



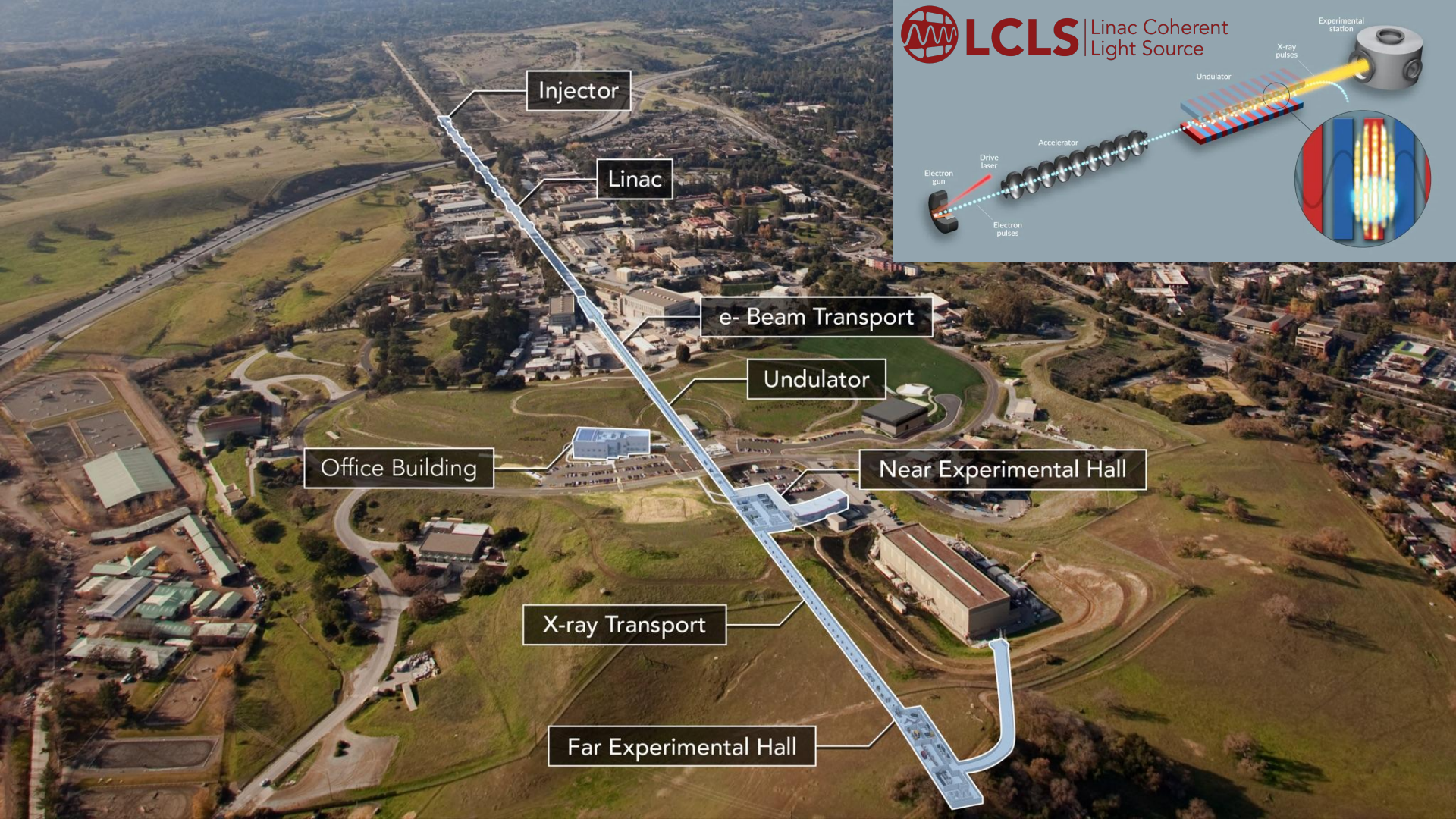
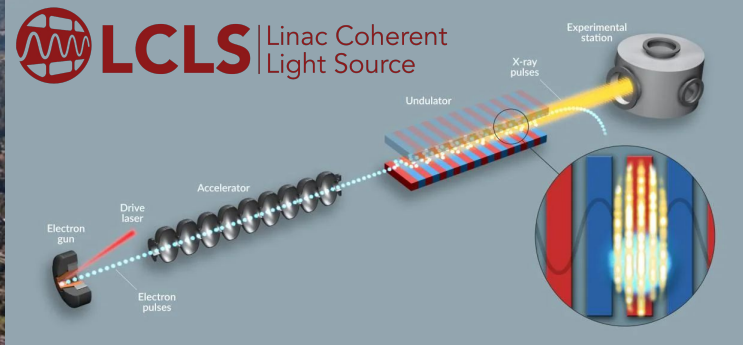
Turning Petabytes into Physics: The case for scalable, user-centric software for next-generation science at LCLS

Jana Thayer/LCLS Data Systems Division Director on behalf of the Data Systems Team & friends
October 21st, 2025

Productive, Performant Software for Large Scale Scientific Data Analysis, SLAC

Many thanks to the people doing the honest work: Matthew Avaylon, Zhantao Chen, Ric Claus, Ryan Coffee, Dan Damiani, Dionisio Doering, Angelo Dragone, Gabriel Dohrlac, Vincent Esposito, Chris Ford, Kevin Dalton, Mikhail Dubrovin, Conny Hansson, Ryan Herbst, Wilko Kroeger, Xiang Li, Doris Mai, Stefano Marchesini, Valerio Mariani, Katalin Mecsesi, Riccardo Melchiorri, Silke Nelson, Chris O'Grady, Amedeo Perazzo, Frederic Poitevin, Omar Quijano, Lorenzo Rota, Amanda Shackelford, Thorsten Schwander, Murali Shankar, Monarin Uervirojnangkoorn, Matt Weaver, Seshu Yamajala, Cong Wang

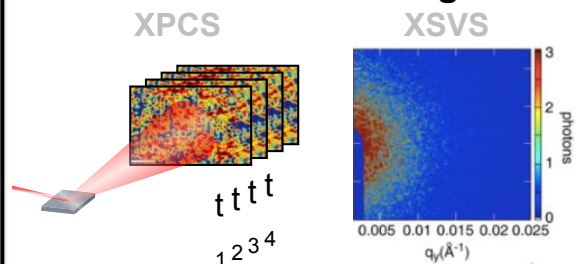
Challenges of Large Experimental Facilities



Visualizing Fundamental Processes at Extreme Timescales

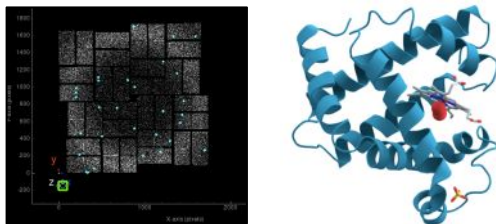
20+ experimental techniques with different workflows, & extreme throughput, and compute needs

Coherent Scattering



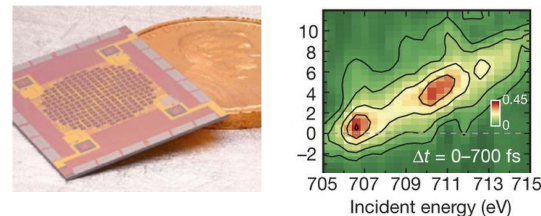
2024: 20 GB/s, 4 TF (reduction), 34 TF (analysis)
2027: 80 GB/s, 34 TF (reduction), 270 TF (analysis)

Nanocrystallography



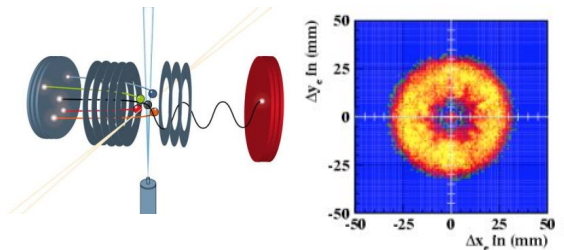
2024: 64 GB/s, 3 TF (reduction), 4 TF (analysis)
2028: 1.2 TB/s, 16 TF (reduction), 20 TF (analysis)

Resonant Inelastic Scattering



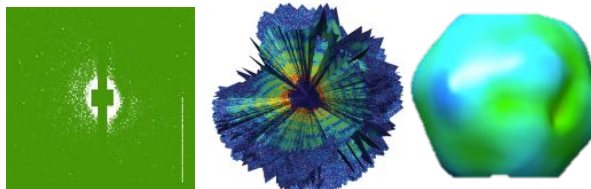
2023: 20 GB/s, 4 TF (reduction), 1 TF (analysis)
2025: 200 GB/s, 40 TF (reduction), 2 TF (analysis)

Coincidence Spectroscopy



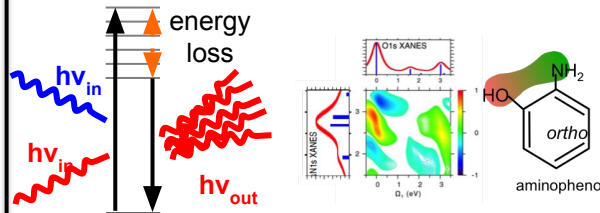
2024: 200 GB/s, <1TF (reduction), <1TF (analysis)

Coherent Imaging

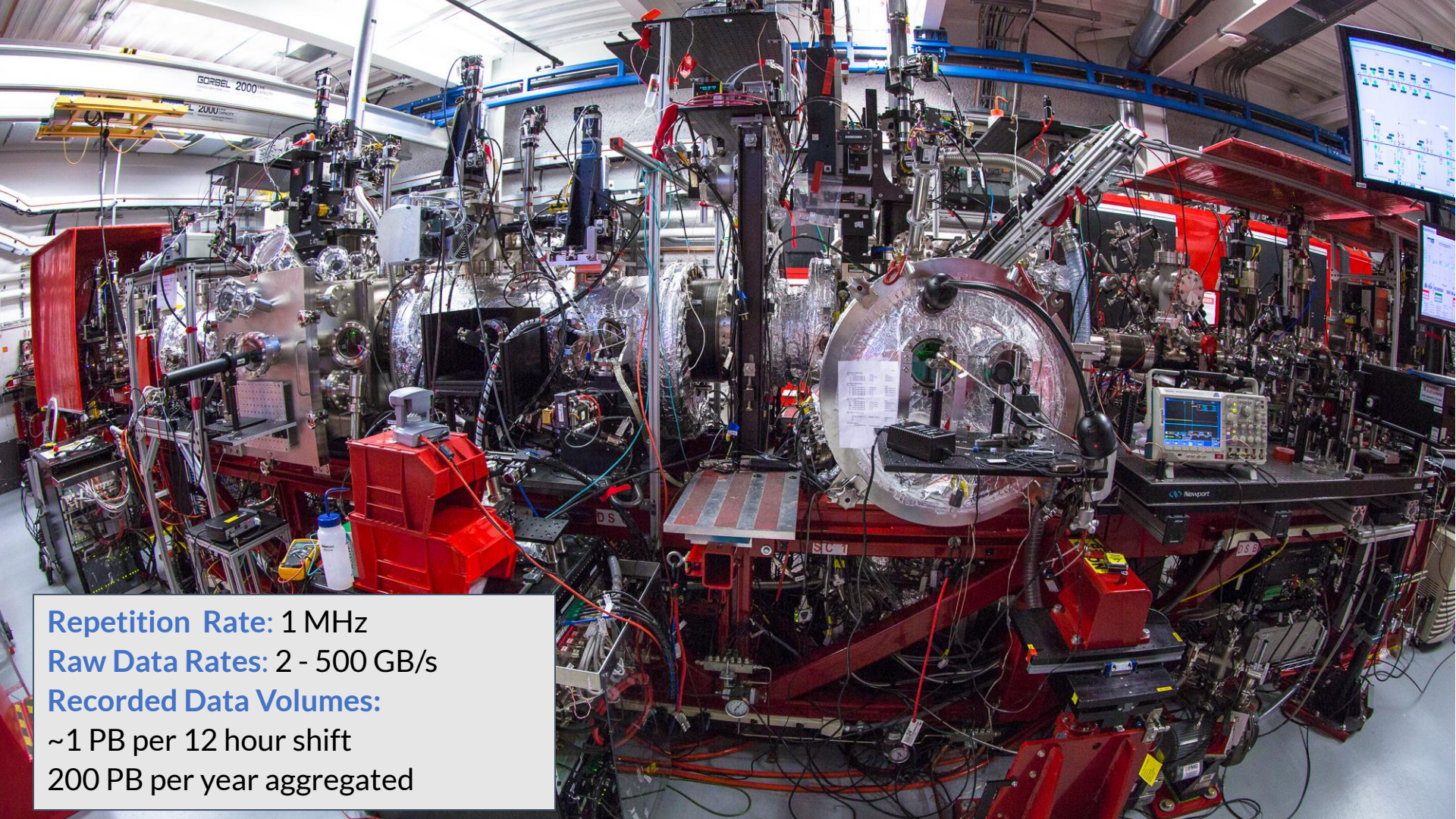


2024: 64 GB/s, 3 TF (reduction), 270 TF (analysis)
2027: 1.2 TB/s, 16 TF (reduction), 1340 TF (analysis)

Nonlinear Spectroscopy



2024: 20 GB/s, 3 TF (reduction), <1 TF (analysis)
2025: 80 GB/s, 16 TF (reduction), <1 TF (analysis)

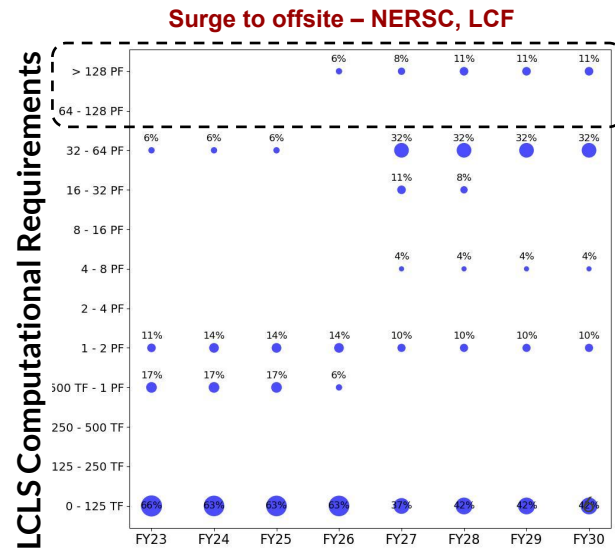
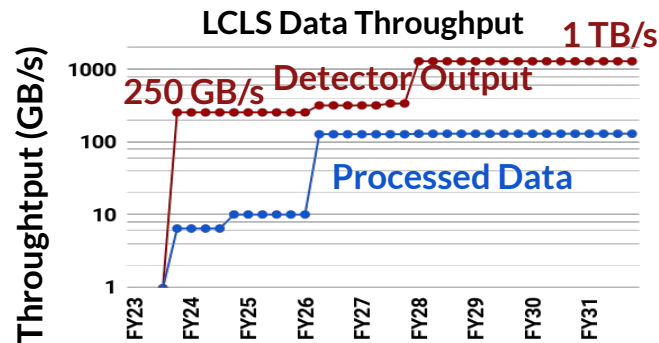


Repetition Rate: 1 MHz
Raw Data Rates: 2 - 500 GB/s
Recorded Data Volumes:
~1 PB per 12 hour shift
200 PB per year aggregated

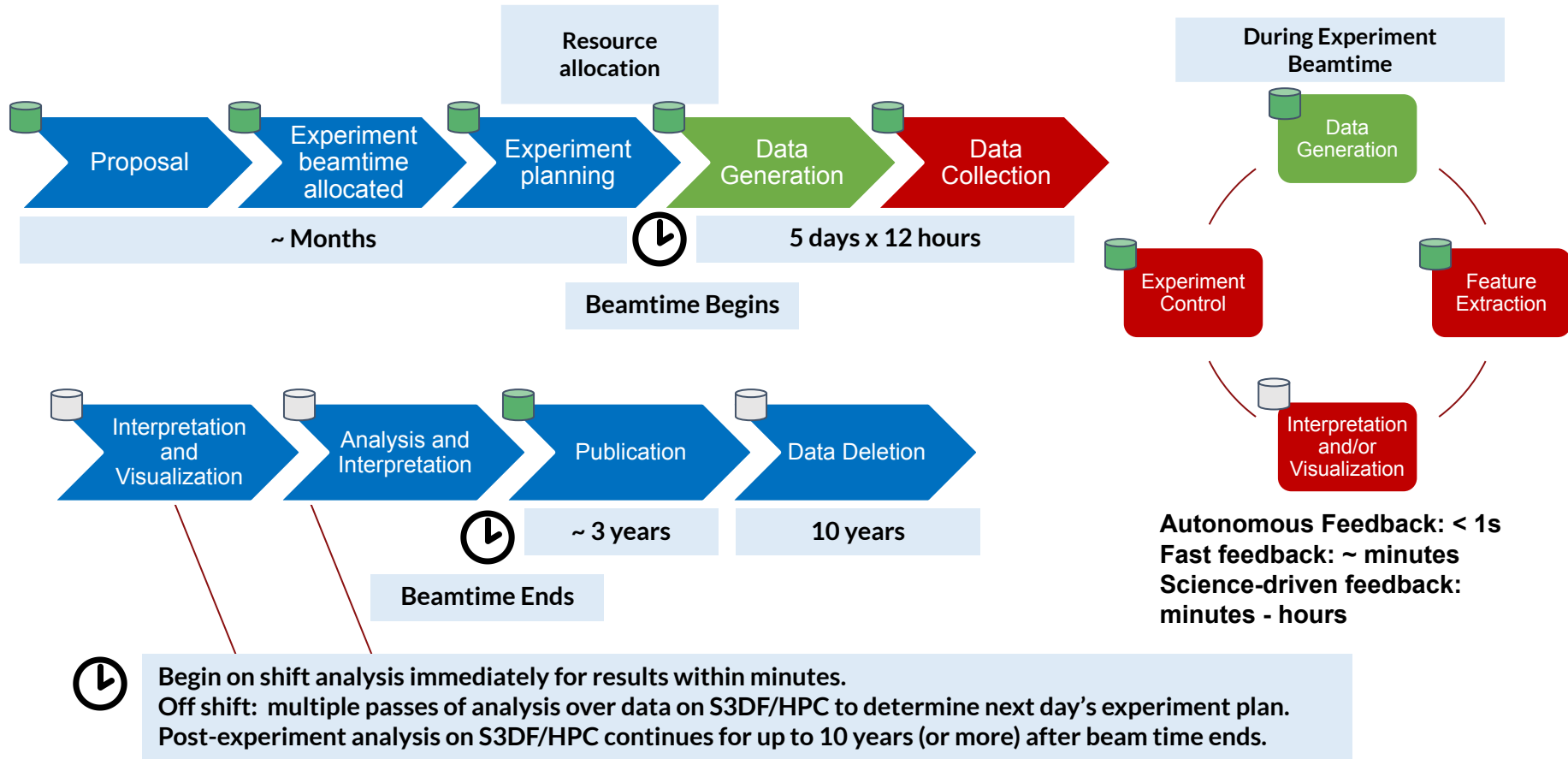
LCLS-II Data System Drivers

- **LCLS-II Upgrade:** greater data *velocity*, *volume*, and *complexity*
 - Data Rates:** 120 Hz to 1 MHz (**10000x**)
 - Raw Data Volumes:** 2 GB/s to 500 GB/s (**100x**)
 - Recorded Data Volumes:** 2 GB/s to 50 GB/s (**10x**)
 - Computational Requirements:** 80% ~1 PF, 20% ~1 ExaFLOP
- **Fast Feedback:** real-time analysis (sec/min) is essential to the users' ability to make informed decisions during experiments.
- **Variability:**
 - **Wide variety of experiments** with turnaround ~days
 - **Large dynamic range:** device readout 0.01 Hz - 1 MHz
 - **Data Complexity:** Variable length data (raw, compressed)
 - **Access patterns** to data vary by experiment and detector
 - Analysis is a mix of **tried-and-true** & **innovative techniques**
- **Time to Science:** **Development cycle** must be fast & flexible
- **No user left behind:** alleviate the pressure on users to gather resources to mount a significant computing effort.

Wide variety of experiments that need to modify analysis during experiments

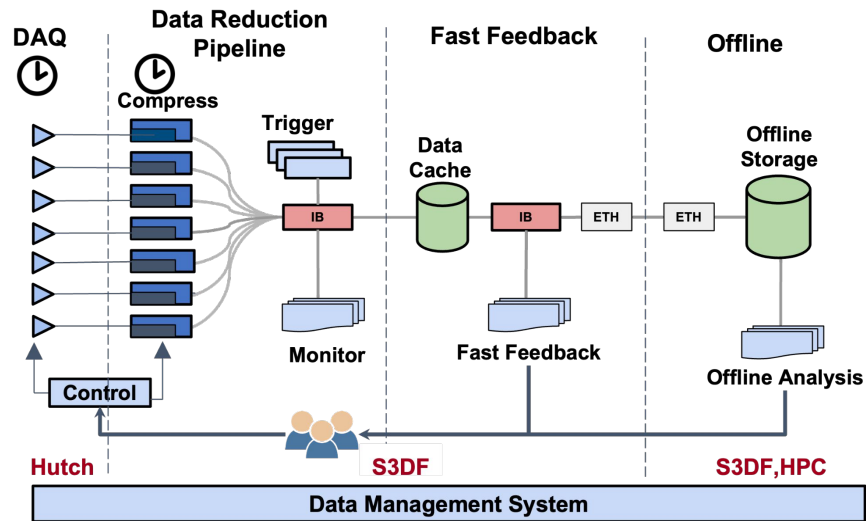


Data Lifecycle Repeated for Hundreds of Experiments / Year



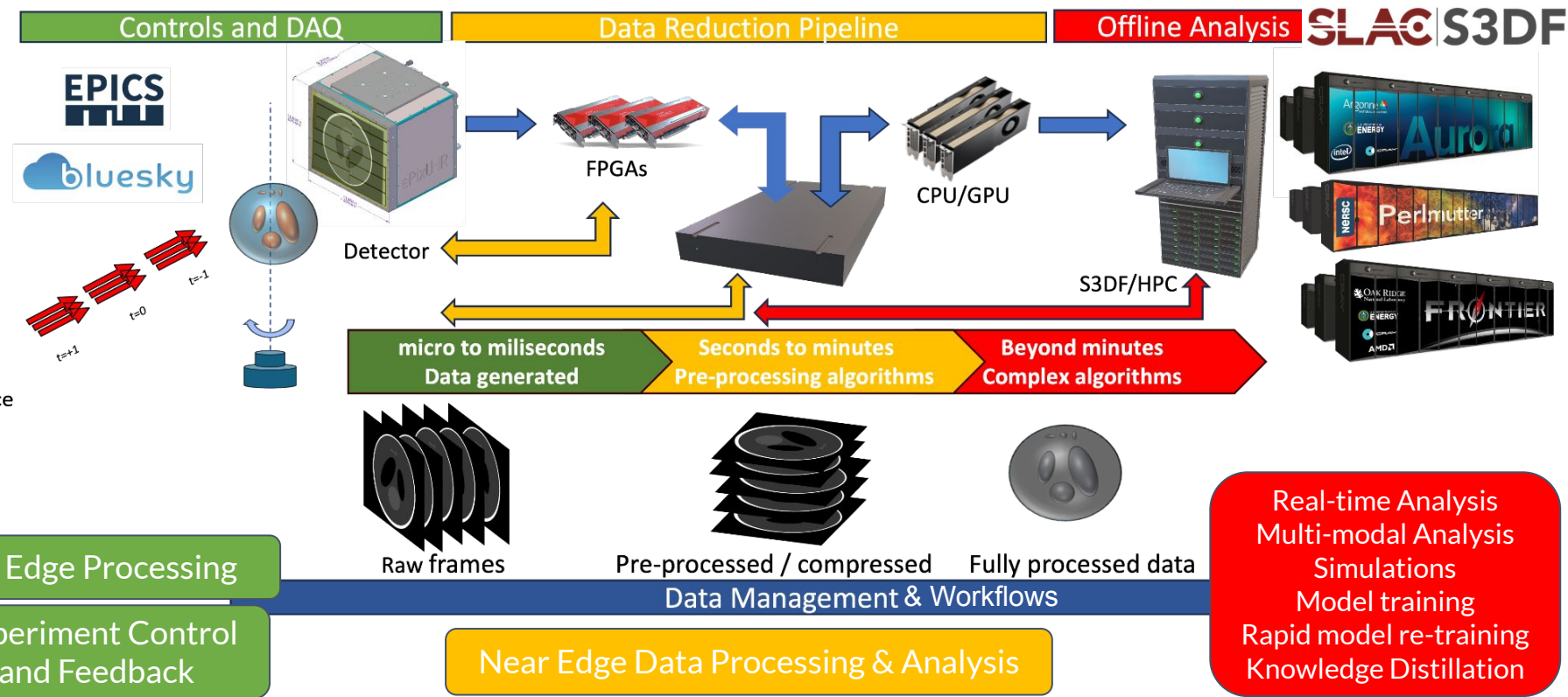
Big Data Handling Strategies for the Linac Coherent Light Source

- Do More With Less Data:
 - Data Reduction/Feature Extraction
 - Implications for Data Format
- Keep Up with the Data Rate
 - Analyze Data at the Rate of Production
 - Actively Monitor Data Quality
 - Performant, Scalable Software
- Seize the Means of Data Production
 - Use Actionable Info to Steer Experiments
 - AI assisted Decision-Making
 - Run algorithms at every layer of the pipeline
- Use Integrated Hardware and Software Infrastructure
 - Local resources and seamless access to remote HPC resources
 - Workflow Orchestration
- Data is a National Resource: Data curation and data management
 - Generate analysis pipeline-ready data products at the point of data collection
 - Automate metadata collection to render data findable and found data useable/reuseable



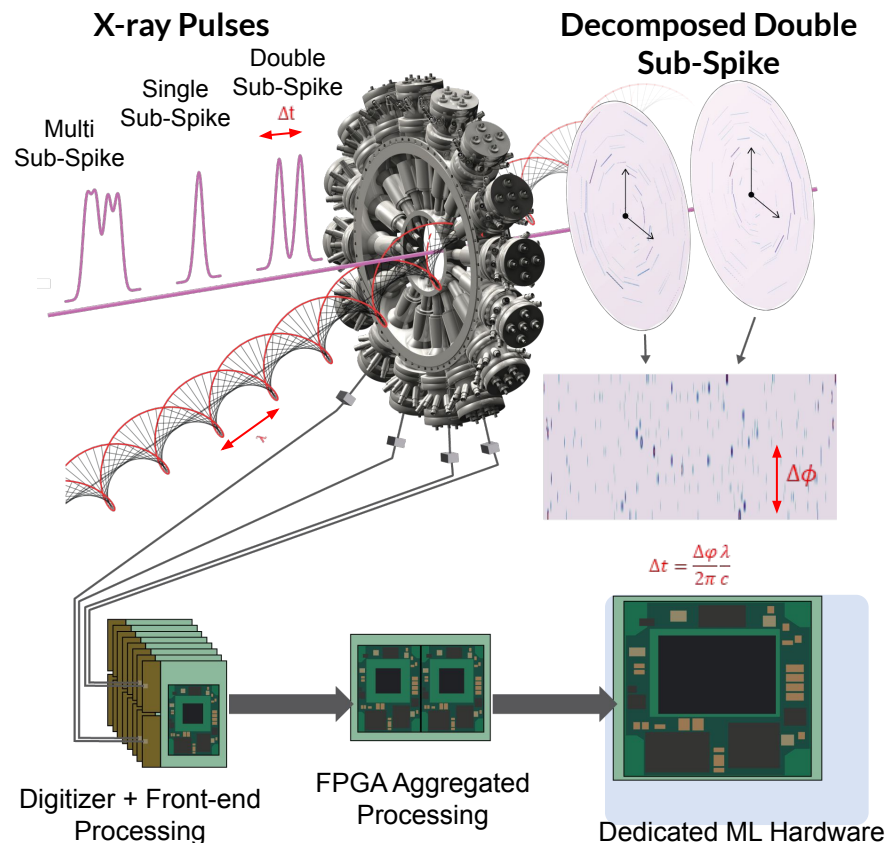
LCLS-II Data System infrastructure for big data handling strategies

LCLS-II Data System is a heterogeneous pipeline with intelligence and computational power at the detector, the data reduction pipeline, local data center and remote HPC.



Life at the Edge

EdgeML for Source to Sink Analysis and Steering



Multi-Resolution COokiebox (MRCO)

- Destination: S3DF
- Computing: EdgeML in FPGA
- Throughput: 200 GB/s \rightarrow 2 GB/s

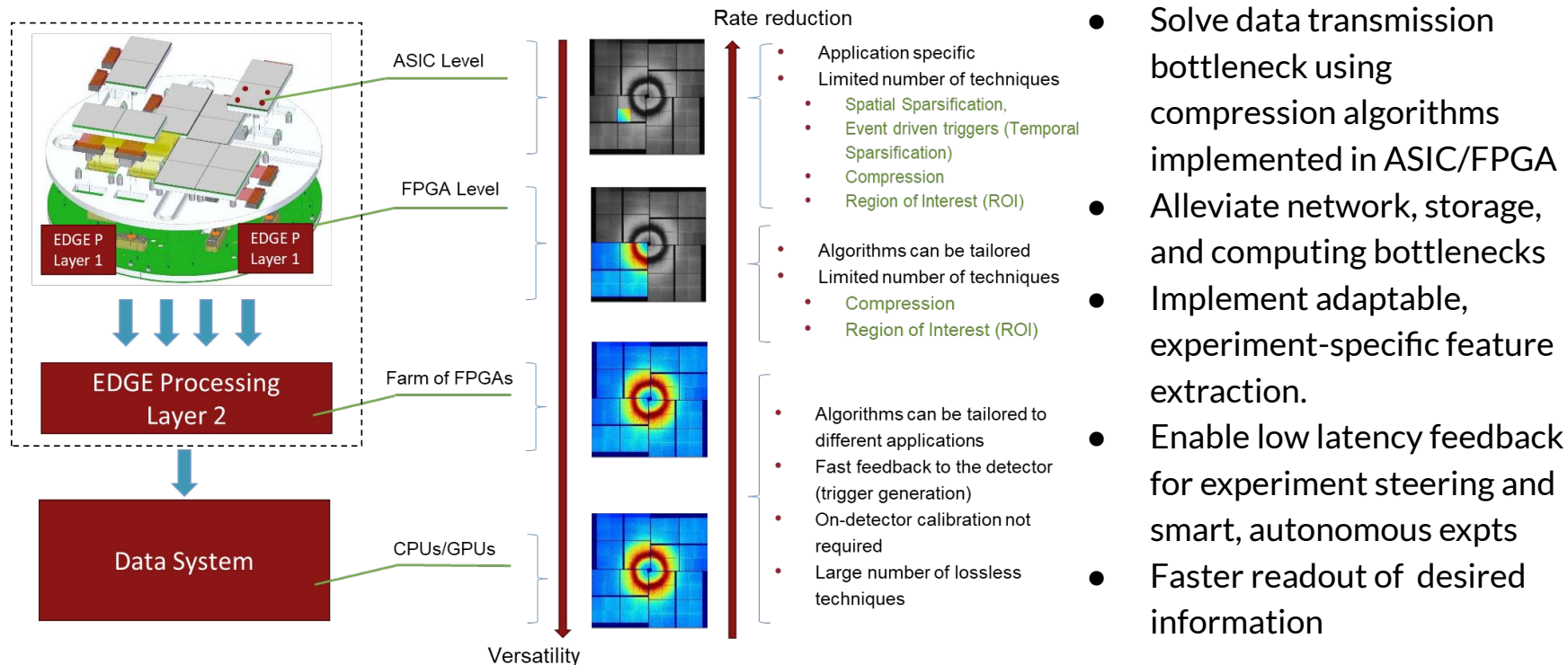
EdgeML for pulse reconstruction, closed-loop control, and AI enabled decision making feedback achieves 30 attosecond resolution

- EdgeML in FPGA differentiates # of sub-spikes/shot, enables veto and live data sorting
- Groq inference throughput > 10 kHz
- Experiment steering

In addition to ASIC/FPGA, novel accelerators possible for inference at edge, but need programmability

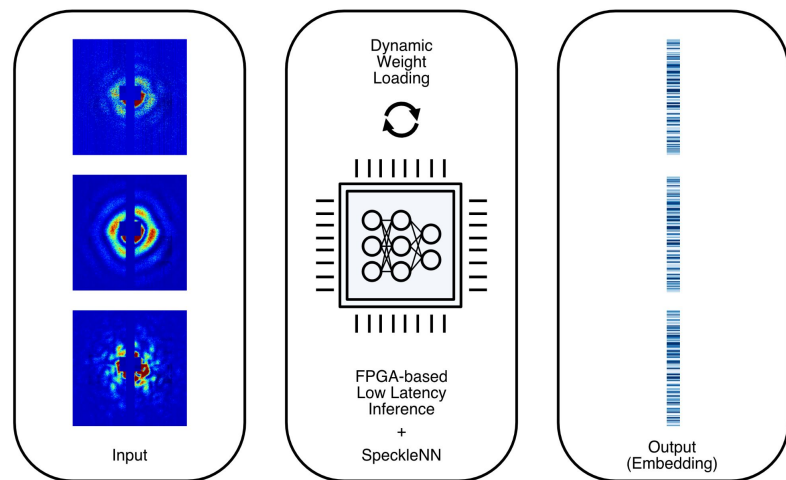
Edge Intelligence: push computing closer to the source

Distribute computing to extract information as data move through the network in the most efficient way.



Daily/weekly adaptability requirements imply the need for programmability & portability of algorithms

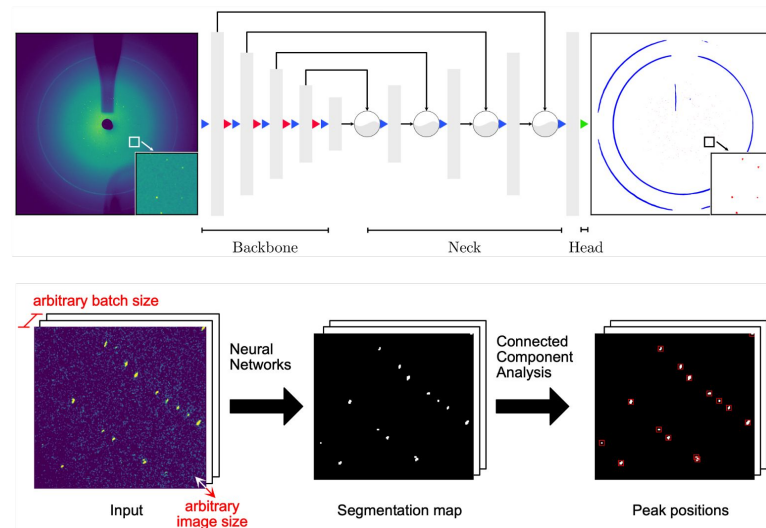
Produce Actionable Information for AI-enabled classification and feature extraction at the Edge



AI Classification of diffraction images in
FPGA for real time vetoing

Wang, C. et al., 2023 (<https://doi.org/10.48550/arXiv.2302.06895>)
Herbst, R., et al. 2023 (<https://doi.org/10.48550/arXiv.2305.19455>)

Cong Wang
Abhilasha Dave



1 MHz Autonomous Bragg Peak Finder
for online feature extraction

Wang, C. et al., 2023 (<https://doi.org/10.48550/arXiv.2303.15301>)

DAQ/Data Reduction Pipeline

Data Reduction Pipeline - Do More with Less Data

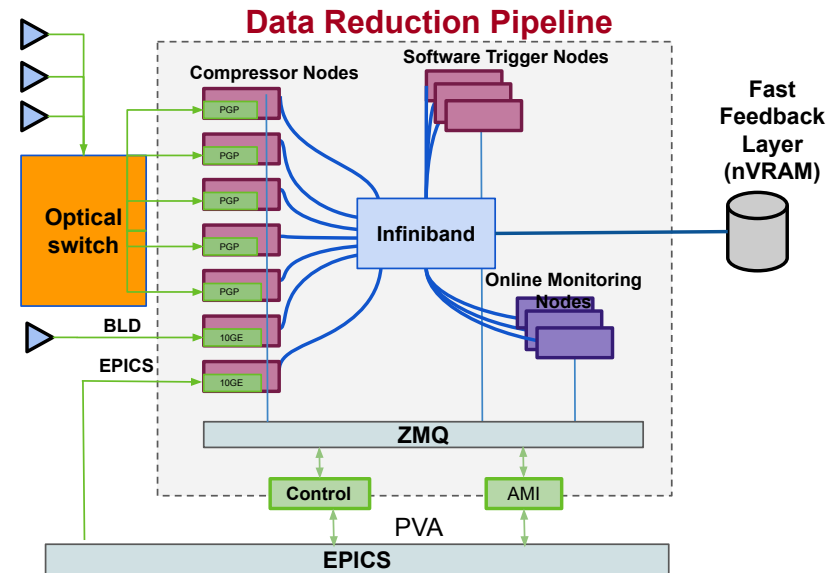
All data are equal, but some are more equal than others

On-the-fly data reduction

- Mitigates network, storage, computing bottlenecks
- Enables streaming to HPC
- Common algorithms, parameterized for adaptability

Software Trigger Nodes perform online event build collecting data from multiple detectors from same event.

Algorithms	TMO	RIX	XPP	TXI	DXS	MFX	CXI
Lossless	X	X	X	X	X	X	X
Veto					X	X	X
SZ Compression		X	X	X	X	X	X
Average image binning			X	X	X	X	
Pixel binning					X	X	X
ROI/Projection			X	X			
Angular integration and pie slicing			X	X		X	
Peak-finding/ threshold	X	X	X	X	X		

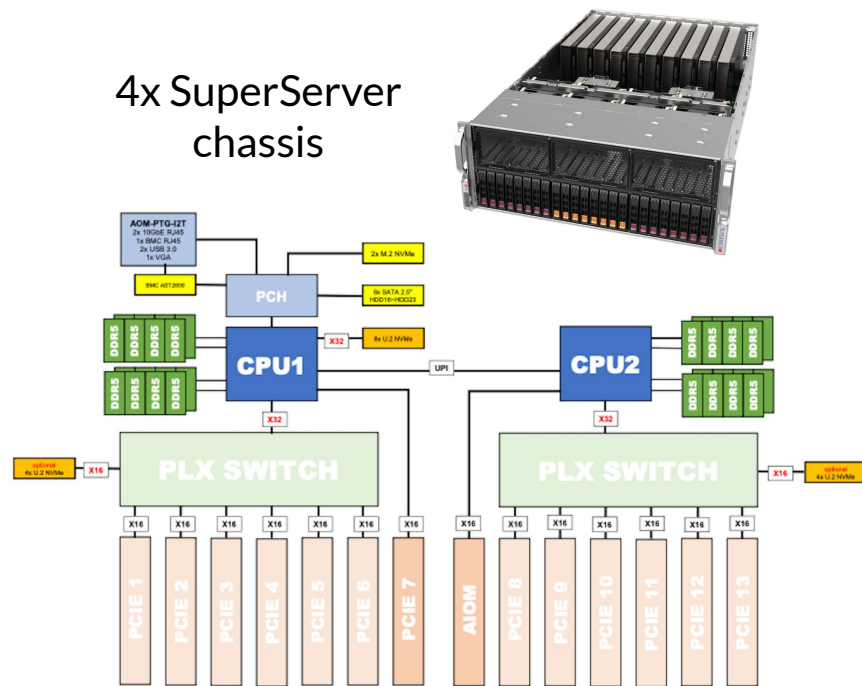


Save reduced data, and a pre-scaled fraction of raw data to allow us to verify that reduction does not affect physics results.

Data Reduction algorithms developed offline, then migrated online → portability is important

Data reduction for large, imaging high-rate area detectors

XPP ePixUHR 4 MP @ 35 kHz for LCLS-II-HE will produce ~ 280 GB/s on 240 fiber pairs

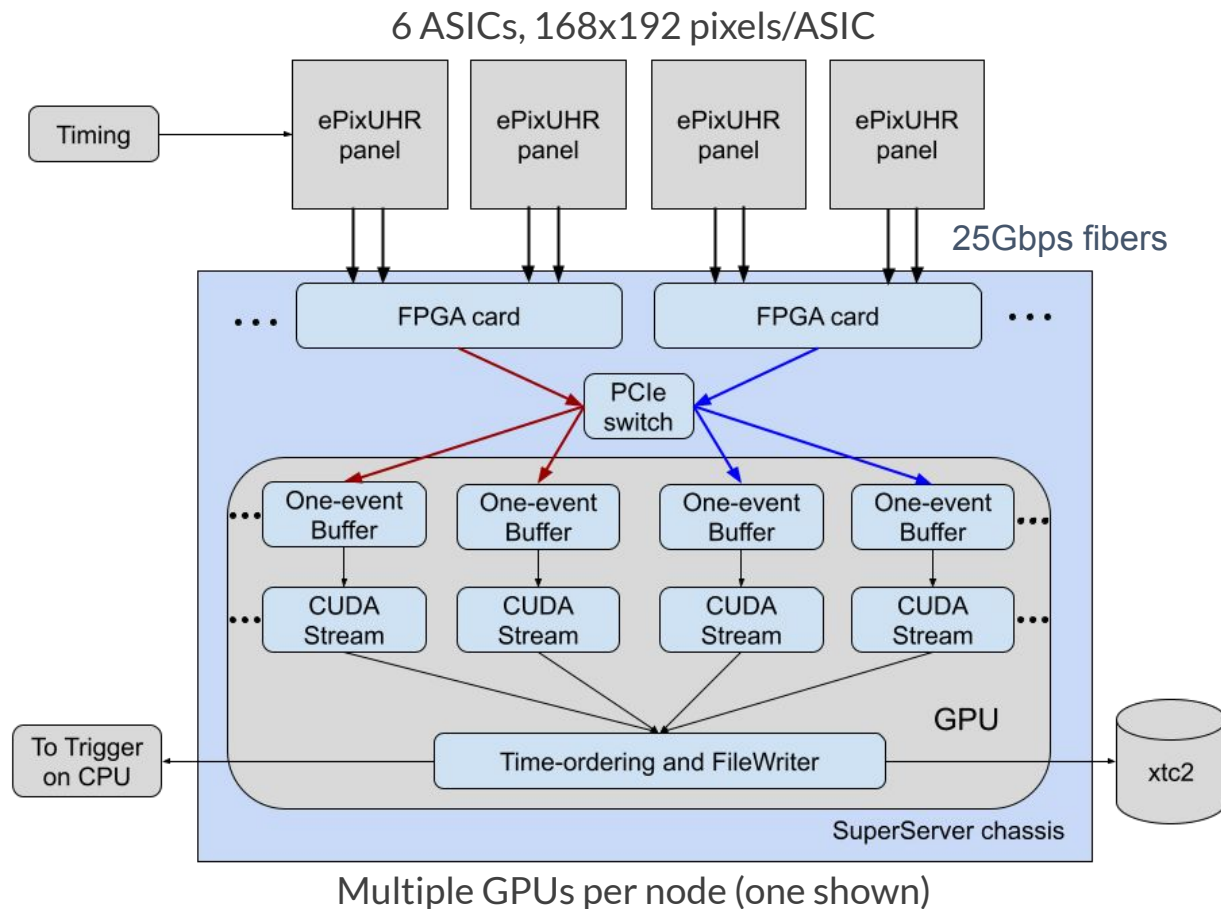


- **Detector protocols:**
 - PGP: SLAC high-speed FPGA protocol, based on Xilinx Aurora
 - UDP unicast/multicast
 - EPICS (TCP)
- Detectors with **widely varying readout rates** (1Hz thru 1MHz)
- **Firmware timestamping**, except for some detectors < 120Hz
- **Deadtime:** disable triggers when buffers are full

4x chassis with 24x FPGA boards, 24x H200 GPUs Worst case: 16 MP@35 kHz

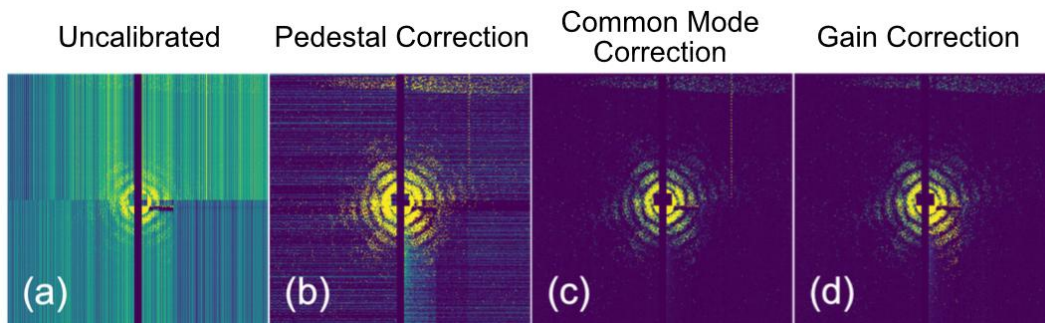
LCLS-II-HE XPP ePixUHR @ 35 kHz - GPU Readout Dataflow

- **Minimal online event-build:** event build a small fraction of events to implement veto
- Use libfabric to **abstract out infiniband/ethernet differences**
- **Reuse buffers:** no malloc/free per-event, but do copy the data from the driver (non-ideal) to get standard format



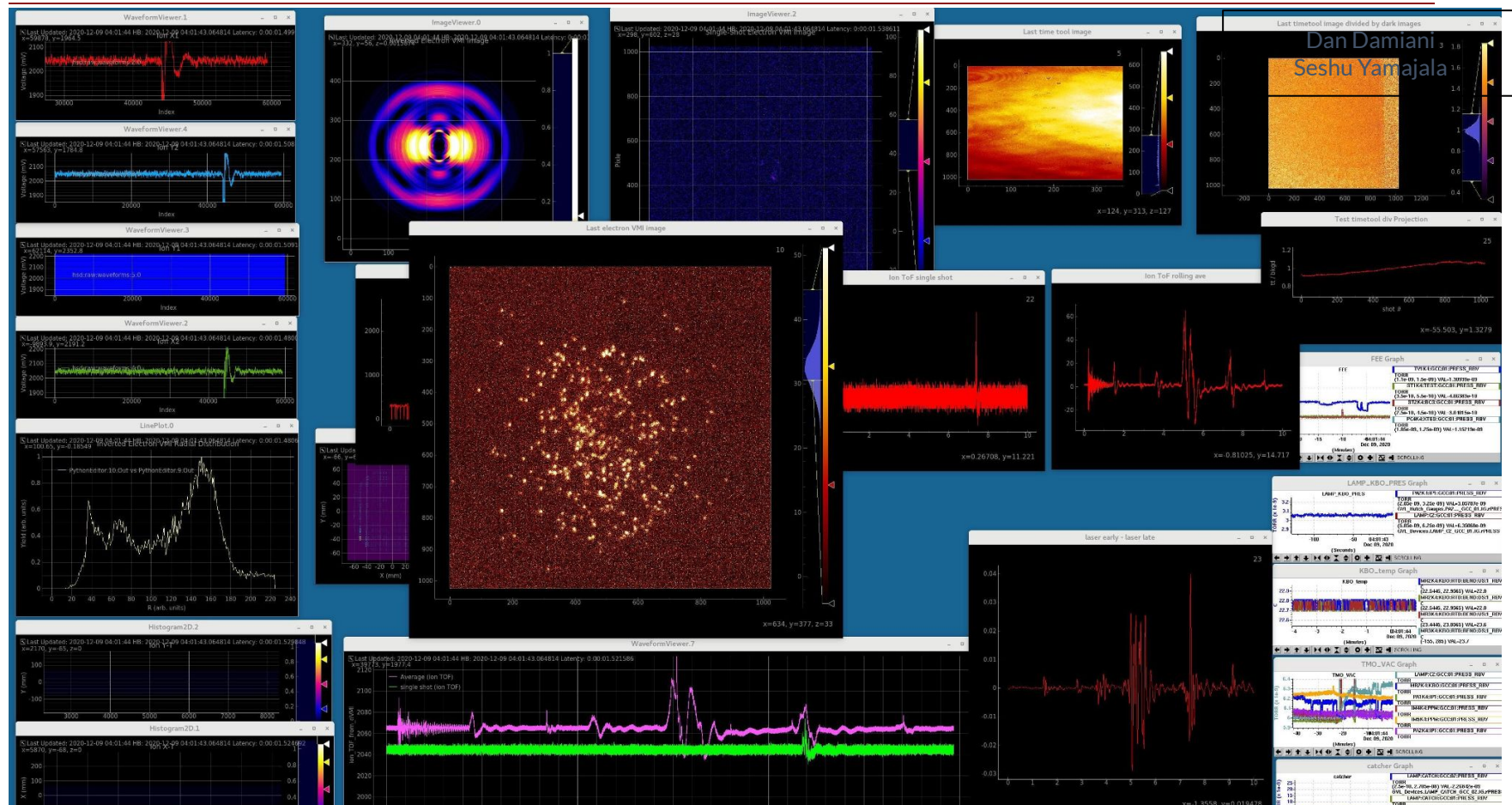
GPU Data Reduction Algorithm Performance

Prerequisite to compression algorithms below is calibration and sometimes binning and shot-filtering.



Algorithm	Rate	Notes
cuSZ lossy compression	91 GB/s	useful for SAXS/WAXS data
LC lossy compression	60 GB/s	4 GPU streams, useful for SAXS/WAXS data
pyFAI angular integration	1.7 GB/s	
Sparse-matrix angular integration	6.6 GB/s	
pyFAI peakfinder8	1.5 GB/s	Issue: requires bkgd calculation from all panels

Keep up with the Data Rate - Scalability and Parallelization



Data Format

xtc2: A Data Format for a Heterogeneous Pipeline

Data format is tightly packed and the same “on the wire” as it is in memory and in the file

- **In-memory format identical to persistent-storage format**
 - allows psana to run from shared memory and embedded in DAQ
- **No serialization required** when data is transmitted to remote machines
- **Natural support for fundamental types** (integers, floats...), arrays of fundamental types and variable-length data
- **Fields that describe the software/version** needed to interpret a particular block of data
- **Separate small metadata from large data** so it can be read more quickly. (e.g. fseek offsets of large data (useful for parallelization) and other small data that can be used to decide whether or not to pay the penalty of fetching the associated large data).
- **Easy python/C++ interface** (with no array copying when accessing)
- Natural support for **variable-length data**
- **Lightweight** (2500 lines of C++) with no dependencies

Raw Data Format: xtc2 allows for variable length data

xtc2 format is the same in-memory and file, enabling streaming

Datagram Header (per-event):

Timestamp (64-bit)

Datagram Type (32-bit)

Contains several nested (or concatenated) XTC headers
(eXtensible Tagged Container):

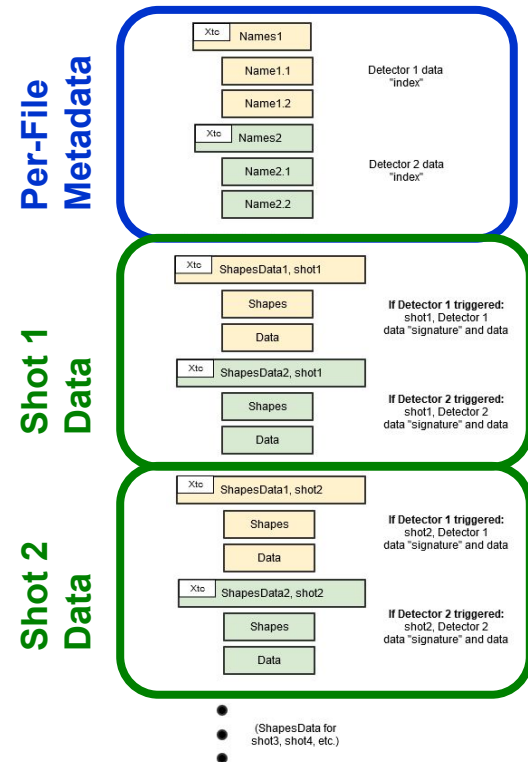
Damage (16-bit)

Src (32-bit)

TypeId (16-bit)

Extent (32-bit)

Payload (as long as Extent)



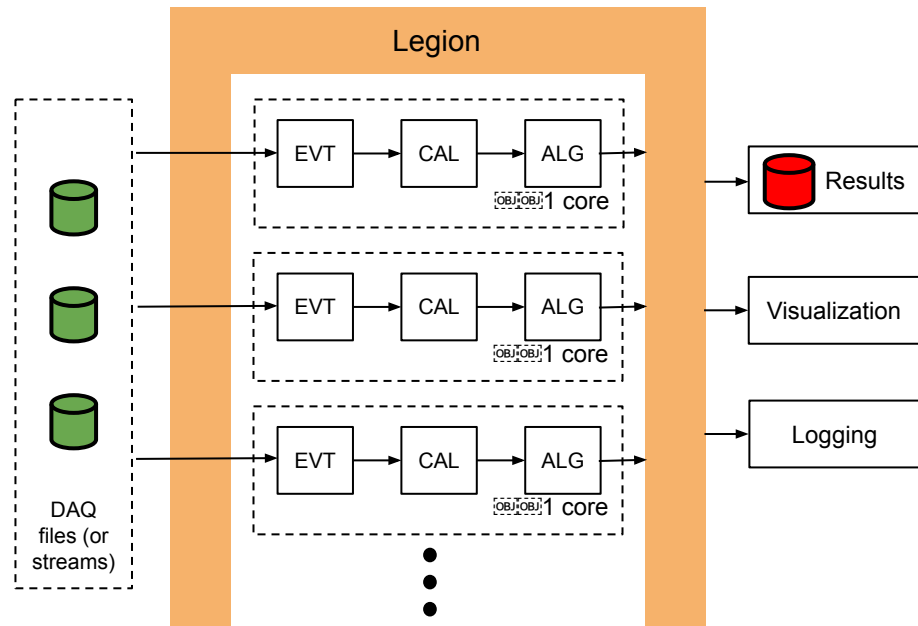
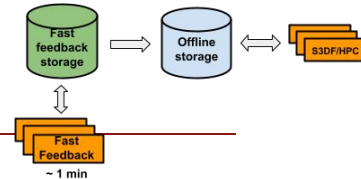
Analysis Framework

Offline Analysis Framework

psana - Photon Science Analysis framework for accessing data

psana provides parallelization, common algorithms, detector corrections, event-building, file format handling, visualization.

- Allows for real-time analysis, whether run online, in the fast feedback, or offline on 1 to 300,000+ cores
- psana uses MPI to provide:
 - a. Overlap of I/O and compute
 - b. Portable performance on new architectures
 - c. “Perfectly parallel” pattern: applications can be scaled, limited by data distribution from filesystem or shared-memory
- LCLS uses conda/spack to create releases that run on S3DF and remote HPC resources



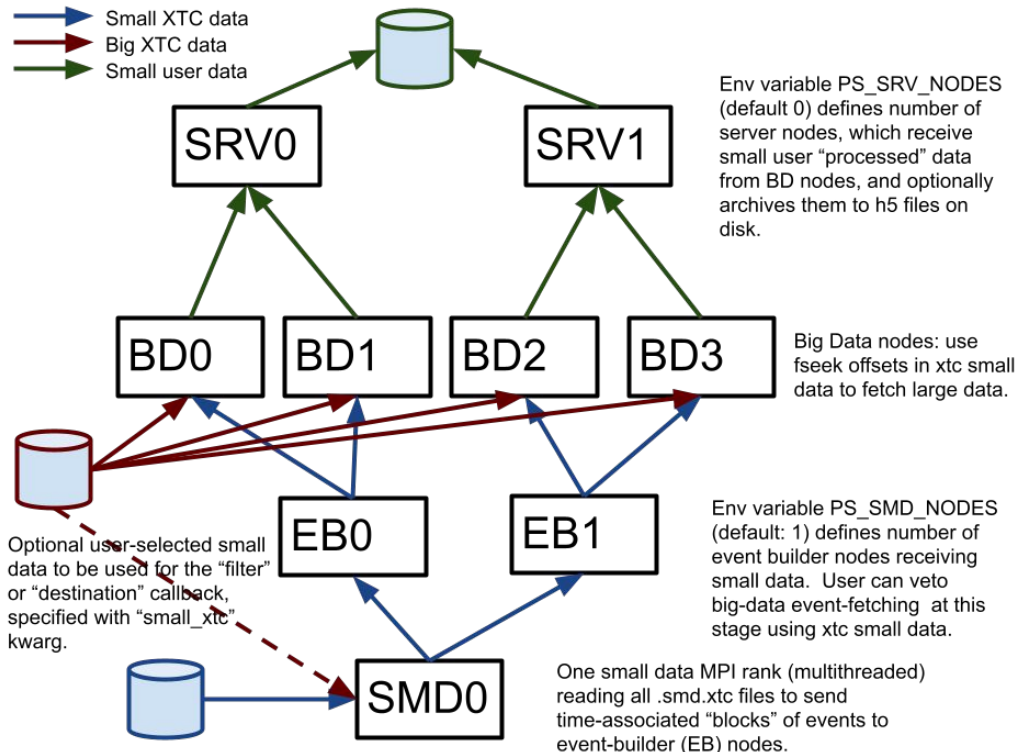
Reads science data from file or stream, distributes one event per core, performs detector calibrations and invokes science specific algorithms

Analysis (“psana”) big ideas

- **Limited data types (float, int, arrays):** raw data translation to Python uses C++ Python extension
- **“Perfectly Parallel”:** different events given to different cores
- **Memory management done with Python** reference counting with zero-copy of array data using `PyArray_SimpleNewFromData()`
- **Two-levels of Python interface:** messy raw data (uncorrected, hidden from user, for experts) and user-level “Detector” interface
- **Small data defined by users** at analysis time for filtering
- Analysis must work on individual area-detector panels, since **detectors have multiple segments**
- Detector **calibrations stored in MongoDB/GridFS** for read-only access at HPC centers
- **Same scripts work everywhere:** real-time and offline analysis
- **Hide unnecessary complexity** from users: MPI parallelization, HDF5 production, detector corrections
- **Integrating detector support** needed to correctly associate high-rate events with slow detectors
- **Live-mode** analyzes data while being written
- **Jump** to selected events (not available in live-mode)
- **Grafana** tools used **to understand/optimize high-rate performance** of analysis workflows

Generating calibrated, event-built HDF5 data products

psana uses MPI to distribute load; HDF5 product is the INPUT to user analysis pipeline



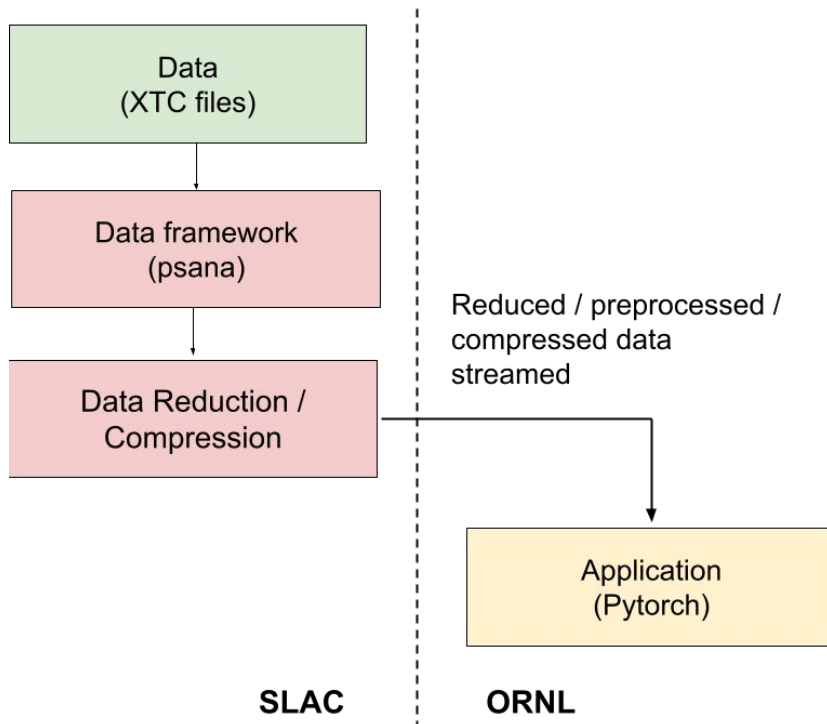
Note that this pattern hides complexity but uses MPI making psana more of an application than a library.

Any workflow built on top of it that thinks in can control distribution of tasks across nodes using MPI is going to have a bad time.

Remote-Location Data Processing: LCLStream

Data Streaming to Remote Facilities: LCLStream

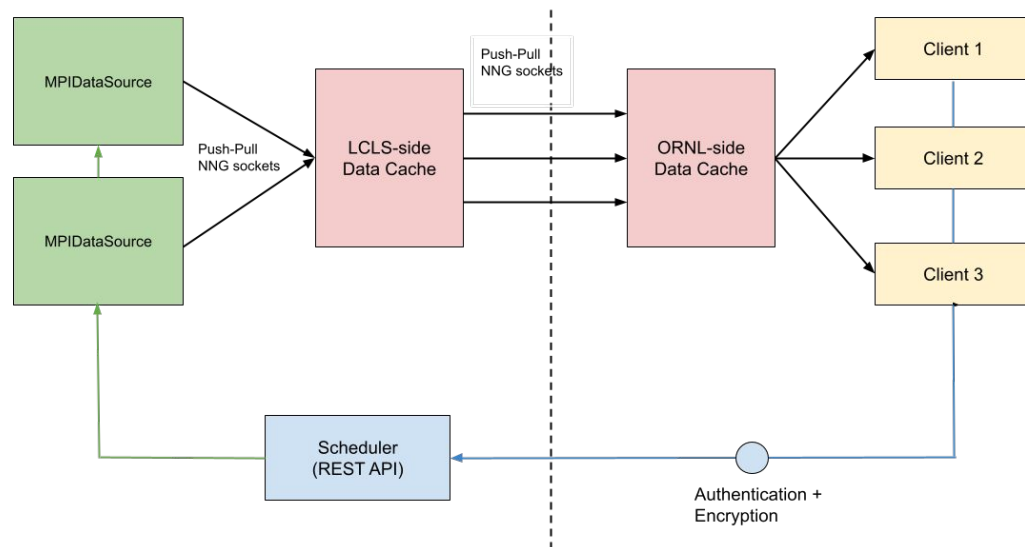
- Streamed data is reduced:
 - Only relevant data is transferred (detectors, hits)
- Streamed data is preprocessed:
 - Ready for scientific interpretation
- Streamed data is compressed:
 - Lossless compression



Remote-Location Data Processing: LCLStream

Data Streaming to Remote Facilities: LCLStream

- Streamed data is reduced:
 - Only relevant data is transferred (detectors, hits)
- Streamed data is preprocessed:
 - Ready for scientific interpretation
- Streamed data is compressed:
 - Lossless compression



Local Heterogeneous Computing using LCLStream

Separate:

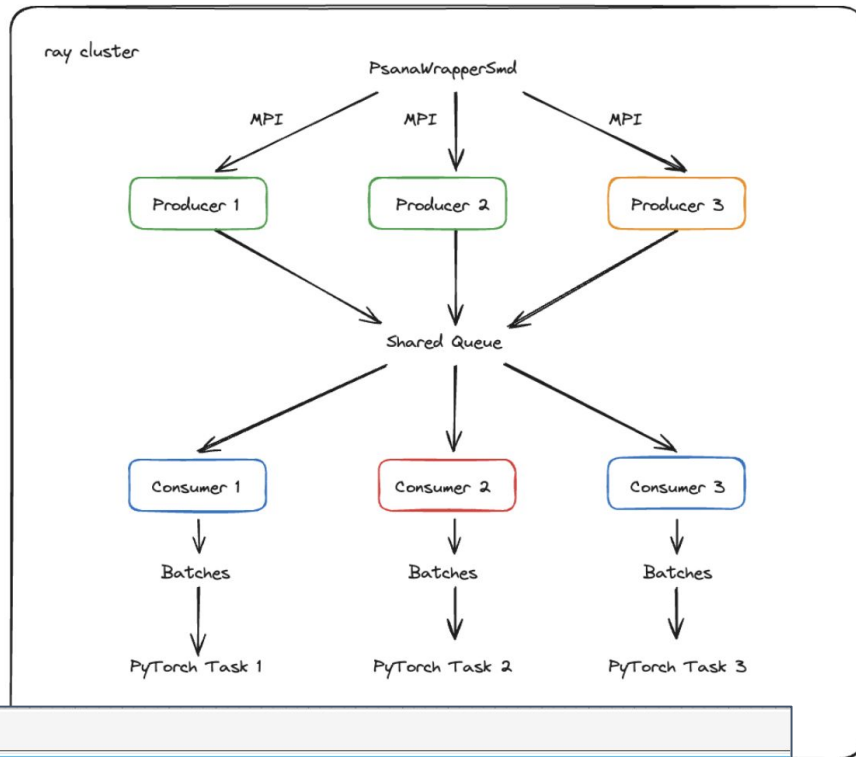
- Data reading (LCLStreamer - Psana)
- Data Processing (Processing code - Ray)

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

Cong Wang

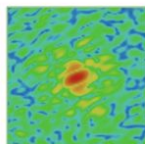


	Source	Target	Source	Target
	Src Server	Target Server	Src Server	Target Server
1	192.168.1.100	192.168.1.101	192.168.1.102	192.168.1.103
2	192.168.1.104	192.168.1.105	192.168.1.106	192.168.1.107
3	192.168.1.108	192.168.1.109	192.168.1.110	192.168.1.111
4	192.168.1.112	192.168.1.113	192.168.1.114	192.168.1.115
5	192.168.1.116	192.168.1.117	192.168.1.118	192.168.1.119
6	192.168.1.120	192.168.1.121	192.168.1.122	192.168.1.123
7	192.168.1.124	192.168.1.125	192.168.1.126	192.168.1.127
8	192.168.1.128	192.168.1.129	192.168.1.130	192.168.1.131
9	192.168.1.132	192.168.1.133	192.168.1.134	192.168.1.135
10	192.168.1.136	192.168.1.137	192.168.1.138	192.168.1.139
11	192.168.1.140	192.168.1.141	192.168.1.142	192.168.1.143
12	192.168.1.144	192.168.1.145	192.168.1.146	192.168.1.147
13	192.168.1.148	192.168.1.149	192.168.1.150	192.168.1.151
14	192.168.1.152	192.168.1.153	192.168.1.154	192.168.1.155
15	192.168.1.156	192.168.1.157	192.168.1.158	192.168.1.159
16	192.168.1.160	192.168.1.161	192.168.1.162	192.168.1.163
17	192.168.1.164	192.168.1.165	192.168.1.166	192.168.1.167
18	192.168.1.168	192.168.1.169	192.168.1.170	192.168.1.171
19	192.168.1.172	192.168.1.173	192.168.1.174	192.168.1.175
20	192.168.1.176	192.168.1.177	192.168.1.178	192.168.1.179
21	192.168.1.180	192.168.1.181	192.168.1.182	192.168.1.183
22	192.168.1.184	192.168.1.185	192.168.1.186	192.168.1.187
23	192.168.1.188	192.168.1.189	192.168.1.190	192.168.1.191
24	192.168.1.192	192.168.1.193	192.168.1.194	192.168.1.195
25	192.168.1.196	192.168.1.197	192.168.1.198	192.168.1.199
26	192.168.1.200	192.168.1.201	192.168.1.202	192.168.1.203
27	192.168.1.204	192.168.1.205	192.168.1.206	192.168.1.207
28	192.168.1.208	192.168.1.209	192.168.1.210	192.168.1.211
29	192.168.1.212	192.168.1.213	192.168.1.214	192.168.1.215
30	192.168.1.216	192.168.1.217	192.168.1.218	192.168.1.219
31	192.168.1.220	192.168.1.221	192.168.1.222	192.168.1.223
32	192.168.1.224	192.168.1.225	192.168.1.226	192.168.1.227
33	192.168.1.228	192.168.1.229	192.168.1.230	192.168.1.231
34	192.168.1.232	192.168.1.233	192.168.1.234	192.168.1.235
35	192.168.1.236	192.168.1.237	192.168.1.238	192.168.1.239
36	192.168.1.240	192.168.1.241	192.168.1.242	192.168.1.243
37	192.168.1.244	192.168.1.245	192.168.1.246	192.168.1.247
38	192.168.1.248	192.168.1.249	192.168.1.250	192.168.1.251
39	192.168.1.252	192.168.1.253	192.168.1.254	192.168.1.255
40	192.168.1.256	192.168.1.257	192.168.1.258	192.168.1.259
41	192.168.1.260	192.168.1.261	192.168.1.262	192.168.1.263
42	192.168.1.264	192.168.1.265	192.168.1.266	192.168.1.267
43	192.168.1.268	192.168.1.269	192.168.1.270	192.168.1.271
44	192.168.1.272	192.168.1.273	192.168.1.274	192.168.1.275
45	192.168.1.276	192.168.1.277	192.168.1.278	192.168.1.279
46	192.168.1.280	192.168.1.281	192.168.1.282	192.168.1.283
47	192.168.1.284	192.168.1.285	192.168.1.286	192.168.1.287
48	192.168.1.			

№	ИД	Имя	Возраст	Пол	Дата рождения	Дата смерти	Дата похорон
26	471 888	СМЕРД	3	м	05/01/2015 09:31:01	04/02/2015 08:38:24	03/01/2015
28	270 554	СМЕРД	3	м	05/01/2015 09:31:01	04/02/2015 08:38:24	03/01/2015
29	036 942	СМЕРД	8	м	05/01/2015 09:31:01	03/02/2015 08:11:48	03/01/2015
29	493 264	СМЕРД	5	м	05/01/2015 09:31:01	03/02/2015 08:11:48	03/01/2015
30	493 264	СМЕРД	5	м	05/01/2015 09:31:01	03/02/2015 08:11:48	03/01/2015



LCLS

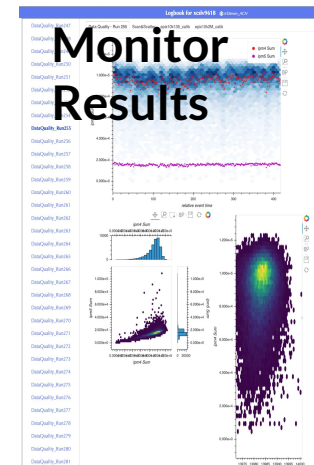
Data
Acqu

XRootD ~ 15 TB/day

Data Transfer Nodes



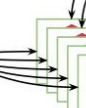
Automated Run Processing



Monitor Results

Users interact with data analysis in real-time

Workflow Coordination



CCLS



Cori Compute Nodes



Analysis results within minutes!

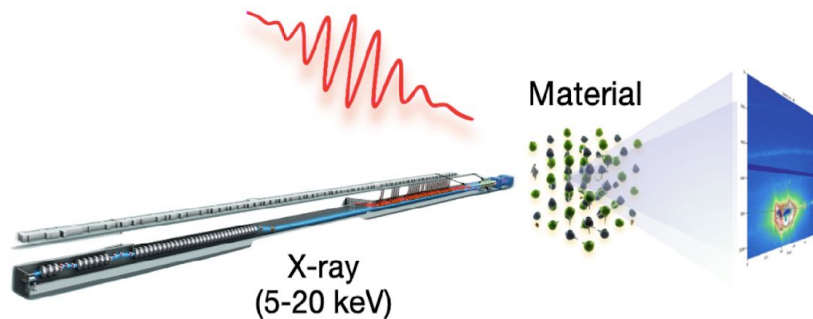
Real-time Analytics:

Using HPC to accelerate analysis for time-resolved experiments

Credit: Quynh Nguyen

cuPyNumeric for performant GPU-accelerated data analysis pipelines
replaces numpy in user code

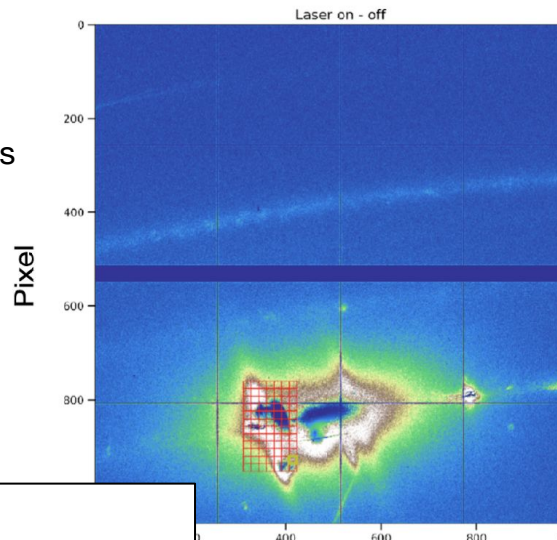
Laser excitation



10+ hours to ~4 minutes
222x Faster



Terabytes/hour
→ 0.1 - 1 Terabytes/s



GPU needs for 33 kHz running

- EOS:
90 nodes, 720 H100s
1060x speedup
- Perlmutter:
256 nodes, 1060 A100s
939x speedup
10+ hours → 90 s

Steps in the pipeline:

- Pedestal subtraction
- Gain correction
- ROI data reduction
- Event filtering (throw out bad events, select for desired qualities)
- Create HDF5/Data product input to user pipeline
- User-defined offline analysis (read from file/decompression penalty)



Dr. Irina Demeshko
Malte Foerster



Seshu Yamalaja
Dr. Quynh L Nguyen

Summary of Challenges and Opportunities

Every time hardware or software updates or a new opportunity (AI/ML) arises, the LCLS Data System platform, infrastructure, and hundreds of user workflows for data reduction, fast feedback, data quality monitoring, decision support, data management and data processing must also change!

Challenges:

- Increased reliance on sophisticated workflows and HPC; some workflows are “canned”, some not
- Growing divide between user capabilities and computing sophistication required to use HPC to execute complex workflows
- Heterogeneous pipelines (ASIC, FPGA, CPU, GPU, TPU, accelerators) help efficiently distribute computing tasks, but complicate performance and portability
- Metadata and data management at high rate with AI/ML require new techniques for data wrangling

Opportunities:

- Wide range of intelligent operations and low code data processing opportunities: Adaptive tagging of code/data; Workflows specified by visual/data-flow languages and/or learned by example; Semantic search and auto-complete for code/data; Speech-directed operations
- Code acceleration methods; advanced heterogeneous memory/processing systems; data file formats; advanced job scheduling (eg to account for I/O, real time queues, etc)
- Cognitive engineering methods can be applied to optimize UI and workflow designs