# Intelligent Decision Assistance Versus Automated Decision-Making: Enhancing Knowledge Work Through Explainable Artificial Intelligence

Max Schemmer
Karlsruhe Institute of Technology
max.schemmer@kit.edu

Niklas Kühl
Karlsruhe Institute of Technology
niklas.kuehl@kit.edu

Gerhard Satzger
Karlsruhe Institute of Technology
gerhard.satzger@kit.edu

## Abstract

*While recent advances in AI-based automated decision-making have shown many benefits for businesses and society, they also come at a cost. It has for long been known that a high level of automation of decisions can lead to various drawbacks, such as automation bias and deskilling. In particular, the deskilling of knowledge workers is a major issue, as they are the same people who should also train, challenge and evolve AI. To address this issue, we conceptualize a new class of DSS, namely Intelligent Decision Assistance (IDA) based on a literature review of two different research streams—DSS and automation. IDA supports knowledge workers without influencing them through automated decision-making. Specifically, we propose to use techniques of Explainable AI (XAI) while withholding concrete AI recommendations. To test this conceptualization, we develop hypotheses on the impacts of IDA and provide first evidence for their validity based on empirical studies in the literature.*

## 1. Introduction

The recent advances in Artificial Intelligence (AI) lead to an increase in automated decision-making [1]. Decisions can be classified into unstructured, semi-structured, and structured decisions [2]. Traditionally, automated decision-making was applied to structured problems, and Decision Support Systems (DSS) enhanced decision-making for unstructured problems [2, 3]. Unstructured tasks were considered too difficult to automate since they require more cognitive flexibility [4]. However, advances in AI, specifically in deep learning, now increasingly enable to automate also more complex cognitive tasks, such as driving a car [5]. Therefore, AI has now the potential to also address semi-structured and unstructured decisions that are far from basic back-office tasks [6]. For example, AI is used to automate loan approval [7], or to conduct recruitment choices [8]—decisions that in the past were unimaginable to automate. Therefore, both the number and complexity of tasks that can be automated increase.

However, it has long been known that increasing automation of decisions can lead to various drawbacks, such as automation bias and deskilling [9, 10]. This is especially challenging since most semi-structured and unstructured tasks are knowledge work incorporating high-stake decision making, e.g. medical diagnosis or jurisdictional decisions. In general, AI for knowledge workers should automate routine and assist knowledge-intensive work with reasoning and other high-level functions [11]. The deskilling of knowledge workers is a major problem, as they are the people who should train, challenge and evolve AI. Knowledge workers create the labels for the AI that is the foundation for its initial training. After changes in the environment of the AI knowledge workers adapt and develop new solutions based on their domain expertise [12]. Furthermore, they should be able to challenge the AI's recommendation, either with regard to performance but also with respect to ethical and fairness concerns. While in many use cases these disadvantages may be negligible there are cases where they must not be ignored. Reasons include, but are not limited to, losing significant competitiveness, e.g. in asset investment strategy decisions, or even potentially harming people, e.g. in medical diagnoses.

Because DSS are explicitly designed to not automate but support decision-makers [13], the initially obvious idea emerges to address these problems by using DSS instead of fully automated systems. However, automation should not be interpreted as a binary state but instead as a continuum [10]. Negative impacts already occur at low automation levels [10]—as positive features of human decision-making are reduced such as human engagement. Therefore, when speaking about automated decision-making, we use the broader understanding of the continuum mentioned above, also including lower automation levels. As many state-of-the-art DSS do include automated, AI-based recommendations [2], they are subject to negative

impacts, like automation bias in the short, reduced engagement in the medium, and deskilling in the long term. Thus, we perceive a major research gap in supporting human decision-making without those downsides, and formulate:

**RQ:** *How can we design AI for decision support without introducing automation disadvantages?*

Based on automation and DSS research, we conceptualize a new class of DSS, *Intelligent Decision Assistance* (IDA), that reduces automation-induced disadvantages while still preserving decision support levels. From the automation literature, we draw the critical evaluation of potential disadvantages of automated decision-making and the awareness of a continuum between full automation and human agency [10]. From DSS literature, we use the concept of guidance [14]. Part of guidance theory is the explainability of DSS [14] which is a traditional topic of IS research [15]. We discuss various combinations of automation levels and explainability and eventually follow the idea of informative guidance as a guidance that foregoes to provide explicit recommendations [16]. In line with this notion, we propose to withhold the AI's decision and let the human "brainstorm" together with the AI by providing techniques from the Explainable AI (XAI) knowledge base [17], such as examples, counterfactuals, or feature importance. After conceptualizing IDA and deriving hypotheses on its impact, we provide first evidence for their validity through a systematic evaluation of empirical studies in the literature. With our work, we contribute to research and practice by conceptualizing a new class of DSS—*Intelligent Decision Assistance*.

## 2. Literature Review

In general, IS are designed to support or automate human decision-making [18]. These two purposes are traditionally analyzed in two different research streams: *decision support* is traditionally covered in DSS literature [19], while *Automation* is mainly addressed in Ergonomics literature [1].

### 2.1. Decision Support Systems

DSS represent an important class of IS that aim to provide decisional advice [13]. In general, "DSS is a content-free expression, which means that there is no universally accepted definition" [2, p. 16]. However, DSS can be used as an umbrella term to describe any computerized system that supports decision-making in an organization [2]. Originally, DSS were defined as supportive IT-based systems, aiming at supporting and improving managerial decision-making [13, 20]. Later

developments in DSS opened the area for application to all levels of an organization [13]. In contrast to other IS, DSS focuses on decision-making effectiveness and decision-making efficiency rather than efficiency alone [21].

In general, the decision-making process consists of three phases that are supported through DSS—the intelligence, design, and choice phase [22]. In the intelligence phase, the decision-maker searches, classifies and decomposes problems [2, p. 48-49]. In the design phase, decision alternatives are derived [2, p. 50]. Finally, in the choice phase, the critical phase of decision making, the decisions are chosen [2, p. 58].

An important concept of decision support is decisional guidance that has a long-lasting history in IS literature [14]. Silver [23] differentiates in the form of guidance, which can be either suggestive, quasi-suggestive, or informative. Suggestive guidance makes judgmental recommendations that can also be a set of alliterative decisions [23, p. 94]. Quasi-suggestive guidance is guidance "that does not explicitly make a recommendation but from which one can directly infer a recommendation or direction" [23, p. 109]. Lastly, informative guidance provides decision-makers only with decision-relevant information without suggesting or implying how to act.

Another form of guidance is the explainability of the DSS [14]. Explainability is a concept with a long tradition in IS [24]. With the rise of expert systems, knowledge-based systems, and intelligent agents in the 1980s and 1990s, the IS community has built the basis for research on explainability [15]. In particular, the research stream of Explainable AI (XAI), which addresses the opaqueness of AI-based systems, is gaining momentum. The term XAI was first coined by Van Lent et al. [25] to describe the ability of their system to explain the behavior of agents. The current rise of XAI is driven by the need to increase the interpretability of complex models [26]. In contrast to interpretable linear models, more elaborate models can achieve higher performance [27]. However, their inner workings are hard to grasp for humans. XAI encompasses a wide spectrum of algorithms. These algorithms can be differentiated by their complexity, their scope, and their level of dependency [17]. The interpretability of a model directly depends on its complexity. Wanner et al. [26] define three types of complexity—white, grey, and black-box models. They define white-box models as models with perfect transparency, such as linear regressions. These models do not need additional explainability techniques but are intrinsically explainable. Black-box models, like neural networks, on the other hand, tend to achieve higher performance

but lack interpretability. Lastly, grey-box models are not inherently interpretable but are made interpretable with the help of additional explanation techniques. These techniques can be further differentiated in terms of their scope, i.e., being global or local explanations[17]: Global XAI techniques address holistic explanations of the models as a whole. In contrast, local explanations function on an individual instance basis. Besides the scope, XAI techniques can also be differentiated with regard to being model-specific or model agnostic.

## 2.2. Automation

Research on automation is an essential part of IS research [28] and has been around for more than a century [4] with the overarching goal to increase the efficiency of work by using automation as a means [29]. In general, humans are performing worse than machines in conducting repetitive tasks and are influenced by cognitive bias [30]. Thereby, automation can reduce human bias-induced errors. Automated decision-making applications are designed to minimize human involvement and relieve humans from exhaustive tasks [31]. Additionally, automation acts as an "talent multiplier" that scales human expertise and frees up human capacity to focus on more valuable work [31].

Traditionally, automation has been seen as a binary state—either none or fully automatic [32]. However, Parasuraman et al. [10, p. 287] define automation as "the full or *partial* replacement of a function previously carried out by the human operator" which implies that automation may occur on different levels. The authors propose a taxonomy of automation and develop ten levels. While humans are responsible for decision making at the first five levels, AI has control at the last five levels up to full autonomy at level ten.

Beyond developing the 10-level taxonomy, Parasuraman et al. [10] provide a four-stage model of automated human information processing consisting of information acquisition, information analysis, decision and action selection, and action implementation. This model allows to precisely specify which stage is automated in the decision process.

Although automation has many advantages, some authors have expressed challenges, such as automation bias or cognitive skill reduction leading possibly to deskilling [33]. In the following, we discuss these disadvantages which essentially represent the problem with current approaches that we want to solve.

In the short-term, automation might lead to *Automation bias* which is the "tendency to use automated cues as a heuristic replacement for vigilant information seeking and processing" [9]—essentially representing an over-reliance on AI recommendations. For this reason, sometimes high levels of automation are not desirable if the automation is not perfectly reliable and recommends wrong decisions [34]. These wrong recommendations then can lead to a negative switch from a previously correct human decision [35].

Furthermore, in the long-term, automation bias can result in *deskilling*, either because of the reduction of existing skills or due to the lack of skill development in general [15, 36]. This attacks the collective intellectual capital that is the key asset of many organizations [6]. Many factors might eventually result in deskilling. One factor is the reduced amount of stored information in memory, and more importantly, the reduced mental capability to store information, when using automation, which is commonly known as the "Google effect" [37]. Users seem to reduce investing energy into storing things that can be easily retrieved [36].

Research shows that human *engagement* in the task is particularly important to keep up the vigilantly [9]. Engagement is a psychological state that is broadly defined as an "individual's involvement and satisfaction as well as enthusiasm for work" [38, p. 269] that could reduce potential deskilling [6]. Exemplary, the danger of deskilling can be highlighted with an intelligent asset solution for financial markets. Thereby, the engagement of the broker in the task will reduce which may lead finally to deskilling. Therefore, within the company implementing that solution, the broker deskills—while brokers from companies not implementing the project stay skilled. In the long-term, the environment may eventually change, for example, because of new regulations. Therefore, existing AI solutions need to be built and trained. One of the most important factors in the development process is domain knowledge which may now be reduced due to deskilling. If other companies did not implement AI, they can build and adapt faster and will, therefore, have competitive advantages .

This long-term disadvantage of automated decision support leads to a discussion of efficiency in the short and long-term in human-AI systems. In the short-term AI might increase performance. However, in the long-term due to deskilling, AI systems will not be effectively further trained and evolved. This potentially results in severe negative long-term effects.

## 3. Conceptualization of Intelligent Decision Assistance

In this section, we use the previously depicted research streams of DSS and automation and synthesize them to conceptualize a solution against the

disadvantages of automation. Subsequently, we discuss three particular techniques of this concept

We see two main dimensions that influence the undesired effects of automation, which we discuss in more detail below: First, the general level of human control and agency [10] and, second, the form and degree of explainability [14].

Which level of human agency in automation should be implemented is a notorious discussion in automation literature [30]. Asatiani et al. [6] have discussed that retaining control of human workers may help to sustain their skill level. Similar, Endsley et al. [39] argued that lower automation levels, in general, can keep them cognitively engaged.

Regarding the second dimension, the literature suggests "that a seamless, collaborative interaction between human agents and automated tools, as opposed to using automation as an isolated "black box", could help to prevent the ill effects of deskilling" [6, p. 6]. As discussed, the research stream of XAI addresses this "black box" issue in AI-based automated decision-making. Recent examples [40, 41] demonstrate the capability of XAI to support end-users in their decision-making. By varying the "degree" of explanations, i.e. the system's transparency [42], we believe different effects on the negative aspects of automation could be influenced. On the one hand, some might argue that more explainability is always better. However, the latest research suggests that a high level of automation paired with high explainability might just result in automation bias [43]. Furthermore, the degree of explainability should be adapted to the profession and experience of the end-user, e.g. novice users might need more intuitive and simpler explanations while data scientists can get the full degree of potential explanations [44]. These examples show that also the degree of explainability needs to be chosen thoughtfully.

As introduced, there are many forms of guidance—suggestive, quasi-suggestive, and informative guidance [16]. Suggestive guidance provides the decision-maker with explicit recommendations and tries to increase the guidance of this recommendation. However, as Parasuraman et al. [10] states, also partially automated systems can lead to automation bias and skill degradation. In contrast, as mentioned, informative decisional guidance is a form of guidance where users do not receive explicit recommendations [16]. We follow this line of reasoning and propose a system could simply withhold its recommendation—although it is aware of that recommendation. Parkes [45] validates that suggestive guidance—which is actually a form of automated decision making—can lead to automation

bias, while informative guidance does not have such effects. Research also shows that the effects of the types of guidance vary depending on the task complexity. Montazemi et al. [46] found that suggestive guidance is better for less complex tasks and informative guidance is better with increasing task complexity. This argument strengthens our derivation. Following this line of thought gives rise to the idea to set the degree of automation to almost zero and withhold explicit AI recommendations while keeping support through explanations up. By doing so, we can minimize the drawbacks of automation while still assisting human decision-making. We are creating intelligent systems that are fully capable of solving issues on their own but use their capabilities to inspire and support instead of automating. Based on the derivation, we name this new class of DSS, Intelligent Decision Assistance (IDA) and define it as follows:

**Definition:** *Intelligent Decision Assistance (IDA) is an AI that a) supports humans, b) does not recommend explicit decisions or actions, and c) explains its reasoning*

Referring to the three phases of decision making—intelligence, design, and choice—we mainly support with this approach the intelligence and to some extent the design phase. In terms of final effects on the human, we derive three hypotheses (engagement, performance, automation disadvantages).

First, IDA provides decision-makers with options to actively engage with the task by interactively requesting explanations, interpreting them and essentially communicating with the AI. As Asatiani et al. [6] [6] have discussed providing explanations instead of using automation as an isolated "blackbox" could result in an engaged human-AI collaboration. Thus, we hypothesize:

**H1:** *IDA increases engagement with the task.*

Beyond that, we hypothesize that IDA should increase human performance. While especially, if the automation is far better than the human, IDA will most likely not exceed automated decision-making, it should still improve the performance by providing guidance and especially insights. Therefore, we formulate:

**H2:** *IDA performance outperforms the human alone.*

Lastly, because IDA does not incorporate higher levels of automation it should reduce automation disadvantages and especially prevent deskilling. Therefore, we formulate the following hypothesis:

**H3:** *IDA reduces automation induced disadvantages.*

In the next section, we are going to test these hypotheses based on empirical studies in the literature.

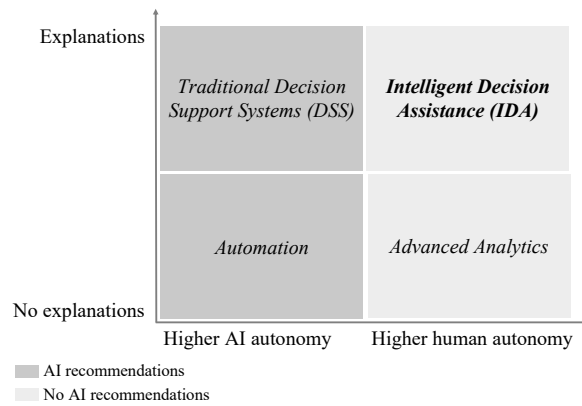Figure 1 depicts IDA in the continuum of both

Figure 1. Positioning of Intelligent Decision Assistance on the two dimensions of explainability and degree of automation

discussed dimensions. We depict different types of systems for decision-making. At a high level of automation and almost no explanations, we position automation [2]. Traditional DSS come also usually with a higher level of automation, through providing explicit recommendation, but additionally provide explanations for the decision-maker. We delimit ourselves from DSS that use AI to transform unstructured data into structured data and DSS that use AI to produce a pre-decision output, e.g. a forecast. As stated, Parasuraman et al. [10] define four stages of automation—information acquisition, information analysis, decision-making, and actions. Following this classification, we focus on the decision-making level. This classification allows us also to differentiate IDA from Advanced Analytics [47]. While advanced analytics may incorporate AI solutions they are always on the information acquisition or analysis level. In contrast, IDA allows the decision-maker to actively engage on the decision level and is positioned in the right top corner of Figure 1 with high explainability and full human autonomy.

Now that we derived, defined, and delimit IDA, we discuss specific explanation techniques that support IDA and consequently pose valid implementation options. Specifically, we discuss feature importance, example-based explanations, and counterfactual explanations. We explain these features based on the example of a loan approval decision-making task.

**Feature importance**: Feature importance is a model-agnostic technique that gives the decision-maker information about the importance of specific data points. Two famous algorithms of feature importance are LIME [48] and SHAP [49]. In a loan approval decision where the banker has information about past credits, expenses, demographics, etc., one could now train artificial

intelligence to make this decision and recommend explicit decisions. In contrast, IDA would withhold the specific AI decision but provide the decision-maker, i.e. the banker, with information on which data was in particular important for the AI's decision. In an IDA this information could now be used for various use cases. Now in the time of big data, e.g. having many information on customers, one particular great use case would be to filter or sort the features in an intelligent way based on the feature importance.

**Example-based explanations**: Example-based explanations provide historical data that is similar to the current instance [50]. Example-based explanations, therefore essentially represent some form of information retrieval. Research in psychology states that humans prefer explanations that show examples [51]. Furthermore, examples can be used within complex tasks [52]. Referring to our loan approval case, the decision-maker would receive information on past approvals that were similar. In an IDA, the decision-maker would get information about similar historical cases that are labeled. Based on these examples, the decision-maker should be able to infer differences or similarities.

**Counterfactual explanation**: Counterfactual explanations give information on what the smallest change would be to get a different AI decision [53]. Counterfactual explanations take a similar form to the statement [54]: "You were denied a loan because your annual income was 20,000. If your income had been 45,000, you would have been offered a loan." In an IDA a counterfactual explanation would look like the following: "Your current annual income is £30,000. If your income would be £45,000, the AI's decision would change." This type of non-intrusive explanation would lead to an increased thought process of the decision-maker.

Figure 2 highlights the idea of IDA for a credit allowance example. On the left side, we display a traditional interface for automated decision-making. On the right side, IDA is visualized. In the traditional interface, the decision-maker gets a specific recommendation. Additionally, the decision-maker gets the available information on the credit applicant, the importance of the features for the decision, and optional explanation options. In contrast, an IDA does not provide a specific recommendation, but rather various XAI techniques that allow the decision-maker to "brainstorm" with the AI.
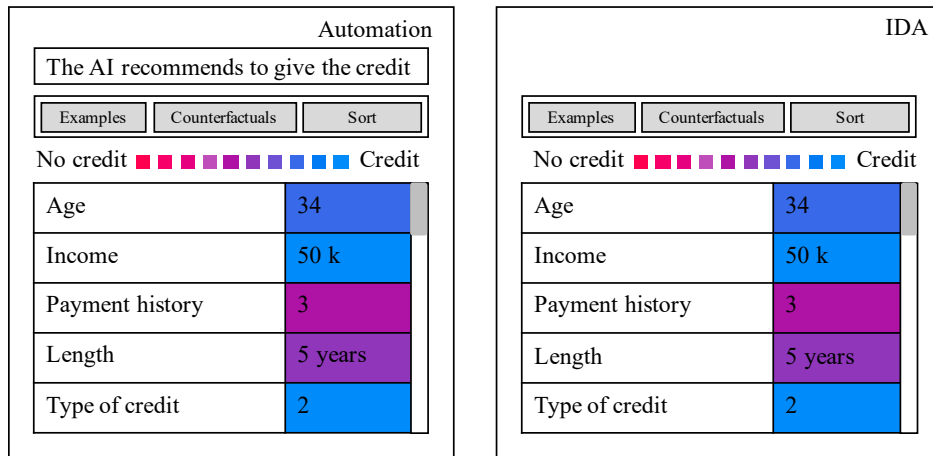
**Figure 2.** Comparison of traditional automated decision-making and Intelligent Decision Assistance (IDA)

## 4. Validation Study

After deriving a conceptualization of IDA, we validate our concept by conducting a literature-based validation study based on the methodology outlined by Brocke et al. [55]. The goal of the study is to find empirical studies that tested variations of automation and explainability and to analyze whether the findings do support our hypotheses above. This means they should address the degree of automation and explainability. For this reason, our search string consists of two main parts. The first reflects XAI, including relevant synonyms, such as "explainable AI" or "interpretability" comprises of "Artificial Intelligence". The second part comprised synonyms of behavioral experiments, e.g., "user study" or "user evaluation". To find the synonyms, we initiated our SLR with an explorative search. The search string was iteratively extended resulting in the following final search string:

*TITLE-ABS-KEY("explainable artificial intelligence" OR XAI OR "explainable AI" OR ( ( interpretability OR explanation ) AND ("artificial intelligence" OR ai OR "machine learning" ) ) ) AND ( "human performance" OR "human accuracy" OR "user study" OR "empirical study" OR "online experiment" OR "human experiment" OR "behavioral experiment" OR "human evaluation" OR "user evaluation")*

Then, we selected an appropriate database. Our exploratory search indicated that relevant work is dispersed across multiple disciplines, publishers, conferences, and journals. For this reason, we chose the SCOPUS database, to ensure comprehensive coverage. Following that, we defined our inclusion criteria. We included every article that (a) conducted empirical research, (b) reported performance measures, (c) focused on an application context where AI supports humans on the decision level, and (d) provided an IDA setting. With our search string defined, we conducted the SLR from January to March 2021. We identified 256 articles through the keyword-based search. As a next step, we analyzed the abstract of each article and filtered based on our inclusion criteria, leading to 61 articles. Afterward, two independent researchers read all articles in detail and applied the inclusion criteria again. Based on these, we conducted a forward and backward search. This led to a total of five articles that were consequently analyzed in-depth to collect data about each experiment. The data collection process was conducted by two independent researchers who discussed and homogenized differences. The main focus of the validation study was to extract the treatments and outcomes of each experiment reported in the studies. For example, if two XAI techniques were used and compared as separate experimental treatments we added two entries into our database. In total, we identified five articles and 12 experiments [40, 56, 57, 58, 59]. In the following, we describe the studies and their results with regard to IDA in detail.

Carton et al. [56] conduct an experiment on online toxicity classification of social media posts. They use feature importance to highlight words that were relevant for the classification. As one condition they have the prediction presence. In their experiment, they find no significant effect of examples. However, they find signs of automation bias:"We find that the presence of a visible model prediction tends to bias subjects in favor of the prediction, whether it is correct or incorrect." [56, p. 101]

Chu et al. [57] conduct an experiment on age guessing supported through AI. They test three different conditions of explanations and the visibility of AI

predictions. The authors found no significant effects of explanations but also signs of automation bias: "The predictions generally help whenever the human is inaccurate [...], but can hurt when the human is accurate and the model is inaccurate [...]." [57, p. 5]

Lai and Tan [40] and Lai et al. [58] refer in their studies also to the ten levels of automation introduced by [10] and test various XAI techniques without ever displaying what the actual AI's decision is on a deception detection task. For example, they highlight all words that were relevant for the decision (unsigned) [58]. Another condition was to colorize this highlight differently depending on the influence of the words (signed). Their results show that signed highlights result in a significant increase in XAI-assisted performance (70.7% for signed, and 60.4% for human performance) [58]. In Lai and Tan [40] they test additionally the influence of example-based explanations with also positive but not significant effects. However, also in Lai and Tan [40] two highlight-based conditions showed significant positive effects in terms of short-term performance.

Lastly, Schmidt and Biessmann [59] conduct two different tasks in their experiment—a book category classification based on their descriptions and a movie rating classification. They test two different XAI algorithms, both feature importance techniques to highlight important words. Both data sets and both XAI algorithms show an increase in IDA performance with one algorithms generating significant results on both data sets.

Table 1. Validation study results

| Source | Engagement | Performance | Automation |
|--------|------------|-------------|------------|
| [56] | No Measurement | No effect | Automation Bias |
| [57] | No Measurement | No effect | Automation Bias |
| [40] | No Measurement | Improvement | No Measurement |
| [58] | No Measurement | Improvement | No Measurement |
| [59] | No Measurement | Improvement | No Measurement |

Table 1 summarizes our results of the validation study. Regarding our first hypothesis (**H1**), we can see that current research fails to provide insights into the effect of IDA on engagement. Regarding **H2**, three papers validated our hypotheses that IDA performance should exceed human performance. Lastly, regarding **H3**, two of the studies showed signs of Automation Bias in the presence of explicit AI recommendations, which is an indicator of potential long-term deskilling effects [15, 36].

## 5. Discussion

Overall, the validation study provides first support for the hypotheses on the impact of IDA and highlights the potential of IDA through five experiments with significant positive effects and none with significant negative effects. Furthermore, the study shows that current research lacks insights on the influence of IDA on engagement which should be addressed in future research.

IDA has of course also limitations. One of them might be the perceived usefulness. Telling the decision-maker that the AI would be theoretically capable of providing them with a recommendation but this recommendation is to withhold may be perceived as annoying for decision-makers, especially if they are under time pressure. Therefore, the advantages of IDA need to be highlighted. One attenuated option could be to show the explanations on default, but the recommendation just on request. Another limitation is the potential high computational costs. Some XAI techniques, e.g. SHAP values [49], are computational inefficient. Therefore, the computational costs, especially in comparison to traditional analytics tools might be much higher. This trade-off has to be determined for individual cases.

We want to clarify that IDA should not be applied in every use case. We explicitly derive this idea for knowledge work and not for repetitive structured work. Especially for jobs where the disadvantages of automation are critical, IDA should be taken into account. Among others, in high stake decision-making such as medicine, law, or human resource. But also in knowledge-intensive areas where the competitive advantage is based on knowledge, such as finance. However, as pointed out by Endsley and Kaber [32], for structured tasks that require low flexibility and have a high system performance, full automation can be the best option.

Additionally, we want to discuss an additional advantage that may have a temporary influence on the adoption of IDA. Paragraph 22 of the GDPR states: "The data subject shall have the right not to be subject to a decision based solely on automated processing [...]." [60] This means that in some cases automated decision-making is simply forbidden. Here the best possible augmentation through IDAs could be a valuable approach.

Furthermore, IDAs could have a positive influence on the fairness of AI-enhanced decision-making. AI algorithms can have biases that can lead to unfair decision-making. With IDAs, we allow people to have full control over the final decision and can thus reduce

bias.

Finally, there are some open questions. Future work should empirically validate whether IDAs prevent deskilling and other automation disadvantages and in contrast increases engagement. Furthermore, one should access the efficiency effects of IDA on human decision-making. For example, Fazlollahi et al. [61] find that decisional guidance increases decision time. However, also direct recommendations may decrease efficiency if they lead to cognitive dissonance and consequently to an in-depth analysis of the decision-maker. The efficiency of IDAs needs to be compared to pure human and automated approaches.

## 6. Conclusion

The main goal of this study was to conceptualize a solution to automation-induced disadvantages, such as automation bias or deskilling. To do so, we initiated our research by conducting a literature review of automation and DSS literature. Based on these two research streams, we conceptualized a new class of DSS, namely *Intelligent Decision Assistance* (IDA). IDA augments human decision-making through Explainable AI (XAI) while withholding explicit AI recommendations. Thereby, IDA aims to provide insight into the data without generating automation disadvantages. Subsequently, we validated our conceptualization by searching for empirical literature which shows first evidence of our hypotheses.

Our contributions are threefold: First, we synthesize the body of knowledge in automation sciences and decision support literature. Second, we conceptualize a new class of systems—IDA—and third, we test three hypotheses regarding the potential of IDA.

Unleashing the potential of IDA requires a multidimensional design process. For this reason, we see the IS research community as the predestined research discipline to advance research in this field. We hope to motivate IS researchers and practitioners to actively participate in the exploration of IDA.

## References

[1] Crispin Coombs, Donald Hislop, Stanimira K Taneva, and Sarah Barnard. The strategic impacts of intelligent automation for knowledge and service work: An interdisciplinary review. *The Journal of Strategic Information Systems*, 29(4): 101600, 2020.

[2] Efraim Turban, Ramesh Sharda, and Dursun Delen. Decision support and business intelligence systems (required). *Google Scholar*, 2010.

[3] George Anthony Gorry and Michael S Scott Morton. A framework for management information systems. 1971.

[4] Mary C Lacity and Leslie P Willcocks. A new approach to automating services. *MIT Sloan Management Review*, 58(1):41–49, 2016.

[5] Carl Benedikt Frey and Michael A Osborne. The future of employment: How susceptible are jobs to computerisation? *Technological forecasting and social change*, 114:254–280, 2017.

[6] Aleksandre Asatiani, Esko Penttinen, Tapani Rinta-Kahila, and Antti Salovaara. Implementation of automation as distributed cognition in knowledge work organizations: Six recommendations for managers. In *40th international conference on information systems*, pages 1–16, 2019.

[7] Infosys. How FinTechs can enable better support to FIs' credit decisioning? 2019.

[8] Edward Tristram Albert. Ai in talent acquisition: a review of ai-applications used in recruitment and selection. *Strategic HR Review*, 2019.

[9] Kathleen L Mosier and Linda J Skitka. Automation use and automation bias. In *Proceedings of the human factors and ergonomics society annual meeting*, volume 43, pages 344–348. SAGE Publications Sage CA: Los Angeles, CA, 1999. ISBN 1541-9312.

[10] Raja Parasuraman, Thomas B. Sheridan, and Christopher D. Wickens. A model for types and levels of human interaction with automation. *IEEE Transactions on Systems, Man, and Cybernetics Part A:Systems and Humans.*, 30(3):286–297, 2000. ISSN 10834427.

[11] Jennifer Adelstein. Disconnecting knowledge from the knower: the knowledge worker as icarus. *Equal Opportunities International*, 26(8): 853–871, 2007.

[12] Lucas Baier, Niklas Kühl, and Gerhard Satzger. How to cope with change?-preserving validity of predictive services over time. In *Proceedings of the 52nd Hawaii International Conference on System Sciences*, 2019.

[13] David Arnott and Graham Pervan. A critical analysis of decision support systems research revisited: the rise of design science. In *Enacting Research Methods in Information Systems*, pages 43–103. Springer, 2016.

[14] Stefan Morana, Silvia Schacht, Ansgar Scherp,

and Alexander Maedche. A review of the nature and effects of guidance design features. *Decision Support Systems*, 97:31–42, 2017.

[15] Christian Meske, Enrico Bunde, Johannes Schneider, and Martin Gersch. Explainable artificial intelligence: Objectives, stakeholders, and future research opportunities. *Information Systems Management*, pages 1–11, 2020.

[16] Mark S Silver. Decisional guidance for computer-based decision support. *MIS quarterly*, pages 105–122, 1991.

[17] Amina Adadi and Mohammed Berrada. Peeking inside the black-box: a survey on explainable artificial intelligence (xai). *IEEE access*, 6: 52138–52160, 2018.

[18] Shoshana Zuboff. Automate/informate: The two faces of intelligent technology. *Organizational dynamics*, 14(2):5–18, 1985.

[19] Daniel J Power. A brief history of decision support systems. *DSSResources. com*, 3, 2007.

[20] Lawrence F Young. Right-brained decision support systems. *ACM SIGMIS Database: the DATABASE for Advances in Information Systems*, 14(4):28–36, 1983.

[21] Gerald E Evans and James R Riha. Assessing dss effectiveness using evaluation research methods. *Information & management*, 16(4):197–206, 1989.

[22] Herbert A Simon. The new science of management decision. 1960.

[23] Mark S Silver. Decisional guidance. *Human–Comouter*, page 90, 2015.

[24] Shirley Gregor and Izak Benbasat. Explanations from intelligent systems: Theoretical foundations and implications for practice. *MIS quarterly*, pages 497–530, 1999.

[25] Michael Van Lent, William Fisher, and Michael Mancuso. An explainable artificial intelligence system for small-unit tactical behavior. In *Proceedings of the national conference on artificial intelligence*, pages 900–907, 2004.

[26] Jonas Wanner, Lukas-Valentin Herm, Kai Heinrich, Christian Janiesch, and Patrick Zschech. White, grey, black: Effects of xai augmentation on the confidence in ai-based decision support systems. 2020.

[27] Erica Briscoe and Jacob Feldman. Conceptual complexity and the bias/variance tradeoff. *Cognition*, 118(1):2–16, 2011.

[28] Ulrich Frank. Increasing the level of automation in organisations: Some remarks on formalisation, contingency and the social construction of reality. *The Systemist*, 20:98–113, 1998.

[29] Katsundo Hitomi. Automation—its concept and a short history. *Technovation*, 14(2):121–128, 1994.

[30] Jeffrey Heer. Agency plus automation: Designing artificial intelligence into interactive systems. *Proceedings of the National Academy of Sciences*, 116(6):1844–1850, 2019.

[31] Jeanne G Harris and Thomas H Davenport. Automated decision making comes of age. *MIT Sloan Management Review*, 46(4):2–10, 2005.

[32] Mica R Endsley and David B Kaber. Level of automation effects on performance, situation awareness and workload in a dynamic control task. *Ergonomics*, 42(3):462–492, 1999.

[33] Lisanne Bainbridge. Ironies of automation. In *Analysis, design and evaluation of man–machine systems*, pages 129–135. Elsevier, 1983.

[34] Nadine B Sarter and Beth Schroeder. Supporting decision making and action selection under time pressure and uncertainty: The case of in-flight icing. *Human factors*, 43(4):573–583, 2001.

[35] Kate Goddard, Abdul Roudsari, and Jeremy C Wyatt. Automation bias: a systematic review of frequency, effect mediators, and mitigators. *Journal of the American Medical Informatics Association*, 19(1):121–127, 2012.

[36] Steve G Sutton, Vicky Arnold, and Matthew Holt. How much automation is too much? keeping the human relevant in knowledge work. *Journal of emerging technologies in accounting*, 15(2): 15–25, 2018.

[37] Betsy Sparrow, Jenny Liu, and Daniel M Wegner. Google effects on memory: Cognitive consequences of having information at our fingertips. *science*, 333(6043):776–778, 2011.

[38] James K Harter, Frank L Schmidt, and Theodore L Hayes. Business-unit-level relationship between employee satisfaction, employee engagement, and business outcomes: a meta-analysis. *Journal of applied psychology*, 87(2):268, 2002.

[39] Mica R Endsley et al. The role of situation awareness in naturalistic decision making. *Naturalistic decision making*, 269:284, 1997.

[40] Vivian Lai and Chenhao Tan. On human predictions with explanations and predictions of machine learning models: A case study on

deception detection. In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 29–38, 2019.

[41] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. Anchors: High-precision model-agnostic explanations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.

[42] Michael Vössing. *Designing Human-computer Collaboration: Transparency and Automation for Intelligence Augmentation*. PhD thesis, Karlsruher Institut für Technologie (KIT), 2020.

[43] Patrick Hemmer, Max Schemmer, Michael Vössing, and Niklas Kühl. Human-ai complementarity in hybrid intelligence systems: A structured literature review. In *PACIS 2021 PROCEEDINGS*, 2021.

[44] Niklas Kuehl, Jodie Lobana, and Christian Meske. Do you comply with ai?—personalized explanations of learning algorithms and their impact on employees' compliance behavior. *ICIS*, 2020.

[45] Alison Parkes. Persuasive decision support: Improving reliance on decision aids. *PACIS*, 4(3): 2, 2012.

[46] Ali R Montazemi, Feng Wang, SM Khalid Nainar, and Christopher K Bart. On the effectiveness of decisional guidance. *Decision Support Systems*, 18 (2):181–198, 1996.

[47] Hugh J Watson. Tutorial: Big data analytics: Concepts, technologies, and applications. *Communications of the Association for Information Systems*, 34(1):65, 2014.

[48] Marco Tulio Ribeiro, Sameer Singh, and Carlos Guestrin. "why should I trust you?": Explaining the predictions of any classifier. In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, San Francisco, CA, USA, August 13-17, 2016*, pages 1135–1144, 2016.

[49] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in NIPS 30*, pages 4765–4774. Curran Associates, Inc., 2017.

[50] Jasper van der Waa, Elisabeth Nieuwburg, Anita Cremers, and Mark Neerincx. Evaluating xai: A comparison of rule-based and example-based explanations. *Artificial Intelligence*, 291:103404, 2021.

[51] Carrie J Cai, Jonas Jongejan, and Jess Holbrook. The effects of example-based explanations in a machine learning interface. In *Proceedings of the 24th International Conference on Intelligent User Interfaces*, pages 258–262, 2019.

[52] Robert Glaser. Intelligence as acquired proficiency. *What is intelligence*, pages 77–83, 1986.

[53] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without opening the black box: Automated decisions and the gdpr. *Harv. JL & Tech.*, 31:841, 2017.

[54] Jakob Schoeffer, Yvette Machowski, and Niklas Kuehl. A study on fairness and trust perceptions in automated decision making. 2021.

[55] Jan vom Brocke, Alexander Simons, Bjoern Niehaves, Bjorn Niehaves, Kai Reimer, Ralf Plattfaut, and Anne Cleven. Reconstructing the giant: On the importance of rigour in documenting the literature search process. 2009.

[56] Samuel Carton, Qiaozhu Mei, and Paul Resnick. Feature-based explanations don't help people detect misclassifications of online toxicity. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 14, pages 95–106, 2020.

[57] Eric Chu, Deb Roy, and Jacob Andreas. Are visual explanations useful? a case study in model-in-the-loop prediction. *arXiv preprint arXiv:2007.12248*, 2020.

[58] Vivian Lai, Han Liu, and Chenhao Tan. " why is' chicago'deceptive?" towards building model-driven tutorials for humans. In *Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems*, pages 1–13, 2020.

[59] Philipp Schmidt and Felix Biessmann. Quantifying interpretability and trust in machine learning systems. *arXiv preprint arXiv:1901.08558*, 2019.

[60] Art. 22 gdpr – automated individual decision-making, including profiling, Jul 2018. URL https://gdpr-info.eu/art-22-gdpr/.

[61] Bijan Fazlollahi, Mihir A Parikh, and Sameer Verma. Evaluation of decisional guidance in decision support systems: an empirical study. *Retrieved November*, 2:2012, 1995.