# Eos

# The Next 100

(Mg,Fe)SiO3

AGU 100

ADVANCING EARTH
AND SPACE SCIENCE

# A GEODATA FABRIC

## FOR THE 21st CENTURY

WE HAVE THE POTENTIAL TO TRANSFORM OUR UNDERSTANDING OF EARTH—IF WE CAN JUST FIGURE OUT HOW TO HARNESS EVER GROWING DATA STREAMS.

By Jeff de La Beaujardière

**THE NATURE** of scientific data has changed considerably since AGU was founded a century ago. Observations then were manual and laborious, data were recorded in paper notebooks and on photographic plates, and the most powerful "computer" was perhaps the Powers Accounting Machine, an electromechanical device for tabulating U.S. Census Bureau data recorded on punched cards. In 1919, the Eddington experiment provided the first observational confirmation of Einstein's theory of general relativity by measuring the deflection of light from stars during a total solar eclipse; the entirety of the published data comprises 20 tables, 2 diagrams, and 1 low-resolution black-and-white photograph [*Dyson et al.,* 1920].

As the digital age arrived and the technology for recording observations progressed, the geosciences and other disciplines began to face the "big data" problem, often characterized by four V words:

• Researchers generate enormous data **volumes** from new observing systems and from simulations run on supercomputers. The sheer size of contemporary data sets makes them expensive to store, difficult to access in more than small subsets, and nearly impossible to ship to all the institutions that would like copies.

• The bewildering **variety** of data sets from sources that use different file formats and organizational systems is a barrier t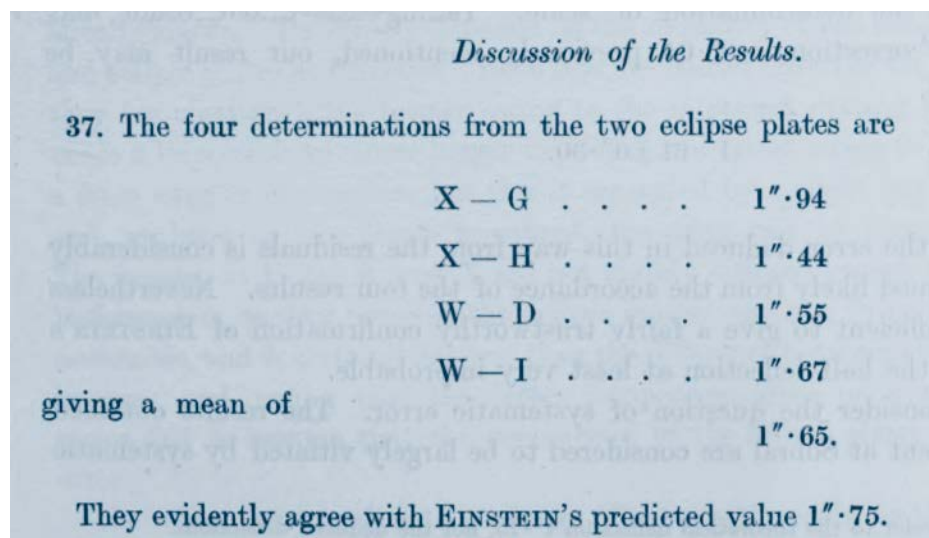o performing interdisciplinary research and analyzing phenomena that have multiple signals observed by different platforms. In addition, there exist to this day myriad small, manually collected observational data sets that are difficult to integrate [*Genova and Horstmann,* 2016].

• Real-time data are being collected, processed, and disseminated at ever increasing **velocities,** particularly in emergency situations such as earthquakes, fires, and other disasters.

• **Variability** caused by surges in data arrival or user requests poses a challenge for facilities, which must have enough capacity to handle the highest loads yet may be idle during down times.

These issues have escalated into major challenges to performing scientific research and providing accessible, usable information to decision-makers. In just 3 years, the NASA–Indian Space Research Organisation Synthetic Aperture Radar (NISAR) Earth-observing satellite is expected to generate some 85 terabytes of data per day. Supercomputers in the "exascale" range, capable of performing at least $10^{18}$ operations per second, are under contract to be deployed by 2021 at Argonne and Oak Ridge national laboratories; these will be 1,000 times more powerful than the petascale machines of only a decade ago.

Unless geoscience research in the next century is to be hamstrung by the very data it is collecting, we'll have to find the answers to two questions: How can we provide storage and access for big data? And



*One of the 20 tables needed in "A Determination of the Deflection of Light by the Sun's Gravitational Field, from Observations made at the Total Eclipse of May 29, 1919," to confirm Einstein's theory of general relativity. Credit: Dyson, et al., 1920, https://doi.org/10.1098/rsta.1920.0009.*

more important, how can we enable "science at scale," such that researchers and other users can work with large, multisource data sets without getting lost in a tangle of incompatible systems?

**The Geoscience Advantage**
Fortunately, we benefit from two crucial facts. First, the geosciences are not alone in facing the big-data problem. In astronomy, the Square Kilometre Array—an enormous radio telescope with installations in South Africa and Australia—is expected to generate 160 terabytes of raw data per second. The genomics community estimates a need for at least 2,000 petabytes of storage by 2025. In the private sector, Facebook had already accumulated at least 300 petabytes of data as of 2014. Our community will therefore be able to leverage work by others in the same predicament.

Second, the very fact that we do geoscience provides a useful organizing framework: Much of our data are, by definition, based on a time and a place. Every Earth observation, every numerical simulation grid point, has an associated temporal range, a position or region on the planet, and possibly an elevation or depth range. Instead of individual files and collections, we could organize these data within a multidimensional "Geodata Fabric." A good analogy is Google Maps, which integrates detailed road, transit, boundary, river, and hiking trail data from multiple sources, along with satellite imagery, elevation data (terrain layer), in situ data (photos and Street View), ancillary information (businesses and facilities), crowdsourced content (reviews of businesses and places), personal annotations (favorite places), real-time information (traffic), and even basic computation (driving directions). Can we do the same for Earth science?

Creation of a Geodata Fabric requires the geosciences to take a huge leap from where we are now, barely past the stage



*The Powers Accounting Machine.*
*Credit: Mahlum/Wikimedia Commons*

of paper maps and guidebooks, with disjointed web servers each providing only a tiny portion of the vast body of environmental data. We need a more unified approach such that each data provider—whether in the atmosphere, land surface, seismology, hydrology, oceanography, or cryosphere domain—can contribute to a shared and commonly accessible framework. The same concept could be extended to domains with differing coordinate systems such as other planets or interplanetary space. Some work that clearly demonstrates the usefulness of the approach has already been done along these lines—notably, the Open Data Cube project in Australia—but the concept must be extended to the entire planet and include more than just satellite imagery.

**Elements of the Solution**
If we want to create this Geodata Fabric, we'll need to rely on automation and standardization at every step as we acquire data; perform initial processing; move them; store them; and enable discovery, access, and analysis. Let's walk through how we might approach four of those steps.
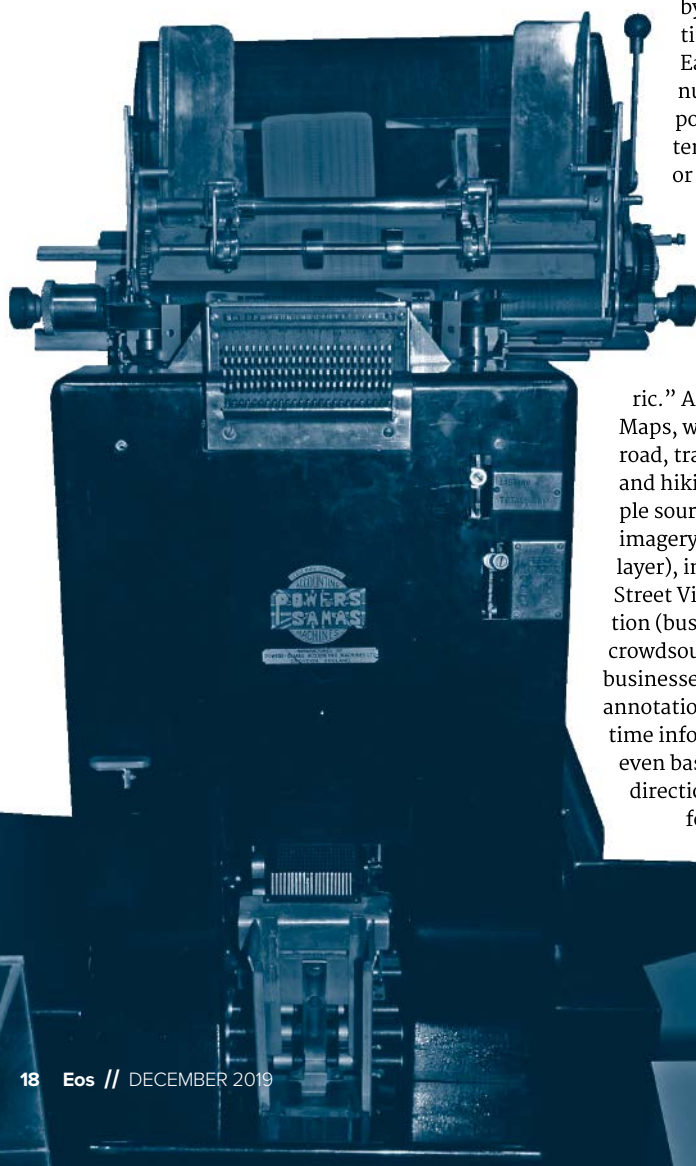
### 1. A New Type of Storage for Big Data
If you were to log on to any computer storing scientific data today, the organization of the information would typically look a lot like that on your personal computer, with a hierarchy of folders and files. However, these systems become sluggish when they contain billions of files. A better solution, known as object storage, has been adopted by most private companies with huge data needs such as Facebook, YouTube, Netflix, Google, and Amazon.

Object storage uses standard disk drives. Instead of a folder hierarchy, however, it simply has a pool of capacity that can be expanded as needed, has customizable metadata for each object, does not need to actively monitor all the files in the system, and recovers much more easily from the failure of individual drives. Because of this, it is simpler to maintain and better suited for large data volumes.

Each item saved in object storage has a unique ID and a user-defined name, either of which can be used to retrieve the object. For example, an early draft of this article was backed up on Google Drive with the ID "1_8Q2pMN6BHI2N8F5qCaVO5ASdGuGGSOy" and the name "MyDrive/work/AGU/Eos _article.docx." Although it *appears* to be in a folder hierarchy for my personal conve-

nience, it is merely a blob of data in a huge storage pool.

Lawrence Berkeley National Laboratory in California is among those institutions recommending the use of object storage for big data [*Lockwood et al.,* 2017]. To make the switch, developers will need to slightly modify software that assumes that data are in a traditional file system.

Object storage can be located at the data collector's facility or in the cloud. We use the cloud daily in our personal lives for such data as email, calendars, documents, and photos, but most geoscience institutions still rely primarily on storage hosted on-site. Cloud storage offers three key advantages: The maintenance of the hardware is outsourced to professionals, the amount of available storage is essentially unlimited, and public access can be easily granted.

Whether on-site or in the cloud, object storage will be necessary to accommodate vast volumes of information in the Geodata Fabric.

## 2. Move the Data Once Or Never

We cannot create a Geodata Fabric if each swath of data is isolated. Instead, we need to perform initial raw data processing at the source and then consolidate the usable information. However, moving huge data sets across the Internet is a slow endeavor. Even with such resources as the high-bandwidth Internet2 connecting some research institutions, overall performance is no better than the slowest link in the path. Traditionally, subsetting services have allowed people to "clip and ship" only what they need, but this doesn't enable large-scale science on multiple large data sets.

Some companies now offer bulk data transfer appliances—literally, disks packed in crates—that can be transported as freight. In 2018, a U.S. Geological Survey facility on the Big Island of Hawaii, threatened by possible lava flow from Kīlauea volcano, used an Amazon Web Services (AWS) Snowball physical transport device to quickly copy critical data to the mainland for safekeeping. In 2017, the commercial satellite company DigitalGlobe was the first to use an even larger AWS capability called Snowmobile. The company copied the contents of 8,700 tapes—nearly 100 petabytes of data—into a 16-meter, fully powered shipping container for transfer to the cloud in a matter of months by truck instead of years via the Internet.

Once the data are consolidated, we can bring the computing to the data such that only the output of the analysis software need be sent to the user.

## 3. The Cloud Advantage

Traditionally, science facilities provide public-facing web interfaces for preconfigured types of analyses, and perhaps allow authorized users to log in directly and run whatever analysis code they prefer. This type of access does not support customized analysis by unknown external users.

The cloud provides several important advantages. Users can operate directly on the data using whatever software they choose. Processing power can be scaled up or down as needed to accommodate large analyses or spikes in demand. Institutions can focus on science instead of operating hardware and can even take advantage of "managed services" ranging from databases to text, image, and video analyzers and, as of earlier this year, entire satellite ground stations.

A potentially revolutionary concept is known as "serverless computing": Instead of keeping a server running (and paid for) constantly to handle only occasional requests, the cloud vendor runs a shared pool of servers on which you can perform brief computations on data as needed, paying only for the amount of time and memory your function uses. This temporary usage is analogous to renting a Zipcar for an

WHETHER ON SITE OR IN THE CLOUD, OBJECT STORAGE WILL BE NECESSARY TO ACCOMMODATE VAST VOLUMES OF INFORMATION IN THE GEODATA FABRIC.



*The Apple II personal computer. Credit: FozzTexx/ Wikimedia Commons*

hour instead of owning a vehicle that spends most of its time in the garage. For organizations that manage large data collections, there are clear benefits to outsourcing infrastructure and focusing only on data stewardship, access, and analysis.

NASA has been a leader in this area, using commercial cloud servers to store, archive, process, distribute, and manage large volumes of Earth-observing mission data—predicted to be over 45 petabytes per year by 2022 and over 245 petabytes total volume by 2025. The agency's Common Metadata Repository and Earthdata Search services are now running on AWS, with more to come. Kathleen Baynes, NASA system architect, told me, "Beyond consolidating software systems and streamlining processes, we anticipate this effort will provide exciting opportunities to further expand the impact of NASA's Earth Science holdings: introducing new paradigms for interacting with data, improving interoperability, facilitating innovative research, and helping to drive from data toward knowledge."

There are some disadvantages to using cloud computing for analysis. Chief among these is the pay-as-you-go model, with costs not fully known in advance. This is a

*The Square Kilometer Array (artist's impression). Credit: Mathieu Isidro (SKA)*



hurdle in particular for government agencies because of legislation that forbids spending more funds than were allocated by Congress, but it can be reduced or eliminated by prepayment, monitoring, throttling usage if necessary, and efficient system design. Storage costs can be minimized by moving infrequently used data to less expensive tiers of storage or even by discarding unimportant data. Egress costs can be minimized by enabling and encouraging users to compute directly on the data. Computing costs can be minimized through more efficient code, using discounted computing time when the vendor has periods of low demand, and using managed services and serverless functions.

Cloud computing is quickly becoming a viable approach for most science being done today and will be essential for moving data out of institutional silos and into a Geodata Fabric.

## 4. The Need for Simplicity

Naturally, researchers tend to store the data they collect in ways that make the most sense for their projects. But if we want to establish an easily accessible and broadly usable Geodata Fabric, we must improve standardization and enable a higher level of abstraction. A user should be able simply to ask for—or directly visualize—a desired data set, time range, and area of interest while software behind the scenes automatically provides what was requested.

The National Science Foundation–funded Pangeo project, for example, is "a community promoting open, reproducible, and scalable science." Pangeo allows users to combine a variety of open-source tools such as the Zarr format to break multidimensional data into chunks, the Dask library to read many chunks simultaneously, the Xarray package to address data at an abstract level, and Jupyter Notebook for customizable web-based analysis workflows. The goal is to enable users to explore data interactively without worrying about storage details and to efficiently perform computations over large data sets.

Elizabeth Maroon, a project scientist at the National Center for Atmospheric Research, told me that "for cutting-edge climate science, we need to use large ensembles and explore increasing resolution, and that makes our data sets huge. We are using the Pangeo tools to calculate ocean circulation metrics in our new high-resolution climate model simulations. Without the parallel computation made rel-

atively easy by the Dask package, these metrics would take prohibitively long to calculate."

As data volumes grow, it is becoming increasingly difficult to have knowledgeable people inspect all the data for interesting phenomena. We need machine learning—large-scale, automated analysis of big data—to become commonplace, which can happen only once we standardize and consolidate data storage and access.

Planetary scientists have demonstrated the ability of machine learning to find novel features in multispectral images from NASA's Mars Curiosity rover [*Kerner et al.,* 2019]. Furthermore, some data are of little long-term value, such as long periods of undersea video showing only murky seabed punctuated by the occasional arrival of an interesting creature. Machine learning may prove useful in isolating the most relevant subsets of the data and allowing humans to decide whether to discard the rest.

### The Future of Geodata

Our capabilities to collect and store data have evolved greatly in the past century, from penciling observations into paper notebooks in the field to using automated sensors injecting information directly into databases. The cost to store a given amount of digital data has decreased thanks to the increasing density with which we can store it, from the $50,000 5-megabyte IBM RAMAC 350 in 1956 to $50 for a 1-terabyte consumer-level hard drive today. Unfortunately, that rate of storage density is flattening while data volumes continue to grow, but new storage technologies will doubtless emerge over the next century. For example, the University of Southampton in the United Kingdom is working on "memory crystals" that could potentially store hundreds of terabytes of information by etching nanoscale-sized structures in quartz [*Kazansky et al.,* 2016].

Whatever the future holds, it is clear that our geoscience data centers must evolve away from traditional silos of in-house systems offering access only to their own data. We must enable science on high-volume, high-variety, high-velocity, and high-variability data by uniting them from multiple sources, standardizing at a higher level of abstraction, and moving the computation to the data. Embracing a concept like the Geodata Fabric would enable researchers to focus more on the science and less on the plumbing.

We are in dire need of better understanding our constantly changing planet, and that requires that we establish better pathways to information.

### References

Dyson, F. W., A. S. Eddington, and C. R. Davidson (1920), A determination of the deflection of light by the Sun's gravitational field, from observations made at the solar eclipse of May 29, 1919, *Philos. Trans. R. Soc. London, Ser. A, 220,* 291–333, https://doi .org/10.1098/rsta.1920.0009.

Genova, F., and W. Horstmann (2016), Long tail of data, e-IRG Task Force Report, Eur. e-Infrastruct. Reflection Group, The Hague, Netherlands, e-irg.eu/documents/10920/238968/ LongTailOfData2016.pdf.

Kazansky, P. G., et al. (2016), Eternal 5D data storage via ultrafast-laser writing in glass, SPIE Newsroom, https://doi.org/10.1117/2 .1201603.006365.

Kerner, H., et al. (2019), Novelty detection for multispectral images with application to planetary exploration, in *Proceedings of the Thirty-First AAAI Conference on Innovative Applications of Artificial Intelligence,* pp. 9,484–9,491, Assoc. for the Adv. of Artif. Intell., Palo Alto, Calif., hannah-rae.github.io/files/Kerner _et_al_2019_IAAI.pdf.

Lockwood, G. K., et al. (2017), Storage 2020: A vision for the future of HPC storage, *Rep. LBNL-2001072,* Lawrence Berkeley Natl. Lab., Berkeley, Calif., escholarship.org/uc/item/744479dp.

### Author information

**Jeff de La Beaujardière** (jeffdlb@ucar.edu), National Center for Atmospheric Research, Boulder, Colo.

▶**Read the full story at bit.ly/Eos -geodata-fabric**

## EMBRACING A CONCEPT LIKE THE GEODATA FABRIC WOULD ENABLE RESEARCHERS TO FOCUS MORE ON THE SCIENCE AND LESS ON THE PLUMBING.

monsit/Depositphotos