# Office of Science Data Management and Sharing Plan requirements

Basic Energy Science Advisory Committee Meeting
DOE Public Access Plan and Data Management Panel Discussion

**Alexander Hexemer,** Advanced Light Source, Lawrence Berkeley National Laboratory

**Nicholas Schwarz**, Advanced Photon Source, Argonne National Laboratory

**Stuart Campbell**, National Synchrotron Light Source II, Brookhaven National Laboratory

**Vivek Thampy**, Stanford Synchrotron Radiation Lightsource, SLAC

**Jana Thayer**, Linac Coherent Light Source, SLAC

# Advances in BES Light Sources

Computing addresses rapid data increases

## More complex experiments

- Multi-modal experiments combine data from multiple samples, techniques, and facilities

- *In situ* and *in operando* experiments require real-time feedback and autonomous control

- Spectroscopy with 1000's of scans in just a few seconds.

## Sources—orders-of-magnitude brighter

- Facility upgrades:
  - NSLS-II  • LCLS-II  • APS-U  • ALS-U • LCLS-II HE

## Detectors—orders-of-magnitude faster

- Faster readout

- Larger arrays

Analyze and reconstruct massive multimodal data volumes

Identify and classify features and patterns across modes

Merge simulation & experiment data to drive experiments and new results

Execute experiments dynamically using real-time reduction and AI/ML

BERKELEY LAB | ADVANCED LIGHT SOURCE

# Scale of the Data Volume

**In 5 years, DOE light sources are projected to**

– generate 1 exabyte of data/year

– 10s-of-petaflop/s to 1-exaflop/s, peak computing power



> **1 exabyte/y = 1.5 million movies every day**
>
> - *Analyze every frame* in near real time; guide experiments
> - Hundreds of experiment types require custom solutions
>
> **1 exaflop/s = 500,000 servers**
>
> - Fast networks (multiple Tbps)
> - Storage
> - Analysis infrastructure

# Changing Landscape for Facilities and Users

compounding the computational and data challenges

The **user community is diverse**: a wide variety of backgrounds and domains

- Varying expectations on types and scales of computing capabilities & services provided by the facilities

There is an **increasing digital divide** within the scientific user community

- Currently, few user groups have the ability to manage/process their data. This challenge will only increase.

Increased interest in **FAIR, open data,** and **data interoperability**

- The role of the facilities is unclear: facilities do not have the infrastructure in place to consistently collect, curate, archive, and disseminate data and metadata at the anticipated scale required

# What does this mean for users?



- New science opportunities

- Take advantage of the wealth of facility data to augment your own science

- Develop & test new algorithms on open and shared data

- Train ML models on large shared data sets

- Use existing ML models to accelerate knowledge extraction from data
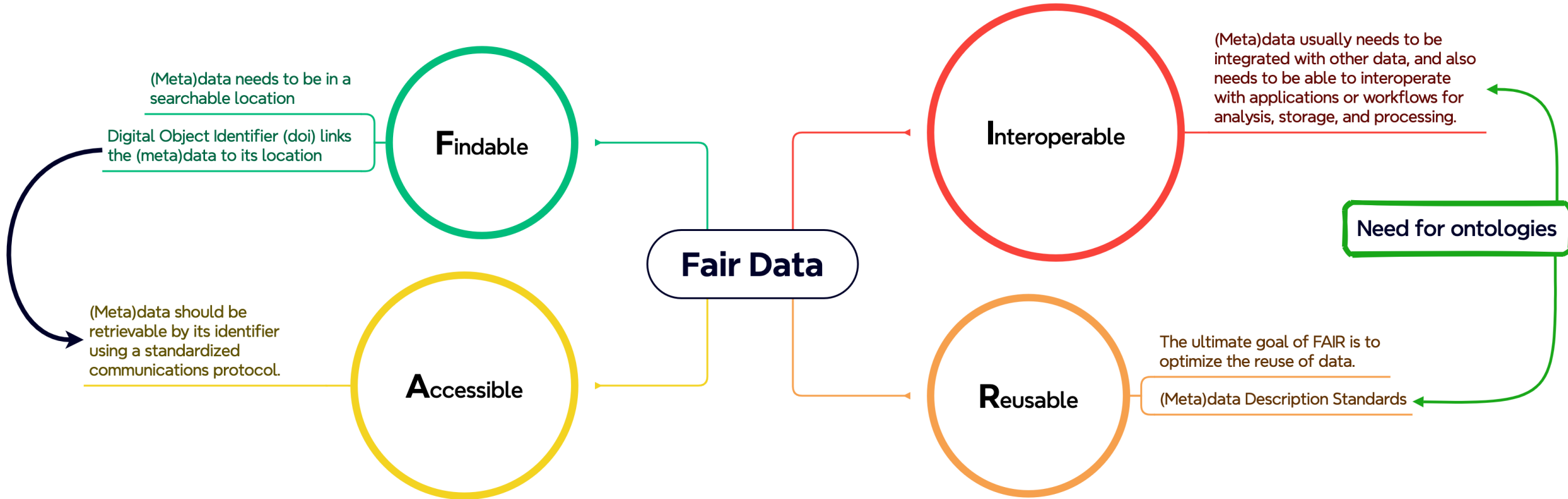
- Make data FAIR

BERKELEY LAB | ADVANCED LIGHT SOURCE

# FAIR Data

**Findable: Data and supplementary materials have sufficiently rich metadata and a unique and persistent identifier.**

**Accessible: (Meta)data are understandable to humans and machines. Data is deposited in a trusted repository.**

**Interoperable: (Meta)data use a formal, accessible, shared, and broadly applicable language for knowledge representation.**
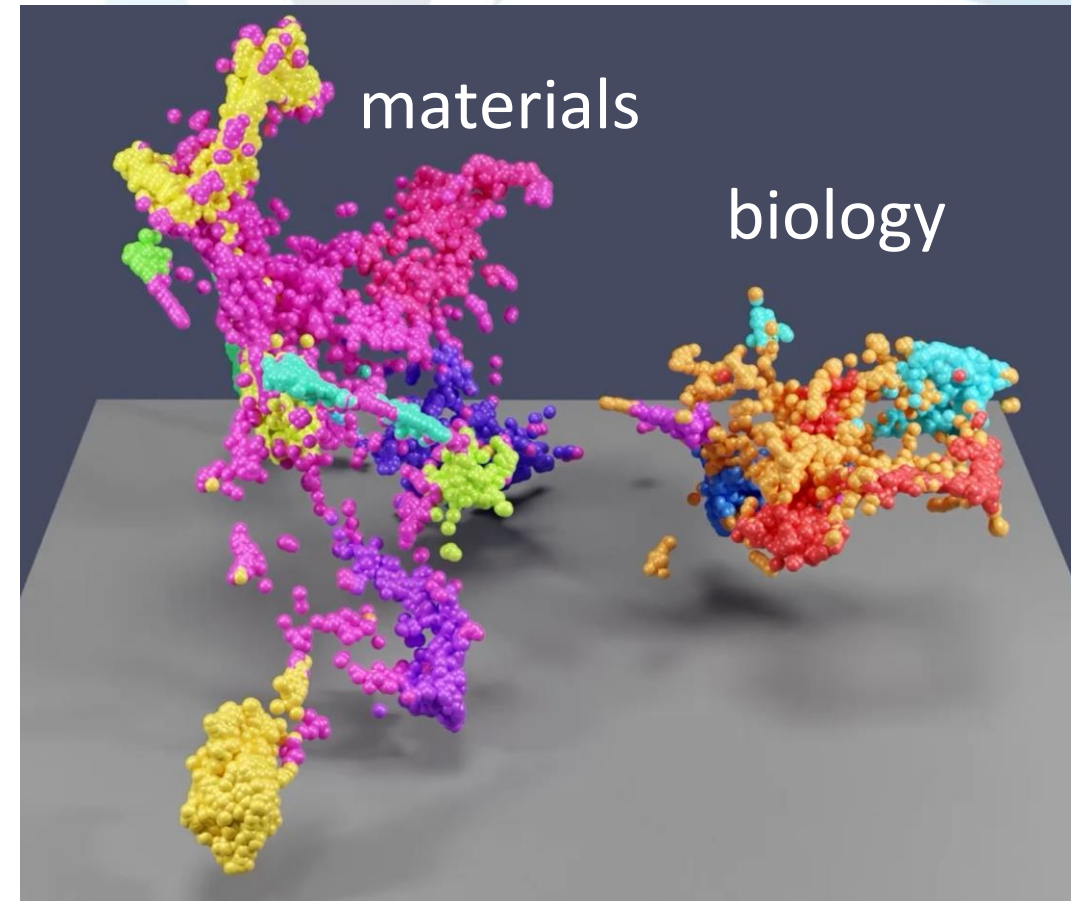
**Reusable: Data and collections have a clear usage licenses and provide accurate information on provenance.**

(Meta)data needs to be in a searchable location

Digital Object Identifier (doi) links the (meta)data to its location

**Findable**

(Meta)data usually needs to be integrated with other data, and also needs to be able to interoperate with applications or workflows for analysis, storage, and processing.

**Interoperable**

**Fair Data**

Need for ontologies

(Meta)data should be retrievable by its identifier using a standardized communications protocol.

**Accessible**

The ultimate goal of FAIR is to optimize the reuse of data.

**Reusable**

(Meta)data Description Standards

# Challenges with FAIR data
*(a scientific user facility perspective)*

- The **ontology** in each field must be well-structured and precise to ensure clear communication and data interoperability.

- Not all the **metadata** of an experiment is available to the user facilities (e.g., material synthesis)

- Data sets can be *very* large & difficult to handle, stored and served

- (Meta)data of a single study can be spread across **multiple facilities**

- What are the implications of **data deletion** (raw or derived data)?

- **Authentication** and authorization challenges across facilities

- Increased **complexity** of working with data



materials

biology

Clustering of ALS publications using a LLM and UMAP to show the wealth of information and diversity of research areas

BERKELEY LAB | ADVANCED LIGHT SOURCE

# Opportunities using FAIR data
## *(a scientific user facility perspective)*

- Data reusability

- Reproducibility of experimental results & analysis

- Development and testing of new algorithms on well-described data

- Common ontologies allows for better cross-facility collaboration

- Seamless integration of data with (HPC) compute resources

- Improved training data quality for ML models with AI ready data

- Opportunity to share data and trained ML models

- Opportunities for Unsupervised and Semi-Supervised Learning

- Using generative AI to create data sets specific to a given experiment

GIWAXS
*Experimental data*

GIWAXS
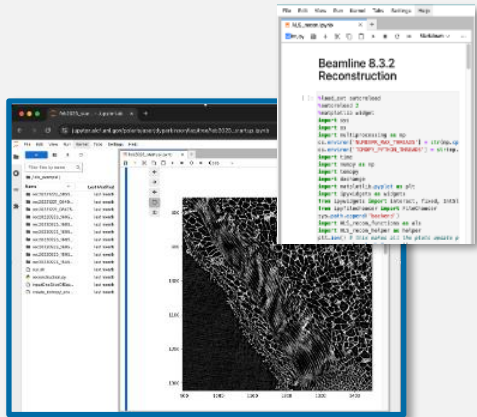*AI-generated data*



**prompt:**
"GIWAXS data with rings and peaks"

BERKELEY LAB | ADVANCED LIGHT SOURCE

# Data Portals and Access across Light Sources

# Working Together Across BES & ASCR Facilities

**ALS** Reconstructions: **ALCF & NERSC**

**APS** On-Demand Workflows: **ALCF**

**LCLS** Data Analytics: **NERSC, ALCF, OLCF**

**NSLS-II** Prototype Pipeline: **ALCF**
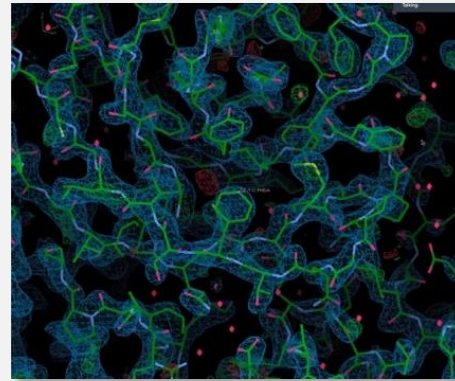
**SNS/HFIR** Data Processing: **OLCF**

- Globus Transfer between ALS, NERSC and ALCF
- Tomography reconstruction on NERSC and using Globus Compute on ALCF
- Results are transferred back to ALS or NERSC

- ALCF Polaris system: continuous on-demand data processing
- Operational workflows for over 10 techniques
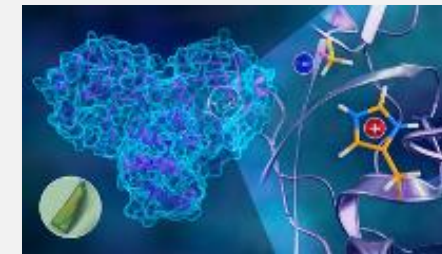- Globus tools provide workflow and web portal services

- Automated SFX pipeline at NERSC; demonstrated during COVID-19 LCLS experiments
- AI training on OLCF Summit+
- Prototype ptychography workflow demonstrated on ALCF Polaris

- Export to file, transfer using Globus
- Processing XPCS using Jupyter Notebooks on ALCF using Globus Flows
- Next, integrate with bluesky/tiled for data access and output

- Web-based data platform, integration with OLCF
- Neutron users perform calculations as part of an analysis workflow
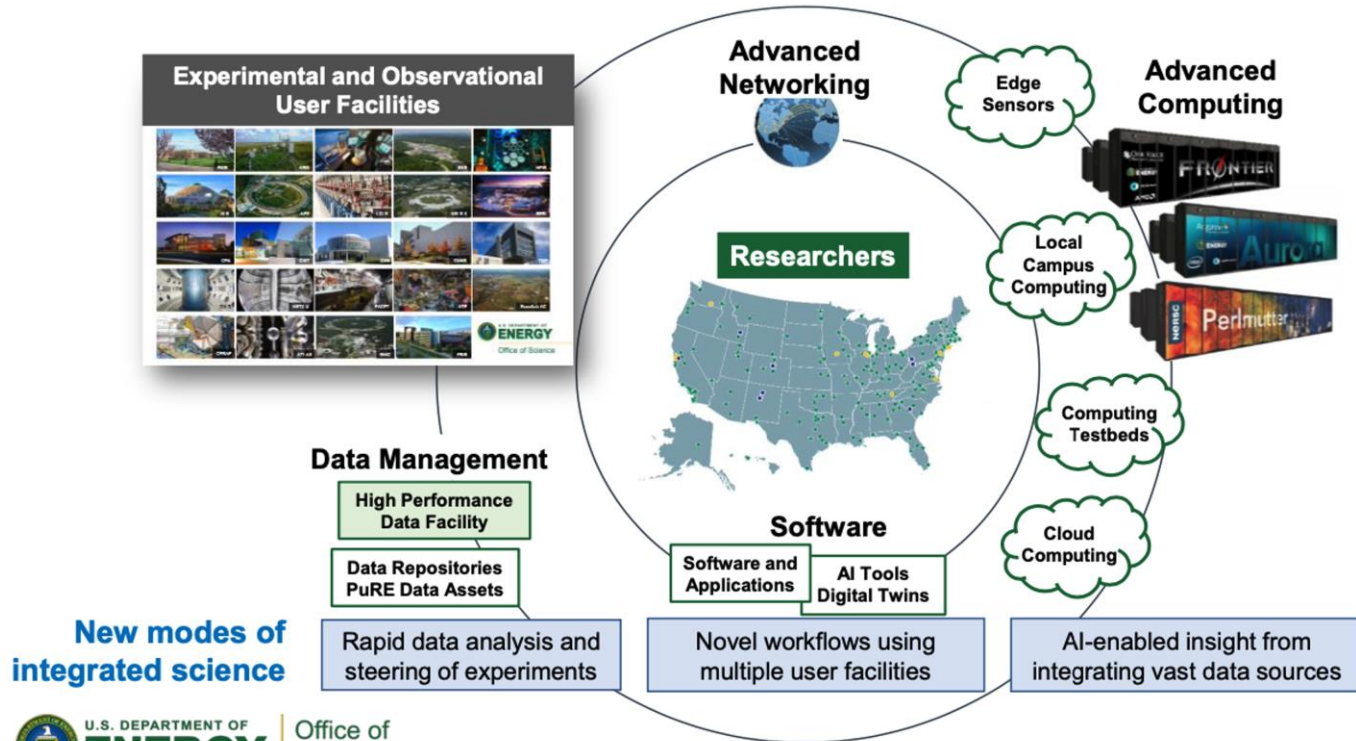- Intra-experiment training to predict protonation states of active protein sites for neutron MX

BERKELEY LAB | ADVANCED LIGHT SOURCE

# Integrated Research Infrastructure (IRI)
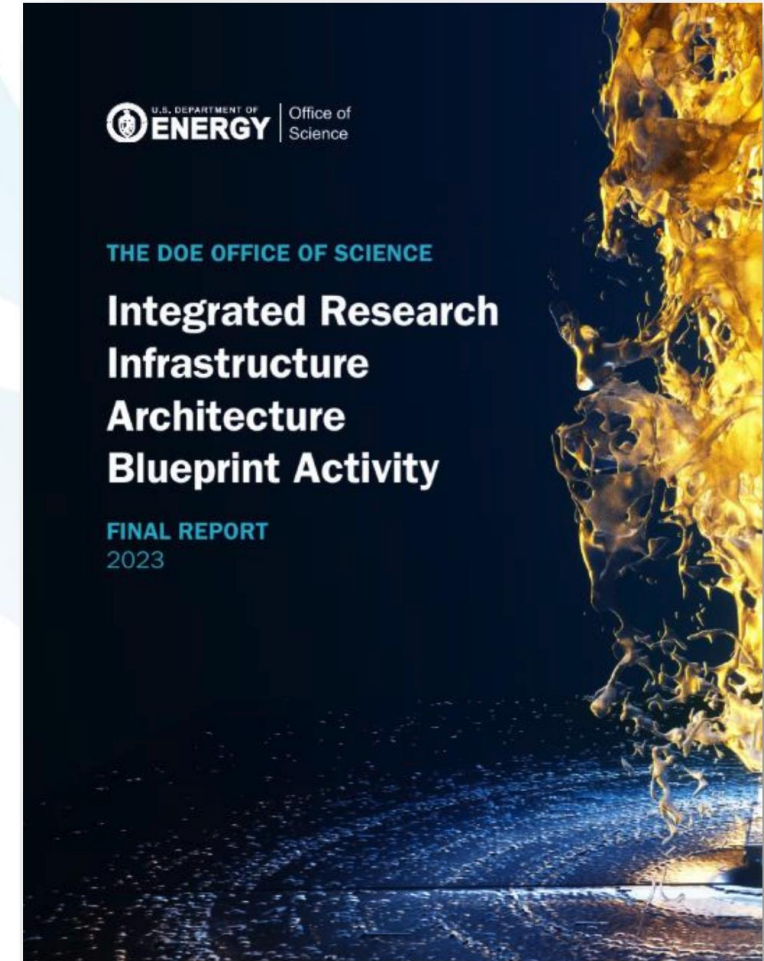
## The BES Light Sources are one of the initial *Pathfinder* projects

- The ALS, APS and NSLS-II are a single combined pathfinder
- Target beamlines and techniques are being finalizes

# Summary of High-Priority Computing Developments across Light Sources

1. **Data management and workflow tools** that integrate beamline instruments with computing & storage, for use during experiment, as well as facile user access for post-experiment analysis.

2. **Real-time data analysis capabilities** to significantly reduce data volumes and provide feedback during experiments, improving data quality and driving autonomous experiments.

3. **On-demand utilization of computing environments** to enable quasi-real-time data processing

4. **Data storage and archival resources** to house the increasing amounts of valuable scientific data produced by the BES Light Sources in a common, smart data portal.

5. **Easy-to-use** solution to provide an inclusive environment for researchers at SUFs

BERKELEY LAB | ADVANCED LIGHT SOURCE