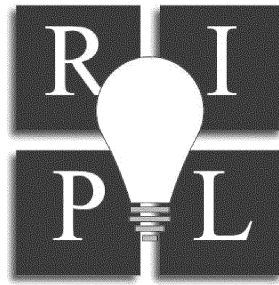


THE JOHN MARSHALL REVIEW OF INTELLECTUAL PROPERTY LAW



THE APPLICATION OF DATA ACCESS POLICIES DESIGNED FOR GENOME- WIDE ASSOCIATION STUDIES TO SMALLER SCALE DATABASES

DONNA M. GITTER

ABSTRACT

Scientific progress thrives with open discussion of new ideas and supporting data. To this end, researchers traditionally publish their results in scientific papers—papers that contain the new ideas and the underlying data supporting those ideas. With the advent of large-scale and high-throughput data analysis, however, the creation of scientific databases have replaced the traditional model. For such publically-funded, data-intensive projects, funding agencies typically require that all relevant data be made available on a publicly accessible website at the time of the paper's publication. Against the backdrop of the public accessibility model used in the 1000 Genomes Project, the author recommends that a modified framework be applied to smaller scale data collection projects. Such a framework could overcome the data producers' concern for protecting the data they have created, thereby encouraging researchers to share data from smaller scale studies.

Copyright © 2011 The John Marshall Law School



*Cite as Donna M. Gitter, The Application of Data Access Policies Designed
for Genome-Wide Association Studies to Smaller Scale Databases,*

10 J. MARSHALL REV. INTELL. PROP. L. 476 (2011).

THE APPLICATION OF DATA ACCESS POLICIES DESIGNED FOR GENOME-
WIDE ASSOCIATION STUDIES TO SMALLER SCALE DATABASES

DONNA M. GITTER

I. INTRODUCTION.....	477
II. THE 1000 GENOMES PROJECT: GOALS AND CHALLENGES	478
III. THE SCIENCE OF THE 1000 GENOMES PROJECT	479
IV. THE DATA ACCESS POLICY OF THE 1000 GENOMES PROJECT	481
V. OBSTACLES TO DATA-SHARING AND POTENTIAL SOLUTIONS OFFERED BY THE TORONTO STATEMENT	483
A. Lack of Recognition of Data Producers	483
B. The Challenge of Crafting Appropriate Publication Policies	485
C. The Complexity of Establishing a Suitable Infrastructure	488
VI. THE NEED FOR RESTRICTED ACCESS DATABASES	489

THE APPLICATION OF DATA ACCESS POLICIES DESIGNED FOR GENOME-WIDE ASSOCIATION STUDIES TO SMALLER SCALE DATABASES

DONNA M. GITTER*

I. INTRODUCTION

Scientific discourse and progress depend upon the open discussion of ideas and full disclosure of supporting facts.¹ This discussion has traditionally occurred through the process of publication, which is the primary means by which researchers achieve recognition for their work.² Traditionally, researchers published papers that combined in one entity both their ideas and the underlying data.³ With the advent of large-scale and high-throughput data analyses, however, the creation of scientific databases replaced the traditional model.⁴ Typically, for such data-intensive projects, funding agencies require that all relevant data must be made available on a publicly accessible website at the time of the paper's publication.⁵

The Human Genome Project ("HGP"), completed in 2003, demonstrated to the scientific community that making data broadly available before publication results in valuable benefits to the public.⁶ This is particularly true where there is a community of scientists who can use the data more quickly than the data producers themselves, and in ways not originally anticipated at the outset of the project.⁷ One successor to the HGP is the 1000 Genomes Project, which provides that project data will be released quickly, prior to publication, into the public domain.⁸

While the open access approach of the 1000 Genomes Project is the norm for large-scale, publicly-funded genomic databases, many smaller projects that likewise produce vast amounts of scientific data nonetheless do not embrace data-sharing.⁹

* © Donna M. Gitter 2011. Associate Professor of Law, Baruch College, New York, New York. J.D., University of Pennsylvania Law School; B.A., Cornell University College of Arts and Sciences. E-mail: Donna.Gitter@Baruch.cuny.edu. I would like to thank two colleagues at The Australian National University College of Law for encouraging me to develop this work: Dr. Matthew Rimmer, Associate Professor, Associate Director of Research, and Associate Director, The Australian Centre for Intellectual Property in Agriculture; and Alison McLennan, Vice Chancellor's Scholar. A version of this work will appear in a book they have edited, entitled *Intellectual Property and Emerging Biotechnologies*, which will be published in 2011. I would also like to thank the editors of the *Review of Intellectual Property Law* of The John Marshall Law School for inviting me to contribute this work to their symposium "Biotechnology and Health-Related Issues in IP Law."

¹ Ewan Birney et al., *Prepublication Data Sharing*, 461 NATURE 168, 168 (2009).

² Paul N. Schofield et al., *Post-publication Sharing of Data and Tools*, 461 NATURE 171, 171 (2009).

³ Birney et al., *supra* note 1, at 168.

⁴ *Id.*

⁵ *Id.*

⁶ *Id.*

⁷ *Id.*

⁸ See *About the 1000 Genomes Project*, 1000 GENOMES, <http://www.1000genomes.org/about#ProjectOverview> (last visited Mar. 25, 2011).

⁹ J.H. Reichman & Paul F. Uhler, *The Public Domain: A Contractually Reconstructed Research Commons for Scientific Data in a Highly Protectionist Intellectual Property Environment*, 66 LAW & CONTEMP. PROBS. 315, 344 (2003).

Quite often, researchers involved in “small science,” meaning research “performed by individual investigators or small and autonomous research groups operating outside large, organized research programs, often with non-federal sources of funding,” are not as influenced by the norms of open access and sharing and instead rely on informal exchanges of data and samples.¹⁰ Scientific researchers, ethicists, lawyers, representatives of funding agencies, and journal editors are all actively engaged in discussions aimed at establishing new scientific norms in order to overcome reluctance to data-sharing in the research community.¹¹ Close investigation of the 1000 Genomes Project provides guidance in terms of applying the data-sharing lessons learned from large-scale, publicly-funded databases to smaller projects.¹²

II. THE 1000 GENOMES PROJECT: GOALS AND CHALLENGES

Launched in January 2008, the 1000 Genomes Project aims to develop a new map of the human genome that will provide an extraordinarily detailed view of biomedically relevant DNA variations among individuals.¹³ A clearer understanding of these small genetic variations will assist in explaining individual differences in susceptibility to disease, responses to drugs, and reactions to environmental factors.¹⁴ The 1000 Genomes Project builds upon the HapMap Project, a large-scale, publicly-funded genomic project that established in 2005 a public database of common genetic variants in human beings.¹⁵ Using HapMap data and related resources, researchers have succeeded in identifying more than one hundred regions of the genome containing genetic variants that are associated with risks of common human diseases such as diabetes, coronary artery disease, and prostate and breast cancer.¹⁶ The goal of the 1000 Genomes Project is to study these genetic differences in greater depth, thereby increasing the ability of researchers to discern the causal relationships between genetic variations and human disease.¹⁷ Ultimately, the 1000 Genomes Project will expedite efforts to diagnose, treat, and prevent human diseases.¹⁸

Like the HapMap Project, the 1000 Genomes Project relies upon a combination of public and private support, along with the expertise of an international consortium of research teams.¹⁹ The research teams for the 1000 Genomes Project will sequence the genomes of at least one thousand people worldwide in order to create a highly detailed database cataloguing human genetic variation.²⁰ As with other large-scale, publicly funded genomic projects, data from the 1000 Genomes Project will be made

¹⁰ See *id.* at 322–23.

¹¹ See Birney et al., *supra* note 1, at 168–69.

¹² See *id.* at 169.

¹³ *International Consortium Announces the 1000 Genomes Project*, NAT'L INSTS. HEALTH (Jan. 22, 2008), <http://www.nih.gov/news/health/jan2008/nhgri-22.htm>.

¹⁴ *Id.*

¹⁵ See generally Int'l HapMap Consortium, *A Haplotype Map of the Human Genome*, 437 NATURE 1299 (2005) (analyzing data collected from phase one of the HapMap Project).

¹⁶ *International Consortium Announces*, *supra* note 13.

¹⁷ See *About the 1000 Genomes Project*, *supra* note 8.

¹⁸ *International Consortium Announces*, *supra* note 13.

¹⁹ *Id.*; see *Three Sequencing Companies Join 1000 Genomes Project*, NAT'L INSTS. HEALTH (June 11, 2008), <http://www.nih.gov/news/health/jun2008/nhgri-11.htm>.

²⁰ *Three Sequencing Companies Join*, *supra* note 19.

available as rapidly as possible to the international scientific community via release into freely accessible public databases.²¹ The quantity of data will be immense, posing great challenges for the experts in the fields of bioinformatics and statistical genetics.²² According to Dr. Gil McVean of the University of Oxford in England, one of the co-chairs of the consortium's analysis group: "At six trillion bases, the 1000 Genomes Project will generate sixty-fold more sequence data over its three-year course than have been deposited into public databases over the past twenty-five years."²³ Certainly, the 1000 Genomes Project Consortium faces a significant challenge in devising the most effective way to analyze this overwhelming quantity of data.²⁴ What is more, the consortium must ensure that its data access policy will make the data available to the maximum number of users, while simultaneously permitting researchers to seek intellectual property protection that will foster innovation.²⁵

III. THE SCIENCE OF THE 1000 GENOMES PROJECT

A gene is the basic physical and functional unit of heredity for humans and almost all other organisms. Genes are made up of DNA, which is composed in part of chemical compounds called bases.²⁶ Genes cause the body to make molecules called proteins that guide bodily functioning.²⁷ In humans, genes vary in size from a few hundred DNA bases to more than two million bases.²⁸ The Human Genome Project has estimated that humans have between 30,000 and 40,000 genes.²⁹

The DNA sequence of any two people is 99.9 percent identical, with very small genetic variations proving quite significant in terms of an individual's susceptibility to disease or response to pharmaceuticals.³⁰ Sites in the DNA sequence where individuals differ at a single DNA base are called single nucleotide polymorphisms ("SNPs").³¹ Sets of nearby SNPs on the same chromosome are inherited in blocks called haplotypes.³² While a haplotype block may contain a large number of SNPs, researchers are able to uniquely identify a haplotype by using a few SNPs, called tag SNPs.³³

²¹ *International Consortium Announces*, *supra* note 13.

²² *Id.*

²³ *Id.*

²⁴ *Id.*; see also *Three Sequencing Companies Join*, *supra* note 19; John M. Conley et al., *Enabling Responsible Public Genomics*, 20 HEALTH MATRIX 325, 329–34 (2010).

²⁵ *Three Sequencing Companies Join*, *supra* note 19; see The 1000 Genomes Project Consortium, *A Map of Human Genome Variation from Population-Scale Sequencing*, 467 NATURE 1061, 1062 (2010).

²⁶ LEROY WALTERS & JULIE GAGE PALMER, *THE ETHICS OF HUMAN GENE THERAPY* 5 (1997).

²⁷ *Id.* at 4.

²⁸ *Id.* at 5–7.

²⁹ Donna M. Gitter, *International Conflicts Over Patenting Human DNA Sequences in the United States and the European Union: An Argument for Compulsory Licensing and a Fair-Use Exemption*, 76 N.Y.U. L. REV. 1623, 1633 (2001).

³⁰ *International Consortium Announces*, *supra* note 13.

³¹ *International HapMap Project*, NAT'L HUM. GENOME RES. INST., <http://www.genome.gov/10001688> (last visited Mar. 25, 2011).

³² *Id.*

³³ *Id.*

The HapMap is a map of the haplotype blocks.³⁴ Its value lies in the fact that the HapMap reduced the number of SNPs required to examine the entire genome for association with a phenotype from the ten million SNPs that exist to roughly 500,000 tag SNPs.³⁵ Thus, researchers can now locate much more efficiently and comprehensively the regions of the genome with genes that affect disease, because it is not necessary to study more SNPs than necessary.³⁶

The HapMap catalogue, however, identifies only those genetic variations that occur in at least five percent of the population.³⁷ The 1000 Genomes Project will use new sequencing technologies to “produce a catalogue of variants present at 1 percent or greater frequency in the human population across most of the genome, and down to 0.5 percent or lower within genes.”³⁸ This will permit researchers to gain a better understanding of some diseases that arise less often.³⁹

The 1000 Genomes Project will also surpass the HapMap Project in that it will map not only SNPs, the single-letter differences in people’s DNA, but also map larger differences in genome structure called structural variants.⁴⁰ Structural variants are rearrangements, deletions, or duplications of segments of the human genome which are believed to be associated with predisposition to conditions such as mental retardation and autism.⁴¹

In order to achieve its scientific goals, the 1000 Genomes Project will sequence the genomes of at least one thousand individuals who gave informed consent for their DNA to be analyzed and placed in public databases.⁴² In fact, the first thousand samples will come from those used for the HapMap and from additional samples in the extended HapMap set.⁴³ These participants are drawn from a geographically wide-ranging group of people of African, East Asian, and European descent.⁴⁴ In addition, an effort began in 2010 to expand data collection to include individuals of Hispanic and African-American descent.⁴⁵ All of the research participants will remain anonymous and will not have any individual medical information collected from them, because the project is intended to catalogue genetic variation rather than to study the disease profile of the individual research participants.⁴⁶ In the future, researchers will be able to use the resource that is developed from the 1000 Genomes Project in order to conduct studies of people affected by various diseases.⁴⁷

The institutions supporting and conducting the genetic sequencing work of the 1000 Genomes Project are the Wellcome Trust Sanger Institute in England, the

³⁴ *Id.*

³⁵ *Id.*

³⁶ *Id.*

³⁷ *Three Sequencing Companies Join*, *supra* note 19.

³⁸ *International Consortium Announces*, *supra* note 13.

³⁹ *Three Sequencing Companies Join*, *supra* note 19.

⁴⁰ *International Consortium Announces*, *supra* note 13.

⁴¹ *Id.*

⁴² *Id.*

⁴³ *Id.*

⁴⁴ *Id.*

⁴⁵ Julia Karow, *1000 Genomes Project to Sequence Nearly 1,000 More Samples by Early 2010; New Samples Collected*, GENOMEWEB (Sept. 8, 2009), <http://www.genomeweb.com/sequencing/1000-genomes-project-sequence-nearly-1000-more-samples-early-2010-new-samples-co>.

⁴⁶ *International Consortium Announces*, *supra* note 13.

⁴⁷ *Id.*

Beijing Genomics Institute in China, and the National Human Genome Research Institute, part of the National Institutes of Health (“NIH”) in the United States.⁴⁸ Most of the funding for the 1000 Genomes Project, which will cost between U.S. \$30 to \$50 million, will come from the existing sequencing budgets of the participating institutes.⁴⁹ In addition, three private U.S. firms that have pioneered development of new sequencing technologies are contributing a significant portion of the sequence data, indicating their belief in the value of the large datasets to be produced and made available through publicly accessible databases.⁵⁰ Indeed, these firms, 454 Life Sciences, a Roche Company; Applied Biosystems, an Applied Biosystems Corp. business; and Illumina Inc., like all other participants in the 1000 Genomes Project, have consented to the open access policies established by the 1000 Genomes Project Steering Committee.⁵¹ These policies include “rapid public release of the data; a prohibition on project participants gaining early access to the data; an intellectual property policy that precludes any participants from controlling the information provided by the project or filing for intellectual property rights in the sequences they produce; regular progress reporting; and coordination of scientific publications with the rest of the consortium.”⁵²

IV. THE DATA ACCESS POLICY OF THE 1000 GENOMES PROJECT

In order to advance research and encourage international scientific collaboration, participants in large-scale, publicly-funded sequencing projects have adopted broad data-sharing policies.⁵³ In February 1996, at the First International Strategy Meeting on Human Genome Sequencing in Bermuda, representatives of laboratories involved in human genome sequencing and of funding agencies unanimously agreed to implement the Bermuda Statement, an international agreement favoring release into the public domain of genetic databases achieved through public funding.⁵⁴ The Bermuda Statement provides that “all human genomic sequence information, generated by centres funded for large scale human sequencing, should be freely available and in the public domain in order to encourage research and development and to maximize its benefit to society.”⁵⁵ One of the fundamental requirements of the Bermuda Statement is the stipulation that sequences longer than 1,000 base pairs must be made freely available to the public, preferably within twenty-four hours of generation.⁵⁶ For the HGP, the highest profile data set rapidly released before publication, the international sequencing centers involved in the

⁴⁸ *Id.*

⁴⁹ Jocelyn Kaiser, *A Plan to Capture Human Diversity in 1000 Genomes*, 319 SCI. 395, 395 (2008).

⁵⁰ *Three Sequencing Companies Join*, *supra* note 19.

⁵¹ *Id.*

⁵² *Id.*

⁵³ Schofield et al., *supra* note 2, at 171.

⁵⁴ See Summary of Principles Agreed at the First International Strategy Meeting on Human Genome Sequencing (held in Bermuda on Feb. 25–28, 1996), http://www.ornl.gov/sci/techresources/Human_Genome/research/bermuda.shtml#1 [hereinafter Bermuda Statement].

⁵⁵ *Id.*

⁵⁶ Bryn Nelson, *Empty Archives*, 461 NATURE 160, 162 (2009).

HGP retained the right to be the first to describe and analyze their complete data sets in peer-reviewed publications, in return for the early release of their data.⁵⁷

According to its proponents, the Bermuda Statement data release policy fosters independent checking of the sequence data by other researchers,⁵⁸ leading to huge gains in life sciences research.⁵⁹ With respect to the HGP, this policy resulted in the publication of new information about thirty genes associated with disease even prior to the publication of the draft sequences.⁶⁰ In addition, a broad data access policy prevents publicly funded large-scale sequencing centers from “establishing a privileged position in the exploitation and control of human sequence information.”⁶¹ What is more, the policy also delivers a symbolic message that “the genome belongs to everybody.”⁶² When sequencers met again in 2003 in Florida, they reaffirmed their commitment to the Bermuda Statement in a statement referred to as the Fort Lauderdale Principles.⁶³ Consequently, organizations such as the National Institutes of Health's National Human Genome Research Institute, the Department of Energy, and the Wellcome Trust require compliance with this policy as a condition of receiving funding.⁶⁴

The 1000 Genomes Project has announced its adherence to the Fort Lauderdale principles on its website.⁶⁵ The site emphasizes that data producers will “release the Project data quickly, prior to publication, in the expectation that they will be valuable for many researchers” and that “data users may use the data for many studies, but are expected to allow the data producers to make the first presentations and to publish the first paper with global analyses of the data.”⁶⁶

In October 2010, the 1000 Genomes Project published in the journal *Nature* an analysis of its completed pilot phase.⁶⁷ After the completion of the pilot phase, the 1000 Genomes Project began full-scale studies, which will require an additional two years.⁶⁸ Data from the pilot studies and the full-scale project are freely available on the project's website, www.1000genomes.org.⁶⁹

The data release policy of the 1000 Genomes Project aims to make the genotype data freely available, while reserving to data producers the right to make the first

⁵⁷ Birney et al., *supra* note 1, at 168.

⁵⁸ David R. Bentley, *Genomic Sequence Information Should Be Released Immediately and Freely in the Public Domain*, 274 SCI. 533, 533 (1996).

⁵⁹ *Id.* at 534.

⁶⁰ Birney et al., *supra* note 1, at 168.

⁶¹ Bermuda Statement, *supra* note 54.

⁶² Eliot Marshall, *Bermuda Rules: Community Spirit, with Teeth*, 291 SCI. 1192, 1192 (2001).

⁶³ See Wellcome Trust, *Sharing Data from Large-Scale Biological Research Projects: A System of Tripartite Responsibility 2* (conference of Jan. 2003), <http://www.genome.gov/Pages/Research/WellcomeReport0303.pdf>.

⁶⁴ Lee Rowen et al., *Publication Rights in the Era of Open Data Release Policies*, 289 SCI. 1881, 1881 (2000).

⁶⁵ *Use of the Project Data, Presentations and Publications, and Authorship*, 1000 GENOMES, <http://www.1000genomes.org/data#DataUse> (last visited Mar. 25, 2011).

⁶⁶ *Id.*

⁶⁷ *1000 Genomes Project Publishes Analysis of Completed Pilot Phase*, NAT'L HUM. GENOME RES. INST. (Oct. 27, 2010), <http://www.genome.gov/27541917>; see *Map of Human Genome Variation*, *supra* note 25.

⁶⁸ *Analysis of Completed Pilot Phase*, *supra* note 67.

⁶⁹ *How to Access 1000 Genomes Data*, 1000 GENOMES, <http://www.1000genomes.org/data#DataAccess> (last visited Mar. 25, 2011).

presentations and to publish the first paper with global analyses of the data.⁷⁰ Examination of the 1000 Genomes Project data release policy is particularly important in light of the recommendation arising from the Fort Lauderdale meeting that rapid prepublication data release be applied to other data sets created primarily as a resource for the scientific community.⁷¹ Experts have observed that data producers are no longer only large genomics centers, and that, increasingly, smaller groups are able to produce a large amount of data.⁷² Although funding agencies anticipated that the policies from the HGP and 1000 Genomes Project would be automatically adopted by these smaller projects, this did not occur.⁷³ For this reason, between eighty to one hundred scientists, ethicists, lawyers, representatives of funding agencies, and journal editors convened at a data release workshop in Toronto in May 2009, hosted by Genome Canada, in order to discuss the extension of such data release policies to fields such as proteomics and biobanking.⁷⁴ The product of that meeting was a set of proposals published in *Nature* in 2009, referred to as “the Toronto Statement,” which aims to establish a set of “best practices” for scientists, both data producers and data users, as well as funding agencies and journal editors.⁷⁵ The Toronto Statement is meant to combat the many factors that inhibit researchers from sharing their data.⁷⁶ Particularly in the realm of small science, these include data producers’ rational desire to be recognized for producing and analyzing data, thereby preserving their ability to win future research grants; the challenge of crafting appropriate publication and intellectual property policies; and the complexity of establishing a suitable infrastructure to house the data.⁷⁷

V. OBSTACLES TO DATA-SHARING AND POTENTIAL SOLUTIONS OFFERED BY THE TORONTO STATEMENT

A. Lack of Recognition of Data Producers

In pursuing career advancement, data producers are understandably reluctant to grant data users access to research results achieved after devoting years of hard work to “planning a project, securing funding and ethics approval, recruiting the participants, collecting the data and materials, performing analyses, managing the collection and infrastructure, controlling technical quality, and generally nourishing the project through waves of funding and maturation.”⁷⁸ In addition, data producers

⁷⁰ *Use of the Project Data*, *supra* note 65.

⁷¹ Wellcome Trust, *supra* note 63, at 2.

⁷² Ciara Curtin, *When to Share*, GENOMEWEB (May 2010), <http://www.genomeweb.com/when-share>.

⁷³ *Id.*

⁷⁴ *Id.*

⁷⁵ Birney et al., *supra* note 1, at 168–70.

⁷⁶ See Donna M. Gitter, *The Challenges of Achieving Open Source Sharing of Biobank Data*, 29 BIOTECH. L. REP. 623, 625 (2010).

⁷⁷ *Id.*

⁷⁸ William W. Lowrance, *Access to Collections of Data and Materials for Health Research* 24 (Mar. 2006), http://www.wellcome.ac.uk/stellent/groups/corporatesite/@msh_grants/documents/web_document/wtx030842.pdf.

fear being “scooped” by data users who may mine data and discover relationships in it that the producer did not discern.⁷⁹ As noted in the Toronto Statement, an emphasis on early data release gives rise to tension between data producers and users.⁸⁰ The former typically wish to publish a first description of a data set, while the latter desire to publish their own analyses of the data.⁸¹ One way to reduce this tension is to ensure proper attribution and recognition for data producers.⁸²

One means of acknowledging the contributions of data producers is to track the usage and citation of data sets through the use of electronic systems similar to those used for traditional publications.⁸³ In order to quantify the usefulness of particular data sets, Professor Boyle has offered the example of a music site associated with Creative Commons,⁸⁴ known as ccMixer.⁸⁵ This site allows users to download music from the site and remix the samples into new tracks, while simultaneously maintaining a record of the credits due to each musician.⁸⁶ A comparable system would permit attribution to the appropriate data producer and also allow universities and funding agencies to track the number of uses, and therefore the value, of a researcher’s data.⁸⁷ This satisfies researchers’ desire for attribution, and assists them in demonstrating their productivity to funding agencies.⁸⁸

Another expert, Myles Axton of *Nature*, has proposed recognizing contributors via “microattribution,” meaning the tracking of researchers’ contributions “down to the very smallest meaningful unit (database record, gene), as opposed to whole papers as is tradition.”⁸⁹ Microattribution could minimize the shortcomings of the Fort Lauderdale Agreement, which Axton has critiqued for failing to reward data producers with credit for uses of their data.⁹⁰ Nor does the Agreement make any provision for crediting data generators, such as genomics “factory” laboratories.⁹¹ In addition, notes Axton, the Fort Lauderdale Agreement is not truly enforceable.⁹²

The Toronto Statement also suggests that funding agencies have a role in fostering data-sharing.⁹³ First, the Statement advises that funding agencies should mandate rapid prepublication data release for projects that “have broad utility, are large in scale, are ‘reference’ in character and typically have community ‘buy-in.’”⁹⁴

⁷⁹ Nelson, *supra* note 56, at 163.

⁸⁰ Birney et al., *supra* note 1, at 169.

⁸¹ *Id.*

⁸² *Id.* at 169–70.

⁸³ Nelson, *supra* note 56, at 163.

⁸⁴ Creative Commons is a nonprofit corporation that aims to facilitate content-sharing in accordance with the law of copyright. See *About Creative Commons*, CREATIVE COMMONS, <http://creativecommons.org/about/> (last visited Mar. 25, 2011).

⁸⁵ Nelson, *supra* note 56, at 163 (citing Professor James Boyle of Duke Law School).

⁸⁶ See *About*, ccMixer, <http://ccmixter.org/about> (last visited Mar. 25, 2011).

⁸⁷ Nelson, *supra* note 56, at 163.

⁸⁸ *Id.*

⁸⁹ Gudmundur A. Thorisson, & Anthony J. Brookes, *Meeting Minutes, IRBW2009 Workshop of May 13–14, 2009*, at 5 (Aug. 27, 2009), <http://www.gen2phen.org/system/files/private/IRBW2009%20meeting%20minutes%20v2.pdf>.

⁹⁰ *Id.*

⁹¹ *Id.*

⁹² *Id.*

⁹³ Birney et al., *supra* note 1, at 169.

⁹⁴ *Id.*

Second, according to the Toronto Statement, funding agencies should make explicit the requirement for prepublication data release and also become actively involved throughout the data release process, as demonstrated by the International HapMap Project and the 1000 Genomes Project, among others.⁹⁵ Among the essential roles for funding agencies is participation in the development of effective consent, security, access, and governance mechanisms that protect research participants while encouraging prepublication data release.⁹⁶

Third, for projects generating large data sets, the Toronto Statement advises that funding agencies require that grant applications include data-sharing plans and undergo peer review, and notes that this is not commonly the case at the present time.⁹⁷ Experts have recommended, in the context of post-publication data release (and it is equally applicable prepublication) that, where they do not yet exist, clear criteria should be developed for reviewers of grants to help them evaluate applicants' plans for sharing data and material.⁹⁸ Among examples of good practice cited in this regard are the NIH,⁹⁹ the U.K. Wellcome Trust,¹⁰⁰ and the U.K. Medical Research Council ("MRC"),¹⁰¹ a government-funded research agency akin to the NIH.¹⁰² Commentators note that, for these agencies, "[d]ata-sharing plans are required in proposals, efforts are made to facilitate sharing, such as putting investigators in touch with repositories and, for some organizations, compliance is an important consideration in funding renewal."¹⁰³

Fourth, in terms of compliance, funding agencies should clearly express the rewards for honoring, as well as the sanctions for ignoring, their policies relating to prepublication data release, and also consistently enforce such policies.¹⁰⁴ One possibility is to condition grant renewals and promotions to some degree on the extent to which the investigator shares data.¹⁰⁵

B. The Challenge of Crafting Appropriate Publication Policies

In light of researchers' desires to protect the data they have gathered and advance their own careers, efforts to create a culture of data-sharing are inextricably linked to the question of the optimal publication policies to support data-sharing

⁹⁵ *Id.*

⁹⁶ *Id.*

⁹⁷ *Id.*

⁹⁸ Schofield et al., *supra* note 2, at 172.

⁹⁹ See generally *Final NIH Statement on Sharing Research Data*, NAT'L. INSTS. HEALTH (Feb. 26, 2003), <http://grants.nih.gov/grants/guide/notice-files/NOT-OD-03-032.html>.

¹⁰⁰ See generally *Policy on Data Management and Sharing*, WELLCOME TRUST (Aug. 2010), <http://www.wellcome.ac.uk/About-us/Policy/Policy-and-position-statements/wtx035043.htm>.

¹⁰¹ See generally *MRC Policy on Data Sharing and Preservation*, MED. RES. COUNCIL, <http://www.mrc.ac.uk/Ourresearch/Ethicsresearchguidance/Datasharinginitiative/Policy/index.htm> (last visited Mar. 25, 2011).

¹⁰² *About Us*, MED. RES. COUNCIL, <http://www.mrc.ac.uk/About/Structure/index.htm> (last visited Mar. 25, 2011).

¹⁰³ Schofield et al., *supra* note 2, at 172.

¹⁰⁴ *Id.*

¹⁰⁵ Birney et al., *supra* note 1, at 169.

efforts.¹⁰⁶ According to the Toronto Statement, although many scholarly works have been published by third parties reporting research findings based upon data sets released before publication, few data producers were affected by these publications that pre-dated their own.¹⁰⁷ Nevertheless, the Toronto Statement proposes that “this ongoing concern is best addressed by fostering a scientific culture that encourages transparent and explicit cooperation among data producers, data analysts, reviewers, and journal editors.”¹⁰⁸

One way of achieving such cooperation is to permit investigators who contribute data to enjoy some period of exclusivity during which they will have the sole right to publish analyses of the data.¹⁰⁹ The Toronto Statement declares that data producers should, as early as possible, and ideally before large-scale data generation begins, produce a citable statement or “marker paper” in which they describe the data set and their intentions in respect of analysis and publication.¹¹⁰ This statement should include “clear details about the data set to be produced, the associated metadata, the experimental design, pilot data, data standards, security, quality-control procedures, expected timelines, data release mechanisms, and contact details for lead investigators.”¹¹¹ In the event that data producers request a protected time period to allow them to be the first to publish the data set, data producers should expect this period to be used only for the publication of a global analysis of the data and to expire within one year.¹¹²

The NIH and the U.K. Medical Research Council have both adopted a period of data exclusivity for data producers.¹¹³ The NIH Policy for Sharing of Data Obtained in NIH Supported or Conducted Genome-Wide Association Studies (“GWAS”) declares that investigators who contribute data to a NIH GWAS data repository will retain the exclusive right to publish analyses of the dataset for a maximum of twelve months following its release via the NIH GWAS data repository.¹¹⁴ During this period of exclusivity, the NIH grants data access through data access committees to other investigators, who may analyze the data, but not submit for publication their analyses or conclusions, until the expiration of the exclusivity period.¹¹⁵ The NIH also expects all investigators who access GWAS datasets to acknowledge in all publications the data producers who conducted the original study, along with the funding organization that supported the work and the NIH GWAS data repository.¹¹⁶ Similarly, the U.K. Medical Research Council provides for a period of interim exclusivity for data producers.¹¹⁷ The MRC data-sharing policy provides that “[a] limited, defined period of exclusive use of data for primary research is reasonable.”¹¹⁸

¹⁰⁶ *Id.* at 169–70.

¹⁰⁷ *Id.*

¹⁰⁸ *Id.* at 169–70.

¹⁰⁹ *Id.* at 170.

¹¹⁰ *Id.*

¹¹¹ *Id.*

¹¹² *Id.*

¹¹³ *MRC Policy, supra* note 101; *see also* Policy for Sharing of Data Obtained in NIH Supported or Conducted Genome-Wide Association Studies, 72 Fed. Reg. 49,290, 49,293 (Aug. 28, 2007).

¹¹⁴ *See* Policy for Sharing, 72 Fed. Reg. at 49,296.

¹¹⁵ *Id.*

¹¹⁶ *Id.*

¹¹⁷ Lowrance, *supra* note 78, at 25.

¹¹⁸ *MRC Policy, supra* note 101.

The U.K. Biobank follows suit, asking researchers to recruit 500,000 participants and collect their data and samples, in return for which the researchers will gain exclusive access for a period of time or for research into certain specialty areas.¹¹⁹ In general, the purpose of the period of exclusivity is to provide the data producer sufficient time to produce, organize, document, verify, and analyze the data in preparation for publication in a scientific journal.¹²⁰ After publication, the data is made available more broadly, either via a publicly available database or via an application system.¹²¹

The Toronto Statement emphasizes that data users should honor the “scientific etiquette” that allows data producers to publish the first global analyses of their data set.¹²² Participants at the Toronto meeting urged, in cases where the publication plans of data users overlap the data producers’ proposed analyses, that the data users approach the producers with the goal of “creating a mutually agreeable publication schedule, such as co-publication or inclusion within a set of companion papers.”¹²³ In addition, data users should acknowledge any data or materials used and the originating sources.¹²⁴ Some of the potential methods to achieve this include the addition of metadata tags linking to data and bioresources, or a digital object identifier for resources in public repositories, which would facilitate the searching of the literature for specific bioresources.¹²⁵ These approaches also offer appropriate recognition of data producers who comply with data release and deposition policies, making it possible to offer incentives such as funding contingent upon such sharing.¹²⁶

Like data users, journal editors and reviewers also have an important role in the data-sharing process.¹²⁷ Participants at the Toronto meeting recommended that journals engage actively in the dialogue about rapid prepublication data release, both in their formal guide to authors and informal instructions to reviewers.¹²⁸ Journal editors should work in tandem with reviewers to consider carefully the specific policies regarding citation and use of the data that may be associated with certain large-scale data sets.¹²⁹ This will help to raise the quality of analysis and also promote fairness in citation of published studies.¹³⁰

Even a decade before the Bermuda principles, the editors of journals, such as *Nucleic Acids Research*, fostered the early development of GenBank and other genomic repositories by requiring researchers to deposit their data there as a precondition for publishing.¹³¹ Newer journals, such as the open-access Public Library of Science journals, have made publication contingent on making the data “freely available without restriction, provided that appropriate attribution is given

¹¹⁹ Lowrance, *supra* note 78, at 25.

¹²⁰ See Reichman & Uhler, *supra* note 6, at 333.

¹²¹ *Id.*

¹²² Birney et al., *supra* note 1, at 169.

¹²³ *Id.* at 170.

¹²⁴ *Id.*

¹²⁵ *Id.* at 169.

¹²⁶ Schofield et al., *supra* note 2, at 171–72.

¹²⁷ Birney et al., *supra* note 1, at 170.

¹²⁸ *Id.*

¹²⁹ *Id.*

¹³⁰ *Id.*

¹³¹ Nelson, *supra* note 56, at 162–63.

and that suitable mechanisms exist for sharing the data used in a manuscript.”¹³² *Nature* journals require authors to “make materials, data and associated protocols promptly available to readers without preconditions.”¹³³

Despite the increased call in many journals for data-sharing, the ability of journals to compel such behavior is limited.¹³⁴ In March 2009, for example, the journal *Epidemiology* declared that “We invite our authors to share their data and computer code when the burden is minimal,” according to an editorial in that issue.¹³⁵ Miguel Hernán, an epidemiologist at Harvard University and a co-author of that editorial, acknowledged the trend toward data-sharing but contended that forcing a sharing requirement on authors “would be suicidal”, especially with unresolved concerns over patient confidentiality, and would likely cause authors to submit their papers elsewhere.¹³⁶ Moreover, many journals have no written policy on the availability of either resources or primary data.¹³⁷

C. The Complexity of Establishing a Suitable Infrastructure

The Toronto Statement also recommends that funding agencies offer infrastructural support by funding the creation and long-term maintenance of databases.¹³⁸ Among the costs associated with creating a suitable infrastructure to house and maintain the data and any associated physical specimens are: the maintenance of physical premises; the development of appropriate information technology; the preparation of data for storage or archiving, such as anonymising the data and documenting the variables; salaries for administrators, managers, and staff; and the creation and maintenance of an accessible database.¹³⁹ These costs are all ongoing ones.¹⁴⁰ Because granting agencies understandably focus primarily on research, they therefore frequently fail to invest in the infrastructural support necessary to support their archiving requirements, in a sort of “tragedy of the commons.”¹⁴¹

The internationally collaborative nature of biomedical research exacerbates this problem by allowing each funding agency to leave the problem to another agency. Yet overcoming the infrastructure issue is crucial in order to optimize the benefits of technological progress in the life sciences.¹⁴² Taking human genome sequencing as one example, one source estimates that the cost of storing all known DNA sequence information in openly accessible databases costs less than one percent of the sum necessary to generate such sequence data.¹⁴³

¹³² *Id.* at 163.

¹³³ *Id.*

¹³⁴ *Id.*

¹³⁵ *Id.*

¹³⁶ *Id.*

¹³⁷ Schofield et al., *supra* note 2, at 171.

¹³⁸ Birney et al., *supra* note 1, at 169.

¹³⁹ Gitter, *supra* note 76.

¹⁴⁰ *Id.*

¹⁴¹ Nelson, *supra* note 29, at 162.

¹⁴² *Id.*

¹⁴³ EUR. LIFE SCI. INFRASTRUCTURE FOR BIOLOGICAL INFO., ELIXIR: DATA FOR LIFE, http://www.elixir-europe.org/bcms/elixir/Documents/Elixir_brochure.pdf (last visited Mar. 25, 2011).

VI. THE NEED FOR RESTRICTED ACCESS DATABASES

Notwithstanding the prescriptions in the Toronto Statement for encouraging data-sharing, not all data is suitable for an open access approach.¹⁴⁴ Indeed, samples and data that retain their link with an individual research participant must be protected with additional procedures.¹⁴⁵ Examples of such data include personal information such as individual-level genotypic and phenotypic data; exposure to drugs and environmental factors; and pedigree data, including information about familial relationships, along with analyses of such data.¹⁴⁶ Experts have offered several technical and policy approaches to preserve patient privacy.¹⁴⁷ They recommend, *inter alia*, the establishment of policies to assess credentials of data users; the execution of clear contracts with data users that define the appropriate use of data; formalization of liability rules for misuse of data; and the use of a technical data management approach that increases the number of research participants whose data will be aggregated where the data is considered more sensitive.¹⁴⁸ Experts note that subjecting data to access controls erects a barrier that limits the number of researchers who will reuse the resources.¹⁴⁹ What is more, researchers may decline to abide by overly complicated data access agreements.¹⁵⁰

Nevertheless, restricted access databases are also useful in terms of building upon the collaborations already established by scientific researchers, particularly those involved in small science.¹⁵¹ One such biobank is the U.K. DNA Banking Network (“UDBN”), which is a secondary biobank, meaning that it aggregates and manages tissue samples and associated data gathered by clinicians who gather the samples in the course of studying particular diseases.¹⁵²

In order to establish UDBN, the U.K. Medical Research Council offered grants to centers that housed DNA collections.¹⁵³ These awards required the collections to be maintained as “shared national resources,” and “made available to collaborators.”¹⁵⁴ In addition, awardees were “required to transfer a portion of each sample” to the network and “to add any genotype data they obtain to the common database.”¹⁵⁵

In granting access to its collections, the UDBN explicitly rejected the unrestricted “open access” model, instead developing a “fair access” regime, which

¹⁴⁴ Policy for Sharing of Data Obtained in NIH Supported or Conducted Genome-Wide Association Studies, 72 Fed. Reg. 49,290, 49,291 (Aug. 28, 2007).

¹⁴⁵ *Id.* at 49,292.

¹⁴⁶ *Id.* at 49,293.

¹⁴⁷ See Bradley Malin et al., *Technical and Policy Approaches to Balancing Patient Privacy and Data Sharing in Clinical and Translational Research*, 58 J. INVESTIGATIVE MED. 11, 15–17 (2010).

¹⁴⁸ *Id.*

¹⁴⁹ *Id.*

¹⁵⁰ Presentation at the 2010 Advances in Genomics Symposium, Paul Flicek, Vertebrate Genomics, Challenges for the Data Management and Analysis of Large-Scale Human Genome Sequencing, <http://www.advances-in-genomics.org/presentations/Flicek.pdf> (last visited Mar. 25, 2011).

¹⁵¹ See generally Martin Yuille et al., *The UK DNA Banking Network: A “Fair Access” Biobank*, 11 CELL TISSUE BANK 241–51 (2010), available at <http://www.springerlink.com/content/11082h68g0645517/fulltext.pdf>.

¹⁵² *Id.* at 241.

¹⁵³ *Id.* at 242–43.

¹⁵⁴ *Id.* at 243.

¹⁵⁵ *Id.*

derives inspiration from the 2003 United Nations Educational Scientific and Cultural Organization (“UNESCO”) International Declaration on Human Genetic Data.¹⁵⁶ The UDBN distributes data via the project website to third party researchers, who apply online for registration.¹⁵⁷ After verifying the researcher’s credentials, UDBN grants access to a restricted area of the website.¹⁵⁸ The third party researchers can then communicate online with the data collectors in order to negotiate collaboration.¹⁵⁹ If a collaborative relationship is successfully negotiated, UDBN permits the data collector to grant the third party access to fuller data.¹⁶⁰

The fair access model builds upon researchers’ willingness to share data with collaborators, and concomitantly denies access to non-collaborators.¹⁶¹ This approach also accepts a collector’s right to “exclusive access to his/her collection for the purposes of the investigational goals stated in the initial collection proposal,” recognizing that granting a “first mover” advantage is likely to motivate scientific discovery.¹⁶² This model also acknowledges that tensions frequently arise when a potential collaborator requests data access from a data producer.¹⁶³ Thus, the restricted access model provides a necessary complement to the open access databank, particularly in the realm of “small science.”

¹⁵⁶ *Id.*

¹⁵⁷ *Id.*

¹⁵⁸ *Id.* at 246.

¹⁵⁹ *Id.*

¹⁶⁰ *Id.*

¹⁶¹ *Id.* at 244–46.

¹⁶² *Id.*

¹⁶³ *Id.*