

September 26, 2024 - NOBUGS

# Streamlining Scientific Discovery with Data Pipelines at the Advanced Photon Source



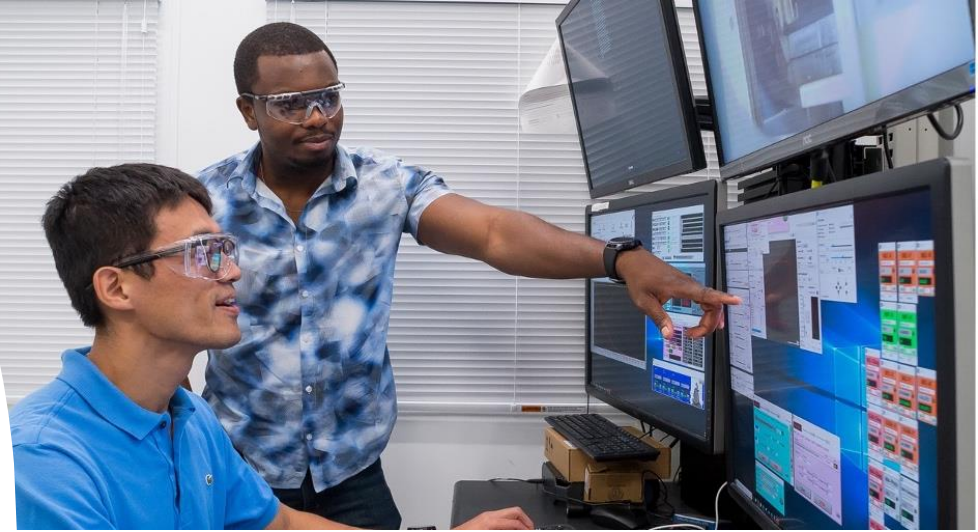
**Hannah Parraga**  
Scientific Software Engineering  
& Data Management Group  
X-ray Science Division  
Advanced Photon Source

R. Chard, J. Hammonds, P. Jemian, S. Henke,  
N. Saint, N. Schwarz, R. Vescovi, S. Veseli



# OVERVIEW

**The Advanced Photon Source is creating automated data processing pipelines leveraging high performance computing to address increasing data needs and enable scientific discovery**



# Scale of the Problem

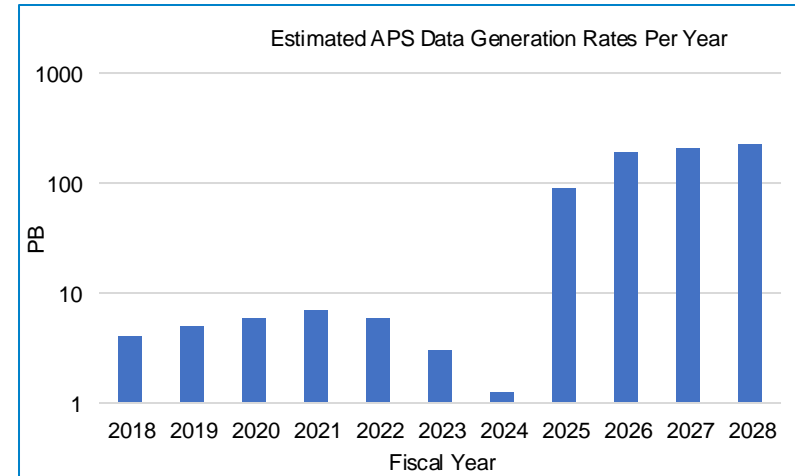
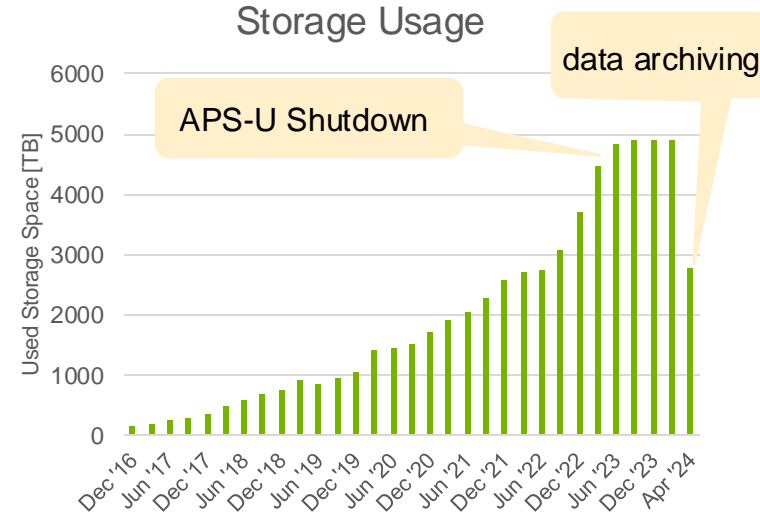
## Multiple order-of-magnitude increase in demand for computing resources

Over the past decade the APS has

- Created over 9000 experiments in the Data Management database
- Used 4.9PB of storage space

Over the next decade the APS will

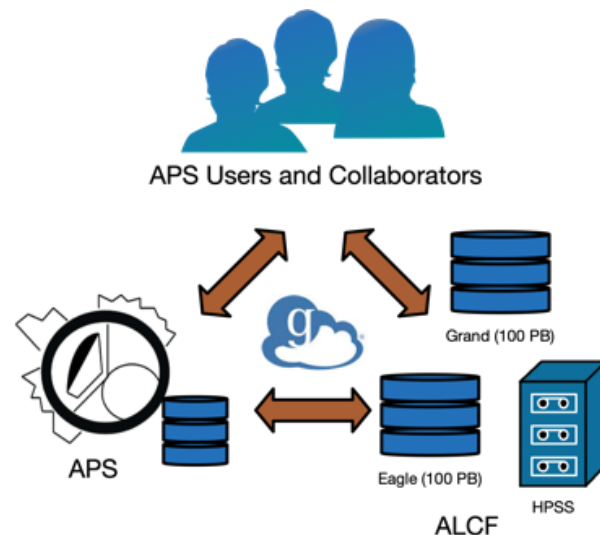
- Generate 100s of petabytes (PBs) of raw data per year
- Require 10s of petaflop/s of on-demand computing power



# The APS Data Management System

## Facility-wide software and hardware system for managing data

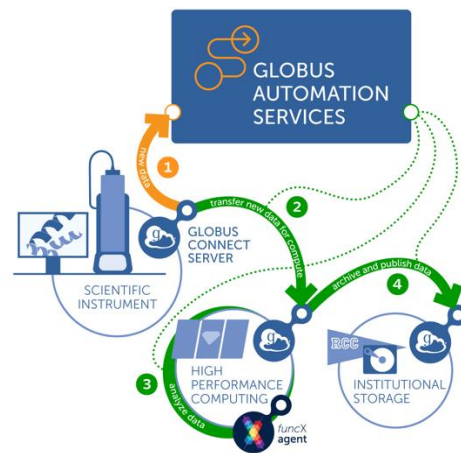
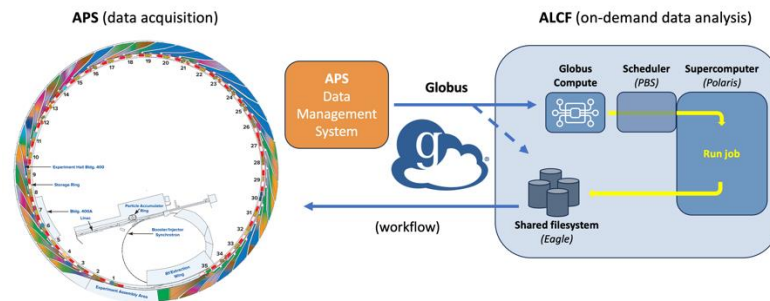
- Tools to automate **transfer** of data, manage storage **access** permissions, and **metadata** catalog
- Workflow tools automate data **processing** via plug-ins
  - Automated or user-initiated jobs
  - Flexible execution of any commands
  - Execute jobs locally or remotely
  - Integrated with Bluesky and Globus Compute



# Bridging Large-Scale Facilities

## Leveraging advances from the ALCF and Globus

- Partnership with Argonne Leadership Computing Facility (ALCF) means APS users do not need to navigate HPC
  - Priority “on-demand” queue for real time processing
  - Service accounts for beamlines
  - Computing allocation for APS
- Globus Compute
  - Function-as-a-service platform for remote job execution
  - Endpoints deployed at ALCF
  - Secure access to data and compute resources



# INTEGRATING CONTROLS WITH ANALYSIS

- Bluesky uses information on experiment conditions and where data is located to inform downstream analysis
- Bluesky plans use the APS Data Management API to launch workflows
- Users automatically get processed data when they run a Bluesky plan without extra effort

[apstools](#) / [apstools](#) / [devices](#) / [aps\\_data\\_management.py](#) 

 prjemian MNT #872

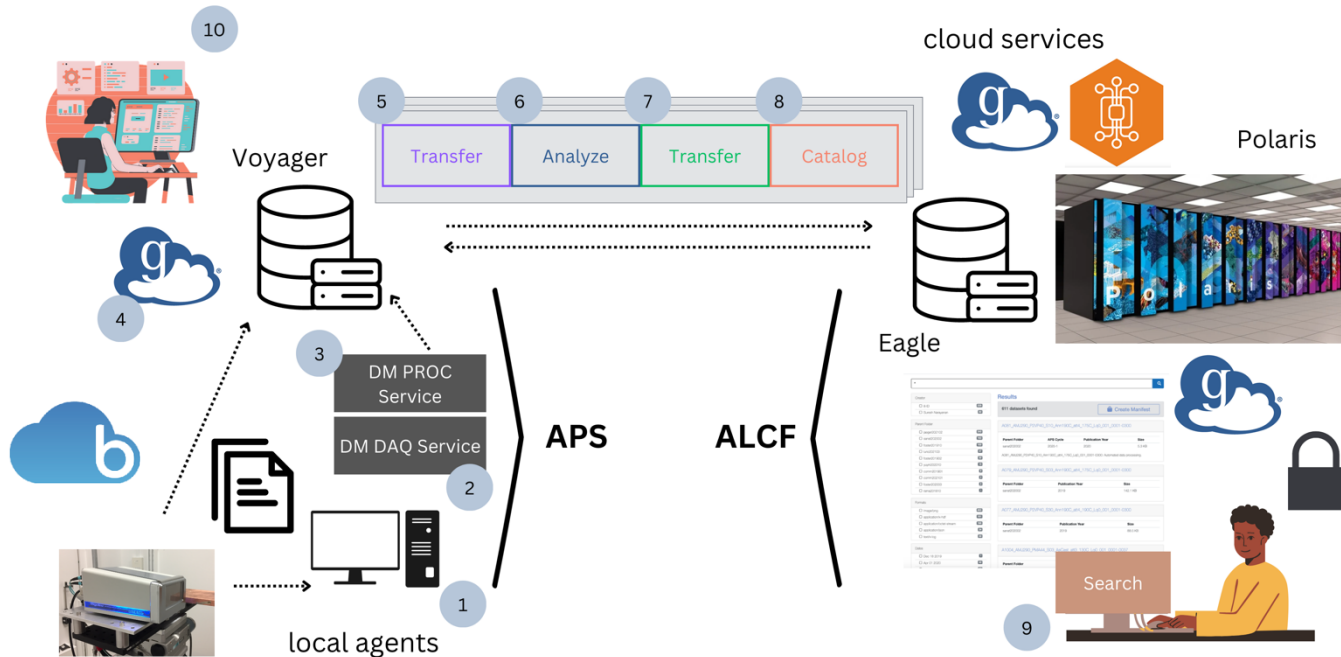
Code Blame 347 lines (287 loc) · 11.4 KB

```
1 """
2 Connect with APS Data Management workflows.
3
4 Example::
5
6     import bluesky
7     from apstools.devices import DM_WorkflowConnector
8
9     RE = bluesky.RunEngine()
10
11     dm_workflow = DM_WorkflowConnector(name="dm_workflow", labels=["DM"])
12     RE(
13         dm_workflow.run_as_plan(
14             workflow="example-01",
15             filePath="/home/beams/S1IDTEST/
16         )
17     )
18 """
```

Bluesky tools for APS by Pete Jemian



# DATA LIFE CYCLE



1. Users fill form about experiment  
 2. Experiment database info populated from form and directory created in central storage

3. Start monitor for files  
 4. Acquire data and writes metadata  
 5. Transfer files to storage  
 6. Select workflow and arguments

7. Data processed locally or with HPC  
 8. Result published to portal  
 9. After time, data archived to tape



# DATA PROCESSING

## Standardized processing for common X-ray techniques

### High-energy Diffraction Microscopy

1ID, 20ID

<https://github.com/marinerhemant/MIDAS>

### Wide Angle and Small Angle X-ray Scattering

1ID, 20ID

<https://github.com/marinerhemant/MIDAS>

### X-ray Photon Correlation Spectroscopy

8ID, 9ID

[https://github.com/AdvancedPhotonSource/boost\\_corr](https://github.com/AdvancedPhotonSource/boost_corr)

### Crystallography

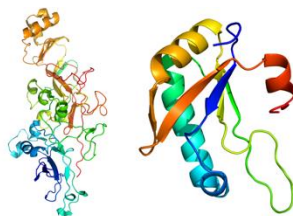
23ID

<https://www.gmca.aps.anl.gov/>

### Laue Micro-diffraction

34ID

Prince *et al.*, 2023



### Ptychography

2ID, 4ID, 12ID, 9ID, 19ID, 26ID, 28ID, 31ID, 33ID

<https://github.com/AdvancedPhotonSource/ptychodus>

### Grazing Incidence X-ray Scattering

9ID

Werzer *et al.*, 2024

### Tomography/Laminography

1ID, 2ID, 2BM, 7BM, 19ID, 20ID, 32ID

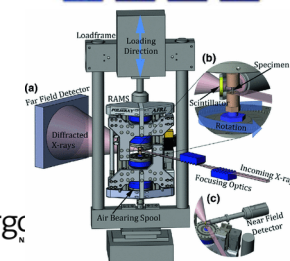
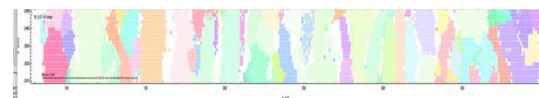
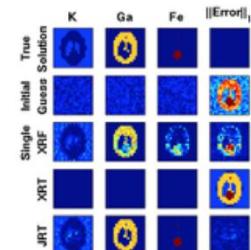
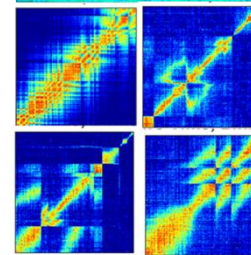
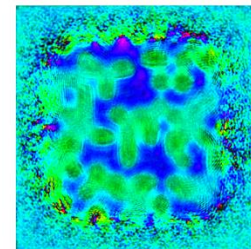
<https://github.com/tomography/tomocupy/>

### X-ray Fluorescence Microscopy

2ID, 19ID

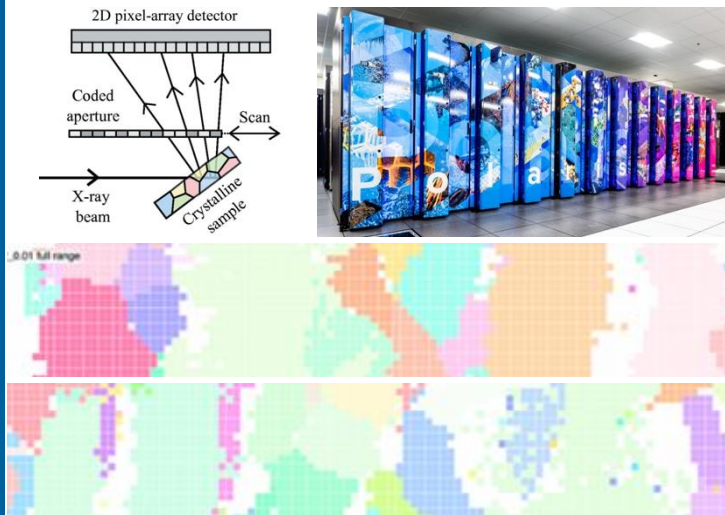
<https://github.com/AdvancedPhotonSource/XRF-Maps>

■ = operational  
■ = in progress





# POLARIS SUPERCOMPUTER ENABLES ON-DEMAND LAUE RECONSTRUCTIONS



Data from the new coded aperture at APS 34-ID-E (top left) is automatically transferred to the Polaris supercomputer (top right) where it is reconstructed on-demand in real-time (bottom)

Michael Prince, Doğa Gürsoy, Dina Sheyfer, Ryan Chard, Benoit Côté, Hannah Parraga, Barbara Frosik, Jon Tischler, and Nicholas Schwarz <https://doi.org/10.1145/3624062.3624613>

## Technical Achievement

Staff at the APS, the Argonne Leadership Computing Facility (ALCF), and Globus have successfully utilized the Polaris supercomputer for automated on-demand data processing of data for the APS-U 3D Micro- and Nano-diffraction (3DMN) Feature Beamline that will be typical in the APS-U Era.

## Significance and Impact

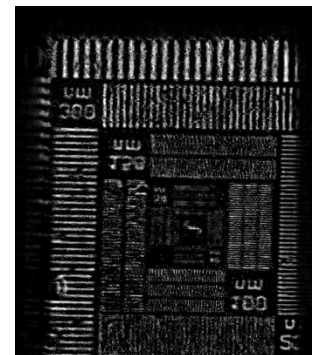
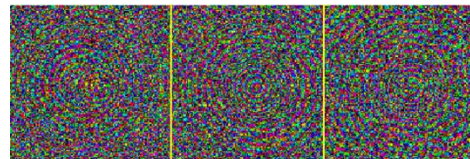
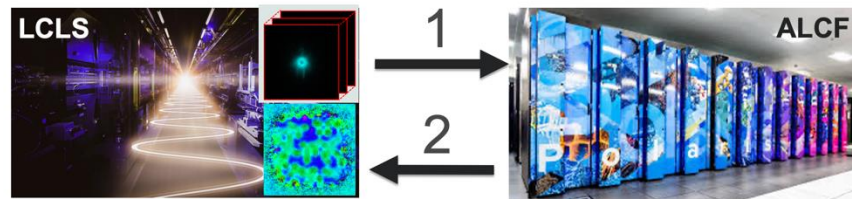
Using Polaris, the APS will be able to reconstruct coded aperture Laue datasets automatically in near real-time, accelerating time to science.

## Research Details

- Continuously utilized 50 nodes on Polaris (~4 petaflop/s) during beamtime
- New parallel GPU-based software implements a new coded aperture Laue reconstruction algorithm
- The APS Data Management System integrates with Globus workflow tools to provide a single end-to-end data pipeline

# PORTABLE CROSS-FACILITY WORKFLOWS FOR X-RAY PTYCHOGRAPHY

- Ptychography data volumes are expected to increase by *orders of magnitude* at leading X-ray research facilities due to next-generation upgrades
- Ptychography benefits from access to GPU computing resources
- We demonstrate cross-facility capabilities by deploying software at the **Linac Coherent Light Source (LCLS)**, packaging data and transferring it to **ALCF** for processing
- Used the ptychodus package

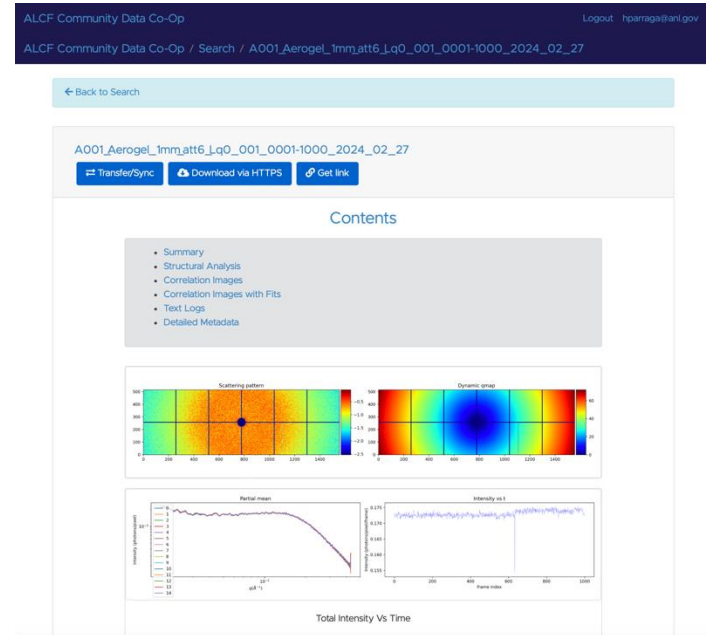


Albert Vong, 2024

# Data Portals

## The Globus Portal Framework Makes APS Data More Accessible to Users

- Built with Globus Django Portal Framework
- **Searchable** data and visualization
- Experiment specific metadata
- **Secure** access to processed data and metadata
  - Access controlled by Globus groups
- **Customizable** for each beamline
  - Close collaboration with beamline staff



# Multi-Tiered Approach

## Utilize local and remote resources

### High-end compute resources

- Large data processing tasks, ML training, post-processing, and data refinement

### Local compute resources

- Perform pre-analysis/data reduction
  - Compression and running ML models
  - Quality control and experiment steering
- May include a GPU workstation at a beamline or the APS computing cluster

#### Argonne Leadership Computing Facility (ALCF)



**Polaris**  
~44 PFLOP/s



**Aurora**  
> 1 EXAFLOP/s

~4 PFLOP/s of Polaris is prioritized for prototype on-demand use by experimental and observational facilities; when Aurora is in User operations, all of Polaris will be prioritized for on-demand use

**Synergy**  
Planning is underway for the next generation on-demand system prioritized for experimental and observational facilities

#### Next Generation Supercomputer

Planning is underway for the next generation leadership class supercomputer

#### Argonne Laboratory Computing Resource Center (LCRC)



**Improv**  
~2.51 PFLOP/s  
825 nodes with 2 AMD EPYC CPUs each

**Bebop**  
~1.75 PFLOP/s  
672 nodes with 36 Intel Broadwell cores each

**Swing**  
~925 TFLOP/s  
48 NVIDIA A100s | 768 AMD EPYC cores

#### Advanced Photon Source (APS)



#### APS general purpose distributed-memory compute cluster

- ~20 TFLOP/s CPU cores

#### ~50 High-Performance Computing Workstations

- Califone – 8 H100s
- Ecto – 8 RTX A6000s
- Refiner – 4 A100s
- Many others...

#### Edge Computing Devices

- 1 x NVIDIA Jetson AGX Orin
- 2 x NVIDIA BlueField-3 DPUs



# FUTURE WORK

- Expand use of data portals and integrate data management features into web UI
- Develop additional workflows and deploy at more beamlines
  - Xia2, GSAS II, combined ptychography/XRF, GIXS and more
- Develop streaming workflows using PvaPy Streaming Framework



# ACKNOWLEDGEMENTS

## ALCF:

- William Allcock
- Benoit Cote
- Jennifer Francis
- Carissa Holohan
- Ryan Milner
- Michael Papka
- Paul Rich
- George Rojas
- Haritha Siddabathuni Som

## DSL:

- Tekin Bicer
- Ian Foster
- Rajkumar Kettimuthu

## Globus:

- Rachana Ananthakrishnan
- Benjamin Blaiszik
- Ryan Chard
- James Pruyne
- Nickolaus Saint
- Rafael Vescovi

# ACKNOWLEDGEMENTS

## 3DMN & Atomic

- Jon Tischler
- Dina Sheyfer
- Wenjun Liu
- Ross Harder
- Wonsuk Cha

## CHEX

- Hua Zhou
- Dillon Fong
- Matthew Highland
- Hawoong Hong
- Stephan Hruszkewycz
- Zhan Zhang

## HEXM

- Jon Almer
- Andrew Chuang
- Peter Kenesei
- Jun-Sang Park
- Victoria Cooley
- Leighanne Gallington
- John Okasinski

## Imaging

- Francesco De Carlo
- Alan Kastengren
- Viktor Nikitin

## Polar

- Daniel Haskel
- Yongseong Choi

- Gilberto Fabbris
- Joerg Stempffer
- David Gagliano
- Jong Woo Kim
- Srtarshi Banerjee

## PtychoProbe & ISN

- Luxi Li
- Volker Rose
- Junjing Deng
- Barry Lai
- Curt Preissner
- Jorg Maser
- Zhonghou Cai
- Sarah Wieghold

- Olga Antipova
- Si Chen
- Yi Jiang
- Fabricio Marin
- Yanqi Luo

## XPCS & CSSI

- Suresh Narayanan
- Joe Strzalka
- Qingteng Zhang
- Eric Dufresne
- Zhang Jiang
- Jin Wang
- Ashish Tripathi
- Peco Myint

# ACKNOWLEDGEMENTS

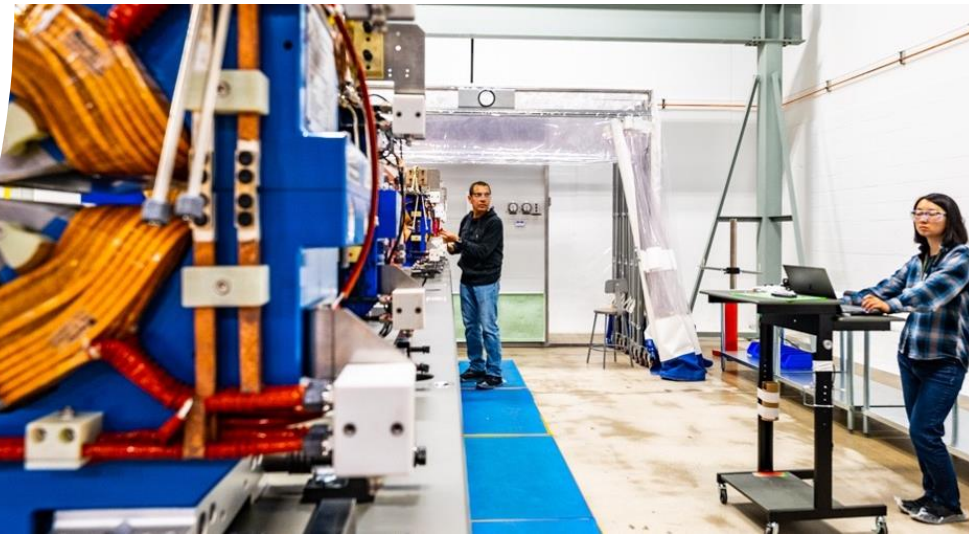
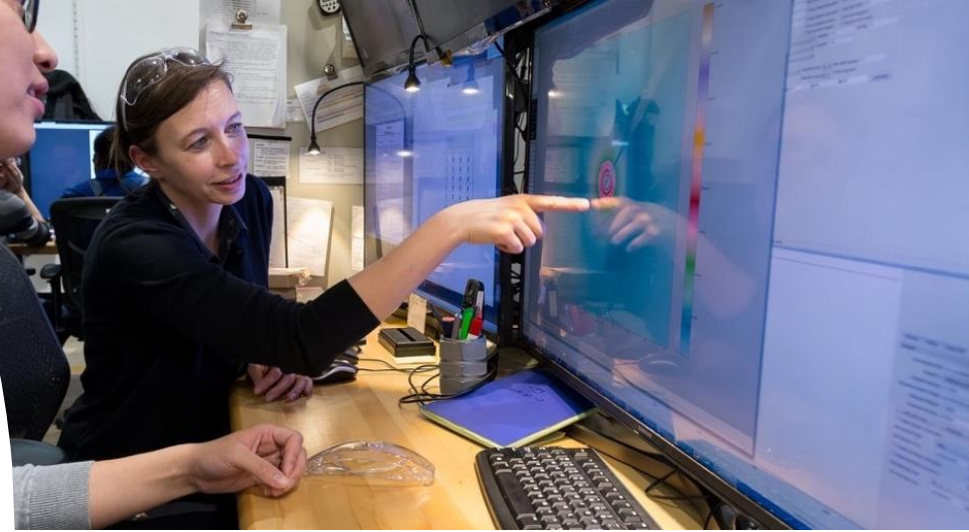
## APS:

- Ryan Aydelott
- Mathew Cherukara
- Miaoqi Chu
- Barbara Frosik
- Arthur Glowacki
- Doga Gursoy
- Tejas Guruswamy
- John Hammonds
- Steven Henke
- Pete Jemian
- Jonathan Lang
- Keenan Lang
- Alex Lavens
- David Leibfritz
- Antonino Miceli
- Tim Mooney
- Tom Naughton
- Michael Prince
- Alec Sandy
- Nicholas Schwarz
- Roger Sersted
- Hemant Sharma
- Kenneth Sidorowicz
- Joseph Sullivan
- David Wallis
- Mary Westbrook
- Max Wyman
- Sinisa Veseli
- Stefan Vogt
- Albert Vong
- Qingping Xu
- Xuan Zhang
- Tao Zhou



# SUMMARY

**The Advanced Photon Source is enabling scientific discovery and addressing increasing data needs by creating automated data processing pipelines leveraging high performance computing**



# REFERENCES

- Siniša Veseli, John Hammonds, Steven Henke, Hannah Parraga, and Nicholas Schwarz. 2023. Streaming Data from Experimental Facilities to Supercomputers for Real-Time Data Processing. In Proceedings of the SC '23 Workshops of The International Conference on High Performance Computing, Network, Storage, and Analysis (SC-W '23). Association for Computing Machinery, New York, NY, USA, 2110–2117. <https://doi.org/10.1145/3624062.3624610>
- Babu, A.V., Zhou, T., Kandel, S. *et al.* Deep learning at the edge enables real-time streaming ptychographic imaging. *Nat Commun* **14**, 7059 (2023). <https://doi.org/10.1038/s41467-023-41496-z>
- R. Pokharel, "Overview of High-Energy X-Ray Diffraction Microscopy (HEDM) for Mesoscale Material Characterization in Three-Dimensions," *Materials Discovery and Design By Means of Data Science and Optimal Learning*, Springer, 2018, 167 - 201. DOI: 10.1007/978-3-319-99465-9\_7
- Z. Li et al., "funcX: Federated Function as a Service for Science," in IEEE Transactions on Parallel and Distributed Systems, vol. 33, no. 12, pp. 4948-4963, 1 Dec. 2022, doi: 10.1109/TPDS.2022.3208767.
- Michael Prince, Doğa Gürsoy, Dina Sheyfer, Ryan Chard, Benoit Côté, Hannah Parraga, Barbara Frosik, Jon Tischler, and Nicholas Schwarz. 2023. Demonstrating Cross-Facility Data Processing at Scale With Laue Microdiffraction. In Proceedings of the SC '23 Workshops of The International Conference on High Performance Computing, Network, Storage, and Analysis (SC-W '23). Association for Computing Machinery, New York, NY, USA, 2133–2139. <https://doi.org/10.1145/3624062.3624613>
- D. Sheyfer, Q. Zhang et al., Phys. Rev. Lett. 125, 125504 (2020)
- Z. Liu, T. Bicer, R. Kettimuthu, D. Gursoy, F. De Carlo, and I. Foster. "TomoGAN: Low-Dose X-Ray Tomography with Generative Adversarial Networks." preprint arXiv:1902.07582 (2019)
- Z. Liu, H. Sharma, J. Park, P. Kenesei, J. Almer, R. Kettimuthu, I. Foster. arXiv <https://arxiv.org/abs/2008.08198>
- S. Veseli, N. Schwarz and C. Schmitz, J. Synchrotron Rad. 25 1574 (2018)
- Q. Zhang et al., J. Synchrotron Rad. (2021)

# REFERENCES

- <https://www.globus.org>
- <https://www.aps.anl.gov/Dynamics-and-Structure>
- <https://www.gmca.aps.anl.gov/>
- <https://acdc.alcf.anl.gov/>
- <https://git.aps.anl.gov/DM/dm-docs>
- <https://github.com/globus-gladier/gladier-xpcs>
- [https://github.com/AdvancedPhotonSource/boost\\_corr](https://github.com/AdvancedPhotonSource/boost_corr)
- <https://github.com/AdvancedPhotonSource/ptychodus>
- <https://github.com/AdvancedPhotonSource/XRF-Maps>
- <https://github.com/Linked-Liszt/laue-parallel>

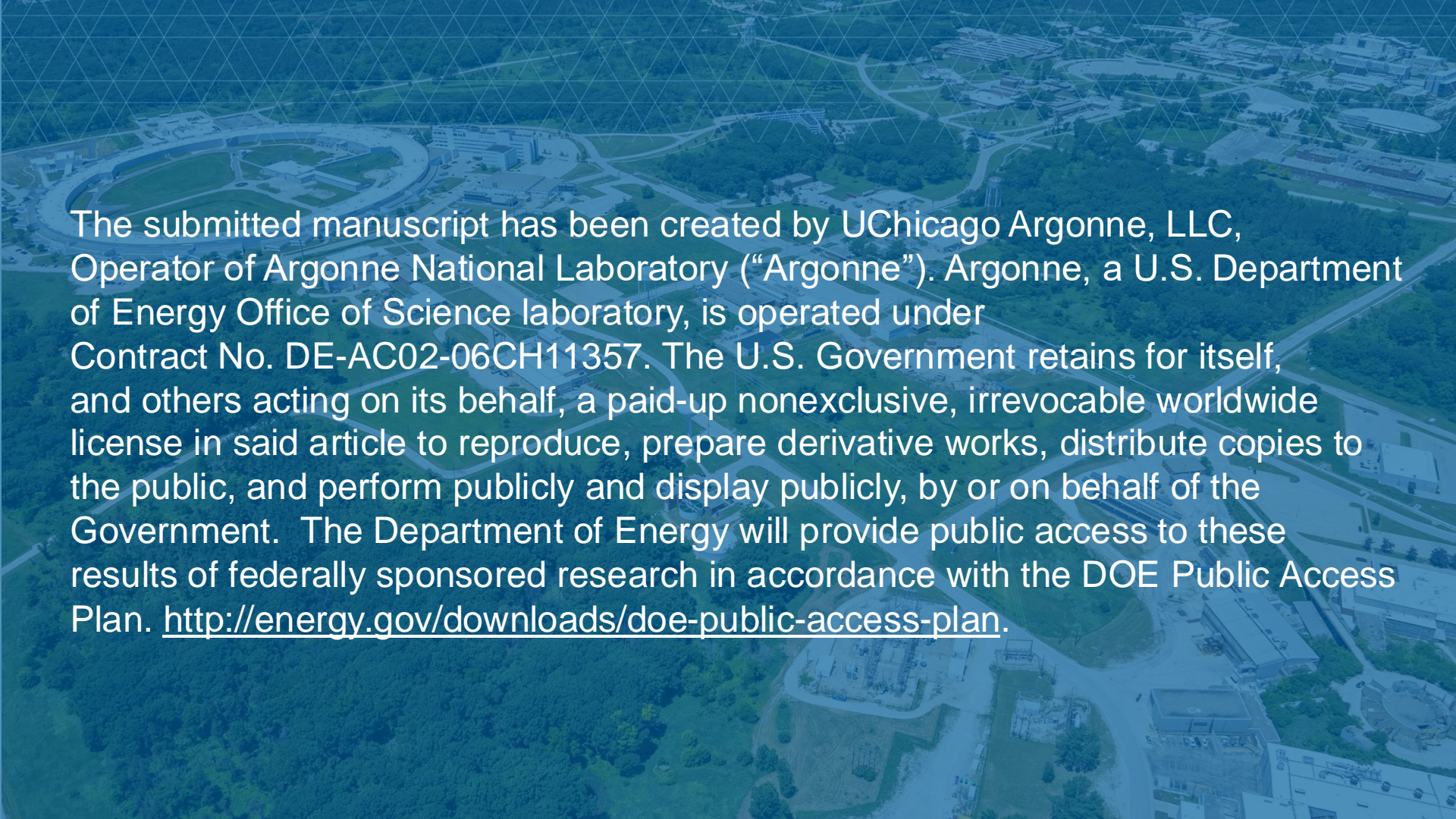


**Argonne**  
NATIONAL LABORATORY



**Advanced  
Photon Source**



An aerial photograph of the Argonne National Laboratory campus, featuring various buildings, parking lots, and green spaces. The image is overlaid with a semi-transparent blue grid pattern. A large block of white text is centered over the image, providing information about the manuscript's creation and public access policy.

The submitted manuscript has been created by UChicago Argonne, LLC, Operator of Argonne National Laboratory (“Argonne”). Argonne, a U.S. Department of Energy Office of Science laboratory, is operated under Contract No. DE-AC02-06CH11357. The U.S. Government retains for itself, and others acting on its behalf, a paid-up nonexclusive, irrevocable worldwide license in said article to reproduce, prepare derivative works, distribute copies to the public, and perform publicly and display publicly, by or on behalf of the Government. The Department of Energy will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan. <http://energy.gov/downloads/doe-public-access-plan>.