

# Challenges and Opportunities

SCIENTIFIC INNOVATION HAS BEEN CALLED ON TO SPUR ECONOMIC recovery; science and technology are essential to improving public health and welfare and to inform sustainability; and the scientific community has been criticized for not being sufficiently accountable and transparent. Data collection, curation, and access are central to all of these issues. For this reason, *Science* has joined with colleagues from our sister publications *Science Signaling*, *Science Translational Medicine*, and *Science Careers* to provide a broad look at the issues surrounding the increasingly huge influx of research data. The entire collection is compiled online at [www.sciencemag.org/special/data/](http://www.sciencemag.org/special/data/). As you will discover, two themes appear repeatedly: Most scientific disciplines are finding the data deluge to be extremely challenging, and tremendous opportunities can be realized if we can better organize and access the data.

Our authors explore data issues that apply to specific fields as well as challenges shared between fields. These articles clearly show that the challenges are difficult and growing. We have recently passed the point where more data is being collected than we can physically store (see Hilbert *et al.*, published online). This storage gap will widen rapidly in data-intensive fields. Thus, decisions will be needed on which data to archive and which to discard. A separate problem is how to access and use these data. Many data sets are becoming too large to download. Even fields with well-established data archives, such as genomics, are facing new and growing challenges in data volume and management. And even where accessible, much data in many fields is too poorly organized to enable it to be efficiently used.

To delve deeper into these issues, *Science* polled our peer reviewers from last year about the availability and use of data. We received about 1700 responses, representing input from an international and interdisciplinary group of scientific leaders. About 20% of the respondents regularly use or analyze data sets exceeding 100 gigabytes, and 7% use data sets exceeding 1 terabyte. About half of those polled store their data only in their laboratories—not an ideal long-term solution. Many bemoaned the lack of common metadata and archives as a main impediment to using and storing data, and most of the respondents have no funding to support archiving.

Many of the responders indicated that they seek or would like additional help in analyzing the data that they had collected. If we can use and reuse scientific data better, the opportunities, as indicated in many examples in this special section, are myriad. Large integrated data sets can potentially provide a much deeper understanding of both nature and society and open up many new avenues of research. And they are critical for addressing key societal problems—from improving public health and managing natural resources intelligently to designing better cities and coping with climate change.

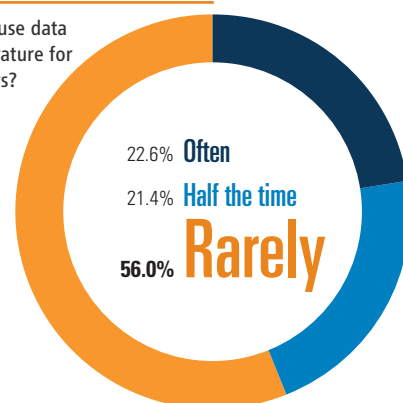
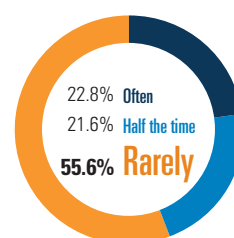
To realize these opportunities, many of the articles in this collection speak of changing the culture of science and the practices of scientists, as well as recognizing the growing responsibility for much better data stewardship. Several of the pieces illustrate steps toward these goals. But it is clear that organized effort and leadership are needed from funders, societies, journals, educators, and individual scientists—and from society at large.

We hope that this collection spurs additional thinking and catalyzes new efforts in dealing with these critical issues. As a start, we invite you to share your thoughts at [talk.sciencemag.org](http://talk.sciencemag.org), where you can also contribute to our poll.

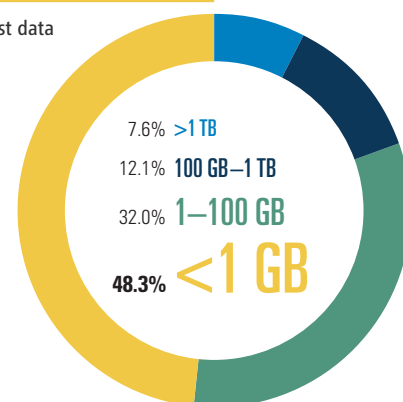
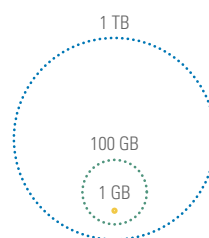
— SCIENCE STAFF

How often do you access or use data sets from the published literature for your original research papers?

From archival databases?

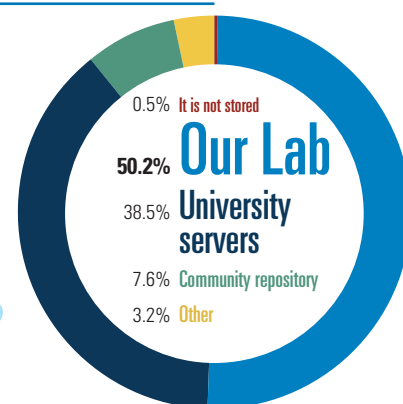


What is the size of the largest data set that you have used or generated in your research?



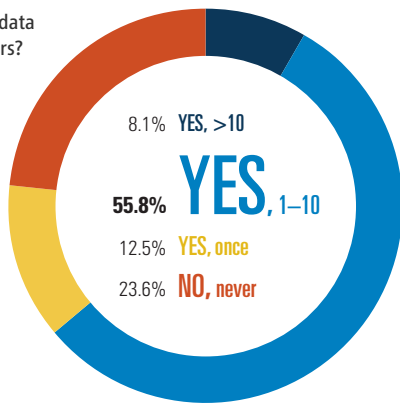
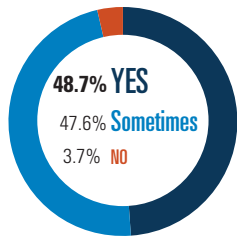
Where do you archive most of the data generated in your lab or for your research?

“Even within a single institution there are no standards for storing data, so each lab, or often each fellow, uses ad hoc approaches.”



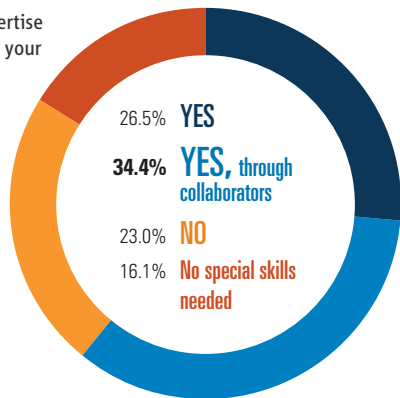
Have you asked colleagues for data related to their published papers?

If you answered yes, have the appropriate data been provided?



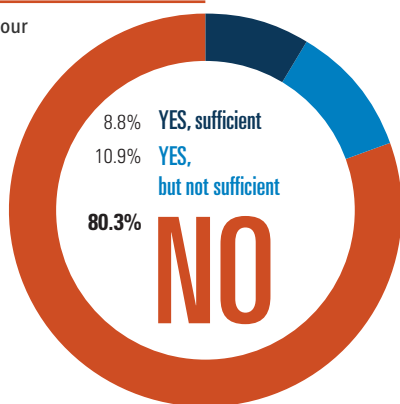
Do you have the necessary expertise in your lab or group to analyze your data in the way you want?

“The next few years [particularly in medicine] the volume of data we need to analyze will expand exponentially.”



Is there sufficient funding for your lab or research group for data curation?

“There are many tales of early archaeologists burning wood from the ruins to make coffee. If we fail to curate the environmental archives we collect from nature at public expense, we essentially repeat those mistakes.”



## CONTENTS

### News

- 694 Rescue of Old Data Offers Lesson for Particle Physicists
- 696 Is There an Astronomer in the House?
- 698 May the Best Analyst Win

### Perspectives

- 700 Climate Data Challenges in the 21st Century  
*J. T. Overpeck et al.*
- 703 Challenges and Opportunities of Open Data in Ecology  
*O. J. Reichman et al.*
- 705 Changing the Equation on Scientific Data Visualization  
*P. Fox and J. Henderl*
- 708 Challenges and Opportunities in Mining Neuroscience Data  
*H. Akil et al.*
- 712 The Disappearing Third Dimension  
*T. Rowe and L. R. Frank*
- 714 Advancing Global Health Research Through Digital Technology and Sharing Data  
*T. Lang*
- 717 More Is Less: Signal Processing and the Data Deluge  
*R. G. Baraniuk*
- 719 Ensuring the Data-Rich Future of the Social Sciences  
*G. King*
- 721 Metaknowledge  
*J. A. Evans and J. G. Foster*
- 725 Access to Stem Cells and Data: Persons, Property Rights, and Scientific Progress  
*D. J. H. Mathews et al.*
- 728 On the Future of Genomic Data  
*S. D. Kahn*

See also:

### Editorial

- 649 *Making Data Maximally Available*  
B. Hanson, A. Sugden, and B. Alberts

### News Focus

- 662 *What Would You Do?*  
J. Couzin-Frankel
- 666 *Will Computers Crash Genomics?*  
E. Pennisi
- 669 *Drag-and-Drop Virtual Worlds*  
R. Service

### Books

- 676 *Bounds and Vision*  
M. A. Porter

### Policy Forum

- 678 *Measuring the Results of Science Investments*  
J. Lane and S. Bertuzzi

### Science Express Research Article\*

- The World's Technological Capacity to Compute, Store, and Communicate Information*  
M. Hilbert and P. López

### Science Signaling\*

- Conquering the Data Mountain*  
N. R. Gough and M. B. Yaffe
- Effective Representation and Storage of Mass Spectrometry-Based Proteomic Data Sets for the Scientific Community*  
J. V. Olsen and M. Mann
- The Potential Cost of High-Throughput Proteomics*  
F. M. White

- Integrating Multiple Types of Data for Signaling Research: Challenges and Opportunities*  
H. S. Wiley

- Setting the Standards for Signal Transduction Research*  
J. Saez-Rodriguez et al.

- Visual Representation of Scientific Information*  
B. Wong

### Science Translational Medicine\*

- Power to the People: Participant Ownership of Clinical Trial Data*  
S. F. Terry and P. F. Terry

- Electronic Consent Channels: Preserving Patient Privacy Without Handcuffing Researchers*  
R. H. Shelton

### Science Careers\*

- More Than Words: Biomedical Ontologies Provide New Scientific Opportunities*  
C. Wald

- Surfing the Tsunami*  
E. Pain

- Sharing Data in Biomedical and Clinical Research*  
K. Travis

\*These items, plus a related podcast and online discussion, are available at [www.sciencemag.org/special/data/](http://www.sciencemag.org/special/data/)