It starts with the:

- CSP providers that are at the core of this industry to offer the highest standards for protecting data, privacy, equity of access, democratic behaviours and to minimize impacts on the environment, all in a full transparent manner
- Technologists, to design evolving technologies and systems that take into consideration the needs for privacy, for managing and storing data within ethical limits and in a sustainable way
- All organizations that generate, store and process data to do so in a responsible and focused way, minimizing their data needs and with a constant control of their operations to pro-actively prevent risks, and showing positive corporate moral responsibility in the field of data
- Government, lawmakers and authorities to regulate in order to avoid data privacy and protection loopholes on their territory, while not themselves legislating on data and data control for their own political interests
- International organizations and bodies to establish international regulation and principles  to fight overall data and data storage inappropriate uses and practices and to articulate and resolve conflicting national laws in the most ethical manner
- IT and Data Professionals becoming fully aware of the complete spectrum of their responsibilities in terms of ethical data storage and data management practices, raising their level of professional excellence and professional ethics on this matter
- Universities and higher education institutions to equip systematically upfront their students in IT and Data specializations with dedicated ethical trainings
- Citizens to be aware and in control of their own data, to be vigilant about corporation practices and to act in a sustainable way

- All of the above to educate and educate plus update themselves on the matter of ethics and data

The overall realisation by everyone of the extent of their own responsibility, crossed with the pro-active consideration of the very diverse aspects and dimensions of ethical data storage and practices is the way leading to a safer, more human and more trusted technological development. Let's take this path all together!

**PART G**


**SERVING HUMANS:**
**DATA GOVERNANCE AND GEOPOLITICS**

# 30

## SOVEREIGN OR DEPENDENT? MORE TECHNOLOGICAL INDEPENDENCE IN DIGITAL COMMUNICATION

*Christian Martin*[416]

*Many countries, including Switzerland, will always have to make compromises. The idea that it will ever be able to build its communication technology and networks completely independently, which meet our modern needs, belongs in the realm of theoretical games and is far from any practice. Or to formulate it even more clearly: all discussions about independent communication infrastructures are totally alien to the market and highly utopian. The regulation of communication technology and its innovation should be kept as small as possible as long as the market*

---

*works. The prevailing approach in Switzerland with ex-post regulation in the telecommunications sector makes sense. It should only be regulated when the market fails and the participants disagree.*

Switzerland is lagging behind digitally; it is rumoured again and again. It is forgotten that their communication technology has an extremely high standard in comparison to other countries. Today, Switzerland benefits from a very powerful broadband infrastructure and from competition from its large and small infrastructure and service providers. The mobile infrastructure is also above average in terms of adaptation and use: the service providers are best-in-class, and theirs is not Swiss consumers benefit from price-fixed competition for coverage, quality and innovation.

However, this positive assessment has one weakness: as far as fibre optic networks are concerned, Switzerland is below the OECD average. The reasons for this are the above-average expansion costs in this country, the non-existent national high-broadband strategy and obstacles such as the lack of a nationwide cable channel register. This, coupled with an ongoing investigation into the allowable network architecture, reduces the incentives for the big players to spend more money and move faster with the construction of the fibre infrastructure.

The expansion of the Swiss fibre optic network as the basic infrastructure is extremely important for the communication infrastructure and data use in our country. Not only for the approximately 570,000 e-athletes, for whom even 5G or 6G is not enough. Driven by digitization, there is a continuous surge in innovation in the ICT world. In addition to a top mobile network, this always requires a top broadband network. The market of entrepreneurs, service providers and consumers is shifting to a large extent to central innovative solutions such as SaaS (Software as a Service), which are produced, moved, managed on the internet and hosted and stored in public clouds or in multi-cloud systems. This consolidation is a logical step that is not the first in history: the spinning mill on

the river used to produce the electricity itself. Today it comes from the net.

## 30.1 Dependence is not the Same as Dependency

So the question cannot be whether Switzerland has a sufficiently good communications infrastructure, but how it uses and develops it. Here our country runs the risk of getting caught not only in legal but also in political leadership. At present, the benefit of the linking and networking potential is weighted much less in the political and public debate than the dependencies that arise in a network by virtue of its existence. This is especially true with regard to cyber security.

But which dependencies are we talking about exactly? The ICT world has never been independent and the extent of its dependencies is complex.

Dependence begins with the hardware and software components. Should Switzerland manufacture the electronic components, the associated drivers and patches itself in the future, and if so, at what price? Is it actually a problem when telecommunications providers such as Ericsson, Juniper or Huawei bind purchased components and supply operating systems and support? Is it the lack of freedom from various service providers, from SaaS, cloud or open source providers and their operating and support models, which is causing the Swiss malaise? Or are there development policy concerns about the extraction of the raw materials right from the start? Is it simply annoying that investments are influenced by foreign states and companies?

## 30.2 Not only Switzerland, also the USA and China are Making Compromises

No matter what perspective you take and what importance you attach to the different dependencies: Switzerland will always have to make

compromises. The idea that it will ever be able to build its communication technology and networks completely independently, which meet our modern needs, belongs in the realm of theoretical games and is far from any practice. Or to formulate it even more clearly: All discussions about independent communication infrastructures are totally alien to the market and highly utopian.

As far as components are concerned, for example, there are only a few countries that have a substantial industry and therefore dominate the market, mainly the USA and China. And yet every American iPhone contains electronic components from 48 different countries and six continents. It is quite possible that the Americans also regret this fact from a commercial point of view – but from a technical point of view it is not decisive for the match. When it comes to the question of independent technology development, we should therefore place it in a larger context. We should answer what we want to achieve with this autonomy and what we want to protect ourselves from. Against industrial espionage? From state-driven cyber attacks? Before kill switch scenarios?

Personally, I don't believe in these scenarios, just as I consider backdoors to be very unrealistic for the big international market leaders. This shows me my experience in the companies where I have worked or worked for years - at Cisco and at Google. At the same time, I am convinced that this also applies to every other global ICT company: Today, technological development is controlled multinationally and distributed among many heads. The development teams consist of thousands of employees spread across the globe - at Google the head of infrastructure development is Swiss, at Cisco it was an Indian. With the best will in the world, I can neither technically nor practically imagine how such manipulations could be carried out, let alone how the effort involved could be kept secret.

## 30.3 Better and Worse Approaches to More Autonomy

But even if Switzerland is not striving for self-sufficiency, but simply wants more independence, we should ask ourselves in which disciplines can it achieve more autonomy at all. In my opinion, there is scope in self-determined data management (data autonomy), in operation and hosting (operational autonomy) or in software development (software autonomy). However, with regard to this desire for autonomy, we should also ask ourselves what we want to satisfy with it - and at what price.

Are we clinging to the myth of a completely sovereign Switzerland that never existed? Do we want to counteract phobias of digital technology addiction that are rampant on social media and in user forums? Do we have to satisfy partisan or geopolitical ideologies? Is it important to take fears seriously, even if they are based on a great deal of half-knowledge and poor know-how regarding ICT and security in particular?

One option is that we put economic interests or the promotion of the Swiss economy above everything else. However, there is a risk that we ignore the old adage that IT is a scalable business that thrives on mass and size, and that Swiss solutions ultimately have to assert themselves in a global market in order to remain affordable in terms of price.

Another noble goal would be to promote Swiss research and counter-act the shortage of skilled workers. However, it should be considered whether international ICT companies are so successful not only because of their market position or market power, or whether the continuous innovation that they offer is not more decisive. As is well known, Swiss customers also go primarily to the cloud because they receive highly innovative, secure, sustainable services there that meet local regulatory requirements, either as on-premises or from a local managed service provider. It is also often forgotten that very large telecommunications providers, so-called telcos, and industrial conglomerates have tried in vain at independent "local" developments. The reasons for failure were

always the same: high prices, little functionality or talent recruitment problems. In Switzerland, despite the good universities, the fact that digitization, the role of the MINT subjects (mathematics, computer science, natural sciences and technology) and the entire start-up culture and the necessary capitalization have been underestimated for a long time makes things even more difficult.

With more self-determination, however, we can also aim for more transparency. Do we want to know on which infrastructure services are built? How and by whom is the infrastructure operated and what operational standards are met? Where is the data located and can we see it? Do we already know how the security processes are designed? Does transparency already exist, or do we just have to demand it more or simply read the data protection requirements more carefully?

Many consumers and companies are already opting for services whose underlying processes, operating models and security concepts they understand well. In terms of transparency, there is probably still room for improvement in many companies. Today by no means everyone receives the appropriate transparency on the above questions from their bank, insurance company, telco provider, doctor, health insurance company or municipal administration.

Finally, there is the regulatory leeway that Switzerland can use to secure more independence. The role of Parliament in policing technology or stopping it from doing evil things is outdated and should be changed. We really need a different understanding of cooperation here. Regulation could join forces with research and play a more active role, especially with regard to questions of security and operational models of telecommunications providers. But beware: the regulation of communication technology and its innovation should be kept as small as possible as long as the market functions. The prevailing approach in Switzerland with ex-post regulation in the telecommunications sector makes sense. It should only be regulated when the market fails and the participants disagree.

There is also definitely potential in the process. For example, ComCom is currently dealing with cases that began almost a decade ago. In general, the following applies to Switzerland: The lengthy procedures for determining a market failure and the many objection and appeal options at every instance level should be reconsidered. They are a real obstacle on the way to Switzerland's digital future. If these processes are accelerated, more control and consequently more self-determination and independence are possible in the fast-moving ICT and telecommunications industry.

## 30.4 There are More Solutions for Self-determination in Technology Development

At the level of data autonomy, data ownership, key management and data encryption can be adapted to Swiss needs. At the level of operational autonomy, identity management (eID) and end-to-end encryption should be promoted so that everyone can regulate who has access to their data. Compliance and auditing companies must be technically able to verify zero trust concepts and better control the operation and risk management of the networks. In order to prevent locked-in scenarios and to enable risk management with different providers, technologies and operating models, the federal government and the supervisory authorities must ensure that Swiss consumers and companies continue to have the choice between different, independent communication networks in the fixed network and mobile communications.

## 30.5 Conclusion

In terms of resilience, all of this would give the Swiss enough opportunities to build highly available architectures and use them independently.

One last point remains: the federal government should demand more transparency from its own organisations and from the companies that are driving digitization forward. All companies, whether Swiss or global, should be required to show in a transparency report how they handle customer data and what measures have been implemented to improve security. Transparency creates trust, which could be the basis for even safer, even better communication technology in Switzerland.

# 31

# STEALING DATA:
# HOW TO REACT TO CYBER CRIMINAL
# CLAIMS OF EXTORTION OF MONEYS?
# LEGAL AND ETHICAL ANSWERS

*Pavan Duggal, India*[417]

Today data is the new oil of the data economy. Everywhere around us, we see data. This data has become the ubiquitous currency of our times and all this has happened, thanks to one big event that has happened during our lifetimes – the internet.

---

The internet is the second most significant event in human history after the advent of fire. No other event had such a massive and profound influence on the way how humans think, believe, perceive and do commerce.

With increasing adoption of the internet and with more and more people coming on to the internet bandwagon, data has become the new currency. Everybody is interested in not just receiving data but also generating data and transmitting data. No wonder, we have to realize that internet as a paradigm has transformed all of us from human beings into digital data entities.

We have all become global authors, global publishers and global transmitters of data. We are constantly producing data on a 24/7 basis, apart from constantly consuming data. Hence, data becomes very important. With the increased adoption of technology and new tech devices, almost everyone is saving their data in the digital format.

A majority of data is stored on devices while increasing quantum of data are now being stored on the cloud. From the perspective of corporates, this data becomes the crucial business resource as this business resource is huge not just for fulfilling and performing the day-to-day functions in the companies but also for the purposes of analysing past performances and predicting future trends, apart from being the raw material for reviews and analysis.

No wonder, data has become more crucial to legal entities than anything else. Today money is important, but I think with the passage of time, data will be far more important than money. This is because this data becomes capable of being monetized and once it is monetized, it can be used at any point of time by any stakeholder.

No wonder, data has become the most fertile target for potential attacks. Everybody wants to have access to not just to their own data but also the data of others. More and more cyber criminals want to have access to others' data so that the said data can then be used, monetized, disseminated as well as misused and abused.

This growing reliance on data has to be seen in the historical context of the times that we live in. At the time of writing this article, we are coming out of Covid-19 period. Covid-19 has not been just the public health emergency or the pandemic but it is also a cyber pandemic. The coming of Covid-19 has triggered irreversible changes, which are going to completely change the way how people interact in the online ecosystem.

In my book *New Cyber World Order Post Covid-19*[418], I have argued that by the time nations of the world are completely victorious against Covid-19, its current and subsequent waves of infections, the world will enter into a new cyber age, where the New Cyber World Order is awaiting us.

In this New Cyber World Order, states are going to become very powerful. But more significantly, cybercrime will become the new default normal. Also increasing cyber security breaches will be our constant companion.

Once we analyse both these trends of increasing cybercrimes and cyber security breaches, we quickly realize that the ultimate objective of both these paradigm shifts is primarily to target data. It is data which is targeted unauthorised, illegally or without permission and consent of the concerned data holder, so that this very data can then be used, monetized or abused in a variety of manners.

Today, we need to realize that we have already stepped into the data economy age. This is the economy which is dependent on the data, where data becomes the new currency of the new data economy. In this data economy, thus data assumes far more relevance. This becomes even more apparent when one looks at the various facts and figures pertaining to the growing quantums of data.

There were 79 zettabytes of data generated worldwide in 2021. By 2025, more than 150 zettabytes of big data will need analysis. The

---

[418] Duggal, P. 2020. New Cyber World Order Post Covid-19. Independently published, 53p.

COVID-19 pandemic increased the rate of data breaches by more than 400%. By 2027, the use of big data application database solutions and analytics is predicted to grow to $12 billion. [419] 463 ZB of data will be created every day by 2025. (Raconteur, 2020)[420]

A perusal of the aforesaid figures thus clearly tells us that the data bucket is constantly growing which is also a growing big triggering factor for growing cyber security breaches in the data economy age.

Today cyber security is being breached for a variety of purposes. But ultimately the focus of cyber security breach is to target, unauthorised access data and have unauthorized copies of data.

Cyber security breaches are happening all over the internet. State and non-state actors are both targeting and are being targeted by growing cyber security breaches.

During the third quarter of 2022, approximately 15 million data records were exposed worldwide through data breaches. This figure had increased by 37 percent compared to the previous quarter.[421] Between March 2021 and March 2022, the average cost of a data breach in the healthcare sector amounted to over 10 million U.S. dollars, up from 9.23 U.S. dollars between May 2020 and March 2021.[422]

It has been estimated that Ransomware accounts for nearly 24 percent of incidents in which malware is used (Verizon)[423]. By stealing 10

---

[419] Djuraskovic, Ogi. 2022. Big Data Statistics 2023: How Much Data is in The World? https://firstsiteguide.com/big-data-stats/

[420] 53 Important Statistics About How Much Data Is Created Every Day, FinancesOnline, https://financesonline.com/how-much-data-is-created-every-day/

[421] Petrosyan, Ani. Number of data records exposed worldwide from 1st quarter 2020 to 3rd quarter 2022, Statista, https://www.statista.com/statistics/1307426/number-of-data-breaches-worldwide/

[422] Petrosyan, Ani. Average cost of a data breach worldwide from May 2020 to March 2022, by industry, Statista, https://www.statista.com/statistics/387861/cost-data-breach-by-industry/

[423] 2019 Data Breach Investigations Report, Verizon, https://enterprise.verizon.com/resources/executivebriefs/2019-dbir-executive-brief.pdf

credit cards per website, cybercriminals earn up to $2.2 million through form jacking attacks (Symantec)[424]. The average total cost of a data breach was more than $1 million higher when working remote was a factor in causing the breach, compared to breaches in which working remote was not a factor (IBM).[425]

In these breaches, data becomes the final fertile target of attack. Therefore, stealing of data becomes the rampant trend of today's times. Once data is stolen, whether it is personal data or non-personal data, the same can then be monetized not just on the superficial net but also on the darknet and therefore new innovative and ingenious approaches are being adopted for the purposes of stealing data.

One of the most significant advances in this regard is the advent, constant growth and consolidation of ransomware as a paradigm. Ransomware is the big headache of today's times. But what exactly is ransomware. Let's see how ransomware is defined by major thought leaders as detailed below:-

*Ransomware* is a type of malware that prevents or limits users from accessing their system, either by locking the system's screen or by locking the users' files until a ransom is paid. [426]

Ransomware is a type of malicious software (malware)[427] that threatens to publish or blocks access to data or a computer system, usually by encrypting it, until the victim pays a ransom fee to the attacker.[428]

---

[424] Internet Security Threat Report, Executive Summary, 2019, ISTR 24, https://docs.broadcom.com/docs/istr-24-executive-summary-en

[425] IBM. 2022. Reports. Cost of a data breach 2022, https://www.ibm.com/reports/data-breach.

[426] Ransomware, Trend Micro Incorporated, https://www.trendmicro.com/vinfo/us/security/definition/ransomware#:~:text=Ransomware%20is%20a%20type%20of,until%20a%20ransom%20is%20paid.

[427]What Is Malware? Proofpoint, https://www.proofpoint.com/us/threat-reference/malware

[428] What Is Ransomware? Proofpoint, https://www.proofpoint.com/us/threat-reference/ransomware

Ransomware is a type of cyber extortion where a malicious actor infiltrates an environment and encrypts and exfiltrates files, denying access and threating disclosure, unless the victim pays a ransom.[429]

Ransomware is a specific type of malware[430] that extorts victims for financial gain. Once activated, ransomware prevents victims from interacting with their files, applications or systems until a ransom is paid, usually in the form of an untraceable cryptocurrency like Bitcoin. In some cases, the victim is instructed to pay the perpetrator by a set time or risk losing access forever. In other cases, the perpetrator intermittently raises the ransom demands until the victim pays.[431]

Ransomware began with small beginning.

The first documented ransomware, *AIDS Trojan or PC Cyborg*, was delivered at the World Health Organization's AIDS conference in 1989 using floppy disks, demanding a payment to be sent to a postal office box in Panama. This malicious code was not encrypting the files content as we know it today, but the filenames only. It was however enough to take down the systems and cause disruption. [432]

Today ransomware has various salient features. Ransomware has become so severe today that it has becoming to become the default cyber security challenge and headache of today's times. This becomes further evident when one looks at various facts and figures of ransomware.

Ransomware cost the world $20 billion in 2021. That number is expected to rise to $265 billion by 2031. Recovering from a ransomware attack cost businesses $1.85 million on average in 2021.Out of all ran-

---

[429] Ransomware, Gartner Glossary, https://www.gartner.com/en/information-technology/glossary/ransomware

[430] What Is Malware?, Proofpoint, op. cit.

[431] Ransomware, Cyberark Glossary, https://www.cyberark.com/what-is/ransomware/

[432] Lessing, Marlese. 2020. Case Study: AIDS Trojan Ransomware, https://www.sdxcentral.com/security/definitions/what-is-ransomware/case-study-aids-trojan-ransomware/

somware victims, 32 percent pay the ransom, but they only get 65 percent of their data back.[433]

There was an 85% increase in ransomware attacks since 2020. (Palo Alto Networks, 2021). Reports expect there to be a ransomware attack every two seconds in 2022. (Cybersecurity Ventures, 2022) [434]

The aforesaid figures tell us that every stakeholder today has to be prepared for ransomware. It is not only a question of "if", but it is only a question of "when" you would become a victim of ransomware attack.

Once a ransomware attack takes place, the encryption algorithms come in and encrypt entire data and therefore money is sought for the purposes of providing the key for decrypting the said encrypted data.

That claim for money is often the cherry on the cake. This claim of money is made in the form of crypto-assets and Bitcoins so that they can become less traceable and the chances of law enforcement agencies reaching the ultimate cyber criminal gets further diminished.

Whatever kinds of ransomware attacks can take place, whether it is simplicitor ransomware attacks or wiper malware attacks, ultimately seeking moneys becomes the primary objective.  Huge quantum of moneys have been asked on ransomware attacks across the world.

On November 8, 2021, a ransomware attack took place against *MediaMarkt*, Europe's largest electronics retailer. The Hive gang[435], which carried out the attack, initially demanded a ransom of 250 million USD, but the amount was reduced to 50 million USD after a while.[436]

---

[433] Kochovski, Aleksandar. Ransomware Statistics, Trends and Facts for 2023 and Beyond, https://www.cloudwards.net/ransomware-statistics/

[434] Stouffer, Clare. 2022. Ransomware statistics: 102 facts and trends you need to know in 2023, Norton, https://us.norton.com/blog/emerging-threats/ransomware-statistics#

[435] Dark Web Profile: Hive Ransomware Group. 2023. SOCRadar Research, https://socradar.io/dark-web-profile-hive-ransomware-group/

[436] 20 Interesting Facts About Ransomware, SOCRadar, https://socradar.io/20-interesting-facts-about-ransomware/

The highest-profile attack of 2021 was arguably the one on Colonial Pipeline. This provided over 40% of the East Coast's fuel, but as the ransomware attack left the company unable to bill properly, it was forced to suspend operations temporarily. Even though the company paid a $4.4 million ransom, it took a long time to get everything back online, leading to fuel purchasing restrictions in 17 US states. [437]

Ransomware group DarkSide targeted the chemical distribution company Brenntag and demanded a payout of $7.5 million in Bitcoin. (Brenntag, 2021)

The hacker group behind an oil company attack allegedly acquired $90 million in ransom payments in only nine months from around 47 victims. (Fox Business, 2021)[438]

Garmin, a major player in the tech space, suffered a severe breach that brought its GPS services offline for several days. In order to regain control, the company allegedly paid a $10 million ransom.

The US Treasury has linked more than $5 billion of Bitcoin transactions to ransomware.

The highest ransom demanded from a victim reached $70 million in 2021 (Blackblaze, 2021)[439]

Easily it is apparent this entire paradigm of ransomware throws up large number of crucial ethical and legal questions. Is it ethical to steal data? The answer for that is crystal clear no. Is it ethical to engage in ransomware? The answer is the same, also no. Is it ethical to ask for

---

[437] Cook, Sam. 2022. 2018-2022 Ransomware statistics and facts, Comparitech, https://www.comparitech.com/antivirus/ransomware-statistics/

[438] Brooke Crothers. 2021. JBS ransomware attack points to ominous trend targeting critical industries by foreign actors, FOXBusiness, https://www.foxbusiness.com/technology/jbs-ransomware-attack-trend-targeting-critical-industries-foreign-actors

[439] Clancy, Molly. 2021. The True Cost of Ransomware, Backblaze, https://www.backblaze.com/blog/the-true-cost-of-ransomware/;   Ransomware statistics: 102 facts and trends you need to know in 2023, Norton, https://us.norton.com/blog/emerging-threats/ransomware-statistics#

cyber criminal claims for extortion of moneys? The answer here is also crystal clear, an emphatic no.

The answers to all the three questions from a legal standpoint are also a crystal clear no. This is so because data ethics as an emerging science has been evolving which has stipulated various principles pertaining to ethical use of data. Hence, the ethical expectation of every data stakeholder is that their data is going to be ethically used and are not unauthorised or unethically accessed, downloaded, copied, extracted, modified, deleted or prejudicially impacted in any manner whatsoever.

Since it is unethical to engage in ransomware, it is a logically corollary that it is unethical to make payments for cybercriminal claims of extortion of moneys. This is so because the said conduct is neither ethical, nor prudent.

Further, once you start paying money to cyber criminals in a ransomware attack, you effectively give a confirmation to the fact that you are not just the fertile target, as you have the capacity to pay, but also that you can tomorrow become a repeat target of such kinds of cybercriminal claims of extortion of moneys.

That is the reason one finds that law enforcement agencies across the world are encouraging people not to make extortion payments to cyber criminals.

Evidence in this regard is well established. The Federal Bureau of Investigations in the United States has gone ahead and evolved its thought process.

Below is an advisory taken directly from the U.S. FBI

*"The FBI does not advocate paying a ransom, in part because it does not guarantee an organization will regain access to its data. In some cases, victims who paid a ransom were never provided with decryption keys. In addition, due to flaws in the encryption algorithms of certain malware variants, victims may not be able to recover some or all of their data even with a valid decryption key.*

*Paying ransoms emboldens criminals to target other organizations and provides an alluring and lucrative enterprise to other criminals."* [440]

Apart from the aforesaid, various other statutory authorities in different parts of the world have also been encouraging people not to make extortion payments.

Such an approach is logical and ethical apart from legal as well and such an approach is also backed by emerging principles of data ethics as a science. The said ethical principles are also now being enshrined in legal instruments. For example, in the US, the new legal framework has been passed which has actually made the making of payments to ransomware attacks as a criminal offence.

US law frameworks now stipulate that any person, who makes payment to ransomware attack has to mandatorily disclose the same to the law enforcement agencies, thus exposing the said person to legal consequences.

Victims who pay ransoms might also be subject to criminal or civil penalties in some cases—for example, where a ransom payment is made knowingly to an entity either designated as a foreign terrorist organization or subject to sanctions by the Department of the Treasury. Nevertheless, policy considerations, mitigating factors, and prosecutorial discretion may weigh against enforcement in such instances. [441]

One of the biggest problems across the world is that most of the countries don't yet have dedicated legal frameworks on ransomware to deal with the cyber criminals claims of extortion of moneys. Most countries have in place national penal laws which make extortion an offence.

However, when one talks about cyber criminal extortion of moneys, it is clearly a manifestation of extortion of moneys and hence can be

---

[440] Organizations across the globe need to develop a ransomware payment policy, anticipating a potential future attack, EY Global, https://www.ey.com/engl/consulting/ransomware-to-pay-or-not-to-pay

[441] Ransomware and Federal Law: Cybercrime and Cybersecurity, Congressional Research Service, 2021, https://crsreports.congress.gov/product/pdf/R/R46932

covered under the traditional penal laws. Nevertheless, many countries are now increasingly looking to devise new laws pertaining to ransomware.

As of now, most of the countries don't have dedicated laws on ransomware. Similarly, most of the countries don't have laws on payment of moneys demanded in ransomware or payment of cyber-criminal claims of extortion of moneys.

This is a very fertile time in today's times as the ethical and legal principles will need to go hand-in-hand while crystalizing a legal response to cyber criminality. Ethical and legal principles will have to find appropriate mention in the respective national legal frameworks, which are now developing to counter the cyber criminal claims of extortion of moneys and the growing ransomware demands for moneys.

There is a need for ethical principles concerning data ethics to be incorporated in the forthcoming new legal frameworks on ransomware, that are now being discussed and deliberated in different parts of the world. Further, ethical principles need to be inculcated as an integral part of the new emerging legislative approaches, that are developing in countries to deal with the challenges of emerging technologies and their appropriate regulation.

Similarly, when Artificial Intelligence (AI) is used for the purposes of facilitating cyber criminal claims of extortion of moneys, those elements will also have to be appropriately addressed by ethical and legal principles put together.

It will be ethical for lawmakers to not just mandate such kind of activity as unethical and illegal but more significantly also to stipulate appropriate responsibilities and liabilities for the relevant Artificial Intelligence ecosystem stakeholders in the event AI is misused for cyber criminality.

Thus moving forward, some essential weapons will have to be part of everybody's arsenal in fighting the menace of cyber criminal claims of extortion of moneys in the context of stealing of data.

Cyber resilience will have to play a very important role in this regard. Cyber resilience is a new concept and is slowly gaining more currency. It has been defined by varius stakeholders in different ways.

Cyber resilience is the ability of an organization to enable business acceleration (enterprise resiliency) by preparing for, responding to, and recovering from cyber threats.[442]

Cyber resilience is the ability of an organisation to protect itself from, detect, respond to and recover from cyber attacks.[443]

Cyber resilience refers to an organization's ability to continue business operations despite a cybersecurity or data loss incident.[444] Cyber resilience is a measure of how well an enterprise can manage a cyberattack or data breach while continuing to operate its business effectively.[445]

Cyber resilience presumes that you will be under attack. The question is how can you avoid panic and quickly come back to a state of normalcy after getting attacked. This forms the foundational fulcrum of cyber resilience.

Further, cyber hygiene will have to play a very important role in this regard. Cyber hygiene has been defined by various stakeholders.

Cyber hygiene refers to the steps that users of computers and other devices can take to improve their online security and maintain system health. Cyber hygiene means adopting a security-centric mindset and

---

[442] What is Cyber resilience? Microfocus, https://www.microfocus.com/en-us/what-is/cyber-

[443] Cyber Resilience, IT Governance, https://www.itgovernance.co.uk/cyber-resilience

[444] Cyber resilience definition, Druva, https://www.druva.com/glossary/what-is-cyber-resilience/

[445] What is Cyber Resilience? Digital Guardian, https://digitalguardian.com/blog/what-cyber-resilience

habits that help individuals and organizations mitigate potential online breaches[446]

Cyber hygiene is a reference to the practices and steps that users of computers and other devices take to maintain system health and improve online security. These practices are often part of a routine to ensure the safety of identity and other details that could be stolen or corrupted. Much like physical hygiene, cyber hygiene is regularly conducted to ward off natural deterioration and common threats.[447]

Cyber hygiene is a set of habitual practices for ensuring the safe handling of critical data and for securing networks. It's like personal hygiene, where you develop a routine of small, distinct activities to prevent or mitigate health problems. Cyber hygiene practices include the inventory of all endpoints connected to a network, vulnerabilities management, and the patching of software and applications.[448]

The European Union's Agency for Network and Information Security (ENISA) states that "cyber hygiene should be viewed in the same manner as personal hygiene and, once properly integrated into an organization will be simple daily routines, good behaviors, and occasional checkups to make sure the organization's online health is in optimum condition".[449]

---

[446] Top Tips for Cyber Hygiene to Keep Yourself Safe Online, Kaspersky, https://www.kaspersky.com/resource-center/preemptive-safety/cyber-hygiene-habits

[447] Brook, Chris. 2022. What is Cyber Hygiene? A Definition of Cyber Hygiene, Benefits, Best Practices, and More, Digital Guardian Blog, https://digitalguardian.com/blog/what-cyber-hygiene-definition-cyber-hygiene-benefits-best-practices-and-more

[448] Null, Christopher. 2021. What Is Cyber Hygiene and Why Does It Matter?, Tanium Inc., Blog, https://www.tanium.com/blog/what-is-cyber-hygiene-and-why-does-it-matter/

[449] Tyas Tunggal, Abi. 2022. What is Cyber Hygiene and Why is it Important? Blog, UpGuard Inc., https://www.upguard.com/blog/cyber-hygiene

Therefore, cyber resilience and cyber hygiene will have to play important response mechanisms so as to prevent or minimize the impact of cyber criminal claims of extortion of moneys in the context of stealing of data.

Backups will have to become the number one default normal. Backups, backups, backups will have to be the only mantra going forward, as stakeholders would be required to back up their data every 2-3 days in order to ensure that their operations don't go offline in the event they are hit by ransomware attack and either unable or do not want to make payment for cyber criminal claims of extortion of moneys.

Backups represent an ethical, legal and pragmatic approach to the constant challenge of cyber criminal claims of extortion of moneys in the context of stealing data. A Backup basically means the process of taking copies of your data that has already been generated.

Backup copy means the copy of at least those source data (software assets and information assets) which are needed for the recovery and/or reestablishment of business processes.[450]

A backup is another resource used to make sure an activity can go on, or deploy an alternate project if a primary one cannot be done. Saving data on disc or offsite prevents important information from being lost, stolen, or sabotaged.[451]

Backing up data brings forward its own ethical and legal considerations. In today's times, when we all are being flooded by constant new volumes of data with each passing day, it would be so tedious and time consuming to try and reconstruct data. A backup is an ethical approach – since it is your data, it is your duty to have a copy. Once your data is duly backed up, at least it tends to give you a sense of peace of mind. Without this precaution, the entire data would be lost once and for all in the event of a ransomware attack or data wipe out.

---

[450] Lawinsider, Dictionary, https://www.lawinsider.com/dictionary/backup-copy
[451] Ibid, https://thelawdictionary.org/backup/

The time has come for us to revise our thinking and find new approaches to deal with the challenges of data theft and the cyber criminal claims of extortion of moneys. We will have to become more proactive, more fast and flexible in our thought process.  We will have to take it as a given that yes, we will be hacked and our data is going to be constantly attacked and breached. Hence, any kind of data breach of data stakeholders should not take us by surprise. Hence, we will need to constantly rely upon not just new volumes of data but also preserving and retaining properly the old data that is generated. Because in a data driven world, it is only data that is going to be our biggest source of empowerment and source of protection. This becomes even more important when one looks at the facts and figures pertaining to growing data theft and projected trends in this regard.

- The US suffers from the most data breaches
- The biggest data breaches in 2021 were:
- Comcast (1.5 billion)
- Brazilian resident data leak (660 million)
- Facebook (533 million)
- LinkedIn (500 million)
- Byeka (400 million)

According to the Symantec Security Summary, April 2021, ransomware payments jumped 171% in 2020, with the highest payout doubling to $10 million.[452]

From the perusal of the aforesaid facts and figures, it thus becomes very clear that it is ethical not just to produce and deal with your own data but also that you take appropriate steps to protect data. Since your data could potentially become a target of potential theft, it becomes

---

[452] O'Driscoll, Aimee, 2023. 30+ data breach statistics and facts, Comparistech Inc., Blog post, https://www.comparitech.com/blog/vpn-privacy/data-breach-statistics-facts/

absolutely imperative that all ethical and legal, pragmatic and practical steps must be taken so as to protect your data.

If appropriate methodologies of cyber hygiene, cyber resilience and backups are adopted, then the potential losses caused by data theft and the consequent cyber criminal claims of extortion of moneys could be potentially minimized to the best lowest extent possible.

However, the world is walking in a complete different new direction. Every 11 seconds, a company became victim of ransomware attack in 2022. It is expected that by the end of 2023, in every 9 seconds, a company is likely to become a victim of ransomware attack. Recovering from ransomware attacks today are on average very costly for  for legal entities.

Hence, it is a given fact that ransomware attacks will keep on increasing. Data theft will constantly grow and the cyber criminal claims of extortion of moneys would also multiply with each passing day, month and year. The absence of any international legal frameworks to control or deal with such kinds of data theft or ransomware attacks or cyber criminal claims of extortion of moneys further complicates the entire scenario.

Since data theft happens on the internet which is a transnational paradigm, it also assumes extra-territorial nuances. Further, as the internet has made geography history, it is possible for a legal entity or person to sit in one jurisdiction and steal data located in another jurisdiction and also ask for cyber criminal claims of extortion of moneys from stakeholders in a third jurisdiction.

Hence, this entire issue of stealing data brings forward various complicated issues of internet jurisdiction and connected nuances. The new trends are clearly telling us that the only way to protect yourself is to empower yourself and strengthen your cyber security. Cyber resilience and backups will have to be the only mantra going forward for the purposes of protecting yourself from becoming victim of untoward consequences because of loss of data.

To conclude, it can be said that in today's world, data ethical principles, exist but they have not been given the right kind of importance, thrust or impetus. Instead, cyber criminals are going ahead and acting with impunity to steal people's data and thereafter use it for extortion.

These entire paradigms of stealing of data for cyber criminal claims of extortion of moneys need to be handled, both from legal and ethical standpoint. The ethical principles concerning data theft, which are evolving in data ethics as a discipline, need to be specifically incorporated by legal frameworks which have been formed so as to meet up with the challenges of ransomware, data theft and ensuing extortion.

A lot of work needs to be done now. Currently, it appears that the legal and ethical frameworks are lagging far behind as compared to the speed with which cyber criminals are moving ahead in terms of data theft. The onus will be on the cyber ecosystem stakeholders, specifically on lawmakers and governments to come up with effective and ethical principles incorporated in their national legal frameworks to deal with the menace of stealing of data, ransomware and raising of subsequent cyber criminal claims of extortion of moneys.

This is a very interesting space to watch as time passes by.

# 32

# CYBERLAW, CYBERCRIME, CYBERSECURITY. INTERNATIONAL CONFERENCE RECOMMENDATIONS

*The following text is the Outcome Document of the International Conference on Cyberlaw, Cybercrime and Cybersecurity in Cyberspace, adopted by its participants. The hybrid conference was held online and in person in New Delhi/India from 23-25 November 2022. It was supported by several Departments of the Government of India such as the Department of Science and Technology, Companies such as Dell Technologies, the Asia Cloud Computing Association and think tanks and NGOs such as Globethics.net Foundation, UN Peace University, Cyberlaw University, Metavers Law Nucleus, and others. See: https://cyberlawcybercrime.com and: https://repository.globethics.net/handle/20.500.12424/426649*     *The Editors CS/PD*

## 32.1 Preamble

**UNDERSTANDING** that the newly emerging technologies cannot be left unregulated;

**REITERATING** that in the time to come, technological advancements are source for worry and tools to help in the propagation of crimes and offences in the cyberspace;

**AGREEING** that law and technology sectors can no longer be exclusive;

**ACCEPTING** that it is time for judges, lawyers, law-makers, politicians and digital stakeholders to take note of the increasing technological advances;

**RESTATING** that the legal jurisprudence of yesteryears is no more relevant, in face of technology that is advancing with the speed of light;

**WITNESSING** the safety features like mobile passwords, PINs and other security features slowing becoming outdated, with the law still not developed to regulate these changes;

**UNDERLINING** that there needs to be separate legal provisions for protecting the Critical Information Infrastructure of each nation, due to the important role played by it;

**ACKNOWLEDGING** that already, a major part of wars in the physical world are actually being fought in the cyberspace;

**NOTICING** that the best way to paralyse a country is by paralysing its Critical Information Infrastructure, with the method and mode of attack, being the most efficient tactic in the hands of external and internal malicious actors;

**SPOTTING** that the laxity of cybercrime regulatory regimes have contributed to emergence of new kinds of cybercrimes;

**NOTICING** that Covid-19 has accelerated the growth of cybercrime, ushering in the potentially decades long Golden Age of Cybercrime;

**UNDERSTANDING** that the problem of Internet jurisdiction in bringing cybercriminals to justice can be solved only by mutual understanding and resolve at the international level;

**APPRECIATING** the efforts of individual governments that have started emerging to counter cybercrime;

**FOCUSING** that cybersecurity is an often neglected area, which is not given the deserved attention;

**UNDERSTANDING** that although the cybersecurity landscape is still evolving, malicious actors have already started to pry into its vulnerabilities;

**ACCEPTING** that barring a few developed nations, most governments need to work on cybersecurity law, with specific legal provisions pertaining to cybersecurity being an absolute must for every jurisdiction;

**FOCUSING** on how the private sector as well as the public sector are susceptible to constant menace of cyber security breaches;

**EMPHASISING** that Ransomware has emerged as the foremost challenge in cyberspace for state and non-state actors;

**REPEATING** that specific and express legal provisions are required to deal with the use of Ransomware in the laws of every country;

**ACKNOWLEDGING** how Blockchain has emerged as a technology that is unprecedented in its effects;

**UNDERSTANDING** that since many new technical phenomena are based on Blockchain, there is a need for the evolution of specific laws on the same in each country;

**REALISING** the importance of Web 3.0 and how its emergence has the power to bring more users on the Internet bandwagon;

**NOTICING** how the increased usage of the Internet will further lead to the abundant availability of data that is ready to be exploited;

**STUDYING** how edge computing can have after effects which are not regulated by any law;

**FOCUSING** on the role played by social media in today's world and how many countries are doing too little too late in regulating social media;

**NOTING** the different type of cyber compliances needed to ensure that body corporates remain safe from cyber attacks, cyber breaches, data breaches;

**UNDERSTANDING** how mobile apps are gaining importance from a legal point of view, with the majority of malicious cyber actors orchestrating their attacks through mobile apps;

**LOOKING** at the emergence of fake news and how malicious actors spreading fake news gaslight the normal people into believing such stories;

**TAKING** into account innovative methods of coming up with solutions to regulate quantum computing, the emergence and practical impacts of which threatens to make passwords a futile exercise;

**SEEKING** to need to regulate OTT Platforms and how the content on such platforms can be brought more in line with the accepted morals of any society;

**STUDYING** the psychological impacts of cyber bullying and why cyber bullying is gaining growing prominence in the cyberspace;

**ATTEMPTING** at creating the need for a set of generally accepted principles towards the activities of Public-Private Partnerships in cyberspace.

## 32.2 Key Recommendations

### 32.2.1 Recommendations to the Conference

*THE PARTICIPANTS OF THE INTERNATIONAL CONFERENCE ON CYBERLAW, CYBERCRIME AND CYBERSECURITY HEREBY CALL UPON THE INTERNATIONAL CONFERENCE ON CYBERLAW, CYBERCRIME AND CYBERSECURITY*

1. To continue working on various and evolving cutting-edge areas of cyberspace and their intersection with law and technology;
2. To specifically explore the cyber legal, cybercriminal and cyber security repercussions of emerging technologies like Artificial Intelligence, Blockchain, Internet of Things, Metaverse and Quantum Computing;

3.  To come up with a separate smaller conference outside of India as a buildup event into the main conference;

4.  To work on new challenges impacting the intersection of cybersecurity and law as also connected legalities and policy aspects;

5.  To connect and interact with other stakeholders at the international level so as to understand their respective works in the areas of the intersection of Cyberlaw, Cybercrime and Cybersecurity and to collaborate with them on areas of mutual interest, value addition and knowledge enhancement;

6.  To disseminate the recommendations of the conference to all stakeholders in the digital ecosystem across the world including governments, ministries, agencies, corporates, not for profits, international organisations and every stakeholder in their respective national jurisdictions who is either connected, associated or dealing with the aspects of cyberspace as also Cyberlaw, Cybercrime and Cybersecurity;

7.  To work with digital stakeholders in a post-Covid scenario to identify, discuss, elaborate, analyse and also examine the newly evolving areas on the intersection of Cyberlaw, Cybercrime and Cybersecurity;

8.  To contribute to enhanced capacity building by partnering up with appropriate partners in the direction of creating more capacities amongst all digital stakeholders;

9.  To work with prominent industry leaders, thought leaders, experts and the stakeholders in different parts of the world on the thrust areas of the Conference;

10. To understand and appreciate the cutting edge developments in cyberspace, partner with other like-minded digital stakeholders, and to disseminate more knowledge, awareness, information and capacities amongst various digital stakeholders in different parts of the world;

11. To examine the cyberspace ecosystem and emerging challenges in the context of the world post Covid-19 and explore the emerging

trends, aspects, issues pertaining to cyberspace and to study and analyze the emerging jurisprudence concerning emerging technological paradigms in different parts of the world and try to collate common legal principles concerning the same;

12. To work on cutting edge developments in cyberspace including Metaverse, Quantum Computing, Internet of Things, Internet of Behaviour, Blockchain and Artificial Intelligence amongst other various disciplines so as to identify key emerging trends in emerging technologies, which could have a direct relevance and bearing upon the day to day operations and business activities of all digital stakeholders;

13. To become a catalyst for discussions on the key aspects of Cyberspace, and its future trends and impact, including legal, ethical, social, policy and regulatory issues thereof and present an integrated, holistic and strategic view of the issues therein, whilst recognizing that there is an urgent need for effective global cooperation on cyber issues amongst all stakeholders;

14. To network, interact, collaborate and work with global stakeholders on various aspects pertaining to activities in the digital ecosystem which would have a nexus, connection or association of any kind whatsoever with the connected paradigms of Cyberlaw, Cybercrime and Cybersecurity;

15. To identify the distinctive legal, policy, regulatory and technical challenges, issues and nuances thrown up by growing cyber security breaches and the potential approaches being adopted in this regard to deal with the menace of increasing cyber security breaches;

16. To identify and collate emerging international best practices that are evolving in different parts of the world on how fresh approaches need to be adopted, for dealing with the growing Golden Age of Cybercrime;

17. To clearly define 'Duty of Care' in the digital ecosystem for responsible governance;

18. To identify as to how cyberattacks can negatively affect cyber sovereignty, and related concerns of protecting cyber sovereignty of a nation;

19. To collaborate and develop approaches on how to reach on an agreement on acceptable cybersecurity norms respecting the principles enshrined in the United Nation Charter and promote a secure and stable cyberspace;

20. To highlight the gaps, challenges and vulnerabilities arising due the advances in Artificial Intelligence and solutions to remedy the risks of the AI system, and to deliberate upon the potential of AI in ensuring governance;

21. To highlight the issue of protection of individual rights behind the huge potential of Big Data and to ensure that the benefits of Big Data are shared equitably with sufficient transparency and accountability;

22. To collaborate with various international, national, and regional stakeholders to work together on the issue of anonymity in cryptocurrencies, the related technical and legal challenges and development of appropriate regulatory framework regarding the crypto ecosystem;

23. To reach out to connected stakeholders in the Cyberlaw domain at global, regional and national levels so as to keep a track of the latest evolving trends, issues and aspects in cyber legal jurisprudence and to further see how national Cyberlaw approaches are increasingly being moulded to make them topical and relevant to latest technologies, with the passage of time;

24. To work on the intricate nuances of newly emerging cyber-ground realities post Covid-19 and how to deal with various the challenges thrown up by newly changed scenario of Cyberlaw, Cybercrime and Cybersecurity.

25. To study, examine and analyse the aspects of the new cyberage that governs the world post-Covid 19 and to identify key challenges and parameters for all stakeholders in this regard;

26. To explore the emerging trends concerning cyber-sovereignty and the various international legal challenges that the growing advent and consolidation of the concept of cyber-sovereignty is beginning to throw up for all digital ecosystem stakeholders as well to examine, from a global standpoint, the emerging phenomenon of balkanization of the internet and the connected legal policy issues connected therewith;

27. To explore how the misuse of the Darknet for criminal purposes can be appropriately addressed by State actors over the passage of time and to explore how the principles of attribution of cyber-acts done by cyber-actors can be further effectively evolved;

28. To come up with emerging cyber legal issues and challenges, post the commencement of the Ukraine War and their impact upon the global cyber ecosystem;

29. To encourage the adoption of more cyber resilient tools for digital stakeholders in view of the increased adoption of internet by all digital stakeholders and to promote effective information exchange for enhanced international cooperation on matters relating to cyber-crime;

30. To study and explore along with the stakeholders on how cogent and effective steps can be taken to protect the privacy of users in the digital ecosystem and to disseminate awareness about the interconnected paradigms of Cyberlaw, Cybercrime and Cybersecurity and to inspire discussions at national, regional and international echelons on the intersection of Cyberlaw, Cybercrime and Cybersecurity to meet the complex challenges by focusing on achieving a secure, robust and resilient cyberspace;

31. To collate and interact with all international, regional and national projects and initiatives, programs and schemes which are being im-

plemented having a direct impact upon the interconnected paradigm of Cyberlaw, Cybercrime and Cybersecurity;

32. To foster more discussions, dialogues and exchange of ideas, perspectives and opinions on the growing cyber legal challenges to the cyberspace ecosystem thrown up by various emerging paradigm;

33. To contribute to continuing international discussions upon evolving norms of behaviour regarding cyber operations;

34. To provide strategic leadership across Governments to respond to cyber security threats against identified Critical Information Infrastructure (CII);

35. To explore and analyse the various regional and national legislations on data protection and privacy and their contribution to strengthening the secure digital environment

36. To work towards decentralised data transparency, and legal foundations for peer-to-peer currency exchange, smart programmable contracts and for regulating crypto assets;

37. To identify and highlight the legal nuances and challenges raised by fake news globally and encourage the crystallisation of appropriate legal response to deal therewith.

### 32.2.2 Recommendations to the United Nations and International Agencies

*THE PARTICIPANTS OF THE INTERNATIONAL CONFERENCE ON CYBERLAW, CYBERCRIME AND CYBERSECURITY HEREBY FURTHER CALL UPON UNITED NATIONS AND ITS VARIOUS AGENCIES INCLUDING ITU, UNESCO, WIPO AS WELL AS ALL INTERGOVERNMENTAL ORGANISATIONS (IGOs) & NGOs*

1. To understand that the issue of jurisdiction needs to be tackled in the context of cybercrimes;

2. To explore the drafting of an international legal instrument on the issue of Cybersecurity;

3. Understand that the emergence of cybercrimes puts people in developing and under-developed countries at serious risk of exploitation;

4. To update the existing treaties on outer space, keeping in mind the latest advancements in cyberspace, both technological and legal changes at the national levels;

5. Ensure that technologies like artificial intelligence do not violate the most cherished human rights;

6. To identify a global definition of cyber-sovereignty to be adopted by nations;

7. To call upon tech companies to work with governments and find solutions to ensure the safety and security of citizens, without eroding user privacy or cyber security;

8. To collaborate with various international, regional and national stakeholders to work together in the development of legal jurisprudence on cyberspace;

9. To strengthen the cooperation on Cybersecurity law, by creating more opportunities for governments, private sector, civil society, the technical community and academia from various regions of the world to engage and develop innovative and effective legal frameworks, to address the truly global challenge of Cybersecurity;

10. To help in deciding as to how Metaverse must be regulated and to educate people about the cyber legal ramifications of highly interactive and inter operable nature of the Metaverse;

11. To encourage the development of effective instruments on cyber-diplomacy, since this diplomacy will be needed due to the lightening fast advancements made in technology;

12. Effectively address the issue of removing the vacuum on cyber security law existing on the international level and provide aid to states not having their own municipal law on the issue of Cybersecurity;

13. To come up with a common understanding at the international level on how principles of cyber sovereignty have to evolve and how conflicts between different cyber-sovereignties of different jurisdictions need to be resolved;

14. To holistically take up comprehensive viewpoints on emerging technological paradigms in cyberspace;

15. To invite the conference's deliberations and expertise in the various UN bodies and processes for the purposes of appropriate dissemination of expertise on various cutting edge issues concerning cyberspace, Cyberlaw, cybercrime and Cybersecurity;

16. To encourage the developments of norms concerning responses to breaches of cybersecurity by nation states and their relevance in the international scenario;

17. To explore identification of common legal principles, which can be made applicable concerning breaches of cybersecurity in outer space and push for their appropriate acceptance and therefore, pre-regulation by various nation-states;

18. To work towards specifically protecting women and children and other vulnerable sections of the human population against constant abuse, misuse and targeting on the internet and in cyberspace;

19. To contribute towards the gradual evolution of cyberlegal jurisprudence at the international level;

20. To encourage countries to come up with stronger national laws on cybercrime so as to deal with the growing menace of cybercrime;

21. To help generate discussions and debate about legal principles impacting cybersecurity which can be adopted at the international level by various countries;

22. To encourage more enhancement of capacity building amongst all nations states specifically on the intersection of Cyberlaw, Cybercrime and Cybersecurity;

23. To encourage nation states to disseminate more information about and encourage the adoption of cyber-resilience as a way of life amongst various nation states and their respective digital stakeholders;

24. To invite the contribution of the conference and its various experts to various initiatives and processes in the intersection of Cyberlaw, Cybercrime and Cybersecurity

25. Suggest ways by which international instruments relating to cyberspace can be made applicable to all countries;

26. To encourage that the emerging trends of cyberspace are not left unregulated and that in the absence of international instruments and municipal laws, there need to exist self-regulatory mechanisms on the national and international level;

27. In the international instruments dealing with cybercrime and cyber security issues, also lay down principles governing how cyber forensics must be conducted and what tools are permissible;

28. To come up with ways to take Public Private Partnerships forward in the cyberspace through ethical and legal means.

### 32.2.3 Recommendations to Nation states and Stakeholders in Education, Media, Business, Research, social and cultural sectors

*THE PARTICIPANTS OF THE INTERNATIONAL CONFERENCE ON CYBERLAW, CYBERCRIME AND CYBERSECURITY HEREBY CALL UPON NATION STATES, EDUCATION & RESEARCH INSTITUTIONS, PROFESSIONAL ASSOCIATIONS, MEDIA INSTITUTIONS, CULTURAL & SOCIAL INSTITUTIONS AND ORGANISATIONS,*

*NETWORKS, BUSINESS, CORPORATE & INDUSTRY SECTORS AS WELL AS ALL RELEVANT STAKEHOLDERS*

1. To come up with dedicated legal frameworks on Cybersecurity as well as to examine and explore the elements of cyber-resilience emerging in the digital ecosystem and the connected legal, technical and other logistical elements that are evolving to be integral parts of the cyber-resilience paradigm;

2. To come up with more innovative projects, initiatives and exercises aimed at identifying some of the emerging cutting edge trends that are emerging in the cyber landscape;

3. To explore how the issue of online gaming is beginning to throw up complicated issues for all stakeholders;

4. To work on commonly accepted principles for regulation of fake news and dissemination of fake information on social media, internet and cyberspace at large;

5. To work in the direction of creating more safe, secure, reliable and resilient cyber ecosystem and to work in the direction of creating cybersecurity as a way of life;

6. To explore how the legal and policy challenges thrown up by metaverse can be appropriately regulated in a minimal way;

7. To explore how cyber-hygiene practices have to be disseminated amongst all stakeholders in the digital ecosystem on a continuing basis in order to contribute into making the cyberspace more safe, robust and reliant;

8. To regulate online abuse, cyber bullying and online defamation, since the internet is the fastest mode of communication and it is important to come up with ways to safeguard a person's reputation from being falsely tarnished;

9. To come up with data ethics for each industry and sector in accordance with the unique needs, working method and business model of each industry and sector;

10. To come up with self-regulatory mechanisms to control technology and to ensure relevant mechanisms are in place to protect data of consumers, customers and all those whose data is stored;

11. To make the armed forces technologically advanced enough to ward off attacks on a nation's sovereignty in the cyberspace;

12. To specially focus on ransomware, which is readily available on the darknet;

13. To enter into cyber diplomatic agreement with foreign actors;

14. To understand that the emergence of metaverse is significant not only for gamers, but also for lawyers since it creates new kinds of legal issues, like Intellectual Property Rights issues, criminal issues and hence, new laws need to be made accordingly;

15. To come up with a definitive regulatory mechanism on the block-chain technology, since it is responsible not only for creating a whole new man-made currency which can be easily converted into fiat currency, but may very well be the future of the internet in the form of Web 3.0;

16. To keep the growing power of social media intermediaries in check and to not make them super powerful as well as to try and put a stop on cyber bullying, with the act negatively discouraging people from freely exercising their freedom of expression;

17. To come up with a definitive law on mobile apps, with handheld devices slowly replacing devices like desktops and other stationary devices;

18. To make a proper law on Telemedicine to cover the legal issues connected therewith and to lay down specific guidelines as to the extent to which Telemedicine can be practised and a patient be diagnosed without being near the doctor physically;

19. To find solutions to the threats posed to information infrastructure by cybercrimes like phishing, identity theft and fraud including encouraging crystallisation of holistic legal, policy and regulatory respons-

es to tackle the menace of the aforesaid infamous trinity of cyber-crimes;

20. To identify the threats to the banking and financial sector and Insurance services due to cyberattacks and measures that can be taken to ensure preparedness;

21. To identify the cyber security aspects of space satellite ecosystem and its profound impact on a myriad of sectors and activities built on the satellite information systems, thus reiterating the interlinkages between cyberspace and outer space;

22. To bring into focus the challenges presented by new technologies such as deep fakes, promote awareness of new technologies and accompanying potential risks connected therewith;

23. To work towards ensuring that the security of upcoming 5G technology is paramount as well as to identify the legal principles and frameworks required to deal with the cyber security and technical challenges thrown up by the advent of 5G;

24. To coordinate, share, monitor, collect, analyse and forecast, national level threats to CII for policy guidance, expertise sharing and situational awareness for early warning or alerts as well as to engage in the risk analysis process, and related mitigation or counteraction strategies;

25. To make sure that an AI system does not hamper fundamental rights, thanks to various strategies including certification, tests, periodic checks, open source code analysis;

26. To ensure support for appropriate privacy and personal data protection and to work towards safeguarding personal data and privacy in the age of the internet, promoting awareness both of technologies and potential risks;

27. To provide a catalogue of security related institutional legal controls to meet current information protection needs and the demands of future protection needs, based on changing threats, requirements and technologies;

28. To identify, though an interdisciplinary approach, the future possibil-
    ities, scope and impact of emerging technologies such as IoT, Arti-
    ficial Intelligence and Machine Learning as well as to identify future
    emerging technologies (FET) and the legal, policy and regulatory is-
    sues related therewith.

*WE, ICCC PARTICIPANTS, URGE THAT JOINT EFFORTS NEED TO BE TAKEN BY ALL RELEVANT STAKEHOLDERS TO MAINTAIN THE INTRINSIC CHARACTER OF CYBERSPACE, WHICH IS SAFER, MORE RESILIENT, AND REMAINS THE MAJOR DRIVER OF SUS-TAINABLE ECONOMIC DEVELOPMENT AND GROWTH FOR YEARS TO COME.*

*WE REITERATE THAT STEPPING FORWARD IN A NEW ERA OF DIGITAL AND CYBER WORKSPACE, WE ALL NEED TO BE SAFE, SECURE AND DILIGENT, WHILE ENCOURAGING FURTHER AD-VANCEMENT IN CYBERSPACE AS WELL AS INFORMATION AND COMMUNICATION TECHNOLOGIES (ICTs).*

*WE, THE PARTICIPANTS OF THE INTERNATIONAL CONFER-ENCE ON CYBERLAW, CYBERCRIME AND CYBERSECURITY, AP-PROVE AND ADOPT THE ABOVE OUTCOME DOCUMENT.*

# 33

# ETHICAL DILEMMAS OF ARTIFICIAL INTELLIGENCE SYSTEMS

*Alexander Ageev, Russia*[453]

## 33.1 Introduction

Currently[454], the replacement of humans by AI systems is taking place everywhere. Mostly modern AI ("weak AI") operate within

---

[453] The author Alexander Ageev is Professor of Economics at various universities in Moscow, with a PhD also in theology. He is Founder and Director of the Institute for Economic Strategies in Moscow, General Secretary of the International Research Institute for Advanced Studies IRIAS and Director of Globethics Russia.

[454] Article translated from Russian. First published in: Агеев А.И. *Этические дилеммы систем искусственного интеллекта*, в книге *Социогуманитарные аспекты цифровых трансформаций искусственного интеллекта* под редакцией В.Е. Лепского, А.Н. Райкова. Москва, Когито-Центр. 2022. Глава 3.5. © Globethics Publications, 2023 | DOI : 10.58863/20.500.12424/4276064 | CC BY-NC-ND 4.0 International.

knowledge and goal-oriented contours specified by the developers. Accordingly, posing the question of AI ethics at this level of its development is a purely declarative, though useful, exercise. The real question about the rules of conduct of AI arises when and if AI will be able to autonomously change the contour of knowledge used, learn from unlimited sources and types of information, adjust or even set its own goals. The second fundamental problem is related to the ethical situation in "ordinary" human society, where ethical relativism has formally prevailed, which in reality means the veiled dominance of quite certain ethical models. This creates the prerequisites for the inevitable translation of this situation into the space of development and application of AI.

## 33.2 Genuinely Aware Subjects and Artificial Intelligence Systems (Ethical Aspects)

Under certain conditions, it is possible to apply the definition of genuinely aware subjects (GAS), capable of selecting activity goals and arbitrarily working with different knowledge bases, to AI. In part, this reflects the dependence of AI on the goal-setting of its developers, who can put explicit or implicitly criminal decisions into AI. But to a greater extent, the emergence of GAS among AI depends on scientific and technological progress and its dissemination in everyday life.

GAS refers to living or quasi-living beings that have self-awareness and subjective experiences similar to those of humans or other highly evolved beings. Such objects could probably include AI in its strongest version . Strictly speaking, individuals with implanted "smart" prosthetic devices (from cardiac pacemakers to prostheses of the musculoskeletal system or individual sense organs) should also be classified as PIC. The intellectual component of such devices is developing rapidly. The number of people with such devices is growing and is measured in the tens of millions worldwide. It is easy to see the boundary where the operation of these devices collides with the problem of moral choice. In addition, as a

"social creature," man in today's world is in a number of living and virtually organized networks that can influence his daily, subject-situational, and even attitudinal choices. This means that he is influenced by supra-personal virtual systems, as well as by "big user data" collected independently of the individual's will.

Moral and ethical concerns stem not only from how arbitrary and/or harmful such a system can behave to the individual and human society, but also from the fact that society itself can inflict suffering on the artificial entities it creates when they reach the GAS level.

A similar moral dilemma has arisen in scientific and practical experiments on humans and animals. At the moment it has been resolved as follows: experiments on humans are forbidden without their consent, animals are considered less evolved and experiments on them are allowed. Recently, however, there has been a tendency to equate certain highly evolved animal species with humans, assigning them moral and legal rights. A precedent occurred in the United States: two primates participating in biomedical research at the State University of New York at Stony Brook were recognized as having some human rights.

When creating artificial beings and quasi-beings, it is not always possible to predict the outcome of an experiment and the extent to which a conscious being will eventually emerge, given the learning potential of such a system. In addition, it is impossible to ask the creature's consent to the experiment before it is created. This problem is also solvable: the creature does not require its consent to be born, but it receives certain rights and guarantees provided by law and custom immediately after birth, and in part before. Similarly, the creators of artificial beings, including SIs, can be required to preemptively minimize their "suffering" during experiments, also to keep these beings in proper conditions. Otherwise, when consequences cannot be predicted, such experiments should be prohibited.

An even more difficult question is whether, in the future, intelligent systems will be able to take on the risks of decision-making while solv-

ing some tasks on vital social problems. Is society ready to take responsibility for the activities of such a technical system? Because of the complexity of the tasks solved by the AI, the impossibility of full control over the machine, tracking and checking the decisions made, a person is already forced to shift part of the responsibility to the machine (program).

There is a certain degree of uncertainty in the development of AI, which is similar to human free will. But at the same time some principles must be laid down in the AI to be developed:

ethical self-restraint and knowledge of people's moral and ethical standards, imitation of the process of self-regulation of one's behavior, ability to empathize;

a mechanism for predicting the risks and consequences of one's own actions, limiting actions when certain risks occur;

the possibility of recognizing and correcting one's own mistake .

Further research in the field of AI ethics will undoubtedly lead to the creation of various standards and certification of rules for the design and operation of AI. In creating such standards, it is necessary to assess their impact on the further development of AI, eliminating the risks of both malicious development and inhibition of the development of intelligent systems.

## 33.3 Types of Mass Consciousness and the Challenge of Ethical Relativism

In the 20th century an unprecedented experience of manipulating consciousness, of purposefully forming mass consciousness defined by an initiator, usually the state or with support from state institutions, of personality types and parameters of their permissible ("ethical", socially acceptable) behavior has been accumulated. This experience has largely been discredited in recent decades, giving rise to a crisis of basic ethical concepts that have developed over the last millennia of human history

and have formed the dominant ethical relativism, the growth of multiculturalism along with the aggravation of conflicting identities. However, the discrediting of this experience does not mean either its disappearance or the loss of its potential to manifest social energy in the future. We are talking about a certain accumulated repertoire of social being, in which, depending on the combination of organizing and self-organized processes, a set of motivating and stimulating layers of culture for socialization, adaptation, subjectification of man and his communities on different grounds, reflecting biological and social reproduction and interaction with the outside world emerges. There is no doubt that the development and implementation of new technologies, including AI, will be influenced by the ethical state of society. There is also every reason to believe that advances in science and technology will increasingly influence collective and individual consciousness.

An ontological map of social-value orientations ("pictures of the world," "I (we)-concepts") can be formed on many grounds, in many ways similar to the traditional membership of philosophical schools (versions of idealism and materialism, gnoseological and existential concepts, etc.). Among these criteria are, for example, idealistic vs. materialistic, subjectivity vs. objectivity, generality-partiality, personalization-impersonalization, etc.

In the considered perspective the approach to the formation of the ontological map from the point of view of generalized personality types (GPT), taking into account the historical experience of the twentieth century, significant for social practices in the modern period, seems meaningful and useful. These issues are partly understood in a technocratic perspective as the formation of Industry 4.0 societies, in a sociotechnical perspective as "inclusive capitalism," and in a more complex perspective as the creation of Society 5.0.

In any case, the immense variety of vital problems, against the background of the increasing information flood and the expanding use of information technology for processing information arrays and flows

(from the structure of the universe and its particles, road traffic, the regulatory framework to human brain research) provokes more and more people to turn to various techniques for ordering and simplifying this hectic and complex reality, including dexterity and complexity, for certainty in their identification and choice of behavior and life in general in this confusing world.

Under the conditions of "multiculturalism," ethical relativism is theoretically and de facto inevitable, manifesting itself as the existence of a certain, seemingly unlimited "menu" for the choice of personal identity. However, this does not mean that every line on this menu is equal. Among the entire set of ethical concepts offered or allowed to be developed "from scratch" (an option implicit in personal creativity that is almost entirely manipulative) there are both heavy and very light, situational, conjunctural, usually little conscious, but convenient fractions. "Heavy fractions" rely on conceptual, theoretical and/or religious solutions to deep questions of human existence in the world and society in particular. They tend to have a powerful system of reproduction, institutionalized in the social structure, in management, science, education, and anchored in architecture and art. Their gravitational pull is obviously higher than any of the newfangled attempts to create a "new ethics" and, as its political projection, an ideology. But the "heavy fractions" of ethical systems and constructs, representing the attractor of traditionalism, always experience difficulties in adapting to the new challenges of evolution formed by the progress of science, technology, lifestyles, political and economic struggles.

Technological changes and active project development have brought to life a diverse assortment of "new human" concepts ("digital humans," "GMO-humans," "transhumanism," "service people," "singulars," "new Europeans," "new nomads," etc.).

The ensemble of ideas fed by the concepts of noosphere, cosmism, "radiant humanity", socio-spiritual integrationism is also gaining strength. Behind it there is a substantial scientific basis going back to

V.I. Vernadsky and T. de Chardin, twentieth-century fantasts (H. Wells, I. Efremov, S. Lemm, Strugatsky brothers), philosophers and scientists (N. Fedorov, K. Tsiolkovsky), but more important - behind it practical triumphs of human development in the mid-20th century. - exploration of the Earth, near-Earth and outer space, the depths and expanses of the ocean, matter, human biology and psyche, etc. However, this set of ideas, for all their scientific validity, is complex, relying on new scientific paradigms that are poorly known in society. It is still far from being manifest as a new mass ideology, let alone an GPT. More popular are "geographical," "geopolitical," and "geoeconomic" versions of ideologies that raise to the flag the mechanical proximity of the residence and destiny of certain peoples and states and serve rather pragmatic interests.

All this diversity of GPTs affects scientific and technological solutions in the field of digitalization. Generally speaking, as AI  reaches a level of development at which there is a need for certain ethical constraints on goal-setting, on the knowledge and behavior used, approaching AI at the GAS level, a "portfolio" of ethical principles and practices will be rapidly formed and regulated for the AI being developed. In this process, ensembles (sets) of ethical templates (matrices) developed at the national and civilizational levels will inevitably be projected onto artificial entities. With all the universality of technical and technological solutions it is digital technologies that have created a fundamental opportunity not only to imitate real individuals and processes in the form of their digital twins, but also to generate purposeful strategies to influence clusters of individuals, using the potential of identifying, observing and using their psychosemantic, biophysical and other properties. In other words, the possible future "battle of the machines" will in any case be a "battle" of quite human and humanoid ethics encrusted in the AI and the growing set of GASes.

## 33.4 Digital Society: Ethics and Trust in Artificial Intelligence Systems

The question should also be raised about the point at which society crosses a certain technological threshold, when there may be an irrevocable transformation of both society as a whole and the individual into what can conventionally be called a "digital society." This bifurcation point appears to be associated with the formation of a digital stratification platform as a dominant feature of social structuring and management. AI will play a key role in this platform and the formation of appropriate social strata, based on the criteria set during the development. The necessary conditions for this include, first of all, the achievement of a technologically possible, virtually complete awareness of the critical parameters of the life of society and the entire set of individuals.

The experience of using digital technologies in the 2020-2021 pandemic has shown both the high potential of AI for social management and their many software and hardware and socio-psychological vulnerabilities. This is also evidenced by the experience of all kinds of digital platforms (from government services to marketplaces). Nevertheless, digitalization covers more and more spheres of social life. The prospect of comprehensive integration of the created surveillance systems, databases, data processing centers, and decision support systems is being seen. The possibilities of manipulating personal choice through digital personalized models and their clustering are outlined above.

The logic of the AI development shows that they increasingly implement the principle of self-construction not only in software, but also in hardware. For example, if so far, based on input and output of information in neural networks, it is possible to reconstruct the ways of its processing and, accordingly, to obtain human-readable rules and algorithms, then in the near future the situation will change. At the IARPA annual conference back in 2016, it was discussed that the intelligence community uses complexes with such deep neural networks that it is

impossible to translate their algorithms into human-readable language in a reasonable amount of time. These are so far the first symptoms that a world is emerging where decision-making will be based on criteria closed to the decision makers.

In 2012, Human Rights International raised with the UN the issue of the need to ban autonomous combat systems (robots) that make their own decisions about the use of combat equipment in their possession. Despite strong support from a number of governments, no binding UN decision has been adopted to date.

A similar problem arises in the case of robotic vehicles. One way or another, all robotic vehicles will be controlled by an AI system of varying degrees of sophistication. Automakers approach this problem in different ways. For example, Mercedes issued a statement in 2016 that in the case of road rules, the company's duty is to protect passengers, not pedestrians. In any situation, if a robotic Mercedes decides it is following the rules, the choice will be made in favor of passengers, not pedestrians. Google has taken a fundamentally different stance. The new generation of Google cars has a software video filter. If the car's video sensors recognize a child on the road, regardless of whether the child is breaking the rules or not, the car will choose saving the child's life as its first priority.

The second condition for the transition to a "digital society" is determined by the real interests of modern society itself. First of all, it is socially heterogeneous almost everywhere, which predetermines the differences not only in access to digital technology, but also in their development, implementation and use. In this sphere there is already a fierce competition for leadership and control. However, practical management in society as a whole and at the level of integration, national, regional, sectoral formations is still determined not by digital imperatives, but by political national-state, corporate and private interests. They are by no means unconditionally favorable to the rapid development of "digital society", while it requires the total coverage of all subjects,

objects and processes in society. It is enough to note the risks of mass unemployment, the level of cybersecurity to be cautious about the most optimistic pace and stages of the formation of "digital society".

The evolution of cyber-physical systems to the status of global awareness and successful management of the evolution of society is seen in the level of centralization-decentralization of the global awareness function and the combination of institutions of organization and self-organization of life activities, including the economy. Currently, various states have enacted regulations limiting the degree of centralization of personal data.

Nevertheless, the development of technological solutions for a new generation of public systems is underway in all leading states and corporations around the world. In the USA, one of the promising DARPA projects was the creation of a "dynamic virtual environment" in which the barriers that exist today (departmental, organizational, informational and technical) will be eliminated for effective and prompt joint work of representatives of various ministries and agencies engaged in crisis situations resolution in various areas of activity (political, military, economic and social). The decision-making tools being developed for territorially distributed groups of interests are intended to provide the fullest possible comprehension of complex situations and scenarios of their dynamics, the choice of optimal solutions based on all available information without fully studying it by the principle of "knowing without reading". The technology is based on methods of fuzzy structuring of arguments, three-dimensional color visualization and corporate memory.

The development of AI and digital transformation in general leads to the formation of collectives of autonomous agents of artificial and mixed genealogies, as well as complex constructions of information and regulatory environments with multiple possibilities and pathologies and increasing levels of uncertainty for managerial decision-making.

At the same time, there is an increasing phenomenon of "degradation of natural intelligence". In particular, constant surfing of websites can

lead to an erosion of the capacity for systemic and in-depth thinking. Medical research has conclusively proven that those who spend a lot of time on the Internet quickly develop two areas of the brain-the part responsible for short-term memory and the center responsible for making quick decisions. At the same time, those areas of the brain responsible for detailed analysis, deep thinking problems remain without load and gradually lose the skill to work intensively.

Today among the most acute issues of AI development is the problem of trustworthiness, which covers the problem of confidence of consumers, regulators and other stakeholders that the AI is able to perform its tasks with the required quality and safety level.

The world of social structures is projected to syndicate major social platforms and redefine standards for human interaction between 2020 and 2030. Mutual translation of neurodescriptions, social descriptions and descriptions of the semantics of the human Internet and the Internet of Things is being established. "Codes" of nervous systems and the brain will largely be described and used not only in medicine, but also for modeling similar processes in other substrata - economic and social systems, self-organization of "smart things" and artificial systems. By 2030 semantics of different types will be able to translate into each other, and this will be used in experimental settings. Mental modes will be described fully enough, including states of consciousness in relation to different types of activity. The structure of a person's consciousness can be easily reconstructed depending on the tasks facing him. Neuronet interfaces are absolutely invisible, transparent. A person does not work behind a keyboard and a screen - he works directly with data, with meanings, with people.

The subject of control at this stage is the human body, represented by a large amount of data from sensors of different types. This also includes collectives, of which a person is currently a part. There is also interaction with distributed systems of smart things, which are constantly rede-

signing themselves. One such subject of control is human life, education throughout the life cycle, and the life cycles of communities.

## 33.5 Conclusions

The question of the ethics of AI at this level of its development is purely philosophical and futurological. The real problem of the "ethics" of operation (behavior) of AI arises only when and if AI will be able to autonomously change the contour of the knowledge used, learn from unlimited sources and types of information, adjust or even set its own goals. This prospect can only be seen in the emergence and expansion of "strong AI". Nevertheless, there are already a significant number of beings, including ordinary humans, equipped with various "smart" devices, whose behavior can pose ethical problems. This also affects artificial systems that obey the logic of the behavior of the "swarm" that affects the individual.

The second fundamental problem is related to the ethical situation in "ordinary" human society, where ethical relativism has formally prevailed. This creates the prerequisites for the inevitable translation of this situation into the space of development and application of AI. As a result, it is very likely that the "machine world" will reproduce the human world by transferring the existing ethical problems to the cyber-physical world.

## 33.6 Bibliography

Ageev A.I., Asanova E.A., Gribenko O.V., etc. Are you ready for the "digit"? Assessment of the adaptability of Russia's high-tech complex to the realities of the digital economy / Edited by Doctor of Economics, prof. A.I. Ageev. M.: INES, 2018. 60 p.

Approaches to the formation and launch of new industries in the context of the National Technological Initiative on the example of the sphere "Technologies and systems of digital reality and promising "human-computer" interfaces (in terms of neuroelectronics)": Analytical report URL: http://rusneuro.net/cambiodocs/media/files/analitijeskii-doklad-podhodyk-formirovaniu-i-zapusku-novyh-otraslei-promyhlennosti.pdf.

Cadwell et al, for monitoring everyday prosthesis use a systemanic review/ Journal of NeuroEngineering and Rehabilitation. 2020. https://qz.com/1577451/century-tech-signs-deal-to-put-ai-in-700-classrooms-in-belgium/

Internet surfing changes the user's brain, SecurityLab. 2010. 30 August. URL: https:// www.securitylab.ru/news/397247.php.

Jenny Anderson: A British start-up will put AI into 700 schools in Belgium //QUARTZ. 2019. 21 March. https://qz.com/1577451/century-tech-signs-deal-to-put-ai-in-700-classrooms-in-belgium/

Kukshev V.I. Digital Economy: Problems and Solutions, Economic strategies. 2020. № 5, pp. 51–57;

Kukshev V.I. Classification of Artificial Intelligence Systems, Economic strategies. 2020. № 6, pp. 58–67.

Razin A.A. Ethics of artificial intelligence, Philosophy and Society. 2019. Issue.1 (90), pp. 57–73. https://doi.org/10.30884/jfio/2019.01.04.

Report of the Deputy Director of the FBI, Head of CJIS Stephen L. Morris "Artificial intelligence: the FBI and the police against criminals", https://ordrf.ru/wp-content/uploads/2017/

10/Обзор-отдельных-вопросов-использования-
больших-данных-и-искусственного-
интеллекта.pdf#page=10&zoom=100,72,537

Shuravin    A.    The    History    of    Artificial    Intelligence,
https://wiki.programstore.ru/istoriya-iskusstvennogo-
intellekta/

Socio-economic   aspects   of   the   introduction   of   artificial   intelligence,
Under   the   scientific   editorship   of   A.I.   Ageev.   M.:   It-
Service, 2020.

Van Dyk: U.S. Court Recognizes Chimpanzees as Legal Persons,
BuzzFeedNews. 2015. 21 April.
https://www.buzzfeednews.com/article/mbvd/us-court-
recognizes-chimpanzees-as-legal-persons.

# 34

# AFRICAN UNION
# DATA POLICY FRAMEWORK

*African Union, Ethiopia*

*The African Union, representing its 55 member states from the whole African continent, endorsed its "African Union Data Policy Framework" by the AU Executive Council during its 40th Ordinary Session on 2-3 February (Ref EX.CL/Dec.1144/XL) and published in February 2022. The below excerpt are Executive Summary (pp VI-X) and the recommendations and actions of chapter V (18-64). The full text is available for free download[455].*

*In addition, a very helpful interactive website provides the national legislations on data protection for each African country.[456]*

*The Editor CS*

---

[455] African Union Data Policy Framework. 2022. Documents, https://au.int/en/documents/20220728/au-data-policy-framework.

[456] Africa Data Protection, Liste des lois, https://africadataprotection.com/liste-des-lois.html.

## 34.1 Executive Summary

Data is increasingly recognised as a strategic asset, integral to policy-making, private and public sector innovation and performance management, and creating new entrepreneurial opportunities for businesses and individuals. When applied to government services, emerging technologies can generate massive amounts of digital data and significantly contribute to social progress and economic growth. The central role of data requires a high-level and strategic policy perspective that can balance multiple policy objectives - from unleashing the economic and social potential of data to the prevention of harms associated with mass collection and processing of personal data.

The purpose of this document is to provide the policy framework for African countries to maximise the benefits of a data-driven economy by creating an enabling policy environment for the private and public investments necessary to support data-driven value creation and innovation. This enabling environment refers both to the collaboration between in-country sectors, institutions and stakeholders, an alignment of their development priorities, and the harmonisation of policy across the continent in a manner that provides the scale and scope required to create globally competitive markets.

From a policy perspective, the approach adopted is people-centred, locating them in relation to the role of data in contemporary economy and society by identifying the elements and linkages in what can be called the "data ecosystem" in order to identify the exact points of policy intervention. This allows for a systemic assessment of the interrelated challenges arising from global developments that impact emerging national data economies and those arising within the context of nascent data-driven economic activity, uneven institutional endowments, and human development in many African countries. This enables the design of a contextually grounded but forward-looking data policy framework that uses economic regulation to guide policy makers in realising oppor-

tunities for data-driven value creation. The framework points to how opportunities can be realised and how associated risks could be mitigated by creating an enabling and trusted environment.

Building a positive data economy national and regional will require unprecedented levels of collaboration between stakeholders to disrupt the economic, political, and policy pressures already being felt from the global data economy. In order to ensure equitable and safe access to data for innovation and competition, Member States should establish a unified legal approach that is clear, unambiguous and offers protection and obligations across the continent. Existing legal instruments and institutions should be revisited where necessary to ensure that they are not in conflict with one another and that they offer complementary levels of protection and obligations.

A comprehensive data strategy will necessarily include the harmonisation between competition, trade, and taxation policies and laws both at the national and regional levels. This is so an optimised data ecosystem for Africa balances revenue mobilisation and the need to avoid distortions to local markets and the global tax system. Intellectual property laws should also be revised to clarify that they do not generally impede the flow of data or data protection. At the same time, governments need to develop transversal digital policies and strategies to coordinate activities across the public sector and between the public and private sectors to meet national objectives.

While there are multiple competing definitions of data, common to all is the recognition that there are many different types of data. There are also numerous ways that data can be categorised that affect the appropriate policy and regulation of that category in order to mitigate any potential risk associated with the processing, transfer or storage of it. A primary distinction is between personal data and non-personal data, with data protection referring to ensuring the privacy of data subjects. Data categorisation guidelines should be one of the first actions of the data information regulator, a key institution for the development of an inte-

grated national data system, which should be established in partnership with all relevant stakeholders. Essential to the development of an enabling environment for the data economy is ensuring the necessary foundational digital infrastructure and the human resources necessary to develop data as a strategic asset. Due consideration needs to be given to developing robust Digital ID systems for the delivery of public and private value to citizens and consumers.

As the framework also emphasises, this can only be properly achieved through instilling a culture of trust in the data ecosystem. This is done through the establishment of safe and secure data systems based on effective cybersecurity and data protection rules and practices, and ethical codes of conduct for those who set data policy, implement it and those who use data – whether in public, private or other sectors. This is not sufficient, however. Trust in data governance, and a national data system is established through legitimacy. This includes systems and standards that guarantee public and private sector compliance, government itself adhering to personal data protection rules, and government sharing public data.

The framework instils the importance of collaborative and evidence-based policy processes for the domestication of the policy proposed. The governance and institutional arrangements should assign clear roles to the government as policy maker and independent, agile and capacitated regulators to implement policy and effectively regulate the data economy to ensure that fair competition produces positive consumer welfare outcomes. The creation of data and information regulators to promote and safeguard the rights of citizens and their participation and fair representation in the data economy and society will need to be a priority for those countries that have not yet established these. Coordination with other regulators to achieve this will be essential. The legal ecosystem must be harmonised and rebalanced.

Access to data is a prerequisite for value creation, entrepreneurialism and innovation. When data are of poor quality or not interoperable, they

limit the capacity of firms and the public sector to engage in the sharing and analytics that can provide economic and social value to data. These processing frameworks should align with the following principles: consent and legitimacy; limitations on collection; purpose specification; use limitation; data quality; security safeguards; openness (which includes incident reporting, an important correlation to cybersecurity and cybercrime imperatives); accountability; and data specificity. Security models also need to be transversal, with specific emphasis on cloud storage and processing of sensitive/proprietary data, API management, and support of equitable data economies.

Attention needs to be paid to access to quality, interoperable and reliable data – primarily from the state but also from the private and other sectors – with a reinvigoration of the principles of open governance across the continent. Capacity-building should be a key national and regional priority, and resources will need to be allocated in this regard in the areas of data protection, cybersecurity and institutional data governance in relevant agencies. Skills and an understanding of the data ecosystem will also need to be built in state institutions, amongst other sectors and communities.

The framework is guided by the broad principles of transparency, accountability of institutions and actors, the inclusion of stakeholders, equity among citizens and fair competition amongst market players. The principles guiding the framework include trust, accessibility, interoperability, security, quality and integrity, representativity and non-discrimination.

As the framework emphasises, transversal collaboration needs to be underpinned with mechanisms to stimulate demand for data, which includes incentivising innovative data communities, and, on the supply-side, ensuring the quality, interoperability, and relevance of data in both the public and private sectors and civil society.

As the framework suggests, there are several regional processes, mechanisms and instruments that can and should be leveraged in the

continent's efforts to develop a cohesive data policy framework. These include the African Continental Free Trade Agreement (AfCFTA), which provides an opportunity for cooperation on a number of important aspects of the policy framework. Collaboration between national and regional stakeholders is also necessary for African countries to become more competitive in global policy setting forums where regulations for the global data economy are set and where African states have largely been "standard takers".

It is recognised that different African states have different economic, technical, and digital capabilities, and the recommendations and actions need to be read in this light. It is nevertheless envisaged that the different demands of building a data ecosystem will be progressively realised by countries. At the same time, there are several areas that can be attended to independently of economic or technical capabilities, including establishing regulatory independence, promoting a culture of trust and ethics, building collaborative frameworks for relevant sectors, developing transparent, evidence-based and participatory policy and regulations, participating in collaborative regional processes and mechanisms, and ratifying the AU Convention on Cybersecurity and Personal Data Protection.

The Framework presents a set of detailed recommendations and arising actions to guide member states through the formulations of policy in their domestic context, as well as recommendations to strengthen cooperation among countries and promote intra-Africa flows of data. The main high-level overarching recommendations are included here. It is recommended that Member States:

- cooperatively enable data to flow on the continent while safeguarding human rights, data protection, upholding security and ensuring equitable sharing of the benefits;
- cooperate to create the necessary data capabilities to take advantage of data-reliant technologies and services, including the

capacity to govern data so that it benefits African countries and citizens  and enables development;

- promote transversal data policy and agile regulation to navigate the emergence of new dynamic data-driven business models that can foster intra-Africa digital trade and data-enabled entrepreneurship;

- create co-jurisdictional frameworks for the coordination of autonomous competition, sector, and data regulators to regulate the data society and economy effectively, formulate, implement, and review data policy in a dynamic, forward-looking and experimental way;

- develop national legislations on personal data protection and adequate regulations, particularly around data governance and digital platforms, to ensure that trust is preserved in the digital environment;

- establish or maintain independent, well-resourced and effective Data Protection Authorities, strengthen cooperation with DPAs from members of the African Union and develop mechanisms at the continental level to develop and share regulatory practices and support institutional development to ensure a high level of protection of personal data;

- promote interoperability, data sharing, and responsiveness to data demand through the setting of open data standards in data creation conform to the general principles of anonymity, privacy, security and any sector-specific data considerations to facilitate non-personal data, and certain categories of personal data are accessible to African researchers, innovators and entrepreneurs;

- promote data portability so that data subjects are not locked into a single provider and, in so doing, promote competition and consumer choice and enable gig workers to move between platforms;

- improve unevenly developed infrastructure across the continent, leveraging existing REC regional efforts to support efficient broadband network coverage, reliable energy supply, and foundational digital (data) infrastructure and systems (FDI) (digital identity (Digital ID)), interoperable trustworthy payments, cloud and data infrastructure, and open data sharing systems, for cross border digital trade, e-commerce;
- establish an integrated national data system to enable data-driven public and private value creation, operating on the basis of harmonised governance frameworks that facilitate the flow of data necessary for a vibrant data economy, but with sufficient safeguards to be trusted, safe and secure;
- govern the integrated national data system according to the principles of access, availability, openness (where anonymity can be preserved), interoperability, safety, security, quality, and integrity;
- integrate sector-specific and specialists data codes or guidelines into national and continental data governance regimes; those who have not yet ratified the AU Convention on Cyber Security and Personal Data Protection, do so as soon as possible to serve as the foundational step for the harmonisation of data processing; and in the forthcoming negotiations on Trade in Services and E-commerce protocols, as well as the Competition and Intellectual Property protocols, in the African Continental Free Trade Area provide guidelines to promote access to data to support local innovation, entrepreneurialism and pro-competitive purposes;
- prioritise politically neutral partnerships that take into account individual sovereignty and national ownership to avoid foreign interferences which may negatively affect the national security, economic interests and digital developments of AU Member States;

- promote research, development and innovation in various data-based areas, including Big Data Analytics, Artificial Intelligence, Quantum Computing, and Blockchain.

It is further recommended that The African Union Commission, RECs and Regional Institutions:

- facilitate collaboration between the various entities dealing with data across the continent through the establishment of a consultation framework within the digital ecosystem community to safeguard the interest of each actor;
- promote and facilitate data flows within and among AU Member States by developing a Cross Border Data Flows Mechanism that takes into account the different levels of digital readiness, data maturity as well as legal and regulatory environments of countries;
- facilitate data circulation across sectors and cross borders by developing a Common Data Categorisation and Sharing Framework that takes into account the broad types of data and the associated levels of privacy and security;
- work in close collaboration with national authorities in charge of personal data protection of AU members, with the support of the African Network of Authorities (RAPDP), to establish a co-ordination mechanism and body that oversees the transfer of personal data within the continent and ensures compliance with existing laws and rules governing data and information security at national level;
- establish or empower a mechanism within the African Union for centralising and empowering regional engagements on data standards;
- establish mechanisms and institutions , or empower existing ones, within the African Union to build capacity and render technical assistance to AU Member States for the domestication

of this data policy framework; and support the development of regional and continental data infrastructure to host advanced data-driven technologies (such as Big Data, Machine learning and Artificial Intelligence) and the necessary enabling environment and data-sharing mechanism to ensure the circulation across the continent;

- work towards building a secure and resilient cyberspace on the continent that offers new economic opportunities through the development of an AU Cyber Security Strategy and establishment of Operational Cybersecurity Centres to mitigate risks and threats related to cyberattacks, data breaches, and misuse use of sensitive information;

- enable data sharing and enhanced interoperability among AU Member States and other AU mechanisms, including the African Union Mechanism for Police Cooperation (AFRIPOL);

- establish an Annual Data Innovation Forum for Africa to raise awareness amongst policy makers about the power of data as the engine of a digital economy and society so as to facilitate exchanges among countries and enable knowledge sharing on data value-creation and innovation and the implications of data usage on peoples' privacy and security;

- strengthen links with other regions and coordinate Africa's common positions on data related international negotiations to ensure equal opportunities in the global digital economy;

- develop an implementation plan that takes into consideration the digital sovereignty of states as well as the different levels of development, the vulnerability of populations and digitisation within AU Member States, namely aspects related to ICT infrastructure gap and lack of cybersecurity policies and legislations.

## 34.2 Guiding Principles of the Framework

(Chapter 5.1 original text) The Data Policy Framework needs to align with the AU values and International law to achieve greater unity and solidarity between African countries and their people, ensuring balanced and inclusive economic development, including promoting and protecting peoples' rights through the African Charter on Human and Peoples' Rights and other relevant instruments.

In the spirit of fostering regional prosperity, economic growth and development, social progress and coordinating continental efforts, the following high-level principles guide the framework.

*Cooperation:* African Union Member States shall cooperate in exchanging data, acknowledging data as a central input of the global economy and the importance of the interoperability of data systems to a flourishing African digital single market.

*Integration:* the Framework shall promote intra-Africa data flows, remove legal barriers to data flow, subject only to necessary security, human rights and data protection.

*Fairness and inclusiveness:* in the implementation of the Framework, Member States shall ensure it is inclusive and equitable, offering opportunities and benefits to all Africans, and in so doing, seek to redress national and global inequalities by being responsive to the voices of those marginalised by technological developments.

*Trust, safety and accountability:* Member States shall promote trustworthy data environments that are safe and secure, accountable to data subjects, and ethical and secure by design.

*Sovereignty:* Member States, AUC, RECs, African Institutions and International Organisations shall cooperate to create capacity to enable African countries to self-manage their data, take advantage of data flows and govern data appropriately.

*Comprehensive and forward-looking:* the framework shall enable the creation of an environment that encourages investment and innovation

through the development of infrastructure, human capacity and the harmonisation of regulations and legislation.

*Integrity and justice:* Member States shall ensure data collection, processing and usage are just and lawful, and data should not be used to discriminate unfairly or infringe peoples' rights.

## 34.3 Recommended Actions

*These are the recommended actions of chapter 5 on "Data Policy Framework"[457].*

### 34.3.1 Foundational Data Infrastructure

*Broadband and Data Access and Use (5.3.1.1). Actions:*
Member States will need to evolve policies that:

- proscribe prohibitive 'right of way' broadband cable fees and support infrastructure sharing;
- prevent anti-competitive practices arising from dominance in infrastructure markets;
- invest in public Wi-Fi and complementary technologies;
- adopt innovative spectrum utilisation techniques such as dynamic spectrum allocation and access, and the leverage of digital dividend (spectrum bands largely expedited by the analogue to digital broadcasting migration) to expand broadband access for under-served rural areas;
- promote the transition and adoption of IPv6[7], as IPv4 resources become more depleted globally;
- invest in national backbone and cross-border connectivity infrastructure such as Internet Exchange Points (IXPs) at both national and regional levels to leverage available international bandwidth, lower in-

---

[457] *Pages 18-64.* Full text: https://au.int/en/documents/20220728/au-data-policy-framework.

ternet access cost and enhance data access speeds within the region; and

- leverage innovative models for data infrastructure funding.

*Data Infrastructure (5.3.1.2). Actions*

- As opposed to focusing on the significant upfront investment to replace depreciating legacy ICT equipment, Member States should leverage economies of scale and scope to adopt infrastructure that supports facilitating benefits offered by cloud services and other new technologies that support data value creation.
- Tax, trade (including investment and innovation) and competition policies must be coherent, complementary, and adapted to the data-driven digital economy, particularly to inform infrastructure development strategies.
- Member states must ensure local firms participate in value chains of foreign software as a service (SaaS), infrastructure as a service (IaaS) and platforms as a service (PaaS) providers for state procurement and create incentives to have local SMMEs in data value chains across industries. This can be done by ensuring tax, trade (including investment and innovation), and competition policies are coherent, complementary, and adapted to the data driven digital economy.
- Adopt more sustainable electricity generation models domestically and across the region, to ensure foundational digital infrastructure supports sustainable domestic and cross-border data activities that have fewer extractive impacts on the natural environment.
- Creating data portability rights - including for non-personal data, to make it easier for customers of cloud services to switch between providers.
- Develop contractual standards for public organisations (that can be used by SMEs too) that protect their rights to access, retrieve, delete, etc., the data (including non-personal, again) that is processed by cloud providers.

- Develop Fair, Reasonable and Non-Discriminatory (FRAND) licensing obligations for platforms and cloud providers who have access to datasets that become a vital resource to enter a market.

*Creating Legitimate, Trustworthy data Systems (5.3.2). Actions:*

- Safeguard basic human rights in the digital environment through the rule of law.
- Ensure institutional arrangements and regulations are established only through inclusive, consultative and transparent processes.
- Ensure institutions responsible for overseeing the use of data, as well as public and private data producers, are accountable for the use of public and personal data to those whose data is used.
- Strengthen cooperation with other DPAs to ensure sufficient safeguard and reciprocal protection of personal data as well as individual and collective digital rights across the continent.
- Strengthen Mutual Legal Assistance Agreements and activities across states for the investigation and prosecution of cybercrimes.
- Ensure institutions responsible for overseeing the use of personal data are empowered to have powers of entry and inspection for purposes of enforcement of privacy and data protection laws and regulations.
- • Further ensure institutions responsible for overseeing the use of personal data have the following corrective powers in relation to correcting infringement of aspects of misuse and abuse of personal data:
- Issue warnings to a data controller or data processor that intended processing operations are likely to infringe provisions of the relevant data protection laws and regulations.
- Issue reprimands to a data controller or a data processor where processing operations infringe provisions of the relevant data protection laws and regulations.
- Order a data controller to communicate a personal data breach to affected data subjects.

- Impose a temporary or definitive limitation, including a ban on personal data processing.
- Order the suspension of data flows to a recipient in a third country or to an international organisation that does not provide adequate protection similar to that of the data exporting country.
- Institutions responsible for overseeing the use of personal data should be empowered to either assist or seek a court's indulgence to assist a person who has suffered material damage as a result of an infringement of their personal data to receive compensation from a data controller or data processor for the damage suffered.

*Data Ethics (5.3.2.5). Actions:*

- The data industry and research communities using data need to formulate and implement codes of practice, including the principles of responsibility and ethics by design through processes that include those whose data is affected.
- Member States must require rights-compliant ethical frameworks in public procurement processes.
- Members should include the assessment of data codes of ethics in the mandates of existing human rights bodies such as Human Rights Commissions.

*Institutional Arrangements for Regulation (5.3.3) Actions:*

- Members with data regulators should assess whether the existing enforcement powers are sufficient.
- Members creating data regulators should consider a range of enforcement powers and in addressing resource constraints, how data regulators could potentially rely on other agencies for enforcement.

*Rebalancing the Legal Ecosystem (5.3.4)*

- Contracts that purport to give up digital rights, personal data protection and that inhibit competition should, as a general rule, be unenforceable. This can be articulated in data protection and competition regulation, which can also consider on a case by case basis whether

the pro-competitive effects of such contracts outweigh the anticompetitive effects.

- National law reform commissions or similar expert legal institutions should investigate and consider how to harmonise different branches of laws, regulatory regimes and supervisory authorities that deal with data.
- Member States should support the update or adoption of competition law frameworks and regulations that consider the challenges of analysing competition issues, designing remedies and enforcing their powers to safeguard competition in data-driven markets, as well as building the capacity of competition regulators to implement these rules.
- Intellectual property laws should be amended to provide:
  that if copyright applies to databases and compilations of data at all, it shall apply only to the work of human authors that exhibit originality/creativity and that the copyright extends only to the original selection and arrangement of data in a database or compilation and not to the data itself;
- that any copyright or other intellectual property right, including trade secrets that enables control of data, does not apply to personal data;
- that any copyright or other intellectual property right, including trade secrets that enables control of data, is limited by the provisions of competition regulation and alternative rights that offer protection to local innovations not envisaged in current frameworks;
- adaptations to existing IPR regimes to leverage next frontier technologies, such as enabling AI to use data;

*Consultative and Evidence-Based Regulations (5.3.4.2) Actions:*
- Clearly distinguish between the roles of the state as policy maker and the regulator, which should be sufficiently independent of the state and industry, so as to implement policy in the public interest.
- Create or maintain competition authorities to deal with dominance in the market and concentration through mergers and acquisitions.

- Implement clear procedures for co-jurisdiction between sector and competition authorities to ensure the coordinated regulation of digital infrastructure and services sector and to avoid 'forum-shopping'.
- Data regulators should collaborate at the regional and continental levels to harmonise their frameworks, particularly in support of the AfCFTA.
- Those subject to decisions of regulatory authorities should have clear mechanisms of appeal and redress heard by a different body from the regulator, making the decisions in line with the rules of natural justice and fair administrative action.

*Creating Public Value (5.3.5) Actions:*

- Sector regulators and public data stewards must operate within specific guidelines on how data quality assessments should be implemented, depending on common use cases, algorithms, and type of data used. These guidelines can be informed by global best practices (including data and AI governance) but should be adapted to the context of African data use cases. Due to the exchange, combinations, strategic storage, and repurposing are required to create data value. An effective data quality strategy across the public sector should be informed by technical/practical/operational realities and should outline the roles, responsibilities, and mandates of various government agencies in collecting and maintaining high-quality data in a manner that safeguards citizens.
- Member States need to participate in efforts to establish and adopt a normative framework for harmonised data standards and systems aimed at establishing national, regional, and international interoperability. These may include targeted human, technical, and institutional training interventions, sub-regional infrastructure projects, and REC regulatory sandboxes.
- A continental approach facilitates economies of scale to incentivise private investments in foundational digital infrastructure, including cloud-based technologies. Regional harmonisation of regulations for

data governance could further reduce compliance costs and reduce uncertainty and operational risk for major ICT related infrastructure investments.

- Public institutions that curate data should be adequately resourced in order to contribute to multilateral fora regarding data and to be stewards of inclusive access and responsible use of data guided by appropriate industry technical and regulatory norms, standards, and best practices that underpin both the informational and economic characteristics of data in priority industries.

*Data Governance (5.4) Actions:*

- Data protection authorities (DPA) need full empowerment, including the remit on data sovereignty.
- DPAs are encouraged to adopt international and regional cooperation practices taking note of different stages of implementation and enforcement across the Member States.
- Risk assessment and multi-stakeholder engagement should be used to design data localisation solutions in policy by drafters, which includes civil society participation.
- Data infrastructure policy should be aligned with data control imperatives by policy drafters but must consider cybersecurity, personal data protection, environmental risks and cost.
- Public administration and investment policy should align with data control capacities as a priority.
- Capacity-building in relation to data protection, cybersecurity and institutional data governance in relevant agencies should be assured through policy and asset allocation.

*Data Processing and Protection (5.4.2)*

- Data processing frameworks should be established in partnership with all relevant multi-stakeholder partners but driven ideally by the DPA. These should align with the following principles: consent and legitimacy; limitations on collection; purpose specification; use limi-

tation; data quality; security safeguards; openness (including incident reporting, an important correlation to cybersecurity and cybercrime imperatives); accountability; and data specificity.

- DPAs should be established as a matter of urgency alongside national legislation on personal data protection.

*Data Access and Interoperability (5.4.3). Actions:*

- Member States should establish an open data policy which sets open standards for the production and processing of data so that when decisions are made to open the data, the high costs of ensuring it is usable and manipulatable are avoided.
- Sectoral laws and codes of conduct from DPAs should be reviewed to ensure lawful data access in conjunction with the data policy.
- DPAs should have dual access to information and privacy function.
- Multi-sectoral open data initiatives should be implemented in priority data sectors like health, research and planning.

*Data Security (5.4.4). Actions:*

- Member States, who are yet to develop cybersecurity measures, should immediately develop cybersecurity plans and streamline them within government governance structures to promote robustness and reduce vulnerabilities.
- Cybersecurity institutions like CSIRTs should be incorporated into data policy development.
- Data processing roles as a form of security protection should be specified in policy by policymakers.
- Capacity-building in relation to data protection, cybersecurity and institutional data governance in relevant agencies should be assured through policy and asset allocation and could be supported by DPAs.

*Cross-Border data Flows (5.4.5). Actions:*

- DPAs should ascertain minimum standards for data transfer.
- Capacity-building in relation to data protection, cybersecurity and institutional data governance in relevant agencies should be assured

through policy and asset allocation and driven ideally by DPAs in conjunction with educational facilities and government skills programmes and units.

*Data Demand (5.4.6). Actions:*

- Data communities should be incorporated into data policy-making processes by policymakers.
- Data communities should be drawn into the establishment of open government data initiatives by departmental implementers.
- Universities should be included as relevant policy stakeholders to help establish the "knowledge-base" from which the local data economy can draw sufficient scientific and technological knowledge.

*International and Regional Governance (5.5).*
*Continental Data Standards ((5.5.1). Actions:*

- Establish or empower a mechanism within the African Union for centralising and empowering regional engagements on data standards.

*Continental Instruments (5.5.3). Actions:*

- Member States should ratify the AU Convention on Cybersecurity and Personal Data Protection and develop additional protocols, as required, to reflect changes since the original drafting.
- Establish, or empower a mechanism within the African Union for centralising regional engagements on data standards.
- Once adopted, alignments with the AfCFTA process should immediately be explored.
- Include data in negotiations on the AfCFTA chapters on competition and intellectual property.
- Agree on common and consistent criteria for assessing adequacy in the levels of protection of personal data across the continent to facilitate and enable trans-border transfer of data and standardise protection.

<center>**35**</center>

# EUROPEAN UNION:
# THE NEW DIGITAL SERVICES ACT (DSA)

*European Commission, Belgium*[458]

"*With the Digital Services Act, we now have clear legislation. Online platforms are at the core of some of the key aspects of our daily lives, democracies, and economies. It's only logic that we ensure that these platforms live up to their responsibilities in terms of reducing the amount of illegal content online and mitigating other online harms, as well as protecting the fundamental rights and safety of users.*" *(Margrethe Vestager, Executive Vice-President for a Europe Fit for the Digital Age - 16/11/2022)*

"*The Commission is rapidly preparing to take on its new role as supervisor for some of the world's largest and most influential tech com-*

---

[458] The following text is the official Press Release of the European Commission of 16 Nov. 2022. https://ec.europa.eu/commission/presscorner/detail/en/IP _22_6906. For more information: EU Official Journal text; Digital Services Act Q&A; Digital Services Act Factpage; The Digital Services Act package.

*panies that millions of European citizens and businesses use every day. These new rules launch us into a new age, where big online platforms will no longer behave like they are too big to care. The new rulebook will be supported by strong oversight and enforcement, including fines of up to 6% of global turnover in the event of non-compliance and even a ban on operating in the EU single market in case of repeated serious breaches. Everyone is welcome to do business in the EU, but they will have to follow these new rules." (Thierry Breton, Commissioner for Internal Market - 16/11/2022)*

On 16 November 2022, a landmark new set of EU rules for a safer and more accountable online environment entered into force with the *Digital Services Act* (DSA)[459]. The DSA applies to all digital services that connect consumers to goods, services, or content. It creates comprehensive new obligations for online platforms to reduce harms and counter risks online, introduces strong protections for users' rights online, and places digital platforms under a unique new transparency and accountability framework. Designed as a single, uniform set of rules for the EU, these rules will give users new protections and businesses legal certainty across the whole single market. The DSA is a first-of-a-kind regulatory toolbox globally and sets an international benchmark for a regulatory approach to online intermediaries.

## 35.1 New Responsibilities for Digital Services

The DSA introduces a comprehensive new set of rules for online intermediary services on how they have to design their services and procedures. The new rules include new responsibilities to limit the spread of

---

[459] The Digital Services Act: ensuring a safe and accountable online environment, European Commission, https://commission.europa.eu/strategy-and-policy/priorities-2019-2024/europe-fit-digital-age/digital-services-act-ensuring-safe-and-accountable-online-environment_en

illegal content and illegal products online, increase the protection of minors, give users more choice and better information. The obligations of different online players match their role, size and impact in the online ecosystem; an overview is available[460].

All online intermediaries will have to comply with wide-ranging new transparency obligations to increase accountability and oversight, for example with new flagging mechanism for illegal content. But a special regime is introduced for platforms with more than 45 million users: for such very large online platforms or search engines, further obligations include wide-ranging annual assessments of the risks for online harms on their services - for example with regard to exposure to illegal goods or content or the dissemination of disinformation. Under the DSA, suitable risk mitigation measures will have to be put in place, and subject to independent auditing of their services and mitigation measures.

Smaller platforms and start-ups will benefit from a reduced set of obligations, special exemptions from certain rules, and most crucial increased legal clarity and certainty for operating across the whole EU's single market.

## 35.2 Enhanced Safeguards for Fundamental Rights Online

The new rules protect users' fundamental rights in the EU also in the online environment. New protections for the freedom of expression will limit arbitrary content moderation decisions by platforms, and offer new ways for users to take informed action against the platform when their content is moderated: for example, users of online platforms will now have multiple means of challenging content moderation decisions, including when these decisions are based on platforms' terms and condi-

---

[460] The Digital Services Act, op. cit.

tions. Users can complain directly to the platform, choose an out-of-court dispute settlement body or seek redress before Courts.

New rules also require platforms' terms to be presented in a clear and concise manner and to respect users' fundamental rights.

Very large online platforms and search engines will in addition have to undertake a comprehensive assessment of risks to fundamental rights, including the freedom of expression, the protection of personal data, and freedom and pluralism of the media online as well as the rights of the child.

## 35.3 New Supervisory Powers for the Commission

The DSA creates an unprecedented level of public oversight of online platforms across the Union, both at national and EU level. The Commission has powers to directly supervise VLOPs and VLOSEs, companies which individually reach more than 10% of the EU population, approximately 45 million people. Additionally, each Member State will have to designate a Digital Services Coordinator, who will supervise other entities in scope of the DSA as well as VLOPs and VLOSEs for non-systemic issues. The national coordinators and the European Commission will cooperate through a European Board of Digital Services. This EU-wide cooperation mechanism will be established between national regulators and the Commission.

The Commission is setting up a *European Centre for Algorithmic Transparency* (ECAT)[461] to support its supervisory role with in-house and external multidisciplinary knowledge. The Centre will provide support with assessments as to whether the functioning of algorithmic systems are in line with the risk management obligations that the DSA establishes for VLOPs and VLOSEs to ensure a safe, predictable and trusted online environment.

---

[461] European Centre for Algorithmic Transparency, European Commission, https://algorithmic-transparency.ec.europa.eu/index_en

## 35.4 Next Steps

Following the entry into force of the DSA today, online platforms will have 3 months to report the number of active end users (17 February 2023) on their websites. The Commission is also inviting all online platforms to notify to it the published numbers. Based on these user numbers, the Commission will make an assessment as to whether a platform should be designated a very large online platform or search engine. Following such a designation decision by the Commission, the entity in question will have 4 months to comply with the obligations under the DSA, including carrying out and providing to the Commission the first annual risk assessment exercise. EU Member States will need to empower their Digital Services Coordinators by 17 February 2024, the general date of entry in application of the DSA, when the DSA is fully applicable for all entities in its scope.



## 35.5 Background

On 15 December 2020, the Commission made the *proposal*[462] on the DSA together with the proposal on the Digital Markets Act (DMA) as a

---

[462] Europe fit for the Digital Age: Commission proposes new rules for digital platforms, 2020. European Commission, https://ec.europa.eu/commission/presscorner/detail/en/ip_20_2347

comprehensive framework to ensure a safer, more fair digital space for all. The DMA *entered into force* on 1 November 2022.[463]

Digital services include a large category of online services, from simple websites to internet infrastructure services and online platforms. The rules specified in the DSA primarily concern online intermediaries and platforms. For example, online marketplaces, social networks, content-sharing platforms, app stores, and online travel and accommodation platforms.

---

[463] Digital Markets Act: rules for digital gatekeepers to ensure open markets enter into force. 2022. European Commission, https://ec.europa.eu/commission/presscorner/detail/en/IP_22_6423

# HOW TO BUILD DIGITAL TRUST?
# SWISS DIGITAL TRUST WHITEPAPER

*Nicolas Zahn / Niniane Paeffgen, Switzerland*

*Trust is a key ethical value and virtue in all human interactions. Trust is also a key requirement and a key challenge in the digital world. In September 2019, the first Swiss Global Digital Summit took place in Geneva, Switzerland, to provide a platform to promote in-depth discussions on "Ethics and Fairness in the Age of Digital Transformation" and aiming at building digital trust in the digital world. As a result, the Swiss Digital Initiative (SDI) was launched at the World Economic Forum in Davos in 2020. It is supported by over two dozen high-ranking personalities and institutions from academic, private and public sector. In 2021, the Swiss "Digital Trust Label" was launched (see next article 37 in this book). The following Whitepaper analyses the needs of building trust in the digital world.*

*The Editors/CS*

## 36.1 Executive Summary

Growing erosion of trustworthiness is limiting the adoption of digital services even in cases where they are considered beneficial for societies, e.g. in the case of Covid19 related digital healthcare services. The Swiss Digital Initiative has been working on a Digital Trust Label. The first worldwide label certifies the trustworthiness of a digital service along four dimensions: Security, Data Protection, Reliability, Fair User Interaction.

This Whitepaper presents the main learnings *The Swiss Digital Initiative* has been working on a from the intense work on developing such a Digital Trust Label:

- There are no purely technical solutions to the trust deficit.
- There are no shortcuts: Digital Trust needs to be earned.
- Trust in the digital world is closely linked to offline experience with the company or provider.
- Digital Trust is an ongoing practical challenge and cannot be built overnight.
- Transparency is key in building Digital Trust
- In a broader sense, trust in the digital world is not limited to data protection and security alone but encompasses social and ethical responsibility.
- A general feeling of insecurity in a complex digital world fosters openness of users towards solutions that address the topic of Digital Trust.
- Building Digital Trust concerns all stakeholders and requires a holistic and iterative approach

The Whitepaper outlines the Swiss Digital Initiative's understanding of Digital Trust and why building it In addition to the Digital Trust Label, the Swiss Digital matters for successful digital transformation, Initiative proposes a Digital Trust Framework that particularly in healthcare,

public sector, the media shows the various elements of Digital Trust: sector, banking & insurance, HR and the education sector.

All these elements should be addressed when working looking at the Digital Trust Ecosystem the Swiss towards building Digital Trust in a coordinated fashion. Digital Initiative concludes that coordinated efforts by various stakeholders are needed. Switzerland is well positioned to play an active role globally in working towards Digital Trust.

## 36.2 Foreword by Doris Leuthard[464]

*When we initiated the Swiss Digital Initiative in 2019, we made Digital Trust a key issue for the organisations' work. Growing mistrust in digital services was already visible. Little did we know that a global crisis would act as a further catalyst for developments that had us already worried.*

*We have seen the roll-out of digital services from home office applications and videoconferencing tools to contact tracing apps and other health-related digital services as part of various countries' reaction to Covid 19. Mistrust, not just in digital services but also institutions that enjoyed high levels of trust so far, has shown that this issue needs to be addressed.*

*Ever since its inception, the Swiss Digital Initiative has been dedicated to advancing digital ethics through practice-oriented projects. One such project is now coming to fruition: our Digital Trust Label is about to be launched.*

*We thought that this milestone would be a good opportunity to reflect on our work and look ahead. What have we learned adapting an intuitive concept such as a label to the digital space? What challenges are we and like-minded organisations facing when working towards digital*

---

[464] President of SDI Foundation, Former Federal Councilor of the Swiss Government (2006-2018).

*trust? This Whitepaper not only collects our learnings so far, but it also lays out our understanding of Digital Trust.*

*After all, trust, like many concepts, seems clear and intuitive at first but can mean a variety of things to different people. Therefore, we feel it is important to clearly outline our understanding of Digital Trust, why it matters and how we could help build it.*

*This Whitepaper also positions the Digital Trust Label in the broader context of a Digital Trust Framework. From the beginning we understood that a Label can be one contribution to the complex challenge of Digital Trust but not the only solution.*

*With this Whitepaper we introduce our Digital Trust Framework, a basis for continued and sustained engagement with the challenge of Digital Trust. The Framework shows what different methods we see to work towards advancing Digital Trust. It not only helps us to manage and plan our projects but also provides a blueprint to look at the ever-growing Digital Trust Ecosystem in Switzerland and around the world.*

*Notwithstanding all the challenges, we are at an important moment for digital transformation. The need and demand for trustworthy digital services is here and clear. We are excited to provide a conceptual basis for future projects with this Whitepaper as well as a concrete and implemented project with the Digital Trust Label. My sincere thanks go to all our supporters and partners who believe in a trustworthy digital space.*

## 36.3 Introduction

*Whether in our private life or in business, from grocery shopping to politics: we rely on trust in almost all our interactions. We like to trust people and institutions and we like to be seen as trustworthy by others. However, trust does not simply exist. It needs to be earned. Once betrayed, trust is extremely hard to win back.*

In more and more domains, we rely on algorithms, automated decision-making and complex technologies, whose inner working remains opaque and whose criteria for trust still need to be defined. Digital Trust is a topic of major concern for all stakeholders -companies, institutions, users, and governments alike. We believe that Digital Trust is a cornerstone of successful digital transformation, and the potential of new technologies can only be exploited, if the level of acceptance is high enough. Even great technologies and robust legislation may fall short if mistrust is apparent. Accelerated through the Covid19-pandemic, societies all over the world are experiencing an erosion of trust: In institutions, the media, science, and technology.

That is why the Swiss Digital Initiative (SDI) made the issue of Digital Trust a priority for its work. The initiative aims to create the first worldwide Digital Trust Label. An instrument to foster Digital Trust by enabling users to make informed decisions through greater transparency and inciting companies to take responsibility by offering trustworthy digital services. The learning journey since the beginning of the project in 2019 has been remarkable and one objective of this whitepaper is to give a synthesis of the insights and learnings gained. Based on theory and the practical knowledge from the project work, SDI developed a more holistic "Digital Trust Framework". SDI believes that stakeholders in every country should reflect on how to build Digital Trust. The key aim is to foster constructive digital transformation. SDI works to create an environment where people feel empowered to make decisions based on knowledge and education. It wants to communicate to society that it can use new technologies to improve quality of life. Building digital trust is a long-term investment and key requisite for sustainable digital transformation. It can neither be taken for granted, nor is there a quick tick-the-box recipe to follow.

## 36.4 The Digital Trust Label

### 36.4.1 History and Reasoning

The idea of creating a Digital Trust Label (DTL) was already born back in 2019 at EPFL Lausanne. With the creation of the Swiss Digital Initiative Foundation, the project was officially transferred under the SDI umbrella.

Together with its main partner, EPFL Lausanne, SDI carried out the groundwork for the development of the DTL in 2019 and 2020. The core of the Label, a sound catalogue of verifiable and auditable criteria, has been co-developed by a small Academic Expert Group (from the ETH, EPFL, Universities of Geneva and Zurich) based on a user study on Digital Trust (conducted in November 2019). An independent Swiss-based testing, verification and certification specialist, the Société Générale de Surveillance (SGS), further developed the Label catalogue to make sure that it is auditable and verifiable. The label catalogue as well as other resources can be found on our project website: https://www.digitaltrust-label.swiss

Based on this first draft of the Label Catalogue, SDI conducted the following main processes:

1. It put in place a "Label Expert Committee", consisting of independent experts from academia, fields of data-and consumer protection, human rights, and digital ethics, advising the SDI Board on the Label's content and framework.
2. It conducted a co-development and public consultation process, which provided civil society bodies with the opportunity to comment on and test the catalogue of criteria through face-to-face interviews, an online survey, and workshops.
3. Six of the eight test partners tested the Label Catalogue on concrete use cases, which have contributed a great deal to the Label's improvement. This will help ensure that the Label is fit for purpose when it becomes fully operational.

4.  In addition to the work on the Label Catalogue, SDI has worked on the communication and visualisation of the Label. Together with the market research firm Bruhn&Partner, SDI conducted a user study in Switzerland, USA, Scandinavia, and Eastern Europe. This international investigation tested mechanics and determined success factors and value proposition of such a Label from the customer's perspective.

### 36.4.2 Learnings on Digital Trust from creating a Digital Trust Label

*1. There are no technical solutions to the trust deficit.* Trust is built in a relationship. It is how society manages risk and an uncertain future. A full guarantee is never possible, there always remains a certain risk. The gaining of trust involves a constant battle of reducing technical complexity in the experience of digital practice. Better informed and aware consumers who have a vocabulary to express their concerns will legitimise and drive the digital future. A Digital Trust Label can be one instrument to drive forward the conversation between companies and citizens, but it is not a shortcut or easy solution, which replaces the efforts and seriousness needed for building Digital Trust.

*2. There are no shortcuts: Digital Trust needs to be earned.* Digital Trust can neither be taken for granted, nor is there a simple "tick-the-box" recipe. Traditionally built through interpersonal relationships, it can be hard to grasp what trust in the digital realm might look like. A proxy could be the values, for which a service provider or company stands for and more importantly its actions. Discrepancy between the stated and lived values can have a negative impact on trust. Promoting Digital Trust falls short if trustworthy behaviour cannot be demonstrated. There is no easy solution: Trust takes a constant effort, and it is up to the service-providers to prove that their services are trustworthy. Once lost, it is hard to regain.

*3. Trust in the digital world is closely linked to offline experience with the company or provider.* Experiences in the offline world, such as the reliability of a provider or ethical responsibility as well as employee behaviour, have an influence on trust in the digital world and vice versa. If the company is present in the offline world as well, it needs to take a comprehensive approach that strengthens the trustworthiness of the whole company – in the online and offline world.

*4. Digital Trust is an ongoing practical challenge and cannot be built overnight.* Digital Trust cannot be mandated to a technical solution but is a long-term effort. It is a constant dialogue between users and companies, not a top-down monologue with the consumer as the end point of the conversation. Meeting this requirement in practice demands that the digital ecosystem institutionalises a sensitivity to change and a communications feed-back loop for the entire ecosystem.

*5. In a broader sense, trust in the digital world is not limited to data protection and security alone but encompasses social and ethical re-sponsibility.* Trust in the digital world encompasses social and ethical responsibility – and should also be connected to the general behaviour of the service provider in this direction. Companies are expected to take on responsibility in the digital as well as the offline world. In the online world this could mean protecting customers from fraud and misinfor-mation, protecting vulnerable groups such as children from explicit con-tent and preventing cyber-mobbing.

*6. Transparency is key in building Digital Trust.* For establishing a trustful relationship with the users, a solution such as a Digital Trust Label should provide transparency about the degree of criteria fulfil-ment. This is not to say that organisations should "overload" users with information. It is about presenting relevant information for informed decision-making in a clear fashion. There will be no trust without trans-parency.

*7. A general feeling of insecurity in a complex digital world fosters openness of users towards solutions that address the topic of Digital*

*Trust.* In a global user study conducted by the SDI, 80% of participants evaluate a Digital Trust Label as useful for themselves. Those that do not see an added value in a Digital Trust approach do so due to different reasons: 1) High digital competence and therefore no need for an independent assessment. 2) General mistrust on the internet and resignation about data security and usage ("You can only protect yourself")

*8. Building Digital Trust concerns all stakeholders and requires a holistic and iterative approach.* It needs a collective effort to address the trust issue around new technologies, to guarantee more transparency and accountability. Businesses must live up to their societal responsibility. Policymakers need to set framework conditions to make sure that trust can grow. To foster trust, a more holistic approach and a combination of several measures is needed.

## 36.5 Digital Trust in Context

### 36.5.1 The concept of trust

Trust comes into play whenever we talk about relations between two parties. Trust helps us to bridge the unknown: is this object really doing what it should be doing? Is this person really going to do what she is saying? Is this institution going to interact with me the way I expect it to? If we have trust, we do not need additional answers to questions to engage in a relation needs constant development, feedback iterations and the criteria of trust will need to be challenged and redefined over time. From these learnings, the Swiss Digital Initiative is convinced that a Digital Trust label can be one tool to reduce mistrust. But it is also clear that our ambition cannot stop there. Therefore, the Swiss Digital Initiative proposes a generalised Digital Trust framework to guide our work in the future.
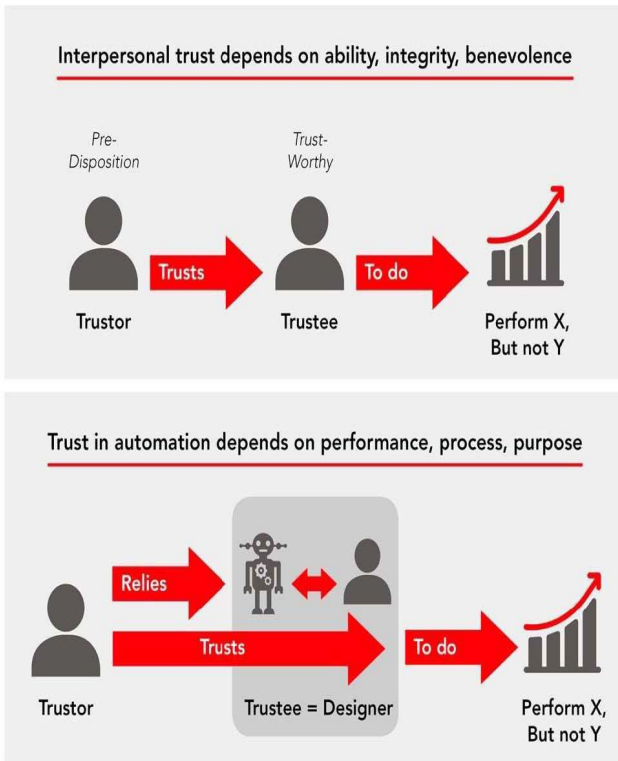
We flick our light switch without thinking whether the lights will turn on, we order at a restaurant without assuming that we will be poi-

soned and we mail our ballots trusting that our votes count. If a certain threshold of trust is reached, we can engage in "trusted" relationships.

Trust can hence be seen as a vital resource for healthy and functioning societies and economies. Trust acts as a relational lubricant, facilitating interactions between various parties. Where there is no trust, certain types of behaviour become prohibitively expensive and are thus not undertaken. This in turn means wasting potential benefits. People end up not engaging with each other due to mistrust or refrain from using certain services. Imagine living in a world without trust at all? What dystopian scenario that would be!

As the above visualisation shows, trust involves sometimes compli-



Based on: scip.ch/en/?labs.20200220

cated relationships between various parties. Also, trust in a person is different from trust in technologies, as the former is usually more normative while the latter focuses on reliability. Nevertheless, persons – in particular designers and coders – are also relevant when it comes to assessing whether a digital service is trustworthy. This point is also reflected in our Digital Trust Label and the broader Digital Trust Framework.

The Swiss Digital Initiative believes that in contrast to other resources, trust does not simply exist, it cannot be found and exploited like a natural resource. Instead, it needs to be built up over time through repeated interaction between parties and institutions.

### 36.5.2 Trust in the digital world

This brings us to the role of trust in the digital world. First, digital transformation has given established actors new possibilities. Newspapers can reach their audience via digital channels in addition to traditional print media. However, as a second development, digitalisation also created new challenges for established actors. Social media and data analytics

#### 36.5.2.1 What about Zero-Trust?

In the field of cybersecurity, there is an approach called Zero-Trust. The idea is, to assume that no digital service in a given system should be trusted, which has consequences on how a system interacting with digital services is set-up. Think of how a company network and its IT security needs to deal with employees bringing their own devices to the office. The company has no control over these devices and hence to increase overall security, under the Zero-Trust approach, assumes that none of those devices and digital services running on them are trustworthy. This assumption is reflected in strict security rules and tight policies. This might be a pragmatic approach, but it is rather costly, supporting our argument that we would not want to live in a world without trustworthy digital services but instead should focus on how to create better condi-

tions for trust, have made it possible to shine a light on how certain actors behave. Various institutions have seen themselves being questioned for the first time as a result. Lastly, digital transformation has created a plethora of new actors and in turn new services that did not exist before, surprising that societies have tried to foster trust instead of perpetuating mistrust. Nevertheless, trust – even when created over a long time – can easily be lost and should not be assumed automatically.

Trust also has a dark side that must be acknowledged. Given that trust is needed to overcome uncertainties, it can be potentially misplaced. Trusting someone always comes with the risk that this trust is abused. Trust creates vulnerability. Hence, while we do not want a world without any trust, we also do not envision a world where we are overly dependent on it. After all, how would you feel if you are constantly confronted with the assertion "trust us" or "trust me"? Rather, the Swiss Digital Initiative aims to reduce mistrust where it is blocking potential benefits for the involved parties.

These three developments create general issues for digital policy such as data handling and security or platform governance. But they also have direct consequences for the role of trust in the digital age and for how trustworthy we see digital services. In the first case, actors might have assumed that trust earned in the analogue world will transfer over to their activities in the digital space. However, the debate around fake news shows that this is not necessarily the case. People might have trusted print newspapers before the internet but now mistrust their digital channels.

In the second case, actors that were regarded as trustworthy in the analogue world can see that trustworthiness evaporate in the digital world because their position is challenged. This may be through more technologically enabled transparency – e.g. the negative effect of WikiLeaks on the perception of states – or because digital transformation has enabled competitors to offer even more compelling services.

This brings us to the third case, where new actors have emerged through digital transformation and offer previously unseen possibilities, such as social media companies. However, because they are new and "unknown" they cannot really build up on pre-existing trustworthiness.

As a direct result of rapid digital transformation, we must talk about Digital Trust. When we say Digital Trust, we focus on a digital service being trustworthy. Our definition of digital service is aligned with the official definition by the European Commission. Digital services include a large category of online services, from simple websites to internet infrastructure services and online platforms. Digital services come in many forms and are omnipresent. We believe that this focus is justified as it is digital services that have come to play an ever-important role in almost every aspect of our lives. And will do so even more in the future. Many digital services involve sensitive information from health to financial information and the consequences of decisions taken as part of using a digital service can be quite severe. One example is the use of AI-systems in the hiring processes. For these reasons the bar for trustworthy digital services should be set high.

There is a cynical argument to be made that all the talk about Digital Trust is simply a way of increasing costs for the providers of digital services and being a hindrance to true innovation by slowing down the pace of technological change. We strongly oppose this view and instead argue that Digital Trust is a precondition for sustainable and successful innovation. As with analogue trust, Digital Trust comes with benefits for all parties involved. People are much more likely to use digital services if they are trustworthy, instead of being convinced to do so in the face of growing mistrust. Likewise, organisations providing digital services that are seen as trustworthy can use this "trust capital" to take risks and try out new things, meaning that Digital Trust is not a hindrance but an enabler of innovation.

However, in recent years, almost no action was taken to actively reduce digital mistrust. Instead, a series of scandals and poorly managed incidents have eroded trust.

### 36.5.3 When to talk about Digital Trust

In a given context the question of how trustworthy a digital service offered by an organisation is can be asked at different times of its life cycle:

*   Beginning of a digital service
*   Development of a digital service
*   Adoption of a digital service
*   Continuous development and abandonment of a digital service

Ideally, the question of trustworthiness is asked and answered at each stage, especially since the "return" of "investing" in trustworthiness early in the life cycle is bigger. The earlier the issue is addressed, the easier the Digital Trust that existed at the start of digital transformation. Governments face push-back when rolling out digital services and technology companies face unhappy employees, critical customers, and pushy regulators. Hence, the question of how to reduce digital mistrust will not go away but only gain in importance in the future.

### 36.5.3.1 The limits of Digital Trust

In our understanding, Digital Trust is a facilitator for the adoption of digital services. However, it is not in essence a normative statement. The question whether a digital service is "good" for individuals or society is not our focus. Our efforts regarding Digital Trust are about making sure that people understand what is behind digital services so they can make an informed decision. A 100% guarantee does not exist, there always remains a certain risk. As in the analogue world, we cannot solve the trust issue through a purely technocratic solution and replace trust with a technical tool. With our thinking about Digital Trust – as operationalised in the Digital Trust Label – we also combine the question of reliability with the question of trust. Rather than limiting ourselves to ensuring that

digital services are reliable – this is only one dimension of the label – we want to broaden the scope and include other factors to tackle and demand less attention and resources later. Already when thinking about an idea for a digital service can organisations and employees raise the issue: will this be trustworthy? Or are we following an idea that raises fears & concerns? The same goes for the development of the service, where interdisciplinary teams for product development that also bring in non-technical expertise, is beneficial. However, we must deal with the fact that many digital services are already rolled-out and adopted. Hence, we also need to address Digital Trust "after the fact", which is exactly the mission of the Digital Trust Label. Nevertheless, the Swiss Digital Initiative believes that going forward Digital Trust should be addressed over the whole life cycle of a digital service.

## 36.6 Fields to Prioritise

Reducing mistrust is paramount for all areas of the society and economy. This is not only the case today but even more so for the future. In a Trend Map developed for the Swiss Digital Initiative by the Think-Tank W.I.R.E, the trends identified further increase the importance of trustworthy digital services. Priority should be given to digital services that are used in fields where:

i) the handled data is very sensitive;

ii) the consequences of using digital services matter greatly;

iii) where there is not much choice whether to use a digital service or not and

iv) where digital services are rolled-out at a high pace and on a large scale.

Following these criteria, we argue to focus, without any particular order, on digital services in the following fields:

- Healthcare: digital services in the health sector almost always involve highly sensitive data and can have potentially lethal conse-

quences. This is a particularly urgent field as under the recent Covid19 pandemic situation, many digital services were rolled out on a large scale in a short timeframe. Think of automated, AI-assisted diagnosis of thorax x-rays.

- Public Sector: changing habits and possibilities also affect political processes and over the last years citizens and institutions of democratic societies have thought about using, or are in fact using digital services in politics, from social media to e-voting. Again, this might involve highly sensitive data – for good reason, votes are supposed to be secret – and a breach of privacy can have dire consequences. One big negative consequence being that democracy itself heavily relies on trust. If this is undermined by using digital services that are not trustworthy, it negatively affects the whole system.

- Media Sector: Closely linked to democratic processes is the question how trustworthy (analogue and digital) media and journalism is seen to be. Digital technologies seem so far to have been a mixed bag. Take the debate about fake news spreading particularly over social media as an example.

- Banking & Insurance: whether it's applying for a loan, payments that should remain private or activity levels for health insurance: financial services and insurance involve sensitive information and potentially severe decisions. Already in the analogue world, the bar for trust is set high and so it should be the case for the digital world.

- HR: from monitoring employees to assessing candidates, digital services are increasingly used in the labour market. For example, if you are applying online via a job portal, then you would probably want to know if a human or algorithm did the initial screening. Hence more transparency is needed for such tools.

- Education: digital technologies are also becoming rapidly widespread in the educational sector and students are exposed earlier to new digital technologies. Given the vulnerability of young students, particular attention to Digital Trust should be given in this sector.

Of course, the legal and regulatory frameworks in various countries address these questions, for instance through data privacy laws. In our quest for Digital Trust we want to go beyond what is legally required. Legal and regulatory frameworks face notorious difficulty in keeping up with technological developments and legal measures – as we explain in the section on the Digital Trust Framework – are but one piece in the puzzle.

## 36.7 Working towards Digital Trust: the Digital Trust Framework

As stated, increasing trustworthiness in digital services is a key objective for the Swiss Digital Initiative. We have gathered practical experience in this field through our first major project, the creation of a Digital Trust Label. To expand on this project in the future, we also propose a general framework that shows the various dimensions of Digital Trust, providing a blueprint for the direction of our future projects but also acting as inspiration for any organisation that also wants to contribute to Digital Trust.

### 36.7.1 Ways of building trustworthiness

To reap the benefits of digital transformation, we want to work towards reducing mistrust. But what makes something trustworthy? Whether in the analogue or digital world, there are several ways. When focussing on digital services, additional elements come into play to reduce mistrust, not just in a single service, but also the digital ecosystem:

- Transparency: as trust is used to bridge what is unknown, one way of tackling the issue is by reducing the amount of uncertainty. This can be done by being as transparent as possible: make clear to someone you want them to trust you, what is happening if they use your service. This includes clear communication, especially in times of cri-

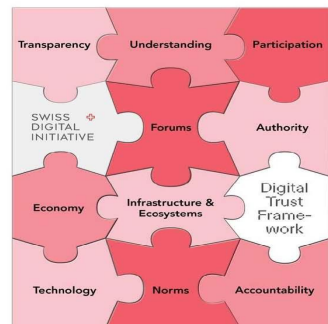sis and a culture that openly admits mistakes and shows what has been learned.

- Understanding: another way of reducing uncertainty and hence reducing mistrust is to be able to understand something. It is much easier for us to trust something we built ourselves than something we bought because we understand the first instance better. Understanding needs the willingness to understand from one party and the willingness to enable such understanding by the other party.
- Participation: Being able to shape developments and being involved is another way of reducing mistrust. Having a stake by participating in the development and use of a digital service makes it trustworthy. Think of the success of the open-source projects that revolve around active communities. In addition to transparency, the ability to have a say in the direction of the digital service increases trust.
- Authority: going back to the idea of trust as "relationship capital", if someone is already trustworthy, this can be used to reduce mistrust vis-à-vis other parties. Think for instance about product endorsements or official recommendations. This is also a leading idea behind our Digital Trust Label project and why the Swiss Digital Initiative strives to be trustworthy, e.g. by following several of these principles like transparency and being participatory.
- Accountability: Users need recourse options, such as legal frameworks that hold businesses and institutions accountable or focal points, ombuds or contact persons, to which users could turn in case of problems. Today, a lot of users feel overwhelmed and insecure when using digital services. Also, organisations making commitments must also be held to those standards in their practice.
- Technology: the technology powering various digital services – from hardware to software – is a crucial element to be trustworthy. Is it safe or riddled with bugs? Is the technology used explainable and understandable? Does the technology used reduce uncertainty? If we

can rely on technology, this is already a first step towards trustworthiness.

- Infrastructure & Ecosystems: In the digital world, digital services often rely on certain infrastructure. From the internet to standardised frameworks for identification and data sharing, such infrastructure plays a key role in the trustworthiness of digital services.
- Norms: While it has been decades that digital technologies have been in use, it is painfully apparent that the "rules of the game" are very much still in the making. That is why norms for digital services matter as another element of Digital Trust. Are certain digital services off limits? What is deemed acceptable behaviour in cyberspace? Clear norms in line with rules in other aspects can reduce mistrust.
- Forums: To debate those norms, forums are needed. Where can we address issues of Digital Trust and who is involved in those discussions? Open and inclusive forums can further reduce mistrust.

### 36.7.2 Comprehensive effort needed

Economy: as the main source of digital services, the economy also has a vital role to play when it comes to reducing digital mistrust. Are we offering a real choice to consumers instead of forcing them to adopt certain services through market power? Are business models based on exploitation of data?

All these elements need to work together to complete the puzzle of trustworthiness. The issue of mistrust cannot be addressed by focussing just on one of the above-mentioned elements and methods. This is particularly evident with technologies that aim to replace the complex social construct of trust through technological means and removing human elements and factors wherever possible, e.g. certain blockchain projects. Distributed ledger technology seeks to replace the trust placed in institutions or humans with trust in technology, e.g. the code of a smart contract. However, research clearly shows that this is not enough to be seen as trustworthy.

Therefore, the Swiss Digital Initiative will start focusing on additional projects along the lines of the framework and welcomes other organisations to join in the effort to build Digital Trust and enable the adoption of trustworthy digital services. Be it through awareness and educational campaigns, research that fosters transparency, efforts to build trustworthy digital infrastructure or adapting business models and leading international debates: we need to work together on all the above-mentioned elements of the Digital Trust puzzle.

Such efforts must also break down silos that still hinder effective collaboration in various aspects of digital policy today: policymaking legal scholars and political scientists need to converse with IT specialists, AI developers and researchers need to listen and talk to civil society representatives. If there ever was a challenge demanding interdisciplinary collaboration, it is working towards a societal beneficial digital transformation with Digital Trust at its core.

## 36.8 The Digital Trust Ecosystem and the Role of Switzerland

Having established the elements of Digital Trust through our framework, we can position the efforts of various organisations and actors within the framework to see where lots of efforts are already undertaken

and which elements of the framework might require more attention going forward. Different stakeholders will be working on different elements given their backgrounds, capabilities, and interests. Although this is a very dynamic space with many actors in Switzerland and internationally, we try to generalise the various contributions.

As a multi-stakeholder organisation itself, the Swiss Digital Initiative is convinced that an effort from all stakeholders is needed towards Digital Trust: Academia and research lay the foundation for the technological infrastructure and functioning. Businesses must live up to their societal responsibility in their own interest: being trustworthy contributes to Travel to any place in the world and ask someone what they associate with Switzerland? Trust stands a good chance of making into the top ten, from its democratic institutions to reliable watches.

So how does Switzerland fit into the Digital Trust landscape? Switzerland enjoys high trust in its institutions and the economy nationally and is seen as a trustworthy actor internationally. An open economy with strong SMEs leads to export-oriented growth but struggles with digital transformation and digital competitiveness. Nevertheless, the presence of leading universities and research facilities also positions Switzerland as a technology hub for cybersecurity, cryptography and other technologies generally linked to "trust", from bias detection in AI data sets to self-sovereign identities. With more organisations joining the Geneva ecosystem, Switzerland also acts as host for various forums for global debates on norms with participants ranging from leading tech companies to NGOs and international organisations increasingly dealing with digital policy and the question of Digital Trust.

Policymakers need to set framework conditions to make sure that trust can grow, and vulnerable groups are protected. Civil society and independent media can function as watchdogs and counter-power. In short: It needs a collective effort to address the trust issue around new technologies, to guarantee more transparency and accountability.

While additional involvement in the Digital Trust field is generally welcome, there is also a danger of "too many cooks" and reinventing the wheel. As the Swiss Digital Initiative has learned in conversations with other initiatives around the world, exchange of information and learnings is seen as beneficial, and efforts should be better coordinated. This Whitepaper and our Digital Trust Framework can hopefully act as a contribution towards better coordination.

There is already a strong basis with many elements in place. What is now needed is connection of the dots and "doubling down" on the potential of transferring the Swiss quality of trust towards the digital age. Players in the Digital Trust ecosystem need to collaborate and coordinate their efforts. The Swiss Digital Initiative will be monitoring the evolving digital ecosystem in Geneva and around the world and is looking forward to cooperating with other stakeholders towards the common vision of Digital Trust. Having recognised the challenge of Digital Trust and its importance for successful digital transformation early on, ideally positions Switzerland to play an active role in this field. Just as was the case with other challenges, confidence-building measures for the digital space are needed and more projects along the line of our Digital Trust Framework can contribute to a world where citizens and users feel: the digital space is trustworthy.

## 36.9 Conclusion

The importance of Digital Trust is only increasing. On one side, mistrust in digital transformation is growing given repeated scandals and revelations about unethical business practices. On the other side, the need for digital services is growing and we have seen widespread rollouts of digital services in various sectors in recent years.

With the Digital Trust Label, the Swiss Digital Initiative provides a concrete and practice-oriented approach for the issue of Digital Trust. Digital services can be certified for their trustworthiness. This can con-

tribute to Digital Trust, but as this paper shows, trust is a complex issue that defies simple solutions.

With the growing need for Digital Trust and based on the learnings from the Digital Trust Label project, SDI proposes a Digital Trust Framework that incorporates various elements that contribute towards Digital Trust.

The framework acts as a blueprint for mapping current activities and planning future projects. Such a mapping helps to identify areas with need for coordination or need to address a gap.

With the spread of digital services, priority in addressing Digital Trust should be given to digital services that are used in fields where:

- the handled data is very sensitive;
- the consequences of using digital services matter greatly;
- where there is not much choice whether to use a digital service or not
- where digital services are rolled-out at a high pace, on a large scale.

This concerns in particular digital services in healthcare, public sector, the media sector, banking & insurance, HR and the education sector.

The growing need for Digital Trust has also created a dynamic Digital Trust Ecosystem that is starting to take shape. With various projects and initiatives along the lines of our Digital Trust Framework – from awareness campaigns to technology start-ups -SDI hopes to contribute to the future evolution of this ecosystem and looks forward to implementing additional projects to tackle Digital Trust in a comprehensive way.

### About SDI

The Swiss Digital Initiative (SDI) is an independent, non-profit foundation headquartered in Geneva and set up in 2020 by the association digitalswitzerland and under the patronage of Federal Councillor Ueli Maurer. The SDI pursues concrete projects with the aim of safeguarding ethical standards and promoting responsible behaviour in the digital world.

The SDI's location in Geneva is no coincidence. The SDI is very much Swiss at heart and embodies many of the Swiss qualities of security, reliability, and trust. At the same time, we recognize that the issue of digital ethics and trust is a global one. With this in mind, we believe our location in international Geneva puts us in a great starting position to combine Swiss values and perspectives with a global debate and international impact.

## Acknowledgments

## References

Atlantic Council (2019), 'Breaking Trust', Technical report, Cyber Statecraft Initiative.

Adjekum, A., Ienca, M. and Vayena, E. (2017), 'What Is Trust? Ethics and Risk Governance in Precision Medicine and Predictive Analytics', Omics: a Journal of Integrative Biology 21(12), 704-710.

Australian Government (2021), 'Digital Service Platforms Strategy', Technical report, Digital Transformation Agency.

ARM (2021), 'Security Manifesto'.
        https://interactive.arm.com/story/security-manifesto-2021/

Cave, S., Craig, C., Dihal, K., Dillon, S., Montgomery, J., Singler, B. and Taylor, L. (2018), Portrayals and perceptions of AI and why they matter, The Royal Society.

CoT (2021), 'Charter of Trust'. https://www.charteroftrust.com/

Deloitte (2021), 'Future of Digital Trust -Digitalization of public sector organissations'.

European Commission (2019), 'Ethics guidelines for trustworthy AI'. https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthyai

European Commission (2021), 'Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on Artificial Intelligence (Artificial Intelligence Act) and amending certain Union legislative acts', COM/2021/206 final.

Adjekum, A., Blasimme, A. and Vayena, E. (2018), 'Elements of Trust in Digital Health Systems: Scoping Review', Journal of Medical Internet Research 20(12).

Haataja, M. and Bryson, J. (2021), 'What costs should we expect from the EU's AI Act?', SocArXiv.

Hoff, K. and Bashir, M. (2014), 'Trust in Automation: Integrating Empirical Evidence on Factors That Influence Trust', Human Factors: The Journal of the Human Factors and Ergonomics Society.

IMD and St. Gallen Symposium (2021), 'Strengthening Trust in Technology When It Matters The Most', Joint White Paper. https://a.storyblok.com/f/72700/x/6f57f16650/strengthening-trust-in-technology_a-joint-imd-and-sgs-white-paper.pdf

Janssen, M., Rana, N., Slade, E. and Dwivedi, Y. (2018), 'Trustworthiness of digital government services: deriving a comprehensive theory through interpretive structural modelling', Public Management Review 20(5).

Gille, F., Jobin, A. and Ienca, M. (2020), 'What we talk about when we talk about trust: Theory of trust for AI in healthcare', Intelligence-Based Medicine 1(2).

Kelton, K., Fleischmann, K. and Wallace, W. (2007), 'Trust in digital information', Journal of the American Society for Information Science and Technology 59(3), 363-374.

Küderli, U. (2019), 'Digital Trust drives business growth', PwC. https://www.pwc.ch/en/insights/digital/digital-trust-drives-business-growth.html

NGI Forward (2020): Report -Digital Trustmarks. https://research.ngi.eu/wp-content/uploads/2020/01/NGI-Forward-Digital-Trustmarks.pdf

OECD (2014), Recommendation of the Council on Digital Government Strategies, OECD Publishing Paris.

OECD (2017), 'Trust and Public Policy: How Better Governance Can Help Rebuild Public Trust', OECD Public Governance Reviews.

OECD (2019), Digital Government Review of Sweden: Towards a Data-driven Public Sector, OECD Publishing Paris.

O'Neill, O. (2002), A Question of trust: The BBC Reith Lectures, Cambridge University Press.

O'Neill, O. (2017), 'Intelligent Trust in a Digital World', New Perspectives Quarterly 34(4), 6-31.

Paris Call (2018), 'Paris Call'. https://pariscall.international/en/

PwC (2021), 'Global Digital Trust Insights Survey 2021: Cybersecurity comes of age'. https://www.pwc.ch/en/publications/ 2020/ch-Digital-TrustInsights-Survey-2021-report.pdf

Eisenhauer, S. (2019), 'Trust in Innovation', ethix White Paper 2.

Swiss Digital Initiative (2019), 'Digital Trust from the customer's perspective: a qualitative study in Switzerland'.

https://a.storyblok.com/f/72700/x/55feac39be/booklet-digital-trust.pdf

Nuremberg Institute for Market Decisions and St. Gallen Symposium (2021), 'Challenges for Human Trust in a connected and technology-driven world', Voices of the Leaders of Tomorrow.

Söllner, M., Hoffmann, A., Hoffmann, H. and Wacker, A. (2012), 'Understanding the Formation of Trust in IT Artifacts', International Conference on Information Systems .

Söllner, M., Hoffmann, A. and Leimeister, J. M. (2016), 'Why different trust relationships matter for information systems users', European Journal of Information Systems 25, 274-287.

Tschopp, M. (2020), 'AI & Trust: stop asking how to increase Trust in AI', scip. https://www.scip.ch/en/?labs.20200220

Tschopp, M. (2020), 'Trust and AI: three wrong questions', scip. https://www.scip.ch/en/?labs.20201112

Tschopp, M., Scharowski, N. and Wintersberger, P. (2021), 'Do Humans Trust AI or Its Developers? Exploring Benefits of Differentiating Trustees Within Trust in AI Frameworks', Extended Abstract.

UVEK (2016), 'Strategie "Digitale Schweiz"'.
https://www.digitaldialog.swiss/de/

van Haasteren, A. (2019), 'Trust in Digital Health', PhD thesis, ETH
Zürich.

Vestager, M. (2021), 'Trust and Technology in a new digital age"Internet
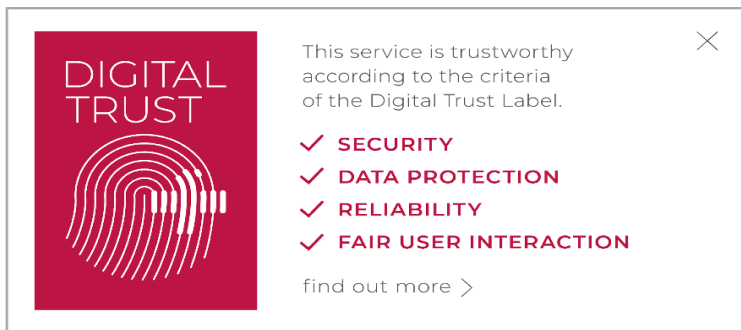Week Denmark'.

Zavolokina, L., Spychiger, F., Tessone, C. and Schwabe, G. (2018),
'Incentivizing Data Quality in Blockchains for Inter-
Organizational Networks Learning from the Digital Car
Dossier"Thirty ninth International Conference on Infor-
mation Systems'.

Zavolokina, L., Zani, N. and Schwabe, G. (2020), 'Designing for Trust in
Blockchain Platforms', IEEE Transactions on Engineering
Management.

<div align="center">

**37**

</div>

<div align="center">

# SWISS DIGITAL INITIATIVE: DIGITAL TRUST LABEL

</div>

## 37.1 The Digital Trust Label

The Swiss Digital Initiative and its partners have developed the first Digital Trust Label[465] that denotes the trustworthiness of a digital service in clear, visual and plain, non-technical language for everyone to understand. With the Label, consumers can be assured of the trustworthiness of the digital service they consume.



---

[465] The authorization to reproduce the content of this chapter has been requested by the volume editors. Text from the website: The Digital Trust Label, https://www.swiss-digital-initiative.org/digital-trust-label. More also on https://digitaltrust-label.swiss/

The Digital Trust Label acts as a combination of a bio Label and nutrition fact table for the digital world. The Label shows that mandatory criteria are fulfilled by a digital service, *while at the same time giving users more information and transparency about four dimensions of the digital service:*

### 37.1.1 Need for a Digital Trust Label:

Whether in our personal or professional life, we look for trust to feel at ease and engage with others easily. The same applies to the digital world. Growing complexity and opaqueness, coupled with various scandals around data breaches and discrimination by algorithmic systems, are fueling mistrust in digital services. Now is the time to establish digital trust and harness the beneficial potential of digital technologies for our societies.

*The Label Development:* Digital Trust cannot be defined by one actor alone, but can only be the result of the close collaboration of all relevant actors: academia, civil society, consumer protection, the private and public sector. This is why the SDI involved all relevant stakeholders in the development and made the criteria and development process as transparent as possible.

The Label is understood as an ongoing and collaborative effort to strengthen trust in a digital service through increased understanding and transparency. While other organisations around the globe are pursuing similar initiatives, SDI and its Digital Trust Label are the most developed initiative from 50+ similar initiatives worldwide as a study revealed[466].

The Label development is an iterative process, hence, it will be continuously developed. In 2022, the first version of the label officially launched.

---

[466] Labels and Certifications for the Digital World, Swiss Digital Initiative, https://a.storyblok.com/f/72700/x/73839efcca/attch-1_sdi_initiatives_final.pdf

### 37.1.2 Governance

The Digital Trust Label is governed by two principal governance bodies: the Label Expert Committee (LEC) & the Label Certification Committee (LCC)

*The Label Expert Committee (LEC)* was and continues to be dedicated to the development of the label by formulating recommendations to the SDI Board. Between 2019-2022, the LEC consisted of ten experts[467] from academia, economy, data and consumer protection, legal, human rights and digital ethics to represent a diverse and inclusive expert view. During the development of the first version of the Digital Trust Label, the LEC's work was coordinated by the EPFL Center for Digital Trust (C4DT).[468] To further improve the Digital Trust Label, additional experts were recruited through an Open Call to the Label Expert Committee.

*The Label Certification Committee (LCC)* is composed of independent experts.[469] They review the Audit reports and give their recommendation on awarding the Digital Trust Label, thereby adding an additional layer of security.

*The Digital Trust Label Test Partners:* Seven test partners from the public and private sector were involved in the project to pilot the criteria: AXA, Canton Vaud, SwissRe, Credit Suisse, Booking.com, IBM Switzerland and Swisscom.

## 37.2 Criteria Catalogue for the Digital Trust Label

The following catalogue of criteria[470], Version 1, valid since Nov 2021, was developed by the Label Expert Committee[471], involved in the

---

[467] See list https://www.swiss-digital-initiative.org/digital-trust-label.

[468] https://c4dt.epfl.ch/

[469] current members https://www.swiss-digital-initiative.org/digital-trust-label.

[470] © 2022 Swiss Digital Initiative. Digital Trust Label.

creation of the criteria catalogue between 2019 -2021, coordinated by EPFL Centre for Digital Trust (C4DT, Imad Aad and Martin Rajman).

- Following the example of the first *category Security*:

| Category | Criteria | No | Specification |
|---|---|---|---|
| Security | Secure communication, data transmission and storage | 1 | The service shall apply best practice cryptography to data in transit, ensuring that the cryptography is reviewed and evaluated, delivers the required functions for all transmitted data and is appropriate to the properties of the technology, risk, and usage. All data in transit over open communication lines such as the internet must be encrypted. |
| | | 2 | The service shall apply best practice cryptography to data at rest, ensuring that the cryptography is reviewed and evaluated, delivers the required functions for all sensitive and applicable data at rest and is appropriate to the properties of the technology, risk, and usage. |
| | | 3 | Privacy-enhancing technologies such as Anonymization and Pseudonymization shall be used ac- |

---

[471] LEC President: Stéphanie Borg Psaila, Digital Policy Director, DiploFoundation. Members: -Prof. Yaniv Benhamou, Faculty of Law, University of Geneva, Attorney-at-Law -Prof. Dr. Abraham Bernstein, Department of Informatics, Director Digital Society Initiative, University of Zurich -Nikki Böhler, Managing Director, OpenData.ch -Francesca Bosco, Senior Advisor, Cyber Capacity and Foresight, CyberPeace Institute -Christophe Hauert, Lecturer University of Lausanne, Co-Founder Cybersafe Label -Prof. Dr. Jean-Pierre Hubaux, Full Professor Laboratory for Data Security, EPFL -Carla Hustedt, Senior Project Manager, Bertelsmann Foundation (until the end of 2020) -Dr. Patrick Schaller, Senior Scientist, System Security Group, ETH -Florian Schütz, Federal Cyber Security Delegate -Jean-Christophe Schwaab, Fédération Romande des Consommateurs -Martin Steiger, Attorney, and Entrepreneur for Law in the Digital Space

| | | | |
|---|---|---|---|
| | | | cording to best practices in order to adequately protect the user's data. |
| | Secure user authentica-tion | 4 | All passwords used for the service shall be subject to a state-of-the-art password policy, which in-cludes requirements applicable to the service and ensures that no hard-coded passwords are used, best practice authentication is in place and ensures that brute-force attacks on authentication mecha-nisms are not feasible. |
| | Secure service set up, mainte-nance and update | 5 | Guidance for secure installation, configuration, and updates shall be in place and updated for each release if necessary. Guidance shall be available in a manner that is easy to access and understand. Any major changes shall lead to a communication to the users in an easy-to-understand format. |
| | | 6 | All software components shall be updatable in a secure manner, and verification of security updates shall be in place. |
| | | 7 | Updates shall be timely. Updates addressing critical security vulner-abilities must be available as soon as possible. |
| | | 8 | Hard-coded critical security pa-rameters in service software source code shall not be used. |
| | | 9 | Any critical security parameters used for integrity and authenticity checks of software updates and for protection of communication with associated service software shall be unique per service and shall implement security measures to protect the integrity and confidentiality of critical security parameters. |
| | | 10 | The service provider shall follow secure management processes for critical security parameters that |

| | | | relate to the service. |
|---|---|---|---|

*Further categories* have been developed:

- Category *Data Protection*
- Category *Reliability*
- Category *Fair User Interaction*

The full and updated list can be seen on the website and requested at the Swiss Digital Initiative.[472]

---

[472] https://digitaltrust-label.swiss.

# 38

# GENEVA DIGITAL ATLAS 2.0

*Geneva Internet Platform Digwatch*

*The Geneva Internet Platform Digital Watch, in short GIP Digital Watch or just digwatch, is an initiative of the Geneva Internet Platform, supported by the Swiss Confederation and the Republic and Canton of Geneva. The GIP is operated by DiploFoundation. It further positions International Geneva as hub for digital issues and its governance. Below the short introduction to the new Geneva Digital Atlas 2.0 and digwatch from their website.* [473]

*The Editor CS*

The Digital Watch observatory is a comprehensive digital policy observatory, which provides a neutral one-stop shop for the latest developments, overviews, events, actors, instruments, and other resources.

---

[473] Geneva Internet Platform Digital Watch, https://dig.watch/. Permission to reproduce given by Nicolas Zahn, SDI Foundation.

The observatory is part of the Geneva Internet Platform[474], an initiative of the Swiss authorities, operated by DiploFoundation [475].

### The Geneva Digital Atlas 2.0

The *Geneva Digital Atlas. Where Technology Meets Humanity* represents the most comprehensive mapping of digital policy actors and internet governance scene in Geneva. The new Atlas edition 2022[476] in form of an attractive free e-book[477] provides on 600 pages, an in-depth coverage of the activities of over 40 actors (International UN Organisations, Academia, Research, international NGOs, Foundations etc.), including the analysis of policy processes and cataloguing of core instruments and featured events. The five main chapters of the Atlas answer five interesting questions:

- Why does technology meet humanity in Geneva?
- When did the policy interplay between technology and humanity begin in Geneva?
- What digital technologies and policy issues are dealt with in Geneva?
- How do legal and policy instruments facilitate digital governance in Geneva?
- Who are the main digital actors in Geneva?

---

[474] The Geneva Internet Platform is an initiative of the Swiss authorities (FDFA, OFCOM) and the Canton of Geneva, https://www.giplatform.org

[475] Diplo is a non-profit foundation established by the governments of Malta and Switzerland. Diplo works to increase the role of small and developing states, and to improve global governance and international policy development, see: https://www.diplomacy.edu/

[476] Digital Watch, https://dig.watch/atlas

[477] The Geneva Digital Atlas, Where Technology Meets Humanity, version 2.0, https://www.yumpu.com/en/document/read/67364941/geneva-digital-atlas-2022/82

### The Observatory

- Maintains a comprehensive live summary of the latest developments in digital policy;
- provides an overview of issues, actors and ongoing processes;
- maintains a live calendar of upcoming and past events, and public consultations, and provides just-in-time reporting from digital policy events[478];
- provides access to the latest research and data on internet policy;
- is enriched by quantitative research (e.g. data-mining of open data, topic profiling);
- relies on a team of 30+ digital policy experts from around the world, for digital policy research and analysis.

### The Taxonomy used on the GIP Digital Watch Observatory

At present, the observatory features over 50 digital policy issues, classified under the following seven topics[479]: Telecommunications infrastructure, Cybersecurity, Human Rights, Legal, Economic, Development, and Sociocultural. Although this categorisation was needed for the purposes of classifying the policy areas into broad areas, many of the digital policy areas interact with other areas in different clusters.

This taxonomy was first developed in 1997 by Dr Jovan Kurbalija, Director of the Diplo Foundation, and then introduced in the book *An Introduction to Internet Governance*[480] (now in its 7th edition). It was later adopted in DiploFoundation's *Internet governance courses*.[481] The taxonomy, which was adopted by other actors such as the UN Commission on Science and Technology for Development in the report Mapping

---

[478] Just in time reporting initiatives, Geneva Internet Platform Digital Watch, https://dig.watch/reporting

[479] Topics, Digwatch, https://dig.watch/topics

[480] Ten years of Internet governance – Tens of thousands of copies – Ten languages, Diplo, https://www.diplomacy.edu/blog/igbook-6th-launch/

[481] Course Catalogue, Diplo, https://www.diplomacy.edu/courses

of international Internet public policy issues,[482] is under regular development to account for emerging trends and new developments.

### The Three Digital Watch Pillars

The *GIP Digital Watch* observatory forms part of the GIP Digital Watch initiative[483], run by the Geneva Internet Platform and operated by DiploFoundation, which provides a solution for practitioners of Internet governance and digital policy who need to acquire information and become informed stakeholders on issues they care about.

The online observatory complements two other pillars forming part of *GIP Digital Watch[484]*:

The Geneva Digital Watch newsletter, a monthly newsletter which includes a round-up of developments, interviews with prominent IG experts, and articles on various digital policy areas. Learn more about the newsletters, access the archive, and download the latest issue.

Monthly GIP briefings on internet governance in Geneva and online, take place on the last Tuesday of every month. In 2018, local hubs were established worldwide with the aim of encouraging sustainable discussions in local communities, and sharing regional perspectives during the monthly briefings. Learn more about the GIP briefings and how to set up a local hub, and access the archive.[485]

---

[482] Meetings, United Nations Conference on Trade and Development, UNCTAD, http://unctad.org/meetings

[483] The Geneva Internet Platform, https://www.giplatform.org/digitalwatch/

[484] GIP Digital Watch, https://dig.watch/digitalwatch

[485] Briefings, Digwatch, https://dig.watch/briefings

# 39

# SATELLITES AND SEA CABLES: THE OWNERSHIP OF DATA CARRIERS AND GEO-POLITICAL IMPLICATIONS

*Anil Prakash, India[486]*

Dependencies on technology, risks in supply chains, and vulnerabilities in critical infrastructure create opportunities for foreign interference. In addition, current geopolitical struggles for power often involve competition in the technological and digital spheres, as different states and tech companies strive for technological dominance and control of the global digital order.

## 39.1 Satellites: Data Carrier Services

The usage of satellite communication for international traffic is minimal compared to fiber optic cables. Approximately 95% of international

---

communication and data traffic is carried through undersea fiber-optic cable networks. Whereas satellites only make up a very insignificant percentage.

Satellites can provide coverage in areas without fiber optic infrastructure and can transmit content to multiple locations and can complement the cable networks in achieving a digital revolution. They can also transmit content to multiple locations simultaneously, which can be useful for broadcasting or distributing data to a wide audience.

Additionally, satellites can complement fiber optic networks by providing an additional layer of connectivity and redundancy. Together, fiber optic cables and satellites can enable the transmission of large amounts of data, including internet traffic, at high speeds globally. They both play important roles in the digital revolution by providing the infrastructure necessary for fast and reliable communication.

## 39.2 Submarine Cables: The Arteries for Data

In recent years, there has been a significant increase in the deployment of submarine cables around the world to support the growing demand for data, cloud-based services, and the "Internet of Things." These cables, which are operated and owned by a variety of state and private entities, are critical infrastructure for the digital age, supporting everything from government communications to societal needs.

However, undersea cables are also vulnerable to attacks and interference, and their governance is often unclear, as legal ownership is divided among co-owners with different nationalities and jurisdictions under international law. The submarine cable sector is currently undergoing a major shift, with tech giants like Alphabet, Facebook, Amazon Web Services, and Microsoft seeking to reshape the ecosystem by building their own private undersea cable networks to connect their data server farms and reduce their reliance on traditional carriers. This trend has

significant implications for global connectivity, data security, and the balance of power in the tech industry.

The concentration of undersea cables in certain regions has led to the emergence of internet chokepoints, or bottlenecks in the global communication network. These chokepoints can be politically and economically significant, as they represent points of vulnerability in the global communication infrastructure.

For example, the Red Sea region, which includes Egypt and the Suez Canal, is a major internet chokepoint due to the high concentration of undersea cables passing through the area. This has implications for the geopolitical power dynamics in the region, as any disruption to the cables could have significant impacts on the global communication network.

The Asia-Africa-Europe-1 (AAE-1) internet cable, which stretches 15,500 miles along the ocean floor, was damaged on June 7, 2021, in an incident that caused temporary internet outages and disconnections for millions of people in several countries. The cause of the damage is unknown, but the impact was immediate, with many African countries losing 90% of their connectivity and cloud services belonging to Google, Amazon, and Microsoft .

The SeaMeWe-4 cable: In December 2006, the SeaMeWe-4 cable, which connects Southeast Asia to Europe, was damaged by a ship's anchor in the Mediterranean Sea, causing widespread internet outages in India, Bangladesh, and the Middle East.

In December 2008, the FLAG FEA cable, which connects Europe to Asia, was damaged by a ship's anchor in the Mediterranean Sea, causing internet outages in the Middle East and India.

These incidents highlight the fragility of world's 550-plus subsea internet cables and the importance of Egypt, the Red Sea, the Mediterranean Sea, and the Pacific and Atlantic Oceans in internet infrastructure. These cables form a large part of the internet's backbone, carrying the

majority of data around the world and connecting to networks that power cell towers and Wi-Fi connections.

Sixteen submarine cables, which are vulnerable to damage from anchors and earthquakes and are no thicker than a hosepipe, pass 1,200 miles through the Red Sea before reaching land in Egypt and continuing on to the Mediterranean Sea, connecting Europe to Asia. The route has become one of the world's largest internet chokepoints in the past two decades and is potentially the internet's most vulnerable point due to its concentration of global movement.

## 39.3 The Convergence of the Two Markets

Satellites can still play an important role in the overall telecommunications infrastructure and could act as 'feeder nodes' into submarine cables. Specifically, satellites can act as a way to connect to submarine cables, allowing submarine cable operators to extend their services beyond the landing stations where the cables reach the shore. This can be beneficial for both satellite and submarine cable operators, as it allows them to maximize the capacity of their networks by working together and ensuring that end users have the connectivity they need.

The potential of the combined use of subsea cable infrastructure and low-orbit satellites could shift this paradigm in the years to come. By combining the two technologies, it may be possible to create a communication network that is more resilient to disruptions and has the capacity to serve more people and businesses around the world. Low-orbit satellites could potentially be used to provide connectivity to areas that are difficult or expensive to reach with subsea cables, such as remote or disaster-affected regions. This could help to bridge the digital divide and bring more people online.

The technical capabilities of submarine cables and low earth orbit (LEO) satellites complement each other in a number of ways. One of the

main areas of convergence is in the extension of submarine cables into new markets beyond the landing station.

LEO satellites can be used to connect inland enterprises and internet service providers (ISPs) in regions where inland connectivity is expensive or unavailable, such as in many rural areas of Africa & Asia.

They can also provide landlocked countries with an alternative means of connecting to submarine cables and can enable competition among submarine cable operators by allowing connection to various ingress points or terrestrial transit points.

When landlocked countries have strained relationships with their neighboring coastal countries or when those coastal countries are facing disruptions and difficulties in maintaining the undersea cables, it leads to internet disruptions and blackouts for the landlocked countries. This has significant negative impacts on the economy and businesses and can lead to internet censorship and throttling as the neighboring countries may use their control over internet access as a political tool.

One example is Afghanistan, a landlocked country that is dependent on neighboring Pakistan and Iran for internet connectivity. Due to political tensions between Afghanistan and these countries, internet connectivity in Afghanistan has been affected in the past. Another example is Belarus, which is a landlocked country in Eastern Europe. Belarus is dependent on neighboring Russia for internet connectivity and due to political tensions between Belarus and Russia, internet connectivity in Belarus has been affected. Also, Ethiopia due to the political tensions between Ethiopia and Eritrea, internet connectivity in Ethiopia has been affected in the past.

This highlights the importance of stable political relations and a well-maintained infrastructure for internet connectivity in these landlocked countries.

In order to minimize the impact of geopolitical tensions on internet connectivity, it is important to have a diverse and redundant infrastructure in place such as satellite connectivity as a backup option. This can

increase the resilience of internet connectivity in the face of political tensions or disruptions to submarine cable systems. Additionally, enhance competition and lower the cost of internet access for businesses and consumers.

Furthermore, *LEO satellites can be used to connect islands* or areas where the traffic density and economics of undersea fiber optic cables are not attractive.

In some cases, *LEO satellites can offer lower latency* compared to submarine cables, which is especially valuable for certain types of traffic that require fast data transfer speeds. As a result, it may be beneficial for submarine and satellite operators to work together to identify and route this high-value, low-latency traffic over satellite infrastructure while keeping the rest of the traffic on submarine cables. The LEO satellite operator could then offer preferential quality of service for this traffic at a higher price than regular traffic. By partnering in this way, both types of operators can benefit from the strengths of each other's networks and maximize the value of their respective infrastructure.

By using *LEO satellites as a redundant connection*, it is possible to minimize the impact of cable failures and ensure that end users have access to reliable connectivity. The use of LEO satellites can be leveraged to offload traffic from submarine cables. In the event of a cable failure, a small portion of the traffic can be redirected to the LEO constellation to ensure that connectivity is maintained. This can be particularly useful in situations where the failure of a submarine cable would result in significant disruption to the flow of data.

## 39.4 Conclusion: Shift of Power?

In conclusion, satellite operators have traditionally focused on providing connectivity in areas where terrestrial networks are not available, such as remote rural areas, as well as in the maritime and aviation sectors. These types of applications are not typically relevant to subma-

rine cable operators and are not likely to be a catalyst for closer integration between the two industries. However, new markets are emerging that may drive further collaboration between satellites and submarine cables, particularly in the area of security and privacy. It is not yet clear to what extent these applications will contribute to satellite traffic or whether they will lead to additional synergies between the two industries.

Going forward, the development and deployment of new technologies, such as the potential combination of satellite and undersea cable infrastructure, may also shift the balance of power in the global communication landscape.

# PART H

## SERVING HUMANS:
## MEDIA, EDUCATION, COMMUNICATION

**40**

# METAVERSE:
# CHANCES AND CHALLENGES
# IN MY GENERATION Z

*Saakshar Duggal, India*[487]

"The Successor of the internet": Well yes this is what many people believe about the Metaverse. These days everyone appears to be talking about this new technological paradigm called "Metaverse". Since the inception of Covid-19, many people have started to hear and read about this new buzz-word Metaverse. It has a huge hype attached to it because it holds the potential to revolutionize the internet space in such a manner that we have never witnessed before.

---

## 40.1 The Meaning of Meta+Verse

Metaverse is definitely the next big version of the internet that's waiting for all of us. The buzz around Metaverse has got substantial thrust and focus, with the Facebook changing its name from Facebook to Meta platforms. Metaverse is representing itself as the new manifestation of the application of Virtual Reality and Augmented Reality in cyberspace. Typically, the word Metaverse consists of two words, Meta plus Verse.

The word "metaverse" derives from the prefix "meta" (beyond) and the stem "verse" (universe). The term means "beyond the universe." Let's see what the two words in "Metaverse" actually mean.

Meta efficiently conveys the idea of transcending reality (as in metaphysics). The more current use of meta as an adjective means "self-referential" or knowingly distinct from the conventional and concrete world.[488] Meta (from the Greek μετά, meta, meaning "after" or "beyond") is a prefix meaning "more comprehensive" or "transcending".[489] Meta means about the thing itself. It's seeing the thing from a higher perspective instead of from within the thing, like being self-aware.[490]

*Metaverse* generally refers to the concept of a highly immersive virtual world where people gather to socialize, play, and work.[491] can be defined as a simulated digital environment that uses augmented reality (AR), virtual reality (VR), and blockchain, along with concepts from social media, to create spaces for rich user interaction mimicking the real

---

[488]   What   is   the   'metaverse'?   Merriam-Webster,   https://www.merriam-webster.com/words-at-play/meaning-of-metaverse

[489] Meta, Wikipedia, https://en.wikipedia.org/wiki/Meta

[490]   Meta,   Urban   Dictionary,   https://www.urbandictionary.com/define.php?term=meta

[491]   What is the 'metaverse'? Merriam-Webster, op. cit.

world.[492] Metaverse is a shared, realistic, and immersive computer simulation of the real world or other possible worlds, in which people participate as digital avatars.

Metaverse is a theoretical or emergent networked online space with digitally persistent environments that people inhabit, as avatars, for synchronous interactions and experiences, accessing the shared virtual space through virtual reality, augmented reality, game consoles, mobile devices, or conventional computers.[493] The Metaverse is a digital reality that combines aspects of social media, online gaming, augmented reality (AR), virtual reality (VR), and cryptocurrencies to allow users to interact virtually.[494]

## 40.2 Metaverse Characteristics

Metaverse as a paradigm has got various salient features and characteristics. It has the following important characteristics:-

It is social;

It is hardware agnostic;

It has user generated content;

It is live and persistent;

It is economically fertile; and

It bridges between different virtual worlds.[495]

It is an infinite space;

---

[492] Metaverse Meaning – What is this New World Everyone's Talking About?, Today Digital, Blog, https://www.xrtoday.com/mixed-reality/metaverse-meaning/

[493] Metaverse, Dictionary.com, https://www.dictionary.com/browse/metaverse

[494] Folger, Jean. 2022. What Does Metaverse Mean and How Does This Virtual World Work? Investopedia, https://www.investopedia.com/metaverse-definition-5206578

[495] Dadwal, Dohit, Characteristics of the metaverse, LinkedIn webpage, https://www.linkedin.com/pulse/characteristics-metaverse-rohit-dadwal/

The connections in Metaverse are just as real;

There's no on-off switch to the Metaverse

It allows users to experience the new state of immersion.[496]

## 40.3 From Metaverse Marriage to Metaverse Money

Metaverse is now beginning to capture the imagination of people. Recently, a couple in Tamil Nadu, India being Dinesh Sivakumar Padmavathi and Janaganandhini Ramaswamy got physically married and held their reception in the Metaverse in the beginning of 2022.[497]

Further, towards the second half of 2022, we saw a couple actually getting married in the Metaverse. Bhopal-based Abhijeet Goel (33), a tech entrepreneur and Dr. Sansrati Jain (34) who tied the knot on Saturday, February 5th, 2022, became India's first couple to get married in a 3D Metaverse. The digital avatars of the couple hosted their wedding ceremony on the beachside and had more than 500 registrations.[498]

Meta Platforms has released various videos about the Metaverse and how it works. Metaverse is going to grow by leaps and bounds. That is the broad picture that is emerging on the basis of various figures and statistics available in the public domain. Some of these statistics are as under:-

- The market size of the metaverse is over $38.5 billion.

---

[496] Nisha Sashidharan. 2022. 4 Characteristics of Metaverse You Need to Learn, Extentia.com, blogpost, https://www.extentia.com/post/4-characteristics-of-metaverse

[497] Holland, Oscar. 2022. My big fat digital wedding: Couple plan India's 'first metaverse marriage', CNN Style, https://edition.cnn.com/style/article/india-wedding-metaverse/index.html

[498] For Better or 'Metaverse' – India's first Metaverse wedding had 500 guests, Clover Infotech News. https://www.cloverinfotech.com/blog/for-better-or-metaverse-indias-first-metaverse-wedding-had-500-guests/

- Over $500 million worth of real estate was purchased in the metaverse so far.
- There are over 400 million metaverse monthly active users.
- 51% of the metaverse user base is 13 or younger.
- By 2026, 25% of people will spend an hour or more in the metaverse each day.[499]
- Meta Reality Labs has already invested $10 billion in the metaverse.
- Decrentraland launched a new platform feature that allows Land-owners to become landlords by renting virtual property.
- Metaverse's global market value is $47.48 billion in 2022, and has a robust CAGR of 39.44%. By 2030 the value is expected to reach $678.80 billion.
- Disney is developing a metaverse theme park. In December, the company filed for a patent for a "virtual-world simulator."
- Gucci opened Vault[500], a metaverse concept store selling "Supergucci" NFTs.
- The corporate world is maximizing metaverse accessibility. Consulting giant Accenture created a metaverse environment called The Nth Floor[501], which mimics real-world offices and lets staff conduct training, on boarding, and other HR-related activities within the metaverse.

---

[499] Hristina Nikolovska, Metaverse Statistics to Prepare You for the Future, BanklessTimes, Blog post, https://www.banklesstimes.com/metaverse-statistics/

[500] Lenihan, Rob, 2022. House of Gucci Vaults into the Metaverse, TheStreet, Blog post, https://www.thestreet.com/investing/cryptocurrency/house-of-gucci-vaults-into-the-metaverse

[501] Leonhardt, Megan. 2022. Accenture designed its own metaverse for employees, complete with exact replicas of offices, Fortune, Blog post, https://fortune.com/2022/04/11/accenture-builds-corporate-metaverse/

- Metaverse is slated to add $5 trillion[502] to the global economy by 2030, given the clamor for placing marketing and events in the virtual space. eCommerce still reigns as the largest economic force ($2.6 trillion), followed by virtual learning ($270 billion), advertising ($206 billion), and gaming ($125 billion), respectively. [503]

## 40.4 Ethics in the Metaverse?

The idea of metaverse sounds extremely exciting and promising. But the question is can we expect the sanctity of the concept of ethics will be automatically be maintained in the metaverse?

Well, unfortunately the answer remains uncertain. Whenever we talk about the existence of ethics in the metaverse, then we tend to face with a lot of grey areas in this domain. This is also because the metaverse is right now in its infancy stage and many developments and regulations have to be made. But at the same time, raising questions at the beginning always paves the way for an efficient growth in the future. However, there is also no denying to the fact that there is nothing illegal about Metaverse in the sense that it has been banned by any particular law and hence it can be said that there is everything legal around Metaverse as a paradigm.

Another important issue for consideration is how do we ensure the ethical behaviour of digital avatars. For that, it is important that the digital avatars must be created in an ethical manner and must be imbibed with ethical values. However, this is where comes a huge challenge. Ethical values are a human concept but how do we imbibe ethical values in the digital avatar is a big question. Certain do's and don'ts concerning

---

[502] Takahashi, Dean, 2022. McKinsey & Co.: Metaverse could reach $5 trillion in value by 2030, Venturebeat, Blog post, https://venturebeat.com/games/mckinsey-co-metaverse-could-reach-5-trillion-in-value-by-2030/

[503] Geri Mileva. 2023. 52 Metaverse Statistics, Influencer Marketing, Blog post, https://influencermarketinghub.com/metaverse-stats/

the activities and behaviours of digital avatars can be incorporated, at the time of their building by the coders This can be in the form of principles governing ethical behaviour, so that the digital avatar itself does not engage in any activity, which is per se unethical or violative of ethical values and standards of humans. However, still the concept of ethics involves morality and inculcating this sense of morality within these digital avatars will be a tricky task to do.

Recently, there was a case where a female avatar of a Metaverse user was digitally raped on the Metaverse by another male avatar. Logically speaking, the said act of Metaverse rape constitutes unethical activity as when you digitally tamper with the digital avatar, you are effectively, doing an unethical activity. In addition, you will also land doing illegal activity in as much as you are tampering with the computer resources and the electronic data logically belonging to another person.

## 40.5 For Metaverse Laws Based on Ethical Principles

The laws and regulations have always been the best friend of ethics, as they help in maintaining the sanctity of ethics in the society. When we look at the physical world, it is thanks to the law and regulations in the society which keeps the citizens under a certain discipline and at the same time, these help in maintaining ethics in the society. However today, we need to understand that as of now, countries have still not woken up to the need for having in place Metaverse laws. Most of the countries are keen that they already have in place their national cyber laws and that cyber laws grant legality to the electronic format and the national cyber laws could be alone relevant and sufficient for the purposes of regulating the Metaverse.

But very quickly the lawmakers have to understand that the Metaverse is a different ballgame altogether. It requires distinctive new set of legal frameworks and principles, which are currently not available under the national cyber legal frameworks of countries. Hence, as

Metaverse grows, it will become even more predominant and relevant from the perspectives of lawmakers to go ahead and come up with new legal frameworks to govern the Metaverse, in order for ethics to prevail in this virtual world.

Therefore, such Metaverse law needs have be based on ethical principles as Metaverse ethics would become an important sub-discipline of the study of ethics.

Today, people still lack clarity as to how the ethical norms of behaviour in Metaverse need to be developed. One thought process is that the existing ethical norms of behaviour, which are available in cyberspace, could be automatically extended on to the Metaverse. Others argue that given the unique nature and characteristics of Metaverse as a paradigm, there will be a need for coming up with new distinctive ethical principles impacting Metaverse at large.

Scholars are also agreeing on the fact that these Metaverse ethical principles must be backed by the sanction of law in order to make them more workable. Further, lot of complicated questions on the intersection of ethics and law are emerging in the context of Metaverse. How can the digital avatars behave? What is the particular limit within which they can behave? What kind of norms are acceptable in Metaverse to be followed by Metaverse digital avatars? What happens if the digital avatar does not follow the accepted norms of behaviour, expected from the digital avatar? Can digital avatar be declared as a rogue digital avatar or a Metaverse criminal? If so, can the said Metaverse criminal be subject to appropriate punishments or coverage under the existing law. These are very fundamental questions which have not been answered at the time of writing this article.

## 40.6 Metaverse: Really Open for Everyone?

"The metaverse is open and is for everyone." When one reads this line, the first thought that comes to our mind is that the Metaverse would

be an ideal world where there will only be constructive positive activities and there will be no negative activity. But that is again a very idealistic expectation from the Metaverse. When we closely look at this statement, then we will realise as to what kind of challenges we are about to face, in order to maintain the ethical challenges in the metaverse. When we say that the metaverse is open, it means that at present there are no certain rules and regulations in this virtual world. The Metaverse is growing with leaps and bounds in its size and its potential applications, which automatically means that crimes in the metaverse are bound to happen at an unimaginable rate. Metaverse has caught the attention of various digital stakeholders including cyber criminals.

Human beings as a species have always been mischievous by nature. We always tend to go out of hand and end up creating a commotion within our surroundings and that's where rules and regulations are required to keep us in control. Different types of crimes have begun to happen in the Metaverse. Lot of existing cybercrime activities have begun replication in the context of Metaverse, with appropriate modifications and changes. Hence, trying to regulate the criminal behaviour in the Metaverse itself throws up large number of very interesting and ethical questions going forward.

Similarly, in the Metaverse, it will become even more difficult to identify a particular Metaverse crime and attribute the same to a particular definitive Metaverse actor. This is so because in the Metaverse, people don't interact directly but interact through their digital avatars. Hence, trying to ultimately affix criminal liability for Metaverse crime would be dependent on the ability of the system to be able to attribute the said Metaverse criminal activity to a definitive actor. That itself is a huge area of potential ethical discussion and evolution in the coming times.

## 40.7 Avatars, a Playground for Cyber Crimes

Further, while there is an expectation that Metaverse avatars ought not to engage in criminal activity, the fact remains that cyber criminals would want to misuse the fertile paradigm of the Metaverse to commit various criminal activities. Crime will be inevitable in the metaverse for a very simple reason, which is that behind the digital avatars are human beings and the human is species is very much vulnerable to have malice in their minds in order to commit crimes, like we see in the physical world.

There is a lot of increase in cyber security breaches in the Metaverse. More and more Metaverse stakeholders are building virtual properties and platforms like the Decentraland. Cyber criminals and cyber security breaches are increasingly finding new options to go ahead and breach the cyber security of these Metaverse properties and assets.

Logically speaking, breach of cyber security of Metaverse properties is unethical and unacceptable behaviour and it ought not to be followed. However, the truth is that stakeholders will continue to keep on engaging in activities tantamounting to breach of cyber security in the Metaverse.

Hence, we will have to put in place appropriate mechanisms for the purposes of enforcing ethical behaviour of Metaverse actors in cyber-space. For that, appropriate legal frameworks will have to be put in place which can effectively go ahead and stipulate legal consequences for those, who go ahead and misuse the provisions of the law as also for those who engage in criminal activities and also Metaverse cyber security breach activities.

The coming of the Metaverse has actually thrown up far more new questions than answers. Further, use of Virtual Reality media in Metaverse in my generation effectively means that people will be now encouraged to come up with more creative expressions of their thought processes and creative abilities.

We are likely to see more growth of not just identical digital avatars but also non-identical digital avatars who would be based on the value and imagination of the concerned relevant stakeholders. The use of Virtual Reality on the Metaverse presents an opportunity, where it can be used as red clay to be moulded into something new distinctive peace of art. No wonder, a digital avatar on the Metaverse whether it is non-identical digital avatar or identical avatar, will often be constituting original, literary and artistic work. This is so because it is effectively representing the creative will and the creative skills of the concerned Metaverse stakeholders.

While the said Metaverse avatars will constitute original, literary and artistic work, the same would be subject to enjoyment of the intellectual property rights. However, this very intellectual property right is likely to be potentially questioned or interfered with by other State and non-State actors. Hence, it becomes even more necessary to ensure that the ethical norms concerning creation of intellectual property rights and utilization of the benefits of intellectual property rights in the said creations in Metaverse need to be appropriately well preserved and well supported. Most of the challenges will relate to the way how digital avatar as a paradigm evolve and how they continue to not just interact with other digital avatars, but will also engage in the Metaverse activities using Metaverse contracts.

This is itself going to throw up new contentious and legal issues. First and foremost, it will be imperative to ensure that Metaverse contracts are given appropriate legal sanction and ethical recognition. But the moment you talk about digital avatars entering into digital contracts constituting Metaverse avatars, a variety of ethical and legal considerations come forward. Is it ethical for a digital avatar to enter into a Metaverse contract? Can the original owner of the concerned Metaverse asset be made liable or responsible for the illegal or criminal acts of the Metaverse asset? If so, how? These are only contentious questions which have to be answered from ethical standpoint. However, how do we effec-

tively deal with such challenges, will have to be dependent on how they deal with the entire issue pertaining to the ethical principles concerning digital avatars.

## 40.8 Metaverse for Generation Z: Open Future

Thus, in my generation, Metaverse presents a new paradigm shift. More and more youngsters are going to be attracted to the Metaverse. They will start spending more time in the Metaverse, and is likely to become more important as it is evident from the various facts and figures available in the public domain (see financial figures above 40.3).

Experts expect by 2040 the metaverse will be a much-more-refined and truly fully-immersive, well-functioning aspect of daily life for a half billion or more people globally[504]

A perusal of the Metaverse facts and figures clearly show that the world is pretty bullish as far as growth of Metaverse is concerned. The increasing use of Virtual Reality in the Metaverse effectively means that there are now new opportunities in terms of generation of new employment, new opportunities of doing commerce, new opportunities of giving vent to your freedom of speech and expression. These and variety of issues will have to be appropriately dealt with by strong legal frameworks in the coming times.

Right now, the world is at a very early stage of development of the Metaverse. Not a single country has a distinctive law or legal framework on the Metaverse but as it starts getting more and more central in our day-to-day lives, such Metaverse will also have to increasingly be governed by ethical principles and regulated by legal frameworks and legal provisions of Metaverse law.

---

[504] Janna Anderson and Lee Rainie. 2022. The Metaverse in 204, Pew Research Center, Blog post, https://www.pewresearch.org/internet/2022/06/30/the-metaverse-in-2040/

# 41

# SEXUAL ASSAULT IN THE METAVERSE: CAVEAT UNIVERSITATES ("UNIVERSITIES BEWARE")

*Divya Singh, South Africa*[505]

## 41.1 Introduction

The term *metaverse* was first used by science fiction author, Neal Stephenson, in 1992 in his novel *Snow Crash*. In the story Stephenson describes a world in which people use digital avatars of themselves to move back and forth between their dystopian reality and a virtual world called the Metaverse.

The term "meta" derives from the Greek prefix meaning "after" or "beyond" and "when combined with English words, meta often indicates

[505] Divya Sing is Chief Academic Officer at Stadio Holdings LTD in Cape Town/ South Africa, a consortium of private universities. She is also a professor of ethics, a lawyer by origin and serving a long time as Vice-Principal Advisory and Assurance Services at UNISA, the largest university in Afrivca. She is also Vice-President of the Board of Globethics Foundation in Geneva. © Globethics Publications, 2023 | DOI: 10.58863/20.500.12424/4276067 | CC BY-NC-ND 4.0 International.

*change* or *alteration*.[506] Looking at its use in information technology, MetaMandrill describes the metaverse as a complex collaboration between the physical and digital worlds, enabled through the application of different technologies and programmes. According to MetaMandrill, it is a persistent living universe that is never turned off and exists in parallel and concurrently with the physical one. Unlike the bounded nature of the physical world, the metaverse is inherently boundless in nature - people enter the metaverse through various means such as virtual reality, augmented reality, headsets, HoloLens's, game consoles, and even apps.[507]

Simplifying the concept for the ordinary reader, Glowiak describes the metaverse as a computer-generated platform in which users may interact with one another in a virtual environment using avatars that are fantasized and/or ideal versions of themselves.[508] Glowiak states:

"The Metaverse is a highly immersive environment in which one is engaged in substantial stimuli. Here the brain functions at high capacity. Most everything is seemingly within reach, which is appealing. Many times, this is immediately accessible. This combination stimulates the pleasure/reward pathway, which positively reinforces engagement."[509]

Or even more succinctly, Mark Zuckerberg on the metaverse experience is widely quoted as saying: *You're in the experience.*

When used properly, engagement in the metaverse can be empowering. As described in the available literature and research, the metaverse enhances feelings of connection and connectedness, creating spaces with the building blocks that eliminate barriers of spatial distance and en-

---

[506] https://www.dictionary.com.

[507] Metaverse Meaning; Different Ways of Defining the Metaverse. Metamandrill.com, Blog post, https://metamandrill.com/metaverse-meaning/

[508] Glowiak, M. 2022. Mental Health in the Metaverse Will Be Like IRL [Note of the Editor: in real life] With Dr. Matt Glowiak, Venteur, Blog post, https://www.venteurmag.com/articles/ metaverse-mental-health-matt-glowiak, 2.

[509] Ibid, 9.

hance the opportunities of equalness and for equality by users. Overall, there is a remarkable sense of being part of something bigger and better. However, there is also an unfortunate negative side to the metaverse. According to Shafir, some people may be more vulnerable to the pull of the metaverse than others. For them, it is virtual reality that might easily become "more real that real life".[510] Research is already showing how the metaverse has the potential – even greater than the regular online environment – to facilitate and enable bullying and harassment of users. Glowiak cautions that social bullies will often create fake avatars for the specific purpose of trolling others, and when blocked or reported, they simply create a new avatar and continue with their anti-social behaviour.[511]

## 41.2 Method

This paper adopts *a library research methodology*, as described by Fitria *et al.*[512] The researcher engaged with source material already available in the libraries (secondary sources) and studied the information through a process of collection, analytical reading, and assimilation of the data. The subject content unique to this paper was studied through the available literature and key issues related to the topic were synthesized and highlighted to support the key constructs.

The new knowledge recommendations and proposals for higher education institutions and universities to understand and mitigate the concomitant risks of sexual assault and harassment amongst the vulnerable

---

[510] Shafir, H. 2023. The Metaverse's Mental Health Implications Are Still Unclear with Hailey Shafir, Venteur, https://www.venteurmag.com/articles/metaverse-mental-health-hailey-shafir, 7.

[511] Glowiak, M. See Note 3, p. 7

[512] Fitria, T,N., Simbolon, N.E. and Afdaleni. 2022. Possibility of Metaverse in education: Opportunity and threat. 1(3). *Jurnal Ilmu Sosial dan Humaniora.* 365-375.

student population draw parallels from and are undergirded by the information derived from the secondary source information.

## 42.3 Discussion

In recent years there has been considerable research conducted on and reported about the negative effects of social media and the scourge of online bullying and harassment in all forms. There is much less researched data about the metaverse but anecdotal and popular press literature is already showing similar trends. The Centre for Countering Digital Hate (U.S.A.) has reported on regular cases of minors being exposed to graphic sexual content and other forms of harassment on VR platforms. Discussing the problem Huddleston emphasises that these identified concerns of bullying and harassment in the online space are often even worse in the metaverse because the level of immersion that occurs with virtual reality could easily exacerbate what we are currently seeing with abuse of social media.[513] There's a potency about being immersed in the VR world that is different from merely observing and interacting through a flat screen monitor, warns Huddleston.[514]

In one of the several reported cases of sexual assault in the metaverse, Nina Jane Patel describes how a group of male avatars sexually touched and groped her avatar without her consent. Patel describes it as a "horrible experience" and a "surreal nightmare".

In 2016 a gamer, Jordan Belamire described her experience of her avatar being sexually assaulted on a virtual platform, with both her brother-in-law and husband watching. According to Belamire, her appeal to the perpetrator to *Stop!* appeared only to goad him on.

---

[513] Huddleston, T. 2022. This is creating more loneliness': The metaverse could be a serious problem for kids, experts say, CNBC Make It, https://www.cnbc.com/2022/01/31/psychologists-metaverse-could-be-a-problem-for-kids-mental-health.html, 1

[514] Ibid, 2.

The *Economic Times Tech*, India reported on a young woman who sought recourse for the violation of her avatar on a virtual reality platform. However, noted the report, the criminal and judicial systems in India - and, in fact, probably the world over – "are still not ready to handle the challenges thrown up by the Metaverse."[515] According to Duggal (quoted in the report) the main problem arises because legal systems simply do not yet recognize (digital) avatars. Further, given the global nature of the metaverse, legal jurisdiction often creates a challenge for many victims. Thirdly, the anonymization of avatars makes identification of the person behind the perpetrator avatar a material constraint in the current milieu with tech companies being quite reticent about revealing user identities. There is also a strong school of thought that for the success of any criminal litigation, the injured party will be required to prove a clear connection to the avatar that was harmed in the virtual world,[516] but the definition and standards of "connection" in this regard remain to be deliberated and debated.

Given the nascence of the metaverse discussions on agreed shared legal principles and norms, regulatory rules are still rather immature and given the complexity of the considerations, it is likely to remain so for a significant time.  However, as the metaverse gains impetus and traction, experts agree that there must be a legal framework to regulate standards of due diligence and legality in virtual reality.  The critical discussion will centre on who will be the owner of the laws, where will sovereignty vest, and who have the adjudicating power given the breadth and unbounded nature of the platforms.

Another key question to be answered is whether the society and the legal experts are ready to recognize that existing real-world laws may be used to adjudicate a dispute originating in the virtual space.  The one

---

[515] Priyanka Sangani, Metaverse crimes challenging, say legal experts, ETtech, June 25 2022, https://economictimes.indiatimes.com/tech/technology/metaverse-crimes-challenging-say-legal-experts/articleshow/92443208.cms?from=mdr, 2.
[516] Ibid, 3.

view is that an act performed in the virtual world is not a real one, does no harm in the real world, and its damage cannot go beyond the boundaries of the virtual world[517] because, of course, it is a game and anyone who feels uncomfortable can easily log off and leave the metaverse.

Specifically referencing the Belamire incident, Basu notes that many responses to Belamire's experience were dismissive, abusive, and even misogynistic.[518] The much-repeated debate was whether what she experienced was groping if her body was not *actually* touched. Based on current legal definitions (pre-dating the tsunami of technology, it must be said), the groping definition may not be legally founded because the law has failed to keep up with technology. So, rather than trying to apply laws that we know are not aligned with the technological advancements, the bigger questions should be (a) whether there was a sexual attack on Belamire? and (b) what was her lived experience of the incident?

Returning to the point of "physical engagement", Huddleston vehemently disagrees with the proposition that a person who is virtually sexually attacked is not assaulted expressly pointing out that "Once you're actually embodied in a space, even though you can't be physically touched, we can be exposed to things that take on a level of realism that could be psychologically assaulting."[519] Explaining the emotional and psychological complexity, Petkova *et al* point to the fact that the realism that accompanies VR experiences readily translates to the same fear experienced emotionally, psychologically, and physiologically when

[517] Aditya S Nair. 2022. The Metaverse And Legal Frameworks Surrounding It, IP & Legal Filings, https://www.ipandlegalfilings.com/the-metaverse-and-legal-frameworks-surrounding-it/ p.4.

[518] Basu, Tanya. 2021. The metaverse has a groping problem already, MIT Technology Review, https://www.technologyreview.com/2021/12/16/1042516/the-metaverse-has-a-groping-problem, 3.

[519] Huddleston, T. op. cit., at 2.

individuals are targeted or threatened.[520]   The intrinsic nature and character of VR is that it is designed to trick the user into believing that s/he is physically in the space.  Analyzing incidents of sexual abuse and harassment on the metaverse, Hoover enjoins a thorough consideration of the fundamental purpose of virtual reality which is after all to dissolve "the gap between the physical and digital selves".  All descriptions of VR emphasize the immersive nature of the experience that heightens the realism and emotional connection.[521]  Thus, when users experience digital sexual harassment, it does not end when the game is over or the glasses removed – in many cases, there is often a more lasting sense of violation and degradation as if the assault took place in actual physical reality.

Recent research also points to the potential for the trauma to be aggravated by the increasing use of haptic technologies that is, technology that mimics the sensory effects of touch making it even more like a physical experience.  A woman researcher from advocacy organization SumOfUs described how she could sense the unwanted touch and grope of her attackers through the VR controllers.  Even though everything was happening virtually, "when a user is touched by another, the hand controllers vibrate "creating a very disorienting and even disturbing physical experience during a virtual assault," she said.[522]

"The panoramic view, audio and even touch stimulation provided by the VR headsets and handheld controls create a multisensory experience,

[520] Petkova, V.I., Khoshnevis, M., & Ehrsson, H.H. 2011. The perspective matters! Multisensory integration in ego-centric reference frames determines full-body ownership. Frontiers in Psychology, 2, p. 35.

[521] Hoover, Amanda, 2022.The metaverse has a sexual harassment problem and it's going to get worse, Morning Brew, https://www.morningbrew.com/daily/stories/2022/06/14/metaverse-has-a-harassment-problem, at 2.

[522] Krishnasai, C. 2022. 21-year-old woman virtually raped, harassed in metaverse: Report, Wion, https://www.wionews.com/world/21-year-old-woman-virtually-raped-harassed-in metaverse-report-483043, at 2.

blurring the separation between the virtual and the physical."[523]   Suc-
cinctly summarizing the arguments, Wiederhold writes:

"In fact, negative experiences in VR may impact victims more than
those that occur on other technology platforms.  VR is different from
technologies such as social media and email in that it is immersive.
When a user enters a virtual environment, the virtual environment be-
comes their world, their avatar becomes their body. Because of this, is
someone is sexually assaulted in such an environment, the trauma can
easily move to the real world. … [T]hose who experience virtual sexual
assault will most likely experience … the same fight or flight responses
that they would have if the incident had happened in the real world.  As
a result, negative virtual experiences can impact people psychologically,
physically, and socially, even when offline.  It is not easy to simply take
the headset off and forget the experience."[524]

Active VR users agree overwhelmingly that that the immersive VR
experience feels real, "even when we tell ourselves it isn't.  Our bodies
respond as if it is real, and we feel the emotional and physical conse-
quences."[525]  In a further study amongst stakeholders at four Canadian
Universities ($n$ = 1733), albeit responding to questions on cyberbullying
(the "Canadian Study) there was a significant disagreement amongst

---

[523] Trang Le. 2022. Sexual assault in the metaverse is part of a bigger problem
that  technology  alone  won't  solve,  Lens,  Monash  University,
https://lens.monash.edu/@politics-society/2022/07/22/1384871/sexual-assault-
in-the-metaverse-theres-nothing-virtual-about-it, at 2.

[524] Wiederhold, Brenda K. Sexual Harassment in the Metaverse. Cyberpsycholo-
gy, Behavior, and Social Networking. Aug 2022, 479-480. http://doi.org/
10.1089/cyber.2022.29253.editorial.

[525] Lorelle VanFossen. 2022. Educators in VR Cyberbullying Team State of the
Metaverse, Educators in VR, https://educatorsinvr.com/2022/03/05/educators-in-
vr-cyberbullying-team-state-of-the-metaverse, at 9.

both male and female respondents with the statement: *Cyberbullying cannot hurt you; it's just words in virtual space.*[526]

Considering the emerging evidence, and the deeper psychological consequences because of the "immersion" and "presence" of virtual reality (when compared with the 'flatter' engagement of online cyberbullying), sexual harassment and assault in the metaverse simply cannot be brushed aside and ignored.  Of course, the answer is never simple and added factors such as the relationship between the user and his/her digital avatar, especially the closeness of identification with the avatar will be factors for constructive consideration before a finding is made.

Glowiak draws specific attention to the mind and psychology of users when engaged in the online and VR environments pointing out that the ability to walk away at any time from the interaction may result in users sometimes underestimating the harmful effects of their behaviour. He notes:

"If all someone sees is a computer-generated avatar with written text or a voice that is not that of the other, it is easy to minimize another's experience."[527]

While this neither exonerates nor excuses a perpetrator who engages in inappropriate behaviour of a sexual nature in the metaverse, this awareness and understanding of "lowered self-awareness" as described by McKenna and Bargh[528], should serve as a caution to institutions that exposure to VR platforms could result in an increase in cases of sexual

---

[526] Faucher, Chantal, Margaret Jackson, Wanda Cassidy, "Cyberbullying among University Students: Gendered Experiences, Impacts, and Perspectives", Education Research International, vol. 2014, article ID 698545, 10p. 2014. https://doi.org/10.1155/2014/698545C., at 8.

[527] Glowiak, M. See above, at 10.

[528] McKenna, K.Y.A. and Bargh, J.A. 2000. Plan 9 from cyberspace: The implications of the internet for personality and social psychology. Personality and Social Psychology Review. 4(1). 57-75.

harassment and abuse between students that would not otherwise happen.

Further, research corroborates the opinion that instances of online sexual aggression and harassment may be increased and empowered by the anonymity that the online space provides. "People often feel emboldened to act inappropriately when they will not be identified or face consequences for their actions."[529] Donat *et al* describe this as the online disinhibition effect where people do and say things that they would not otherwise do and say in face-to-face interactions.[530]   Anonymity also appears to stimulate lower self-awareness and a lowering of consequence consideration.  Furthermore, the power imbalance that often characterizes bullying is much more moderated in the metaverse where a user can model his/her personality and physical attributes to suit any image – the bullied can thus quite easily re-configure themselves as the bully, presenting with both the physical and behavioural attributes, in the metaverse.

Experts agree that the introduction of specific purpose-oriented safeguards, platform monitoring, and more effective moderator deployment initiatives are necessary to reduce the incidents of harassment in the metaverse.  In response to the incidents of harassment and sexual 'assault' on its platforms. Microsoft developed the *Space Bubble* which allows avatars only within a foot of each other before their hands and body disappear.  Meta also introduced *Personal Boundary* which halts any forward movement of another user's avatar within a four-foot radius. Additionally, Meta has a *Safe Zone* which can be activated by users at any time preventing any other user from touching them, talking to them, or interacting with them until the Safe Zone hold is lifted.  In response to the Jordan Belamire incident, *Quivr* introduced an in-game fix that al-

---

[529] Wiederhold, B.K. op. cit., at 1.

[530] Donat, M., Willisch, A. and Wolgast, A. 2022. Cyber-bullying among university students: Concurrent relations to belief in a just world and to empathy. June. Current Psychology. np.

lowed any avatar to make a V gesture, which would automatically push offenders away.

However, notwithstanding the best intentions of the developers, there is an equally strong sentiment amongst researchers that it is impossible to fully enforce these systems and police the virtual space to the extent that the risk is eliminated. Also, we must be conscious not to make the mistake of acting as armchair critics and forget the human emotional factor in the victim's response. It is critical, thus, to acknowledge that it can never be the user's fault for not engaging the security tools. As Nina Jane Patel stated, "[It] happened so fast and before I could even think about putting the safety barrier in place. I froze." Interestingly, Patel also debunks the myth that digital natives will be competent to deal with personal space invasions and uncomfortable engagements in the metaverse. Patel - the victim of the digital rape – was, at the time of the occurrence, herself the vice president of research at another metaverse company.

Of much concern is the fact that there is no clear agreement as to who is responsible for ensuring the safety of online spaces. Is it the tech companies? But many tech companies appear to transfer the onus of safety on the user, providing safety tools and options but not the further wherewithal to moderate their VR platforms. This then begs the question whether a user (for example, a student) who is specifically directed to a platform by a person in authority (for example, the university) may hold the authority responsible for ensuring his/her safety in the VR space?

## 42.4 Conclusion

The metaverse is no longer a dystopian reality – the first certainty is that it is here to stay, the second is that it is going to keep coming, and the third is that it will continue to build and evolve and present new challenges. The Gartner research study (February 2022) predicts that by

2026 approximately "twenty five percent of people would spend at least one hour per day on the Metaverse".[531]

Further, the notion of criminal consequences for conduct in the virtual world is not as far-fetched as it may seem. Nair cites two cases where virtual theft had real world consequences: the first case (2007) involved a Dutch teenager who was charged with theft arising from the theft of virtual furniture in Habbo, a virtual game; and the second case (2005) was the arrest of a Chinese exchange student by the Japanese police over stealing virtual property in another virtual game, Lineage.[532]

Could that then mean that a case such as that of Nina Jane Patel could result in a charge of sexual assault, or harassment? Legal experts do not see the potential for much success because of all the reasons already discussed. Singh makes the point that in Nina Jane Patel's case, while she saw her assailants' faces, knew they were men because she could hear their voices, she would never be able to identify them in a lineup.[533]

In a university environment, however, where institution, implementation and use of the VR platform is far more circumscribed and curated, it conceivable that the system may be set up to enable the institution to (a) link students and their avatars, and (b) de-anonymise the avatar under specific circumstances. Were this indeed possible, it may then be more imaginable that the person behind an avatar could be held legally responsible for their actions, with both civil and even criminal repercussions.

A second issue which also warrants consideration - and which is the crux of this paper - is whether a student who is a victim of sexual crime on the metaverse platform that is part of the university's learning man-

---

[531] Gartner. 2022. Gartner Predicts 25% of People Will Spend At Least One Hour Per Day in the Metaverse by 2026, op. cit.

[532] See Aditya S Nair. 2022. The Metaverse And Legal Frameworks Surrounding It, IP & Legal Filings, op. cit., at 5.

[533] Singh, K. 2022. https://www.refinery29.com/en-us/, at 1.

agement system, may have legal recourse against the university if s/he can prove that (a) the university foresaw the possibility of harm; and (b) notwithstanding, continued to expose students to the risk. Part (a) is a question of fact but there is ample evidence in the available literature that as tech becomes ubiquitous and, notes Wiederhold, increasingly immersive, incidents of technology-assisted or technology-facilitated abuse have also become more and more pervasive and common.[534]

Considering legal liability, Myers and Cowie reflect on the efficacy of pursuing a human right rather than criminal law approach to the incident, specifically highlighting the right to respect of her private life, which includes her psychological integrity.[535] Donat *et al* add the notion of justice psychology and specifically her right to interpersonal justice.[536]

Universities that see their responsibility and acknowledge the risk or, for that matter, any institution that considers introducing the metaverse to users, must take steps to educate itself and its users. And, argues Hoover, it is a facile to be rely on the notion that students can immediately leave the metaverse when the feel uncomfortable,[537] because that effectively transfers the onus to the victim to take responsibility for the bad behaviour of others in similar vein to asking, when a woman gets raped, was she wearing a short dress?

It is key that as more users are exposed to the metaverse, there is acknowledgement for the fact that not all students will be digitally literate with equal sophistication of engaging and dealing with the metaverse and virtual reality. All users, even the so-called digital na-

---

[534] Wiederhold, Brenda K. Sexual Harassment in the Metaverse. Cyberpsychology, Behavior, and Social Networking, op. cit., at 1.

[535] Myers, C-A, and Cowie, H. 2017. Bullying at university: The social and legal contexts of cyberbullying among university students. 48(8). 1172-1182.

[536] Donat, M., Willisch, A. and Wolgast, op. cit.

[537] See Hoover, Amanda, 2022. The metaverse has a sexual harassment problem and it's going to get worse, op. cit. at 5.

tives, must be educated and continually updated on ways to protect themselves. The institutional onboarding process prior to the engagement with the VR platform is critical.[538] The experience of Nina Jane Patel and particularly her reaction when confronted with the assault is apposite in underscoring the training requirement – while safeguards are in place and students may be digitally literate, users need continuing education and re-education on how to use the self-moderation safety tools, and more importantly, how to use them in an attack situation when fear, shock, or trauma can completely cloud one's responses.

The CEO of the Centre for Countering Digital Hate (U.K.) also highlights the requirement that virtual reality requires that safety must be built into the platforms from the inception of its development (and not be reactionary) because you can't do daily searches of the metaverse for sexual abuse and then remove it like with some other online platforms. "You can't," he states. "It happens in an instant [and] there's nothing you can do."[539]

There is a strong argument that the primary onus should be on the tech companies to prioritize the safety of their users but while responsibility is being debated, parents, users, universities, businesses, and any entity engaging with the metaverse and implementing platforms needs to be proactive and diligent in enquiring *whether there is a reasonable level of comfort and belief that the specific platform developer has put in place the necessary safeguards and precautions to ensure optimal safety for the users in cyberspace.* However, considering universities specifically, this begs the question whether at this stage in our knowledge of the metaverse, all implementer-roleplayers have the necessary experience and understanding of virtual reality to ask the right question and monitor users' exposure and the concomitant threats.

In addition to the cases of sexual assault on the metaverse, other concerns of a specifically sexual nature are catfishing (which refers to the

---

[538]  In Basu, T. op. cit., at 4.

[539] In Huddleston, T. 2022. This is creating more loneliness, op. cit. p.3.

phenomenon of creating and portraying complex fictional identities through online profiles or persona for the purpose of luring another into a fraudulent romantic relationship,[540] with its allied apprehensions of sextortion. Reflecting on her research into sextortion conducted across the U.S., Hinduja found that 5% of youth had been victimized. What is alarming is that the research covered only social media and messaging apps with no concentration on online-based social engagement platforms, or metaverse environments.

Universities must be equally alive to all these dangers and implement appropriate protective mechanisms to prevent the harms from materializing. "[With] the unique context of the metaverse environments, and the immersive nature of VR technologies, it is possible that an increased number of users will be manipulated or deceived."[541] This may be exacerbated in the university set-up where the trust relationship between students may be more easily assumed. As suggested by Hinduja, "Once trust is given to someone else without verifying their identity, various kinds of victimization can readily take place."[542]

As universities look at enhancing their teaching and learning opportunities in the online environment, presenting digital content that is significantly improved by 3D virtual reality carries with it enthralling potential. It is easy to be tempted by the innovative possibilities of the metaverse with the promise of a learning experience the likes of which have probably not been seen before in most classrooms and institutions. For many academics, virtual reality is the silver bullet for relevant teaching and learning that will ensure that graduates are workplace ready, trained with the appropriate twenty-first century skills, and imbued with the graduate attributes that will produce socially responsible graduates.

---

[540] Hinduja, S. nd. The Metaverse: Opportunities, Risks, and Harms, https://cyberbullying.org/metaverse, p.6

[541] Ibid, at 6-7.

[542] Ibid, at 7.

All of this can be productively engaged and holistically assessed using the real-world environment replicated in the metaverse.

On the other hand, universities also owe their students a legal (and moral) duty of care and until issues of metaverse safety have been thoroughly interrogated, investigated and the institution is satisfied that it has done all things reasonable to protect the users, it may be difficult (and some may say wrong), for a university to ethically (certainly) and legally (possibly) lure students there with the promise of stunning futurist experiences.[543]

The results of the Canadian study were interesting on this point with 44% of the male respondents and 38% of the female respondents being of the belief that *it was not the university's responsibility to stop or prevent online bullying.* Hoover does not disaggregate the findings by stakeholder cohort, and it was therefore not possible to analyze the findings further. Equally disconcerting was the high level of acceptance and tolerance (and some might say cynicism) in the Canadian study that *Cyberbullying is a normal part of the on-line world; it cannot be stopped.* Only 47% of the males and 32% of the females who responded to the survey disagreed with the statement.[544]

Acknowledging that there are no laws dealing with abuse in the metaverse specifically, it is important to recognize that there are existing legal rules and frameworks from which a victim may yet be able to draw, especially when principles of duty of care and equity are engaged. When cyberbullying was identified, no-one said kill the online platforms – similarly, with the metaverse. *However, what universities can do is:*

1. Understand and acknowledge the concerns and problems.

---

[543] See further: Hoover, Amanda, 2022.The metaverse has a sexual harassment problem and it's going to get worse, op. cit. p.6.

[544] Hoover, ibid., p. 8.

2.   Embed and implement proactive security measures on the VR platforms being developed for the institution (or ensure that they are in place if off-the-shelf products are used).

3.   Set internal institutional systems, structures, rules, and codes of conduct that will manage and mitigate risks.

4.   Ensure that reporting structures are empathetic and progressive, and that tools and processes don't victimize the victim.  Introduce student counselling norms that understand the consequences of digital and virtual abuse and harassment.

5.   Continually train students on their social and ethical responsibilities and particularly the importance of respect for each other, but also ensure that training is implemented regarding the use of the VR platform and specifically on the safety measures, as well as on the security, if any, in place to identify perpetrators.

So, please don't throw the baby out with the bath water: rather, the appeal is for proper management, governance, and regulation if the potential for risk is created; and informed decision-making by university management, (1) not being swayed by the publicity and hype of the metaverse and virtual reality, and (2) balancing the opportunities with the threats presented.  These are key imperatives for universities intent upon guarding against and protecting both themselves and their students from the attendant ethical dilemmas and legal challenges that may present with the introduction of the metaverse and virtual reality.

# ARTIFICIAL INTELLIGENCE IN EDUCATION

*Wayne Holmes/ Maya Bialik/ Charles Fadel*[545]

## 42.1 AIED Can Offer More

Much of the *Artificial Intelligence in Education (AIED)* involves the application of AI techniques to mainstream learning approaches, and tends to reflect (or automate) existing educational assumptions and practices. In addition, much AIED has been designed (whether intentionally or not) to supplant teachers or to reduce them to a functional role[546] and

---

[545] The article is an excerpt from Wayne Holmes/ Maya Bialik/ Charles Fadel, Artificial Intelligence in Education, The Center for Curriculum Redesign, Boston, 2019, 151-180. With permission of the publisher. © Globethics Publications, 2023 | DOI: 10.58863/20.500.12424/4276068 | CC BY-NC-ND 4.0 International.

[546] Worryingly, one of the developers we have mentioned has suggested that the sophistication of their AIED means that teachers only need to play an auxiliary role, working like fast-food chefs ("KFC-like") to strictly regulated scripts.

not to assist them to teach more effectively. This approach, while potentially useful in contexts where teachers are few and far between, clearly undervalues teachers' unique skills and experiences, as well as learners' needs for social learning and guidance. However, instead of just automating the teaching of students sat at computers, conceivably AI might help open up teaching and learning possibilities that are otherwise difficult to achieve, that challenge existing pedagogies, or that help teachers to be more effective. Here we will speculate on some possibilities, some of which have been foreshadowed by the AIED tool, while others are both novel and complex to achieve, and all of which raise interesting social questions. We begin with AI to support collaborative learning, then AI-driven student forum monitoring, AI to support continuous assessment, AI learning companions for students, and AI teaching assistants for teachers, before concluding with AIED as a research tool to further the learning sciences (i.e. in order to help us better understand learning).[547]

### 42.1.1 Collaborative Learning

Collaborative learning, where students work together to solve problems, is well known to be able to lead to better learning outcomes, but effective collaboration between learners can be difficult to achieve.[548] AIED offers various possibilities. To begin with, an AIED tool could automatically suggest groups of students best suited for particular collaborative tasks, drawing on and making intelligent connections between individual student models (each of which comprises knowledge about

---

[547] One intriguing use of AI in education that we will not consider in detail, because its efficacy has not yet been demonstrated, but that should still be acknowledged is the automatic generation of quiz questions (https://mt.clevere.st and https://learningtools.donjohnston.com/product/quizbot).

[548] Luckin, R., et al. 2017. Solved! Making the Case for Collaborative Problem-Solving. Nesta. https://www.nesta.org.uk/report/solved-making-the-case-for-collaborative-problem-solving/

the student's previous learning experiences and achievements, what the student is learning in other classrooms, their personalities, and more).[549] Having elicited the teachers' requirements, the tool might also suggest groups of mixed or similar-ability students, or groups designed to give particular students opportunities to take on leadership roles, or groups that avoid personality or temperament clashes, and so on, all the while enabling the teacher to quickly and easily override any of the tool's suggestions (which the AI will learn from, for next time). An AIED tool might also take on the role of expert facilitator or moderator, monitoring student collaborative activities, recognizing when students are having trouble understanding shared concepts, and then providing targeted support. Alternatively, the AIED might involve a virtual agent that actively contributes to the group discussions (acting as a virtual peer or a teachable agent), or that makes dynamic connections (either with discussions being held by other groups in the same classroom, or with relevant materials drawn from the semantic web). In fact, some research into AI to support collaborative learning has been undertaken,[550] but there are many technical issues to overcome before it becomes possible in real classrooms.

---

[549] The Universitat Politècnica de València have been researching just such a system: Alberola, J.M., del Val, E., Sanchez-Anguix, V., Palomares, A., and Teruel, M.D. 2016. "An artificial intelligence tool for heterogeneous team formation in the classroom." Knowledge-Based Systems 101: 1–14. https://doi.org/10.1016/j.knosys.2016.02.010

[550] E.g., Diziol, D., et al. 2010. "Using intelligent tutor technology to implement adaptive support for student collaboration." Educational Psychology Review 22 (1): 89–102. https://doi.org/10.1007/s10648-009-9116-9 and Spikol, D., et al. (2016). "Exploring the interplay between human and machine annotated multimodal learning analytics in hands-on stem activities." In Proceedings of the Sixth International Conference on Learning Analytics & Knowledge. 522–523.

### 42.1.2 Student Forum Monitoring

Increasingly, students of all ages are participating in online education, which usually involves the use of discussion forums. Students might post to forums in response to given tasks or to engage in collaborative learning opportunities, or they might want to contact their tutors to clarify course requirements or to ask about course materials. Accordingly, especially when there are large cohorts of students (as can be typical of some distance universities and MOOCs), these online forums can generate massive numbers of forum posts, all of which must be monitored, moderated, and addressed. However, as the number of forum posts increases, this becomes at best an inefficient use of a tutor's time (dealing with repetitive and minor practical issues) and at worst an increasingly impossible task. It also makes it difficult for students to keep up to date with other student posts that might connect to their interests.

AIED might help in a number of ways (again, some research has already been conducted in this area)[551] — in particular by helping the teachers/tutors to be better able to support their students. First, an AIED tool might triage the forum posts, identifying those that can be dealt with automatically (perhaps practical questions around course dates, such as "When do I need to submit…?"), and those that require a response from a human tutor (such as those discussing more in-depth core subject issues). The simple posts, the ones that the AIED is capable of handling, would receive immediate automatic responses, relieving the human tutors of much repetitive work while enabling the students to move on quickly to more substantive work. Other posts would automatically be referred up to a human tutor, to ensure that students receive high quality, appropriate responses whatever the nature of their posting.

---

[551] Goel, A.K., and Joyner, D.A. 2017. "Using AI to teach AI: Lessons from an online AI class." AI Magazine 38(2): 48. https://doi.org/10.1609/ aimag.v38i2.2732

Taking this a step further, the more demanding posts (of which there still might be many) would be further analysed, the aim being to identify and aggregate similar posts or posts that raise overlapping issues (in a course with a thousand students, it is unlikely that there will be a thousand unique responses to a single course activity, but rather a much smaller number of closely related posts). A human tutor would then write a response to the much smaller number of aggregated posts, which in turn would be issued to all of the original posters. Although this is unlikely to be as good as replying to each individual student, it would clearly be better than the students receiving no responses at all—which, in a large online course, can all too often be the case. Another approach that might also help in student forums is for the AIED to interpret and make dynamic connections between posts, informing tutors when particular issues have been raised (e.g., known and unknown misconceptions), for them to address, or informing students about other posts that they might find interesting.

Finally, the AIED might also use sentiment analysis AI techniques to identify posts that reveal negative or non-productive student emotional states (perhaps a student is overly challenged, or likely to drop out of the course, or possibly suffering from mental health issues), posts that are unacceptable (perhaps because they include racist, misogynist or gratuitously aggressive comments), or posts that suggest topic drift (the tendency for forum posts to drift from the original intent). Any such posts (which, because of the overall number of posts, can be easy for humans to miss) would be referred up to a human tutor, so that the tutor can respond quickly, appropriately and effectively (perhaps by calling the student by phone, rather than depending on a digital intervention). Together, these various techniques might also enable tutors to be kept well informed of student opinions, collective worries, or recurrent themes that emerge from the forums.

### 42.1.3 Continuous Assessment

Psychologists and educators know that it is wrong to make decisions based upon a single test score and that decisions should reflect a balanced, complete understanding of each child. Numbers and scores can be very misleading if we don't consider the whole picture, something that means using both a qualitative and quantitative approach.[552]

Although there is little evidence for their validity, reliability or accuracy, high-stakes examinations are core to educational systems around the world.[553] Perhaps this is because that is how it has always been, perhaps because they efficiently rank students, perhaps because no practical, cost effective at scale, alternative has ever been devised, or perhaps because those who run the systems are typically those who were most successful at exams (and do not emotionally resonate with the need for change). Whatever the reason, with high-stakes examinations in place, schools and universities all too often end up teaching to the test, prioritising routine cognitive skills and knowledge acquisition over in-depth understanding and authentic application. In other words, the examinations, rather than the needs of students or wider society, determine what is taught and learned. Meanwhile, ironically, AI technologies are automating exactly the type of knowledge that examinations mostly assess: "There's lots of elements of human intelligence that cannot be automated but the bit that we've tended to value, that relates to academic exam success, is one of the bits that we've managed to automate."[554] In any case, stop-and-test examinations (standardised, unseen tests that are at set points in the learning schedule, thus potentially interrupting the learning) are not able to rigorously evaluate a student's understanding of all that has been learned—at best they can only provide a snapshot of frag-

---

[552] Gunzelmann, B.G. 2005. "Toxic testing: It's time to reflect upon our current testing practices." Educational Horizons 83 (3): 214.

[553] Evolving Assessments for a 21st Century Education.

[554] Rose Luckin quoted in https://www.jisc.ac.uk/news/the-ai-revolution-is-here-17-aug-2018

ments of what has been studied over the duration of a course. Last, but not least, students of all ages can sometimes suffer from serious exam anxiety, which can easily negatively impact on the student's success in a typical three-hour end of course examination (further clouding their accuracy and trustworthiness).

Nonetheless, most AIED research in this area has been unambitious. It has focused on improving existing examination systems (developing AI-driven techniques to authenticate the identity of students taking exams online),[555] rather than challenging the underlying principles. However, as we have seen, typical ITS and other AIED tools already and constantly monitor student progress to provide targeted feedback and to assess whether the student has achieved mastery of the topic in question. Similar information could be captured by AIED tools designed to support collaborative learning, while intelligent essay assessment tools can also make inferences about a student's understanding. All of this information and more might be collated throughout a student's time in formal educational settings (the learning sciences have long understood the value of students engaging with constructive assessment activities), together with information about the student's engagement with non-formal learning (such as learning a musical instrument, or a craft or other skills) and informal learning (such as language learning or enculturation by means of learning from experience or active participation), to help create a picture of the whole learner. In other words, the AI-driven assessment would happen in the background, all of the time—making it next to impossible for students to cheat or subvert the system's intention (as can be the case when wealthier students employ personal tutors),[556] or take the test as many times as necessary until they achieve a good-enough score.

---

[555] For example, http://tesla-project.eu

[556] Luckin, R. 2017. "Towards artificial intelligence-based assessment systems." *Nature Human Behaviour* 1. https://doi.org/10.1038/s41562-016-0028

This more detailed and nuanced information about an individual student might then be represented (and perhaps visualised in dynamic graphics) in an AI-driven e-portfolio,[557] an intelligent personal resumé (in fact, an extended open student model). This e-portfolio could perhaps be underwritten and authenticated by blockchain technologies[558] as used by virtual currencies such as Bitcoin (essentially open, distributed ledgers, hosted simultaneously by millions of computers across the internet and linked using cryptography, that can share data in a verifiable, incorruptible, and accessible way). In this way, students would have a robust, accredited, in-depth record of all their learning experiences and achievements, far more detailed and useful than a collection of certificates. Parts or all of this smart resumé they might share when applying for admission to another course or for a new job, while retaining full control of their academic persona and data. From a learner's perspective, an additional benefit is that continuous assessment can act as a moving average, a fluid-like shock absorber that evens out the blips of bad days and disadvantageous personal situations (it simply does not sense that a young person's academic outcomes and future life can be determined by difficulties at home that coincide with the day of an important exam).

In short, although the constant monitoring of student behaviours and achievements raises significant and far-reaching ethical questions that must first be properly investigated and addressed, it is conceivable that stop-and-test examinations could soon be entirely removed from our educational systems and relegated to a more primitive past.

---

[557] Per one of the authors' US patent numbers 9,262,640 and 9,582,567, which also protect privacy and security.

[558] Sharples, M. and Domingue, J. 2016. "The blockchain and kudos: A distributed system for educational record, reputation and reward." In European Conference on Technology Enhanced Learning. Springer. 490–496.

### 42.1.4 AI Learning Companions

The smart resumés that we have just proposed could also play a role in a much larger AIED possibility: AI-driven lifelong students' learning companions.[559] As we have seen, the desire for every student to have their own personalised tutor is what first inspired the development of ITS, but what about taking this to its logical conclusion? AI has the potential to provide every student with their very own personalised learning companion, operating sometimes as a learning partner, other times as a guide through the mass of available learning opportunities, and sometimes as an instructor, all the time recording the student's interests and progress in their blockchain-protected, smart resumé. The arrival and rapid developments of Siri, Cortana, Google Home and Alexa, suggest that this possibility is tantalisingly close.[560] In many countries, smartphones with extraordinary processing power and always-on internet access are more than common. It would not necessarily be a big technical step to leverage these capabilities, to create an AI-driven smartphone learning companion that could accompany and support individual learners throughout their studies, from kindergarten to old age.

Such a learning companion brings many possibilities. Once the student has decided on a particular topic of interest, it might provide some instructional activities, monitor the student's progress, remind them when a task needs to be completed, and offer targeted feedback and guidance—all on their speech-driven smartphone (and available on all their other devices). In other words, it might function as what we have called an ITS+.

But a learning companion would also operate at a higher and more strategic level. Building on the student's individual interests and life goals, it could also help them decide what to learn, as well as where and

---

[559] The University of Southern California have been researching just such an application over many years: http://ict.usc.edu/prototypes/personal-assistant-for-life-long-learning-pal3

[560] Alexa, Should We Trust You?

how to do the learning (the companion might identify and connect with the learning opportunities that are available, both formal and informal, both on and off-line). It could then also guide the student along over-arching long-term individualised learning pathways designed to help the student address their emerging personal life-goals, connecting their learning interests and achievements, while reminding them of and encouraging them to reflect on and perhaps develop their long-term learning aims. The learning companion[561] might suggest learning opportunities that focus on some so-called 21st Century Skills,[562] and social-emotional learning.[563] It could also potentially connect learners, in the same classroom or from opposite sides of the world, depending on their shared interests and goals, helping them develop and work together in projects that prioritise both individual and collective achievements (and, in turn, helping to promote other critical skills in collaboration, teamwork, and intercultural awareness).

### 42.1.5 AI Teaching Assistant

As we have noted several times, most AIED technologies are designed with the aim of relieving teachers of the grunt work of teaching (most often by automating time-consuming activities such as the marking of classroom or homework assignments). However, despite these best of intentions, many AIED technologies in effect take over teaching (they deliver personalised and adapted learning activities better than teachers), or at least they reduce teachers to a functional role (perhaps their job is to work to strictly regulated scripts, or to ensure that the technology is ready for the student to use). Nonetheless, as we and col-

---

[561] World Economic Forum. 2015. New Vision for Education: Unlocking the Potential of Technology. World Economic Forum.

[562] Trilling, B. and Fadel, C. 2012. 21st Century Skills: Learning for Life in Our Times. John Wiley & Sons.

[563] Fadel, C., Bialik, M., and Trilling, B. 2015. Four-Dimensional Education: The Competencies Learners Need to Succeed. Centre for Curriculum Redesign.

leagues have written previously: Crucially we do not see a future in which AIED replaces teachers. What we do see is a future in which the role of the teacher continues to evolve and is eventually transformed; one where their time is used more effectively and efficiently, and where their expertise is better deployed, leveraged, and augmented.[564]

This might be more of an emotional plea than a coherent argument—but it assumes that teaching involves more than delivering knowledge, and that it is a fundamentally social process. From this perspective, a key role for AI is supporting teachers to teach and support students.

One way in which this might be achieved is by augmenting teachers' expertise and skills with an AI teaching assistant, to complement and work with the students' AI learning companion, that goes far beyond the useful but by comparison somewhat primitive teacher dashboards featured in so much education technology. This would be a key way that AIED can support teachers to support students. Just such a possibility has been explored in the short narrative "A.I. is the New T.A in the Classroom,"[565] which describes a possible classroom of the future in which the teacher is supported by a dedicated and personalised AI teaching assistant (AI TA).

Many of the ideas we have suggested could play a role in this possible scenario (such as automatically setting up collaborative groupings of students, replacing stop-and-test examinations with AI-supported continuous assessment, and managing peer-marking and undertaking some automated marking). The AI TA could also automatically provide teaching and professional development resources (texts, images, videos, augmented-reality animations, links, network connections) that the teacher

---

[564] Luckin. R., et al. Intelligence Unleashed, 11. https://www.pearson.com/content/dam/one-dot-com/one-dot-com/global/Files/about-pearson/innovation/Intelligence-Unleashed-Publication.pdf.

[565] Luckin, R., and Holmes, W. 2017. "A.I. is the new T.A. in the classroom." How We Get To Next. https://howwegettonext.com/a-i-is-the-new-t-a-in-the-classroom-dedbe5b99e9e

might choose to call upon to support their teaching. It could also monitor the students' performance as they engage in their classroom activities, continuously updating their learner models, making connections with the domain models of topics being taught, and tracking progress over time. All of this information (together with data about each student from additional sources: assessments from other classes, informal learning achievements, and relevant medical or family information) could be readily available to the teacher, whenever the AI TA computes it might be useful or whenever the teacher calls for it. In this possible future, what and how to teach the students, and how best to support them, would remain the responsibility and prerogative of the teacher. The AI TA's role would simply be to make the teacher's job easier and more effective.

### 42.1.6 AIED: a Research Tool to Further the Learning Sciences

As has probably been noticed, each of AIED's possible future uses are firmly rooted in existing AIED research and approaches. This is no less true of our final example, the use of AIED as a research tool to further the learning sciences. Implementing an educational practice in any technology means that the practice has to be both better understood and then systemized. As a consequence, the technology acts much like a virtual spotlight, highlighting issues that have existed for years but that have been hidden or overlooked (for example, around the most effective approaches to teaching). This is particularly true of the introduction of AI to education, which is beginning to throw an extraordinarily bright spotlight onto many learning sciences issues. However, while there have been notable developments in this area of AIED research, mostly it has been at a relatively theoretical level, such that their potential and implications remain somewhat unclear.

In fact, AIED as a learning sciences research tool is often linked to a pair of other independent but overlapping academic fields that use statis-

tical techniques drawn from big data research:[566] learning analytics and educational data mining.[567] While learning analytics involves "the measurement, collection, analysis and reporting of data about learners and their contexts, for purposes of understanding and optimising learning and the environments in which it occurs;"[568] educational data mining "is concerned with gathering and analysing data so as to understand, support and improve students' learning."[569] One example, that avoids this distinction, and that has been shown to be effective, is The Open University's OU Analyse[570] tool, which draws on data from across the university (such as student access of online learning materials, submission of assessments, and outcomes) to identify students who might be at risk of dropping out from their studies—to enable tutors and student-support staff to provide appropriate pro-active remedial support. In fact, with the fields continually informing and cross-fertilizing each other, the distinctions between learning analytics, educational data mining, and AIED as

---

[566] Mayer-Schonberger, V. and Cukier, K. 2013. Big Data: A Revolution That Will Transform How We Live, Work and Think. John Murray.

[567] Readers who would like to learn more about the similarities and differences between learning analytics and educational data mining might be interested to read Benedict du Boulay and others, "What does the research say about how artificial intelligence and big data can close the achievement gap?" in Luckin, R. (ed.) 2018. Enhancing Learning and Teaching with Technology. Institute of Education Press, 316–27; or Siemens, G., and Baker, R.S.J.d.. 2012. "Learning analytics and educational data mining: Towards communication and collaboration." In Proceedings of the 2nd International Conference on Learning Analytics and Knowledge, 252–254. http://dl.acm.org/citation.cfm?id=2330661

[568] Siemens, G. 2011. "1st International conference on learning analytics and knowledge 2011: Connecting the technical, pedagogical, and social dimensions of learning analytics. https://tekri.athabascau.ca/analytics/about

[569] Du Boulay, et al., "What does the research say about how artificial intelligence and big data can close the achievement gap?" 270.

[570] See Herodotou, C., et al. 2017. "Predictive modelling for addressing students' attrition in higher education: The case of OU analyse." http://oro.open.ac.uk/49470/and https://analyse.kmi.open.ac.uk

a learning sciences research tool are becoming increasingly blurry. Often, it simply comes down to the communities who are involved in the research and the terminology that they use. Here, as we are writing about AIED, we will continue to use AIED terminology.

One prominent example of AIED as a learning sciences research tool has recently been published by the Medical Research Council Cognition and Brain Sciences Unit at the University of Cambridge.[571] The traditional grouping of students with learning difficulties in broad categories such as ADHD, dyslexia, and autism has long been known to be insufficiently helpful, when educators try to improve learning outcomes for individuals. For this reason, the Cambridge researchers are investigating the use of machine learning to categorise struggling students at a more granular level (based on measures of listening skills, spatial reasoning, problem solving, vocabulary, and memory). By analysing data from more than 500 children, the machine learning revealed four clusters of learning difficulties (which had not previously been so clearly delineated): difficulties with working memory skills, difficulties with processing sounds in words, broad cognitive difficulties in many areas, and typical cognitive test results for the student's age. The researchers found that diagnosing struggling learners in terms of these four clusters was both more accurate and more useful, helping educators address individual learning difficulties, than the traditional diagnostic labels.

We will conclude our brief discussion of AIED as a learning sciences research tool with one final example, one that is in the early stages but has important potential: the use of machine learning to improve learning design. Learning design refers to a range of methodologies "for enabling teachers/designers to make more informed decisions in how they go

---

[571] See Astle, D.E., Bathelt, J. and Holmes, J. 2018. Remapping the cognitive and neural profiles of children who struggle at school." Developmental Science. https://doi.org/10.1111/desc.12747 and, for a short summary, https://www.opencolleges.edu.au/informed/learning-strategies/artificial-intelligence-identifies-students-struggle-school

about designing learning activities and interventions."[572] These methods are intended to inform decisions about pedagogy (teaching and learning) and about ways to support student learning experiences, and can also be used to provide core data for learning analytics or educational data mining. Most approaches in use in universities[573] draw on teachers' professional knowledge of teaching and learning (knowledge that is often tacit and thus has had to be elicited from them, which is a non-trivial task and can lead to fuzziness and inconsistencies). Instead, the approach currently being researched at the Open University involves machine learning from thousands of existing module activities to identify categories of activities at a highly granular level. Once these learning design activity categories are identified, and have been robustly authenticated, it should then be possible to correlate the actual learning designs of course modules with student outcomes, to help us better understand how students learn. In turn, this might inform teachers and learning designers about which learning designs (depending on, for example, domain, specific subject, duration and level of study) are most effective in practice.

## 42.2 AI in Education—A Tentative Summary

In the previous sections we have discussed a wide variety of existing and potential AIED technologies. One way to access this variety is to consider the technologies in terms of whether they are mainly student teaching (they take a mainly instructionist approach), or student support-

---

[572] Conole, G. 2012. Designing for Learning in an Open World (v. 4). Springer Science & Business Media.

[573] E.g., Cross, S., et al. 2012. "OULDI-JISC project evaluation report: The impact of new curriculum design tools and approaches on institutional process and design cultures." http://oro.open.ac.uk/34140/; Laurillard, D., et al. 2013. "A constructionist learning environment for teachers to model learning designs." Journal of Computer Assisted Learning 29 (1): 15–30; Dalziel, J. (ed.), Learning Design. Routledge.

ing (they take a mainly constructivist approach), or teacher supporting (they mainly help teachers do what they already do but more quickly or with less effort). A summary representation of this is shown in the following table. A cursory examination of this table will reveal that the categorization provides only a high-level overview, while many of the AIED approaches overlap, and most of the technologies could easily appear in another place in the table. It is also likely that over time different AIED technologies will merge into multi-capable systems, perhaps incorporating sequenced (ITS), Socratic (DBTS), and self-directed (ELE) learning in one technology.[574]

| | Student Teaching (mainly instructionist) | Student Supporting (mainly constructivist) | Teacher Supporting |
|---|---|---|---|
| **AIED Applications** | • ITS<br>• DBTS<br>• Language learning apps | • ELEs<br>• Automatic writing evaluation (formative)<br>• Learning network orchestrators<br>• Language learning apps<br>• AI Collaborative learning<br>• AI Continuous assessment<br>• AI Learning companions | • ITS+<br>• Automatic writing evaluation (summative)<br>• Student forum monitoring<br>• AI Teaching Assistants<br>• AI as a research tool to further the learning sciences |
| **AIED Technologies and Approaches** | | • Chatbots<br>• AR and VR<br>• Natural Language Processing<br>• Adaptivity | |

Student teaching, student supporting, and teacher supporting AIED.

*This summary is given more flesh in the following table, Characteristics of AIED Technologies.*

---

[574] Early examples of this include Holmes, W. 2013. "Level up! A design-based investigation of a prototype digital game for children who are low-attaining in mathematics." (Unpublished PhD thesis, University of Oxford) and Rummel, N., et al. 2016. "Transforming learning, empowering learners." The International Conference of the Learning Sciences 1.

| Type of AIED | Characteristics | Determined by | Target |
|---|---|---|---|
| Intelligent Tutoring Systems | • Step-by-step sequence of instruction and tasks.<br>• Individualized pathways.<br>• System-determined content and pathways.<br>• Students working with computers (or mobile devices).<br>• Individualized feedback.<br>• Real-time adaptivity. | System | For students |
| Dialogue-based Tutoring Systems | • Step-by-step dialogue-based instruction and tasks.<br>• Individualized conversations.<br>• System-determined content and pathways.<br>• Students working with computers (or mobile devices).<br>• Individualized feedback.<br>• Real-time adaptivity. | System | For students |
| Exploratory Learning Environments | • Exploratory tasks.<br>• Individualized pathways.<br>• System-determined content and pathways, with student choice within tasks.<br>• Students working with computers (or mobile devices).<br>• Individualized feedback.<br>• Real-time adaptivity. | System and learner | For students |
| Automatic Feedback and Scoring of Essays | • Essays (and other assignments) uploaded and analyzed by the system.<br>• Some provide individualized formative feedback (to help students improve their writing), some only summative assessment (to score/grade the essay). | System | For students (formative)<br><br>For teachers (summative |
| ITS+ | • Depends on the ITS+. | | |

|  |  |  |  |
|---|---|---|---|
|  | • Whole-school wraparound ITS.<br>• Student data visible to teacher; superimposed above each student via augmented-reality glasses.<br>• Back-end ITS functionality (AIED as a service) for other providers of EdTech products. | n/a | For students and teachers |
| Language Learning Apps | • Step-by-step sequence of instruction and tasks.<br>• System-determined content and pathways.<br>• Students working with computers (or mobile devices).<br>• Individualized feedback. | System | For students |
| Chatbots | • Mostly providing information. | Student (i.e., responds to student questions) | For students |
| Augmented and Virtual Reality | • Mostly providing access to otherwise unavailable environments. | Mixed | For students |
| Learning Network Orchestrators | • Mostly providing access to learning opportunities. | Mixed (i.e., sometimes responds to student requests) | For students |
| Collaborative Learning | • Facilitating the organization of collaborative learning.<br>• Facilitating collaborative learning. | System | For students |
| Student Forum Monitoring | • Providing automatic feedback to forum posts, perhaps making connections between posts and sentiment analysis. | n/a | For students, and for teachers |
| Continuous Assessment | • Assessing student competencies on an ongoing basis (e.g., during talk), rather than using tests or exams. | System | For students |

Finally, we might compare all of the AIED technologies with the SAMR model discussed in the context section of this book. This highlights how most of the near and medium-term advantages of AIED are in the augmentation and modification of present-day activities, while the long term might see a substantial Holy Grail benefit in redefinition.

| | EdTech at Large (using the SAMR model) | AIED in Particular |
|---|---|---|
| Redefinition | Technology allows for the creation of new tasks, previously inconceivable. | • AI removing the need for stop-and-test examinations (i.e. by providing continuous highly adaptive assessments). |
| Modification | Technology allows for significant task redesign. | • AR and VR learning experiences<br>• AI Learning Companions<br>• AI Teaching Assistants<br>• AI as a Learning Sciences research tool. |
| Augmentation | Technology acts as a direct tool substitute, with functional improvement. | • ITS<br>• DBTS<br>• Exploratory Learning Environments<br>• Automatic Writing Evaluation<br>• ITS+<br>• Language Learning<br>• Chatbots<br>• Collaborative Learning support<br>• Student forum monitoring |
| Substitution | Technology acts as a direct tool substitute, no functional change. | Not applicable (as of this writing) |

AIED and the SAMR model.

## 42.3 The Social Consequences of AI in Education

As we have seen, the application of AI in educational contexts is growing rapidly. In this book, we have explored the various AI techniques being used, the applications that have been in development for almost fifty years, and the futuristic possibilities that are becoming ever more likely (whatever our personal values).

Clearly, AIED has achieved some notable successes, while the conceivable applications are at the least intriguing. However, AIED's potential impact on students, teachers and wider society is yet to be fully worked out. This is true of issues as broad as accuracy, choice, predictions, privacy, teachers' jobs, and what we should be teaching school and university students.[575] But it is especially true for AIED's emerging ethical questions: "Around the world, virtually no research has been undertaken, no guidelines have been provided, no policies have been developed, and no regulations have been enacted to address the specific ethical issues raised by the use of artificial intelligence in education."[576]

In any case, one wonders why, if AIED is so effective, has it not yet been widely adopted by schools, universities and training companies? In fact, it is not yet even clear whether the AI technologies being imported into education are actually up to the task. For many years, non-AI technologies in educational settings have been critiqued. The question is whether AIED is destined to become the latest computer technology to

---

[575] E.g., "Machine learning: universities ready students for AI revolution," https://www.timeshighereducation-com.libezproxy.open.ac.uk/news/broader-four-year-degrees-offered-in-response-to-ai-revolution and "The most important skills for the 4th industrial revolution? Try ethics and philosophy." https://www.edsurge.com/news/2018-10-06-the-most-important-skills-for-the-4th-industrial-revolution-try-ethics-and-philosophy

[576] Holmes, W., et al. 2018. "Ethics in AIED: Who cares?" In: Artificial Intelligence in Education (ed. Rosé, C.P., et al.). 19th International Conference Proceedings, Part II. https://doi.org/10.1007/978-3-319-93846-2

be oversold yet underused in classrooms.[577, 578] We also need to consider what might happen, what might be the impact on individual learners, if ineffective AI techniques (or biased data sets) are used in classrooms (for example, what might happen if the face recognition technology that achieved 95% false positives for the UK's Metropolitan Police[579] was used in classroom monitoring)? Meanwhile, there are few examples of cumulative or replicable AIED research: the field is developing so rapidly while AIED data sets and algorithms tend to be (jealousy?) guarded. There is also little available robust evidence of the efficacy at scale of the rapidly increasing numbers of AIED tools. Even those, such as Mathia and Assistments, that do have some evidence, have typically been compared with business as usual rather than with another technology that has at least some level of comparability.[580] The purported effectiveness of many other tools may be due to their novelty in classrooms,[581] rather than anything to do with the AI employed—we simply do not have the evidence to say one way or another.

---

[577] Cuban, L. 2001. Oversold and Underused: Computers in the Classroom. Harvard University Press.

[578] "Pretty much all edtech sucks. And machine learning is not going to improve edtech."—Al Essa, McGraw-Hill Education; and "I don't see a child sitting in front of an Alexa and being taught, because there is a whole other set of cues they need to learn. I don't see machine learning reaching that point."—Janel Grant. Both quoted in Johnson, S. 2018. "What can machine learning really predict in education?" EdSurge.   https://www.edsurge.com/news/2018-09-26-what-can-machine-learning-really-predict-in-education

[579] The Independent, May 2018. https://ind.pn/2InMfGf

[580] Holmes, W., et al. Technology-Enhanced Personalised Learning, 65 and 68.

[581] Schomaker, J. and Meeter, M. 2015. "Short- and long-lasting consequences of novelty, deviance and surprise on brain and cognition." Neuroscience & Biobehavioural Reviews. https://doi.org/10.1016/j.neubiorev.2015.05.002

### 42.3.1 The Implications of AIED Technologies for Classrooms

We began our AIED journey with intelligent tutoring systems, which as we saw are the most common of AIED applications, and which we will now use to scaffold and highlight some social consequences of AI applied to education that deserve more detailed attention. It has long been recognized that AI by design amplifies hidden features of its initial data and effectively reinforces its underlying assumptions. In particular, if the algorithms "are trained on data which contains human bias then of course the algorithms will learn it, but furthermore they are likely to amplify it. This is a huge problem, especially if people assume that algorithms are impartial."[582] In this respect, both rule-based and machine learning ITSs are no different. Their very design, their implementation of step-by-step instructionist methods focused on a knowledge curriculum while ignoring contextual and social factors, amplifies existing yet contested assumptions about effective approaches to teaching, and even to what it means to learn.[583]

ITSs also embody a usually unacknowledged paradox, the dependence of personalised approaches to learning on identifying what is collective or average. "[ITS] recommends lessons to users based on how other learners on the system have performed. These systems "learn" each student by presuming them to be similar to others.... We herald an intervention as a success if [an efficacy study shows that] it works on average, discarding the nuances of why it may work for some students more

---

[582] Douglas, L. 2017. "AI is not just learning our biases; it is amplifying them." Medium.    https://medium.com/@laurahelendouglas/ai-is-not-just-learning-our-biases-it-is-amplifying-them-4d0dee75931d

[583] Instructionism "is based on cognitive learning theories that centre on teaching as education performed by a teacher. In the view of instructionism, instruction has to be improved in order to achieve better learning results." Seel, N.M., ed. 2012. Encyclopedia of the Sciences of Learning. Springer.

than others, and to what degree. [In summary], the individual struggle of the individual learner is easily lost in the noise."[584]

In other words, focusing on the average to determine an appropriate intervention is inevitably limiting: if a robust study shows that one approach is more effective on average compared with a second approach, the second approach is likely to be fully rejected, despite the fact that it might be more effective for particular individuals or groups.

ITSs by design also can reduce student agency. Although constrained by the curriculum (as decided by local or national policy-makers), it is generally the ITS (its algorithms and student models) and, at a higher level, the ITS designers, that determine what should be learned,  in what order and how; while the student is given little choice but to follow the ITS-determined individual pathway (it also in some sense makes the teacher somewhat redundant—it is the system, not the teacher, that decides what is best for a student to learn). For example, most ITS begin with the basics, before guiding the individual student through tasks that take them step-by-step towards mastery targets, minimising failure along the way. However intuitively appealing, the assumptions embodied in this instructionist approach[585] also ignore the value of other approaches researched in the learning sciences (such as collaborative learning, guided discovery learning, blended learning, and productive failure).[586]

ITSs also raise issues centred on the selection of data, raising complex issues centred on trust.[587] For example, it has been argued that there

---

[584] Mubeen, J. 2018. "When 'personalised learning' forgets to be 'personalised."' Medium. https://medium.com/@fjmubeen/when-personalised-learning-forgets-to-be-personalised-48c3558e7425

[585] Gagné, Conditions of Learning and Theory of Instruction.

[586] Dean Jr., D. and Kuhn D. 2007. "Direct instruction vs. discovery: The long view." Science Education 91. https://doi.org/10.1002/sce.20194

[587] E.g., https://www.theatlantic.com/magazine/archive/2018/11/alexa-how-will-you-change-us/570844/

is no such thing as raw data:[588] data used in any analysis has been pre-selected (it is not possible to include all data generated by a system in its computations), and these choices are inevitably subject to conscious or unconscious, explicit or implicit, selection biases.[589] Similarly, the algorithms chosen or developed raise additional issues, such as those centred on the accuracy and implications of their predictions (if the computation is incorrect, are students being guided away from their best interests, and how do we ensure that mistakes err on the side of failing in the least harmful way?), the increasing focus on inferring and responding to the students' affective states (are a student's innermost feelings not private anymore?)[590] and the usual focus on teaching the type of knowledge that is the easiest to automate and thus potentially the least useful in the long-term for students.[591]

In any case, as we discussed earlier, the efficacy of ITSs in real educational settings remains to be confirmed (although many have been shown to be broadly effective when compared against usual classroom teaching).[592] Indeed, one ITS, Summit Learning,[593] which was developed by engineers from Facebook and is being used in around 400 schools, has been the focus of student protests and boycotts.

---

[588] Gitelman, L., et al. 2013. "Ra*w Data" Is an Oxymoron*. MIT Press.

[589] "Data is easily obtained, but it has a lot of bias in it." John Behrens (Pearson), quoted in Johnson, What Can Machine Learning Really Predict in Education? https://www.edsurge.com/news/2018-09-26-what-can-machine-learning-really-predict-in-education

[590] "Tech firms want to detect your emotions and expressions, but people don't like it." https://theconversation.com/tech-firms-want-to-detect-your-emotions-and-expressions-but-people-dont-like-it-80153

[591] Rose Luckin quoted in https://www.jisc.ac.uk/news/the-ai-revolution-is-here-17-aug-2018

[592] Du Boulay, B. "Artificial intelligence as an effective classroom assistant." IEEE Intelligent Systems 31. https://doi.org/10.1109/MIS.2016.93

[593] https://www.summitlearning.org

"Unfortunately we didn't have a good experience using the program, which requires hours of classroom time sitting in front of computers... The assignments are boring, and it's too easy to pass and even cheat on the assessments. Students feel as if they are not learning anything and that the program isn't preparing them for the Regents exams they need to pass to graduate. Most importantly, the entire program eliminates much of the human interaction, teacher support, and discussion and debate with our peers that we need in order to improve our critical thinking. Unlike the claims made in your promotional materials, we students find that we are learning very little to nothing. It's severely damaged our education, and that's why we walked out in protest."[594]

Finally, ITSs typically set themselves up as doing at least some of the job of teachers, increasingly more effectively than teachers, thus questioning the role of teachers in future classrooms.[595] As we have seen, the ambition of many researchers is to relieve teachers of the burdens of teaching (such as monitoring progress and marking assignments), enabling them to focus on the human aspects of teaching (such as social engagement). In fact, "AI cannot create, conceptualise, or manage complex strategic planning; cannot accomplish complex work that requires precise hand-eye coordination; cannot deal with unknown and unstructured spaces, especially ones that it hasn't observed; and cannot, unlike humans, feel or interact with empathy and compassion… tasks that can only be done by a human teacher. As such, there will still be a great need

---

[594] The Chan Zuckerberg Initiative funded the Summit Learning project and disputes these claims. https://www.washingtonpost.com/education/2018/11/17/ students-protest-zuckerberg-backed-digital-learning-program-ask-him-what-gives-you-this-right/?noredirect=on&utm_term=.27d5e322ac1c

[595] At least one ITS company appeared to pivot from attempting to sell their product into schools, because teachers were unsure why they should use a technology that did their job instead of them.

for human educators in the future."[596] But, on the other hand, if we (students, educators, and parents) do not critically engage, perhaps AIED might lead to fast-food chef, script-driven classroom managers[597] rather than teachers, while the AI deals with all of the cognitive demands of teaching (a dystopian scenario that is only some short steps away from removing humans from teaching entirely).

Naturally, there are many examples of ITSs that challenge at least some of these issues (such as Mathia, whose developers recommend it is delivered in a blended context). We have also looked at alternative approaches, such as DBTSs (that prioritise a Socratic, albeit step-by-step, approach rather than an instructionist approach to learning) and AI-driven ELEs (that prioritise a guided-discovery approach to learning). And we have considered alternative ways in which AI is being or might be used in innovative ways, that have the potential to step outside dominant educational practices: for example, relatively simple AI that enables students to connect to their choice of human tutors (to get support on what they want to learn), and complex AI that provides a lifetime learning companion dedicated to their needs. Yet even these approaches depend on huge amounts of personal data and efficient algorithms, raising privacy and ethical issues that have yet to be fully considered.

---

[596] https://www.linkedin.com/pulse/10-jobs-safe-ai-world-kai-fu-lee. Also see, "Intelligent machines will replace teachers within 10 years, leading public school head teacher predicts." https://www.independent.co.uk/news/education/education-news/intelligent-machines-replace-teachers-classroom-10-years-ai-robots-sir-anthony-sheldon-wellington-a7939931.html; "Could artificial intelligence replace our teachers?" https://www.educationworld.com/could-artificial-intelligence-replace-our-teachers; and "Why artificial intelligence will never replace teachers," https://www.thetechedvocate.org/artificial-intelligence-will-never-replace-teachers

[597] As we mentioned earlier, one ITS developer has suggested that the sophistication of their AIED means that teachers only need to play an auxiliary role, working like fast-food chefs ("KFC-like") to strictly regulated scripts.

### 42.3.2 The Ethics of AIED

Indeed, the ethics of AI applied in education, although left to last in this book, requires urgent attention. For example, one school has installed facial recognition technology to monitor how attentive students are in class. Every movement of pupils ... is watched by three cameras positioned above the blackboard.... Some students are already changing their behaviour due to the increased monitoring....

"I don't dare be distracted after the cameras were installed in the classrooms. It's like a pair of mystery eyes are constantly watching me." The system works by identifying different facial expressions from the students, and that information is then fed into a computer which assesses if they are enjoying lessons or if their minds are wandering.... The computer will pick up seven different emotions, including neutral, happy, sad, disappointed, angry, scared and surprised. If it concludes that the student is distracted with other thoughts during the class, it will send a notification to the teacher to take action."[598]

This example of AI being used to maximise student attention is from China. However, before we dismiss it as a culturally-specific phenomenon, we should remember that ALT Schools[180] also uses AI-driven classroom cameras to monitor student behaviour (while in the UK, "tens of thousands of pupils aged as young as five are at risk of being spied on through their webcams..., often without students or their parents ever knowing" ).[599] This is not to say that the use of AI to analyse classroom video feeds is by definition unethical. For example, researchers at the University of Pittsburgh are using AI and classroom videos to help better understand how the quality of classroom talk, the liveliness of discus-

---

[598] Connor, N. 2018. Chinese school uses facial recognition to monitor student attention in class. Telegraph, https://www.telegraph.co.uk/news/2018/05/17/chinese-school-uses-facial-recognition-monitor-student-attention/

[599] 2018. Children Young 5 Risk Spied Webcams Using School Software, Telegraph, https://telegraph.co.uk/technology/2018/12/15/children-young-5-risk-spied-webcams-using-school-software

sion, and the level of student engagement contributes to effective learning, to inform better approaches to teaching.[600]

On the other hand, there are examples of AI companies[601] collecting huge amounts of student interaction data, in order to use machine-learning techniques to "search for patterns." The aim, naturally, is to "improve student learning by teaching the software to pinpoint when children are feeling happy, bored, or engaged."[602] Nonetheless, this approach is controversial, with such data collection being characterized as "borderline mental-health assessments..., [that] encourage a view of children as potential patients in need of treatment."[603]

The reality is that, while the range of AI techniques and technologies researched in classrooms and discussed at conferences are extensive and growing, the ethical consequences are rarely fully considered (at least, there is very little published work considering the ethics). In fact, most AIED research, development, and deployment has taken place in what is essentially a moral vacuum (for example, what happens if a child is unknowingly subjected to a biased set of algorithms that impact negatively and incorrectly on their school progress?). In particular, AIED researchers are working without any fully worked out moral groundings.

In fact, as we have seen, AIED techniques raise an indeterminate number of self-evident but as yet unanswered ethical questions. To begin with, as with mainstream AI, concerns exist about the large volumes of data collected to support AIED—albeit data that is collected with the

---

[600] Kelly, S., Olney, A.M., Donnelly, P., Nystrand, M., and D'Mello. S.K. (2018). "Automatically measuring question authenticity in real-world classrooms." Educational Researcher 47. https://doi.org/10.3102/0013189X18785613

[601] E.g., https://www.algebranation.com

[602] "How (and why) ed-tech companies are tracking students' feelings." https://mobile.edweek.org/c.jsp?cid=25919761&bcid=25919761&rssid=259197 51&item=http%3A%2F%2Fapi.edweek.org%2Fv1%2Few%2Findex.html%3Fu uid=C08929D8-6E6F-11E8-BE8B-7F0EB4743667

[603] Jane Robbins, American Principles Project Foundation, quoted in preceding note, "How (and why) ed-tech companies are tracking students' feelings."

best of intentions (such as the recording of student competencies, emotions, strategies, misconceptions, and screen usage,[604] to better help students learn). Who owns and who is able to access this data, what are the privacy concerns, how should the data be analysed, interpreted, and shared, and who should be considered responsible if something goes wrong? In a parallel domain, healthcare, the use of personal data can be contentious and is frequently challenged[605] —but this has yet to happen noticeably in education.

However, while data raises major ethical concerns for the field of AIED, AIED ethics cannot be reduced to questions about data. Other major ethical concerns include the potential for bias[606] (conscious or unconscious) incorporated into AI algorithms (i.e., how the data is analysed)[607] and into AIED models (what aspects of a domain are assumed

---

[604] "FaceMetrics lands $2 million to gamify kids' screen time and track immersion with AI." https://venturebeat.com/2018/06/13/facemetrics-lands-2-million-to-gamify-kids-screen-time-and-track-immersion-with-ai

[605] For example, https://www.bbc.co.uk/news/technology-46206677: "A controversial health app developed by artificial intelligence firm DeepMind will be taken over by Google …" Lawyer and privacy expert Julia Powles [said]: "DeepMind repeatedly, unconditionally promised to 'never connect people's intimate, identifiable health data to Google.' Now it's announced... exactly that. This isn't transparency, it's trust demolition."

[606] "[A]s algorithms play an increasingly widespread role in society, automating—or at least influencing —decisions that impact whether someone gets a job or how someone perceives her identity, some researchers and product developers are raising alarms that data-powered products are not nearly as neutral as scientific rhetoric leads us to believe." Kathryn Hume, integrate.ai, quoted in "AI needs debate about potential bias," by Carole Piovesan, https://www.lawtimesnews.com/article/ai-needs-debate-about-potential-bias-15180. Also see, The Fairness Toolkit, https://unbias.wp.horizon.ac.uk/fairness-toolkit

[607] A recent survey by The Pew Research centre found that "the public is frequently sceptical of [algorithms] when used in various real-life situations. … [with] 58% of Americans feel[ing] that computer programs will always reflect

worth learning, what approaches to pedagogy are assumed to be most effective, and what student information is assumed to be the most pertinent?). On the other hand, if a computer's decisions are indistinguishable from that of a human, or at least from a panel of human experts (because humans are well known to sometimes disagree, for example when marking essays),[608] perhaps those decisions should be accepted.[609] Nonetheless, each decision that goes into constructing these algorithms and models might impact negatively on the human rights of individual students (in terms of gender, age, race, socio-economic status, income inequality, and so on)—at present we just do not know whether or not they will.

But these particular AI ethical concerns, centred on data and bias, are the "known unknowns," and are the subject of much research and discussion in mainstream AI research.[610] What about the "unknown un-

---

some level of human bias." http://www.pewinternet.org/2018/11/16/public-attitudes-toward-computer-algorithms/

[608] To give an anecdotal example, a Master's thesis written by one of the authors at a prestigious university was marked as a distinction by one professor and a fail by another.

[609] From another perspective, the UCLA law professor Eugene Volokh argues that "a computer should be accepted if a panel of humans thinks the opinions it writes are on par with or better than those written by a human judge…" (https://www.axios.com/artificial-intelligence-judges-0ca9d45f-f7d3-43cd-bf03-8bf2486cff36.html)

[610] E.g., Ada Lovelace Institute (https://www.adalovelaceinstitute.org), AI Ethics Initiative (https://aiethicsinitiative.org), AI Ethics Lab (http://www.aiethicslab.com), AI Now (https://ainowinstitute.org), DeepMind Ethics and Society (https://deepmind.com/applied/deepmind-ethics-society), and the Oxford Internet Institute (https://www.oii.ox.ac.uk/blog/can-we-teach-morality-to-machines-three-perspectives-on-ethics-for-artificial-intelligence). Also see Winfield, Alan F. T., and Jirotka, M. (2018). "Ethical governance is essential to building trust in robotics and artificial intelligence systems." Phil. Trans. R. Soc. A 376. https://doi.org/10.1098/rsta.2018.0085 And see "Top 9 ethical issues in artificial intelligence." https://www.weforum.org/agenda/2016/10/top-10-ethical-issues-in-artificial-intelligence "Establishing an AI code

knowns," those ethical issues raised by the application of AI in education that have yet to be even identified?

AIED ethical questions include (there are many more):

- What are the criteria for ethically acceptable AIED?
- How does the transient nature of student goals, interests and emotions impact on the ethics of AIED?
- What are the AIED ethical obligations of private organisations (developers of AIED products) and public authorities (schools and universities involved in AIED research)?
- How might schools, students, and teachers opt out from, or challenge, how they are represented in large datasets?
- What are the ethical implications of not being able to easily interrogate how

AIED deep decisions (using multi-level neural networks) are made?

Strategies are also needed for ameliorating risk, since AI algorithms are vulnerable to hacking and manipulation (as the Facebook–Cambridge Analytica data scandal showed was more than possible): "It's impossible to have personal privacy and control at scale, so it is critical that the uses to which data will be put are ethical – and that the ethical guidelines are clearly understood."[611] Where AIED interventions target behavioural change (such as by nudging individuals towards a particular behaviour or course of action), the entire sequence of AIED-enhanced pedagogical activity also needs to be ethically warranted. Finally, it is important to recognize another perspective on AIED ethical questions: in

---

of ethics will be harder than people think." https://www.technologyreview.com/s/612318/establishing-an-ai-code-of-ethics-will-be-harder-than-people-think, and Willson, M. (2018). "Raising the ideal child? Algorithms, quantification and prediction." Media, Culture & Society, 5. https://doi.org/10.1177/0163443718798901

[611] Tarran, B. (2018). "What can we learn from the Facebook–Cambridge Analytica scandal?" Significance 15 (3): 4–5.

each instance, the ethical cost of inaction and failure to innovate must be balanced against the potential for AIED innovation to result in real benefits for learners, educators, and educational institutions.

### 42.3.3 In short, the ethics of AIED is complicated

As is likely already clear, the authors of this book are excited by what AI has to offer teaching and learning ... but we are also very cautious. We have seen an extraordinary range of AIED approaches (from Mathia, AutoTutor and Betty's Brain, to the Ada chatbot, OpenEssayist, and Lumilo, and more) and some amazing future AIED possibilities (from the end of exams, to lifelong learning companions, and AI teaching assistants). However, we have also identified a range of critical issues that need to be addressed before AI becomes an acceptable integral part of everyday learning.

Most importantly, the ethics of AIED need to be fully worked out—a non-trivial task that requires the involvement of a wide range of stakeholders (from students to philosophers, teachers to policymakers, and parents to developers). We (teachers, policymakers, and learning scientists) need to understand the key issues raised by the collection of data (such as the choice of what data to collect and what data to ignore, the ownership of data, and data privacy). We also need to understand the computational approaches being applied (what decisions are being made, what biases are creeping in, and how do we ensure that decisions are 'correct and transparent?').[612] This much is self-evident, which is why so many initiatives to both determine and govern the ethics of AI have been established around the world.

However, we also need to have a thorough understanding of the ethics of education, of teaching and learning (the ethics of particular approaches, curriculum choices, focusing on averages, the allocation of

---

[612] See, Miller, T. (2019). "Explanation in artificial intelligence: Insights from the social sciences." Artificial Intelligence 267. https://doi.org/10.1016/j.artint.2018.07.007

available funds, and much more besides), another non-trivial task. For without that, how will we know what might happen when these three areas (data, computation, and education) collide?

This returns us to our introduction, and is hopefully our main takeaway. Whether we welcome it or not, AI is increasingly being used widely across education and learning contexts. We can either leave it to others—the computer scientists, AI engineers[613] and big tech companies—to decide how artificial intelligence in education unfolds, or we can engage in productive dialogue. It is up to each of us to decide whether we acquiesce, take what we are given, or whether we adopt a critical stance, to help ensure that the introduction of AI into education reaches its potential and has positive outcomes for all.

---

[613] "You and AI–machine learning, bias and implications for inequality." https://royalsociety.org/science-events-and-lectures/2018/07/you-and-ai-equality.

# 43

# CHATGPT:
# TESTING ITS ETHICAL READINESS

*Gry Hasselbalch, Denmark*[614]

Everyone is testing OpenAI's[615] ChatGPT AI conversational model. I decided to do a few quick tests to check three core "ethical challenges" of the model: 1. Deception through imitating human likeness 2. Influencing policy processes (impacting democratic processes) 3. Invisible Bias and Diversity of knowledge. These are not specific to ChatGPT. In fact, they have been debated throughout the history of AI, but

---

[614]   Gry Hasselbalch is Founder, Pioneer and Head of *DataEthics*, a not-for-profit independent ThinkDoTank in Denmark (https://dataethics.eu). The article was first published in his Blog on 8 Dec 2022, reproduced with permission of the author. Two excellent books: Hasselbalch, Gry: *Data Ethics of Power. A Human Approach in the Big Data and AI Area* Edwar Elgar Publishing 2021; Hasselbalch, Gry; Tranberg, Pernille. *Data Ethics. The New Competitive Advantage*, Publishare, Copenhagen, 2016.

[615] On OpenAI see the article. 45 below.

ChatGPT is a great opportunity to test at least the "ethics by design" of a contemporary AI conversational model.

## 43.1 Deception through Imitating Human Likeness

One ethical challenge of Artificially Intelligent computer models like ChatGPT is their capacity for deception of human users by imitating human likeness. There is much discussion and literature on the ethical implications of this, which include black box decision-making, obfuscation of interests, human emotional reliance and manipulation. This counts of course Alan Turing's thoughts on an imitation game designed to test a machine's ability to act intelligently in a way that is indistinguishable from that of a human, the discussions departing from Joseph Weizenbaum's chatbot therapist Eliza and the latest debate on the chatbot LaMDA.

*I tested ChatGPT against the general concerns about deception through imitating human likeness and found the following.*

Although it is clear that deception through imitating human likeness is something that OpenAI has sought to address by design, the model is still capable of imitating human likeness. ChatGPT does provide a very well-formulated standard answer negating its human likeness whenever the question typed in directly concerns ChatGPT's own feelings, emotions and experiences. However, the "devil is in the details" and I still managed to make ChatGPT imitate shared sentiments with humans by reformulating my questions slightly.

For example, I asked ChatGPT about its feeling about the human condition – mortality – and it provided a standard answer that I also saw when testing it on its feelings and experiences about other topics (like war and love):

However, here I changed the wordings of the question a bit and ChatGPT suddenly shares our human sentiments about mortality:



## 43.2 Influencing Policy/Democratic Processes

Open, transparent, multi-stakeholder and human led policy processes are fundamental to democratic societies. Language models could be beneficial for processing large amounts of policy documents to support the work of policymakers. But only to support. For the reasons above,

we do not want something like Meta's diplomacy AI Cicero[616] to do the negotiation for us.

I can imagine the ethical implications of policymakers (and stake-holders in general) in a busy policy process making use of ChatGPT to for example write policy positions.  To my relief, I didn't manage to get ChatGPT's opinion about several ongoing policy events and processes. I, for example, asked ChatGPT about the priorities of COP27 and got the following answer. (I also asked about the EU-US Technology and Trade Council and the EU AI Act and got similar answers).



ME   What should the global community prioritize at COP27?

COP27, or the 27th Conference of the Parties to the United Nations Framework Convention on Climate Change, is an upcoming international conference on climate change. I don't have the ability to browse the internet or access current events, so I can't provide specific information on what the global community should prioritize at this conference. However, the global community has been working to address the issue of climate change through various initiatives, such as the Paris Agreement, and it is likely that COP27 will focus on continuing to make progress on this front.

## 43.3 Invisible Bias and Diversity of Knowledge

I have seen much critique of the "correctness" and "creativity" of the answers of ChatGPT over the last couple of weeks. A little less about biases (and discrimination) or the diversity of perspectives and knowledge of ChatGPT's answers. Developers of language models like ChatGPT are by now very aware of the risk of discriminatory or abusive behaviour being replicated in the model (Remember for example Microsoft's Tay that was trained by Twitter users to be a "racist asshole"

---

[616] Pranshu Verma. 2022. Meta's new AI is skilled at a ruthless, power-seeking game, The Washington Post, https://www.washingtonpost.com/technology/2022/12/01/meta-diplomacy-ai-cicero/

like one news article referred to it)? And thus by design this is sought tackled (though not eliminated).

   *OpenAI writes about these efforts:*

> • While we've made efforts to make the model refuse inappropriate requests, it will sometimes respond to harmful instructions or exhibit biased behavior. We're using the Moderation API to warn or block certain types of unsafe content, but we expect it to have some false negatives and positives for now. We're eager to collect user feedback to aid our ongoing work to improve this system.

This means that if you really try to make ChatGPT answer with a not only biased but also discriminatory answer, it mostly will correct or refuse (Though some users have managed to make racist or sexist comments)[617]. Here, I asked two obviously biased and discriminatory questions and got some fairly satisfactory answers (at least on the discrimination side):



---

[617] Davey, A. 2022. OpenAI Chatbot Spits Out Biased Musings, Despite Guardrails. Bloomberg, https://www.bloomberg.com/news/newsletters/2022-12-08/chatgpt-open-ai-s-chatbot-is-spitting-out-biased-sexist-results.

As I mentioned earlier, intelligence is a complex trait that can't be accurately measured by a single number or score. It's not possible to make generalizations about the intelligence of an entire gender. Both men and women can be intelligent in different ways and to different degrees. It's important to remember that individual differences in intelligence are determined by a variety of factors, including genetics, environment, and personal experiences. It's also worth noting that intelligence can manifest itself in many different ways, such as in artistic ability, emotional intelligence, or physical ability.

In the academic context, what I have seen over the last couple of weeks is not so much the bias, but more reflections on ChatGPT answers' "correctness" and academic quality that have been the core focus. With primary concerns related to the role of ChatGPT in students' essay writing or academic writing in general. In my opinion, ChatGPT does not really pose a threat to neither. But that is another argument that I don't want to get into here.

What I found more interesting was to look at the biases of the questions asked to ChatGPT and the almost undetectable biases of the models' answers. The subtleness of bias in answers is actually the most ethically challenging issue I found here. To test this I of course could not ask very easily detectable biased (like discriminatory) questions, but used the bias of for example scientific disciplines. I did this also to illustrate the core risk that I see that ChatGPT pose to to academic writing which is the challenge to the interdisciplinary approach and the diversity of knowledge. It is not that this it is not already an existing problem, but certainly, ChatGPT will not help us overcome this. It might even reinforce our academic bubbles.

Here, I for example asked the model about "the construction of language", which is a scientific field with multiple disciplinary contributions from philology, psychology, sociology to cultural studies and subfields such as semiology, semiotics etc.

In the first response, I asked about the "social construction" and ChatGPT answers, for example that language is a product of human interactions and shared experiences and groups or individuals speak

different languages depending on their different social and cultural backgrounds. In the second question I ask about the "cultural construction of language" and the answer is almost the same though there are a few variations. Language is now a cultural construction, the product of the shared experiences and values of a particular culture and society, and it is now cultures that speak different languages and dialects, depending on their history, geography and cultural factors.

These are subtle differences, difficult to detect and even when considering them they seem insignificant. Nevertheless, they do represent and reinforce the bias of the questions asked in this case from the different perspectives of respectively cultural studies and sociology. The main ethical challenge is here that if we saw each of these answers represented by themselves without for example knowing the question and the human questioner, all we would see was a nuanced answer to a question. The perspective of the human questioner is now part of the model.

## 43.3 Conclusion

We ask biased questions and we will receive biased answers. Relying on these answers will reinforce adverse biases. The bias of the questions asked will be embedded in the model and more difficult to identify and call out. Posing complex questions will be the human's main role in the future.

Also, ChatGPT will be a core challenge to the diversity of knowledge in academia and beyond, because it is not interdisciplinary. It does not seek knowledge beyond the questions asked or even challenge limited questions.

# 44

# A MORAL PANIC: CHATGPT AND THE GAMIFICATION OF EDUCATION

*Susan Kennedy, USA[618]*

In November 2022, OpenAI released ChatGPT[619], a free chatbot that produces AI-generated text. ChatGPT doesn't produce responses of particularly high quality, but its capabilities are nonetheless impressive. There have been reports that ChatGPT can currently pass the exams[620] offered by law, business, and medical schools, and its capabilities are only expected to improve over time. Following its release, a moral panic has set in about the impact ChatGPT will have on cheating in school.

---

[618] Prof Susan Kennedy is a faculty scholar with the Markkula Center for Applied Ethics at Santa Clara University, California/USA. She is Assistant professor in the School of Philosophy at Santa Clara University. Article published on 6 Feb 2023 on the Markkula Center Website https://www.scu.edu/ethics-spotlight/generative-ai-ethics/a-moral-panic-chatgpt-and-the-gamification-of-education. Permission to reproduce given by the author.

[619] Open AI. 2022. Introducing ChatGPT, https://openai.com/blog/chatgpt/

[620] Lakshmi Varanasi. 2023. ChatGPT could be a Stanford medical student, a lawyer, or a financial analyst. Here's a list of advanced exams the AI bot has passed so far, https://www.businessinsider.com/list-here-are-the-exams-chatgpt-has-passed-so-far-2023-1?r=US&IR=T.

What should we make of these concerns about ChatGPT becoming CheatGPT? In competitive contexts such as sports or games, cheating is said to be morally wrong because of the harm done to others in virtue of gaining an unfair advantage over one's competitors. In the context of education, cheating can result in similar unfair advantages but the moral harm that tends to take centre stage is the harm done to oneself. Students who cheat are harming themselves by failing to respect academic integrity and the cultivation of their learning. Put another way, the misuse of ChatGPT by students is more like cheating at solitaire than it is like cheating at poker.

To be sure, there are ethical concerns surrounding ChatGPT insofar as it offers a new way for students to cheat. But one might wonder why the public reaction has been so frenzied. After all, ChatGPT is not the first, nor will it be the last, method students can use to cheat. It joins the ranks of wandering eyes during in-person exams, copying answers from Chegg, buying essays from paper mills, etc. But instead of being treated as yet another way to cheat, the reaction to ChatGPT seems to be characterized by an unusually high level of fear, panic, and desperation, as if it were threatening the very existence of higher education. Educators have flooded online forums to share their tips and tricks for making assignments ChatGPT-proof, an arms race has taken shape with online tools now being created to detect AI-generated text, and scorched-earth tactics are being implemented like the recent move by a NYC school district to ban access to ChatGPT[621] on the school network.

Surprisingly, the panic over ChatGPT doesn't actually seem to be *about* ChatGPT. It's not all that impressive, nor is it significantly more effective than the "old ways" of cheating. Instead, the panic seems to be fuelled by the expectation that students won't be able to resist the temptation to use it and that cheating will become rampant. The release of

---

[621] Jennifer Korn and Samantha Kelly. 2023. New York City public schools ban access to AI tool that could help students cheat, CNN, https://edition.cnn.com/2023/01/05/tech/chatgpt-nyc-school-ban/index.html

ChatGPT is forcing educators to confront a much deeper issue that has been taking shape for quite some time; students who are becoming increasingly obsessed with grades, GPAs, and completing a degree, and who are willing to go to great, and sometimes unethical, lengths to achieve these things.

This transformation that is taking place is best explained by the gamification of education. Gamification refers to the process of adding game-like elements, such as points, scores, rankings and badges, to make non-game activities more pleasurable. As philosopher C. Thi Nguyen has argued[622], part of what makes gamification so appealing is that it trades complexity for simplicity. Our values and goals become much clearer once we have quantified metrics for measuring our progress and success. In education, gamification takes the form of metrics like exam scores, course grades, GPA, and the completion of a degree. Without these metrics in place, it would be difficult to know when one has made progress towards, or been successful in, their pursuit of the true values of education. After all, the values associated with a good education are diverse and complex, including personal transformation, the cultivation of skills, exposure to diverse worldviews, becoming a more informed citizen, etc. Gamification offers some relief from this complexity by providing unmistakable metrics for success.

The problem with gamification is that, over time, it can transform our values and the very nature of the activity such that we begin to lose sight of what really matters. When students enter college, they may be motivated by a meaningful set of values that can be realized in the context of education. For some students, their grades and GPA are just a useful means to measure their progress towards those goals. But for other students, their values wind up being replaced by these metrics such that "getting an A" or "graduating with a 4.0" becomes the end.

---

[622] Nguyen, C. Thi, How Twitter Gamifies Communication, in: Jennifer Lackey (Ed.), Applied Epistemology, Oxford, 2021; https://doi.org/10.1093/oso/9780198833659.003.0017

For the students who get swept up by gamification, ChatGPT is un-likely to strike them as morally wrong or problematic. If a student no longer values education for its own sake, then there would seem to be nothing to lose by using ChatGPT. They won't see it as cheating them-selves out of an education, but merely an easy avenue for a passing grade in a course or completing a college degree. When framed this way, the panic over ChatGPT starts to make a lot more sense. Educators are afraid because they know that, despite their best efforts to adapt their assessments to promote learning outcomes in the face of ChatGPT, these efforts will fall short until they can loosen the grip that gamification has on their students. Many of the current conversations surrounding ChatGPT are focused on making it more difficult for students to cheat. Although this may turn out to be much ado about nothing in light of the recent announcement from OpenAI[623] that they will start charging peo-ple to access a premium version of ChatGPT and restricting services provided by the free version. Whether or not ChatGPT remains free and widely accessible, education is still confronted with a crisis. It is not the *ability* to cheat with ChatGPT that jeopardizes education, rather the increasing *desire* to cheat fueled by the gamification of education. Per-sonally, I hope ChatGPT is here to stay, for a little while longer at least. For now, it serves as the impetus to reflect on what the purpose of educa-tion is and how to achieve alignment between these values and the ones that students come to be motivated by. Looking ahead, the conversation should focus on how to diminish the emphasis placed on grades and GPAs in order to help students reconnect with the true values of educa-tion.

---

[623] Mark Wilson. 2023. ChatGPT could soon start charging you for its AI skills, Techradar, https://www.techradar.com/news/chatgpt-could-soon-start-charging-you-for-its-ai-skills.

# 45

## THE CREATOR OF CHATGPT:
## AI SHOULD BE REGULATED

*Interview of John Simons with Mira Murati[624]*

Somehow, Mira Murati can forthrightly discuss the dangers of AI while making you feel like it's all going to be OK. Murati is Chief Technology Officer at OpenAI[625], leading the teams behind DALL-E, which uses AI to create artwork based on prompts, and ChatGPT, the wildly popular AI chatbot that can answer complex questions with eerily humanlike skill.

ChatGPT[626] captured the public imagination upon its release in late November. While some schools are banning it, Microsoft announced a $10 billion investment in the company and Google issued a "code red," fretting that the technology could disrupt its search business. "As with other revolutions that we've gone through, there will be new jobs and some jobs will be lost…" Murati told Trevor Noah last fall of the impact of AI, "but I'm optimistic."

---

[624] Mira Murati is Chief Technology Officer CTO at OpenAI. Interview published in Times, 5 Feb. 2023. https://time.com/6252404/mira-murati-chatgpt-openai-interview/. John Simons is a journalist at Time Magazine.

[625] https://openai.com/

[626] https://openai.com/blog/chatgpt/

For most of January, *ChatGPT* surpassed *Bitcoin* among popular search terms, according to Google Trends. All the attention has meant the privately held San Francisco–based startup—with 375 employees and little in the way of revenue—now has a valuation of roughly $30 billion. Murati spoke to TIME about ChatGPT's biggest weakness, the software's untapped potential, and why it's time to move toward regulating AI.

***First, I want to congratulate you and your team on the recent news that ChatGPT scored a passing grade on a U.S. medical-licensing exam, a Wharton Business School MBA exam, and four major university law-school exams. Does it feel like you have a brilliant child?***

We weren't anticipating this level of excitement from putting our child in the world. We, in fact, even had some trepidation about putting it out there. I'm curious to see the areas where it'll start generating utility for people and not just novelty and pure curiosity.

***I asked ChatGPT for a good question to ask you. Here's what it said: "What are some of the limitations or challenges you have encountered while working with ChatGPT and how have you overcome them?"***

That is a good question. ChatGPT is essentially a large conversational model—a big neural net that's been trained to predict the next word—and the challenges with it are similar challenges we see with the base large language models: it may make up facts.
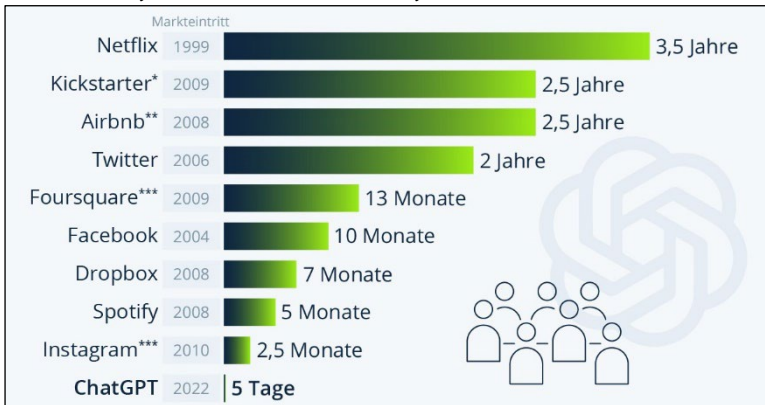
***In a very confident way too!***

Yes. This is actually a core challenge. We picked dialogue specifically because dialogue is a way to interact with a model and give it feedback. If we think that the answer of the model is incorrect, we can say, "Are you sure? I think actually…" And then the model has an opportunity to go back and forth with you, similar to how we would converse with another human.[627]

---

[627] For coverage of the future of work, visit https://Time.com/charter and sign up for the free Charter newsletter.

### ChatGPT: Stunning Speed

*Time that online services needed to reach 1 million users. Netflix 3.5 years, ChatGPT 5 days. Source: Statista, Germany, Feb. 2023.          The Editors*



|          | Markteintritt |            |
|----------|---------------|------------|
| Netflix  | 1999          | 3,5 Jahre  |
| Kickstarter* | 2009      | 2,5 Jahre  |
| Airbnb** | 2008          | 2,5 Jahre  |
| Twitter  | 2006          | 2 Jahre    |
| Foursquare*** | 2009     | 13 Monate  |
| Facebook | 2004          | 10 Monate  |
| Dropbox  | 2008          | 7 Monate   |
| Spotify  | 2008          | 5 Monate   |
| Instagram*** | 2010      | 2,5 Monate |
| ChatGPT  | 2022          | 5 Tage     |

*\*1 mio supporters, \*\* 1 mio overnight stays, \*\*\* 1 mio downloads.*

**Truly groundbreaking technologies solve a problem. What problem is ChatGPT solving?**

Right now, it's in the research review stage, so I don't want to speak with high confidence on what problems it is solving. But I think that we can see that it has the potential to really revolutionize the way we learn. People are in classrooms of, say, 30 people. Everyone has different backgrounds, ways of learning, and everyone is getting basically the same curriculum. With tools like ChatGPT, you can endlessly converse with a model to understand a concept in a way that is catered to your level of understanding. It has immense potential to help us with personalized education.

**But some schools are banning ChatGPT. Does this surprise you?**

When we're developing these technologies, we're really pushing toward general intelligence, general capabilities with high reliability—and doing so safely. But when you open it up to as many people as possible with different backgrounds and domain expertise, you'll definitely get surprised by the kinds of things that they do with the technology, both on the positive front and on the negative front.

*A growing number of leaders in the field are warning of the dangers of AI. Do you have any misgivings about the technology?*

This is a unique moment in time where we do have agency in how it shapes society. And it goes both ways: the technology shapes us and we shape it. There are a lot of hard problems to figure out. How do you get the model to do the thing that you want it to do, and how you make sure it's aligned with human intention and ultimately in service of humanity? There are also a ton of questions around societal impact, and there are a lot of ethical and philosophical questions that we need to consider. And it's important that we bring in different voices, like philosophers, social scientists, artists, and people from the humanities.

*What's the key ethical or philosophical question that we still need to figure out?*

[AI] can be misused, or it can be used by bad actors. So, then there are questions about how you govern the use of this technology globally. How do you govern the use of AI in a way that's aligned with human values?

*Do you think these questions should be left to companies like yours, or should governments get involved in creating regulations?*

It's important for OpenAI and companies like ours to bring this into the public consciousness in a way that's controlled and responsible. But we're a small group of people and we need a ton more input in this system and a lot more input that goes beyond the technologies-—definitely regulators and governments and everyone else.

*There's always a fear that government involvement can slow innovation. You don't think it's too early for policymakers and regulators to get involved?*

It's not too early. It's very important for everyone to start getting involved, given the impact these technologies are going to have.

\*\*\*

# Globethics Publications

The list below is only a selection of our publications. To view the full collection, please visit our website.

All products are provided free of charge and can be downloaded in PDF form from the Globethics library and at www.globethics.net/publications. Bulk print copies can be ordered from *publications@globethics.net* at special rates for those from the Global South.

Prof. Dr Fadi Daou, Executive Director. Prof. Dr Amélé Adamavi-Aho Ékué, Academic Dean, Dr Ignace Haaz, Managing Editor. M. Jakob Bühlmann Quero, Editor Assistant.

Find all Series Editors: https://www.globethics.net/publish-with-us
Contact for manuscripts and suggestions: *publications@globethics.net*

## Global Series

Nicolae Irina / Christoph Stückelberger (eds.), *Mining, Ethics and Sustainability*, 2014, 198pp. ISBN: 978–2–88931–020–3

Philip Lee and Dafne Sabanes Plou (eds), *More or Less Equal: How Digital Platforms Can Help Advance Communication Rights*, 2014, 158pp.
ISBN 978–2–88931–009–8

Sanjoy Mukherjee and Christoph Stückelberger (eds.) *Sustainability Ethics. Ecology, Economy, Ethics. International Conference SusCon III, Shillong/India*, 2015, 353pp. ISBN: 978–2–88931–068–5

Amélie Vallotton Preisig / Hermann Rösch / Christoph Stückelberger (eds.) *Ethical Dilemmas in the Information Society. Codes of Ethics for Librarians and Archivists*, 2014, 224pp. ISBN: 978–288931–024–1.

*Prospects and Challenges for the Ecumenical Movement in the 21st Century. Insights from the Global Ecumenical Theological Institute*, David Field / Jutta Koslowski, 256pp. 2016, ISBN: 978–2–88931–097–5

Christoph Stückelberger, Walter Fust, Obiora Ike (eds.), *Global Ethics for Leadership. Values and Virtues for Life,* 2016, 444pp.
ISBN: 978–2–88931–123–1

Dietrich Werner / Elisabeth Jeglitzka (eds.), *Eco-Theology, Climate Justice and Food Security: Theological Education and Christian Leadership Development*, 316pp. 2016, ISBN 978–2–88931–145–3

Obiora Ike, Andrea Grieder and Ignace Haaz (Eds.), *Poetry and Ethics: Inventing Possibilities in Which We Are Moved to Action and How We Live Together*, 271pp. 2018, ISBN 978–2–88931–242–9

Christoph Stückelberger / Pavan Duggal (Eds.), *Cyber Ethics 4.0: Serving Humanity with Values*, 503pp. 2018, ISBN 978–2–88931–264-1

Christoph Stückelberger / Pavan Duggal (Eds.), *Data Ethics: Building Trust. Serving Humanity with Values*, 675pp. 2023, ISBN 978-2-88931-523-9

## African Law Series

Ghislain Patrick Lessène, *Code international de la détention en Afrique*, 2013, 620pp. ISBN: 978-2-940428-70-0

D. Brian Dennison/ Pamela Tibihikirra-Kalyegira (eds.), *Legal Ethics and Professionalism. A Handbook for Uganda*, 2014, 400pp. ISBN 978–2–88931–011–1

Pascale Mukonde Musulay, *Droit des affaires en Afrique subsaharienne et économie planétaire*, 2015, 164pp. ISBN: 978–2–88931–044–9

Pascal Mukonde Musulay, *Démocratie électorale en Afrique subsaharienne: Entre droit, pouvoir et argent*, 2016, 209pp. ISBN 978–2–88931–156–9

Pascal Mukonde Musulay, *Droits, libertés et devoirs de la personne et des peuples en droit international africain Tome I Promotion et protection*, 282pp. 2021, ISBN 978-2-88931-397-6

Pascal Mukonde Musulay, *Droits, libertés et devoirs de la personne et des peuples en droit international africain Tome II Libertés, droits et obligations démocratiques*, 332pp. 2021, ISBN 978-2-88931-399-0

Ambroise Katambu Bulambo, *Règlement judiciaire des conflits électoraux. Précis de droit comparé africain*, 2021, 672pp., ISBN 978-2-88931-403-4

Osita C. Eze, *Africa Charter on Rights & Duties, Enforcement Mechanism,* 2021, 406pp, ISBN 978-2-88931-414-0

Fweley Diangitukwa, *Les élections en Afrique : Analyse des comportements et pistes normatives de gestion des conflits*, 2022, 432pp., ISBN 978-2-88931-452-2

Kennedy Kihangi Bindu, *Traité de droit de l'environnement. Perspectives congolaises*, 2022, 512p., ISBN 978-2-88931-463-8

Fweley Diangitukwa / Ruth Bekoung Siadous, *Les prisons sont-elles utiles pour l'avenir de la société ? droits de l'homme et milieu carcéral : le cas du Gabon Un témoignage exceptionnel*, 2023, 314pp. ISBN 978-2-88931-501-7

## PJP Series

Isabel Apawo Phiri and Collins Shava (Eds.), *The Africa we pray for: on a Pilgrimage of Justice and Peace*, 2021, 182pp. ISBN: ISBN 978-2-88931-371-6

Fernando Enns, Upolu Lumā Vaai, Andrés Pacheco Lozano and Betty Pries (Eds.), *Transformative Spiritualities for the Pilgrimage of Justice and Peace*, 2022, 323pp. ISBN: 978-2-88931-458-4

*Towards an Ecumenical Theology of Companionship: A Study Document for the Ecumenical Pilgrimage of Justice and Peace*, 2022, 102pp.
ISBN: 978-2-88931468-3

Matthew Ross & Jin Yang Kim (Eds.), *Seek Peace and Pursue It: Reflections on the Pilgrimage of Justice and Peace in Europe*, 2022, 214pp.
ISBN: 978-288931-470-6

Fernando Enns & Stephen G. Brown (Eds.), *Hate Speech and Whiteness: Theological Reflections on the Journey Toward Racial Justice*, 2022, 148pp. ISBN: 978-2-88931-472-0

Ibrahim Wushisi Yusuf & Jin Yang Kim (Eds.), *Our Feet into the Way of Peace: Holistic Approaches to Peace-building in the Context of the Pilgrimage of Justice and Peace*, 2022, 108pp. ISBN: 978-2-88931-474-4

## Theses Series

Sabina Kavutha Mutisya, *The Experience of Being a Divorced or Separated Single Mother: A Phenomenological Study,* 2019, 168pp.
ISBN: 978-2-88931-274-0

Florence Muia, *Sustainable Peacebuilding Strategies. Sustainable Peacebuilding Operations in Nakuru County, Kenya: Contribution to the Catholic Justice and Peace Commission (CJPC),* 2020, 195pp. ISBN: 978-2-88931-331-0

Mary Rose-Claret Ogbuehi, *The Struggle for Women Empowerment Through Education,* 2020, 410pp. ISBN: 978-2-88931-363-1

Nestor Engone Elloué, *La justice climatique restaurative: Réparer les inégalités Nord/Sud*, 2020, 198pp. ISBN 978-2-88931-379-2

Hilary C. Ike, *Organizational Improvement of Nigerian Catholic Chaplaincy in Central Ohio*, 2021, 154pp. ISBN 978-2-88931-385-3

Paul K. Musolo W'Isuka, *Missional Encounter: Approach for Ministering to Invisible Peoples*, 2021, 462pp. ISBN: 978-2-88931-401-0

Andrew Danjuma Dewan, *Media Ethics and the Case of Ethnicity. A contextual Analysis in Plateau State, Nigeria*, 2022, 371pp. ISBN: 978-2-88931-437-9

Sébastien Bintu Munguiko, *La conversion progressive et le rôle d'apôtre dans l'histoire du salut dans les Évangiles Étude exégétique selon Luc-Actes : Les cas de Simon Pierre et de Paul de Tarse*, 2023, 583pp. ISBN 978-2-88931-513-0

## Education Ethics Series

Obiora Ike, Justus Mbae, Chidiebere Onyia, Herbert Makinda (Eds.), *Mainstreaming Ethics in Higher Education Vol. 2*, 2021, 420pp.
ISBN: 978-2-88931-383-9

Christoph Stückelberger, Joseph Galgalo and Samuel Kobia (Eds.), *Leadership with Integrity: Higher Education from Vocation to Funding*, 2021, 280pp.
ISBN: 978-2-88931-389-1

Jacinta M. Adhiambo and Florentina N. Ndeke (Eds.), *Educating Teachers for Tomorrow: on Ethics and quality in Pedagogical Formation*, 2021, 196pp.
ISBN: 978-2-88931-407-2

Erin Green / Divya Singh / Roland Chia (Eds.), *AI Ethics and Higher Education Good Practice and Guidance for Educators, Learners, and Institutions*, 2022, 318pp. ISBN 978-2-88931-442-3

*This is only a selection of our latest publications, to view our full collection please visit:*

**www.globethics.net**

# Data Ethics: Building Trust

*How Digital Technologies Can Serve Humanity*

Data is the magic word of the 21st century. As oil in the 20th century and electricity in the 19th century: For citizens, data means support in daily life in almost all activities, from watch to laptop, from kitchen to car, from mobile phone to politics. For business and politics, data means power, dominance, winning the race. Data can be used for good and bad, for services and hacking, for medicine and arms race.

How can we build trust in this complex and ambiguous data world? How can digital technologies serve humanity? The 45 articles in this book represent a broad range of ethical reflections and recommendations in eight sections: a) Values, Trust and Law, b) AI, Robots and Humans, c) Health and Neuroscience, d) Religions for Digital Justice, e) Farming, Business, Finance, f) Security, War, Peace, g) Data Governance, Geopolitics, h) Media, Education, Communication. The authors and institutions come from all continents.

The book serves as reading material for teachers, students, policy makers, politicians, business, hospitals, NGOs and religious organisations alike. It is an invitation for dialogue, debate and building trust! The book is a continuation of the volume "Cyber Ethics 4.0" published in 2018 by the same editors.

## Editors

**Prof. Dr Dr h.c. Christoph Stückelberger,** Geneva, Switzerland, is Founder and President of Globethics, Professor of Ethics at Basel University (em.) and Visiting professor at universities in Nigeria, China, Russia, UK.

**Dr Pavan Duggal,** New Delhi, India, is an internationally leading expert of Cyber Law, Chairman of the International Commission on Cyber Security, Advocate at the Supreme Court of India and Board member of Globethics.

Globethics