

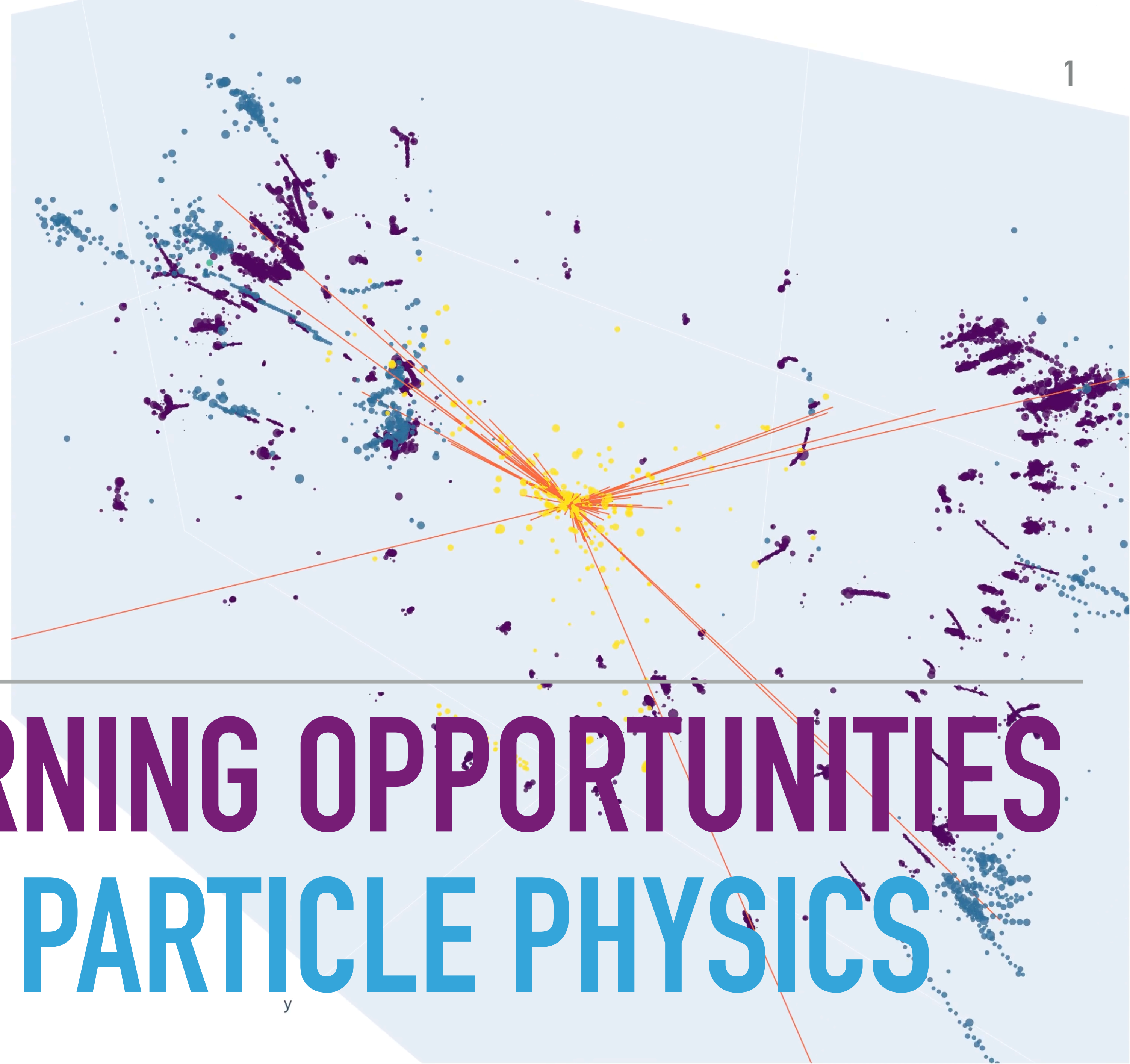
JAVIER DUARTE

ICML

JULY 24, 2024

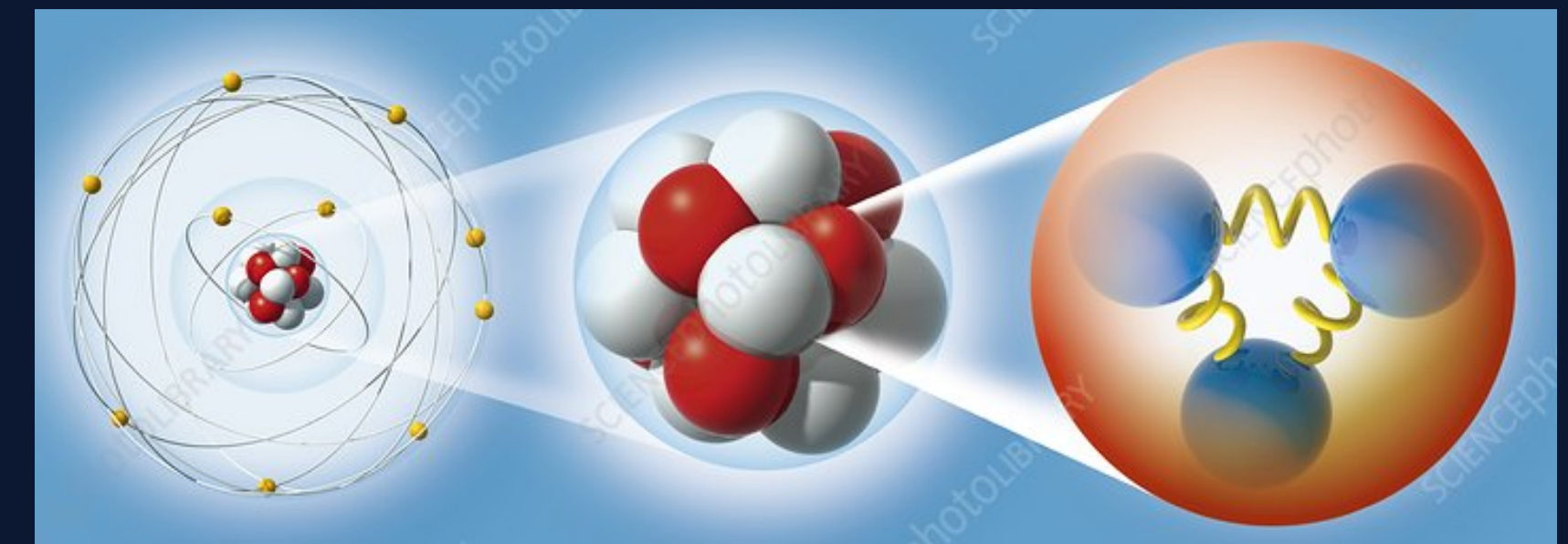
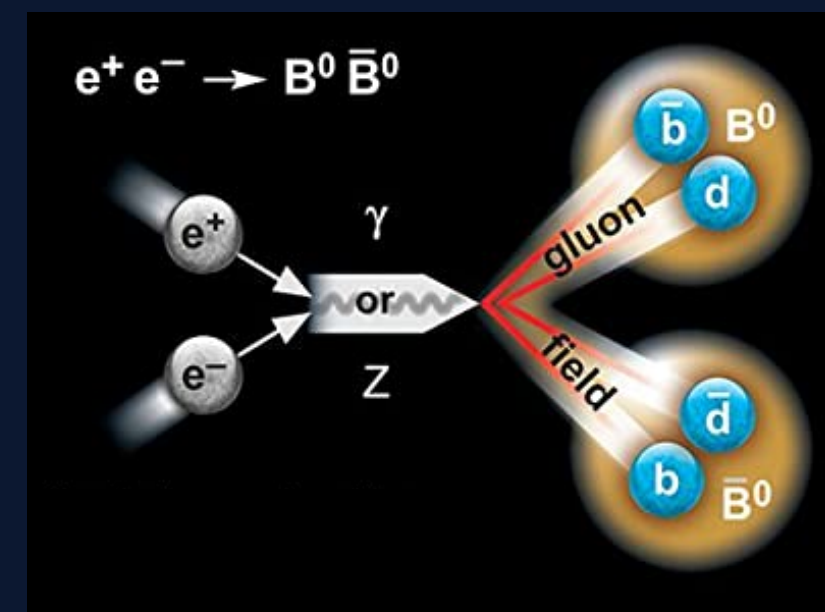
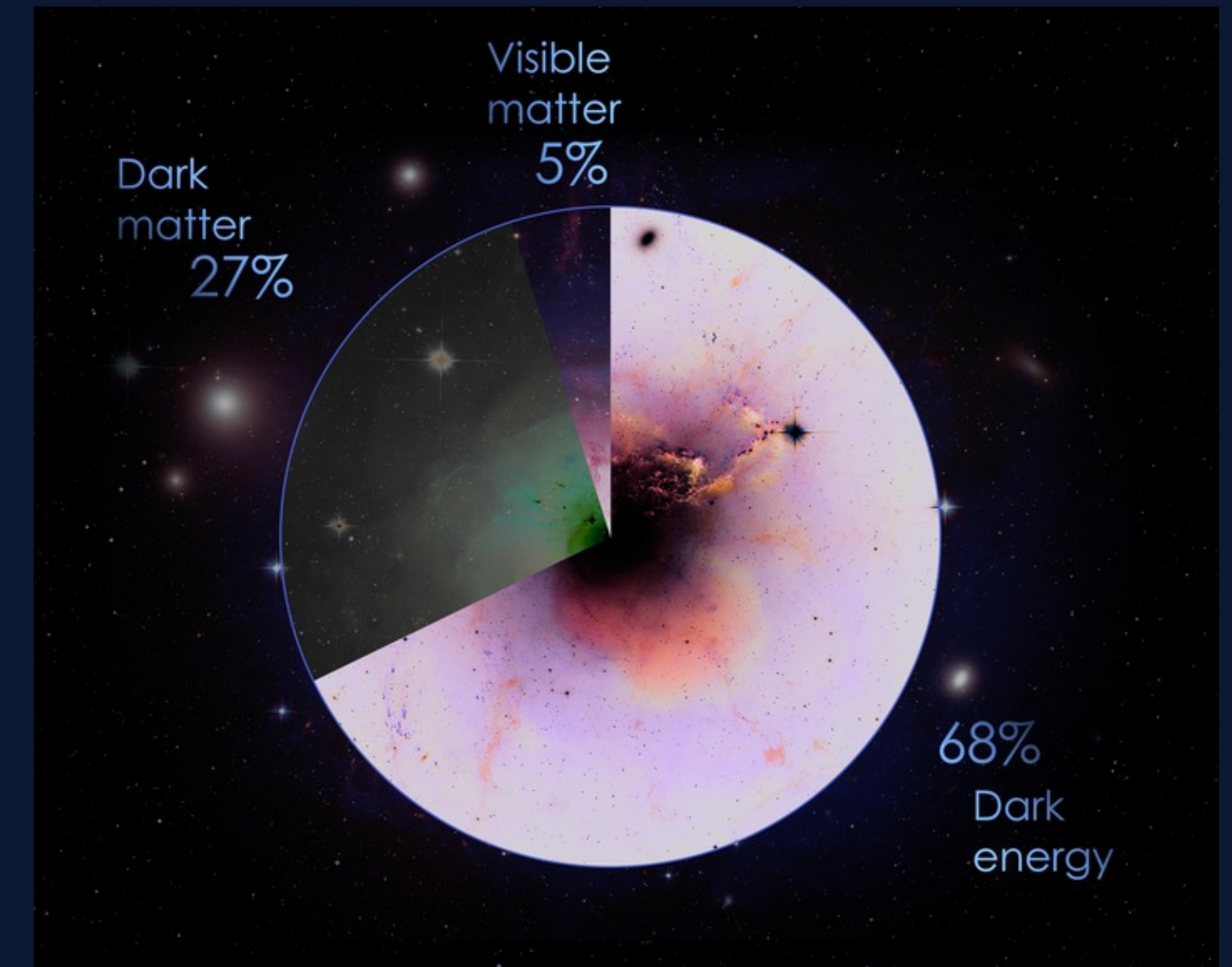
---

# MACHINE LEARNING OPPORTUNITIES FOR NEXT GEN PARTICLE PHYSICS



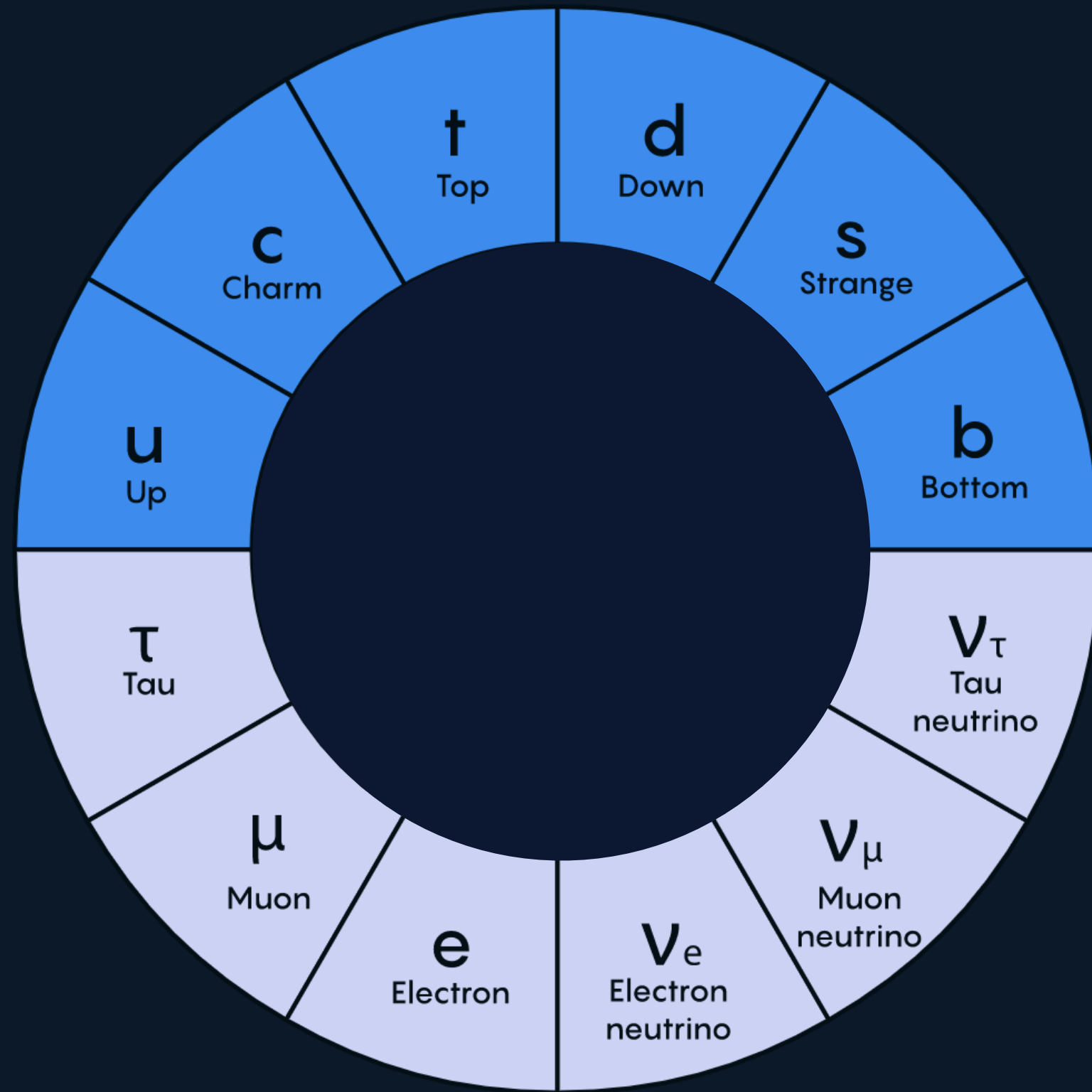


- ▶ What is our universe made of?
- ▶ What are the smallest building blocks of nature?
- ▶ How do they interact with each other?
- ▶ Is our universe stable?

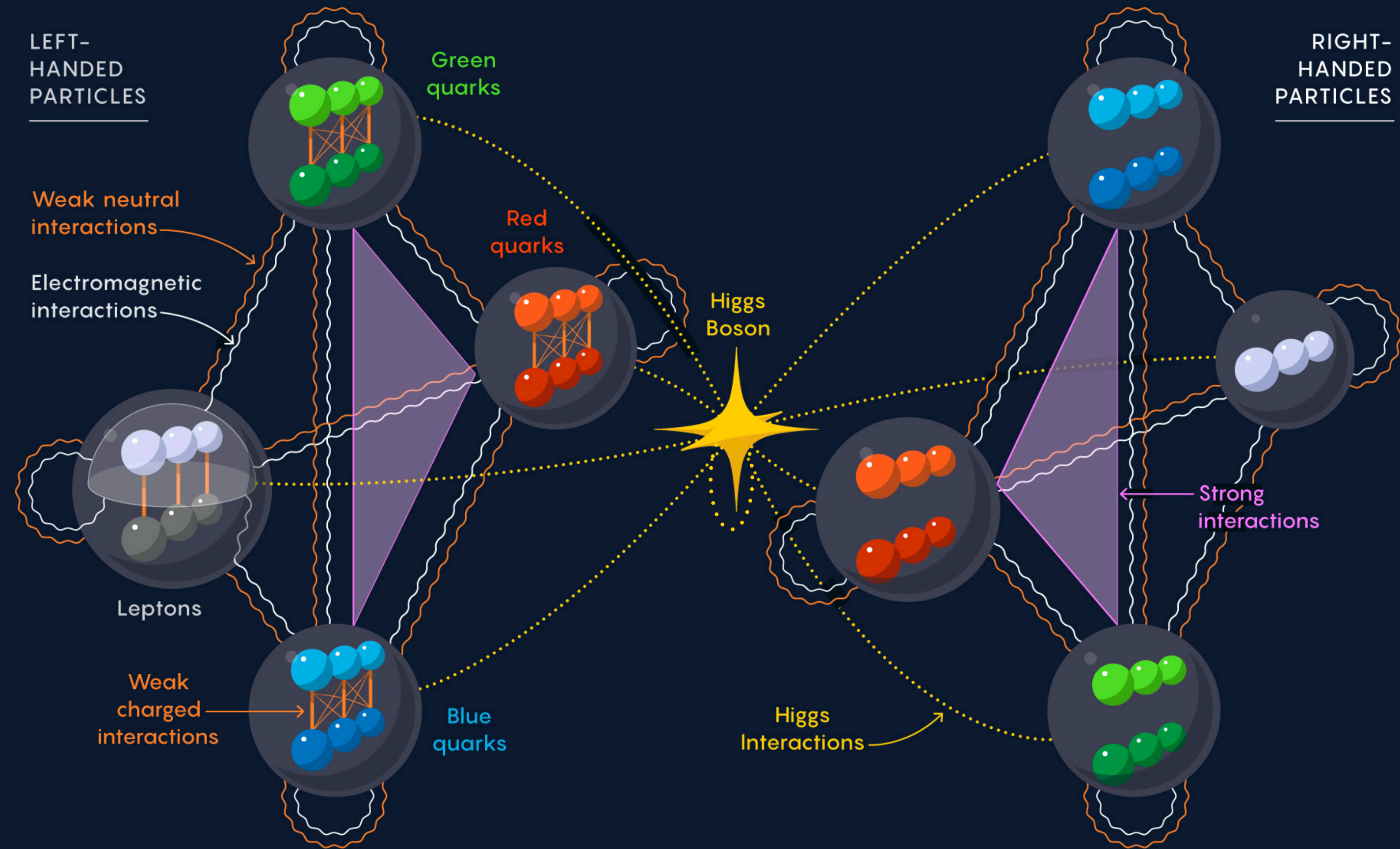




# THE STANDARD MODEL



FERMIONS (MATTER)  
● QUARKS ● LEPTONS



- ▶ But there has to be more to it! SM does not answer all our questions
- ▶ Higgs is a **centerpiece**: Mechanism by which particles **acquire mass**
- ▶ How do we study these microscopic building blocks?



# WHY HIGH ENERGY?

- ▶ High energies  $\leftrightarrow$  short length & time scales
- ▶ Collisions at the highest energy possible today let us recreate conditions 0.1 ns after the Big Bang!

# HISTORY OF THE UNIVERSE

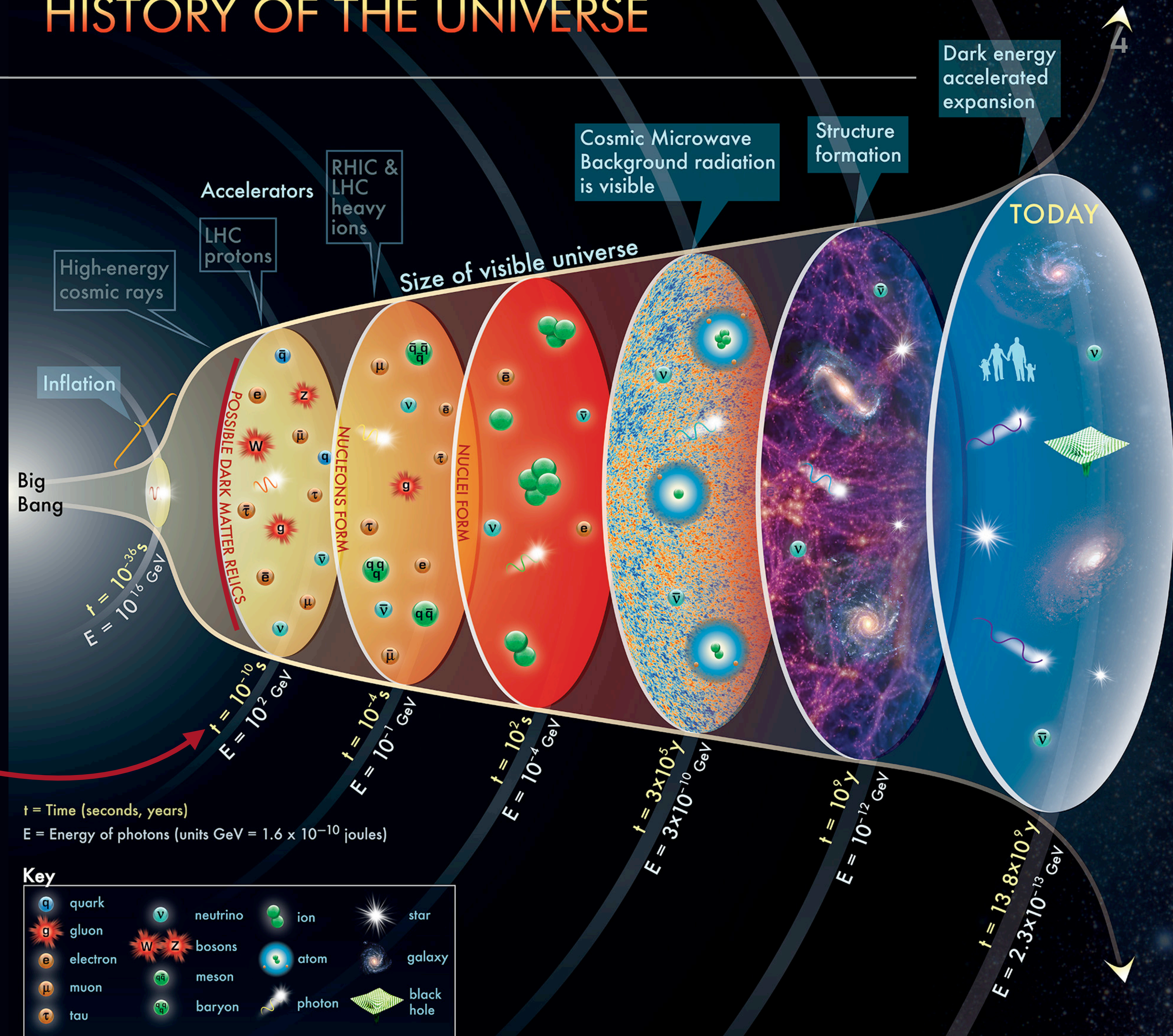


Image: <https://particleadventure.org/history-universe>

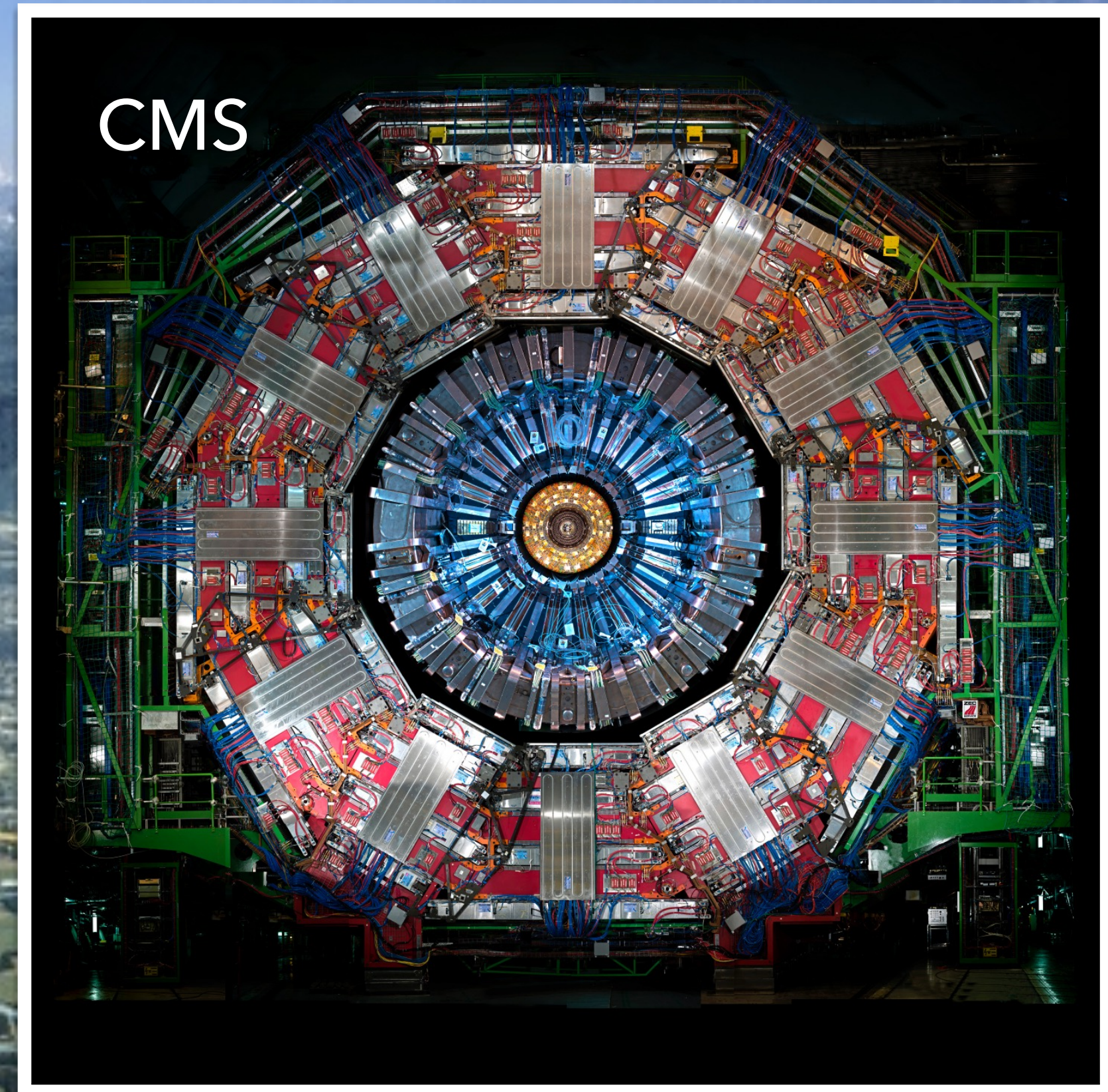
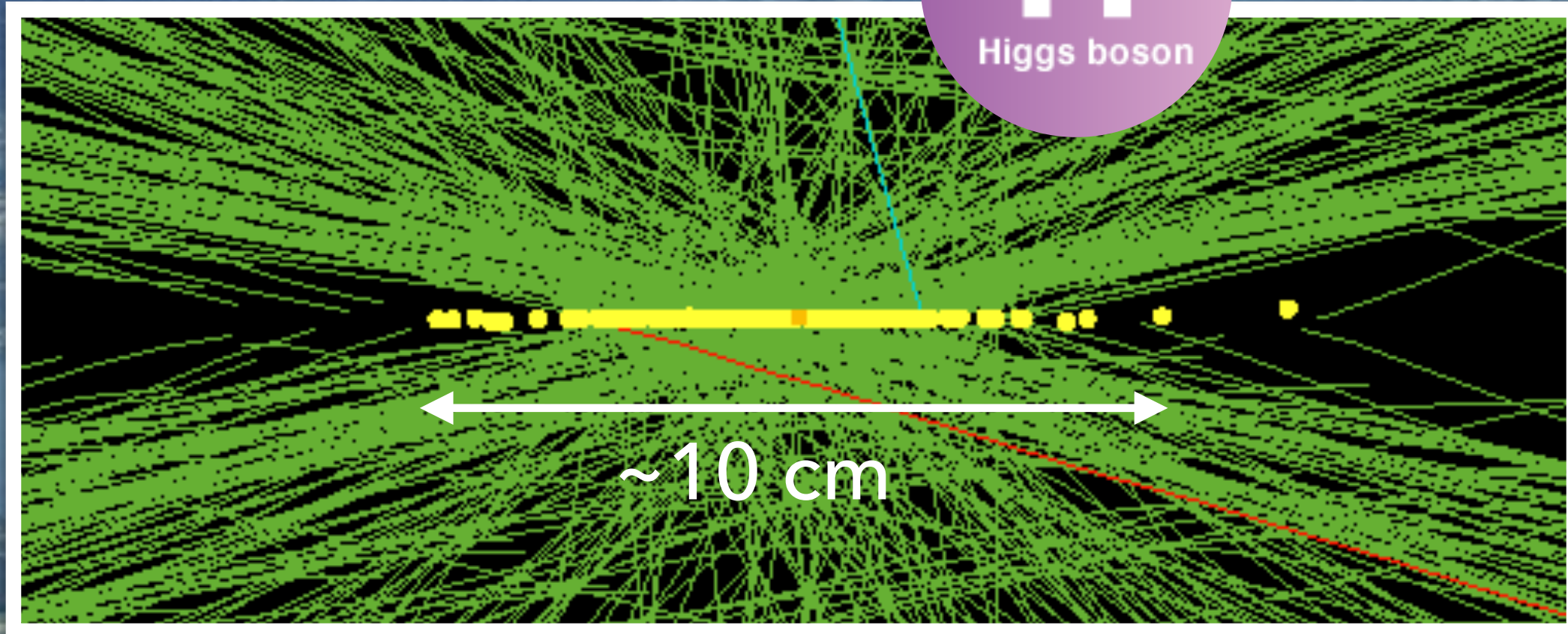
The concept for the above figure originated in a 1986 paper by Michael Turner.

Particle Data Group, LBNL © 2015

Supported by DOE



# THE LARGE HADRON COLLIDER



proton-proton collider @ 13 TeV center-of-mass energy

4 interaction points

40 million collisions / second

**Higgs boson** produced 1/10 billion collisions (every 4 minutes)

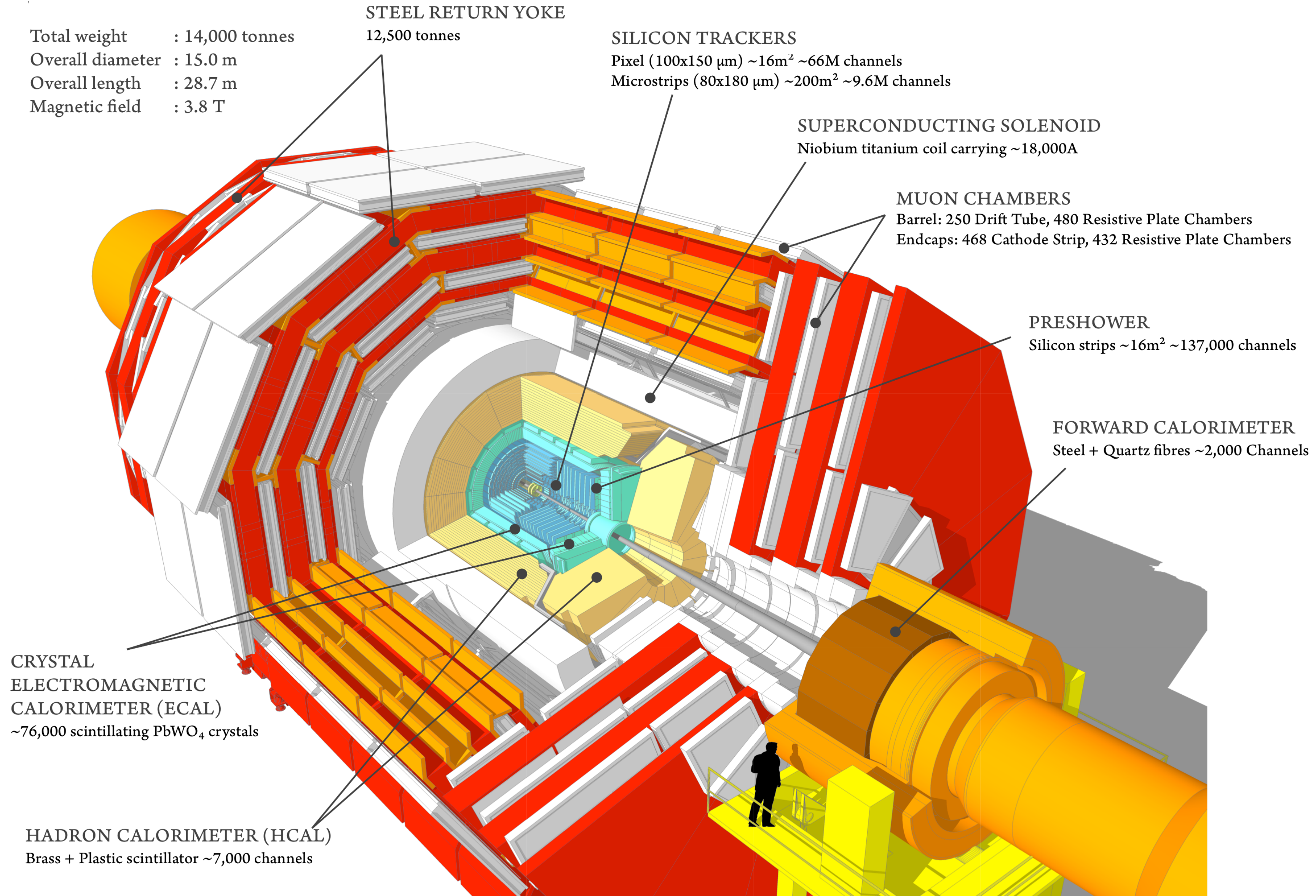
analyze ~1000 collisions / second

LHC 27 km

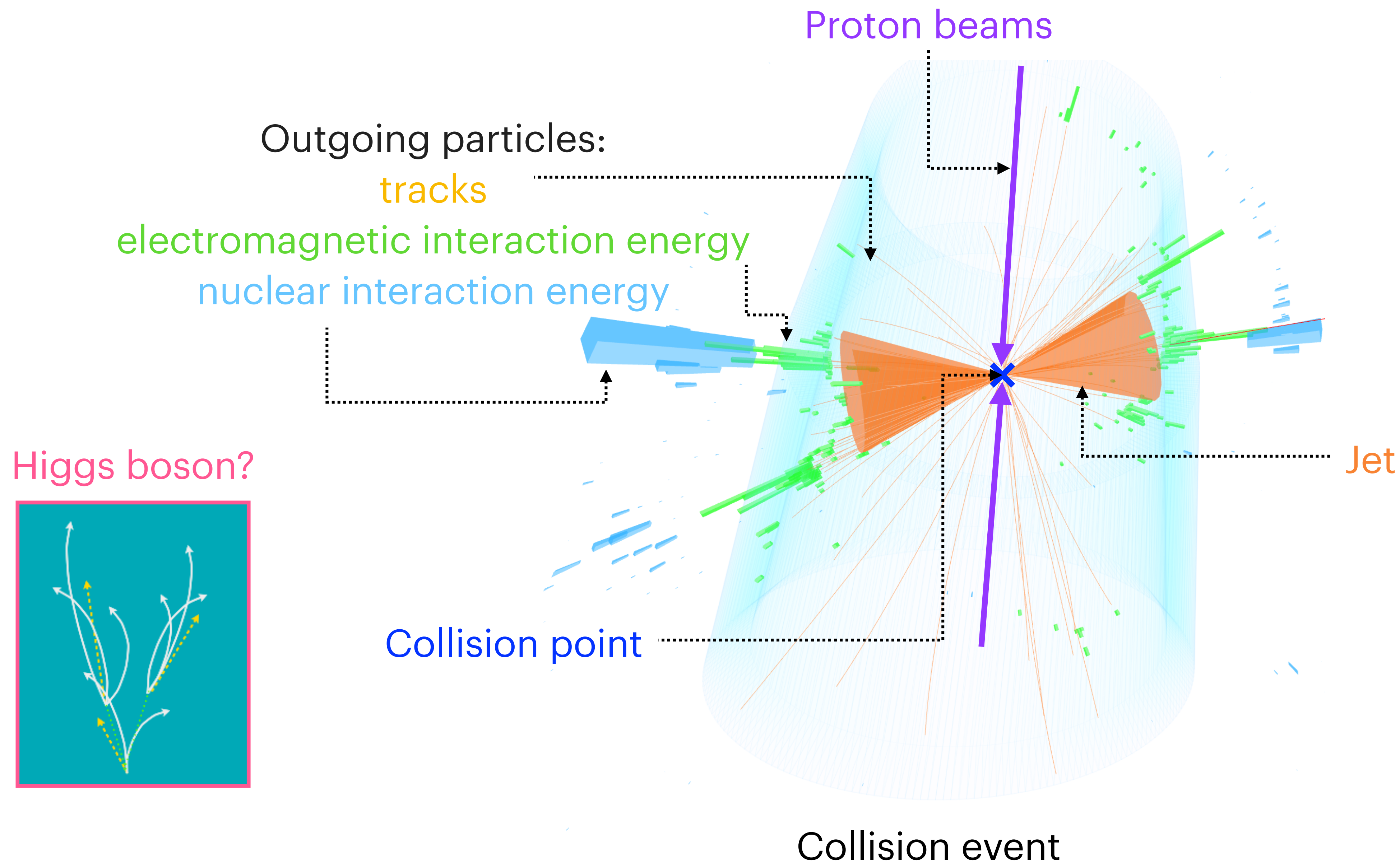


► Specialized components to measure different particles

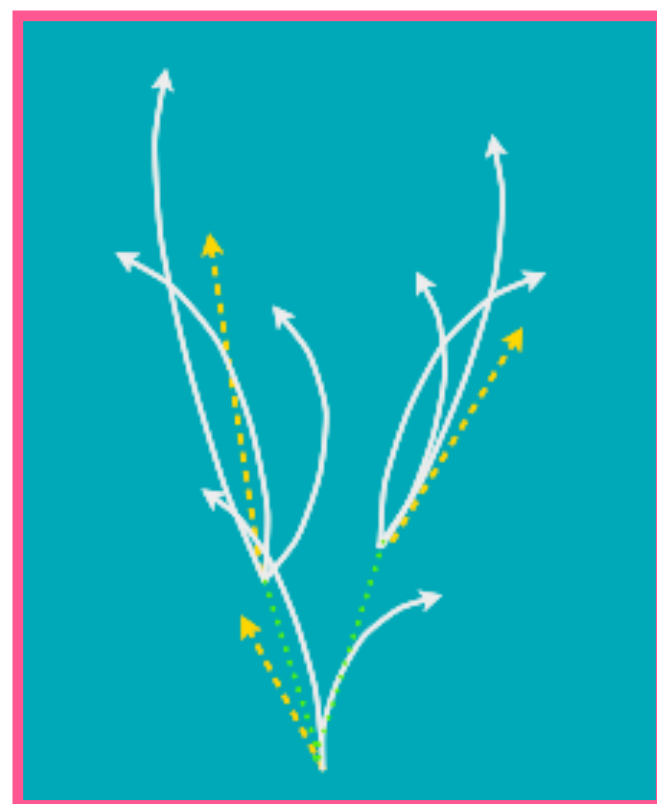
► 100 million channels







Higgs boson?





- ▶ Machine learning has changed the way we do particle physics
  - ▶ It is an essential and versatile tool that we use to improve existing approaches
  - ▶ It enables fundamentally new approaches
- ▶ Highlight a few active areas of R&D, with public datasets and opportunities to collaborate!





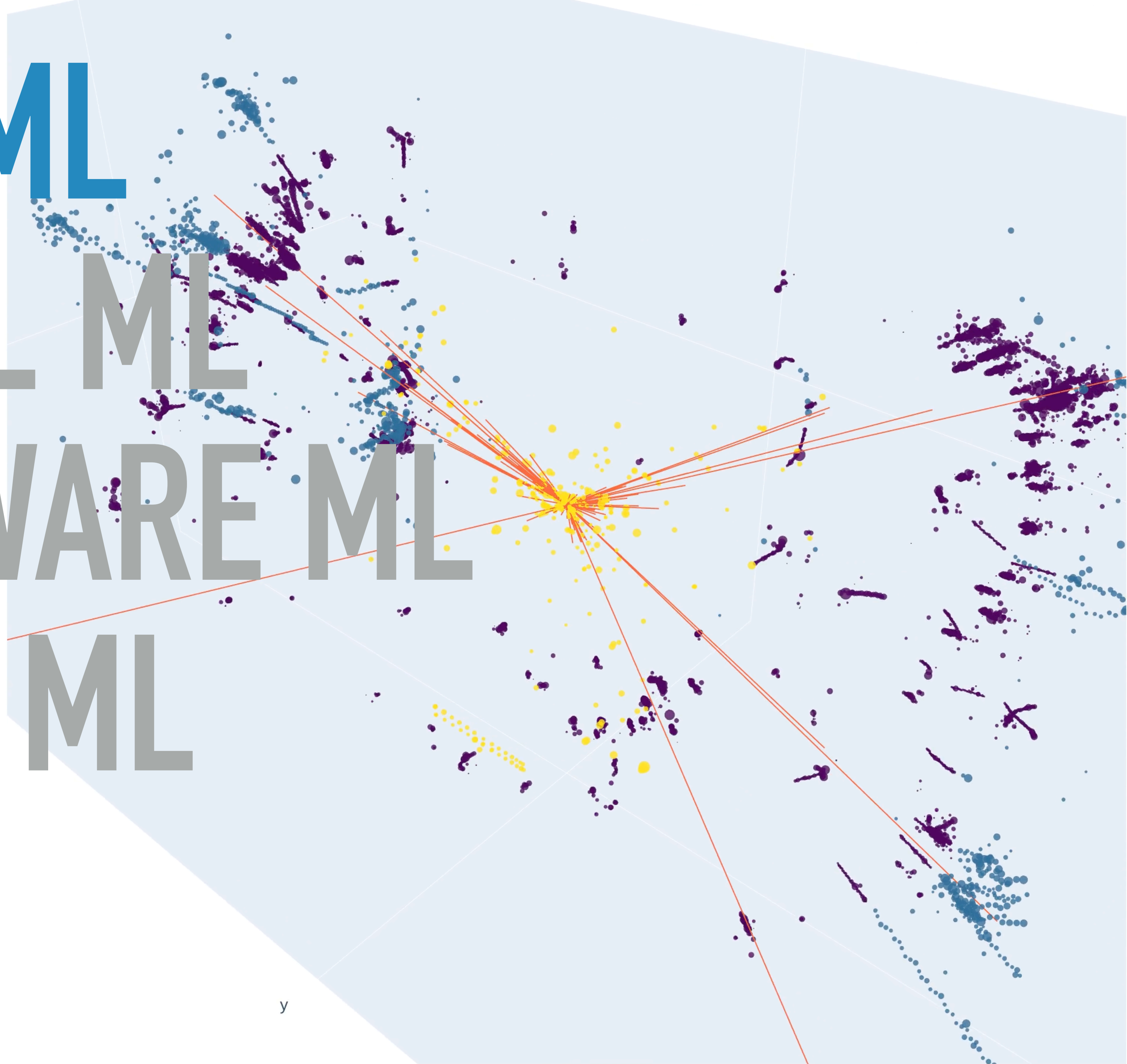
**ULTRAFAST ML**

MULTIMODAL ML

PHYSICS-AWARE ML

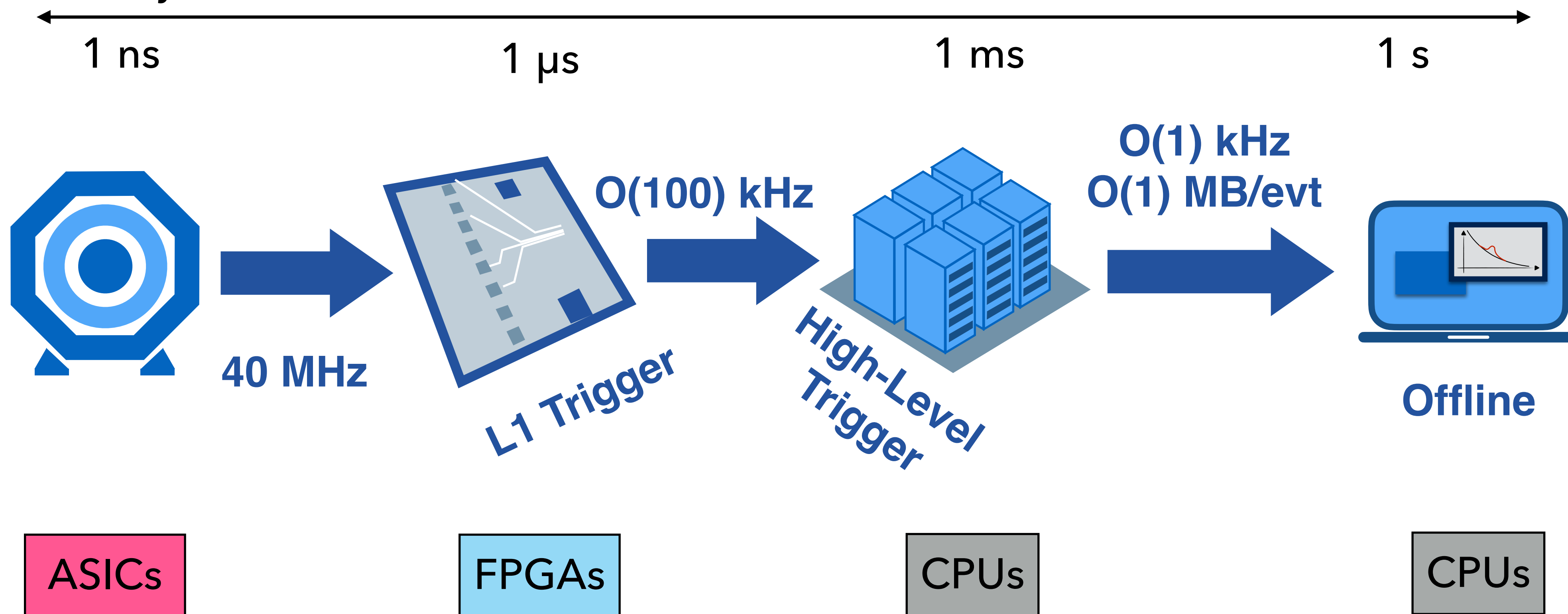
GENERATIVE ML

OUTLOOK





Compute  
Latency



## Challenges:

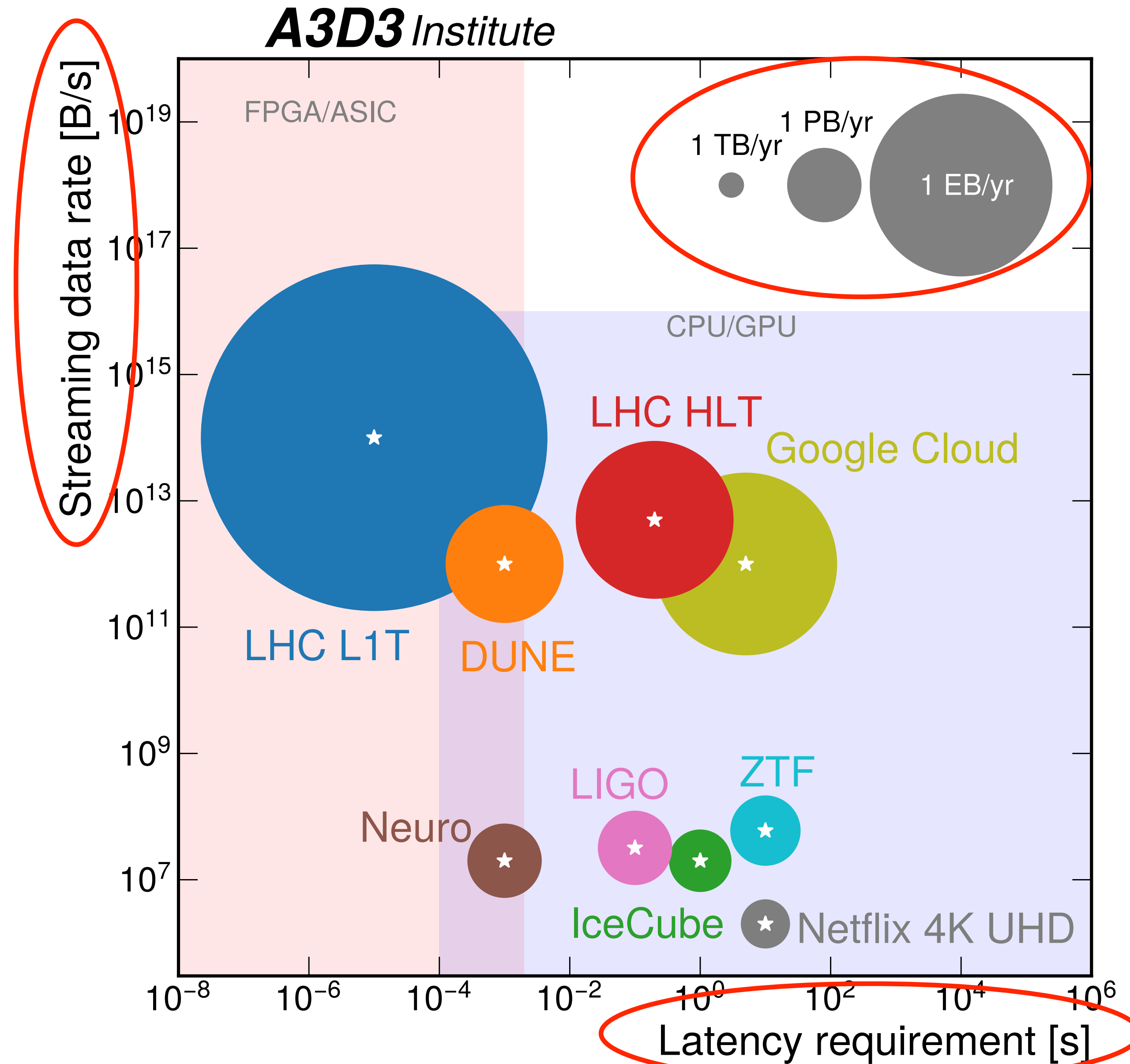
Each collision produces  $O(10^3)$  particles

The detectors have  $O(10^8)$  sensors

Extreme data rates of  $O(100 \text{ TB/s})$

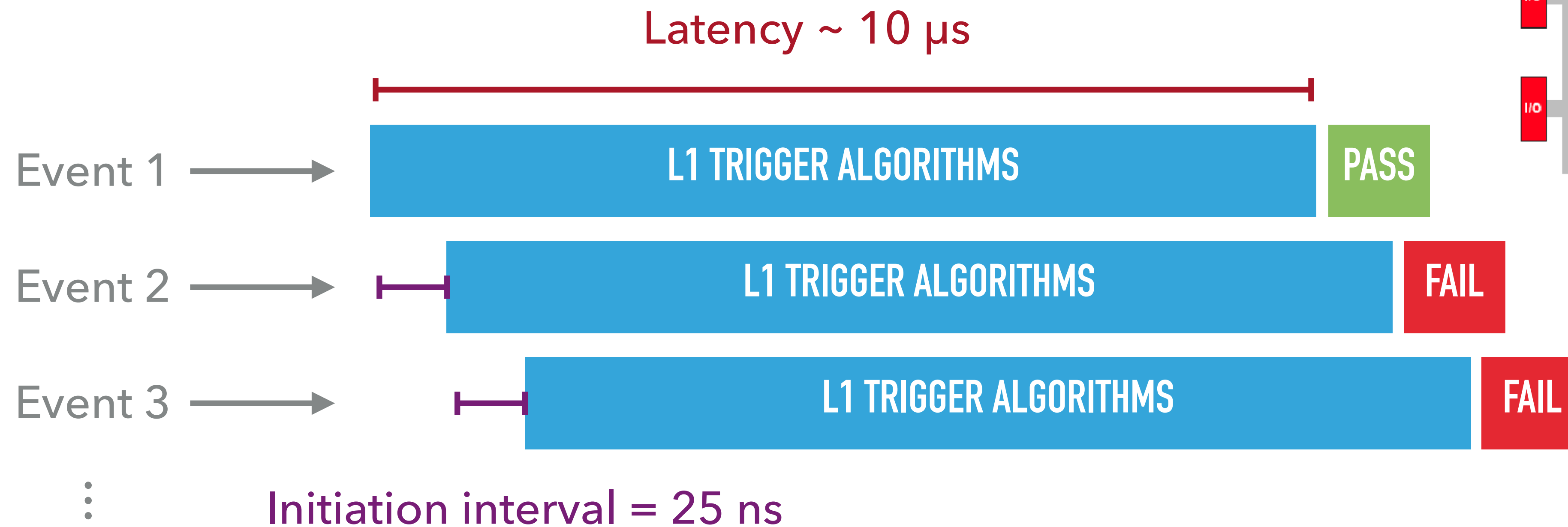
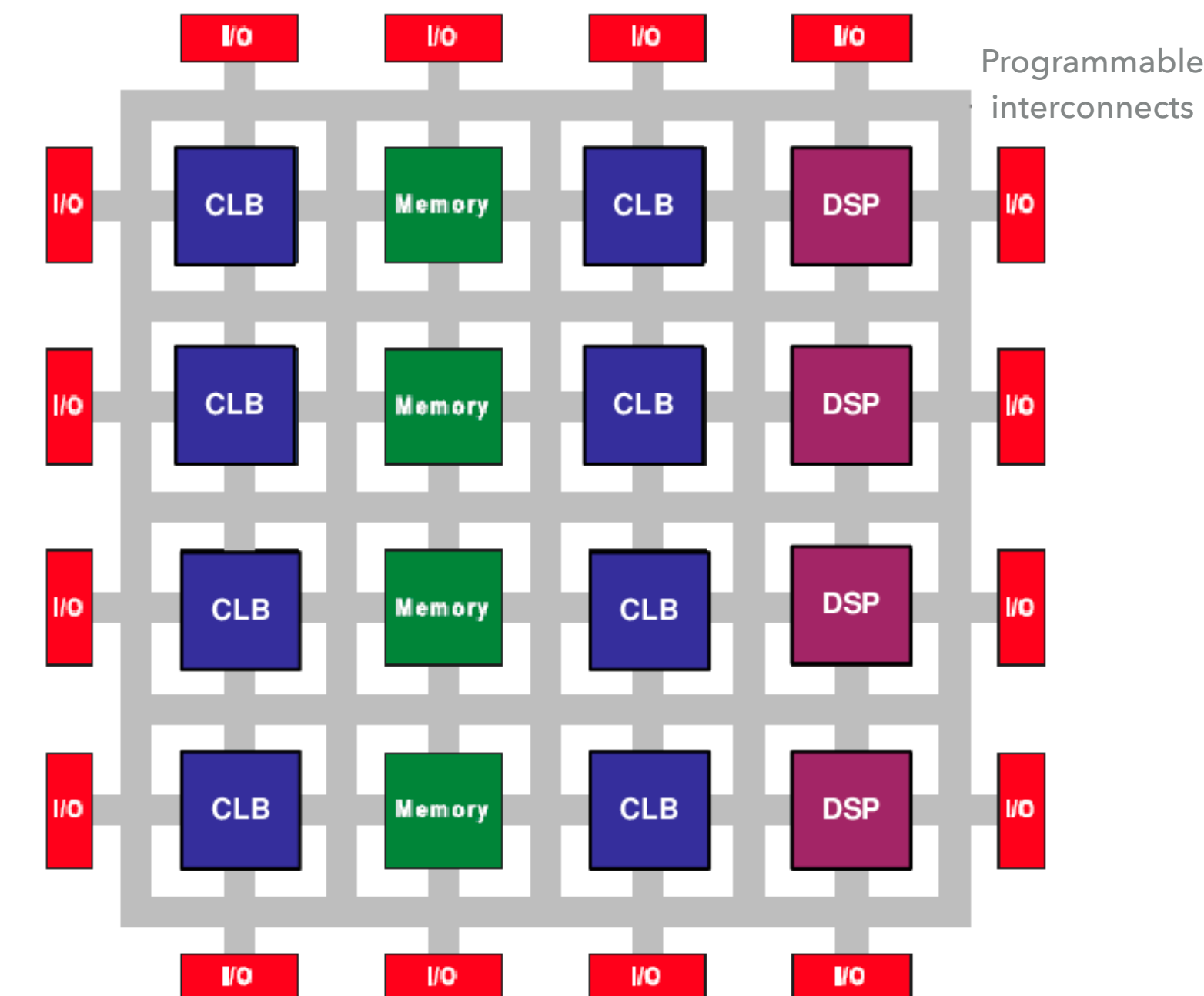
Exabyte-scale  
datasets





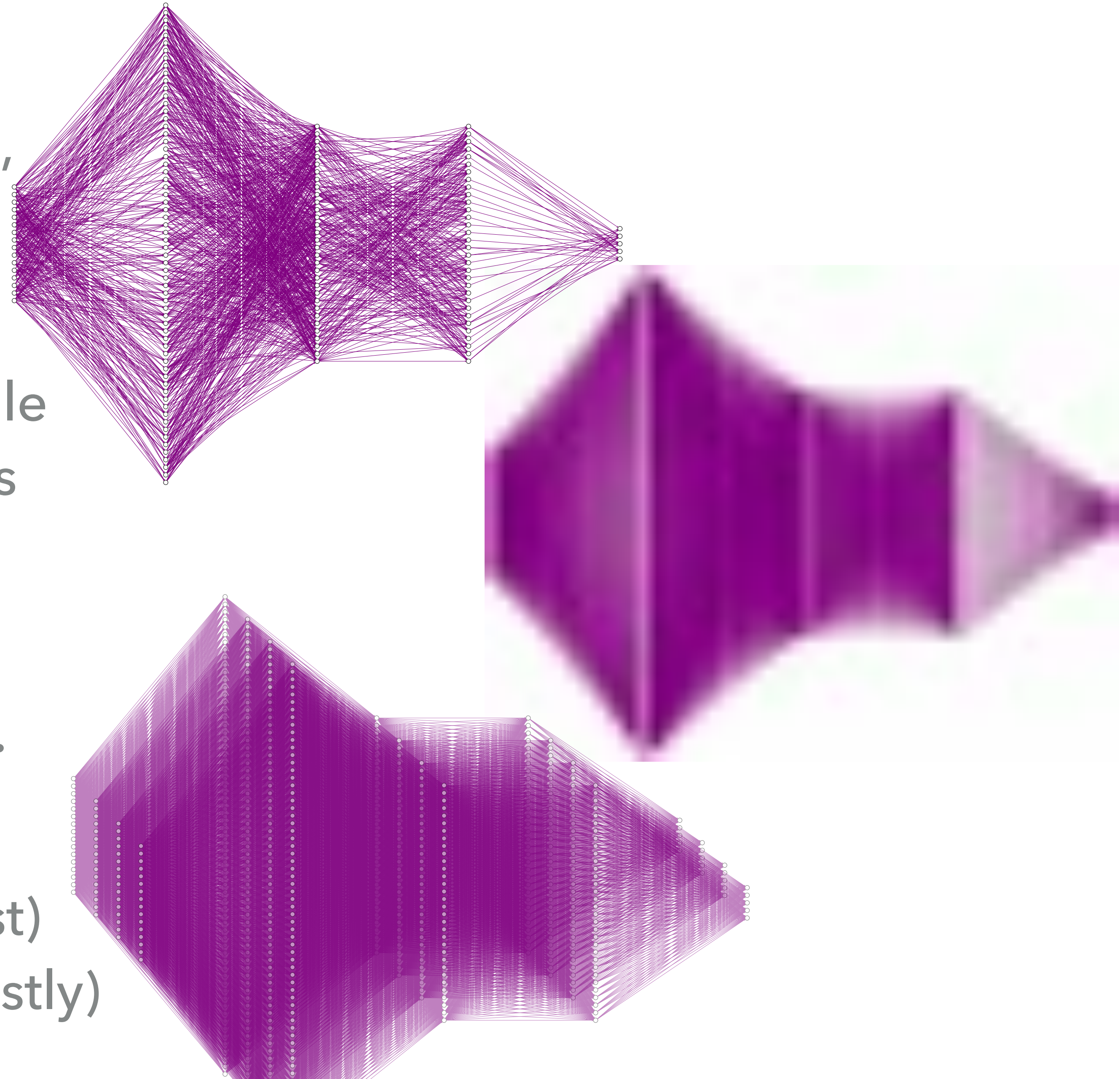


- ▶ Reconstruct all events and reject 98% of them in  $O(10) \mu\text{s}$ 
  - ▶ Algorithms have to be  $<1 \mu\text{s}$  and process new events every 25 ns
- ▶ Latency necessitates all **FPGA** design
  - ▶ Algorithms have to fit on  $<1$  FPGA
- ▶ How can we satisfy these constraints?



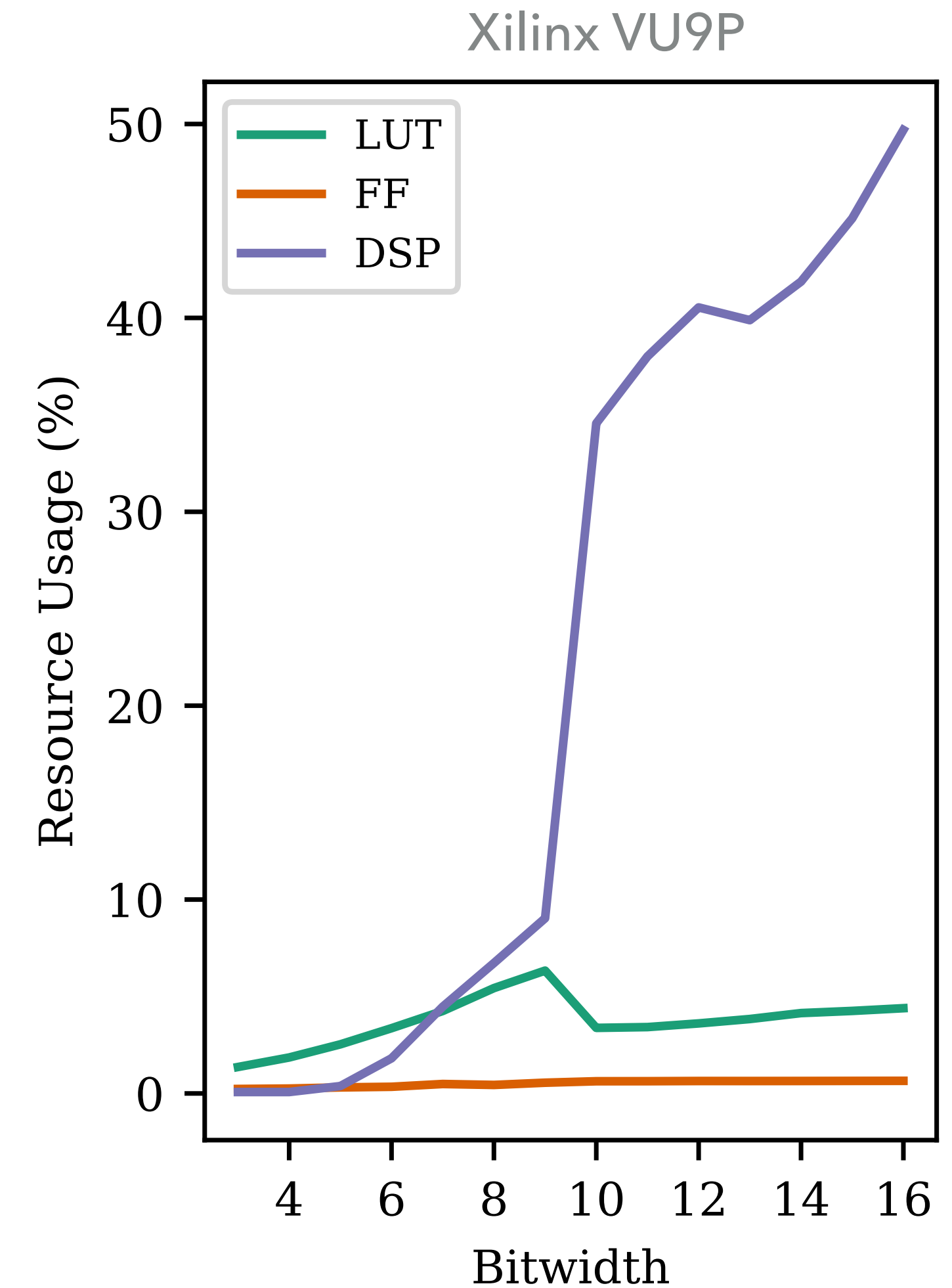
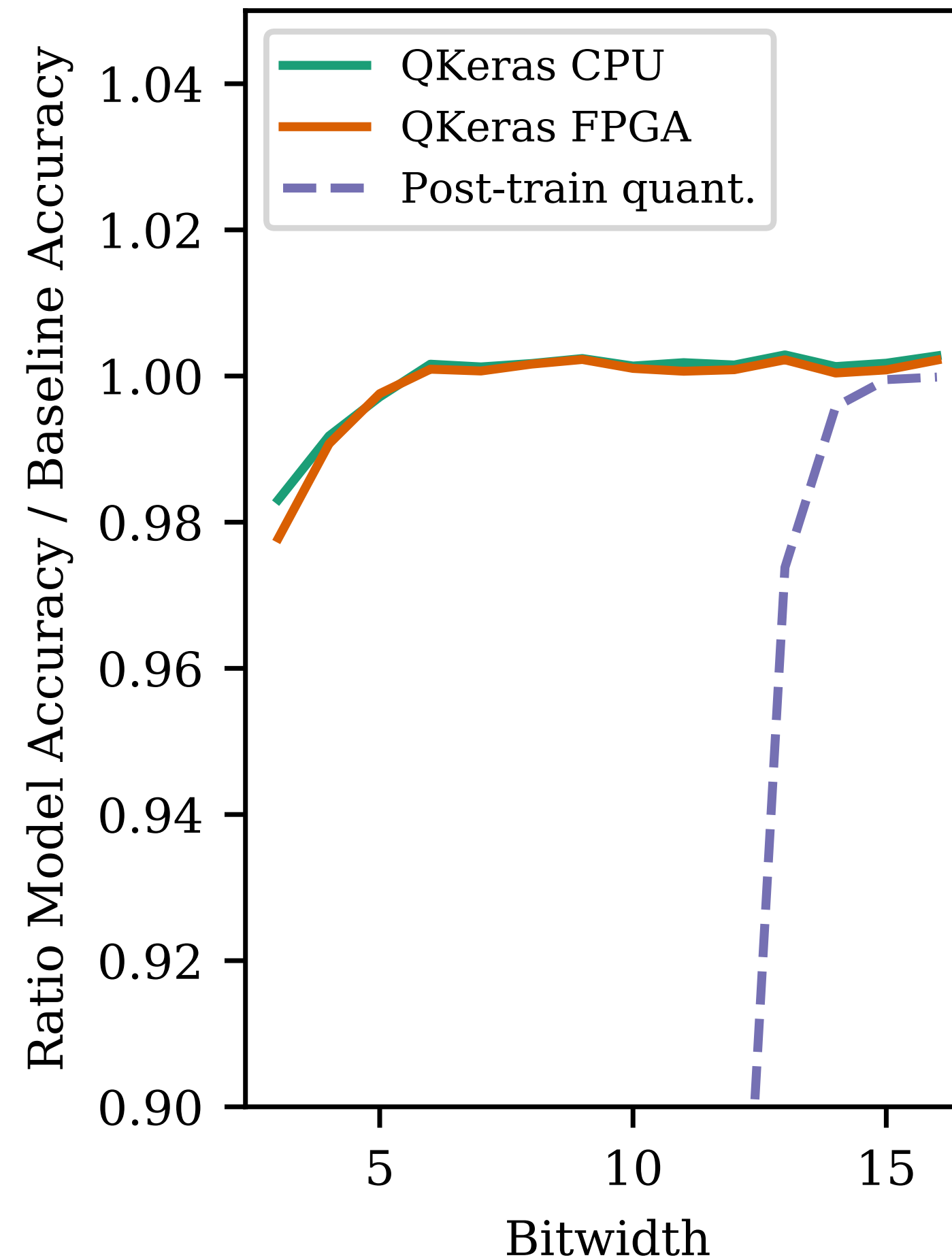
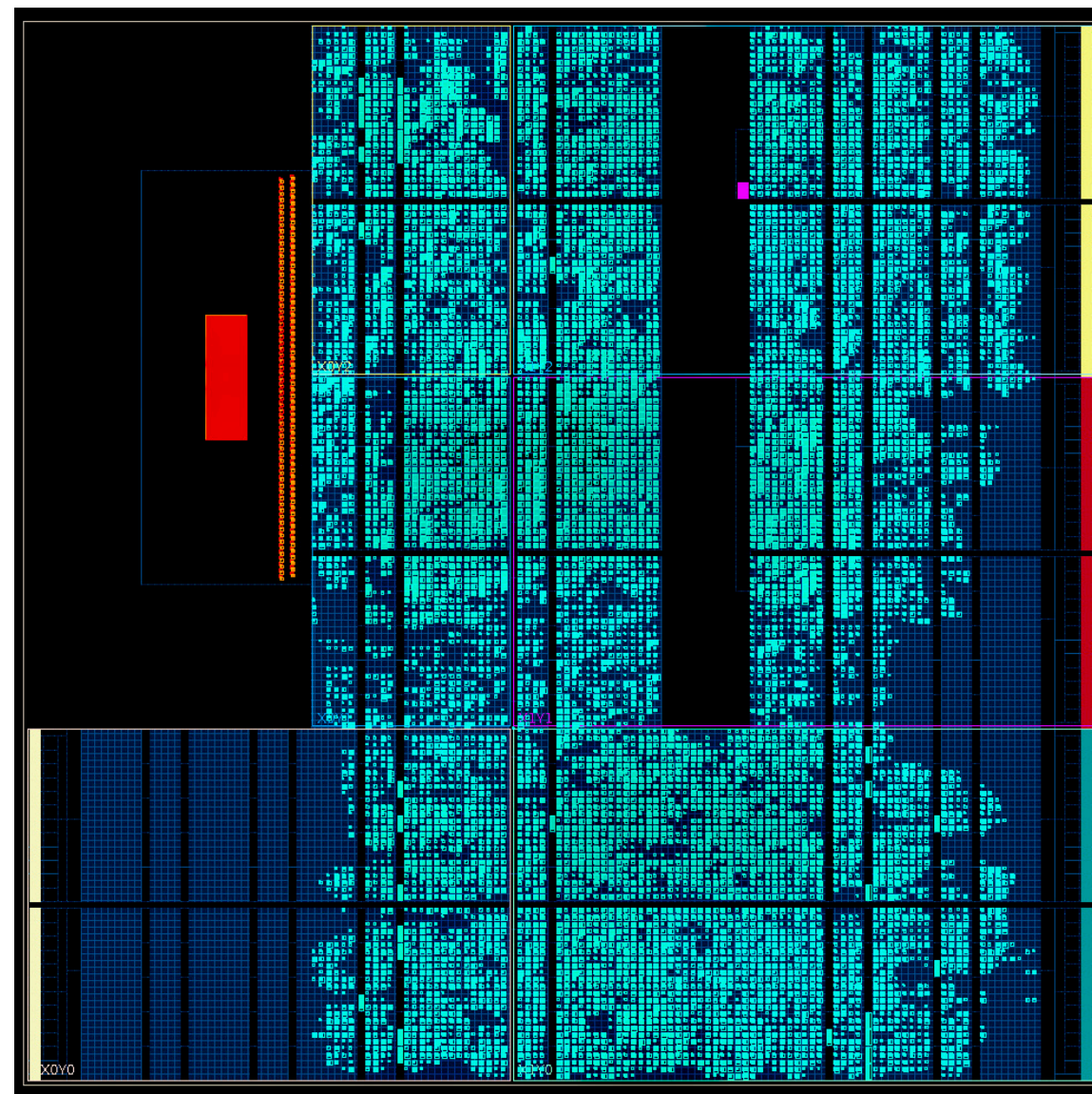


- ▶ **Codesign:** intrinsic development loop between ML design, training, and implementation
- ▶ Pruning
  - ▶ Maintain high performance while removing redundant operations
- ▶ Quantization
  - ▶ Reduce precision from 32-bit floating point to 16-bit, 8-bit, ...
- ▶ Parallelization
  - ▶ Balance parallelization (how fast) with resources needed (how costly)



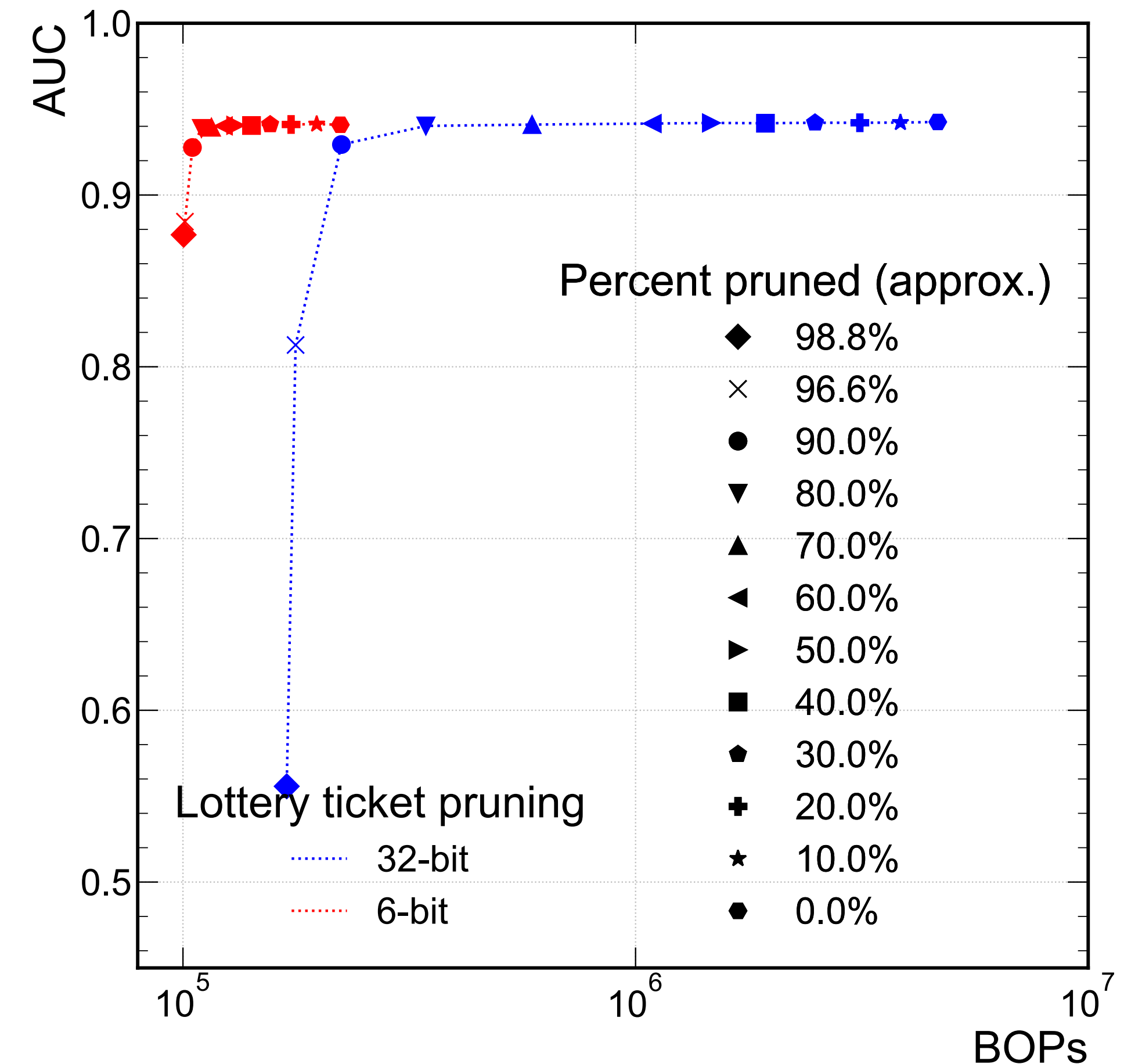


- ▶ Small NN benchmark correctly identifies particle "jets" 70-80% of the time
- ▶ Full performance with 6 bits instead of 14 bits
- ▶ Much smaller fraction of resources





- ▶ Quantization-aware pruning (QAP): iterative pruning further reduces hardware complexity of a quantized model
- ▶ After QAP,  $50 \times$  reduction in bit operations compared to the 32-bit, unpruned model



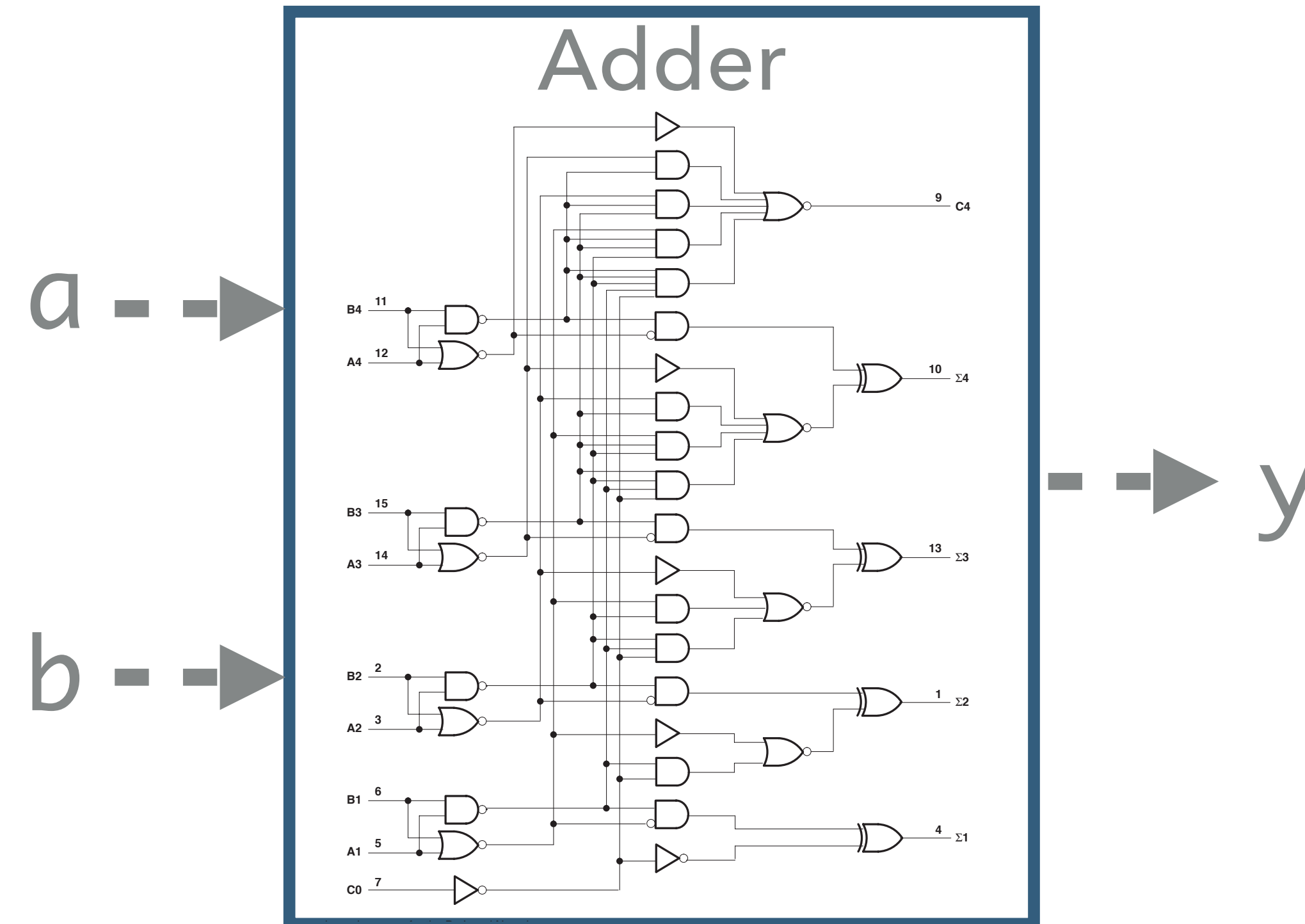
Bit operations (BOPs) definition:

[arXiv:1804.10969](https://arxiv.org/abs/1804.10969)



- ▶ Say you want to program an "adder" function on an FPGA

```
module adder(  
    input wire [4:0] a,  
    input wire [4:0] b,  
    output wire [4:0] y  
);  
    assign y = a + b;  
endmodule
```

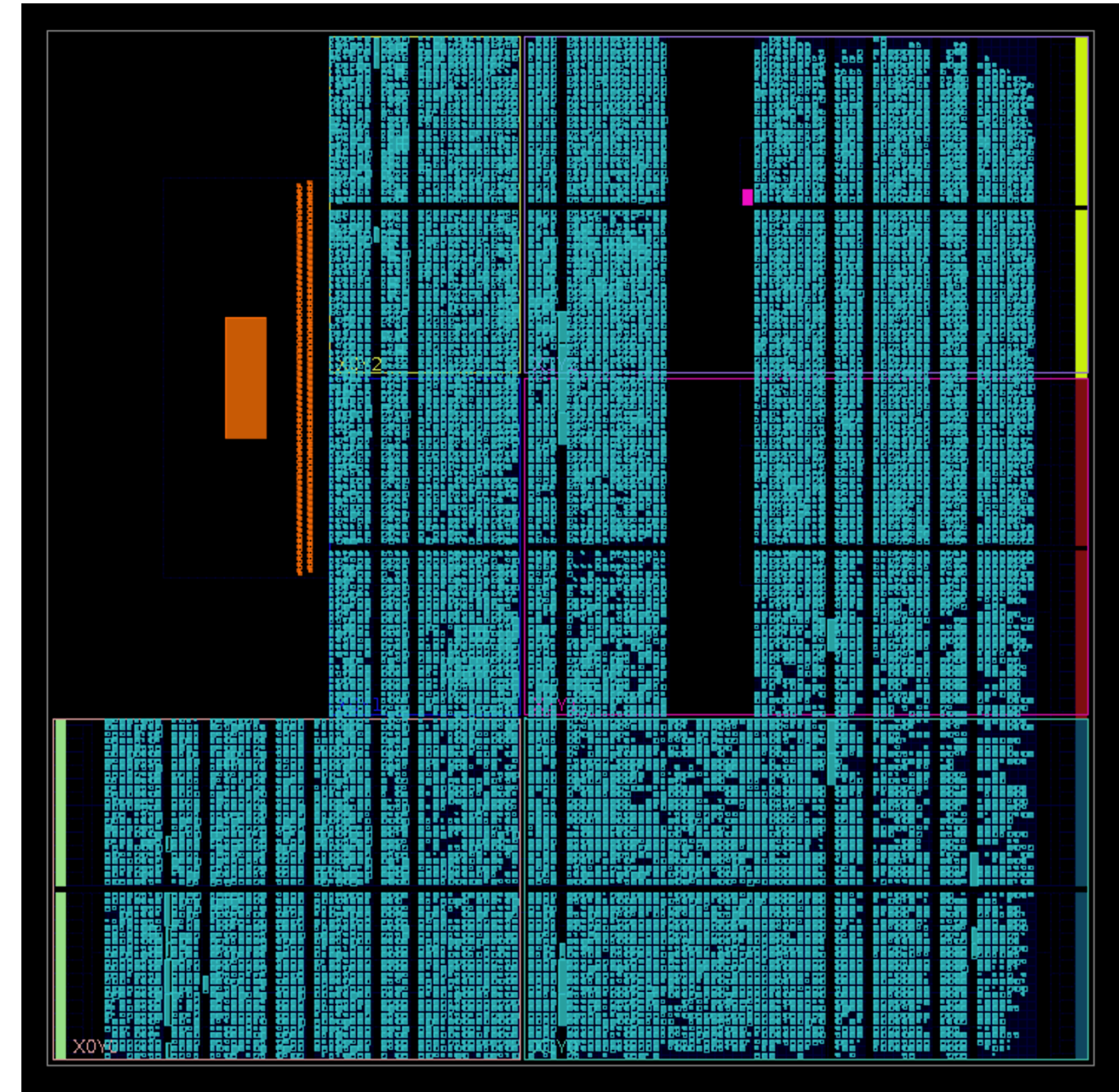
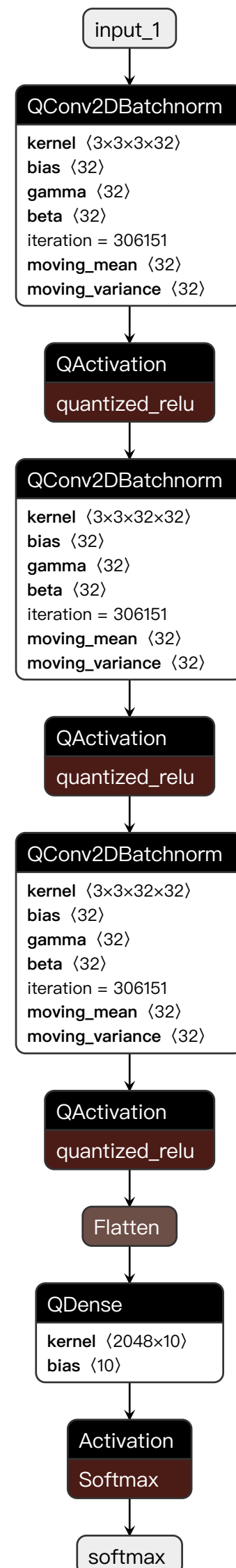


Synthesis

- ▶ Register transfer-level (RTL) code is "synthesized" into gates



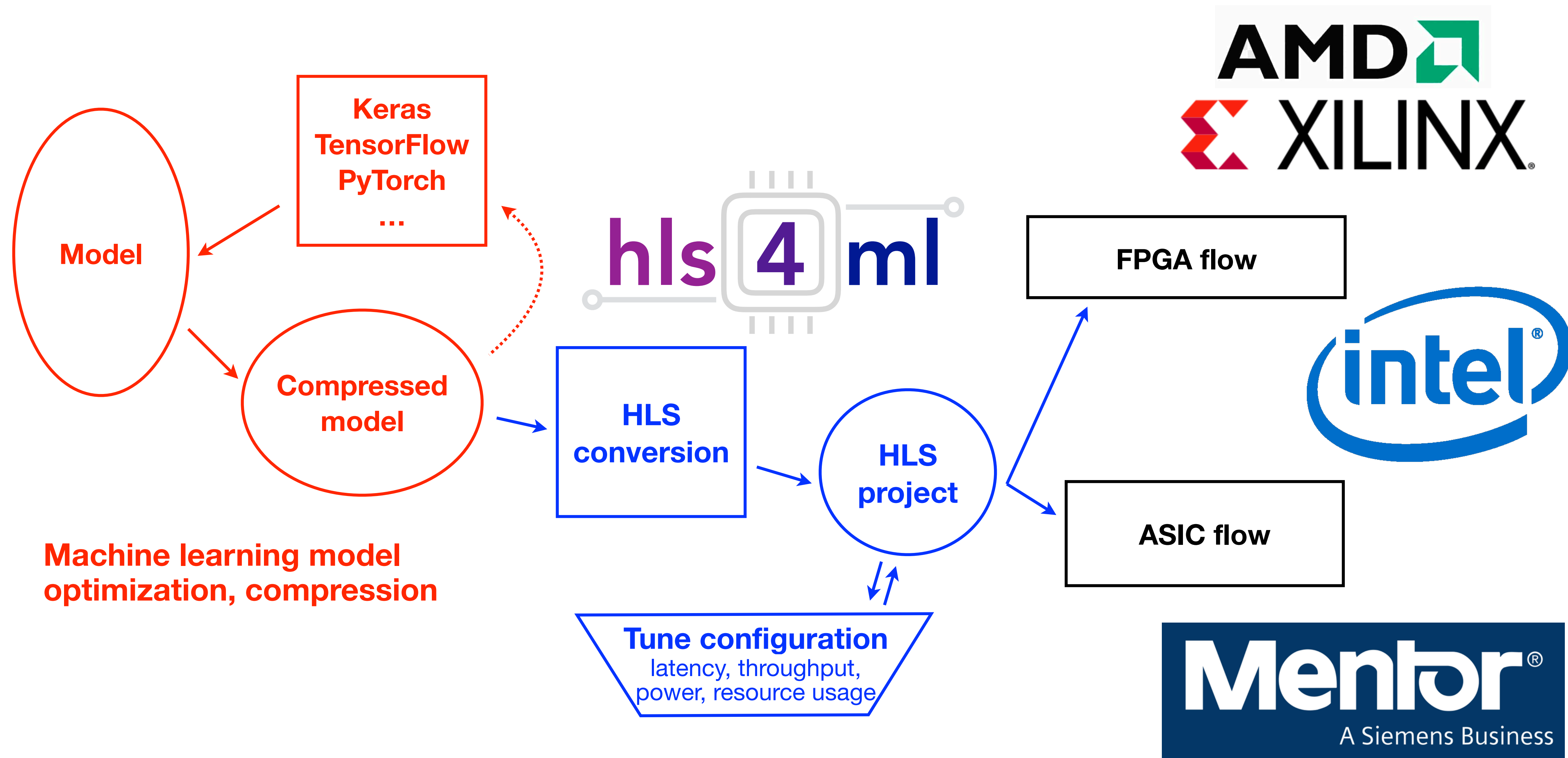
- ▶ What if instead we specify an AI model (e.g., in [QONNX](#))



High-Level Synthesis

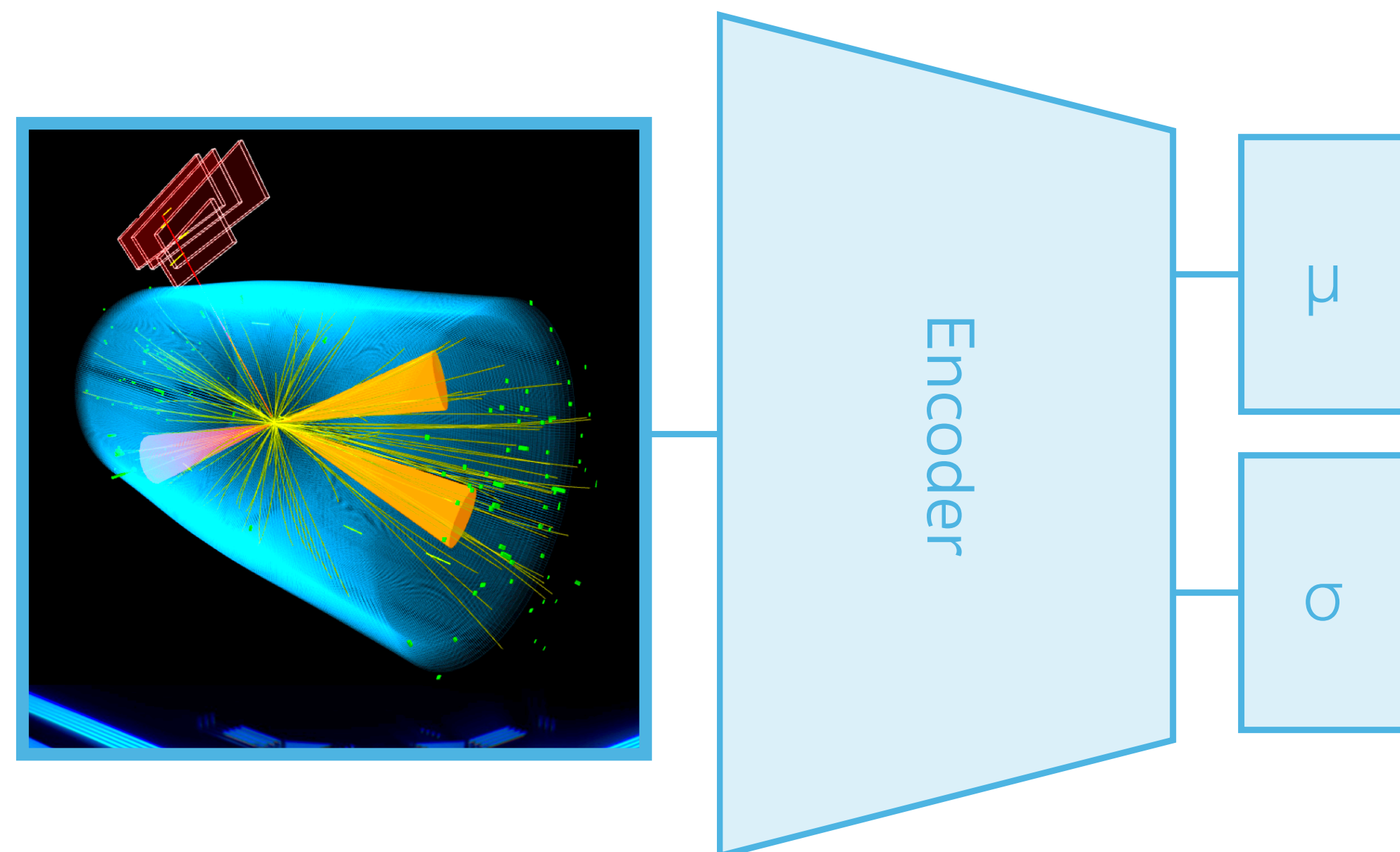


- ▶ [hls4ml](#) for scientists or ML experts to translate ML algorithms into RTL firmware





- ▶ Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious signal events will be discarded
- ▶ Can we use unsupervised algorithms to detect non-SM-like anomalies?
  - ▶ Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference
  - ▶ Variational autoencoders (VAEs): model the latent space as a probability distribution; possible to detect anomalies purely with latent space variables

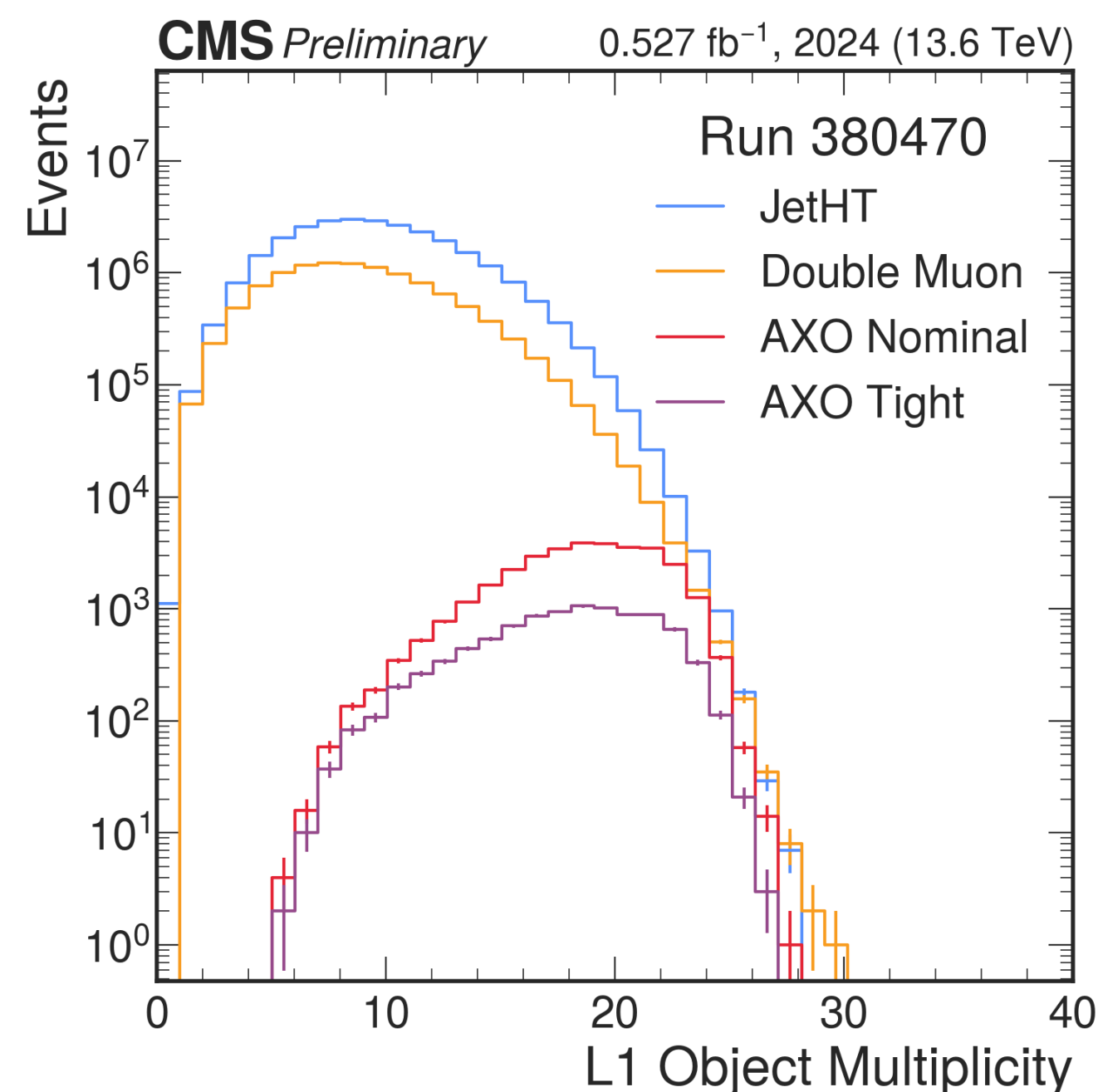
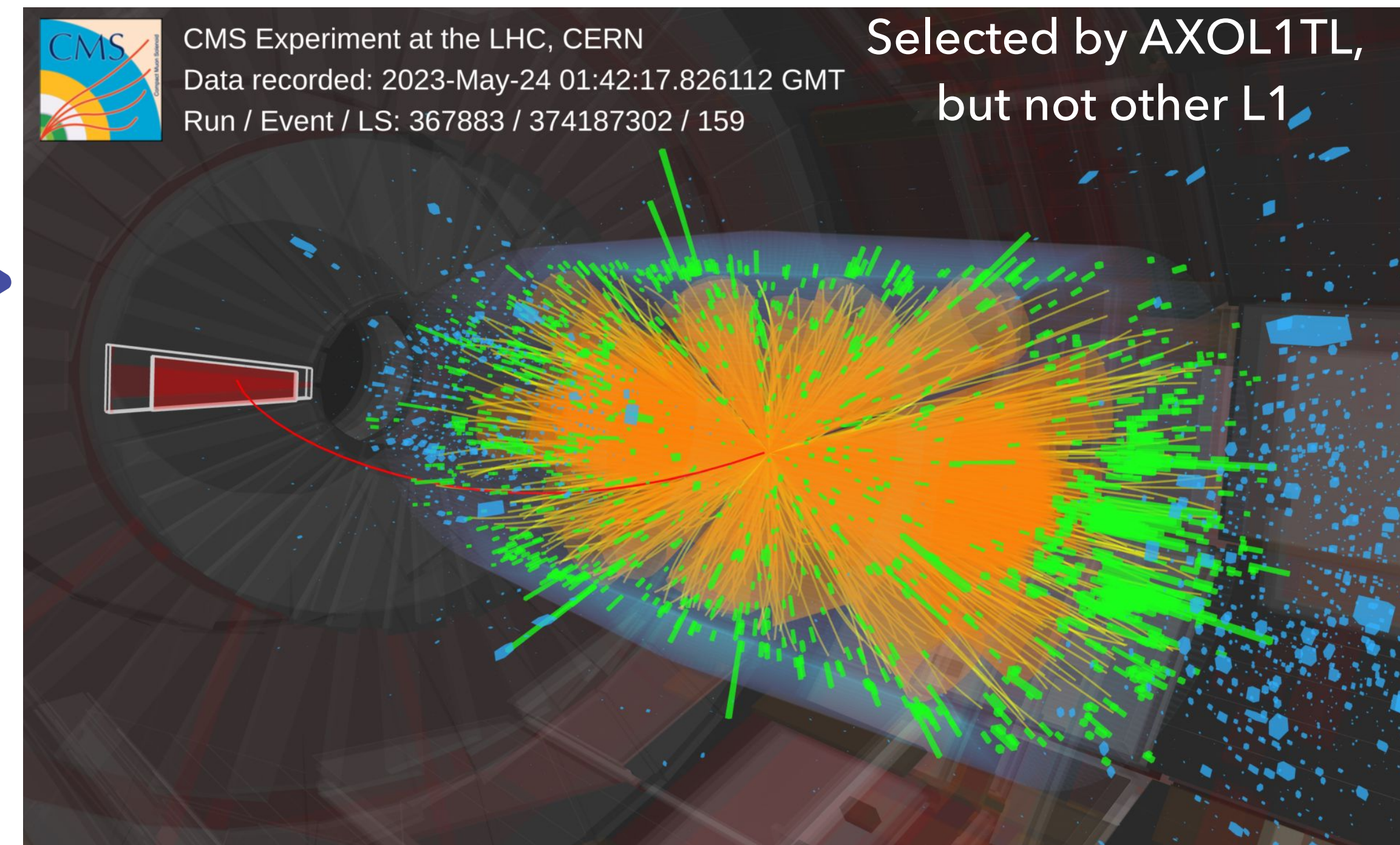
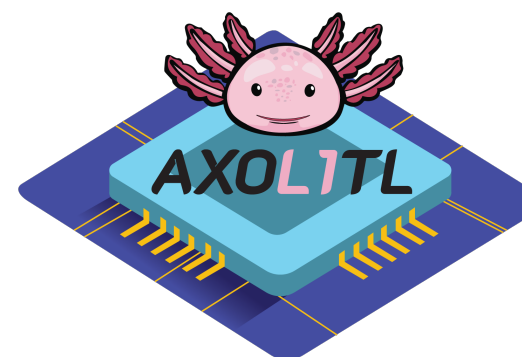


Key observation: Can build an anomaly score from the latent space of VAE directly!  
No need to run decoder!

$$R_z = \sum_i \frac{\mu_i^2}{\sigma_i^2}$$

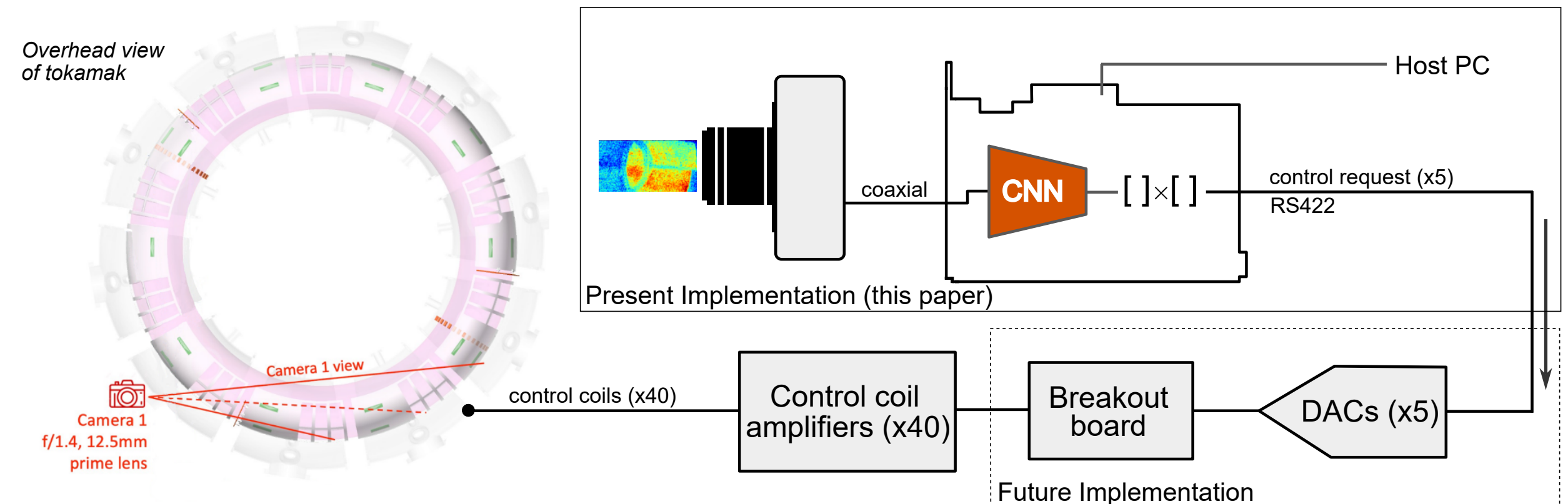
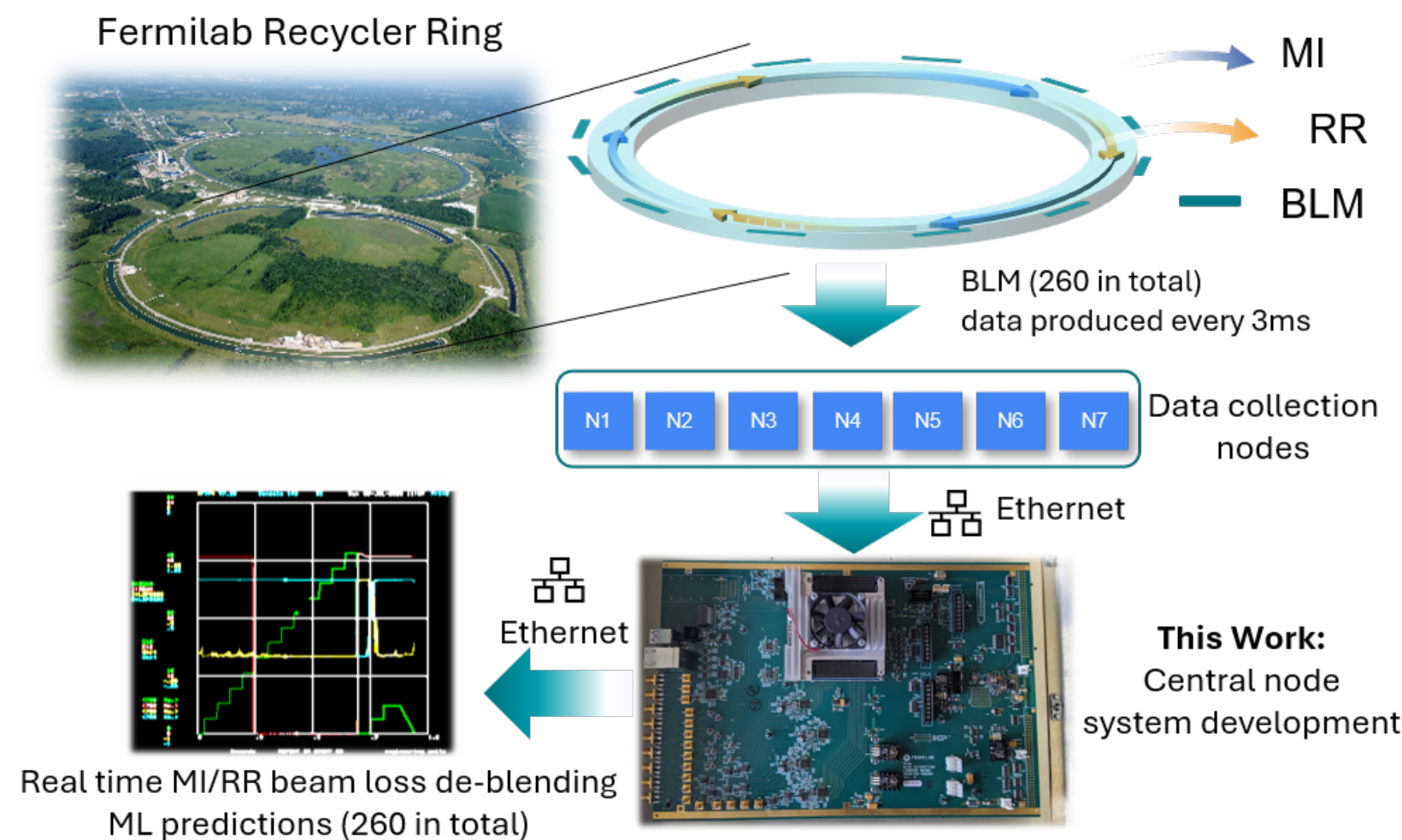
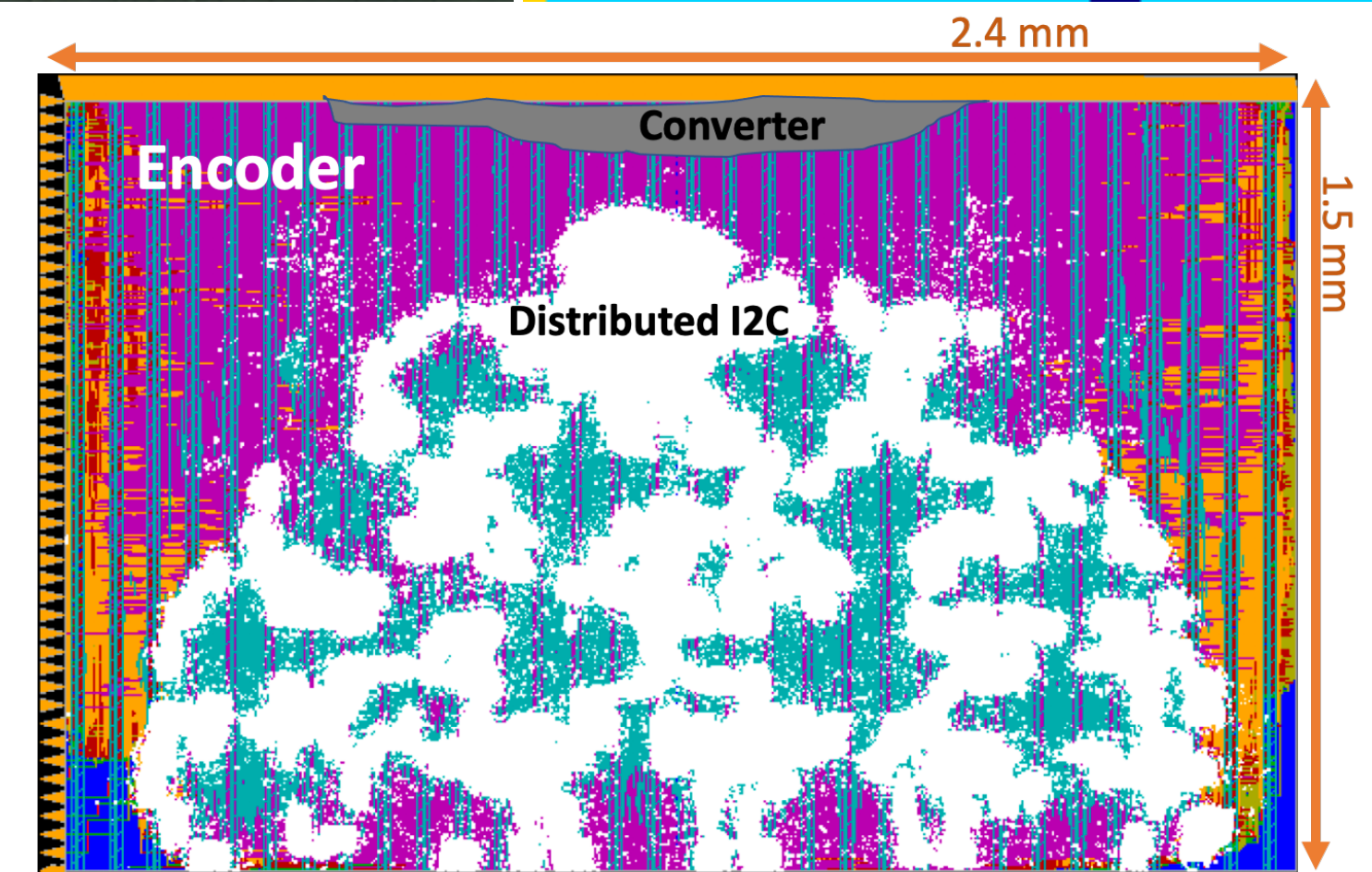


- ▶ AXOL1TL anomaly detection algorithm for the trigger based on a variational autoencoder
- ▶ Selects unique events relative to existing triggers
- ▶ Preference for high multiplicity events



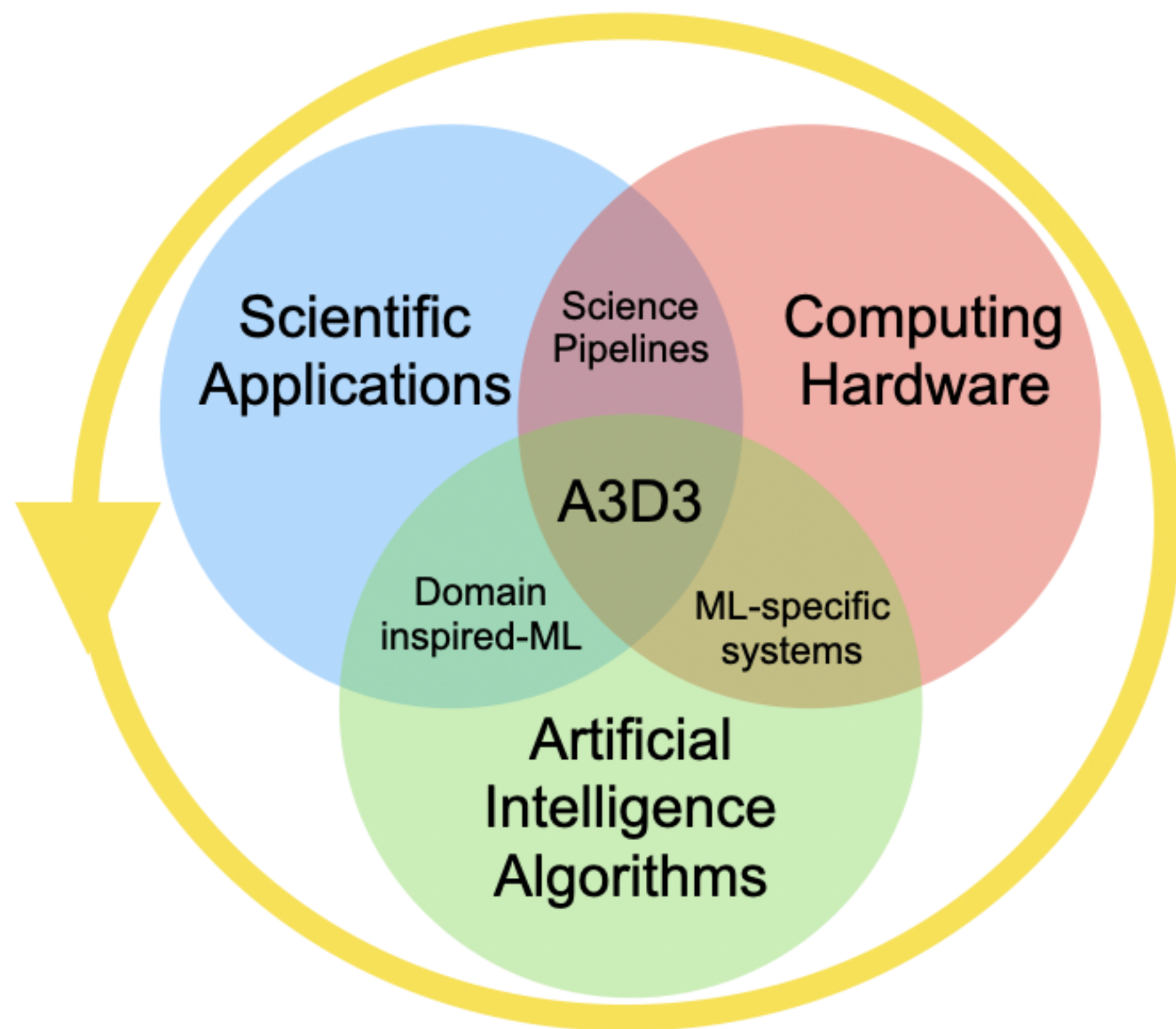


- ▶ Though [hls4ml](https://github.com/hls4ml) developed for particle physics, has seen widespread use for
  - ▶ Self-driving cars [[arXiv:2205.07690](https://arxiv.org/abs/2205.07690)]
  - ▶ Fusion devices [[arXiv:2312.00128](https://arxiv.org/abs/2312.00128)]
  - ▶ Steering particle beams [[arXiv:2011.07371](https://arxiv.org/abs/2011.07371), [arXiv:2311.05716](https://arxiv.org/abs/2311.05716)]
  - ▶ Data compression at the edge [[arXiv:2105.01683](https://arxiv.org/abs/2105.01683)]





- ▶ Tightly coupled organization of domain scientists, computer scientists, and engineers that unite three core components which are essential to achieve real-time AI to transform science: AI techniques, Computing Hardware, Scientific Applications
- ▶ Check the [a3d3.ai](http://a3d3.ai) for events and more information!



[PHY-2117997](http://PHY-2117997)





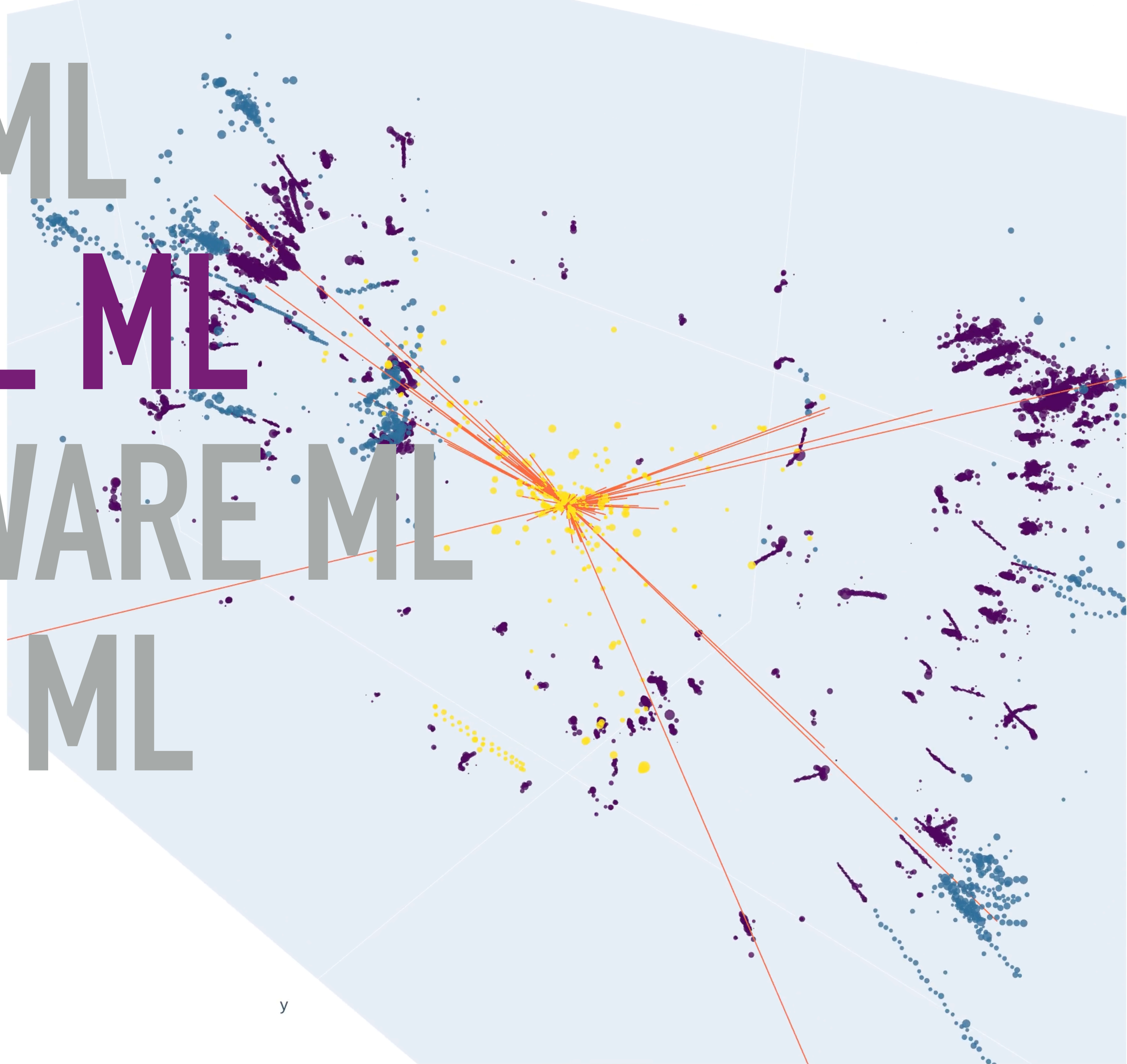
ULTRAFAST ML

**MULTIMODAL ML**

PHYSICS-AWARE ML

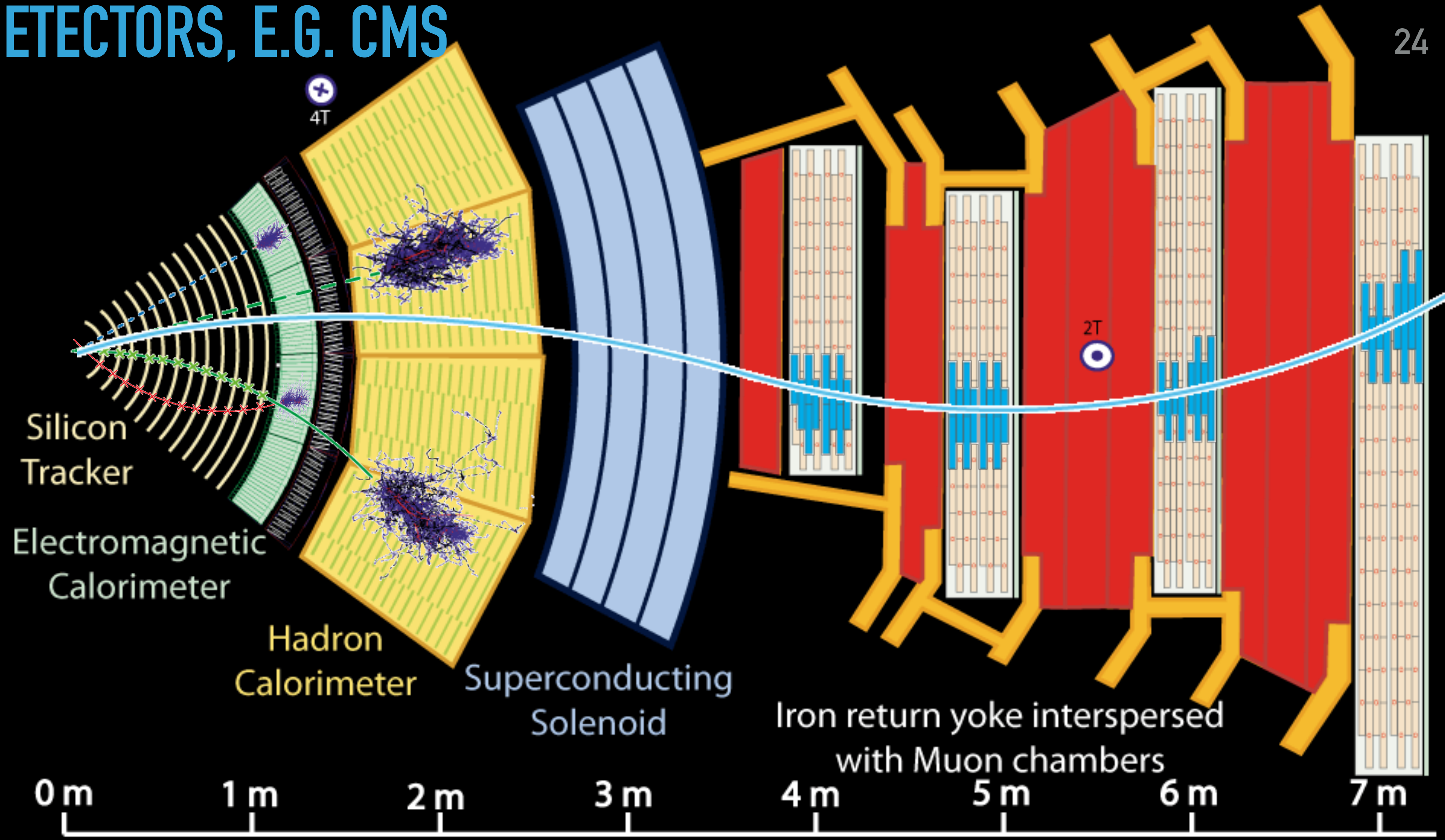
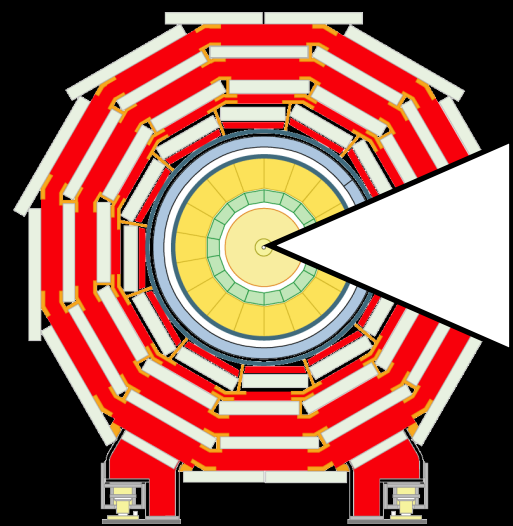
GENERATIVE ML

OUTLOOK





# MULTILAYERED DETECTORS, E.G. CMS



Current and future multilayered detectors...

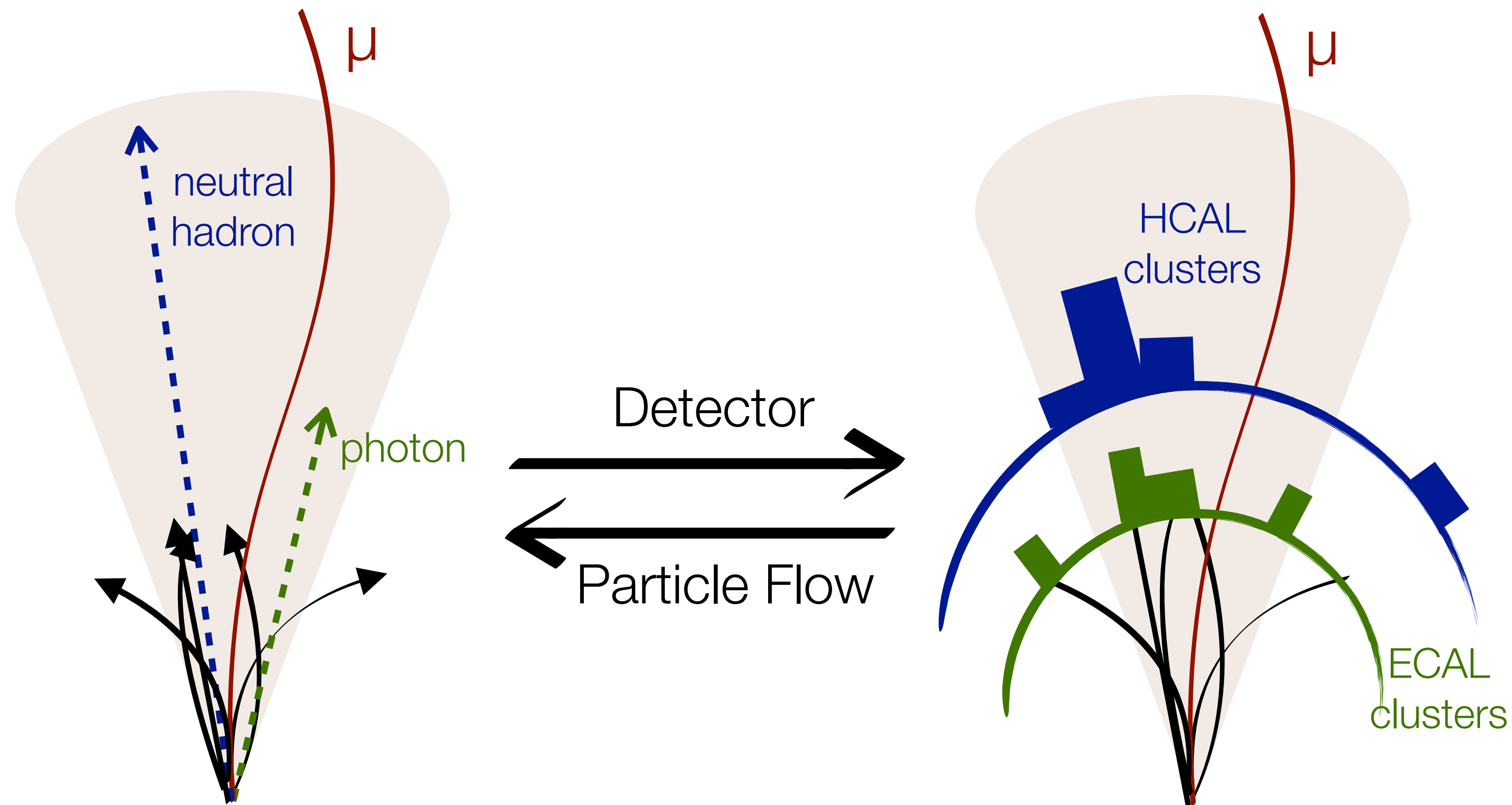
Require complex pattern recognition

Key:

- Muon
- Electron
- Charged Hadron (e.g. Pion)
- Neutral Hadron (e.g. Neutron)
- Photon

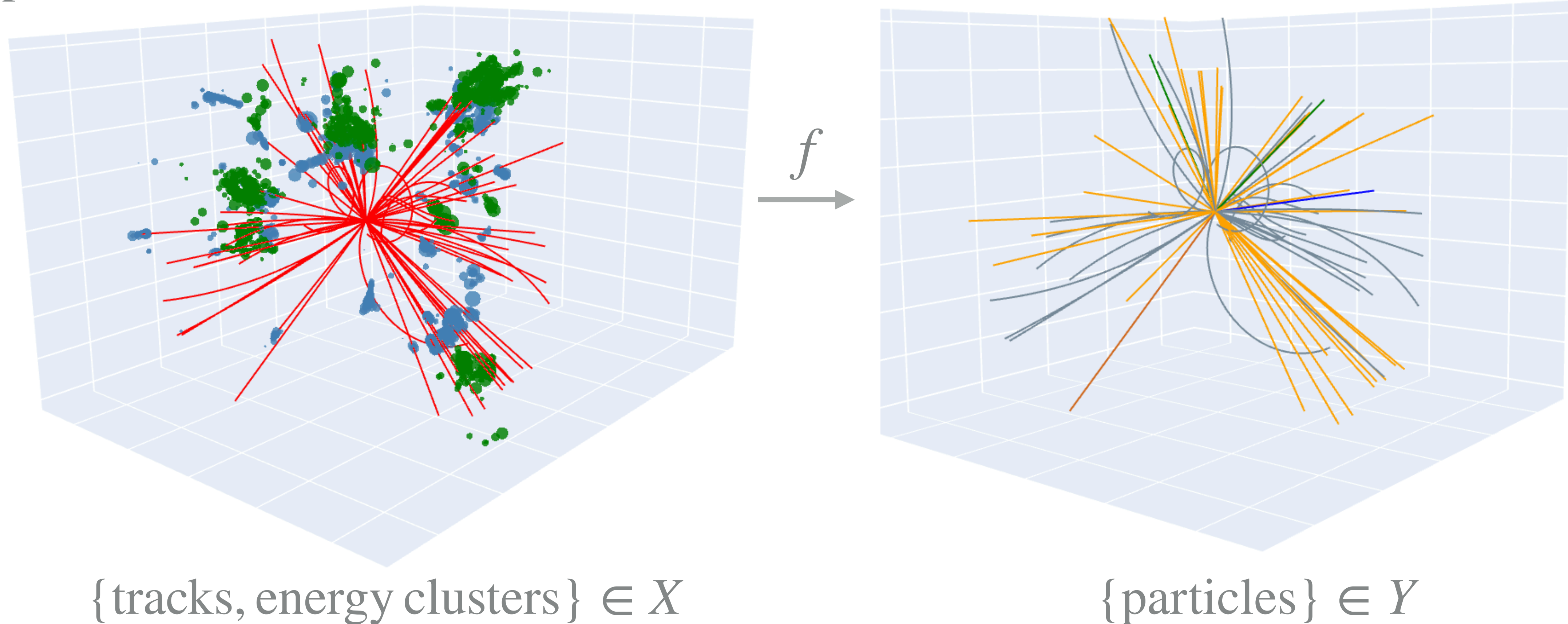


- ▶ Particles interact with detector, leaving energy deposits and tracks
- ▶ Combination of multimodal information from complementary detectors to produce particle-level interpretation of the event based on complex, hand-tuned heuristics





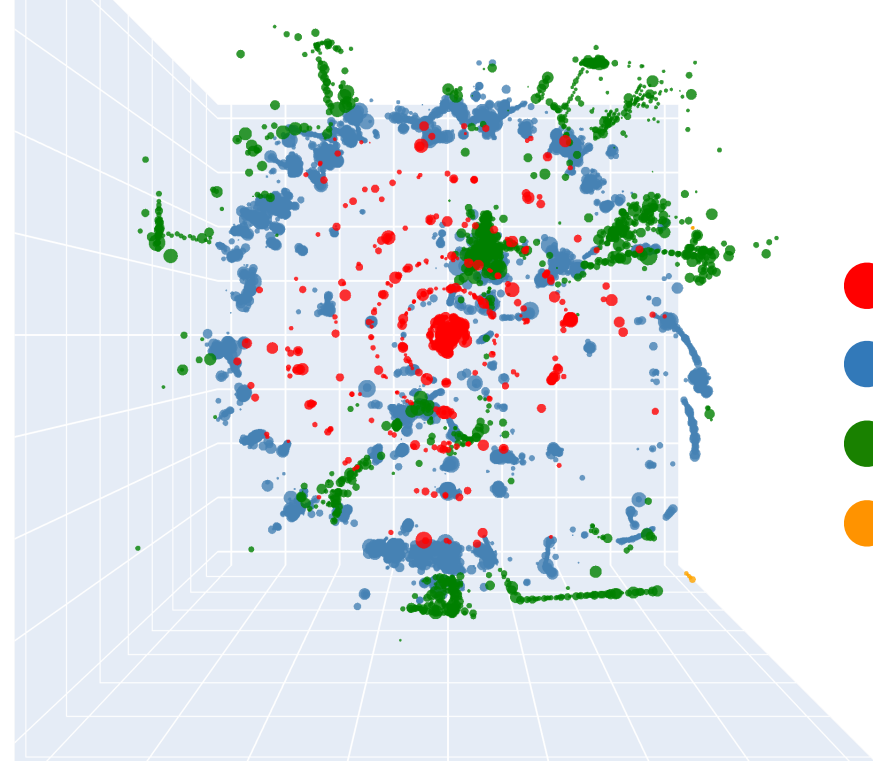
- ▶ Can we instead formulate PF as an ML task (naturally “tunable” through re-training and portable to new hardware)?
- ▶ Learn a “set-to-set” function  $f: X \rightarrow Y$ , where  $\{\text{tracks, energy clusters}\} \in X$  and  $\{\text{particles}\} \in Y$





~100k / event

Raw detector hits

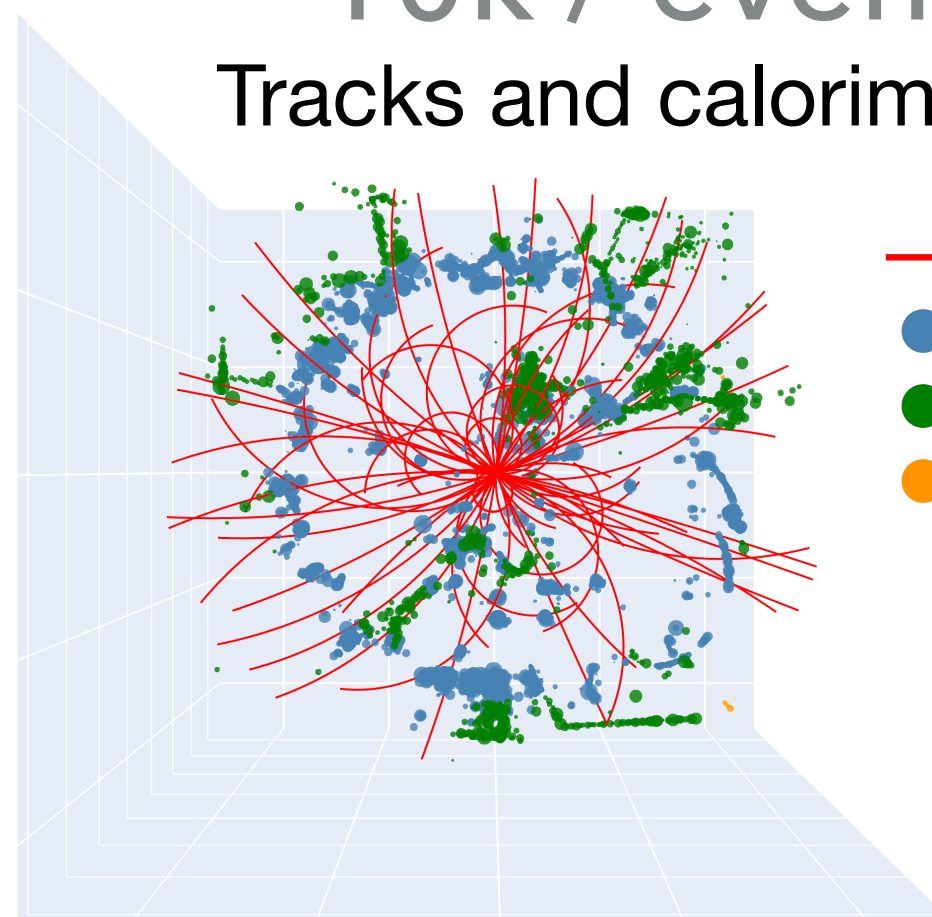


- Raw tracker hit
- Raw ECAL hit
- Raw HCAL hit
- Raw Muon chamber hit

Charged particle tracking

~10k / event

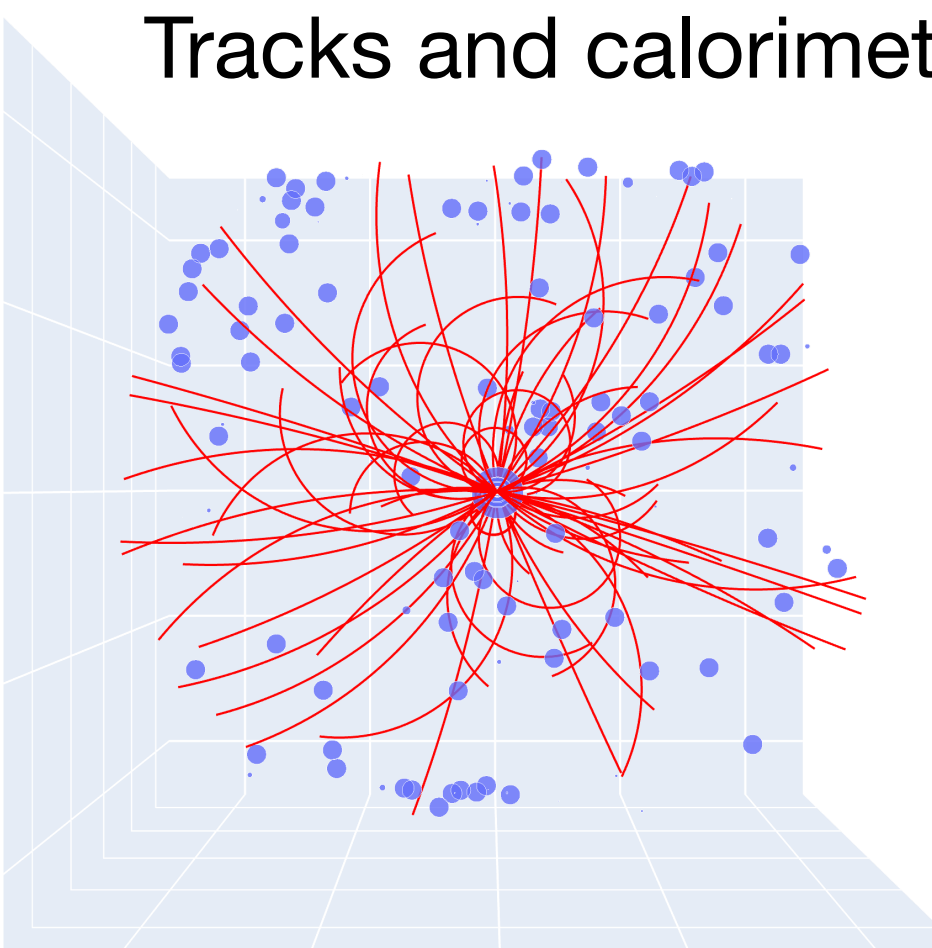
Tracks and calorimeter hits



- Track
- Raw ECAL hit
- Raw HCAL hit
- Raw Muon chamber hit

Calorimeter clustering

Tracks and calorimeter clusters

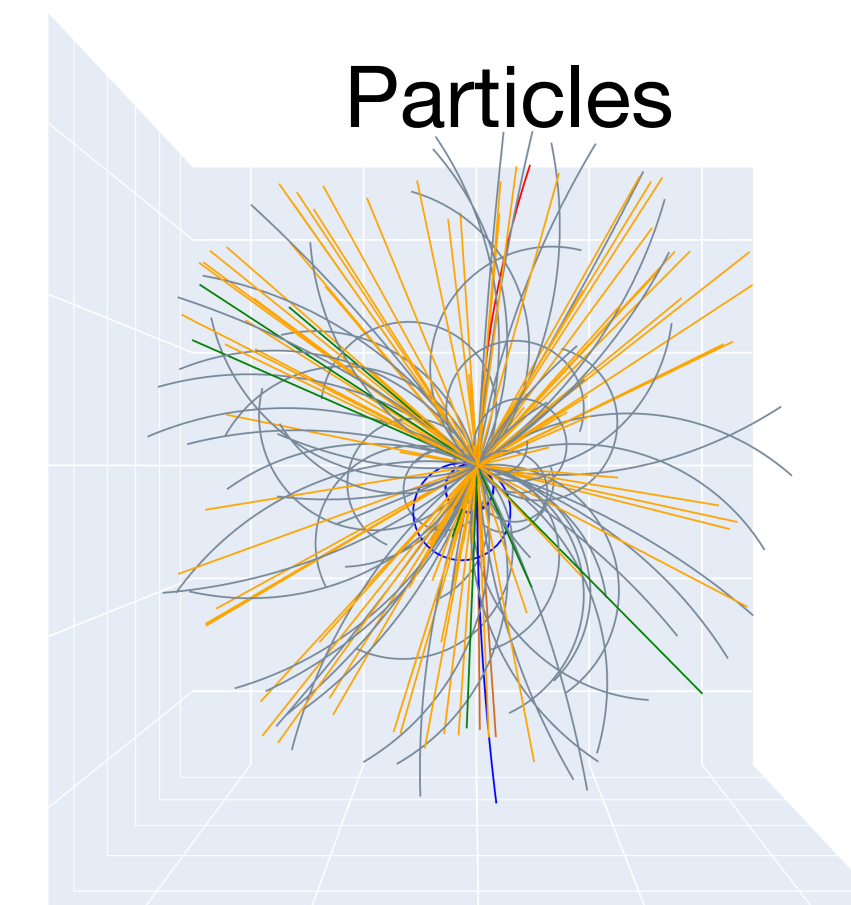


- Track
- ECAL or HCAL cluster

~300-500 / event

▶ 2.5 TB, 6 million events total

Particles



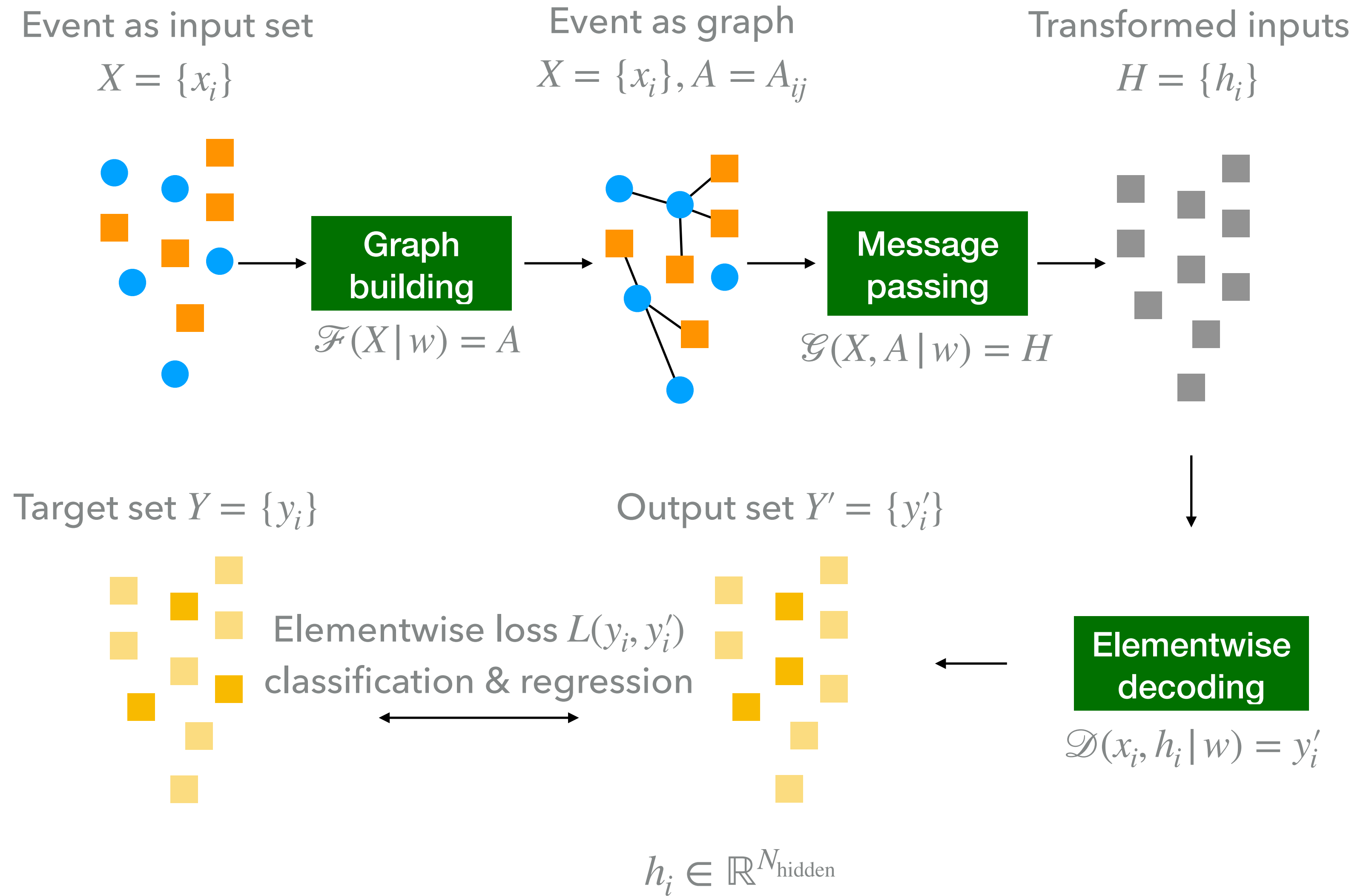
- Charged hadron
- Photon
- Neutral hadron
- Electron
- Muon

Cluster-based ML particle-flow reconstruction

~100-300 / event



- ▶ Convert input set to a locally, **sparsely connected** graph
- ▶ Message-passing NN to transform features
- ▶ Decode transformed inputs elementwise
- ▶ (During training) Compare to target set, optimize weights

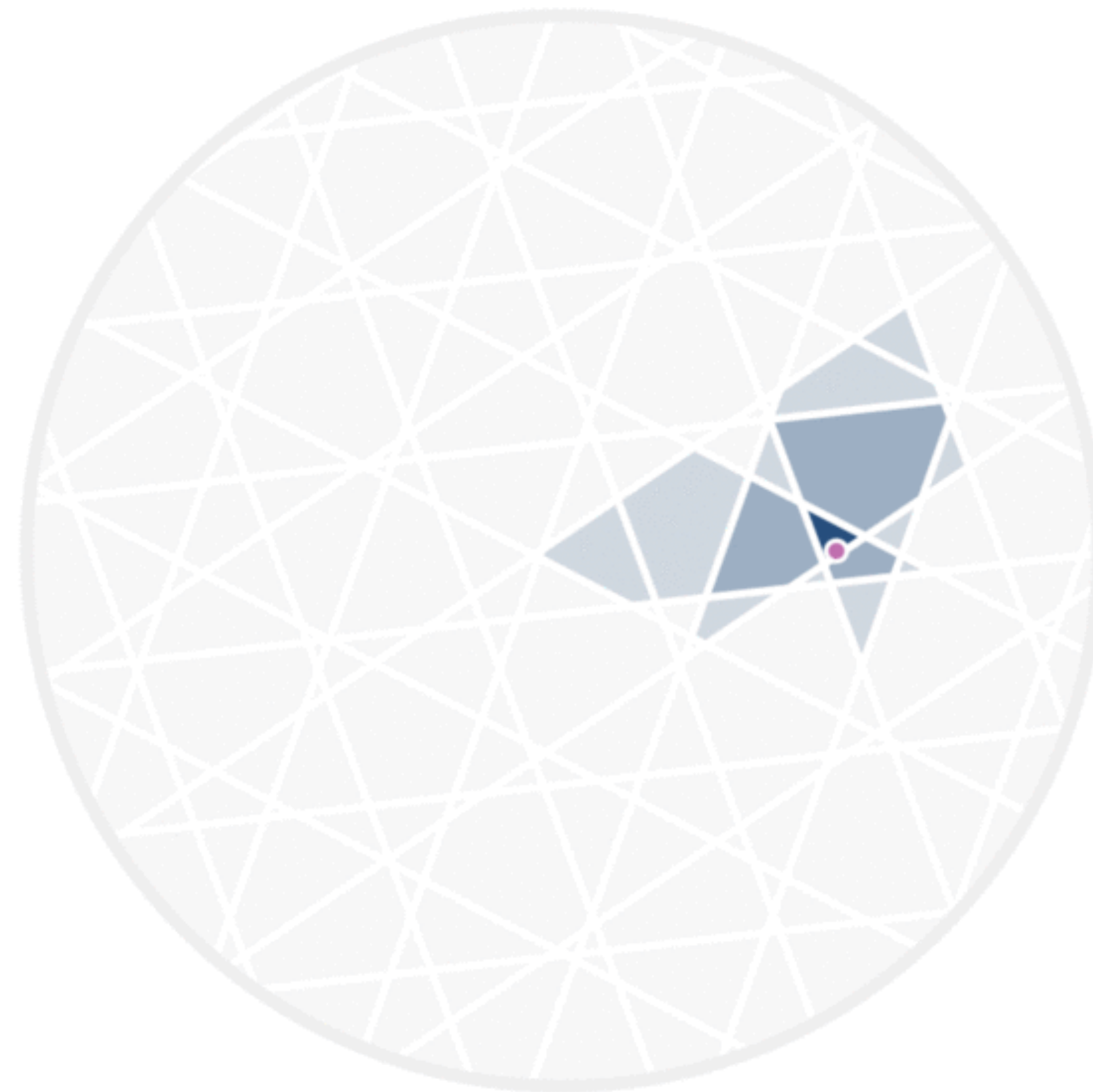


Trainable neural networks:  $\mathcal{F}, \mathcal{G}, \mathcal{D}$

- Track, ■ Calorimeter cluster, ■ Encoded element
- Target (predicted) particle, ■ No target (predicted) particle



Locality-sensitive hashing reduces graph-building complexity



5 random hash functions, where  
darkest blue = 5 hash collisions,  
lightest blue = 3 hash collisions



- ▶ Hyperparameter optimization requires large compute

JURECA Supercomputer at Jülich Supercomputing Centre

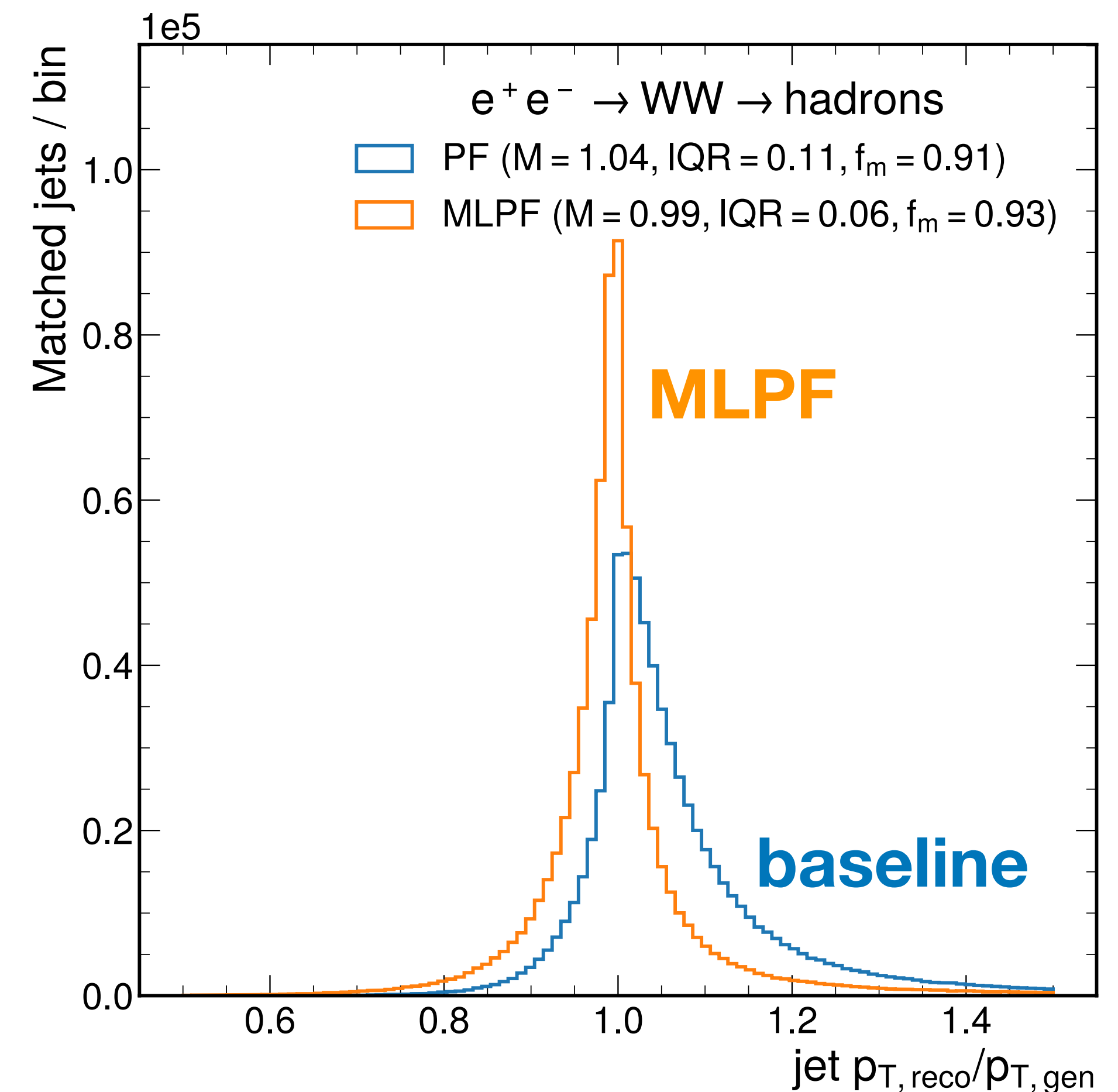
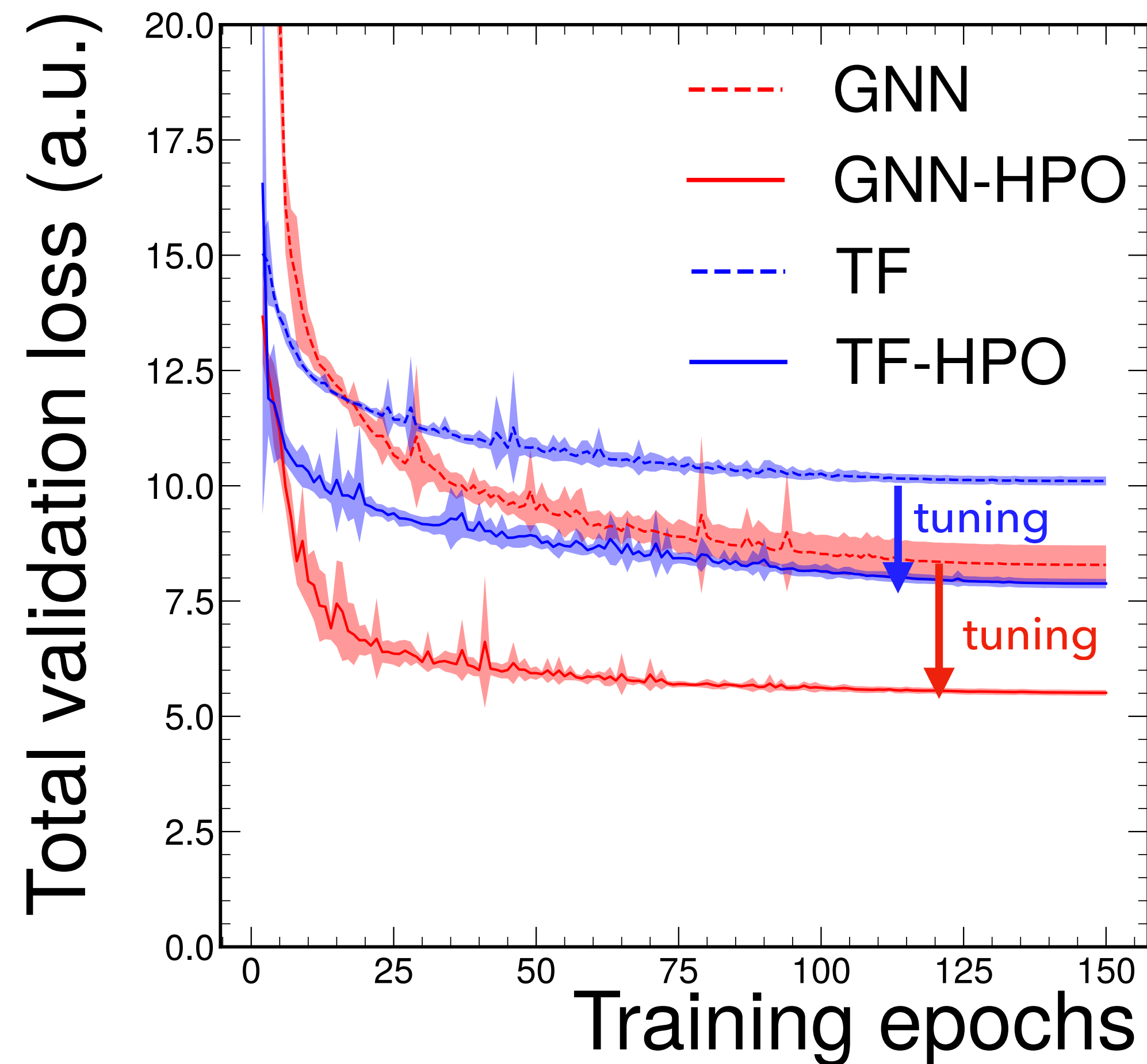


Voyager Supercomputer at San Diego Supercomputer Center



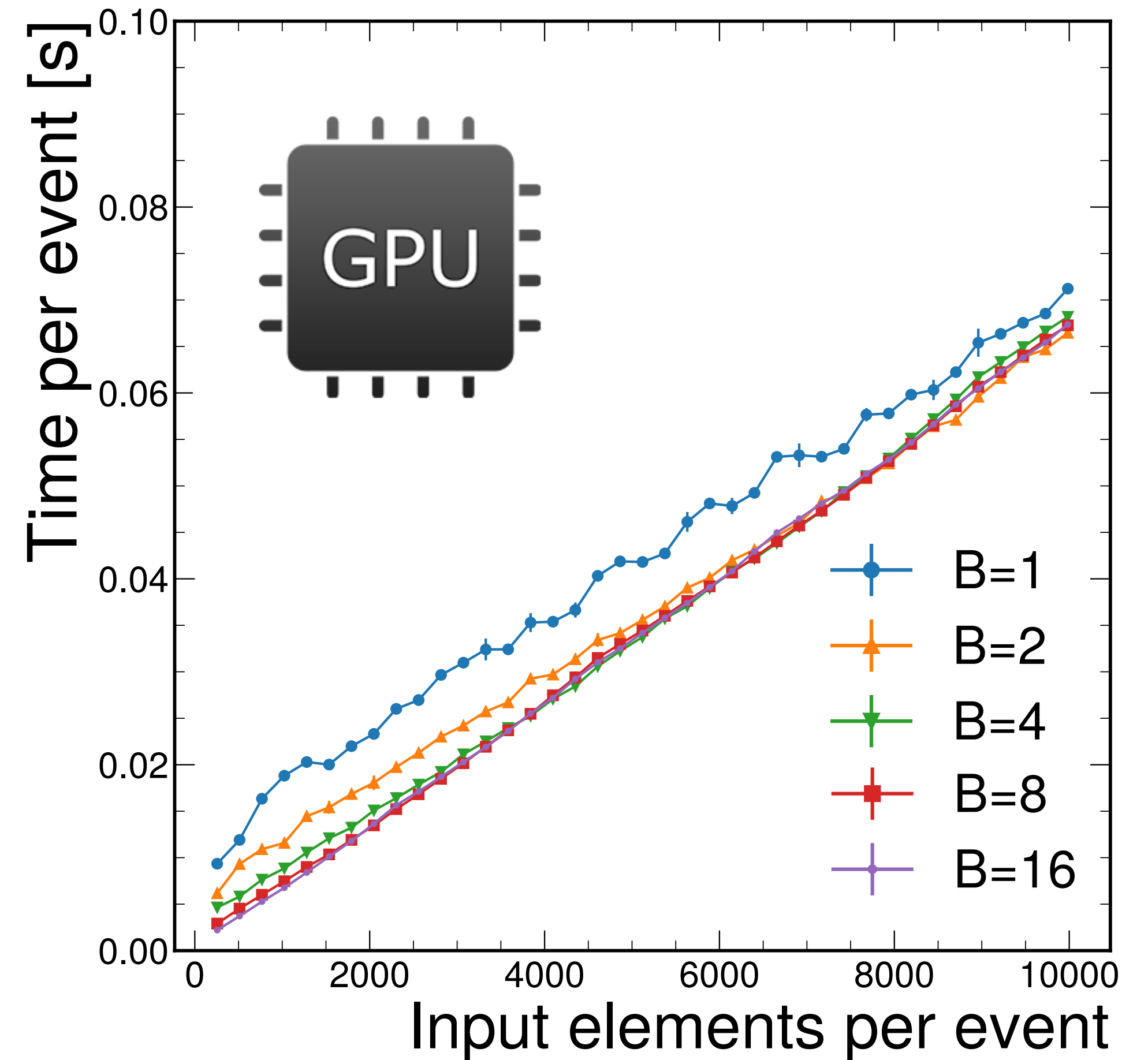
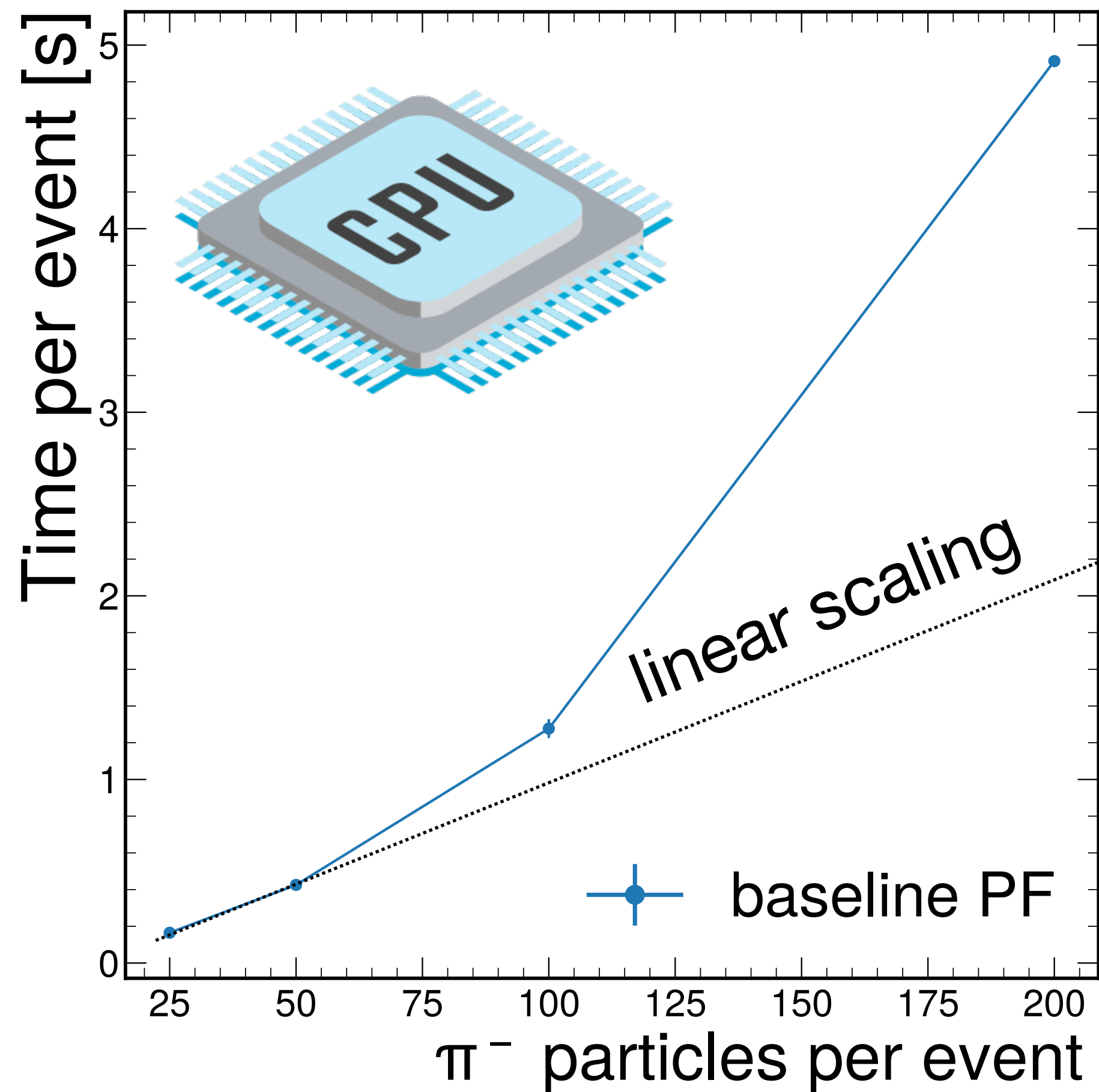


- ▶ Tuning improves particle-level performance dramatically
- ▶ Though we optimize a particle-level loss, also achieve better energy resolution

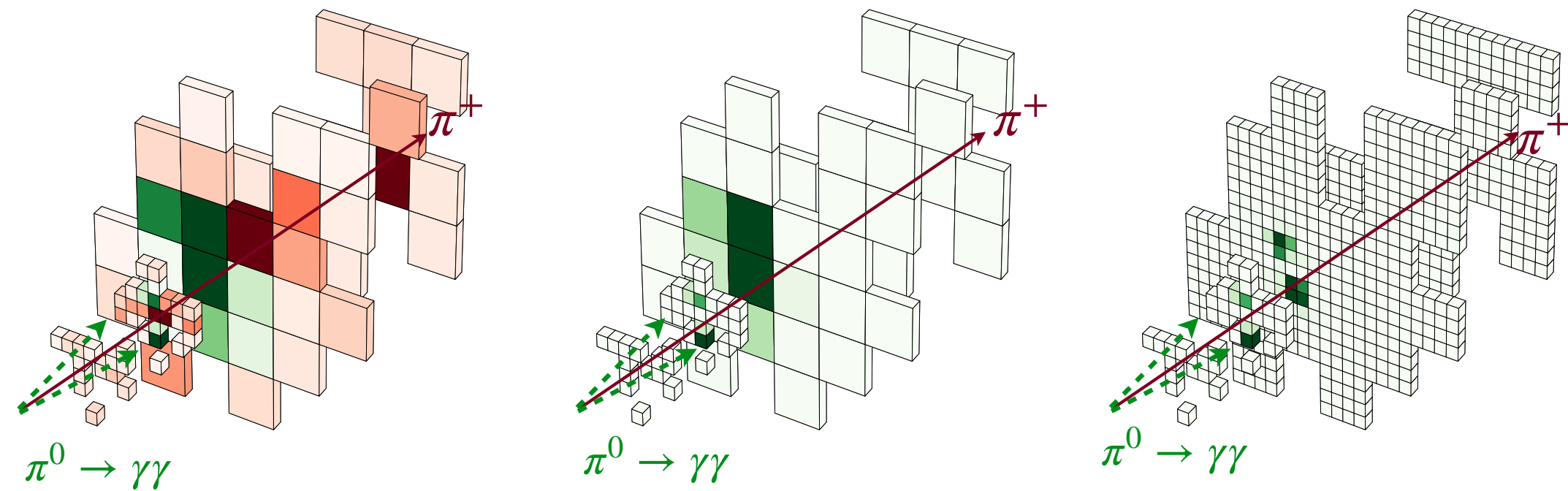




- ▶ Baseline algorithm runs only on CPU, scales ~quadratically, runs in seconds
- ▶ ML model scales linearly, runs in milliseconds on a consumer 8 GB GPU







## Towards a Computer Vision Particle Flow <sup>★</sup>

Francesco Armando Di Bello<sup>a,3</sup>, Sanmay Ganguly<sup>b,1</sup>, Eilam Gross<sup>1</sup>, Marumi Kado<sup>3,4</sup>, Michael Pitt<sup>2</sup>, Lorenzo Santi<sup>3</sup>, Jonathan Shlomi<sup>1</sup>

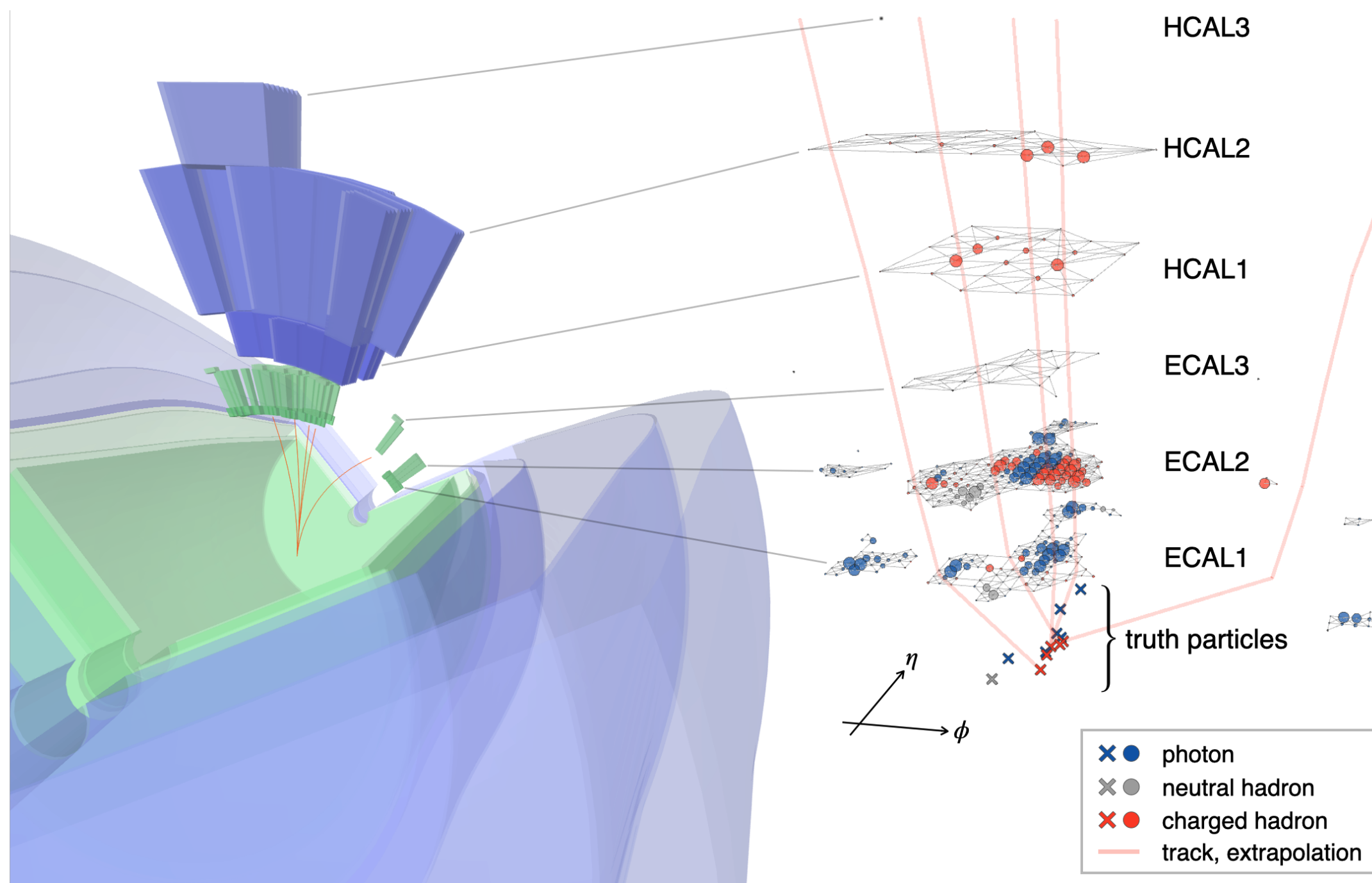
<sup>1</sup>Weizmann Institute of Science, Rehovot 76100, Israel

<sup>2</sup>CERN, CH 1211, Geneva 23, Switzerland

<sup>3</sup>Università di Roma Sapienza, Piazza Aldo Moro, 2, 00185 Roma, Italy e INFN, Italy

<sup>4</sup>Université Paris-Saclay, CNRS/IN2P3, IJCLab, 91405, Orsay, France

[Eur. Phys. J. C 81, 107 \(2021\)](#)



## Reconstructing particles in jets using set transformer and hypergraph prediction networks

Francesco Armando Di Bello<sup>1,a</sup>, Etienne Dreyer<sup>2,b</sup>, Sanmay Ganguly<sup>3</sup>, Eilam Gross<sup>2</sup>, Lukas Heinrich<sup>4</sup>, Anna Ivina<sup>2</sup>, Marumi Kado<sup>5,6</sup>, Nilotpal Kakati<sup>2,c</sup>, Lorenzo Santi<sup>6</sup>, Jonathan Shlomi<sup>2</sup>, Matteo Tusoni<sup>6</sup>

<sup>1</sup> INFN and University of Genova

<sup>2</sup>Weizmann Institute of Science

<sup>3</sup>ICEPP, University of Tokyo

<sup>4</sup>Technical University of Munich

<sup>5</sup>Max Planck Institute for Physics

<sup>6</sup>INFN and Sapienza University of Rome

[Eur. Phys. J. C 83, 596 \(2023\)](#)



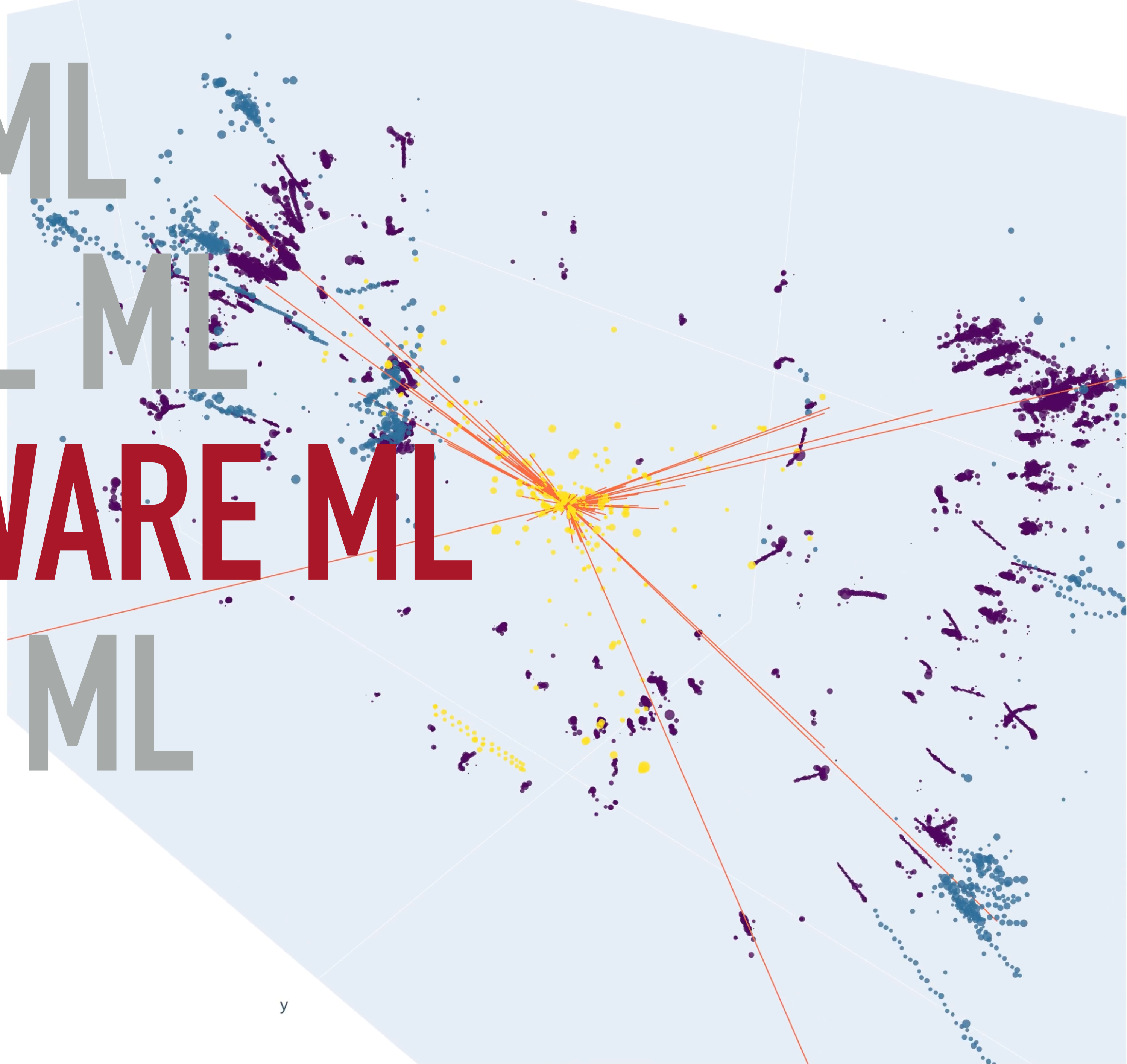
ULTRAFAST ML

MULTIMODAL ML

**PHYSICS-AWARE ML**

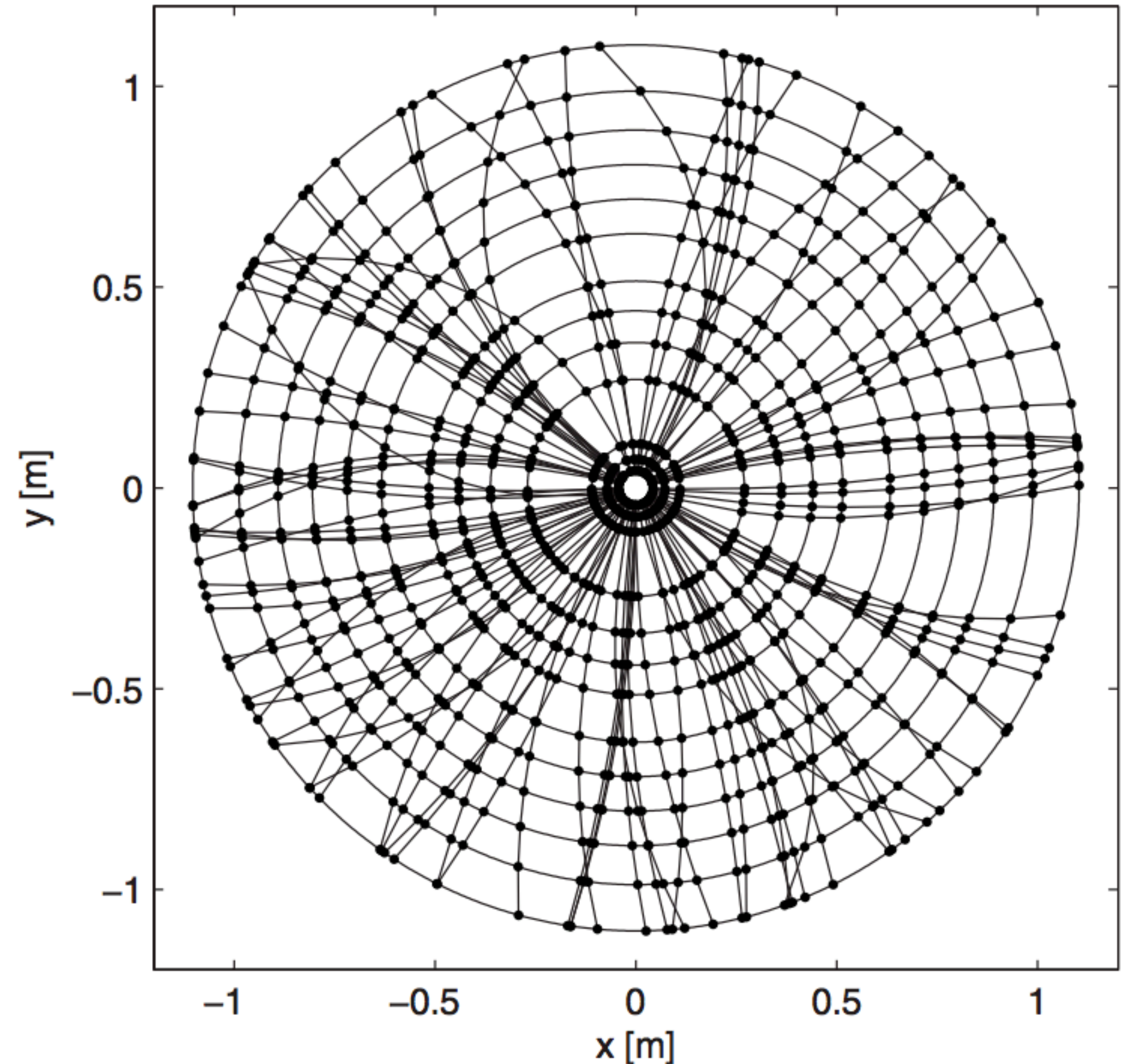
GENERATIVE ML

OUTLOOK



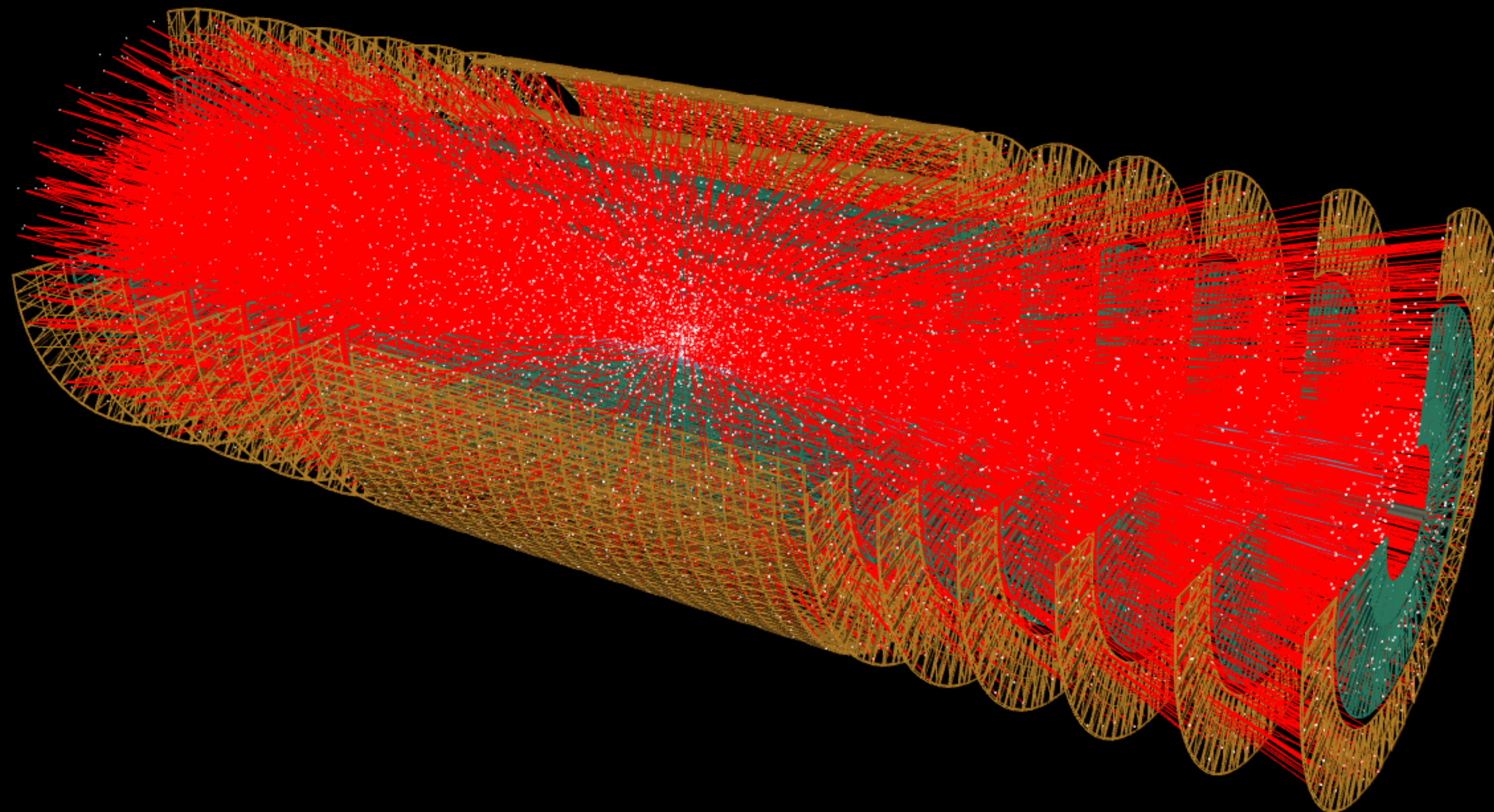


- ▶ Particle tracking is a classic pattern recognition task
- ▶ From a set of hits sampled sparsely in 3D, reconstruct the helical trajectories of particles
- ▶ Traditional algorithms scale worse than  $N^2$  in the number of hits  $N$
- ▶ How about GNNs and transformers?



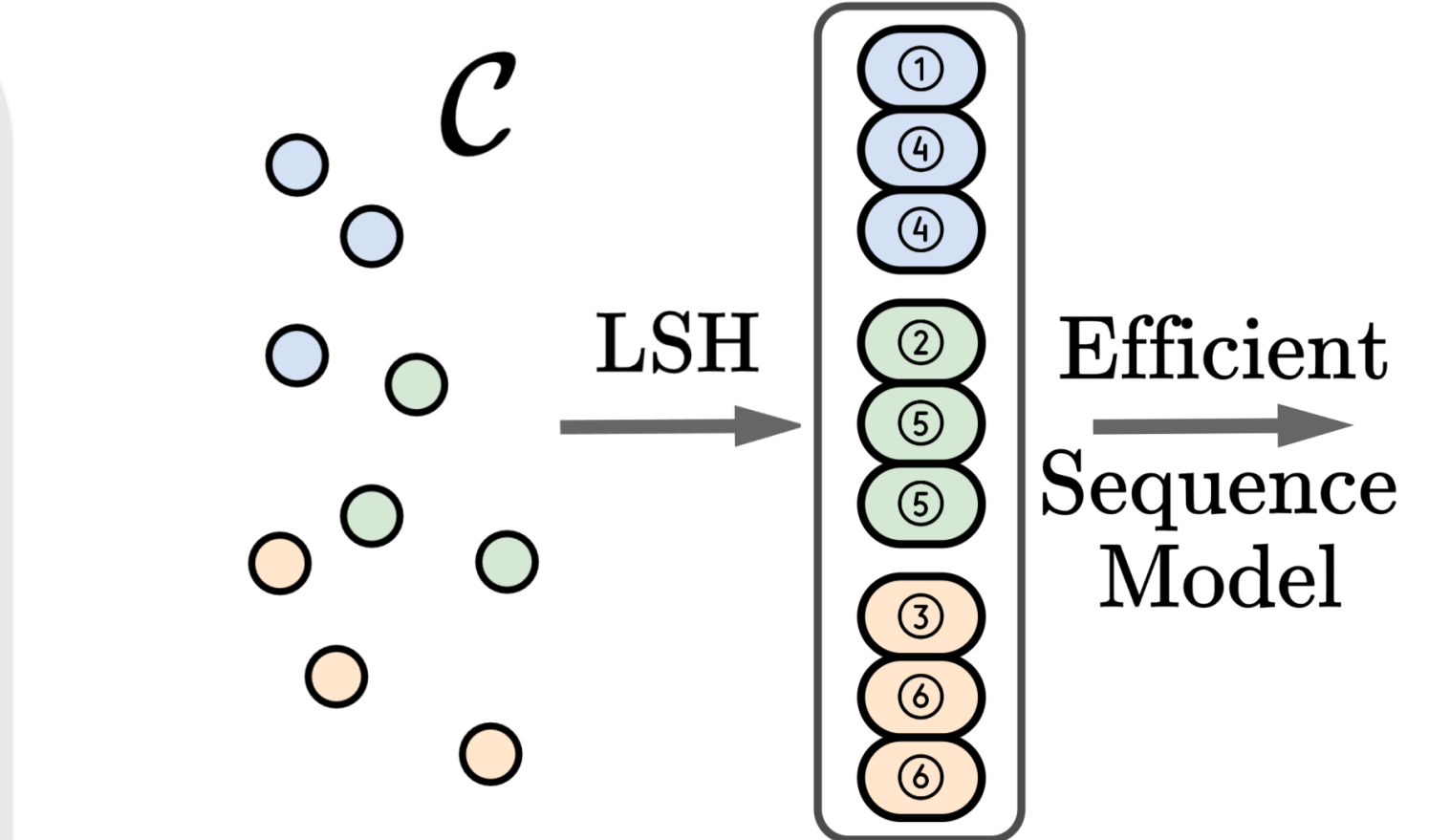
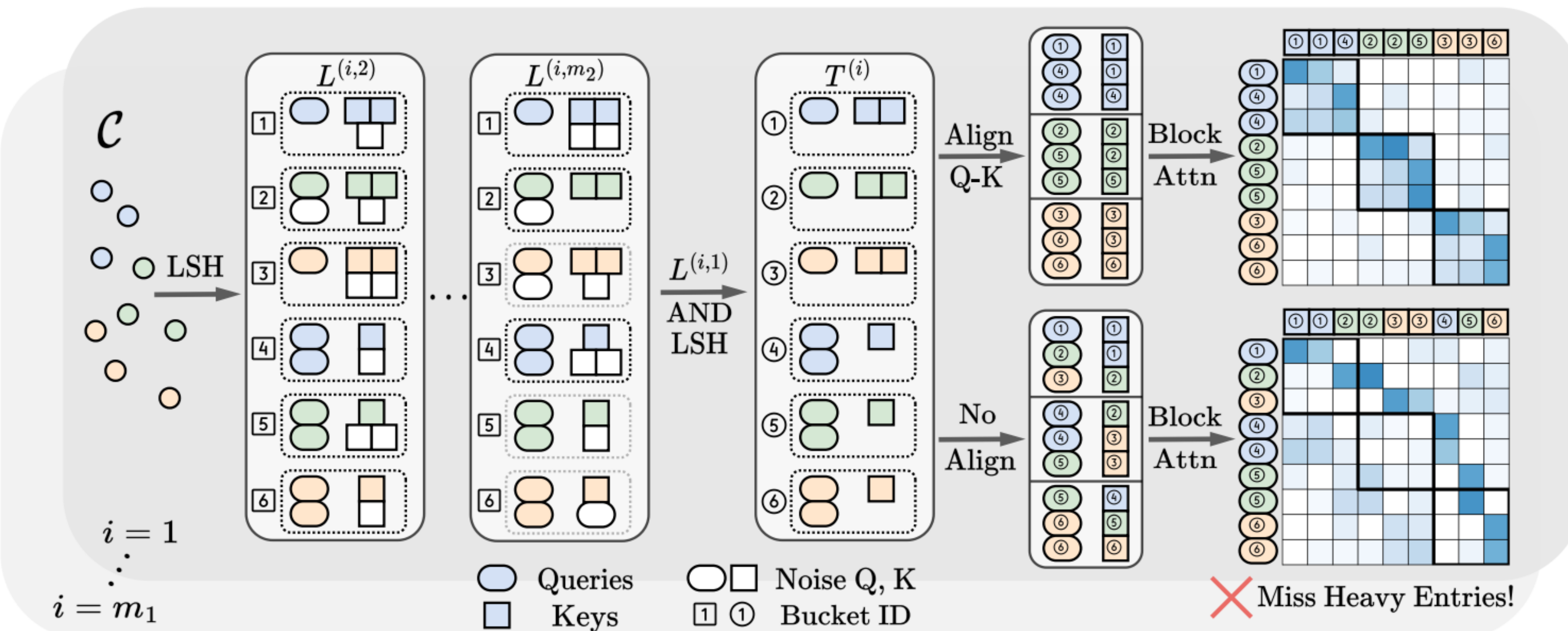
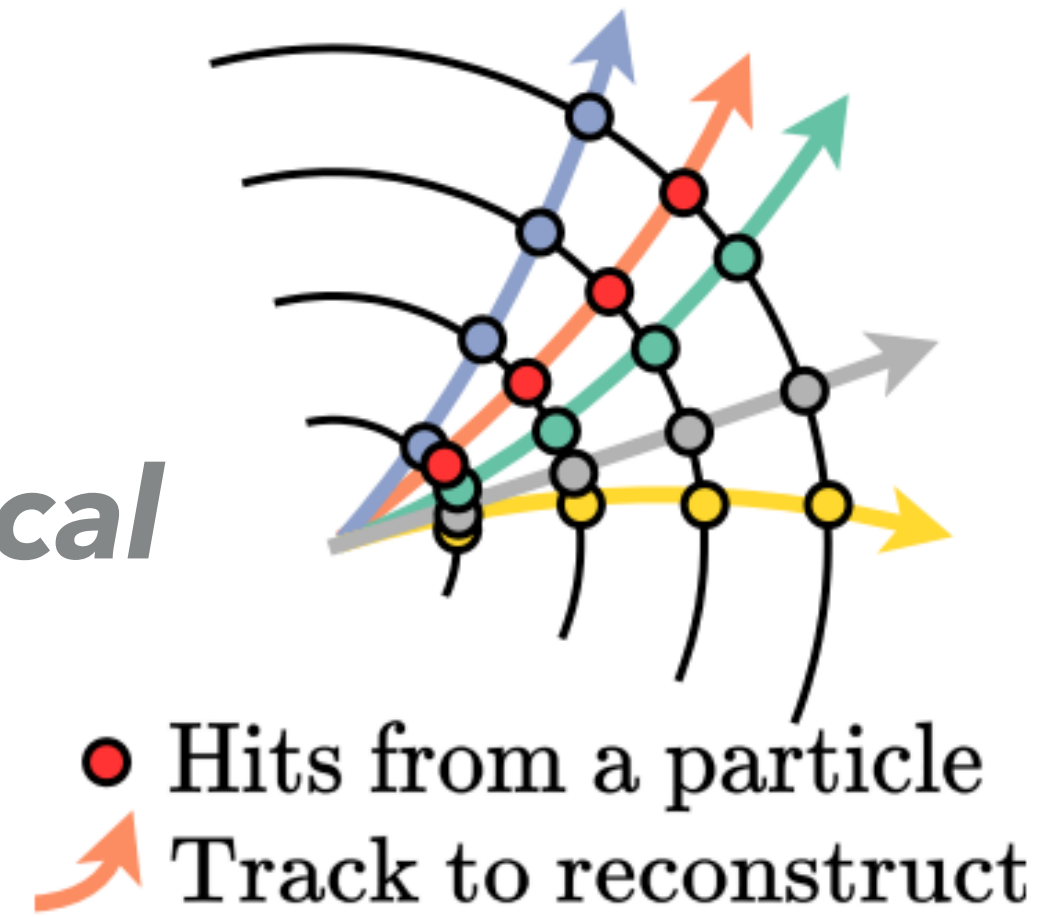


- ▶ More realistic public TrackML dataset used for 2018 Kaggle competition has  $O(100k)$  hits and  $O(10k)$  tracks per event
  - ▶  $O(N^2)$  attention or graph building is too slow!





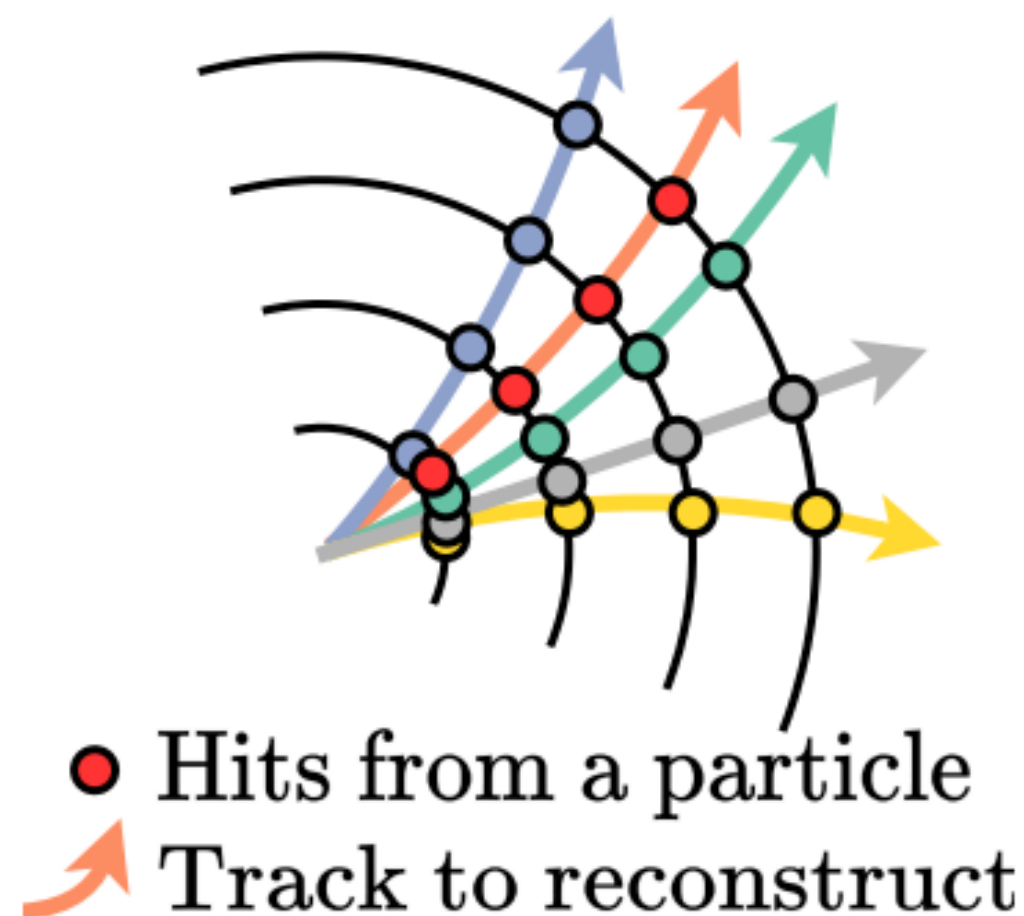
- ▶ HEPT: an efficient point transformer based on OR & AND LSH
  - ▶ No graph construction; only regular computations
  - ▶ Assign hash codes using OR & AND E<sup>2</sup>LSH; **physics-aware local inductive bias**: nearby hits in detector share 1D hash codes
  - ▶ Sort items based on hash codes; block-diagonal attention



Key idea: HEPT projects point clouds to 1D sequences using **physics-aware local inductive bias**

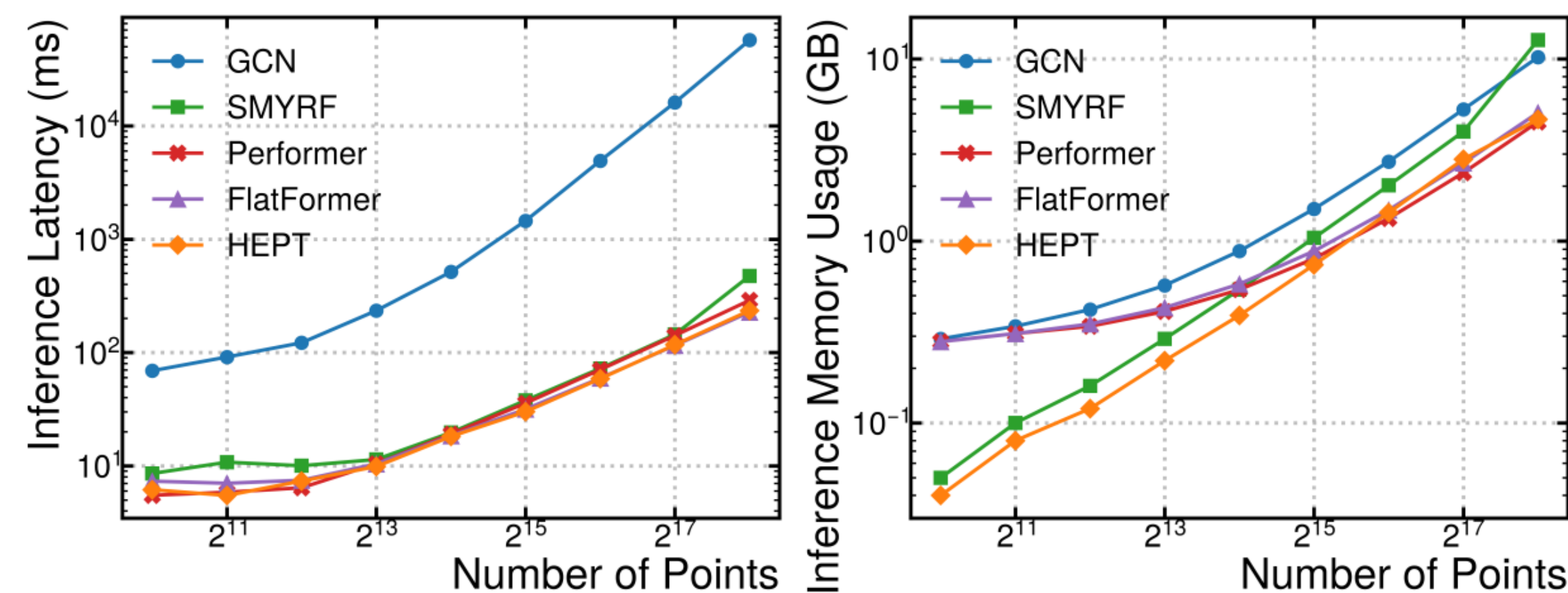


- ▶ Tracking as a representation learning task: learn close embeddings for hits originating from the same particle



	Tracking-6k (AP@k)	Tracking-60k (AP@k)
Random	5.88	5.71
SOTA GNNs	<b>91.00<sup>‡</sup></b>	<b>90.89<sup>‡</sup></b>
Reformer	72.37	<u>72.47</u>
SMYRF	72.98	71.18
HyperAttn	71.49	70.22
Performer	73.17	72.07
FLT	72.55	71.45
ScatterBrain	73.35	72.06
PointTrans	72.33	70.81
FlatFormer	<u>74.22</u>	70.23
GCN	79.61	75.38
DGCNN	<b>90.74</b>	<b>88.66</b>
GravNet	90.11	87.99
GatedGNN	80.98	78.42
Performer- $k_{\text{HEPT}}$	71.97	69.20
SMYRF- $k_{\text{HEPT}}$	83.19	71.04
FlatFormer- $k_{\text{HEPT}}$	88.18	85.06
HEPT	<b>92.66<sup>†</sup></b>	<b>91.93<sup>†</sup></b>

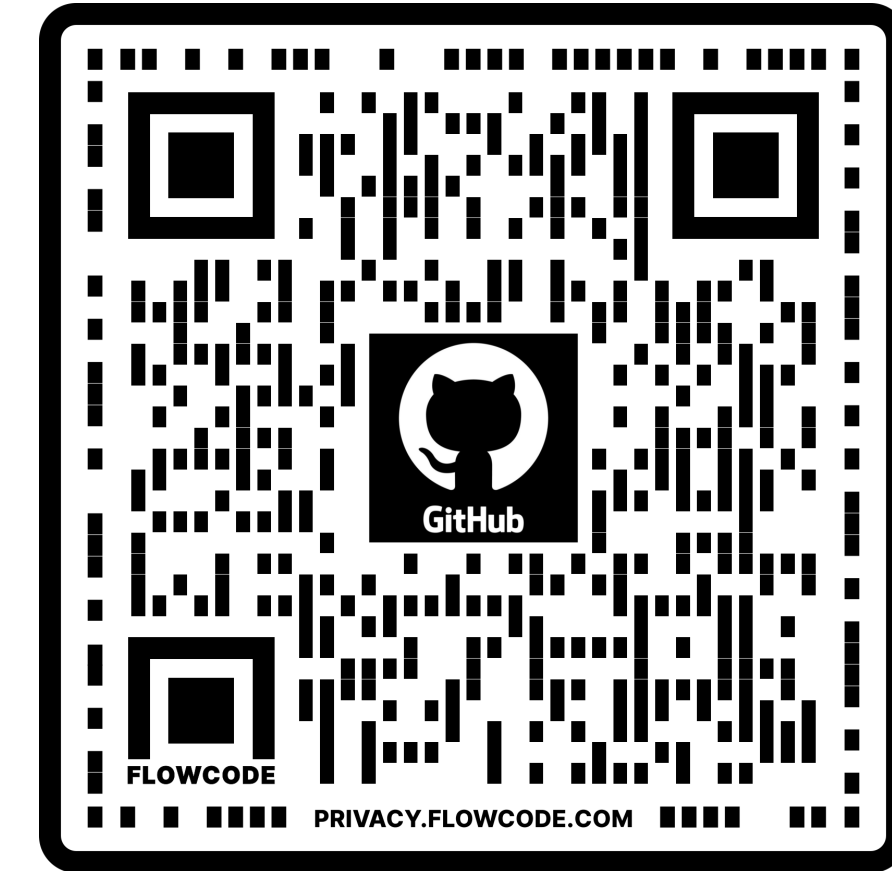
- ▶ HEPT achieves SOTA performance and achieves over 100x speedup on GPUs compared to GNNs on Tracking-60k (60k hits/event)



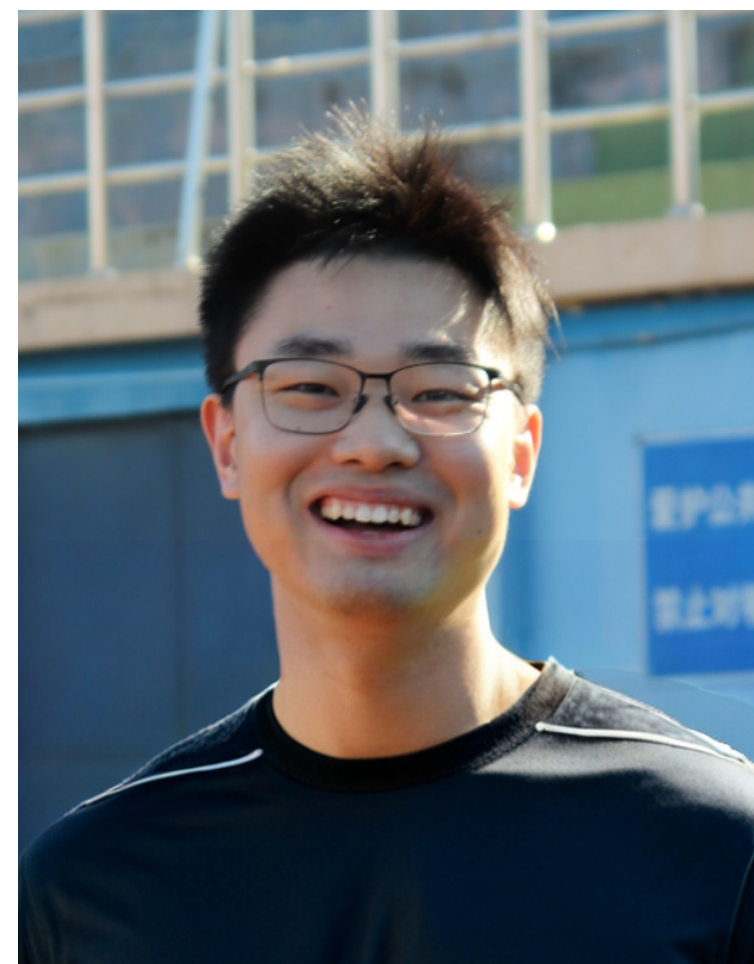


# HEAR MORE TOMORROW!

- ▶ **Oral Session:** Thursday 10:45, Lehar 1-4
- ▶ **Poster Session:** Thursday 11:30, Hall C 4-9 #407
- ▶ **Paper:** [arXiv:2402.12535](https://arxiv.org/abs/2402.12535)
- ▶ **GitHub:** <https://github.com/Graph-COM/HEPT>



Siqi Miao<sup>1</sup>



Zhiyuan Lu<sup>2</sup>



Mia Liu<sup>3</sup>



Javier Duarte<sup>4</sup>

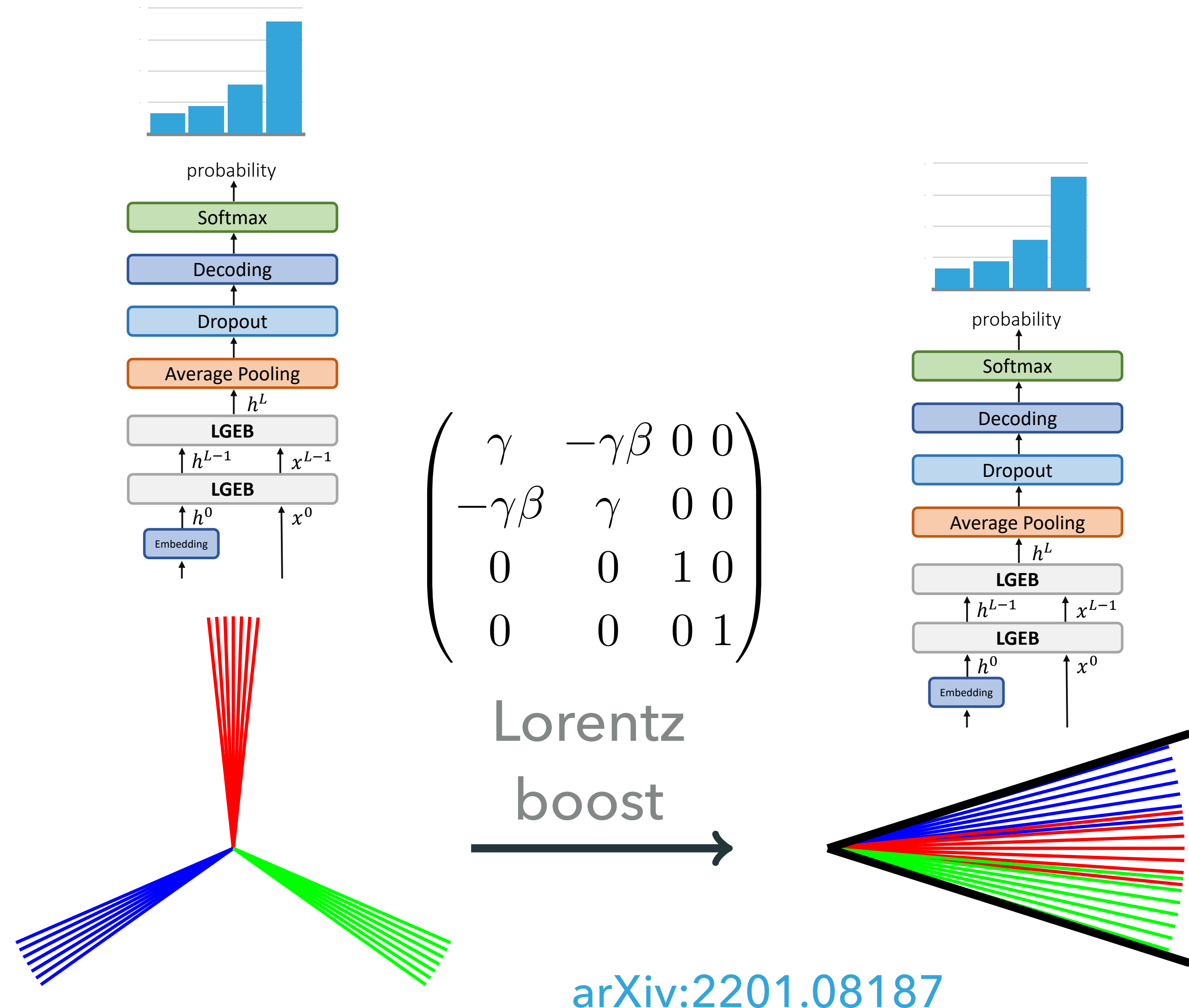


Pan Li<sup>1</sup>





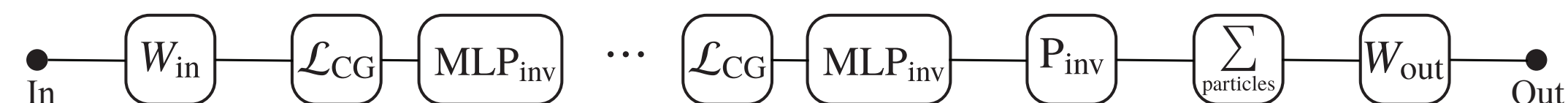
- ▶ Lorentz symmetry: physics is the same no matter which reference frame we consider
- ▶ Lorentz-invariant networks:
  - ▶ Boosting all particles into a new reference frame should give the same result
- ▶ Lorentz-equivariant networks:
  - ▶ Boosting all particles into a new reference frame should give an output that transforms the same way





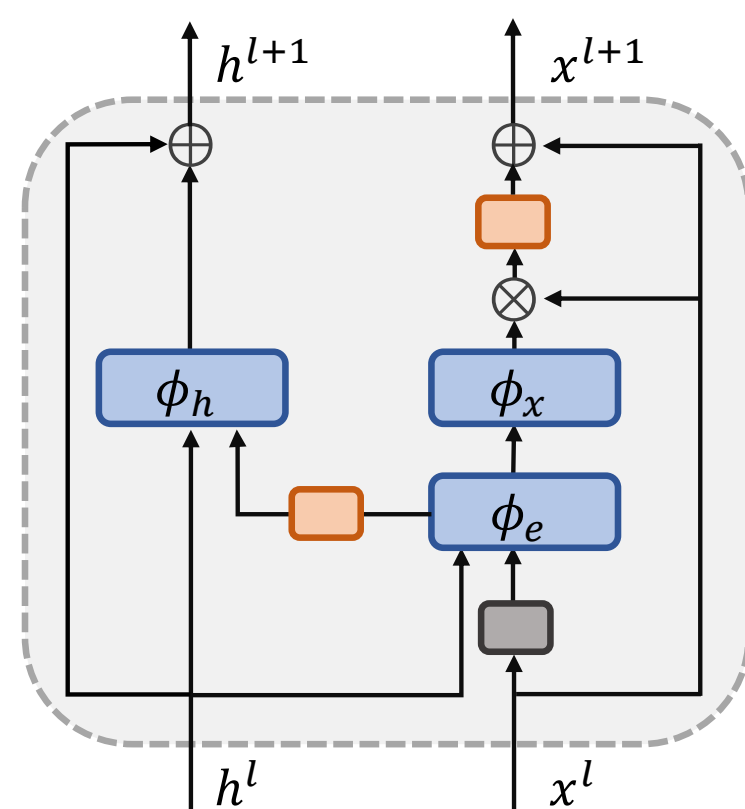
- ▶ Lorentz Group Network [[arXiv:2006.04780](https://arxiv.org/abs/2006.04780)]:

nonlinearity is tensor product followed by Clebsch-Gordan (CG) decomposition



- ▶ LorentzNet [[arXiv:2201.08187](https://arxiv.org/abs/2201.08187)] uses structured message passing based on Lorentz scalars and vectors

- ▶ PELICAN [[arXiv:2307.16506](https://arxiv.org/abs/2307.16506)] builds invariants and covariants based on pairs of inputs



MLP
  Sum Pooling
  Minkowski Norm & Inner Product

**Lorentz Group Equivariant Block (LGEb)**

**Invariants** (output for PELICAN's "scalar form")

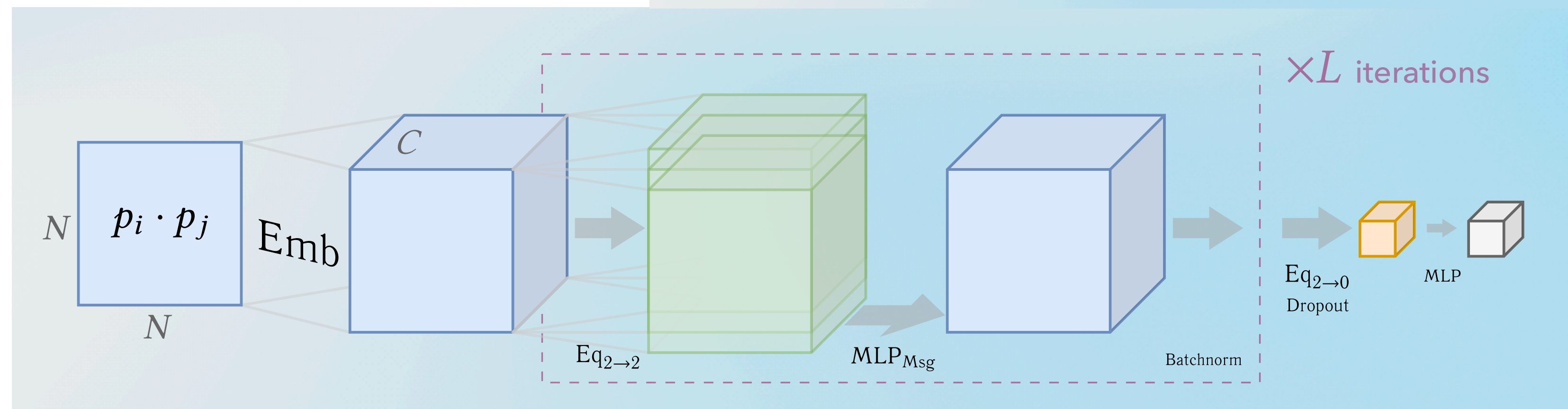
$$I_k(p_1, \dots, p_N) = I_k(\{d_{ij}\}_{i,j}), \quad d_{ij} \equiv p_i^\mu p_{j,\mu}$$

Lorentz-invariant

**Covariants** (output for PELICAN's "vector form")

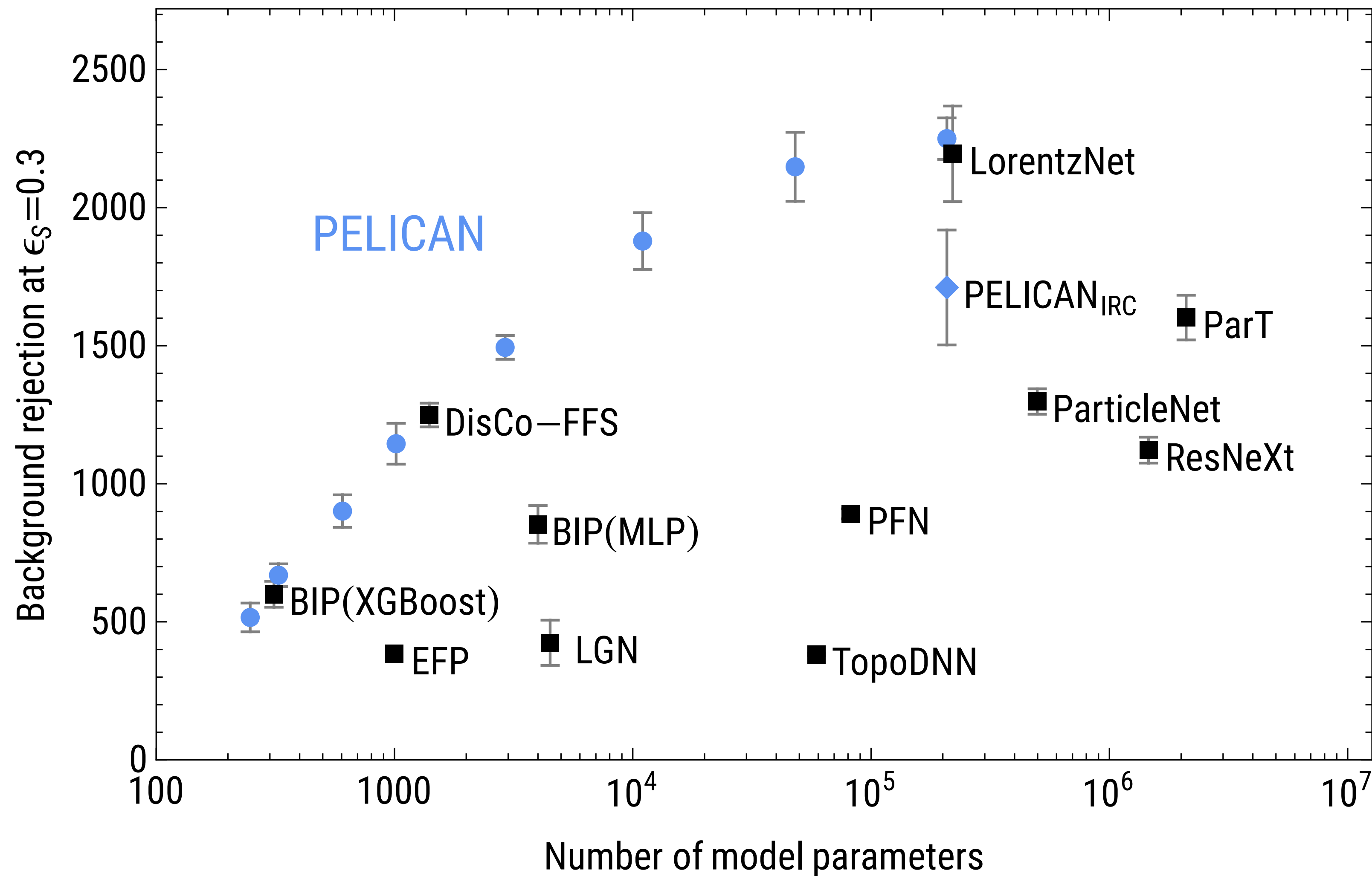
$$F^\mu = \sum_k I_k(p_1, \dots, p_N) p_k^\mu$$

Lorentz-covariant



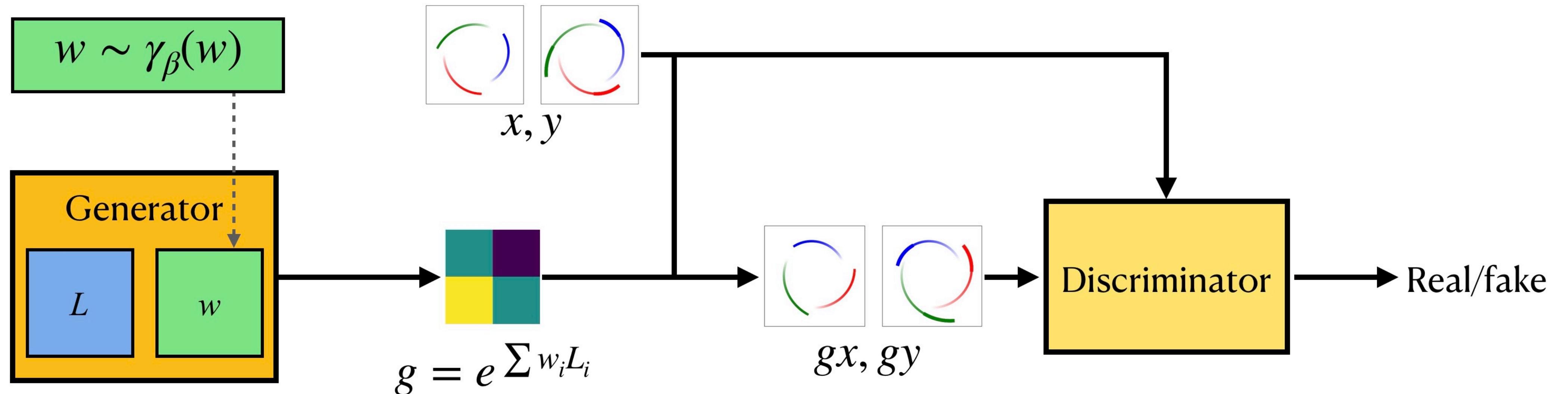
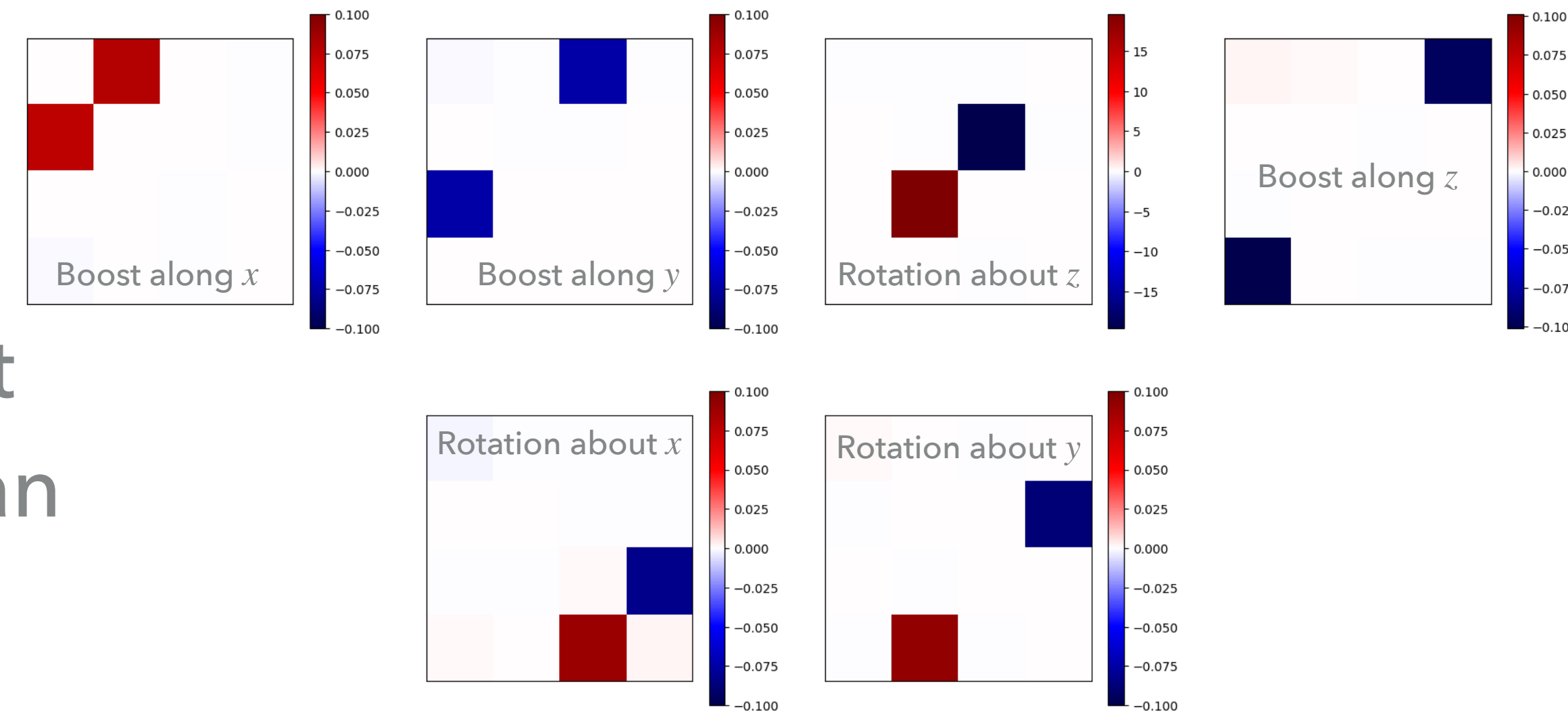


- ▶ PELICAN traces Pareto optimal boundary of performance and model complexity





- ▶ Symmetries are fundamental principles in particle physics
- ▶ LieGAN learns a continuous Lie group that preserves the original data distribution; can discover symmetries present in the data!
- ▶ Discovers approximate  $SO(1, 3)^+$  symmetry in particle physics dataset





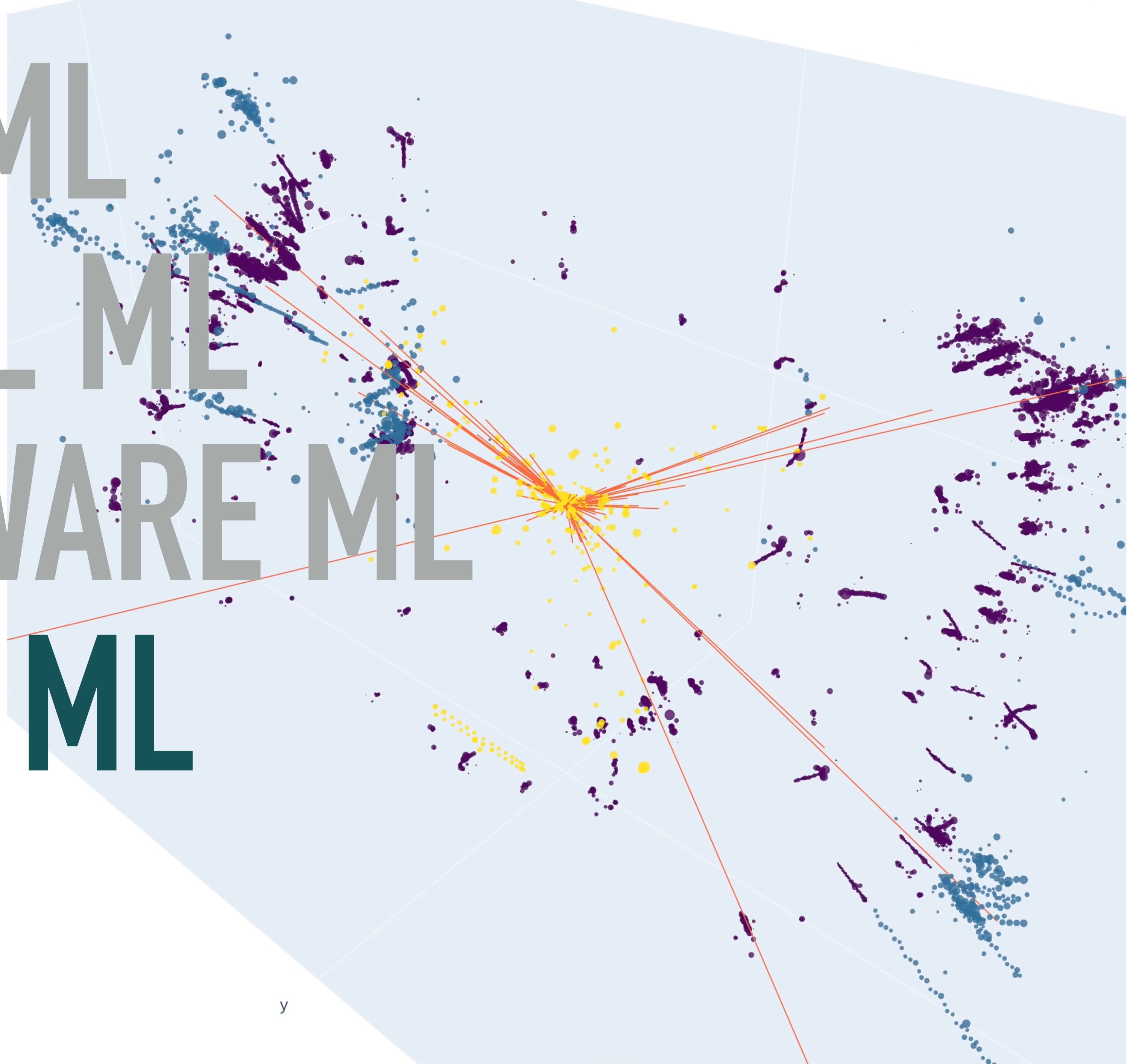
ULTRAFAST ML

MULTIMODAL ML

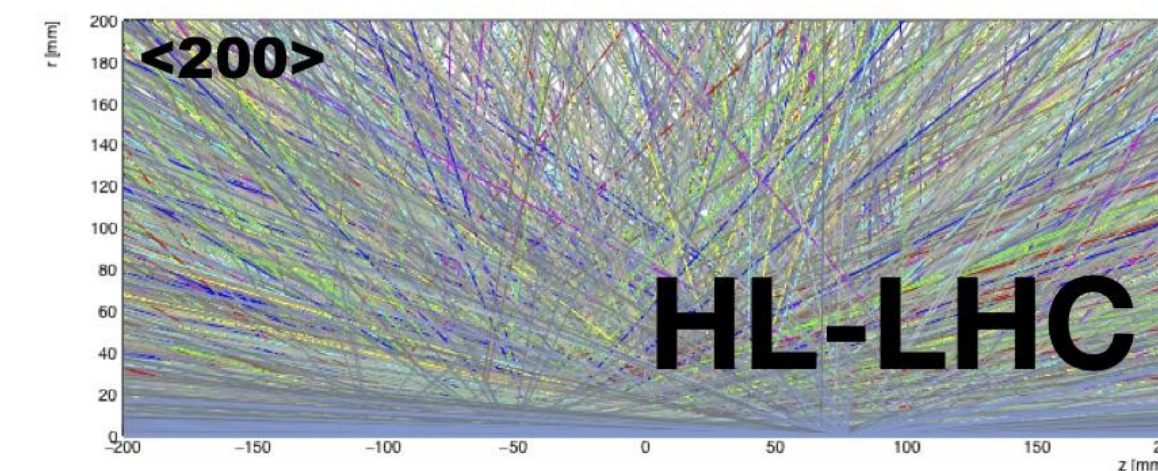
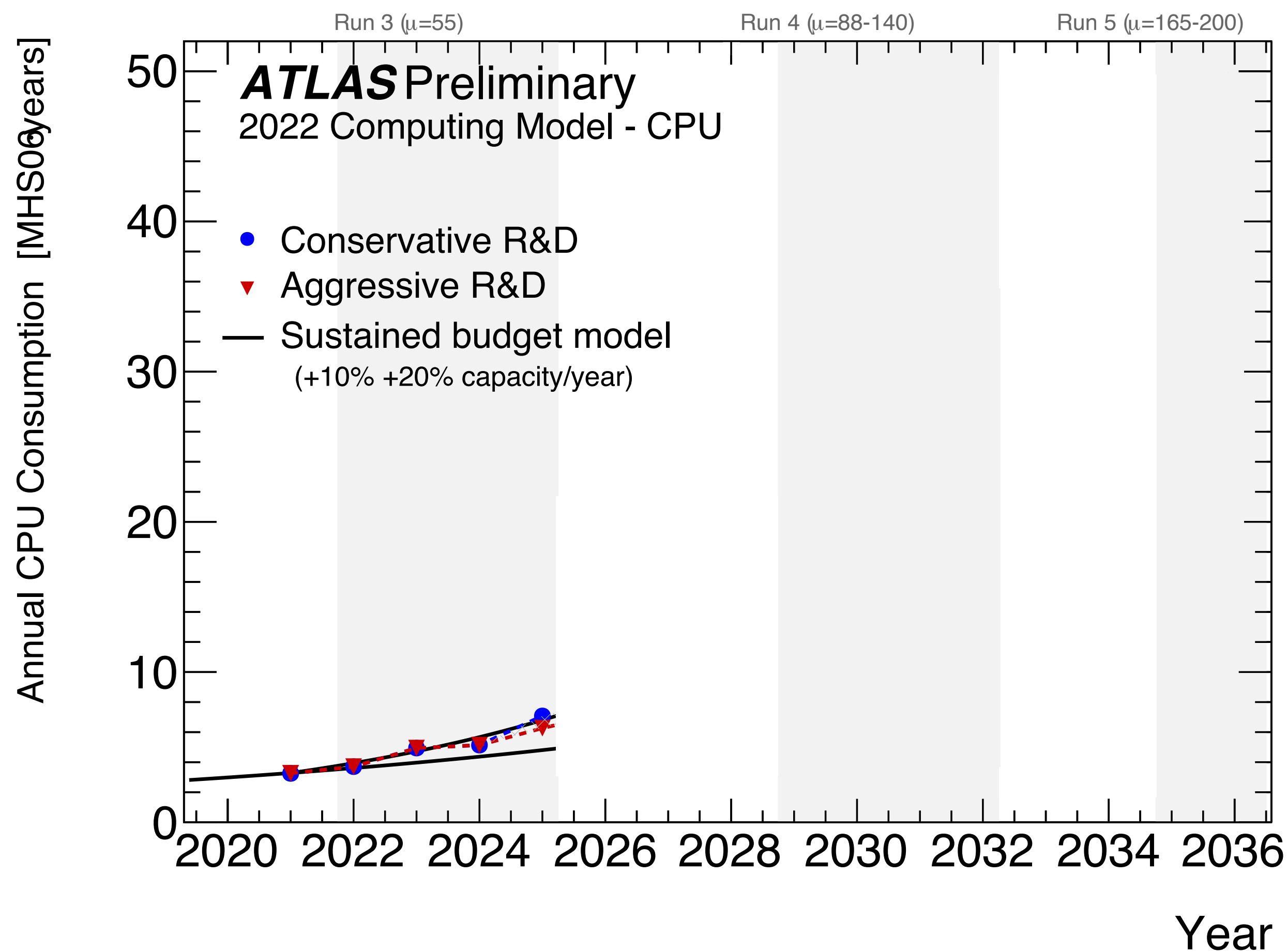
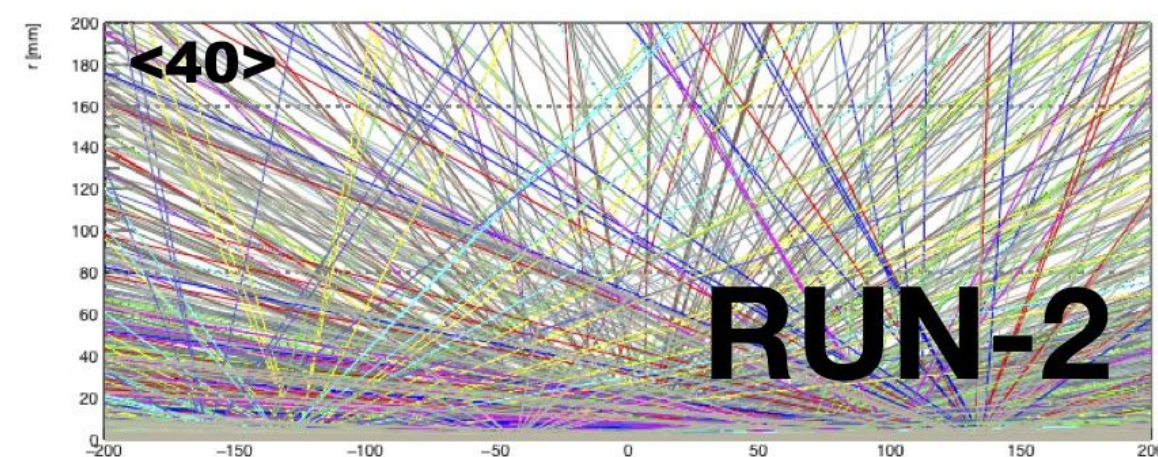
PHYSICS-AWARE ML

**GENERATIVE ML**

OUTLOOK



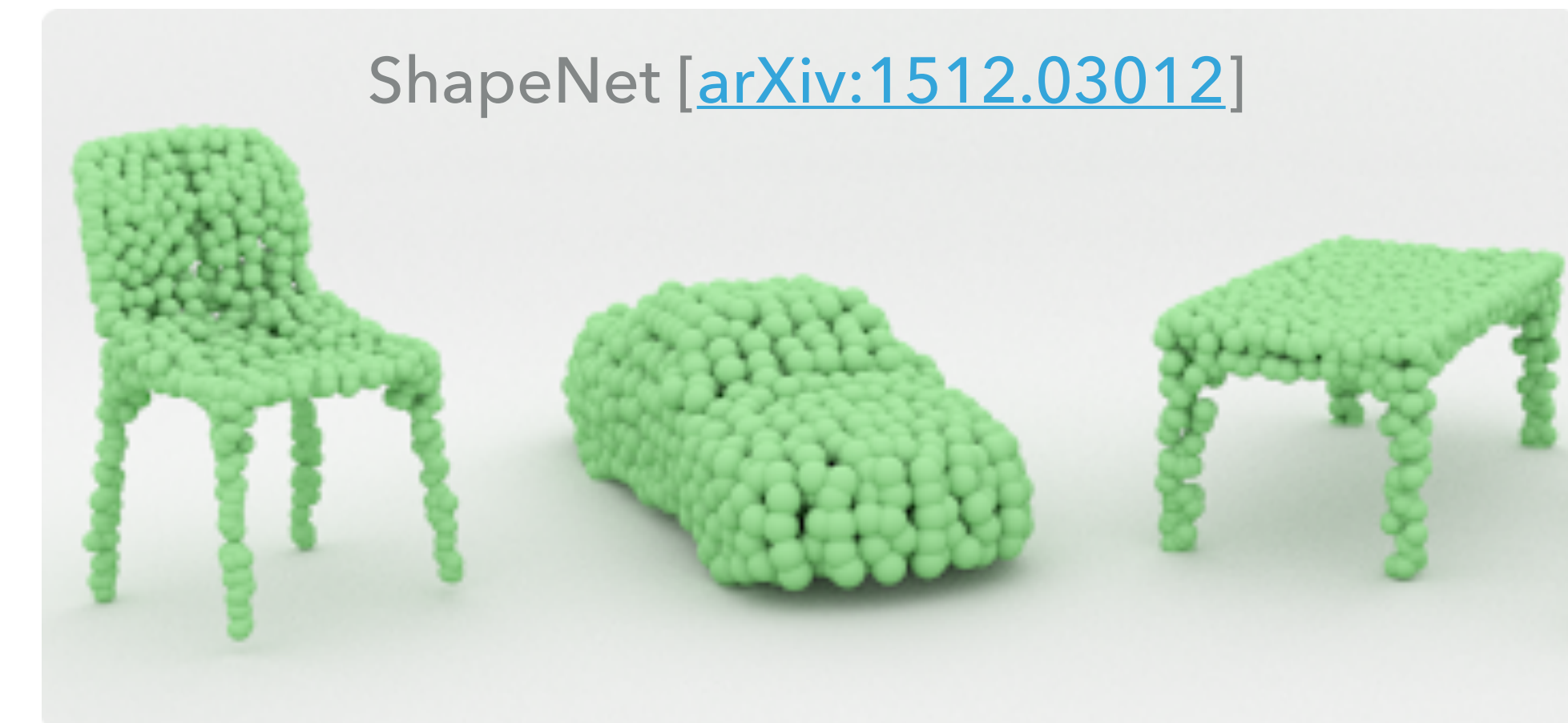
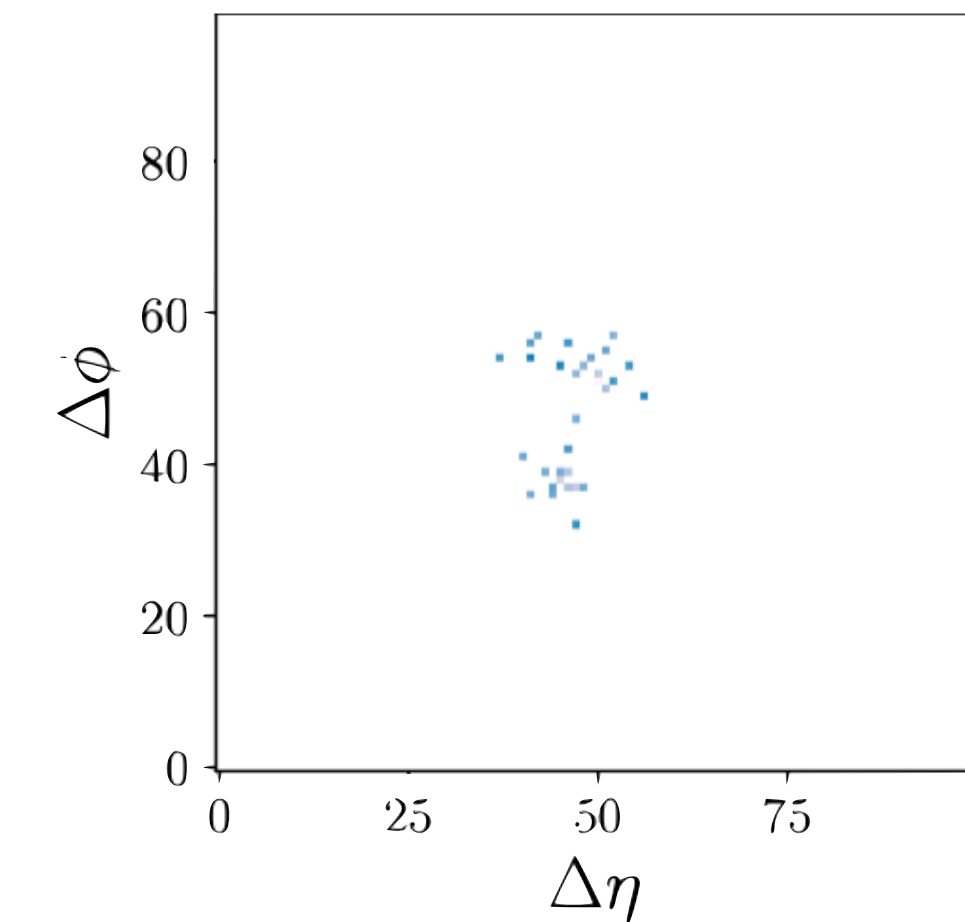




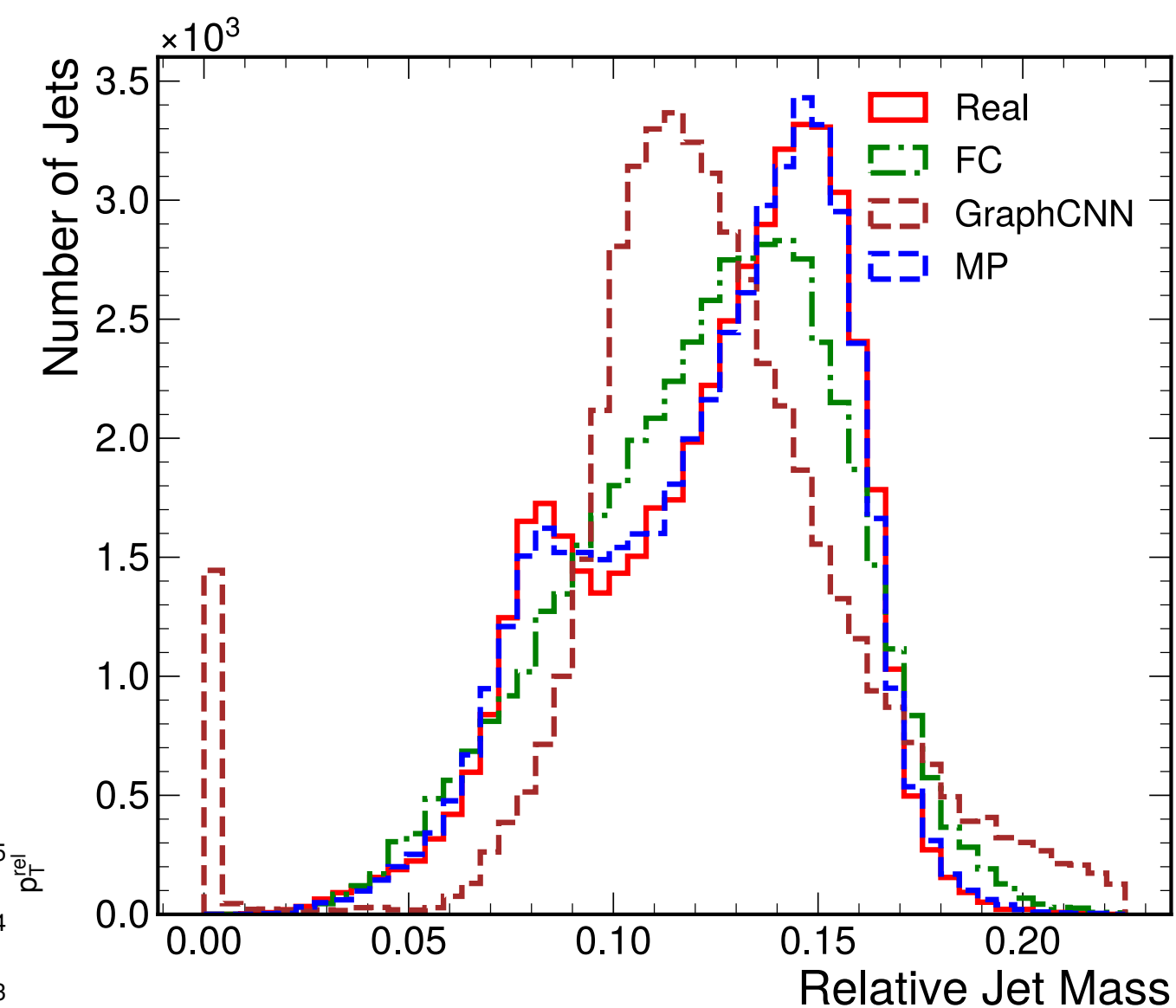
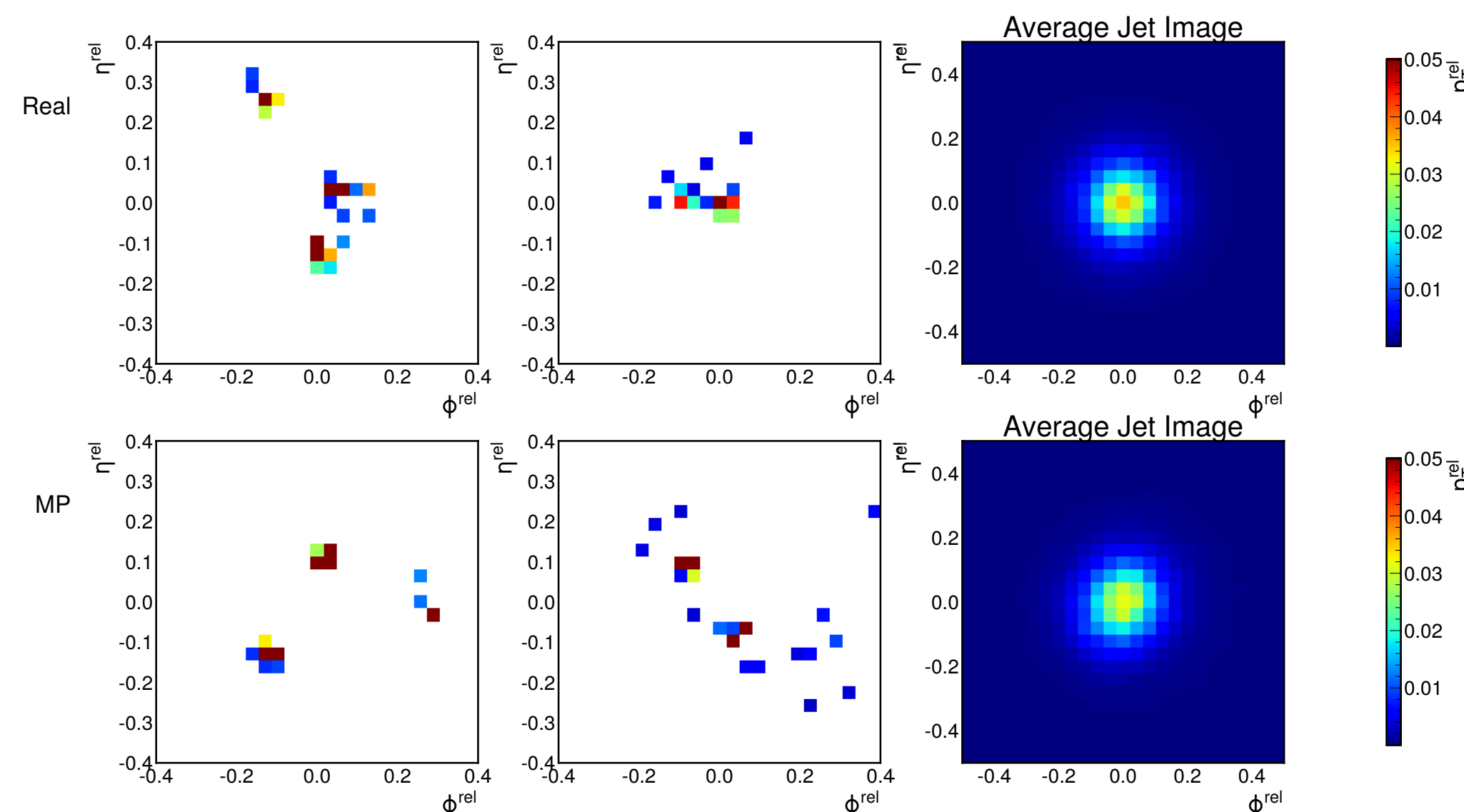
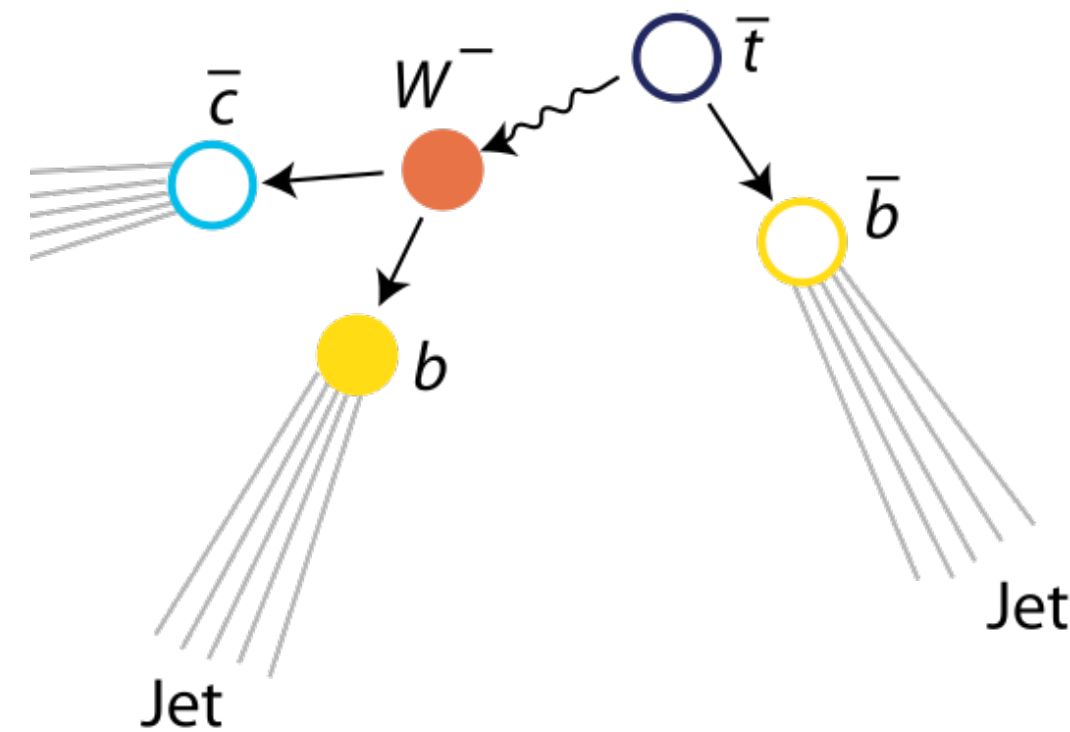
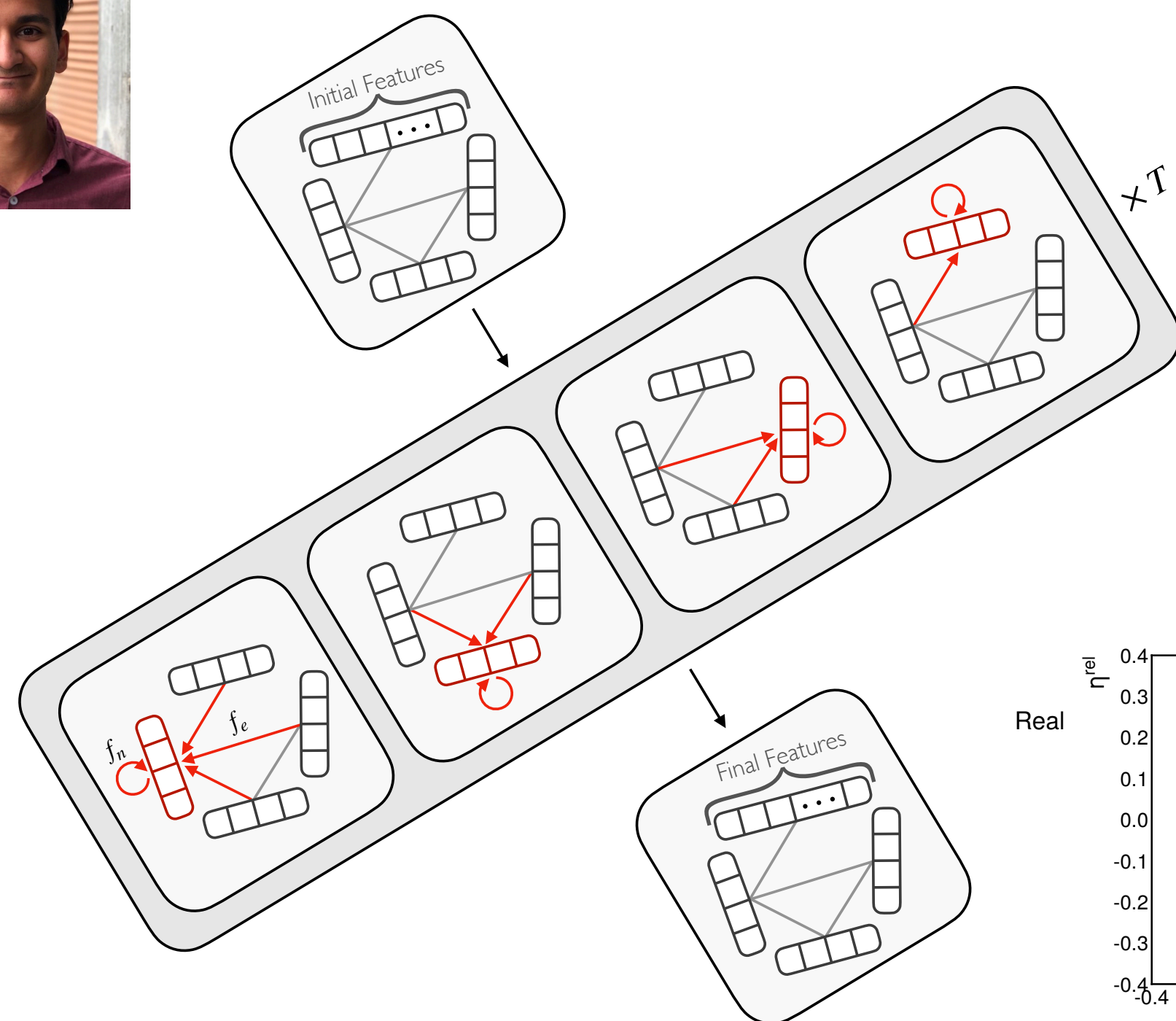
- ▶ Simulation is a key driver of CPU needs for the HL-LHC
- ▶ ML can be used to “short cut” simulation



- ▶ Public dataset for benchmarking generative models focusing on “particle cloud” representations instead of image representations
  - ▶ Similar idea to ShapeNet [[arXiv:1512.03012](https://arxiv.org/abs/1512.03012)]
  - ▶ Consists of up to 150 particles per jet with 3 features:  $(p_T^{\text{rel}}, \eta^{\text{rel}}, \phi^{\text{rel}})$
  - ▶ Available on Zenodo [[doi:10.5281/zenodo.4834875](https://doi.org/10.5281/zenodo.4834875)]
- ▶ Part of **FAIR4HEP** [[fair4hep.github.io](https://fair4hep.github.io)]





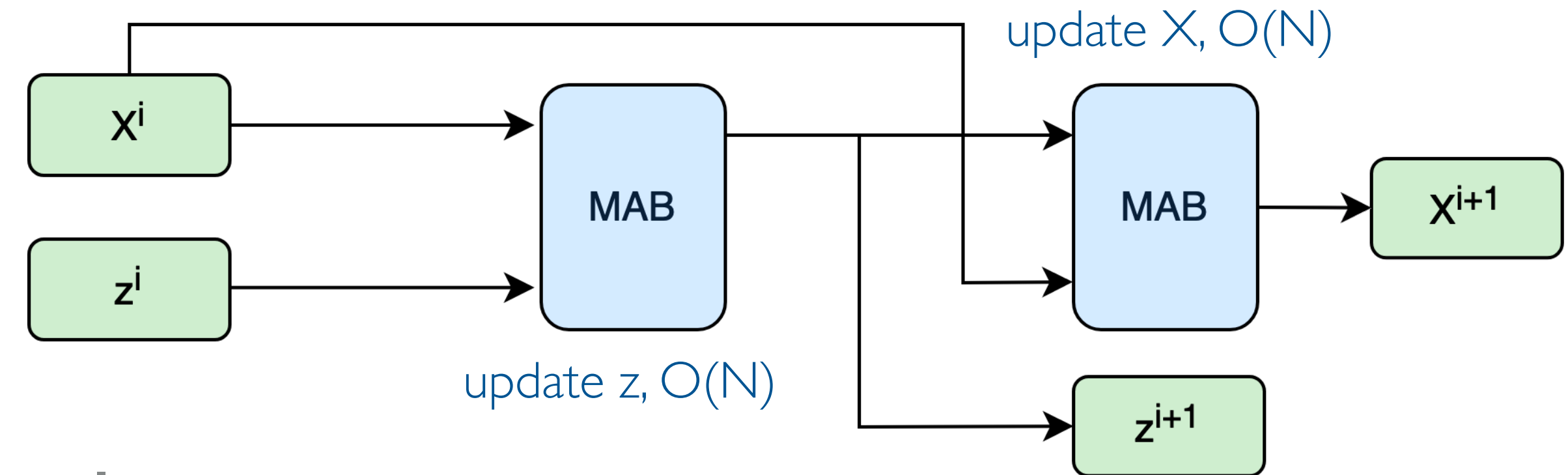


- ▶ Message-passing GAN can generate realistic particle jets
- ▶ Outperforms existing point cloud GANs



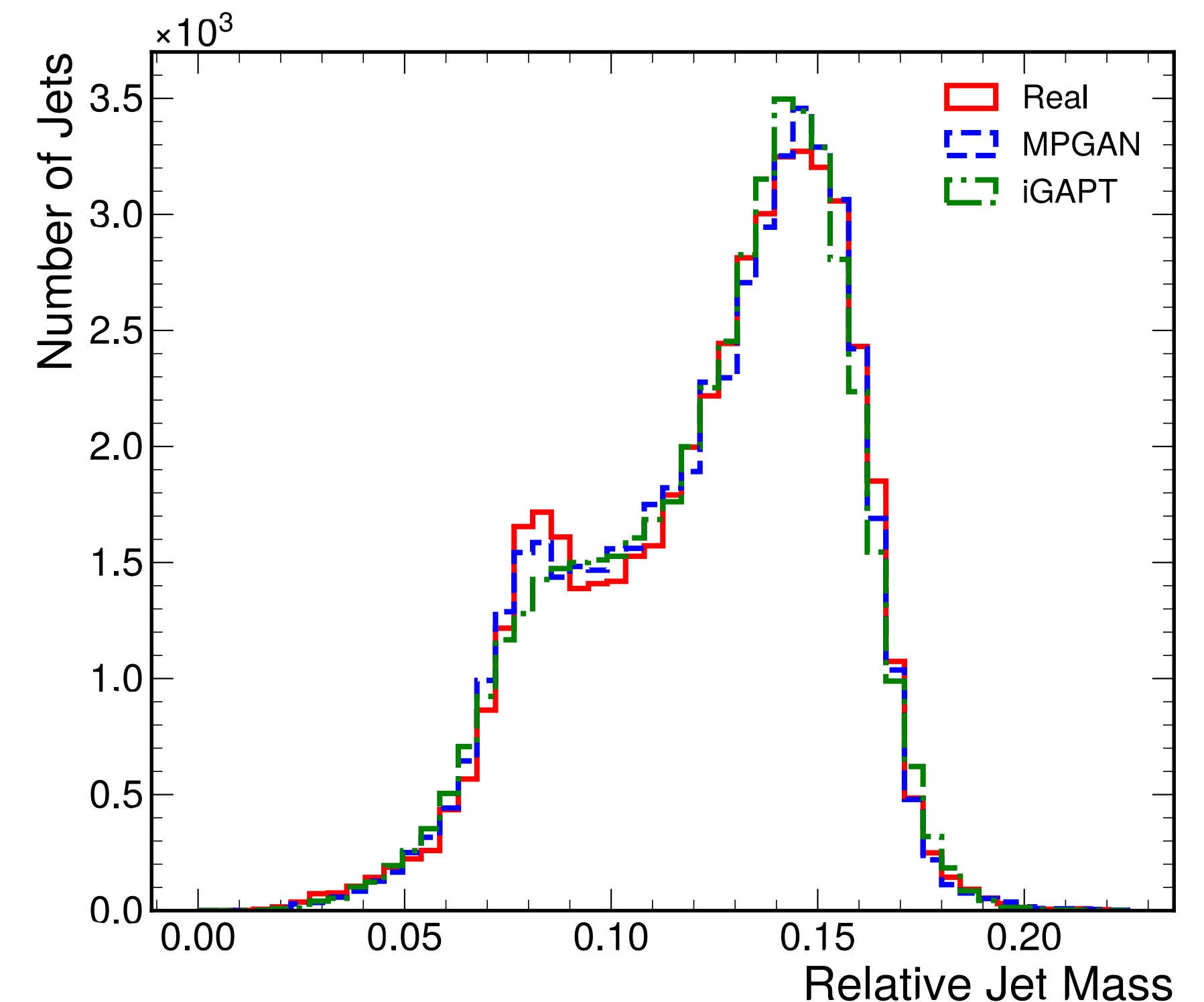
- ▶ Induced graph attention particle transformer (iGAPT) more efficient than MPGAN, but how to compare them?

$N \times 3$  particles  
jet features



- ▶ Metrics developed to comprehensively evaluate generative AI models in HEP
- ▶ **Fréchet Physics Distance:** Inspired by Fréchet Inception Distance with physics-based features

	FPD $\times 10^3$	Inference time ( $\mu s$ ) per jet
Truth	$0.08 \pm 0.03$	
MPGAN	<b><math>0.30 \pm 0.06</math></b>	41
GAPT	$0.66 \pm 0.09$	9





# IMPACT OF JETNET

► Sparked significant R&D into equivariant models, **diffusion models**, and more!

## CaloClouds: Fast Geometry-Independent Highly-Granular Calorimeter Simulation

Erik Buhmann<sup>1</sup>, Sascha Diefenbacher<sup>1,2</sup>, Engin Eren<sup>3</sup>, Frank Gaede<sup>3,4</sup>, Gregor Kasieczka<sup>1,4</sup>, Anatolii Korol<sup>3\*</sup>, William Korcari<sup>1</sup>, Katja Krüger<sup>3</sup>, Peter McKeown<sup>3</sup>

## Pay Attention to Mean-Fields during Particle Cloud Generation

**Benno Käch**  
Deutsches Elektronen-Synchrotron DESY  
Germany  
benno.kaech@desy.de

**Isabell-A. Melzer-Pellmann**  
Deutsches Elektronen-Synchrotron DESY  
Germany  
isabell.melzer@desy.de

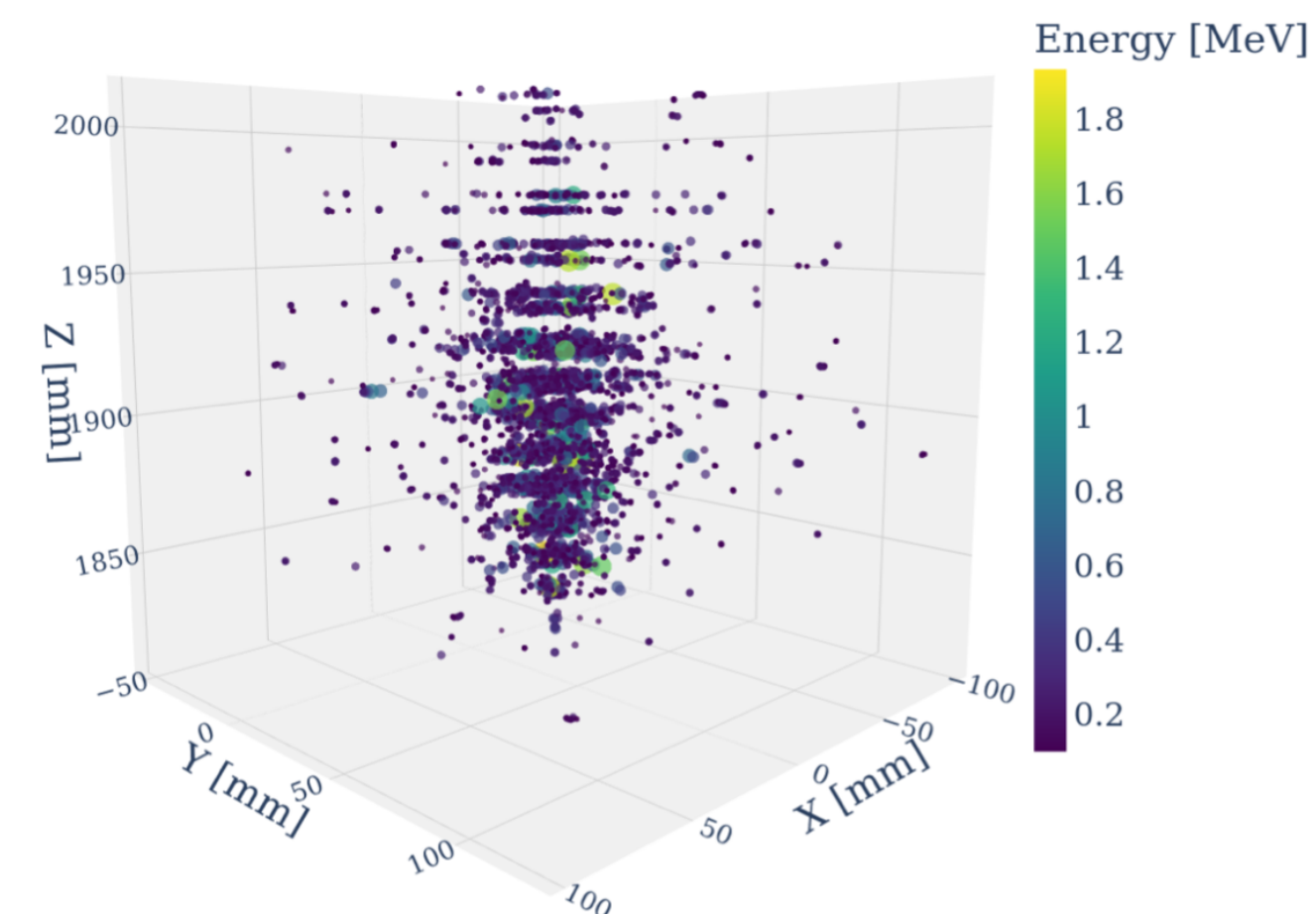
SciPost Physics

Submission

## PC-JeDi: Diffusion for Particle Cloud Generation in High Energy Physics

Matthew Leigh, Debajyoti Sengupta, Guillaume Quétant, John Andrew Raine, Knut Zoch, and Tobias Golling,

University of Geneva



SciPost Physics

Submission

## EPiC-GAN: Equivariant Point Cloud Generation for Particle Jets

Erik Buhmann<sup>1\*</sup>, Gregor Kasieczka<sup>1,2</sup> and Jesse Thaler<sup>3,4</sup>

<sup>1</sup> Institut für Experimentalphysik, Universität Hamburg, Germany

<sup>2</sup> Center for Data and Computing in Natural Sciences (CDCS), Hamburg, Germany

<sup>3</sup> Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, Massachusetts, United States

<sup>4</sup> The NSF AI Institute for Artificial Intelligence and Fundamental Interactions

\* erik.buhmann@uni-hamburg.de

MIT-CTP 5519

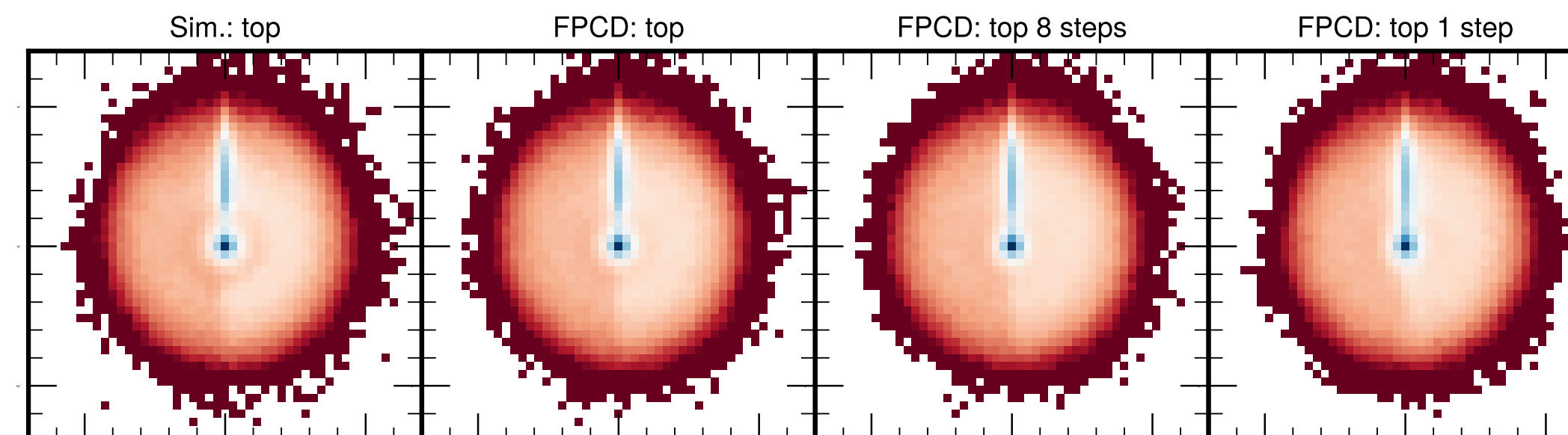
## Fast Point Cloud Generation with Diffusion Models in High Energy Physics

Vinicius Mikuni<sup>1,\*</sup>, Benjamin Nachman<sup>2,3,†</sup> and Mariel Pettee<sup>2,‡</sup>

<sup>1</sup> National Energy Research Scientific Computing Center, Berkeley Lab, Berkeley, CA 94720, USA

<sup>2</sup> Physics Division, Lawrence Berkeley National Laboratory, Berkeley, CA 94720, USA

<sup>3</sup> Berkeley Institute for Data Science, University of California, Berkeley, CA 94720, USA





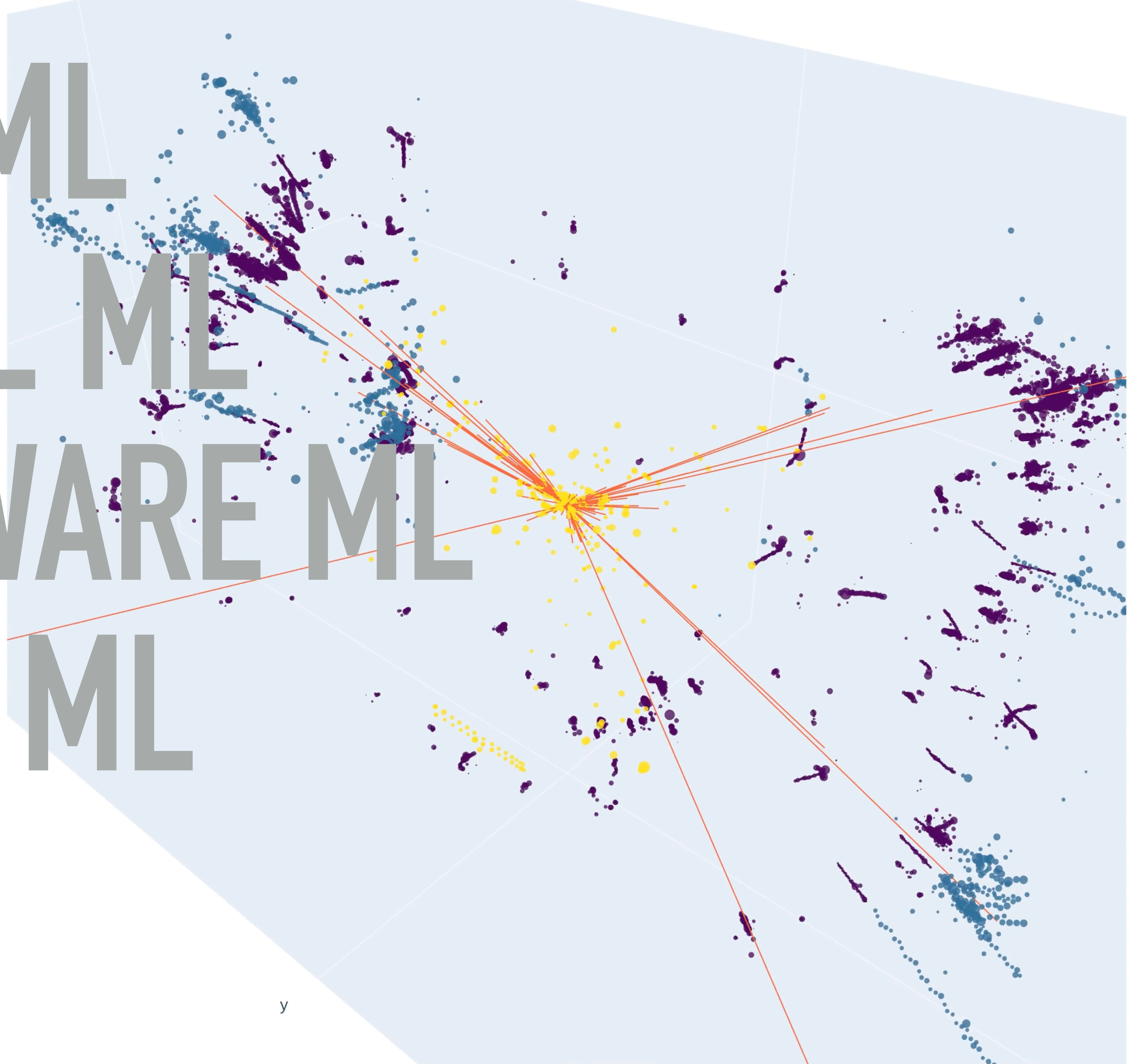
ULTRAFAST ML

MULTIMODAL ML

PHYSICS-AWARE ML

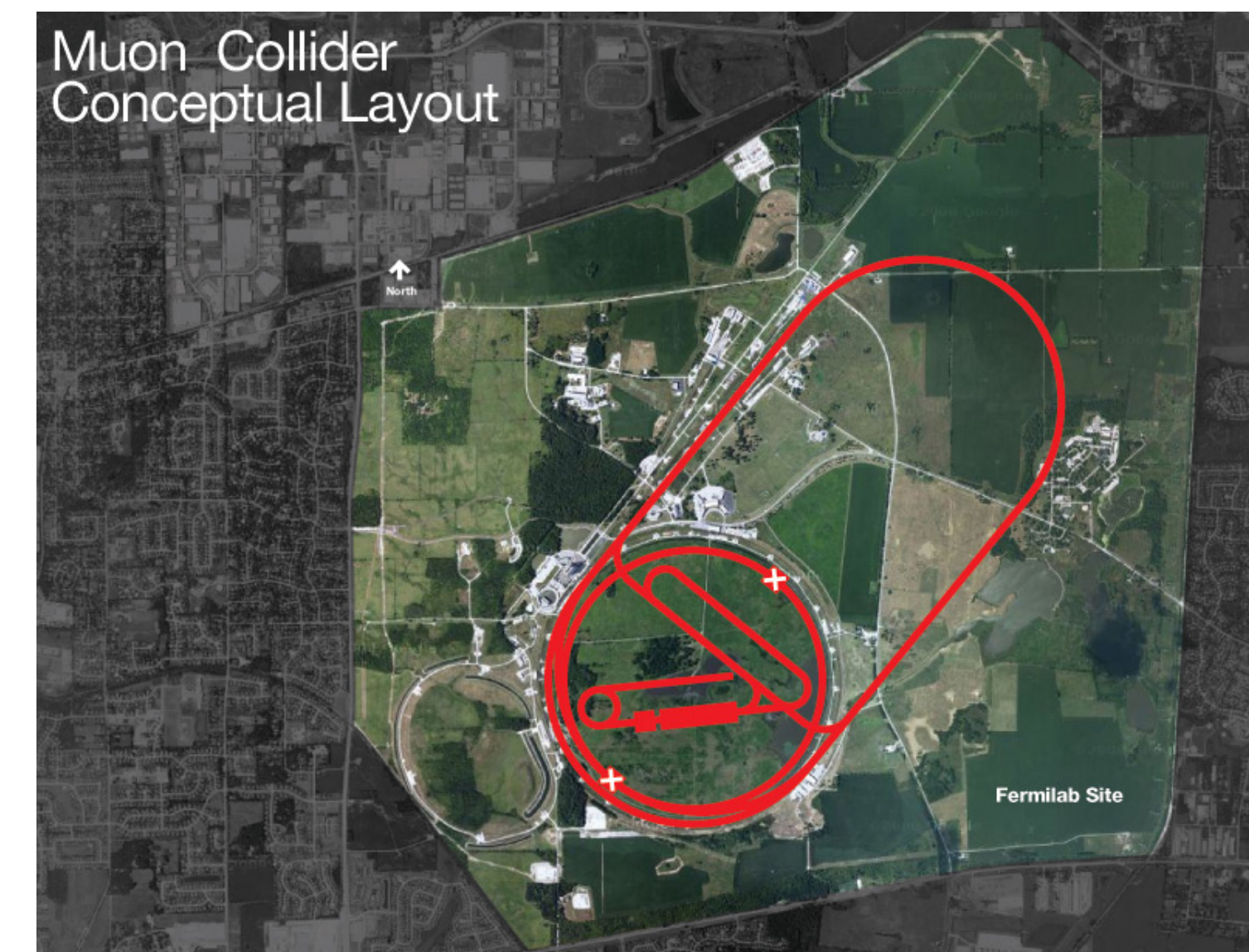
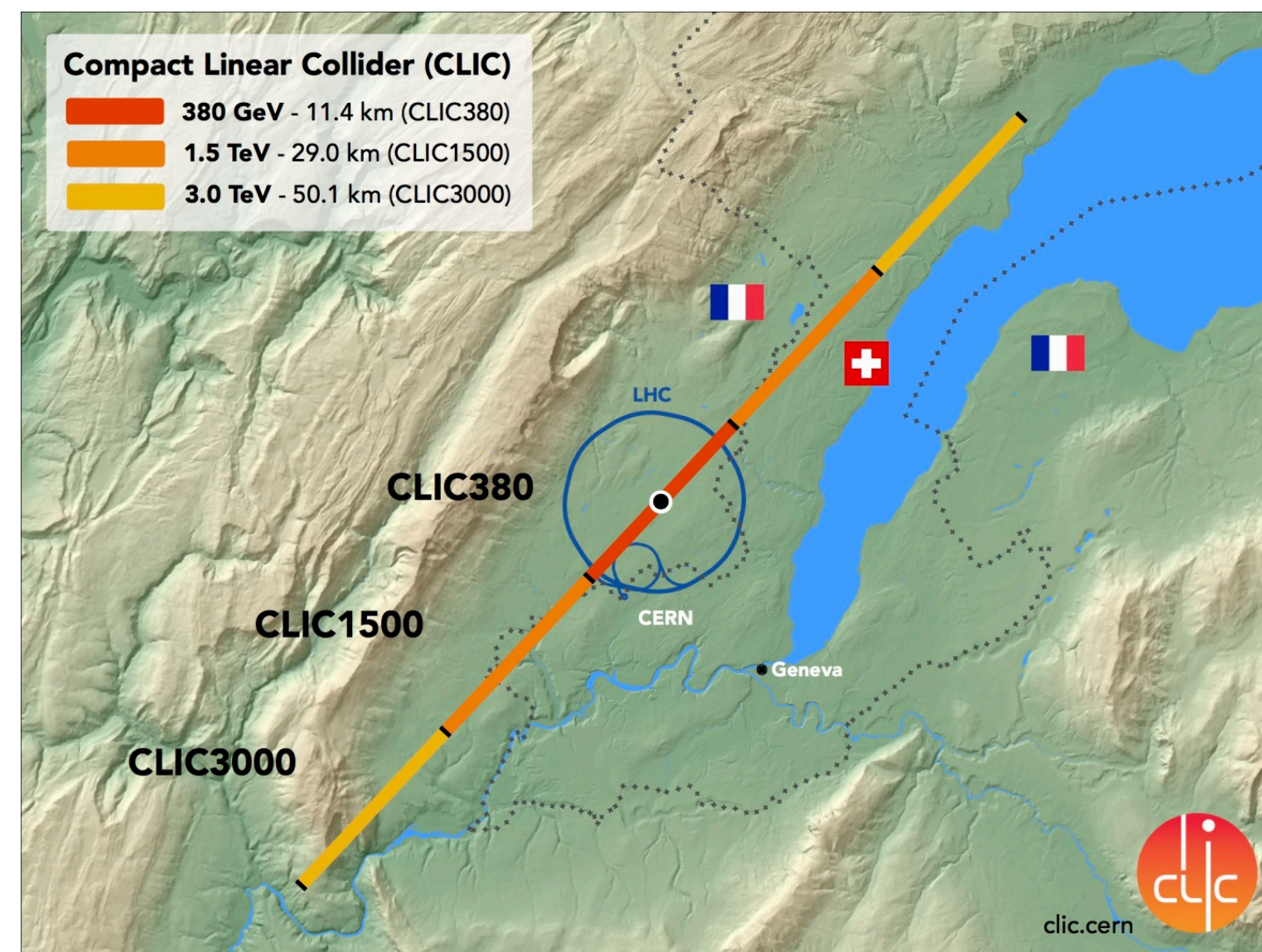
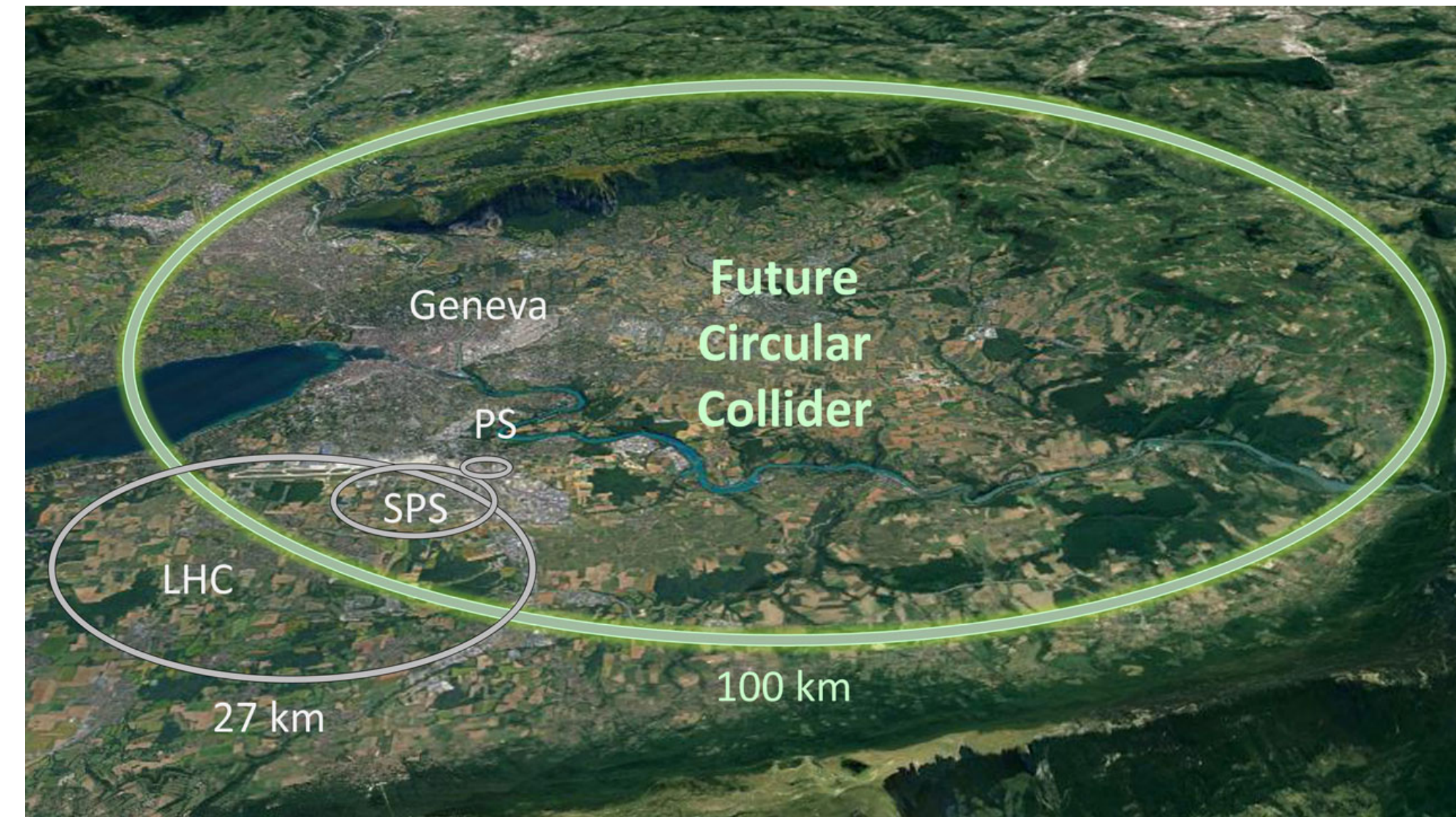
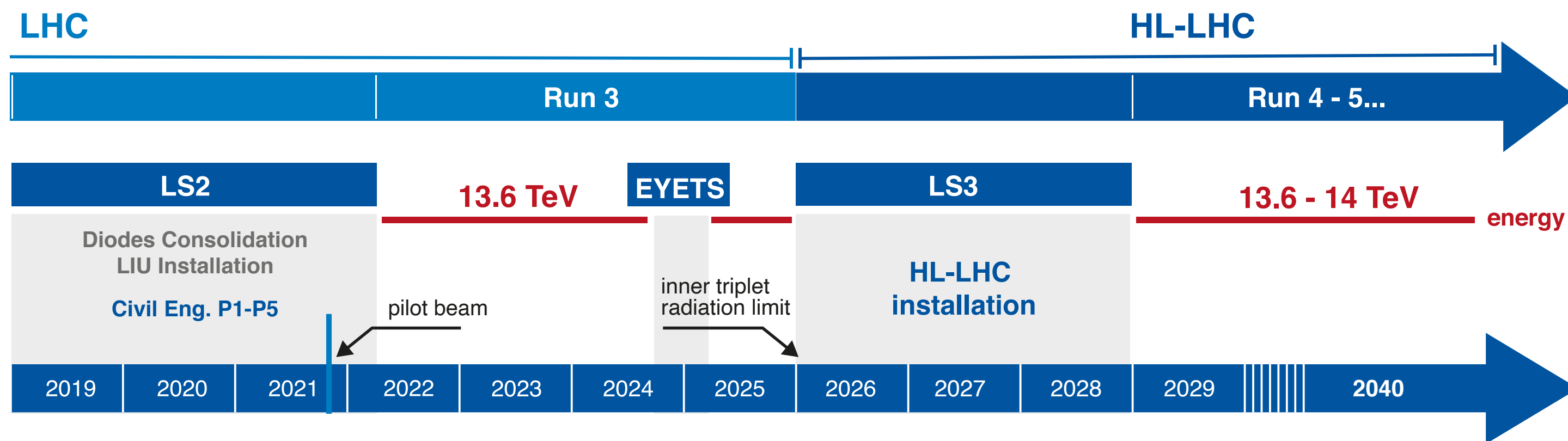
GENERATIVE ML

**OUTLOOK**



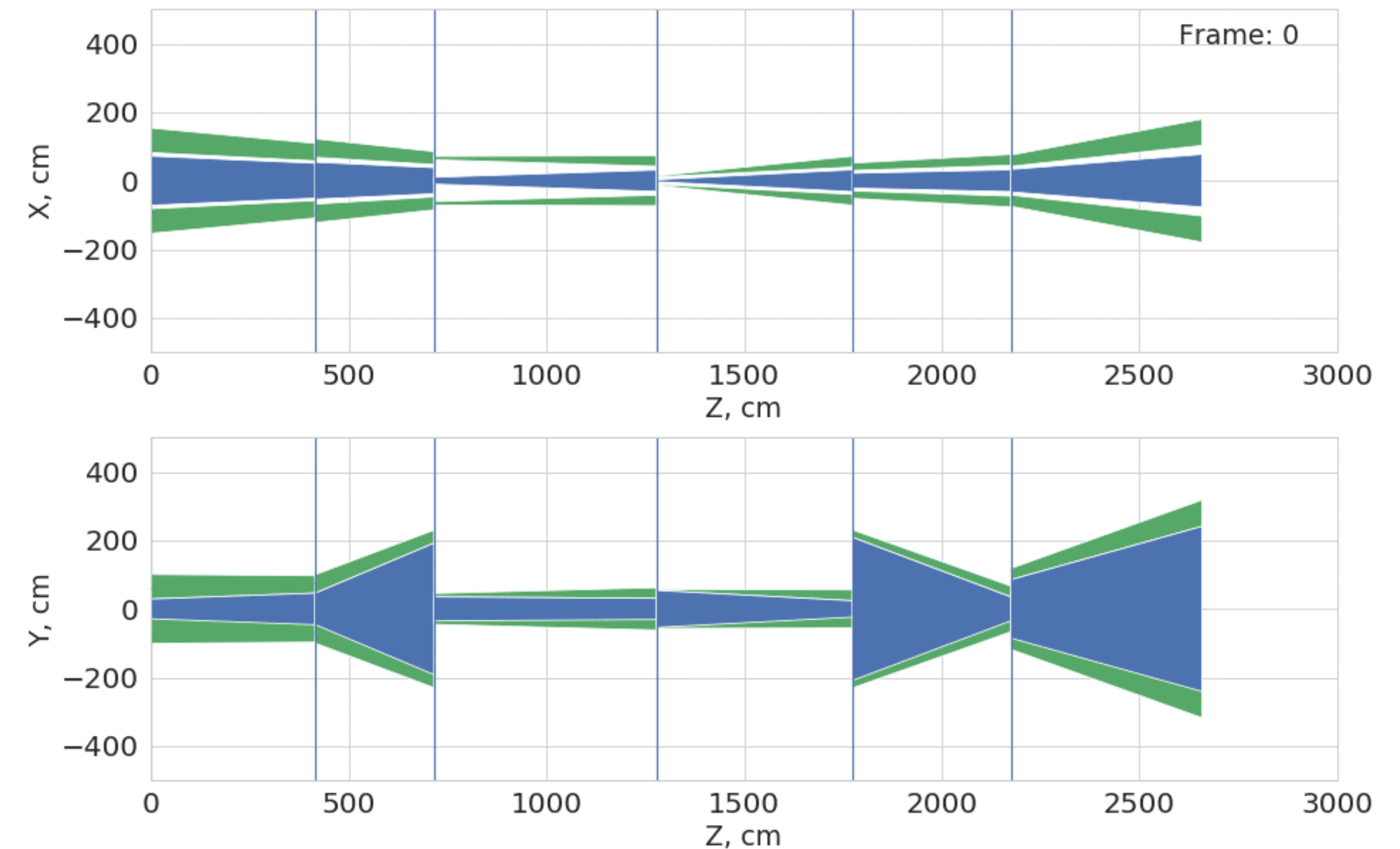
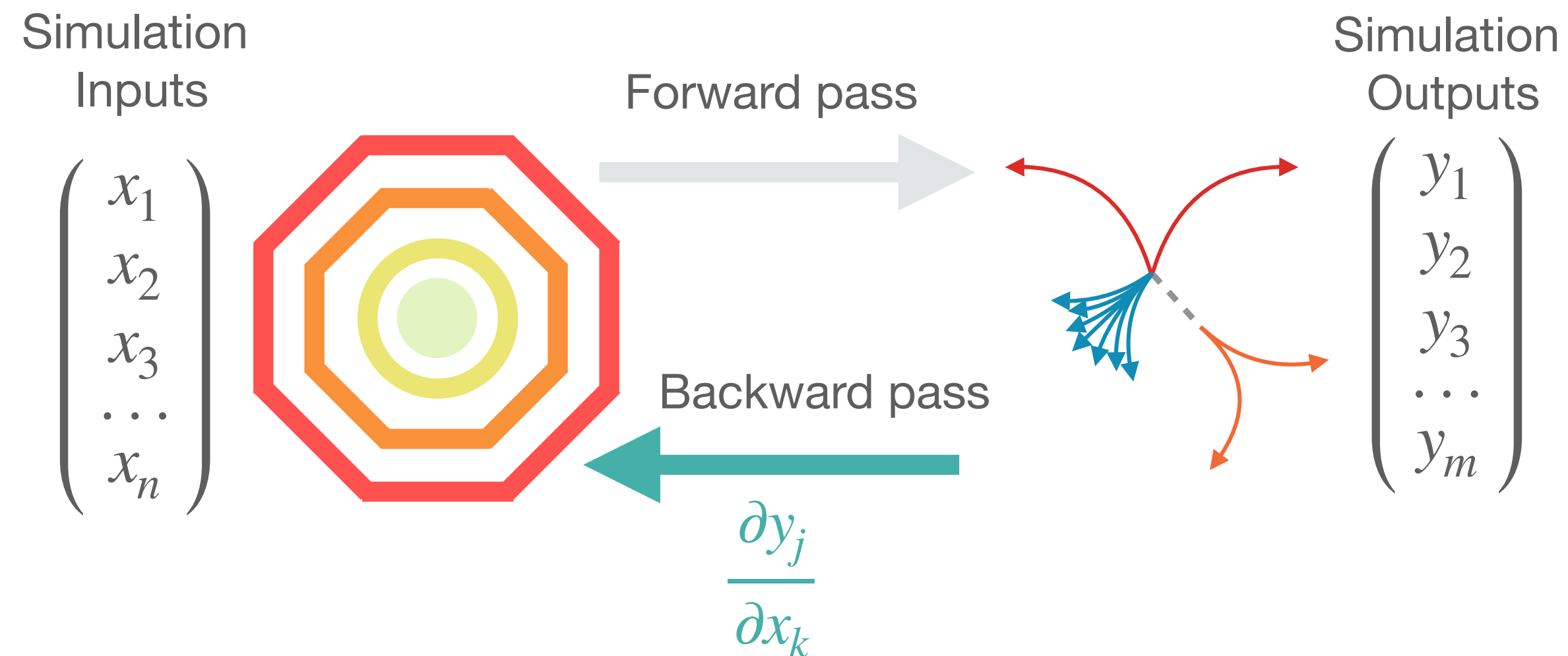


▶ We have been considering what should come after the HL-LHC...

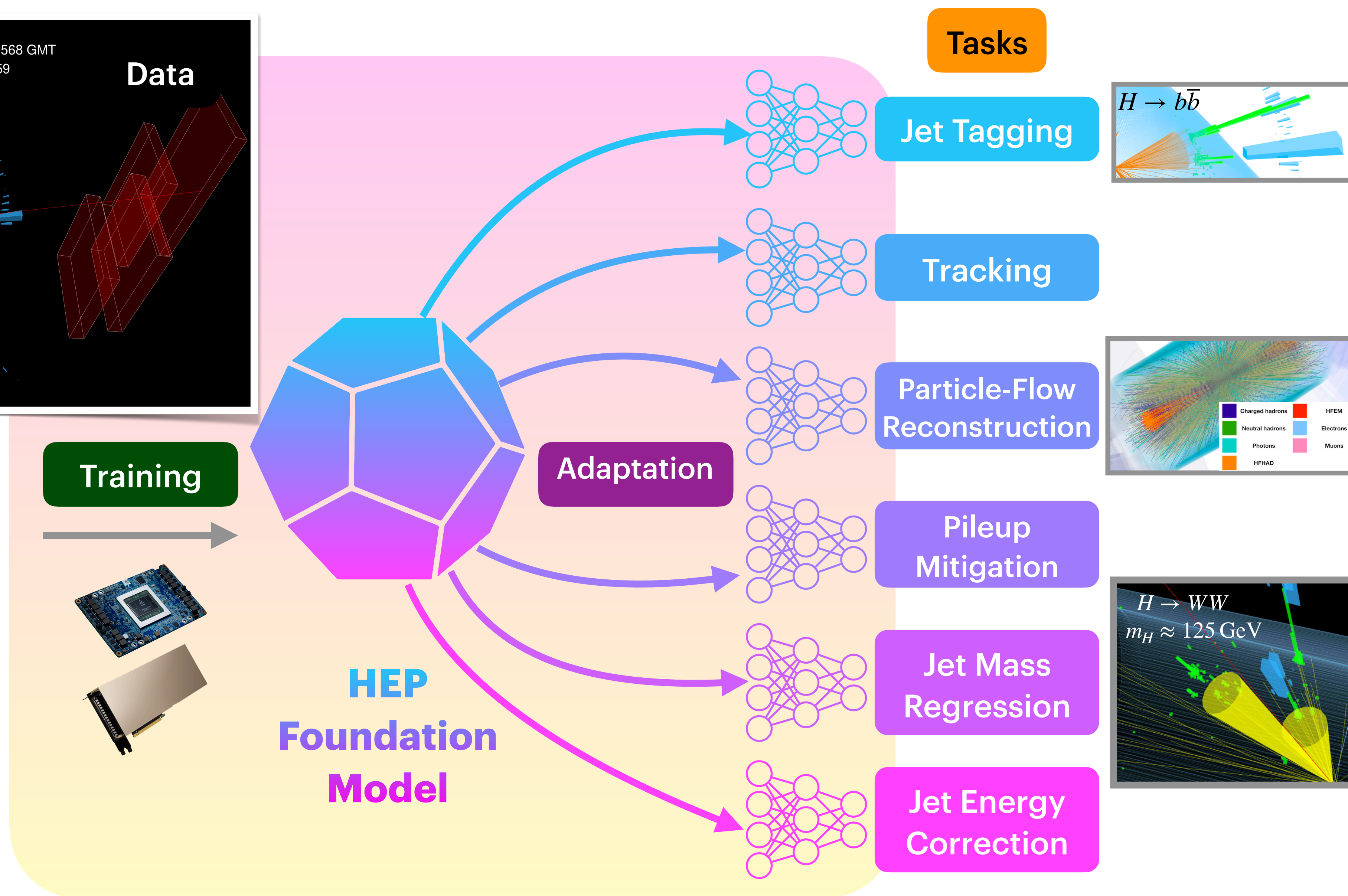
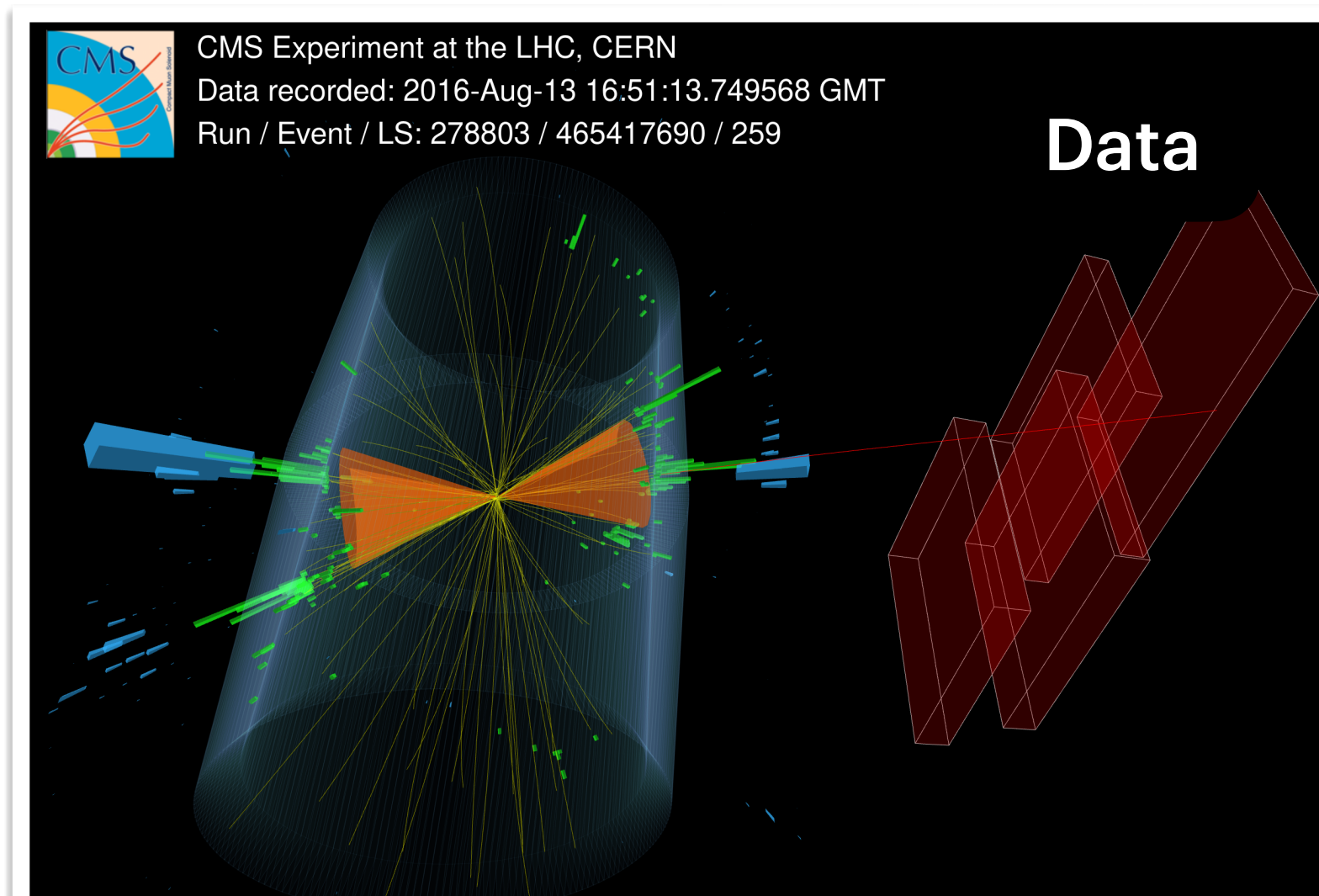




- ▶ ML has the promise to help us optimize the design of future colliders and detectors
  - ▶ Need all aspects of simulation chain to be implemented with differentiable programming [[arXiv:2002.04632](https://arxiv.org/abs/2002.04632)]
- ▶ Check out [Differentiable Almost Everything Workshop](#)



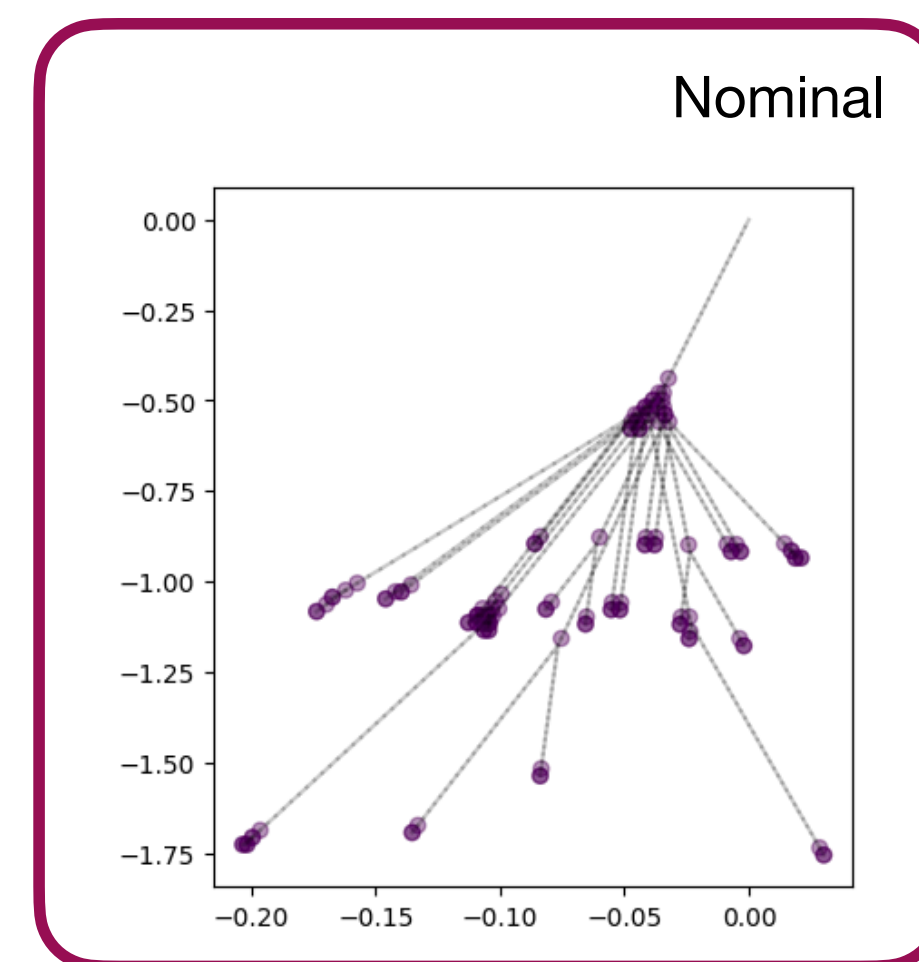
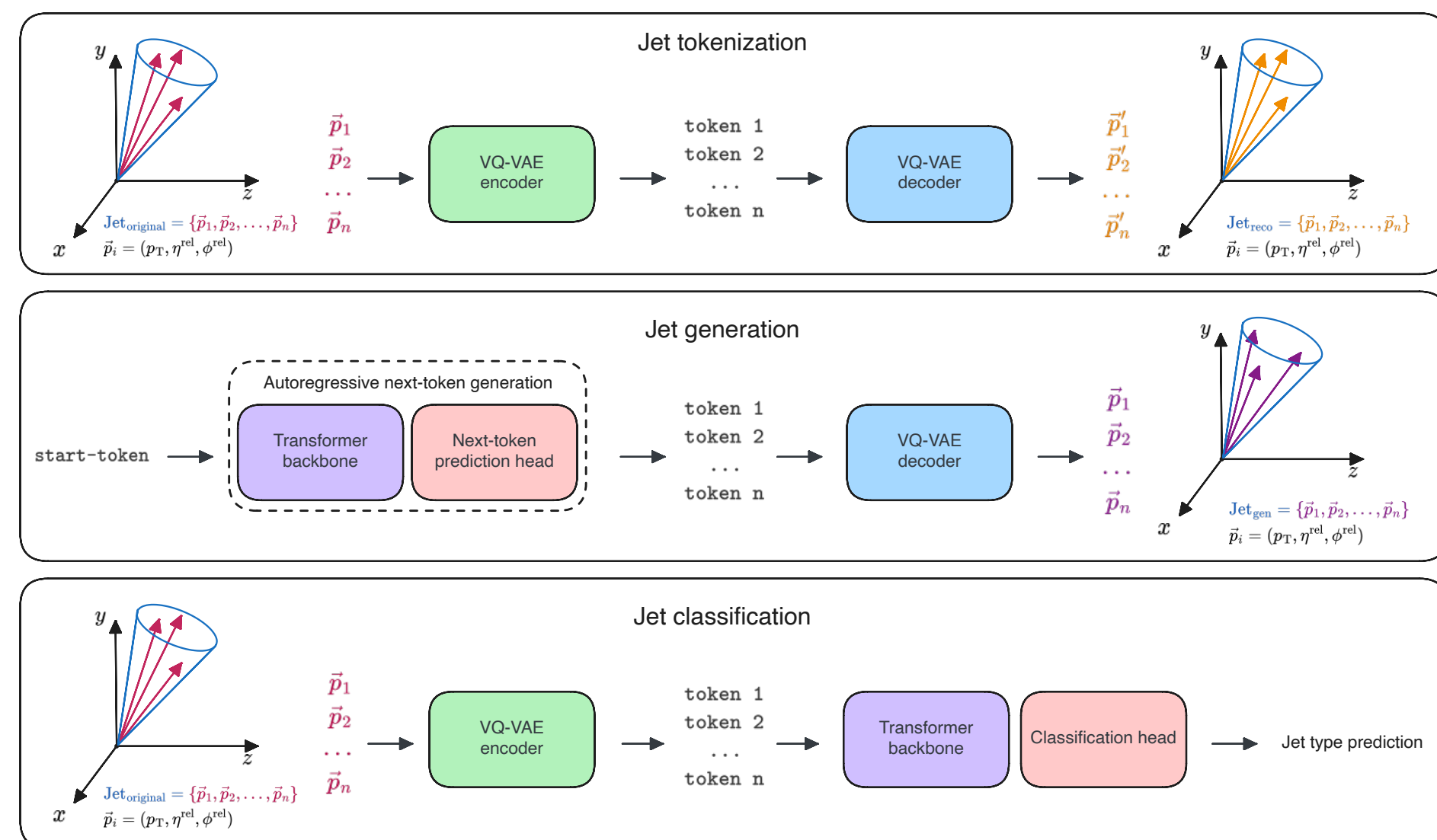
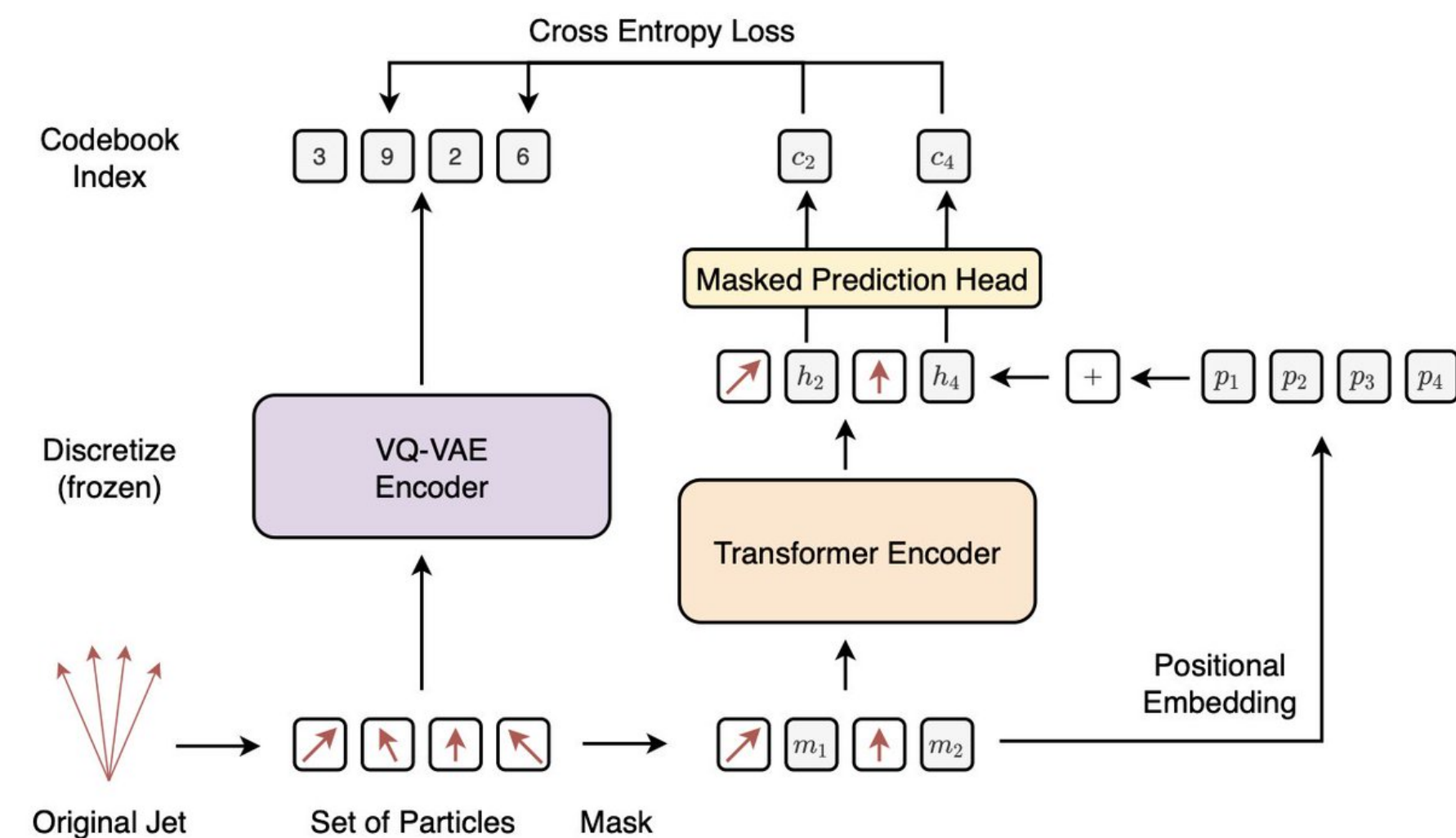
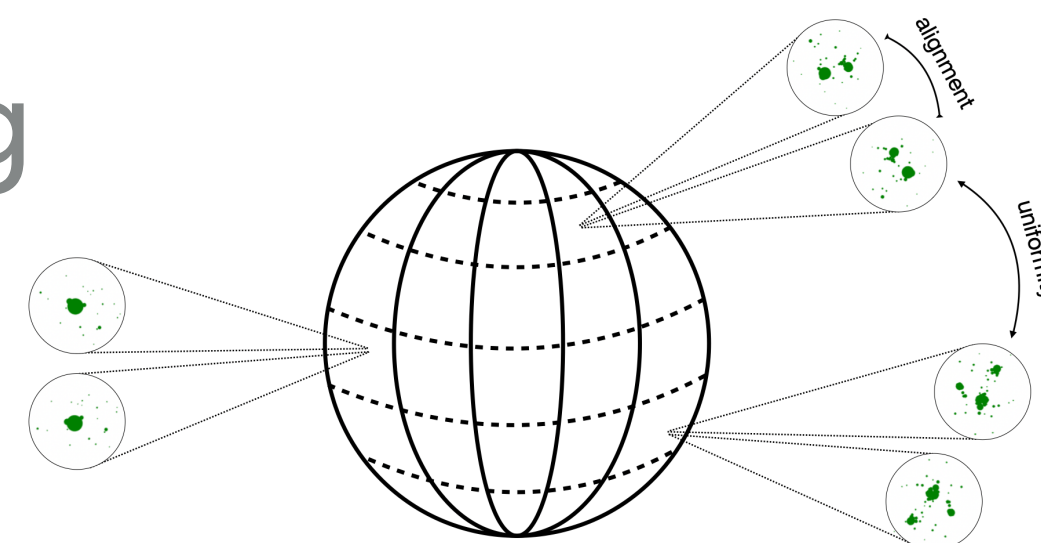




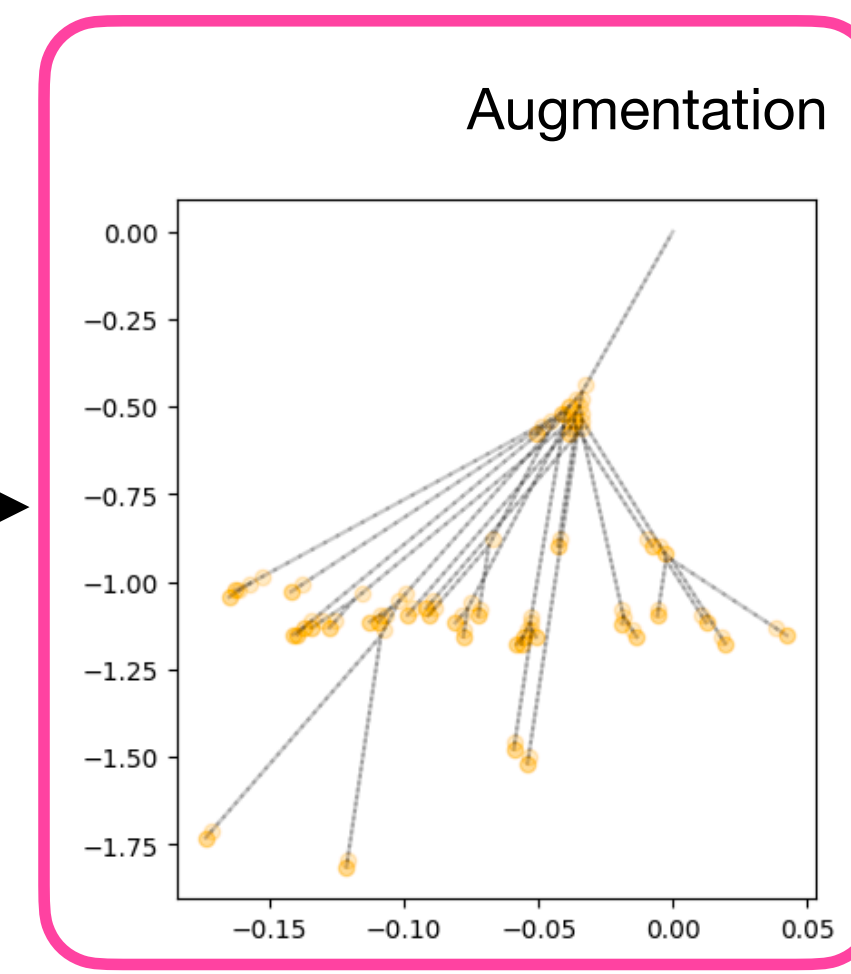


► Many studies implementing CS- and physics-inspired pre-training strategies

- JetCLR [[arXiv:2108.04253](https://arxiv.org/abs/2108.04253)]
- Masked particle modeling [[arXiv:2401.13537](https://arxiv.org/abs/2401.13537)]
- Resimulation [[arXiv:2403.07066](https://arxiv.org/abs/2403.07066)]
- GPT [[arXiv:2403.05618](https://arxiv.org/abs/2403.05618)]



Augment  
Re-shower

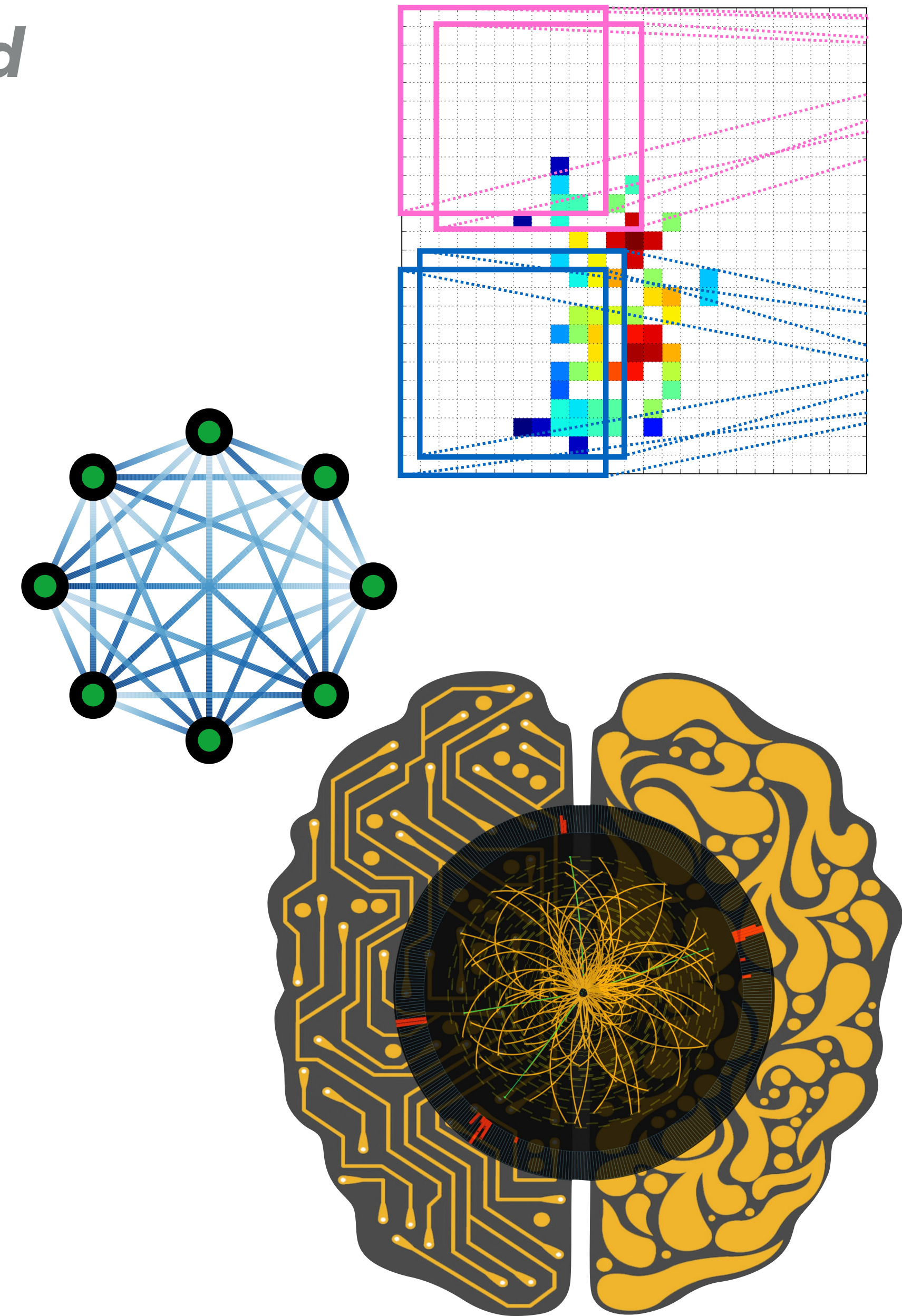


6

Initial parton fixed,  
only re-do shower  
+hadronization  
**Toy simulation**  
Credit to:  
M Kagan  
L Heinrich

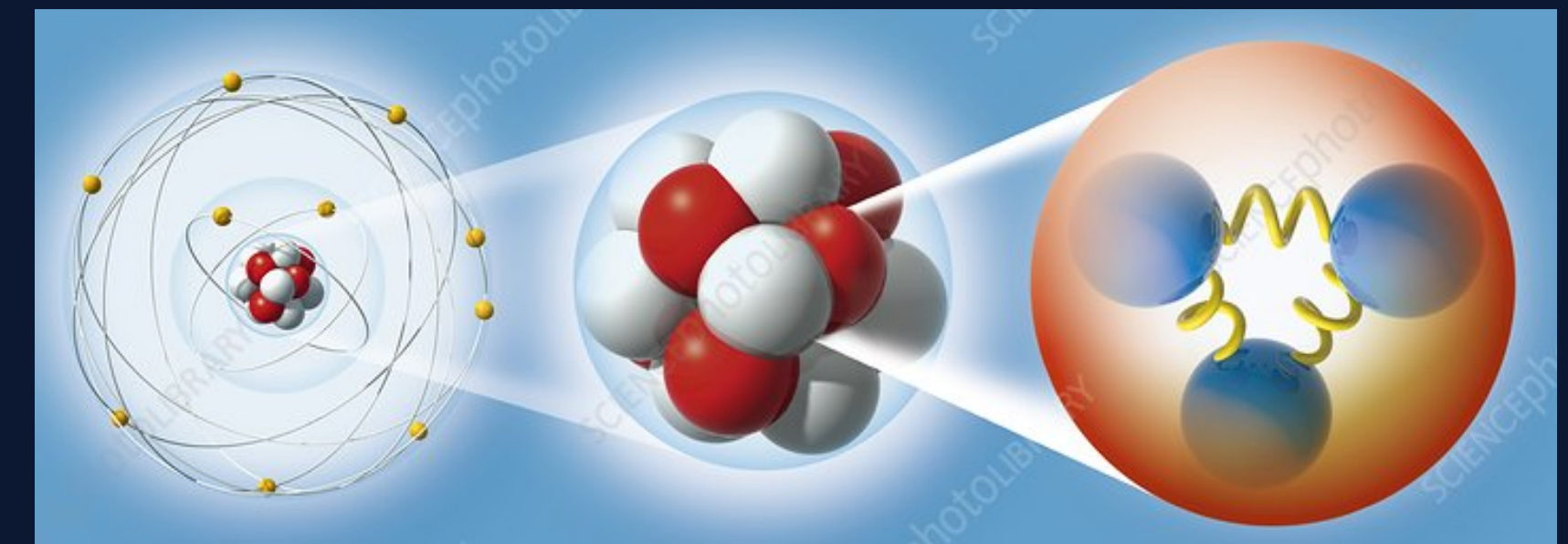
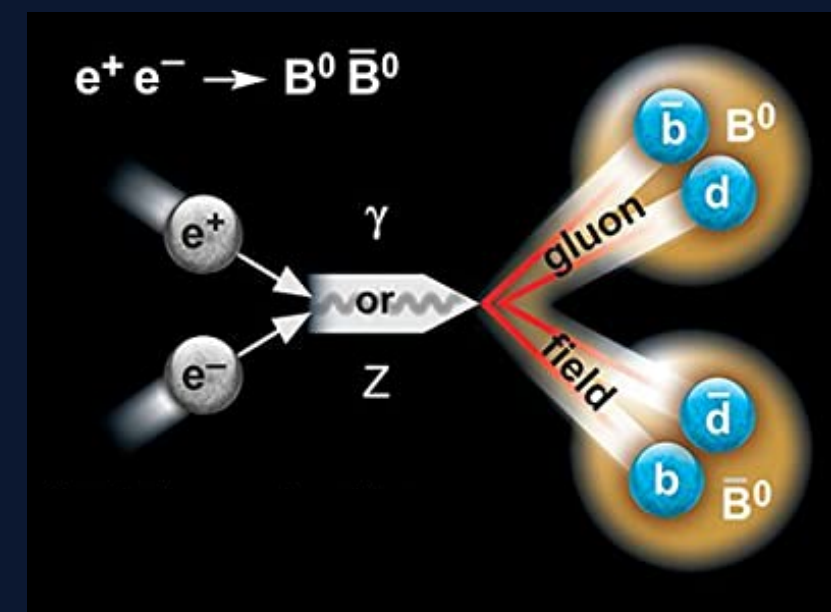
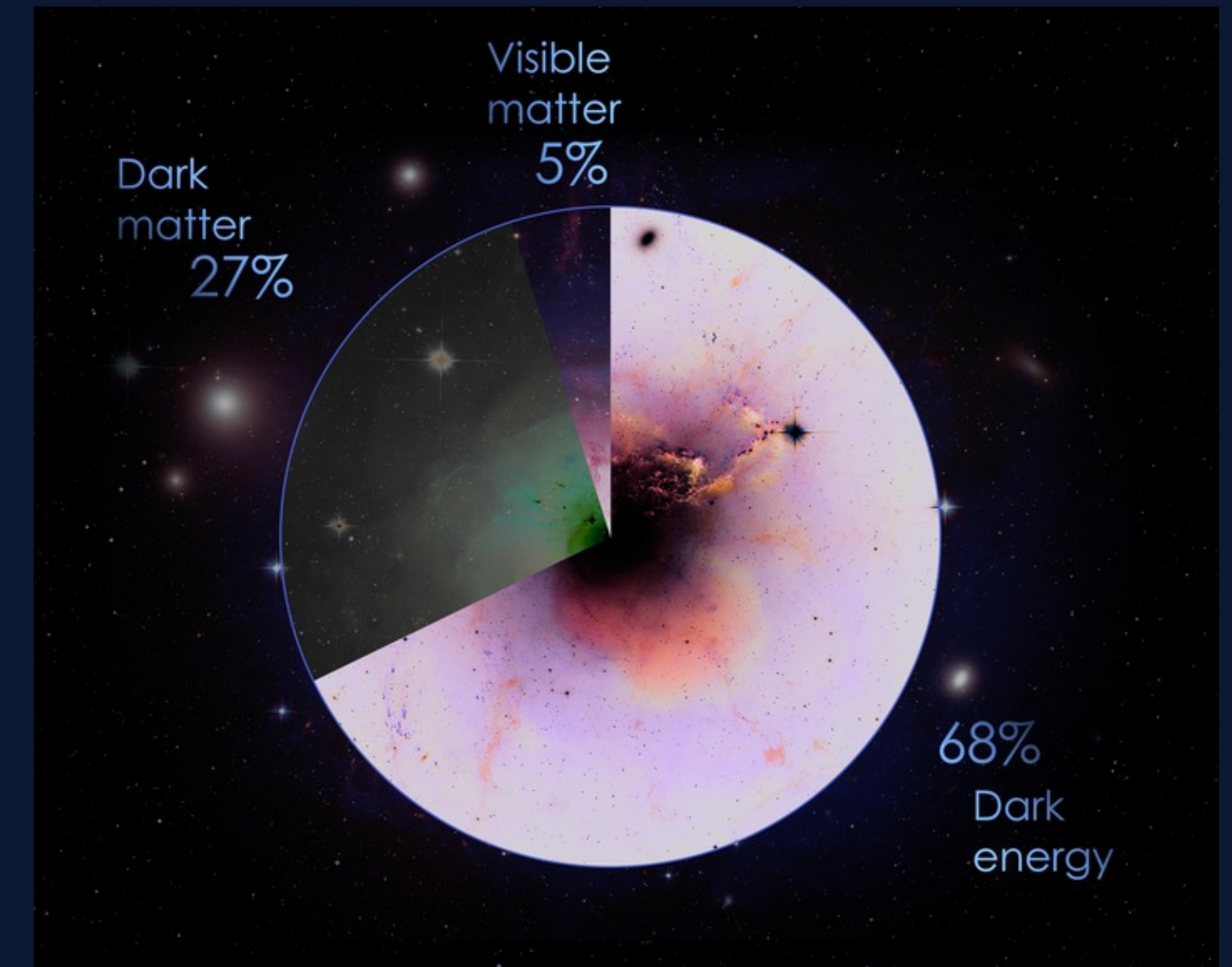


- ▶ Dizzying array of *ML opportunities, innovations, and applications* in *particle physics experiments*
- ▶ ML can help us solve major challenges for the next generation of particle physics experiments
- ▶ So we can (hopefully) get answers to **our big questions**

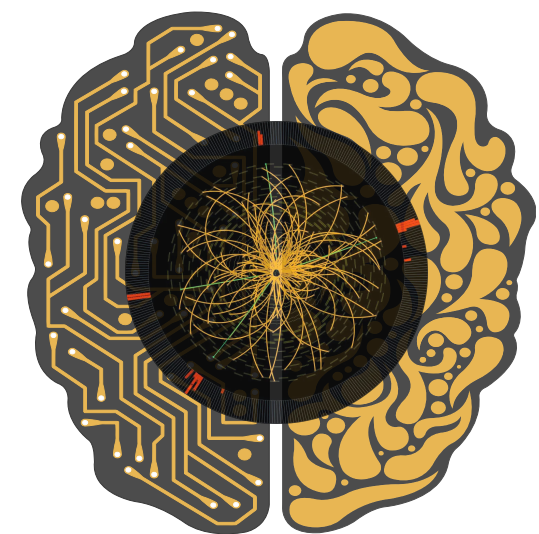




- ▶ What is our universe made of?
- ▶ What are the smallest building blocks of nature?
- ▶ How do they interact with each other?
- ▶ Is our universe stable?

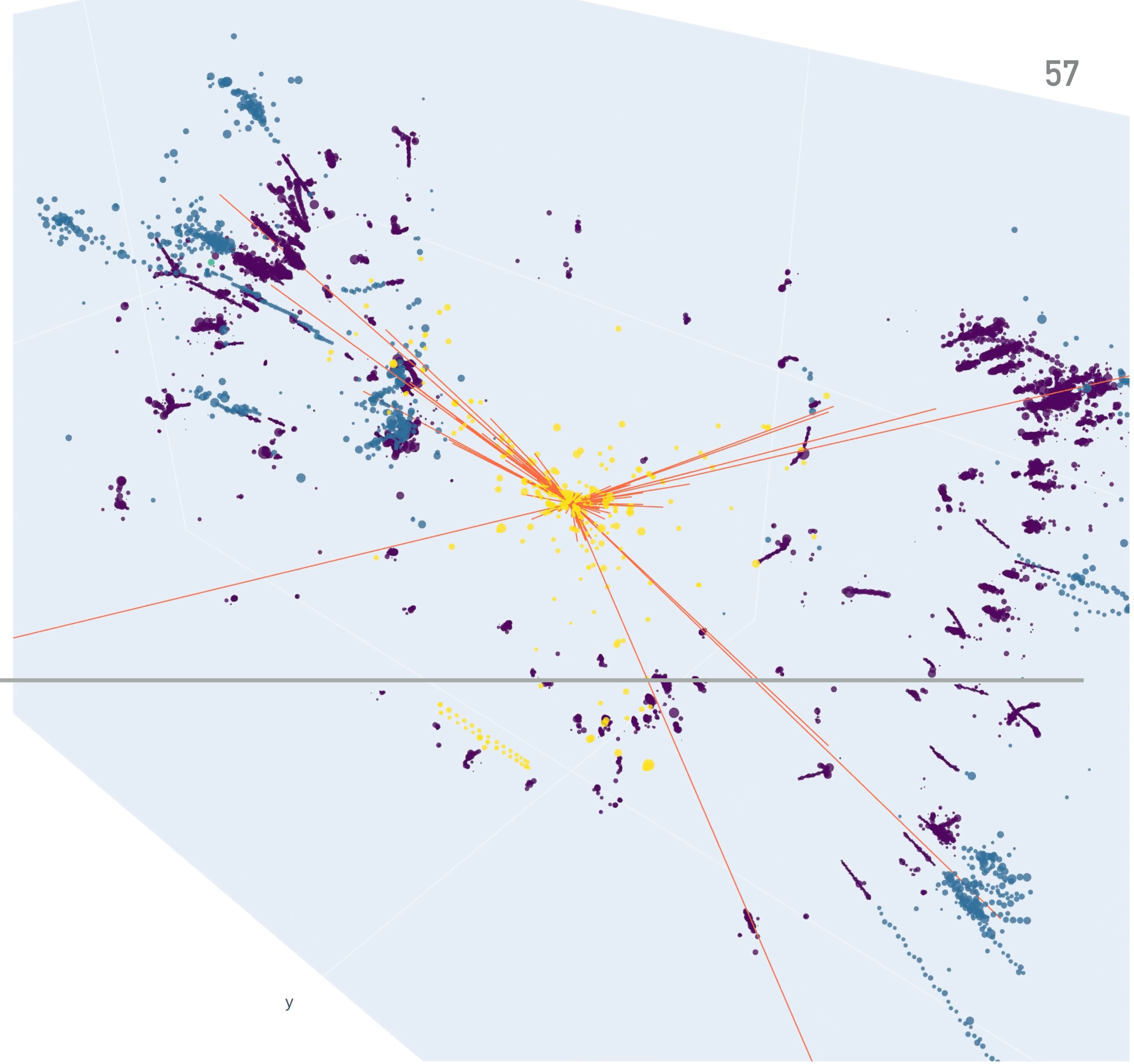






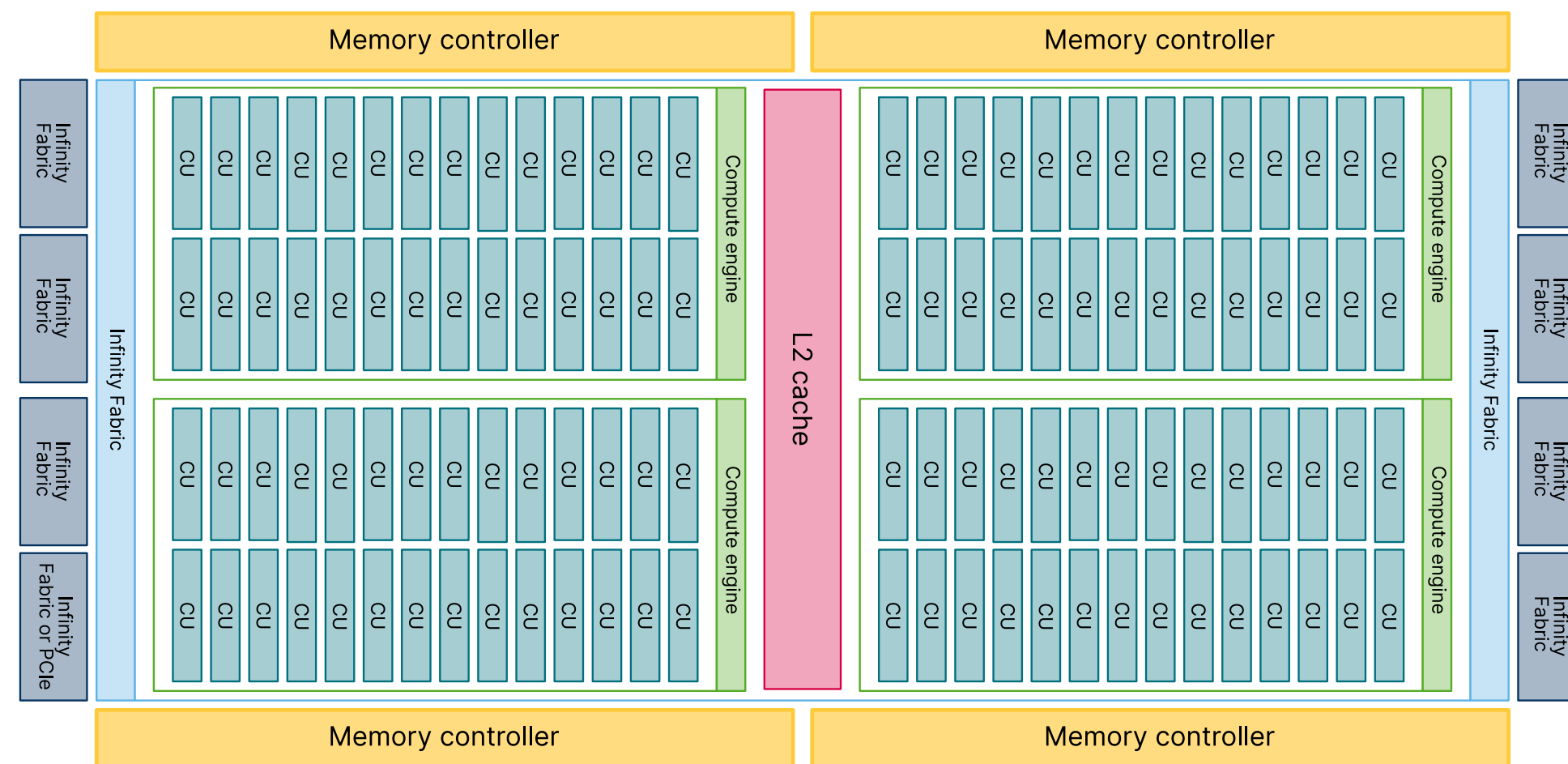
JAVIER DUARTE  
ICML 2024  
JULY 24, 2024

# BACKUP

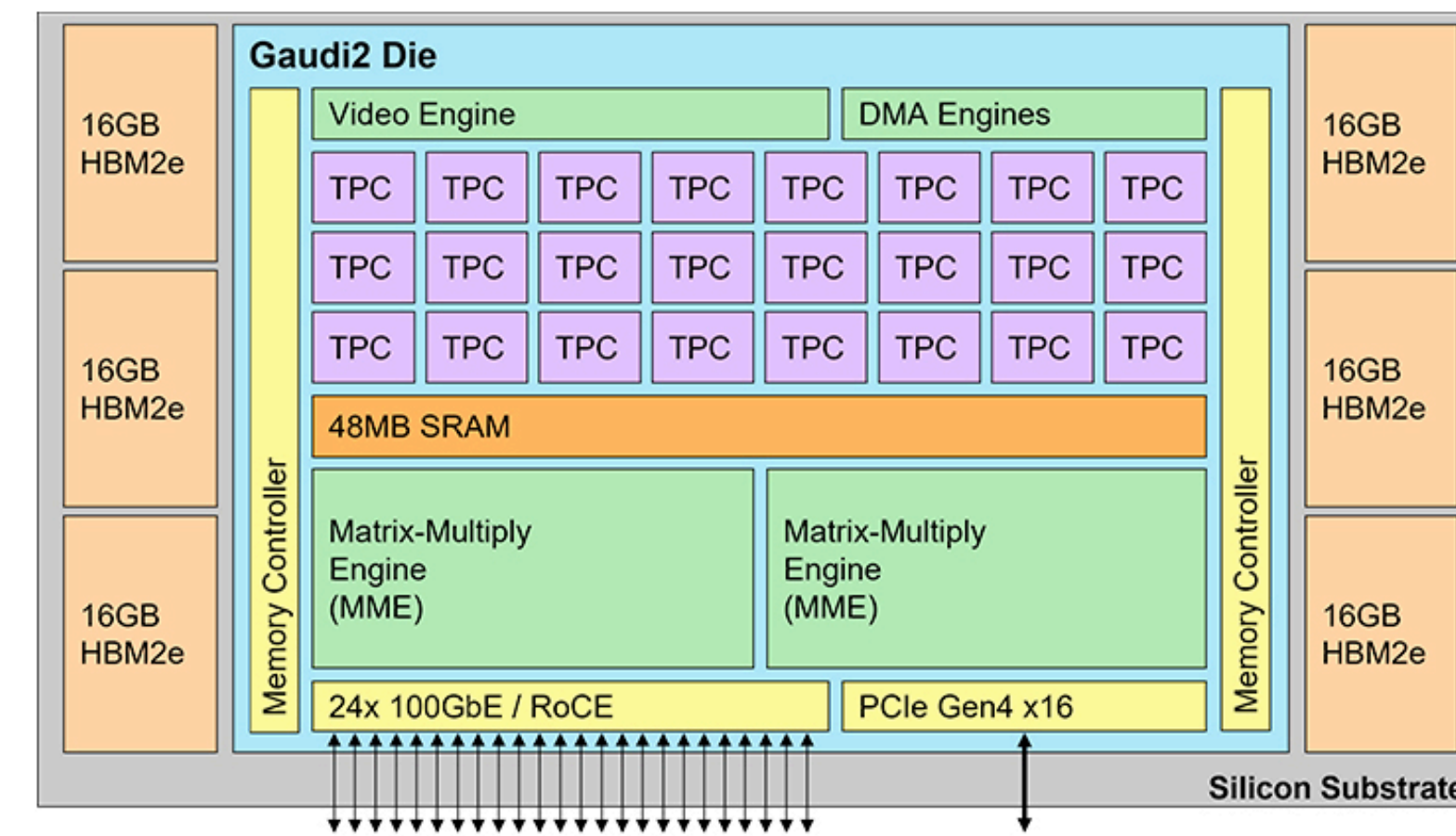




## The HPC AI chip landscape is diversifying



AMD MI250X GPU



Intel Gaudi2 deep learning processor

... we need flexible and portable codes to make use of these resources in the near future!



# PORTABILITY

Portable on CPU,  
Nvidia & AMD GPU,  
Intel Habana Gaudi  
chips

