WHEN LESS IS MORE:

# GarNet with Attention

## Towards a Lightweight Graph Neural Network for Reconstruction

**Uzziel Perez** on behalf of
Miriam Calvo Gomez, Xavier Vilasis Cardona (La Salle), et al.

*CPAN, COMCHA, Valencia, Spain*
*November 20, 2025*

# GarNet Collaborators
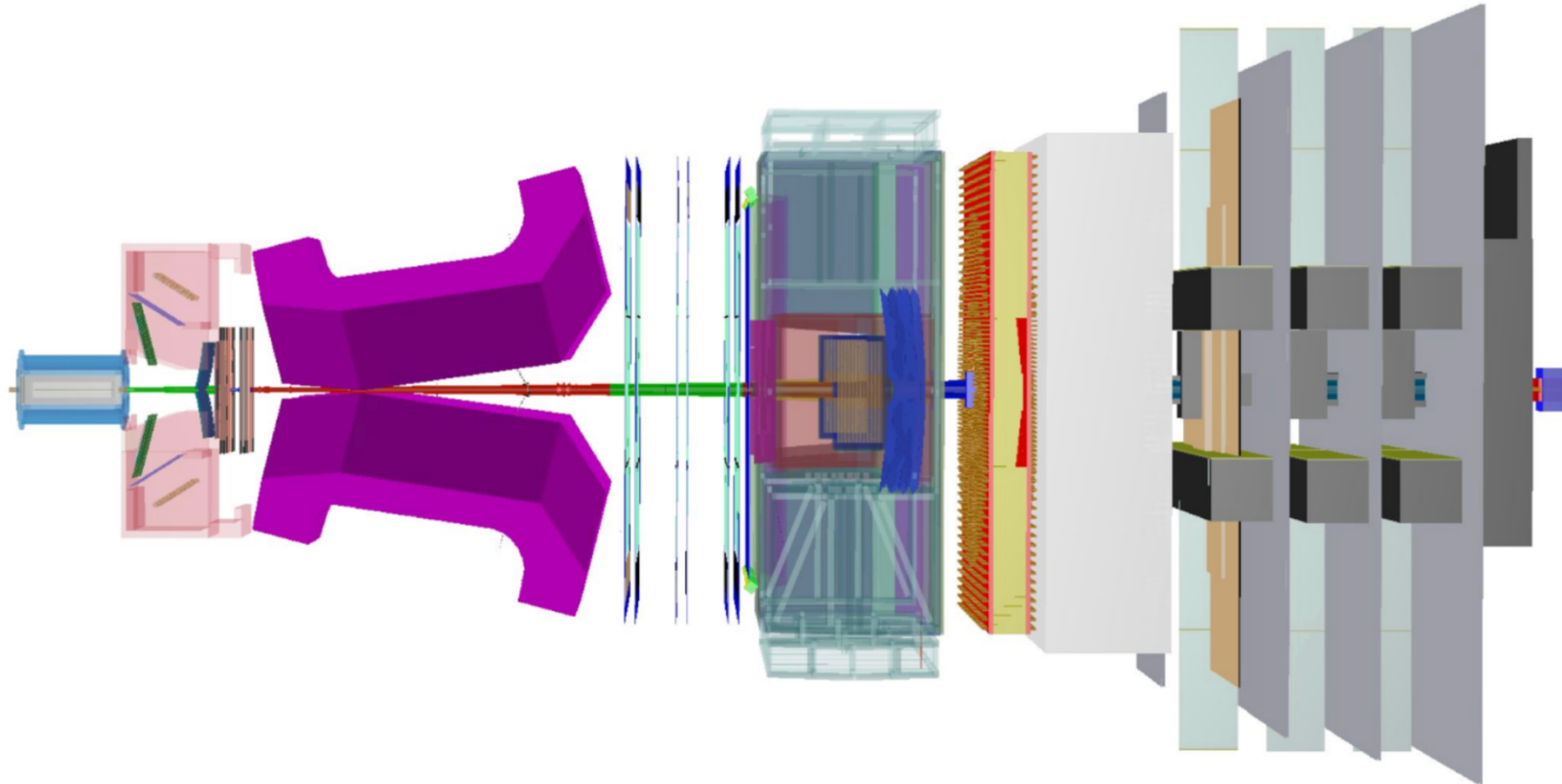
Growing List of Collaborators



Uzziel Perez, Miriam Calvo Gomez, Xavier Vilasis Cardona (La Salle URL, Spain),
Felipe Luan Souza de Almeida (University of Barcelona, Spain),
Justin Bartz, Matthew Rudolph, Wren Vetens, (Syracuse University, USA),
Rafael Silva Coutinho (Centro Brasileiro de Pesquisas Físicas, Brazil)
Katya Govorkova (MIT), Ke Wei (Wuhan University, China)
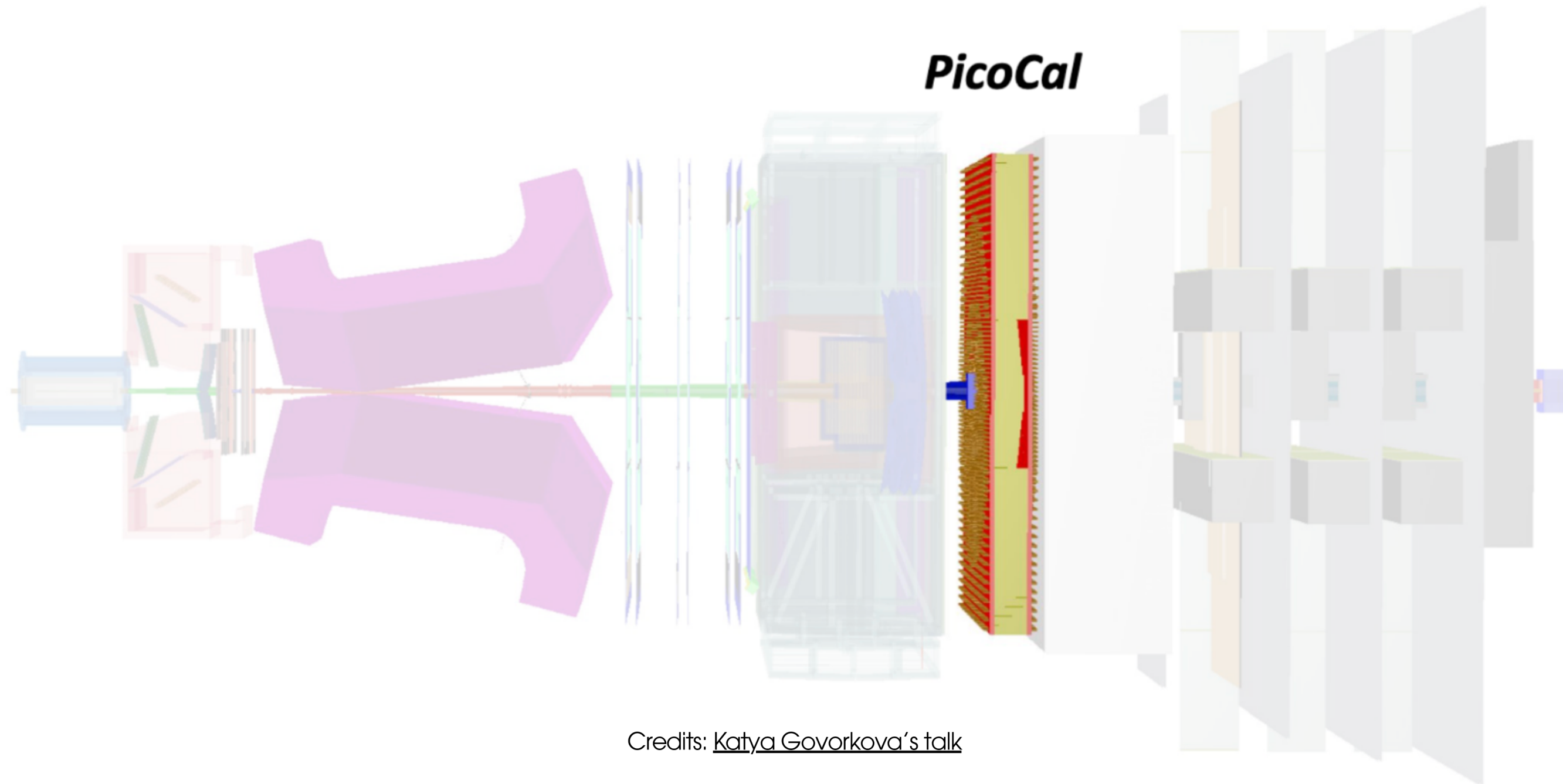
# Introduction

The **LHCb Upgrade II** redesigns LHCb to operate at 5x the instantaneous luminosity with a data rate of 200 Tb/s

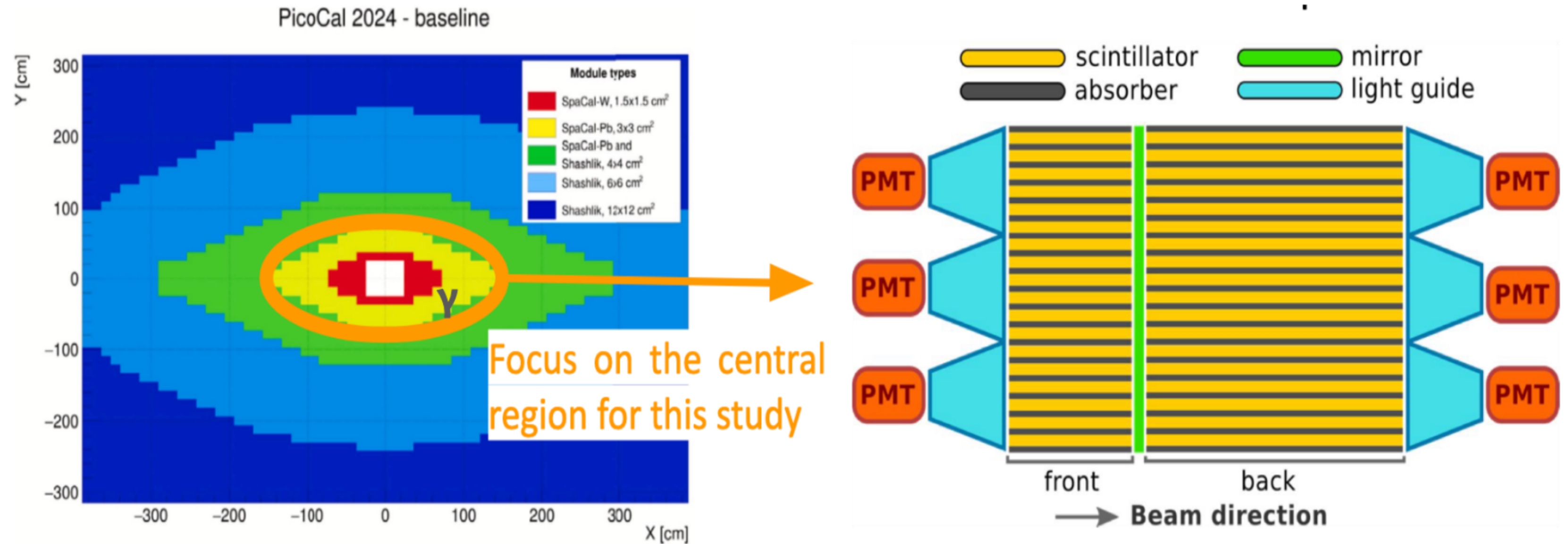Credits: Katya Govorkova's talk

# Next Gen: PicoCal Detector

The **PicoCal** is the next-generation of the electromagnetic calorimeter for γ, e⁻, pion (neutral) reconstruction, which includes timing information of **O(10) ps** precision
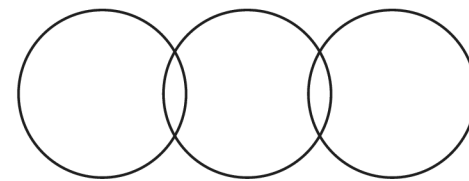


Credits: Katya Govorkova's talk

# Next Gen: PicoCal Detector

The central region will be replaced with radiation-tolerant **SpaCal Modules** which have W/Pb absorbers and crystal/plastic scintillating fibers (LHCb-TDR-026)
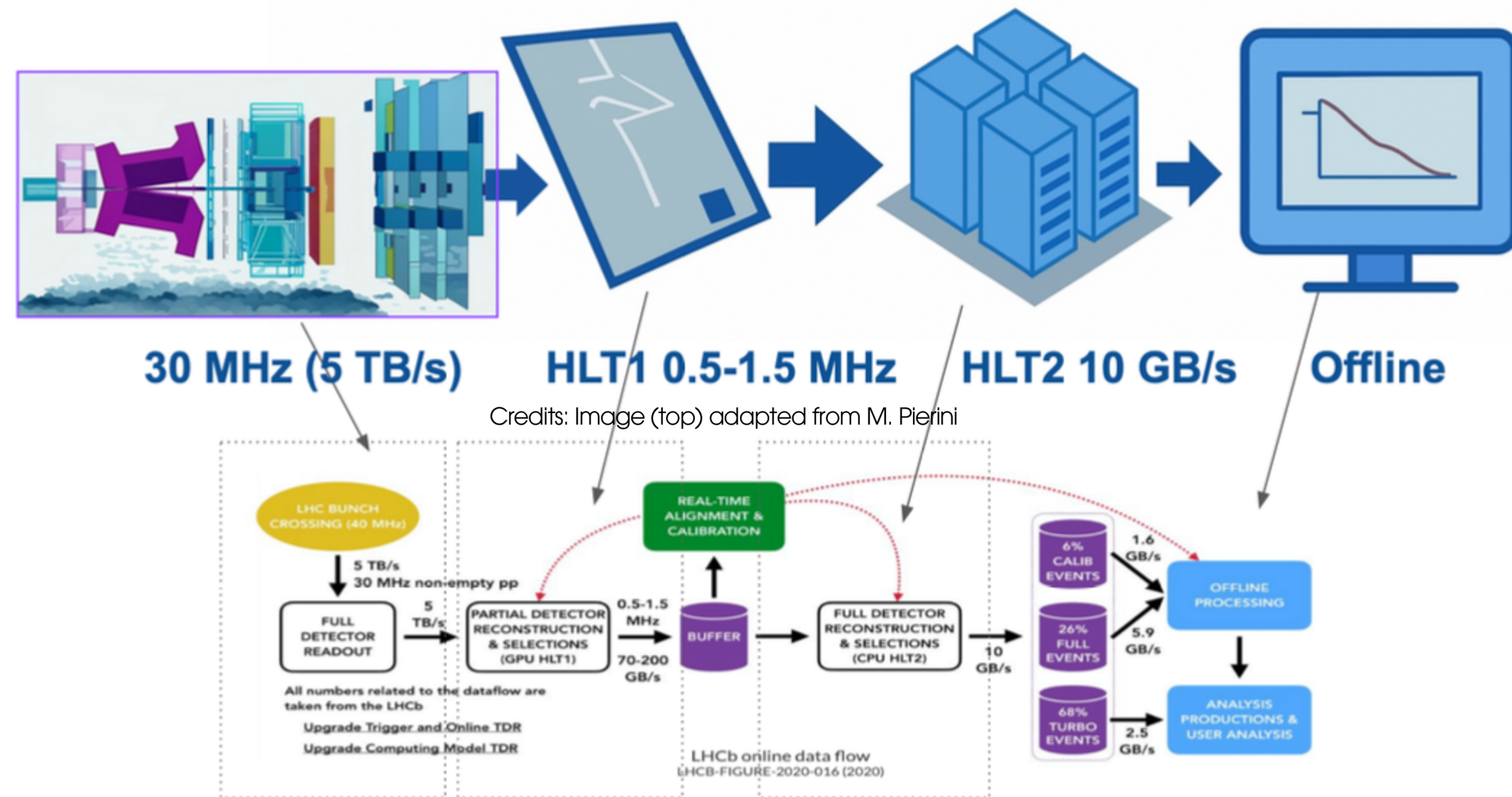


Focus on the central region for this study

Credits: Katya Govorkova's talk

# What are the challenges for Real-Time Reconstruction?

# Latency and Throughput

What are the requirements and Bottlenecks?
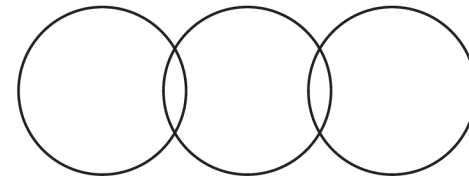


**30 MHz (5 TB/s)**  **HLT1 0.5-1.5 MHz**  **HLT2 10 GB/s**  **Offline**

Credits: Image (top) adapted from M. Pierini

→ LHC bunch crossings occur every **25 ns** → Latency requirement: **~10 μs** w/ buffers absorbing **~0.1-1 ms**

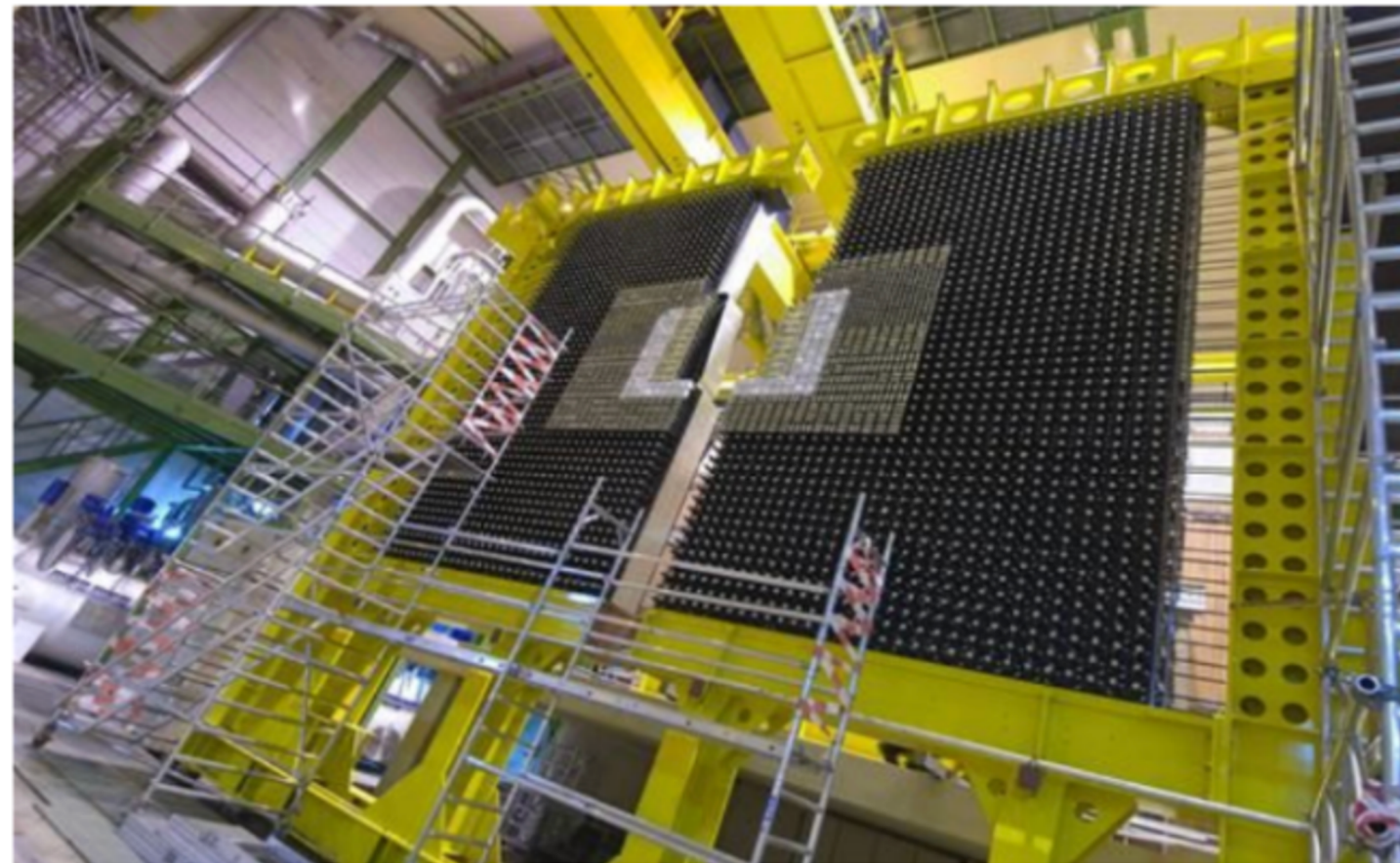→ HLT2 has a 10GB/s throughput bottleneck. To avoid backpressure → have low latency reconstruction

# Evolving Reconstruction Algorithm

# Evolution of Reconstruction Algorithms

→ The current Graph Clustering Algorithm is less than 60% faster than the legacy Cellular Automata
→ GNNs seem like natural successors for future reconstruction



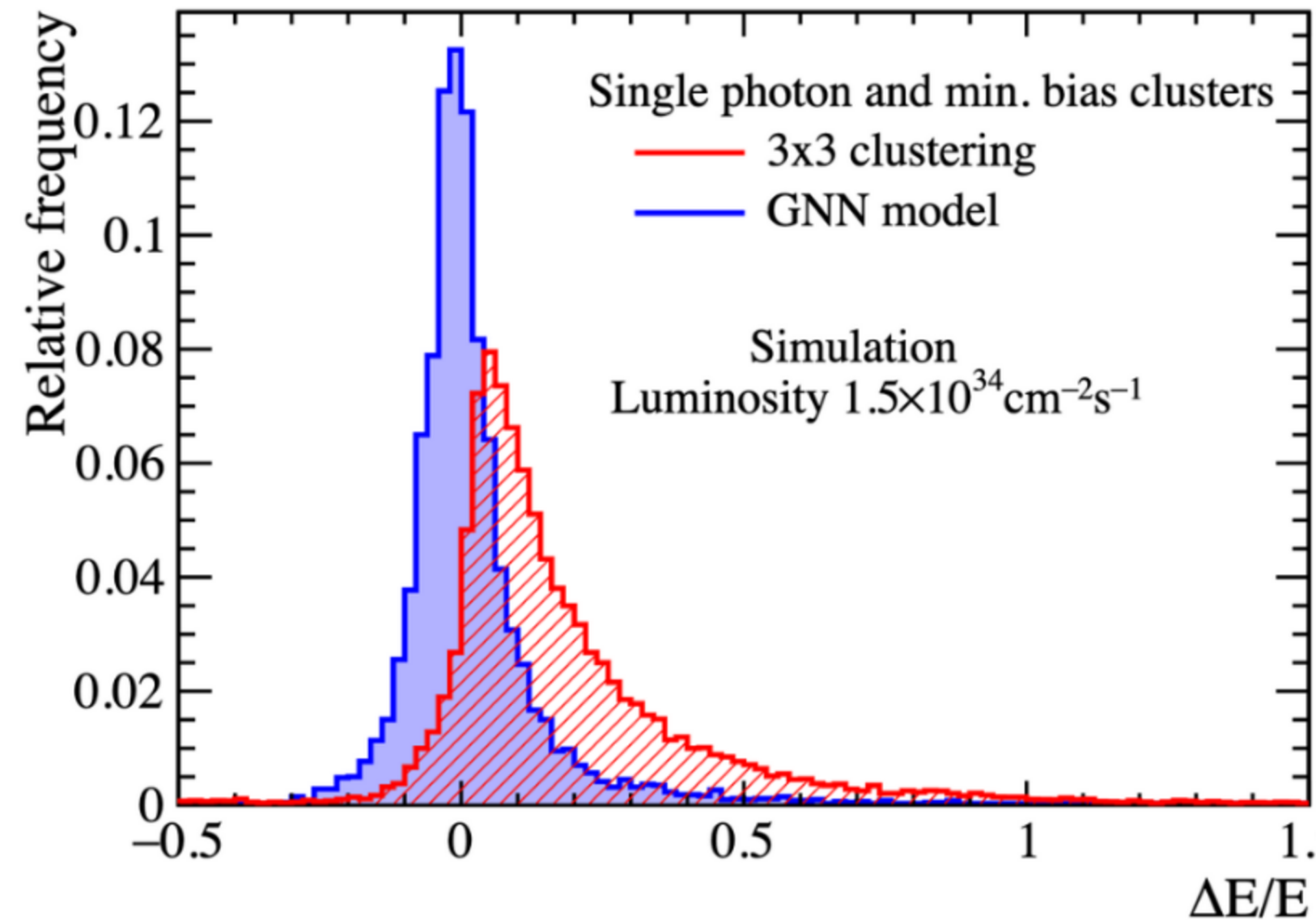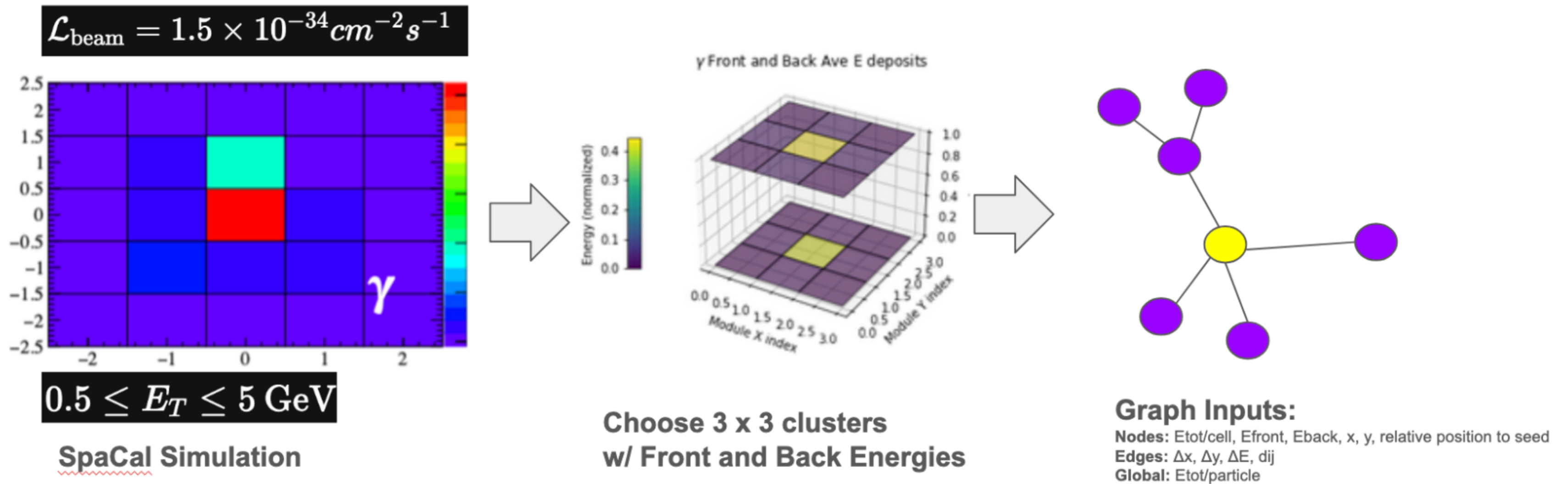| Cellular Automaton | Graph Clustering | Graph Neural Networks |
| PAST | PRESENT | FUTURE |

# Why Graph Neural Networks?

→ Keeps the graph structure of clustering but make the aggregation rules learnable
→ More adaptable for handling irregular detector geometry

# Data Preprocessing

→ Spacal Simulation with Single Photons (particle gun) and minbias clusters
→ Raw PicoCal Data converted to KNN-based graph → node (E, position), edge (spatial links), and global (seed position) features
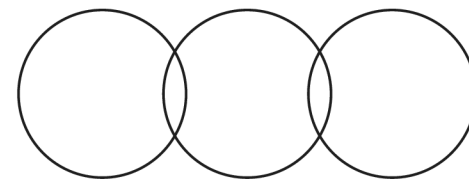


$$\mathcal{L}_{beam} = 1.5 \times 10^{-34} cm^{-2} s^{-1}$$

$\gamma$

$$0.5 \leq E_T \leq 5 \, \text{GeV}$$

**SpaCal Simulation**

y Front and Back Ave E deposits

Energy (normalized)

Module X index

Module Y index

**Choose 3 x 3 clusters
w/ Front and Back Energies**

**Graph Inputs:**
**Nodes:** Etot/cell, Efront, Eback, x, y, relative position to seed
**Edges:** Δx, Δy, ΔE, dij
**Global:** Etot/particle

100k in FULL ECAL, 12k in the Spacal Region

# GNN Fundamentals

→ Input data is also a graph
→ Initial node features updated by message passing layers
→ Nodes are updated by applying a FF-NN on a previous state and received messages
→ Encoder-Decoder is a typical GNN architecture

# GNN Flavours

Full Message Passing
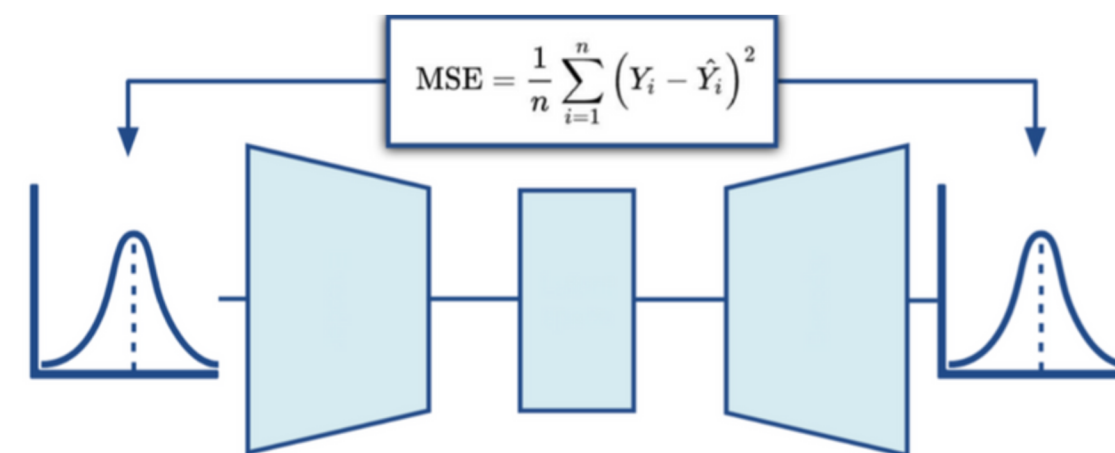
GarNet

# GNN Flavours

In a Nutshell

### With Edges ($\Delta x$, $\Delta y$, $\Delta E$)



$$\text{MSE} = \frac{1}{n}\sum_{i=1}^{n}\left(Y_i - \hat{Y}_i\right)^2$$

Encoder-Processor-Decoder

### No Edges



$$\text{MSE} = \frac{1}{n}\sum_{i=1}^{n}\left(Y_i - \hat{Y}_i\right)^2$$

GarNet Layers 1 + 2 + Global Mean Pool

# GNN w/ Full Message Passing

Full Message Passing → GarNet (see Felipe's Talk!)

→ Encoder: Projects all heterogeneous features into a common embedding space

→ Message Passing: Iterative information propagation between neighbor elements

→ Decoder: Aggregates learned representation to predict shower energy

# GNNMP Training

## Full Message Passing

→ **MSE Loss:** Mean squared error between predicted and true energy

→ **Backpropagation:** Adam Optimizer

→ **Early Stopping:** Stops training if validation loss does not improve for a certain number of epochs

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^{n} \left( Y_i - \hat{Y}_i \right)^2$$

# GNNMP Energy Resolution

## Superior energy resolution

→ Note: 3x3 Clustering resolution of Cellular Automata and Graph Clustering are equivalent

→ GNNMP Model better

# Time-expensive

Too close to the ~1 ms latency benchmark requirement for a single-photon cluster!

## Training Time per Epoch

--- Avg: 321.90s

**~321.9 s or ~5 hours total**

## Batch Inference Time Distribution

**~0.04 s**

❌ Per event latency $= \frac{t_{batch}}{batch\_size} = \frac{0.04}{64} = 0.000625$ s $= 0.625$ ms

❌ Throughput $= \frac{batch\_size}{t_{batch}} = \frac{64}{0.04} = 1600$ events/s $\approx 1.6$ kHz

Disclaimer: lxplus CPU times (Xeon Silver 4216), DDP-gloo with 4 processes

# CMS Reconstruction

Taking inspiration from two papers from CMS on *distance-weighted graph neural networks* and their *FPGA implementations,* we experimented with a similar variant dubbed as **GarNet with Attention**

# GarNet with Attention

→ Explicit edge features, i.e. relative distances and energies removed
→ **Encoder-Processor-Decoder** replaced with simpler **2 GarNet layers** and a **Global Mean Pool**
→ GarNet Layer: Learned aggregators + distance attention to predict incident particle energy from 3×3 front/back cell energies

# How does the attention mechanism work in GarNet?

# Attention Mechanism

→ Each node represents a detector cell
→ **Attention weights**: Each node compared to another and finds "similarities", $\alpha_{ij} = \texttt{softmax}(q_i \cdot k_j)$
→ **Output:** weighted sum of all node features $h_i = \Sigma_j \alpha_{ij} \cdot v_j$
→ The model learns to focus on relevant patterns across the detectors



Your GarNetLayer – Compact Multi-Head Self-Attention (K=4 heads)

# Attention Mechanism

→ **Nodes:** Chaos of fires and spills
→ **Attention weights:** Priority scores
→ **Output**: weighted response strategy $h_i = \Sigma_j \alpha_{ij} \cdot v_j$
→ Prioritizes which chaos to fix first and the **Global Mean Pool** gives a big picture

# Attention Mechanism

→ ✅ Superior energy resolution despite removing edge features

# GarNet Time Complexity

→ ✅ ~8x speedup: Training (from 4 hours down to 30 minutes) and Inference Time (0.05 s to 0.007 s)



**Training Time per Epoch**

~53.22 s or ~33.7 minutes

- - - Avg: 53.22s

**Batch Inference Time Distribution**

~0.007 s

✅ Per event latency $= \frac{t_{batch}}{batch\_size} = \frac{0.005}{64} = 0.000078125$ s $= 0.07812$ ms

✅ Throughput $= \frac{batch\_size}{t_{batch}} = \frac{64}{0.005} = 12800$ events/s $\approx 12.8$ kHz

Disclaimer: lxplus CPU times (Xeon Silver 4216), DDP-gloo with 4 processes

# Head-to-Head vs GNNMP

→ ✅ Removing explicit edge features → less mathematical operations
→ GNNMP Encoder-processor-Decoder → 2 GarNet Layers and Global Mean Pool

Disclaimer: lxplus CPU times (Xeon Silver 4216), DDP-gloo with 4 processes

# Moving Forward

# Active Efforts

Promising Initial Results from:

→ Distillation and Quantization providing additional speedup with and superior resolution (c/o Irvin Chacon/ F. Siles of Univ. of Costa Rica)

→ Conversion to ONNX Format with even more speedup (c/ Ronald Caravaca/ F. Siles of Univ. of Costa Rica)

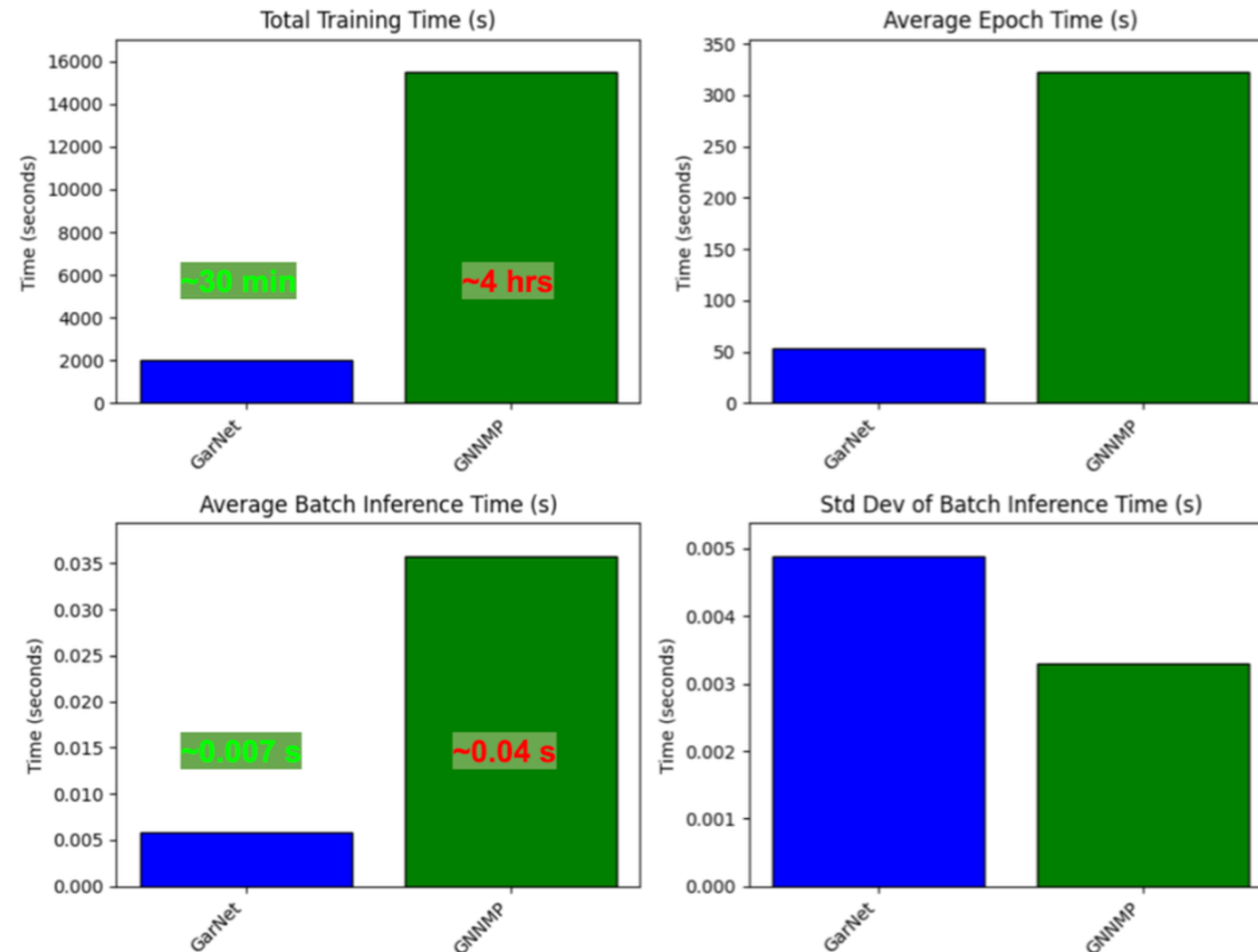→Testing with Allen Framework for Run 5 with the PicoCal (in collaboration with  C. Agapopolou, G. Khreich, A. L. Salvia, J.F. Marchand, et al. of the ODISSEE Project )

# Distillation

Promising Initial Results from:

→ Distillation involves a smaller student network that learns from the output of a teacher network.
→ Loss function is a composite of student learning (λ) on its own and learning from teacher (β)
→ Credits to Irvin Chacon! Stay tuned in another conference for the actual results.
→ We get additional speedup and resolution improvements



Student learns from this output

Total Loss = **β**∗MSELossScaled + **λ**∗MSE Loss, **β=0.3, λ = 0.7**

Benchmark: CPU - MacOS M3 chip, using torch.utils.benchmark.Timer with single thread

# Summary

Promising Initial Results from:

→ Clearer path towards HLT1 reconstruction with Graph Neural Networks with Allen Framework with PicoCal
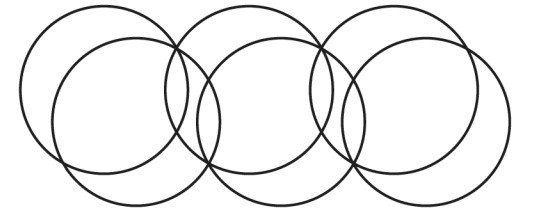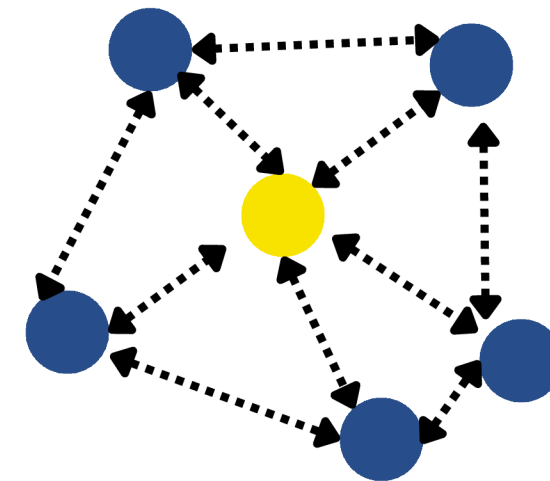
→ We begin with a lightweight version that can be further compressed/ distilled if necessary achieving ~30-100x speedup with our initial results using lxplus CPUs and  (stay tuned for details)

→ Update Testing with Allen Framework for Run 5 with the PicoCal and get more accurate throughput/ latency numbers for Runs 4 and 5

# THANK YOU!



Credits: Nuria Valls-Canudas

# Key Design Principles

## Efficient Architecture for PD

```
> lscpu
Architecture:            x86_64
  CPU op-mode(s):        32-bit, 64-bit
  Address sizes:         46 bits physical, 48 bits virtual
  Byte Order:            Little Endian
CPU(s):                  28
  On-line CPU(s) list:   0-27
Vendor ID:               GenuineIntel
  Model name:            Intel(R) Xeon(R) Silver 4216 CPU @ 2.10GHz
    CPU family:          6
    Model:               85
    Thread(s) per core:  1
    Core(s) per socket:  1
    Socket(s):           28
    Stepping:            7
    BogoMIPS:            4199.76
```

# Future Work and Contact Information

Email
uzzie.perez@cern.ch