# Evolving role of HDF5 at the upgraded Advanced Photon Source

**Tejas Guruswamy**
Detectors Group, Advanced Photon Source
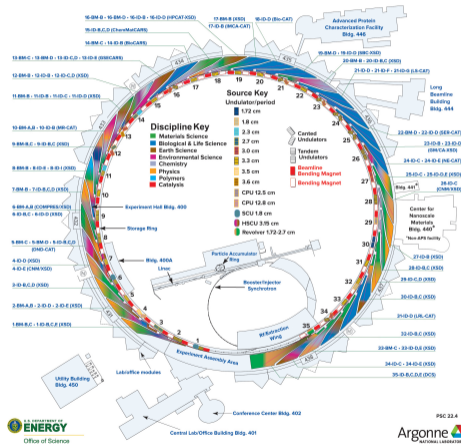
APS-U Beamline Data Pipelines project

HDF5 User Group Meeting 2023
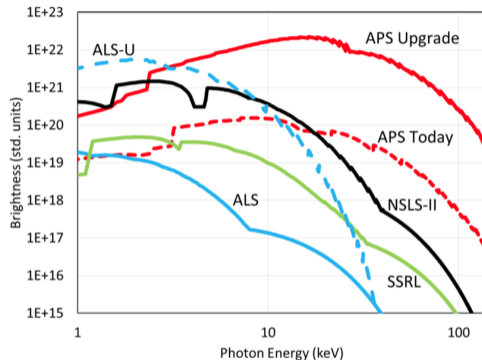
# The Advanced Photon Source

- The **Advanced Photon Source** is a hard X-ray (6 to 150 keV photons) synchrotron light source

- Electrons accelerated to 7 GeV (99.999999% of $c$) generate intense beams of photons for imaging, diffraction, spectroscopy

- Around 70 parallel experimental stations each with unique specialties operating 24/7, 9 months a year

- $10^9$ (one billion) times brighter than a medical X-ray

- >5,000 worldwide users each year from physics, materials science, chemistry, biology, defence, nuclear safety, microfabrication, electronics, pharma . . .



ARGONNE NATIONAL LABORATORY 400-AREA FACILITIES
ADVANCED PHOTON SOURCE
(Beamlines, Disciplines, and Sources)
ADVANCED PROTEIN CHARACTERIZATION FACILITY
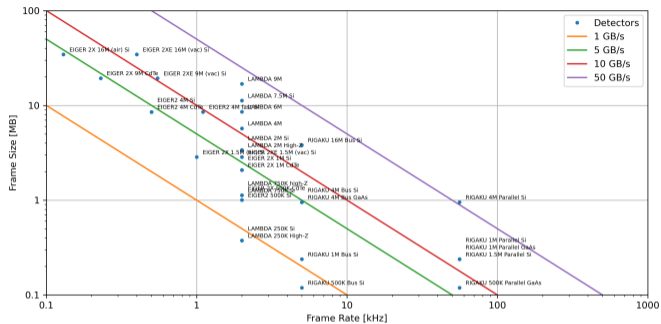CENTER FOR NANOSCALE MATERIALS

# The APS Upgrade (Apr 2023 - Apr 2024)

- Facility undergoing major upgrade, replacing all storage ring magnets
- **500 times increase** in coherence and X-ray brightness at most energies
- Enables better measurements of dilute samples, faster scanning, nanoscale focusing
- Entire beamlines, especially detectors, being upgraded to take advantage of increased photon flux
- New network, data storage and compute infrastructure throughout APS

# Increased data rate & size

- Order-of-magnitude increase in data generated (**<10 PB/year to >100 PB/year**)

- Updated science goals including real-time feedback, use of ALCF compute (exascale Aurora/Polaris), FAIR data policies

- Overhaul of **batch and stream** processing, analysis, visualization for most scientific techniques (*Beamline Data Pipelines* project)
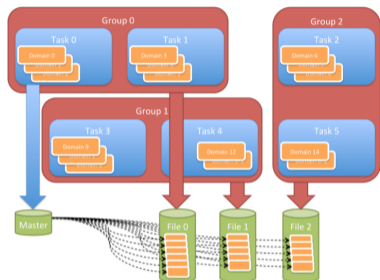
# Data format and management tools

- HDF5 preferred as on-disk data format for all new/upgraded data workflows
- Combine all data and metadata in stable format for archiving and offline (re-)analysis
- Leveraging existing solutions where possible
  - **EPICS areaDetector** plugin(s) for direct to HDF5 data acquisition and network streaming of structured data with same compression algorithms as HDF5
  - **Data Exchange** [https://dxfile.readthedocs.io/] and **NeXus** [http://www.nexusformat.org/] for public, documented schema matching needs of technique
  - **bluesky** experiment orchestration and metadata collection, including **tiled** UI for slicing and exporting HDF5-backed data
  - **globus** cross-institution data management and HPC workflow execution

Argonne NATIONAL LABORATORY

U.S. DEPARTMENT OF ENERGY

# Improving parallel and GPU use



From Mark Miller, 2017
https://www.hdfgroup.org/2017/03/
mif-parallel-io-with-hdf5/

- Parallel read and (sometimes) write important to keep up with data rates
- *External multiple file links, MPI, or other parallelism?*
- Increasing use of GPU-based analysis; experimenting with NVIDIA Rapids libraries, limited success so far
- *Accelerate HDF5 loading to GPU, network streams?*
- Some techniques, e.g. XPCS, require data layouts efficient on FPGA, CPU and GPU
- *Sparse and/or optimized data layouts for diverse architectures?*
- *Enable AI and data mining to efficiently navigate archived data?*