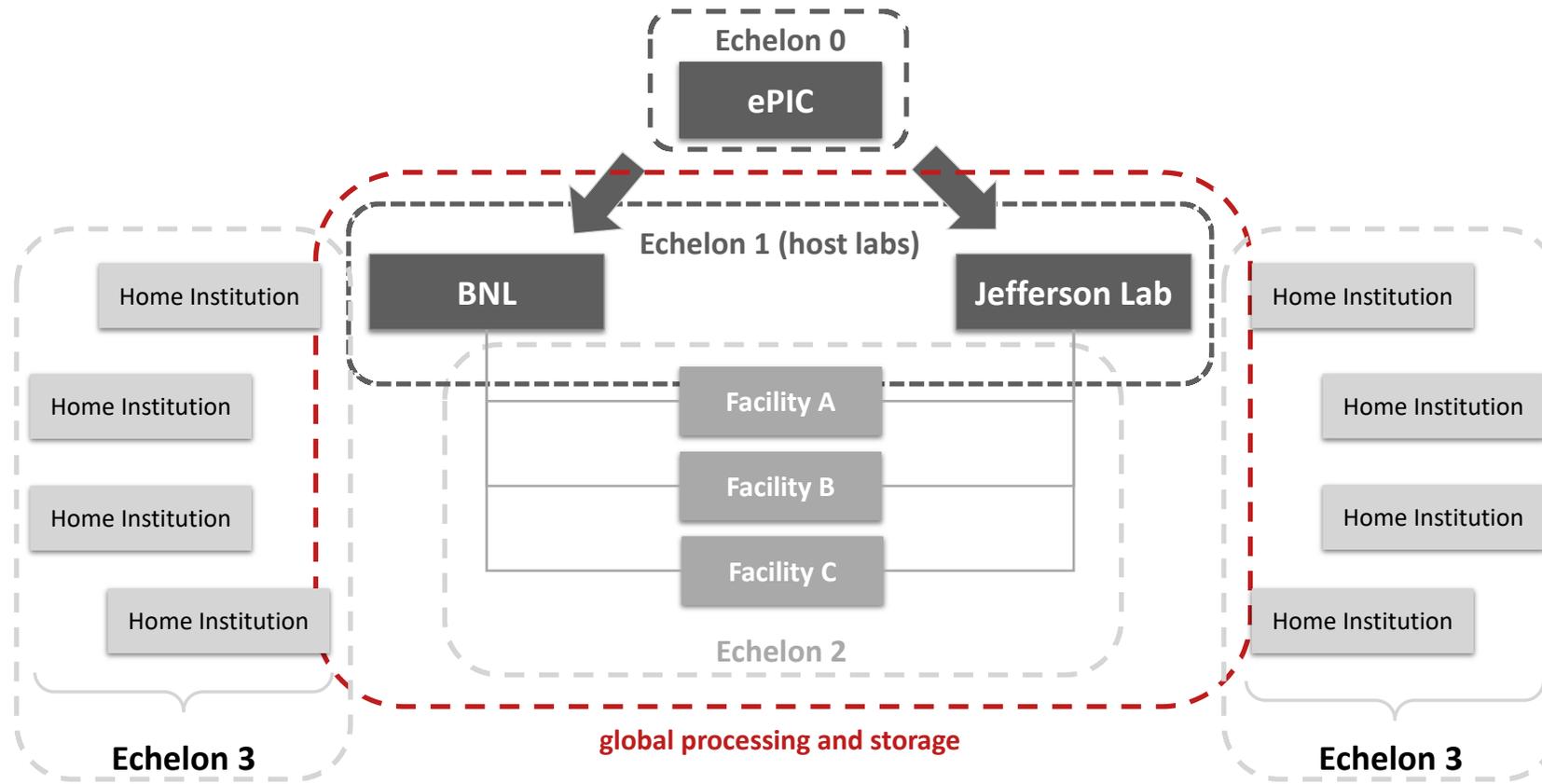


Echelon 2 Roles, Categories, and Resource Requirements



Markus Diefenthaler (Jefferson Lab)

The ePIC Streaming Computing Model

ePIC Software & Computing Report

<https://doi.org/10.5281/zenodo.14675920>

The ePIC Streaming Computing Model Version 2, Fall 2024

Marco Battaglieri¹, Wouter Deconinck², Markus Diefenthaler³, Jin Huang⁴, Sylvester Joosten⁵, Dmitry Kalinkin⁶, Jeffery Landgraf⁴, David Lawrence³ and Torre Wenaus⁴
for the ePIC Collaboration

¹Istituto Nazionale di Fisica Nucleare - Sezione di Genova, Genova, Liguria, Italy.

²University of Manitoba, Winnipeg, Manitoba, Canada.

³Jefferson Lab, Newport News, VA, USA.

⁴Brookhaven National Laboratory, Upton, NY, USA.

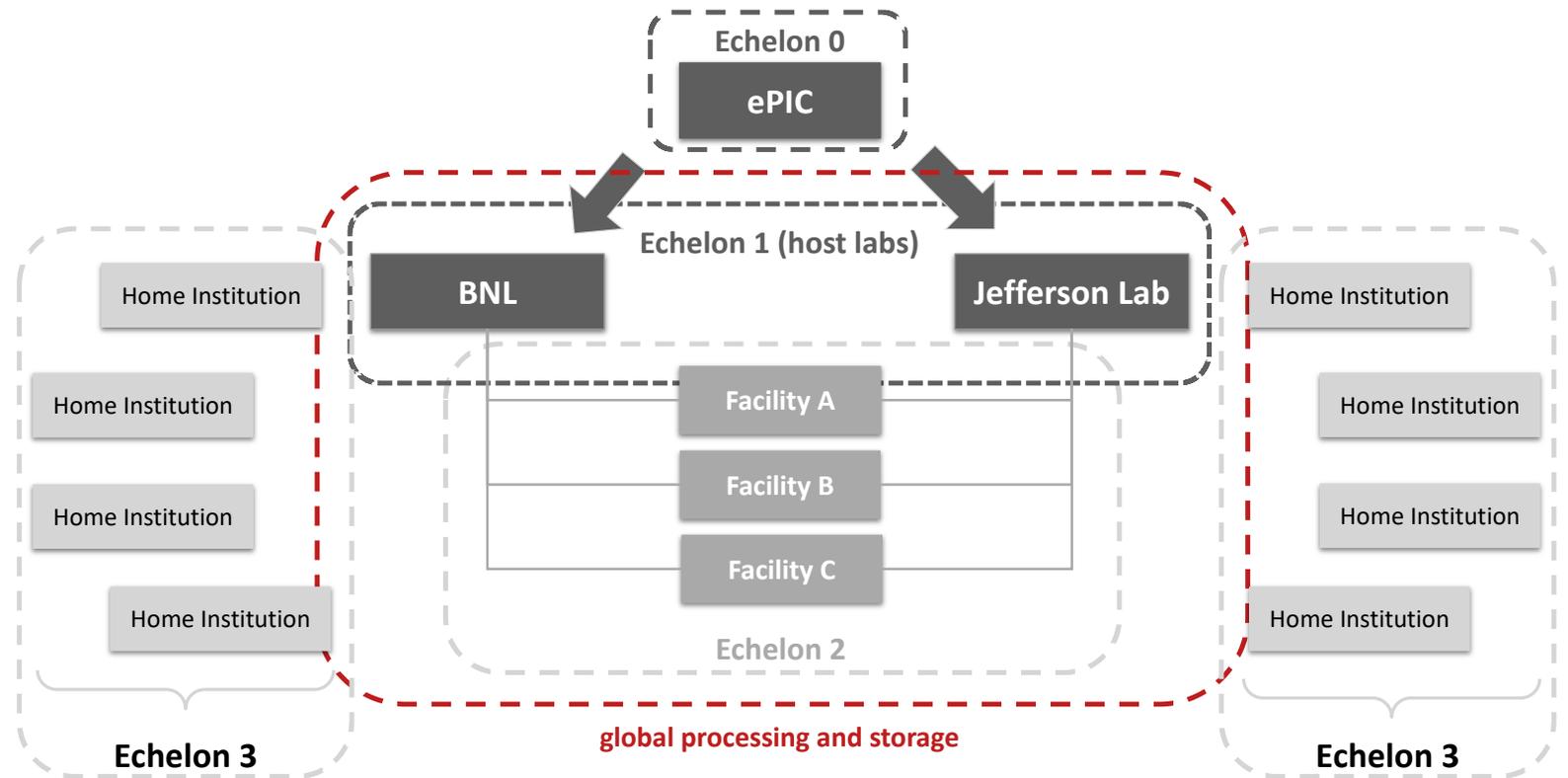
⁵Argonne National Laboratory, Lemont, IL, USA.

⁶University of Kentucky, Lexington, KY, USA.

Abstract

This second version of the ePIC Streaming Computing Model Report provides a 2024 view of the computing model, updating the October 2023 report with new material including an early estimate of computing resource requirements; software developments supporting detector and physics studies, the integration of ML, and a robust production activity; the evolving plan for infrastructure, dataflows, and workflows from Echelon 0 to Echelon 1; and a more developed timeline of high-level milestones. This regularly updated report provides a common understanding within the ePIC Collaboration on the streaming computing model, and serves as input to ePIC Software & Computing reviews and to the EIC Resource Review Board. A later version will be submitted for publication to share our work and plans with the community. **New and substantially rewritten material in Version 2 is dark green.** The present draft is preliminary and incomplete and is yet to be circulated in ePIC for review.

1



We developed the ePIC Streaming Computing Model to accelerate the pace of discovery and enhance scientific precision through improved management of systematic uncertainties. The model is documented in a detailed report and was reviewed during the 2023 and 2024 ECSAC reviews.

Distributed Computing for ePIC

- **Echelon 1** sites uniquely perform the **low-latency streaming workflows**:
 - Archiving and monitoring of the streaming data, prompt reconstruction and rapid diagnostics.
- Apart from low-latency, **Echelon 2** sites fully participate in use cases and **accelerate** them.
- Establishing EIC International Computing Organization (EICO):



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		✓	✓	

Substantial role for Echelon 2 in preliminary resource requirements model

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

Echelon 2 Resource Needs (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

O(1M) core-years to process a year of data:

- Even with performance gains over the years, the required processing scale remains substantial.
- Highlights the need to leverage distributed and opportunistic resources from the outset.

~350 PB to store data of one year.

The resource estimates are based on a luminosity of $L = 1 \times 10^{33} \text{ cm}^{-2} \text{ s}^{-1} = 1 \text{ kHz}/\mu\text{b}$. There is an ongoing discussion regarding the EIC Strategy and the luminosity during the first years of running. I propose that we update our resource estimates once there is clarity on the EIC Strategy. For now, we will keep our estimate.

Echelon 2 Roles

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		✓	✓	

Substantial role for Echelon 2 in preliminary resource requirements model

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

Distributed computing provides flexibility, and this flexibility will influence the fraction of use cases handled at Echelon 2. I propose that, for the purpose of estimating resource requirements, we keep the currently assumed fraction.



Networking Estimates

Echelon 0: The raw data from the ePIC Streaming DAQ (Echelon 0) will be replicated across the host labs (Echelon 1). At the highest luminosity of $1e34$, the data stream from the ePIC Streaming DAQ is estimated at 100 Gbit/s. Consequently, Echelon 0 requires an outgoing network connection of at least 200 Gbit/s.

Echelon 1: Each Echelon 1 facility has similar requirements, as it will receive up to 100 Gbit/s of raw data and will share this data with Echelon 2. In addition, Echelon 1 will send a small amount of monitoring data, approximately 1 Gbit/s, back to Echelon 0. Echelon 1 will also receive calibration and analysis data from various Echelon 2 nodes at a comparable rate of about 1 Gbit/s.

Echelon 2: The network connection requirements for Echelon 2 facilities will depend on the proportion of raw data they intend to process. For the 10% of Echelon 1 scenario, a network connection of 20 Gbit/s would be required.

Echelon 2 Types

Use Case	Echelon 1	Echelon 2a	Echelon 2b	Echelon 2c
Streaming Data Storage and Monitoring	✓			
Alignment and Calibration	✓	✓		
Prompt Reconstruction	✓			
First Full Reconstruction	✓	✓		
Reprocessing	✓	✓		
Simulation	✓	✓	✓	
Physics Analysis	✓	✓	✓	
AI Modeling and Digital Twin	✓	✓		✓

The requirements for Echelon 2 sites depend on the use cases. I propose defining 2–3 categories of Echelon 2 sites and then specifying the requirements for each category.