# NEXUS:
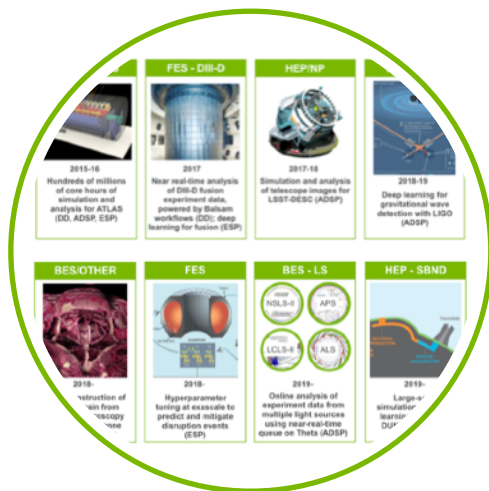# INTEGRATED RESEARCH INFRASTRUCTURE AT ARGONNE

Tom Uram, turam@anl.gov
Argonne Leadership Computing Facility
ALCF Webinar, March 2024

Argonne
NATIONAL LABORATORY

**ALCF Support
of Experimental
Scientific Computing**



**Integrated
Research
Infrastructure**



**Argonne
Nexus**

# ALCF Systems Evolution



**≥2 EF**

**44 PF**

**15.6 PF**

**11.7 PF**

**10 PF**

**557 TF**

**5.7 TF**

IBM BG/L
2004

**Intrepid**
IBM BG/P
2007

**Mira**
IBM BG/Q
2012

**Theta**
Intel-Cray XC40
2017

**Crux**
HPE-AMD

**+**

NVIDIA
DGX A100
2020

**Polaris**
HPE-AMD/
NVIDIA
2021

**Aurora**
Intel-HPE

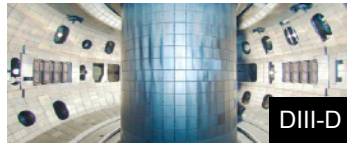**JLSE (2013)**          **AI Testbed (2020)**   **Edge Testbed (2021)**
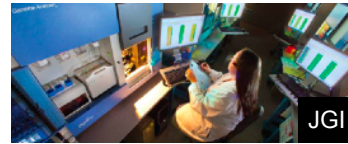
# DOE EXPERIMENTAL USER FACILITIES

- DOE operates 24 experimental user facilities
- Similar to the computing facilities, some of them are undergoing upgrades
- Their data rates and their computing needs will increase accordingly

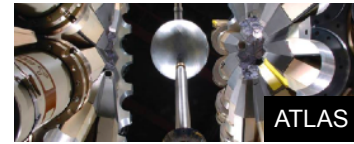

FNAL AC

DIII-D

JGI

NSTX-U

TMF

SNS

CNMS

ATLAS

RHIC

NSLS-II

CNM

SSRL

EMSL

HFIR

CFN

ATF

LCLS

FRIB

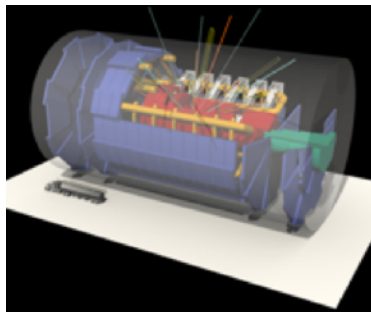CEBAF

ARM

CINT

FACET-II

APS

LCLS

# ALCF SUPPORT FOR EXPERIMENTAL SCIENTIFIC COMPUTING

**Allocation programs: Director's Discretionary, Early Science Program, ALCF Data Science Program**

**Technologies: Workflows, Scheduler, Globus**

## HEP - ATLAS



**2015-16**

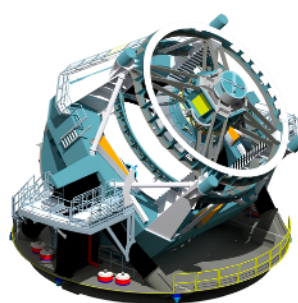**Hundreds of millions of core hours of simulation and analysis for ATLAS (DD, ADSP, ESP)**

## FES - DIII-D



**2017**

**Near real-time analysis of DIII-D fusion experiment data, powered by Balsam workflows (DD); deep learning for fusion (ESP)**
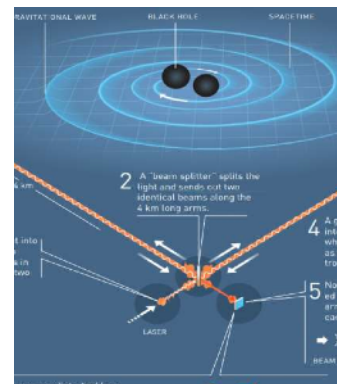
## HEP/NP



**2017-18**

**Simulation and analysis of telescope images for LSST-DESC (ADSP)**

## HEP - LIGO



**2018-19**

**Deep learning for gravitational wave detection with LIGO (ADSP)**

# ALCF SUPPORT FOR EXPERIMENTAL SCIENTIFIC COMPUTING

**Allocation programs: Director's Discretionary, Early Science Program, ALCF Data Science Program**

**Technologies: Workflows, Scheduler, Globus**



### BES/OTHER

2018-

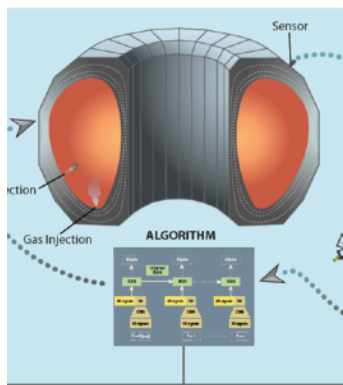**3D Reconstruction of mouse brain from APS imagery at Argonne (DD, ADSP, ESP)**
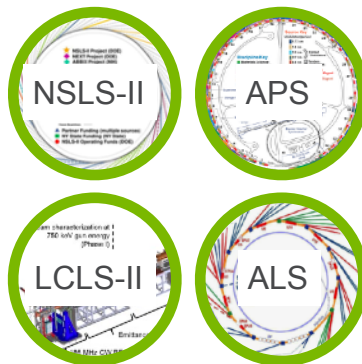
### FES

2018-

**Hyperparameter tuning at exascale to predict and mitigate disruption events (ESP)**
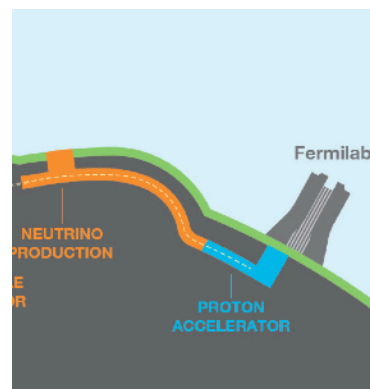
### BES - LS

2019-

**Online analysis of experiment data from multiple light sources using near-real-time queue on Theta (ADSP)**
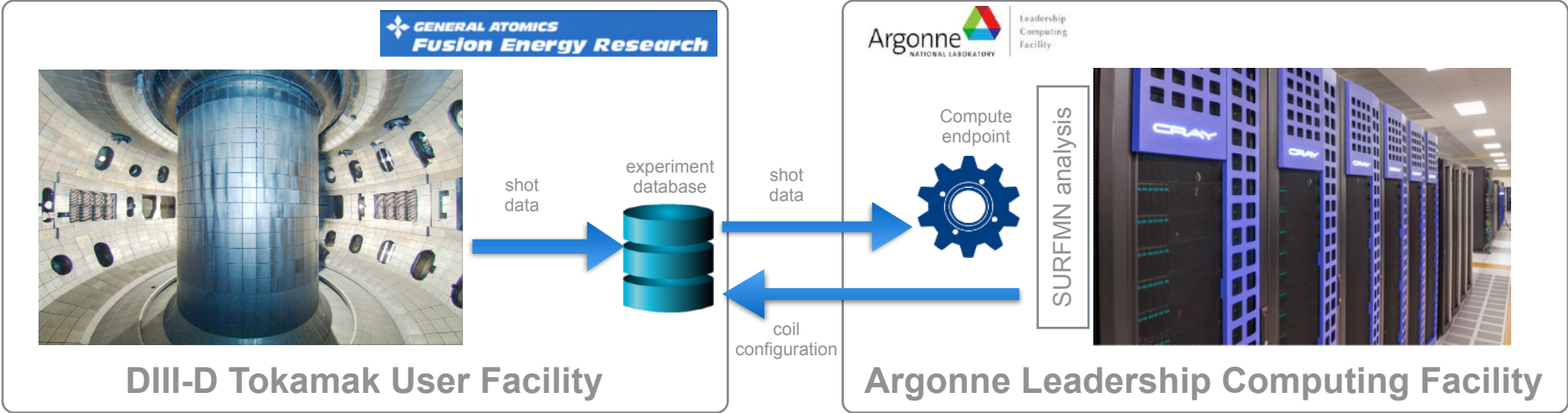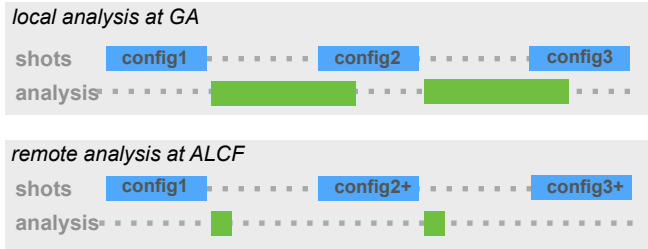
### HEP - SBND

2019-

**Large-scale simulation and deep learning for SBND/ DUNE (DD, ADSP)**

# AUTOMATIC BETWEEN-SHOT ANALYSIS OF DIII-D EXPERIMENTAL DATA



DIII-D Tokamak User Facility
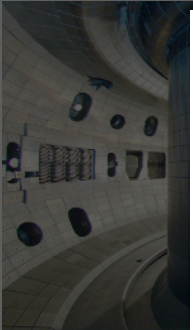
Argonne Leadership Computing Facility

- ‣ Scientists configure experimental "shots" every 20 minutes
  - A shot is an attempt to magnetically confine high temperature plasma
  - The timing/current of magnetic coils are configured to control the plasma during a disruption to avoid damage to the containing vessel (applicable to DIII-D and future reactors)
  - Analyses indicate how to optimize coil configuration for confinement
- ‣ Each shot triggers an automatic, near real-time analysis job at ALCF
- ‣ GA scientists integrate analysis results into configuration for next shot
- ‣ Analysis at ALCF enables more complex analyses (16x resolution) to be completed faster, improving the accuracy of results and allowing analyses to inform every shot instead of every other



*Faster analysis time allows analysis results to be integrated into magnet configuration for subsequent shots. Higher resolution analyses improve configuration accuracy.*

M. Kostuk, T. Uram, et al, 2nd IAEA Technical Meeting on Fusion Data Processing, Validation, and Analysis, 2018

# Automatic Between-shot Analysis of DIII-D Experimental Data



GENERAL ATOMICS
**Fusion Energy Research**

Argonne NATIONAL LABORATORY — Leadership Computing Facility

**DIII-** ... **ng Facility**

local analysis at GA

shots: **config1** · · · · · · **config2** · · · · · · **config3**

analysis: · · · · · ▮▮▮▮▮ · · · · ▮▮▮▮▮ · · · ·

remote analysis at ALCF

shots: **config1** · · · · · · **config2+** · · · · · · **config3+**

analysis: · · · · ▮ · · · · · ▮ · · · · · · ·

*Faster analysis time allows analysis results to be integrated into magnet configuration for subsequent shots. Higher resolution analyses improve configuration accuracy.*

- Scientists conf...
  - A shot is...
  - The timi...
    disrupti...
    reactors...
  - Analyses...
- Each shot trigg...
- GA scientists integrate analysis results into configuration for next shot
- Analysis at ALCF enables more complex analyses (16x resolution) to be completed faster, improving the accuracy of results and allowing analyses to inform every shot instead of every other

**config3**

**config3+**

*Faster analysis time allows analysis results to be integrated into magnet configuration for subsequent shots. Higher resolution analyses improve configuration accuracy.*
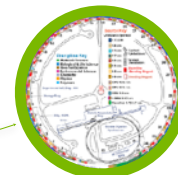
# ALCF SUPPORT OF LIGHT SOURCE COMPUTING

- Remote computing from three light sources at ALCF
- Computing demands are increasing
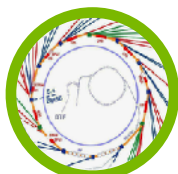
**ALCF Theta (11.7 PetaFLOPs)**
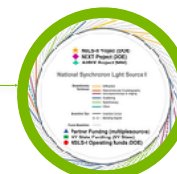


**APS**

XPCS workloads reaching 0.1PF by 2024

XPCS workloads reaching 0.1PF by 2024

XPCS workloads reaching 2.5PF by 2024

**ALS**

**NSLS-II**

# NEAR REAL-TIME PROCESSING OF WORKLOADS FROM THREE LIGHT SOURCES AT ALCF

**Experiment**
- Transfer 40GB input dataset from APS, ALS, NSLS-II
- Analyze data in near real-time with XPCS-Eigen* using **backfill queue** on Theta
- Transfer results to originating light source

**Results**
- Continuously executed for **48+ hours**
- Transferred 23TB input data from APS/ALS/NSLS-II to ALCF
- Analyzed 500+ datasets
- Transferred 179GB output data from ALCF to APS/ALS/NSLS-II



* https://github.com/AdvancedPhotonSource/xpcs-eigen

# AUTOMATED ANALYSIS BETWEEN TWO LIGHT SOURCES AND THREE COMPUTE FACILITIES

- Individual jobs and data transfers scheduled with single Python call

- XPCS dataset (878MB) transferred from APS to Theta/Summit/Cori (left), ALS to Theta/Summit/Cori, and APS/ALS to Theta/Summit/Cori

- Pool of nodes maintained for fast injection of jobs on arrival, and immediate return of results to originating site

# INTEGRATED RESEARCH INFRASTRUCTURE BLUEPRINT ACTIVITY REPORT (2023)

**THE DOE OFFICE OF SCIENCE**

**Integrated Research Infrastructure Architecture Blueprint Activity**

**FINAL REPORT**
2023

## IRI Science Patterns (3)

**Time-sensitive pattern** has *urgency*, requiring real-time or end-to-end performance with high reliability, e.g., for timely decision-making, experiment steering, and virtual proximity.

**Data integration-intensive pattern** requires combining and analyzing data from multiple sources, e.g., sites, experiments, and/or computational runs.

**Long-term campaign pattern** requires sustained access to resources over a long period to accomplish a well-defined objective.

## IRI Practice Areas (6)

**User experience practice** will ensure relentless attention to user perspectives and needs through requirements gathering, user-centric (co)-design, continuous feedback, and other means.

**Resource co-operations practice** is focused on creating new modes of cooperation, collaboration, co-scheduling, and joint planning across facilities and DOE programs.

**Cybersecurity and federated access practice** is focused on creating novel solutions that enable seamless scientific collaboration within a secure and trusted IRI ecosystem.

**Workflows, interfaces, and automation practice** is focused on creating novel solutions that facilitate the dynamic assembly of components across facilities into end-to-end IRI pipelines.

**Scientific data life cycle practice** is focused on ensuring that users can manage their data and metadata across facilities from inception to curation, archiving, dissemination, and publication.

**Portable/scalable solutions practice** is focused on ensuring that transitions can be made across heterogeneous facilities (portability) and from smaller to larger resources (scalability).
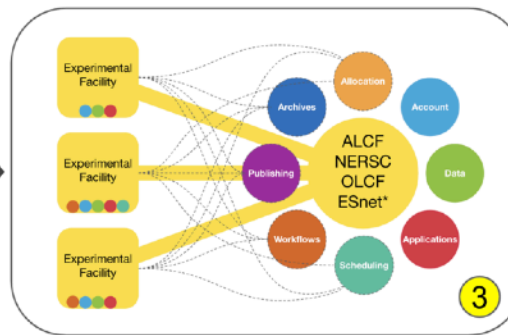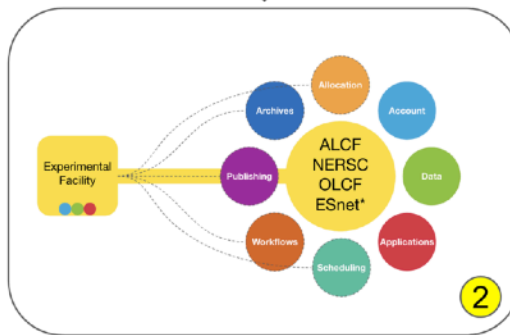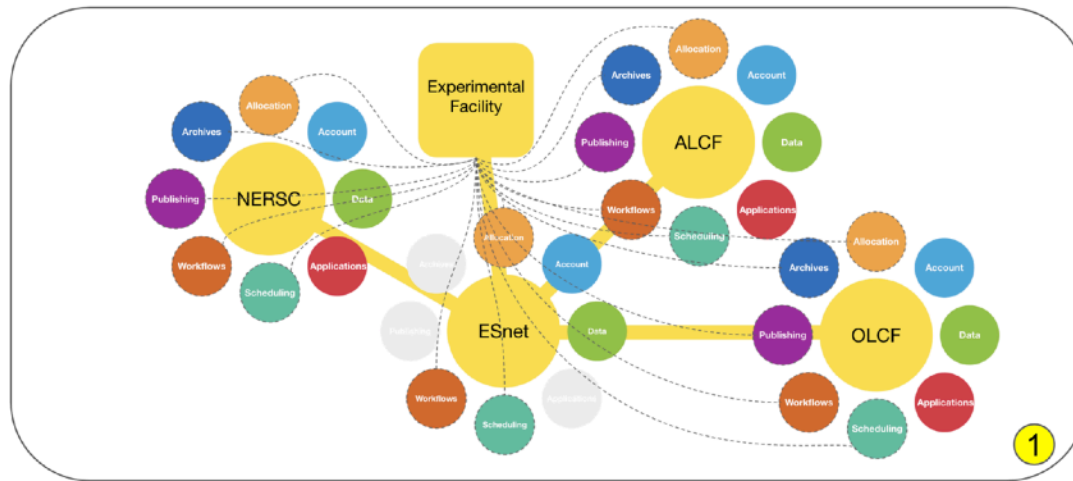
Argonne
NATIONAL LABORATORY

# IRI White Paper March 2021

Depiction of the integration of experimental facilities with computational facilities, across the range of services provided, in contrast with the one-to-one approach required today.

1. Today, an experimental facility must arrange separate bespoke interactions with individual HPC/HPN facilities.

2. A future paradigm with common interfaces could simplify integration of an experimental facility with multiple HPC/HPN facilities.

3. In turn, these common interfaces could support expansion and integration **across multiple experimental facilities and HPC/HPN facilities**.



| | |
|---|---|
| AL | Allocations |
| AC | Accounts |
| DA | Data |
| AP | Applications |
| SC | Scheduling |
| WF | Workflows |
| PB | Publication |
| AR | Archiving |

Toward a Seamless Integration of Computing, Experimental, and Observational Science Facilities: A Blueprint to Accelerate Discovery
(www.osti.gov/servlets/purl/1863562)

**DOE's Integrated Research Infrastructure (IRI) Vision:**
*To empower researchers to meld DOE's world-class research tools, infrastructure, and user facilities seamlessly and securely in novel ways to radically accelerate discovery and innovation*



Experimental and Observational User Facilities

Advanced Networking

Edge Sensors

Advanced Computing

Researchers

Local Campus Computing

Computing Testbeds

Data Management

High Performance Data Facility

Data Repositories PuRE Data Assets

Software

Software and Applications

AI Tools Digital Twins

Cloud Computing

**New modes of integrated science**

Rapid data analysis and steering of experiments

Novel workflows using multiple user facilities

AI-enabled insight from integrating vast data sources

U.S. DEPARTMENT OF ENERGY | Office of Science

# Nexus

Pioneering new approaches to integrating scientific facilities, supercomputing capabilities and data technologies.

Computational scientific research is evolving rapidly with faster data acquisition rates, larger datasets, and increasingly complex processing workflows. Advanced research instruments like the X-ray light sources across the U.S. Department of Energy's (DOE) national laboratory system produce vast amounts of data. Automating these interconnected research processes is critical to fully utilize the power of supercomputing and leading-edge data storage and technologies to drive breakthrough science

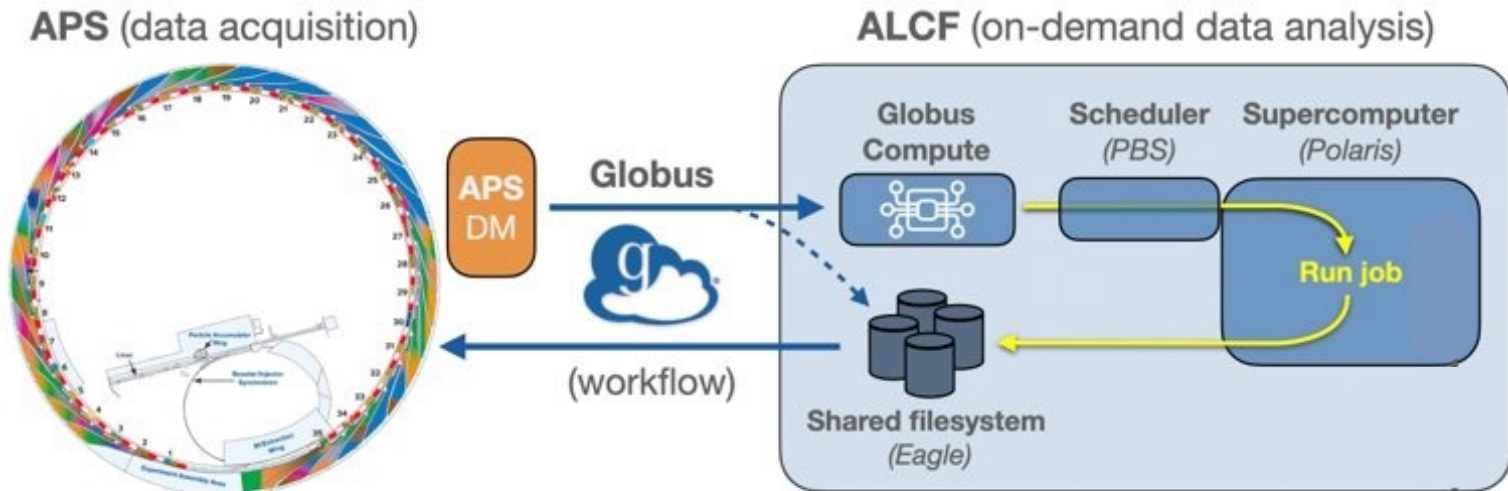| Nexus |
| --- |
| About Nexus |
| Nexus Projects          > |
| Nexus Publications |

SHARE

# ARGONNE NEXUS

- The **Nexus** initiative at Argonne enables experimental facilities to leverage supercomputing facilities for experiment-time data analysis
- Nexus goals align with the DOE vision for integrated research infrastructure (IRI)
  - Simplify access to DOE computing resources (accounts, job submission)
  - Automate data transfer and workflows triggered by experiment events
  - Provide robust, generalizable solutions, not tailored to a single experiment/facility
  - Support interactive inspection and provenance of data and derived products
  - Publish data for community access, citation, and archiving

Argonne
NATIONAL LABORATORY

# ARGONNE NEXUS SERVICES

- Demand Queue
    - Deployed now on Polaris: reduces queue wait time for experiment-time analysis; backfilled by preemptable jobs
- Service Accounts
    - Provides experiment-specific accounts for running analyses in a controlled environment
- Eagle Data Sharing (100PB filesystem)
    - Users can define collections of data to share publicly or with designated Globus users, without an ALCF account
    - Leveraged in ALCF-HEP cosmology data sharing portal
- Globus Infrastructure
    - Distributed compute, transfer, and web-based monitoring
- ALCF Community Data Cooperative (ACDC)
    - Layer atop Eagle data sharing to enable more sophisticated metadata-based navigation and search
    - Metadata extraction and capture services
- Dedicated Web Applications
    - Foundation for user-driven access to data and analysis
    - Developing in collaborations with APS and Argonne-HEP for extension to other science areas
    - Data catalogs resident on Eagle filesystem, searchable via user-provided metadata
    - Community analysis of data supported via back-end integration of workflows on Polaris
    - Reusable web components can be deployed and customized for future science domains

# ARGONNE NEXUS: LIGHTSOURCE AUTOMATION



- Integration with the data management (DM) system at APS **allows the workflow to begin as soon as data is obtained**
- Workflow moves data from the APS beamline to ALCF and submits job to demand queue on Polaris
- Results are written to Eagle, where they're reachable via Jupyter, and also returned to APS for evaluation

**One-time configuration at ALCF**

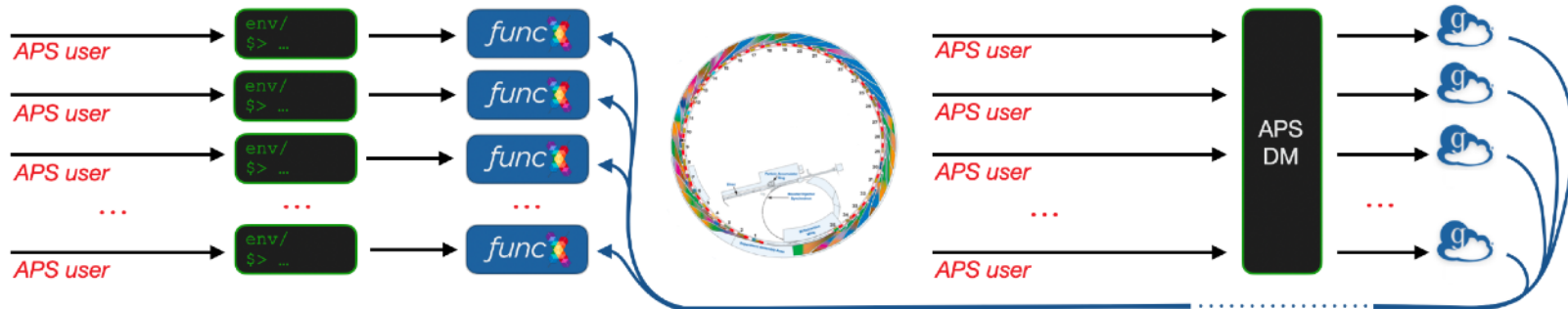**APS experiments**

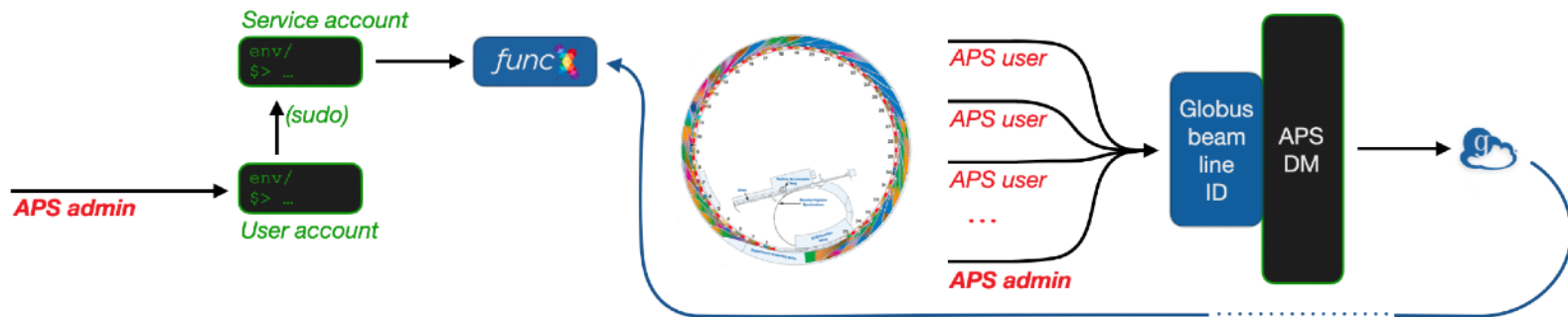**Using user accounts**

**Humans involved in the workflows**

APS requests user accounts at ALCF

ALCF creates user accounts

APS users setup their **personal environment** and funcX endpoint

APS users log into DM system using **personal credentials**

DM system starts Globus flows involving **specific funcX endpoints**

APS user

APS user

APS user

. . .

APS user

APS DM

**One-time configuration at ALCF**

**APS experiments**

**Using a service account per beam line**

**No human involved in the workflows**

APS requests service account at ALCF

ALCF creates service account

APS admin user setups a **shared environment** and funcX endpoint

APS users log into DM system using **shared Globus identity**

DM system starts Globus flows involving **shared funcX endpoint**

Service account

(sudo)

APS admin

User account

APS user

APS user

APS user

. . .

APS admin

Globus beam line ID

APS DM

# NEXUS SUPPORTS SEVERAL APS BEAMLINES

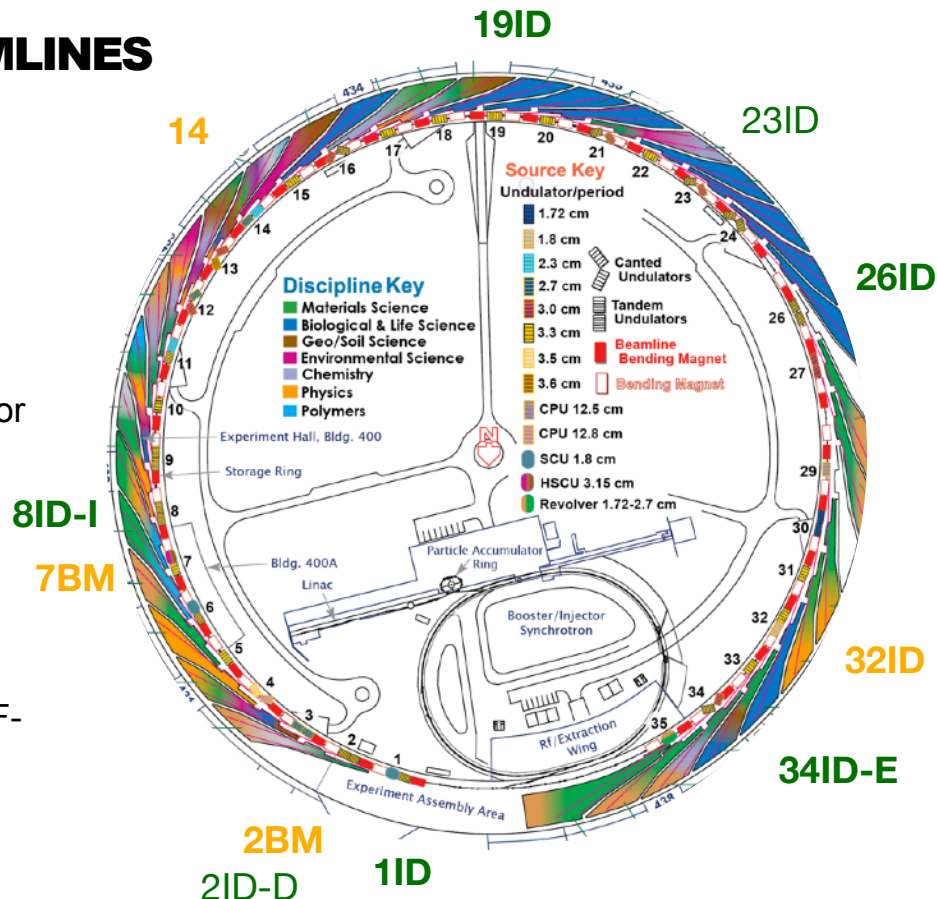**Several APS beamlines have run production experiments with Globus Compute**
- 8ID-I: X-ray Protein Cystal Spectroscopy (XPCS)
- 19ID: Serial Crystallography (SSX) at the structural biology center (SBC), using DIALS for crystallography analysis
- 1ID: High energy diffraction microscopy using MIDAS for tomography and ptychography
- 34ID-E: Laue Micro-diffraction spectroscopy
- 26ID: Ptychography with Ptychodus and Tike

**Other beamlines are running but have not yet run in production**
- 23ID: GM/CA using crystfel toolset (+AlphaFold)
- 2ID-D: X-ray fluorescence image processing using XRF-Maps

**Ramping up**
- 14: BioCARS using crystfel
- 2BM, 7BM, 32ID: Tomography using tomocupy



**running, in production**
running, ready for production
**starting up**

# Leveraging Globus Tools

**Compute**: A managed service that implements a universal computing fabric

Local **Transfer** and **Compute** agents provide global footprint for actions

**Auth** provides federated identities and distributed authorization with delegation

**Flows** orchestrates actions across the computing continuum

**UIs** to monitor, manage, inspect it all

# DIII-D ANALYSIS AT THREE DOE COMPUTE FACILITIES (WIP)

- Experiment shot occurs every 20 minutes
- Simulation of plasma under experimental conditions is initiated based on current experiment parameters
- Flow can be initiated manually or **when shot data is available**; data is transferred to target compute site
- Workflow **conditionally falls back** to a secondary site
- Currently adapting workflow to read directly from the MDSPlus database at DIII-D

**Argonne:** Christine Simpson, Tom Uram, Bill Allcock, Mike Papka
**DIII-D/General Atomics**: Mark Kostuk, Sterling Smith, Juan Colmenares, David Schissel

DIII-D

Start Flow

PBSPro sched

ALCF

Slurm sched

NERSC

LSF sched

OLCF

time intensive

data intensive

long term

## Monitor runs

Sort ∨

- ⊘ FlowSucceeded
- Transfer_Out — ActionCompleted (9 seconds)
- Transfer_Out — ActionStarted (0 milliseconds)
- Postprocessing — ActionCompleted (19 seconds)
- Postprocessing — ActionStarted (0 milliseconds)
- IonOrb — ActionCompleted (52 seconds)
- IonOrb — ActionStarted (0 milliseconds)
- Transfer_In — ActionCompleted (26 seconds)
- Transfer_In — ActionStarted (34 milliseconds)
- FlowStarted

Argonne
NATIONAL LABORATORY

# MULTI-FACILITY CLIMATE DATA ANALYSIS

- Search for data, resolve to computing facility (ALCF, OLCF), analyze locally
- Choice of computing facility driven by user input, data locality, or resource availability
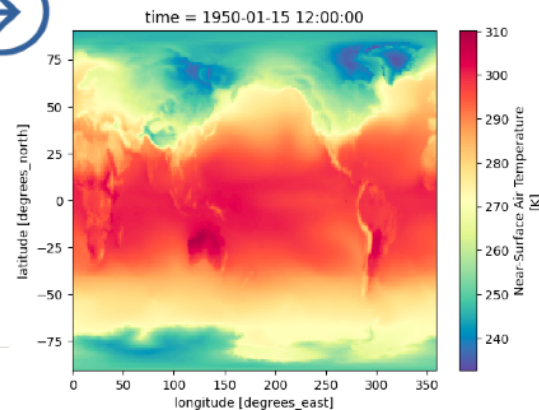


Monitor and manage runs



Dynamically generated web interface
to start analysis

**Argonne:** Ian Foster, Mike Papka, Max Grover, Scott Collis, Tom Uram, Christine Simpson, Bill Allcock, Benoit Cote, Ryan Chard
**UChicago:** Kyle Chard, Nick Saint
**JLab:** Amitoj Singh, Xinxin Mei, Chris Larrieu, Bryan Hess
**ORNL:** Forrest Hoffman, Nathan Collier

- ACDC provides a data publication interface backed by Eagle filesystem
- Data can be tagged with application metadata to facilitate searching
- User-defined plots and visualizations can be included via templates
- One-click access to transfer datasets via Globus

http://acdc.alcf.anl.gov

## Data from MiraTitanU Snapshots

The Mira-Titan Universe simulation suite was carried out on Mira, a supercomputer at the Argonne Leadership Computing Facility, and Titan, at the Oak Ridge National Laboratory. The simulations cover a range of cosmological models including models with a dynamical dark energy equation of state parameterized via $w_0$ and $w_a$. Each simulation covers a $(2.1 Gpc)^3$ volume and evolves $3200^3$ particles. We provide outputs for 27 redshifts, between z=4 and z=0, including halo information and down-sampled particle information.

Please select one or more models from the list below, then select all the relevant redshifts and data products. The SubmitTransfer button will indicate the number and overall size of the selected files that you aim to transfer. This button will lead you to the Globus interface. The Search box at the top allows you to narrow the model selection by specifying a model number or a numerical value for any cosmological parameter.

Search: [                    ]

| | Model | $\Omega_{cdm}$ | $\omega_b$ | $\omega_v$ | h | $\sigma_8$ | $n_s$ | $w_0$ | $w_a$ |
|---|---|---|---|---|---|---|---|---|---|
| ☐ | M000 | 0.2200 | 0.02258 | 0.0 | 0.7100 | 0.8000 | 0.9630 | -1.0000 | 0.0000 |
| ☐ | M001 | 0.3276 | 0.02261 | 0.0 | 0.6167 | 0.8778 | 0.9611 | -0.7000 | 0.6722 |
| ☐ | M002 | 0.1997 | 0.02328 | 0.0 | 0.7500 | 0.8556 | 1.0500 | -1.0330 | 0.9111 |
| ☐ | M003 | 0.2590 | 0.02194 | 0.0 | 0.7167 | 0.9000 | 0.8944 | -1.1000 | -0.2833 |
| ☐ | M004 | 0.2971 | 0.02283 | 0.0 | 0.5833 | 0.7889 | 0.8722 | -1.1670 | 1.1500 |
| ☐ | M005 | 0.1658 | 0.02350 | 0.0 | 0.8500 | 0.7667 | 0.9833 | -1.2330 | -0.0445 |
| ☐ | M006 | 0.3643 | 0.02150 | 0.0 | 0.5500 | 0.8333 | 0.9167 | -0.7667 | 0.1944 |
| ☐ | M007 | 0.1933 | 0.02217 | 0.0 | 0.8167 | 0.8111 | 1.0280 | -0.8333 | -1.0000 |
| ☐ | M008 | 0.2076 | 0.02306 | 0.0 | 0.6833 | 0.7000 | 1.0060 | -0.9000 | 0.4333 |
| ☐ | M009 | 0.2785 | 0.02172 | 0.0 | 0.6500 | 0.7444 | 0.8500 | -0.9667 | -0.7611 |

- HEP-funded project to develop analysis capabilities coincident with cosmology data sharing portal
- Portal currently exposes 150TB of HACC cosmology simulation data from 5 simulations, searchable via metadata
- Users transfer results to their home institutions via Globus sharing, without requiring an ALCF account
- New Globus Compute-based capabilities will allow in-place analysis and visualization

data intensive

long term

PI: Katrin Heitmann          http://cosmology.alcf.anl.gov

Argonne NATIONAL LABORATORY

# ALCF IRI RESOURCES TODAY - EDITH

- Edith (aka Edge Dev Testbed) was a test and development system for Polaris
    - Has been used in ALCF/APS collaboration for years already
- The testbed consists of four physical servers, each containing:
    - (2) 2.4GHz AMD EPYC 7532 32-core Processors
    - (2) NVIDIA A100 GPUs
    - 512 GB RAM
- All ALCF fileystems mounted
    - Access to all data at ALCF (on-node, via Globus transfers)
    - Constrains testbed operations (e.g. no root access)
- Networking
    - Login node accessible to internet
    - Outbound connectivity via HTTP proxy or SSH tunnel
- Scheduling
    - Uses PBSPro, similar to other ALCF systems
- Workflows
    - Globus Compute possible today

# ALCF IRI RESOURCES FUTURE - POLARIS

- Polaris is a hybrid CPU/GPU leading-edge pre-Aurora testbed system

- 56 nodes available in **on-demand queue** for experimental data analysis workloads

- ALCF plans to make Polaris available as a future IRI resource

| | |
|---|---|
| Platform | HPE Apollo 6500 Gen 10+ |
| System Peak | 44 PF DP |
| Peak Power | 1.8 MW |
| Total System Memory | 280 TB CPU DDR4 + 87.5 TB GPU HBM2 |
| System Memory Type | DDR, HBM |
| Node Performance | 78 TF DP |
| Node Processors | (1) AMD EPYC "Milan" 7543P CPU (4) NVIDIA HGX A100 GPUs |
| System Size (nodes) | 560 |

- ALCF is looking to expand into more experimental and observational science areas, and would like to work with you

- The DOE-led IRI effort is formulating what future infrastructure should look like, and will be better informed if you describe your science case

- Thank you to the many people who have contributed to IRI, at DOE; at ALCF, HPDF, OLCF, NERSC, ESnet; the 150+ participants in the Blueprint Activity representing 28 user facilities

- Thank you to the many people who have contributed to ALCF efforts, at ALCF, APS, ALS, NSLS-II, DIII-D, CPAC

**Questions?**

Argonne
NATIONAL LABORATORY