



Responsible, practical genomic data sharing that accelerates research

James Brian Byrd¹, Anna C. Greene², Deepashree Venkatesh Prasad³,
Xiaoqian Jiang⁴ and Casey S. Greene^{3,5}✉

Abstract | Data sharing anchors reproducible science, but expectations and best practices are often nebulous. Communities of funders, researchers and publishers continue to grapple with what should be required or encouraged. To illuminate the rationales for sharing data, the technical challenges and the social and cultural challenges, we consider the stakeholders in the scientific enterprise. In biomedical research, participants are key among those stakeholders. Ethical sharing requires considering both the value of research efforts and the privacy costs for participants. We discuss current best practices for various types of genomic data, as well as opportunities to promote ethical data sharing that accelerates science by aligning incentives.

Metadata

The data that describe the data. For genomic samples, this could be how the sample was processed, the platform that was used to assay the sample, characteristics about the conditions in which the sample was obtained or any other elements that provide context to the genomic data in question.

¹Department of Internal Medicine, Medical School, University of Michigan, Ann Arbor, MI, USA.

²Alex's Lemonade Stand Foundation, Bala Cynwyd, PA, USA.

³Childhood Cancer Data Lab, Alex's Lemonade Stand Foundation, Philadelphia, PA, USA.

⁴School of Biomedical Informatics, The University of Texas Health Science Center at Houston, Houston, TX, USA.

⁵Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA, USA.

✉e-mail: greenescientist@gmail.com

<https://doi.org/10.1038/s41576-020-0257-5>

Genomics has a robust culture of data sharing. We are now nearing the two-decade mark of strong expectations for sharing genome-wide transcriptomic assays and associated metadata¹. This wealth of data has enabled new approaches that rely on the analysis of very large collections of public data by investigators who were not involved in the original data collection^{2–7}. It is also possible to assay genotypes^{8,9}, methylation¹⁰ and many other features of a sample at a genome-wide level, which presents considerable opportunities for secondary analysis.

With proof-of-concept studies showing the potential to uniquely identify an individual in ever-widening types of detailed data sets, the sharing process has become murkier^{11,12}. As expression profiling has switched from array-based profiling to sequencing-based profiling, the reidentification risk from human-derived samples has also increased^{13–15}. For genetic data, the risk of reidentification has led to controlled-access sharing, which is mediated via services such as the [database of Genotypes and Phenotypes \(dbGaP\)](#)¹⁶. However, genotype-related data that contain aggregated estimates, such as variant-level association statistics, pose some risk that individuals could be reidentified¹⁷.

Investigators, funders and other stakeholders supporting responsible data sharing must consider both the risks and benefits to participants as well as other individuals who could be affected positively or negatively by sharing a research data set in different ways. In addition to ethical concerns, it is important to consider the impact of data sharing practices on the overall research ecosystem. Genomic profiling technologies are now ubiquitously available and are becoming widely used in fields with different cultures of sharing. Funders and publishers must balance multiple considerations to develop appropriate policies. For example, adding data

sharing requirements, particularly as unfunded mandates, could hamper the establishment of a pro-sharing culture by creating resentment around reuse¹⁸. However, early genomic scientists recognized the potential for high-dimensional profiling to lead to irreproducible results and spurious findings if source data were not shared¹⁹. Funders and publishers ultimately must take steps to foster a robust, responsible data sharing culture to support rigorous research with high-dimensional genomic profiling technologies.

Investigators who have shared data well increase the impact of their research: publications linked to a data repository or persistent identifier are more cited²⁰. In this Review, we first outline types of data, metadata and frameworks for sharing. We next describe the steps that researchers can take to assess risks and responsibly share data derived through genome-wide profiling technologies. We discuss the rationale for specific data sharing practices. For some data types, there are no widely recognized single point-of-truth repositories, and these principles can guide researchers' current decision-making. For data types with widely used repositories, we provide more detailed guidance. We extensively cover privacy challenges posed by individual-level data derived from human samples because these data pose the most substantial challenges, but we recognize that many types of genomic data, such as those derived from model organisms, pose little to no risk and should be publicly shared in appropriate repositories. Although we focus on genomic profiling, the underlying principles apply to other data-intensive research projects as well. We also note the roles that other stakeholders, including funders and publishers, can play in the process to enhance the pace of discovery, ultimately helping patients. We identify practical changes that could better

High-dimensional profiling

Assays of samples that produce many measurements for each sample. Genomic profiling technologies are high-dimensional ones. For example, assaying the expression level of all protein-coding genes in the genome characterizes each sample in approximately 20,000 dimensions. Genotyping of single-nucleotide polymorphisms can produce more than one million dimensions for each human sample.

Single point-of-truth repositories

Repositories designed to store the archival form of a data set and that assign a unique identifier. Investigators are responsible for all aspects of data provenance until data are put into a single point-of-truth repository, at which point the repository becomes responsible for these.

Intermediate data

Results between raw data and the desired final representation for reporting. For example, in an analysis to identify differentially expressed pathways from RNA sequencing reads, gene expression estimates and differential expression p values could both be considered intermediate results.

align researcher incentives and support the efficient enforcement of sharing for valuable research products.

What are research data?

Research data in genomics are of many different classes and types. We can divide data by the types of biomolecules that they represent. For example, certain assays measure RNA in a sample²¹, and others measure DNA²², protein²³, or metabolite²⁴ content. We can also divide data by the type of measurement technology used to gather them. For example, RNA assays could be based on microarray or sequencing profiling²⁵. A sample itself could be derived from a single organism or many²⁶: it could be a cell line with a treatment²⁷, a human tissue sample²⁸ or a population of organisms gathered from an ocean location²⁹. For the purposes of this Review, we consider genomic data to be those that include the potential to profile the genes or gene products for most of an organism's genes or a collection of organisms' genes.

We also consider derived data that are intermediate between the raw data produced by an instrument and a finding to be research data. In the terminology used in this Review, we consider read files produced by an RNA sequencing (RNA-seq) experiment and represented in FASTQ format³⁰ to be raw data, we consider gene expression estimates to be intermediate data and we consider the findings to be plots, figures and underlying statistics produced by analysis of the gene expression data. There could be multiple intermediate data representations between raw data and a finding. Researchers sequencing paired tumour and normal samples to identify somatic and germline variants would be likely to produce FASTQ files for each tumour and normal sample, variant call format (VCF) files for each sample, separate mutation annotation format (MAF) files for the germline and somatic variants, and, finally, summary results and figures. In this case there could be hundreds of intermediate VCF files and two separate MAF files between the raw data and the findings. We provide recommendations for how investigators can select which items from raw data to findings should be archived and how they can best be shared.

An increasingly common type of derived data is a model produced by machine-learning methods applied to genomic data. Researchers can download publicly available data or process data associated with their study, analyse those data with neural networks^{31,32} or other approaches^{33,34}, and then use those models to either infer something about the biological system that generated the data⁶, to better understand the methods themselves³⁵ or to develop a deeper understanding of a related disease or process⁷. Machine-learning models can often be repurposed in much the same way as underlying data. For example, Gulshan et al.³⁶ took a model trained on generic images and fine-tuned it to detect diabetic retinopathy. In genomics, Kelley et al.³⁷ demonstrated that a model trained on a collection of data from certain cell types could quickly and accurately be adapted to a new cell type. Because machine-learning models are executable, they can also be automatically tested³⁸. New repositories, such as Kipoi, have been designed to support and automatically test such models³⁹, providing

downstream researchers with a library of working models.

Throughout this Review we maintain these distinctions between raw data, intermediate data and findings, and provide specific sharing recommendations for each. We also discuss why certain data are more or less likely to identify a study participant and how sharing is controlled for certain high-risk data. In the interests of providing a review that is as broadly applicable as possible, we also describe the principles that underlie specific recommendations. For data modalities that are either not discussed within this Review or that are developed in the future, we expect that these principles can be applied to develop an appropriate sharing plan.

What are research metadata?

Research metadata are the data that describe research data. If a biospecimen's genomic sequence data are represented in raw form by a FASTQ file, information about the biospecimen is metadata. This could include, for example, a coded identifier, the tissue from which a biospecimen was taken, information about the handling of the biospecimen, information extracted from an electronic health record describing the individual from which the biospecimen was taken and more.

We divide our consideration of research metadata into information about the subject of study, which we term 'sample metadata', and information about a sample's handling and processing, which we term 'handling metadata'. This framing is aligned with how the influential minimum information about a microarray experiment (MIAME)¹ recommendations can be applied to non-microarray settings. It is also aligned with how these types of resources are represented in major databases: for example, in the BioSample database, frequently reused biospecimens such as cell lines or references are designated a single, reusable identifier with additional sample metadata⁴⁰. For derived data, the sample's metadata would often remain unchanged whereas the handling metadata would differ based on the computational processing steps; however, this distinction begins to blur for intermediate forms that integrate multiple samples, such as machine-learning models.

Metadata are provided with a level of detail that can be high or low. The fields that are included as metadata can enhance or reduce the level of detail. For example, a hypothetical sample⁴¹ could be described as 'tumour' or as 'tumour from an 18-month-old male'. The latter has additional age and gender information, which are akin to additional fields. The level of specificity for each field also affects the level of detail of the metadata: the same sample could be described as 'malignant peripheral nerve sheath tumour from an 18-month-old male'.

Metadata can be structured or unstructured. Structured metadata could be represented as a tab-delimited text file containing a unique identifier and experimental factor ontology (EFO)⁴² terms relevant to a sample or its handling. Unstructured metadata could be a paragraph in a manuscript describing the experiment. In our example above, derived from Kudesia⁴¹, 'malignant peripheral nerve sheath tumour from an 18-month-old male' is an unstructured description of a sample. Databases

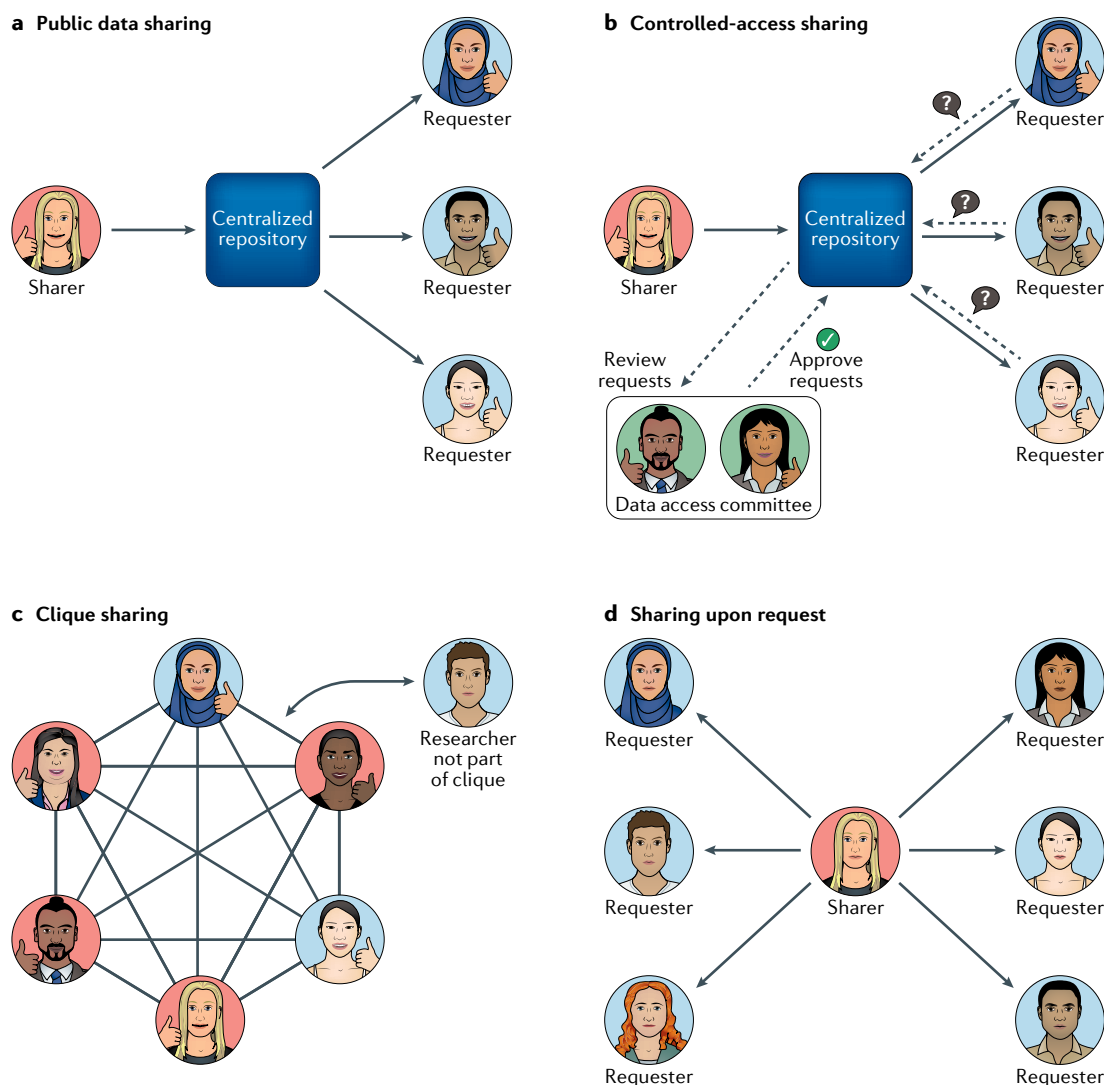


Fig. 1 | Diverse types of data sharing. **a** | In public sharing, researchers make data broadly available without restrictions on use. **b** | In controlled-access data sharing, researchers place some conditions on access and reuse but ideally do not discriminate on the basis of individual projects or proposers. **c** | Clique sharing occurs when researchers form consortia and share within the consortia (data sharers within the clique are shown with a pink background) but have very restrictive policies for external sharing, which hampers engagement and impact. **d** | Sharing upon request offloads the burden for negotiating sharing procedures until there is demand, but when there are multiple requests this approach can become time consuming to manage.

designed to store research data often include fields that allow highly structured information, such as ontology terms that apply to a sample, to be provided alongside fields that are relatively unstructured. For example, the EFO term for malignant peripheral nerve sheath tumour is EFO:0000760, age is EFO:0000246 and male, which is included in EFO from the phenotype and trait ontology (PATO), is term PATO:0000384. The meta-data that describe most repository-stored genomic data are available with some structured and other relatively unstructured elements.

How are data shared?

Genomic data are shared in many ways. We distinguish between public, controlled-access, clique and upon-request sharing approaches (FIG. 1). Data are also

shared on many different platforms, from investigator-specific solutions to those that are purpose-built for a data type to general-purpose repositories that support many data types.

Public data sharing (FIG. 1a) occurs when data are released for reuse without barriers (beyond any applicable ethical considerations and laws, with which the user is expected to be familiar). This level provides the lowest barrier to entry for reuse as researchers can probe the data to gain an understanding of its characteristics. Public data sharing combined with detailed sample and handling metadata can allow researchers to answer numerous questions. The Cancer Genome Atlas (TCGA) data set provides somatic mutation, gene expression estimates, a limited set of clinical metadata and certain other profiling information, which were made available

in a fully-open form and available for publication by anyone after an embargo period⁴³. TCGA has become a remarkably successful example of a public, reusable data resource laying the groundwork for numerous discoveries⁴⁴. At a smaller scale individually — although covering more biological samples — microarray gene expression data sets are also publicly shared in data type-specific repositories, such as *ArrayExpress*⁴⁵ and the *Gene Expression Omnibus (GEO)*⁴⁶.

Controlled-access sharing (FIG. 1b) occurs when data are available for reuse if some fixed criteria are met. These criteria may include a review of protocols, a commitment to use data only for health-related research or other elements that affect how one obtains and uses the data but are not applied differently to different requesters. This level usually provides a modest barrier to entry for reuse efforts and is currently the favoured approach for de-identified genomics data that pose significant reidentification concerns. We discuss such data sets as high risk. The UK Biobank⁴⁷ is an example of a resource that is made available under such criteria. A similar effort is underway in the United States via the All of Us project⁴⁸. Making data sets available in this way allows data set developers to confirm that adequate oversight structures are in place for research that could potentially lead to reidentification of a study participant.

Clique sharing (FIG. 1c) and sharing upon request (FIG. 1d) occur when investigators join a consortium or make individual arrangements to share data. These mechanisms place substantial burdens on data requesters, and those within the clique or who hold data sets can select which requesters will be disadvantaged. Data ostensibly made available upon request are not widely shared in practice^{49,50}. In these cases, the data sharing decisions at each point come down to individual scientists. There may be a mismatch between researchers' perceptions of their own sharing behaviour and their practices. Even when the commitment to share is strong, failure to quickly deposit data in a repository may degrade the investigator's ability to share as personnel come and go from the laboratory, as data are likely to be managed less reliably than they would be in established repositories. Earlier-career scientists report being the most enthusiastic about sharing, and senior researchers report the most reticence⁵¹. In the same survey, early-career researchers report worse sharing behaviours than more senior ones⁵¹; however, Campbell et al.⁵² made data requests and found better sharing behaviours among early-career scientists. These seemingly contradictory results suggest that early-career researchers may hold themselves to a higher standard for sharing. For the purposes of this Review on behaviours supporting an ecosystem that accelerates discovery, we focus on public or controlled-access sharing because of the considerable limitations of clique-based and request-based approaches.

Although the type of sharing influences the extent to which sharing efforts will enhance the impact of the work, it is not the sole factor. For example, Learned et al.⁵³ describe efforts to access and compile publicly available genomic data into a reusable resource for the paediatric cancer community. Even among public data,

the authors found barriers to using some of the data: samples that were mislabelled, purportedly uploaded data that were missing or, in certain cases, a requirement that they would have to use a proprietary cloud platform for analysis at a substantial cost. In subsequent sections we describe potential risks as well as principles and practices that can help investigators maximize the impact of their data through effective sharing.

Data have variable levels of risk

Although we focus a considerable amount of attention in this Review on the risks associated with sharing certain data, in many cases sharing data poses little to no risk. Many experiments involve genomic assays of model organisms, cell lines, environmental samples or agricultural subjects. In other cases, the measurement technology may not be capable of revealing individual characteristics or the assay may provide information that is transient and thus poses little risk. Other data clearly identify the individual from which the data were derived, either through the data themselves or the metadata that describe them.

Data that accurately describe a person for long periods of time typically carry a greater privacy risk compared with information that is only transiently true. For example, the sequence of our genome is with us for our lifetime whereas triglyceride levels may fluctuate with fasting. The risk of reidentification is also related to the extent to which the data modality uniquely identifies individuals. The idea of an equivalence class can help to develop an intuitive understanding of risk: consider an equivalence class to be the number of people for whom a set of values would be true. A measure of the risk of reidentification, given those values, can be considered to be $1 / (\text{number of people in that equivalence class})$ ^{54,55}. In general, the richer the data elements, the smaller the equivalence classes. Transformations of the data can alter the size of equivalence classes; using the decade of life rather than age increases the size of many equivalence classes. However, the effect is not uniform across the data set: equivalence classes can remain very small for those at the extremes of age. Although it is not possible to exhaustively enumerate data types and their associated risk levels, we provide certain examples (TABLE 1) and a fuller discussion of risk levels in the following subsections.

Other types of data encountered in genomic research could also pose risks when shared for reasons other than identification. Certain data, such as the genome sequences of particular pathogens, could pose biosafety concerns. Data that inadvertently disclose the location of endangered species could facilitate poaching. We expect these cases to be rare. In the absence of a clear overriding concern of this type, data not derived from participants should be considered low risk.

Genomic variants are one path to risk. Certain types of genomic data, such as those directly assaying numerous variants across the genome, cannot be de-identified. For other data types, de-identification can be attempted but may not succeed, and as with other data types the key points to consider are the duration and uniqueness.

Table 1 | Genomic data types and levels of risk

Data type	Usual risk level	Sharing to minimize risk
RNA-seq reads of model organisms	None	Public access
Whole-genome sequencing reads of endangered species	Usually none, although location metadata could put species at risk	Public data but controlled-access metadata
RNA-seq reads of human tissue samples	High	Public gene expression estimates Controlled-access for sequencing reads
Whole-exome sequencing reads of cancerous tissue samples	High	Public access for somatic variant data, but controlled access for germline variant data Potential summary-level queries of germline variants
Exome sequencing of non-cancerous human tissue samples	High	Public summary-level information aggregated across many individuals
High-density DNA methylation array of human tissue	High	Remove data from probes that contain common variants before public sharing Controlled access for full data set

RNA-seq, RNA sequencing.

Certain types of genomic data are designed to reveal many of an individual's genetic variants: whole-genome sequencing, high-density genotyping array profiling and whole-exome sequencing. Germline genetic variants accurately describe a person for long periods of time and, with modest numbers of variants, produce very small equivalence classes. Genomic data sharing beacons were an attempt to share only limited, summary-level genomic information to control risks, but equivalence classes are so small that querying beacons for modest numbers of variants was often sufficient to reidentify an individual as a member of a beacon⁵⁶. The clearest avenue to risk is with high-density germline variant calls⁵⁷. Even noisy variant information can be readily cross-referenced with study participants to reidentify an individual⁵⁸. In addition, systems for storing genomic data have, at times, permitted queries of the database using uploaded sequences. Such systems make it possible to find individuals related to an unknown person, given that unknown person's DNA sequence. Law-enforcement entities have used these systems to solve previously unsolved cases, including that of the Golden State Killer⁵⁹. One database has sought to use an opt-in preference from data contributors to control what can be searched; however, a court has recently ruled that with a search warrant, a police agency can search this database without regard for the opt-in preference of the data contributors⁶⁰. The extent to which data can be accessed and obtained in this manner depends greatly on the legal jurisdictions that apply.

Sequencing-based assays are one avenue of risk. Sequencing-based assays can reveal the genetic variants that characterize an individual, even if that was not an intended portion of the experiment. Sequencing cancer genomes with the goal of identifying somatic variants reveals both germline and somatic variants. Even if the goal is to simply measure gene expression with RNA-seq, an experiment on normal human tissue that captures a large fraction of messenger RNA and long non-coding RNAs with high sequencing depth is likely to contain sufficient sequencing depth to call genetic

variants¹³⁻¹⁵. The RNA isolation strategy and sequencing depth will affect the windows of the genome in which variants could be revealed. For certain body sites, a substantial fraction of metagenomic reads intended to measure our microbiomes align to a human reference⁶¹. On the other hand, highly targeted sequencing technologies may assay only small portions of the genome. The key question in each case is whether or not the technology reveals enough variants to identify an individual¹².

Array-based assays can also reveal genetic variants. This can occur intentionally: single-nucleotide polymorphism (SNP) genotyping arrays are specifically designed to capture inter-individual differences in the allele present at a locus. This can also occur unintentionally: DNA methylation profiling with dense arrays can reveal genotypes at roughly 1,000 loci⁶², those for which some people have genetic variants that directly overlap with the profiled positions. In many cases, data from microarray-based transcriptomic profiling technologies are currently considered low risk. For any array-based technology, the more of the genome that is assayed and the more sensitive probes are to short mismatches, the more risk there is of revealing genomic variants.

Especially in the case of genomic data, the probability of reidentification is not static over time and changes based on what other resources are available. Genetic measurements of many individuals provide sufficient information to design artificial queries against data resources that could reveal alleles of interest⁶³. As our understanding of the interrelatedness of genotype and molecular phenotypes grows, it will become easier to identify alleles that underlie high-dimensional data that do not directly measure genotypes⁶⁴. As more data are made available, it becomes easier to find individuals who are closely enough related to a target individual to identify that participant. The observation that certain genetic variants affect gene expression has led to reports of a related risk for gene expression microarray data, but the accuracy of the imputed genotypes is currently relatively low⁶⁴. We find the considerations in [NOT-OD-19-023](#) from the US National Institutes of

Direct identifiers

Information that is replicable, distinguishable and knowable, and that can identify individuals uniquely.

HIPAA privacy rule

The standards for privacy of individually identifiable health information introduced in the Health Insurance Portability and Accountability Act (HIPAA) of 1996. The rule introduces the concepts of expert determination and 'safe harbour' as a means of de-identifying data.

De-identification

As defined by the Health Insurance Portability and Accountability Act (HIPAA), data that have been processed by the expert determination method or the 'safe harbour' method.

Safe harbour

A Health Insurance Portability and Accountability Act (HIPAA)-designated method of de-identification that relies on the removal of identifiers of the individual, or of relatives, employers or household members of the individual. To achieve this method of de-identification, 18 different types of identifiers including e-mail addresses, social security numbers, all elements of dates directly related to an individual, except year for individuals 89 years of age and younger, and many other elements must be removed.

Creative Commons Public Domain Dedication

(CC0). A licence designed to allow a data generator to waive all rights to the extent allowable by law, enabling any recipient to reuse the content to which it is applied without asking permission or meeting other terms. The current version of the licence is 1.0 and is sometimes referred to as CC0 1.0.

Creative Commons Attribution

(CC BY). A licence designed to enable reuse and sharing as long as the person sharing provides appropriate credit, a link to the licence and a notice of whether or not any changes were made. The current version of the licence is 4.0 and is sometimes referred to as CC BY 4.0.

Health (NIH) for genomic summary results (GSR) to be particularly helpful for data with theoretical risks but limited current danger⁶⁵. This policy favours broad sharing except in the case of "studies for which there are particular sensitivities, such as studies including potentially stigmatizing traits, or with identifiable or isolated study populations".

Metadata can confer risk. Submitters should supply metadata at the highest level of detail that is ethically and legally feasible. Certain identifiers are direct identifiers. Others may not be direct identifiers but may produce small enough equivalence classes to make reidentification possible. Although defining which entities or research projects are covered by the Health Insurance Portability and Accountability Act (HIPAA) of 1996 is beyond the scope of this Review, the law defines useful concepts regarding data sharing and privacy, particularly as it relates to metadata. The HIPAA privacy rule provides two approaches for de-identification of a data set: expert determination and the 'safe harbour' method⁶⁶. The expert determination method requires that a person with appropriate knowledge certifies the risk of reidentifying an individual as 'very small'. The safe harbour method requires removal of 18 HIPAA-specified potentially identifying pieces of information from the data. These types of identifiers pose an avenue of risk and include specific geographic locations tied to an individual, absolute dates and times, and other elements. In these cases, it can be helpful to remove absolute dates and times and replace the date and time fields with intervals. In any case where certain metadata fields introduce risk, we recommend that these fields should be separated — low-risk elements are shared openly whereas high-risk fields are shared only via controlled access in accordance with legal and ethical guidelines.

Machine-learning models can confer risk. Machine-learning models are an emerging form of derived research data that often poses little to no risk. Models trained on publicly available data do not pose a risk above and beyond the data themselves. Models with few parameters relative to the number of subjects also pose less risk. However, models with many parameters that are trained on individual-level genomic data or metadata could reveal detailed information about study participants. Certain attacks have been described that are capable of extracting substantial information about training examples from models or, in certain cases, even the predictions from models⁶⁷. In some cases, models can be trained using techniques such as differential privacy that allow investigators to manage this risk^{68,69}. Such techniques should be considered if sensitive data from human study participants are used during model training. We recommend that high-dimensional models trained on sensitive data without any form of protection should be treated as high risk.

The principles that guide best practices

Data sharing is simply a means to an end. The goal of research with genomic data is often to improve human health or to better understand a biological process. For

such research, stakeholders often include foundations and their donors, taxpayers, study participants who are each dedicating personal or financial resources to these ends and patients who could someday benefit from the research. Participants in clinical trials overwhelmingly want their data to be shared with other academic researchers⁷⁰. Researchers generating genomic data should be driven to responsibly advance the aims of these stakeholders as well as their own. We begin from the premise that the goal of sharing is to enhance the overall pace of research in an ethical manner.

Where feasible, data should be shared through data type-specific repositories that are widely used within a field. Existing data type-specific repositories are ideal data warehouses because they have the following four properties. First, they support publicly available or controlled-access sharing, thus increasing the speed at which data can be requested and obtained. Second, they provide long-term access to the data through provision of a persistent ID, such as a digital object identifier (DOI), and archiving. Third, they lower the costs of research by making large collections of similar data available in a consistent place, which can reduce redundant work and encourage the generation of new hypotheses from secondary analyses. Last, they allow data to be cited, which lets scientists generating data accrue credit for sharing data sets⁷¹. For controlled-access data sets, these repositories provide a consistent request approach. In certain circumstances, particularly early in the development of a data modality, there may be no such repository. In these cases, investigators should choose the last-resort option of placing data in general purpose archiving platforms, such as [Figshare](#) or [Zenodo](#), along with metadata that precisely describe the included files and their format. For data that cannot be publicly shared due to privacy concerns, [Synapse](#) provides a similar general-purpose archiving platform that supports controlled-access sharing.

Principles that should guide sharing of data with reduced risk.

The lowest risk data, including those derived from model organisms or experiments not involving humans, should be maximally shared with minimal restrictions. Investigators should apply a licence to public data sets to provide certainty that they can be reused: Creative Commons Public Domain Dedication (CC0) allows for data to be freely used, and Creative Commons Attribution (CC BY) allows reuse as long as the data sources are attributed. Failure to apply a licence can create substantial barriers to reuse for other researchers^{72,73}. In countries that separate the copyright status of facts from those on creative works, it is possible that much genomic data already fall into the public domain but applying a CC0 licence makes the intent to promote reuse clear. We recommend CC0 for all public data. Certain licences create particular challenges for reuse efforts⁷⁴. Additionally, academic norms require attribution, so CC BY adds barriers but is unlikely to change behaviour. Finally, in the event that someone violates a CC BY licence it seems unlikely that investigators would pursue legal action to enforce a citation requirement. For these reasons, we suggest that CC0 is the most appropriate choice for genomic data that are intended to be public.

Box 1 | Repositories for sharing genomic data

Repository selection process

There are numerous data type-specific repositories for sharing genomic data. Investigators should prioritize repositories that are likely to be maintained: these include those built by the National Center for Biotechnology Information (NCBI) in the United States and the European Bioinformatics Institute (EBI) in Europe. We discuss major NCBI and EBI repositories below. For data modalities for which no data type-specific repositories are maintained by these or similar organizations, investigators should prioritize repositories that are well-adopted within the focused research community. In such cases, and also in cases for which no such repositories exist, investigators should archive their data in general purpose repositories.

The NCBI SRA

The [Sequence Read Archive \(SRA\)](#) supports read-level sequencing data. This includes RNA sequencing (RNA-seq), chromatin immunoprecipitation followed by sequencing (ChIP-seq), assay for transposase-accessible chromatin sequencing (ATAC-seq), whole-exome sequencing, whole-genome sequencing and the results of other such assays. The repository is primarily intended to support US National Institutes of Health (NIH)-funded research; however, data sets less than 1 TB in size may be uploaded without cost. Only data that are intended to be public should be uploaded directly to the SRA submission system. This repository still holds controlled-access data, but the upload process is managed via the [database of Genotypes and Phenotypes \(dbGaP\)](#), which is discussed below.

The EBI ENA

The [European Nucleotide Archive \(ENA\)](#) also supports public read-level sequencing data and the same data types as the SRA. This database shares a data model with the SRA, including the BioProject and BioSample concepts. Uploaded public data are mirrored across both systems. We recommend that investigators in Europe or those without NIH funding use the ENA as a mechanism to publicly disseminate sequencing data. Controlled-access data should be uploaded to the [European Genome-phenome Archive \(EGA\)](#).

The NCBI dbGaP

The dbGaP is designed to support the controlled-access sharing of genomic data for NIH-funded projects. Access is managed through a data access committee (DAC). The dbGaP shares data from genetic association studies, studies of methylation and other individual-level data that are high risk. In the case of raw sequencing data, the data are housed at the NCBI SRA but access and upload are managed through the dbGaP.

The EBI EGA

The EGA holds individual-level genetic data similarly to the dbGaP. Access is also managed by a DAC. The EGA also holds controlled-access sequencing data. We recommend it as the primary choice for investigators working in Europe.

The NCBI GEO

The [Gene Expression Omnibus \(GEO\)](#) holds array-based profiling data intended for public release as well as summary-level data from sequencing experiments. Essentially, if the results of an assay can be expressed as a level of observation for a gene, probe or other such entity, then it can be housed here. For microarray-based experiments, raw data should be uploaded directly to the GEO. For sequencing data, raw data should be uploaded to the SRA and summary-level data should be supplied to the GEO.

The EBI ArrayExpress

[ArrayExpress](#) is the EBI repository that parallels the GEO. For many years, ArrayExpress also imported nearly all GEO data, making it the easiest way to gain access to high-throughput profiling data. Data sets that were imported in the past still exist and have E-GEO- prefixes on their identifiers. However, ArrayExpress no longer imports GEO data so investigators seeking to query these resources will need to query both or a meta-repository that contains the contents of both.

Kipoi

The [Kipoi](#) model repository is a recently developed repository for storing machine-learning models that operate over genomic sequences. The models are regularly tested and paired with an application program interface (API) that facilitates their reuse. We recommend that investigators developing models compatible with Kipoi upload them there. For models not yet compatible with Kipoi, we recommend that investigators use general purpose repositories.

Sample metadata should be provided in as structured a form as possible. Unstructured text elements should be used only when a structured representation is not supported by the database. Well-structured metadata maximize the value for downstream use and also make it easier to verify that metadata do not inadvertently reveal a participant's identity. In data type-specific databases (BOX 1), structured fields are often in place for metadata related to sample handling. Some, such as ArrayExpress, provide entries for commonly used protocols that can be reused⁴⁵. Using existing entries makes it easier to add new experiments and allows subsequent users to select all experiments that follow a specific protocol. For handling metadata,

unstructured fields should be used sparingly and may not need to be used at all for very common analytical strategies.

Principles that should guide sharing of data with elevated risk. The vast majority of clinical trial participants favour data sharing despite potential privacy risks⁷⁰. These privacy risks of data sharing scale according to the chance of one or more parties reidentifying a person in the data set and, also, the potential consequences of reidentification. Successful de-identification is key to reducing the chance of reidentification, and investigators should take care to avoid identifier leakage, which is a particular risk with metadata elements.

The 2015 Institute of Medicine consensus report entitled ‘Sharing Clinical Trial Data: Maximizing Benefits, Minimizing Risk’⁵⁴ is a uniquely comprehensive discussion of the risks of data sharing, and the steps that can be taken to mitigate those risks. Although not specific to genomic data, much of the report applies to genomic data. Among the principles identified in the consensus report is that the context must be considered when thinking about risks of data sharing. If data sharing is via controlled access, then the risks of data sharing are mitigated to some extent. The privacy risks of sharing data sets that focus on rare diseases are generally greater than for common diseases, but not necessarily too great to undertake. To identify risks that could be deemed acceptable requires selecting a way to measure the reidentification risk and selecting an appropriate threshold of risk, and finally measuring the risk in the actual data to be shared. The report encourages investigators to consider the maximum risk to an individual when calculating the risks of publicly shared data sets, and the average risk to individuals for controlled-access data sets.

The risk–benefit ratio of data sharing will look different to different study participants because of varying levels of tolerance for risk and individual reasons for participating in the study. Consent to share de-identified data for secondary analyses can be obtained by design. This approach demonstrates the highest regard for study participants’ interest in the issue of data sharing⁵⁴. However, other, less clear forms of consent language have also been used, with varying degrees of consideration for the privacy of the participants. The approach that is most invasive of participants’ privacy is neither to obtain consent for data sharing up front nor to notify the participants that the de-identified data are being shared. Researchers owe it to their participants to make sure that the impact of the data is maximized within ethical and legal constraints. We recommend that researchers ensure that informed consent language explicitly allows for data to be shared and to “promote research initiatives at other institutions” to maximize the impact of participants’ data⁷⁵.

In many cases it is possible to produce low-risk derivatives or views of high-risk data that retain much of the utility while mitigating much of the risk (FIG. 2). Methods include presenting only summary-level data (FIG. 2a), and potentially adding noise (FIG. 2b). The Exome Aggregation Consortium (ExAC) and Genome Aggregation Database (gnomAD) browsers focus on germline exome and whole-genome sequencing data, and yet are relatively low risk to participants, even though the underlying data are not, by providing summary information and limiting the complexity of queries^{76,77}. Other methods of risk mitigation include redacting data (FIG. 2c) or generating synthetic data that preserve certain statistical properties (FIG. 2d). Given that participants often want their data shared, researchers should aim to identify methods to share valuable derivatives while guarding participant privacy, such as the step of removing human reads performed by the Human Microbiome Project before public sharing.

Investigators who wish to maximize the impact of their research projects should always share findings-level

data publicly unless they pose some risk. In many cases it is also possible to responsibly share intermediate-level data publicly as well. Public sharing reduces the barrier to entry for reanalysis and reduces the chance that a request for data will be received years after the work is done: such requests can be time consuming to answer, and the risk of data being irretrievable increases over time. Finally, data that cannot be responsibly shared in a public manner should be shared through a controlled-access repository.

Privacy is a non-renewable resource

Data have been said to be the new oil, the fuel that will power the economic engine of the twenty-first century^{78–80}. However, the metaphor is imperfect; in stark contrast to oil, data are not lost when shared and are not destroyed when used. On the other hand, privacy is a resource that can be lost, and once it is lost, it cannot be regained. Although fully open data sharing would be ideal from the perspective of the pace of scientific discovery, it is important to consider the privacy costs of sharing study participants’ data.

In general, measurements that are transient are of a lower risk than information that rarely or never changes. For example, sharing metadata that reveals participants’ white blood cell count — which can transiently increase for many reasons — would impose less of a privacy risk than sharing participants’ HIV status. Certain measurements also pose additional concerns: HIV infection has unfortunately been the focus of stigmatization. Information associated with social stigma has a greater risk when sharing data.

The potential for someone to cross-reference information in a de-identified database with other data sources expands the possible threats to privacy⁸¹. Suppose a person tweets that she is proud to have volunteered for a clinical study at a medical school on a particular date, but chooses not to disclose which study. A data analyst may accidentally or intentionally become aware that a particular row of data fits that person’s data due to the date of the tweet. Cross-referencing risk and rare observation risk can be interactive. In a different example, if the date of a visit to the research facility is shared, and the research community is aware that only two families in the United States have a particular disorder, the participant’s home state and decade of life could be sufficient to identify the participant. In general: the rarer a measurement, the more risk it poses for privacy.

For newly designed studies, researchers should plan for sharing at the outset. Ultimately, we need to be guided by the feelings of participants, and most participants do want to see data sharing among academic scientists⁷⁰. Still, careful consideration should be given to what certain technologies, especially sequencing-based technologies, can reveal. In many cases it may not currently be possible to reidentify individuals from a certain data type, but this is a function of other data available for cross-referencing, the computational methods and hardware available, and other factors: future risks for reidentification are difficult or impossible to predict. Consent forms should clearly discuss how data will be shared and the known associated risks, including the caveat that for

genomic data there is significant risk that data that are not currently identifiable can become so in the future. Mechanisms for dynamic consent may be helpful in some regards⁸², but the control that they promise must be carefully considered alongside the potential future risk of reidentification due to new analytical methods.

Making repositories the single point of truth

In software engineering, the concept of a single point of truth can reduce errors⁸³, and similar considerations emerge for research involving genomic data. Data and metadata accumulate during the course of a study, and, ideally, they are stored in one place with one set

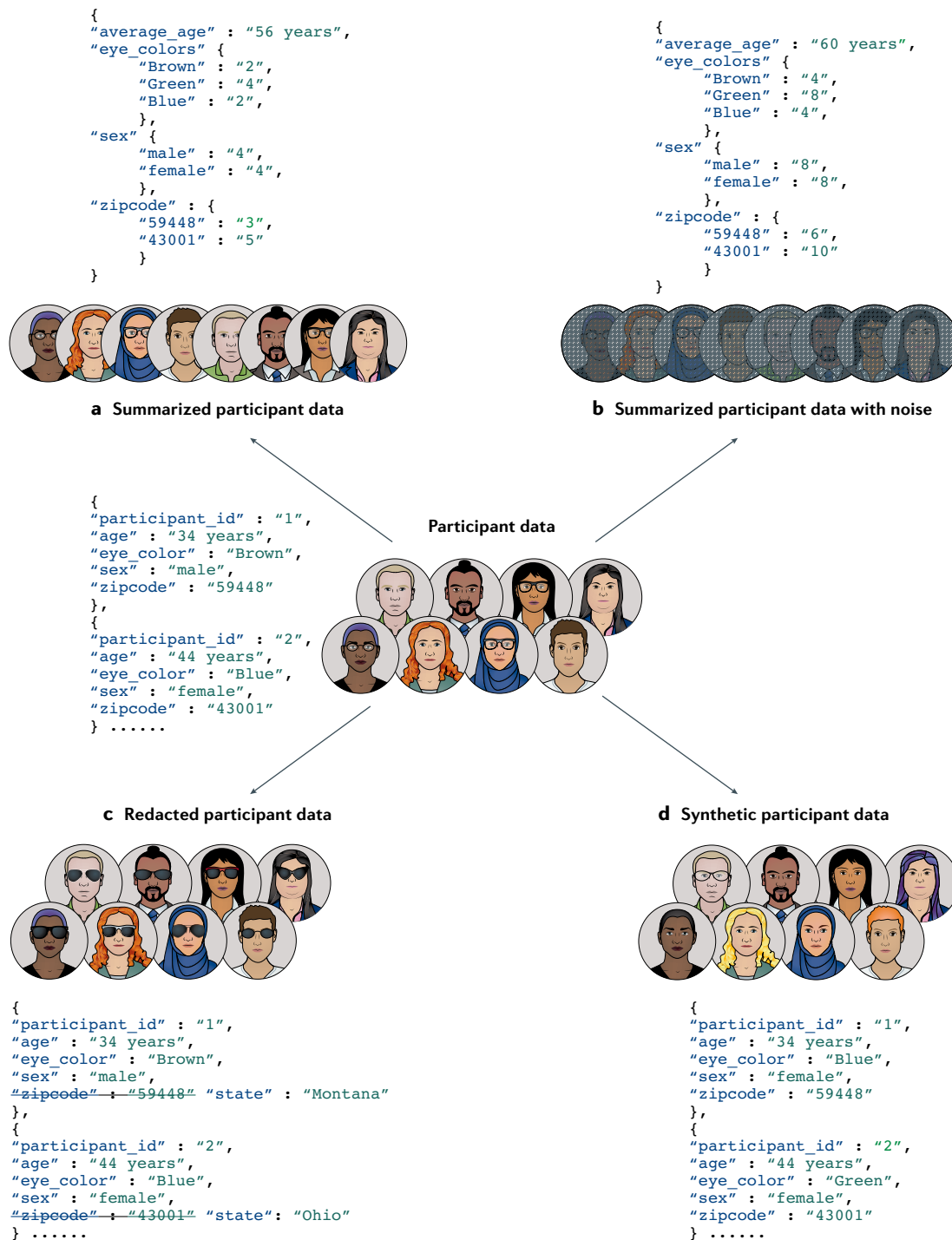


Fig. 2 | **Strategies for de-risking data.** Participant data (centre) can be modified or reported in certain ways to minimize risk. **a** | The data can be reported only at the summary level. **b** | Those summaries can include added noise to make it difficult or impossible to determine the membership of an individual in an aggregate membership. **c** | The data can have identifying fields redacted. **d** | In certain cases, the data can be replaced with entirely synthetic data that have many of the same statistical properties but none of the original individuals.

of metadata descriptors. At this stage it is particularly important for scientists to have procedures in place to track the single point of truth for the data and metadata.

Depositing data into a repository as soon as possible offloads responsibility to the repository and prevents knowledge about the data, including metadata, from atrophying⁸⁴. Depositing data in an accepted repository during a study reduces the risk of turnover leading to lost critical knowledge: with the passage of time, scientists generating data may not remember where the data sets are located and the details describing the data⁷⁵. The repositories also frequently support versioning, allowing researchers to track the state of data over time. Repositories typically do not require that data are made public immediately after they are added: most allow investigators to deposit the data and release them once they are suitably complete and validated for public use.

The concept of a single point of truth also has implications for efforts to construct study-specific data portals or 'data commons'. For such efforts, it is helpful to first deposit data and metadata in data type-specific repositories that are widely used by the biomedical community and then to construct the metadata summaries and derivative files made available on a data commons from these single points of truth.

Repositories for sharing high-risk data. For genomic data, the primary repositories for sharing high-risk data support controlled access. Genetic data, raw RNA-seq reads from human samples and other related data types can often be shared through the same repositories as low-risk data, but with an access control mechanism. As an example, consider the National Center for Biotechnology Information (NCBI) [Sequence Read Archive \(SRA\)](#): access for certain data sets is controlled by the dbGaP. For this database, access is controlled by a data access committee (DAC). Investigators who wish to use such data submit a project description, and the request is submitted by an institutional signing official. This confirms that the host institution is aware of the research and has given ethical approval. The DAC examines the project description and assesses the extent to which the described analysis aligns with the consent that was granted. If an investigator's access is approved, they are then able to access the data.

What if there are no standard repositories? In some cases, there will be no standard repository for the data type. For example, there is not currently a controlled-access repository for machine-learning models trained on clinical data that may leak information about individuals. If there are no standard repositories for the data type, investigators may consider a controlled-access general purpose repository: two primary repositories for public access data are Figshare and Zenodo, and one of the primary such repositories for controlled access data is Synapse, produced by Sage Bionetworks. All general-purpose repositories somewhat hamper reuse. It is harder for users to perform consistent analyses across the contents of the repository, and more onus is on the uploaders to fully document their data formats, metadata and other elements. Because this form of sharing requires more effort

from both sharers and requesters, it should be only used in the case of last resort.

Benefits that accrue to good sharers

Sharing research outputs benefits the scientific community and increases transparency with the public, who predominantly fund the work, through taxpayer dollars as well as charitable giving to non-profit funders⁸⁵. Sharing research outputs promotes reproducible science, with fewer unintentionally duplicate studies allowing research dollars to be put to maximal use. Effective sharing should also accelerate the pace of discovery. Even though sharing benefits the community, it is not necessarily apparent to the scientists generating the data how sharing can benefit them and their careers directly, and this, in particular, is crucially important to address in order to increase their willingness to share high-quality data.

Science progresses by building upon the work of others. Sharing outputs openly leads to better utility and visibility of the research, which leads to more citations of that work^{71,86}. For example, publications with preprints are more cited than those without preprints⁸⁷ and publications with data in openly accessible repositories are more cited when compared with those without accessible data⁸⁸.

To empower the sharing ecosystem, researchers recently created awards to recognize those who share data as well as those who reanalyse publicly available data. The Research Symbiont Awards founded by J.B.B. are given annually to researchers who share data beyond the expectation of their field⁸⁹. The companion award to the Research Symbiont is the Research Parasite Award, founded by C.S.G., which honours those who conduct rigorous secondary analysis of existing data⁹⁰. The goal of these awards is to publicly celebrate those who are committed to sharing and reusing data in a way that contributes to a greater understanding of the world around us.

Open data empower researchers with the ability to pool data, effectively increasing the sample size for appropriately powered studies⁹¹. Furthermore, open data facilitate linking — for example, genomic and epigenomic data with clinical and environmental exposure data — for a greater understanding of disease biology⁹². To further illustrate the power of open data, Milham et al.⁹¹ recently compared the publications resulting from the use of the International Neuroimaging Data-sharing Initiative (INDI) repository by those who contributed data with those who did not. They found that 90.3% of publications resulting from reanalysis of the data in the repository were authored by teams without any data contributions, suggesting that clique/consortium models that only allow access to the data for those who contribute are missing out on bringing new expertise and collaborators into their field who are able to reanalyse the data with fresh perspectives⁹¹.

Funding practices that support sharing

Funders of biomedical research can play a large role in shifting scientific sharing practices. In the absence of sharing requirements, researchers can be reticent to share, but sharing mandates can increase data sharing

FAIR data

Data that are findable, accessible, interoperable and reusable; however, there is no precise definition for each of these criteria, so this is an aspirational goal as opposed to a specific standard.

prevalence^{93–95}. Funders should promote a culture of sharing, and in particular a data sharing culture that builds upon FAIR data standards: ensuring that data are findable, accessible, interoperable and reusable⁹⁶. Barriers to sharing among researchers are multifactorial. Some barriers are practical: researchers may lack the time, funding or understanding of how and where to share. Others are cultural, and may include the lack of adoption in a field, concerns with data misuse or reproducibility and disincentives for sharing related to the potential loss of future publications derived from the data set^{97,98}. We strongly recommend that funders require that data are deposited into standard repositories that provide identifiers to enable output tracking. However, we expect that this alone will not be enough because the quality of data sharing can vary widely⁵³. Hence, there must also be practical ways in which funders can incentivise a greater researcher focus on effective sharing, which we describe next.

A fundamental challenge with incentivising greater sharing is that resources, including data, may not be obviously valuable until a major discovery is made from them (FIG. 3a). However, once a discovery is made, credit for the discovery accrues to the researchers who made that discovery and not necessarily to those who built and publicly shared the resources that enabled it (FIG. 3b). This practice disadvantages sharers: those who share well would do better to hold on to resources and only trade them in the context of a negotiated contract that provides part of the share of the discovery credit (FIG. 3c). Funders have the ability to break this state of poor incentives by considering the applicants' track record of sharing by asking reviewers to consider the evidence of prior sharing. In particular, manuscripts written by unrelated groups using the shared data set can provide *prima facie* evidence of the sharing reputation of the researcher under consideration for funding. If funding decisions are positively influenced by a strong track record, the reputational benefit for sharing can have concrete value that supersedes the value of refusing to share (FIG. 3d). Rewarding open sharing by assessing sharing reputations in funding decisions has the potential to reduce the friction of contract negotiation and accelerate the pace of discovery. Alex's Lemonade Stand Foundation, a leading funder of paediatric cancer research in the United States, is one of the few funders requiring and reviewing prior sharing histories as part of resource sharing plans for all grant applicants, where resource sharing is inclusive of all research outputs, including data^{90,99}.

When funders collectively require and review sharing plans, they provide an amplified voice to this issue that helps to shape sharing practices in the long term. To increase transparency and compliance in data sharing, funders should consider releasing the sharing plans to the public so that the scientific and lay communities know what was promised to be shared, especially when the projects are publicly funded, such as work supported by the NIH⁹⁷. Funders should also require clear statements of when data will be made available.

Although it is important for funders to ask for resource sharing plans, it is also equally important that

funders support the budgeting of reasonable costs for sharing. Sharing effectively requires knowledge, time and money, and funders must be willing to support these costs in order to ensure compliance with sharing policies. For example, Couture et al.⁸⁴ found that compliance with data sharing mandates, despite being higher than that without sharing requirements, is still low: 26% of data were recovered even when required to be shared by a funder mandate. Funders must provide monetary support for high-quality data deposition so that the community does not end up with 'data dumpsters' containing data that are difficult to use due to lack of metadata or meaningful documentation¹⁰⁰.

Funders should also promote the use of university libraries as a resource for the development and implementation of data sharing plans and may consider supporting infrastructure grants that allow for the hiring of personnel devoted to data management or, where needed, support repository formation and/or maintenance^{101,102}. Funders may also consider offering or funding research data management training workshops¹⁰¹. Funders should consider supporting the use of existing tools for the creation of data management plans, including California Digital Library's DMPTool¹⁰³ and Digital Curation Centre's DMPonline¹⁰⁴, which provide templates for data sharing plans⁷⁵.

In summary, funder policies and practices have the potential to dramatically shift the data sharing landscape. Funders should make clear through their actions and funding decisions that they value all research outputs, including data sets, as important scientific contributions^{105,106}. For this to be feasible, unique research outputs should have persistent identifiers that allow them to be cited, highlighting the key importance of sharing via repositories that we emphasize in this Review. Additional open science practices, such as research output sharing, open access publishing and preprinting¹⁰⁵, can help to support this transition. Ultimately, funders should move to establish funding policies based in part on a past track record of effective sharing: this promotes the proactive sharing of high-quality outputs to create an ecosystem where researchers compete to share the highest quality data possible by the most effective method possible.

Publishing practices that support sharing

Journals play a key role in requiring microarray-based gene expression data to be made available at the time of publication¹⁰⁷. Publishers must similarly require that data described in publications are made available. Reviewers should be asked specifically if any data or data sets should be made available. Before an article is published, journal staff should check not only that an accession number is present but also that the accession number resolves to a resource that contains the data described in the published work⁵³. This would avoid certain cases where data that are shared are not as they are described⁵³.

The complement to requiring data availability is ensuring that usage is responsible. Investigators have published research¹⁰⁸ using controlled-access data resources such as the UK Biobank where the research

questions were at best tangentially related to the underlying data access request¹⁰⁹. Journals should require investigators using controlled-access data resources to provide the description of the proposed work as supplementary materials. Reviewers should be asked whether

the study in question aligns with the proposed work. Editors should also use their expert judgement during the editorial review process to assess the extent to which the work described in the manuscript aligns with the underlying request. Journals should refuse to publish

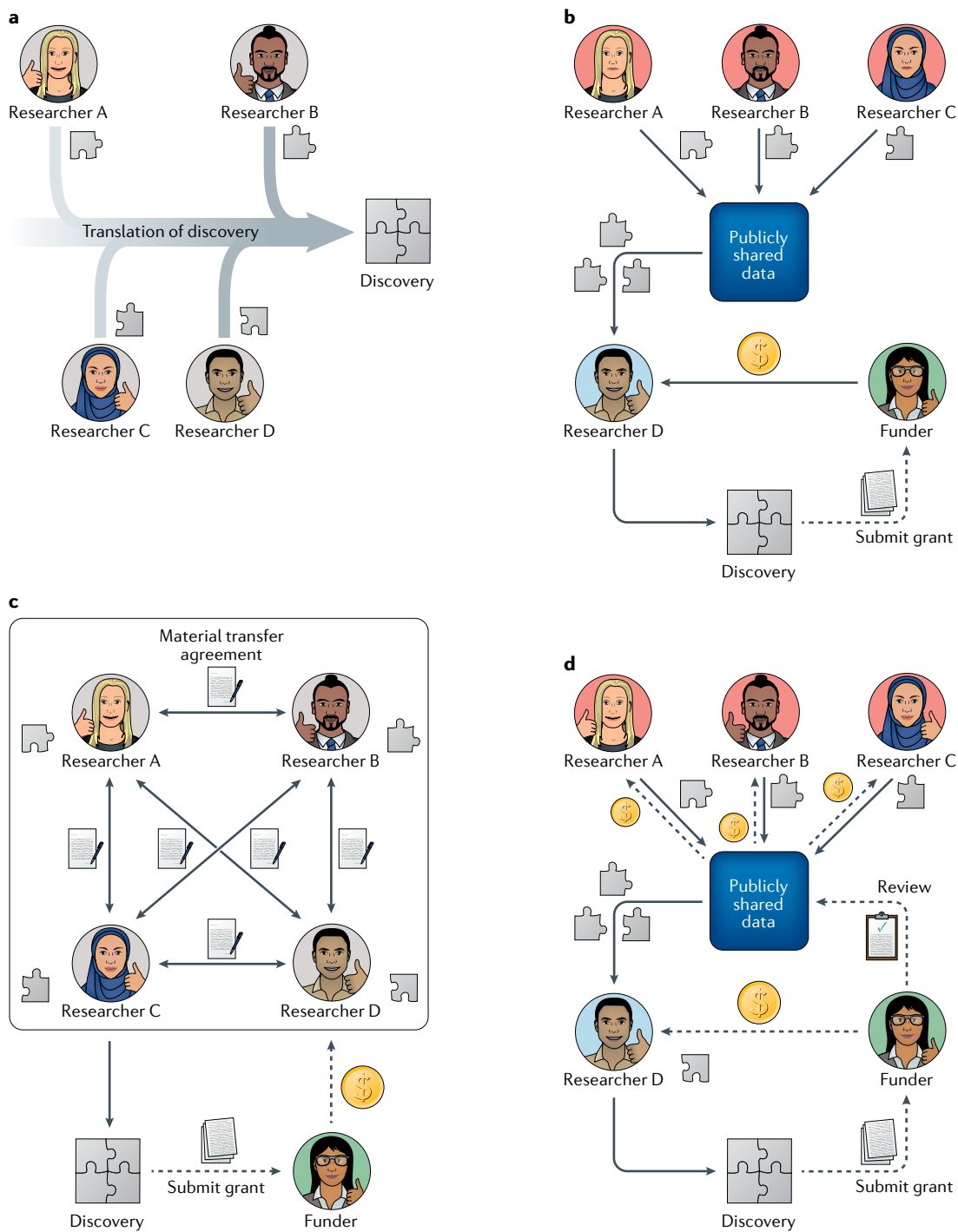


Fig. 3 | How funders consider resources can affect sharing cultures. a | Researcher work products from multiple groups need to be combined to produce a discovery that improves human health. **b** | In a system of open sharing, if funders allocate credit without considering sharing behaviour, much of the credit and funding can accrue to the researcher who brings the final component that enables translation. **c** | Researchers can restrict sharing and negotiate agreements through cliques to enhance the equity of credit distribution, but negotiating agreements is time consuming and may delay or prevent advances. **d** | Funders who consider the value of shared resources when assessing impact provide a benefit not only to the researcher bringing the final component but also all others on the value chain.

work if the data were obtained under pretences that do not match the results.

Perspectives

Investigators must simultaneously balance the wishes of participants to participate in impactful research with the informed risks that participants take in doing so. For genomic data in particular, the risks of participation are not static over time. Our understanding of underlying biological mechanisms, the presence of other complementary data types and the power of our analytical approaches all affect the risk of reidentification. Research is needed on processes that can generate derivatives that maximize reuse value while mitigating the reidentification risk for as long as is possible. Still, because perfect risk reduction is likely to be impossible, researchers should not consent participants under promises that genomic data will be made de-identifiable. Certain efforts are underway to create computing environments that expose data for analysis but that limit risk, but guidance from the trajectory of beacons^{110,111} to reidentification⁵⁶ suggests that technical solutions may be insufficient. In an era when we can expect those interested in reusing data to aim to train high-parameter machine-learning models, investigators should take guidance in designing consent processes from the limited number of efforts that intended to publicly release variant-level data. For the 1000 Genomes Project^{112,113} and the Harvard Personal Genome Project¹¹⁴, participants consented to have their germline genetic data openly shared. In a pilot programme in Texas, many patients with cancer elected to have both germline and somatic variants shared openly¹¹⁵. It is clear that at least some are willing to participate in research, even if this leads to the public release of their germline genetic variants. Even for projects where the primary sharing mechanism is intended to be controlled access,

investigators may wish to offer participants the opportunity to become 'data donors' whose data would be publicly shared.

Researchers recruiting participants must also make every effort to ensure that data sharing and consent processes do not marginalize certain participants or groups of individuals. The overwhelming presence of individuals of European descent in genetic databases has been widely documented^{116,117}. A fuller communication of the potential risks of participating could discourage individuals from certain groups, particularly those who have been minoritized, from participating. Researchers have a responsibility to make sure that benefits of research accrue broadly to society: an increased proportion of individuals who decline to participate in genomic research should not be an acceptable excuse for disparities in the extent to which research benefits the members of that group.

Researchers who generate genomic data can take certain steps to make those data as impactful as possible: adding key metadata elements, sharing the data with the fewest restrictions possible and putting data in data type-specific repositories. However, creating a responsible culture of data sharing that accelerates research is more than just the responsibility of those who generate data in the course of their research. For controlled-access human study participants' data, those analysing the data have a responsibility to do so in accordance with the consent of participants and supplied study plans. Journals have a responsibility to decline to publish analyses that are not conducted in accordance with ethical research practices. Funders have a responsibility to support ethical research in diverse populations while preferentially supporting those who have established exemplary records of generating widely reused resources.

Published online 21 July 2020

- Brazma, A. et al. Minimum information about a microarray experiment (MIAME)—toward standards for microarray data. *Nat. Genet.* **29**, 365–371 (2001).
This paper describes an effort to catalogue what elements must be shared for a genome-wide assay of gene expression to be suitable for reuse and reanalysis.
- Myers, C. L. et al. Discovery of biological networks from diverse functional genomic data. *Genome Biol.* **6**, R114 (2005).
- Mostafavi, S., Ray, D., Warde-Farley, D., Grouios, C. & Morris, Q. GeneMANIA: a real-time multiple association network integration algorithm for predicting gene function. *Genome Biol.* **9**, S4 (2008).
- Huttenhower, C. et al. Exploring the human genome with functional maps. *Genome Res.* **19**, 1093–1106 (2009).
- Lee, I. et al. Predicting genetic modifier loci using functional gene networks. *Genome Res.* **20**, 1143–1153 (2010).
- Tan, J. et al. Unsupervised extraction of stable expression signatures from public compendia with an ensemble of neural networks. *Cell Syst.* **5**, 63–71.e6 (2017).
- Taroni, J. N. et al. MultiPLIER: a transfer learning framework for transcriptomics reveals systemic features of rare disease. *Cell Syst.* **8**, 380–394.e4 (2019).
- Ragoussis, J. Genotyping technologies for genetic research. *Annu. Rev. Genomics Hum. Genet.* **10**, 117–133 (2009).
- Ng, P. C. & Kirkness, E. F. in *Genetic Variation: Methods and Protocols* (eds Barnes, R. M. & Breen, G.) 215–226 (Humana, 2010).
- Beck, S. & Rakyen, V. K. The methylome: approaches for global DNA methylation profiling. *Trends Genet.* **24**, 231–237 (2008).
- Harmanci, A. & Gerstein, M. Analysis of sensitive information leakage in functional genomics signal profiles through genomic deletions. *Nat. Commun.* **9**, 1–10 (2018).
- Gürsoy, G., Brannon, C. M., Navarro, F. C. P. & Gerstein, M. FANCY: fast estimation of privacy risk in functional genomics data. Preprint at *bioRxiv* <https://doi.org/10.1101/775338> (2020).
- Piskol, R., Ramaswami, G. & Li, J. B. Reliable identification of genomic variants from RNA-seq data. *Am. J. Hum. Genet.* **93**, 641–651 (2013).
- Brouard, J. S., Schenkel, F., Marete, A. & Bissonnette, N. The GATK joint genotyping workflow is appropriate for calling variants in RNA-seq experiments. *J. Anim. Sci. Biotechnol.* **10**, 44 (2019).
- Deelen, P. et al. Calling genotypes from public RNA-sequencing data enables identification of genetic variants that affect gene-expression levels. *Genome Med.* **7**, 30 (2015).
- Mailman, M. D. et al. The NCBI dbGaP database of Genotypes and Phenotypes. *Nat. Genet.* **39**, 1181–1186 (2007).
- Homer, N. et al. Resolving individuals contributing trace amounts of DNA to highly complex mixtures using high-density SNP genotyping microarrays. *PLoS Genet.* **4**, e1000167 (2008).
- Longo, D. L. & Drazen, J. M. Data sharing. *N. Engl. J. Med.* **374**, 276–277 (2016).
- Perou, C. M. Show me the data! *Nat. Genet.* **29**, 373–373 (2001).
- Colavizza, G., Hrynaskiewicz, I., Staden, I., Whitaker, K. & McGillivray, B. The citation advantage of linking publications to research data. *PLoS One* **15**, e0230416 (2020).
- Perou, C. M. et al. Molecular portraits of human breast tumours. *Nature* **406**, 747–752 (2000).
- Clarke, L. et al. The 1000 Genomes Project: data management and community access. *Nat. Methods* **9**, 1–4 (2012).
- Aebbersold, R. & Mann, M. Mass spectrometry-based proteomics. *Nature* **422**, 198–207 (2003).
- Trivedi, D. K., Hollywood, K. A. & Goodacre, R. Metabolomics for the masses: the future of metabolomics in a personalized world. *N. Horiz. Transl. Med.* **3**, 294–305 (2017).
- Marioni, J. C., Mason, C. E., Mane, S. M., Stephens, M. & Gilad, Y. RNA-seq: an assessment of technical reproducibility and comparison with gene expression arrays. *Genome Res.* **18**, 1509–1517 (2008).
- Handelsman, J. Metagenomics: application of genomics to uncultured microorganisms. *Microbiol. Mol. Biol. Rev.* **68**, 669–685 (2004).
- Subramanian, A. et al. A next generation connectivity map: L1000 platform and the first 1,000,000 profiles. *Cell* **171**, 1437–1452.e17 (2017).
- Konecny, G. E. et al. Prognostic and therapeutic relevance of molecular subtypes in high-grade serous ovarian cancer. *J. Natl. Cancer Inst.* **106**, dju249 (2014).
- Zinger, L. et al. Global patterns of bacterial β -diversity in seafloor and seawater ecosystems. *PLoS One* **6**, e24570 (2011).
- Cock, P. J. A., Fields, C. J., Goto, N., Heuer, M. L. & Rice, P. M. The sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic Acids Res.* **38**, 1767–1771 (2009).

31. Tan, J., Hammond, J. H., Hogan, D. A. & Greene, C. S. ADAGE-based integration of publicly available *Pseudomonas aeruginosa* gene expression data with denoising autoencoders illuminates microbe–host interactions. *mSystems* **1**, e00025–15 (2016).
32. Zhou, J. & Troyanskaya, O. G. Predicting effects of noncoding variants with deep learning-based sequence model. *Nat. Methods* **12**, 931–934 (2015).
33. Zhou, W. & Altman, R. B. Data-driven human transcriptomic modules determined by independent component analysis. *BMC Bioinformatics* **19**, 327 (2018).
34. Stein-O'Brien, G. L. et al. Decomposing cell identity for transfer learning across cellular measurements, platforms, tissues, and species. *Cell Syst.* **8**, 395–411. e8 (2019).
35. Way, G. P., Zietz, M., Rubinetti, V., Himmelstein, D. S. & Greene, C. S. Compressing gene expression data using multiple latent space dimensionalities learns complementary biological representations. *Genome Biol.* **21**, 109 (2020).
36. Gulshan, V. et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *JAMA* **316**, 2402–2410 (2016).
37. Kelley, D. R., Snoek, J. & Rinn, J. L. Basset: learning the regulatory code of the accessible genome with deep convolutional neural networks. *Genome Res.* **26**, 990–999 (2016).
38. Beaulieu-Jones, B., Greene, C. Reproducibility of computational workflows is automated using continuous analysis. *Nat. Biotechnol.* **35**, 342–346 (2017).
39. Avsec, Z. et al. The Kipoi repository accelerates community exchange and reuse of predictive models for genomics. *Nat. Biotechnol.* **37**, 592–600 (2019). **This paper describes one of the first repositories for machine-learning models and uses continuous integration to verify that the models are reusable and interoperable.**
40. Barrett, T. et al. BioProject and BioSample databases at NCBI: facilitating capture and organization of metadata. *Nucleic Acids Res.* **40**, D57–D63 (2012). **This paper shows that the BioProject and BioSample databases are one of the primary ways in which genomics metadata are stored and accessed.**
41. Kudesia, S., Bhardwaj, A., Thakur, B., Kishore, S. & Bahal, N. Primary MPNST in childhood — a rare case report. *J. Clin. Diagn. Res.* **8**, FD01–FD02 (2014).
42. Malone, J. et al. Modeling sample variables with an experimental factor ontology. *Bioinformatics* **26**, 1112–1118 (2010).
43. Wang, Z., Jensen, M. A. & Zenklusen, J. C. A practical guide to The Cancer Genome Atlas (TCGA). *Methods Mol. Biol.* **1418**, 111–141 (2016).
44. Park, Y. & Greene, C. S. A parasite's perspective on data sharing. *Gigascience* **7**, gyl129 (2018).
45. Rustici, G. et al. ArrayExpress update—trends in database growth and links to data analysis tools. *Nucleic Acids Res.* **41**, D987–D990 (2013).
46. Barrett, T. et al. NCBI GEO: archive for functional genomics data sets—update. *Nucleic Acids Res.* **41**, D991–D995 (2013).
47. Sudlow, C. et al. UK Biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015). **This paper shows that the UK Biobank is a remarkable example of sharing high-risk data in a manner that has accelerated health research.**
48. National Institutes of Health. All of us. *NIH* <https://allofus.nih.gov/> (2020).
49. Savage, C. J. & Vickers, A. J. Empirical study of data sharing by authors publishing in PLoS journals. *PLoS One* **4**, e7078 (2009).
50. Wood, B. D. K., Müller, R. & Brown, A. N. Push button replication: is impact evaluation evidence for international development verifiable? *PLoS One* **13**, e0209416 (2018).
51. Tenopir, C. et al. Changes in data sharing and data reuse practices and perceptions among scientists worldwide. *PLoS One* **10**, e0134826 (2015).
52. Campbell, H. A., Micheli-Campbell, M. A. & Udyawer, V. Early career researchers embrace data sharing. *Trends Ecol. Evolution* **34**, 95–98 (2019).
53. Learned, K. et al. Barriers to accessing public cancer genomic data. *Sci. data* **6**, 98 (2019). **This contribution notes how not all public data sharing is equal, and the implementation greatly affects how reusable and interoperable data are.**
54. Institute of Medicine. *Sharing Clinical Trial Data: Maximizing Benefits, Minimizing Risk* (National Academies, 2015). **This book discusses the risks and benefits associated with sharing and how we can balance them.**
55. Malin, B. A. An evaluation of the current state of genomic data privacy protection technology and a roadmap for the future. *J. Am. Med. Inform. Assoc.* **12**, 28–34 (2004).
56. Shringarpure, S. S. & Bustamante, C. D. Privacy risks from genomic data-sharing beacons. *Am. J. Hum. Genet.* **97**, 631–646 (2015).
57. Erlich, Y., Shor, T., Pe'er, I. & Carmi, S. Identity inference of genomic data using long-range familial searches. *Science* **362**, 690–694 (2018).
58. Gürsoy, G., Harmandi, A., Green, M. E., Navarro, F. C. P. & Gerstein, M. Sensitive information leakage from functional genomics data: theoretical quantifications & practical file formats for privacy preservation. Preprint at *bioRxiv* <https://doi.org/10.1101/345074> (2018).
59. Kaiser, J. We will find you: DNA search used to nab Golden State Killer can home in on about 60% of white Americans. *Science* <https://doi.org/10.1126/science.aav7021> (2018).
60. Hill, K. & Murphy, H. Your DNA profile is private? A Florida judge just said otherwise. *The New York Times* <https://www.nytimes.com/2019/11/05/business/dna-database-search-warrant.html> (5 Nov 2019).
61. Lloyd-Price, J. et al. Strains, functions and dynamics in the expanded human microbiome project. *Nature* **550**, 61–66 (2017).
62. Philibert, R. A. et al. Methylation array data can simultaneously identify individuals and convey protected health information: an unrecognized ethical concern. *Clin. Epigenetics* **6**, 28 (2014).
63. Edge, M. D. & Coop, G. Attacks on genetic privacy via uploads to genealogical databases. *eLife* **9**, e51810 (2020).
64. Schadt, E. E., Woo, S. & Hao, K. Bayesian method to predict individual SNP genotypes from gene expression data. *Nat. Genet.* **44**, 603–608 (2012).
65. National Institutes of Health. Update to NIH management of genomic summary results access. *NIH* <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-19-023.html> (2018).
66. US Department of Health and Human Services. Methods for de-identification of protected health information in accordance with the Health Insurance Portability and Accountability Act (HIPAA) privacy rule (HHS, 2020).
67. Shokri, R., Stronati, M., Song, C. & Shmatikov, V. Membership inference attacks against machine learning models. *Proc. IEEE Symp. Security Privacy* <https://doi.org/10.1109/SP.2017.41> (2017).
68. Abadi, M. et al. Deep learning with differential privacy. *Proc. ACM Conf. Comput. Commun. Security* <https://doi.org/10.1145/2976749.2978318> (2016).
69. Beaulieu-Jones, B. K. et al. Privacy-preserving generative deep neural networks support clinical data sharing. *Circ. Cardiovasc. Qual. Outcomes* **12**, 159756 (2019).
70. Mello, M. M., Lieou, V. & Goodman, S. N. Clinical trial participants' views of the risks and benefits of data sharing. *N. Engl. J. Med.* **378**, 2202–2211 (2018).
71. Furman, J. L. & Stern, S. Climbing atop the shoulders of giants: the impact of institutions on cumulative research. *Am. Econ. Rev.* **101**, 1933–1963 (2011).
72. Oxenham, S. Legal maze threatens to slow data science. *Nature* **536**, 16–17 (2016). **This paper discusses how licensing of data is important, and choosing no licence or a restrictive licence can slow reuse efforts dramatically.**
73. Himmelstein, D. S. et al. Systematic integration of biomedical knowledge prioritizes drugs for repurposing. *eLife* **6**, e26726 (2017).
74. Hagedorn, G. et al. Creative Commons licenses and the non-commercial condition: implications for the re-use of biodiversity information. *ZooKeys* **150**, 127–149 (2011).
75. Mannheimer, S., Pienta, A., Kirilova, D., Elman, C. & Wutich, A. Qualitative data sharing: data repositories and academic libraries as key partners in addressing challenges. *Am. Behav. Sci.* **63**, 643–664 (2019).
76. Lek, M. et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
77. Karczewski, K. J. et al. Variation across 141,456 human exomes and genomes reveals the spectrum of loss-of-function intolerance across human protein-coding genes. Preprint at *bioRxiv* <https://doi.org/10.1101/531210> (2019).
78. ANA Marketing Maestros. Data is the new oil. *ANA Marketing Maestros* https://ana.blog.com/maestros/2006/11/data_is_the_new_oil (2006).
79. European Commission. Meglena Kuneva — European Consumer Commissioner — keynote speech — roundtable on online data collection, targeting and profiling (EC, 2009).
80. Microsoft. Qi Lu: Build 2016. *Microsoft* <https://news.microsoft.com/speeches/qi-lu-build-2016/> (2016).
81. Narayanan, A. & Shmatikov, V. Robust de-anonymization of large sparse datasets. *Proc. IEEE Symp. Security Privacy* <https://doi.org/10.1109/SP.2008.33> (2008).
82. Kaye, J. et al. Dynamic consent: a patient interface for twenty-first century research networks. *Eur. J. Hum. Genet.* **23**, 141–146 (2015).
83. Holzmann, G. J. Points of truth. *IEEE Softw.* **32**, 18–21 (2015). **This paper discusses the principle of a single point of truth in programming, which is a concept that should also be widely considered in data storage and sharing.**
84. Couture, J. L., Blake, R. E., McDonald, G. & Ward, C. J. A funder-imposed data publication requirement seldom inspired data sharing. *PLoS One* **13**, e0199789 (2018). **This paper reports that only around a quarter of source data sets were shared, even when funder mandates required it.**
85. Mervis, J. Data check: U.S. government share of basic research funding falls below 50%. *Science* <https://doi.org/10.1126/science.aal0890> (2017).
86. Piwowar, H. A., Day, R. S. & Fridsma, D. B. Sharing detailed research data is associated with increased citation rate. *PLoS One* **2**, e308 (2007). **This paper demonstrates that publications linked with a public data set accrue more citations than those without accessible data.**
87. Fraser, N., Momeni, F., Mayr, P. & Peters, I. The effect of bioRxiv preprints on citations and altmetrics. Preprint at *bioRxiv* <https://doi.org/10.1101/673665> (2019).
88. Piwowar, H. A. & Vision, T. J. Data reuse and the open data citation advantage. *PeerJ* **1**, e175 (2013). **This report details factors that support reuse and examines reuse over long time intervals. Many data sets still accrue reuse citations 5 years after the initial publication.**
89. Byrd, J. B. & Greene, C. S. Data-sharing models. *N. Engl. J. Med.* **376**, 2305–2306 (2017).
90. Greene, C. S., Garmire, L. X., Gilbert, J. A., Ritchie, M. D. & Hunter, L. E. Celebrating parasites. *Nat. Genet.* **49**, 483–484 (2017).
91. Milham, M. P. et al. Assessment of the impact of shared brain imaging data on the scientific literature. *Nat. Commun.* **9**, 2818 (2018).
92. Joly, Y., Dyke, S. O. M., Knoppers, B. M. & Pastinen, T. Are data sharing and privacy protection mutually exclusive? *Cell* **167**, 1150–1154 (2016).
93. Levenstein, M. C. & Lyle, J. A. Data: sharing is caring. *Adv. Methods Pract. Psychol. Sci.* **1**, 95–103 (2018).
94. Federer, L. M. et al. Data sharing in *PLOS ONE*: an analysis of data availability statements. *PLoS One* **13**, e0194768 (2018).
95. Nuijten, M. B. et al. Journal data sharing policies and statistical reporting inconsistencies in psychology. *Collabra Psychol.* **3**, 31 (2017).
96. Wilkinson, M. D. et al. Comment: the FAIR guiding principles for scientific data management and stewardship. *Sci. Data* **3**, 160018 (2016).
97. Terry, R. F., Littler, K. & Olliaro, P. L. Sharing health research data — the role of funders in improving the impact. *F1000Research* **7**, 1641 (2018).
98. Stuart, D. et al. Whitepaper: practical challenges for researchers in data sharing. *Nat. Res.* <https://doi.org/10.6084/M9.FIGSHARE.5975011.V1> (2018).
99. Teytelman, L. No more excuses for non-reproducible methods. *Nature* **560**, 411 (2018).
100. Merson, L., Gaye, O. & Guerin, P. J. Avoiding data dumpsters—toward equitable and useful data sharing. *N. Engl. J. Med.* **374**, 2414–2415 (2016).
101. Berghmans, et al. Open data: the researcher perspective — survey and case studies. *Mendeley Data* <https://doi.org/10.17632/bwrnfb4bvh.1> (2017).
102. Popkin, G. Data sharing and how it can benefit your scientific career. *Nature* **569**, 445–447 (2019).
103. DMPTool. *California Digital Library* <https://dmptool.org/> (2020).

104. DMPonline. *Digital Curation Center* <https://dmponline.dcc.ac.uk/> (2020).
105. Kiley, R., Peatfield, T., Hansen, J. & Reddington, F. Data sharing from clinical trials—a research funder’s perspective. *N. Engl. J. Med.* **377**, 1990–1992 (2017).
106. Piwowar, H. Altmetrics: value all research products. *Nature* **493**, 159 (2013).
107. Ball, C. A. et al. Submission of microarray data to public repositories. *PLoS Biol.* **2**, e317 (2004).
108. Hill, W. D. et al. Genome-wide analysis identifies molecular systems and 149 genetic loci associated with income. *Nat. Commun.* **10**, 5741 (2019).
109. UK Biobank. The relationship of cognitive function and negative emotions with morbidity and mortality: an aetiological investigation (Biobank, 2015).
110. Fiume, M. et al. Federated discovery and sharing of genomic data using beacons. *Nat. Biotechnol.* **37**, 220–224 (2019).
111. Global Alliance for Genomics and Health. A federated ecosystem for sharing genomic, clinical data. *Science* **352**, 1278–1280 (2016).
112. Siva, N. 1000 Genomes Project. *Nat. Biotechnol.* **26**, 256 (2008).
113. Auton, A. et al. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
114. Ball, M. P. et al. Harvard Personal Genome Project: lessons from participatory public research. *Genome Med.* **6**, 10 (2014).
115. Becnel, L. B. et al. An open access pilot freely sharing cancer genomic data from participants in Texas. *Sci. Data* **3**, 160010 (2016).

116. Hindorff, L. A. et al. Prioritizing diversity in human genomics research. *Nat. Rev. Genet.* **19**, 175–185 (2018).
117. Sirugo, G., Williams, S. M. & Tishkoff, S. A. The missing diversity in human genetic studies. *Cell* **177**, 26–31 (2019).

Acknowledgements

This work was funded in part by grants from Alex’s Lemonade Stand Foundation (CCDL), the US National Institutes of Health (NIH) (K25 HL128909, R01 CA237170 and R01 HG010067) and the Gordon and Betty Moore Foundation (GBMF 4552).

Author contributions

J.B.B., A.C.G., D.V.P. and C.S.G. researched data for the article. J.B.B., A.C.G., X.J. and C.S.G. substantially contributed to discussion of content. J.B.B., A.C.G. and C.S.G. wrote the article. All authors reviewed/edited the manuscript before submission.

Competing interests

A.C.G. is an employee of a funder, Alex’s Lemonade Stand Foundation. As an author, A.C.G. participated in all aspects of conceptualization, design, preparation of the manuscript and the decision to publish. D.V.P. is an employee of a funder, Alex’s Lemonade Stand Foundation. As an author, D.V.P. participated in preparation of the manuscript and the decision to publish. C.S.G. is the Director of the Alex’s Lemonade Stand Foundation’s Childhood Cancer Data Lab. As an author, C.S.G. participated in all aspects of conceptualization, design, preparation of the manuscript and the decision to publish.

Peer review information

Nature Reviews Genetics thanks O. Hofmann and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.

Publisher’s note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

RELATED LINKS

ArrayExpress: <https://www.ebi.ac.uk/arrayexpress/>
Database of Genotypes and Phenotypes (dbGaP): <https://www.ncbi.nlm.nih.gov/gap/>
European Genome–phenome Archive (EGA): <https://www.ebi.ac.uk/ega/home>
European Nucleotide Archive (ENA): <https://www.ebi.ac.uk/ena>
Figshare: <https://figshare.com/>
Gene Expression Omnibus (GEO): <https://www.ncbi.nlm.nih.gov/geo/>
Kipoi: <https://kipoi.org/>
NOT-OD-19-023: <https://grants.nih.gov/grants/guide/notice-files/NOT-OD-19-023.html>
Sequence Read Archive (SRA): <https://www.ncbi.nlm.nih.gov/sra>
Synapse: <https://synapse.org>
Zenodo: <https://zenodo.org/>

© Springer Nature Limited 2020