# Real-Time Data Analysis at NERSC:
## a Trial Run of Nascent Exascale Experimental Data Analysis
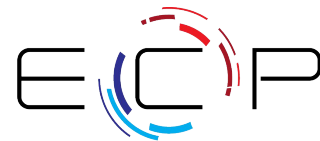
Johannes Blaschke[1], Aaron Brewster[1], Daniel Paley[1], Derek Mendez[1], Nicholas Sauter[1], Wilko Kroeger[2], Murali Shankar[2], Bjoern Enders[1], Deborah Bard[1]
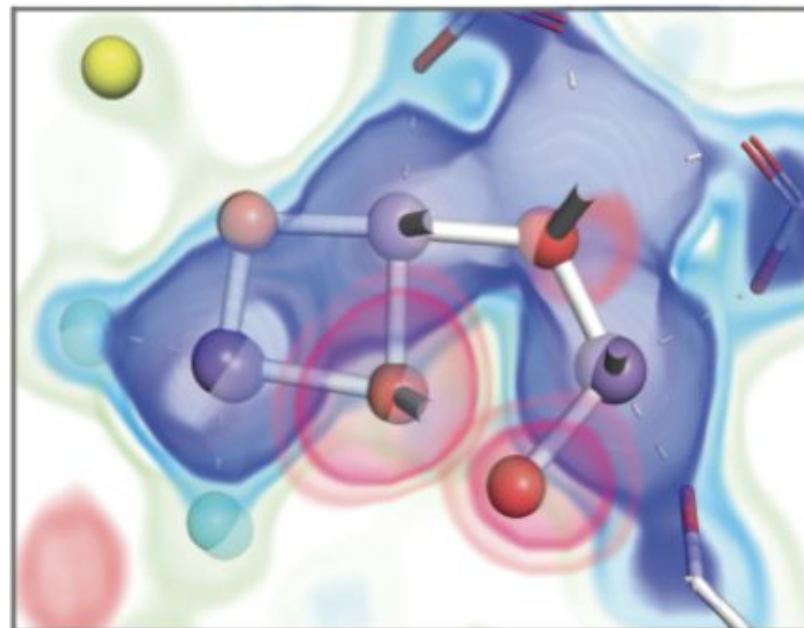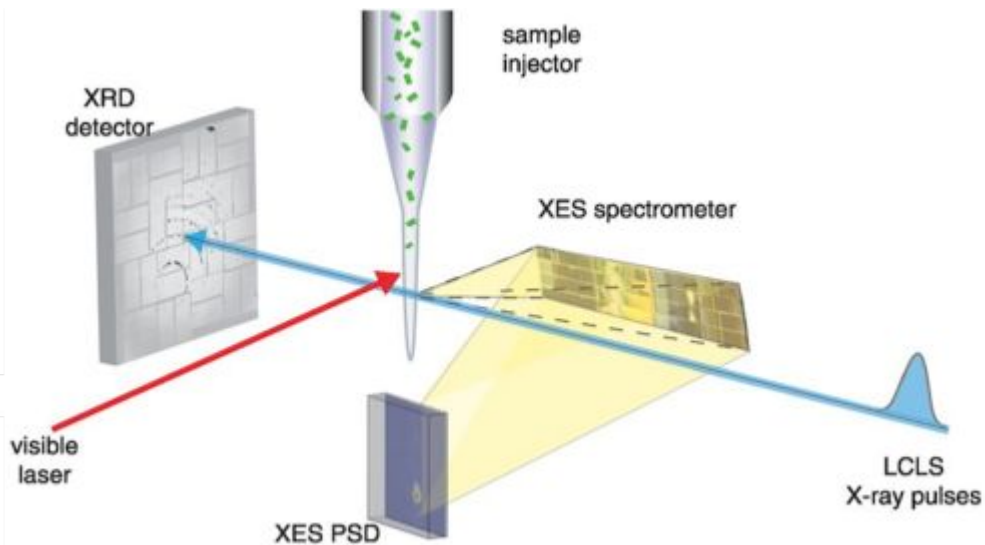
[1]Lawrence Berkeley National Lab
[2]SLAC National Accelerator Lab

CUG, May 5, 2021

# Serial Crystallography
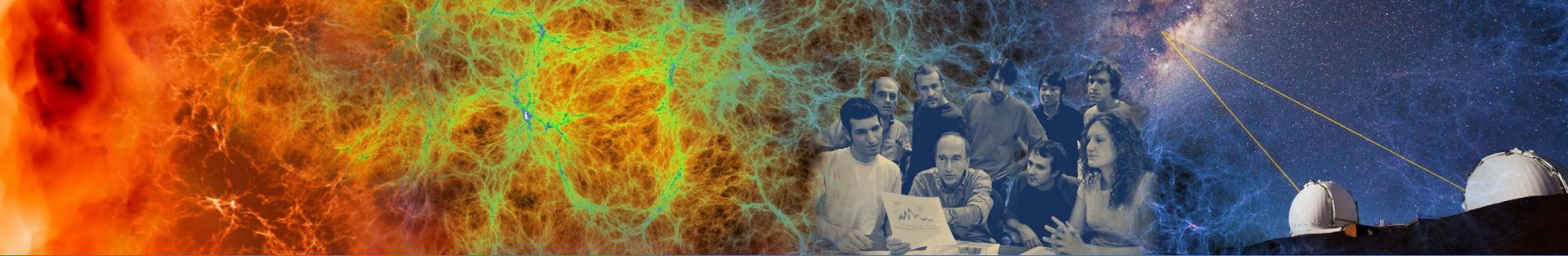


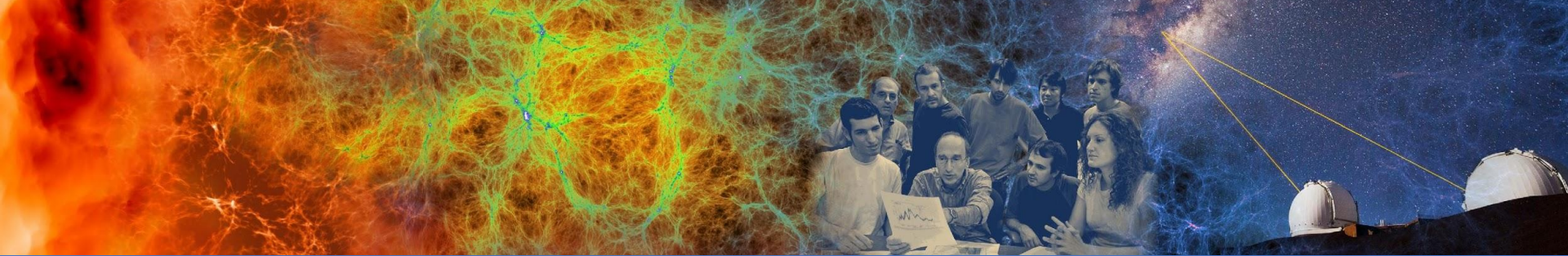PSII 2Fo-Fc and Fo-Fc at 2.07 Å, S3 state

# When should I move onto the Next Sample?

- Beamtime is scarce!

- Critical live feedback:
  - Does the beam hit the sample?
  - Do we see crystals?
  - Does the data make sense?
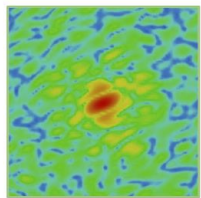  - What is the quality of the data?

- Can I move on to the next sample?

# Experimentalists Are In The Driver's Seat
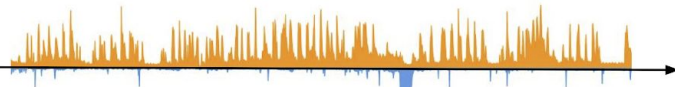Live Data Analysis for Experiments in 2020, and Beyond!

# Deploying CCTBX at NERSC

LCLS

Data Acquisition → Spinning Disk

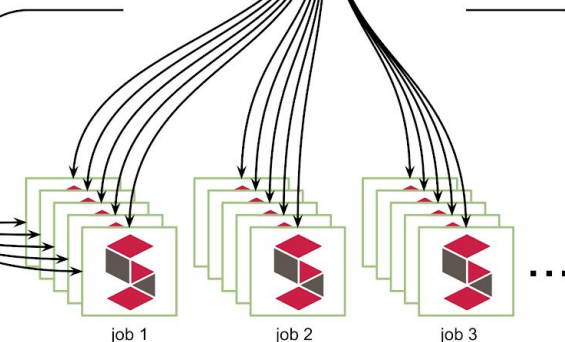XRootD ~ 15 TB/day

Data Transfer Nodes

NERSC

SCRATCH, CFS, DataWarp

Users interact with data analysis in real-time

Workflow Coordination

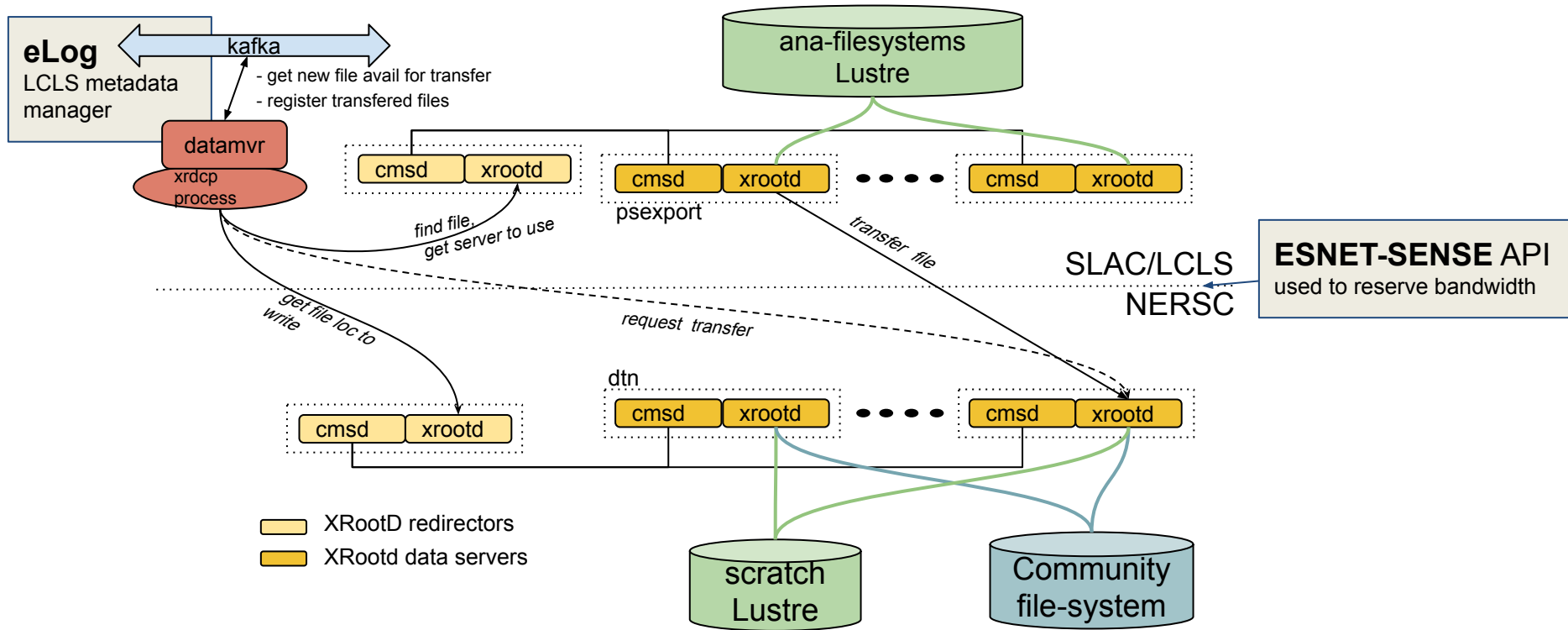MySQL

Spin

Cori Compute Nodes

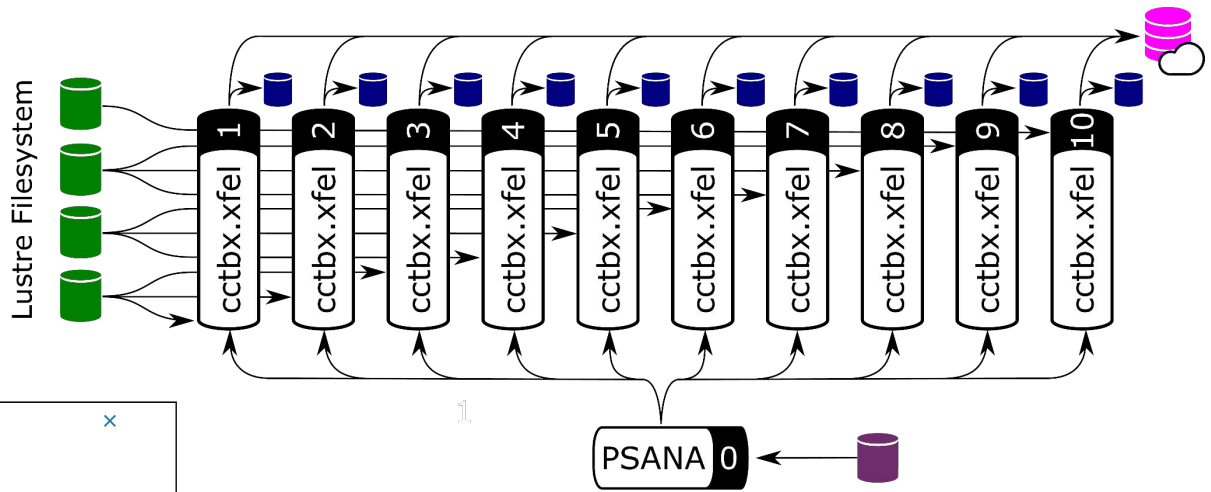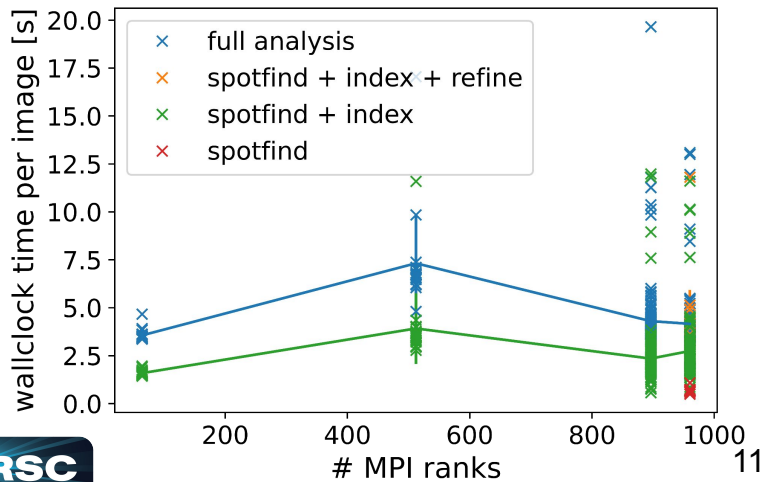job 1    job 2    job 3    ...

LCLS    NERSC

# Data Movement XRootD clusters

# Data Analysis

- Data analysis follows sequential stages:
  - spotfinding
  - indexing Bragg spots
  - model refinement
  - integrating Bragg spots

# How's the Computation Weather Today?

- Computational Weatherplot:
  - Each line shows work done by one MPI rank
  - There is no "*the* `cctbx.xfel` workload"

Start-Up and I/O (PSANA)
Spot Detection (DIALS)
Indexing (DIALS)
Refinement (DIALS)
Integrating (DIALS)

# How's the Computation Weather Today?

- Computational Weatherplot:
  - Each line shows work done by one MPI rank
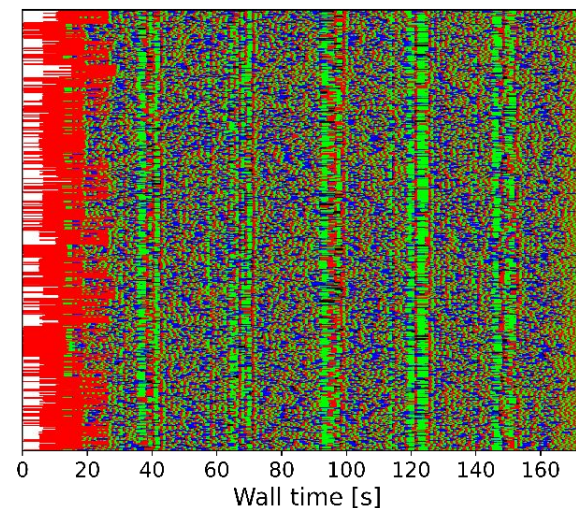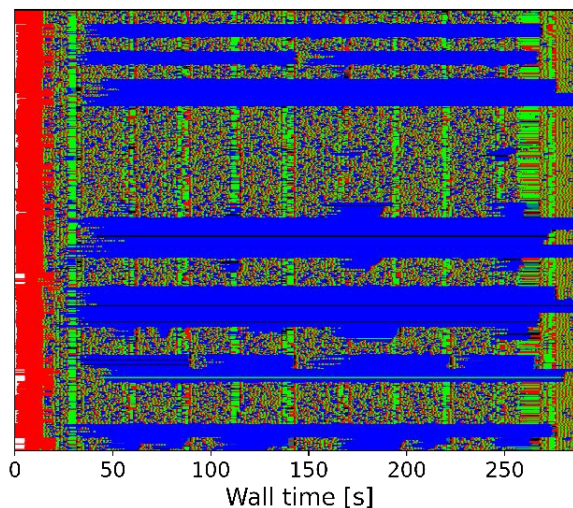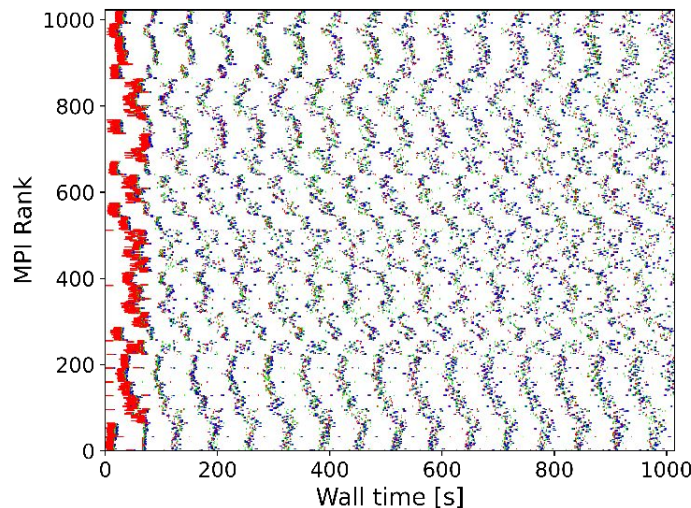  - There is no "*the* `cctbx.xfel` workload"



Start-Up and I/O (PSANA)
Spot Detection (DIALS)
Indexing (DIALS)
Refinement (DIALS)
Integrating (DIALS)

# How's the Computation Weather Today?

- Computational Weatherplot:
  - Each line shows work done by one MPI rank
  - There is no "*the* `cctbx.xfel` workload"



Start-Up and I/O (PSANA)
Spot Detection (DIALS)
Indexing (DIALS)
Refinement (DIALS)
Integrating (DIALS)



Issue: MPI-Communication Bound

Issue: I/O contention

# How's the Computation Weather Today?

- Computational Weatherplot:
  - Each line shows work done by one MPI rank
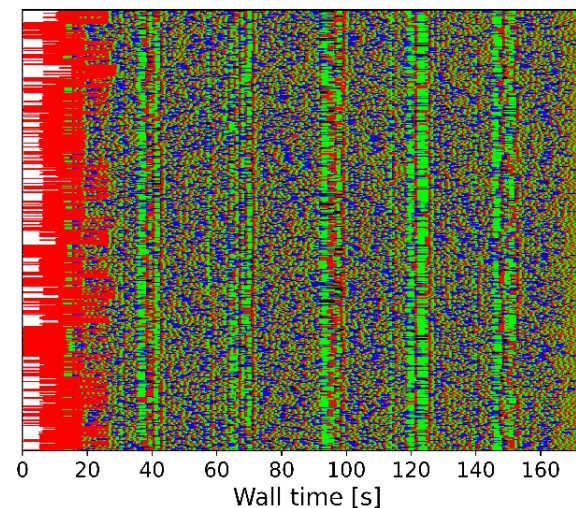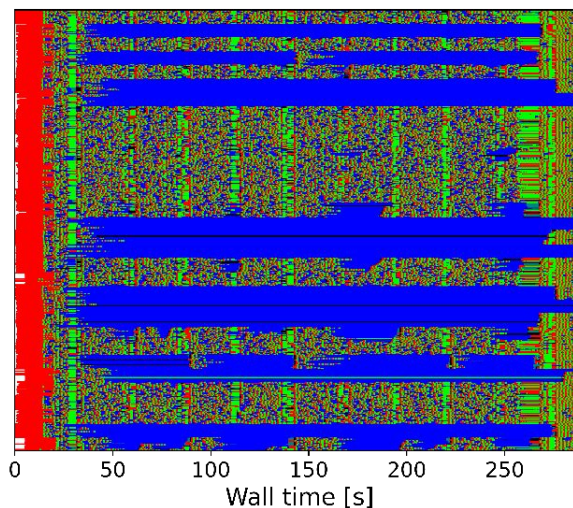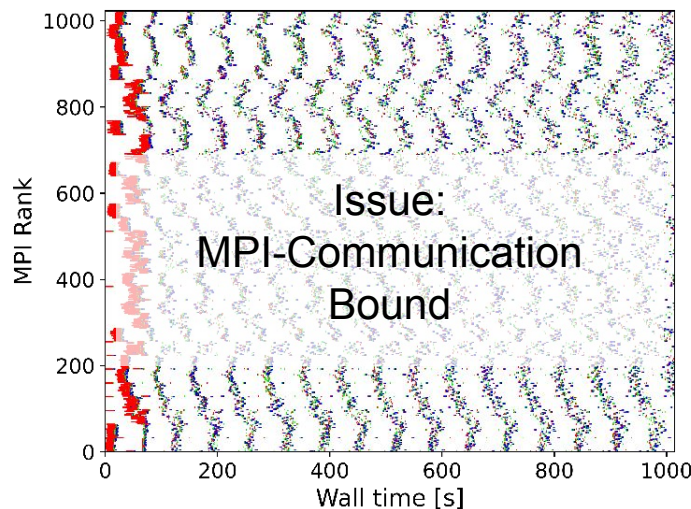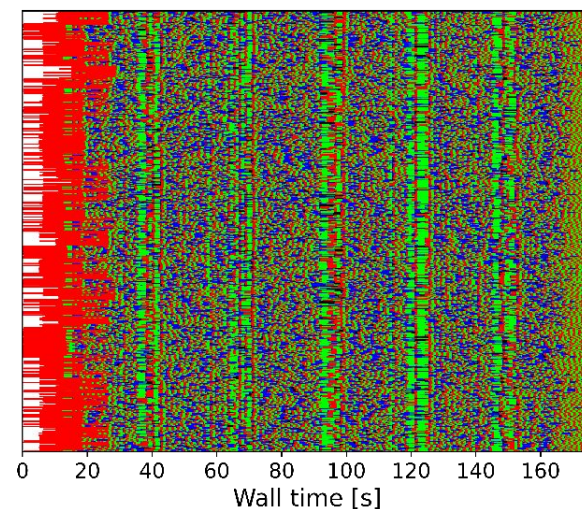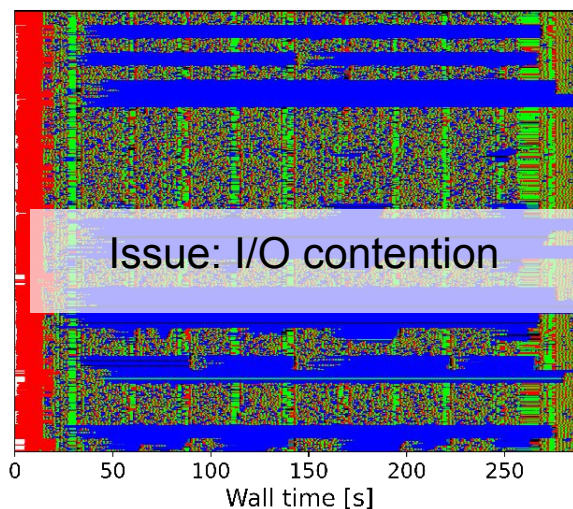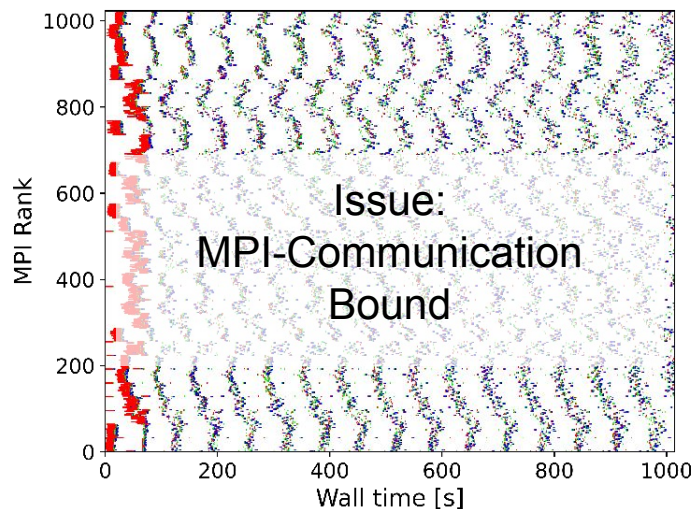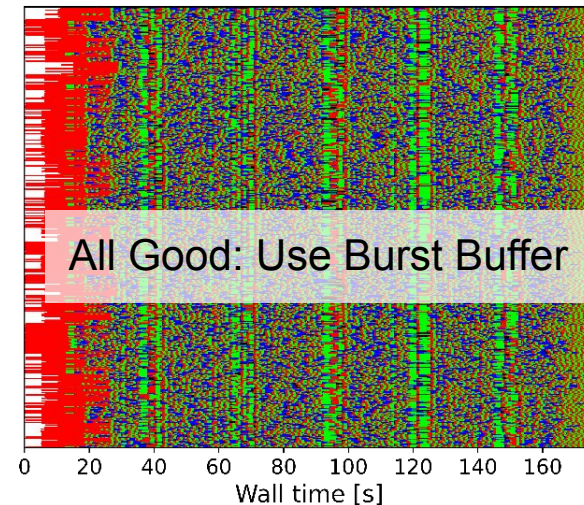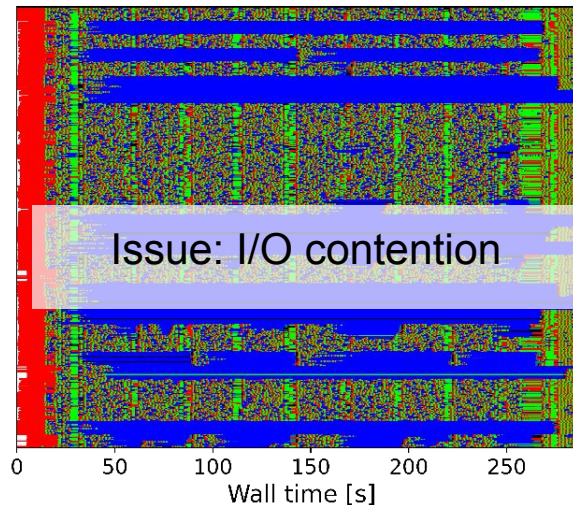  - There is no "*the* `cctbx.xfel` workload"



Start-Up and I/O (PSANA)
Spot Detection (DIALS)
Indexing (DIALS)
Refinement (DIALS)
Integrating (DIALS)

Issue:
MPI-Communication
Bound

Issue: I/O contention

All Good: Use Burst Buffer

# Workflow Coordination using NERSC Spin

- Home-grown workflow manager `cctbx.xfel`
  - mySQL database hosted on Spin (NERSC microservice platform)
  - Each worker commits progress to DB
  - `cctbx.xfel` determines new analysis runs and "assembles" jobs (input files, job scripts, ...)
  - `cctbx.xfel` monitors slurm and DB, reporting live progress



datamvr

`cctbx.xfel` checks status of file transfers

slurm connector submits new jobs (**multiple** users can submit)

**multiple** users query database

**each rank** commits status to database:
1. spot-finding rate
2. indexing rate
3. crystal parameters

NX/Login Nodes          Spin          Compute Nodes

# XFEL as a Proxy for HPC Data Analysis
Why is XFEL relevant to other "Data Analysis for Science" projects?

# Challenge 1: Urgent Computing Resources

■ : collecting/transferring
■ : processing
▯ : no live results

Data Collected (run number) →

Reservation with 3 nodes, each run takes 1 unit of time to collect and process ⇒ can only process 3 yellow squares at once

Wallclock time →

NeRSC

BERKELEY LAB
Bringing Science Solutions to the World

U.S. DEPARTMENT OF ENERGY | Office of Science

# Challenge 1: Urgent Computing Resources

■ : collecting/transferring
■ : processing
▒ : no live results

Data Collected (run number) →

Reservation with 3 nodes, each run takes 1 unit of time to collect and process ⇒ can only process 3 yellow squares at once

Wallclock time →

NeRSC

BERKELEY LAB
Bringing Science Solutions to the World

U.S. DEPARTMENT OF ENERGY | Office of Science

# Challenge 1: Urgent Computing Resources

: collecting/transferring

: processing

: no live results

Reservation with 3 nodes, each run takes 1 unit of time to collect and process ⇒ can only process 3 yellow squares at once

Data Collected (run number) →

Batch reprocessing
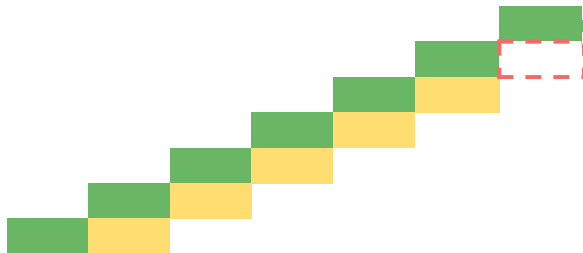
Wallclock time →

NeRSC

BERKELEY LAB
Bringing Science Solutions to the World

U.S. DEPARTMENT OF ENERGY | Office of Science

# Challenge 1: Urgent Computing Resources



: collecting/transferring

: processing

: no live results

Reservation with 3 nodes, each run takes 1 unit of time to collect and process ⇒ can only process 3 yellow squares at once

Falling behind on live fast-feedback

Batch reprocessing

Batch reprocessing

Data Collected (run number) →

Wallclock time →

NeRSC

BERKELEY LAB
Bringing Science Solutions to the World

U.S. DEPARTMENT OF ENERGY | Office of Science

# Challenge 1: Urgent Computing Resources



- **Reservation:**
  - ○ 32 - 64 Haswell nodes for live data processing
  - ○ Can be used for preemptible jobs in the future (avoid idle nodes)

- **Realtime QOS:**
  - ○ Flexibly add up to 20 Haswell nodes for reprocessing

# Challenge 2: High-Speed Data (Network and I/O)

- In data analysis workflows, file systems and network can become bottlenecks

- I/O Optimization:
  - Optimize python logger for high-frequency parallel I/O
  - Write logs to Burst Buffer

- Experience:
  - Transfers ran smoothly, can switch redirect destination
  - **FS performance limited the transfer rate**

- Improvements:
  - Use SSD storage at LCLS to speed up transfers
  - Improve write performance at NERSC
  - Allow users to initiate the remote transfers
  - Better monitoring and alerting

# Challenge 2: High-Speed Data (Network and I/O)

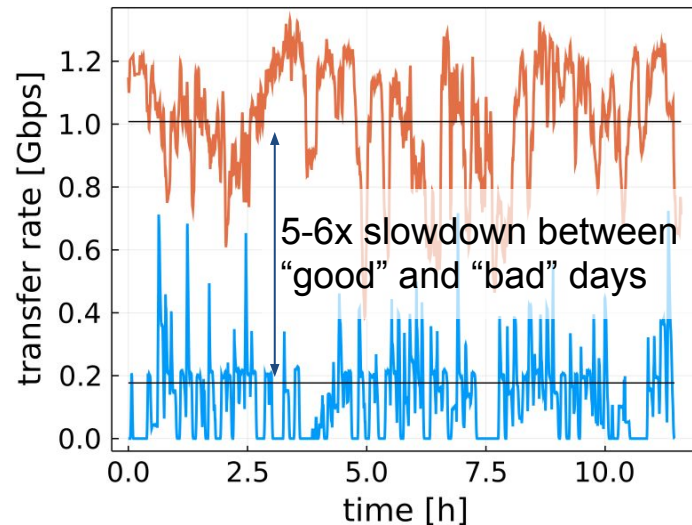- In data analysis workflows, file systems and network can become bottlenecks

- I/O Optimization:
  - Optimize python logger for high-frequency parallel I/O
  - Write logs to Burst Buffer

- Experience:
  - Transfers ran smoothly, can switch redirect destination
  - **FS performance limited the transfer rate**

- Improvements:
  - Use SSD storage at LCLS to speed up transfers
  - Improve write performance at NERSC
  - Allow users to initiate the remote transfers
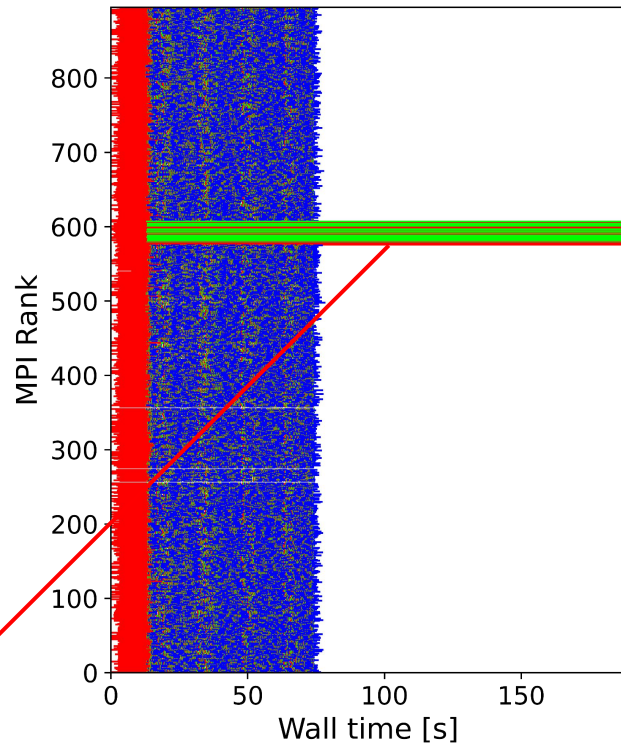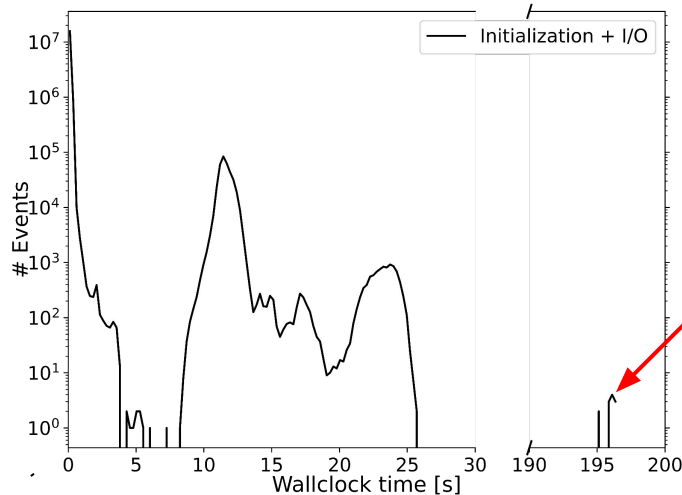  - Better monitoring and alerting



5-6x slowdown between "good" and "bad" days

# Challenge 3: Realtime Monitoring and Workflow Coordination

- Need to identify and deal with variable performance (e.g. rank getting "stuck" on I/O)

- Weatherplots good for identifying load imbalance

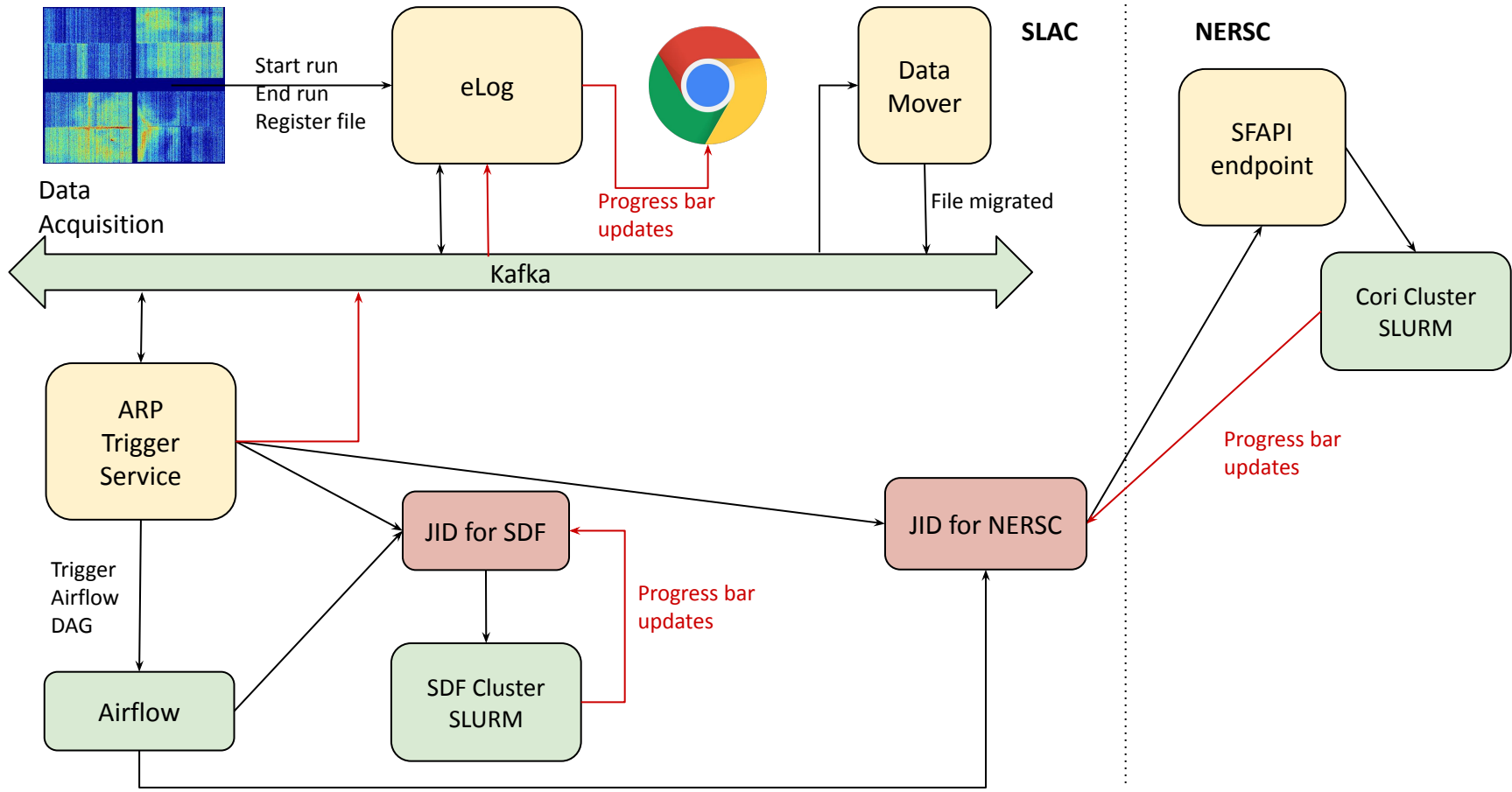- Ongoing research: How to integrate with workflow manager? How to automate?

# Looking Forward:
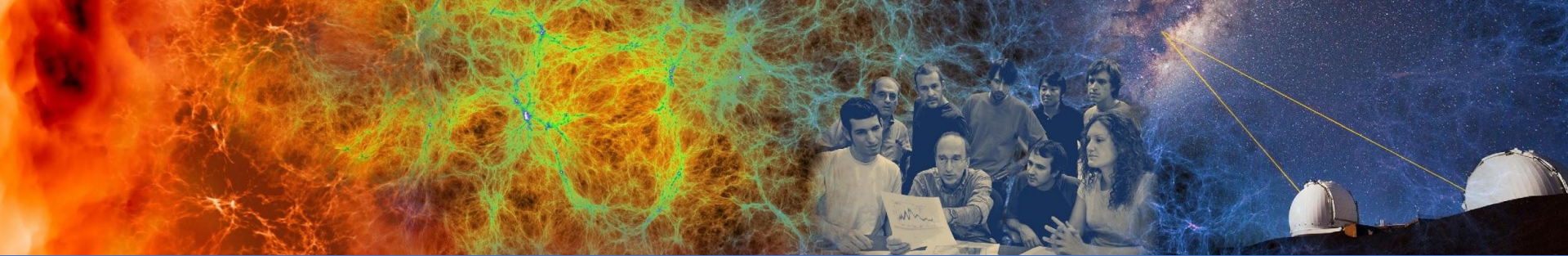# (Standardized) Facility APIs

# What can an API do?

> **Vision: all NERSC interactions are callable;**
> **backend tools assist large or complex operations.**

**Endpoints prototyped or in prep:**

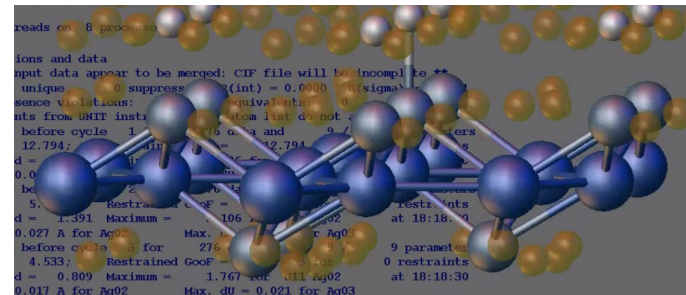| | |
|---|---|
| `/account` | data about the user's projects, roles, groups and usage information. |
| `/compute` | run batch jobs, query job and queue statuses on compute resources. |
| `/task` | get info about asynchronous tasks (eg. from `/compute` or `/storage`). |
| `/status` | query the status of NERSC component system health |
| `/storage` | move data with Globus or between NERSC storage tiers |
| `/reservations` | submit and manage future compute reservations (in prep) |
| `/utilities` | traverse the filesystem, upload and download small files, and execute commands on NERSC systems |

# Conclusion

# Successful Realtime Data Analysis at NERSC



- Live Feedback:
  - **10 mins from end of run to the molecular structures**
    - Enable real-time feedback to beamline staff
  - No babysitting from NERSC staff needed

- XFEL Flexes the following "HPC Muscles":
  - Urgent Computing Resources
  - High-Speed Data (Network and I/O)
  - Realtime Monitoring and Workflow Coordination

- Beamtime is scarce! Fast feedback is critical! [github.com/cctbx/cctbx_project](github.com/cctbx/cctbx_project)
  - G. Winter *et al*. DIALS: implementation and evaluation of a new integration package. *Acta Crystallogr D Struct Biol* **74**, 85-97 (2018)

Structure of Tethrene
(determined during LV95 beamtime)