


# The LCLStream: Multi-Institutional Data Analysis in Heterogeneous Computing Environments

---

Valerio Mariani / LCLS Data Analysis Department  
SLAC National Accelerator Laboratory

# LCLStreamer (Summit Plus)



LEADERSHIP  
COMPUTING  
FACILITY

ABOUT OLCF ▾

OLCF RESOURCES ▾

R&D ACTIVITIES ▾

SCIENCE AT OLCF ▾

FOR USERS ▾

COMMUNITY ▾



# SUMMITPLUS

[HOME](#) / [SUMMIT PLUS](#)

## Summit begins a new era

Five years after its debut as the fastest supercomputer in the world, Summit remains a powerful and reliable instrument for scientific discoveries in artificial intelligence, energy, climate, health, and other areas with a direct impact on national security and global welfare. The Department of Energy is extending Summit operations through October 2024, enabling researchers to pursue projects on one of the world's leading AI-enabled open science supercomputing platforms.

### KEY LINKS

- [Submit a proposal for SummitPLUS](#)
- [2023 Notable System Changes – OLCF User Documentation](#)

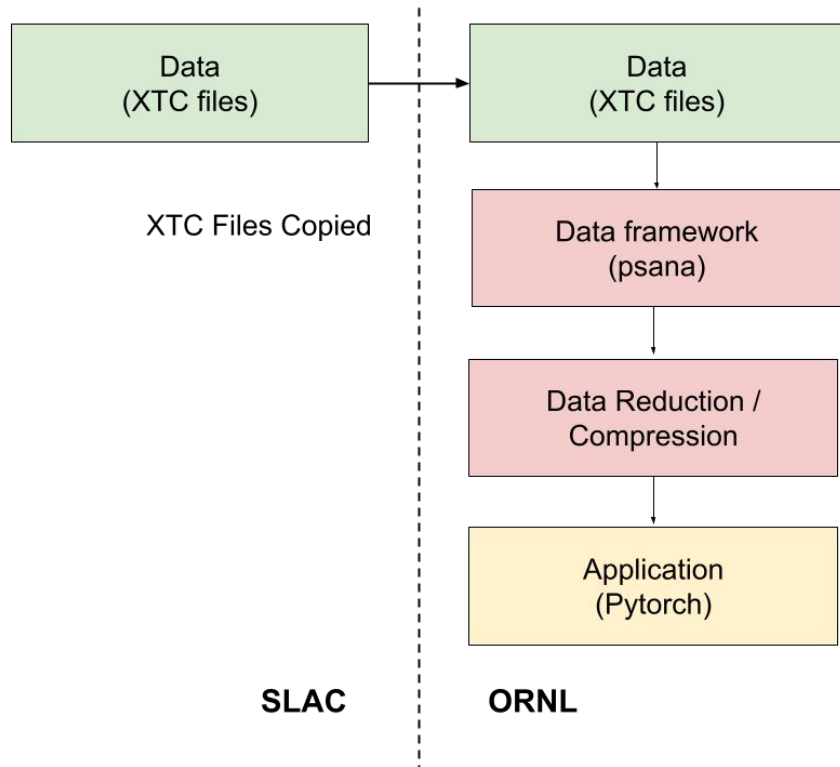
LibreOffice

▲ TOP

# LCLStream (Summit Plus)

## Data Streaming to Remote Facilities: Copying Files

- XTC files are huge
  - All of LCLS's data in raw format
- Psana needed at remote location (might be difficult to set up and run, different architecture):
  - To read the data
  - To interpret the data (Calibration / preprocessing)

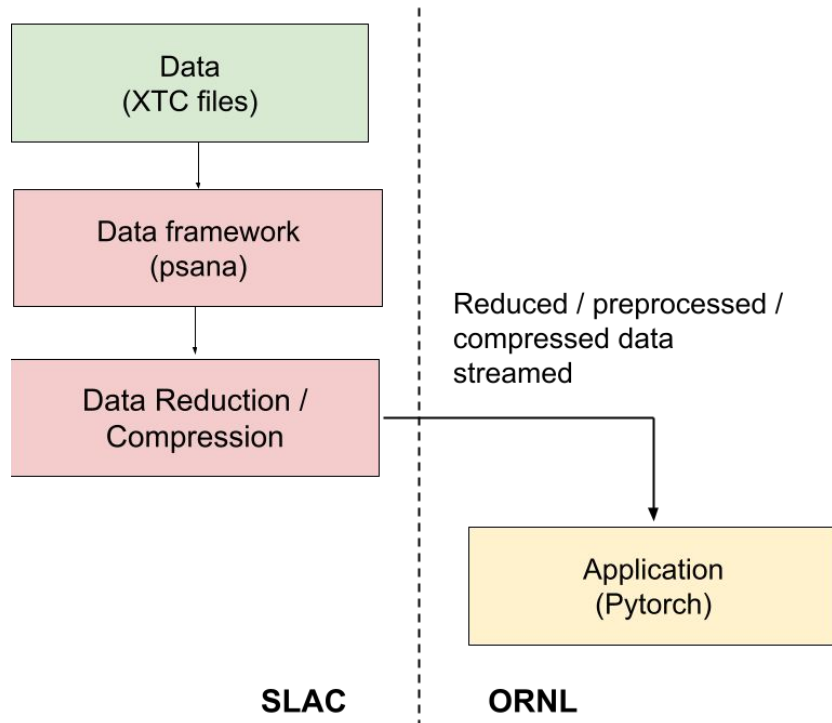


# LCLStream (Summit Plus)

---

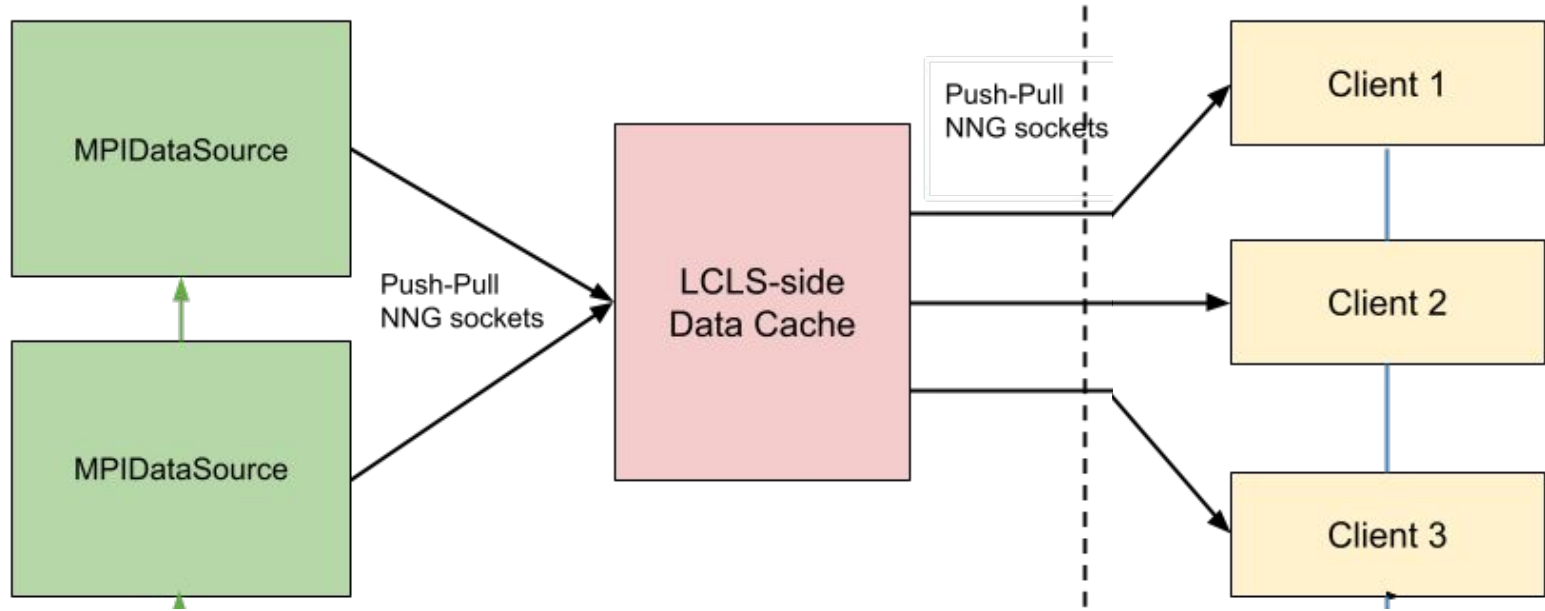
## Data Streaming to Remote Facilities: LCLStream

- Streamed data is reduced:
  - Only relevant data is transferred (detectors, hits)
- Streamed data is preprocessed:
  - “Science Ready”
- Streamed data is compressed:
  - Lossless compression

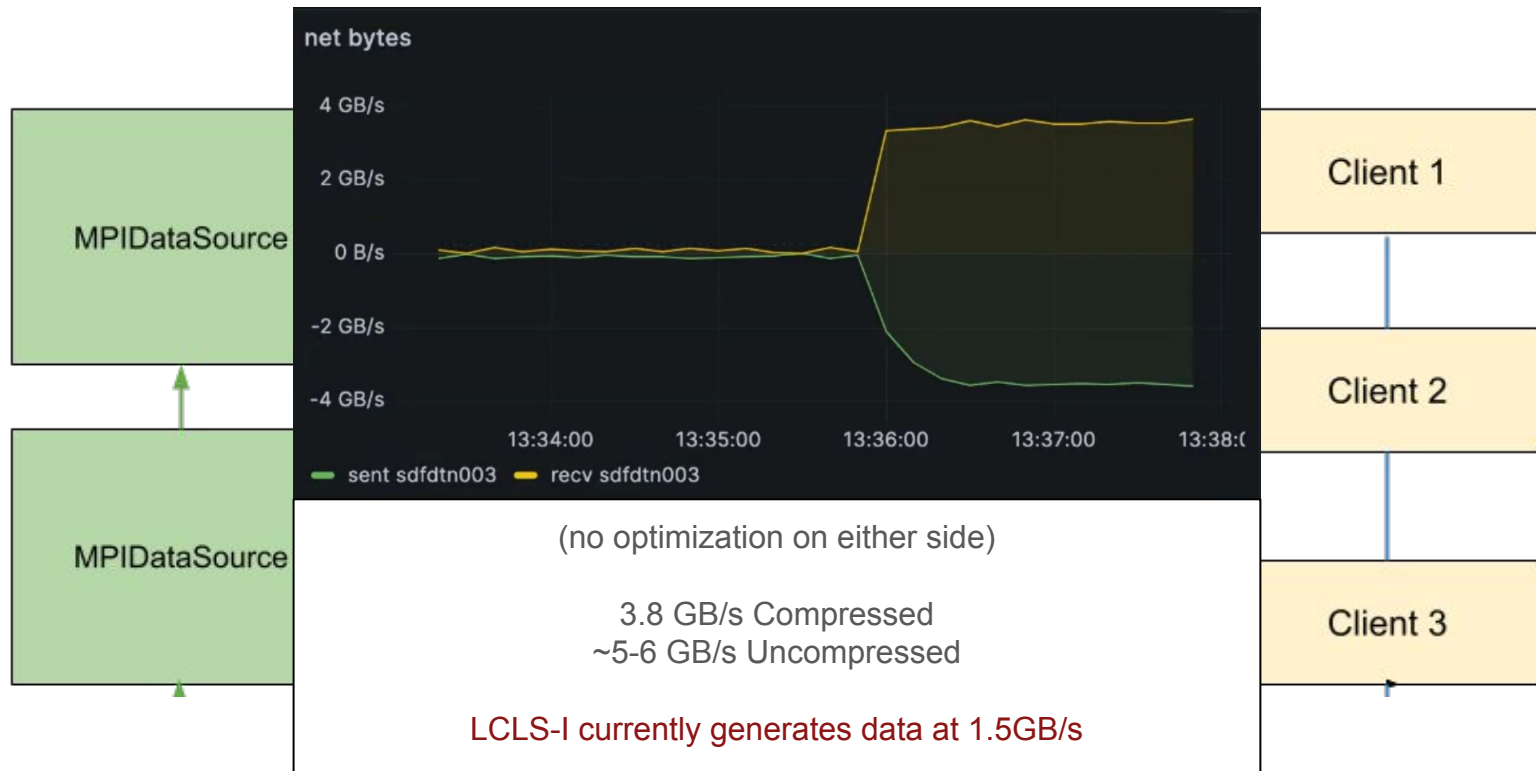


# Remote-Location Data Processing: Streaming Demo

---



# Remote-Location Data Processing: Streaming Demo





# Remote-Location Data Processing: LCLStream

---

## HTTP Server

- REST API

```
{  
  exp: "xpptut15"  
  run: 670  
  access_mode: "smd"  
  detector_name: "epix10k2M"  
  mode: "raw"  
  addr: "tcp://134.79.23.43:5000"  
  img_per_file: 1  
}
```

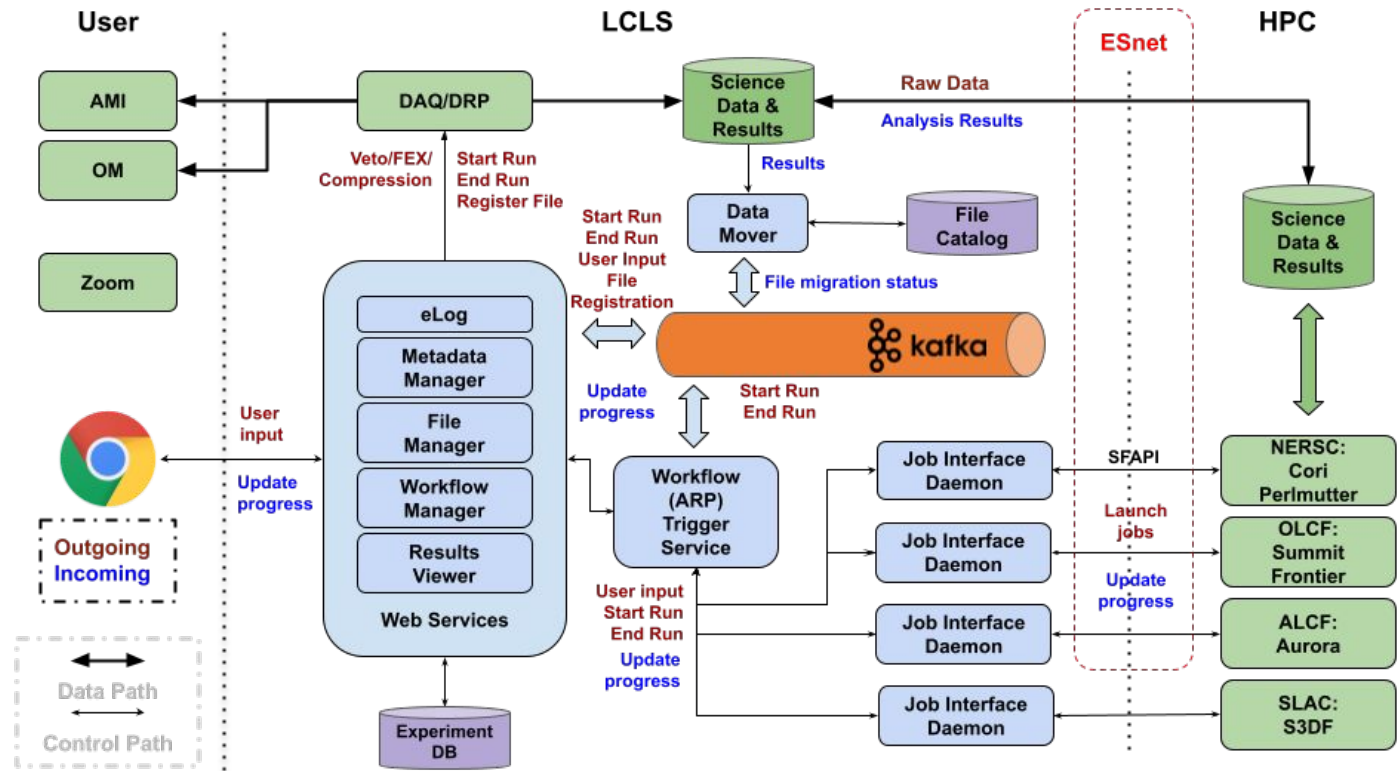


## Streaming Data

- Streamed binary data (No file!)
  - Internal structure: HDF5 file
  - Minimal changes to applications
- ```
data_buffer = socket.recv()  
data = h5py.File(data_buffer)
```
- (then same as before....)
- Compression / filters: HDF5 plugins

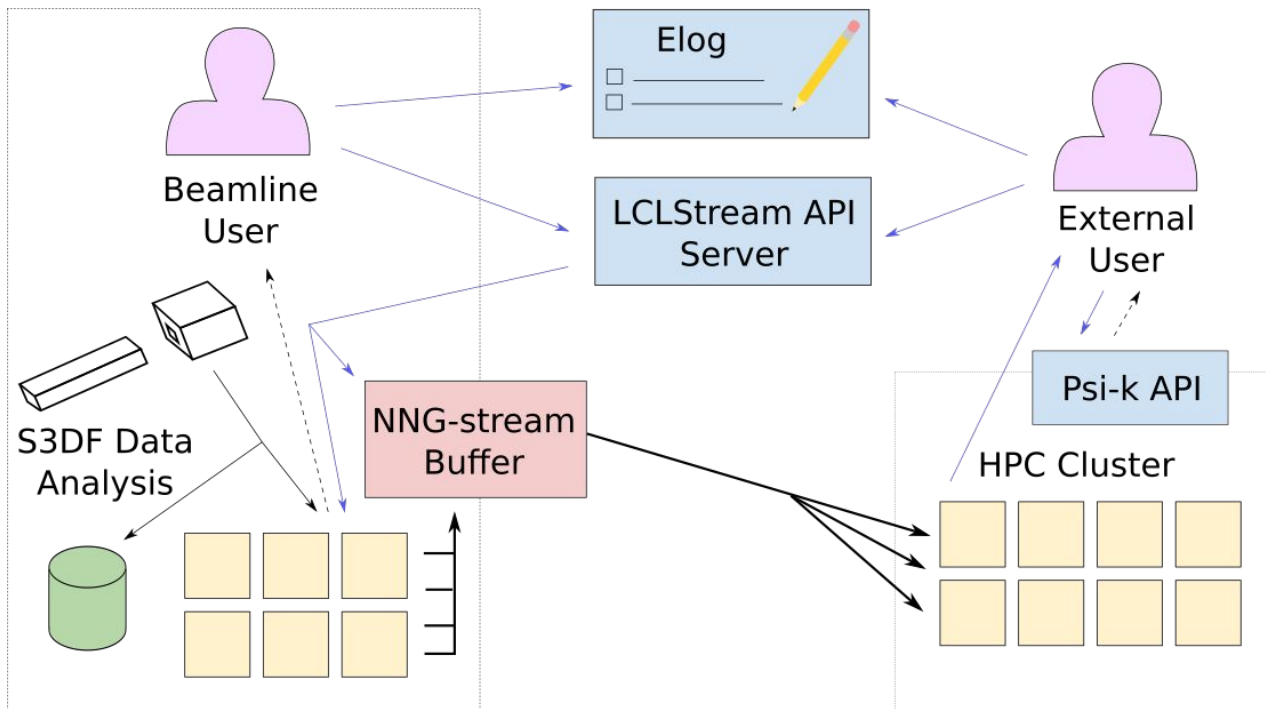


# LCLS: Automated Run Processing (ARP)

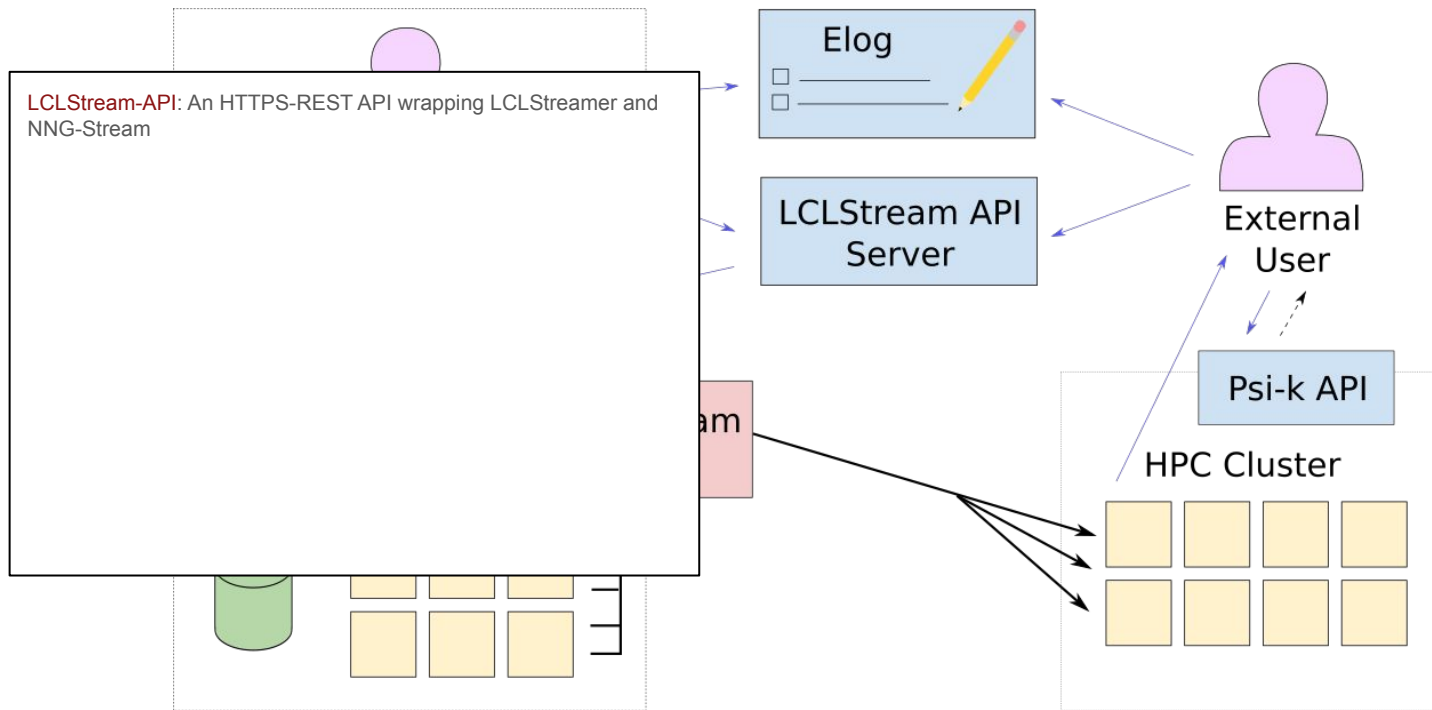




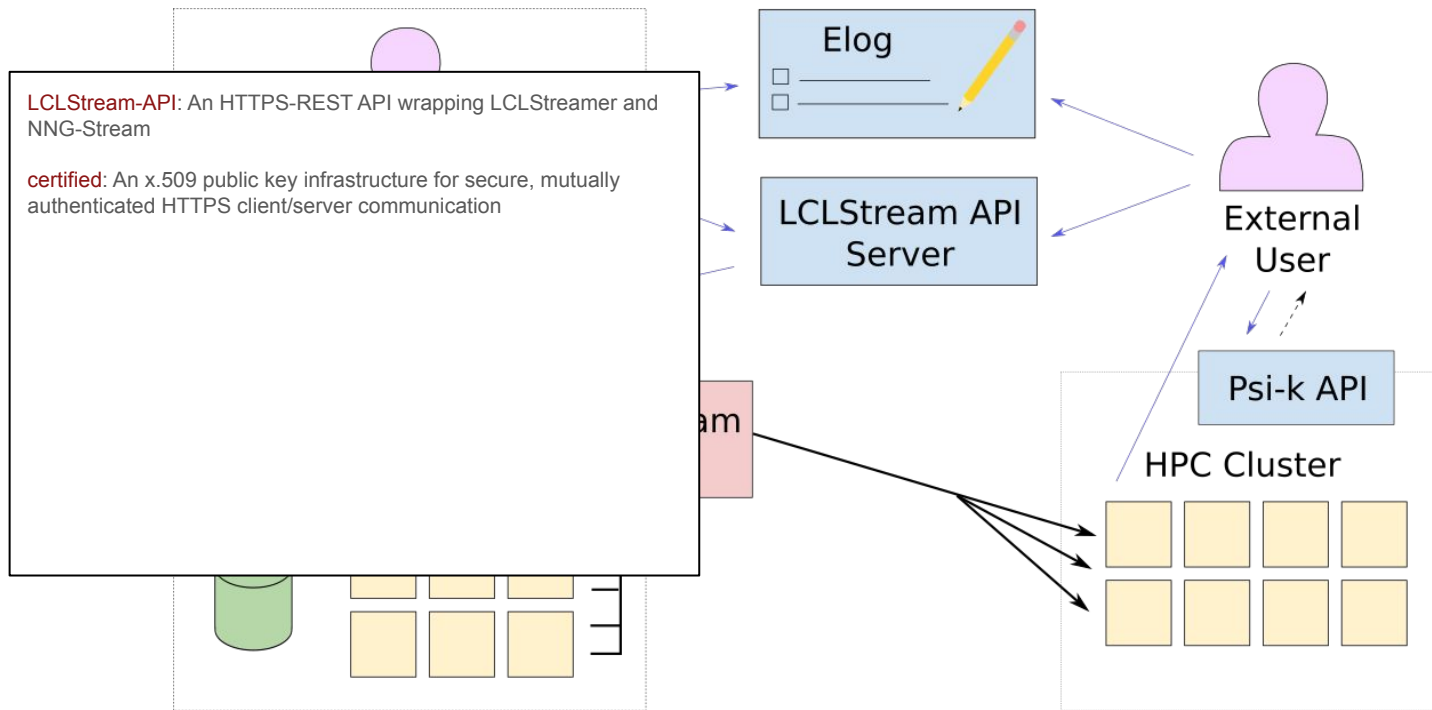
# The LCLStream Ecosystem



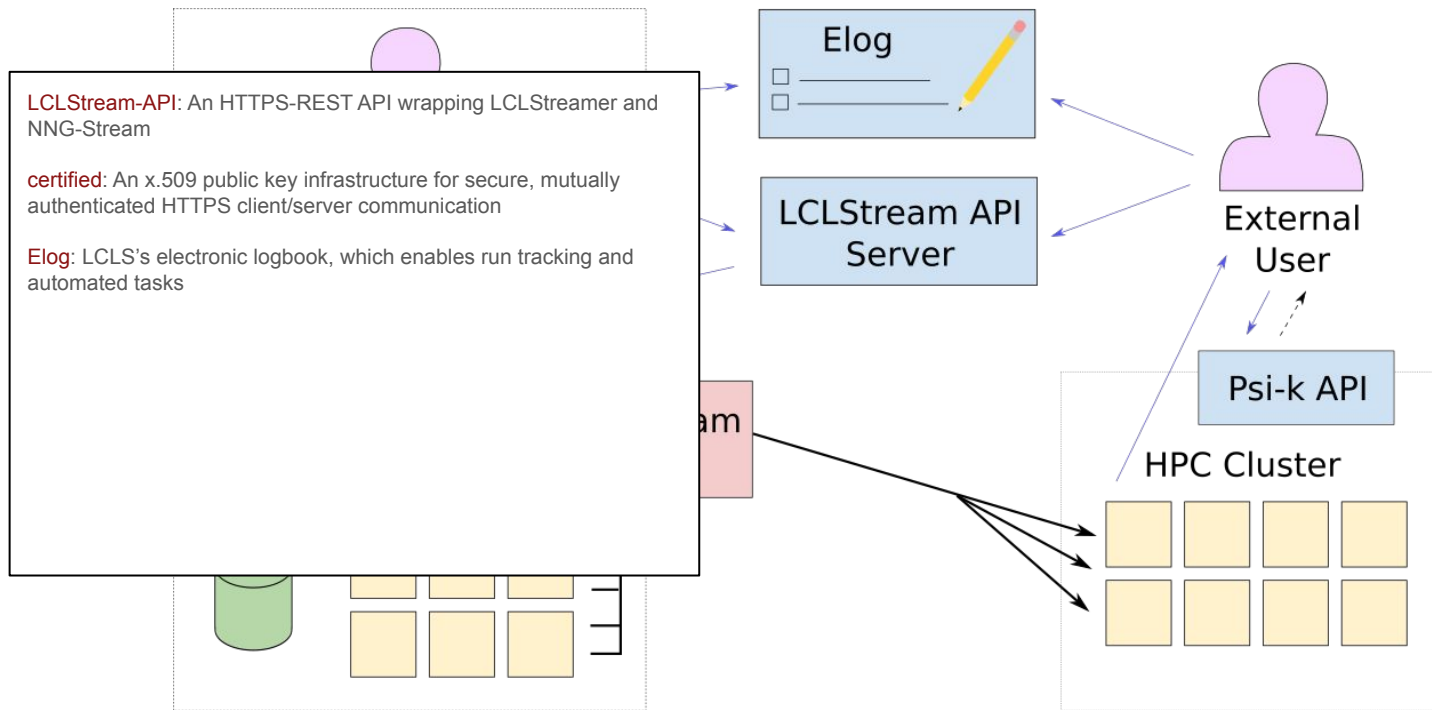
# The LCLStream Ecosystem



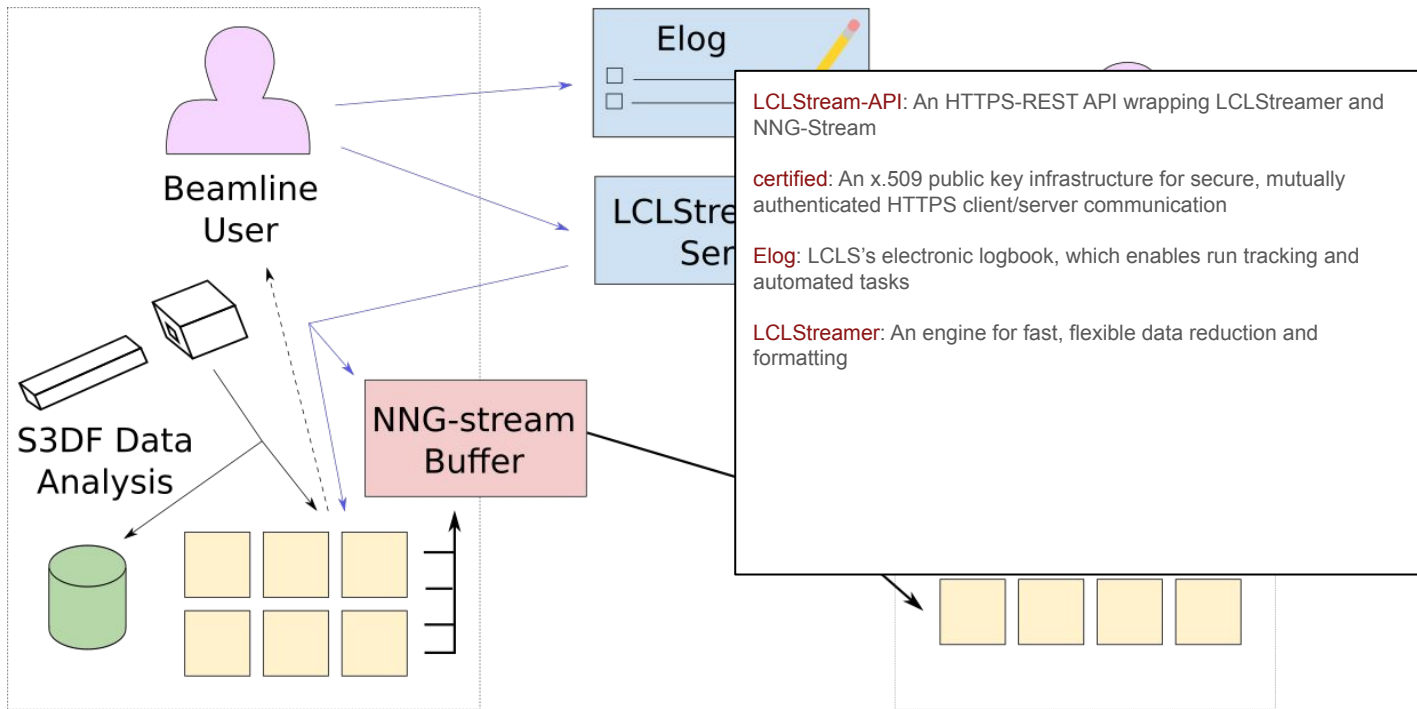
# The LCLStream Ecosystem



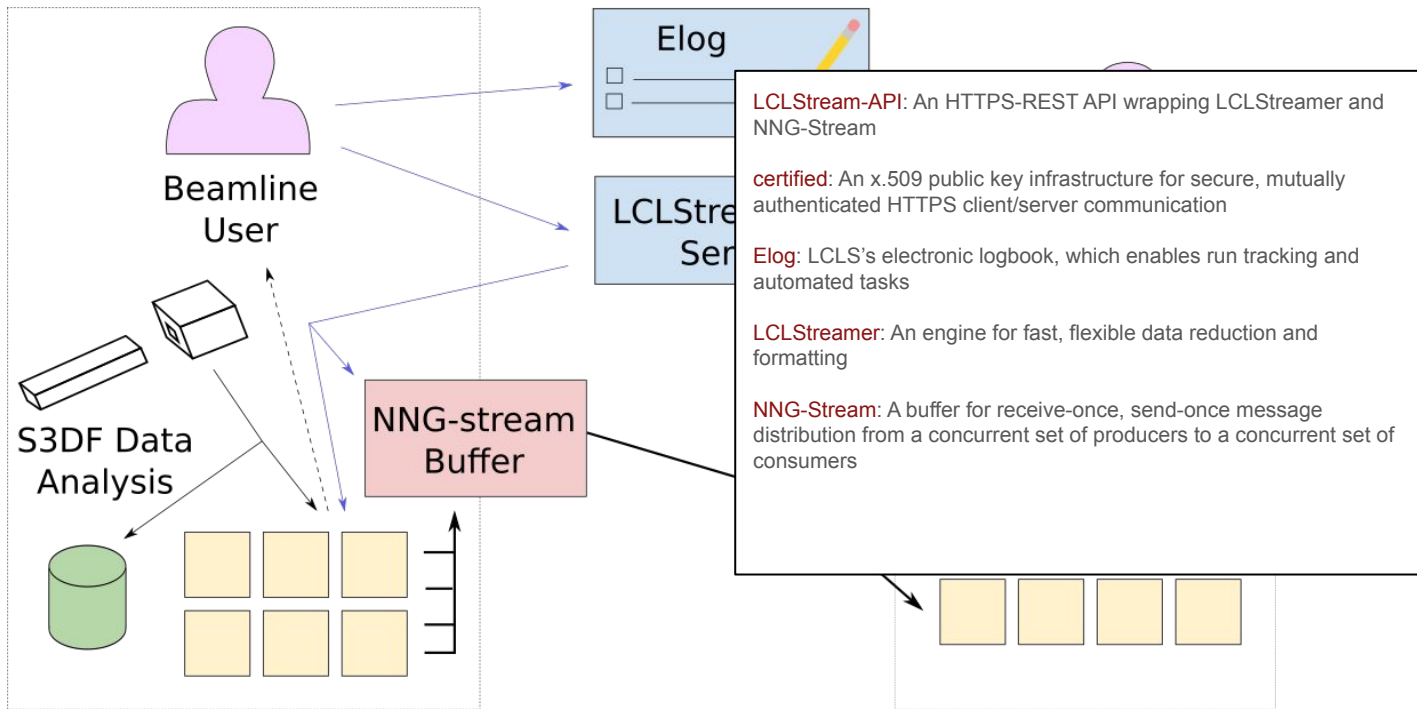
# The LCLStream Ecosystem



# The LCLStream Ecosystem

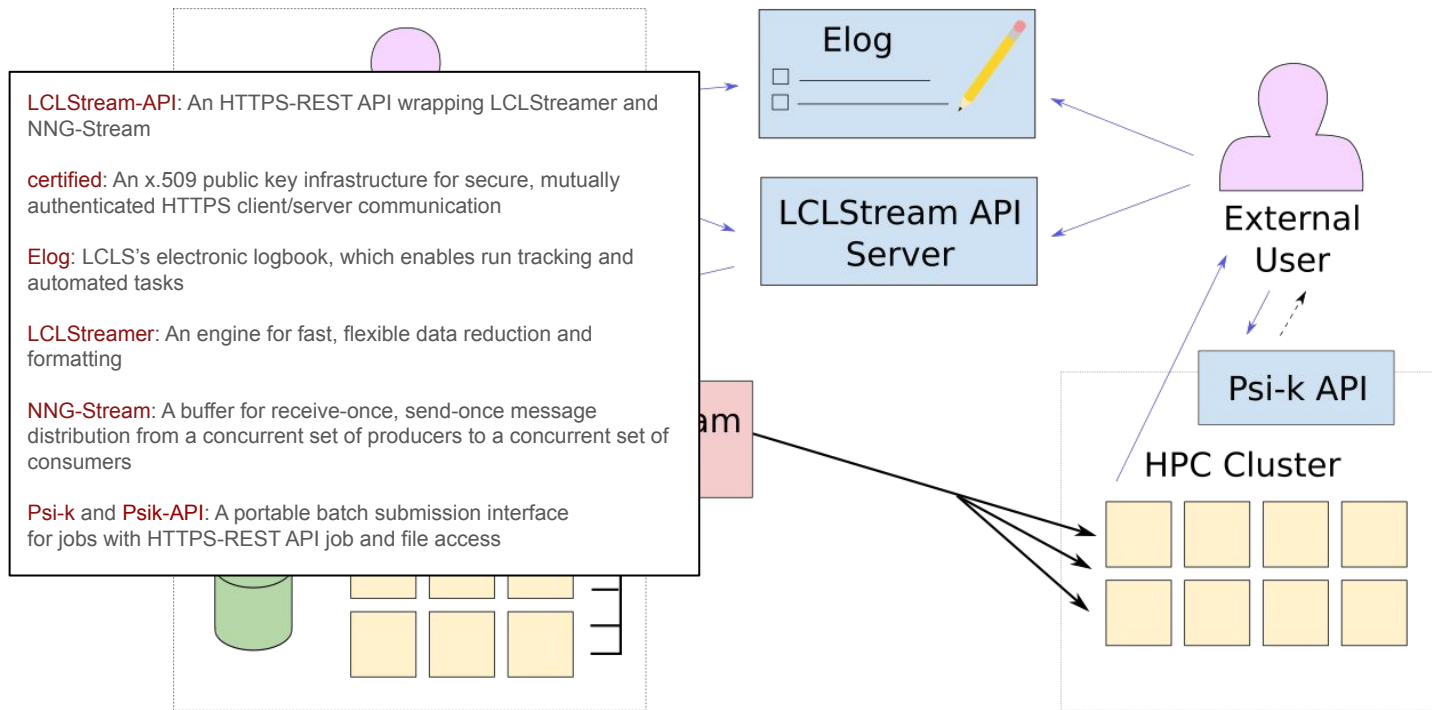


# The LCLStream Ecosystem





# The LCLStream Ecosystem



# CCTBX / SIMPLON

---



- LCLStreamer is modular and configurable
- Data sent in format required by the consumer
- Example: Collaboration with A. Brewser and D. Mittan-Moreau
- Streaming data from LCLS to CCTBX running at NERSC
- Data Format: Simplon (DECTRIS)

# CCTBX / SIMPLON

---



- LCLStreamer is modular and configurable
- Data sent in format required by the consumer
- Example: Collaboration with A. Brewser and D. Mittan-Moreau
- Streaming data from LCLS to CCTBX running at NERSC
- Data Format: Simplon (DECTRIS)

# Local LCLStream: Heterogeneous Computing

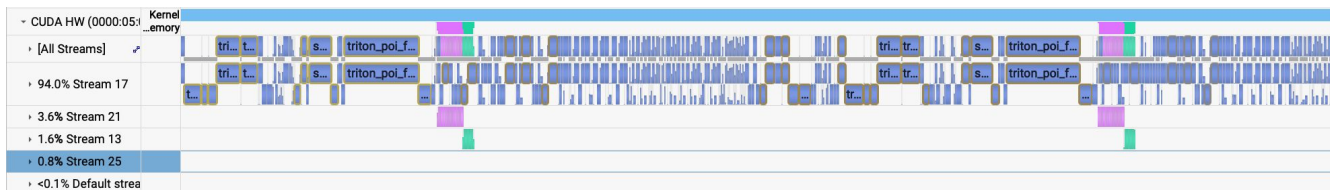
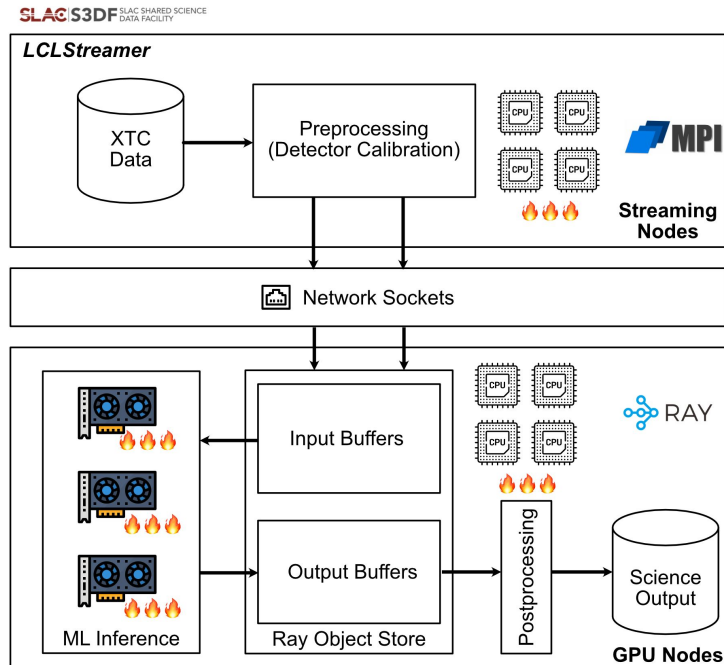
Separate:

- Data reading (LCLStreamer - Psana)
- Data Processing (Processing code - Ray)

Psana: optimized for heavily parallelized processing of single events

GPU: Batches of events in contiguous memory

LCLStream can bridge the gap and optimize GPU usage



# LCLStream: Multiple-lab Collaboration

---

## **SLAC National Accelerator Lab - LCLS**

Valerio Mariani  
Katalin Mecseki  
Cong Wang  
Murali Shankar  
Wilko Kroger  
Jana Thayer

## **Oak Ridge National Lab**

David Rogers  
Tom Beck

## **Lawrence Berkeley National Lab**

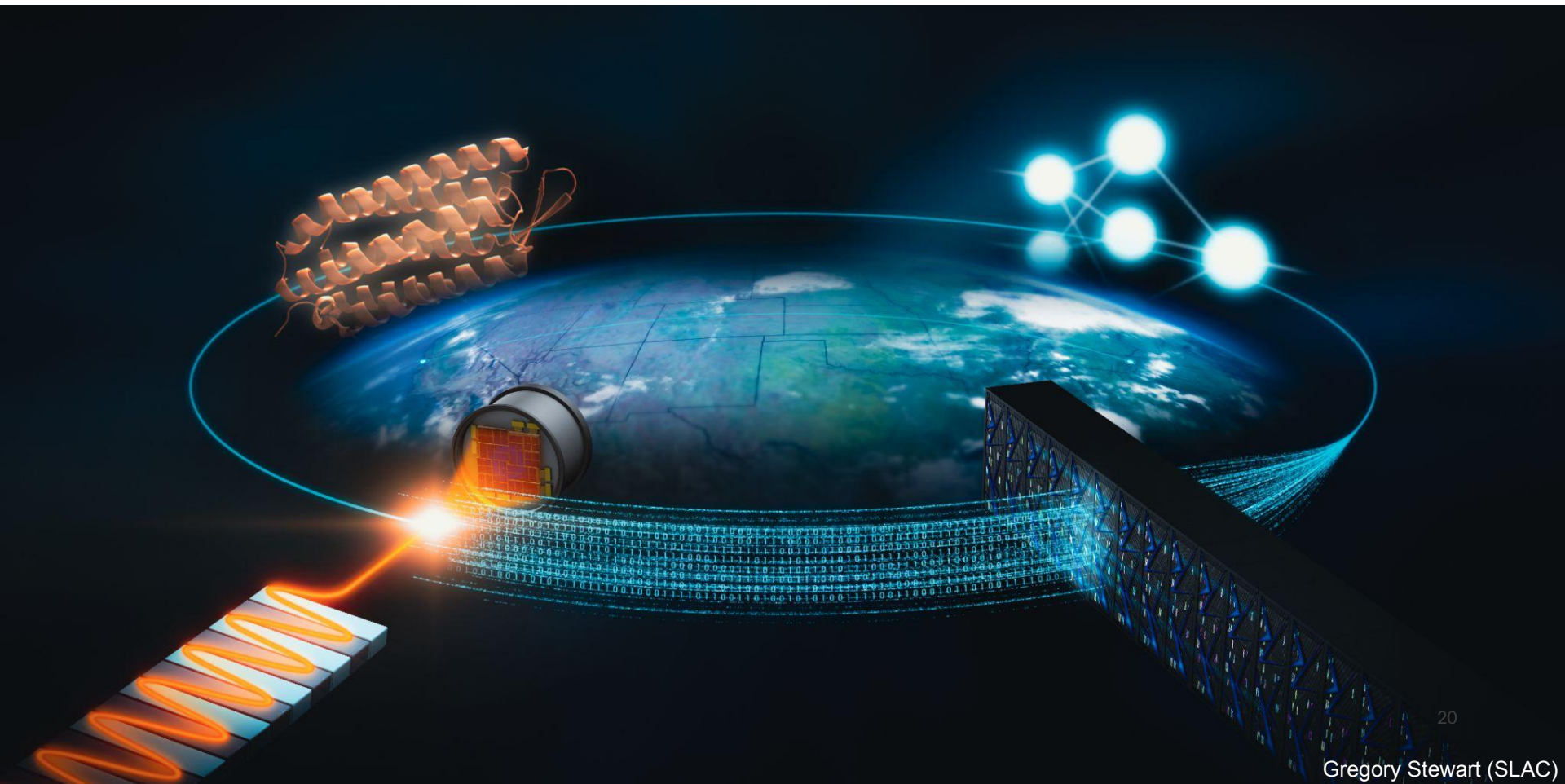
David Mittan-Moreau  
Aaron Brewster

## **National Energy Research Scientific Computing Center**

Johannes Blaschke

THANK YOU!!

Valerio Mariani (valmar@slac.stanford.edu)





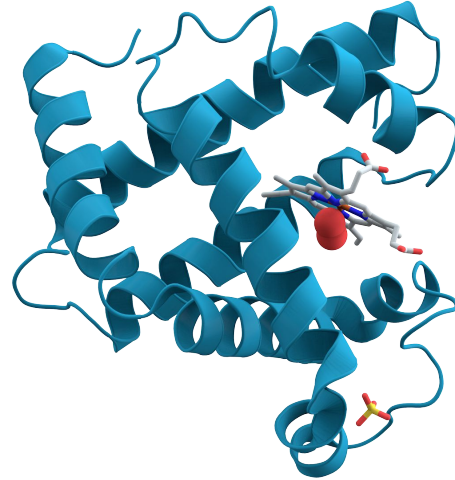


# Serial Crystallography (SFX)

---

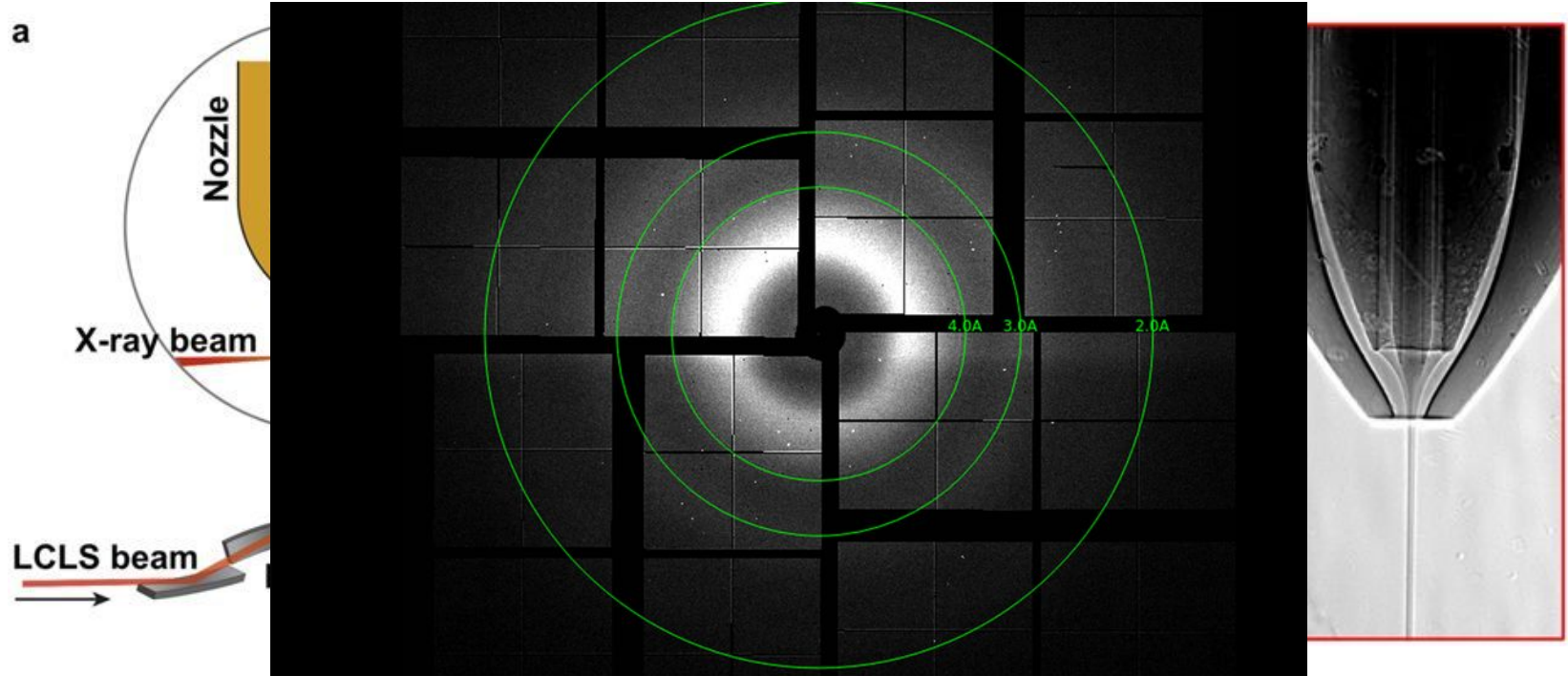


Protein Crystals

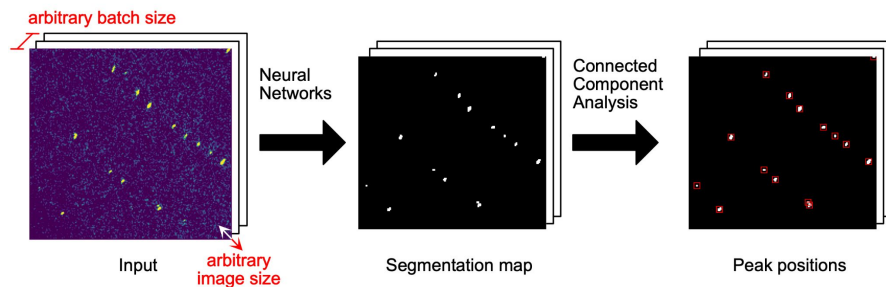


Protein Structure

# Serial Crystallography (SFX)



# PeakNet: A 1 MHz Autonomous Bragg Peak Finder



PeakNet is a deep neural network for

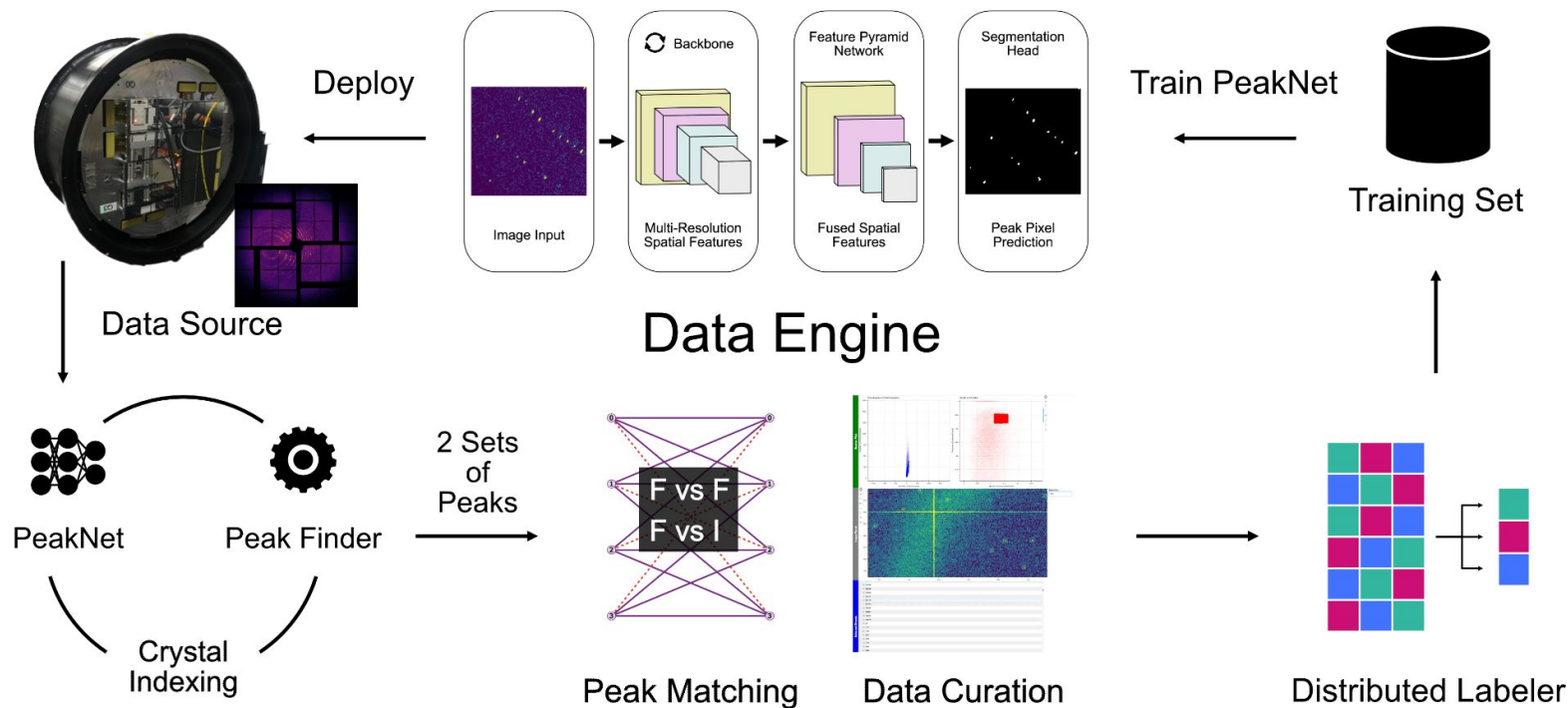
- Autonomous Bragg peak detection in real-time
- Adapts in real-time to shot-to-shot background changes without manual tuning
- Supervised learning (labelled data)

Autonomous pixel segmentation into

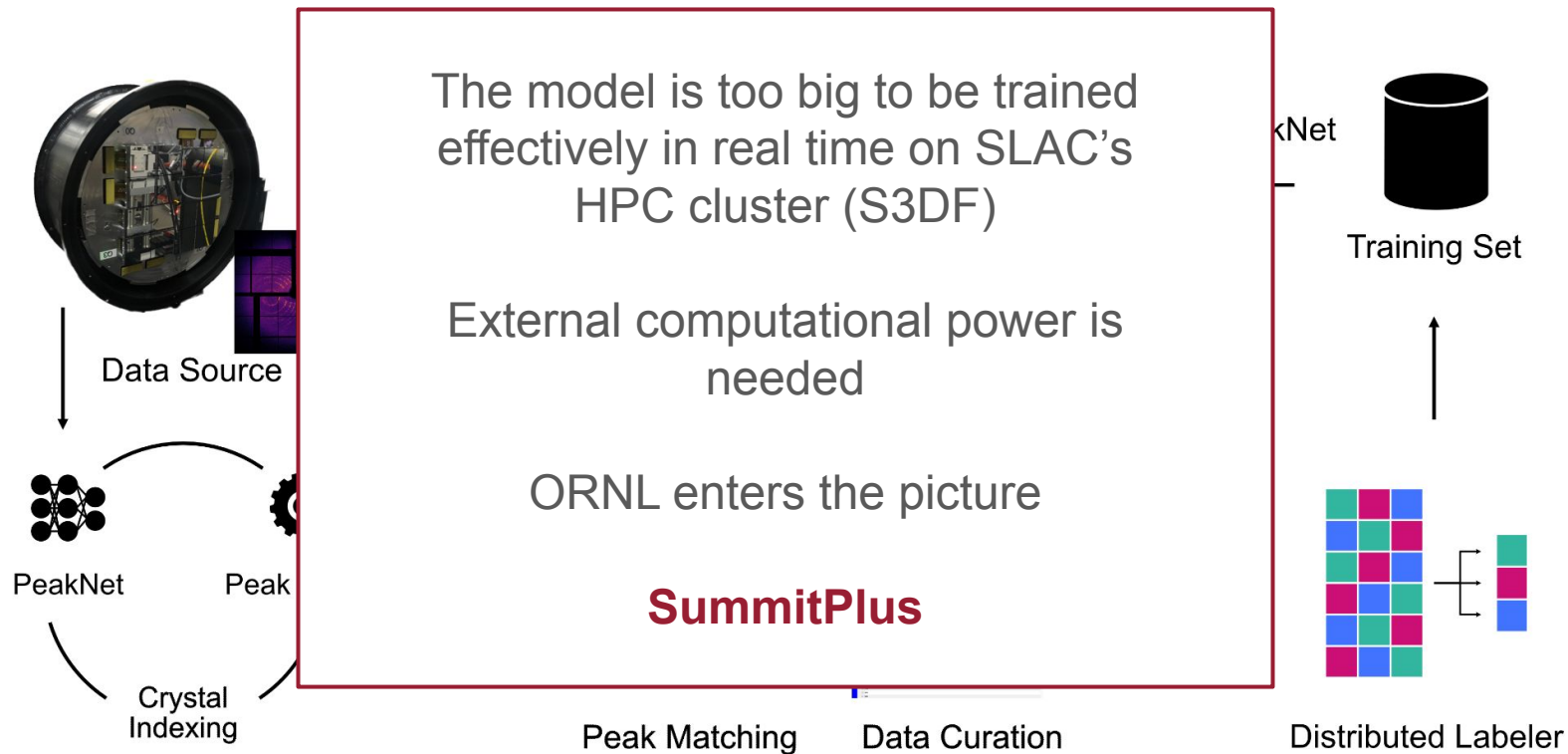
- Artifact scattering
- Background
- Bragg peaks

With no user parameter tuning.

# PeakNet: A 1 MHz Autonomous Bragg Peak Finder



# PeakNet: A 1 MHz Autonomous Bragg Peak Finder





The diagram illustrates the Digital Twin architecture. On the left, a large grey arrow labeled "Data" points into a central box. Above this arrow are logos for NERSC, OLCF, and Argonne National Laboratory, along with a file icon. Inside the central box, the text "Digital Twin" is displayed. To the right, a vertical stack of components is shown: a grey trapezoid labeled "uncompress", a grey rectangle labeled "image", a blue trapezoid labeled "encoder", and a blue rectangle labeled "feature". A red double-headed arrow connects the "feature" block to a blue cube labeled "Model of the Sample". Below the cube are two circular nodes, one yellow and one blue, both with an 'X' inside. A grey arrow labeled "simulated" points from the blue node to the right. A black line connects the "feature" block to the "simulated" output. On the far left, a large grey arrow labeled "Model" points outwards from the central box.

.mtz

# ILLUMINE - SLAC-led 5 light source + neutron source \$10M, 5Y effort for Experiment Steering Infrastructure

A modular framework to close the loop between fast analysis, machine-assisted decision-making, and data acquisition to drive experiments on the timescales of seconds, minutes, or hours

