

EIC E0-E1 mini-workshop

Input to topical discussions

Torre Wenaus

Nuclear and Particle Physics Software (NPPS) Group Leader, BNL Physics Dept
ePIC Deputy S&C Coordinator

EIC E0-E1 mini-workshop
Apr 2 2025
BNL

Topics I'm addressing

- Software stack and computing infrastructure prototypes
 - Workflow management - near term objectives in two areas
 - for streaming
 - for other use cases (e.g. production)
- Computing services at the two labs and their roles in prototyping and production
 - NP Rucio at BNL
 - NP PanDA at BNL
 - PanDA-accessible processing resources at JLab?
- Data Center infrastructure to support Echelons 0 and 1

Workflow management near term objectives: AID2E

- The first use case for NP PanDA@BNL and NP Rucio@BNL has been the AID2E project in scalable distributed AI-based EIC detector design optimization
 - Bayesian optimization across several design parameters and detector subsystems, making for complex optimizations with large processing demands
- PanDA and its iDDS fine-grained workflow management component bring the ‘scalable distributed’ capability
 - Building on prior work on scalable ML services like hyperparameter optimization
- Rucio brings the data/metadata management needed to scale up the full simulation that’s the basis for the detailed design optimization
- We now have both PanDA & Rucio at BNL! Yay and thank you!
- Dealing with the Rucio gotcha of still requiring X509. Increases complexity and startup time for new users significantly.

Workflow management near term objectives: ePIC streaming R&D

- ePIC has a fairly well developed view of its streaming computing model expressed in the Oct 2024 Version 2 of the report describing it
- attention is turning to R&D testbeds to explore aspects of it
- One testbed defined in the last few weeks is workflow management, which will be addressed by a PanDA + Rucio (BNL instances) testbed focused (at least initially) on Echelon 1 near real time processing
- Most interestingly, encompassing the unique ePIC aspects of its streaming computing model
 - Data arrival as 'super timeframes' from DAQ, Rucio registration, archiving, Rucio-triggered PanDA processing of newly arrived data
 - Streaming processing as distinct from conventional run/job/file based processing
 - Low-latency workflows for continuous monitoring of the accelerator and detector status
 - Detector/data state model describing what is on the stream at any given time
 - Orchestration of calibration workflows as well as collision data
- Context for the R&D is the [WFMS requirements](#), draft document out this week
- With the infrastructure established for AID2E, we have what we need to get going!
- We're eager to explore joint work with JLab, e.g. combining the testbeds

Workflow management near term objectives: ePIC production

- At the same time that WFMS testbed R&D for ePIC streaming workflows has been gestating, ePIC production folks have begun to take an interest in how PanDA might be useful to ePIC production in the near term
- If there is such interest, it is at least as high a priority as the testbed (higher, I think)
- NB Rucio is already in use in production
- Anil Panta at JLab and Sakib Rahman have begun to look at PanDA integration with ePIC production
- Interesting twist: it means PanDA@BNL working with Rucio@JLab (the instance used in ePIC production), as well as with Rucio@BNL for testbed R&D
- Motivations
 - Monitoring is a principal one. Anil already likes the PanDA monitor a lot :-)
- Approach
 - Treating PanDA as an alternative submission target to HTCondor. 'PanDA as the batch system' enables monitoring and the rest
 - Requires PanDA submitting workloads to OSG glide-ins, demonstrated a few months ago by Xin
- Important we advance and support this with priority!

Computing services & Roles

- NP Rucio at BNL
 - It's there, it's complete now, it works, it has the storage endpoint(s) we need for now
 - It would be wonderful if we got rid of X509
- NP PanDA at BNL
 - It's there, it's complete, it works, it's in use (AID2E now, ePIC WFMS testbed soon, ePIC production soon?)
 - Fully token based!
 - The accompanying OpenSearch/Grafana analytic services are needed and appreciated
- PanDA-accessible processing resources at JLab?
 - Explore this in the context of PanDA for ePIC production?
 - AID2E (JLab is an AID2E collaborator)? We've asked JLab if they're interested in enabling JLab Rucio for AID2E PanDA processing
- Data Center infrastructure to support near term activities
 - We now have the services we need, as we scale our activities resource availability will grow as an issue
 - EIC processing purchase through program development will help a lot
 - We haven't budgeted AID2E funds for processing, at least at BNL. Using DOE HPCs is part of the project (CSI contribution)

Storage for ePIC Simulation Production Campaigns

XRootD is the official data access protocol for ePIC simulation production. Backend storage choices should be easily scalable and work well with Rucio.

- **Current Storage Resources**

- JLab:
 - Lustre: Used for main simulation outputs; 500 TB limit.
 - ZFS: Used for small log files; 500 TB limit.
 - Tape 1 PB available for long-term archiving based on data retention policy (to be finalized in April, 2025)
 - Write backs handled through X509 authentication in Rucio.
- BNL
 - dCache: Not used for production campaigns but integrated as an RSE with JLab Rucio

- **Storage Consumption Statistics**

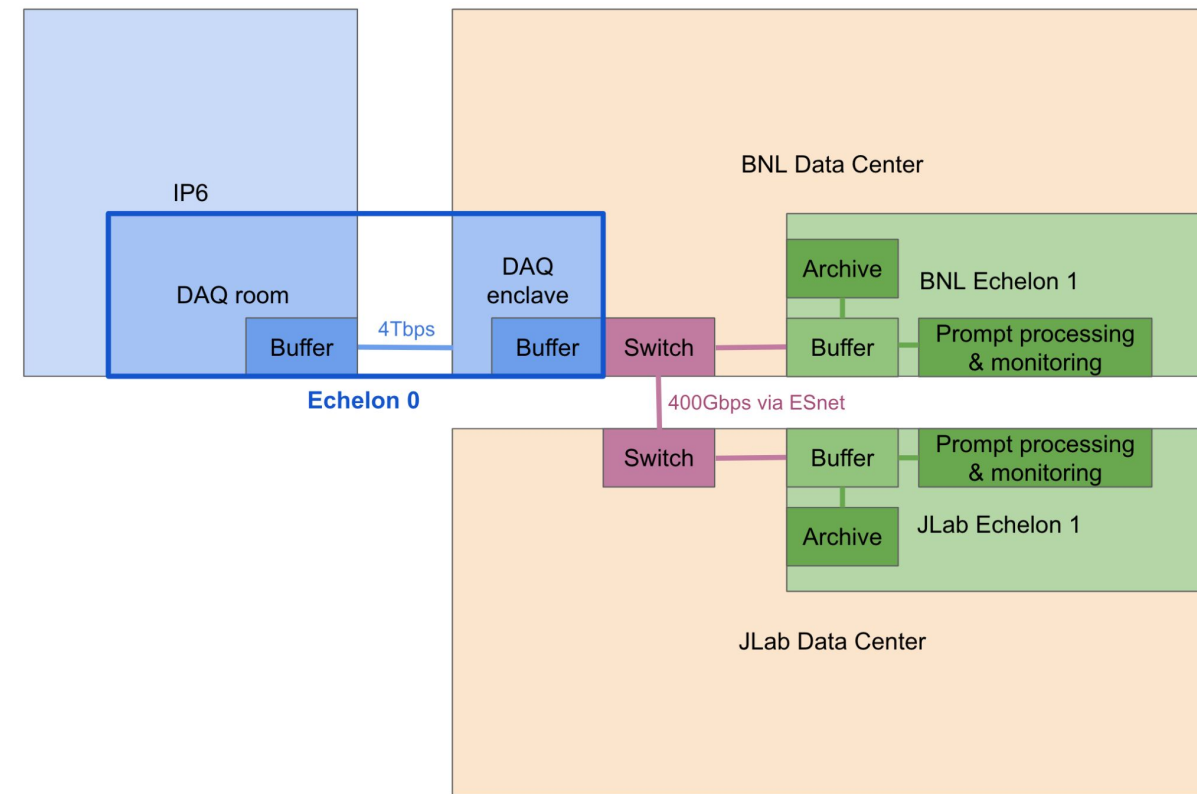
- 572 TB of output files on Lustre (currently exceeding soft quota). Older campaign files to be moved to tape storage as space fills up (data retention policy to be finalized by April 2025).
- January and February, 2025 campaigns: 37 and 42 TB respectively
- April 2025 and near term: Expect at least 2x storage consumption per campaign if running both with and without backgrounds.

Supplementary

End of discussion-seeding materials, the rest is to pull up as needed in the discussion

Echelon 0 (DAQ) to Echelon 1

- We began last fall to work out how data will flow from Echelon 0 to Echelon 1 for prompt consumption
- It's all under discussion; nothing I say here is decided, it's merely my perspective
- DAQ expected to leverage data center resources through a 'DAQ enclave' at SDCC
- Symmetric BNL, JLab Echelon 1 facilities downstream of DAQ are equally capable of performing Echelon 1 workflows
 - The model expresses the desired symmetry
 - Exact E1 roles will be decided by ePIC + facilities
- What's in the data stream sent from DAQ?
 - **Time frames**, each containing all detector data in a time interval, are built in DAQ
 - Time frames are aggregated in **super time frames** (STFs) which are sent out of DAQ to E1
 - Together with a small non-event data component (e.g. slow controls)
 - Data arrives at Echelon 1s in a ~1 week deep disk buffer
 - Workflows consume data from that buffer, first and foremost:
 - Archiving the full stream to tape (at least in early years)
 - Prompt processing, monitoring for a rapid view of data/detector integrity and quality



Time Frames and Super Time Frames

- Each time frame aggregates all detector data within a time window of $\sim 0.6\text{ms}$
- Super Time Frame (STF) is a contiguous set of ~ 1000 time frames
 - Within a STF the TFs are time-ordered, as required for reconstruction
 - No overlaps between time frames, so cannot reconstruct the edges
 - Make them large enough that losses from the edges are negligible
- This $O(1\text{s})$, 2GB STF data unit is an appropriate granularity for Echelon 1 processing
 - Short enough for prompt processing
 - Long enough for tractable bookkeeping and file size
- The STFs are *not* (required to be) time-ordered (facilitated by no overlaps)
 - Friendly to distributed computing and parallel orchestration
- **STF is the atomic unit for ePIC data processing**
 - (Such is the present ~ 4 month old thinking)

Streaming orchestration at Echelon 1s

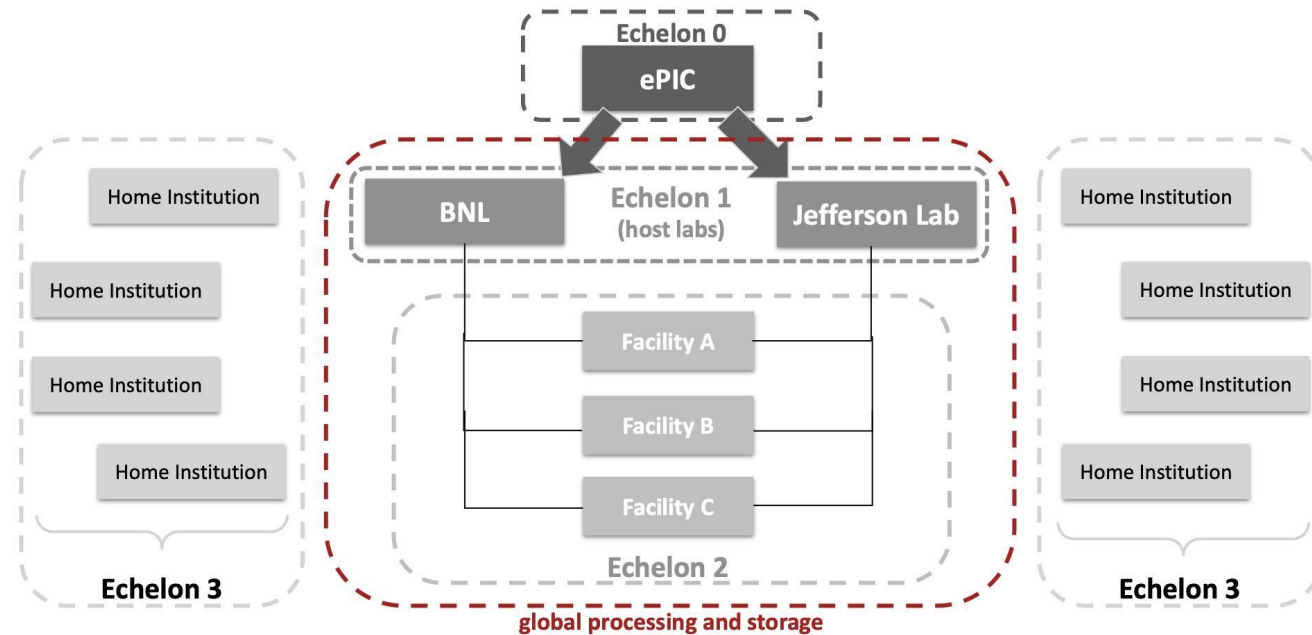
- Each E1 receives and archives 100% of the raw data, packaged identically: each E1 has a full set of super timeframe (STF) replicas
- Prompt processing is the other primary Echelon 1 role
- How will prompt processing be performed? What distributed computing systems will manage the processing at Echelon 1s and downstream?
 - The data management system is decided: Rucio
 - The workflow/workload management system evaluation and decision will take place (we expect) over the next 6mo or so
 - ePIC is very good at doing a considered process of requirements-based surveying, evaluating, and choosing. They've chosen the elements of their software stack this way
 - Key4HEP framework with Jana2 replacing Gaudi; Podio based EDM, etc.
 - Their decision process doesn't demand prototyping, but for PanDA, we do plan to prototype as available effort permits
 - As well as addressing the requirements document ePIC will produce (for which the [CMS PanDA eval Q&A](#) is useful input)
 - We are planning for success: demonstrate PanDA + Rucio working together as designed, for the first time outside ATLAS
 - (in Rubin they don't work together; Rubin middleware sits in the middle)

Prototyping prompt streaming with Rucio + PanDA

- Rucio is operating at JLab, BNL close to completing an instance for EIC use
 - Rucio used in ePIC production jobs for the first time this week
- PanDA prerequisite has been a PanDA instance at BNL, now in place thanks to Xin, PanDA team, SDCC
 - (Parenthetical: it's through the PanDA@BNL for NP work that we set up OSG + GlideinWMS support, essential to CMS and spotted by them immediately in the last PanDA community meeting, towards their [PanDA evaluation](#))
 - Client #1 for BNL PanDA is our AID2E AI-based detector design optimization project funded by DOE NP
 - ePIC PanDA is riding on its coat tails
- Other workflow management system participants in the evaluation TBD
 - DIRAC(x) is the natural, expected as part of the requirements eval, probably not prototyping?
 - NPPS also works/develops with DIRAC(x), we did the Rucio integration for Belle II; we have the most DIRAC experience in ePIC but we're not going to work two evaluations!
- ePIC streaming computing model WG met last week with two long-time friends of this room
 - Cedric assessed Rucio for the STF-based streaming model, all seems good
 - Dan Van Der Ster assessed object stores as the basis for STF storage; (Ceph) S3 seems good
- For the next weeks: develop a prototyping plan!

The ePIC streaming computing model

- The two-host-lab organization motivates the ‘butterfly’ model: BNL and JLab are symmetric peers
 - They avoided ‘Tier’ because 0 and 1 levels differ from their LHC meaning
- Globally distributed community and compute
 - Like LHC, distributed computing is essential
 - Pledged and opportunistic: 80% OSG today
- ePIC’s plan is to draw on existing experience and tools from LHC & elsewhere
 - while addressing the unique aspects of its streaming computing model



Echelon 0: ePIC experiment, DAQ system

Echelon 1: Two host labs, two primary ePIC computing facilities

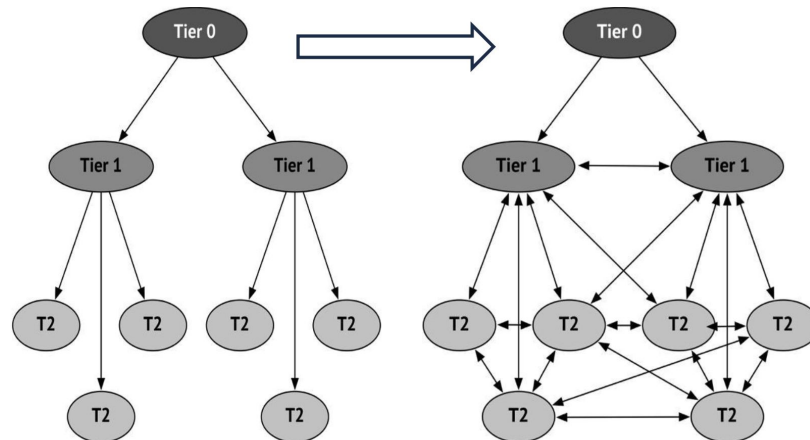
Echelon 2: Global contributions leveraging commitments to ePIC computing from universities and labs domestically and internationally

Echelon 3: Supporting the analysis community where they are at their home institutes, primarily via services hosted at E1s/E2s

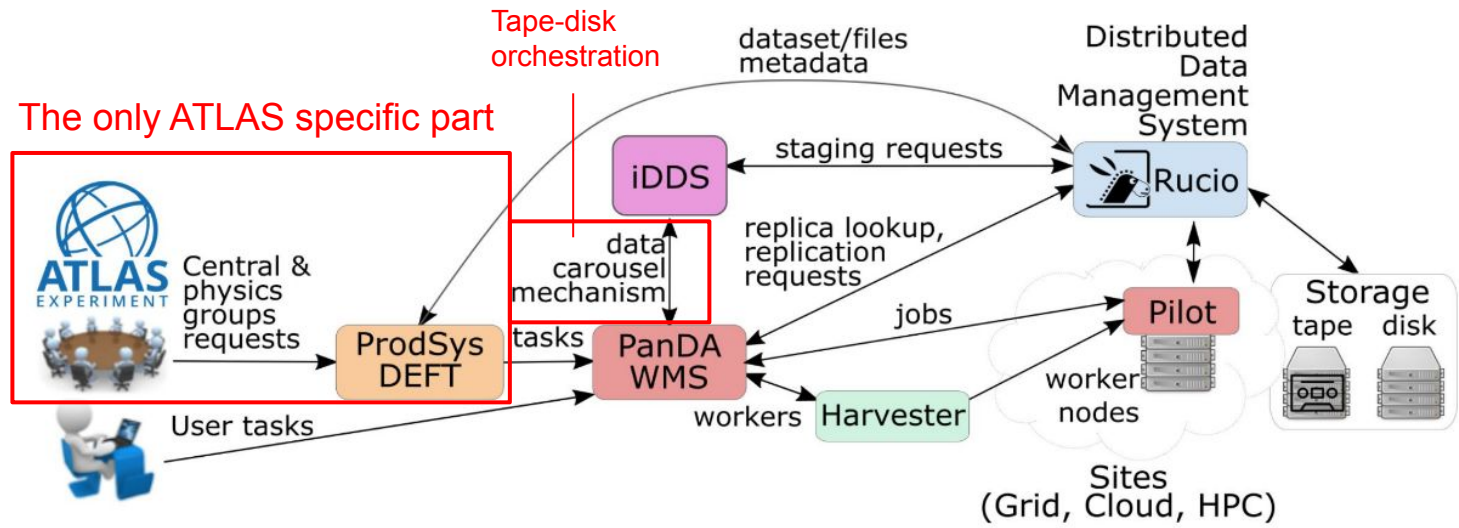
Lesson from LHC: Mesh outperforms hierarchy

Fully interconnected facilities flexibly and efficiently serve many roles

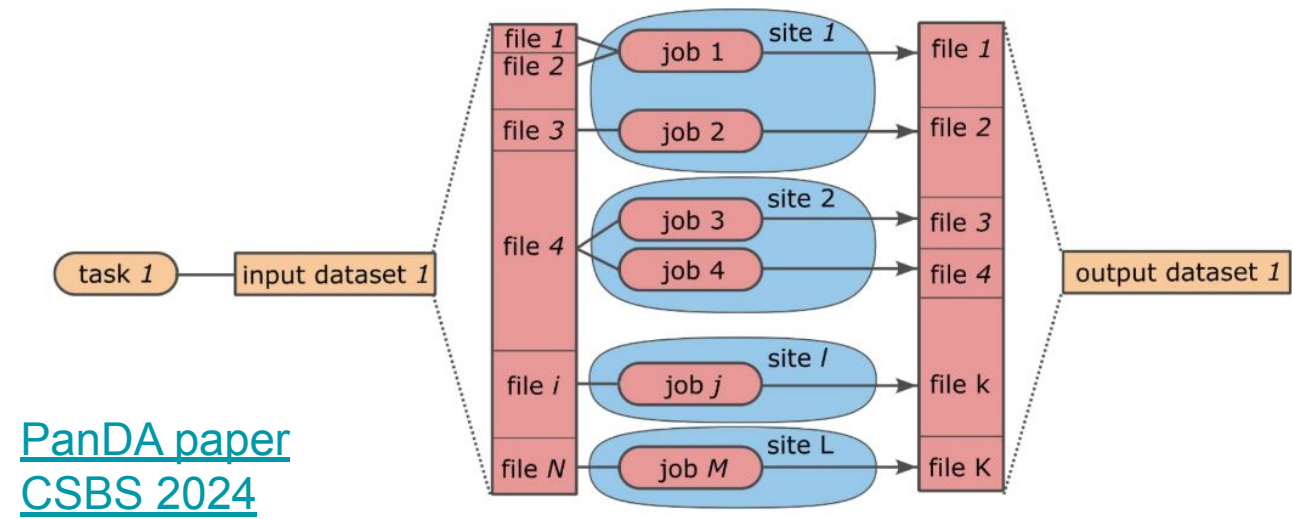
A hierarchy is *more complex* and less efficient & flexible



PanDA is a good fit, including fine grained orchestration

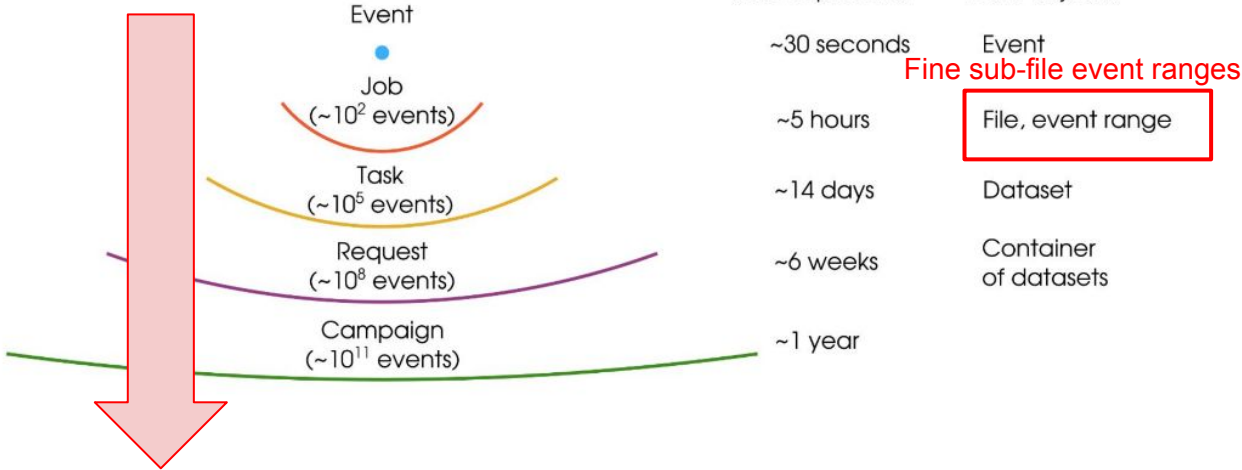


- **DEFT**: Database Engine For Tasks
- **PanDA**: Production AND Distributed Analysis System
- **Harvester**: resource-facing service between the PanDA and collection of pilots
- **Pilot**: the execution environment on a worker node **Fine grained workflow orchestration**
- **iDDS**: Intelligent Data Delivery System
- **Rucio**: Distributed Data Management System



[PanDA paper CSBS 2024](#)

Granularity suited to the use case



Computing use cases and their Echelon distribution

| Use Case | Echelon 0 | Echelon 1 | Echelon 2 | Echelon 3 |
|---------------------------------------|-----------|-----------|-----------|-----------|
| Streaming Data Storage and Monitoring | ✓ | ✓ | | |
| Alignment and Calibration | | ✓ | ✓ | |
| Prompt Reconstruction | | ✓ | | |
| First Full Reconstruction | | ✓ | ✓ | |
| Reprocessing | | ✓ | ✓ | |
| Simulation | | ✓ | ✓ | |
| Physics Analysis | | ✓ | ✓ | ✓ |
| AI Modeling and Digital Twin | | ✓ | ✓ | |

Prompt = rapid low-latency processing

Prompt processing of newly acquired data typically begins in seconds, not tens of minutes or longer

| Assumed Fraction of Use Case Done Outside Echelon 1 | |
|---|-----|
| Alignment and Calibration | 50% |
| First Full Reconstruction | 40% |
| Reprocessing | 60% |
| Simulation | 75% |

- **Echelon 1s uniquely perform the low-latency streaming workflows consuming the data stream from Echelon 0**
 - Archiving, monitoring, prompt reconstruction, rapid diagnostics
- There’s been discussion over whether Echelon 2s have a role in low-latency streaming processing
 - We shouldn’t exclude it in the design, but doing it in year 1,2,... is very unlikely
 - Message from reviewers (e.g. Simone Campana) was don’t plan to do this unless you really have to
- **Ensure the E1s have sufficient processing power for the low-latency workflows**
- **Apart from low-latency streaming, Echelon 2s are full participants in the use cases**
 - Resource requirements model assumes a substantial role for Echelon 2

Computing resource needs and the implications

- See [Markus Diefenthaler's talk this week](#) 24-25 for the numbers behind the numbers
- O(1M) core-years to process a year of data, above ATLAS scale today
 - Optimistic constant-dollar performance gains would reduce the numbers about 5x
 - Based on current LHC measure of 15%/yr
 - But the trend is towards lower gains per year
- Whatever the gains over time, the processing scale is substantial
 - Motivates attention to leveraging distributed and opportunistic resources from the beginning
- ~400PB/yr storage scale, also above ATLAS accumulation rate today
 - Archival storage (probably tape) will play a role
 - Motivates attention to data carousel type orchestration

Estimated needs for ePIC Phase I nominal year (circa 2034)

| Processing by Use Case [cores] | Echelon 1 | Echelon 2 |
|---------------------------------------|----------------|----------------|
| Streaming Data Storage and Monitoring | - | - |
| Alignment and Calibration | 6,004 | 6,004 |
| Prompt Reconstruction | 60,037 | - |
| First Full Reconstruction | 72,045 | 48,030 |
| Reprocessing | 144,089 | 216,134 |
| Simulation | 123,326 | 369,979 |
| Total estimate processing | 405,501 | 640,147 |

| Storage Estimates by Use Case [PB] | Echelon 1 | Echelon 2 |
|---------------------------------------|------------|------------|
| Streaming Data Storage and Monitoring | 71 | 35 |
| Alignment and Calibration | 1.8 | 1.8 |
| Prompt Reconstruction | 4.4 | - |
| First Full Reconstruction | 8.9 | 3.0 |
| Reprocessing | 9 | 9 |
| Simulation | 107 | 107 |
| Total estimate storage | 201 | 156 |

More info

- [ePIC collaboration meeting in Jan 2025](#)
 - [Markus Diefenthaler's S&C report](#)
- [EIC website](#)
- [EIC/ePIC github](#) (ePIC + common software)
- [ePIC Streaming Computing Model Report](#) (currently V2)
- [ePIC streaming computing model WG meeting notes](#)
- [Science Requirements and Detector Concepts for the Electron-Ion Collider: EIC Yellow Report](#)