

# **Data Dissemination and Cyberinfrastructure Ad Hoc Committee Initial Report**

**6 September 2017**

## **1 Introduction**

To facilitate the assessment of Ocean Observatories Initiative (OOI) data quality by the science community, and to accelerate the integration of OOI infrastructure usage into project proposals and scientific publications, the NSF Ocean Observatories Initiative Facility Board (OOIFB) established the Data Dissemination and Cyberinfrastructure (DDCI) ad hoc committee which is tasked with identifying near-term obstacles to the enhanced delivery of data to the scientific community and providing recommendations for removing these obstacles.

The DDCI comprises the following individuals:

- Timothy Crone, LDEO (co-Chair)
- James O'Donnell, UConn (co-Chair)
- Brian Glazer, UH
- Orest Kawka, UW
- Stephanie Petillo, WHOI
- Mary Jo Richardson, TAMU
- Richard Signell, USGS
- Derrick Snowden, NOAA IOOS
- Larry Atkinson, OOIFB Chair (ex officio)

The DDCI met in-person at NSF headquarters on 18 July 2017, and has had several Webex/conference calls prior to and following this meeting. NSF program manager Lisa Clough was present for the meeting and all calls. Other representatives from NSF including Rick Murray, Bob Houtman, and Rachel Shackleford were present for part or all of the in-person meeting. Annette DeSilva (UNOLS/URI) facilitated meetings and calls and took meeting notes.

During the meeting and the calls, the committee spent a significant amount of time learning about the current state of the OOI cyberinfrastructure and how data is currently handled and disseminated. We heard presentations from Mike Vardaro of the Rutgers Data Team and Ivan Roderio of the Rutgers CI Team. We heard a presentation on ERDDAP by Rich Signell, and a presentation on the OOI high-definition camera system (CAMHD) by Tim Crone. We had lively and informative discussions about the needs of the scientific community, potential new modes of data access for the science user, and discussions regarding potential improvements to the management structure of the CI.

This report is a summary of the committee's findings representing our views at this stage of our efforts. The committee expects to continue our work on this problem, to meet again in the future, and to refine our views and recommendations as we learn more about the current state of the system and receive input from operators and the scientific community.

The findings and recommendations in this report are broken down into two sections. In the first section, we detail our recommendations for short-term adjustments to the current cyberinfrastructure priorities and management structure that we believe can be reasonably accomplished in the next few months. In the second section, we detail our recommendations for "OOI 2.0" which include longer-term recommendations that should be considered as the next phase of OOI operations is planned and a new Cooperative Agreement (CA) for the management and operation of OOI is formed and executed.

## 2 Near-Term Recommendations

The committee has several near-term recommendations to facilitate data dissemination in the coming months, which the committee thinks can be reasonably accomplished before the transition to OOI 2.0. These recommendations are:

1. Prioritize the development and public release of the uFrame-powered ERDDAP server.
2. Accelerate the ingestion of backlogged data.
3. Identify a single individual who will be responsible for improving data access for the scientific end user and who has the authority to define both CI and Data Team priorities.

### 2.1 ERDDAP

ERDDAP is a free and open-source Java “servlet” which for the non-expert can be thought of as a kind of specialized web server that excels in serving and converting disparate scientific datasets using a uniform interface. ERDDAP is focused primarily on serving tabular or time-series datasets which are stored on the server as static NetCDF files, and it can serve raw (processed) data in a large number of formats as well as generate plots and maps of requested data. ERDDAP has a standard browser interface that facilitates searching for, converting, and plotting data, but ERDDAP is built on a RESTful API, meaning that the server does not store browser state and all information about every request is contained in the URL of each request. This makes it easy to automate searching for and using data in other applications like Python or MATLAB, and makes it easy for users to build their own custom interfaces if they so wish.

ERDDAP is a server framework that allows anyone with data to serve to serve their data by running their own ERDDAP server. Many dozens of organizations are now running ERDDAP servers to serve their scientific data, and ERDDAP is on its way to becoming a de facto standard in the Oceanographic community. Many oceanographers are already familiar with the ERDDAP interface and have already developed their own tools to work with data served by such systems.

The committee believes that ERDDAP has the potential to serve most of the OOI data in an efficient and useful manner and that the deployment of an ERDDAP system that works on top of uFrame could greatly expand OOI data availability for the scientific community. The committee recommends that the development of the ERDDAP system be made a top priority.

To expedite the development of the ERDDAP system, the committee recommends that the ERDDAP development team be provided with the access they need to complete this task as quickly and as efficiently as possible. At a minimum, the ERDDAP team should be given read access to the production Tomcat logs. Another suggestion for speeding up the development of ERDDAP is to reduce the deployment timeline from two weeks to a few days, specifically in support of the ERDDAP team to accelerate the development of this system.

The committee also recommends that the ERDDAP developers begin or continue to interact with other OOI developers such as the CGSN who have developed internal ERDDAP systems, and that they ensure that the ERDDAP data sets are well-described using best practices for international standards. For example, it would be best if the OOI CI way of publishing moored buoy data via ERDDAP were similar to or even identical to the OceanSITES way of publishing these data. OceanSITES is a 15-year-old program that publishes moored buoy data, has international buy in, and is largely viewed as the best practice for formatting this type of data.

### 2.1 Data ingestion

Ingestion backlogs are an area of concern in terms of data availability for the scientific community. The committee recommends that data ingestion remain a top priority for the CI and Data Teams. The

committee notes that although the M2M ingestion system appears to be promising, the MIOs are not currently using it, and may in fact not be authorized to use it. Also, it is not clear how a distributed ingestion model can work. The committee recommends that the CI and Data teams continue to focus on data ingestion using a centralized model with the understanding that the MIOs may not be involved in the near-term.

### *2.3 Data Delivery Manager*

It is the committee's view that the separated organizational structure of the CI and Data team has led to roadblocks in terms of the effective and efficient dissemination of data to the scientific community, and in terms of the ability for the scientific community to provide input on decisions made by the administrators of the system. The committee believes that the two teams would benefit from the establishment of an OOI Data Delivery Manager with authority over the priorities of the OOI CI and Data Teams. The primary goal of the Data Delivery Manager should be to deliver data to the scientific users in a way that works for those users, and to be responsive to the users' needs and input, with the scientific users defined as the *customers* inside the OOI "business model". With oversight over the CI and Data Teams, the Data Delivery Manager can help reorient the focus to delivery of data to the scientific user in the most efficient and effective manner possible.

## **3 Long-Term Recommendations (OOI 2.0)**

The committee has several longer-term recommendations to facilitate data dissemination as OOI 1.0 transitions to OOI 2.0. These recommendations are:

1. Assess the future viability of uFrame.
2. Place a primary focus on the scientific user base.
3. Consider partnerships for providing remote compute capability for larger OOI datasets.
4. Maintain a Data Delivery Manager in OOI 2.0.

### *3.1 Assess the future viability of uFrame*

The committee and nearly everyone consulted by members of the committee have serious concerns about the uFrame system. One primary concern is that uFrame in effect places a "black box" or at best a "gray box" in the processing pipeline, and it is difficult for end users to fully ascertain how data products are generated from raw data. It is difficult if not impossible for users to run custom processors using different calibrations or other processing parameters to experiment and troubleshoot. This lack of obvious transparency has caused many members of the scientific community to express healthy levels of skepticism regarding the data pipeline.

Another primary concern is the apparent lack of documentation for uFrame and the proprietary nature of many components of the uFrame codebase. In addition to the obvious transparency issues when dealing with closed-source code, the proprietary aspect of the software may become a budget issue in the future. If there are no funds for planned product improvements or if Raytheon is unwilling to make uFrame open source, then OOI could be locked in with a Raytheon product for the foreseeable future. Even if Raytheon does release the source code, there is no guarantee that the current CI team (or the new team if that changes) will have the skills to maintain and modify what would become a fork of the Raytheon product into the public domain.

Despite these concerns, the committee notes that the uFrame/Cassandra database model offers some advantages that *may* not be easily replicated using a simpler file based system. The first among these is that uFrame stores instrument raw data in the database, and applies processors to the data upon data request. This model would theoretically allow users to apply custom processor files to the raw data to generate alternative data products during queries, however it is not clear if this capability has been

realized. Currently changing processors appears to be a long and complex process which regular users do not have easy access to. Another advantage is that the current system is capable of ingesting, processing, and serving data from the Cabled Array in real-time, which provides substantial scientific value.

For OOI 2.0, the committee recommends that uFrame be evaluated in terms of the issues listed above, and that potential alternatives be considered. Any replacement systems considered should not descope the capabilities of the CI, and specifically should maintain and preferably extend the “compute on query” aspect of the system, and should maintain the real-time ingestion/processing/service capability of Cabled Array data.

### *3.2 Place a primary focus on the scientific user base*

The OOI has enormous potential for outreach and education, for use by the general public, the media, and by students of all ages. However, the viability of the observing system during these early years of operation will be dependent on proposal pressure from scientists in the community to use and expand OOI assets, and on the publication of peer-reviewed journal articles based on OOI data. Indeed, a primary goal of this committee is to accelerate the availability of OOI data for scientists and thus expand the use of these data for science. For this reason, the committee recommends that efforts to improve the user experience on the OOI data portal, and the expanded availability of data through systems such as ERDDAP or the M2M interface be focused on the needs of the working scientists. Based on our discussions, our view is that scientists have needs and requirements that are quite different than the casual user, and can be summarized by this list of questions a scientist is likely to ask when looking to obtain data:

1. What data is available? What instruments are working and which ones are not? Scientists need an easy to see overview of the entire system, which instruments are working, which ones are not, and why not, and which ones will be working in the future.
2. How good are the data? Are the metadata flags easy to understand and are they well incorporated into the data provided? Scientists need to know how reliable the data they obtain is.
3. Where are the data? Are the data easy to download in easy to use formats? Can the data be downloaded by clicking a link instead of by waiting for an e-mail to arrive? Scientists may want to see plots of data in real time, but in most cases scientists will want to download data in some sort of table format that allows them to do their own processing and visualization using the software tools of their choice.

Efforts to improve the UI of the OOI should be focused on how working scientists actually use data.

### *3.3 Consider partnerships for providing remote compute capability for larger OOI datasets*

For some OOI data the download model for data access is simply not viable. In particular, the hydrophone and the HD video datasets are so large that researchers cannot hope to download these datasets within any reasonable timeframe. For example, the HD video dataset comprises nearly 7000 high-resolution video files totaling approximately 85 TB in size. Not only might it take many weeks or months to download the entire dataset, but most researchers would struggle to find the space to store such a large amount of data locally.

For this reason, the committee recommends that collaborations and or partnerships be sought to provide combined compute and storage capability for these large datasets. One potential partner is XSEDE, which oversees a consortium of some of the country’s largest supercomputer operators. XSEDE may be able to provide hosting and access to these data using some novel funding model. Other possibilities include the development of commercial partners such as Amazon or Google who may be willing to host these large datasets at affordable rates, or Calit2 which has expressed interest in hosting OOI data. The possibilities

for partnerships abound, and for some of the data in the OOI system, a cloud-based solution is the best way to accelerate data access for the scientific community.

### *3.4 Maintain a Data Delivery Manager in OOI 2.0*

A person or small group should be retained in a management position at the level of the MIOs that has the primary goal of overseeing data delivery to the scientific community, is responsive to the needs of the scientific community, and has oversight authority over all management components of the cyberinfrastructure system and administration so that decisions about the cyberinfrastructure can be made with the needs of the scientific community at the forefront.

## **4 Summary**

In summary, the recommendations of the DDCI Committee over the short term are:

1. Prioritize the release of the OOI ERDDAP server by empowering the ERDDAP development team with all needed access to the production server and by shortening of the deployment cycle timeline.
2. Accelerate the ingestion of backlogged data.
3. Identify a single individual who will serve as OOI Data Delivery Manager and will be responsible for improving data access for the scientific end user and who has the authority to define both CI and Data Team priorities.

To help guide the formation of the CA for OOI 2.0, the recommendations of the DDCI Committee are:

1. Assess the future viability of uFrame.
2. Place a primary focus on the scientific user base for data delivery.
3. Consider partnerships for providing remote compute capability for larger OOI datasets.
4. Maintain the position of OOI Data Delivery Manager in OOI 2.0.