



JAVIER DUARTE
LISHEP SESSION C
JULY 6, 2021

MACHINE LEARNING FOR (EXPERIMENTAL) HIGH ENERGY PHYSICS

NEURAL NETWORKS AND CELLULAR AUTOMATA IN EXPERIMENTAL HIGH ENERGY PHYSICS

B. DENBY

Laboratoire de l'Accélérateur Linéaire, Orsay, France

Received 20 September 1987; in revised form 28 December 1987

Within the past few years, two novel computing techniques, cellular automata and neural networks, have shown considerable promise in the solution of problems of a very high degree of complexity, such as turbulent fluid flow, image processing, and pattern recognition. Many of the problems faced in experimental high energy physics are also of this nature. Track reconstruction in wire chambers and cluster finding in cellular calorimeters, for instance, involve pattern recognition and high combinatorial complexity since many combinations of hits or cells must be considered in order to arrive at the final tracks or clusters. Here we examine in what way connective network methods can be applied to some of the problems of experimental high energy physics. It is found that such problems as track and cluster finding adapt naturally to these approaches. When large scale hard-wired connective networks become available, it will be possible to realize solutions to such problems in a fraction of the time required by traditional methods. For certain types of problems, faster solutions are already possible using model networks implemented on vector or other massively parallel machines. It should also be possible, using existing technology, to build simplified networks that will allow detailed reconstructed event information to be used in fast trigger decisions.

INTRODUCTION

- ▶ Particle physics has been linked to machine learning and neural networks since the 1980s!

Computer Physics Communications 49 (1988) 429–448
North-Holland, Amsterdam

429

NEURAL NETWORKS AND CELLULAR AUTOMATA IN EXPERIMENTAL HIGH ENERGY PHYSICS

B. DENBY

Laboratoire de l'Accélérateur Linéaire, Orsay, France

Received 20 September 1987; in revised form 28 December 1987

Within the past few years, two novel computing techniques, cellular automata and neural networks, have shown considerable promise in the solution of problems of a very high degree of complexity, such as turbulent fluid flow, image processing, and pattern recognition. Many of the problems faced in experimental high energy physics are also of this nature. Track reconstruction in wire chambers and cluster finding in cellular calorimeters, for instance, involve pattern recognition and high combinatorial complexity since many combinations of hits or cells must be considered in order to arrive at the final tracks or clusters. Here we examine in what way connective network methods can be applied to some of the problems of experimental high energy physics. It is found that such problems as track and cluster finding adapt naturally to these approaches. When large scale hard-wired connective networks become available, it will be possible to realize solutions to such problems in a fraction of the time required by traditional methods. For certain types of problems, faster solutions are already possible using model networks implemented on vector or other massively parallel machines. It should also be possible, using existing technology, to build simplified networks that will allow detailed reconstructed event information to be used in fast trigger decisions.

INTRODUCTION

- ▶ Particle physics has been linked to machine learning and neural networks since the 1980s!
- ▶ In recent years, the use of ML has expanded into new territory

Computer Physics Communications 49 (1988) 429–448
North-Holland, Amsterdam

429

NEURAL NETWORKS AND CELLULAR AUTOMATA IN EXPERIMENTAL HIGH ENERGY PHYSICS

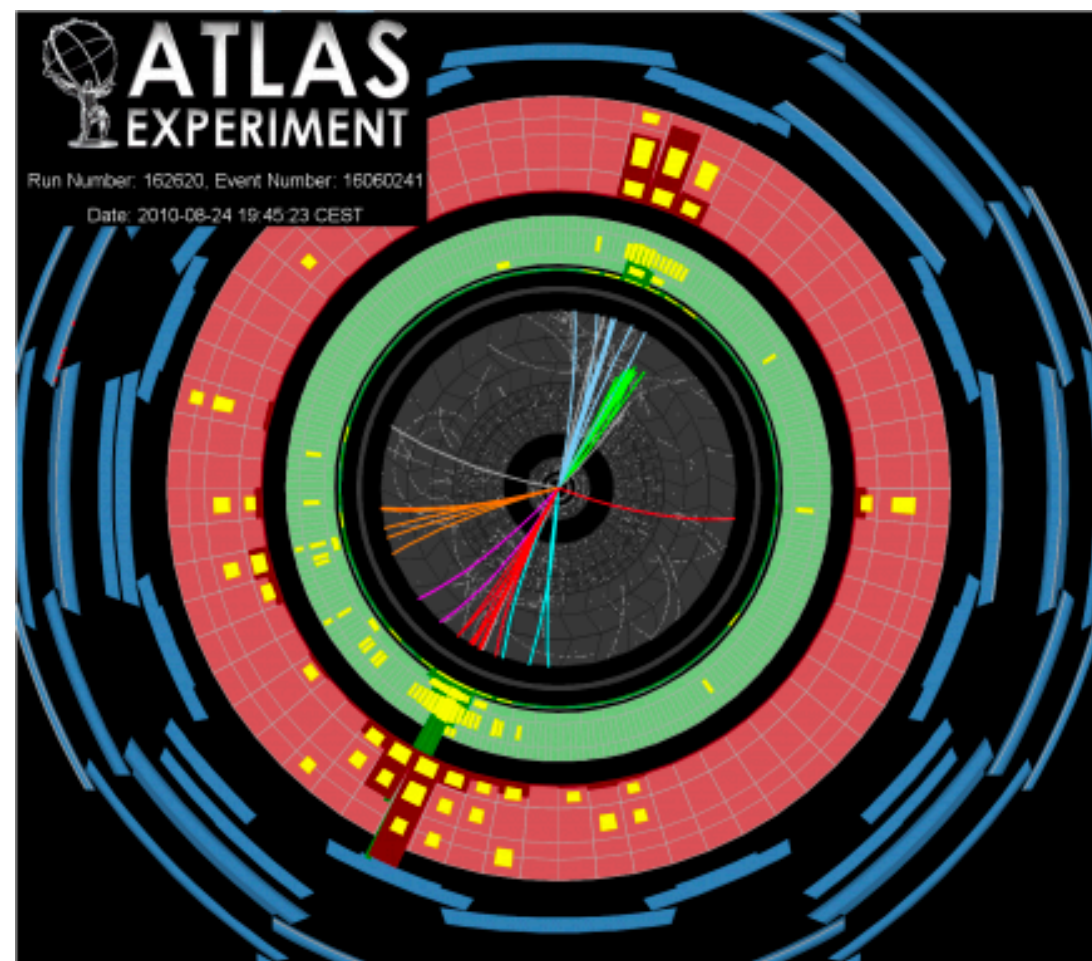
B. DENBY

Laboratoire de l'Accélérateur Linéaire, Orsay, France

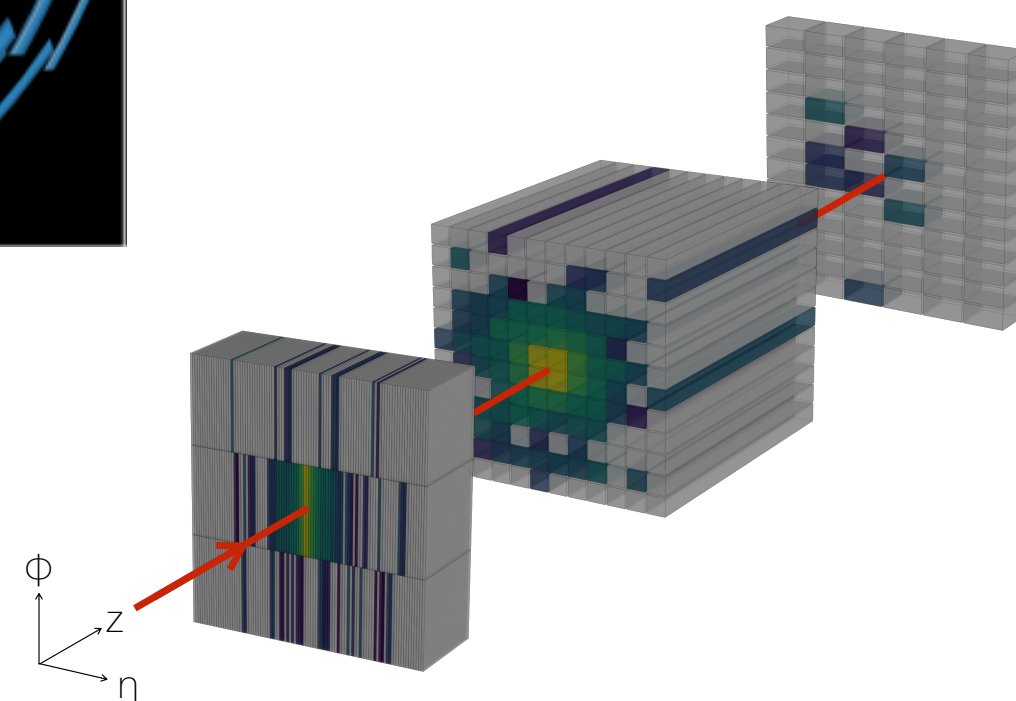
Received 20 September 1987; in revised form 28 December 1987

Within the past few years, two novel computing techniques, cellular automata and neural networks, have shown considerable promise in the solution of problems of a very high degree of complexity, such as turbulent fluid flow, image processing, and pattern recognition. Many of the problems faced in experimental high energy physics are also of this nature. Track reconstruction in wire chambers and cluster finding in cellular calorimeters, for instance, involve pattern recognition and high combinatorial complexity since many combinations of hits or cells must be considered in order to arrive at the final tracks or clusters. Here we examine in what way connective network methods can be applied to some of the problems of experimental high energy physics. It is found that such problems as track and cluster finding adapt naturally to these approaches. When large scale hard-wired connective networks become available, it will be possible to realize solutions to such problems in a fraction of the time required by traditional methods. For certain types of problems, faster solutions are already possible using model networks implemented on vector or other massively parallel machines. It should also be possible, using existing technology, to build simplified networks that will allow detailed reconstructed event information to be used in fast trigger decisions.

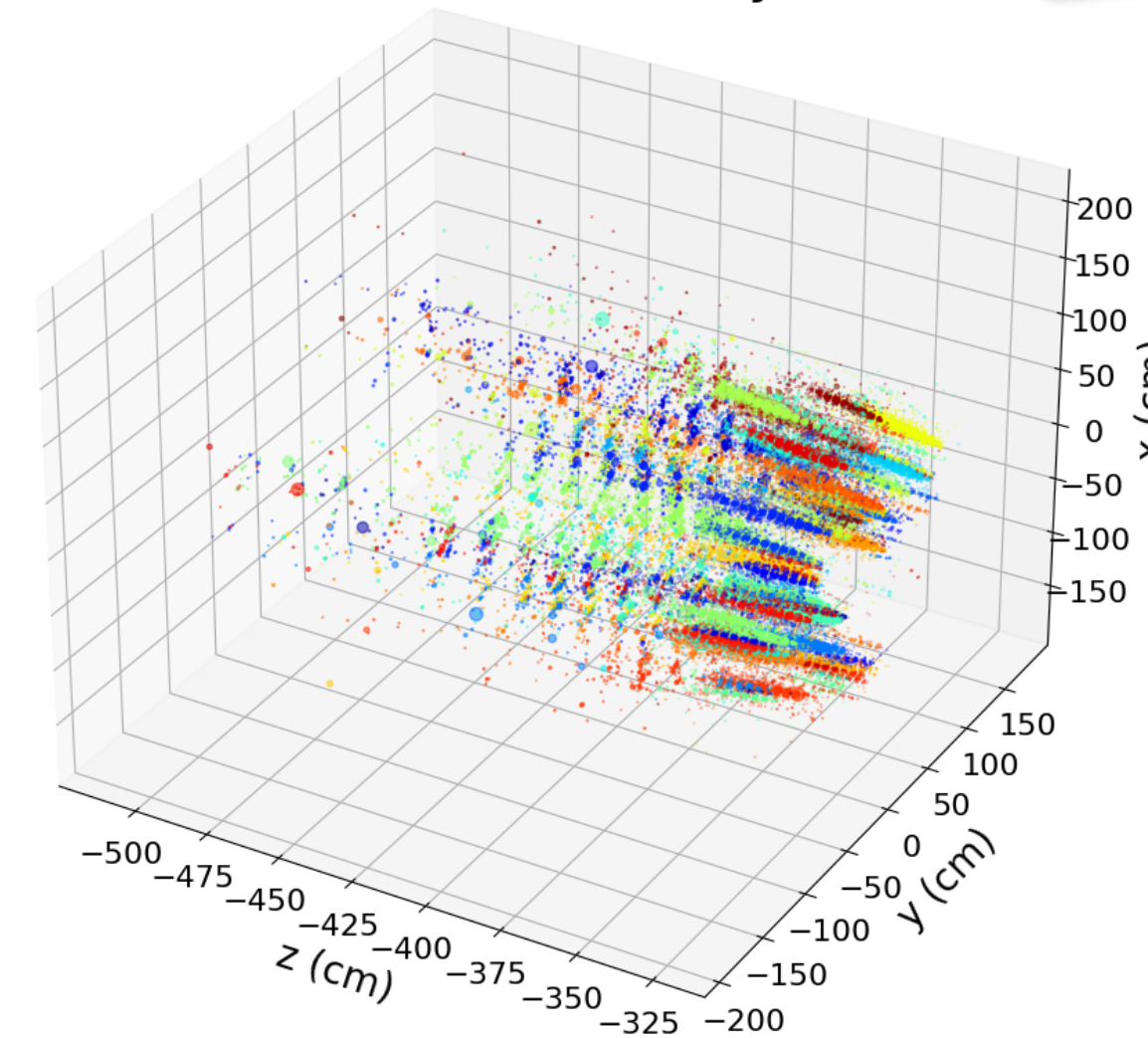
ML for "jet tagging"



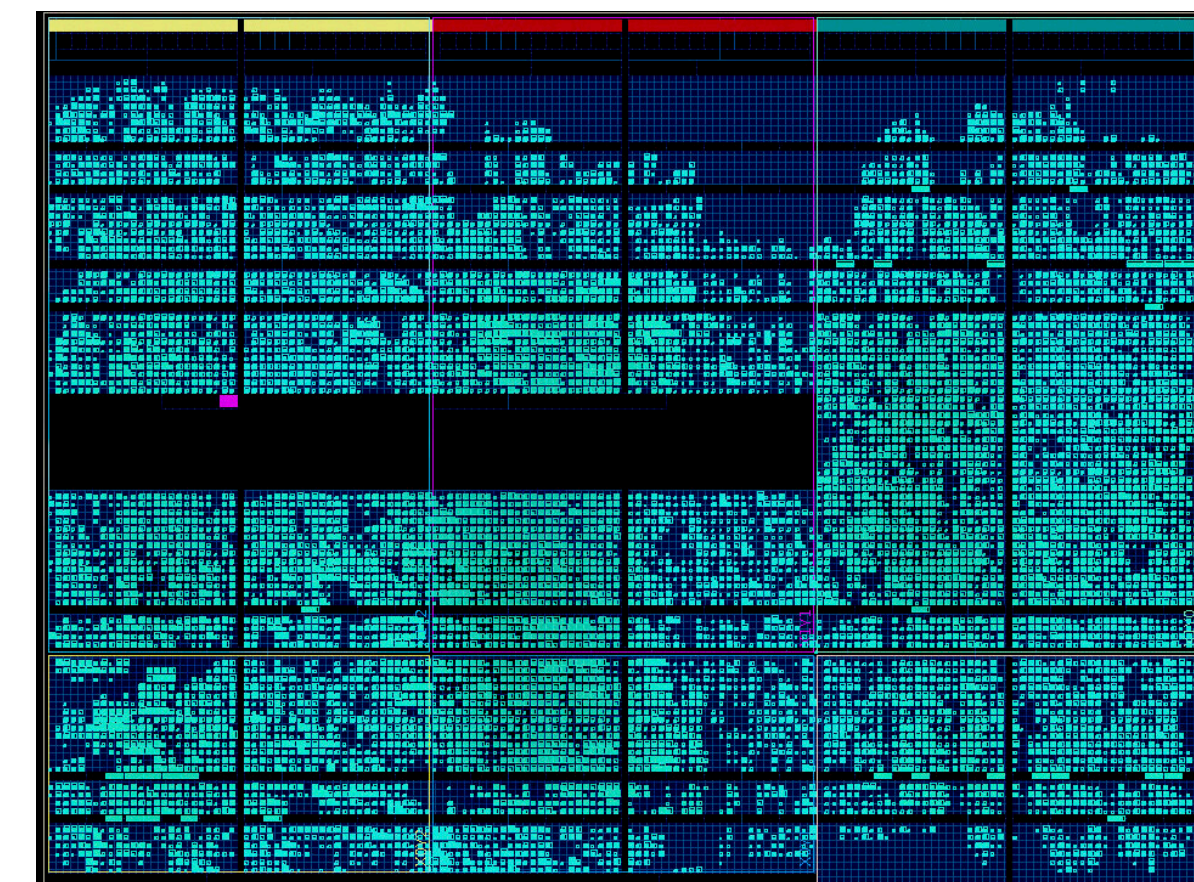
ML for reconstruction



CMS Phase-2 Simulation Preliminary



Fast ML for trigger



ML for generation/simulation

INTRODUCTION

- ▶ Particle physics has been linked to machine learning and neural networks since the 1980s!
- ▶ In recent years, the use of ML has expanded into new territory
- ▶ Broadly speaking, we're interested in two advantages from ML: sensitivity to physics and computational performance

Computer Physics Communications 49 (1988) 429–448
North-Holland, Amsterdam

429

NEURAL NETWORKS AND CELLULAR AUTOMATA IN EXPERIMENTAL HIGH ENERGY PHYSICS

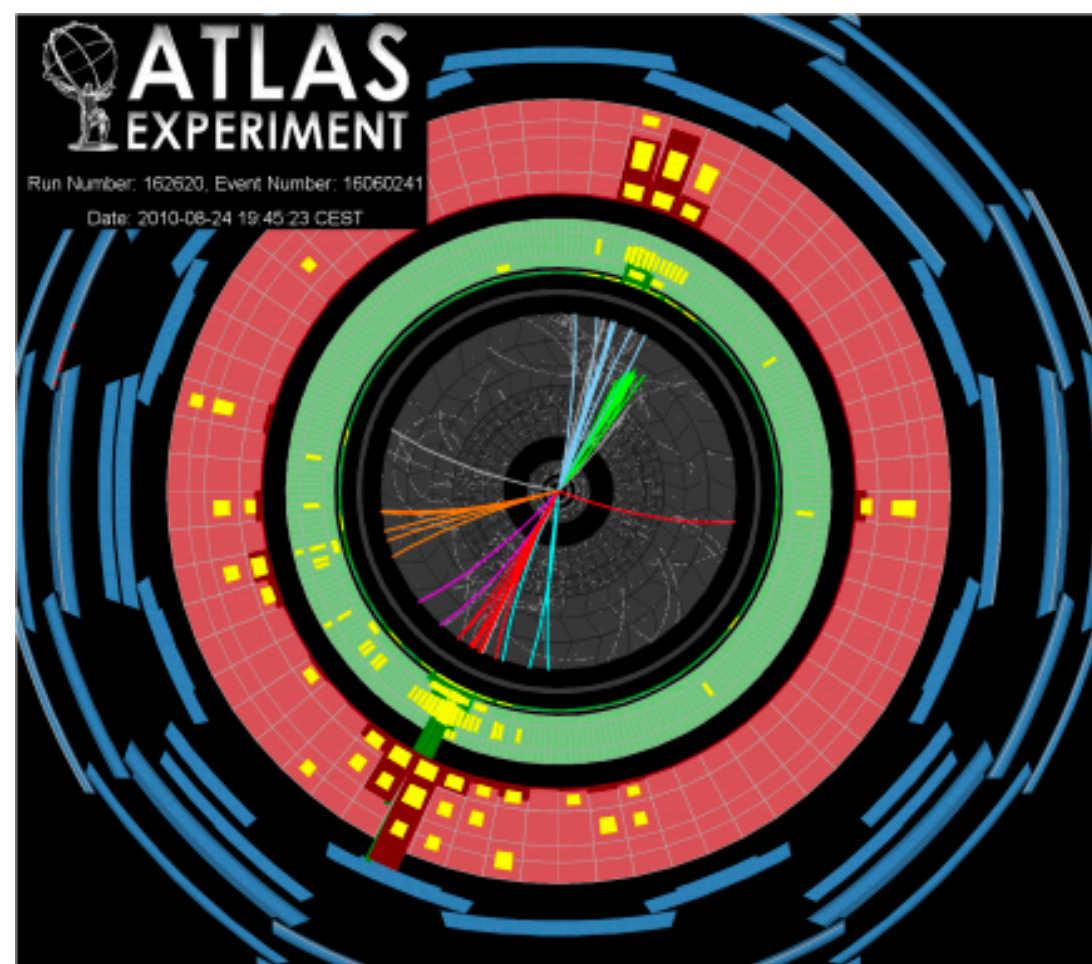
B. DENBY

Laboratoire de l'Accélérateur Linéaire, Orsay, France

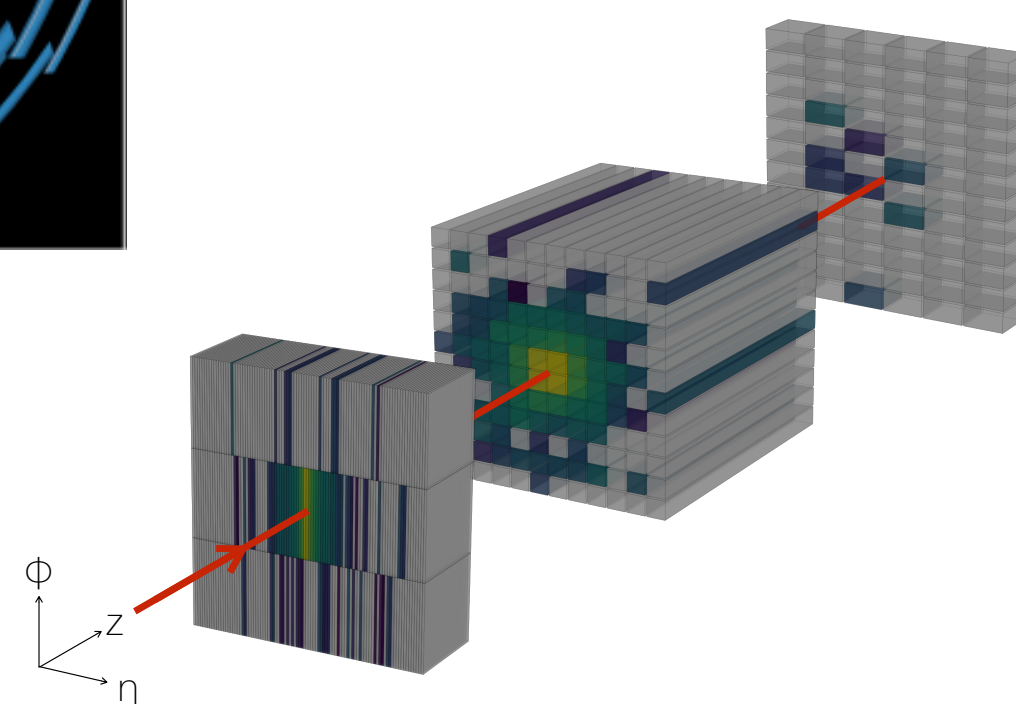
Received 20 September 1987; in revised form 28 December 1987

Within the past few years, two novel computing techniques, cellular automata and neural networks, have shown considerable promise in the solution of problems of a very high degree of complexity, such as turbulent fluid flow, image processing, and pattern recognition. Many of the problems faced in experimental high energy physics are also of this nature. Track reconstruction in wire chambers and cluster finding in cellular calorimeters, for instance, involve pattern recognition and high combinatorial complexity since many combinations of hits or cells must be considered in order to arrive at the final tracks or clusters. Here we examine in what way connective network methods can be applied to some of the problems of experimental high energy physics. It is found that such problems as track and cluster finding adapt naturally to these approaches. When large scale hard-wired connective networks become available, it will be possible to realize solutions to such problems in a fraction of the time required by traditional methods. For certain types of problems, faster solutions are already possible using model networks implemented on vector or other massively parallel machines. It should also be possible, using existing technology, to build simplified networks that will allow detailed reconstructed event information to be used in fast trigger decisions.

ML for "jet tagging"

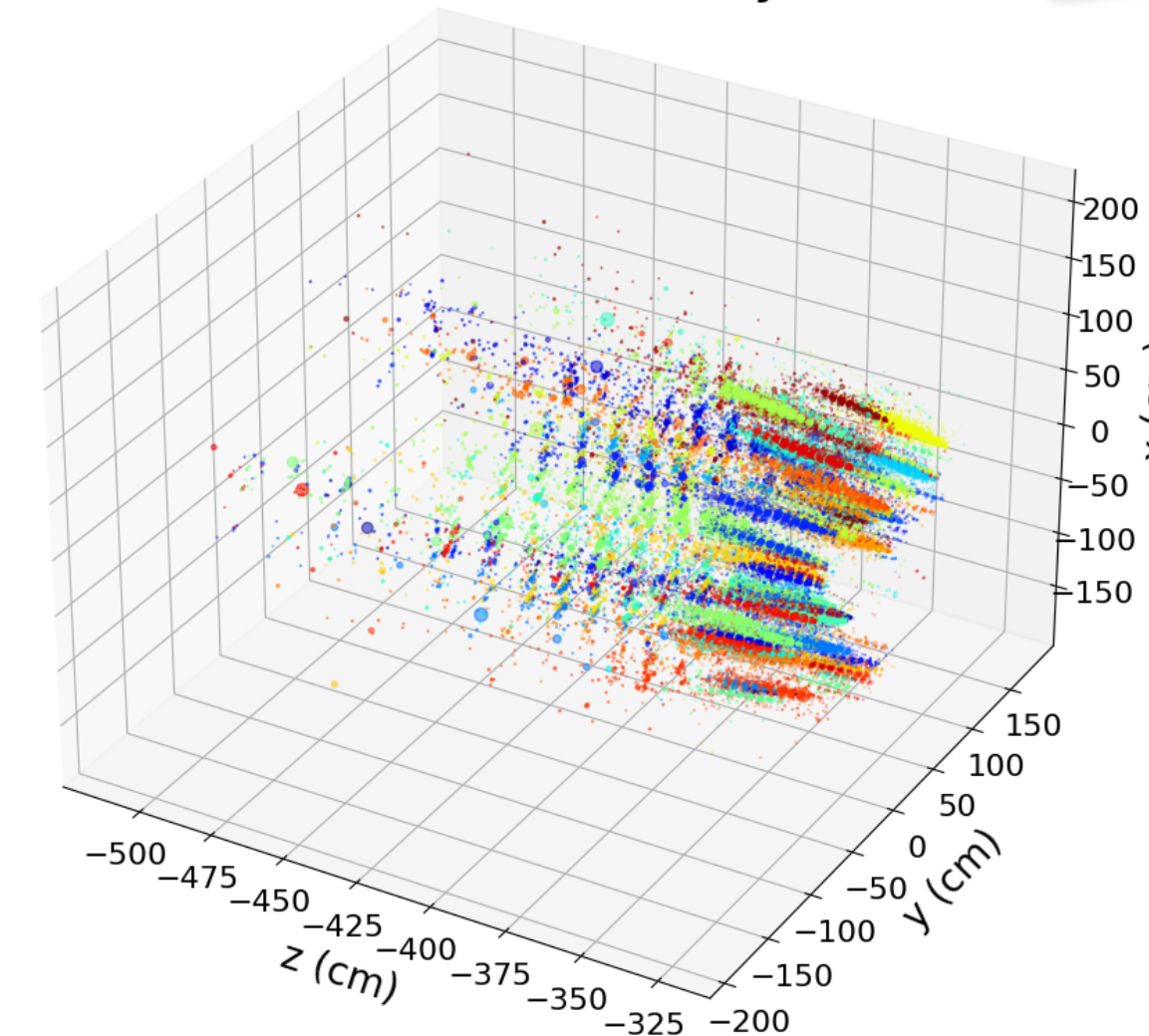


ML for reconstruction

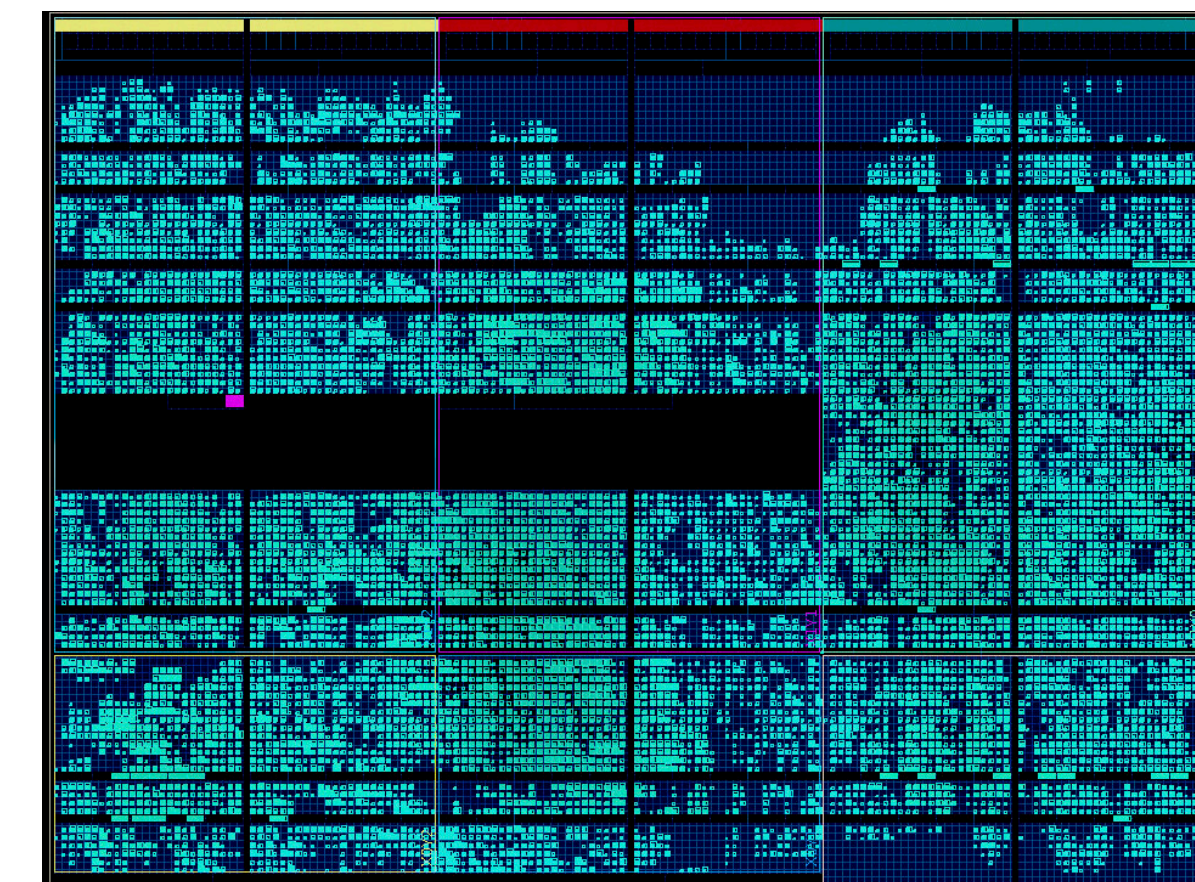


ML for generation/simulation

CMS Phase-2 Simulation Preliminary



Fast ML for trigger



- ▶ Deep (machine) learning is the use of **structured neural networks** with **many hidden layers** as generic functions to approximate the optimal solution for a given task

- ▶ Deep (machine) learning is the use of **structured neural networks** with **many hidden layers** as generic functions to approximate the optimal solution for a given task
- ▶ Why use deep neural networks in particle physics?

- ▶ Deep (machine) learning is the use of **structured neural networks** with **many hidden layers** as generic functions to approximate the optimal solution for a given task
- ▶ Why use deep neural networks in particle physics?
 - ▶ They work!

- ▶ Deep (machine) learning is the use of **structured neural networks** with **many hidden layers** as generic functions to approximate the optimal solution for a given task
- ▶ Why use deep neural networks in particle physics?
 - ▶ They work!
 - ▶ Lots of data & simulation to train them

- ▶ Deep (machine) learning is the use of **structured neural networks** with **many hidden layers** as generic functions to approximate the optimal solution for a given task
- ▶ Why use deep neural networks in particle physics?
 - ▶ They work!
 - ▶ Lots of data & simulation to train them
 - ▶ Fast ways to train them



TensorFlow



Keras



PyTorch

- ▶ Deep (machine) learning is the use of **structured neural networks** with **many hidden layers** as generic functions to approximate the optimal solution for a given task
- ▶ Why use deep neural networks in particle physics?
 - ▶ They work!
 - ▶ Lots of data & simulation to train them
 - ▶ Fast ways to train them
 - ▶ Possibly gain new insights...



TensorFlow

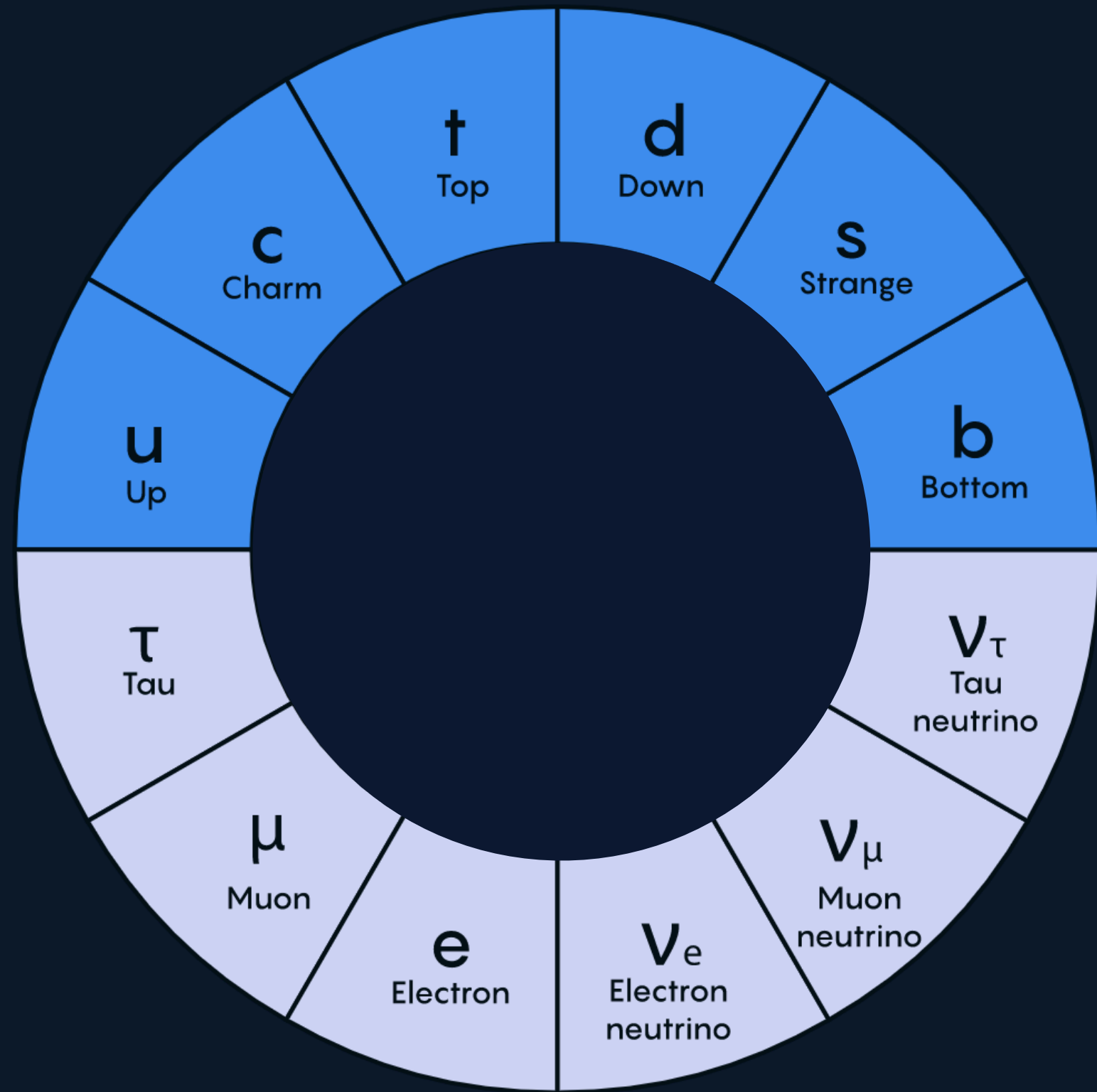


ML FOR JET TAGGING

ML FOR GEN/SIM

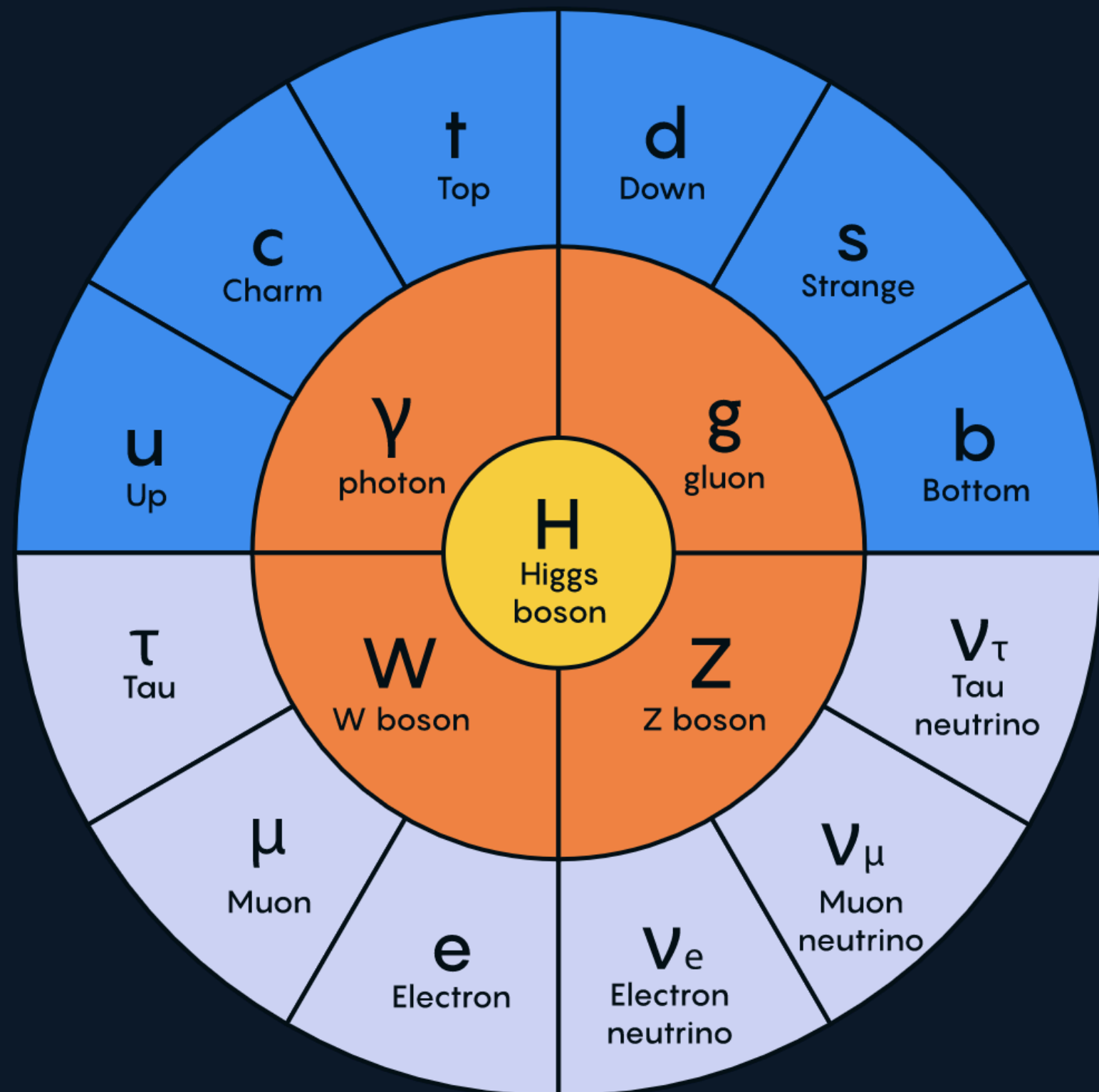
FAST ML FOR TRIGGER

HIGGS BOSON IN THE STANDARD MODEL



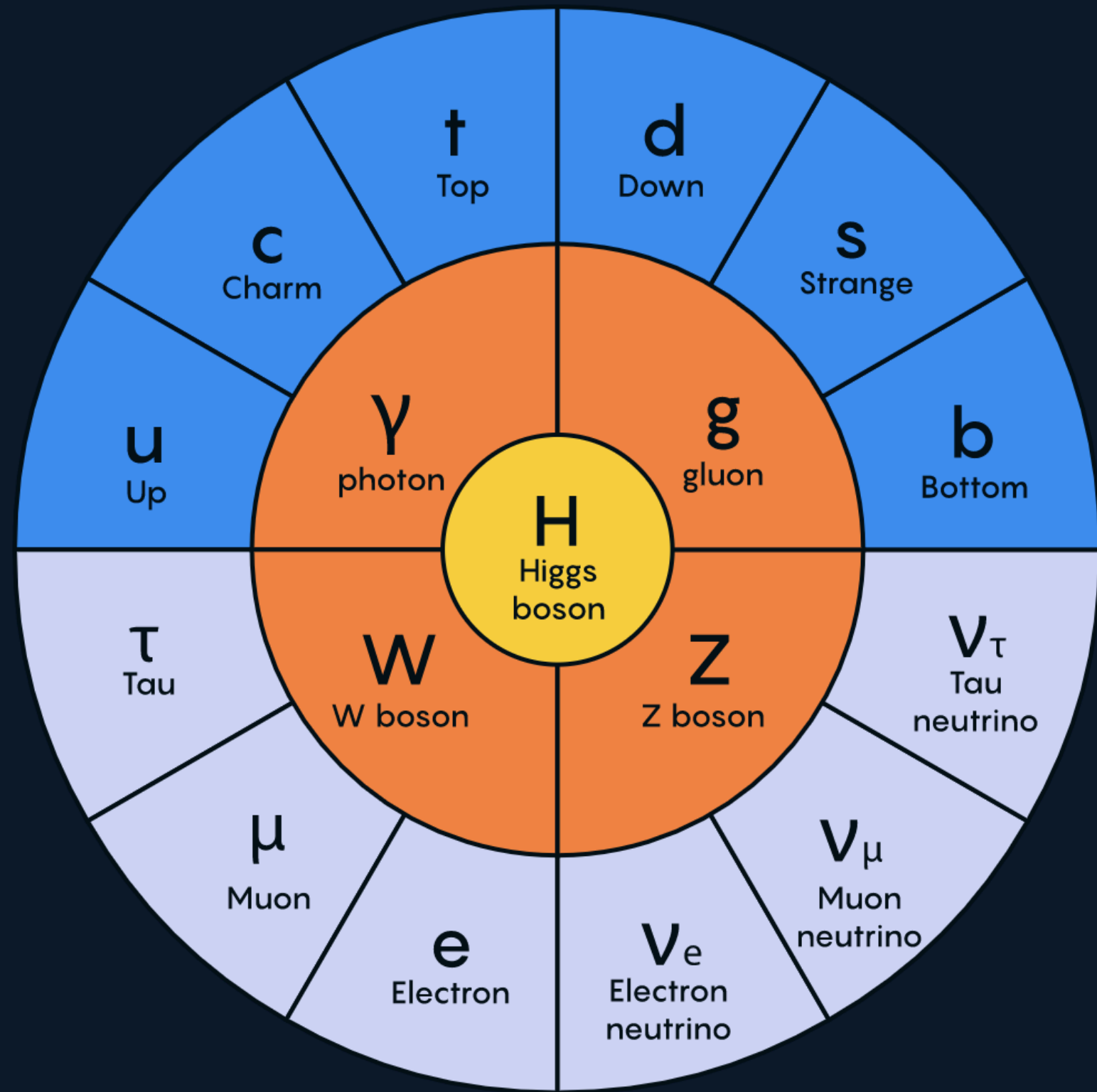
FERMIONS (MATTER)
● QUARKS ● LEPTONS

HIGGS BOSON IN THE STANDARD MODEL

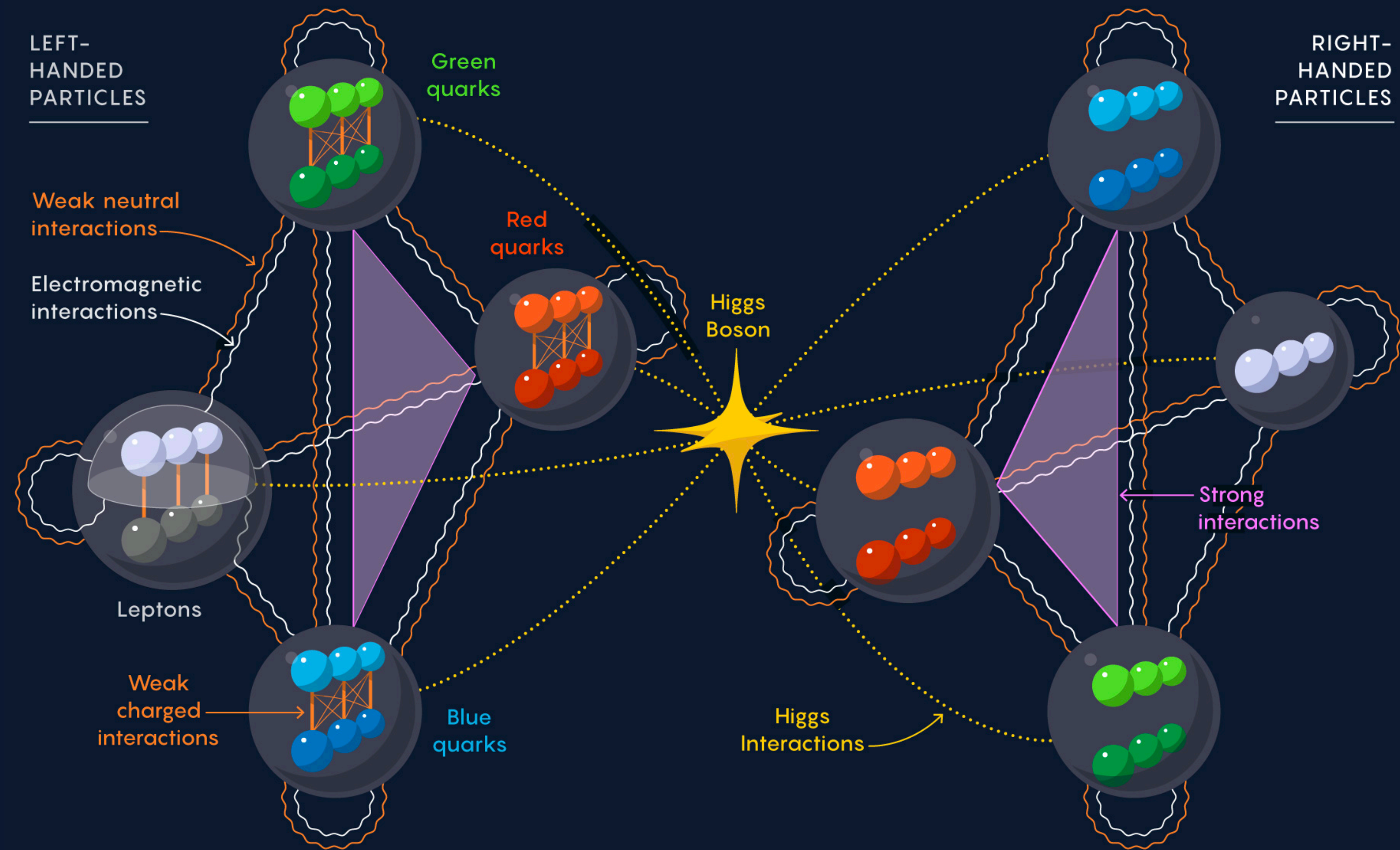


FERMIONS (MATTER) **BOSONS (FORCE CARRIERS)**
● QUARKS ● LEPTONS ● GAUGE BOSONS ● HIGGS BOSON

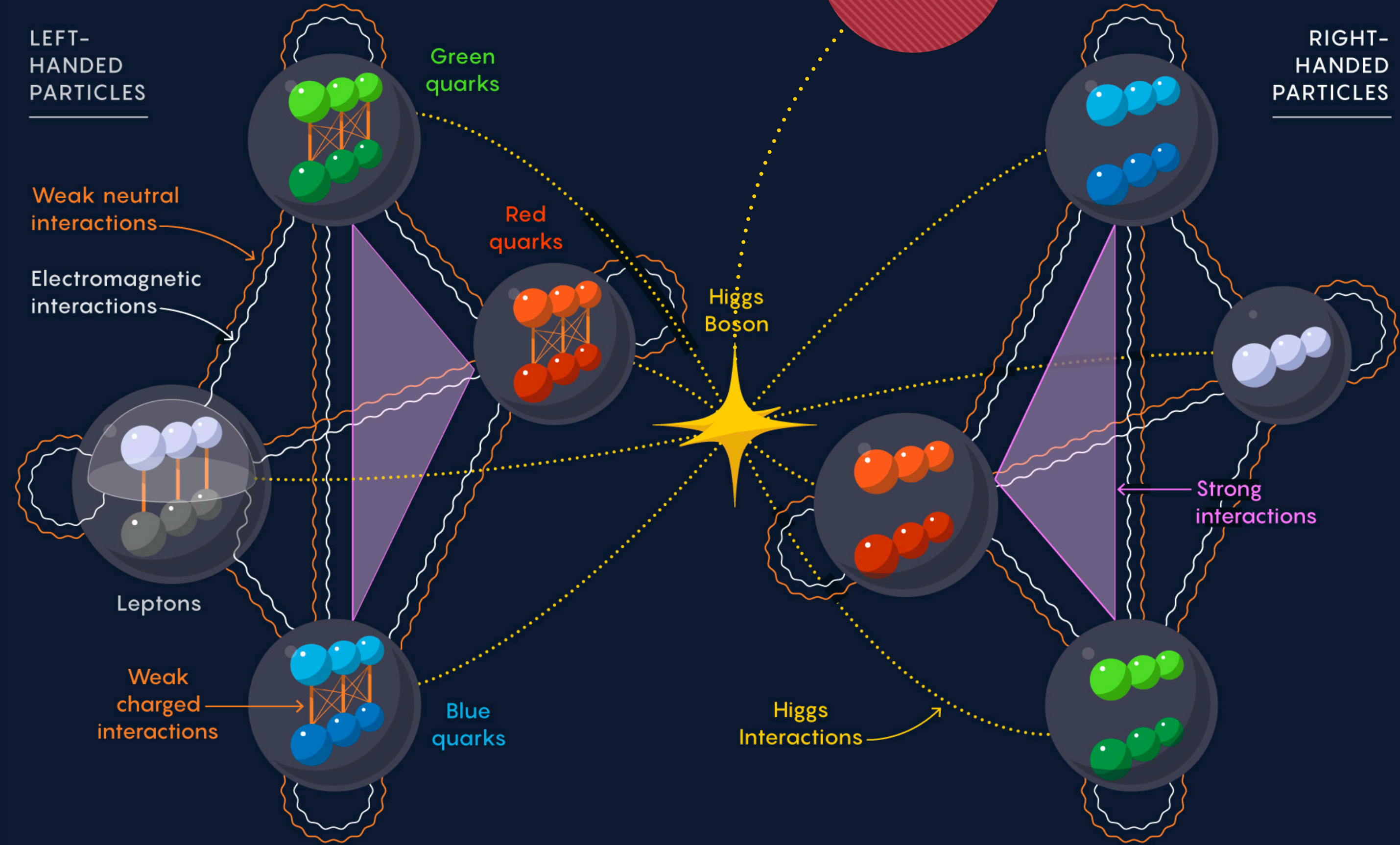
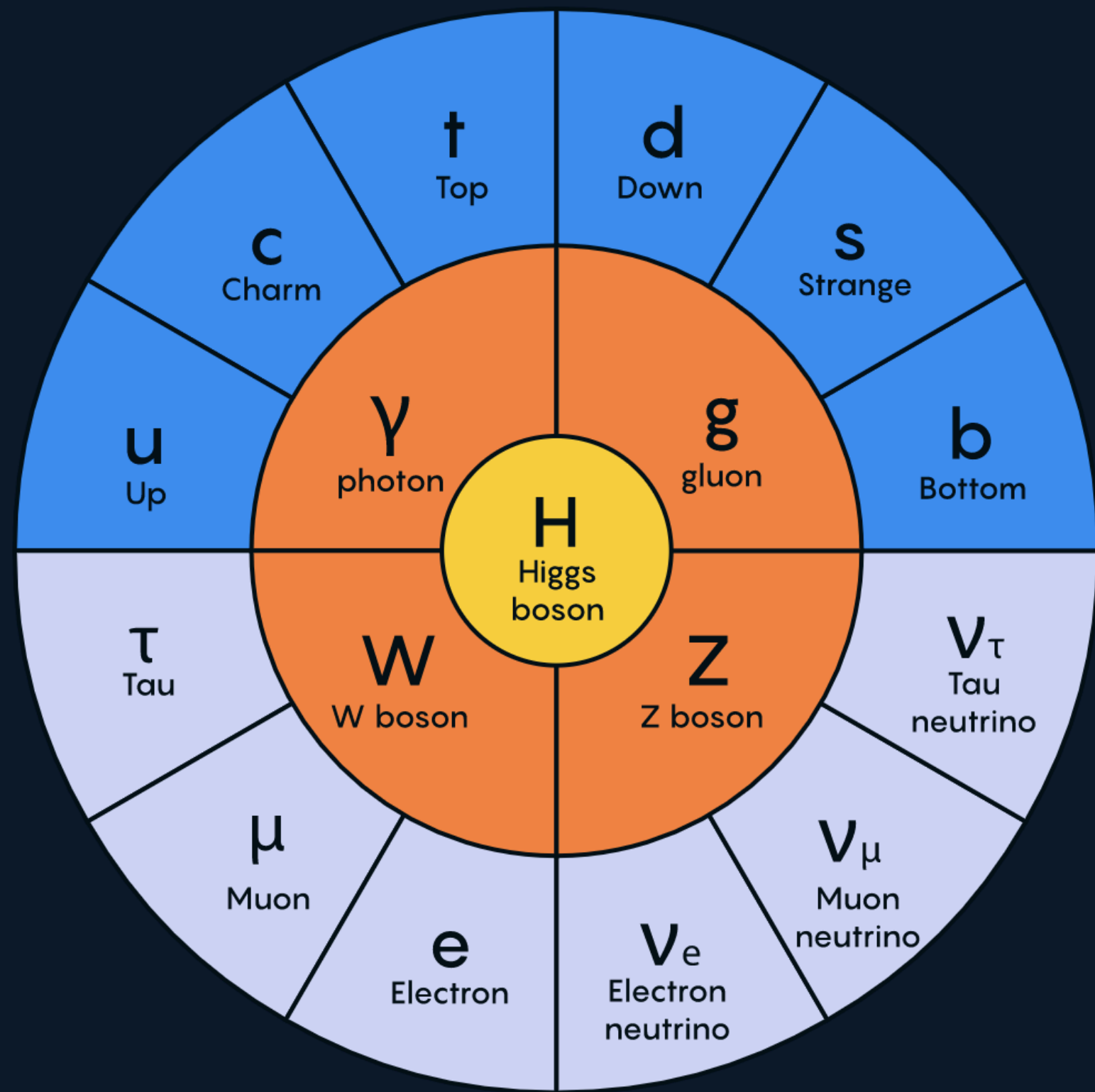
HIGGS BOSON IN THE STANDARD MODEL



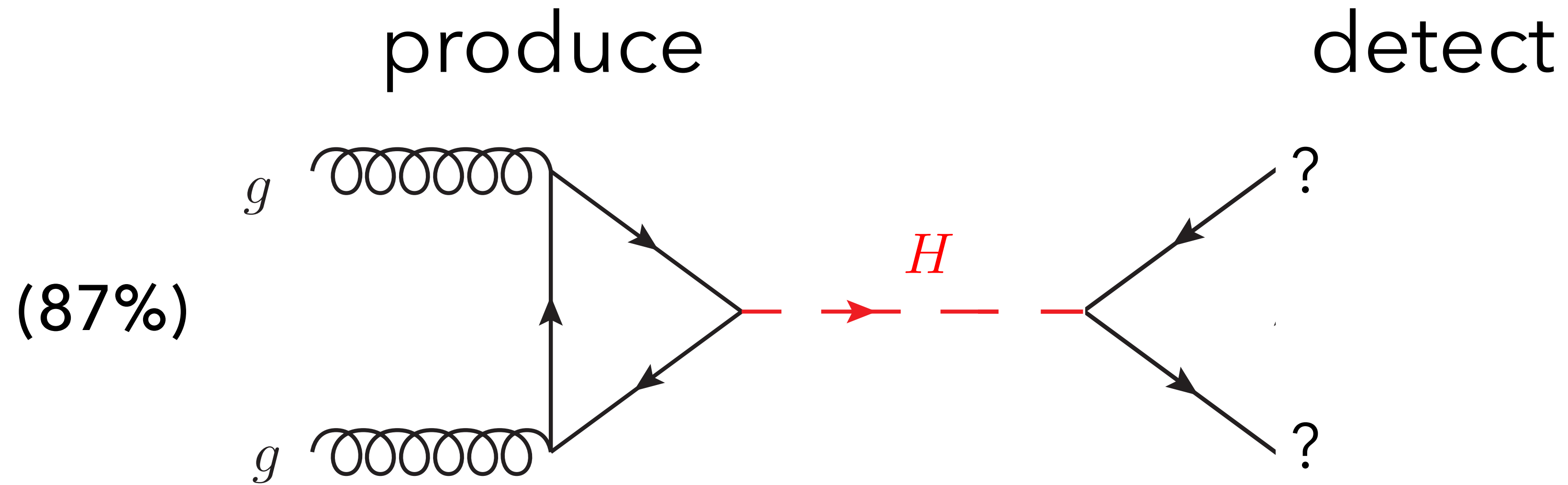
FERMIONS (MATTER) BOSONS (FORCE CARRIERS)
● QUARKS ● LEPTONS ● GAUGE BOSONS ● HIGGS BOSON



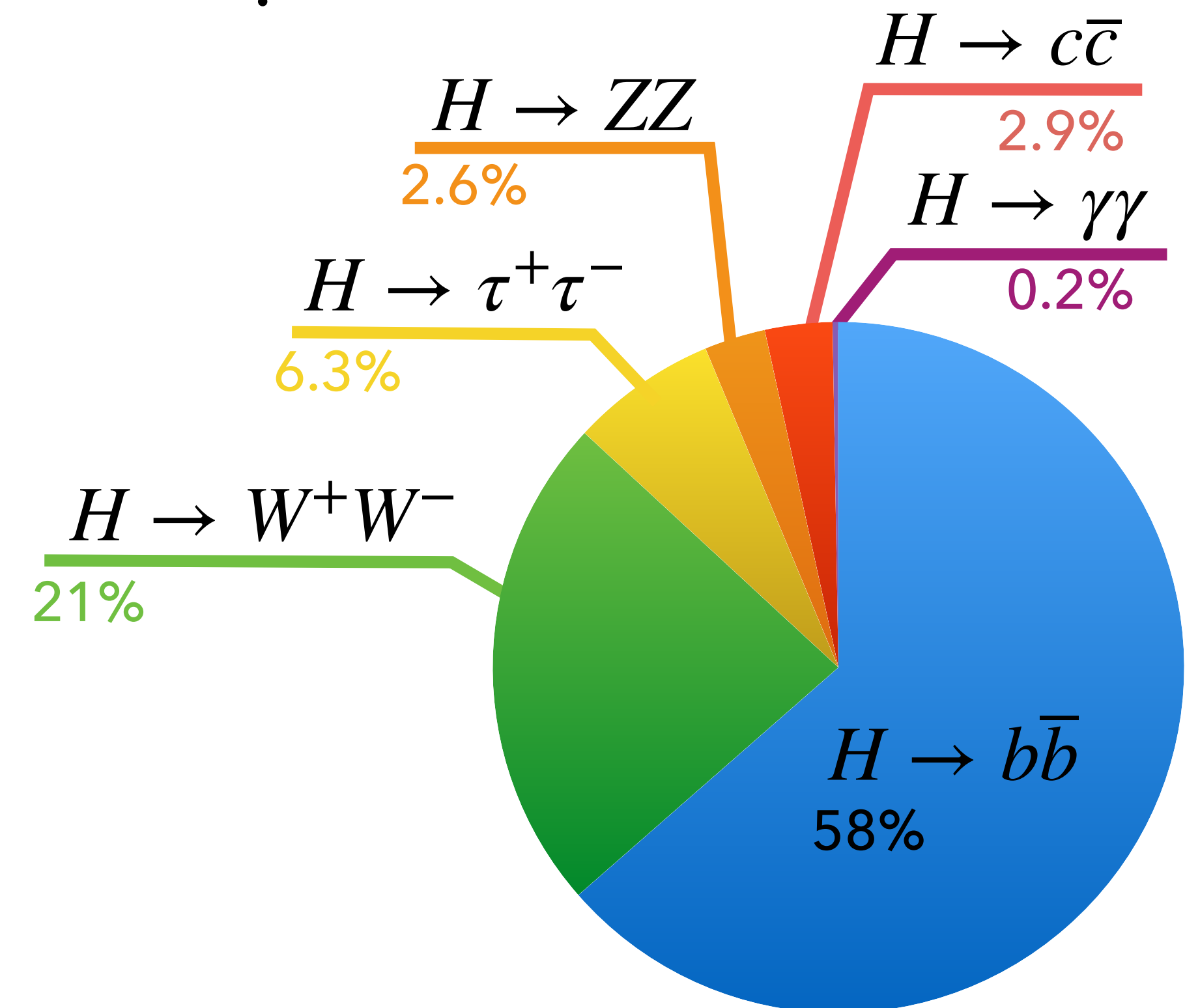
HIGGS BOSON IN THE STANDARD MODEL



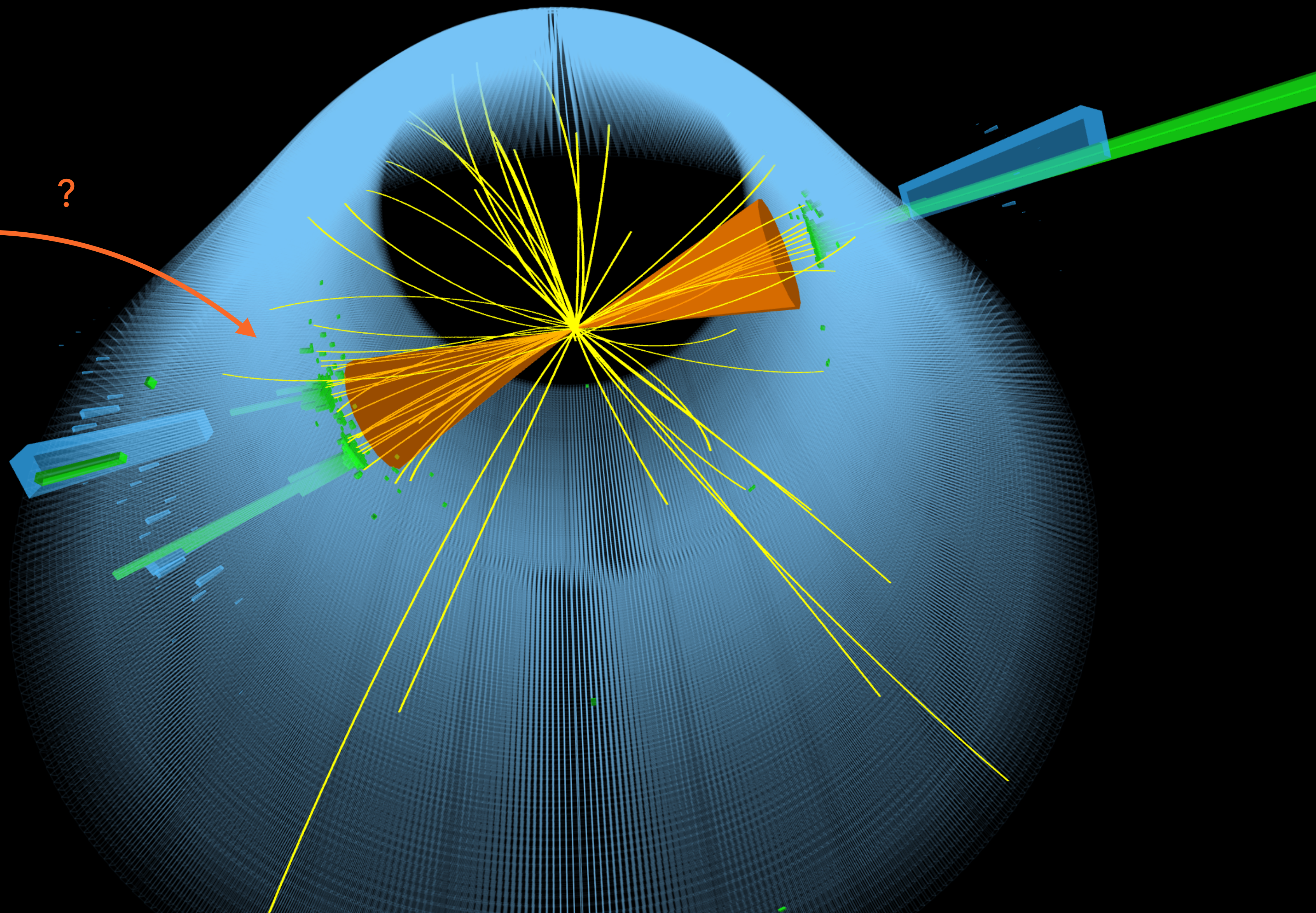
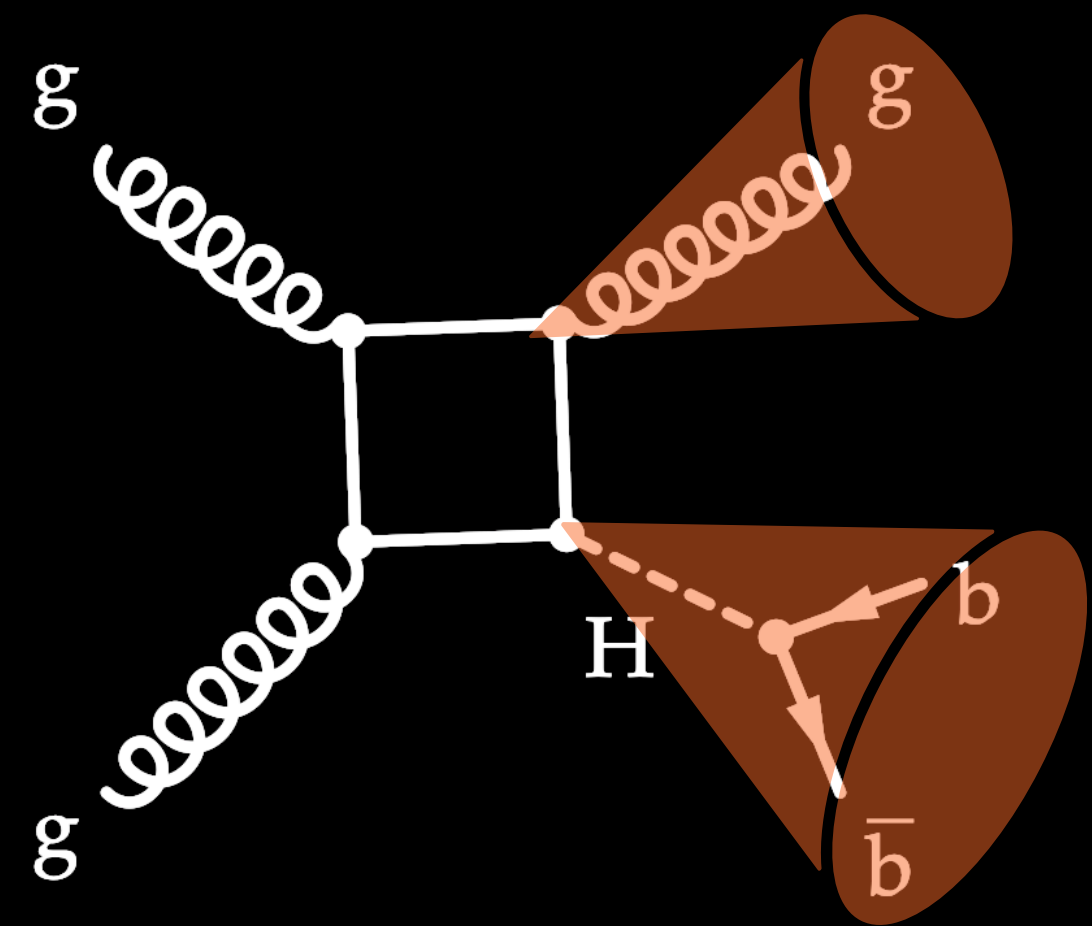
- ▶ Higgs boson is the **centerpiece**: all particles interact with it
- ▶ May be a link to new particles or interactions



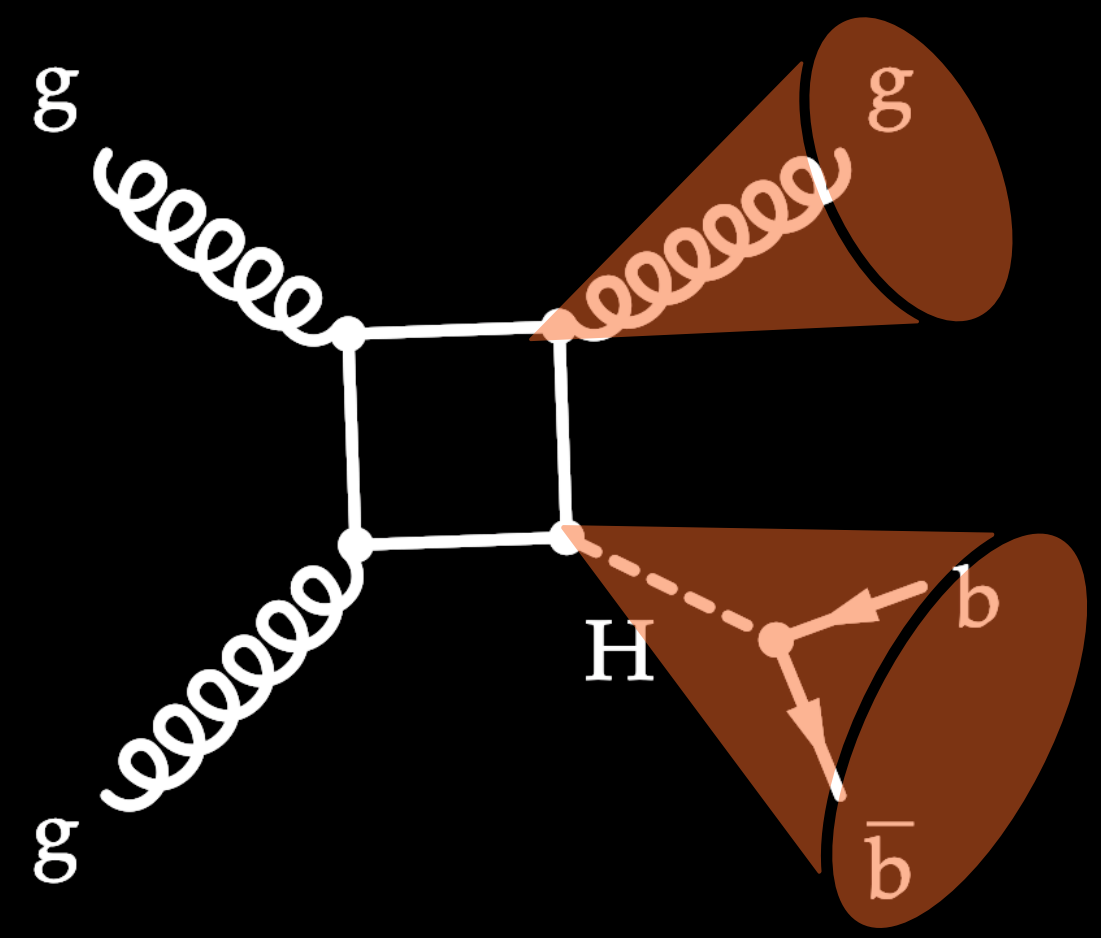
- ▶ $H \rightarrow \gamma\gamma$ and $H \rightarrow ZZ$ have small rate but are very "clean"
- ▶ $H \rightarrow b\bar{b}$ is large, but more difficult due to large backgrounds



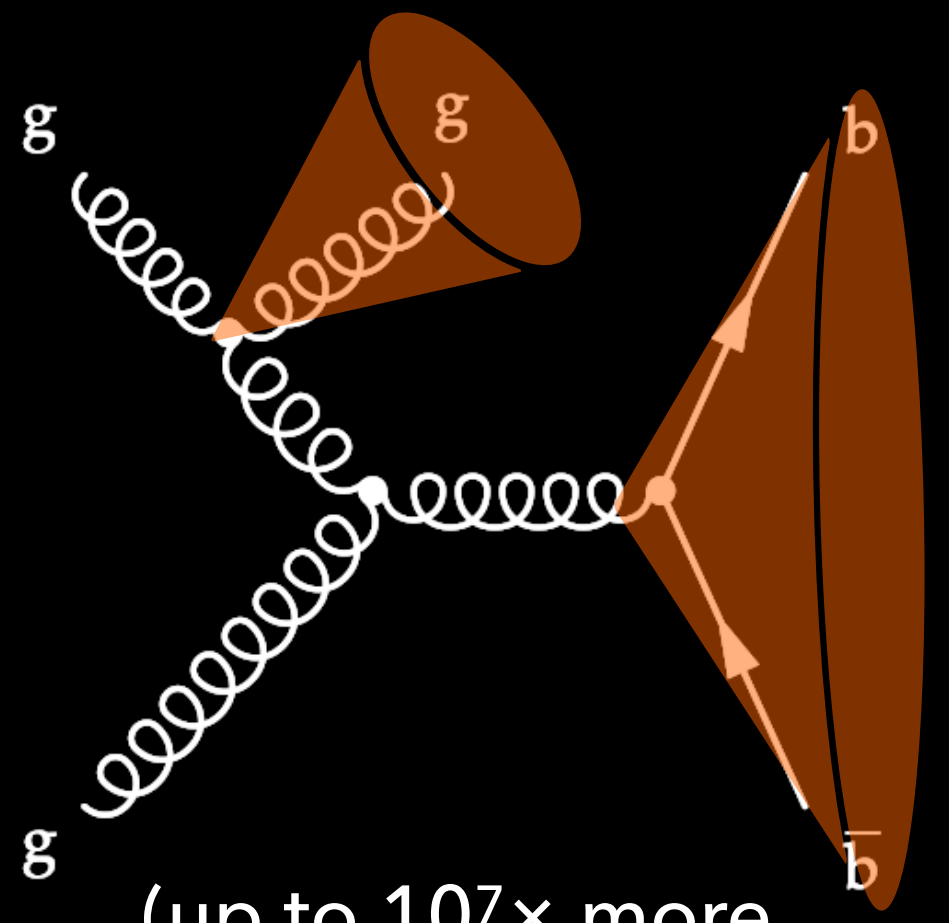
Signal:



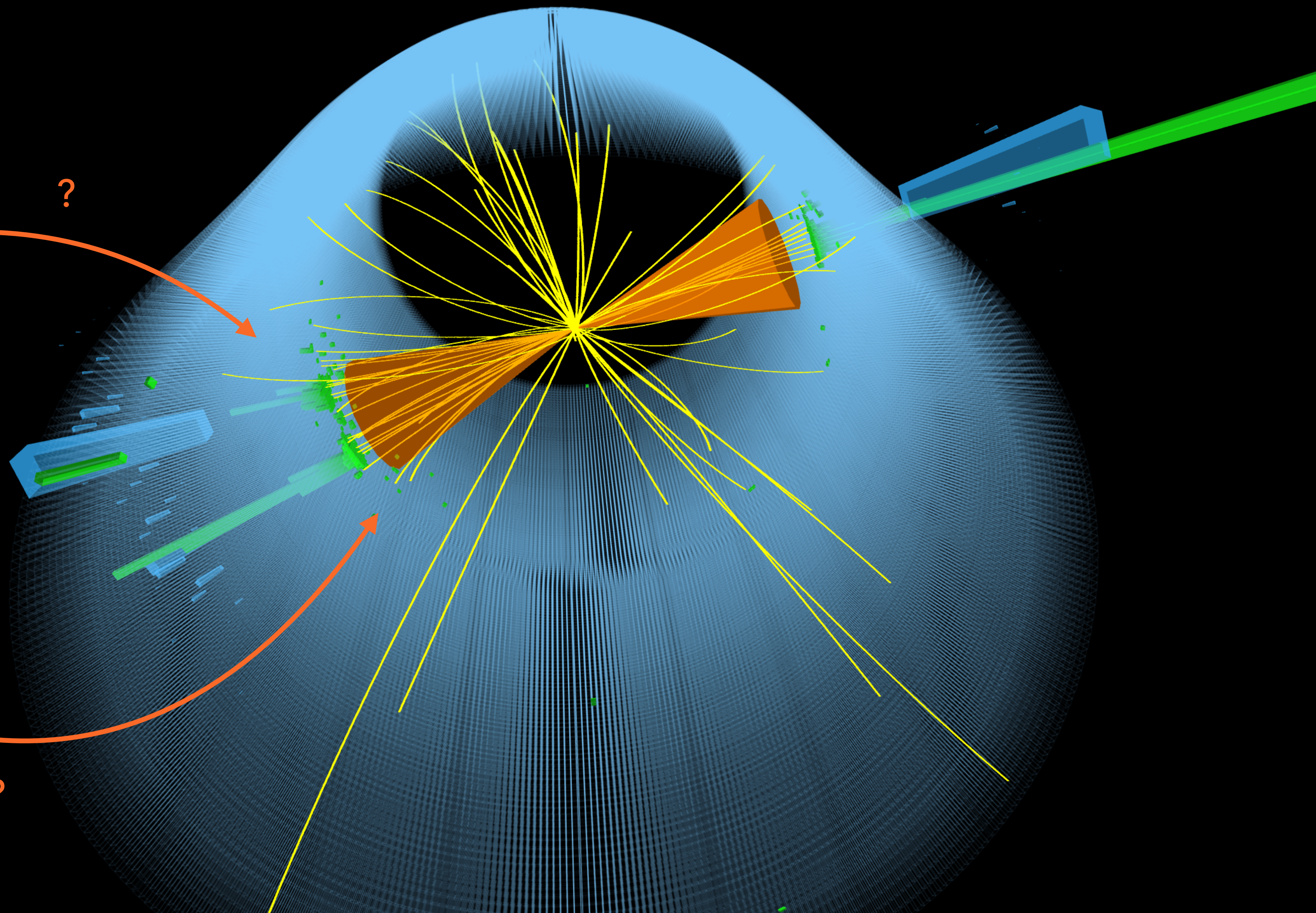
Signal:



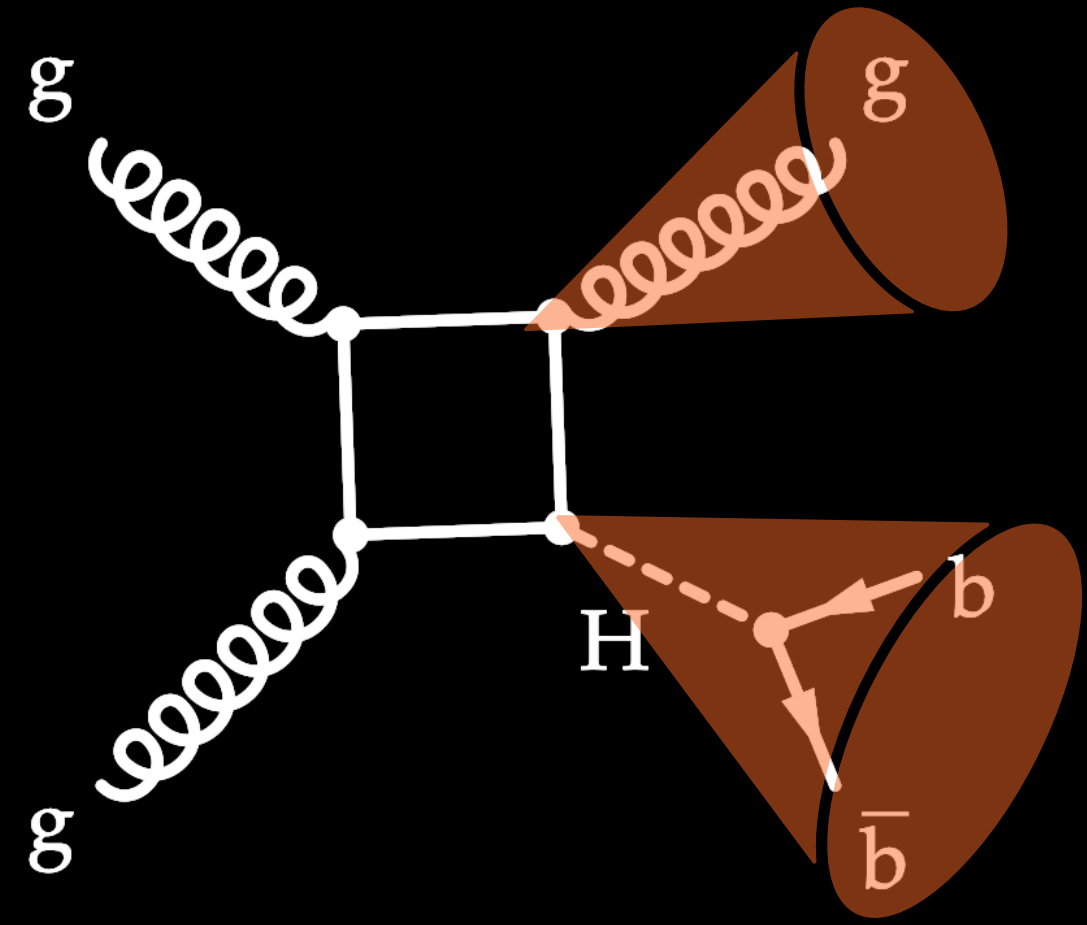
Background:



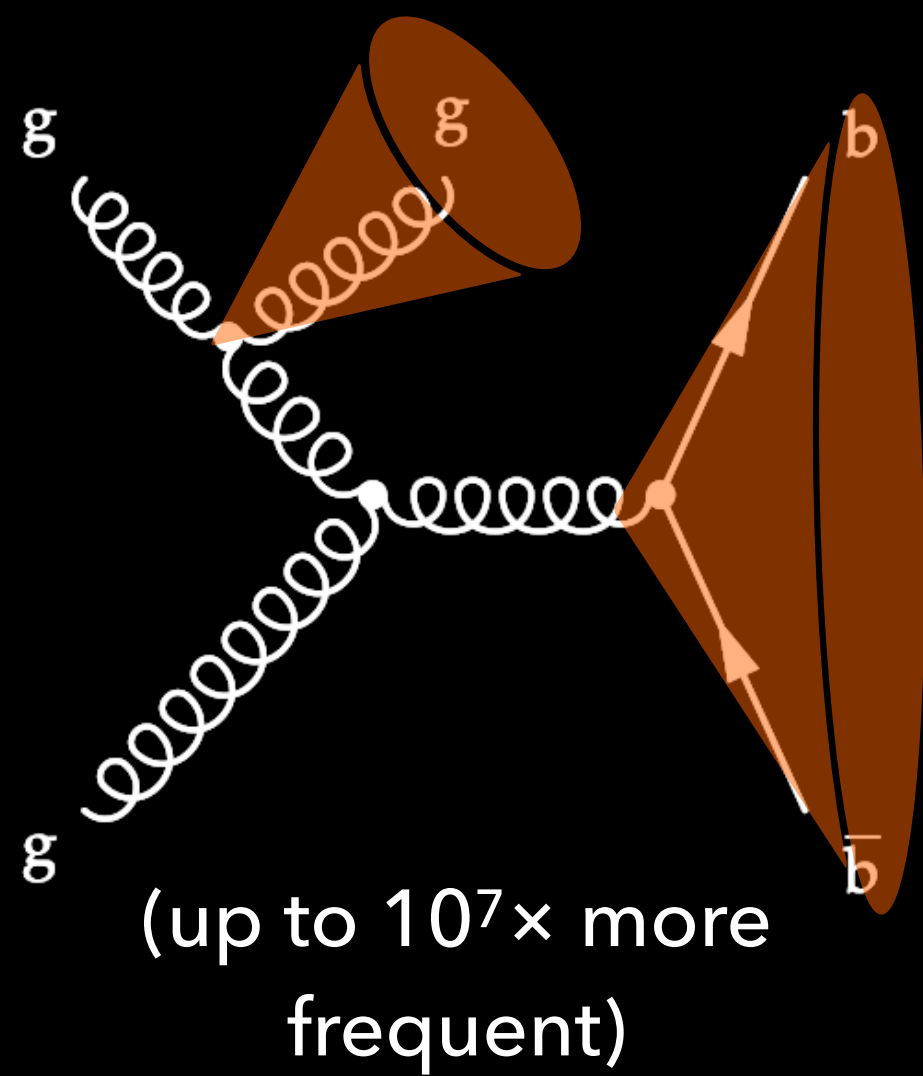
(up to $10^7 \times$ more frequent)



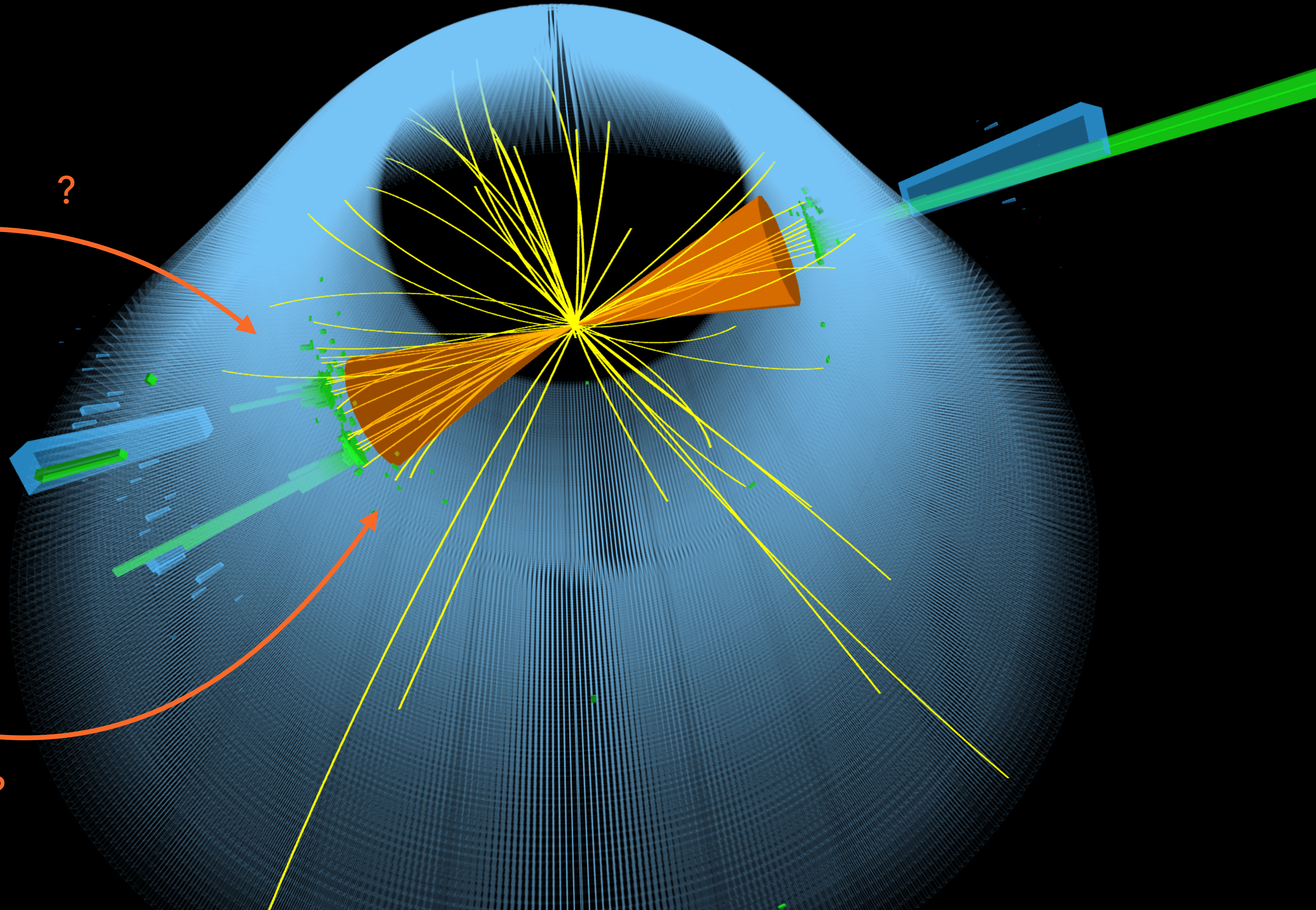
Signal:



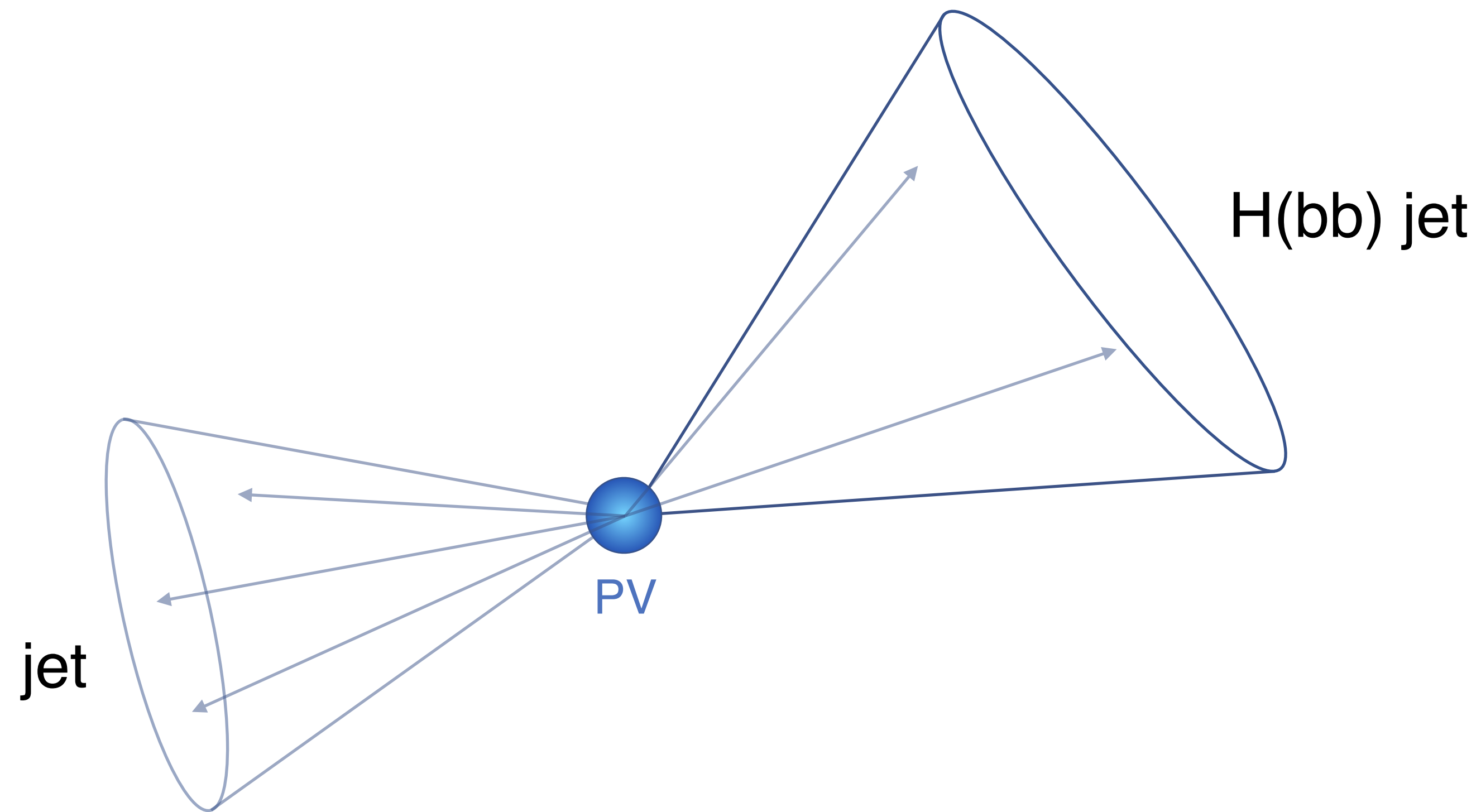
Background:



Can machine learning help us?

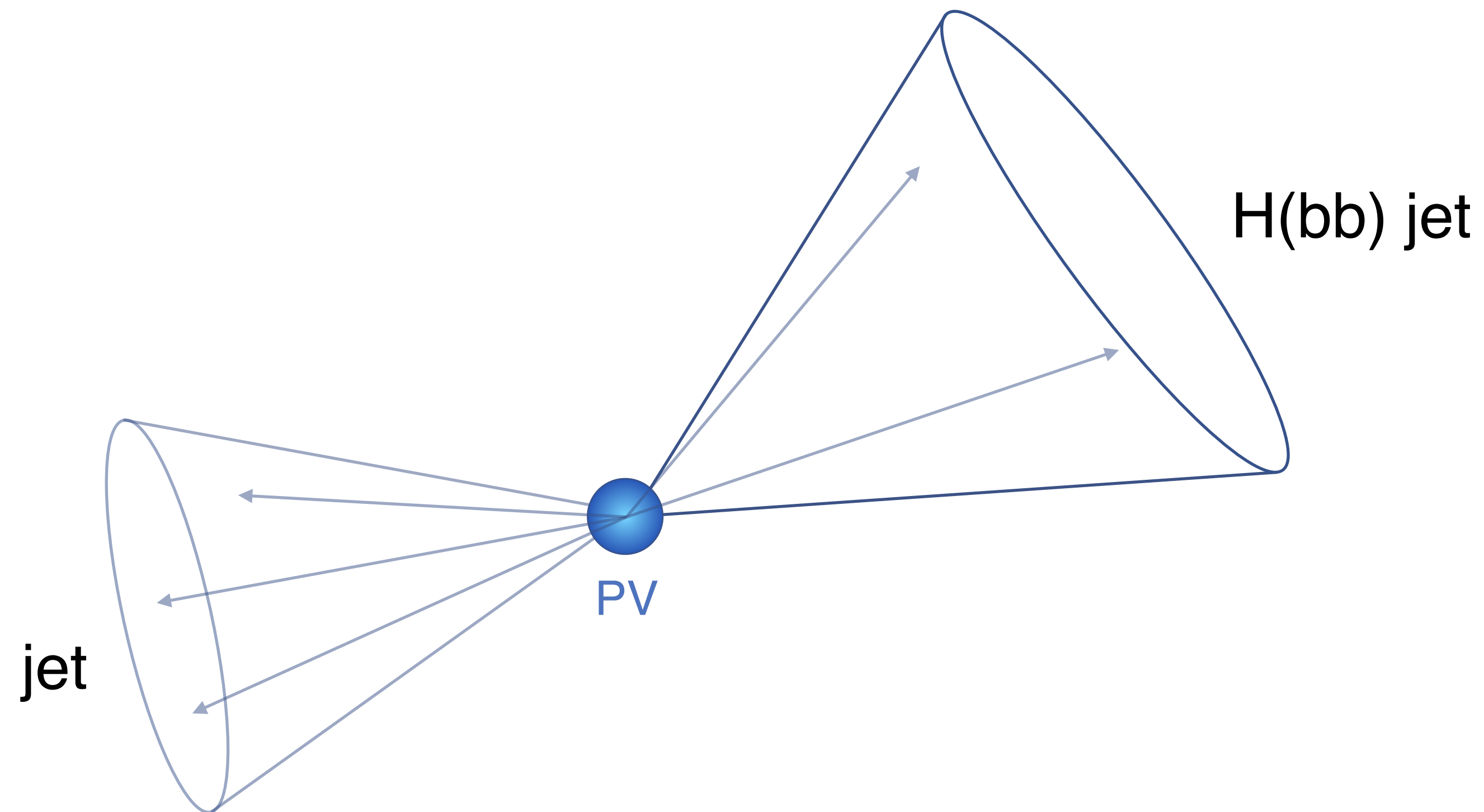


b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!



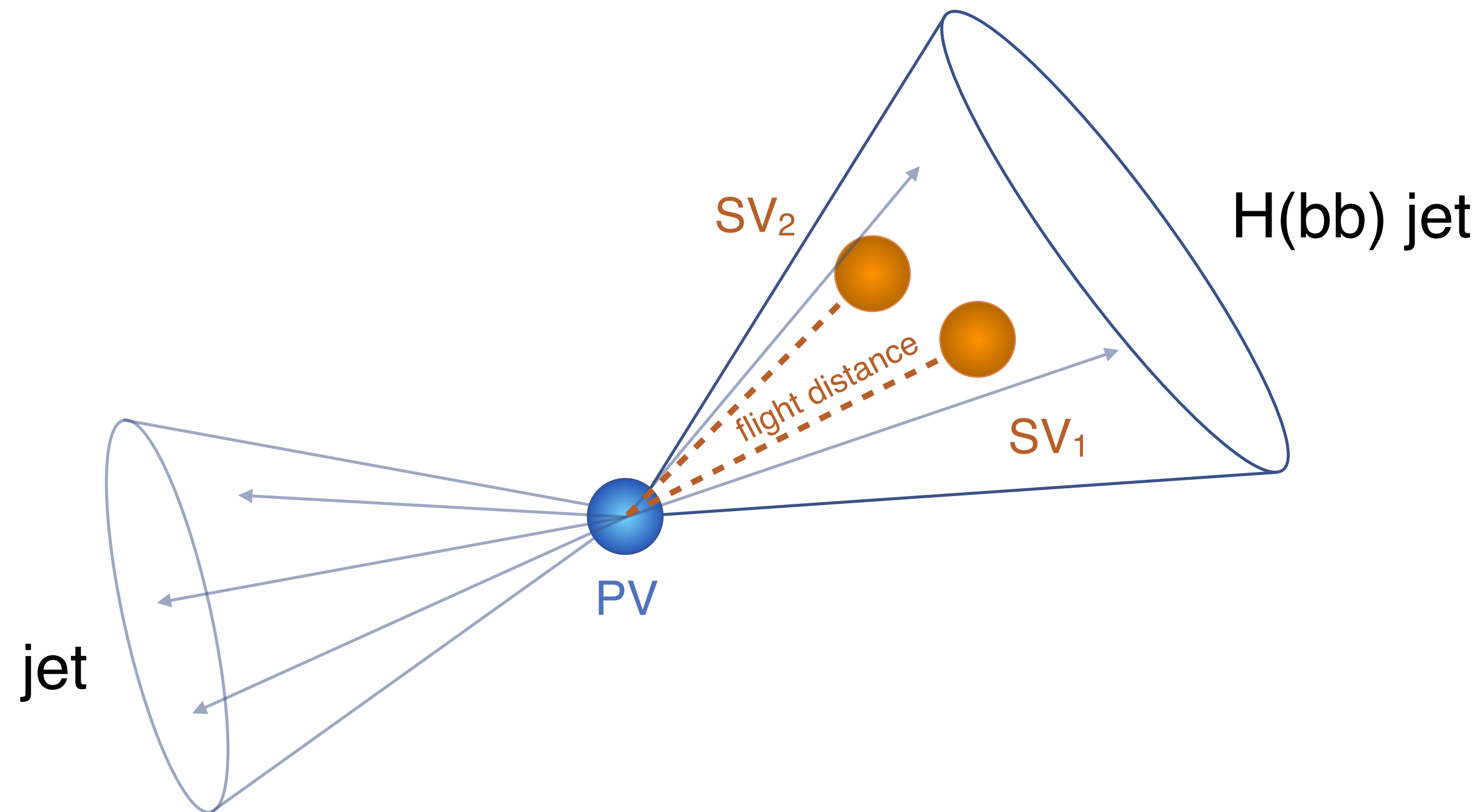
b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!

► Handles:



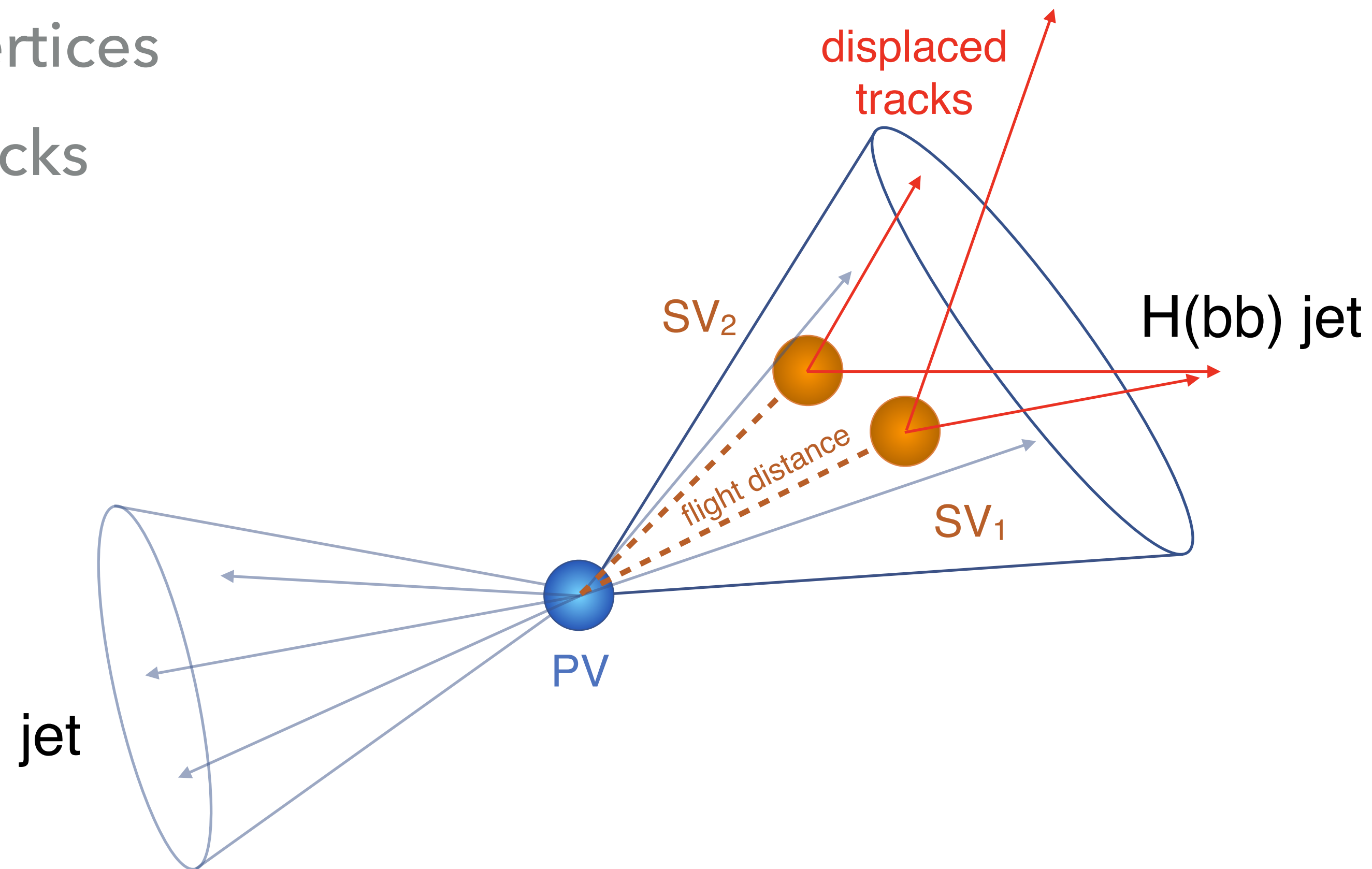
b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!

- ▶ Handles:
- ▶ secondary vertices



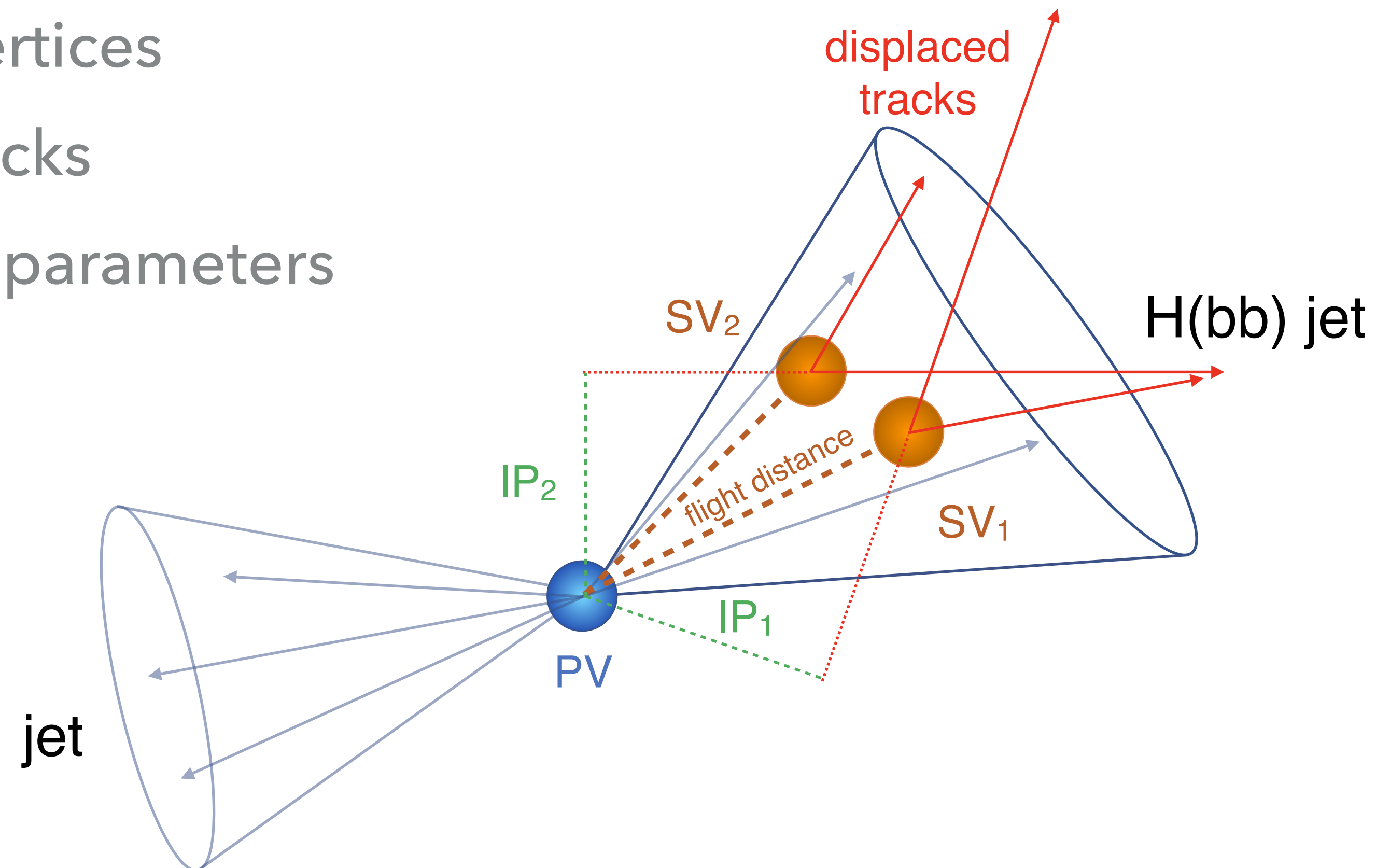
b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!

- ▶ Handles:
- ▶ secondary vertices
- ▶ displaced tracks



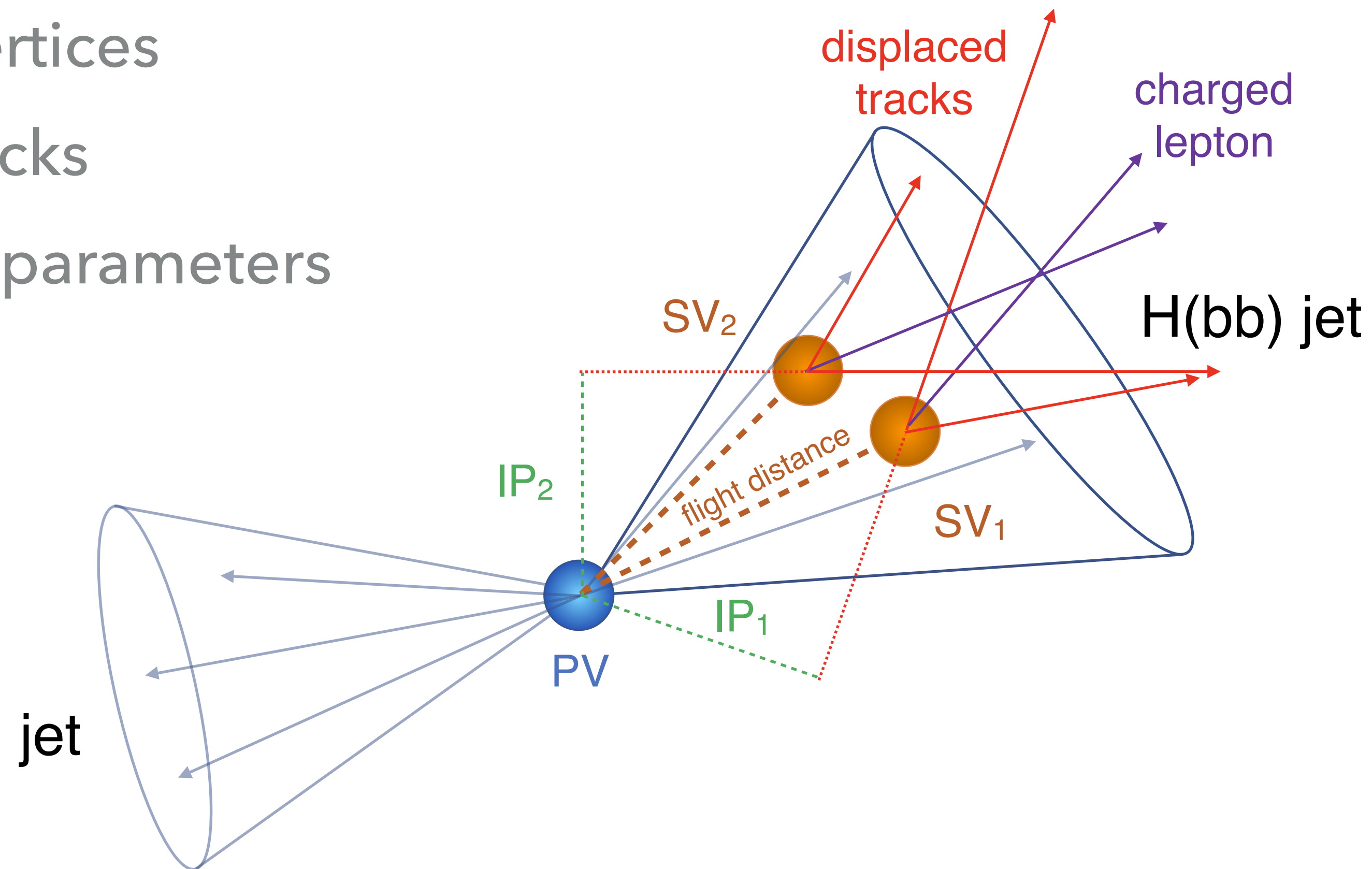
b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!

- ▶ Handles:
 - ▶ secondary vertices
 - ▶ displaced tracks
 - ▶ large impact parameters



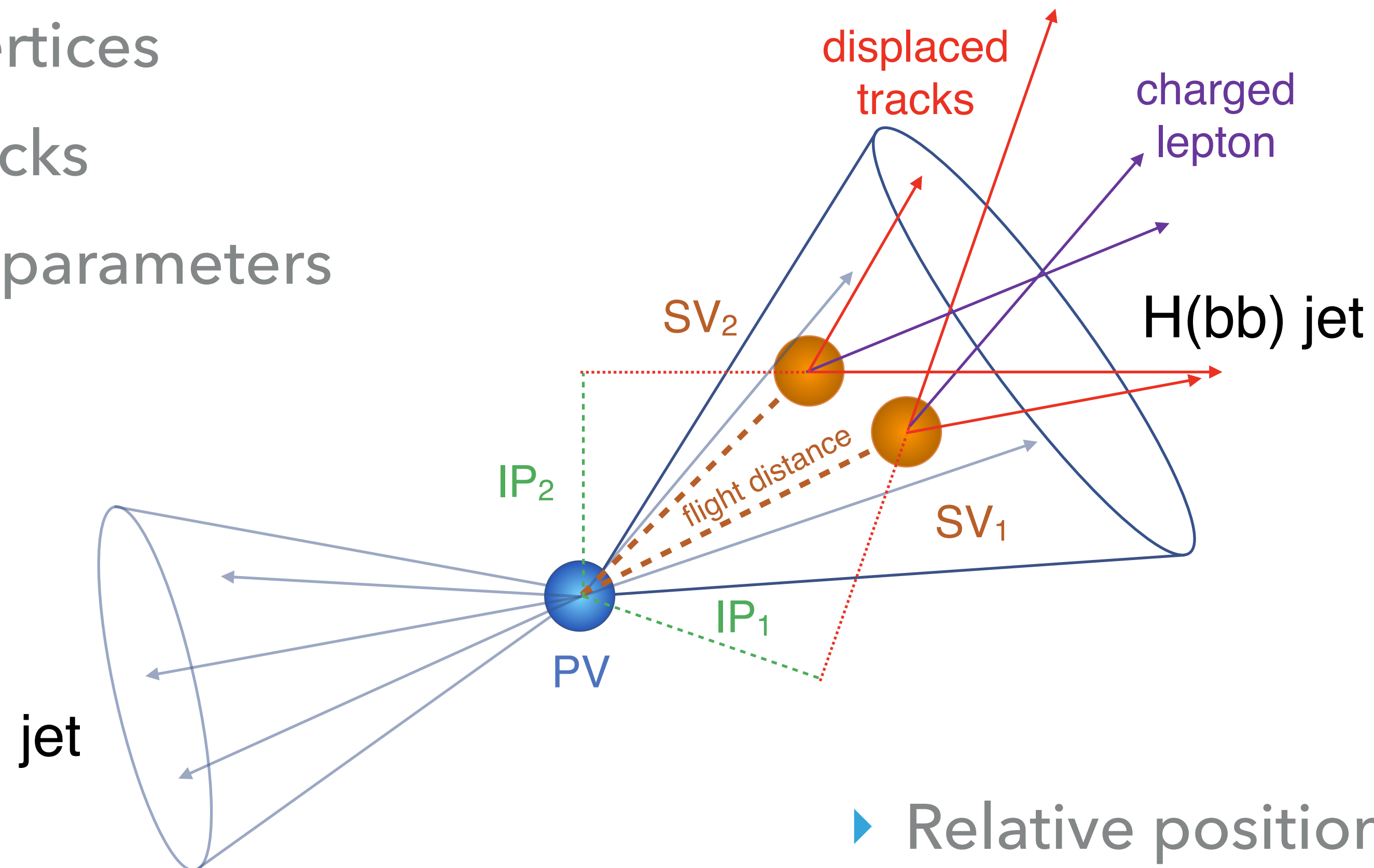
b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!

- ▶ Handles:
 - ▶ secondary vertices
 - ▶ displaced tracks
 - ▶ large impact parameters
 - ▶ soft leptons



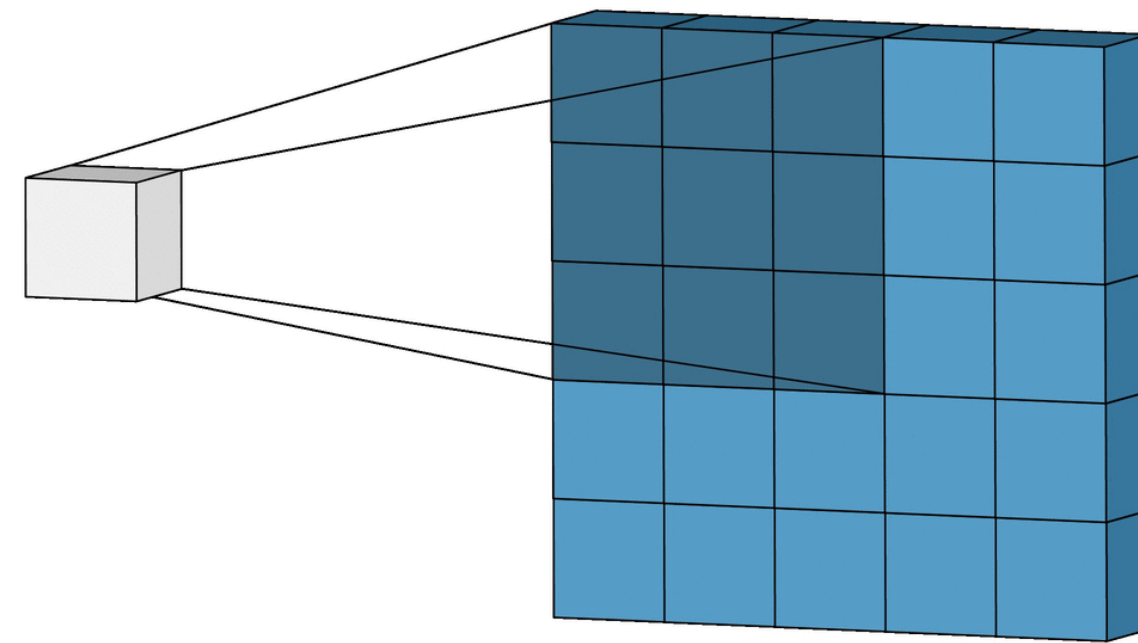
b hadrons have long lifetimes:
travel $O(\text{mm})$ before decay!

- ▶ Handles:
 - ▶ secondary vertices
 - ▶ displaced tracks
 - ▶ large impact parameters
 - ▶ soft leptons



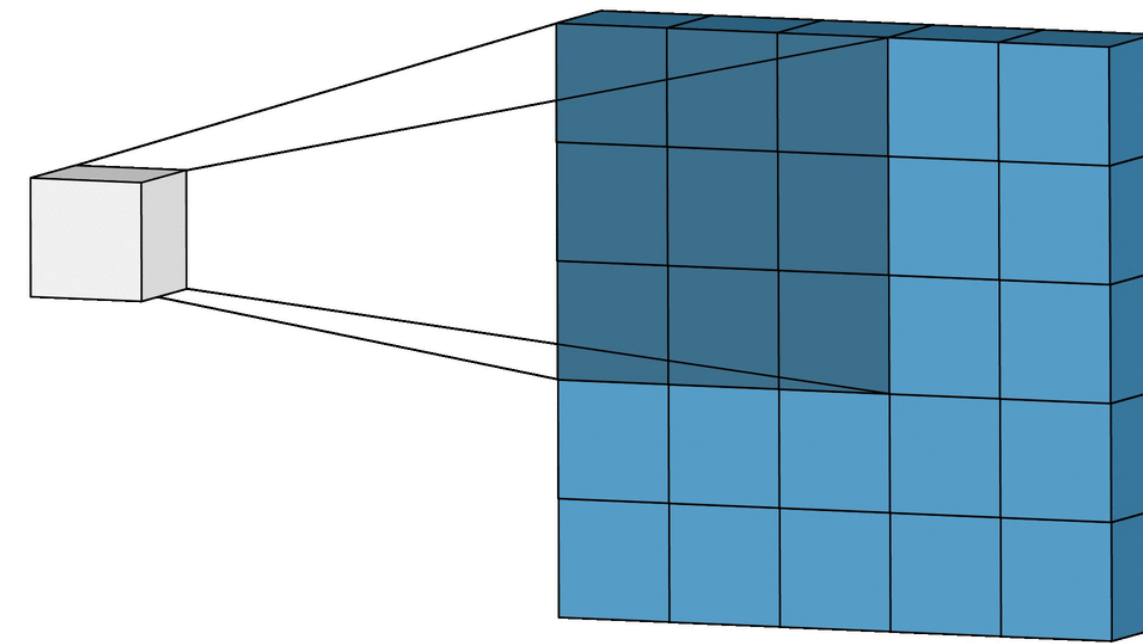
- ▶ In deep learning, tailoring algorithms to the structure (and symmetries) of the data has led to groundbreaking performance

- ▶ In deep learning, tailoring algorithms to the structure (and symmetries) of the data has led to groundbreaking performance
 - ▶ CNNs for images

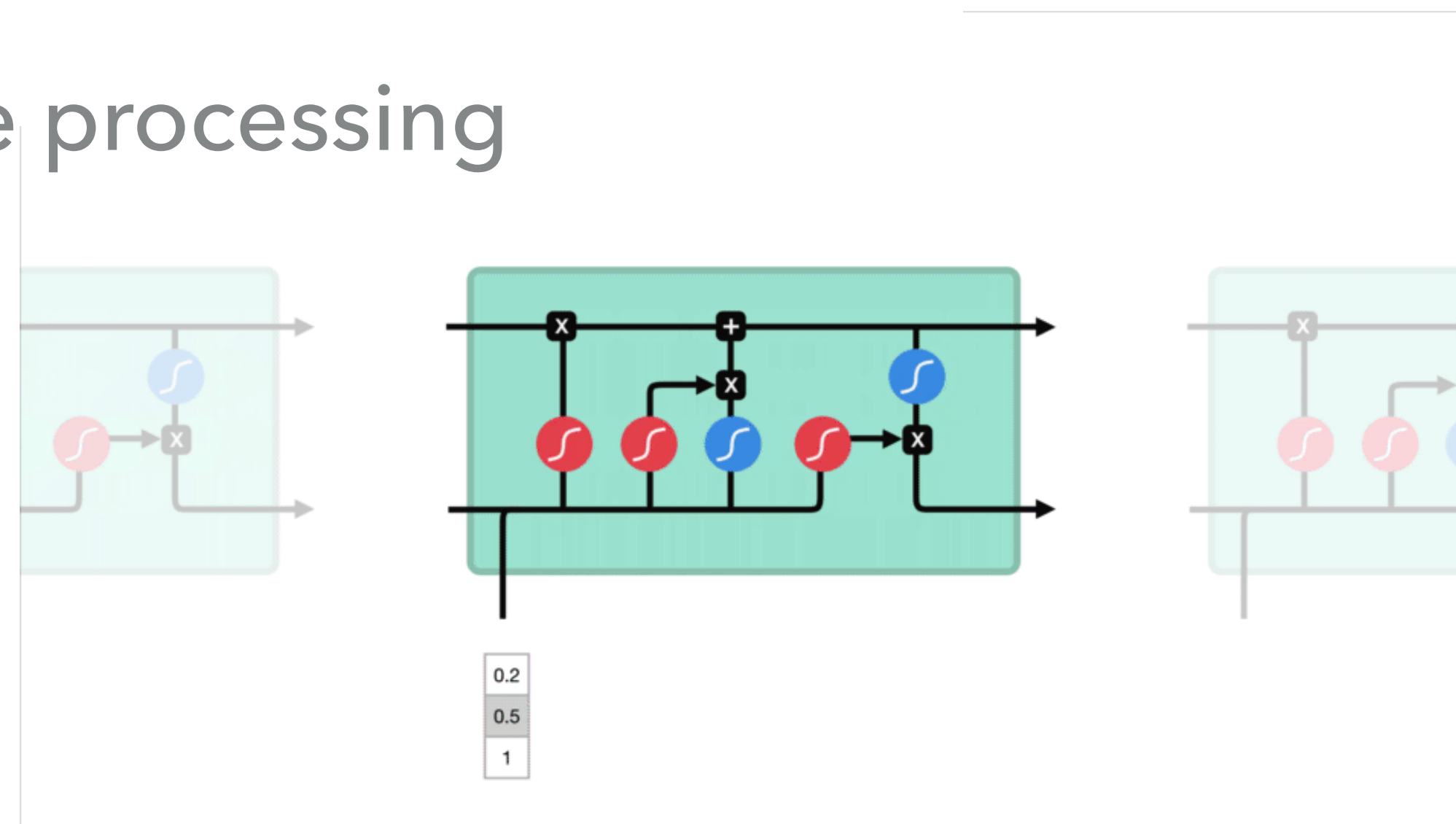


- ▶ In deep learning, tailoring algorithms to the structure (and symmetries) of the data has led to groundbreaking performance

- ▶ CNNs for images

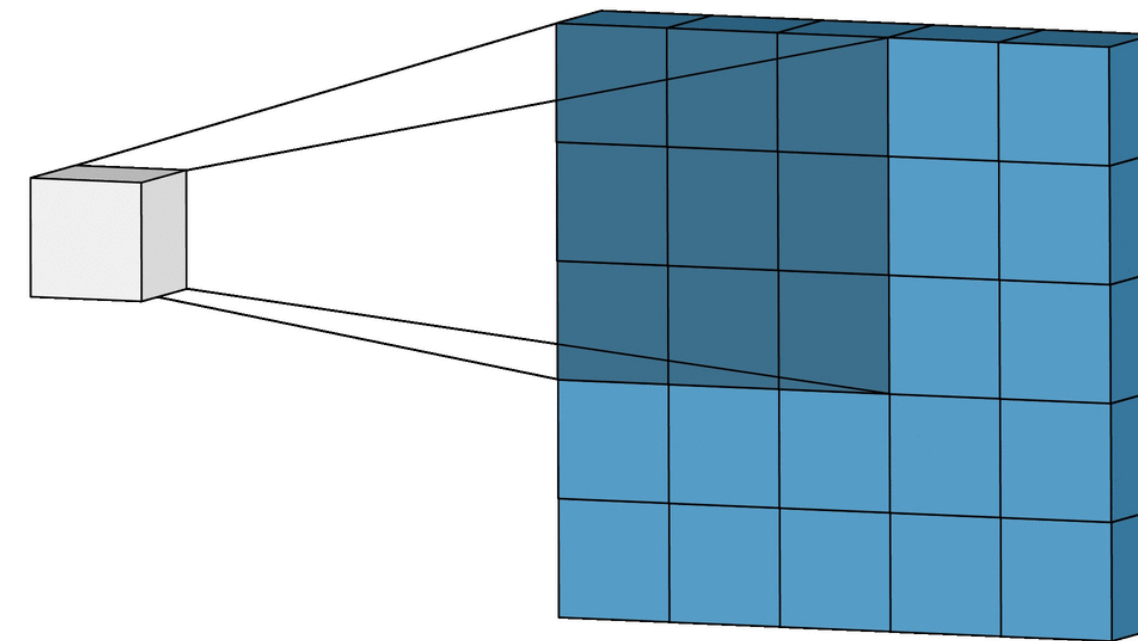


- ▶ RNNs for language processing

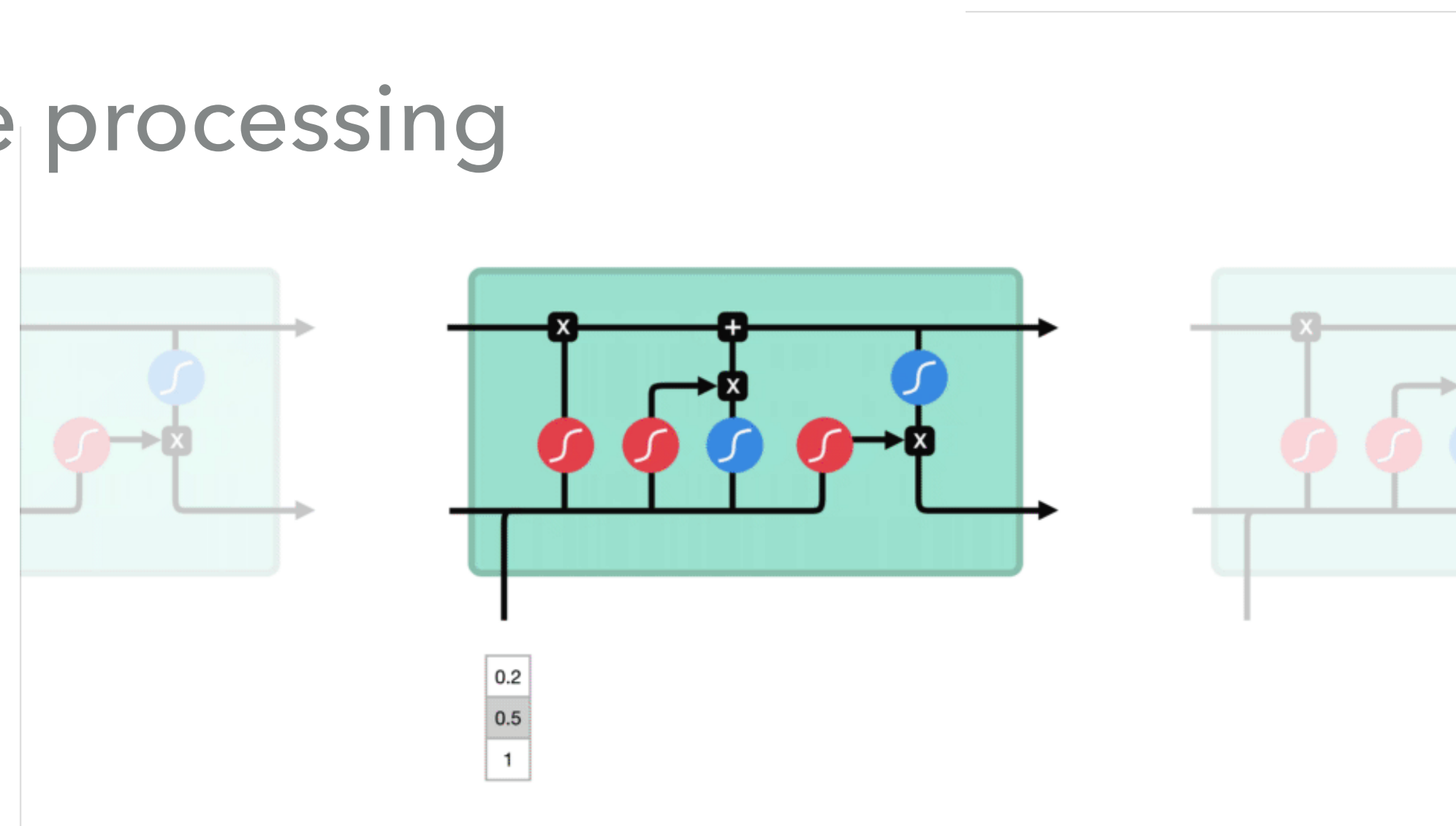


- ▶ In deep learning, tailoring algorithms to the structure (and symmetries) of the data has led to groundbreaking performance

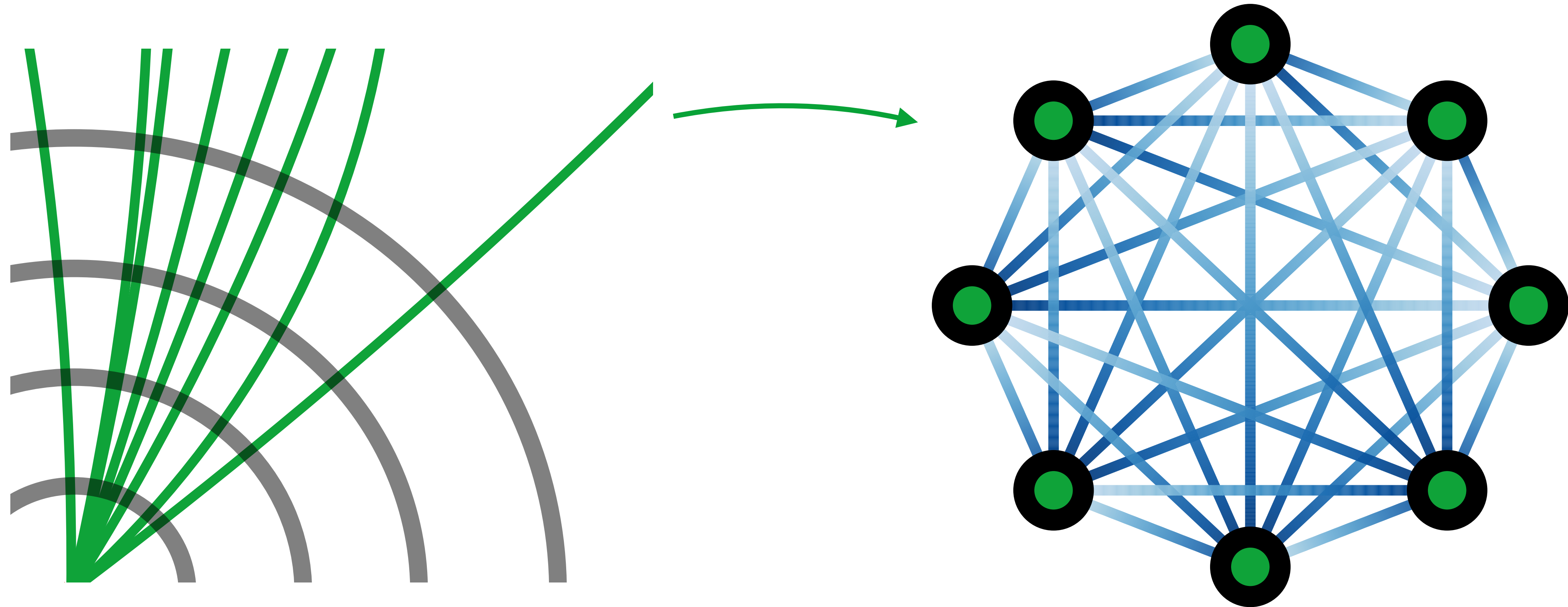
- ▶ CNNs for images



- ▶ RNNs for language processing

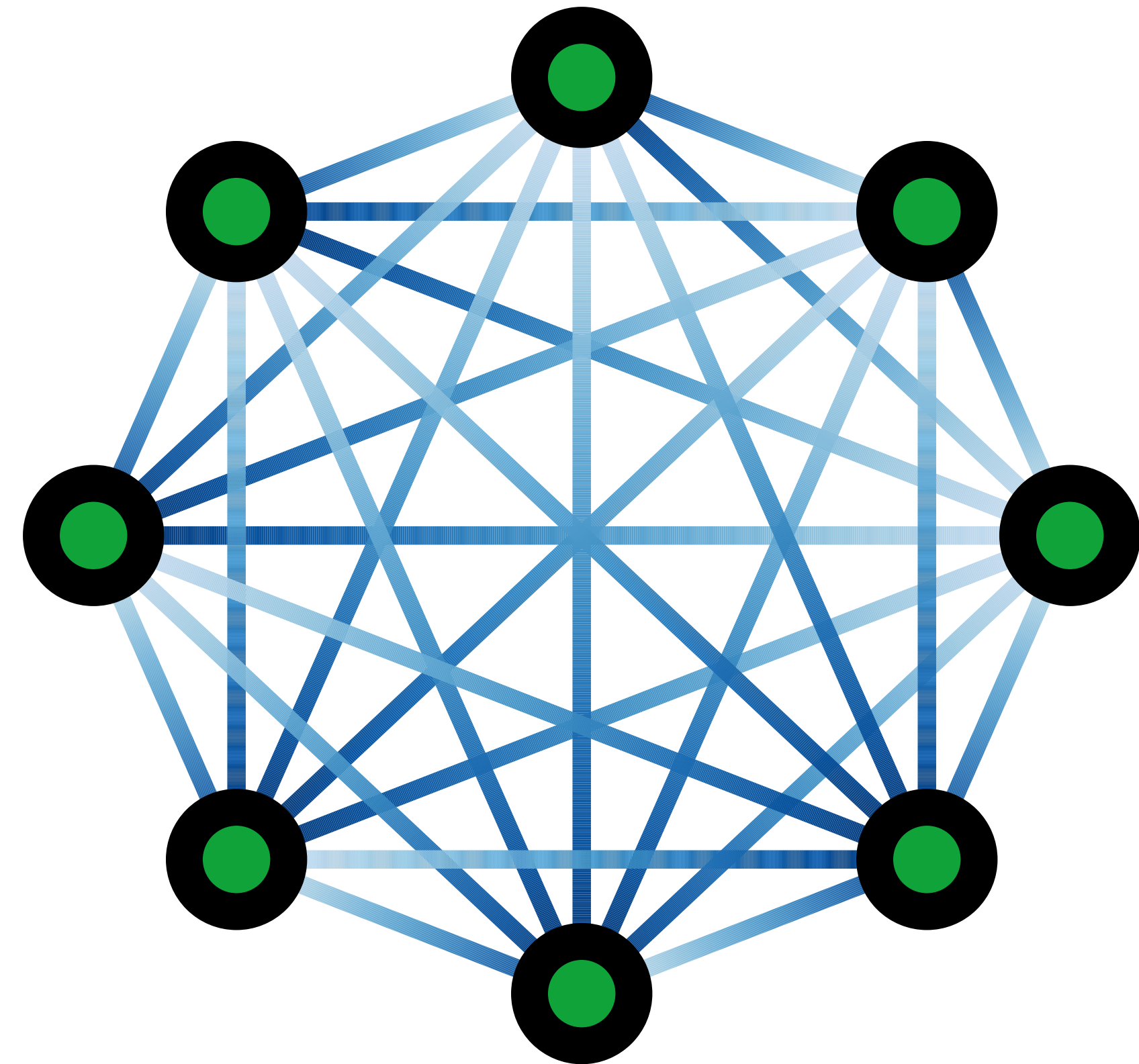


- ▶ What about high energy physics data like jets?



- ▶ Node features \mathbf{v}_i : particle 4-momentum

$$p = [E, p_x, p_y, p_z] \equiv [p_T, \eta, \phi, m]$$

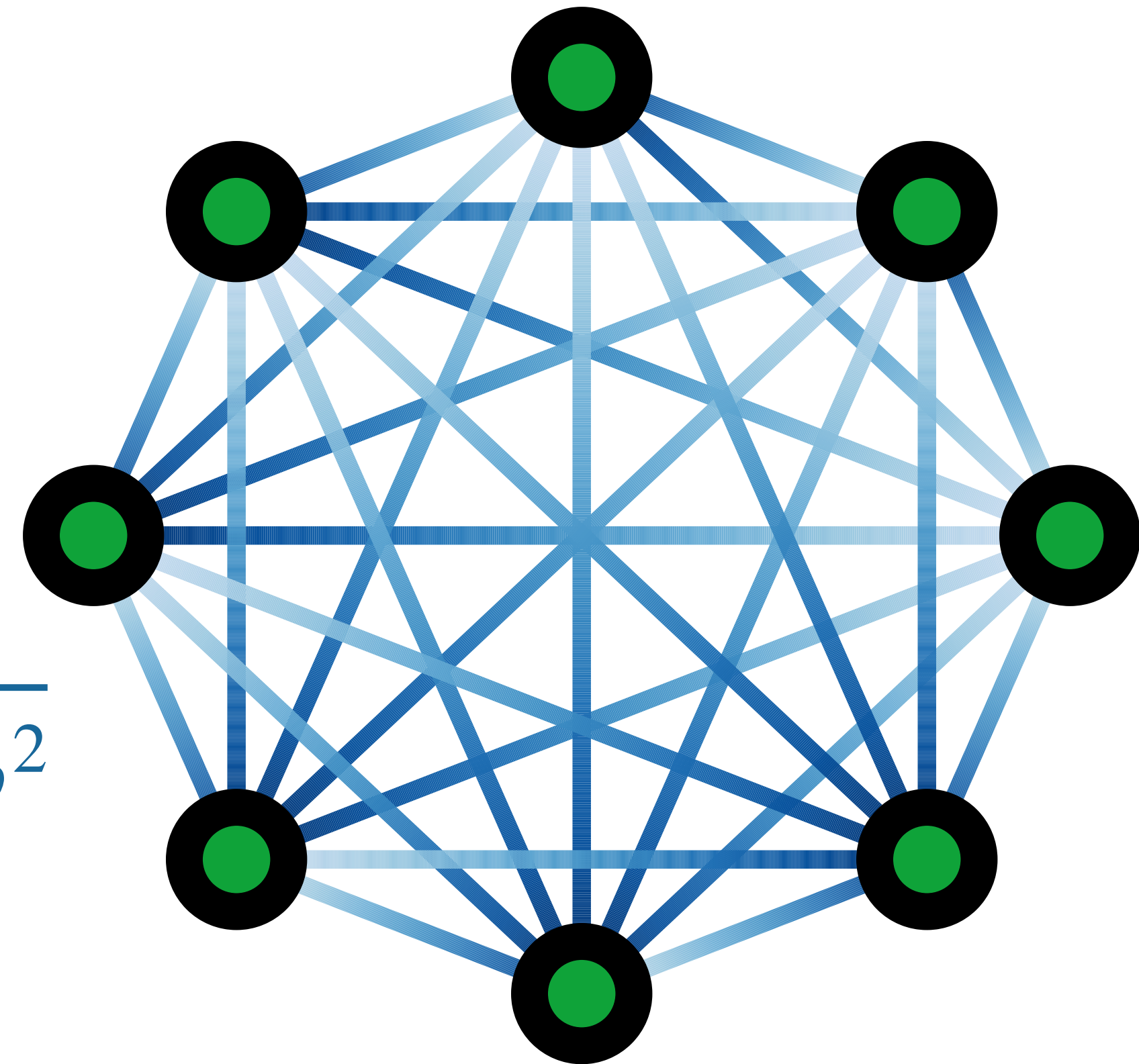


- ▶ Node features \mathbf{v}_i : particle 4-momentum

$$p = [E, p_x, p_y, p_z] \equiv [p_T, \eta, \phi, m]$$

- ▶ Edge features \mathbf{e}_k : pseudoangular distance between particles

$$\Delta R = \sqrt{\Delta\eta^2 + \Delta\phi^2}$$

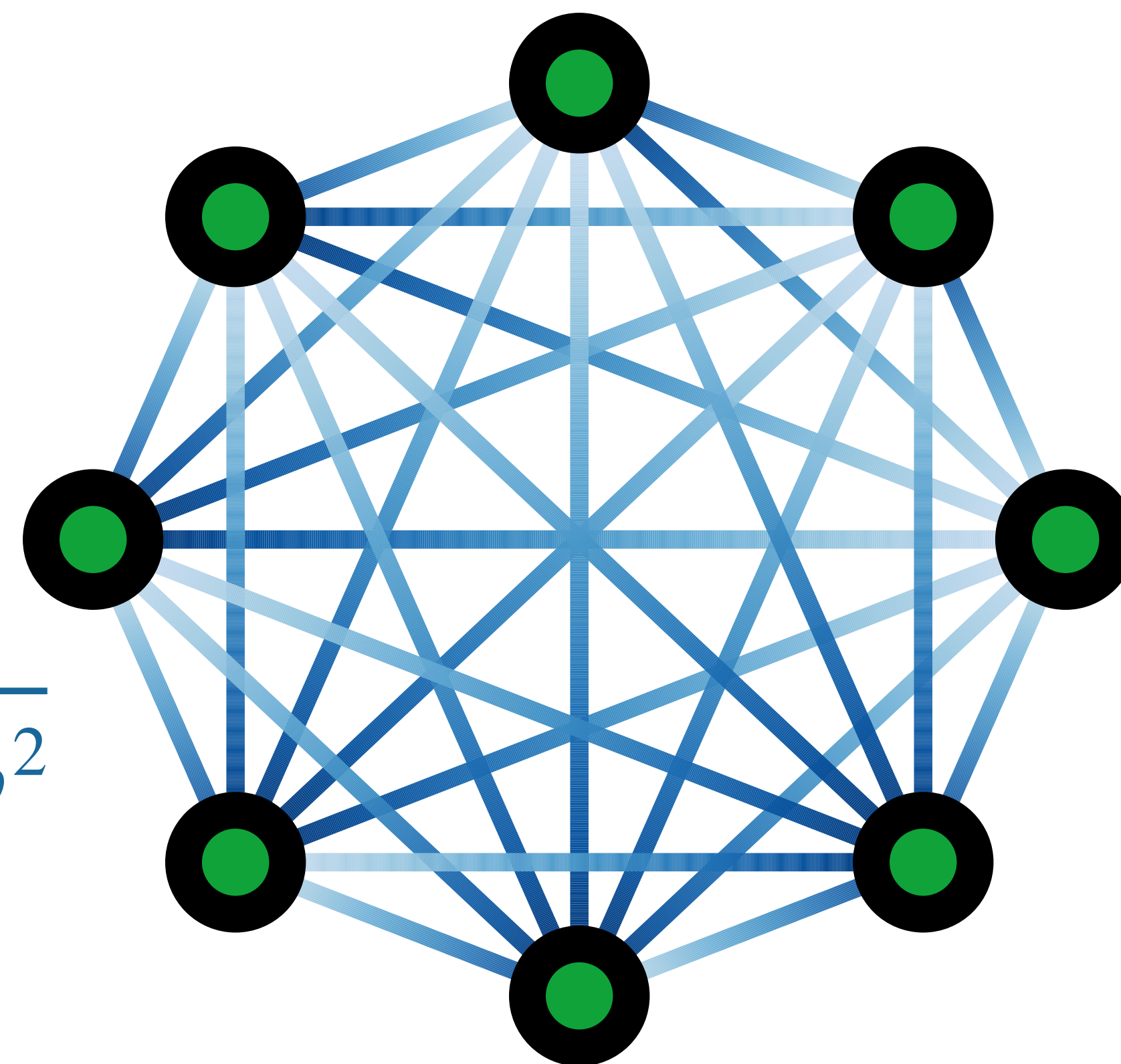


- ▶ Node features \mathbf{v}_i : particle 4-momentum

$$p = [E, p_x, p_y, p_z] \equiv [p_T, \eta, \phi, m]$$

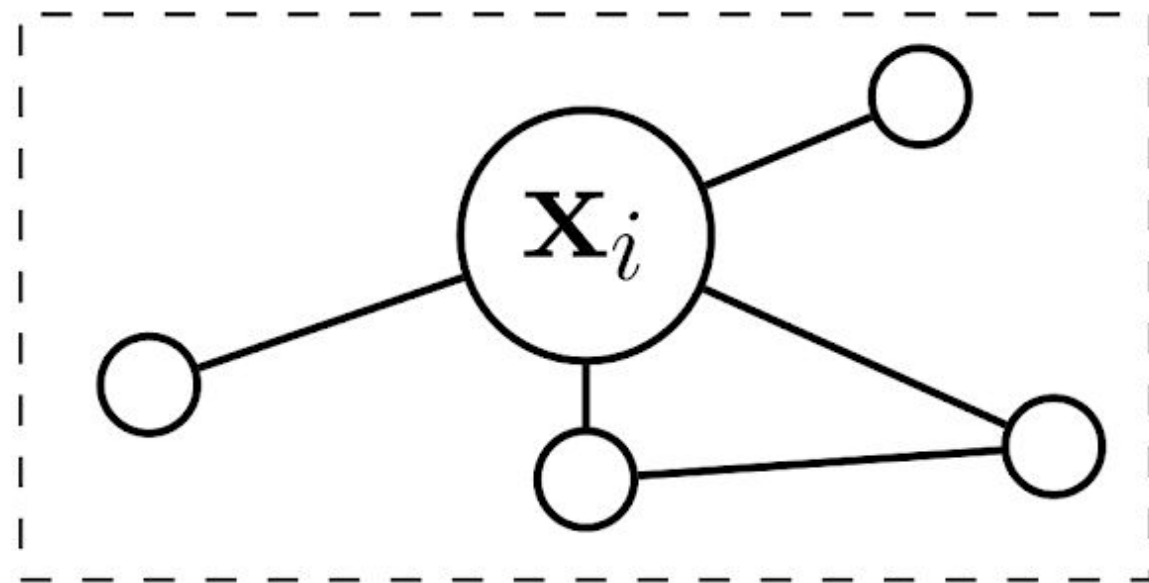
- ▶ Edge features \mathbf{e}_k : pseudoangular distance between particles

$$\Delta R = \sqrt{\Delta\eta^2 + \Delta\phi^2}$$

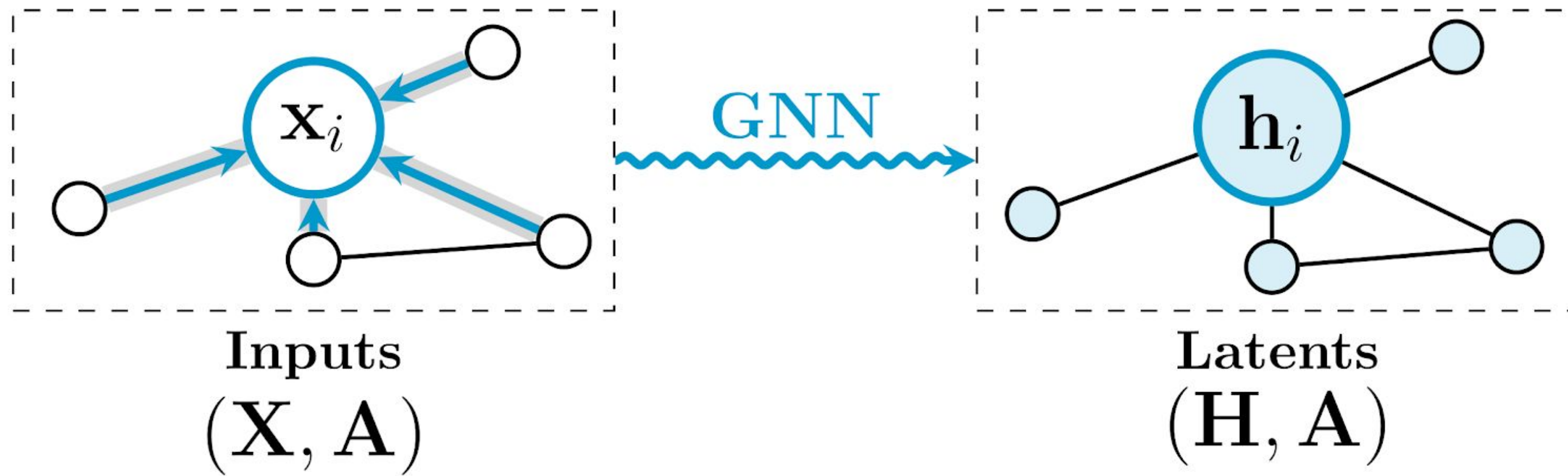


- ▶ Graph (global) features \mathbf{u} : jet mass

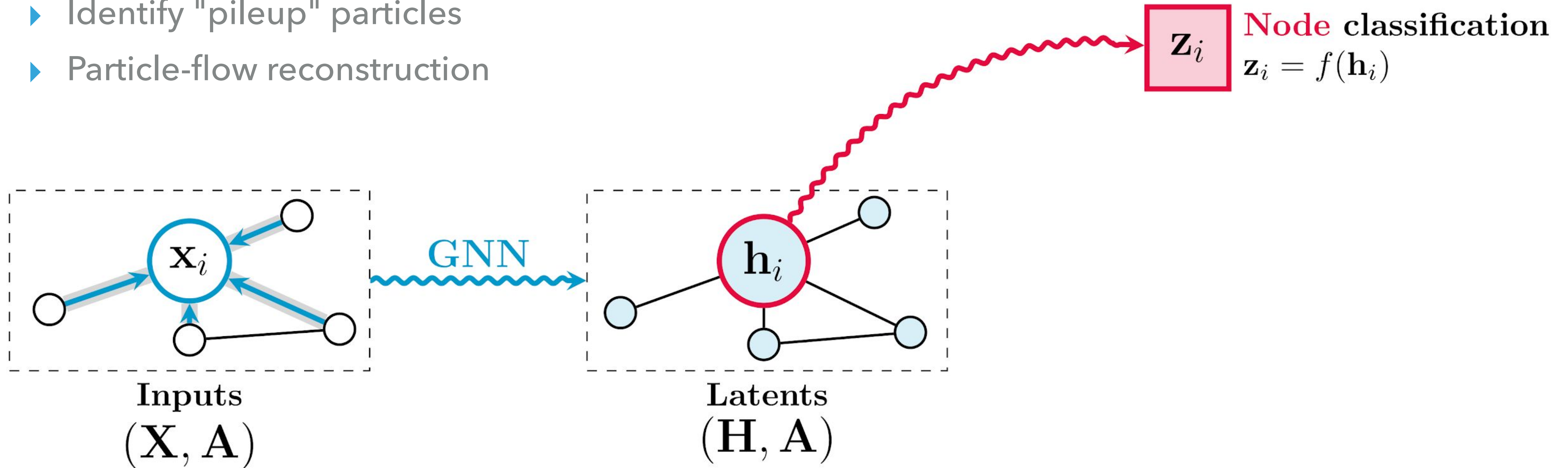
$$m = \sqrt{\sum_{i \in \text{jet}} E_i^2 - p_{x,i}^2 - p_{y,i}^2 - p_{z,i}^2}$$



Inputs
(\mathbf{X}, \mathbf{A})



- ▶ Node-level tasks
 - ▶ Correct cluster energies
 - ▶ Identify "pileup" particles
 - ▶ Particle-flow reconstruction



GRAPH NEURAL NETWORKS

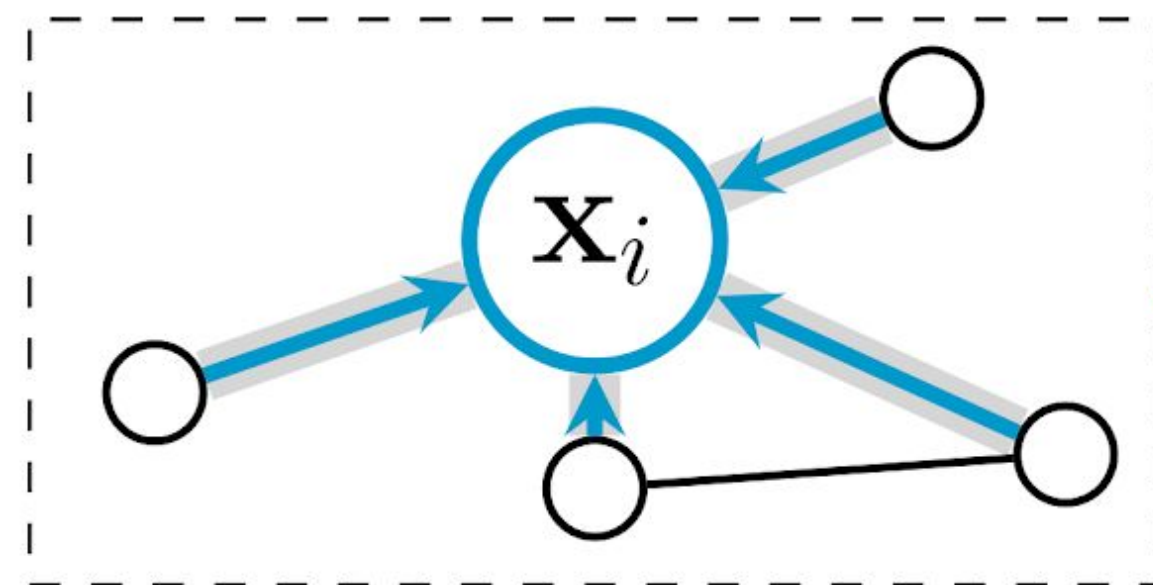
Source: <https://youtu.be/uF53xsT7mjc> 11

▶ Node-level tasks

- ▶ Correct cluster energies
- ▶ Identify "pileup" particles
- ▶ Particle-flow reconstruction

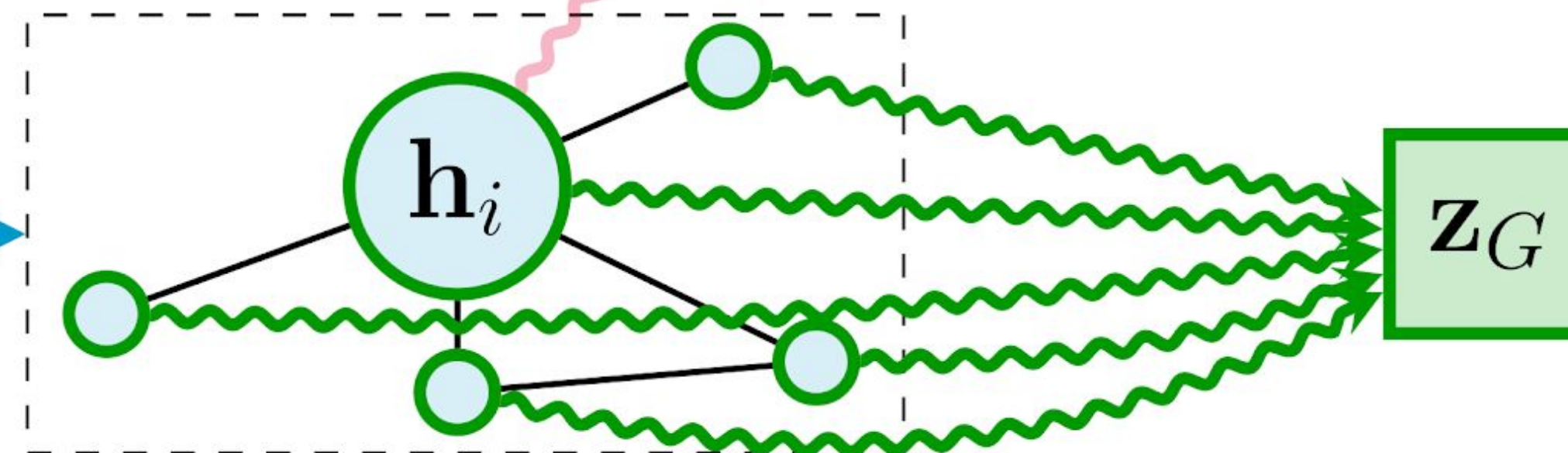
▶ Graph-level tasks

- ▶ **Jet tagging**
- ▶ Estimate shower energy
- ▶ Signal-to-background event discrimination

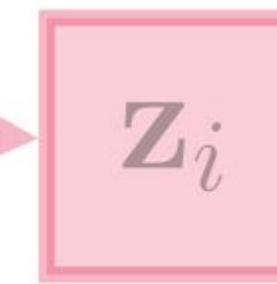


Inputs
 (\mathbf{X}, \mathbf{A})

GNN



Latents
 (\mathbf{H}, \mathbf{A})



Node classification
 $\mathbf{z}_i = f(\mathbf{h}_i)$



Graph classification
 $\mathbf{z}_G = f(\bigoplus_{i \in \mathcal{V}} \mathbf{h}_i)$

GRAPH NEURAL NETWORKS

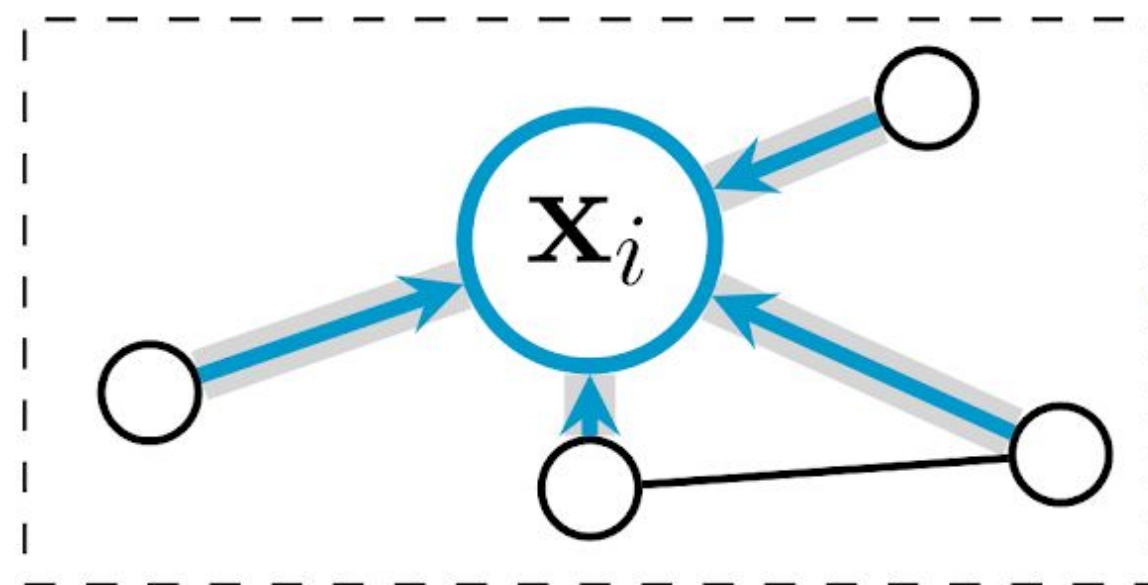
Source: <https://youtu.be/uF53xsT7mjc> 11

▶ Node-level tasks

- ▶ Correct cluster energies
- ▶ Identify "pileup" particles
- ▶ Particle-flow reconstruction

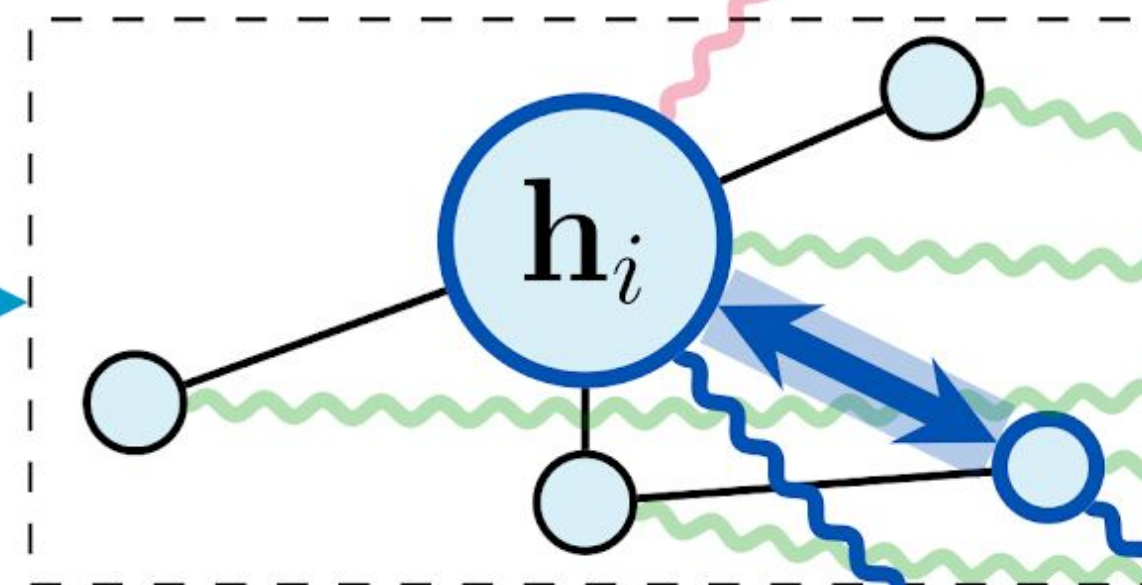
▶ Graph-level tasks

- ▶ **Jet tagging**
- ▶ Estimate shower energy
- ▶ Signal-to-background event discrimination



Inputs
 (\mathbf{X}, \mathbf{A})

GNN



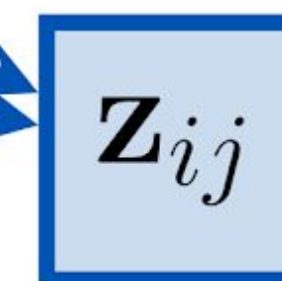
Latents
 (\mathbf{H}, \mathbf{A})



Node classification
 $\mathbf{z}_i = f(\mathbf{h}_i)$



Graph classification
 $\mathbf{z}_G = f(\bigoplus_{i \in \mathcal{V}} \mathbf{h}_i)$

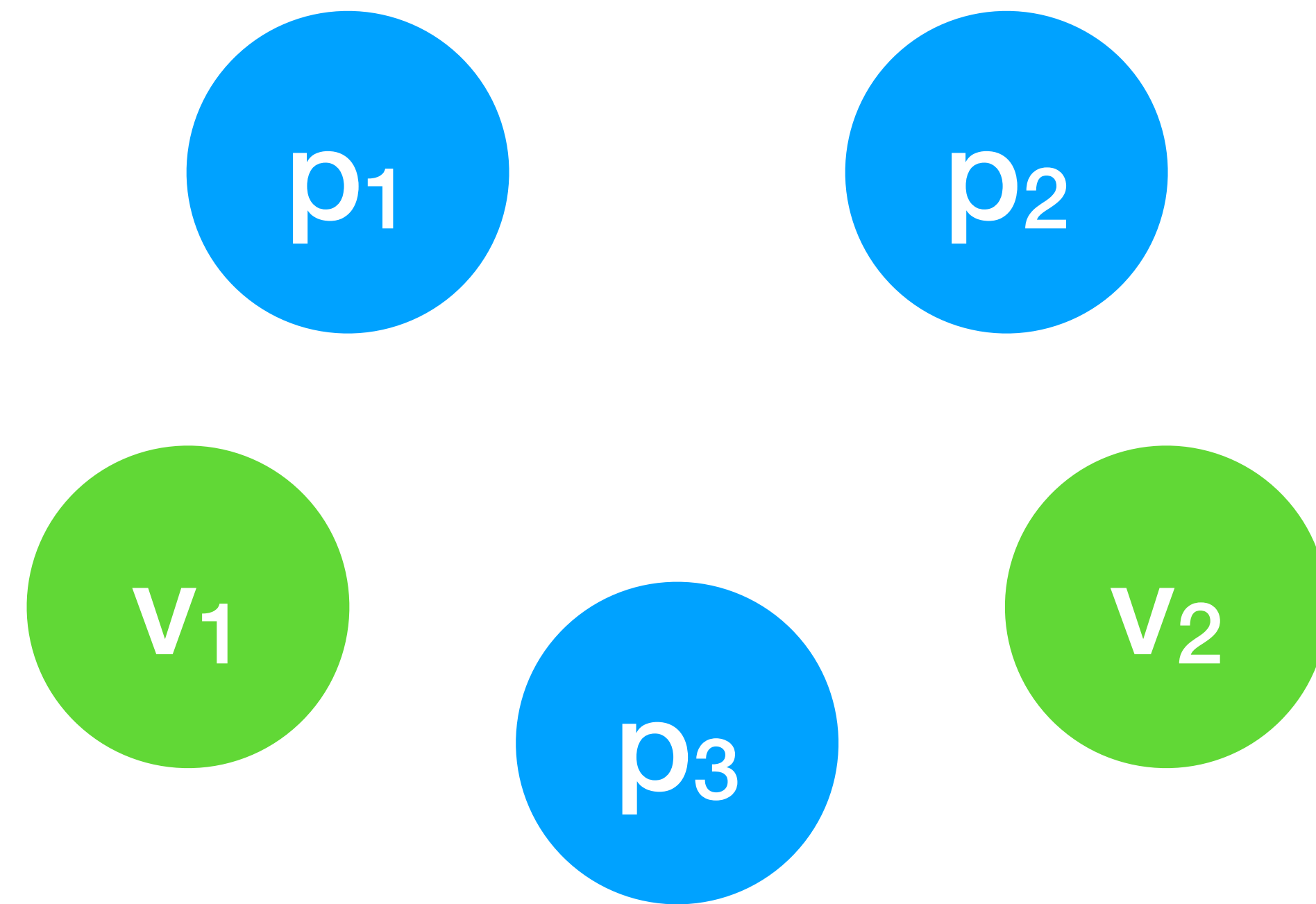


Link prediction
 $\mathbf{z}_{ij} = f(\mathbf{h}_i, \mathbf{h}_j, \mathbf{e}_{ij})$

▶ Edge-level tasks

- ▶ Identify track segments
- ▶ Estimate track parameters
- ▶ Secondary vertex reconstruction

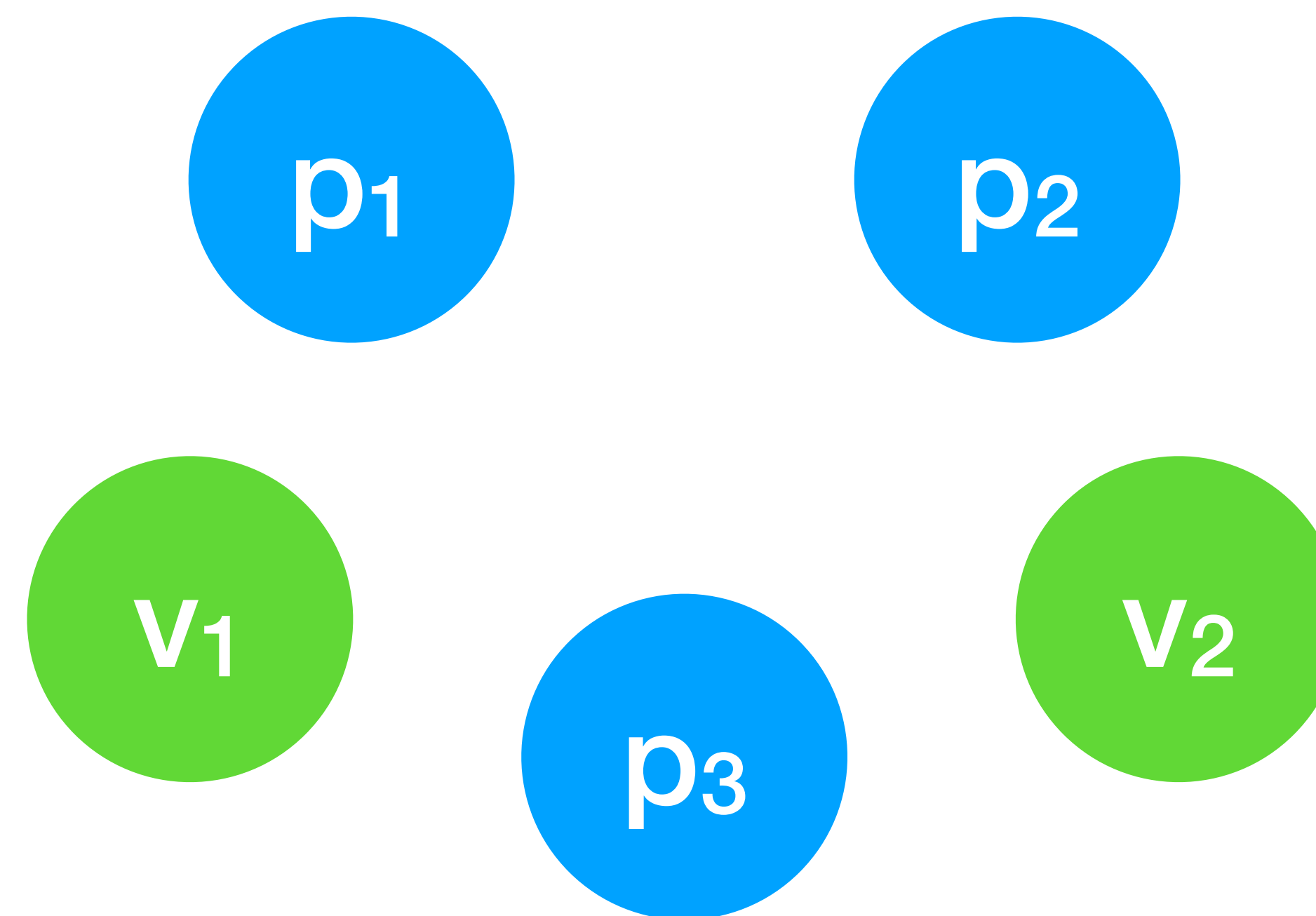
$$p_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, d_{3D}, \text{cov}(p_T, p_T), \dots]$$



$$v_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, n_{\text{tracks}}, \cos \theta_{PV}, \dots]$$

- ▶ Particles (i.e. tracks) and vertices are two separate inputs with different feature vectors (*heterogenous graph*)

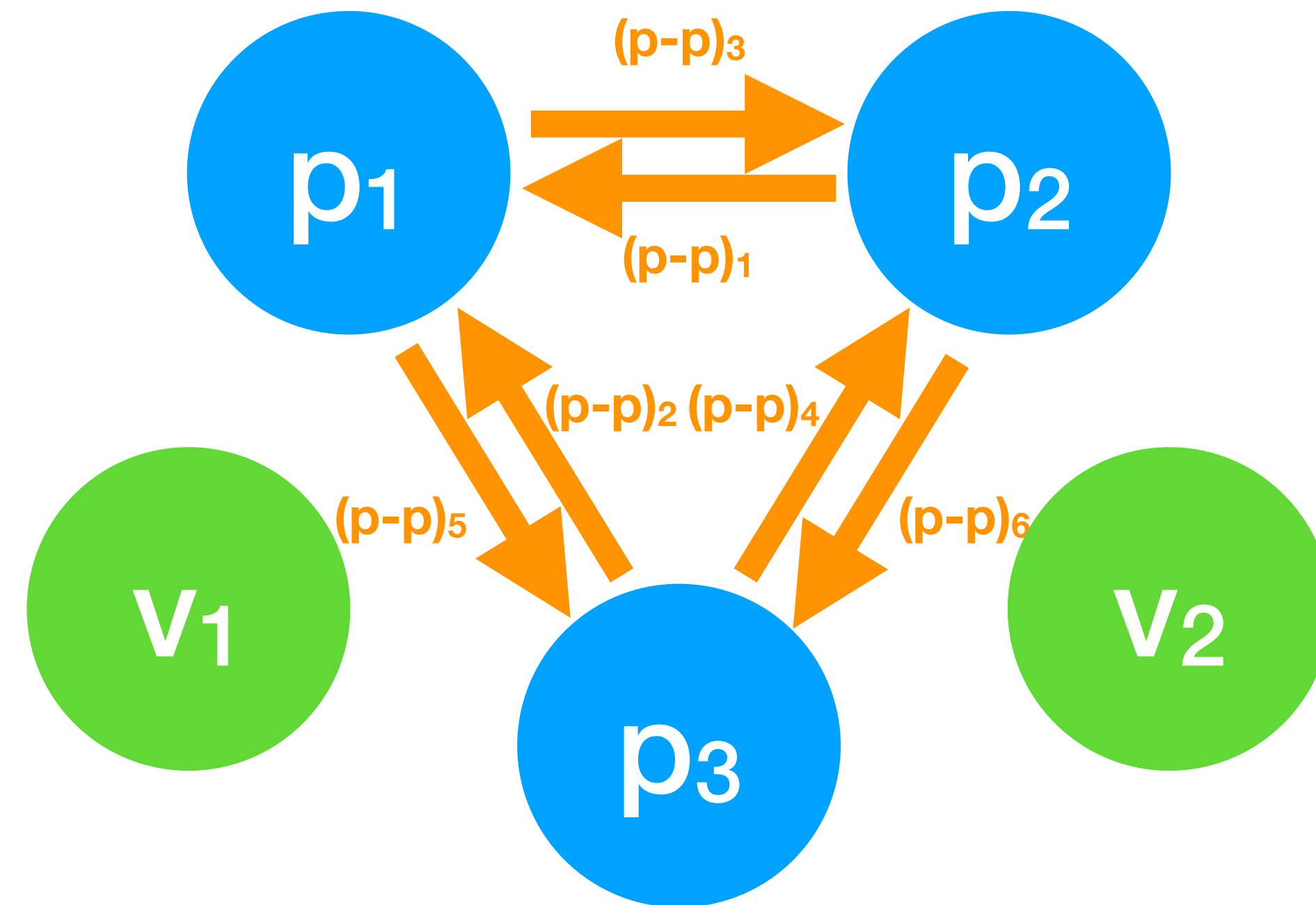
$$p_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, d_{3D}, \text{cov}(p_T, p_T), \dots]$$



$$v_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, n_{\text{tracks}}, \cos \theta_{PV}, \dots]$$

- ▶ Particles (i.e. tracks) and vertices are two separate inputs with different feature vectors (*heterogenous graph*)
- ▶ GNNs typically consider a *homogenous graph* (e.g. particle-particle graph)

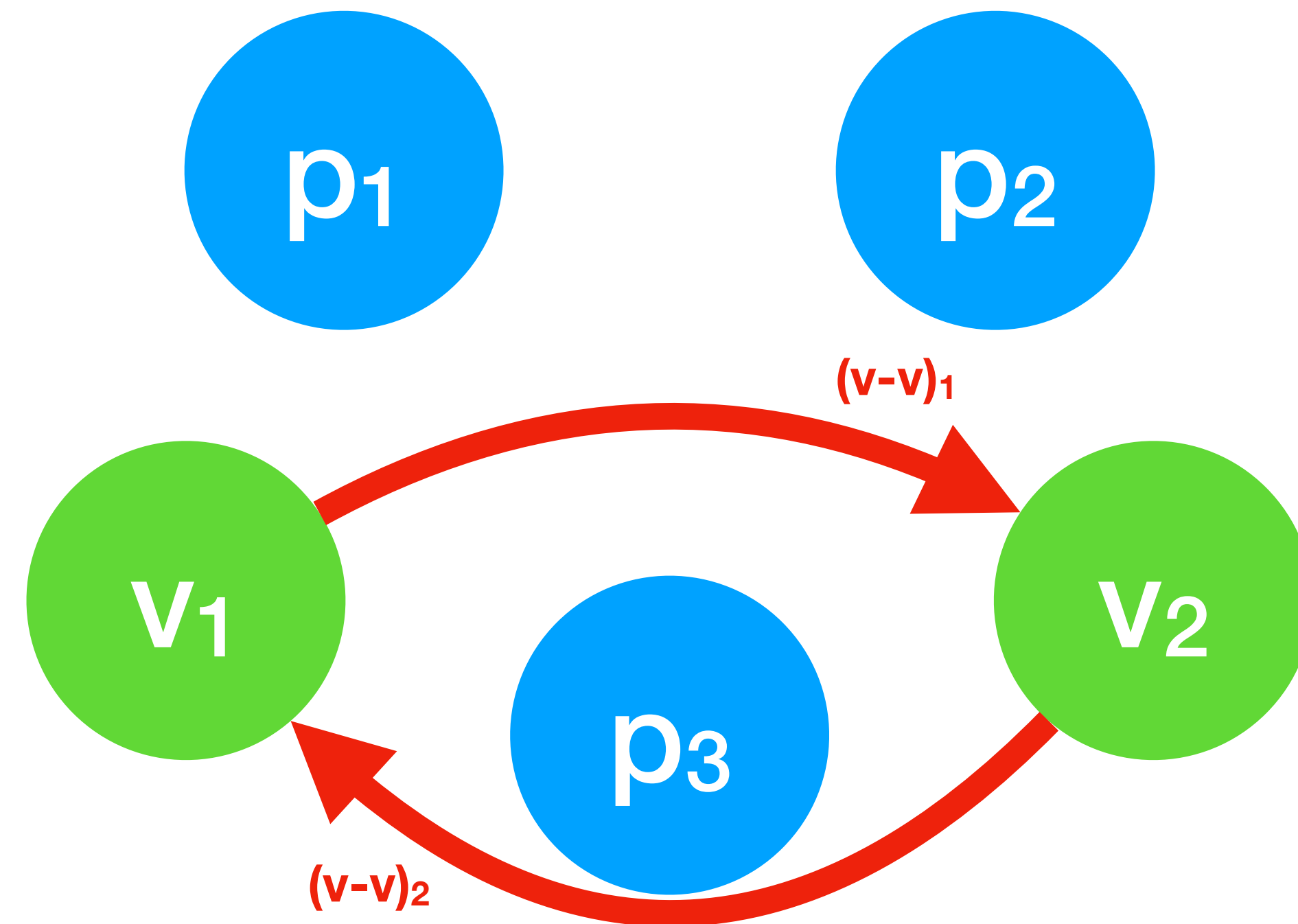
$$p_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, d_{3D}, \text{cov}(p_T, p_T), \dots]$$



$$v_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, n_{\text{tracks}}, \cos \theta_{\text{PV}}, \dots]$$

- ▶ Particles (i.e. tracks) and vertices are two separate inputs with different feature vectors (*heterogenous graph*)
- ▶ GNNs typically consider a *homogenous graph* (e.g. particle-particle graph)
- ▶ Vertex-vertex graph can also be considered

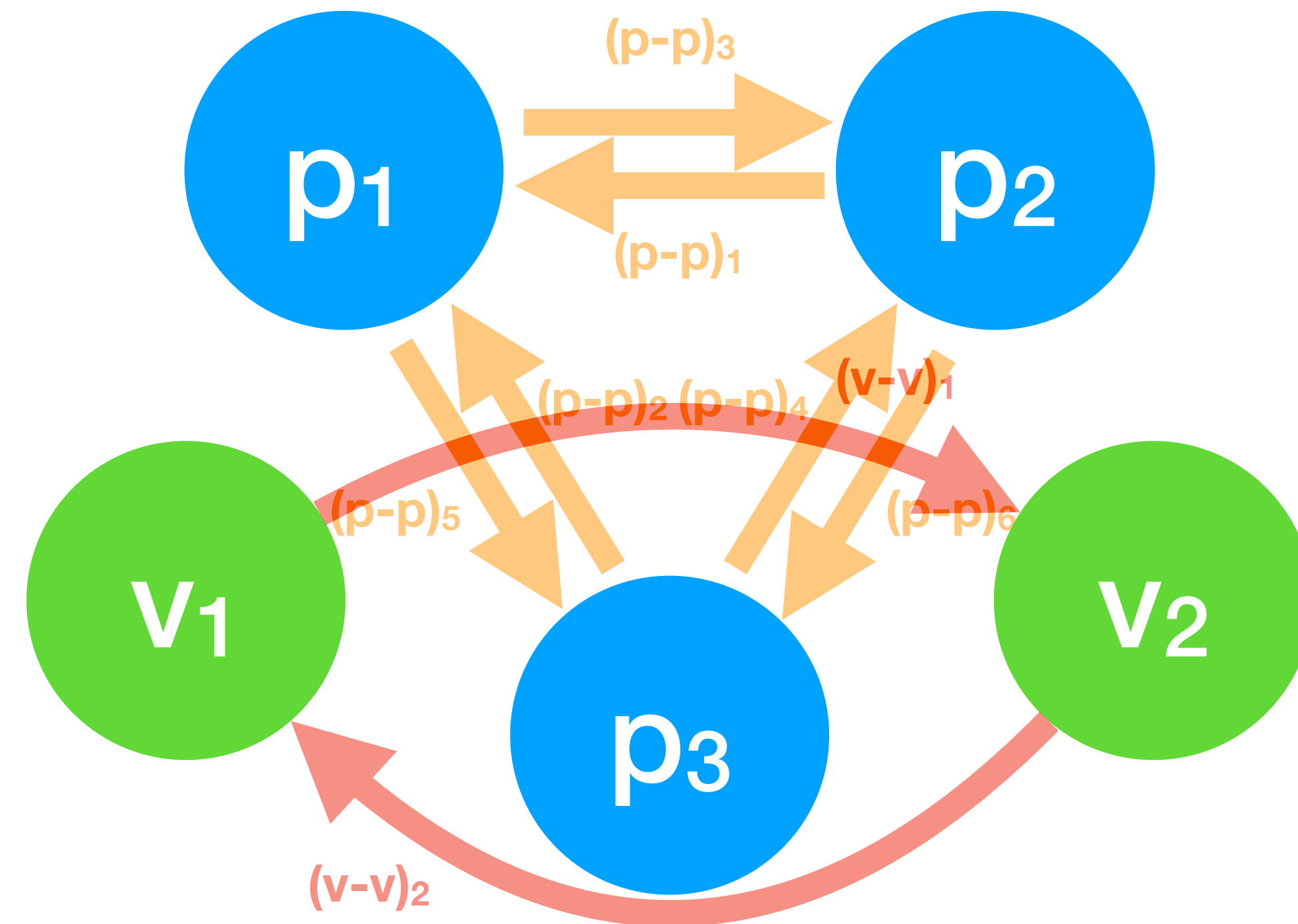
$$p_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, d_{3D}, \text{cov}(p_T, p_T), \dots]$$



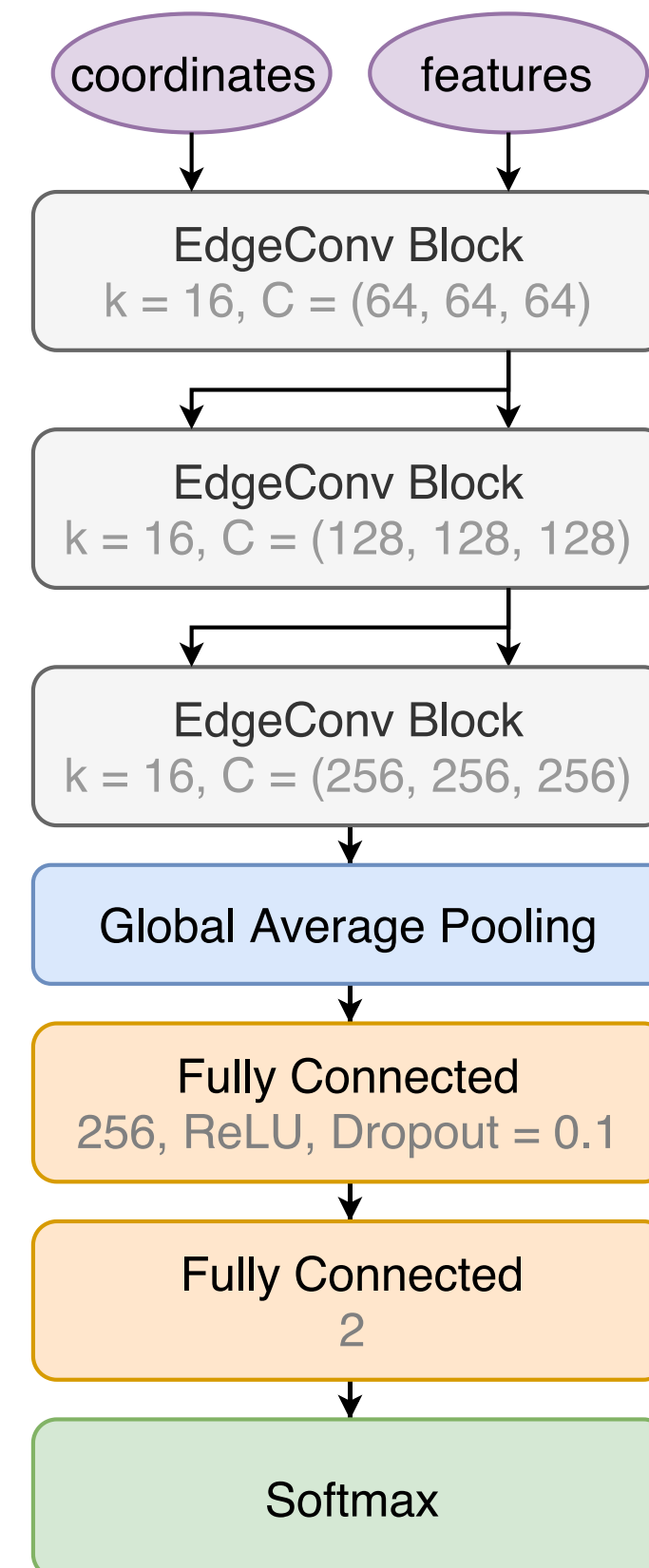
$$v_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, n_{\text{tracks}}, \cos \theta_{\text{PV}}, \dots]$$

- ▶ Particles (i.e. tracks) and vertices are two separate inputs with different feature vectors (*heterogenous graph*)
- ▶ GNNs typically consider a *homogenous graph* (e.g. particle-particle graph)
- ▶ Vertex-vertex graph can also be considered
- ▶ Combined GNN can consider both by constructing two separate graphs

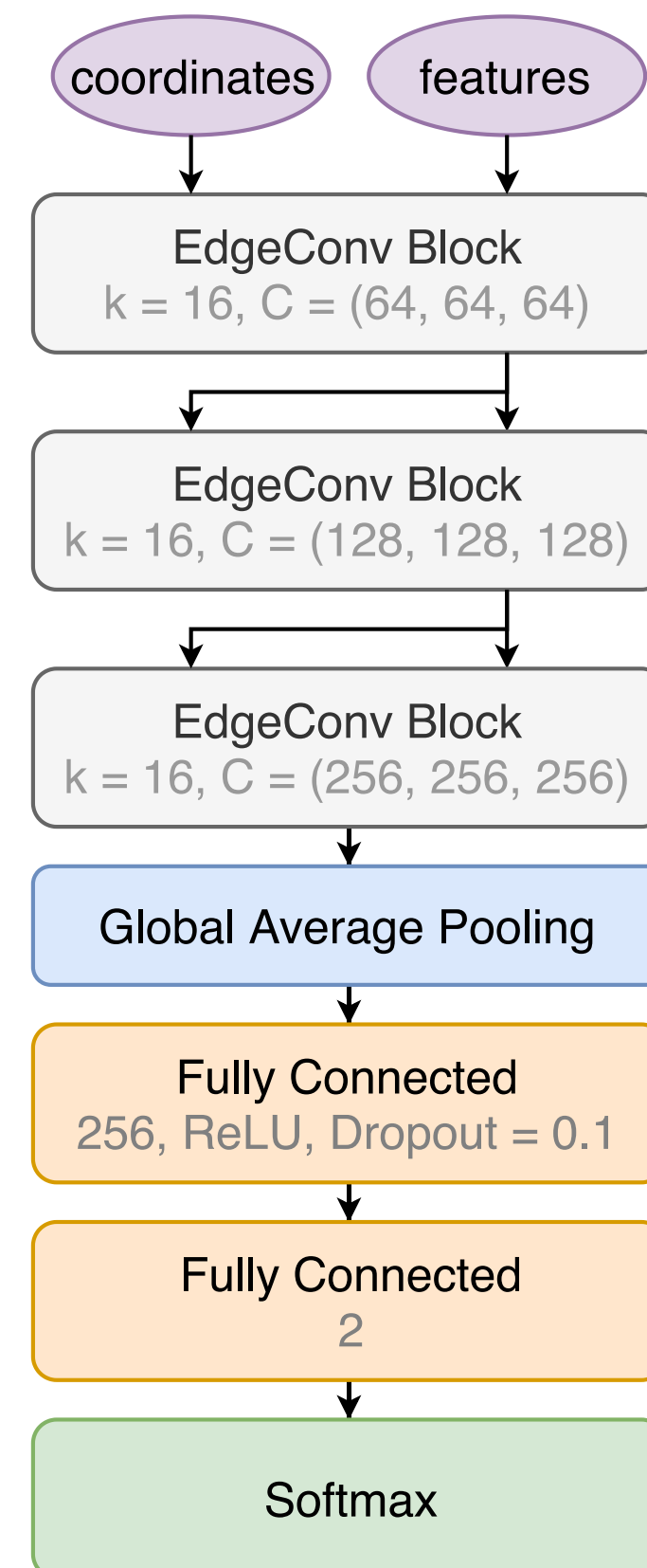
$$p_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, d_{3D}, \text{cov}(p_T, p_T), \dots]$$



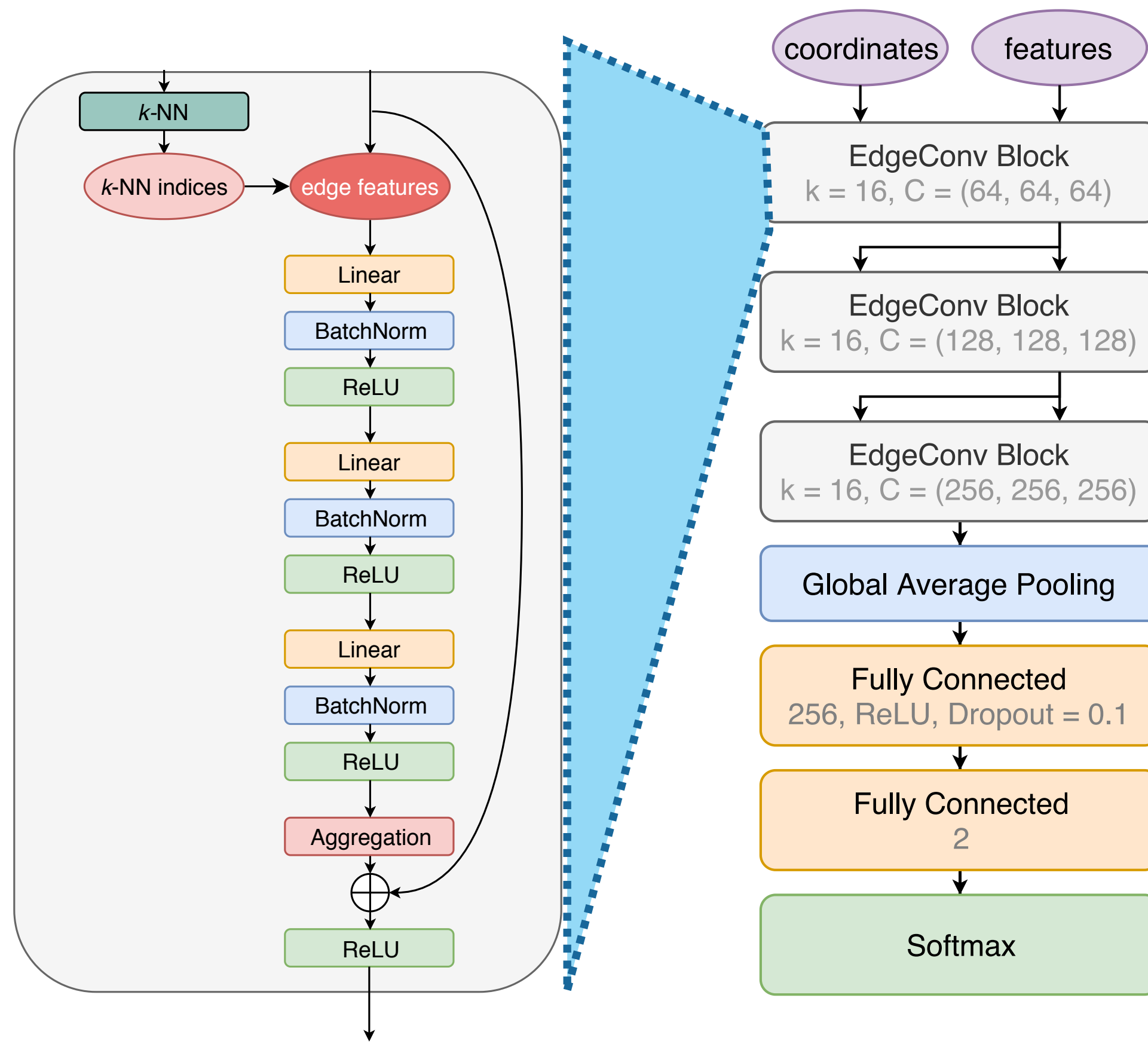
$$v_i = [p_T^{\text{rel}}, \phi^{\text{rel}}, \eta^{\text{rel}}, \dots, n_{\text{tracks}}, \cos \theta_{PV}, \dots]$$



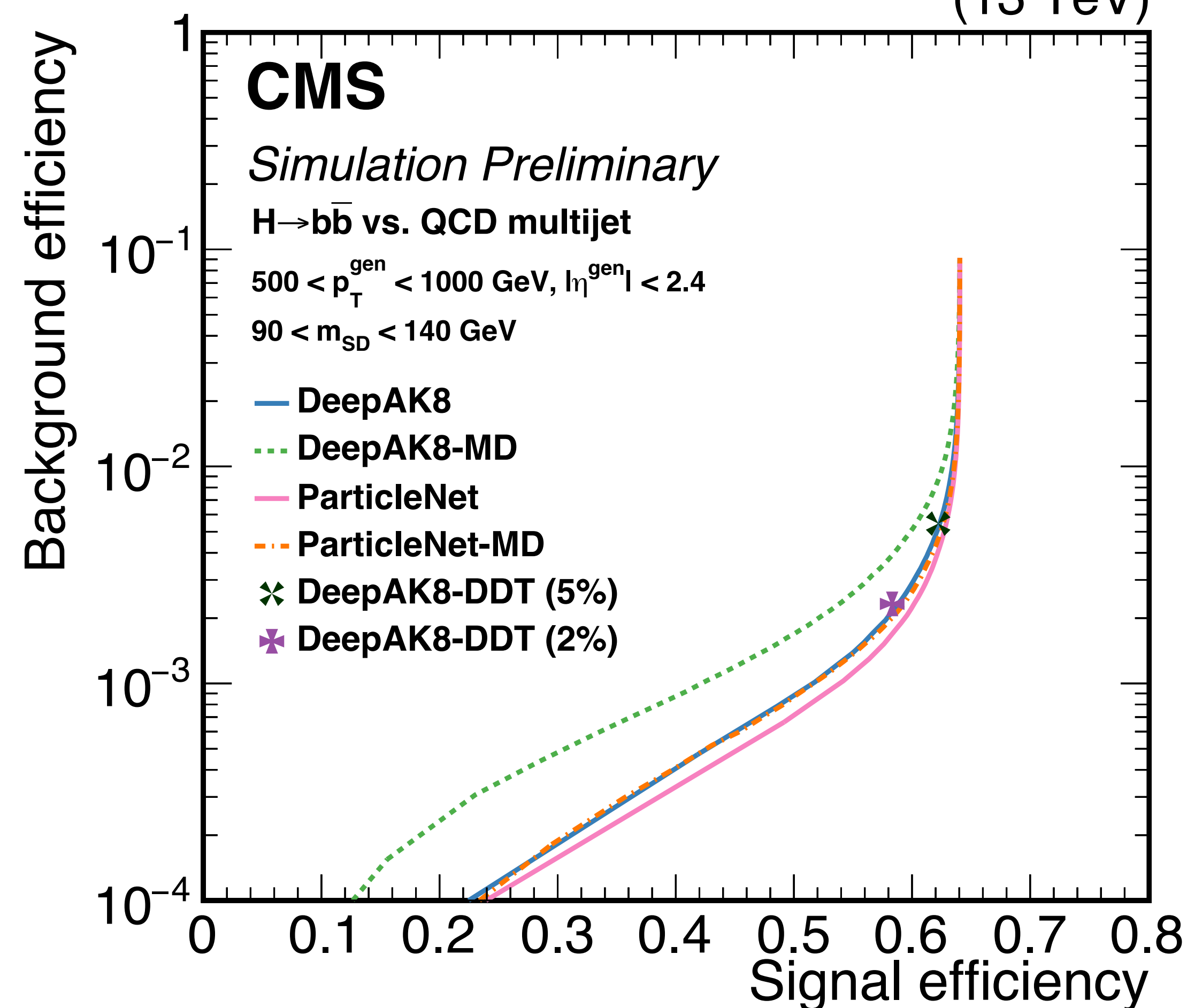
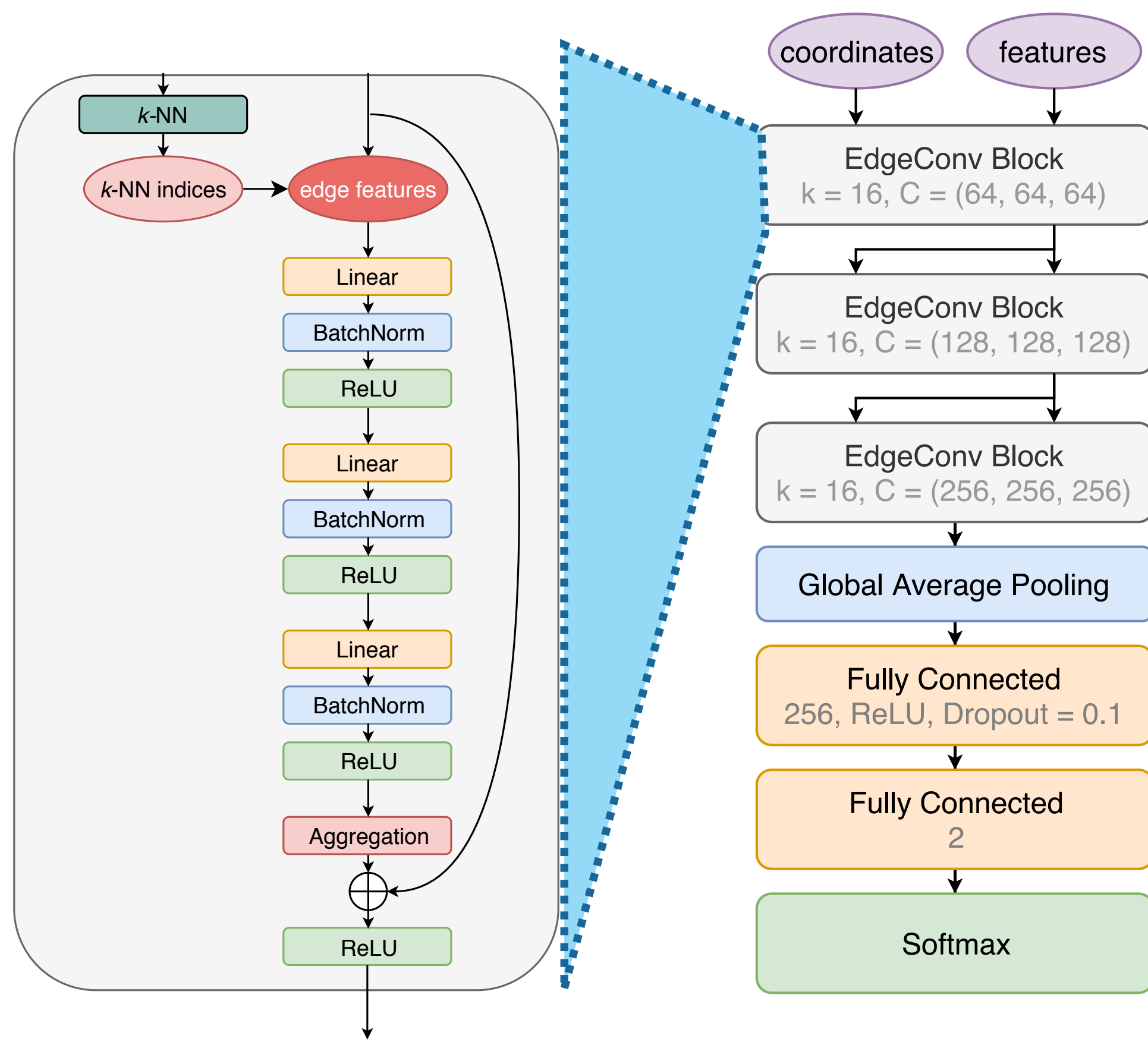
- ▶ ParticleNet, using “dynamic edge convolutions:” graph is constructed based on “closeness” in a latent space



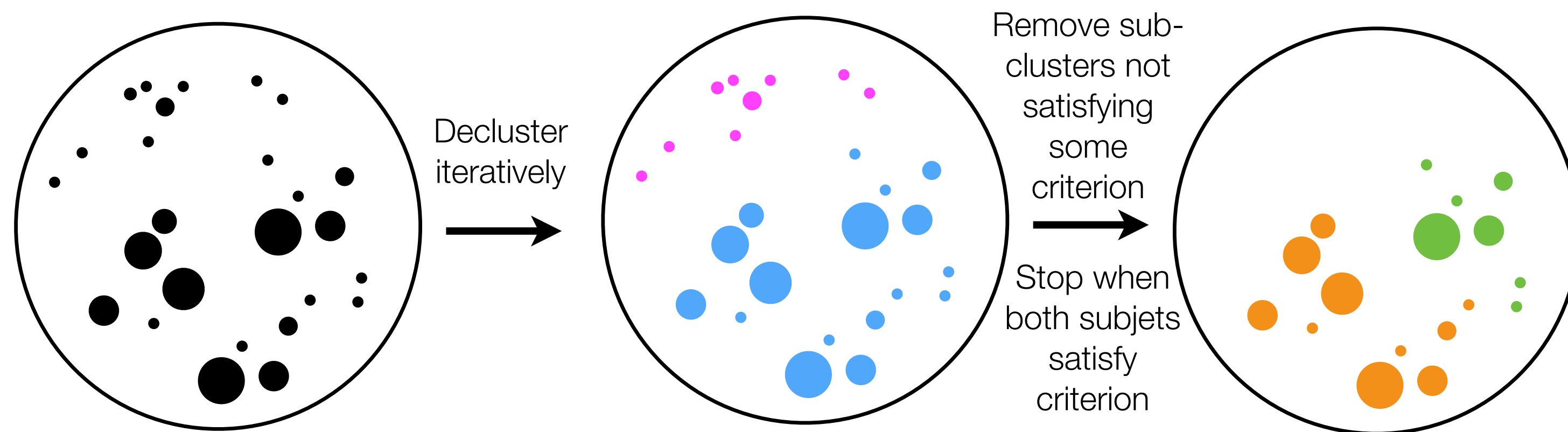
- ▶ ParticleNet, using “dynamic edge convolutions:” graph is constructed based on “closeness” in a latent space



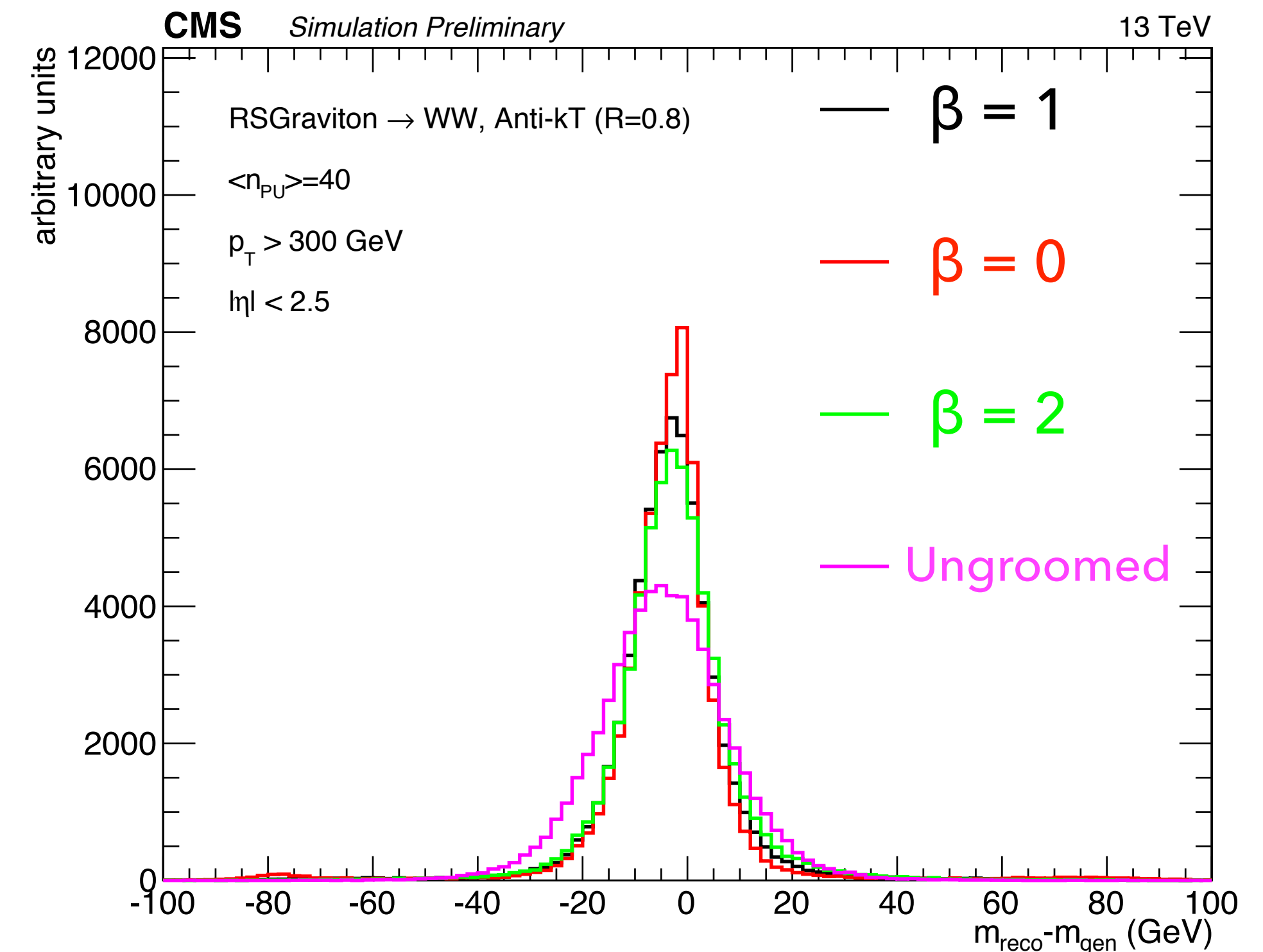
- ▶ ParticleNet, using “dynamic edge convolutions:” graph is constructed based on “closeness” in a latent space
- ▶ Identifies H(bb) with a true positive rate of over 50% and a false positive rate of 0.1%



- ▶ An important property used to analyze Higgs boson jets is the invariant mass
 - ▶ Provides good separation between W/Z/H-jets and q/g jets
 - ▶ Grooming removes soft and wide-angle radiation (soft drop is CMS standard)
 - ▶ Can we do better with ML?



Soft Drop Condition:
$$\frac{\min(p_{T1}, p_{T2})}{p_{T1} + p_{T2}} > z_{\text{cut}} \left(\frac{\Delta R_{12}}{R_0} \right)^\beta$$



CMS: $z_{\text{cut}} = 0.1, \beta = 0$

- ▶ Reuse ParticleNet architecture with a target of the “true” jet mass
- ▶ Special training samples incorporate $X \rightarrow bb$, $X \rightarrow cc$, $X \rightarrow qq$ with varying X mass in $[15, 250]$ GeV

$$M_{\text{target}} = \begin{cases} M_{\text{SD}}^{\text{gen}} & \text{if jet is QCD} \\ m_X \in [15, 250] \text{ GeV} & \text{otherwise} \end{cases}$$

- ▶ Minimize loss function:

$$L(y, y^p) = \sum_{i=1}^n \log \cosh(y_i^p - y_i)$$

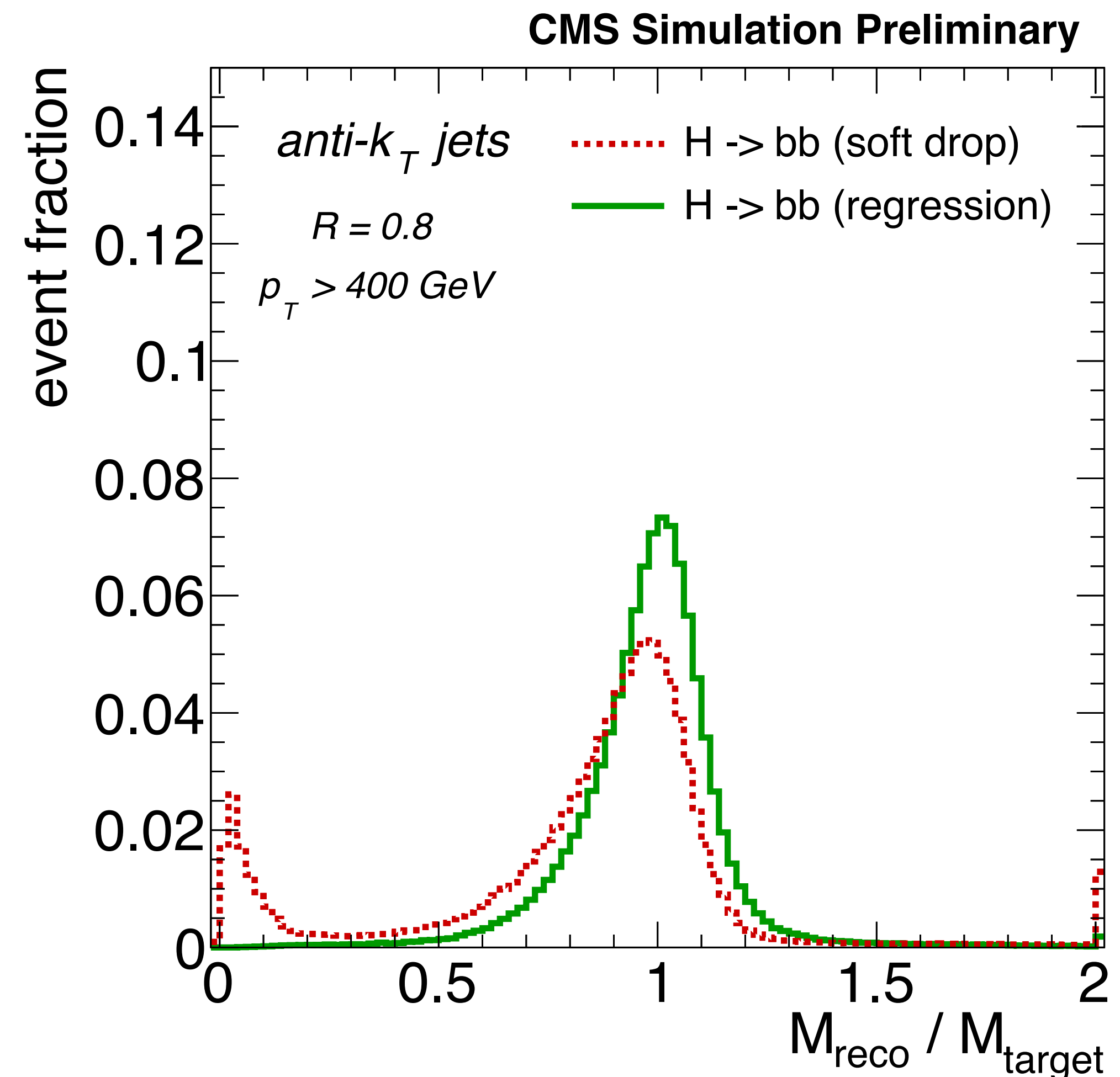
- ▶ Reuse ParticleNet architecture with a target of the “true” jet mass
- ▶ Special training samples incorporate $X \rightarrow bb$, $X \rightarrow cc$, $X \rightarrow qq$ with varying X mass in $[15, 250]$ GeV

$$M_{\text{target}} = \begin{cases} M_{\text{SD}}^{\text{gen}} & \text{if jet is QCD} \\ m_X \in [15, 250] \text{ GeV} & \text{otherwise} \end{cases}$$

- ▶ Minimize loss function:

$$L(y, y^p) = \sum_{i=1}^n \log \cosh(y_i^p - y_i)$$

- ▶ Substantial scale and resolution improvement



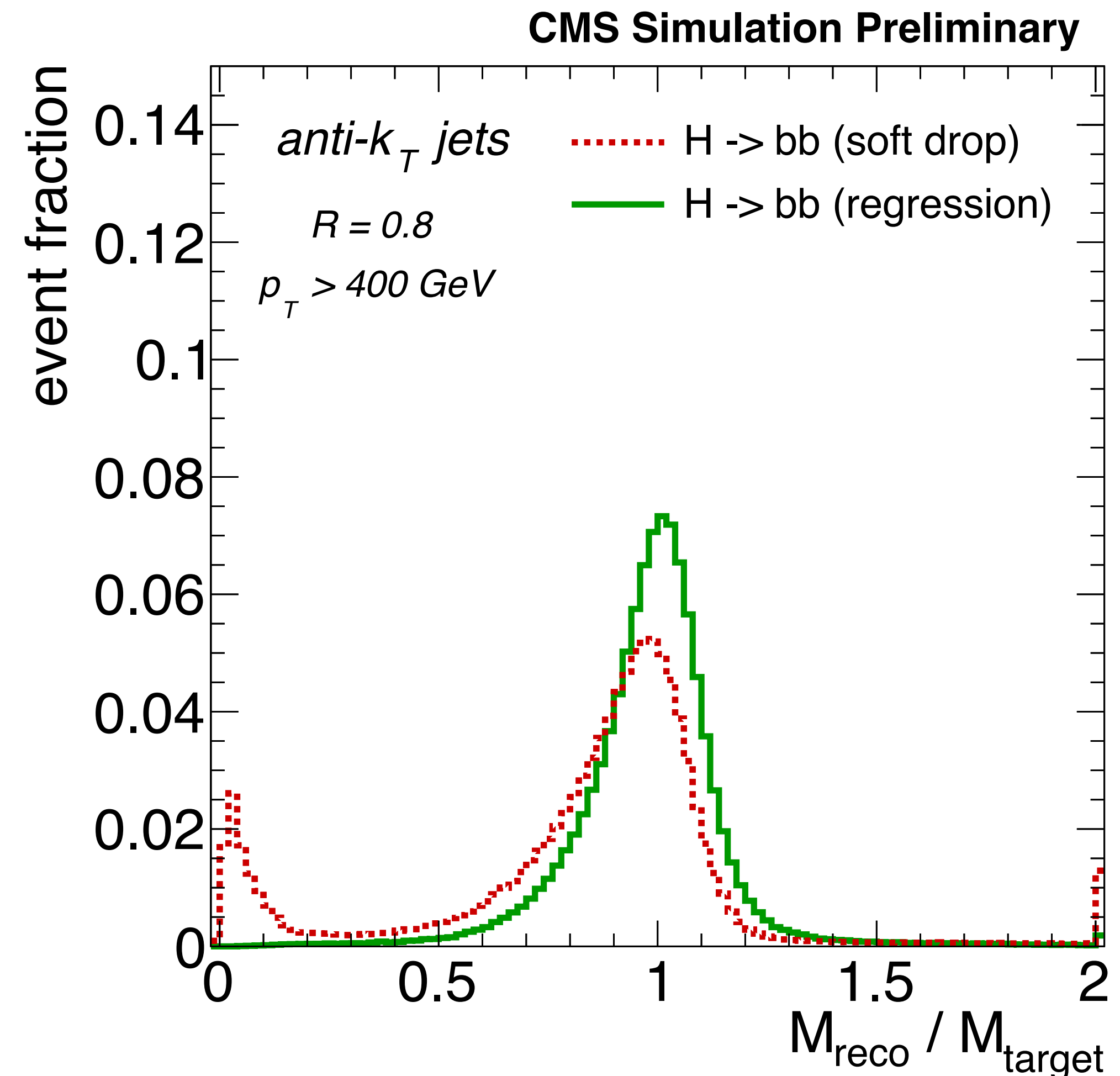
- ▶ Reuse ParticleNet architecture with a target of the “true” jet mass
- ▶ Special training samples incorporate $X \rightarrow bb$, $X \rightarrow cc$, $X \rightarrow qq$ with varying X mass in $[15, 250]$ GeV

$$M_{\text{target}} = \begin{cases} M_{\text{SD}}^{\text{gen}} & \text{if jet is QCD} \\ m_X \in [15, 250] \text{ GeV} & \text{otherwise} \end{cases}$$

- ▶ Minimize loss function:

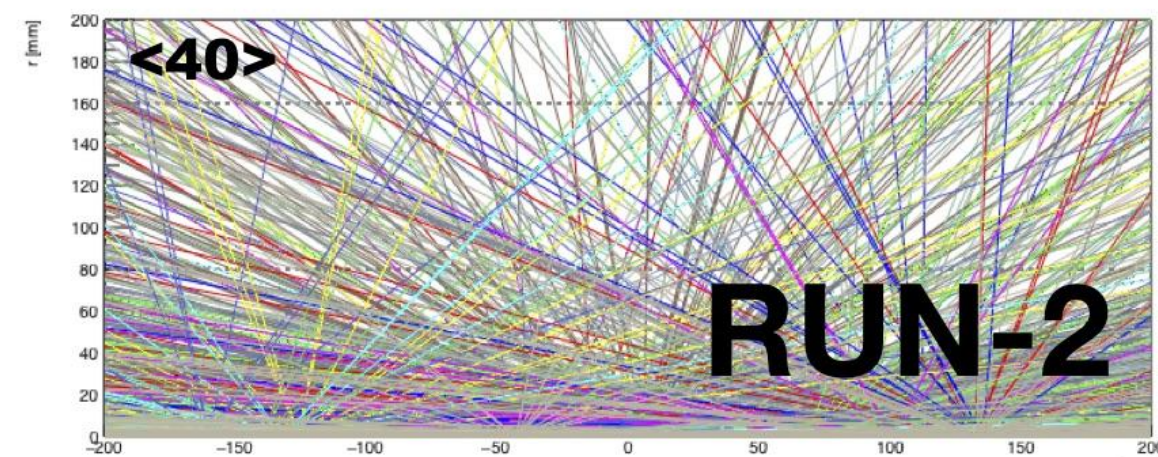
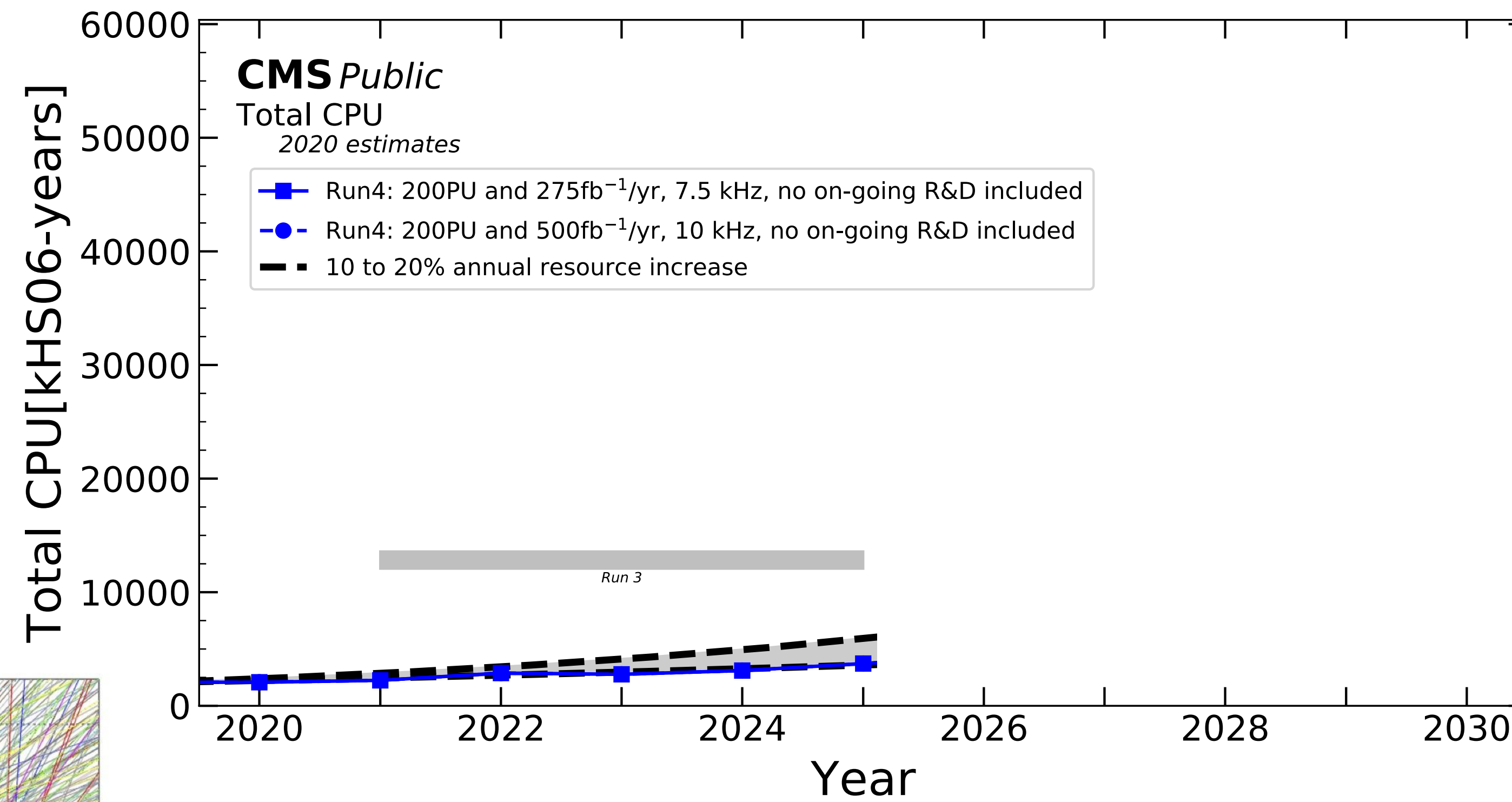
$$L(y, y^p) = \sum_{i=1}^n \log \cosh(y_i^p - y_i)$$

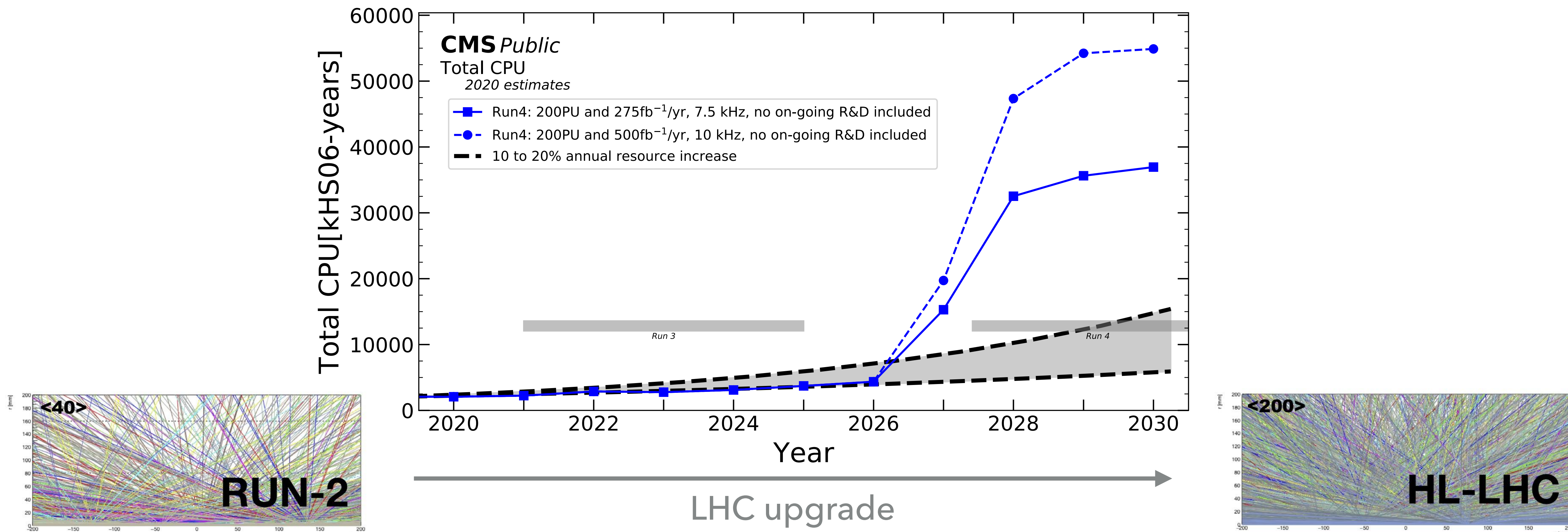
- ▶ Substantial scale and resolution improvement
 - ▶ Can increase sensitivity by 20-25% to rare Higgs boson signals like HH, VBF, ...



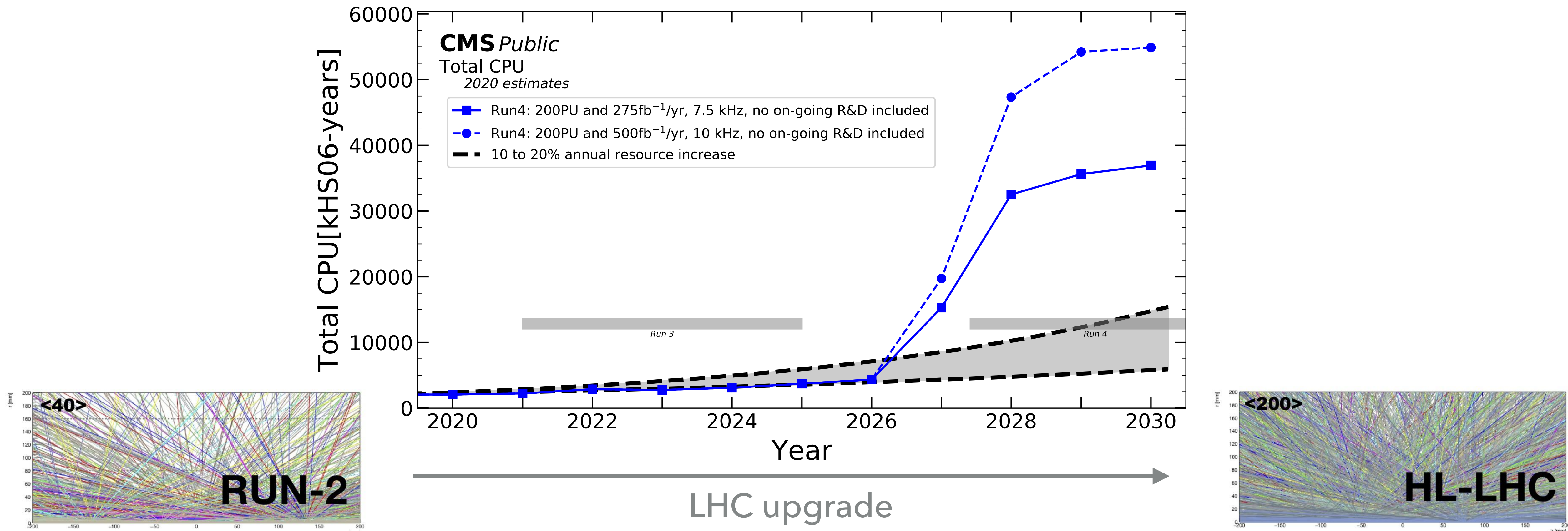


ML FOR TAGGING
ML FOR GEN/SIM
FAST ML FOR TRIGGER

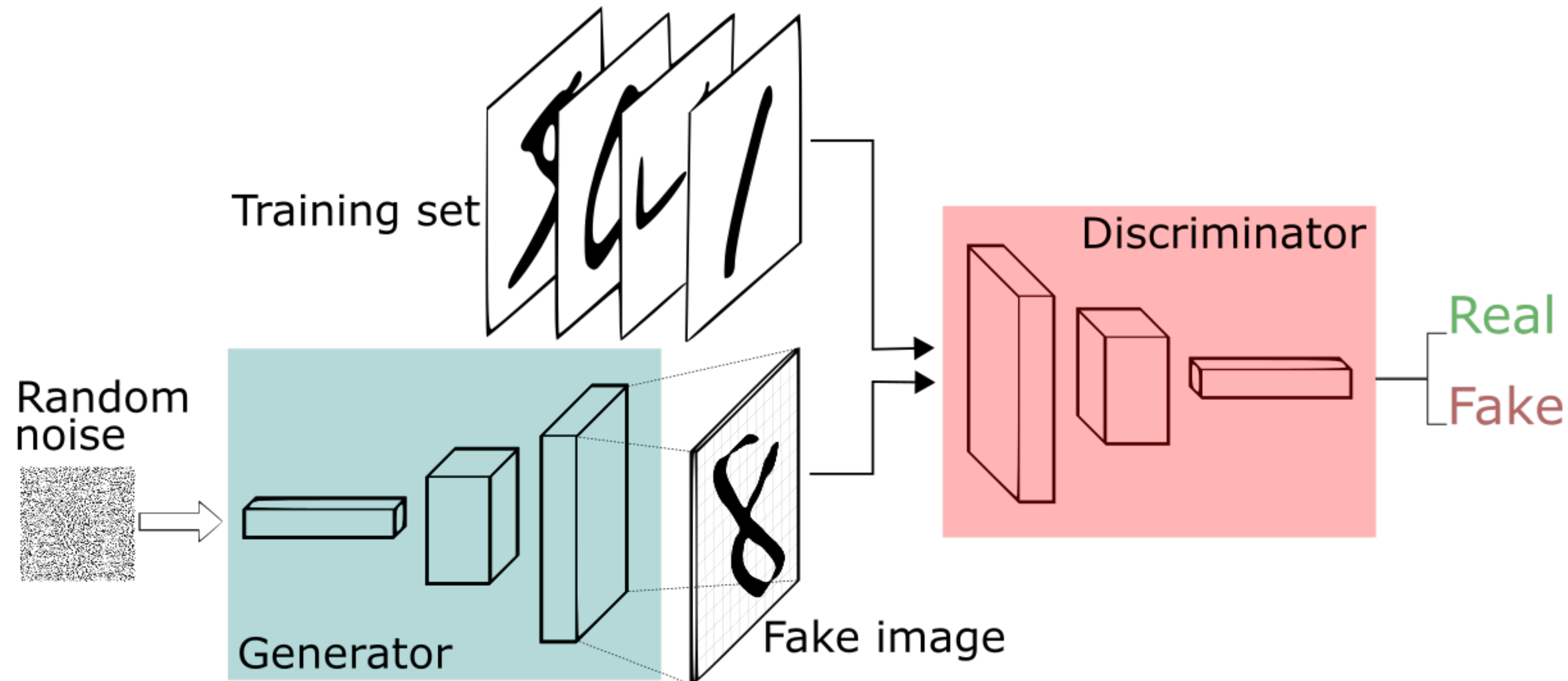




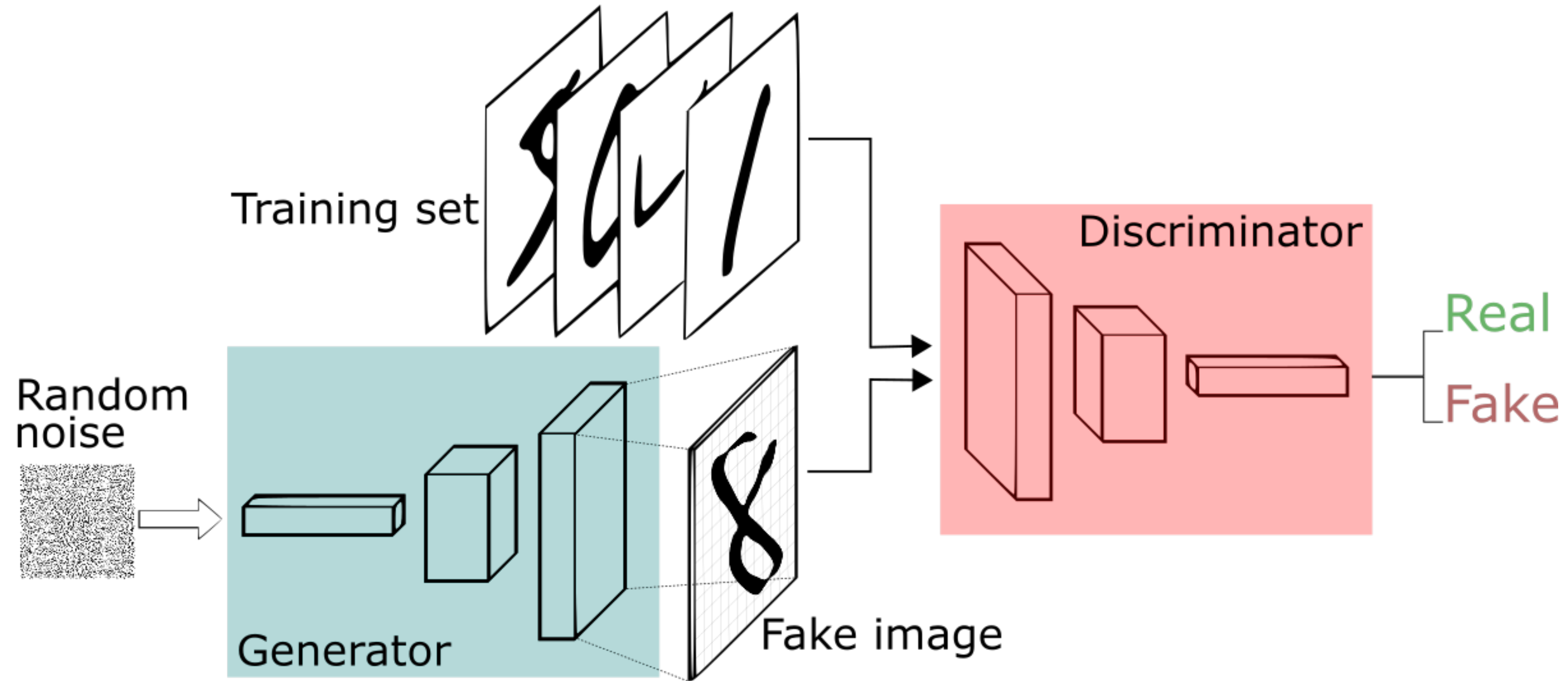
- ▶ Computing demands increase nonlinearly with increasing “pileup”



- ▶ Computing demands increase nonlinearly with increasing “pileup”
- ▶ Need more processing power (or smarter algorithms like **deep learning**) to keep up with demands



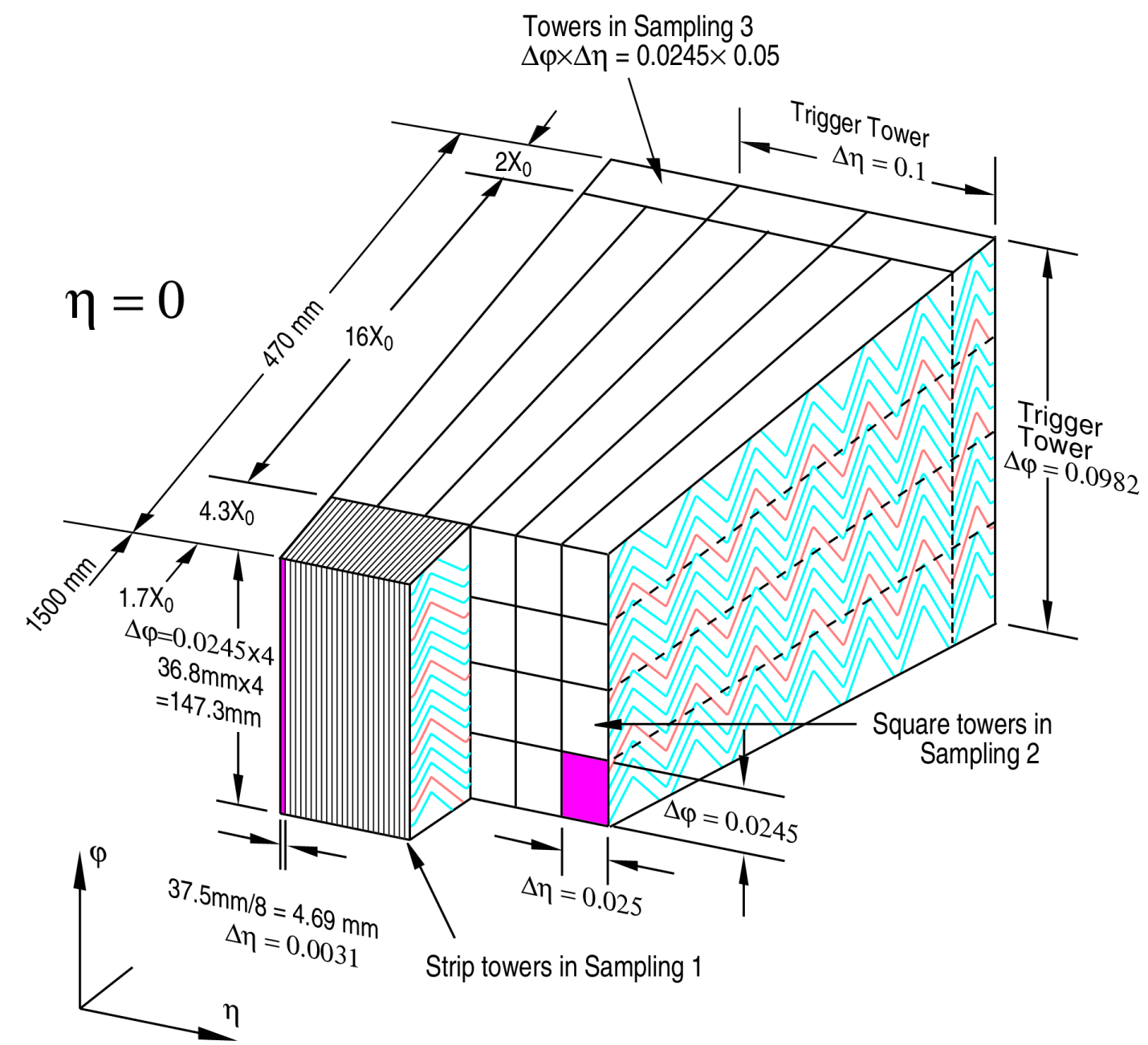
- ▶ Train two neural networks in tandem:
 - ▶ one to generate realistic "fake" data
 - ▶ the other to discriminate "real" from "fake" data

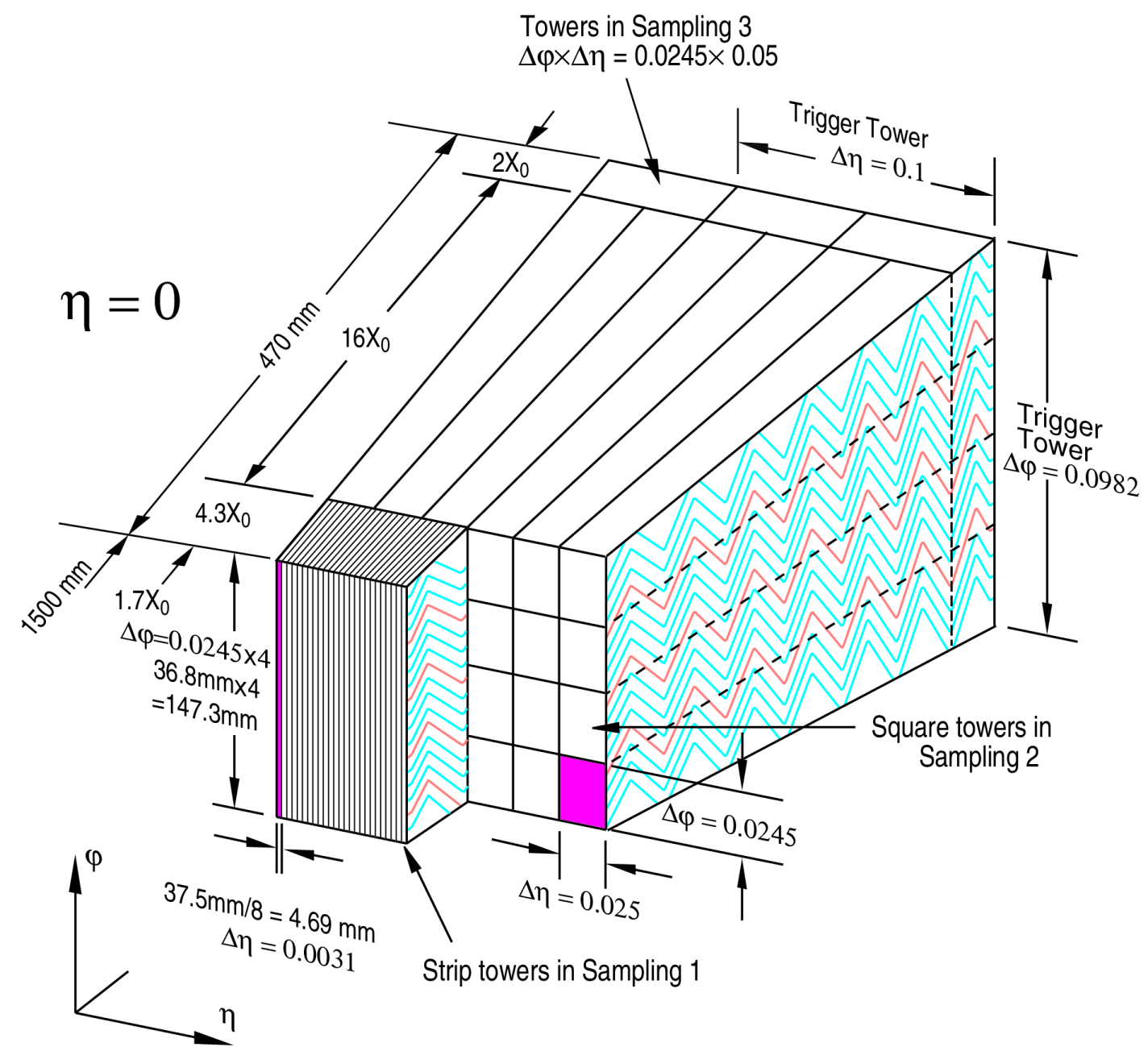


thispersondoesnotexist.com

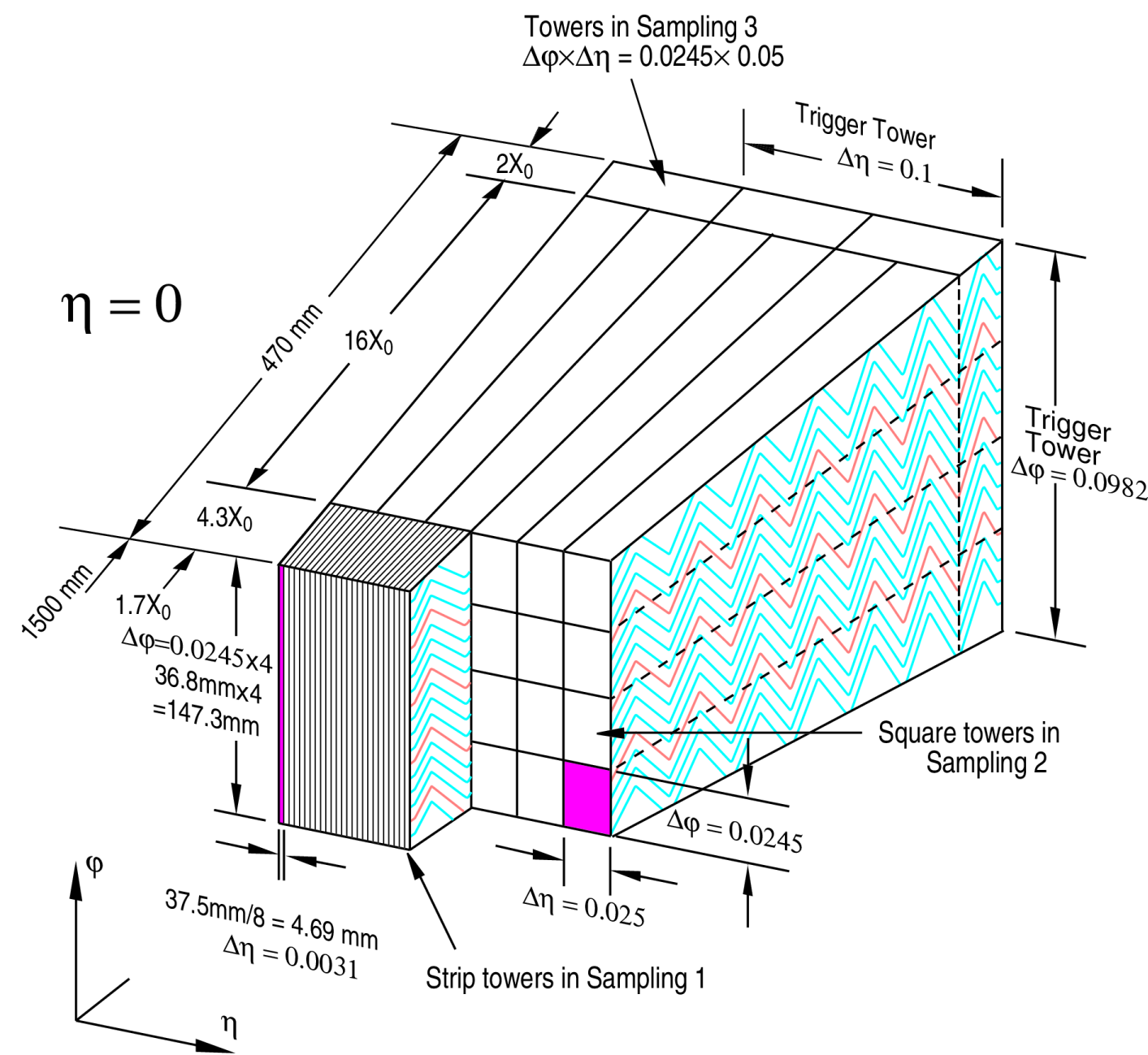


- ▶ Train two neural networks in tandem:
 - ▶ one to generate realistic "fake" data
 - ▶ the other to discriminate "real" from "fake" data

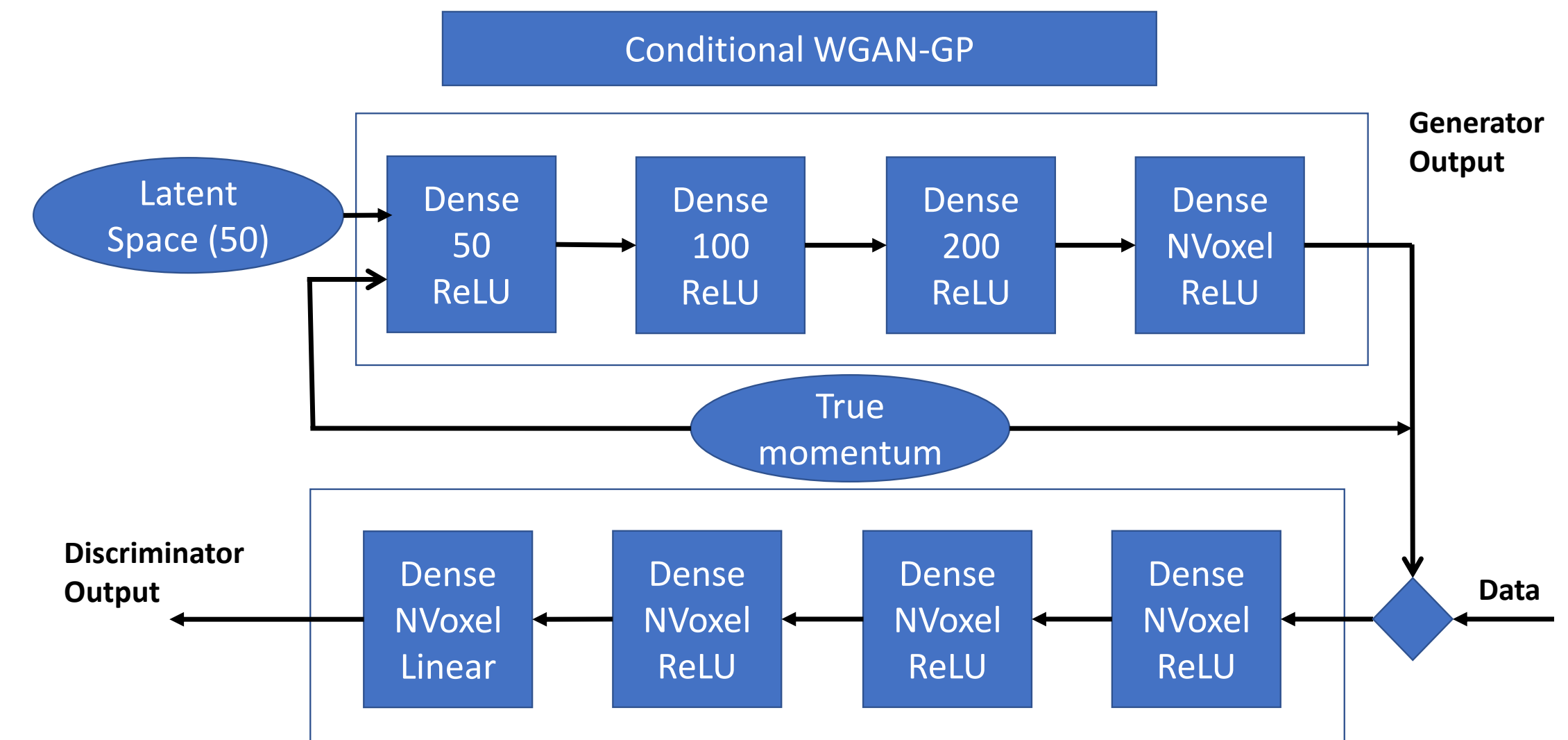




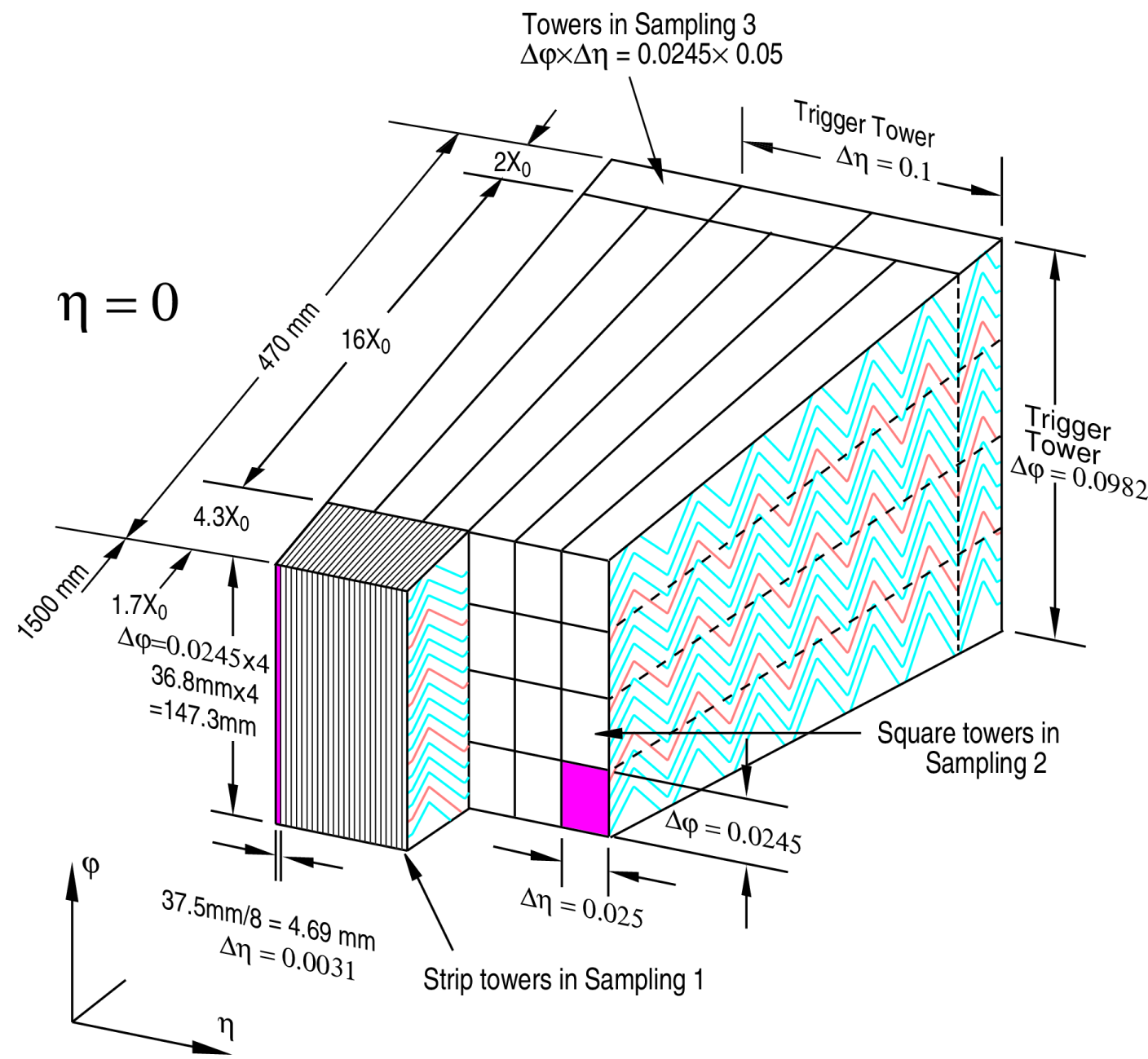
- ▶ Geant4-based ATLAS simulation of the full calorimeter is slow; can a GAN replace this?



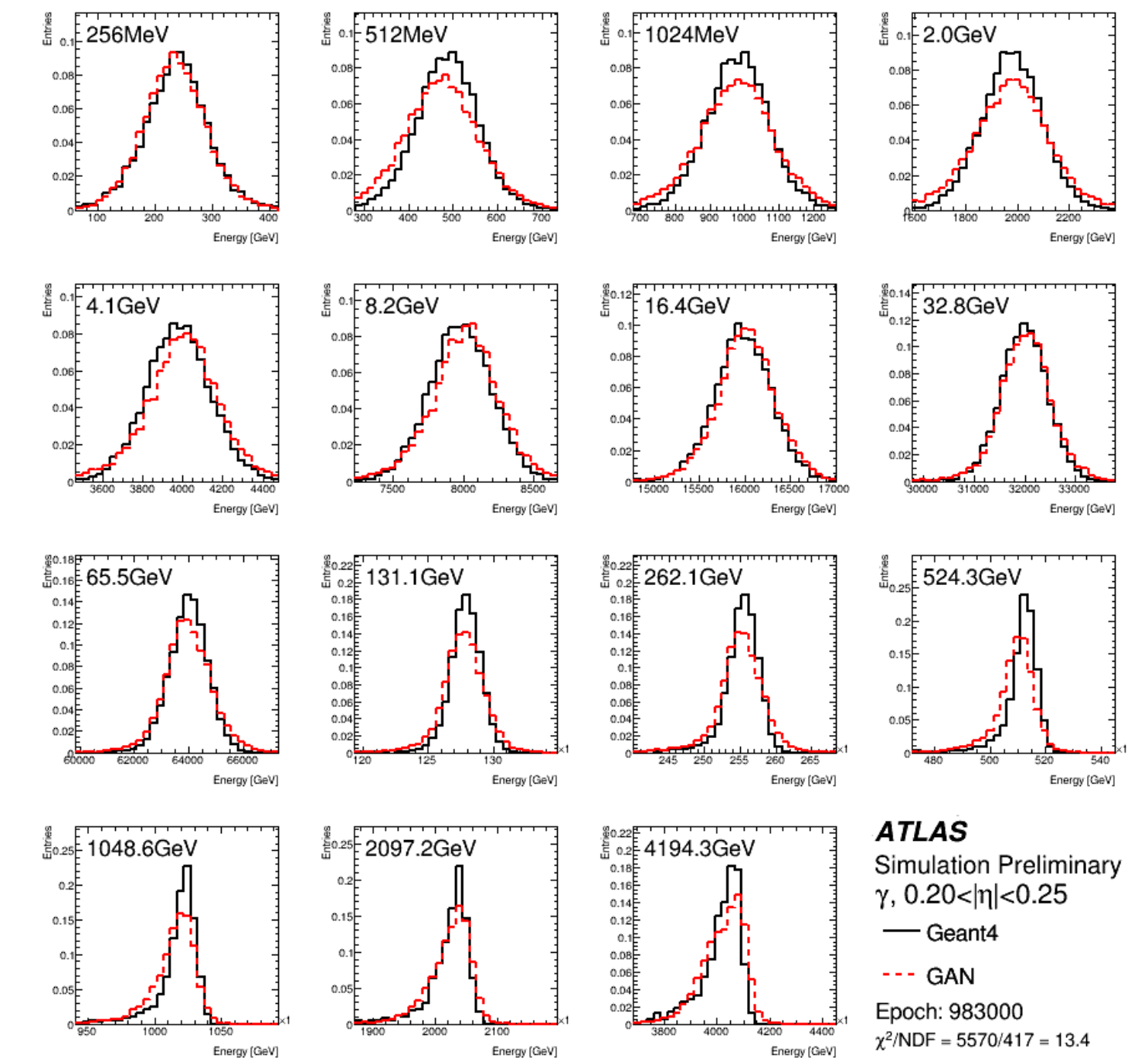
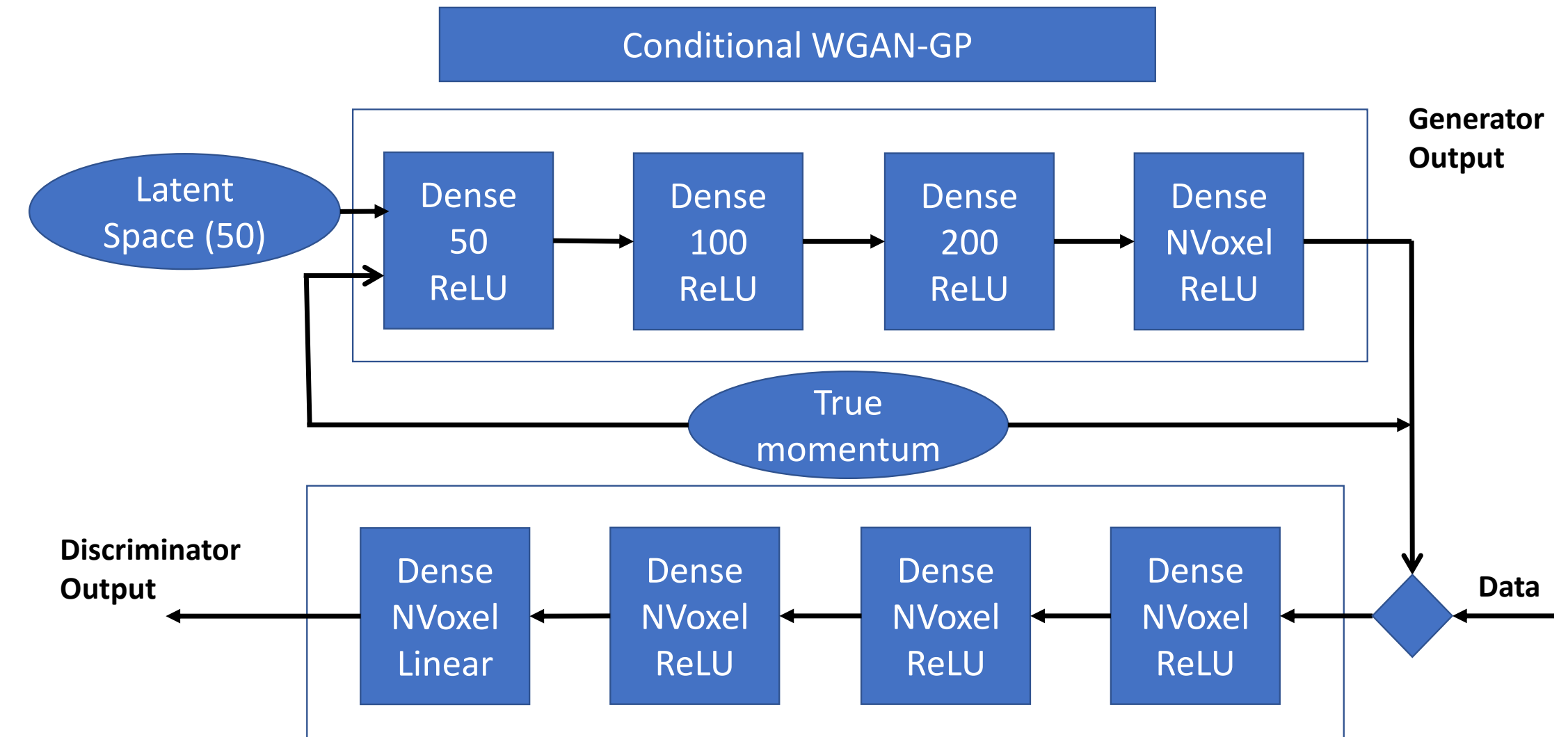
Voxelization



- ▶ Geant4-based ATLAS simulation of the full calorimeter is slow; can a GAN replace this?
- ▶ 300 GANs trained to parametrize the detector response to photons, electrons and pions



Voxelization

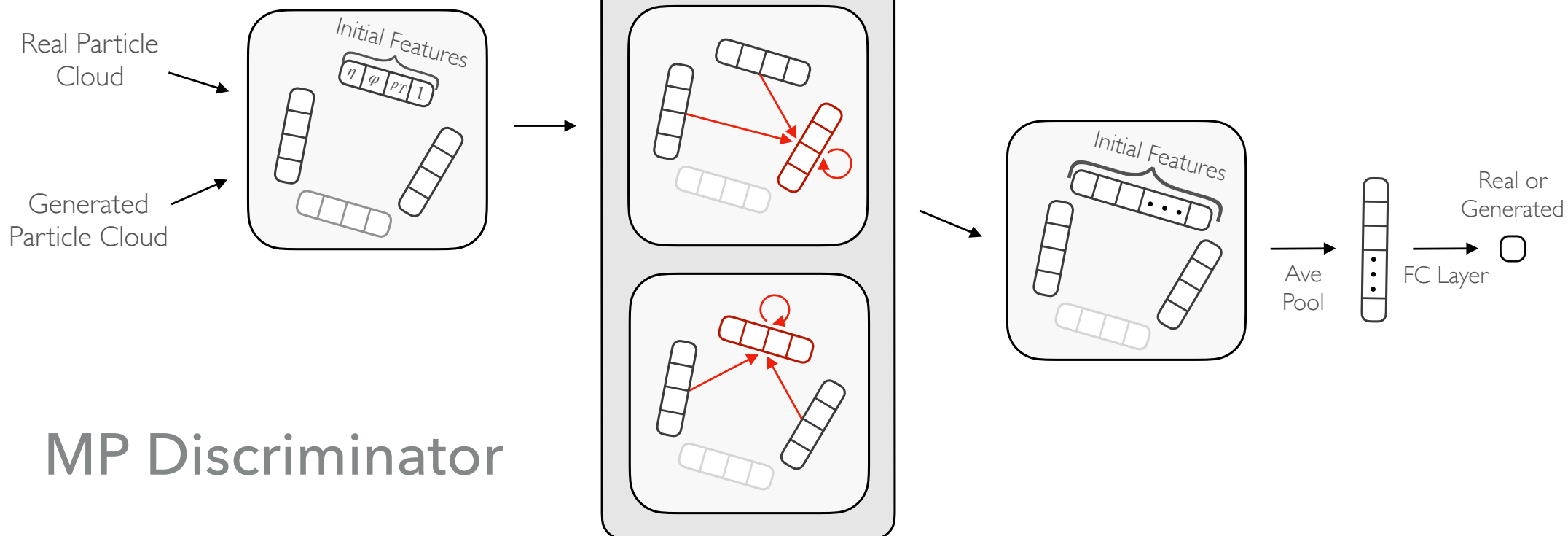
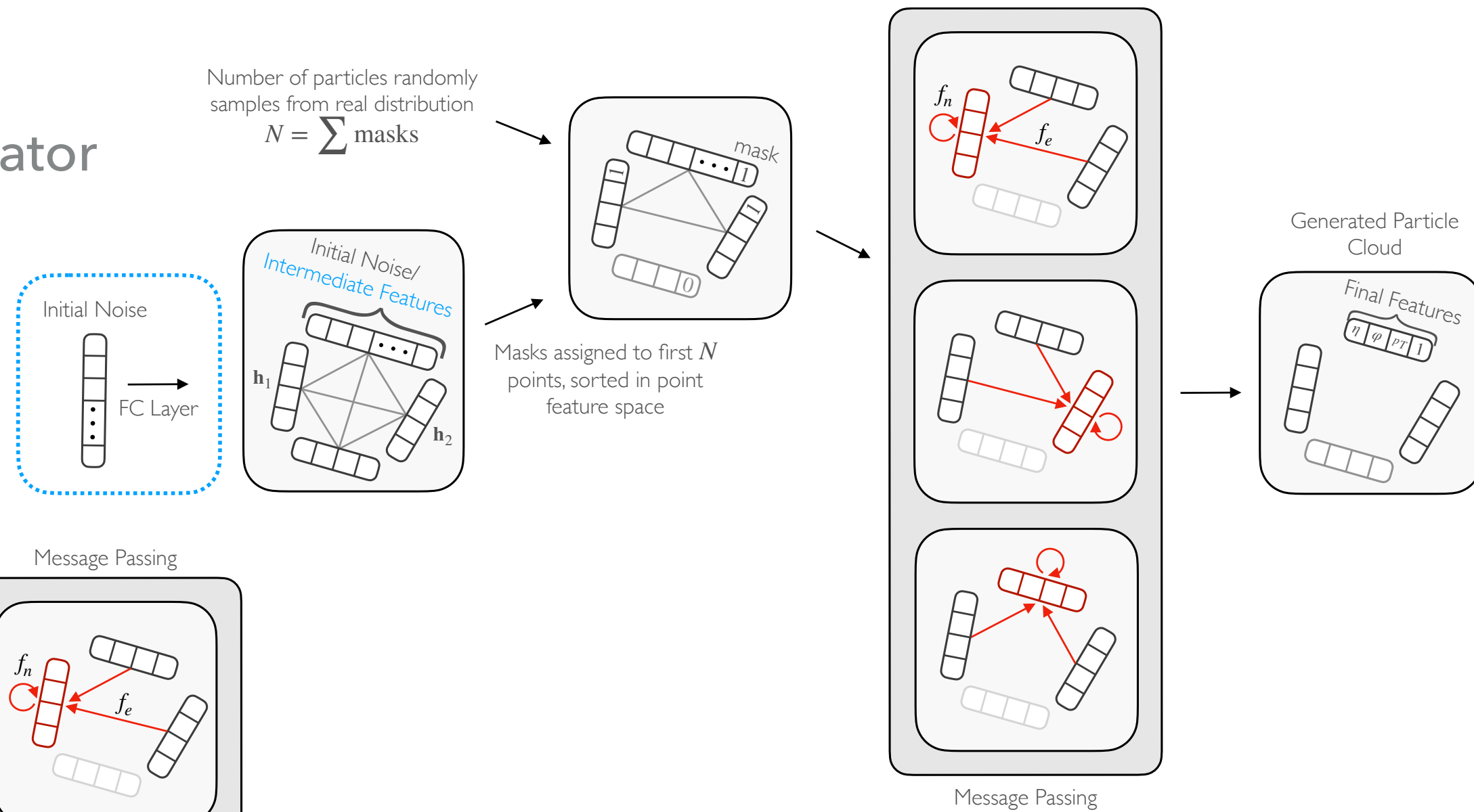


ATLAS
Simulation Preliminary
 γ , $0.20 < |\eta| < 0.25$
— Geant4
- - - GAN
Epoch: 983000
 $\chi^2/NDF = 5570/417 = 13.4$

- ▶ Geant4-based ATLAS simulation of the full calorimeter is slow; can a GAN replace this?
- ▶ 300 GANs trained to parametrize the detector response to photons, electrons and pions
- ▶ Good agreement between the GAN and Geant4 both for single-particle showers and jets

- As an alternative to voxelization, a graph-based GAN can be used to generate jets as particle clouds

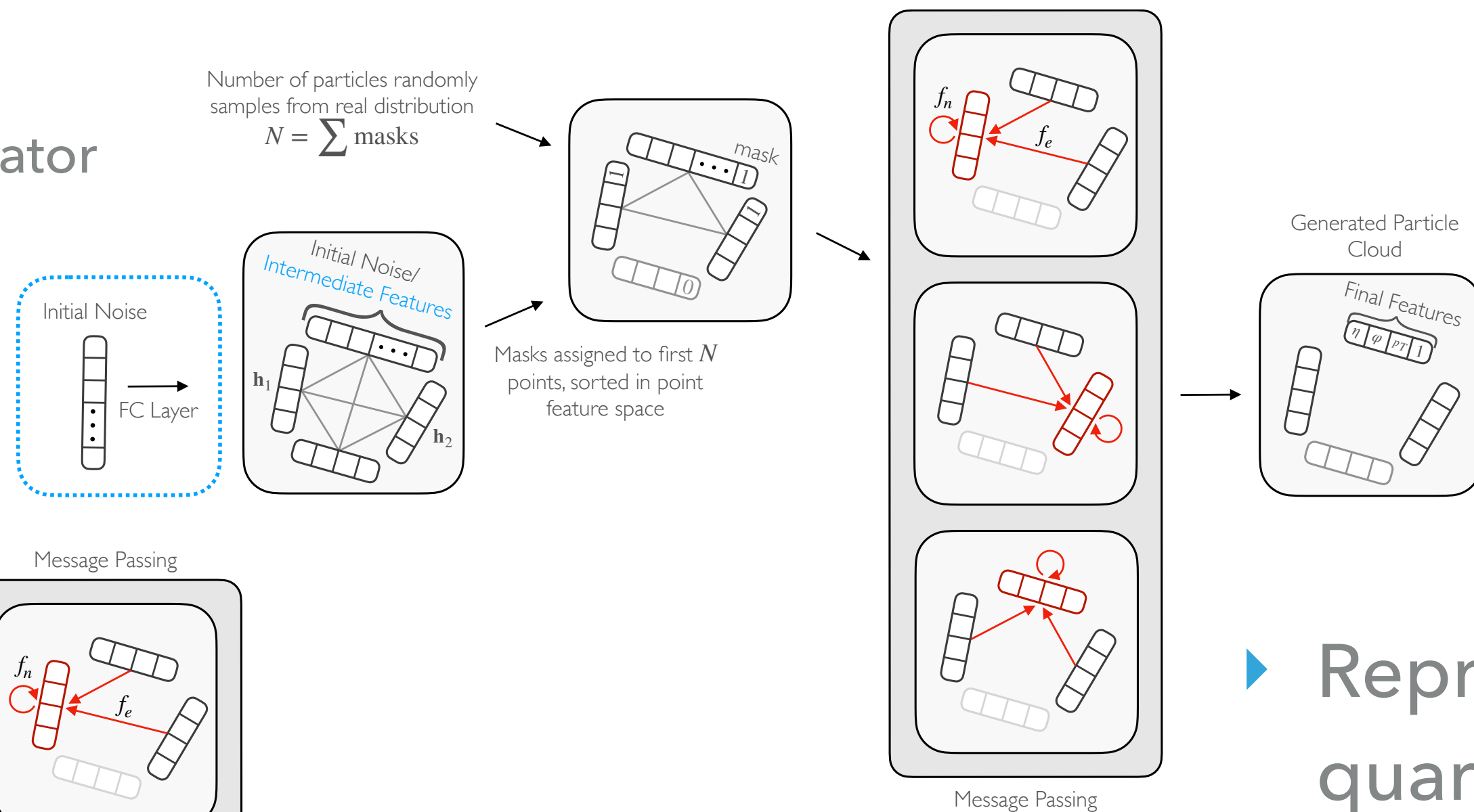
MP-(LFC) Generator



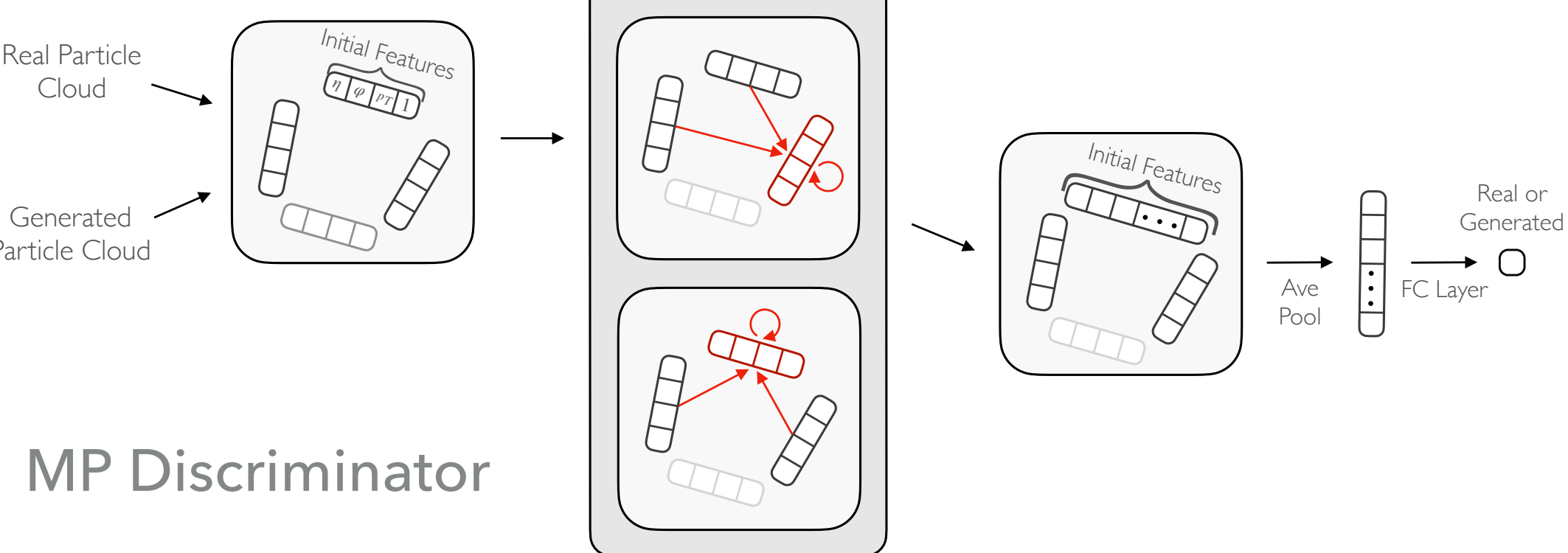
MP Discriminator

- As an alternative to voxelization, a graph-based GAN can be used to generate jets as particle clouds

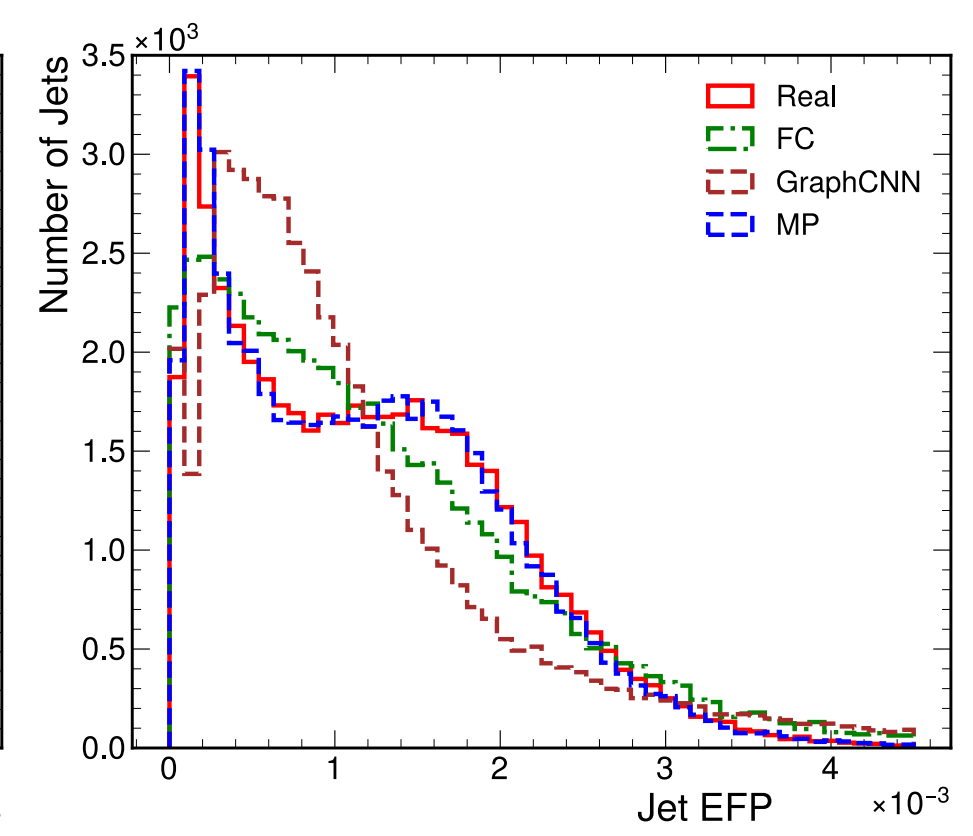
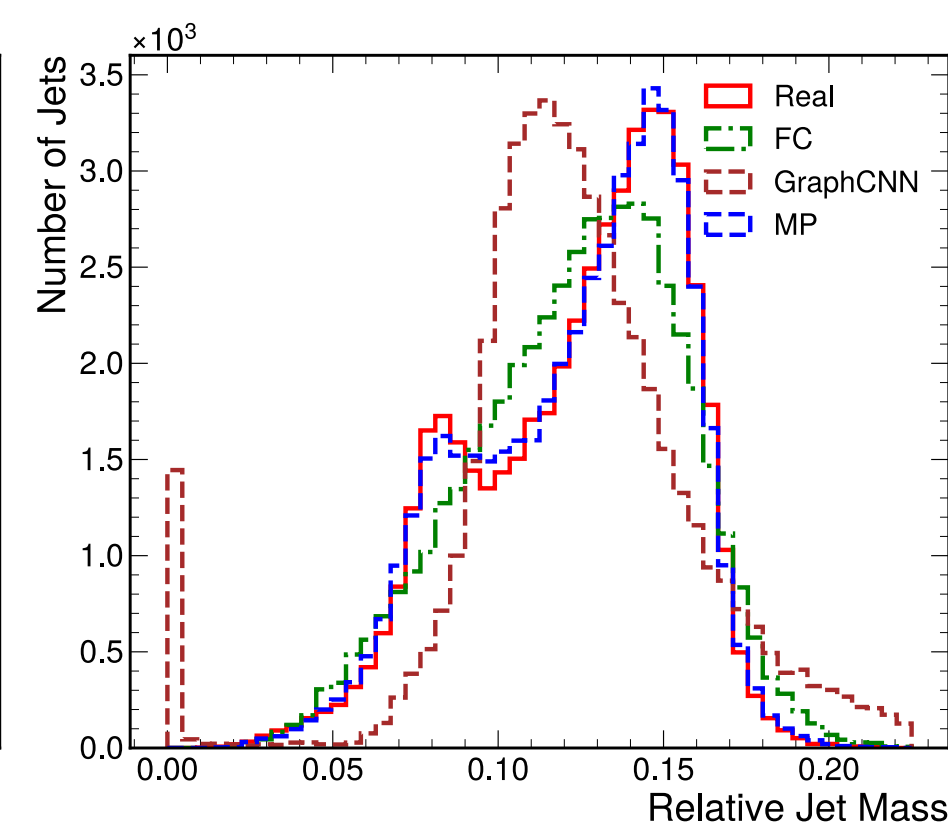
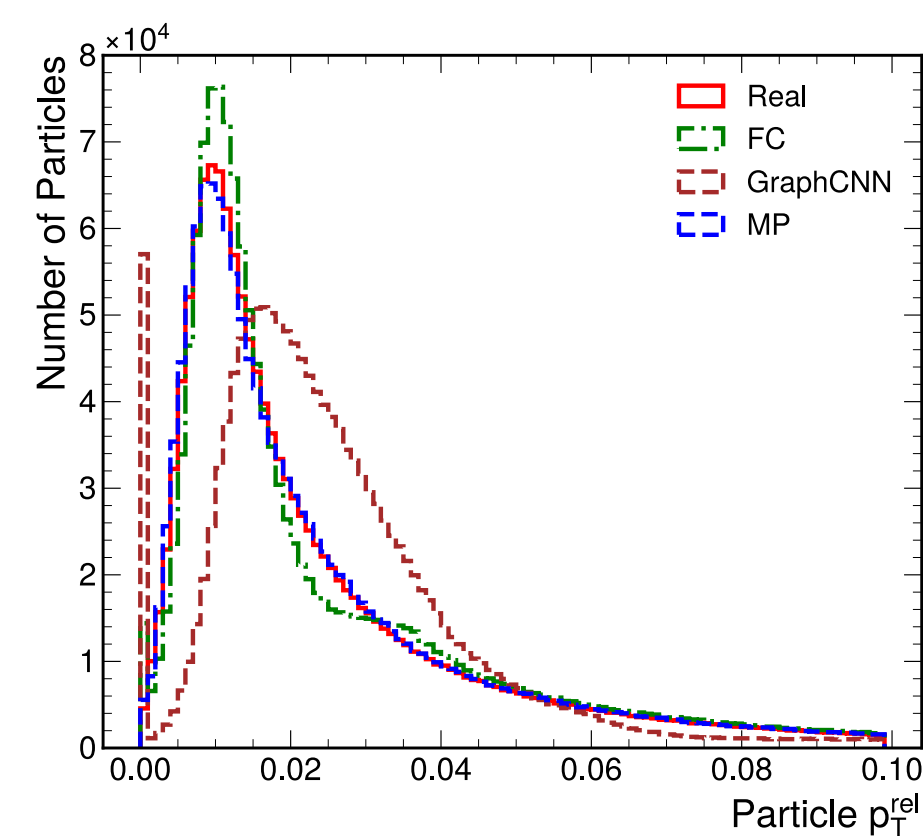
MP-(LFC) Generator



- Reproduces nontrivial properties like top quark jet mass and energy-flow polynomials



MP Discriminator

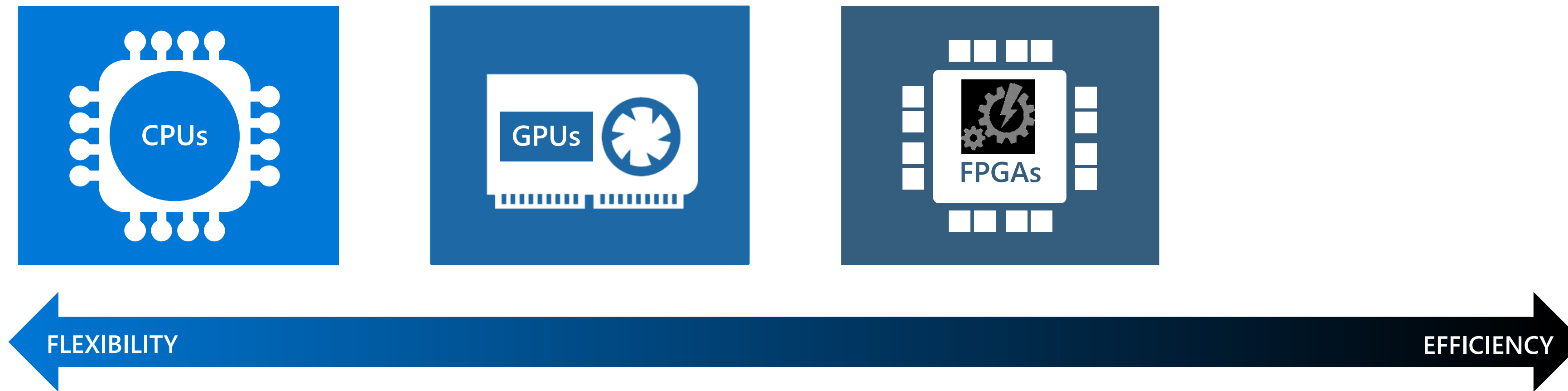


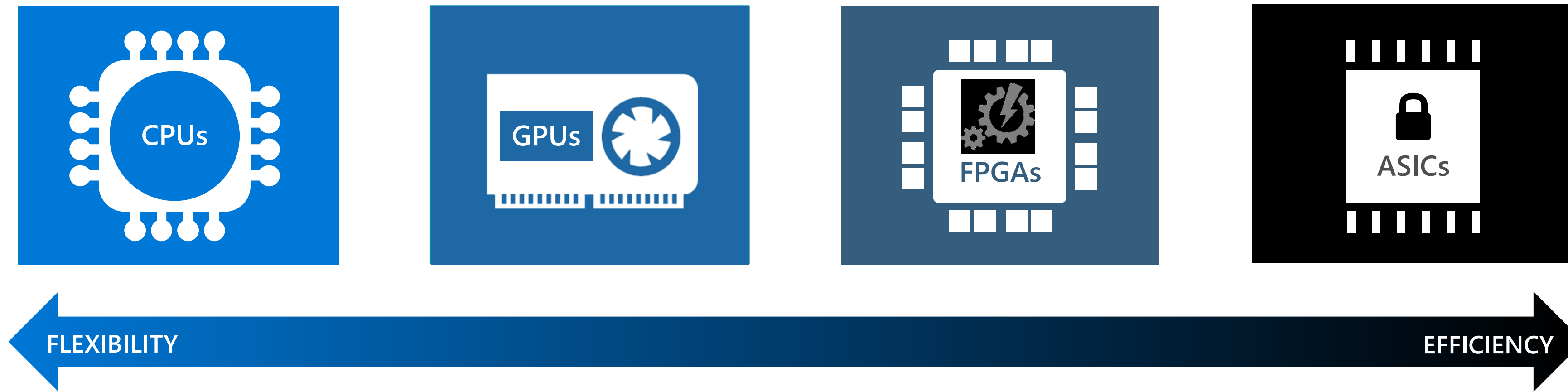


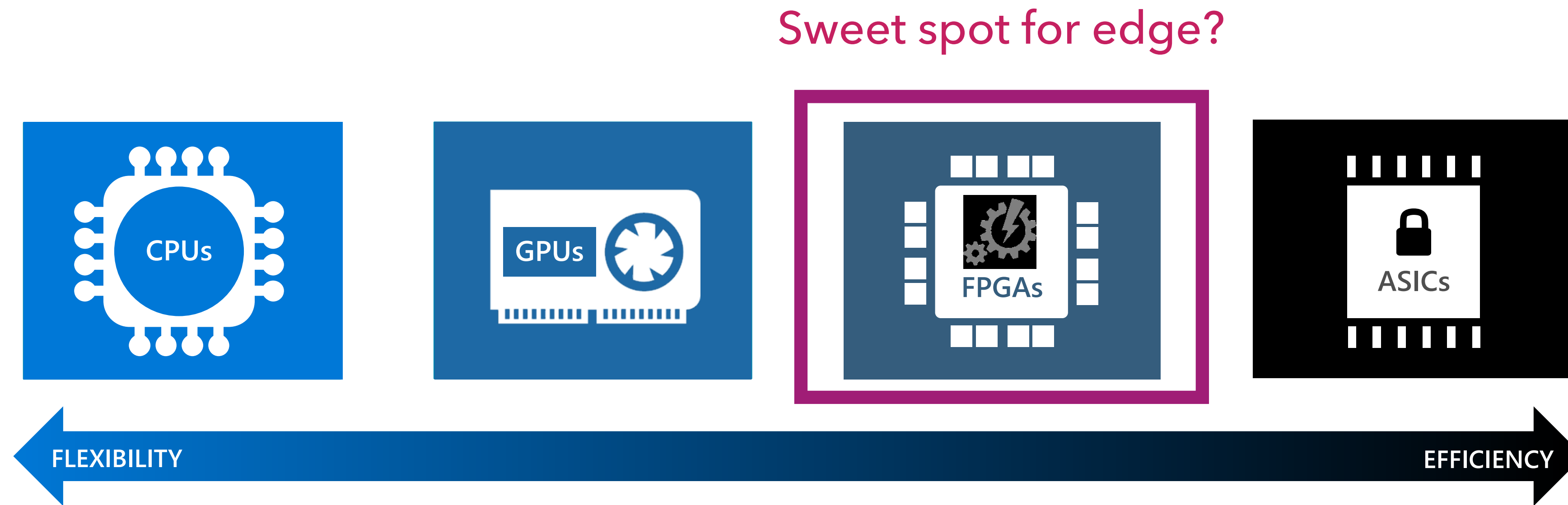
ML FOR TAGGING
ML FOR GEN/SIM
FAST ML FOR TRIGGER



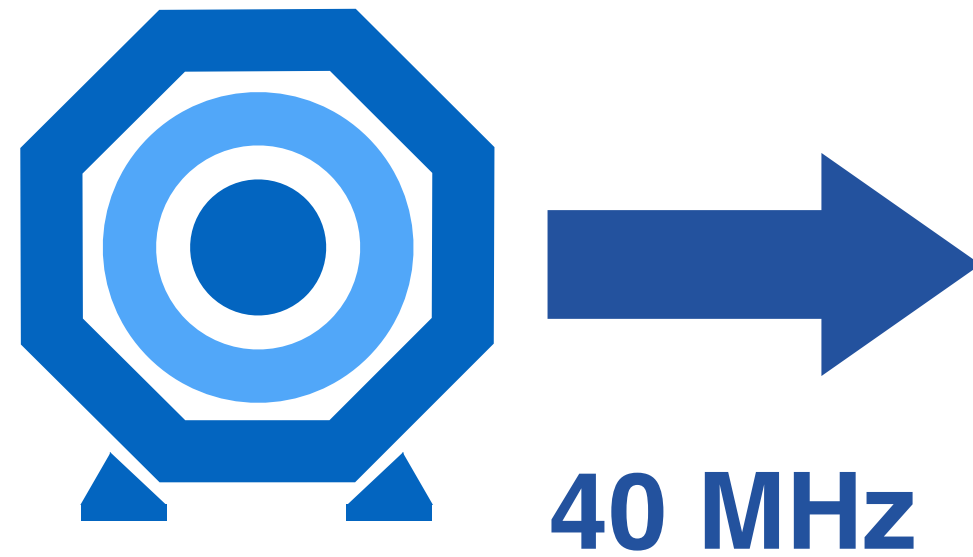
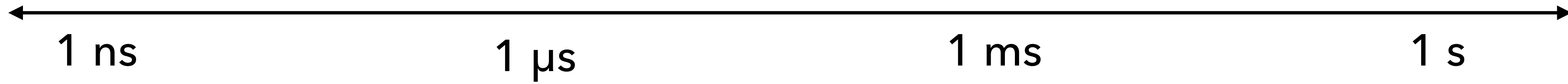








Compute
Latency



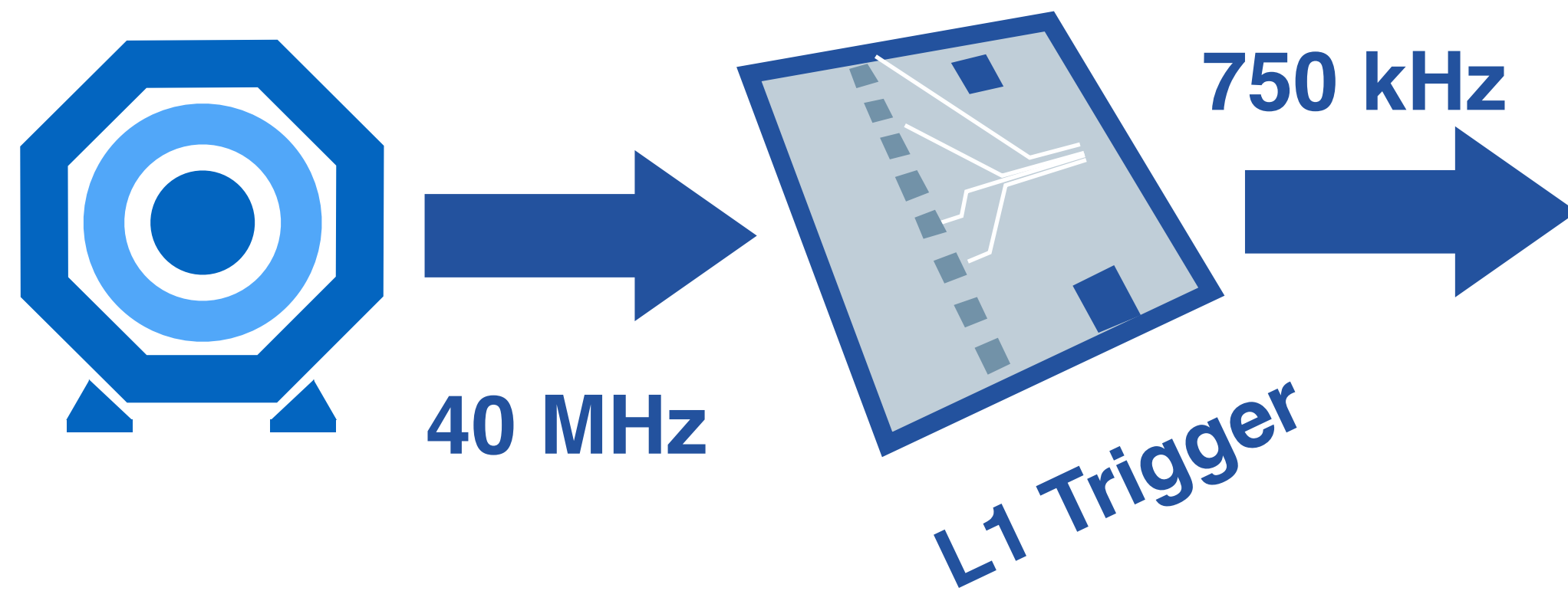
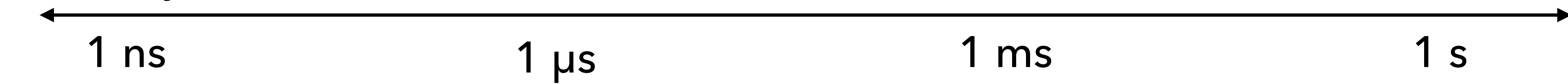
Challenges:

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

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

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

Compute
Latency



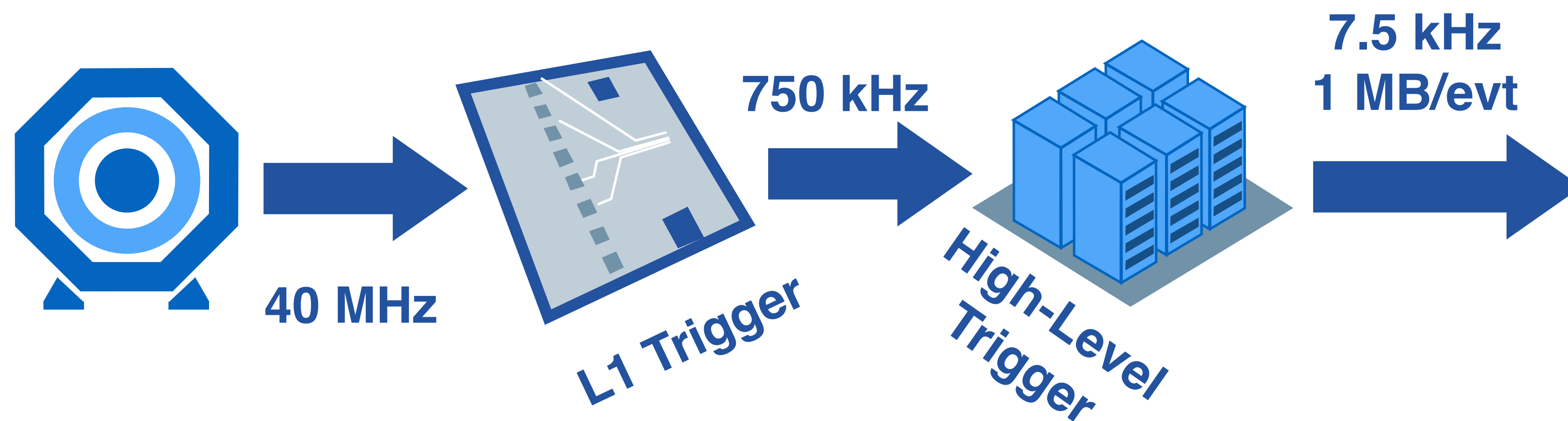
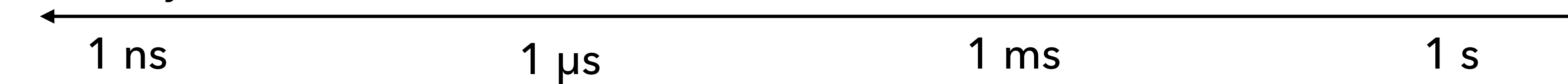
Challenges:

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

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

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

Compute
Latency



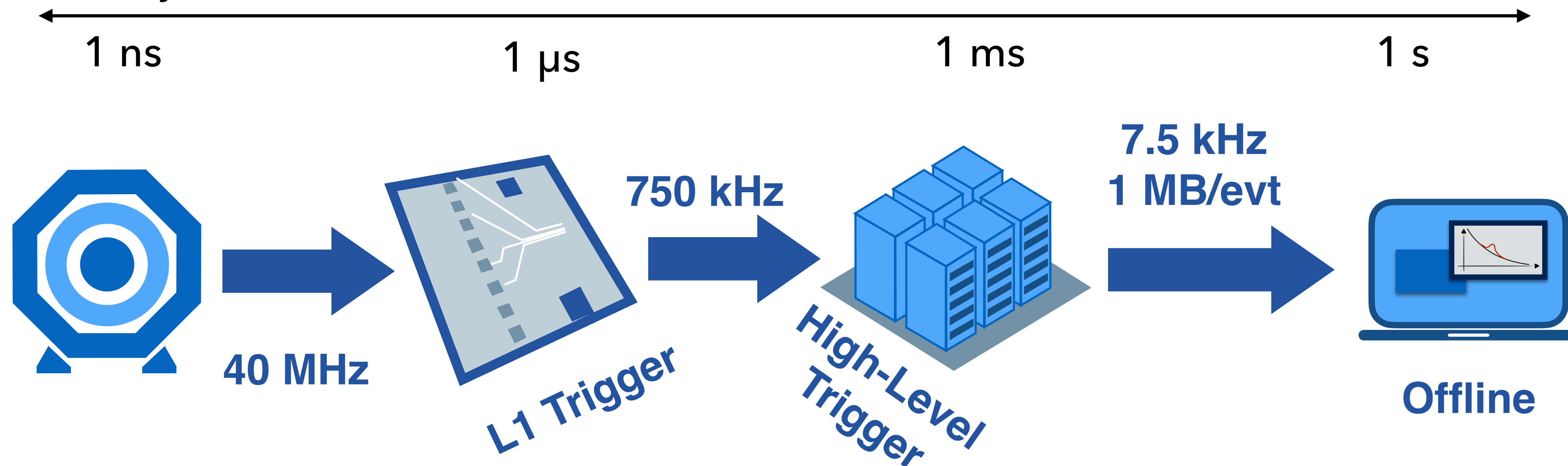
Challenges:

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

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

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

Compute
Latency



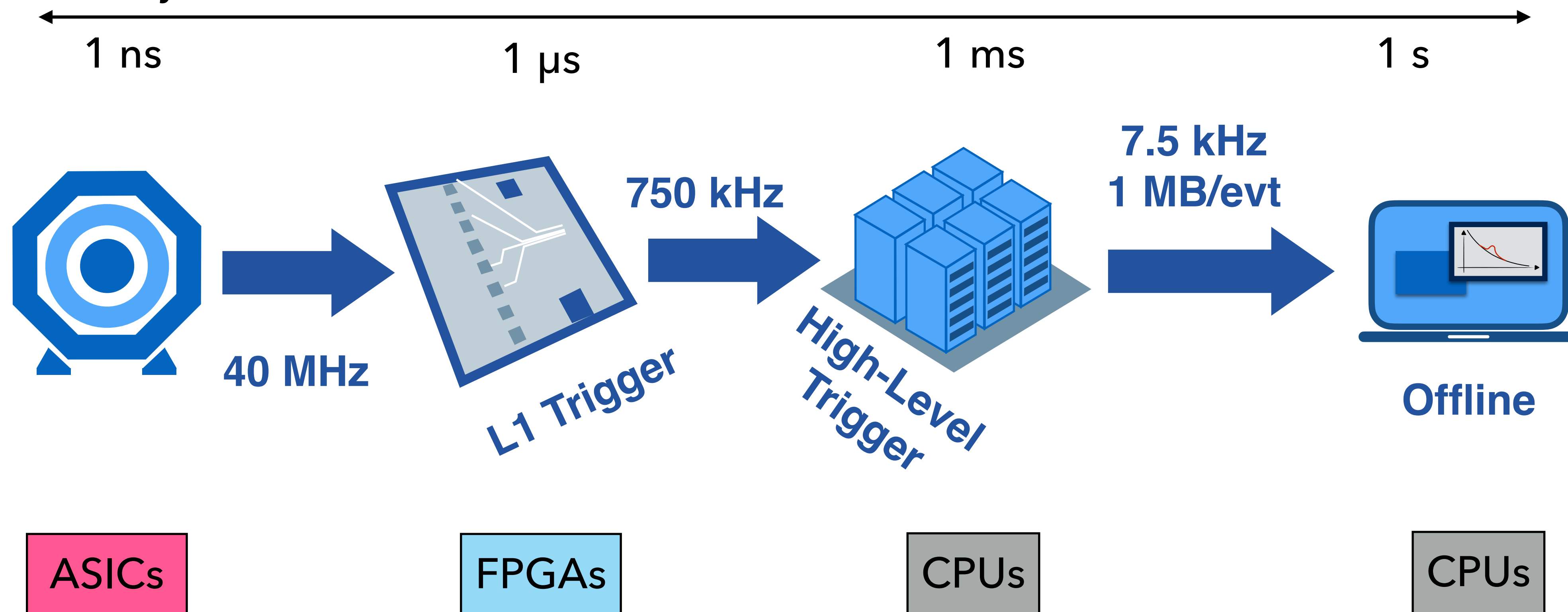
Challenges:

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

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

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

Compute
Latency

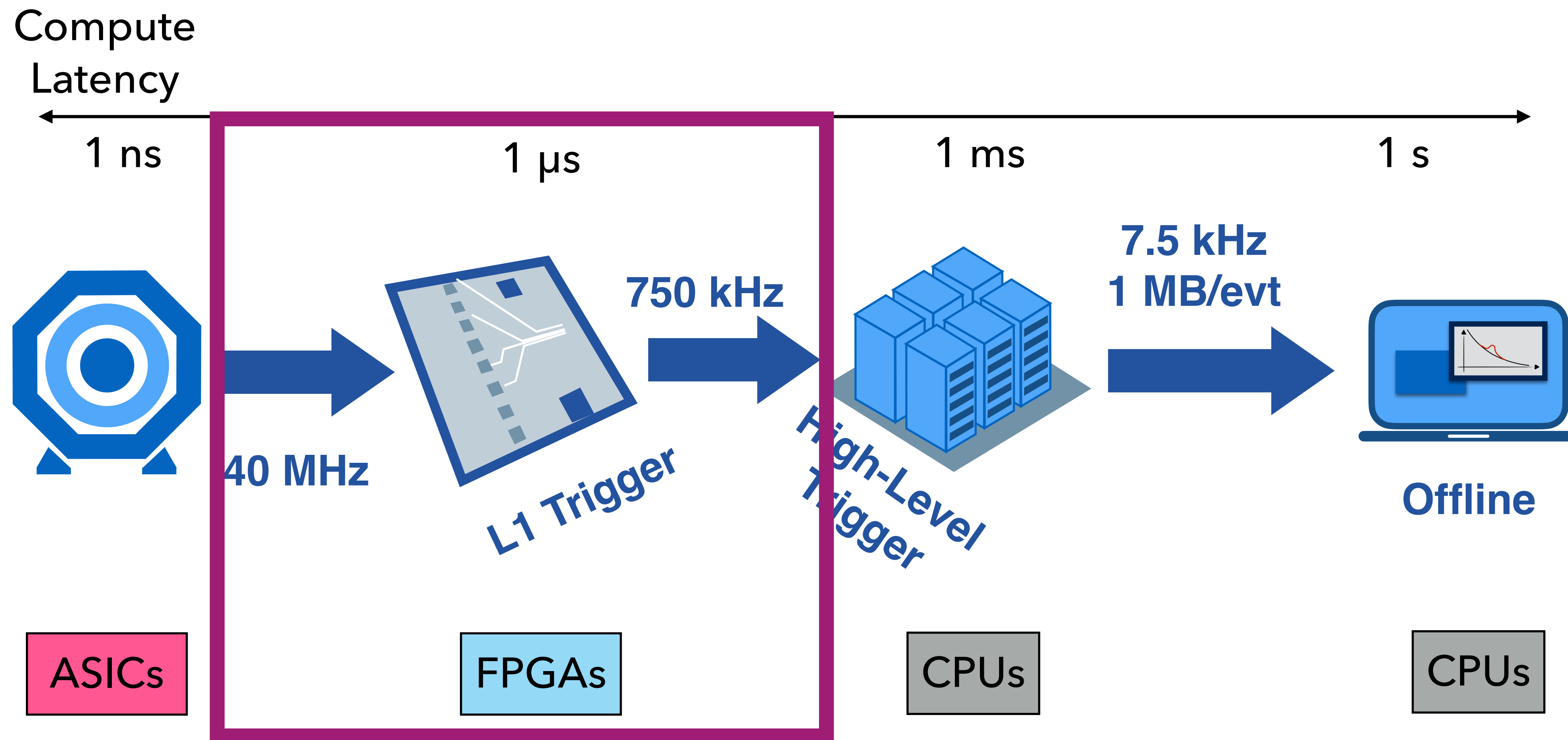


Challenges:

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

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

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



Challenges:

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

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

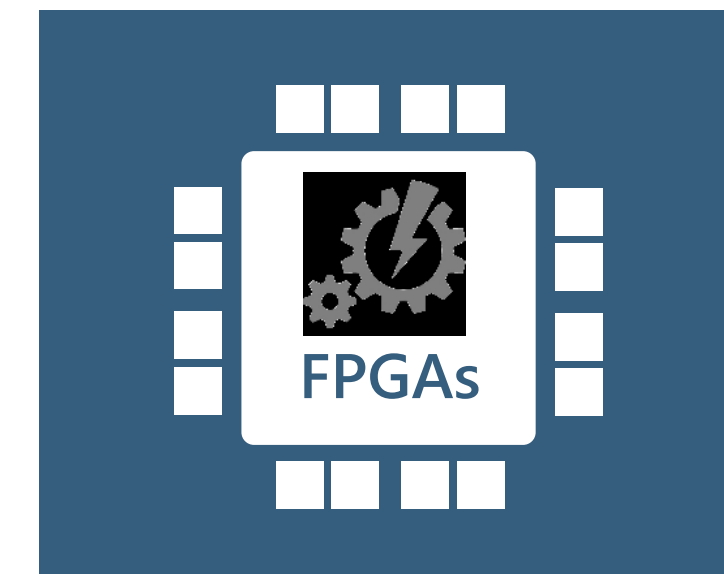
Extreme data rates of $O(100 \text{ TB/s})$

- ▶ 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;
```

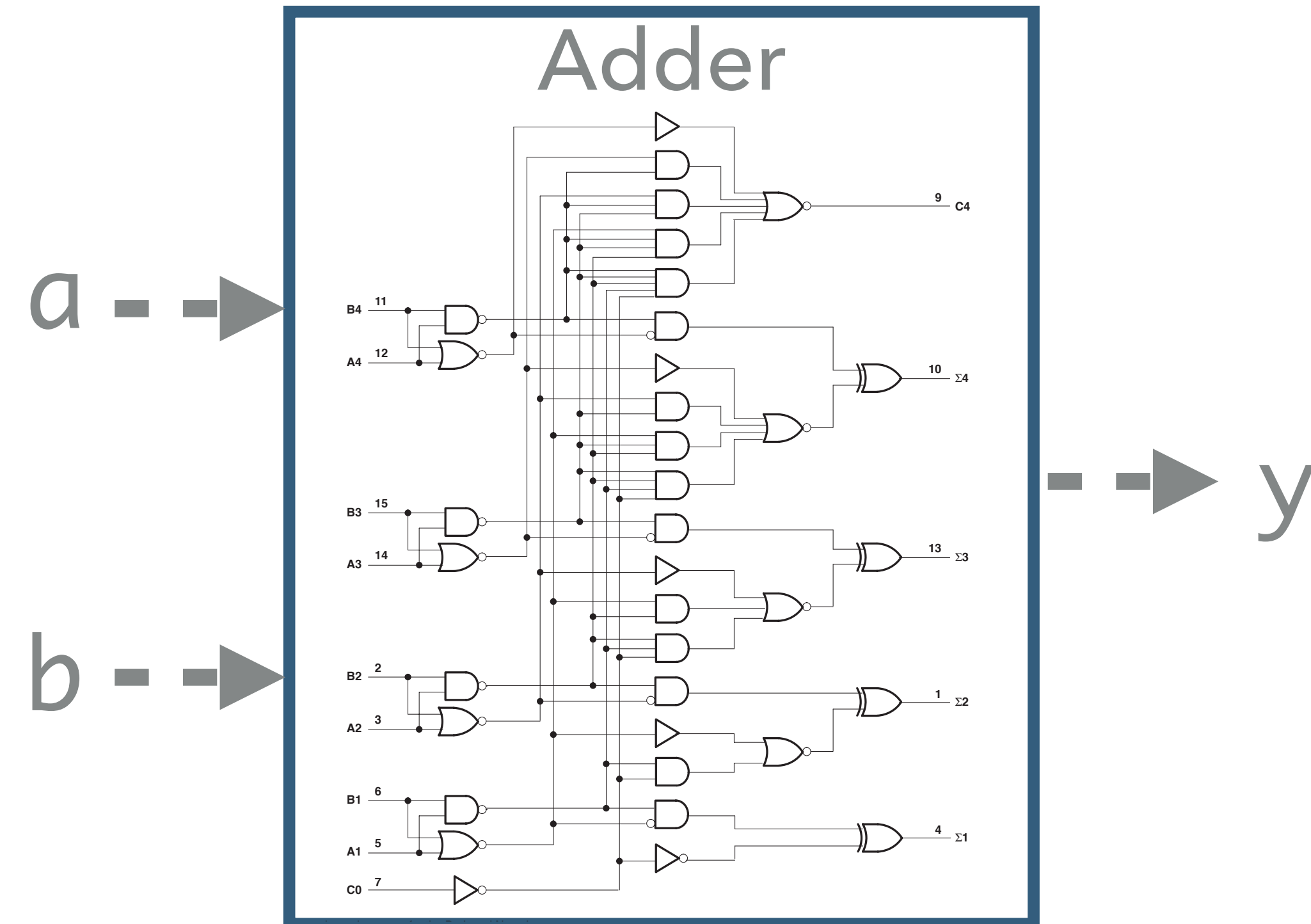
endmodule

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



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

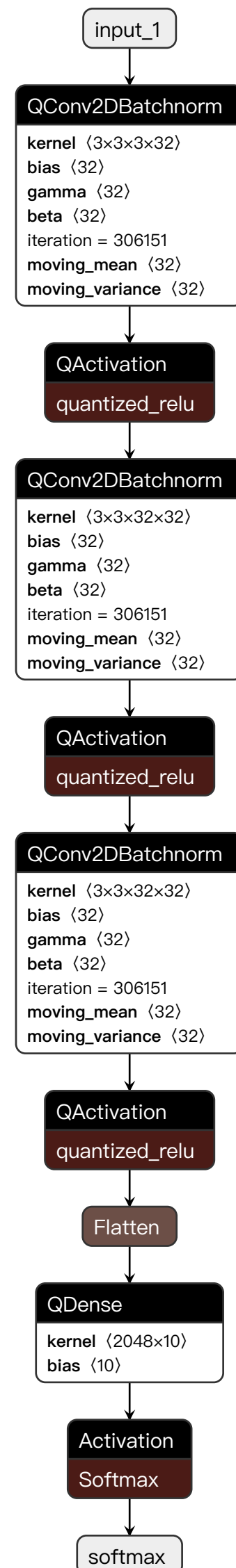
```
module adder(  
    input  wire [4:0] a,  
    input  wire [4:0] b,  
    output wire [4:0] y  
);  
    assign y = a + b;  
endmodule
```



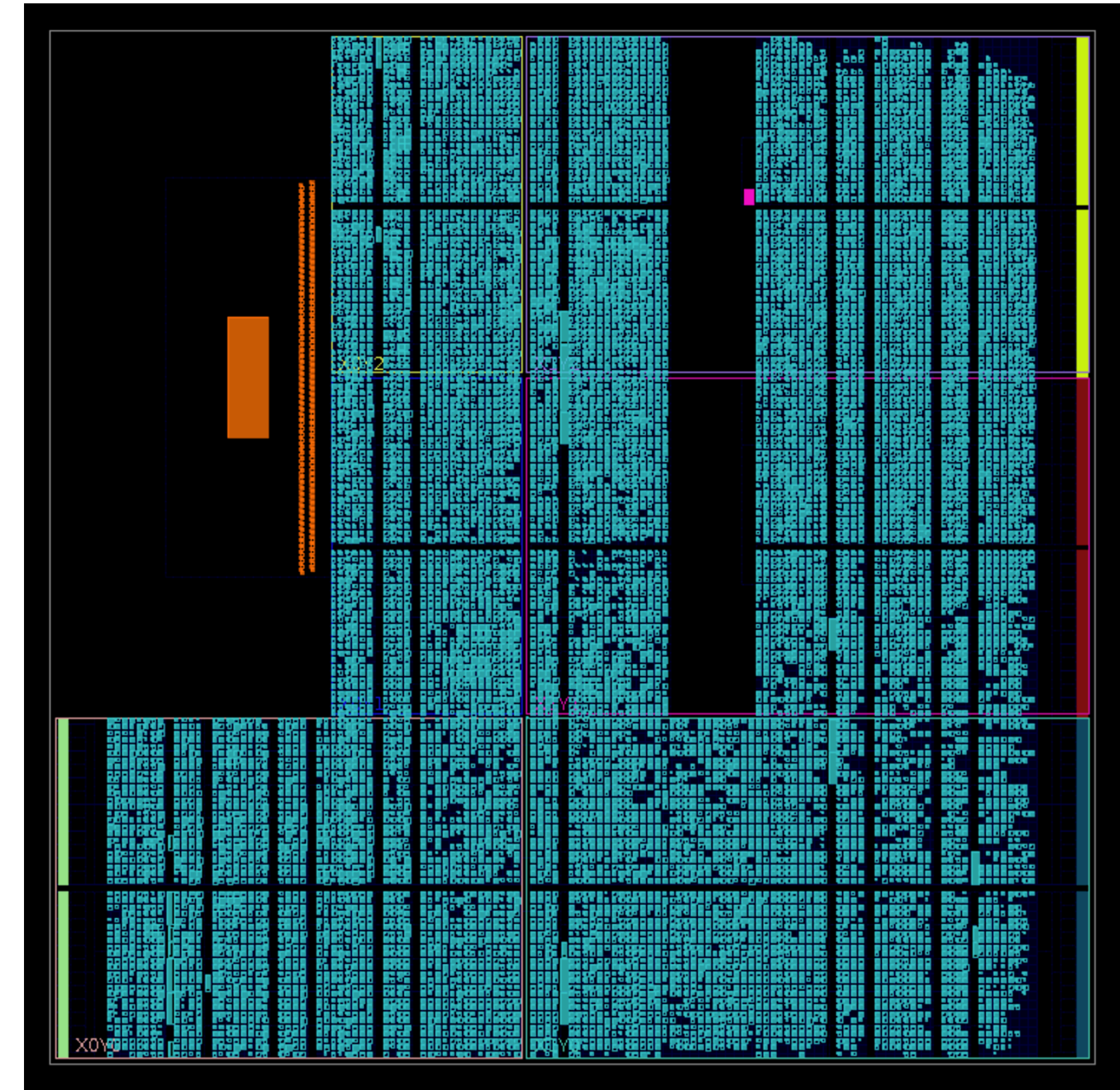
- ▶ Register transfer-level (RTL) code is "synthesized" into gates

Synthesis

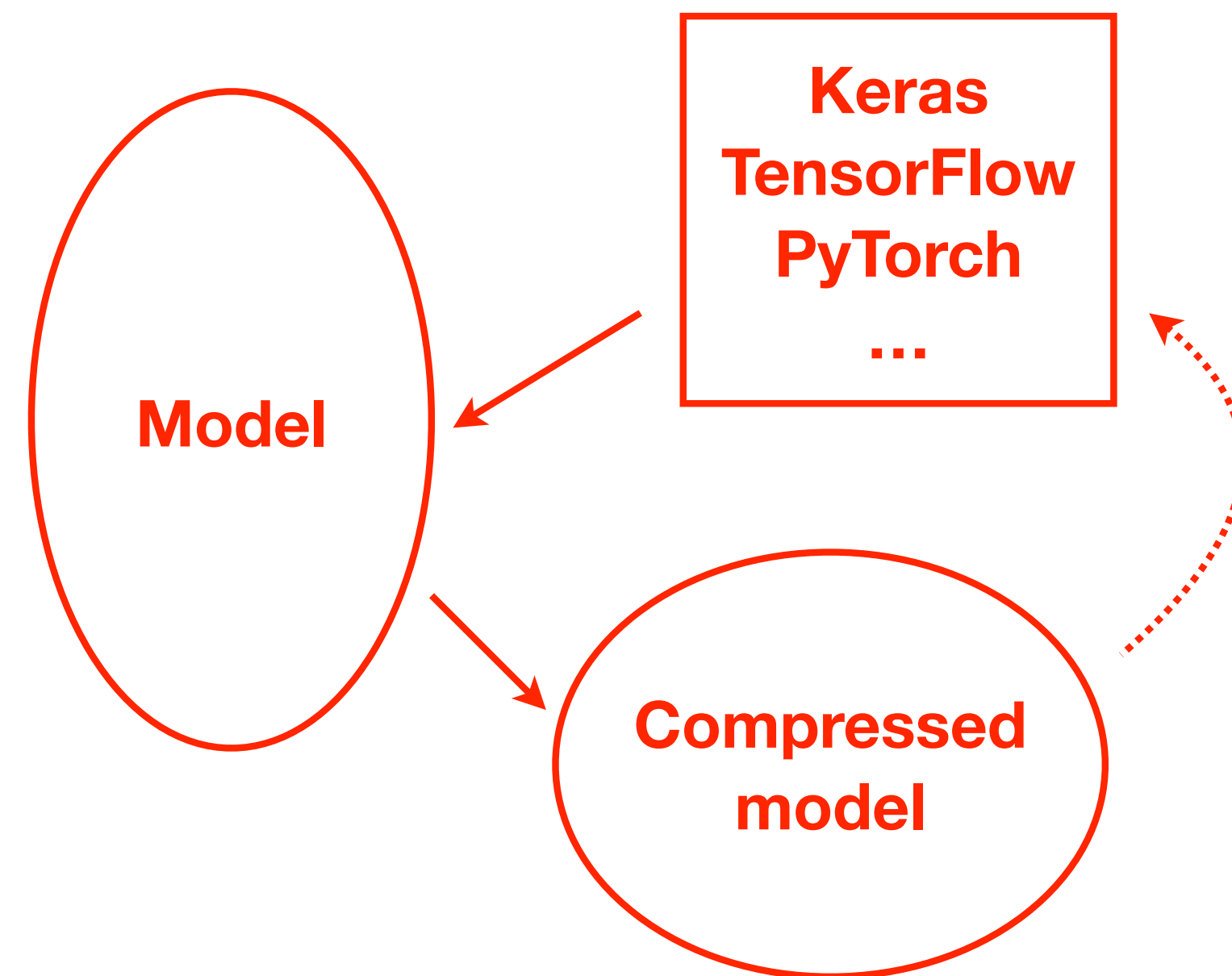
- ▶ What if instead we specify an AI model



High-Level Synthesis

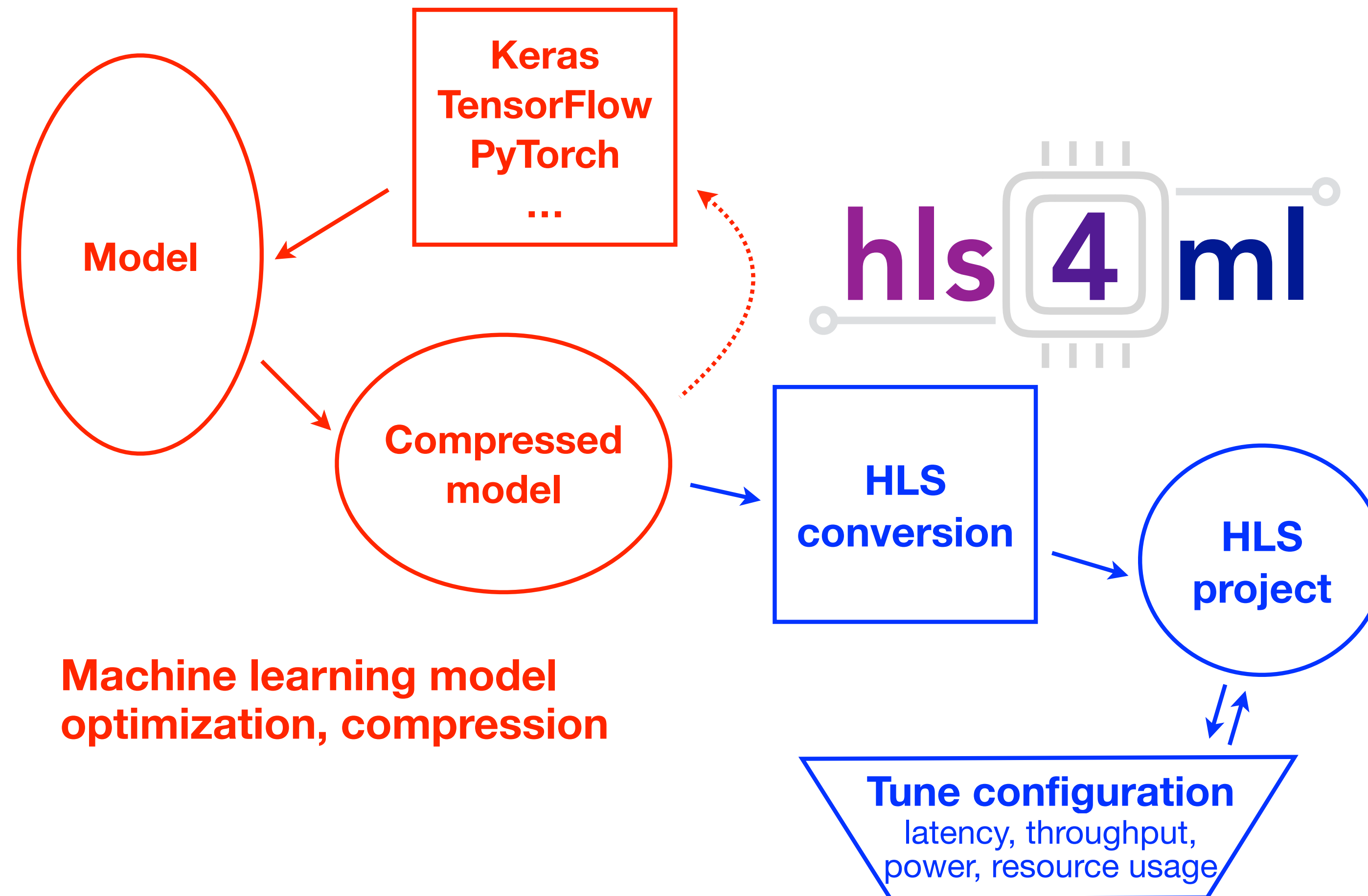


- ▶ [hls4ml](#) for scientists or ML experts to translate ML algorithms into RTL firmware

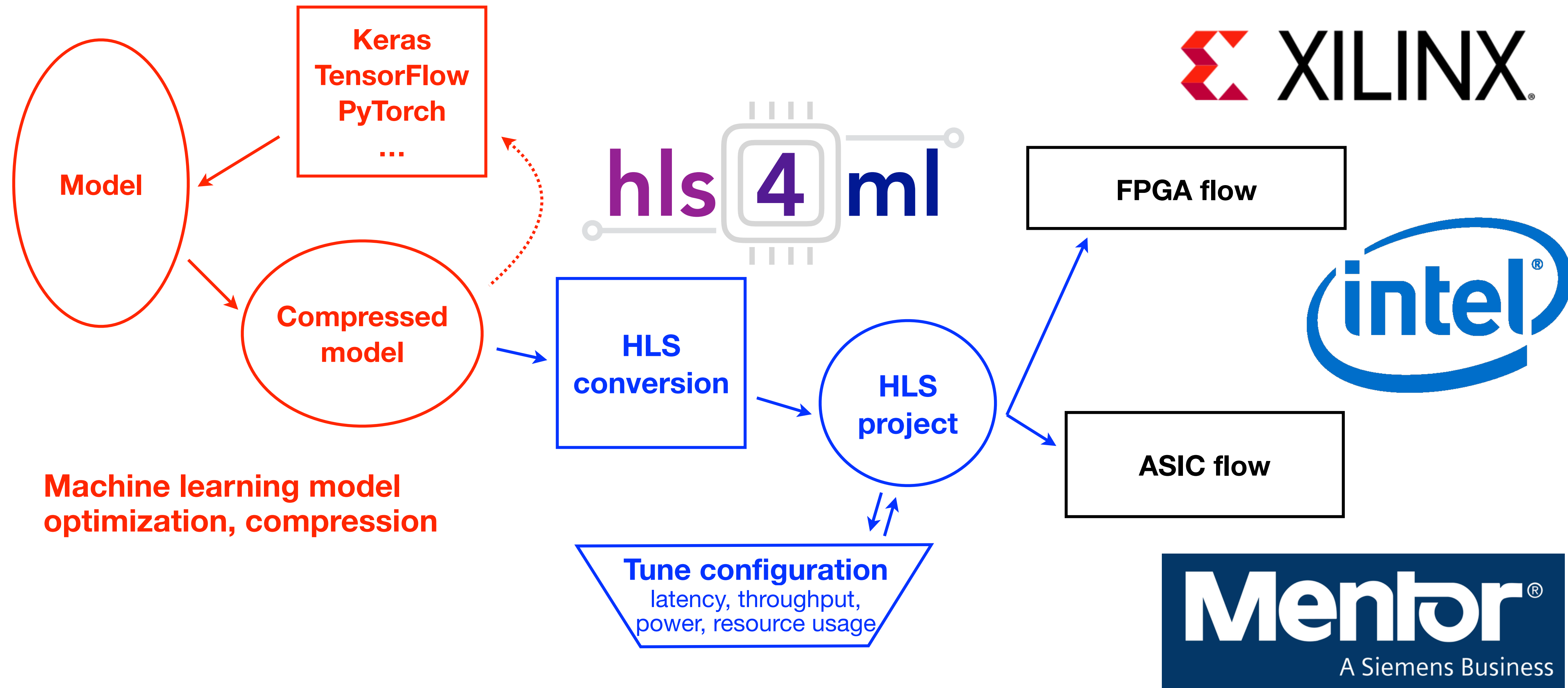


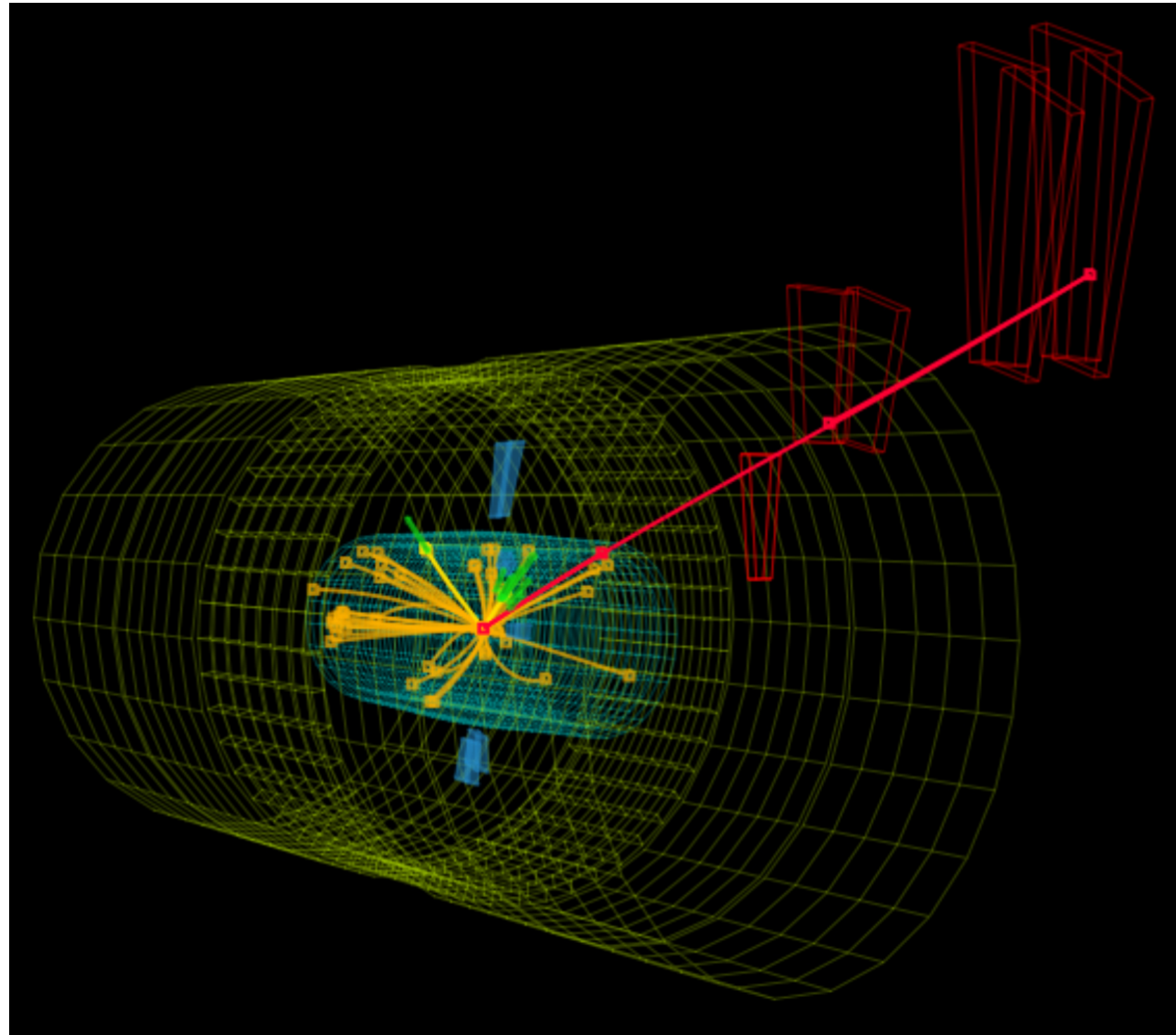
**Machine learning model
optimization, compression**

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

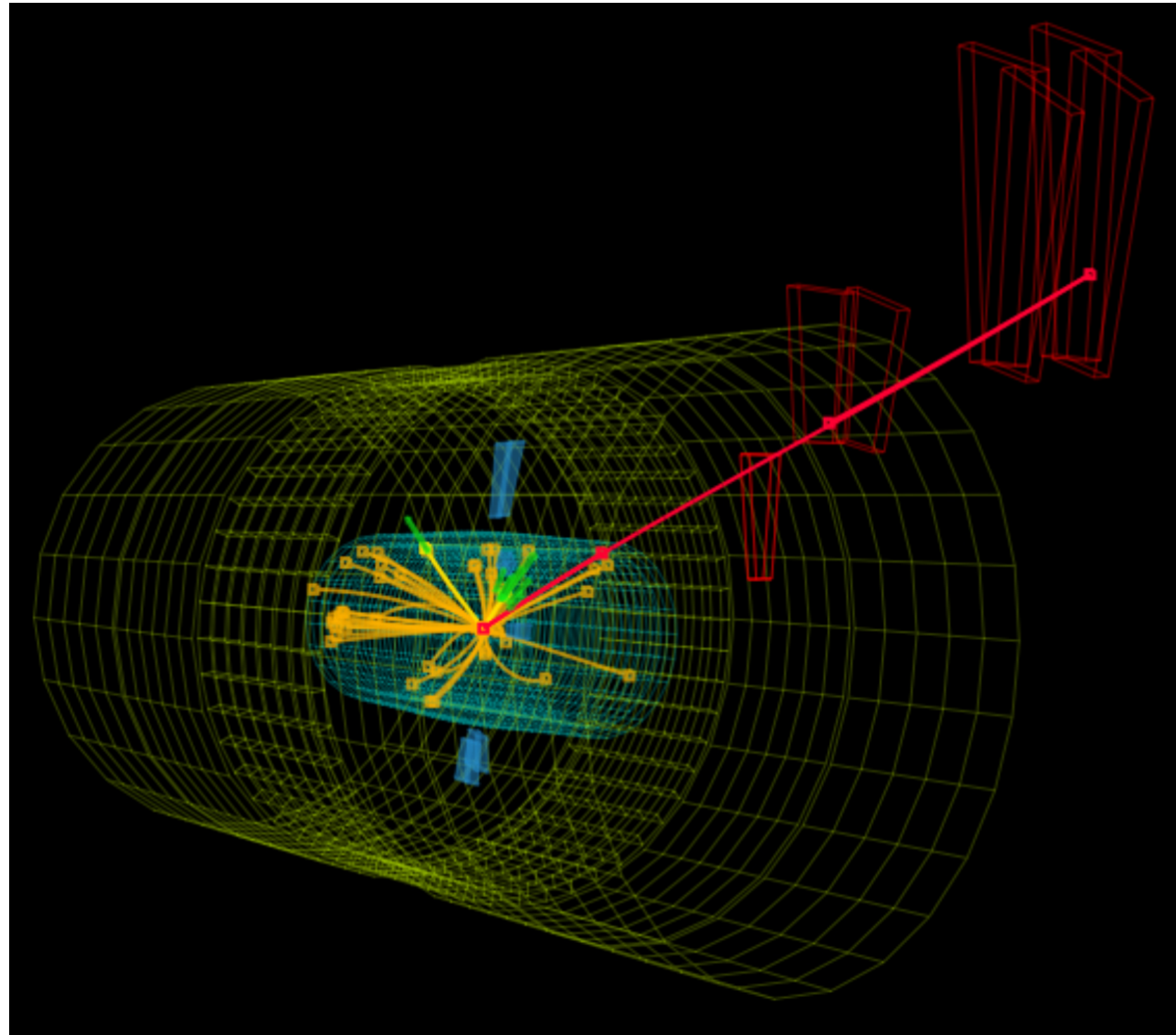


- ▶ [hls4ml](#) for scientists or ML experts to translate ML algorithms into RTL firmware

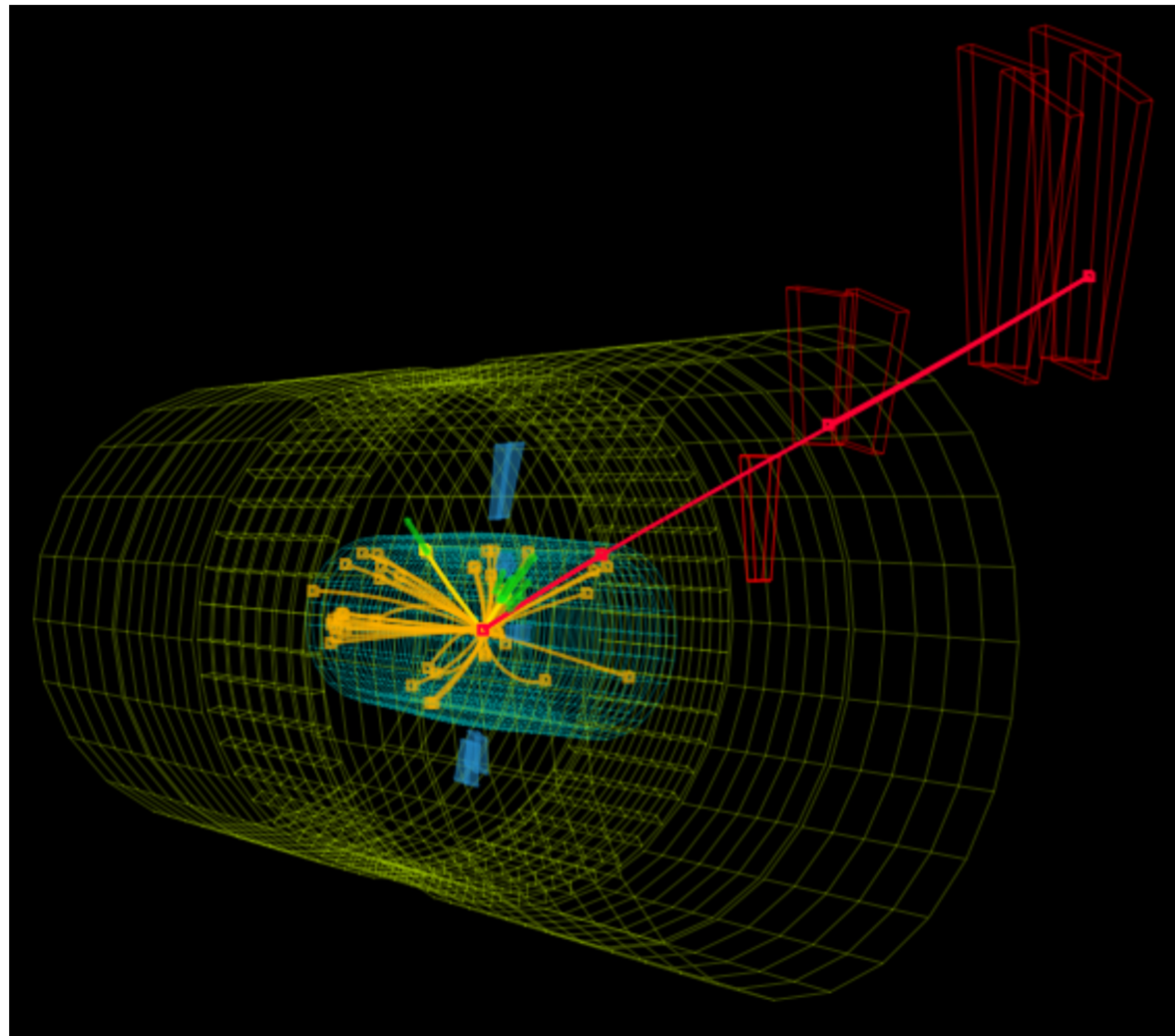




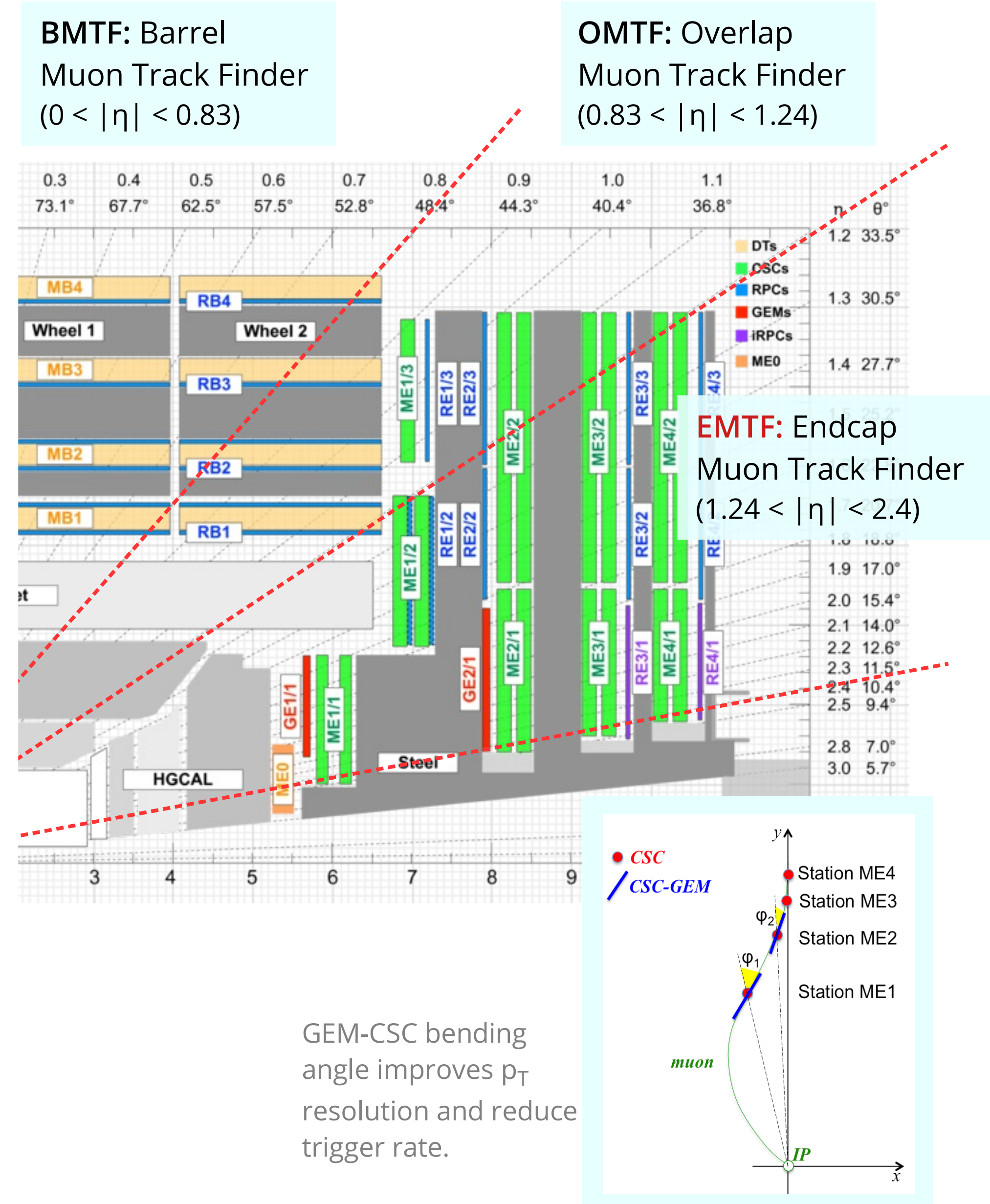
- ▶ Goal: determine muon p_T in endcap based on info. available in the L1 trigger

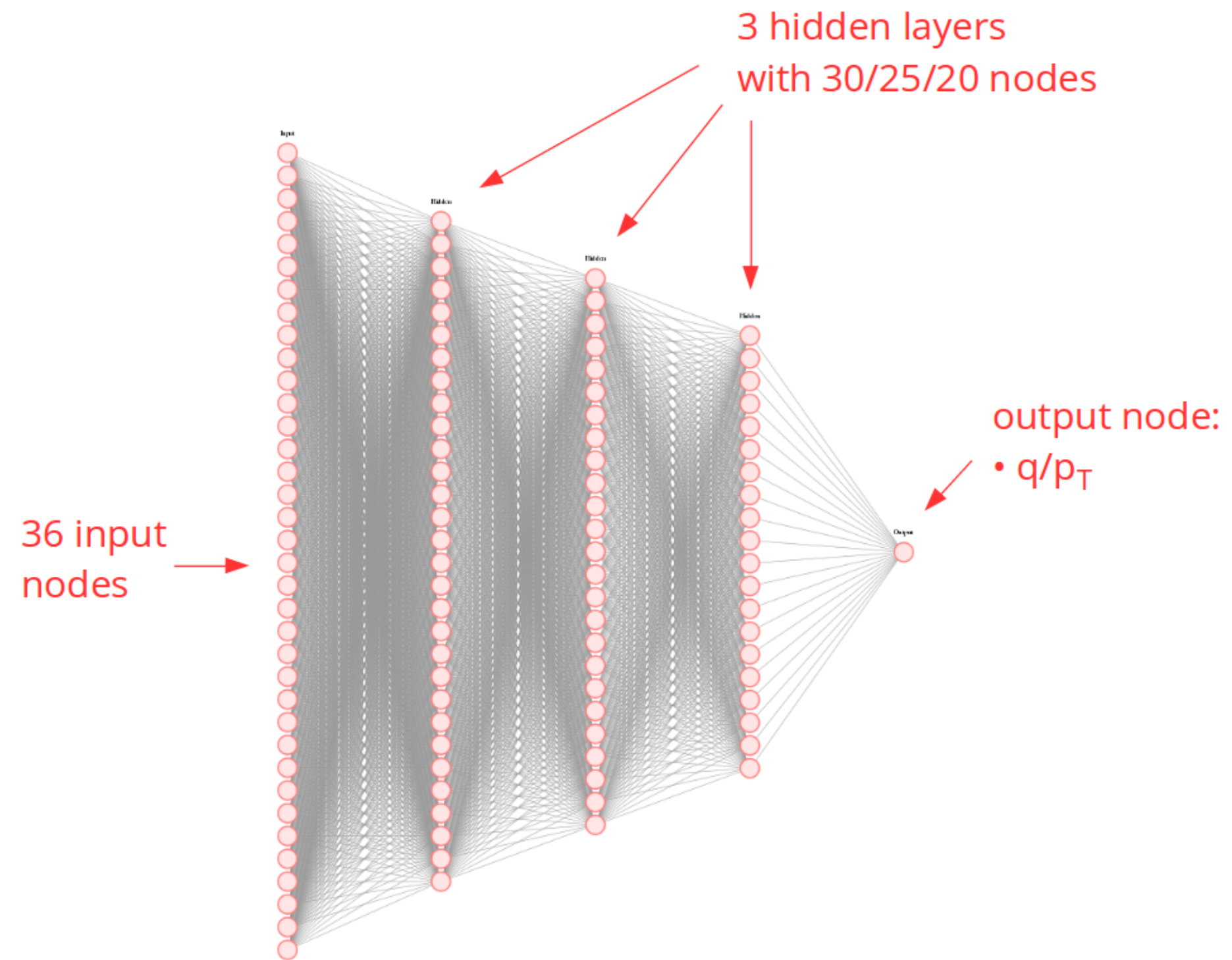


- ▶ Goal: determine muon p_T in endcap based on info. available in the L1 trigger



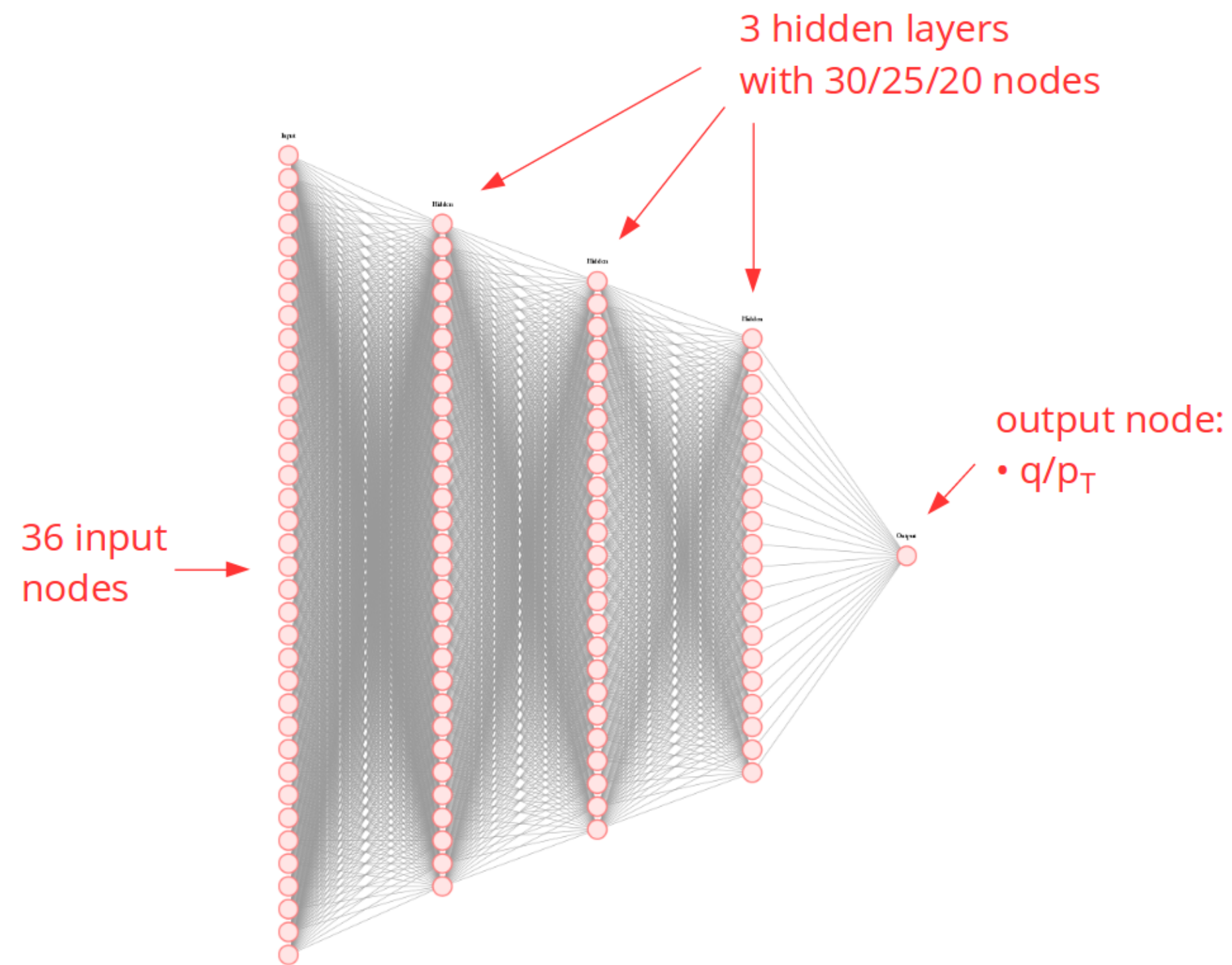
- ▶ Challenges:
 - ▶ Non-uniform magnetic field with little bending
 - ▶ Large background from multiple sources
- ▶ EMTF++ has to evolve to
 - ▶ Incorporate new muon detectors
 - ▶ Improve efficiency, redundancy, p_T resolution, timing
 - ▶ Maintain the same trigger threshold at higher pileup



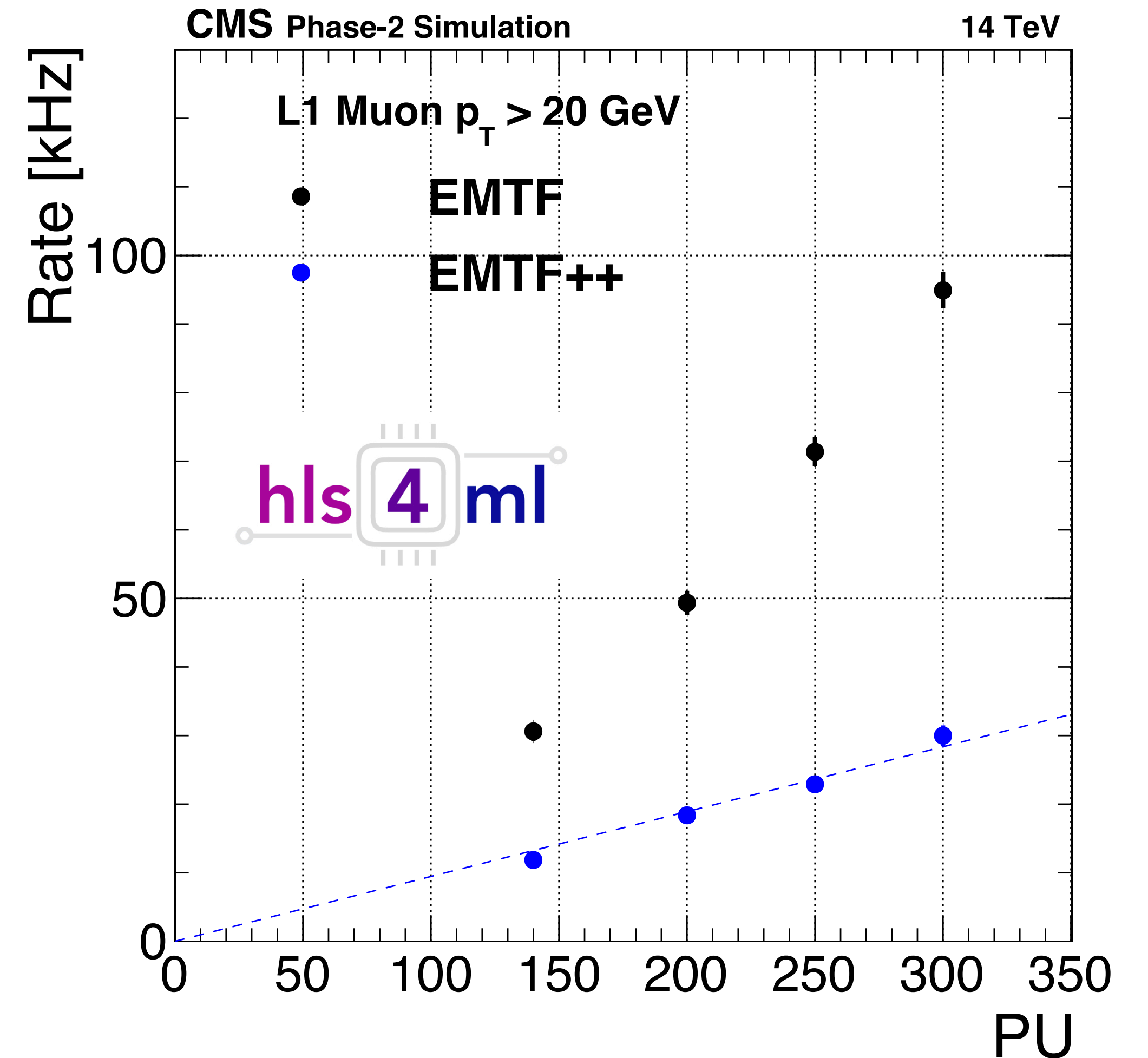


	ME1/1	ME1/2	ME2	ME3	ME4	RE1	RE2	RE3	RE4	GE1/1	GE2/1	ME0
ϕ	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
θ	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
bend	✓	✓	✓	✓	✓							✓
quality	✓	✓	✓	✓	✓							✓
time												

- ▶ NN regresses muon p_T based on 36 inputs

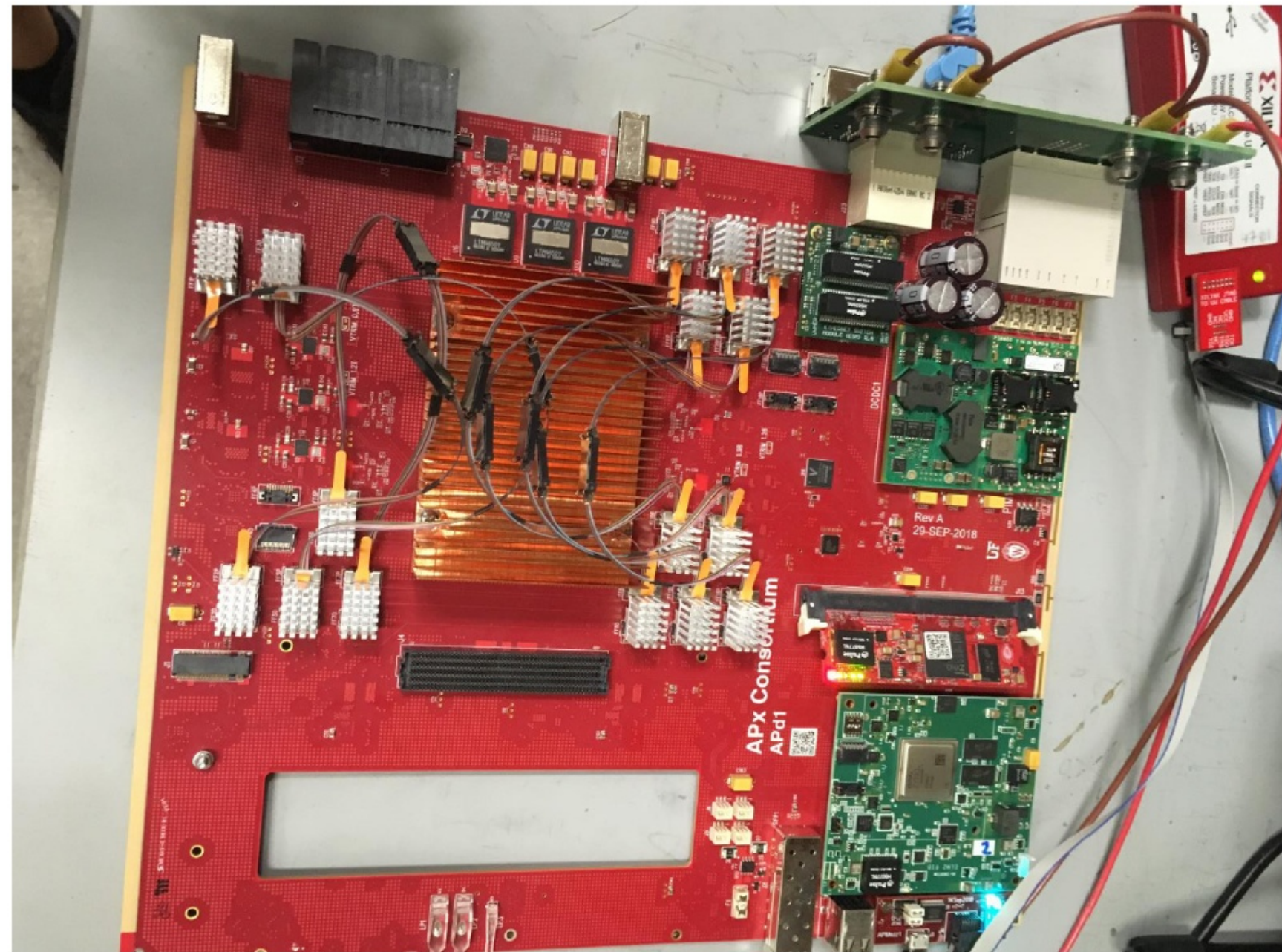


	ME1/1	ME1/2	ME2	ME3	ME4	RE1	RE2	RE3	RE4	GE1/1	GE2/1	ME0
ϕ	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
θ	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
bend	✓	✓	✓	✓	✓							✓
quality	✓	✓	✓	✓	✓							✓
time												

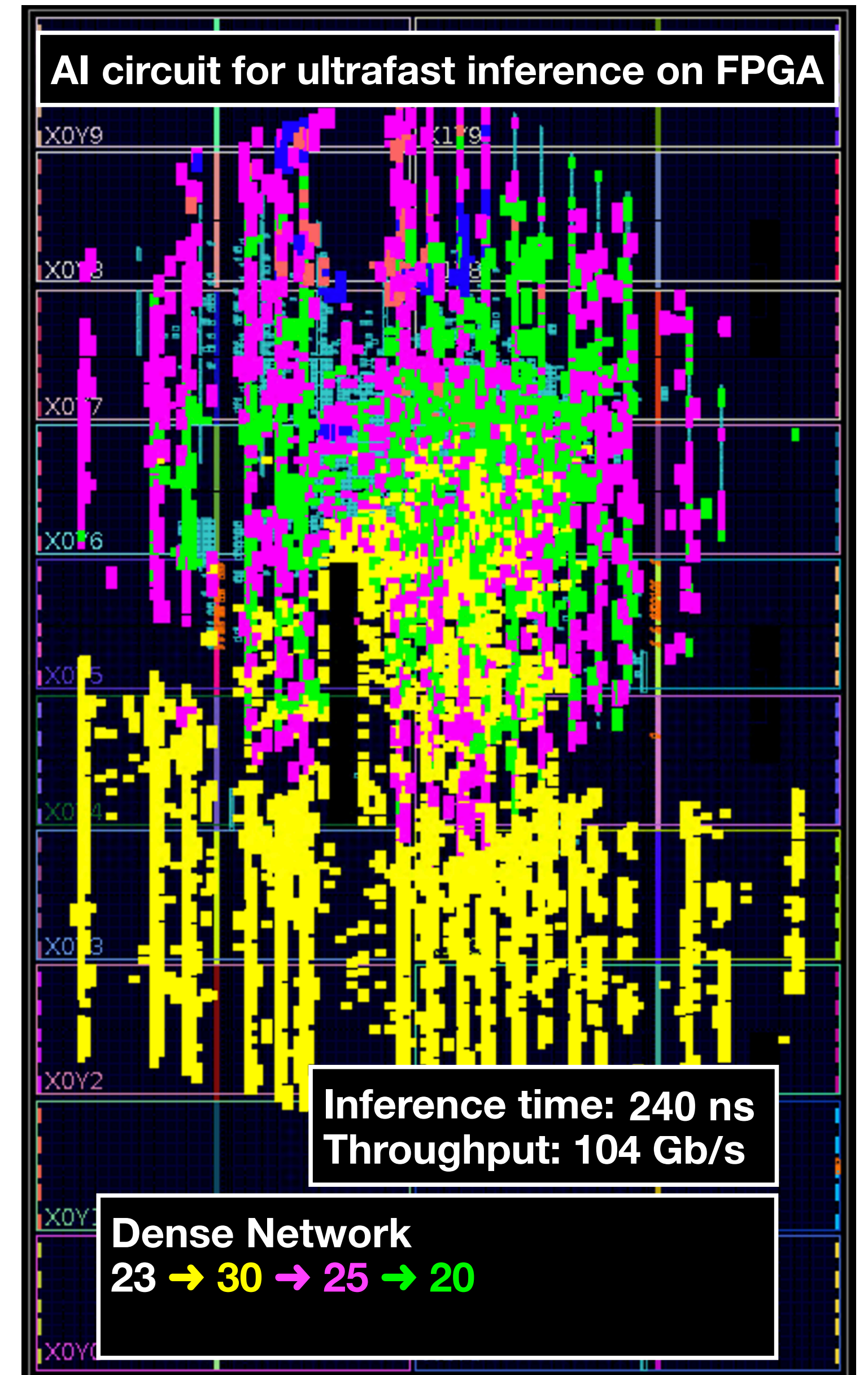


- ▶ NN regresses muon p_T based on 36 inputs
- ▶ 3× reduction in the trigger rate for NN!

Algorithm (target FPGA)	LUT	Flip-flop	Block RAM	DSP
NN + EMTF (VU9P)	28%	8%	30%	30%

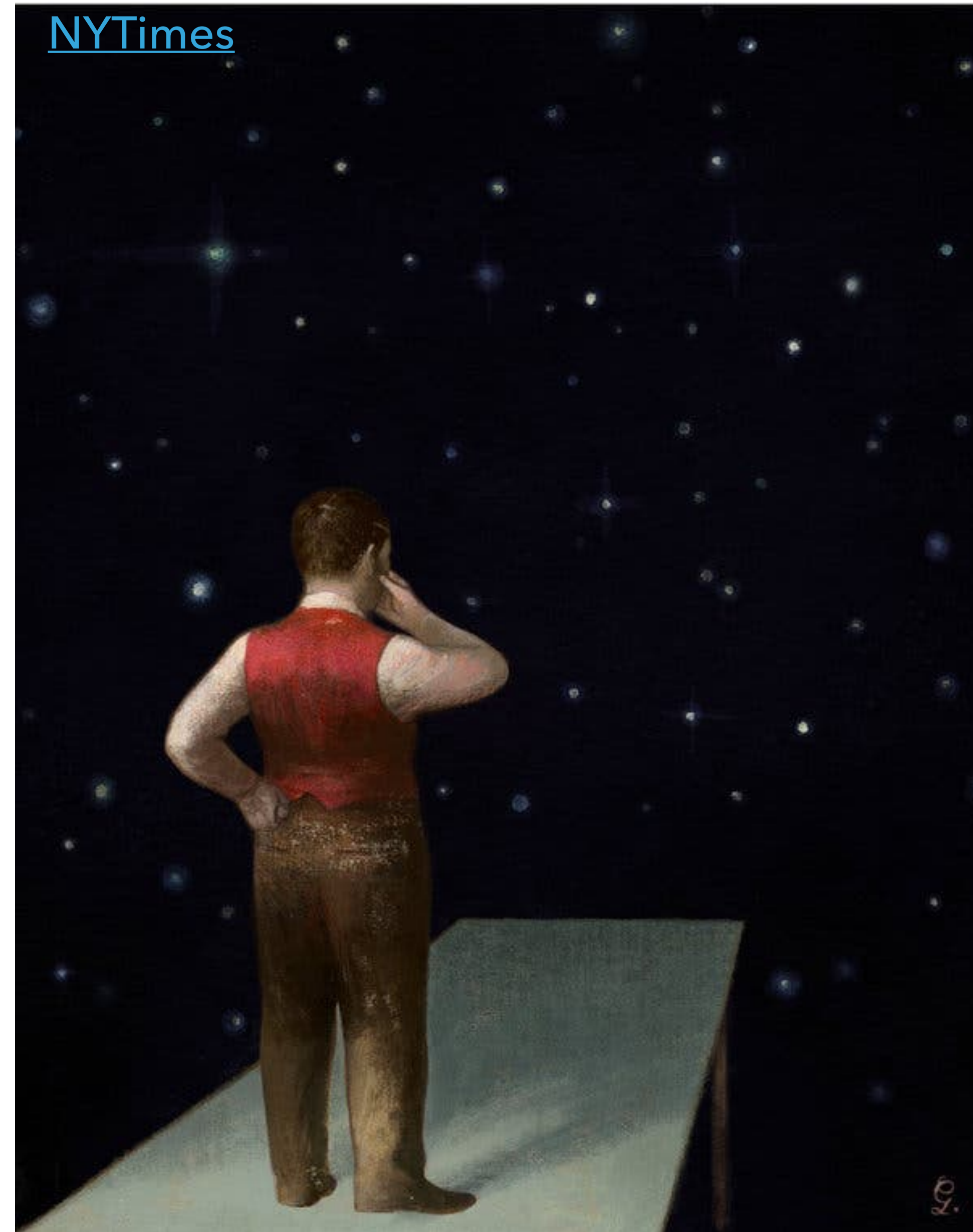


- ▶ Fits within L1 trigger latency (240 ns!) and FPGA resource requirements (less than 30%)



- ▶ Modern ML is the latest tool in the arsenal of HEP that has a wide range of applications
 - ▶ Jet tagging/regression, event reconstruction, anomaly detection, trigger, data compression, generation/simulation
- ▶ We have only scratched the surface of what is possible in the future with ML
 - ▶ Improvements in physics sensitivity, detector design, automatic calibrations, reducing time/cost of data analysis

- ▶ Modern ML is the latest tool in the arsenal of HEP that has a wide range of applications
 - ▶ Jet tagging/regression, event reconstruction, anomaly detection, trigger, data compression, generation/simulation
- ▶ We have only scratched the surface of what is possible in the future with ML
 - ▶ Improvements in physics sensitivity, detector design, automatic calibrations, reducing time/cost of data analysis
- ▶ With upcoming data at the LHC and beyond, we will explore the **edge** of the unknown in particle physics with cutting-edge ML





JAVIER DUARTE
LISHEP SESSION C
JULY 6, 2021

BACKUP