# Webinar Goals



**1** Review & discuss project objectives & proposed plan

**2** Background on other model data repositories and approaches

**3** Discuss & improve the feedback form, and discuss answers
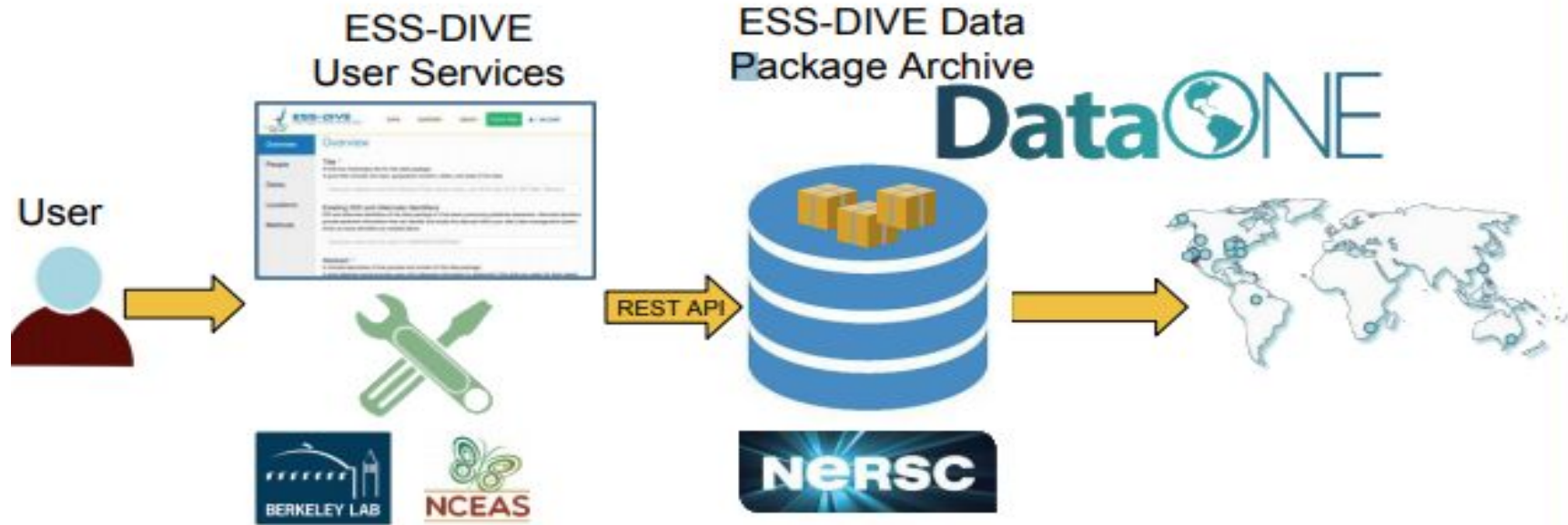
# ESS-DIVE: Current Functionality for Model Data Archiving



CREATE SINGLE OR MULTIPLE DATA PACKAGES

# ESS-DIVE: Current Functionality for Model Data Archiving

First model dataset on
ESS-DIVE

# ESS-DIVE: Current Functionality for Model Data Archiving

Opportunities for improvement!

# ESS-DIVE: Current Functionality for Model Data Archiving

**Problem Statement**

Model data storage is limited
- currently only archiving a limited set of small-sized model outputs
- infrastructure limitations on data size:
  - upload limits 2GB/file on portal and 1GB/file on API
  - architecture limits how much data ESS-DIVE can store and serve
- web interface not the best tool for uploading/downloading large datasets
- API helps but there are still physical limitations

No community consensus yet on what to archive, standards, storage space needed, etc.
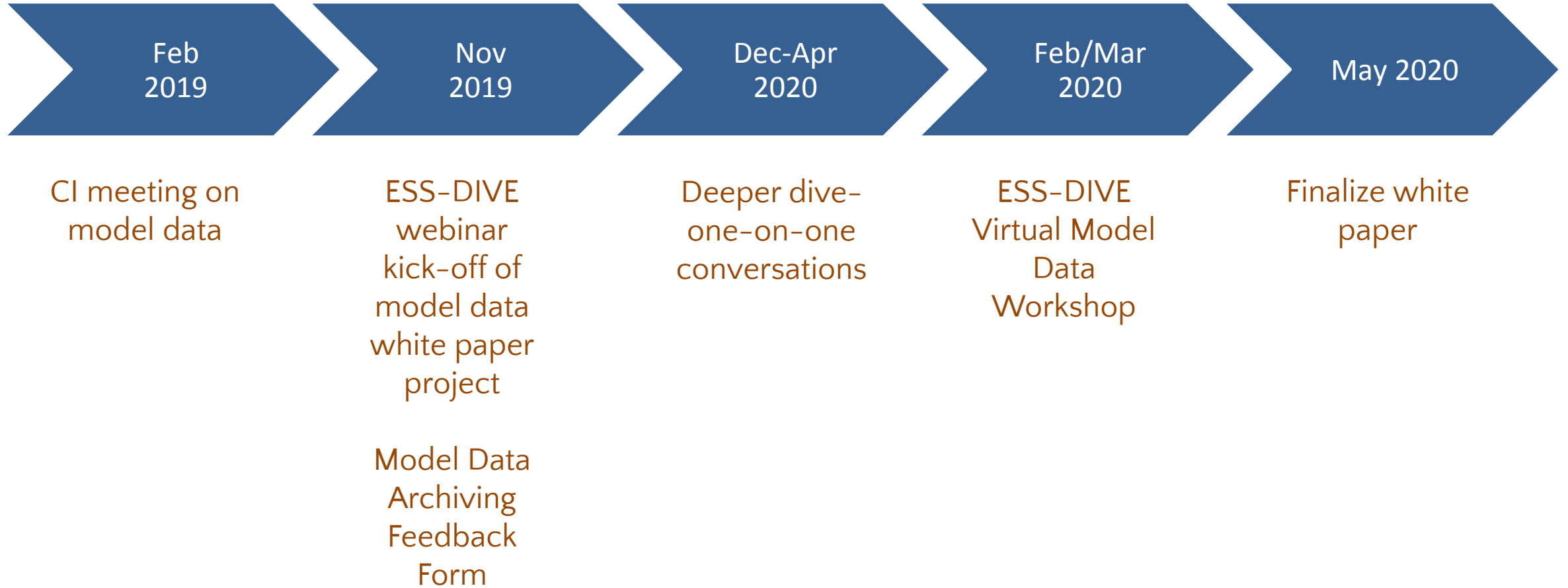
# Main objectives of this project

**Assess:**
- What model data should be archived, purposes of storage, storage capacity needed
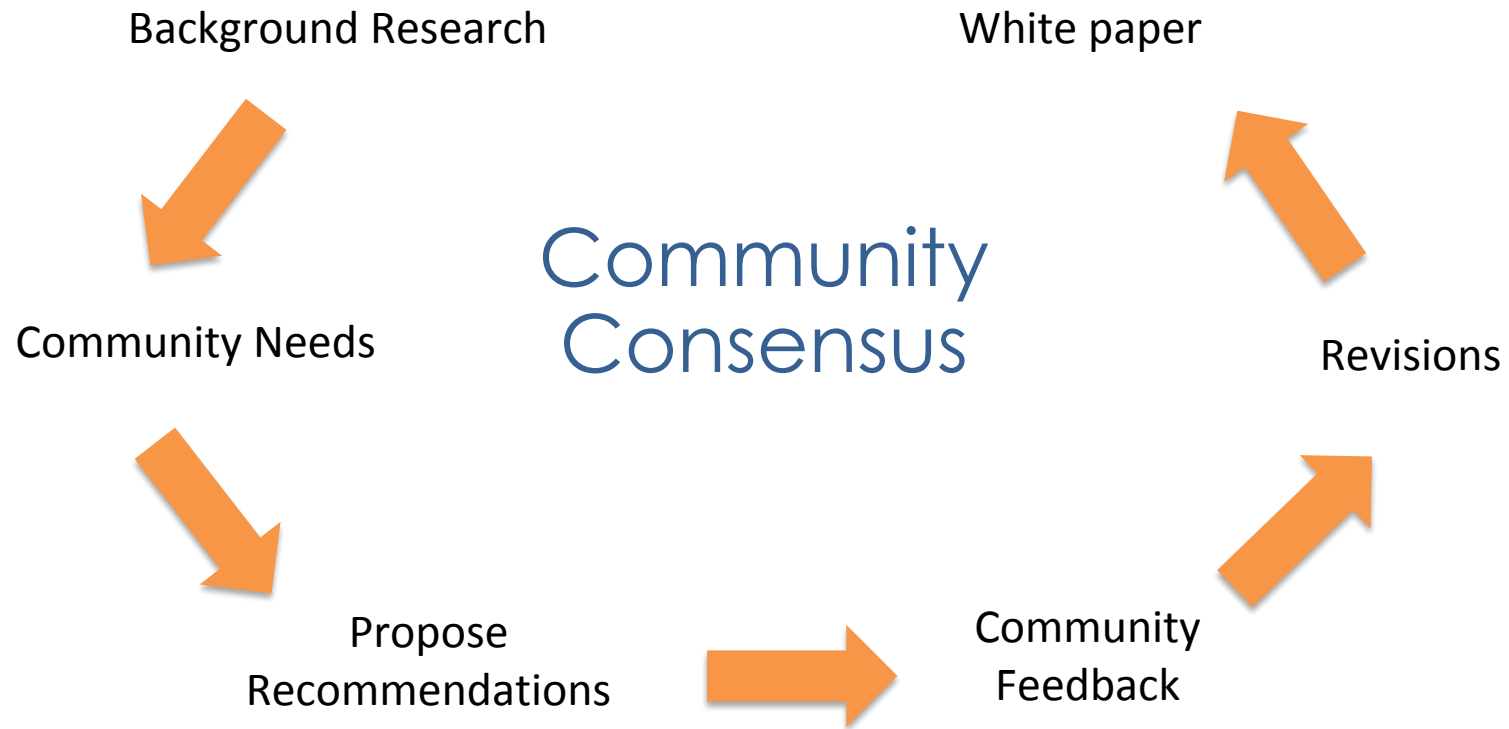- Best approach to store data

**Deliverable:**
- White paper describing data storage capabilities ESS modeling projects need based on community feedback and a few use-cases, and potential storage options

# Timeline

| Feb 2019 | Nov 2019 | Dec-Apr 2020 | Feb/Mar 2020 | May 2020 |
|----------|----------|--------------|--------------|----------|
| CI meeting on model data | ESS-DIVE webinar kick-off of model data white paper project

Model Data Archiving Feedback Form | Deeper dive-one-on-one conversations | ESS-DIVE Virtual Model Data Workshop | Finalize white paper |

# Our Process

ESS-DIVE

Background Research

White paper

Community Needs

## Community Consensus

Revisions

Propose Recommendations
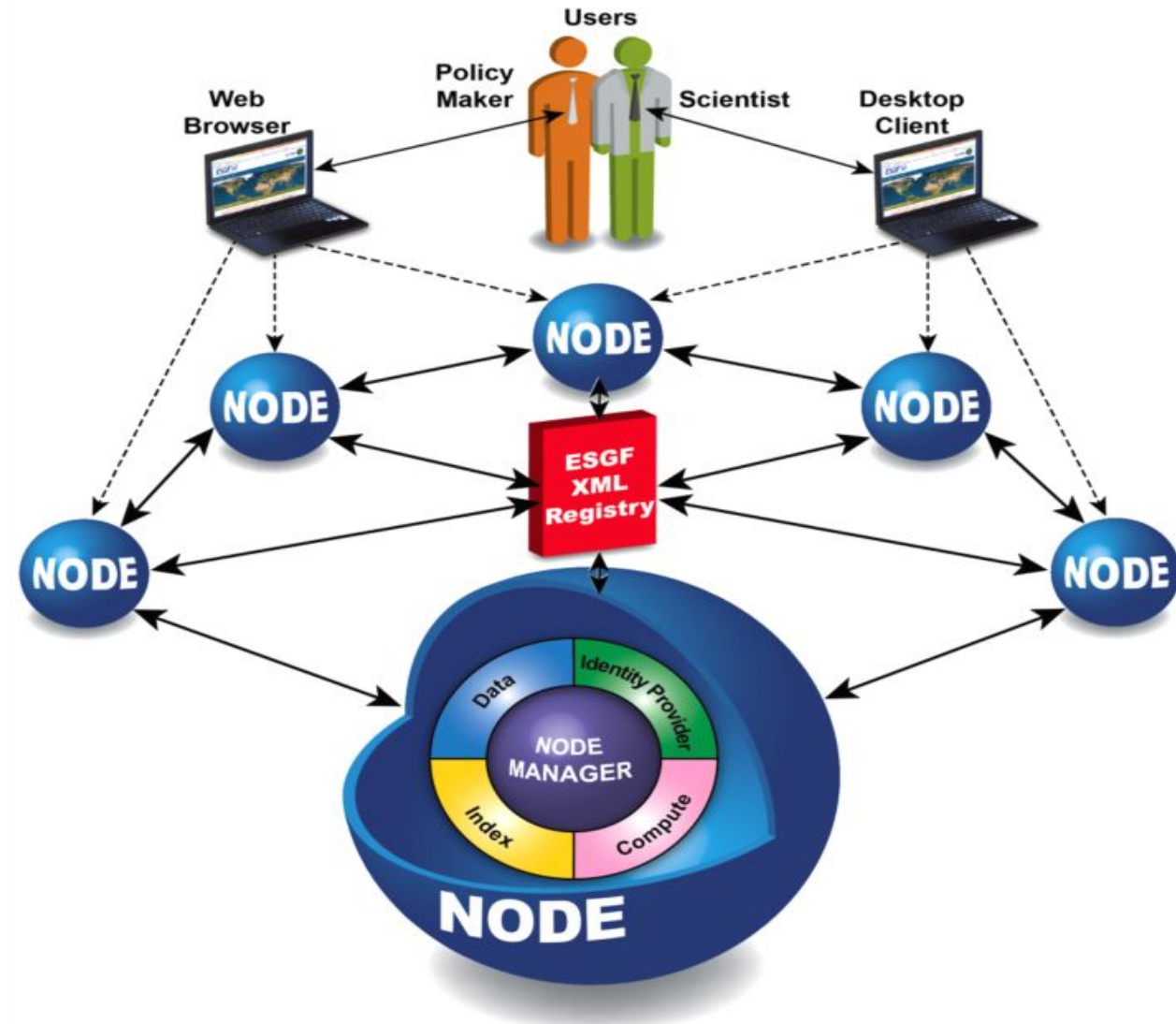
Community Feedback

# Examples of cloud-based web storage for model & observational data

- Earth System Grid Federation (ESGF)

- NASA's Earth Observing System Data and Information System (EODIS)

- NCAR's Earth Observatory Laboratory (EOL) Data Archive

- NCAR's Research Data Archive (RDA)

- CYVERSE

# ESGF

- Federated system for storing and serving data from multiple locations and sources
- Developed for sharing climate model (CMIP) data
- Currently stores some ESS modeling data
- NOT guaranteed to be a long-term host of data, to follow digital library standards, or to assign DOIs to data

# EOSDIS





FIGURE 1. HISTORICAL AND PROJECTED CUMULATIVE ARCHIVE VOLUME IN EOSDIS. (YEARS RUN FROM OCTOBER TO SEPTEMBER.)

# Earth Observatory Lab (EOL) Data Archive

# Research Data Archive (RDA)



**Browse the RDA**

There are 696 public datasets in the CISL RDA. You can begin browsing the datasets by choosing one of the facets in the menu to the left. Facet descriptions are given below, along with the number (in parentheses) of datasets in each.

**Variable / Parameter** (696)
A variable or parameter is the quantity that is measured, derived, or computed – e.g. the data value.

**Type of Data** (696)
This refers to the type of data values – e.g. grid (interpolated or computed gridpoint data), platform observation (in-situ and remotely sensed measurements), etc.

**Time Resolution** (298)
This refers to the distance in time between discrete observation measurements, model product valid times, etc.

**Platform** (661)
The platform is the entity or type of entity that acquired or computed the data (e.g. aircraft, land station, reanalysis model).

**Spatial Resolution** (341)
This refers to the horizontal distance between discrete gridpoints of a model product, reporting stations in a network, measurements of a moving platform, etc.

**Topic / Subtopic** (696)
Topic and subtopic are high-level groupings of parameters – e.g. Atmosphere (topic), Clouds (subtopic of Atmosphere).

**Project / Experiment** (159)
This is the scientific project, field campaign, or experiment that acquired the data.

**Supports Project** (51)
This refers to data that were acquired to support a scientific project or experiment (e.g. GATE) or that can be used as ingest for a project (e.g. WRF).

**Data Format** (695)
This refers to the structure of the bitstream used to encapsulate the data values in a record or file - e.g ASCII, netCDF, etc.

**Location** (108)
This the name of the (usually geographic) location or region for which the data are valid.

# CYVERSE (originally iPlant)



- Data storage geared specifically towards data analysis
- Interactive, web-based analytical platform
- Cloud computing, analysis and storage
- Support services for scaling up computational algorithms & on how to use CI

Any thoughts so far?

# (DRAFT) Feedback form for developing the ESS-DIVE model data repository

- Currently comprised of 20 questions to:
  - inventory models and assess their specific data storage needs
  - evaluate what's worth archiving and for how long
  - get recommendations for archiving protocols and storage options
- Let's review for completeness and start discussing answers!

# ESS-DIVE: Vision

# Next steps

**1** Synthesize our discussion today: compile preliminary poll of responses to feedback form, and email link for additional comments on it from the ESS community.

**2** Revise feedback form and distribute to everyone in ESS community.

**3** Connect for follow-up discussions.

**4** Email any questions or more ideas to Maegen ([mbsimmonds@lbl.gov](mailto:mbsimmonds@lbl.gov)) and Bill ([wjriley@lbl.gov](mailto:wjriley@lbl.gov)).