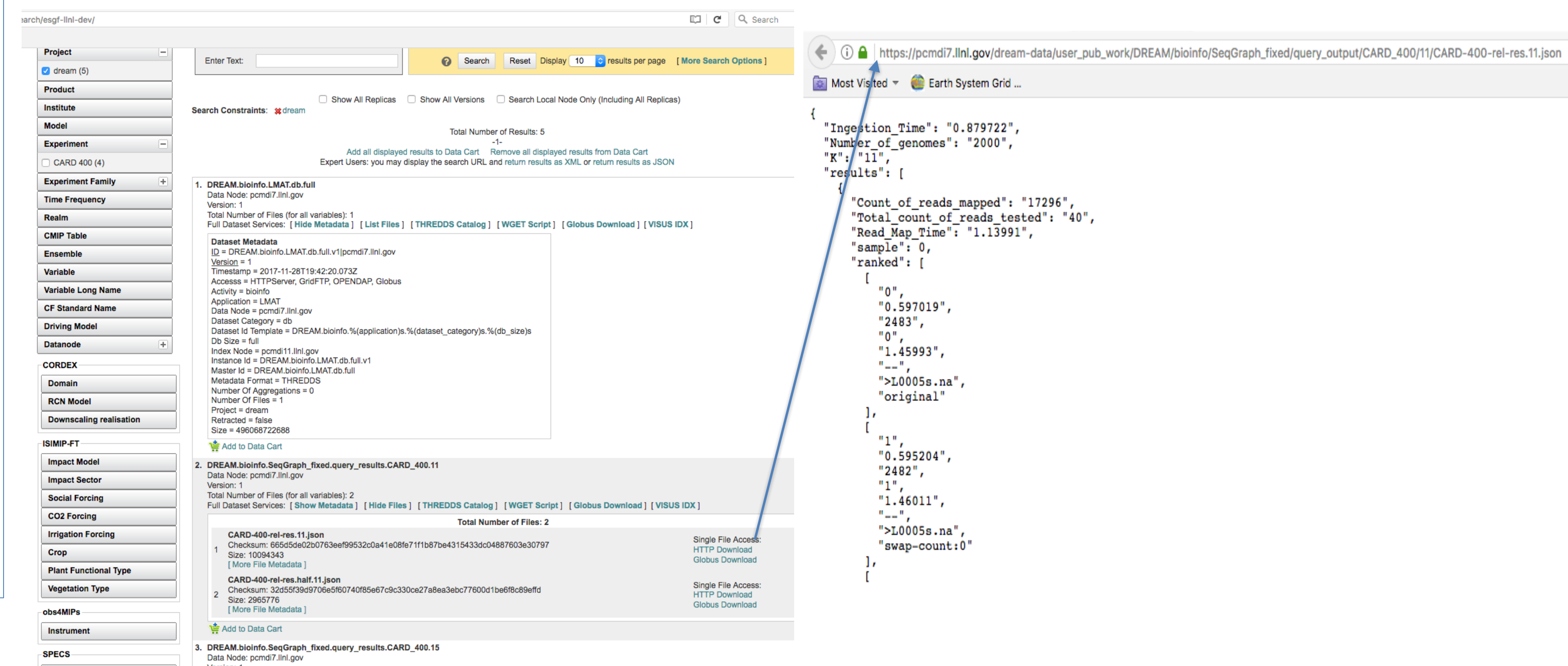# DREAM Data Services for Biological Data and Beyond

## Sasha Ames, Luca Cinquini, and Dean Williams -  DOE LLNL and NASA JPL

- One of the many goals of the DREAM project is to explore enhancements needed for ESGF to enable publishing and access to data in scientific fields other than climate/weather
- THREDDS data server (TDS) has been very effective for serving NetCDF data published to ESGF
- Need a service more specific for alternate data (e.g. ASCII-based) in other domains
  - Example: FASTA format used in bioinformatics to represent genomic and protein sequences
- This service will allow a variety of content types to interoperate properly with a user's web browser



### New tools:

- get_meta.sh
- create_dset.py
    - scripts to generate a generic mapfile with consistent directory_format to dataset_id mapping
- dream-data
  - Service supplies correct mime times for published files
  - eg. images, .pdf, .json
  - Flask module
  - https://github.com/ESGF/esgf-dream-data-service

### Steps taken for publishing to service:

- Modify esg.ini – use /dream-data/ for HTTPServer
- esg.dream.ini – use multiple_handler to accept non-netcdf data
- Generate mapfile with get_meta.sh + create_dset.py
- Reconfigure DRS when sub-project config changes, repeat

esg.dream.ini`



Benefit to input4MIPs: Example MPI-M "dataset" has .pdf, README, and code files

### Future work:

- Add random access (FASTA and JSON)
- Content listing in service
- Support formats with metadata to be extracted (need to determine which)
- JSON input integrated into publisher for attribute and map specification
- Manage config externally (retire ini) using CV services and support mutable DRS projects