# **POLICY**FORUM

## MEGASCIENCE

# **Omics Data Sharing**

Dawn Field,<sup>1\*</sup>†‡ Susanna-Assunta Sansone,<sup>1,2</sup>† Amanda Collis,<sup>3</sup>† Tim Booth,<sup>1</sup> Peter Dukes,<sup>4</sup> Susan K. Gregurick,<sup>5</sup> Karen Kennedy,<sup>6</sup> Patrik Kolar,<sup>7</sup> Eugene Kolker,<sup>8</sup> Mary Maxon,<sup>9</sup> Siân Millard,<sup>10</sup> Alexis-Michel Mugabushaka,<sup>11</sup> Nicola Perrin,<sup>12</sup> Jacques E. Remacle,<sup>7</sup> Karin Remington,<sup>13</sup> Philippe Rocca-Serra,<sup>12</sup> Chris F. Taylor,<sup>12</sup> Mark Thorley,<sup>14</sup> Bela Tiwari,<sup>1</sup> John Wilbanks<sup>15</sup>

Development of high-throughput genomic and postgenomic technologies has caused a change in approaches to data handling and processing (1). One biological sample might be used to generate many kinds of "big" data in parallel, such as genome sequence (genomics), patterns of gene and protein expression (transcriptomics and proteomics), and metabolite concentrations and fluxes (metabolomics). Extensive computer manipulations are required for even basic analyses of such data; the challenges mount further when two or more studies' outputs must be compared or integrated.

Grassroots movements (2–5), efforts including the Science Commons, which is initiating an open-access data protocol (6), as well as topdown (funder-led) efforts (see table, page 235), have led to a range of policies for data management and sharing. A recent European Science Foundation consultation exercise confirmed a lack of explicit, well-documented data-sharing policies for most funding agencies in European countries (7). If we are to avoid squandering the immediate and extended value of big data, a focused strategy will be pivotal.

Early policies were driven by the need to manage long-term data sets (those accrued over 30 or more years), such as those in the social and environmental sciences. More recently, policies have emerged in response to increased funding for high-throughput approaches in major 'omics fields. The European Commission has invited the member states to develop policies to implement access, dissemination, and preservation for scientific knowledge and data (8).

Beyond public and private funding agen-

\*Full author affiliations are available on *Science* Online. †These authors contributed equally to this article. ‡Author for correspondence. E-mail: dfield@ceh.ac.uk cies, regulatory agencies such as the U.S. Food and Drug Administration (FDA) (9), European Medicines Agency (EMEA) (10), and U.S. Environmental Protection Agency (EPA) (11) are also working to define guidelines to facilitate electronic submission of traditional and 'omics data types. These, as well as industry guidelines, are beyond the scope of this document, but much could be learned from an exchange of ideas and practices (12).

The policies listed here share common principles. They aim to protect cumulative data outputs. All recognize data as a public good and data sharing as a way to accelerate subsequent exploitation. On a practical level, all acknowledge the right of first use for data providers and the right to appropriate accreditation. Likewise, these policies have been generated through the same basic process (table S1) (13).

Despite these commonalities, there is still room for heterogeneity, as expected, given the different types of communities served by each funder and the data types they generate. Care must be taken, though, that these differences do not impede seamless interoperability. The path a funding agency takes in supporting its data policy largely reflects the relative emphasis placed on managing versus sharing data. A focus on managing is often accompanied by an institutional infrastructure. Such centralization provides economy of scale, institutional memory, and reusable capability, but it also incurs a substantial direct cost that may compete with research funding (14). The UK Natural Environment Research Council (NERC) sustains a system of national data centers and has invested in the NERC Environmental Bioinformatics Centre (NEBC) to cover 'omics data (15, 16). Similarly, the UK Economic and Social Research Council provides a central data service for social scientists (17). Policies that focus on sharing tend to place more responsibility on researchers. For example, the UK Biotechnology and Biological Sciences Research Council (BBSRC) is supporting its data-sharing policy through funds that allow researchers to develop their own solutions from the bottom up.

Massive-scale raw data must be highly structured to be useful to downstream users. Standardized solutions are increasingly available for describing, formatting, submitting,

# Data sharing, and the good annotation practices it depends on, must become part of the fabric of daily research for researchers and funders.

and exchanging data (18, 19). These reporting standards include minimum information checklists, ontologies, and file formats. Minimum information checklists are simple, structured documents that reflect the consensus view of a community on the information to report about particular kinds of biological studies or instrument-based assays. Ontologies provide terms needed to describe the minimal information requirements. File formats define a shared syntax to transmit and exchange standardized information.

There are now an escalating number of community-developed checklists, ontologies, and file-format projects, a positive sign of community engagement. But this proliferation brings with it new sociological and technological challenges-creating interoperability and avoiding unnecessary overlaps and duplication of efforts. These projects largely focus on a particular technology or a specific biological knowledge domain (e.g., ontologies for anatomy, gene functions, or the environment) and are by nature fragmented and not designed to be interoperable. A range of activities are fostering harmonization and consolidation of these standards for checklists (5), ontologies (4), and representation of information in electronic formats (2, 3).

Many large coordinative initiatives (20–23) are working to address the problem of archiving and integrating data. The ELIXIR project (22) aims to construct and operate a common, sustainable bioinformatics research infrastructure to support the life sciences across Europe. The Infrastructure for Spatial Information in the European Community (INSPIRE) directive requires that Europe binds together its geospatial data into portals (23). Widely useful are initiatives like the Digital Curation Centre (DCC), which tracks data standards, documents best practice, and has published a data life-cycle model to underpin long-term datapreservation policies (24).

# **Achieving Adherence**

Community adherence would be automatic if guidelines aligned with prevailing scientific culture and (emergent) practice. However, there is often a gulf or even outright resistance (25, 26).

Policies that stipulate public data release, especially of prepublication data, raise researchers' concerns about loss of intellectual ownership—for example, by compromising

9 OCTOBER 2009 VOL 326 SCIENCE www.sciencemag.org Published by AAAS

<sup>&</sup>lt;sup>1</sup>U.K. Natural Environment Research Council (NERC), Environmental Bioinformatics Centre. <sup>2</sup>European Molecular Biology Laboratory (EMBL) Outstation, The European Bioinformatics Institute (EBI). <sup>3</sup>U.K. Biotechnology and Biological Sciences Research Council. <sup>4</sup>U.K. Medical Research Council. <sup>5</sup>U.S. Department of Energy. <sup>6</sup>Genome Canada and Wellcome Trust Sanger Institute. <sup>7</sup>Unit for Genomics and Systems Biology, European Commission. <sup>8</sup>Seattle Childrens Hospital. <sup>9</sup>Marine Microbiology Initiative, Gordon and Betty Moore Foundation. <sup>10</sup>U.K. Economic and Social Research Council. <sup>11</sup>European Science Foundation. <sup>12</sup>The Wellcome Trust. <sup>12</sup>U.S. National Institute of General Medical Science, NIH. <sup>14</sup>NERC. <sup>15</sup>Science Commons.

Funding body	COUNTRY	YEAR	POLICY INFORMATION
Economic and Social Research Council (ESRC)	UK	(1994) 2000	www.esrcsocietytoday.ac.uk/ESRCInfoCentre/Images/DataPolicy2000_tcm6-12051.pdf
Natural Environment Research Council (NERC)	UK	(1996) 2008	www.nerc.ac.uk/research/sites/data/policy.asp
National Science Foundation (NSF)	US	2001	www.nsf.gov/pubs/2001/gc101/gc101rev1.pdf
National Institute of Health (NIH)	US	2003	http://grants.nih.gov/grants/policy/data_sharing/
Gordon and Betty Moore Foundation (GBMF)	US	(2005) 2008	www.moore.org/docs/GBMF_Data%20Sharing%20Philosophy%20and%20Plan.pdf
Genome Canada	Canada	(2005) 2008	www.genome can a da.ca/medias/PDF/EN/Data Release and ResourceSharingPolicy.pdf
Medical Research Council Data Sharing and Preservation Policy (MRC)	UK	2006	www.mrc.ac.uk/Ourresearch/Ethicsresearchguidance/Datasharinginitiative/Policy/index.htm
Biotechnology and Biological Sciences Research Council (BBSRC)	UK	2007	www.bbsrc.ac.uk/publications/policy/data_sharing_policy.html
Wellcome Trust	UK	2007	www.wellcome.ac.uk/About-us/Policy/Policy-and-position-statements/WTX035043.htm
Department of Energy (DOE)	US	2008	http://genomicsgtl.energy.gov/datasharing
European Commission	Europe	NA	Issued Communication calling for uniform policies across Member Nations http://ec.europa.eu/research/science-society/document_library/pdf_06/communication-022007_en.pdf
European Science Foundation	Europe	NA	Researchers are expected to follow the policies of the national agencies that directly provide research funding.

chances to publish, to commercialize aspects of funded work, or to collaborate with industry. Public release of 'omics data has also been complicated by the increasing use of human subjects (27) in medical-related studies and the resulting ethical issues. Funding agencies must allay fears that data could be reused without permission or due recognition by clarifying the agency's expectations. There is currently no large-scale infrastructure ready to support data citations, but interest in this issue is growing (28).

Researchers may be limited in their ability to comply by inadequate resourcing; time-inefficient data management at the local or community level; or a lack of tools, databases or informatics expertise. Researchers must now incorporate the cost of this type of essential work into research grants effectively and consistently, and an expert pool of scientists with the requisite skills must be developed, as well as a community of biocurators (29, 30). Mechanisms for crediting data generators when their data sets are published or reused would help justify making the data public in the mind of the researcher, especially if funding decisions took into account prior good practice.

Collecting, holding, and disseminating electronic data are substantial undertakings, if considered at the global level. If policies are to be successful, information superhighway infrastructure must be built. This must involve the creation and adoption of appropriate standards that enable electronic data to be shuttled around, tools for doing the actual task, and world-class database infrastructure to hold the collective submissions. Journals, for example, will only require compliance with reporting standards when appropriate standards-compliant software tools and public repositories become available (*31*). An exemplar project already exists, the Investigation/Study/Assay (ISA) Infrastructure, which is developing standards to enable freely available tools that encompass several 'omics technologies and facilitate curation and reporting at the community level (*3*, *32*). Lack of funding for these activities has already been highlighted (*33*, *34*), and new ways of balancing streams of funding for the generation of novel data versus the protection of existing data must be found.

### The Future

We recommend that a single, brief, high-level consensus guideline serve as a template for policy documents at the funder, community, and project levels. At its heart should be the public and timely release of data. It should be based on the principle that funders and the research community must work together to develop best practice. On enforcement of policy, we suggest that, in addition to mandating the inclusion of data-sharing plans in grant applications, deposition of supporting (or ideally, all) data in appropriate databases be the rule within a specified time period in accordance with international standards. This would uphold and extend the model of "accession number for publication" that has worked well for DNA sequence data (27). "Appropriate" databases, by definition, should be secure, should be publicly accessible, and ought to have a long-term **Examples of data policies from major funding agencies in the United States and United Kingdom.** Funders are listed by the first year in which they made their policy public (in parentheses if a newer version of the policy exists). The NSF document is the Grant General Award document, rather than a formal policy. The DOE example is a program-level policy, as an agency-level policy does not yet exist. NA, not applicable.

funding horizon. This allows reviewers to focus on the science, while creating a simple way to check compliance via a URL. When funders do not have a suitable database or repository to endorse, they should attempt to find or fund one (14).

We created the BioSharing Web site to centralize and to give a higher profile to bioscience data policies and standards (35). It offers a focal point for stakeholders in data policy (i) by providing a "one-stop shop" for those seeking data policy documents and information (including information about the standards and technologies that support them) and (ii) by encouraging exchange of ideas and policy components among funders, and between funders and potential fundees. For example, a recent post covers the "Toronto" (36) and "Rome" datasharing meetings (37) that aimed to build upon the highly influential Bermuda Principles (38) and the Fort Lauderdale report (39). Ideally, this hub could spark the formation of a Bio-Sharing Consortium that would work at the global level to build essential linkages between funders and awardees and among the main research groups.

#### References and Notes

1. Big Data special issue. Nature 455, 1 (2008).

Downloaded from www.sciencemag.org on June 18, 2010

# **POLICY**FORUM

- 2. A. R. Jones et al., Nat. Biotechnol. 25, 1127 (2007).
- 3. S. A. Sansone et al., OMICS 12, 143 (2008).
- B. Smith et al., Nat. Biotechnol. 25, 1251 (2007).
- 5. C. F. Taylor et al., Nat. Biotechnol. 26, 889 (2008).
- 6. "Protocol for implementing open access data," http:// sciencecommons.org/projects/publishing/open-accessdata-protocol.
- 7. European Science Foundation (ESF), Shared Responsibilities in Sharing Research Data: Policies and Partnerships, Report of an ESF-Deutsch Forschungsgemeinschaft workshop, Padua, Italy, 21 September 2007 (ESF, Strasbourg, France, 2008).
- 8. European Commission (EC), "On scientific information in the digital age: Access, dissemination and preservation"; http://ec.europa.eu/research/science-society/document\_library/pdf\_06/communication-022007\_en.pdf
- 9. FDA, "Genomic data submission"; www.fda.gov/Drugs/ ScienceResearch/ResearchAreas/Pharmacogenetics/ ucm083641.htm.
- 10. EMEA, Guideline on Pharmacogenetics Briefing Meetings; www.emea.europa.eu/pdfs/human/ pharmacogenetics/2022704en.pdf.
- 11. EPA, Potential Implications of Genomics for Regulatory and Risk Assessment Applications at EPA; www.epa.gov/ osa/genomics.htm.

- 12. "Pistoia vision," www.pistoiaalliance.org/.
- 13. Organisation for Economic Co-operation and Development, OECD Principles and Guidelines for Access to Research Data from Public Funding (OECD, Paris, 2007); www.oecd.org/dataoecd/9/61/38500813.pdf.
- 14. B. Tiwari, D. Field, J. Snape, Nature 439, 912 (2006).
- 15. D. Field, B. Tiwari, J. Snape, PLoS Biol. 3, e297 (2005).
- 16. D. Field et al., Nat. Biotechnol. 24, 801 (2006).
- 17. Economic and Social Data Service, www.esds.ac.uk/.
- 18. D. Field, S. A. Sansone, OMICS 10, 84 (2006).
- 19. Standardizing data. Nat. Cell Biol. 10, 1123 (2008).
- Cancer Biomedical Informatics Grid, (caBIG), National 20. Cancer Institute, NIH; http://cabig.cancer.gov.
- 21. Biomedical Informatics Research Network, www.nbirn.net/.
- 22. ELIXIR, http://www.elixir-europe.org.
- 23. EC, "INSPIRE Directive," http://inspire.jrc.ec.europa.eu/ index.cfm.
- 24. DCC, www.dcc.ac.uk/.
- 25. C. Thomas, Science 324, 1632 (2009).
- 26. S. Wiley, Scientist 23, 33 (2009).
- 27. E. Pennisi, Science 324, 1000 (2009).
- 28. Earth System Science Data, www.earth-system-sciencedata.net/.
- 29. D. Howe et al., Nature 455, 47 (2008).
- 30. International Society for Biocuration, www.biocurator.org

- 31. H. Barsnes et al., Nat. Biotechnol. 27, 598 (2009).
- 32. Investigation/Study/Assay (ISA) Infrastructure for Manag-
- ing Experimental Metadata, http://isatab.sf.net. 33. C. Brooksbank, J. Quackenbush, OMICS 10, 94 (2006).
- 34. Z. Merali, J. Giles, Nature 435, 1010 (2005).
- 35. Biosharing, http://biosharing.org/.
- 36. Toronto International Data Release Workshop Authors,
- Nature 461, 168 (2009). 37. P. N. Schofield et al., Nature 461, 171 (2009).
- 38. "Summary of principles agreed at the First International Strategy Meeting on Human Genome Sequencing, Bermuda, 25 to 28 February 1996 (Human Genome Organisation, Singapore, 1996); available at www.ornl.gov/sci/ techresources/Human\_Genome/research/bermuda.shtml.
- Sharing Data from Large-Scale Biological Research 39. Projects, A System of Tripartite Responsibility, Fort Lauderdale, FL, 14 and 15 January 2003 (Wellcome Trust, 2003); available at www.wellcome.ac.uk/stellent/groups/ corporatesite/@policy\_communications/documents/ web\_document/wtd003207.pdf.

#### Supporting Online Material

www.sciencemag.org/cgi/content/full/326/5950/234/DC1

10.1126/science.1180598

# **GENOMICS** Genome Project Standards in a **New Era of Sequencing**

P. S. G. Chain, 12,3\*†§ D. V. Grafham, 4†§ R. S. Fulton, 5† M. G. FitzGerald, 6† J. Hostetler, 7† D. Muzny,<sup>8</sup>† J. Ali,<sup>9</sup> B. Birren,<sup>6</sup> D. C. Bruce,<sup>1,10</sup> C. Buhay,<sup>8</sup> J. R. Cole,<sup>3</sup> Y. Ding,<sup>8</sup> S. Dugan,<sup>8</sup> D. Field,<sup>11</sup> G. M. Garrity,<sup>3</sup> R. Gibbs,<sup>8</sup> T. Graves,<sup>5</sup> C. S. Han,<sup>1,10</sup> S. H. Harrison,<sup>3\*</sup> S. Highlander,<sup>8</sup> P. Hugenholtz,<sup>1</sup> H. M. Khouri,<sup>12</sup> C. D. Kodira,<sup>6\*</sup> E. Kolker,<sup>13,14</sup> N. C. Kyrpides,<sup>1</sup> D. Lang,<sup>12</sup> A. Lapidus,<sup>1</sup> S. A. Malfatti,<sup>12</sup> V. Markowitz,<sup>15</sup> T. Metha,<sup>6</sup> K. E. Nelson,<sup>7</sup> J. Parkhill,<sup>4</sup> S. Pitluck,<sup>1</sup> X. Qin,<sup>8</sup> T. D. Read,<sup>16</sup> J. Schmutz,<sup>17</sup> S. Sozhamannan,<sup>18</sup> P. Sterk,<sup>11</sup> R. L. Strausberg,<sup>7</sup> G. Sutton,<sup>7</sup> N. R. Thomson,<sup>4</sup> J. M. Tiedje,<sup>3</sup> G. Weinstock,<sup>5</sup> A. Wollam,<sup>5</sup> Genomic Standards Consortium Human Microbiome Project Jumpstart Consortium, ‡ J. C. Detter<sup>10</sup>†‡

or over a decade, genome sequences have adhered to only two standards that are relied on for purposes of sequence analysis by interested third parties (1, 2). However, ongoing developments in revolutionary sequencing technologies have resulted in a redefinition of traditional whole-genome sequencing that requires reevaluation of such standards. With commercially available 454 pyrosequencing (followed by Illumina, SOLiD, and now Helicos), there has been an explosion of genomes sequenced under the moniker "draft"; however, these can be very poor quality genomes (due to inherent errors in the sequencing technologies, and the inability of assembly programs to fully address these errors). Further, one can only infer that such

draft genomes may be of poor quality by navigating through the databases to find the number and type of reads deposited in sequence trace repositories (and not all genomes have this available), or to identify the number of contigs or genome fragments deposited to the database. The difficulty in assessing the quality of such deposited genomes has created some havoc for genome analysis pipelines and has contributed to many wasted hours. Exponential leaps in raw sequencing capability and greatly reduced prices have further skewed the time- and cost-ratios of draft data generation versus the painstaking process of improving and finishing a genome. The result is an everwidening gap between drafted and finished genomes that only promises to continue (see

<sup>1</sup>U.S. Department of Energy Joint Genome Institute. <sup>2</sup>Lawrence Livermore National Laboratory. <sup>3</sup>Michigan State University. <sup>4</sup>The Sanger Institute. <sup>5</sup>Washington University School of Medicine. <sup>6</sup>The Broad Institute. <sup>7</sup>J. Craig Venter Institute. <sup>8</sup>Baylor College of Medicine. <sup>9</sup>Ontario Institute for Cancer Research. <sup>10</sup>Los Alamos National Laboratory. <sup>11</sup>Natural Environmental Research Council Centre for Ecology and Hydrology. <sup>12</sup>National Center for Biotechnology Information. <sup>13</sup>Seattle Children's Hospital and Research Institute. <sup>14</sup>University of Washington School of Medicine. <sup>15</sup>Lawrence Berkeley National Laboratory. <sup>16</sup>Emory GRA (Georgia Research Alliance) Genomics Center. <sup>17</sup>HudsonAlpha Institute. <sup>18</sup>Naval Medical Research Center.

\*Full affiliations are available on Science Online. †Finishing in the Future Working Group members. ‡These authors contributed equally to organizing this work. \$Authors for correspondence. E-mail: pchain@lanl.gov (P.S.G.C.); dg1@sanger.ac.uk (D.V.G.)

### More detailed sequence standards that keep up with revolutionary sequencing technologies will aid the research community in evaluating data.

the figure, page 236); hence, there is an urgent need to distinguish good from poor data sets.

The sequencing institutes and consortia whom we represent believe that a new set of standards is required for genome sequences. The following represents community-defined categories of standards that better reflect the quality of the genome sequence, based on our understanding of the technologies, available assemblers, and efforts to improve upon drafted genomes. Due to the increasingly rapid pace of genomics, we avoided rigid numerical thresholds in our definitions to take into account products achieved by any combination of technology, chemistry, assembler, or improvement and/or finishing process.

Standard Draft: minimally or unfiltered data, from any number of different sequencing platforms, that are assembled into contigs. This is the minimum standard for a submission to the public databases. Sequence of this quality will likely harbor many regions of poor quality and can be relatively incomplete. It may not always be possible to remove contaminating sequence data. Despite its shortcomings, Standard Draft is the least expensive to produce and still possesses useful information.

High-Quality Draft: overall coverage representing at least 90% of the genome or target region. Efforts should be made to exclude contaminating sequences. This is still a draft assembly with little or no manual review of the product. Sequence errors and misassem-