# LCLS Realtime Analysis Needs at NERSC

**Christopher O'Grady, LCLS Data Systems**

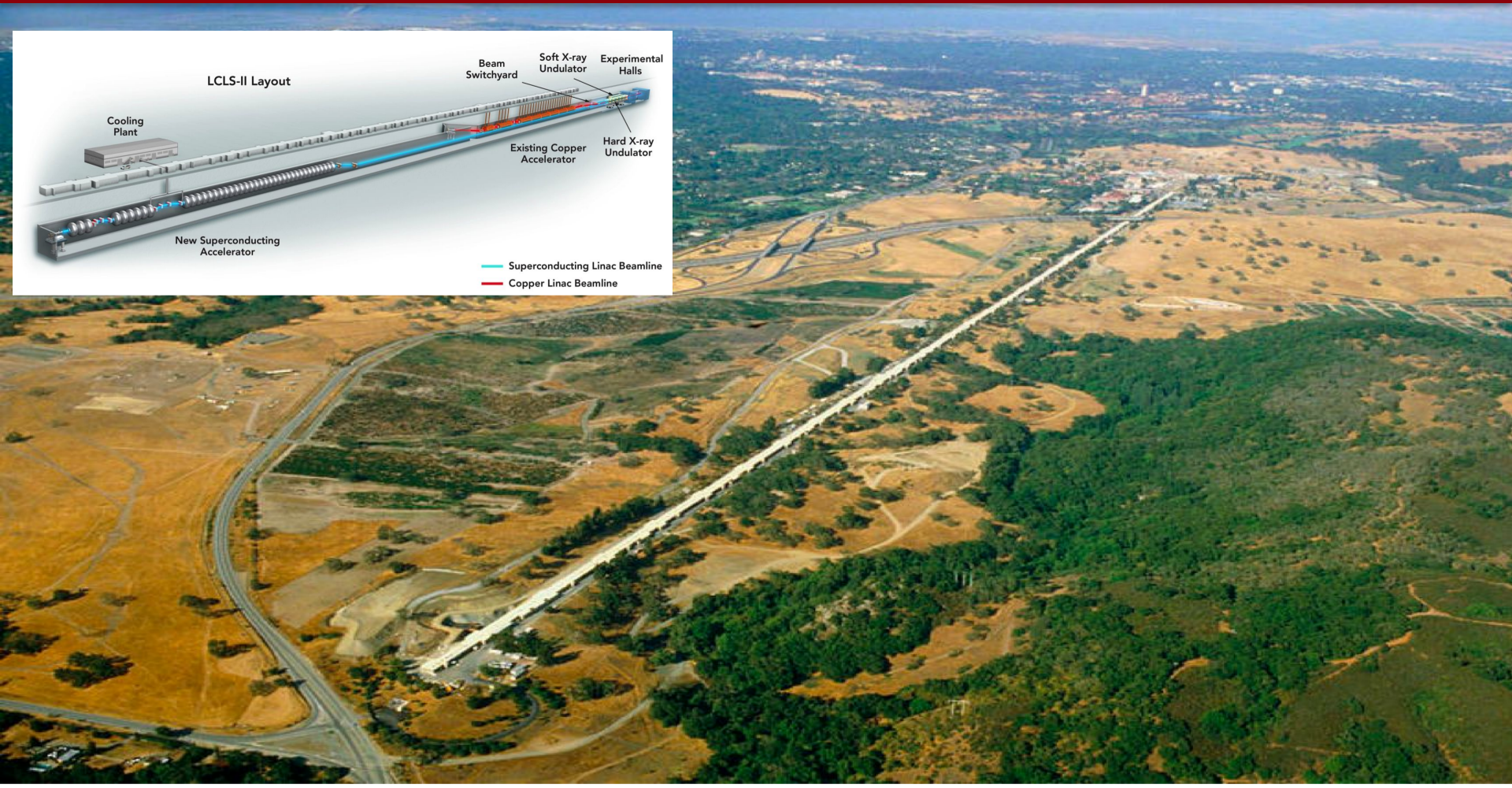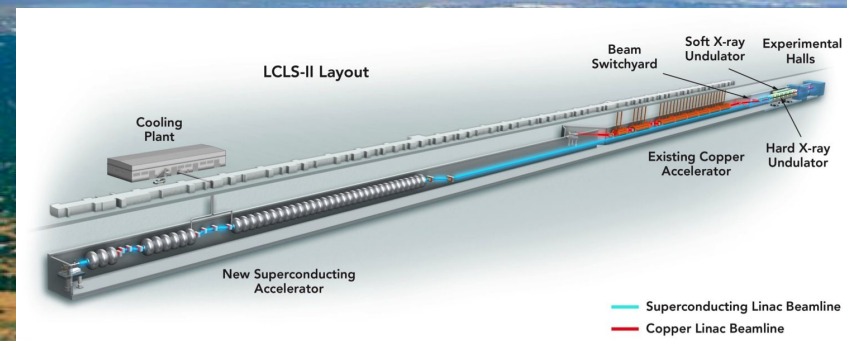**July 8, 2020**

# Linac Coherent Light Source
### … the world's first "hard x-ray" laser

LCLS Injector
(Sector 20)

LCLS Linac
(Sectors 21-30)

LCLS Beam
Transport

LCLS
Undulator Hall

LCLS Near
Experimental Hall

LCLS Office
Building (901)

LCLS X-ray Transport/
Optics/Diagnostics

Endstation
Systems

Endstation
Systems

LCLS Far Experimental
Hall (underground)

**LCLS operates 24 hours/day with 95% beam availability and delivers pulses at 120 Hz**

**LCLS-II, a major (~ B$) upgrade to LCLS is currently underway. Online in 2021.**

# LCLS Realtime Requirements

- ~$1B facility runs 24/7
- **1MHz, 20GB/s** in 2021: requires supercomputers.
- Experiments change significantly multiple times per week
- Realtime data analysis feedback is critical for running experiment
  - ~1s latency for subset of data (before data reaches disk)
  - Few-minute latency for all data (from disk)
- I am here to discuss the **few-minute latency (from disk)** which I will **(loosely) call "realtime"**

**LCLS-II instrument development (underway)**

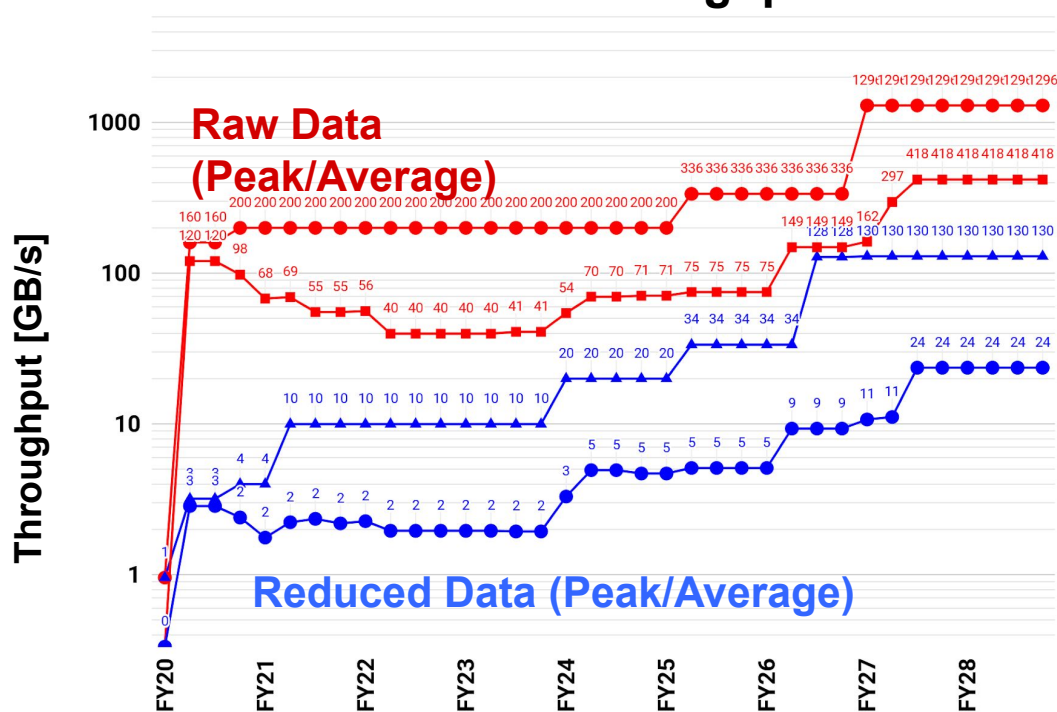**LCLS Hard X-Ray instrument suite, and plans for LCLS-II-HE**



LCLS–II and –HE require a new suite of X-ray instruments, detectors, and data systems, consistent with the leap from 120 Hz to 1 MHz

**LCLS-II will increase data throughput by three orders of magnitude by 2025**

# Throughput Requirements



## LCLS Data Throughput
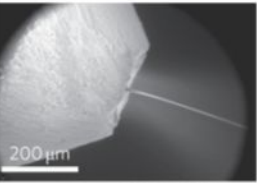
Actual Data Rates lower than peak rates because of:

- **DRP** (shown Reduced Data chart is after data reduction pipeline)

- **Actual utilization** (shown Average Data chart is after adjusting for expected utilization)
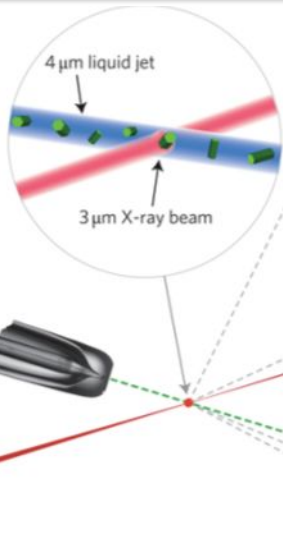
**Data reduction is e$$ential**

# LCLS-II Data System Architecture: Nanocrystallography Example

**SLAC**

## Experiment Description



4 µm liquid jet

3 µm X-ray beam

Gas dynamic virtual nozzle

LCLS beam

200 µm

- Individual nanocrystals are injected into the focused LCLS pulses
- Diffraction patterns are collected on a pulse-by-pulse basis
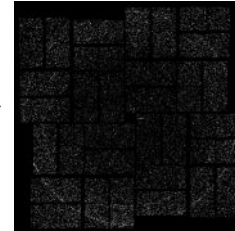- Crystal concentration dictates "hit" rate

**Multi-megapixel detector**



**8 kHz in 2024 (4 MP)**
**40 kHz in 2027 (16 MP)**

**60 GB/s**
**1 TB/s**

**X-ray diffraction image**



**Data Reduction**
- Remove "no hits"
- >10x reduction

**3 TFlops**
**16 TFlops**

**6 GB/s**
**100 GB/s**

**Intensity map from multiple pulses**



**Interpretation of system structure / dynamics**



**Data Analysis**
- Bragg peak finding
- Index / orient patterns
- Average
- 3D intensity map
- Reconstruction

**4 PFlops**
**20 PFlops**

## Data reduction mitigates storage, networking, and processing requirements

## Existing LCLS Computing

- Local hardware:
    - 40 16-core nodes for realtime analysis
    - 80 12-core nodes for offline analysis
    - 7PB Lustre filesystems
- Analysis pattern is embarrassingly-parallel MPI python (scaled to 300,000 cores at NERSC via EXAFEL project)
- LSF batch system (now moving to SLURM)
- Adding more computing ("SDF", shared with all of SLAC) but won't be enough.

- Three levels of job priorities:
  - Running experiment (highest priority)
  - Experiment that will run in 12 hours (second highest)
  - Standard offline analysis
- Each of the 3 priorities can preempt the lower-priority jobs
- Our preemption is imperfect:  lower-priority jobs are suspended but use memory/swap

# Current NERSC Possibilities (my best understanding)

**SLAC**

- **Reservations**
  - Need >1 day advance notice? While useful, LCLS is too dynamic: e.g. accelerator or expt breaks, or job takes longer than expected
- **"Realtime QOS"**
  - Dedicated resources that are idle when not being used. Inefficient, but very useful for smaller users.
  - LCLS has been approved for 20 nodes
- **"Flex" queue**
  - Jobs that can checkpoint (e.g. density functional theory codes like VASP, Quantum Espresso…)
  - Used by NERSC to chop big jobs in small pieces to "fill in the cracks"
- **DMTCP** (https://www.nersc.gov/assets/Uploads/Checkpoint-Restart-20191106.pdf)
  - A work in progress by Zhengji Zhao and others

- "Realtime QOS" is inefficient, so not an option for larger efforts like LCLS
- I've been told "suspended jobs" (remain in memory/swap) is not an option at NERSC
- My **best guess**:
  - **Flex queue is closest**: NERSC system is already preempting checkpointable jobs, which receive a discount
  - **Expand flex-queue idea: a "high-priority queue"** where LCLS pays a premium to be able to preempt flex-queue jobs that can checkpoint (VASP, Quantum, Espresso, DMTCP?)

- How to guarantee **enough low-priority jobs**?
  - No guarantees, but my understanding is VASP etc. is a large fraction of what NERSC does (numbers, anyone?). DMTCP will help.
  - If it looks like it's going to be too small perhaps we can augment with a reservation 1 day in advance, if the "weather prediction" is good enough?
- Only **one level of preemption**
  - Can avoid this if LCLS uses DMTCP so we can preempt our own jobs?