

Supporting Experimental Science at NCSA

Margaret Johnson, Assistant Director



ILLINOIS

NCSA | National Center for
Supercomputing Applications

Why I am here

- I have an interest in supporting experimental science, especially those which follow a HTC model
- Learn more about the HTC community ecosystem and what other institutions are doing
- Give you a flavor of the open science activities at NCSA at campus and national levels

Support at all scales

- NCSA has a role in projects which require support at the scale of a center
- Portfolio of services and capabilities
 - Direct support for experiments
 - Reusable middleware tools and frameworks
 - Technology services which are not purely big-scale MPI
- Mix of innovation and strategic and practical reuse of technologies and processes
- NCSA mission includes strategic support for campus
 - Collaborate with UIUC Technology Services and Engineering IT to deliver a portfolio of IT services supporting research at Illinois

Direct Support of Experiments

- Sky Surveys
 - Dark Energy Survey (DES)
 - LSST
- Earth Science
 - Terra Data Fusion Project
- Crop Science and Genomics
 - Transportation Energy Resources from Renewable Agriculture Phenotyping Reference Platform (TERRA-REF)
- etc.

Dark Energy Survey Data Management

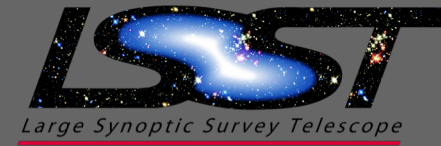


- DES is a 5.5-year mission designed to measure cosmic acceleration → requires large, well-calibrated, uniformly processed, well-described dataset
- NCSA is an institutional partner in the DES collaboration and provides data processing, management, quality assessment, and distribution for the experiment
- Raw data is reduced nightly + reprocessed every year for annual data releases
- Integrate existing infrastructure
 - Production data processing framework built on HTCondor
 - Use dedicated nodes in the Illinois Campus Cluster (ICC), as well as FermiGrid and Blue Waters (via glideins)
 - Centralized data management with location, metadata, provenance (and science catalogs) managed in Oracle database
 - Raw and derived data products on the NCSA Storage Condominium with large VM edge infrastructure for efficient bulk data movement, software packaging, and data analysis



AST-1138766

LSST Data Facility (LDF)



- NCSA is the host institution for the LDF, which will be responsible for transferring, storing, processing and disseminating the LSST data.
- Data will be processed (and reprocessed) on many timescales – real-time, daily, and annually for generating data releases to the community.
- NCSA will also host the US Data Access Center. LDF will manage the data at NCSA and Chile, and distribute to CC-IN2P3 for DR and satellite computing and to facilities with agreements.
- We will also manage computing, storage and networking, as well as provide facility services like network scanning and identity management.

Terra Data Fusion Project



- Collaborative effort by UIUC Department of Atmospheric Sciences (Di Girolamo), NCSA, NASA, and HDF Group
- Brings together ~1PB of data from 5 instruments on NASA's EOS Terra satellite (20 years of data) to create mission-scale value-added dataset
- Goal is to produce basic fusion products and facilitate use by the community in generating and disseminating new data products through use of national computing and data services.
 - Challenges: transfer data from multiple data centers, variation in file formats and metadata
 - Resulted in Basic Fusion dataset of ~2.5 PB
- Data products are residing on nearline at NCSA. Current effort is focused on building the data distribution and analysis infrastructure.



- Data management is collaborative effort between University of Arizona, University of Illinois, and NCSA.
- ~2TB of data are collected per day from gantry field scanner with multiple instruments scanning research field in Arizona, plus aerial scanning and ground vehicles.
- Goal is high-throughput phenotyping. Measure plant traits, produce a large reference dataset for crop genomics, and facilitate use of reference data by the community in accessing, analyzing, and contributing new data products.
- Data are transferred to NCSA for processing and archiving (5PB).
- Data and computation pipeline developed by NCSA to collect, transfer, process and distribute large volumes of crop sensing and genomic data. NCSA runs an instance of Clowder server for data management and visualization.

Software Tools and Frameworks

- Innovative Software and Data Analysis (ISDA) group develops tools for data curation, data management, and visualization through collaborations with researchers at Illinois and nationwide
- Leverage solutions to create scalable and reusable software tools and frameworks for data analysis. Broad participation in NSF DIBBs.
- Clowder
 - A data management system designed to support any data format and multiple research domains
 - Scalable - handles PBs of files
 - Supports both automatic metadata extraction and manual annotation of data elements
 - We are considering Clowder as a metadata backend for creating a cloud modeling training set based on Terra fusion data
 - We are considering Clowder for LSST to manage auxiliary instrument data





Architectural Vision for Research Cyberinfrastructure

Discipline Specific Environments

Science Portals

Applications & Frameworks

Research Facilities

Community Metadata Improvement #1443062

Computer Aided Discovery in Geo #1442997

Local Spectroscopy Data Infra. #1640899

Mobile Sensor Data #1640813

Archiving U-Series Geochronologic Data #1443037

Middleware & Analytics Libraries #1443054

Ocean Cloud Commons #1640775

Modular Eng/Sci Cyber-platform #1443027

Whole Tale #1541450

ClearEarth #1443085

Continuous Capture of Metadata #1640575

Materials Engineering Data Lab #1640867

STORM #1443046

Triple Gateway #1443040

Nanocomposite Resource #1640840

Virtual Data Collaboratory #1640834

Vizier #1640864

Transient Data Access #1443083

LearnSphere #1443068

Collaborative Workflow Design #1443069

Spatial Data Synthesis #1443080

Brown Dog #1261582

HUBzero Geospatial #1261727

Virtual Info. Fabric Infrastructure #1640818

Clowder #1835834

Pacific Research Platform #1541349

SciServer #1261715

User Driven Architecture #1443070

Scalable Data Delivery Platform #1541318

SEAD #0940824

Gravitational-Wave Workflow #1443047

DataONE



XSEDE

Extreme Science and Engineering Discover



Science Grid



OSiRIS: Ceph and SDN #1541335

Data Exacell #1261721

4CeeD #1443013

MBDH #1550320

North East Data Exchange #1640831

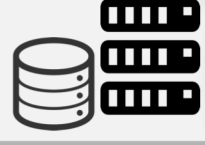
Aristotle Cloud Federation #1541215

NSF Resources

Commercial Resources

Resources

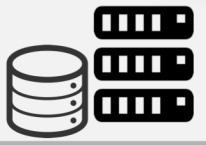
International Resources



OSN #1747552



Confidential Social Science Data #1443014



Integrative Services

Resources

Illinois Campus Cluster Program (ICCP)

- Shared HPC computing and storage services for Illinois faculty, research groups, and campus departments
 - Invest in node(s)/disk(s) or
 - Pay for on-demand use of cycles/storage (Research Computing as a Service (RCaaS), Research Storage as a Service (RSaaS))
 - Administrative costs are subsidized by campus
- ICCP is operated under the leadership of Deputy CIO for Research IT, John Towns
 - Dual role as NCSA Executive Director For Science And Technology and XSEDE Project Director
 - NCSA provides the technical personnel to operate the cluster
- Currently 36 investors (including DESDM and ATLAS Midwest Tier-2), many more users.
- Collaborative pilot project in 2018 to set up a HTC “Condor Annex” with 400 repurposed nodes.

NCSA Storage Condominium

- Reliable mid-scale data storage
- 7 PB usable storage and growing (3 PB more coming soon)
- Supports multiple access methods including GridFTP, NFS and native GPFS clients
- Edge infrastructure - offers a virtual machine capability to support various utility services for projects
 - E.g., a recent test suggests we can support 1 GB/s transfers to S3 bucket → load 1PB of data in 15 days

Current Specific Activities to Integrate into OSG

- Investigating StashCache to reduce data staging
 - DES single image processing (calibrations)
 - DES coadd processing (gather all the overlapping images),
 - Terra basic fusion data analysis
- Status:
 - A data origin server has been set up in the Kubernetes Federation. Imported ~50GB of test data (DES calibrations) so far.
 - We registered [/cvmfs/desdm.osgstorage.org](https://cvmfs.desdm.osgstorage.org) for our data.
 - A submit node with a flocking arrangement to OSG collectors has been set up.
 - 2 PB Terra fusion data are being migrated from nearline to the NCSA storage condo that will support a data origin server.

Current Specific Activities to Integrate into OSG

- To be worked out:
 - DESDM tends to run long jobs – concern about preemption. Will need changes to our framework and the way we package jobs. We are investigating ways to shorten jobs (i.e., parallelization) to fit better into OSG model.
 - DESDM relies on Oracle for workload management. How to handle passing Oracle credentials securely using OSG?
 - Setting up the edge infrastructure to push Terra data volume into a StashCache environment. Fused data products are very large (avg. 25 GB) and expectation is that only elements of the images will be loaded for further processing.
- Also - investigating Rucio with CC-IN2P3 for managing data distribution and long-term storage for LSST

Summary

- NCSA is directly supporting experimental science and providing research computing services at both campus and national scales.
- We have a broad interest in collaboration with OSG and the high-throughput community.
- We're just getting started.
- I'm here to learn.



Thanks!



ILLINOIS

NCSA | National Center for
Supercomputing Applications