# Exploring Approaches for Large Data in Seismology: User and Data Repository Perspectives

Javier Quinteros[*1], Jerry A. Carter[2], Jonathan Schaeffer[3], Chad Trabant[2], and Helle A. Pedersen[3,4]

## Abstract

**New data acquisition techniques are generating data at much finer temporal and spatial resolution, compared to traditional seismic experiments. This is a challenge for data centers and users. As the amount of data potentially flowing into data centers increases by one or two orders of magnitude, data management challenges are found throughout all stages of the data flow.**

**The Incorporated Research Institutions for Seismology—Réseau sismologique et géodésique français and GEOForschungsNetz data centers—carried out a survey and conducted interviews of users working with very large datasets to understand their needs and expectations. One of the conclusions is that existing data formats and services are not well suited for users of large datasets. Data centers are exploring storage solutions, data formats, and data delivery options to meet large dataset user needs. New approaches will need to be discussed within the community, to establish large dataset standards and best practices, perhaps through participation of stakeholders and users in discussion groups and forums.**

## Introduction

New methods of measuring ground motion are significantly reducing the cost of data collection. Two new types of equipment strongly contribute to this cost reduction. The first is nodal data, that is, data from experiments with a very high number of observation points using low-cost sensors, are now becoming common. As an example, more than 5200 high-frequency sensors were deployed over a period of six months in and around Long Beach, California (e.g., Lin *et al.*, 2013). The second is distributed acoustic sensing (DAS) technology, using fiber-optic cables, which is currently being tested and deployed in many locations. As an example, Jousset *et al.* (2018) deployed a 15 km long fiber-optic cable layout on Reykjanes Peninsula, southwest Iceland, with a distance of 4 m between sampling points, to study structural features in the Reykjanes Oblique Rift.

Both types of equipment can generate tens to hundreds of terabytes of data, at a small fraction of the cost of using traditional seismometers and geophones for equivalent data volumes. As the price of collecting data becomes drastically cheaper, a corresponding and dramatic increase in the volume of data being collected provides a potential scientific bonanza.

As dataset sizes increase into the range of tens to hundreds of terabytes; however, barriers to storing, transporting, and processing the data begin to appear. The largest collections of openly available seismological data (e.g., as managed by Incorporated Research Institutions for Seismology [IRIS],

Réseau sismologique et géodésique français [RESIF], and GEOForschungsNetz [GEOFON]) measure the volume of their decades-spanning repositories in hundreds of terabytes, yet some new (e.g., DAS) experiments are gathering data volumes well in excess of a hundred terabytes—a significant fraction of the total data volumes presently stored at the seismological data centers. One might suggest that just a small fraction of the data that are collected should be preserved as was suggested decades ago when broadband digital data was first introduced, but tremendous scientific value has been found in the continuous data that was recorded. Storing these data in perpetuity is a significant and daunting challenge, as is delivering datasets exceeding a few tens of terabytes. High-performance and high-throughput computational resources (HPC and HTC) are increasingly necessary to process these data, and these resources are neither typically provided by the repositories, nor are they collocated with the repositories.

This explosion in data volumes is just beginning, and the data centers that have traditionally been the repositories of data

1. Helmholtz-Zentrum Potsdam Deutsches GeoForschungsZentrum GFZ, Telegrafenberg, Potsdam, Germany; 2. Incorporated Research Institutions for Seismology (IRIS), Seattle, Washington, U.S.A.; 3. Univ. Grenoble Alpes, CNRS, IRD, INRAE, Météo France, OSUG, Grenoble, France; 4. Univ. Grenoble Alpes, Univ. Savoie Mont Blanc, CNRS, IRD, UGE, ISTerre, 38000 Grenoble, France

*Corresponding author: javier@gfz-potsdam.de

for the entire community are not only being asked to host these data for the wider research community but also to accommodate access to the computational resources that are needed to process these data sets. Appropriate data management practices by the data centers must address large data transport, reduced-volume derivative data products, increased access to HPC and HTC, and data formats that are HPC and HTC-friendly.

IRIS, RESIF, and GEOFON, all of which are dedicated to providing free and unrestricted access to their data holdings, have attempted to identify the needs of the community and to look for common solutions and best practices for managing very large data sets, while maintaining their traditional data services. We begin by presenting the results of soliciting user expectations for large data services through a survey of large data providers and users. Next, we describe the challenges posed by large data from its submission to a data center, to archiving, format considerations, data distribution issues, and finally to the diversity of use cases that data centers provide and will need to provide to serve the research community. The conclusions of this article provide some possible solutions and strategies for dealing with the challenges of accommodating large data sets. As a community, we must also recognize the environmental impact of storing, processing, and transporting large datasets (e.g., carbon footprint due to energy consumption, water usage, pollution from backup generators). The strategy that is eventually implemented by the seismological data centers must take this impact into account and minimize it when at all possible.

## User Expectations

To better serve their users, data centers evaluate the data requests that they receive to improve their existing services, and interact with the scientific research and education communities to discover future trends. This provides insight into popular data selection parameters, the data formats requested, and the services that the community uses.

### User survey

A survey was conducted of selected and self-identified large data users; 37 responses were received. The survey responses indicate that researchers anticipate the following:

- Large volume use, from 1 to 300+ terabytes, of traditional data (e.g., broadband seismic) and newer data types such as nodal deployments or DAS.
- Use of both existing data access tools and mechanisms (web services, Python-based clients, miniSEED), and newer data access and processing such as Hierarchical Data Format version 5 (HDF5), Zarr, xarray, especially for larger data sets.

The respondents reported the maximum size of the datasets they expect to be working within the next 3–5 yr. We classified

them in three different categories of data volume: 21 small (1–9 TB), 11 medium (10–50 TB), and 5 large (50+ TB). We consider the small volume data users relatively well served by current data center capacities, especially, as they anticipated using well-established data formats, processing tools, and access mechanisms. Medium volume data users, however, are observed to split their data requests into many small requests. What we summarize below are the survey results primarily for the medium and large data users. The results for four important variables are shown in Figure 1.

What data type are the raw data (e.g., broadband, nodal, DAS)?

The medium data users primarily identified broadband seismic data, nodal data, and a bit of DAS, whereas large data users indicated mostly DAS, with some nodal and broadband seismic.

From which data centers do you expect to request large data volumes?

Users from all categories indicated using a wide variety of data centers, which we interpret to indicate they will access data from wherever they can get it, with no particular preference.

In which data format(s) would you prefer to work with large data (miniSEED, seismic analysis code [SAC], PASSCAL HDF5 [PH5], Adaptable Seismic Data Format [ASDF], HDF5, etc.)?

Both medium and large data users anticipated a use of a variety of data formats, but primarily miniSEED, HDF5 (PH5 and ASDF), with a few cases of Zarr (see Data and Resources). However, for users planning to work with 20+ TB, the option of miniSEED reduces considerably, in comparison to HDF5-based formats or other less standard formats.

Would you process the data with standard codes? With your own code? Using third party frameworks and libraries (e.g., ObsPy)?

Users from all categories indicated that the most common codes they expect to use would mostly be their own code, supported by ObsPy (Beyreuther et al., 2010) or MATLAB (see Data and Resources).
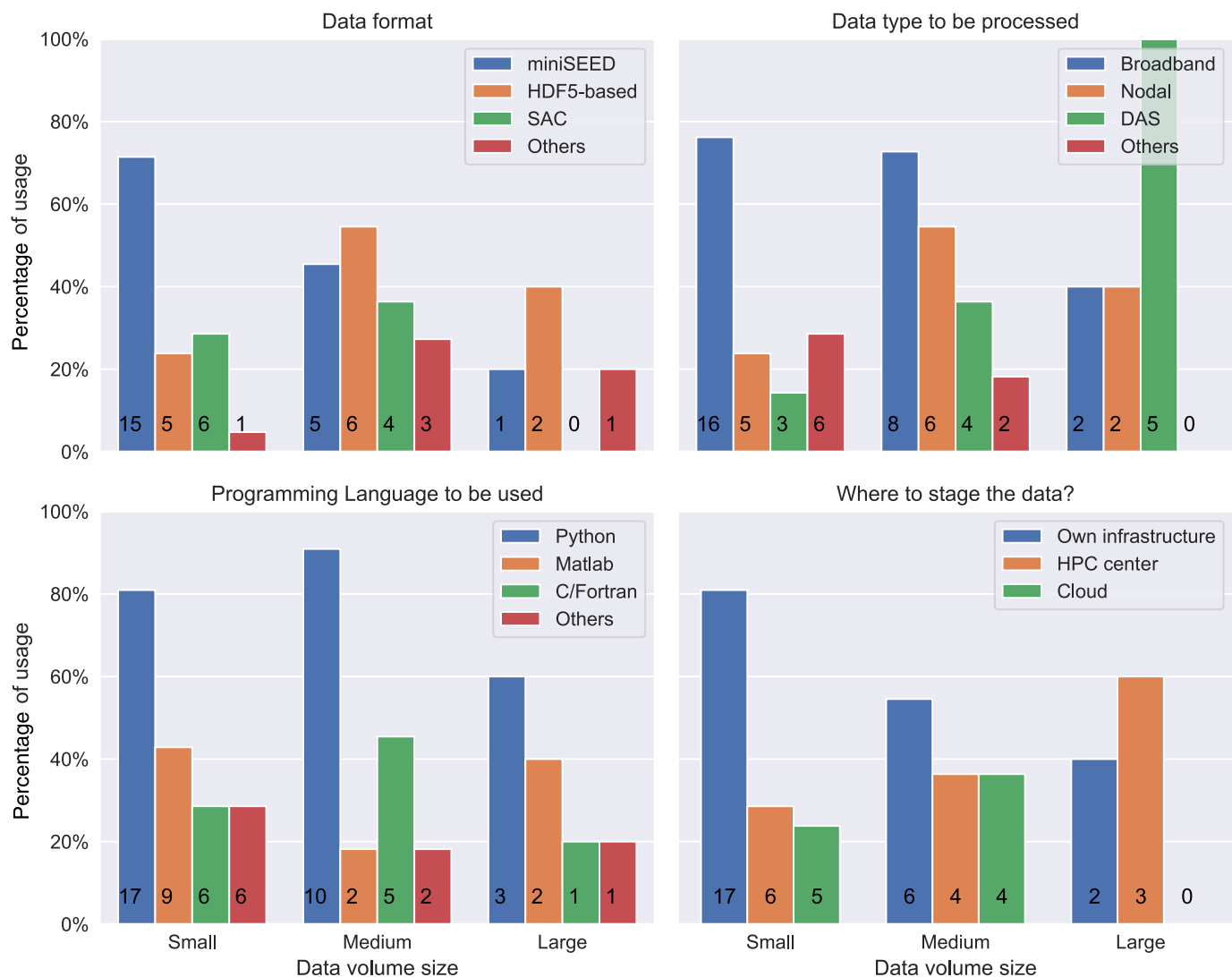
Which programming language(s) do you, or would you expect to, use for data processing and analysis?

Python is the most indicated programming language of choice, followed by MATLAB, C/C++, FORTRAN, and Julia.

Where would you prefer the data to be delivered: your own compute infrastructure, a cloud system, an HPC center?

Many medium and large users indicated that they would like to transfer data to their own computer infrastructure or their institutional infrastructure. Many also indicated a preference to deliver data to an HPC center, one of whom anticipated relying on commercial cloud providers.

Would you be willing to pay for the data processing resources (e.g., compute), if the data are available within a commercial cloud or other environment that is not free to use?

**Figure 1.** Different aspects of users' data workflows by size of datasets to be requested. Counts at the base of bars are the number of respondents. Smaller datasets tend to be requested in miniSEED format, but HDF5-based formats are preferred for larger datasets. Requests to get nodal experiment data will generate medium size datasets, whereas users plan to get very large datasets from distributed acoustic sensing (DAS) experiments. Python is the preferred programming, independent of the dataset size. Small datasets are expected to be staged in the user infrastructure, whereas larger datasets in some high-performance computation resources (HPC) facility. Only a few users mentioned their will to stage data in a cloud service. SAC, seismic analysis code. The color version of this figure is available only in the electronic edition.

Many users are willing to pay (some dependent on grant allowance) for processing. However, as in previous questions, among users working with 20+ TB, there seems to be a clear trend in favor of paying. Some pointed out that it should be an option to transfer the data to a location where they would not be required to pay for processing.

Please briefly summarize your current large data processing workflow, trouble points identified or foreseen, and how it can be improved.

Most users indicated that they use largely ad hoc data processing workflows, with a few identifying the need for parallel processing and processing-ready data access such as from HDF5 containers. The common challenge for many is efficient access to large volume of data, either due to access interface issues or limited transfer capability.

Which software would you use to access or download large data volumes (libraries, applications, etc.)?

Users of all scales report that they expect to use the existing tools and mechanisms of web services in Obspy, or directly in their own codes (in Python or Julia). A few report their anticipated use of larger platforms, such as Pangeo, direct access to object storage (e.g., S3), specialized libraries (e.g., xarray; Hoyer and Hamman, 2017), and advanced data containers such as Zarr.

When asked for other comments or their vision for large seismic data processing in the future, users identified the need to develop access and transfer mechanisms for large volumes of data. In addition, they noted the need for advanced data formats appropriate for large data, the desire for derived (reduced volume) data products for easier use, their preference for compute resources in the same system as the data (to avoid transfer), and cloud ready seismic processing software. Finally, the users expressed their desire for continued collaboration between international seismological data centers so that researchers can discover and access large volumes of data wherever it is available.

## Diversity of data requests

From the survey and from the experience of the IRIS, GEOFON, and RESIF data centers, requests for medium and large data sets span a range of time and spatial scales. On one end of the spectrum are studies that require hours or days of high sample-rate data collected from densely spaced sensors, and at the other end of the spectrum are studies that require decades of broadband data from stations distributed around the globe.

The survey clearly demonstrates that the DAS experiments represent the largest volume datasets; however, both the survey and data center experience show that medium and large data volumes from traditional seismic stations are increasingly in demand. Two typical uses of such medium and large datasets include studies based on cross-correlation techniques and studies using machine learning (ML) methods. In such cases, large datasets are created from years-long time series from a large number of individual seismic stations either at global, regional, or local scales, and the data needed for a comprehensive study may even be distributed across several data centers on different continents. In contrast, requests for node data or other experimental data with high-frequency sampling are often based on data from a single data center. Users may want access to the entire dataset or to a specific time or spatial slice. DAS datasets can be much larger than anything that traditional data centers supporting the scientific research and education communities have experience in managing. Similar to nodal data, researchers may want access to the entire dataset or to subsets of the full dataset; they may also want to process those data without transporting them or need to have them delivered to an HPC center rather than to their home institutions. Because of these processing needs, data centers should be aware of the data formats that are efficient for high-performance computation and be able to deliver data in some of those formats.

The management and organization of data and user services need to be flexible enough to handle not only the requests for medium and large datasets, but also for the diversity of standard data requests for a small number of seismic stations and short-time periods (small datasets). Data centers need to encourage and enhance services associated with data requests, using criteria such as station location, instrument type,
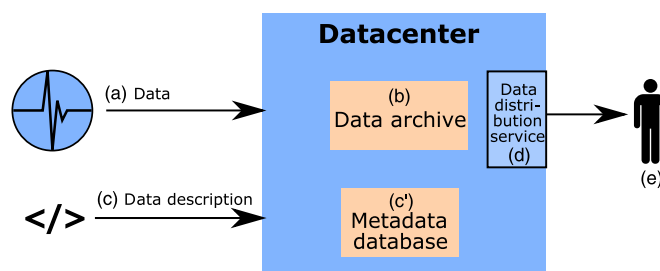


**Figure 2.** Schematic view of the different parts of the data flow in a seismological data center. The inputs are (a) data and (c) data description (metadata), submitted by the producer. Both inputs are managed in separate workflows to be hosted in (b) data archive and (c') metadata database. Data distribution services (d) need access to both repositories to handle the (e) final user's requests. Challenges regarding large data in (a–e) are described in this article. The color version of this figure is available only in the electronic edition.

sampling rate, data quality, data repository location, or more advanced parameters associated with earthquake parameters, propagation path, local geology at the seismic station, etc. MUSTANG (Casey et al., 2018) and WFCatalog (Trani et al., 2017) are some examples of services that allow users to preselect data based on quality metrics.

## The Challenges of Large Data

The seismological community has a long tradition of collaborative work between data centers and researchers. In Europe, data centers collaborate through Observatories and Research Facilities for European Seismology, whereas global collaboration in the seismological community is through the International Federation of Digital Seismograph Networks (FDSN). One of the most important achievements for the discipline is the establishment of well-defined standards under the coordination of the FDSN. These include not only the data formats, but also detailed specifications for providing data and metadata services and how Digital Object Identifiers (DOIs) should be used to identify seismic networks. Standards approved by the global community within the FDSN allow seismological data centers to support the FAIR data principles (Findable, Accessible, Interoperable, and Repeatable; Wilkinson et al., 2016). This global collaboration has contributed to a greater sharing of data among all seismologists and is a model for other disciplines.

Specifications regarding formats and services were designed in a totally different landscape of operational limitations and data usage, compared to what the community will face from now on. As soon as the amount of data increases one or two orders of magnitude, technical problems are expected. From the data center perspective, large data management challenges are found in the following stages of the data flow: (1) data submission, (2) data archival, (3) creation of

TABLE 1
**Summary of Data Formats Mentioned in the Survey and Some of Their Relevant Advantages and Disadvantages**

| Format | Comments Regarding Use in Large Data Management |
| --- | --- |
| miniSEED<br>Standard for the Exchange of Earthquake Data (SEED) (2012) | miniSEED is an international seismological standard for data exchange that is also commonly used for archiving. It was designed for independent, single time series, usually very small chunk sizes, and not for processing. |
| PH5<br>(Hess *et al.*, 2018) | PH5 is an HDF5-based format. A rich toolbox has been provided by Portable Array Seismic Studies of the Continental Lithosphere (PASSCAL), but data subsetting is not possible without high-level translation, and streaming the base format is not supported. |
| ASDF<br>(Krischer *et al.*, 2016) | ASDF is an HDF5-based format. Data subsetting is not possible without high-level translation, and streaming the base format is not supported. |
| Zarr<br>(see Data and Resources) | Zarr is a format and a python library to manage numerical arrays. Data subsetting is allowed, and it is a good design for HPC. Zarr has not been used much in the geophysical community and is only for python. |
| ADIOS2<br>(Godoy *et al.*, 2020) | ADIOS2 is a framework for data IO management; it is not broadly known or used in seismology. |

proper metadata, (4) data distribution, and (5) data usage (Fig. 2).

## Data submission

In the cases of large-$N$ or DAS experiments, the volume of data submitted to a data center could be hundreds of terabytes. Here, it is important to make a distinction between temporary and permanent experiments. For temporary experiments, data are usually acquired in the field and either transmitted or physically transported to the data center. Field data storage devices can be mounted on the data center's network and the data ingested into the system. In the case of permanent deployments, data are usually transmitted to the data center in real time by means of protocols like SeedLink. However, for high-volume data, such as from DAS experiments, real-time transmission of the full-resolution data may not be practical. One option for DAS installations is to preprocess the data at the station, to reduce the data volume before transmitting them to a data center. Raw data would still be available for later delivery to a data center or, should no data center be capable of handling the volume, a suitable archive found by the operator of the DAS installation.

## Data archive planning

Designing an efficient data management system for seismological data requires consideration of several aspects: the storage hardware, the data format(s), and the way data will be accessed. These three factors are connected; storage system performance varies depending on user access requirements, which are supported by the format of the data. Requests that require gathering small snippets of data from many files require higher performance storage systems than request for data in a single file.

The internal data management of seismological data centers is not specified by FDSN; data centers are free to manage their data as they see fit and store the data they have in their repositories and archives in whatever format suits their use cases. The data provisioning system, fdsnws-dataselect (see Data and Resources), must deliver data in the miniSEED format, but the data centers are free to use different data formats to archive data and perform the conversion on-the-fly before sending them to users.

To date, the miniSEED format has been very popular as a storage format among data centers, because it is highly compatible with the software ecosystem in the community. The main advantages are that there is no need to convert on-the-fly when sending most data to users, and that datasets can be formed by a concatenation of miniSEED records that can be streamed as they are retrieved from the storage media. The ability to stream data directly avoids the need to stage the dataset on the server side before sending it. The disadvantage is that miniSEED has no inherent indexing capability; users typically use the file system as a means to organize and locate data, which may require a huge number of files for medium and large datasets.

For medium and large datasets, the HDF5 format is popular among users of HPC. Two popular HDF5-based formats in our community are PH5 (Hess *et al.*, 2018) and ASDF (Krischer *et al.*, 2016). Initial experimentation has indicated that a reduction in storage space is possible using PH5, likely due to the longer data segments being compressed. A redesign of the PH5 format to make it more versatile is expected to be completed in 2021. The new format is expected to be usable not only for nodal data, but also for DAS, magnetotelluric, geodetic, and other data types. A generic HDF5-based solution could have some advantages not only for data centers but also for users needing to process their large datasets in an HPC-friendly format. The main disadvantage for data centers is that HDF5 formatted data cannot currently be streamed on the fly as they are read by the data provisioning systems (see the Data Distribution section). However, there is a rich ecosystem of
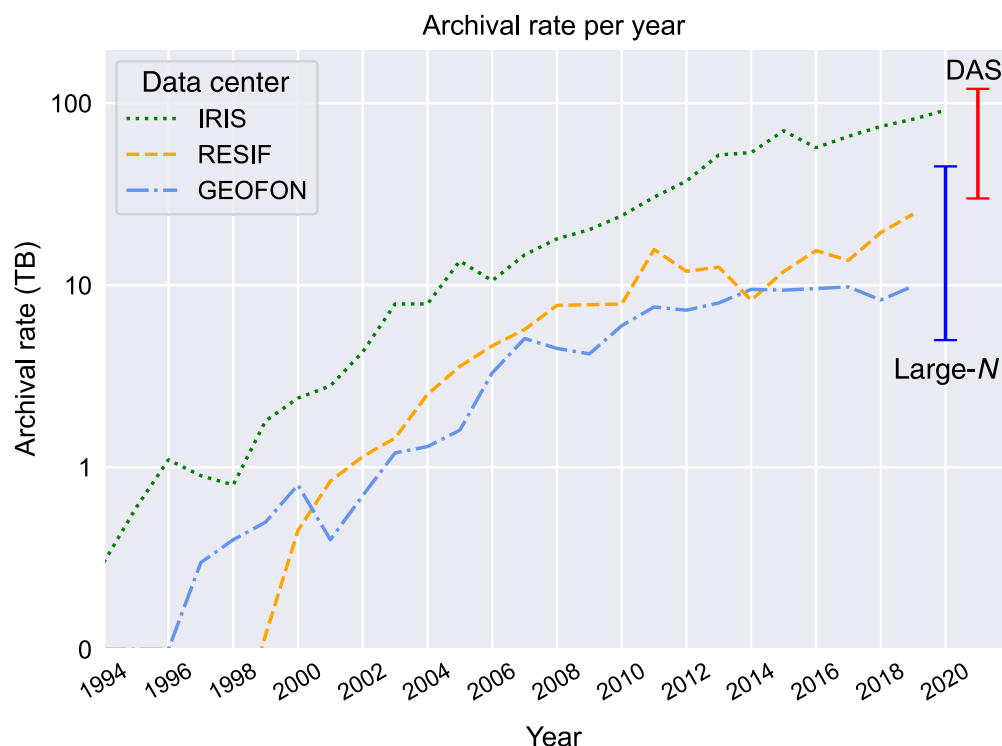
**Figure 3.** Evolution of the data archived yearly at each data center (green, yellow, and blue curves). The error bars on the right show the expected size range of a single dataset for a large-*N* (blue) and a DAS (red) experiment. It can be seen that a single large dataset could be equivalent in size to a whole year of the typical datasets currently being archived. IRIS, Incorporated Research Institutions for Seismology. The color version of this figure is available only in the electronic edition.

the renewal of the hardware (e.g., hard drives) every 5–7 yr, and the regular procedures to ensure that data are still readable.

One way to optimize data storage is to utilize tiered storage strategies; keep data that are rarely accessed in slower, cheaper systems, and data that are requested frequently on systems that provide rapid access. Managing how data are tiered could potentially be dynamic, based on past access patterns. However, these patterns are not easily discerned.

Some data centers have already tested commercial cloud services as a storage system, to evaluate the feasibility of migrating their archives to the cloud. Despite the technical advantages of this approach, there are financial drawbacks that deserve serious consideration, particularly, for large data centers that hold petabytes of data in their repositories, that require substantial compute resources to manage their data, and that deliver petabytes of data each year to their users.

An important consideration regarding choice of storage system is whether the data will be accessed directly for processing, in addition to supporting general extraction systems. Direct data access adds a new dimension to the usage patterns and presents a challenge for managing access permissions, request rates, response times, etc. Direct data access for processing should not be allowed to deteriorate services associated with standard data requests.

interesting open-source tools supporting HDF5, such as the THREDDS data server (Unidata, 2020) or Highly Scalable Data Service (HSDS) distribution services (see Data and Resources).

Other file formats were mentioned in the user's survey, but data center experience with these formats is limited. The desirable features of formats like Zarr, which allow dynamic data chunking, have not been thoroughly tested in the community. In Table 1, we summarize some of the file formats and their characteristics. The different data format features must be explored to optimize the trade-offs between efficient data compression, fast data extraction, and streaming.

Increases of volume also imply new strategies for choice of physical storage. If data centers are expected to increase their archive sizes by at least one order of magnitude (or even more), it might not be feasible to keep all of the data online (Fig. 3). Offline storage should probably be considered inside the data center or provided by external services. This will raise the storage capacity, but data would have to be staged on demand, and with limited size and retention time. Considering the costs, one should consider that the long-term archival of data implies not only the cost of the storage units, but also the backup strategy,

### Data description
There exist comprehensive metadata standards in seismology that are well suited for large nodal data sets or other classical seismic data sets. However, newer, emerging data types, such as, DAS, are not as well handled by these standards in which the concepts of a network, station, location, and a channel are different, and the fundamental instrument responses are new. Therefore, one of the most difficult tasks for data centers attempting to identify and implement standard services is to standardize DAS metadata through the definition of new metadata schemas. StationXML, the international standard for seismic metadata, is not a good fit with the technical features of
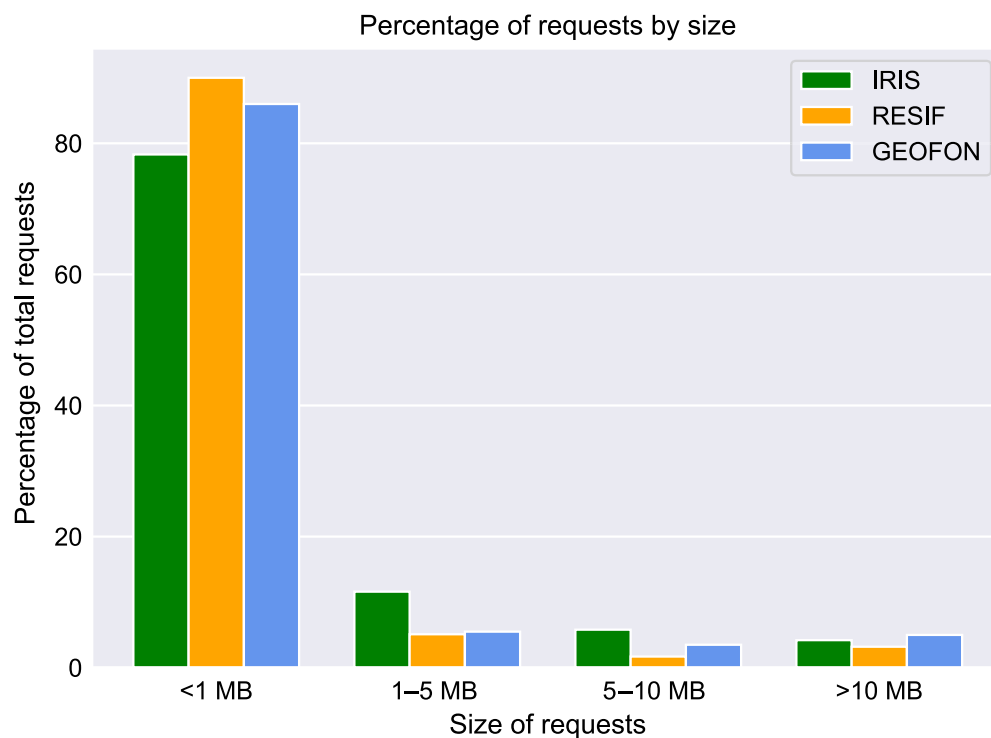
## Percentage of requests by size



**Figure 4.** Percentage of requests per request size and data center. Statistics have been calculated from April 2020 to September 2020. 90% of the requests are less than 5 MB. However, some users request intermediate size volumes through thousands of small data requests. The dataset sizes needed by users could be much larger than what is shown here. The color version of this figure is available only in the electronic edition.

deployment to be a single-station–multiple-channels configuration. We note that the International Seismological Centre (ISC) applies a 1 km station rule for the International Registry of Stations that is generally followed by the community. Quoting ISC: "Because of the need for accurate station positions for hypocenter location programs, a new international code is assigned, if a station is moved more than one (1) kilometer from the previous location (see Data and Resources)". Although, having locations codes more than a kilometer away from the station does not technically violate the rule about stations that move, we believe that it would be extremely confusing to have channels with a location up to 50 km away from their associated station. We also note that some large array deployments (like the Norwegian Seismic Array Network, see Data and Resources) identify each array element with different station names.

DAS interrogators and the generated datasets. There is a need to tackle this in a smart way, to make these data available and interoperable with all popular tools used within the community, which expect typical seismic waveforms with some minimum set of requirements on their metadata, and the possibility of mixing different data types.

A nontechnical (and still pending) issue is how to name DAS data sampling points. If the first approach to standardize and distribute DAS experiment data is to generate derived products, as suggested by many users, there are some details to be considered. Should DAS experiments be registered by the FDSN to get a formal network code? If this is not the case, how could a user identify or select the required sampling points? To have a network code registered would allow the dataset to be uniquely distinguished from other experiments. This would also allow one to formally assign a DOI and the proper citation to the experiment, as agreed and expected in the community.

At levels finer than the network, should each sampling point be considered a station, a location, or a channel? If we consider that each fiber could have thousands of sampling points, with a spatial separation as small as of 2–4 m and a length of tens of kilometers, a station designation seems to be the best fit. As stations have a unique position, we cannot consider such a

A reasonable approach would be to define a station for each sampling point, similar to the situation in nodal experiments in which each node is typically assigned a station code. These naming conventions will be fundamental for interoperability between data centers, as well as a homogeneous definition of the subsets independent of who is archiving the data.

As described earlier, simple information such as latitude, longitude, and orientation are already a challenge, whose solution is not trivial and could potentially make a difference in the results, if data should be automatically processed. It seems unavoidable to perform some postprocessing of the data, to calculate the exact location of each sampling point. This should be based on the Global Positioning System information from the interrogator and the detection of an active source (e.g., "tap test") during the deployment phase. Standard and generic codes, benchmarks, and guidelines provided to/by the community would be desirable, and information about the postprocessing needs to integrate the metadata.

Some of these topics are being discussed within community user groups, such as the DAS research coordination network (DAS RCN, see Data and Resources), and hopefully some proposals to standardize the description of DAS experiments will

emerge. The needs and related challenges from the irruption of new technologies, such as DAS interrogators, could be the triggering factor, to discuss the long-term evolution of our current standards (e.g., StationXML) and to potentially explore more generic formats like SensorML, which is widely used in Geophysics and generic enough to encompass all that StationXML can express.

## Data distribution

Currently, the FDSN-standard suite of web services (fdsnws-station and fdsnws-dataselect) are used by many researchers to discover and download data. From April to September 2020, 90% of the fdsnws-dataselect requests received delivered less than 5 MB of data (Fig. 4). The miniSEED format along with the fdsnws-dataselect webservice is well designed for this purpose. However, users may make many small requests, to optimize their data retrieval success and assemble the complete dataset that they desire on their computers. The size of the datasets needed by the community is probably larger than is observed in Figure 4. It is still unclear what the download pattern for large-$N$ or DAS data would be, but downloads of entire datasets are expected.

The distribution of large data sets (50+ TB) is challenging due to bandwidth constraints. As mentioned previously, one of the main results of the survey we conducted was that the full resolution data are not needed for most studies, particularly for DAS. Instead, a set of products derived from the original data would satisfy most users. The derived data products are often a significantly reduced volume, making them more feasible for archiving and distribution by existing data centers and systems. Although, the derivative products are useful in reducing the burden of transferring bulky data, the need to transfer large volumes is not eliminated for all users.

Another significant challenge is the lack of available services to handle the transmission of large volumes. The current data transfer standard in global broadband seismology (fdsnws-dataselect), defines a synchronous web service interface that accepts arbitrary data selections and is expected to begin returning data with little or no delay. It has become a highly successful mechanism, but although, allowed, requesting large data volumes with this service has the risk of causing connection timeouts. It is difficult to efficiently resume transfers following broken connections, which are increasingly likely with very large volumes. Also, maintaining acceptable performance can require significant software and system engineering by the data centers.

One tool, IRIS's ROVER (see Data and Resources), has overcome the issue of dropped connections using a smart client to manage the synchronous data transfer in miniSEED and constructing the final dataset on the client side. This approach requires a relatively complex client, but is a robust and efficient method of requesting large datasets.

Future distribution capabilities worth exploration are the use of GridFTP (Allcock *et al.*, 2005), or rsync (see Data

and Resources) for efficient bulk file transfer. But these methods can only be used with existing files and not arbitrary data selections. Another possibility is the extension of web service interfaces, such as fdsnws-dataselect, to allow asynchronous, or batch, data transfer to a site (e.g., HPC center) identified by the requestor. In other cases and depending on the archive format, some middleware solutions, such as THREDDS data server (Unidata, 2020) or HSDS (see Data and Resources), could potentially be useful solutions for subselecting and transforming large data volumes.

## Using the data

Considering the classical approach of data usage by seismologists, some important requirements have been identified in the previous sections, in particular, the need for standard data and metadata formats, and correct and complete metadata. Currently, only miniSEED is an FDSN standard within the community, but, this format is inefficient and cumbersome for processing large datasets. Non-FDSN standard formats (e.g., HDF5 based) are presently being used by researchers working with medium and large datasets (20+TB), and they would like to get data in this format from the data center provisioning systems. The definition of an additional FDSN standard aimed at large datasets would spur the development of community processing tools, and would greatly further the cooperation and sharing of resources between data centers.

The user's need to process data in a remote environment (e.g., cloud computing, containers) means that other factors should also be considered by data centers. For instance, even if a data center can provide a container or Jupyter notebook (see Data and Resources) with direct processing of the data, solutions need to be developed to limit restricted dataset access to authorized users. This is not a minor issue, because most of these large datasets go under an embargo period of some years before being opened to the public. This means that, depending on the user, only part of the archive should be accessible to any given user from the computing interfaces. A similar issue is related to the data organization (format and structure). The data center can expose (or let the user access) what exists in the storage system, so the user will, in principle, be able to use only the storage format. It will be difficult (and in some cases impossible) to stage a new dataset in another format to be accessed by the user. Therefore, the data format and organization needs to be known by the user. Because each data center is free to organize the data in whatever way is optimal to them, data processing software cannot easily be used across different data center computational facilities.

The transparent federation of data is probably the biggest challenge for data centers, because it includes not only all the improvements and developments mentioned in this article, but also the edge case in which multiple large datasets are requested from different data centers. In such a case, even running the code on top of the data is not a solution. That would only save the

transfer from that data center, but there is still the problem of requesting all the other data. Some public initiatives exist in the United States and Europe, to propose solutions to the problem of federated data (e.g., Jetstream cloud environment from the Extreme Science and Engineering Discovering Environment, and the European Open Science Cloud). However, there are still no solutions available that cover all the aspects mentioned earlier.

## Conclusions

Despite present initiatives trying to tackle the problem of big seismological datasets and how they can be archived, delivered, and processed, data centers are facing a situation similar to that from the mid-80s. Namely, the amount of data generated was much bigger than the data centers could safely and sustainably archive. At that moment, the temporary solution was to keep only time windows of interest (e.g., waveforms related to an event) and discard the rest of the data after some time. What has been learned from this experience is that discarding data is not the best option to consider, as illustrated by present day research in data mining, new types of seismic signals, and imaging techniques using seismic noise.

With the present increase of data volumes come additional problems associated with the internal organization of data within the data centers and meeting the need for data formats that are well adapted for storing and slicing big datasets. Along with these problems, there is a need for adapting present data request tools or developing new ones.

At present, the IRIS, GEOFON, and RESIF data centers recognize that it is not possible for them to permanently archive very large datasets, such as DAS, due to prohibitive costs of long-term archiving. Safe storage (involving multiple copies) of such large datasets available online is presently not possible either. Therefore, the seismological data centers should explore a strategy for delivering datasets through automatic asynchronous services to designated destinations, such as computing facilities for which data preprocessing or processing could occur. For very large datasets, a common strategy is needed for storing only standard data products, such as spatially and/or temporally downsampled data, and/or small windows of highly sampled data.

Ideally, in the long term, international standards and products should be developed within the framework of the FDSN. As a preparatory phase, we suggest to undertake the following actions:

- Encourage a wide and continuous international collaboration on DAS data. This would include dedicated workshops and special sessions in conferences involving data centers and scientific users who could define a limited set of standard products for DAS data. One natural forum would be a strong international presence in the DAS RCN, which involves both scientific users and data center operators and managers.
- Extend the present discussion between our three data centers to other interested data centers. One option is to set up dedicated technical workshops on very focused subjects,

including solutions to keep the raw data safely for future use. Our community should in these workshops interact with data centers from other disciplines to understand similarities and differences of use cases and technical constraints, to benefit from past experience over a broad range of scientific areas, and to explore possible common technical solutions.

This preparatory work hopefully will lead to new community standards and user services that are well adapted to large datasets.

## Data and Resources

The Standard for Exchange of Earthquake Data (SEED) manual can be downloaded from http://www.fdsn.org/pdf/SEEDManual_V2.4.pdf (last accessed November 2020). Information about Highly Scalable Data Service (HSDS) can be read from https://www.hdfgroup.org/solutions/highly-scalable-data-service-hsds/ (last accessed November 2020). Documentation of Zarr can be found at https://zarr.readthedocs.io/ (last accessed November 2020). Information about the distributed acoustic sensing (DAS) research coordination network (DAS RCN) can be found at https://www.iris.edu/hq/initiatives/das_rcn (last accessed November 2020). Specifications of the International Federation of Digital Seismograph Networks (FDSN) web services can be found at http://fdsn.org/webservices/FDSN-WS-Specifications-1.2.pdf (last accessed November 2020). Information on the Norwegian Seismic Array Network (NO) can be found at http://www.fdsn.org/networks/detail/NO/ (last accessed November 2020). Details of the International Registry of Seismographic Stations can be read from http://www.isc.ac.uk/registries/registration/ (last accessed November 2020). Documentation of ROVER can be found at https://iris-edu.github.io/rover/ (last accessed November 2020). Rsync information is available at https://rsync.samba.org/ (last accessed November 2020). Information about Jupyter notebooks is available at https://jupyter.org/ (last accessed November 2020). Plots were made using Matplotlib (Hunter, 2007). The Portable Array Seismic Studies of the Continental Lithosphere facility (PASSCAL) is a facility of the Incorporated Research Institutions for Seismology (IRIS) consortium that provides instrumentation for National Science Foundation, Department of Energy, and otherwise funded seismological experiments around the world. The Matlab available at https://www.mathworks.com/products/matlab.html (last accessed January 2021).

## Acknowledgments

seismological facility for the advancement of geoscience (SAGE) is operated by Incorporated Research Institutions for Seismology (IRIS) on behalf of the National Science Foundation (NSF) through Award EAR-1851048. The Réseau sismologique et géodésique français (RESIF) data center is operated on behalf of the RESIF consortium, a national research infrastructure coordinated by Centre National de la Recherche Scientifique (CNRS) and which receives funding, personnel and other from all the consortium members. Résif-SI is also supported by the System of Observation and Experimentation for Research and the Environment (SOERE) and by the Ministry of Ecological Transition. Résif-SI has benefited from the Résif-CORE project (11-EQPX-0040), funded by the national French funding program Investissements d'Avenir and managed by the French National Research Agency (ANR).

# References

Allcock, W., J. Bresnahan, R. Kettimuthu, and M. Link (2005). The globus striped GridFTP framework and server, *SC'05: Proc. of the 2005 ACM/IEEE Conference on Supercomputing*, Seattle, Washington, 12–18 November 2005, doi: 10.1109/SC.2005.72.

Beyreuther, M., R. Barsch, L. Krischer, T. Megies, Y. Behr, and J. Wassermann (2010). ObsPy: A Python toolbox for seismology, *Seismol. Res. Lett.* **81,** no. 3, 530–533, doi: 10.1785/gssrl.81.3.530.

Casey, R., M. E. Templeton, G. Sharer, L. Keyson, B. R. Weertman, and T. Ahern (2018). Assuring the quality of IRIS data with MUSTANG, *Seismol. Res. Lett.* **89,** no. 2A, 630–639, doi: 10.1785/0220170191.

Godoy, W. F., N. Podhorszki, R. Wang, C. Atkins, G. Eisenhauer, J. Gu, P. Davis, J. Choi, K. Germaschewski, K. Huck, *et al.* (2020). ADIOS 2: The adaptable input output system. A framework for high-performance data management, *SoftwareX* **12,** 100561, ISSN 2352-7110, doi: 10.1016/j.softx.2020.100561.

Hess, D., N. Falco, rsdeazevedo, damhuonglan, and K. Jacobs (2018). PIC-IRIS/PH5: v4.1.2 (Version v4.1.2), Zenodo, doi: 10.5281/zenodo.1284569.

Hoyer, S., and J. Hamman (2017). xarray: N-D labeled arrays and datasets in Python, *J. Open Res. Software* **5,** no. 1, 10, doi: 10.5334/jors.148.

Hunter, J. D. (2007). Matplotlib: A 2D graphics environment, *Comput. Sci. Eng.* **9,** no. 3, 90–95.

Jousset, P., T. Reinsch, T. Ryberg, H. Blanck, A. Clarke, R. Aghayev, G. P. Hersir, J. Henninges, M. Weber, and C. M. Krawczyk (2018). Dynamic strain determination using fibre-optic cables allows imaging of seismological and structural features, *Nat. Commun.* **9,** 2509, doi: 10.1038/s41467-018-04860-y.

Krischer, L., J. Smith, W. Lei, M. Lefebvre, Y. Ruan, E. Sales de Andrade, N. Podhorszki, E. Bozdağ, and J. Tromp (2016). An adaptable seismic data format, *Geophys. J. Int.* **207,** no. 2, November 2016, 1003–1011, doi: 10.1093/gji/ggw319.

Lin, F.-C., D. Li, R. W. Clayton, and D. Hollis (2013). High-resolution 3D shallow crustal structure in Long Beach, California: Application of ambient noise tomography on a dense seismic array, *Geophysics* **78,** no. 4, Q45–Q56, doi: 10.1190/GEO2012-0453.1.

Standard for the Exchange of Earthquake Data (SEED) (2012). *Standard for the Exchange of Earthquake Data*, SEED Reference Manual, SEED format version 2.4., International Federation of Digital Seismograph Networks, Incorporated Research Institutions for Seismology (IRIS), USGS.

Trani, L., M. Koymans, M. Atkinson, R. Sleeman, and R. Filgueira (2017). WFCatalog: A catalogue for seismological waveform data, *Comput. Geosci.* **106,** no. September 2017, 101–108, doi: 10.1016/j.cageo.2017.06.008.

Unidata (2020). THREDDS data server (TDS) version 4.6.2 [software], UCAR/Unidata, Boulder, Colorado, doi: 10.5065/D6N014KG.

Wilkinson, M., M. Dumontier, I. Aalbersberg, G. Appleton, M. Axton, A. Baak, N. Blomberg, J. W. Boiten, L. B. da Silva Santos, P. E. Bourne, *et al.* (2016). The FAIR guiding principles for scientific data management and stewardship, *Sci. Data* **3,** Article number: 160018, doi: 10.1038/sdata.2016.18.