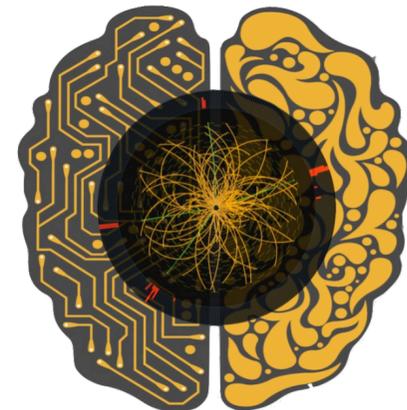


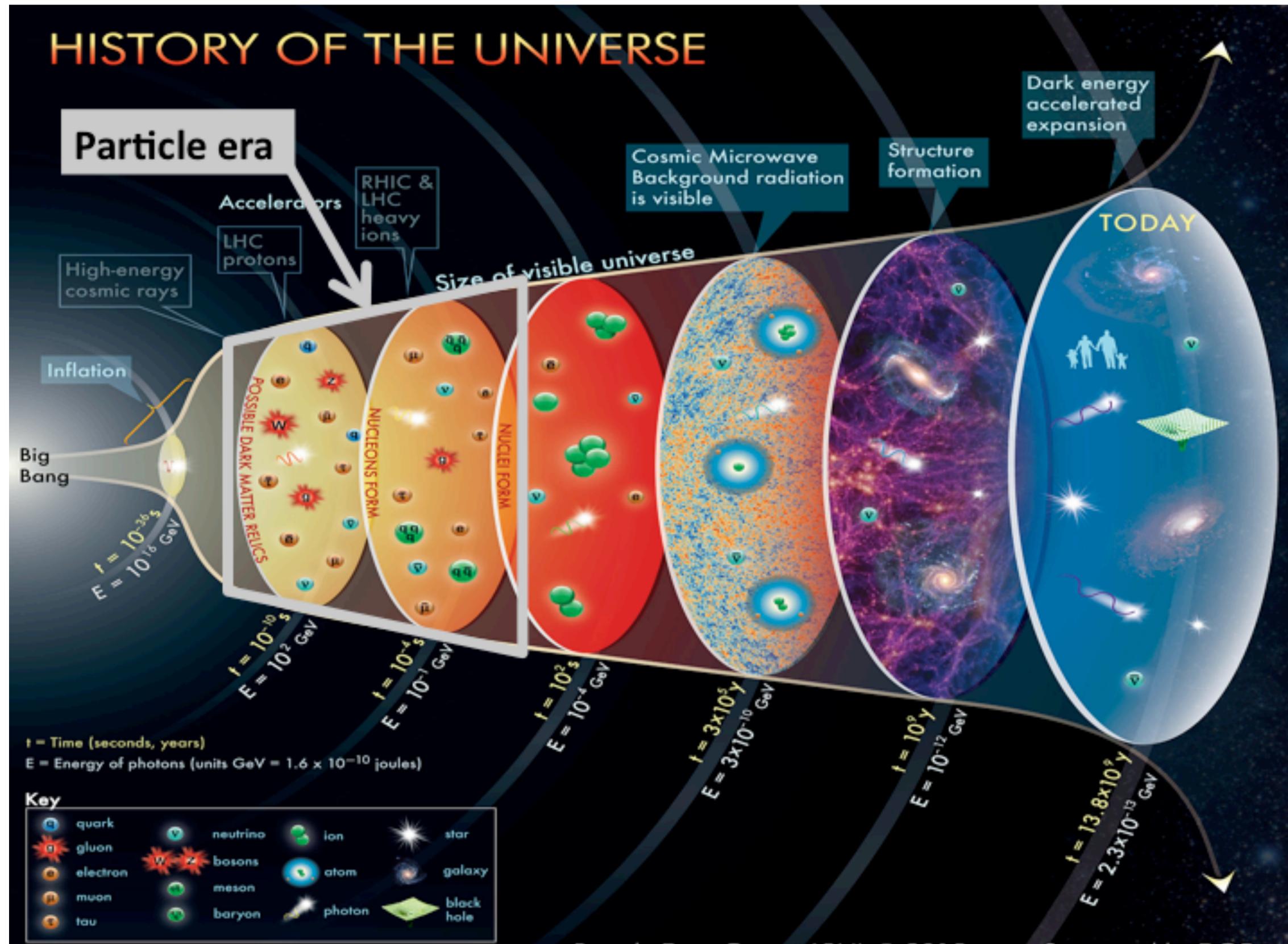
# Artificial Intelligence Accelerated Discoveries At the Large Hadron Collider

Mia Liu  
Purdue University

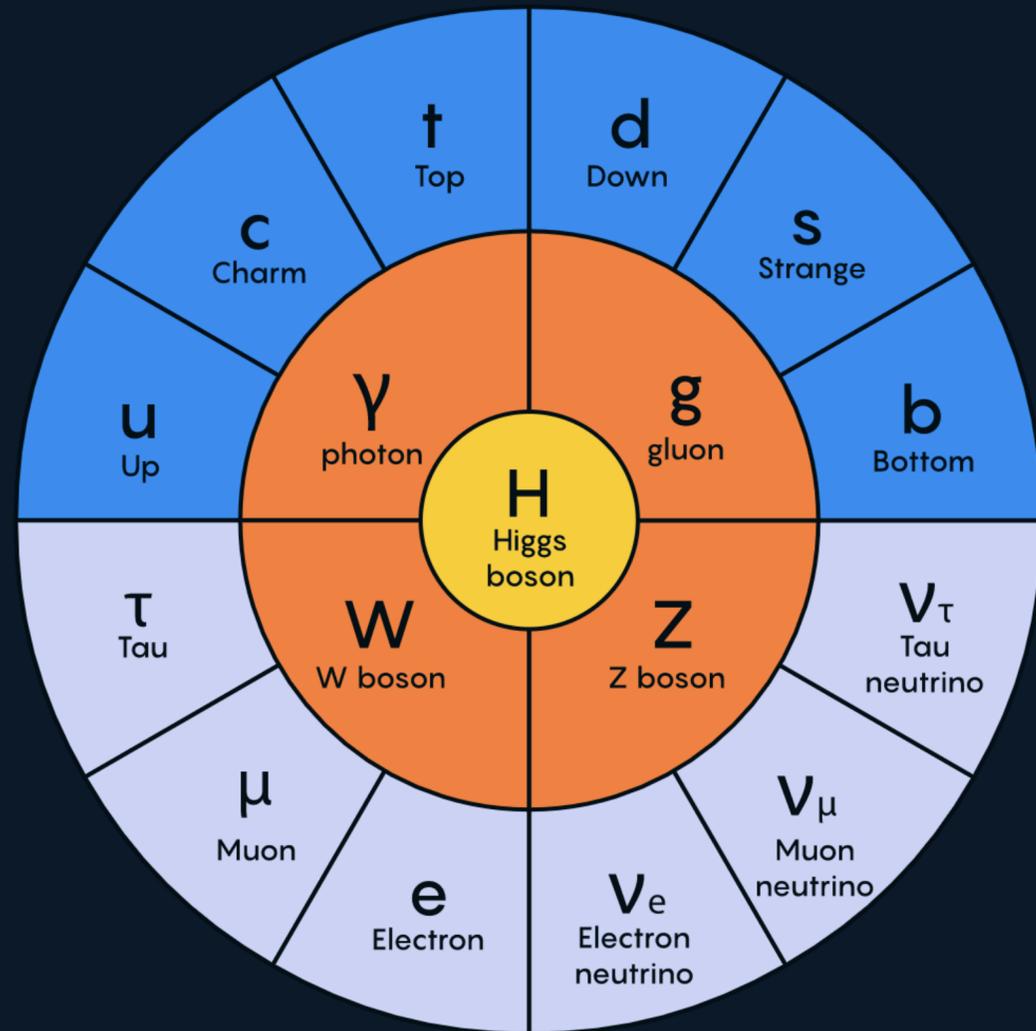
Oct 18, 2021  
A3D3 pre kick-off seminar



# Recreating the universe after the big bang<sub>2</sub>

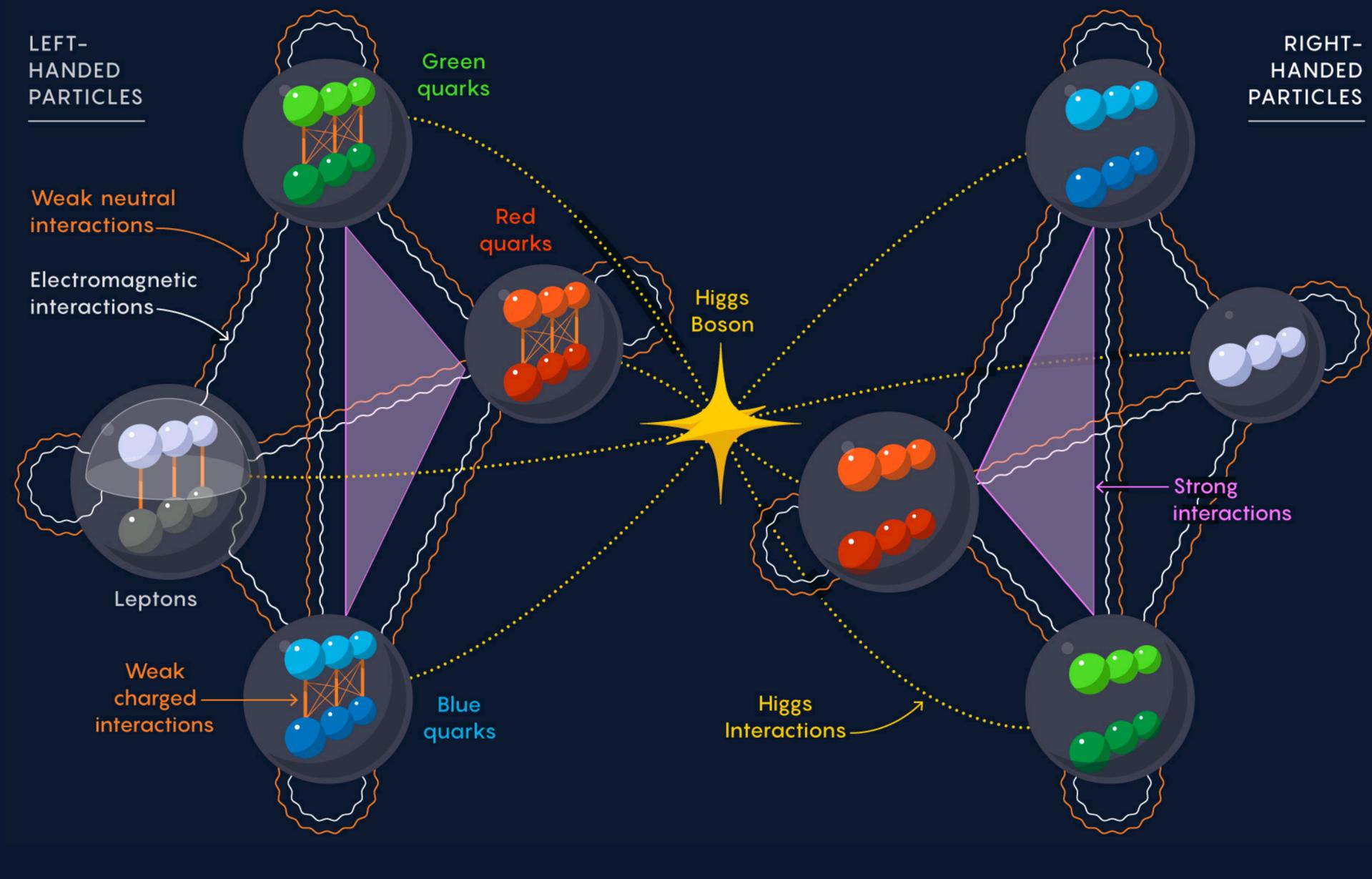


# The Standard Model of Fundamental Particles <sub>3</sub>



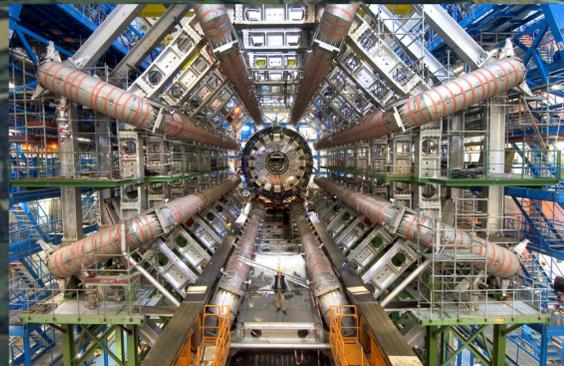
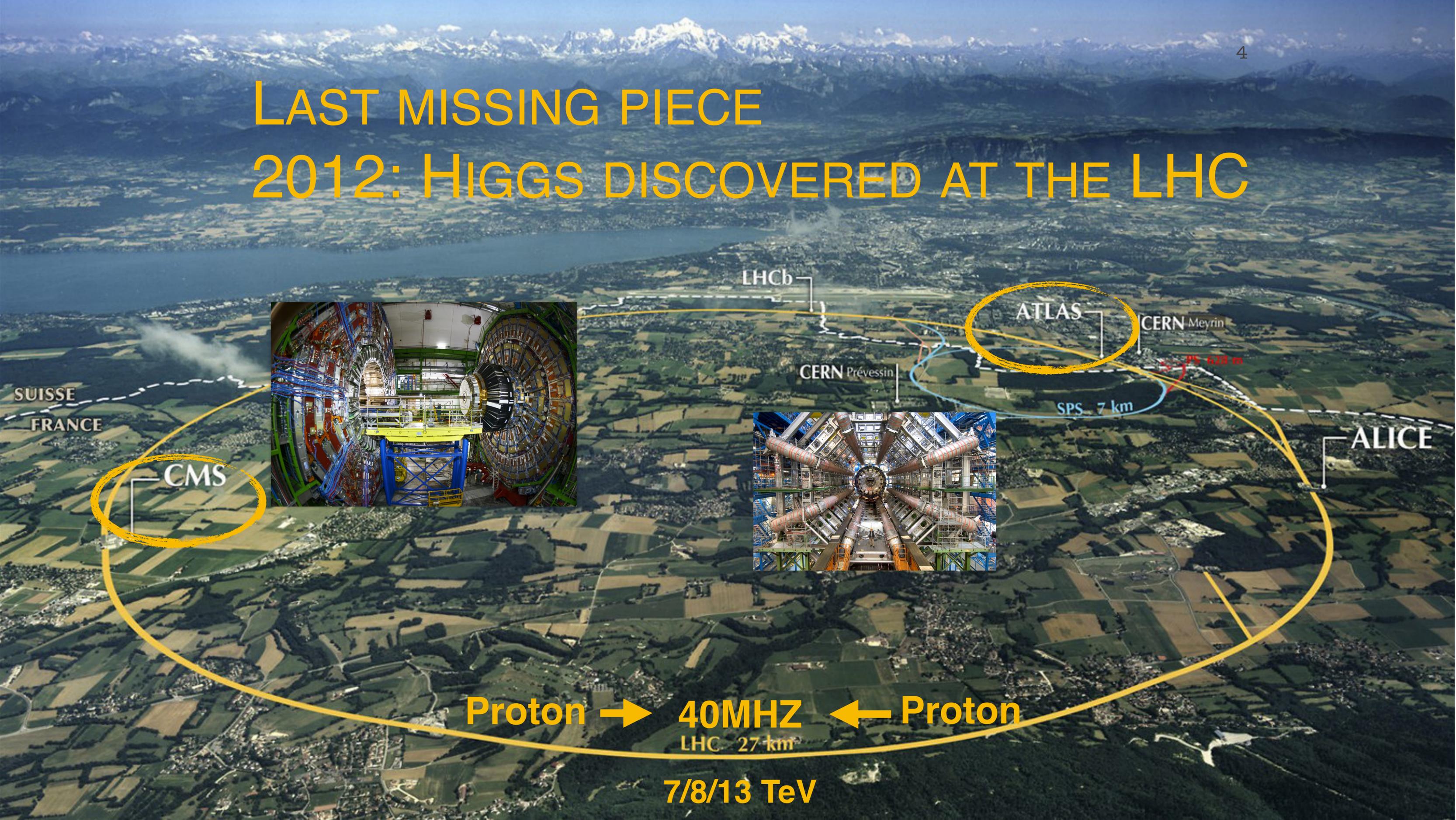
FERMIONS (MATTER)      BOSONS (FORCE CARRIERS)

● QUARKS   ● LEPTONS   ● GAUGE BOSONS   ● HIGGS BOSON



# LAST MISSING PIECE

## 2012: HIGGS DISCOVERED AT THE LHC



SUISSE  
FRANCE

CMS

LHCb

ATLAS

CERN Meyrin

CERN Prévessin

SPS 7 km

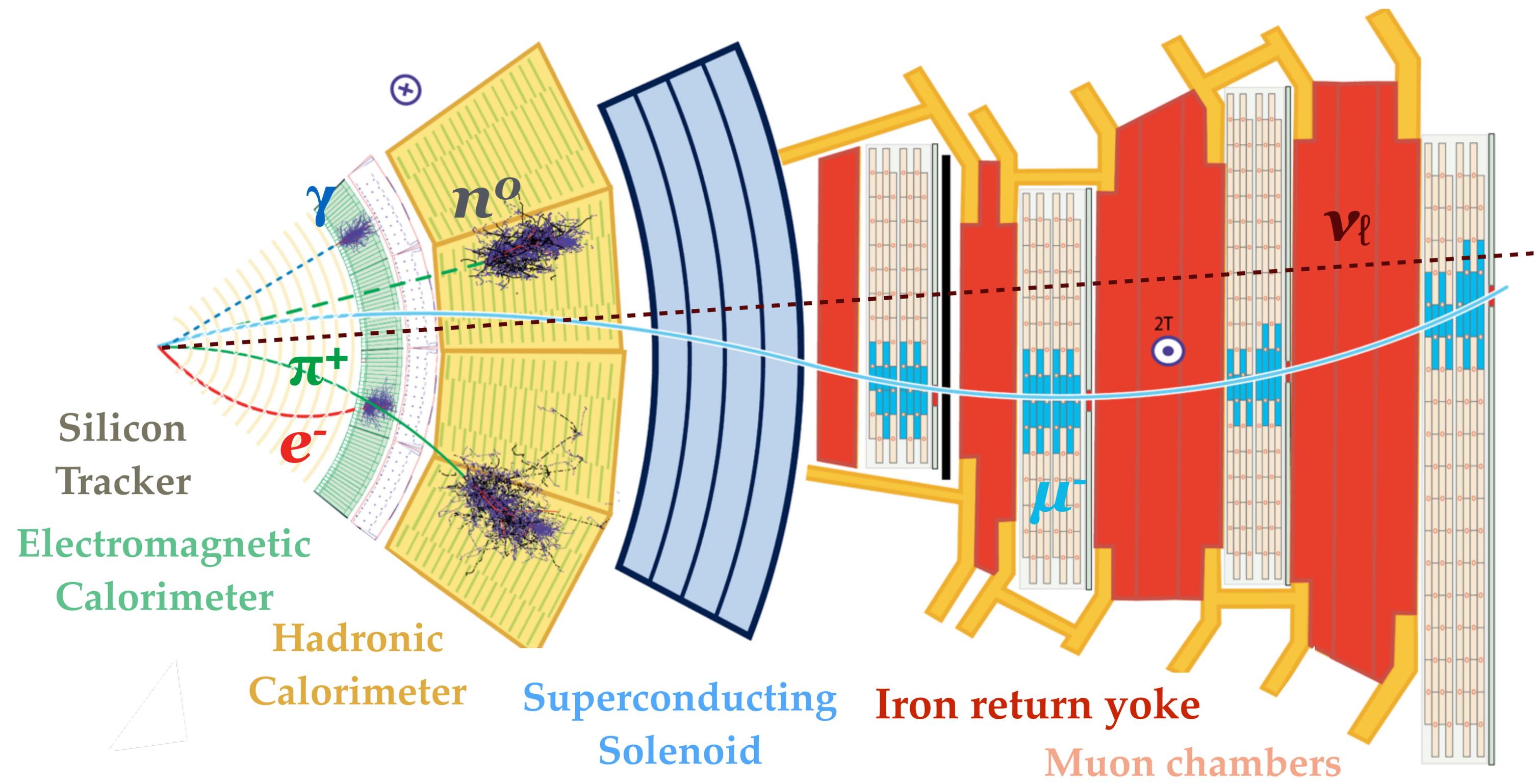
PS 6.28 km

ALICE

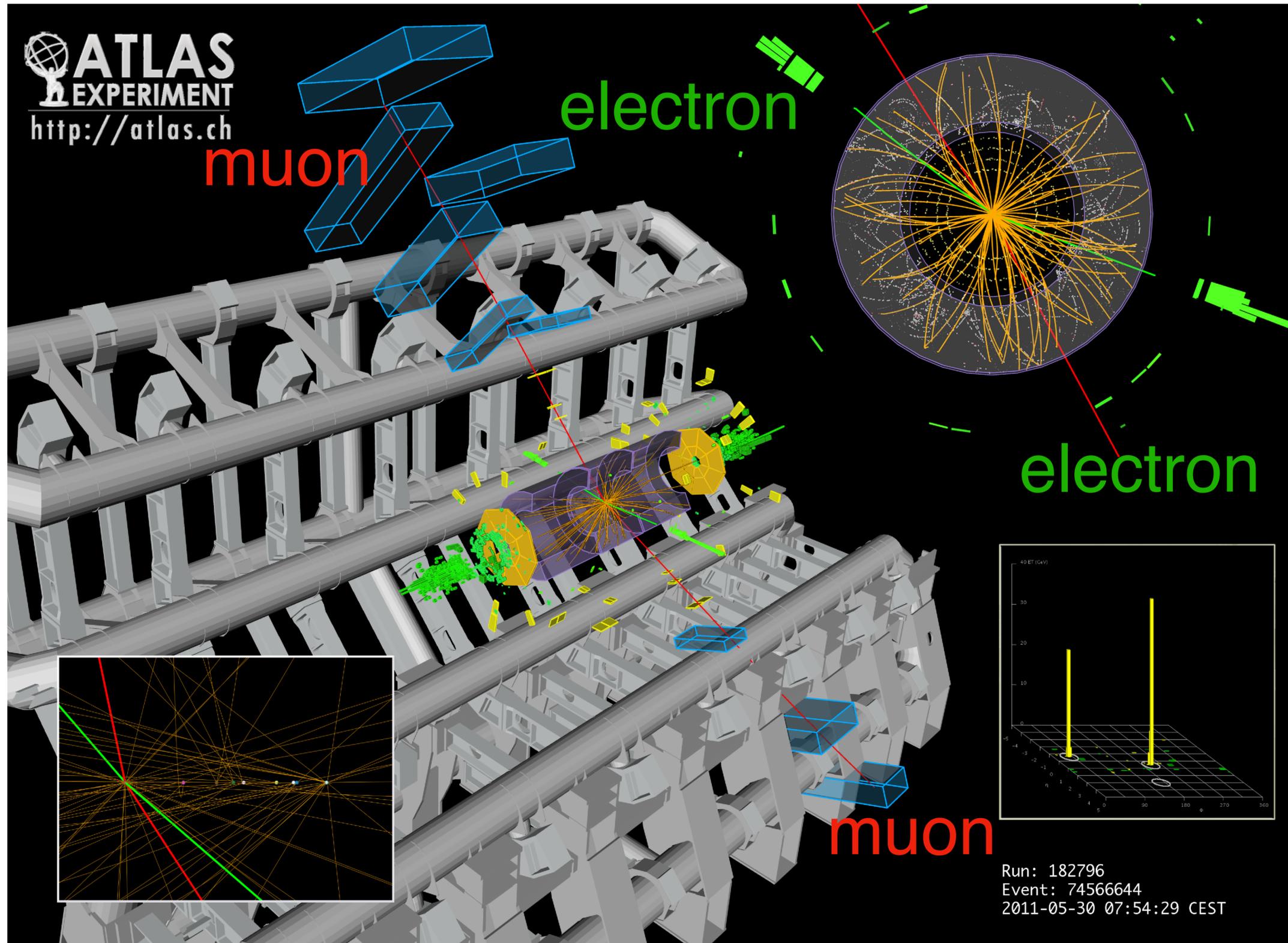
Proton → 40MHZ ← Proton  
LHC 27 km

7/8/13 TeV

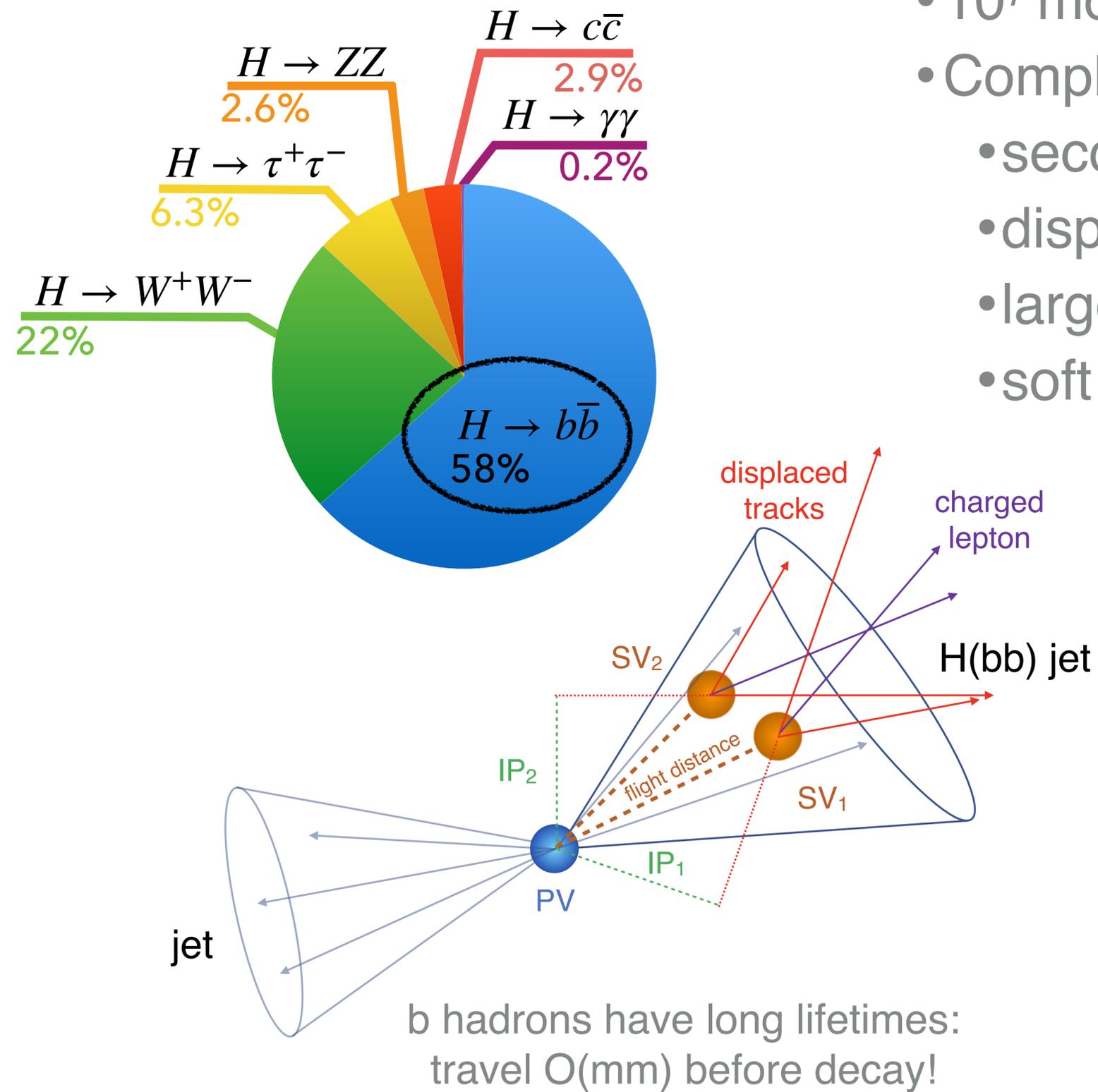
# Particle Detection in CMS



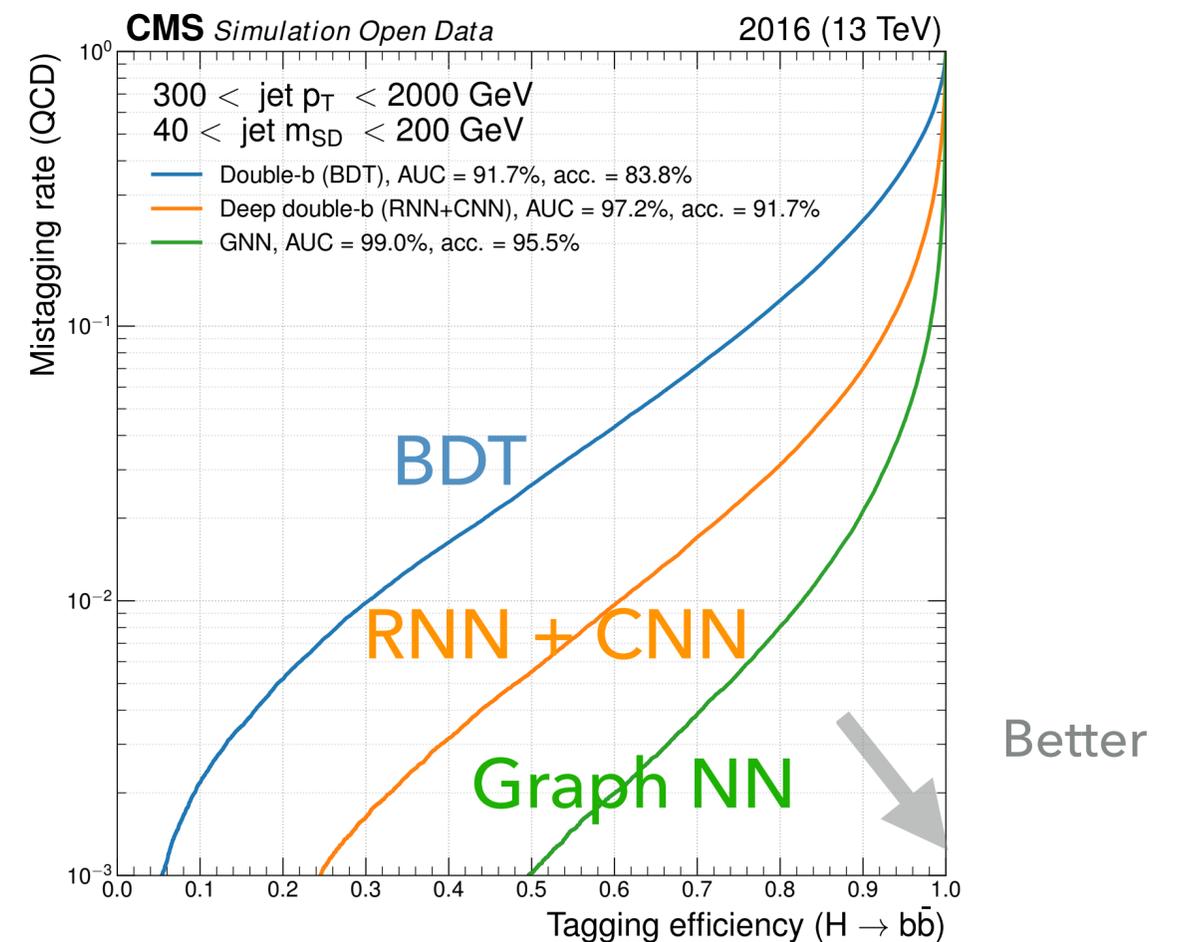
# Discovering the Higgs boson ( $H \rightarrow ZZ \rightarrow ee\mu\mu$ )



# Tagging boosted Higgs with Machine Learning 7



- $10^7$  more frequent background
- Complex signature:
  - secondary vertices
  - displaced tracks
  - large impact parameters
  - soft leptons



# Puzzles

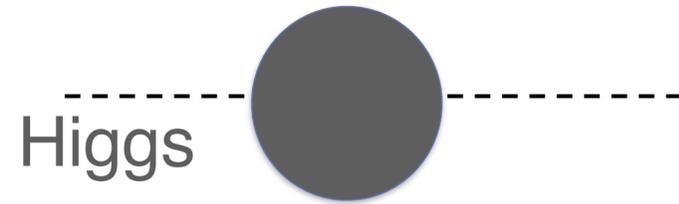
Fine tuning  
Dynamical origin  
...

Experimental:  
Dark matter/dark energy  
Not in SM

Neutrinos

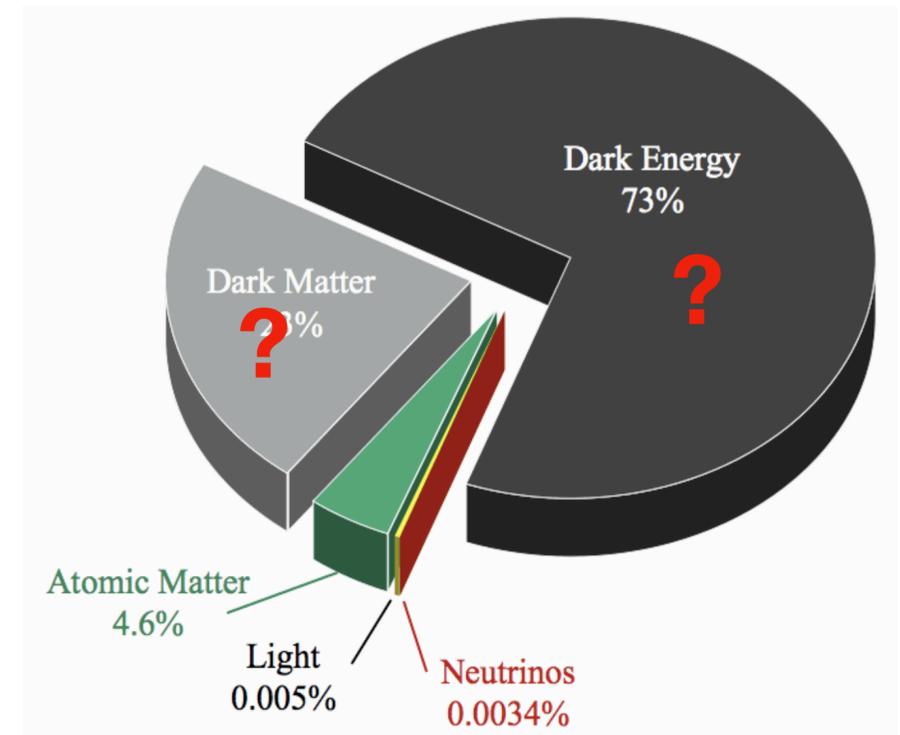
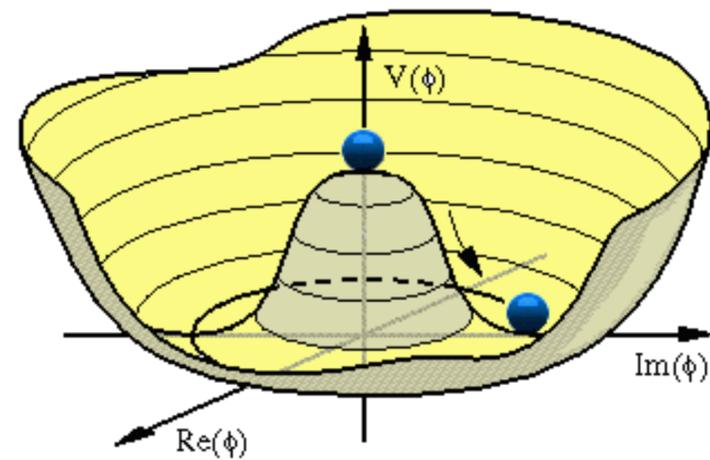
Anomalies: Muon g-2, LHCb  
lepton flavour universality

SM particles

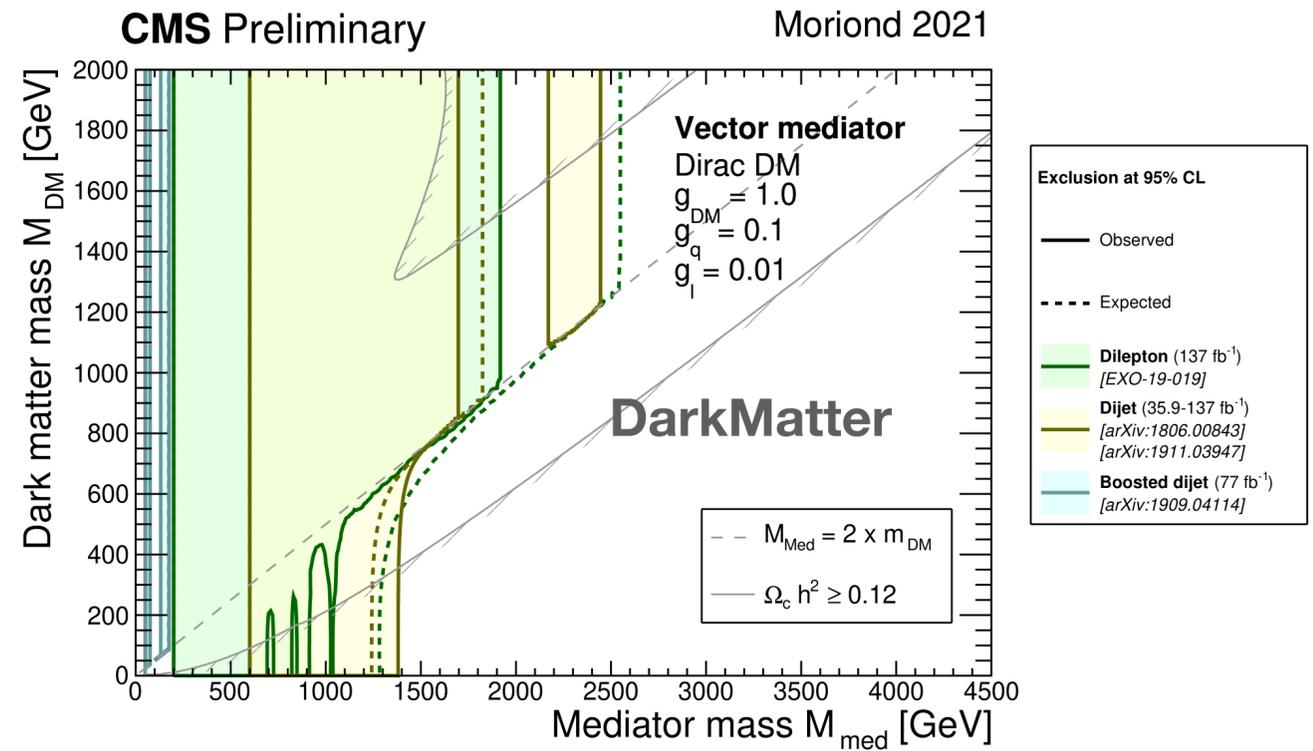
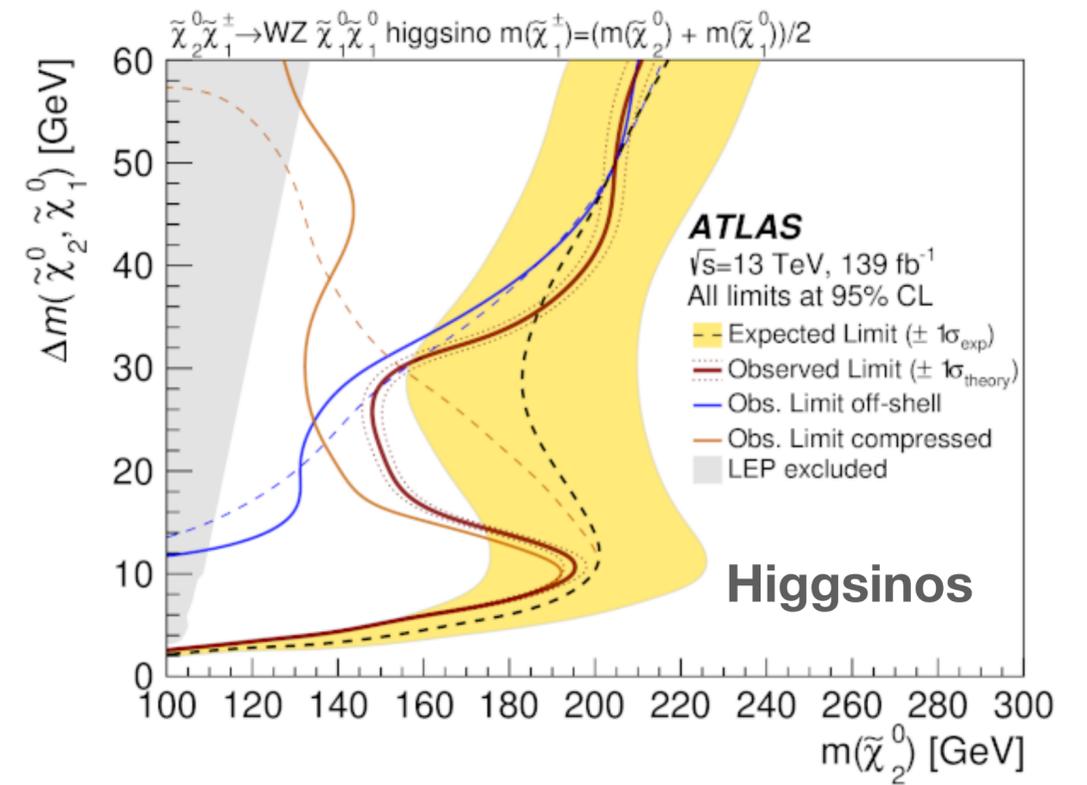
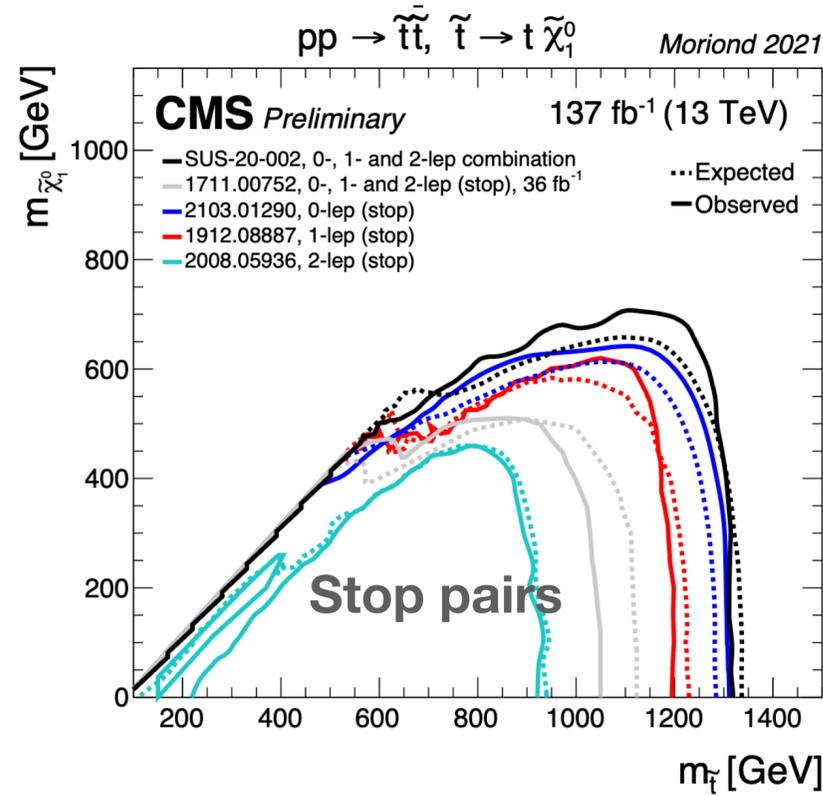
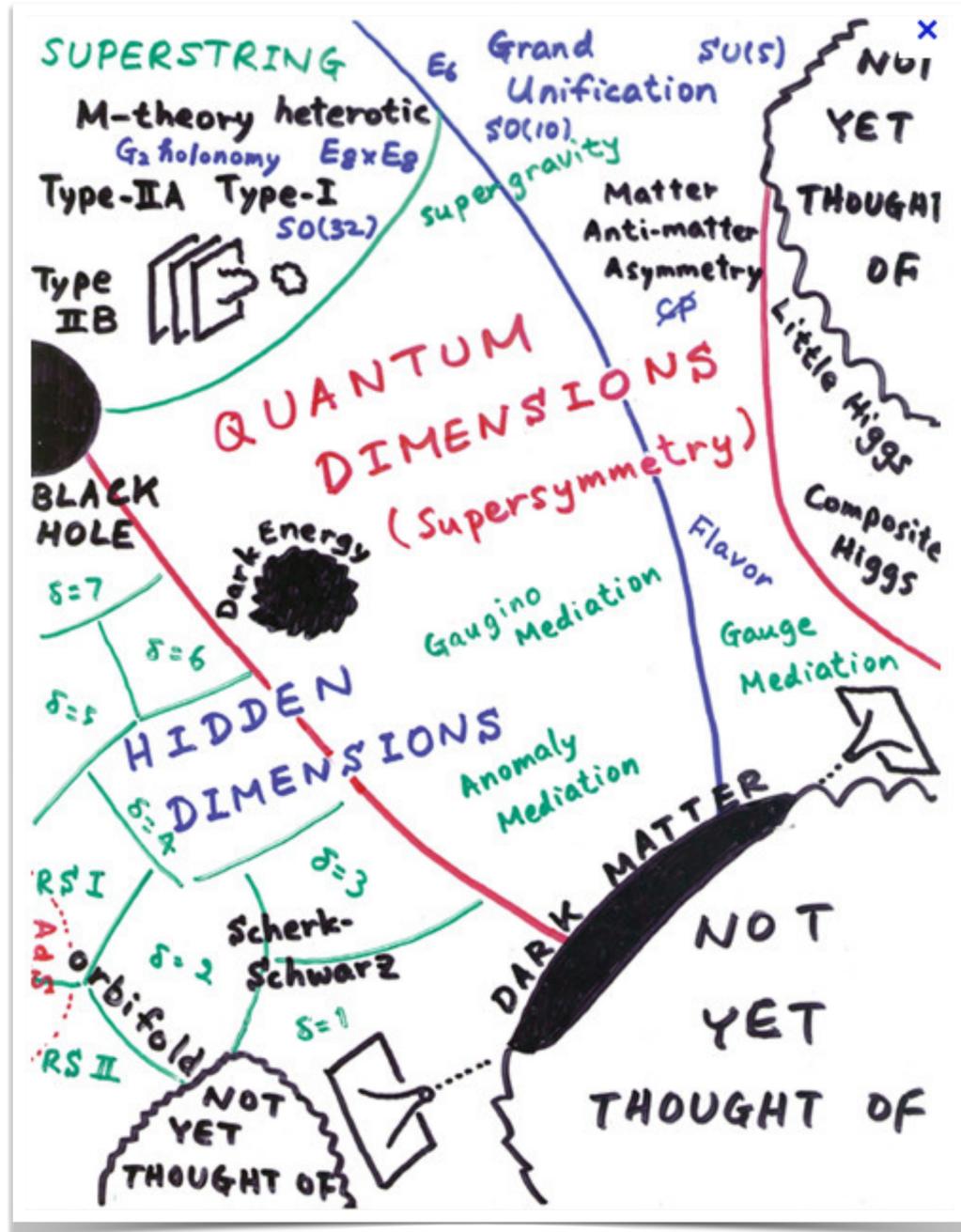


$$m_{obs}^2 = m_{bare}^2 - \frac{\lambda_f^2}{8\pi^2} \Lambda^2$$

125 GeV       $\epsilon * 10^{19}$  GeV?       $10^{19}$  GeV

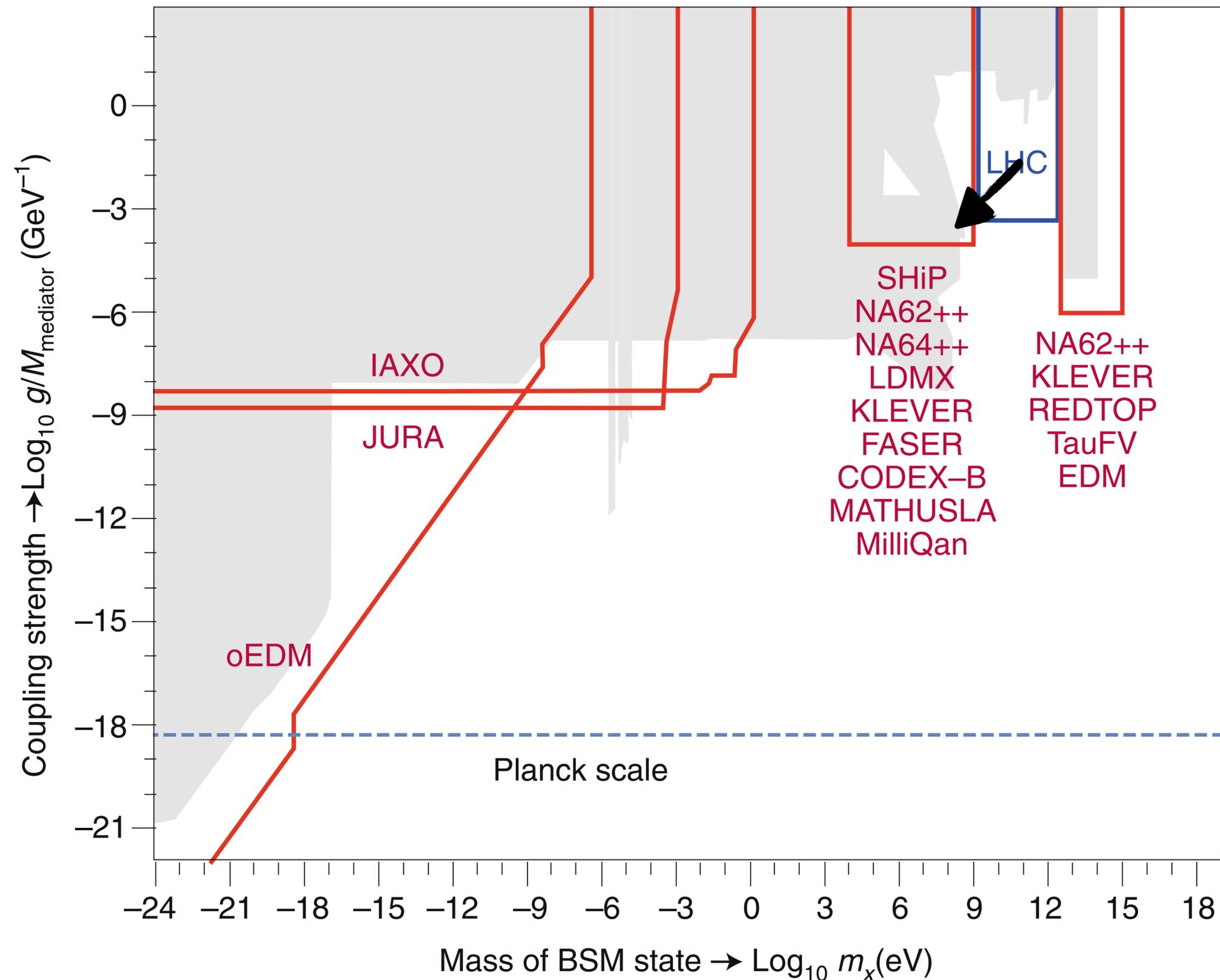


# Extensive searches performed at the LHC



# Improving search space at the LHC

10



**Low mass, weakly coupled BSM particles:**

**Unconventional signatures:**

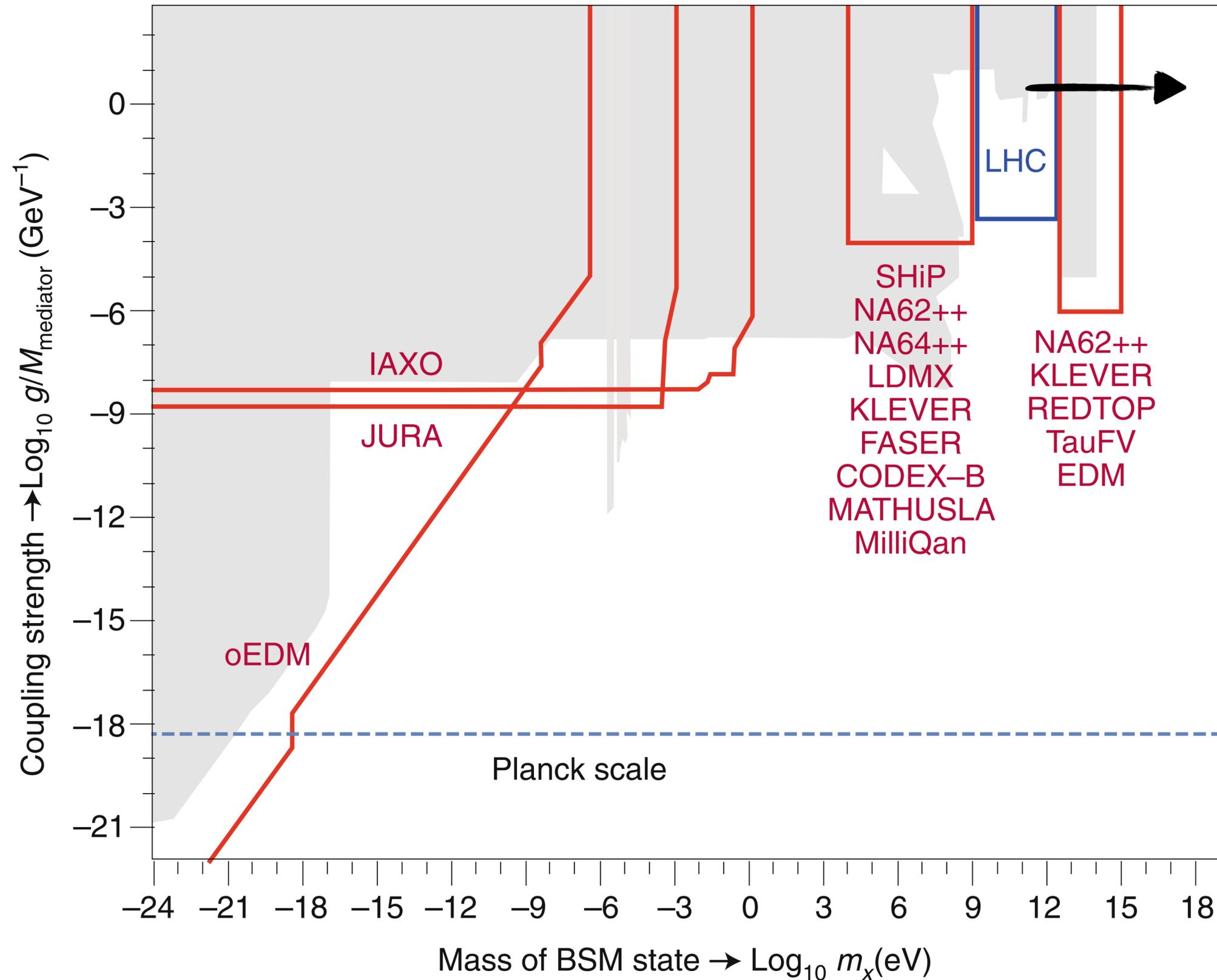
Long lived particles that decay outside beam pipe

**Many are limited by online reconstruction/filtering (trigger):**

Currently relying on additional object present in the event

# Improving search space at the LHC

11



**Improve high energy probe:**

**high pT object detection.**  
e.g. boosted Higgs/W/Z/top  
etc.

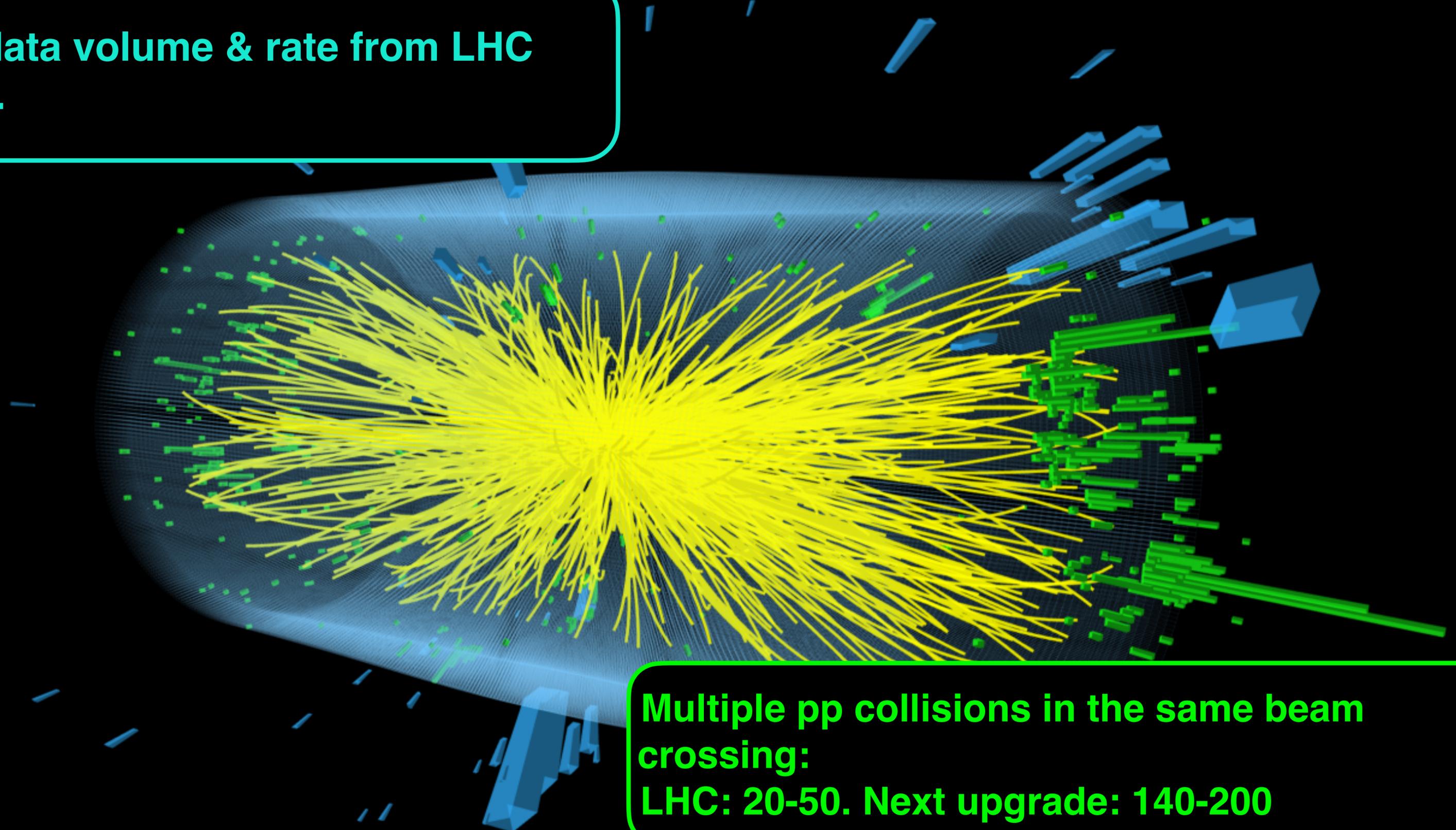
**Loop contributions to ultra-rare processes.**e.g lepton flavor violation

- **soft signatures that are limited by online reconstruction/filtering (trigger) as well**

# The Fast and Furious

12

**Extreme data volume & rate from LHC collisions.**



**Multiple pp collisions in the same beam crossing:  
LHC: 20-50. Next upgrade: 140-200**

# Capturing raindrops

13

Now: The LHC



LHC Run-4



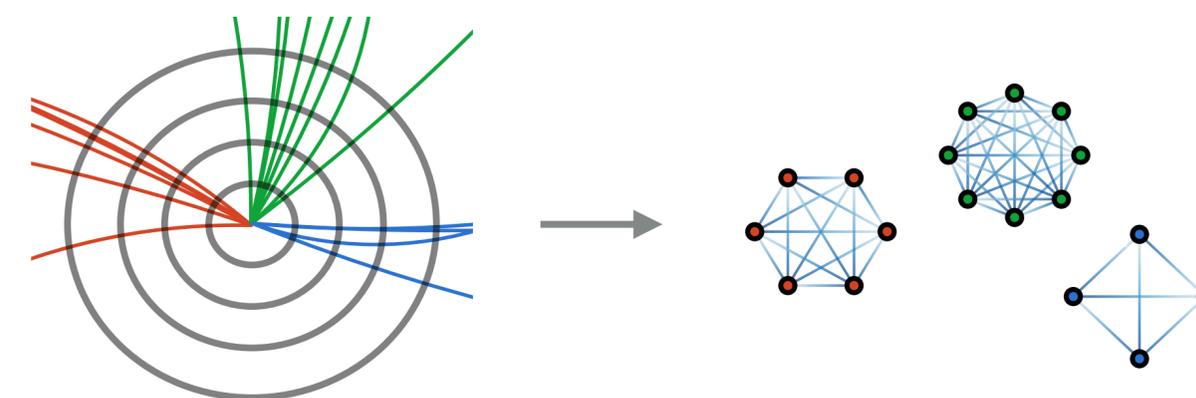
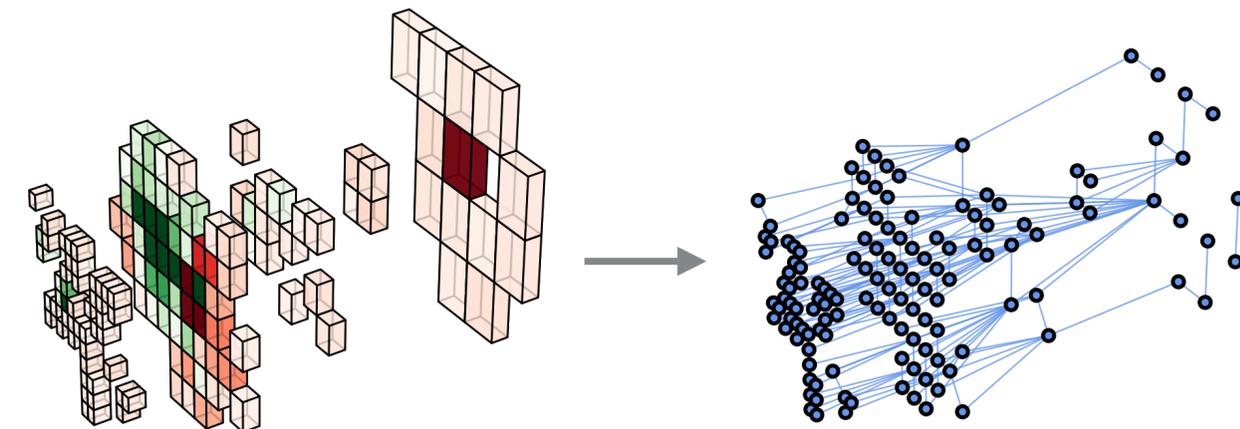
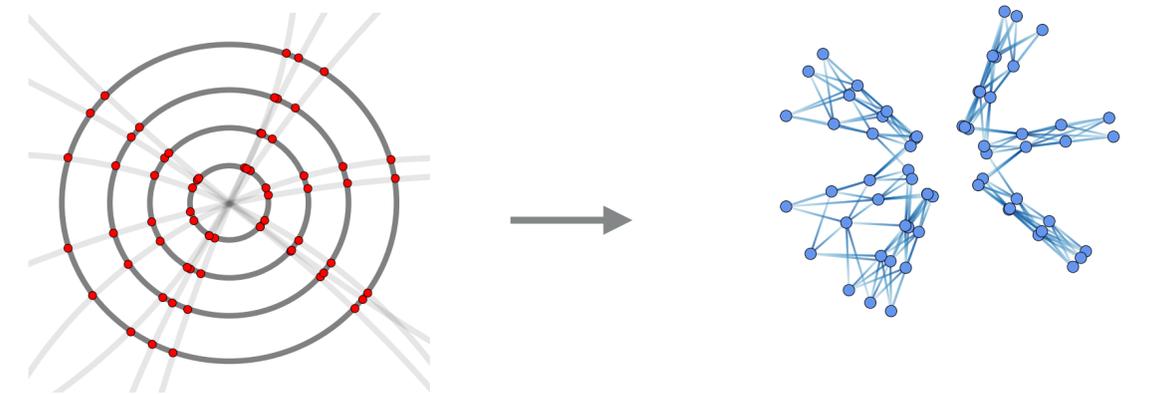
Faster and fine granularity detectors to collect higher pileup collisions  
ML to maintain/improve detection of (un)conventional signatures

## Represent data as nodes and edges

appropriate representation of particle physics data:  
irregular, structural, relational

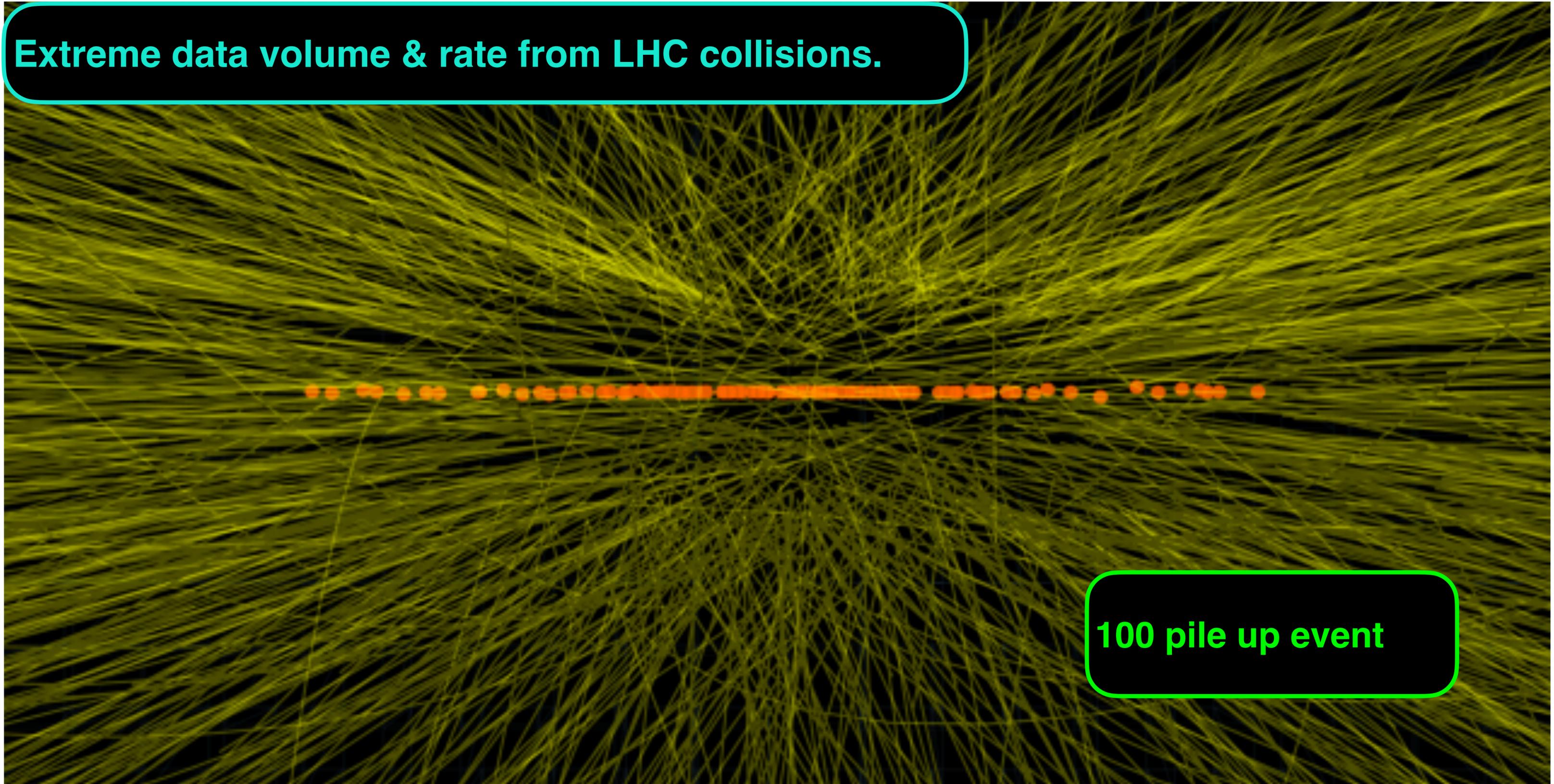
## GNN developments driven by social media, how to adapt/adopt to domain knowledge

- node classification: Pileup mitigation
- graph classification:  $\tau \rightarrow 3\mu$
- edge classification: charged particle tracking
- identifying subgraphs/sets: particle flow



# The Fast and Furious

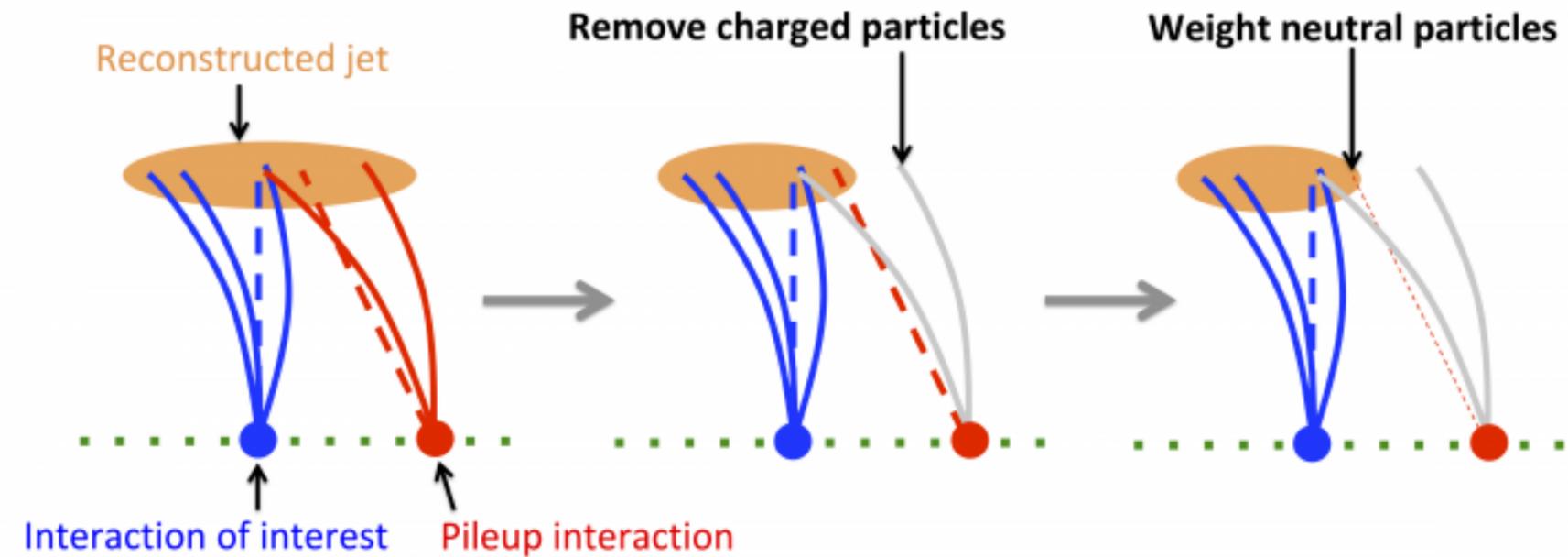
**Extreme data volume & rate from LHC collisions.**



**100 pile up event**

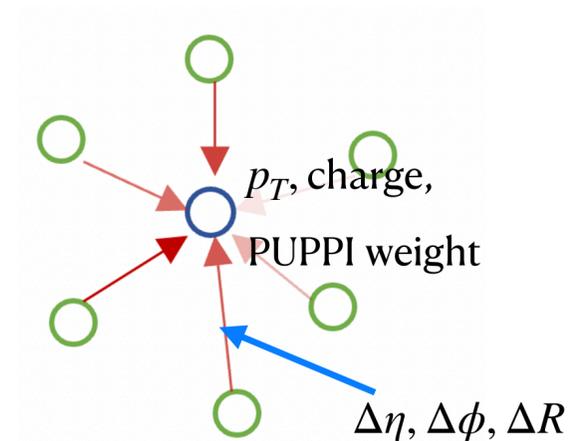
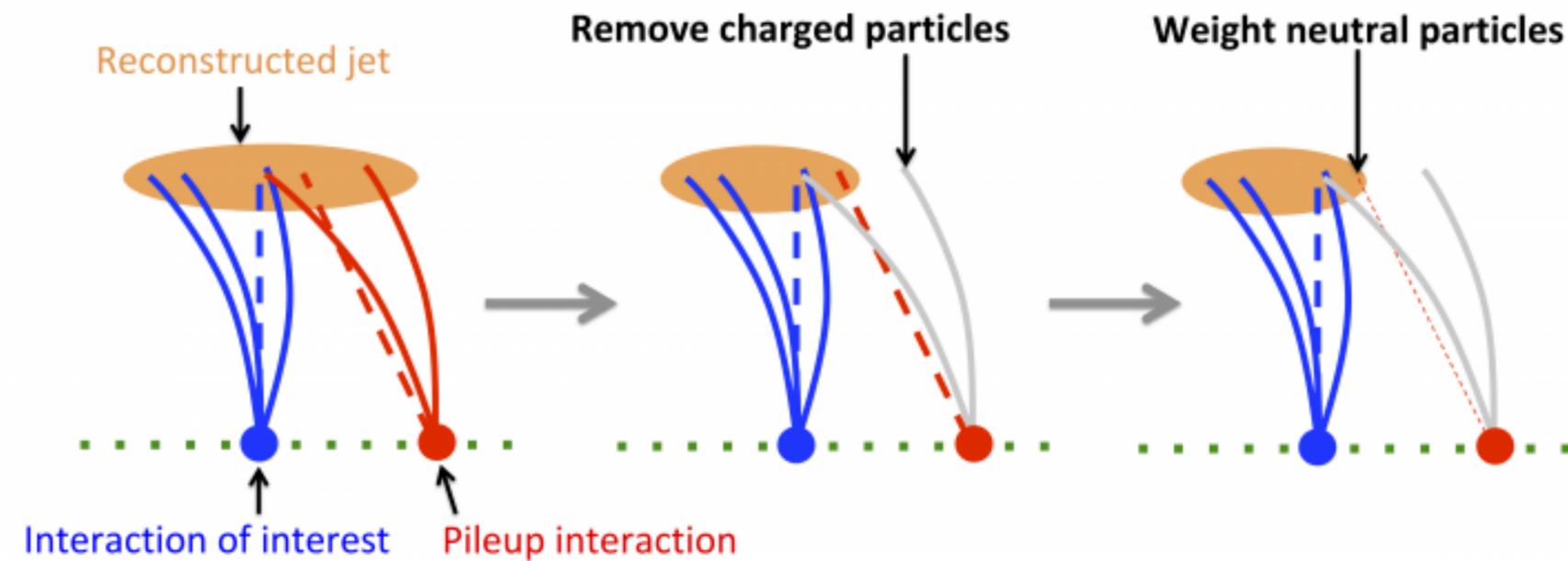
Essential to the LHC physics program.

- “Noise” needs to be subtracted from interaction of interest

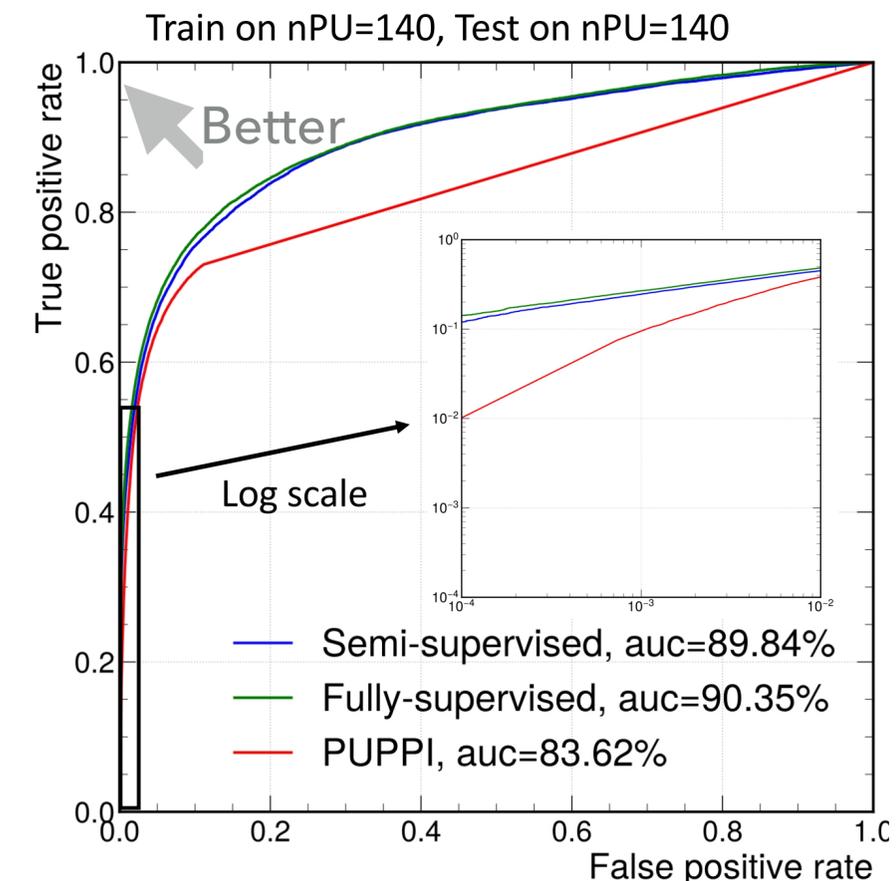


Trained on labeled charged particles and applied to neutral ones  $\rightarrow$  can learn from data.

- Outperforms expert feature method (PUPPI), comparable to fully supervised method.
- Presented at BOOST 2021, Papers coming soon.
- Next: Apply to CMS simulation & data. Neutral particle vertex association for the forward region. Physics with PU interactions?

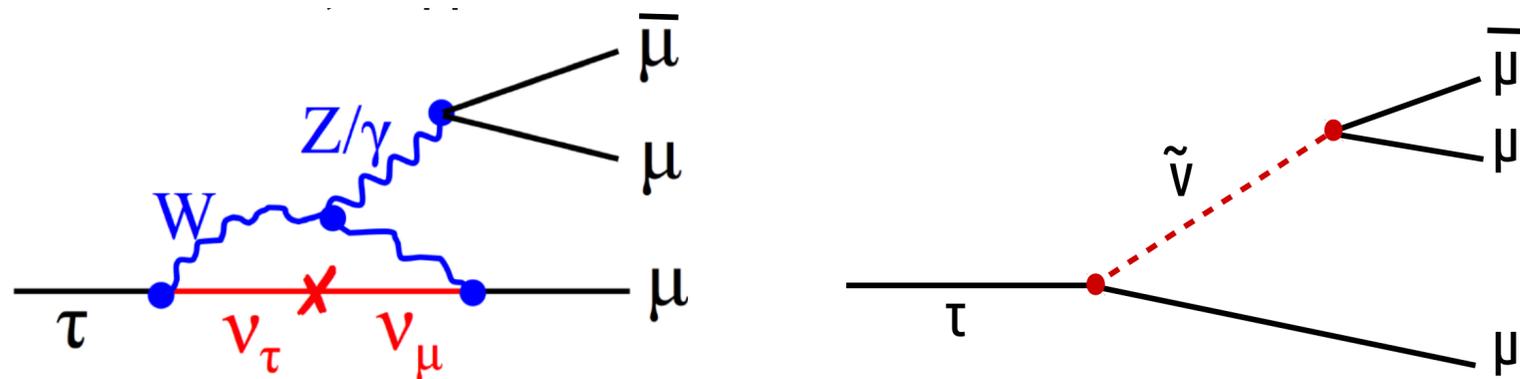


Neighboring structure

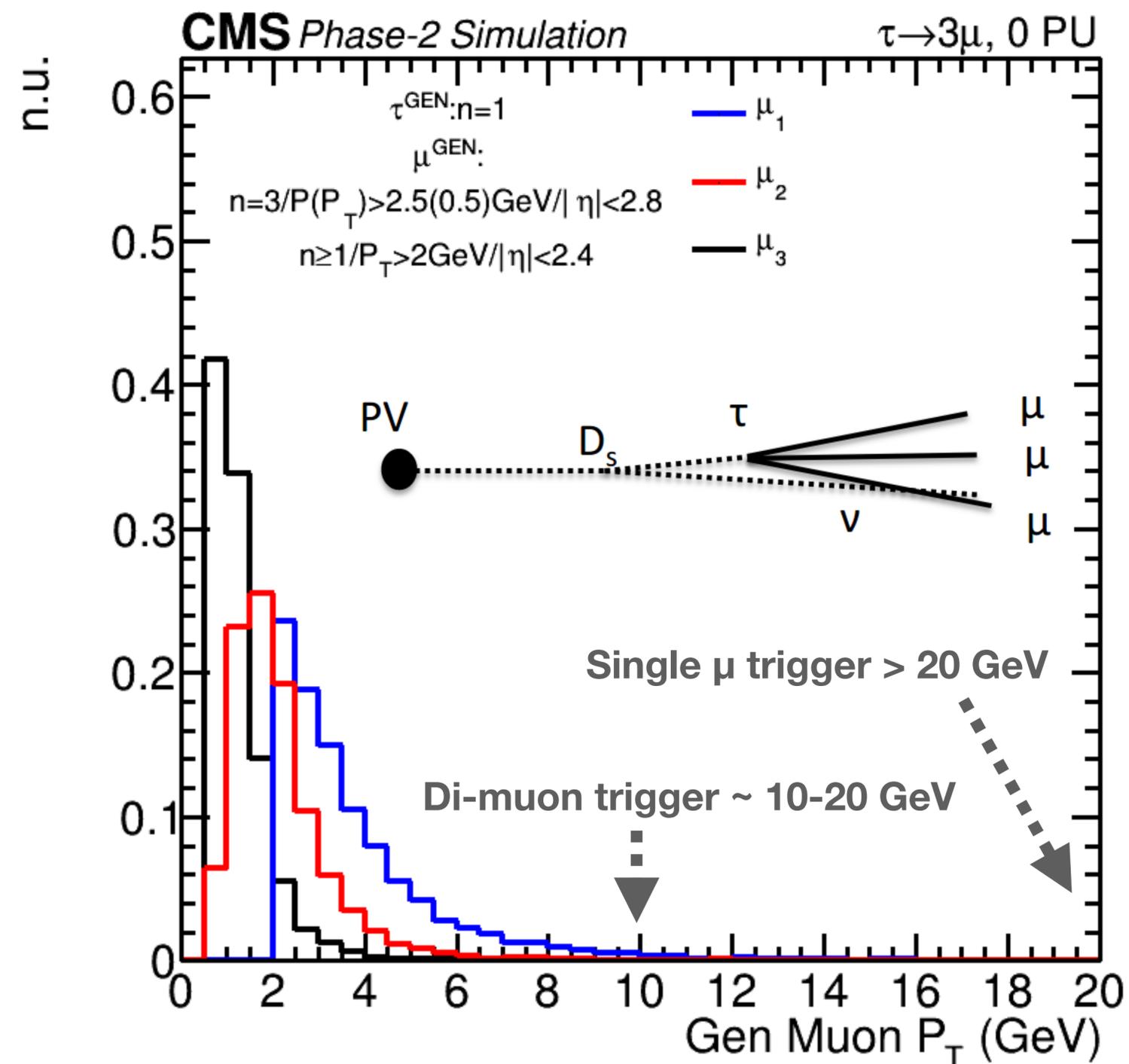


# Graph Classification $\tau \rightarrow 3\mu$

18



- **Detecting lepton flavor violating  $\tau \rightarrow 3\mu$  decays**
  - Very (super ultra) rare in the SM, Neutrino oscillations: BR  $\sim O(10^{-14})$
  - BSM e.g. R-parity violating SUSY enhanced:  $O(10^{-8})$
- **Collimated, low  $p_T$  and very forward muons**
- Next upgrade uses tracking information, **only 25% efficient** (already 10 times better than current method)



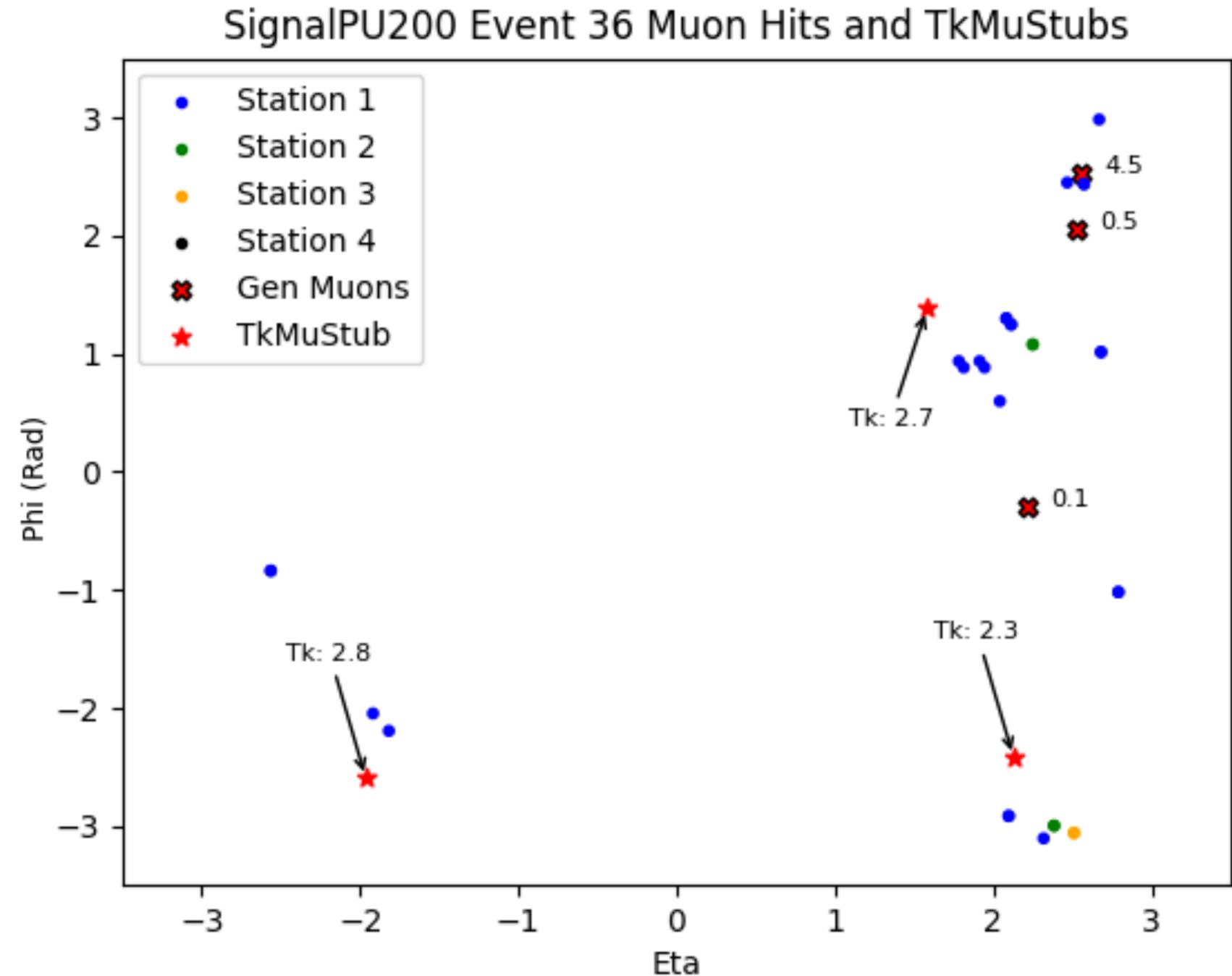
**Graph Inputs:** Muons hits (Coordinates and bending angle), represented as nodes in graph.

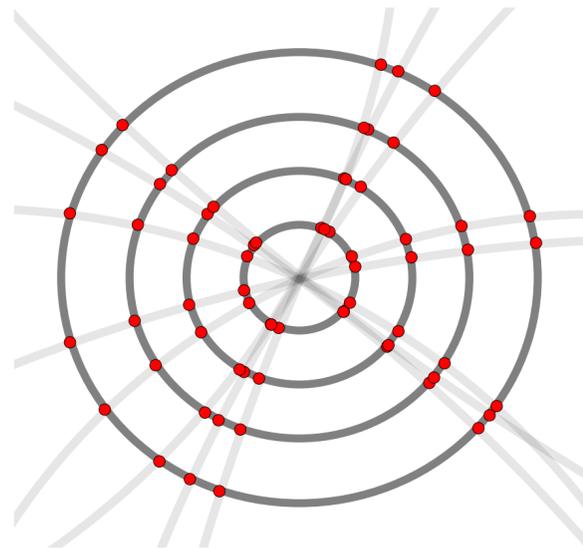
**Attention based GNNs** for information aggregation between local nodes and global node

**Training setting:**  $\tau \rightarrow 3\mu$  signals mixed with PU backgrounds vs pure PU sample

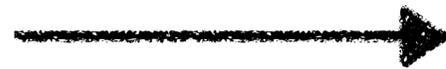
**> 90% efficient within trigger bandwidth.**

**End-to-end approach.** Model interpretations, Model adaption (to other signals & other experiments), Data augmentation. Regression of  $\tau/\mu$  kinematics

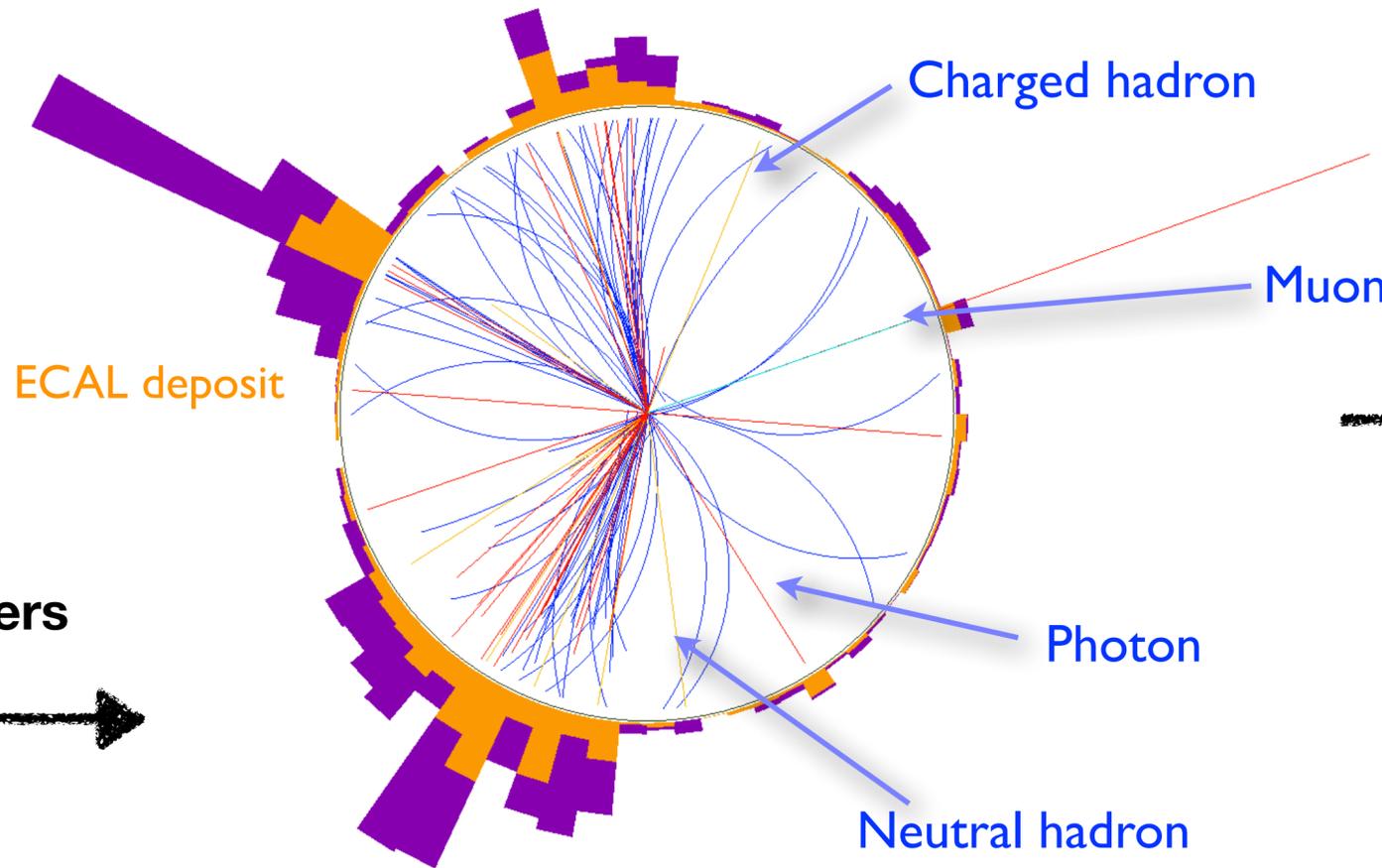




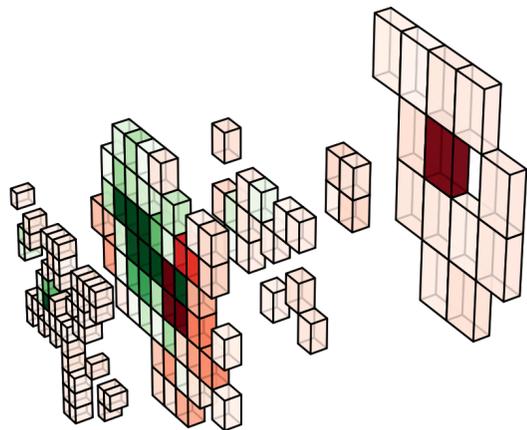
**Charged particle tracks**



HCAL deposit



**Connecting the dots**



**Energy clusters**



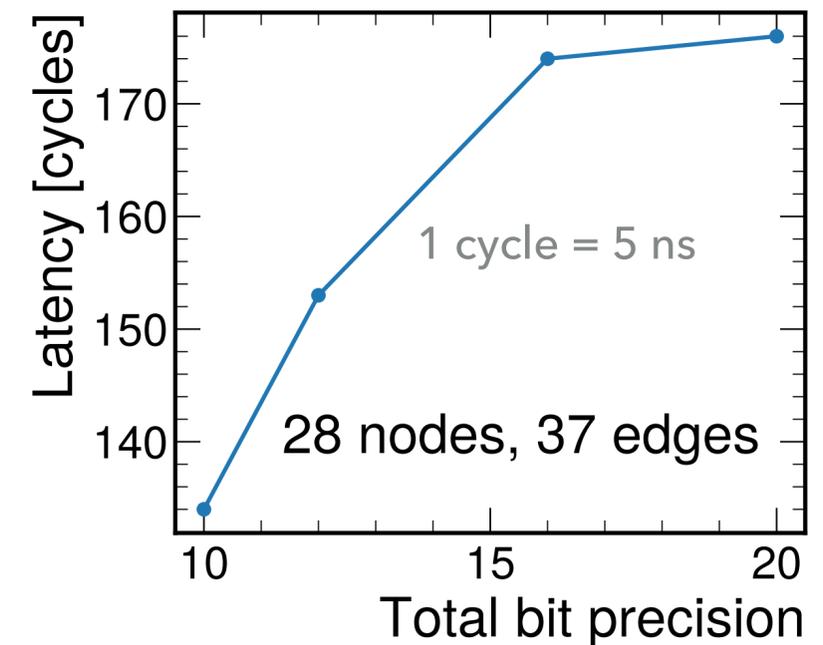
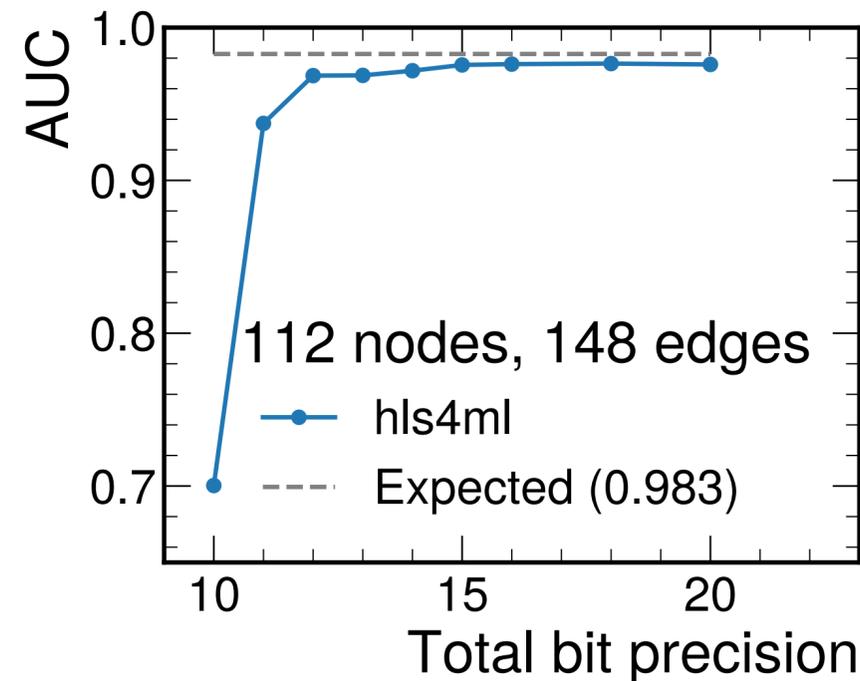
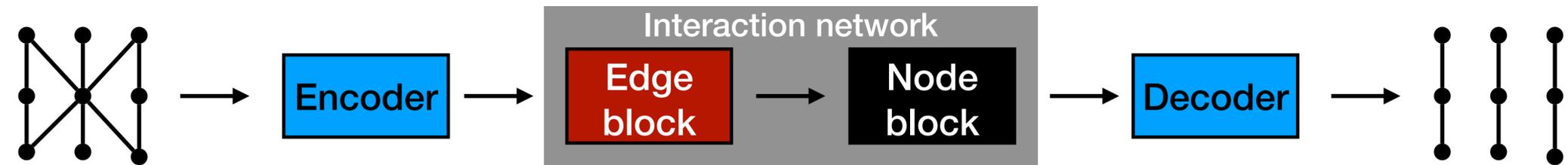
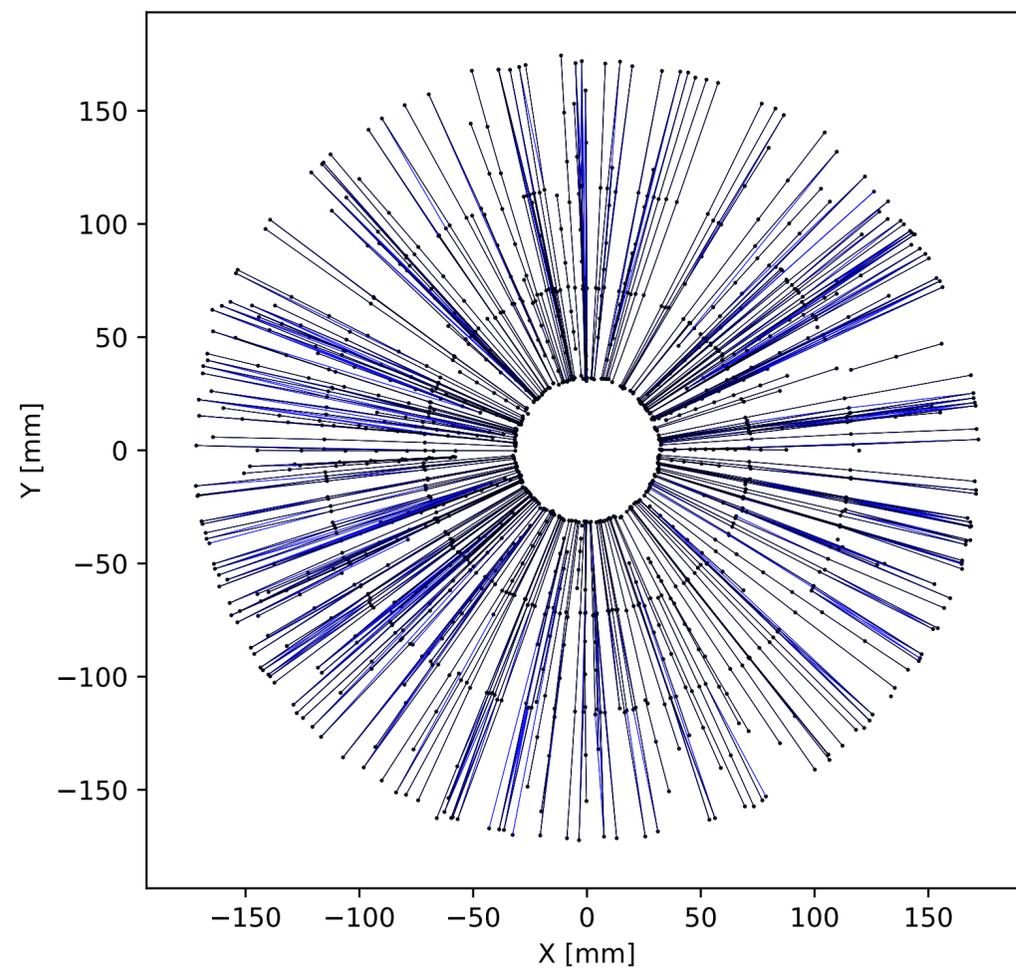
**Find the clusters**

**Particle Flow Reconstruction**

**Analysis Reconstruction**

- Higgs?
- Top Quark Pair?
- W boson?
- Z boson?
- Multiple boson?
- ⋮
- ⋮
- ⋮
- ⋮
- New Physics?

# Speed up tracking with GNN

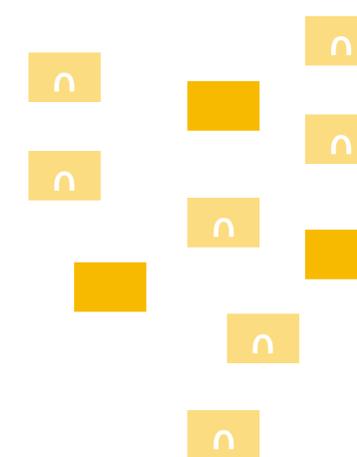
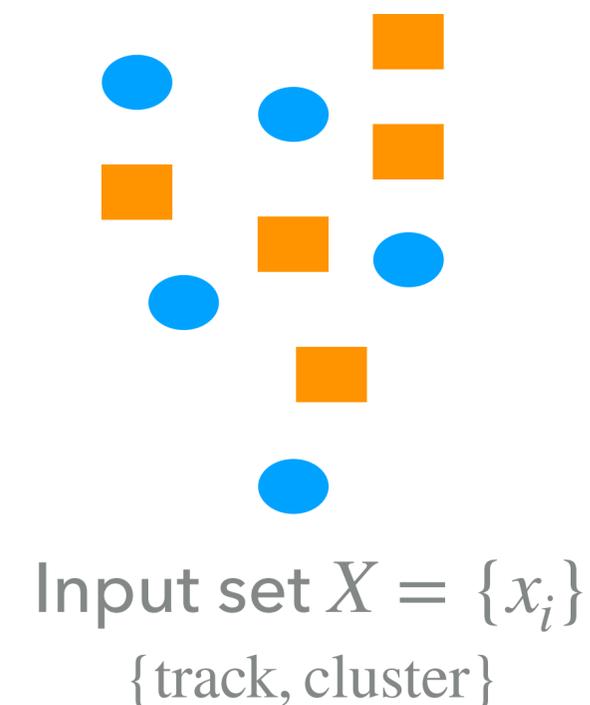
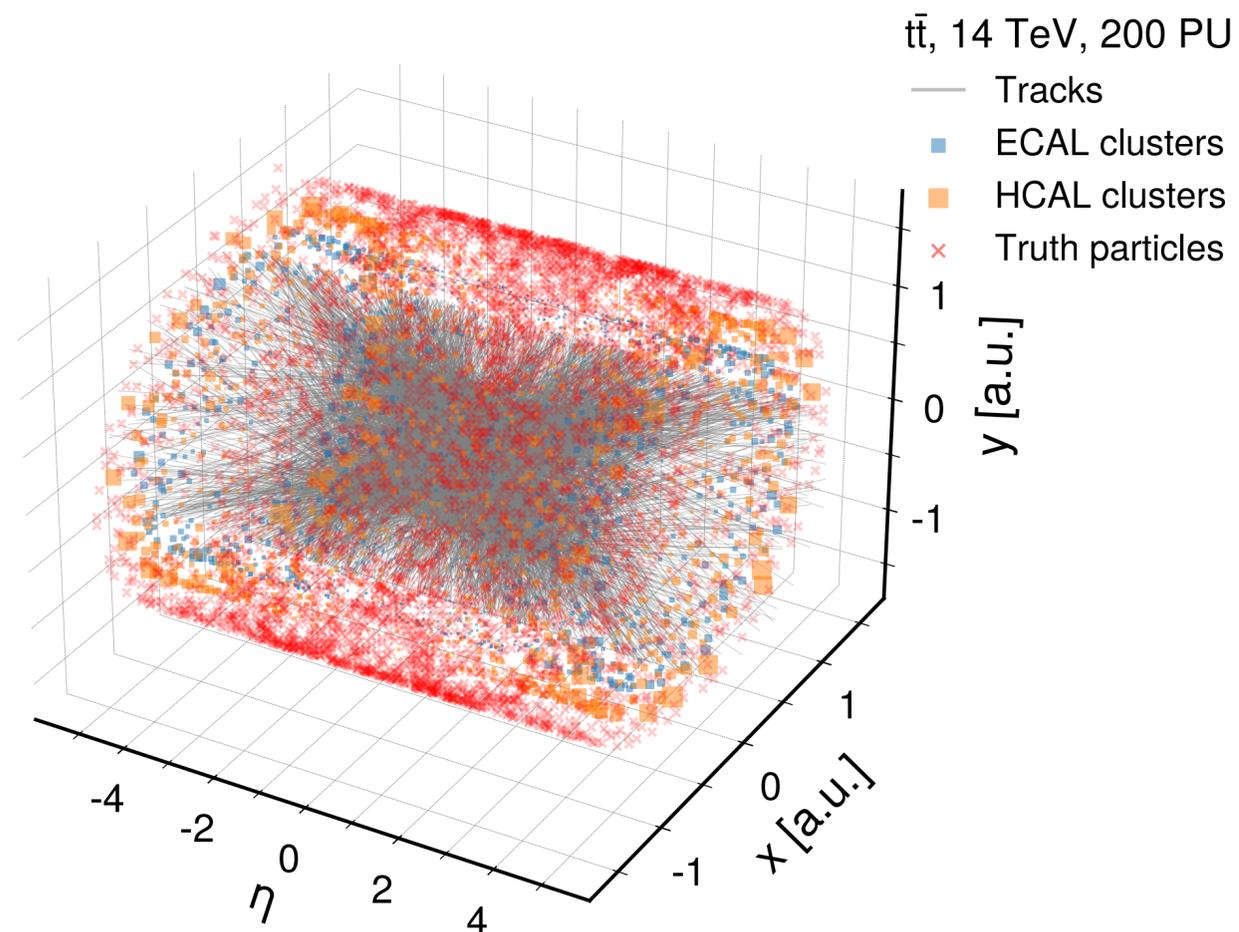


- Traditional Kalman Filter: computational expensive, serial operations
  - Efficient tracking with ML.
- GNN: Identifying edges compatible with real particle trajectories
- Implemented a GNN on an FPGA for particle tracking on small sectors
  - Found it can complete a tracking task in  $< 1 \mu\text{s}$ .

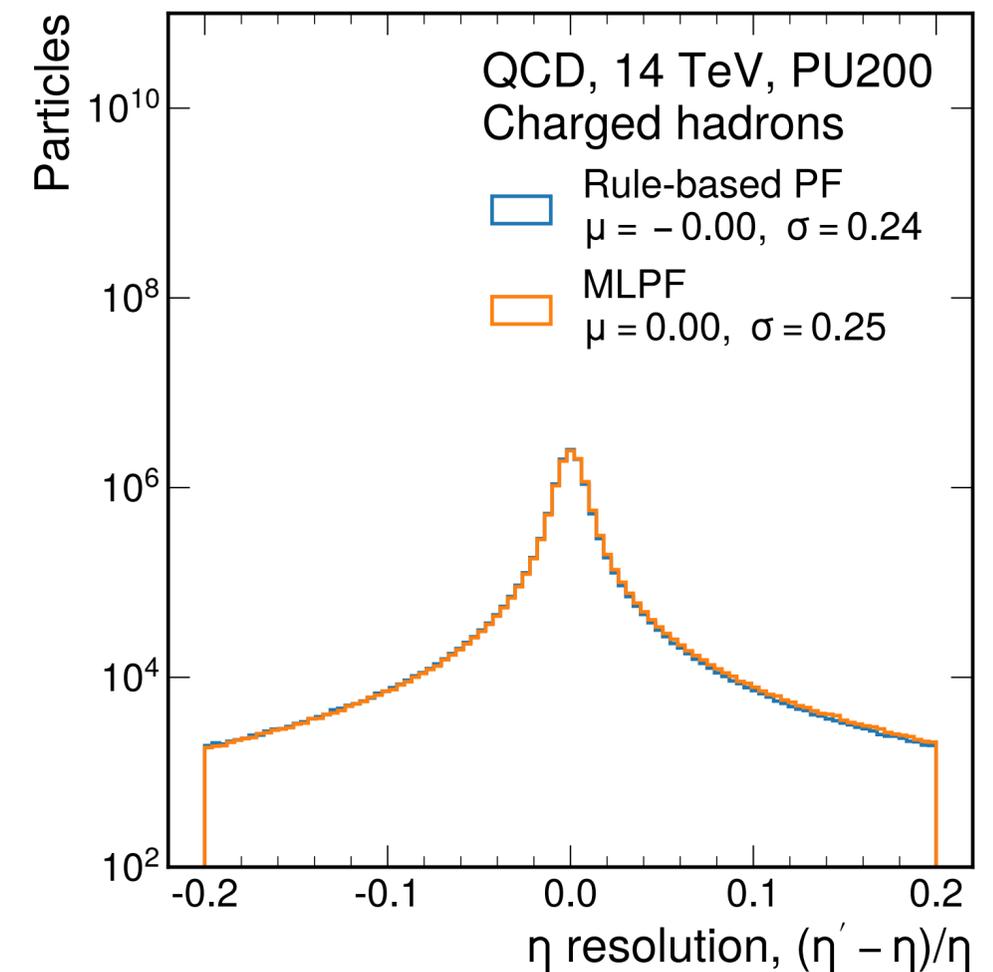
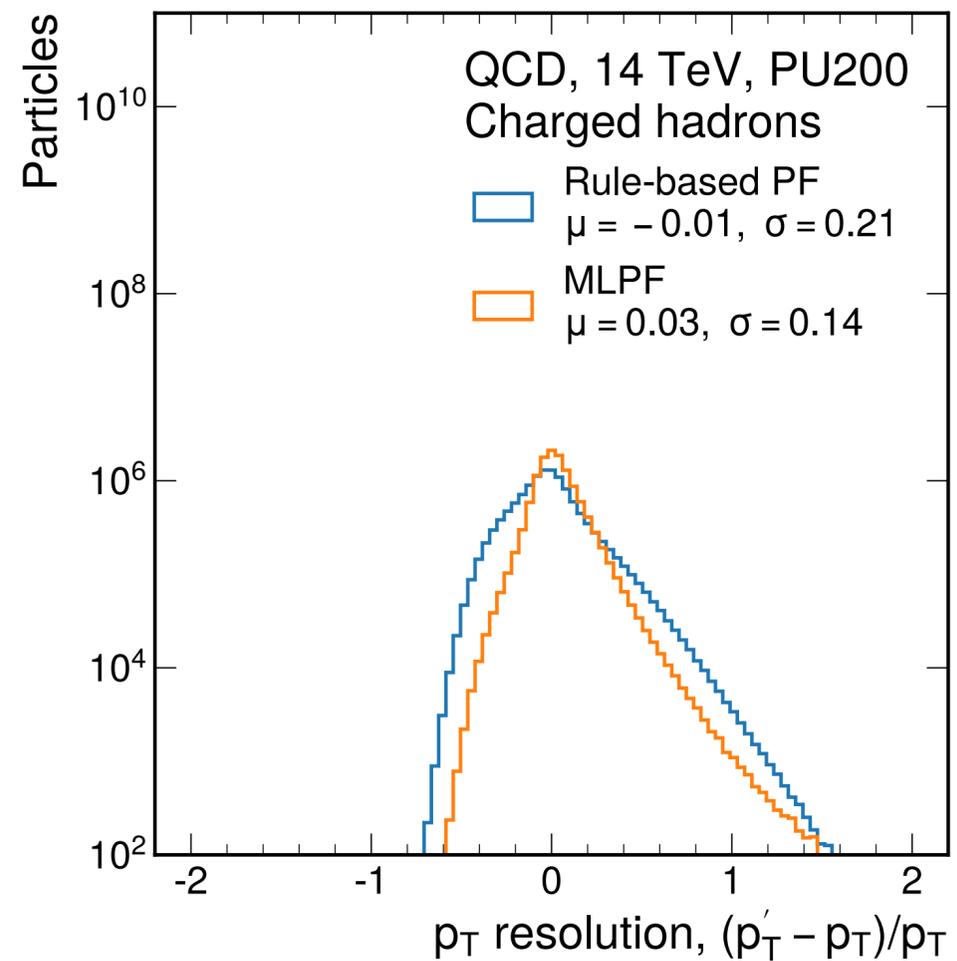
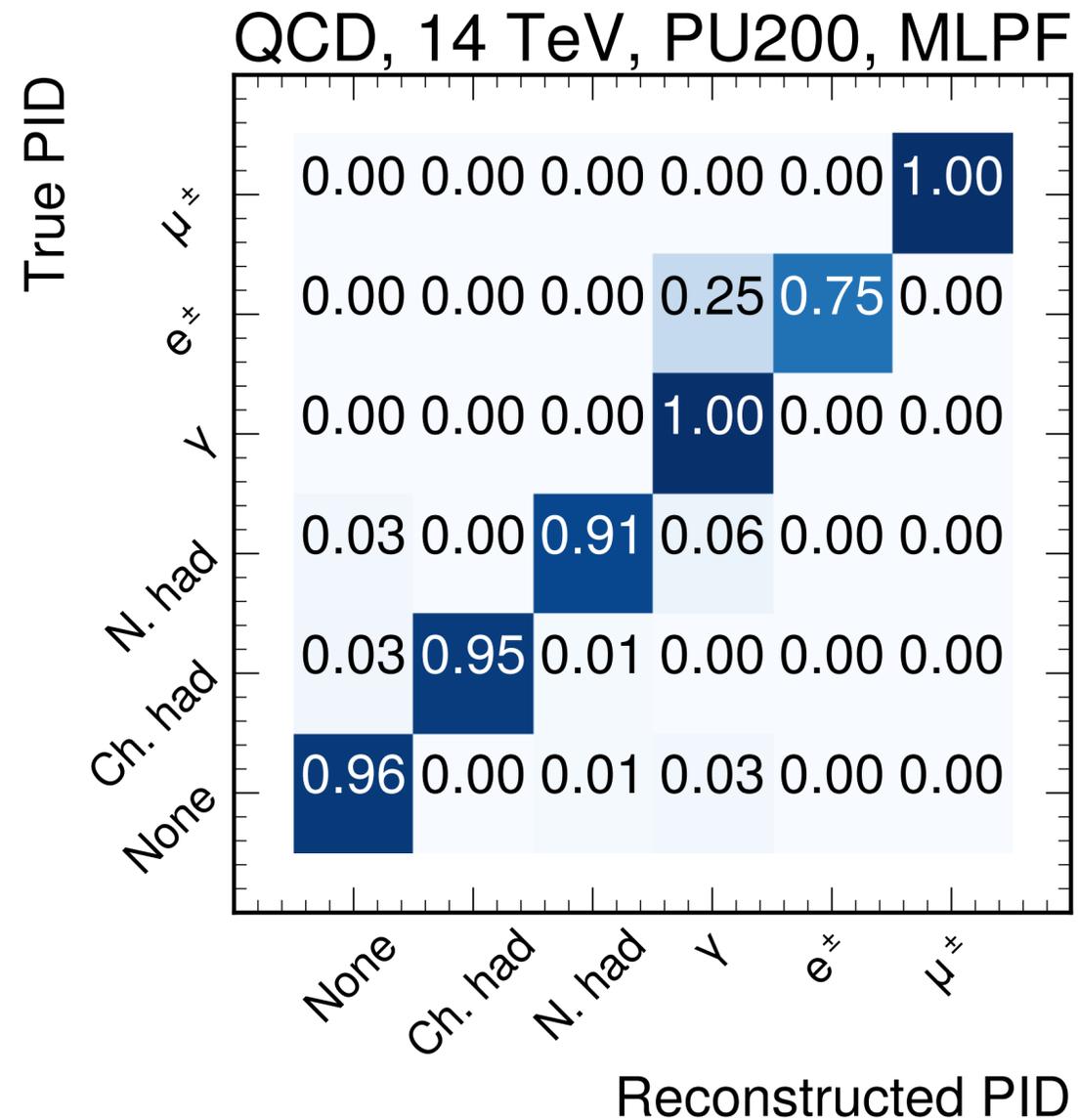
- Input set (or point cloud or heterogeneous graph) of tracks,
- Target set of **truth particles**  $Y = \{y_i\}$
- Goal: construct a mapping  $\mathcal{U}(X) = Y' \sim Y$  that minimizes some distance

$$||Y - Y'|| \equiv \sum_{i \in \text{event}} L(y_i, y'_i),$$

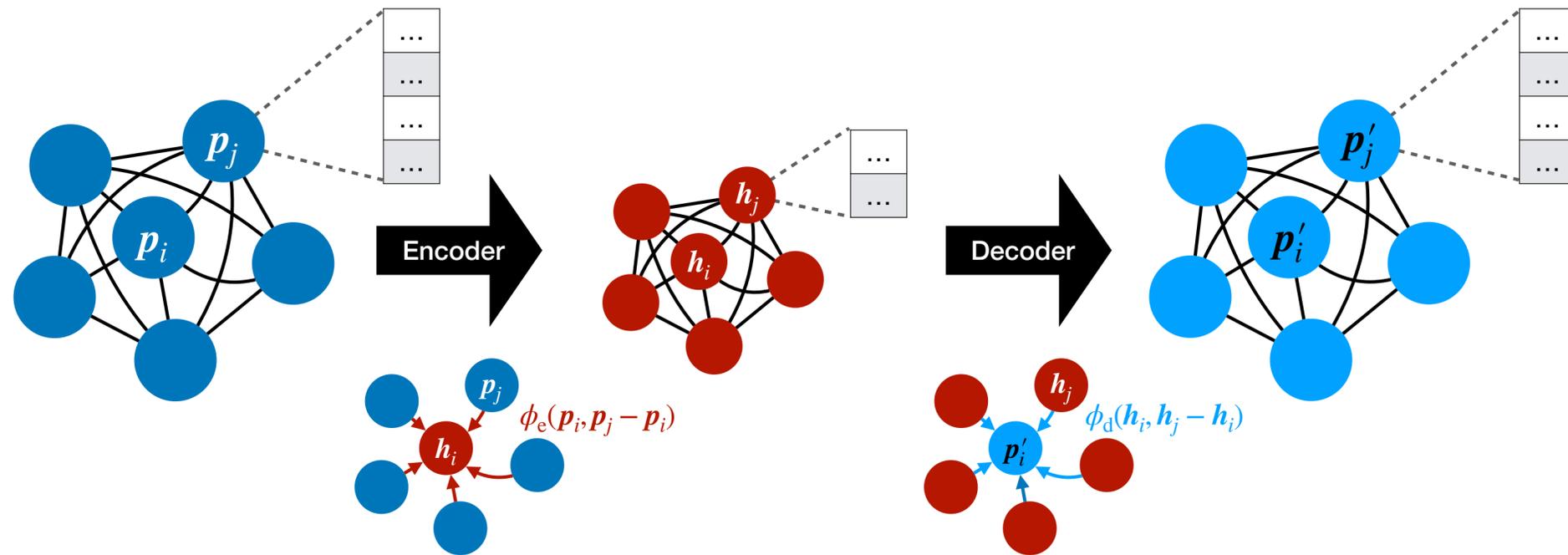
$$L(y_i, y'_i) \equiv \text{CLS}(c_i, c'_i) + \alpha \text{REG}(p_i, p'_i),$$



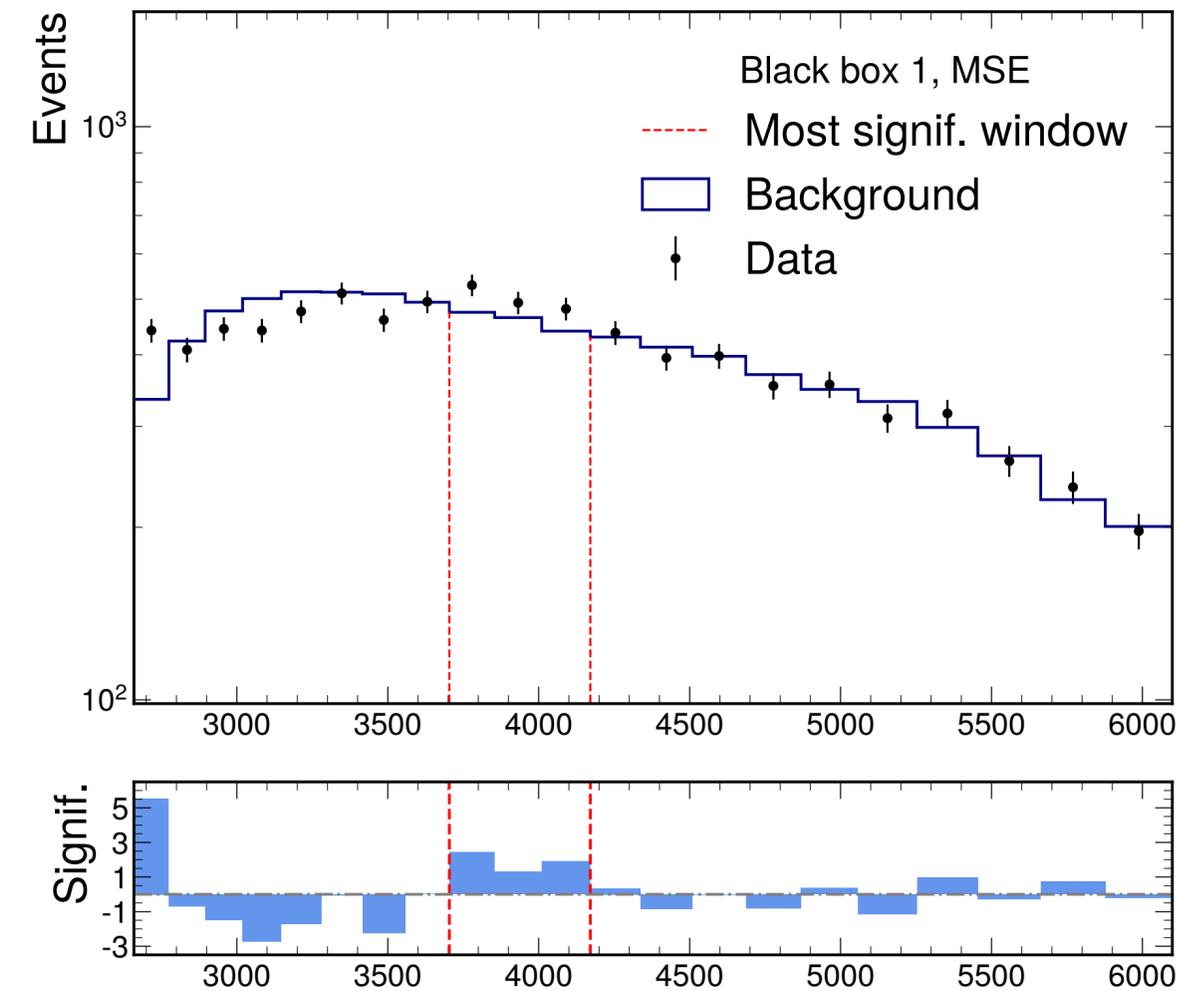
Target set  $Y = \{y_i\}$   
 {charged hadron, neutral hadron,  $\gamma$ ,  $e^\pm$ ,  $\mu^\pm$ }



- MLPF same or better than rule-based PF algorithm



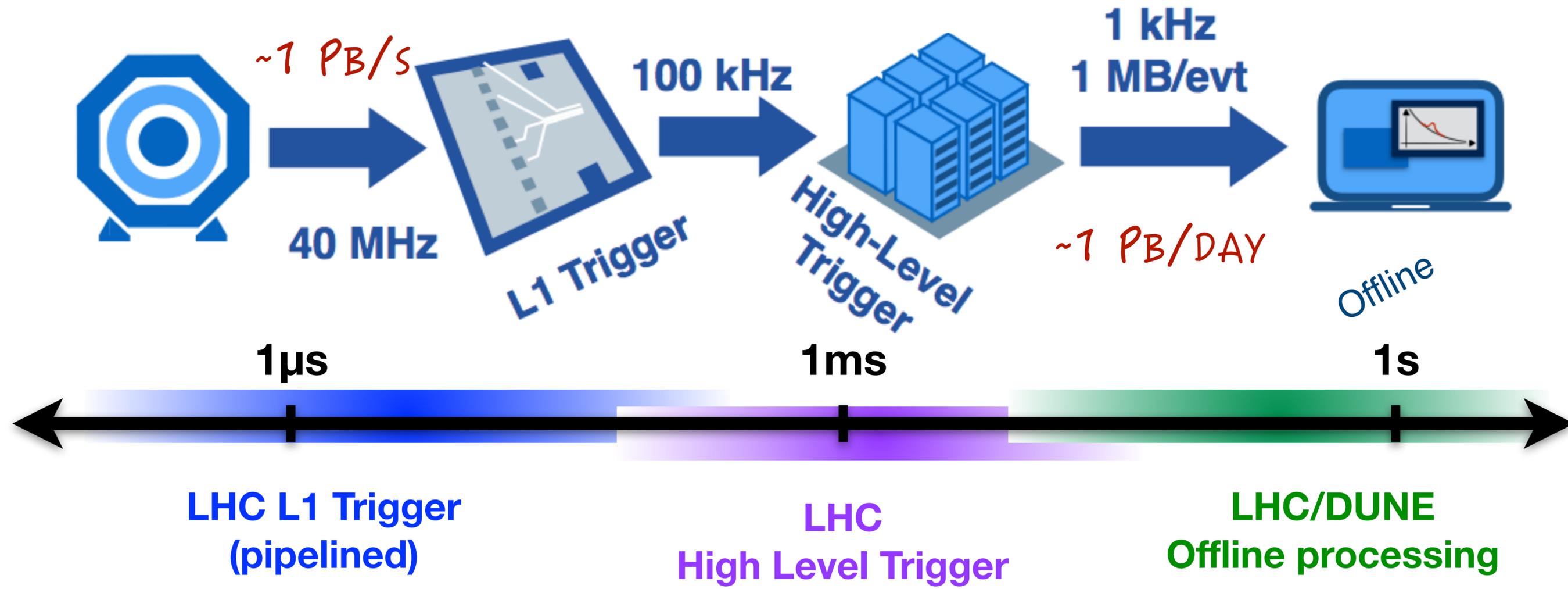
[arXiv:2101.08320](https://arxiv.org/abs/2101.08320)



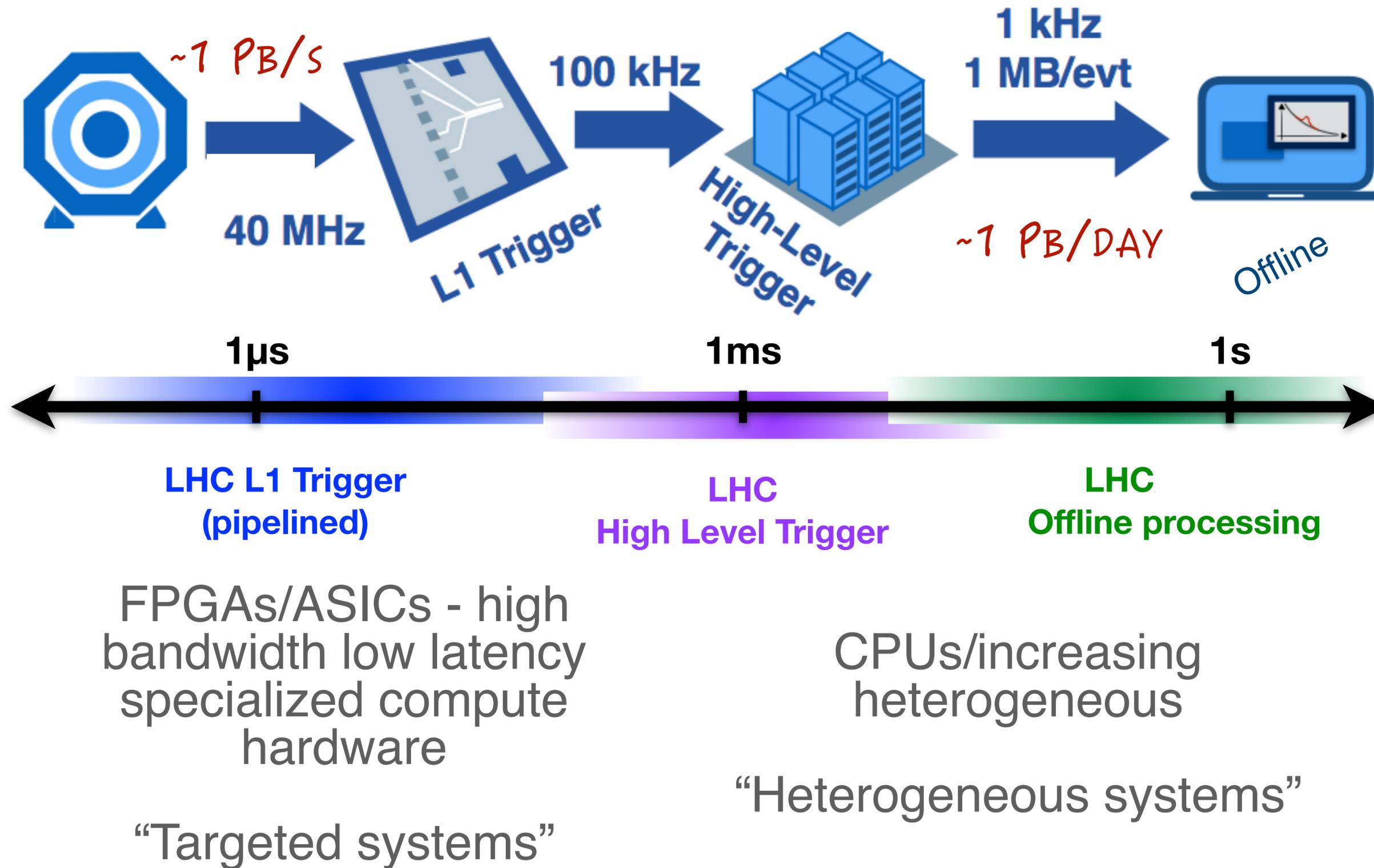
- **Autoencoders** compress data to a smaller representation and reconstruct it
  - Apply it to particles (E,  $\mathbf{p}$ ): train autoencoder on known SM data
- Found to be effective on a mock dataset with a signal injected at 3.8 TeV

The reality:  
Have to be fast enough

# From Collisions to Discoveries

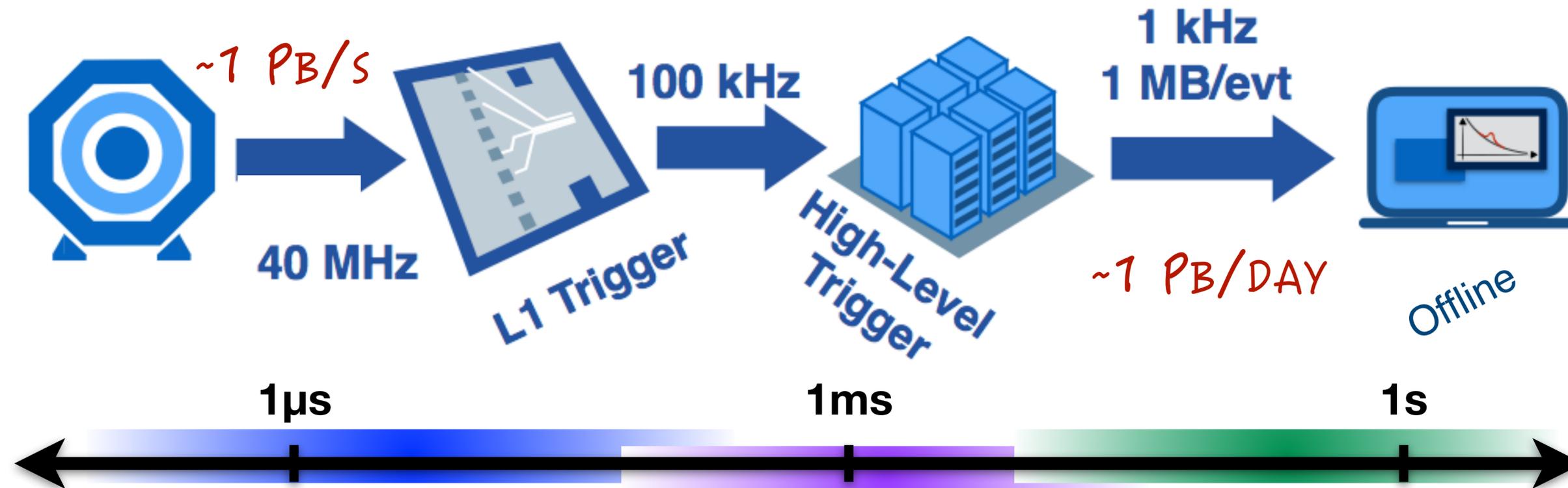


# Accelerate ML landscape



# Products developed in Fast ML

28



LHC L1 Trigger  
(pipelined)

LHC  
High Level Trigger

LHC  
Offline processing

## High Level Synthesis 4 Machine Learning

**Applications in CMS:** muon pT in endcap L1T, VBF tagger in L1T, BSM AE in L1T, ASIC data compression, AE in HGCAL, b-tagging/MET in L1T

## ML-as-a-service (SONIC):

Cost-effective, flexible, scalable, reduced software stack maintenance

Applications in CMS data processing workflows (ML or CUDA)

# Summary

---

29

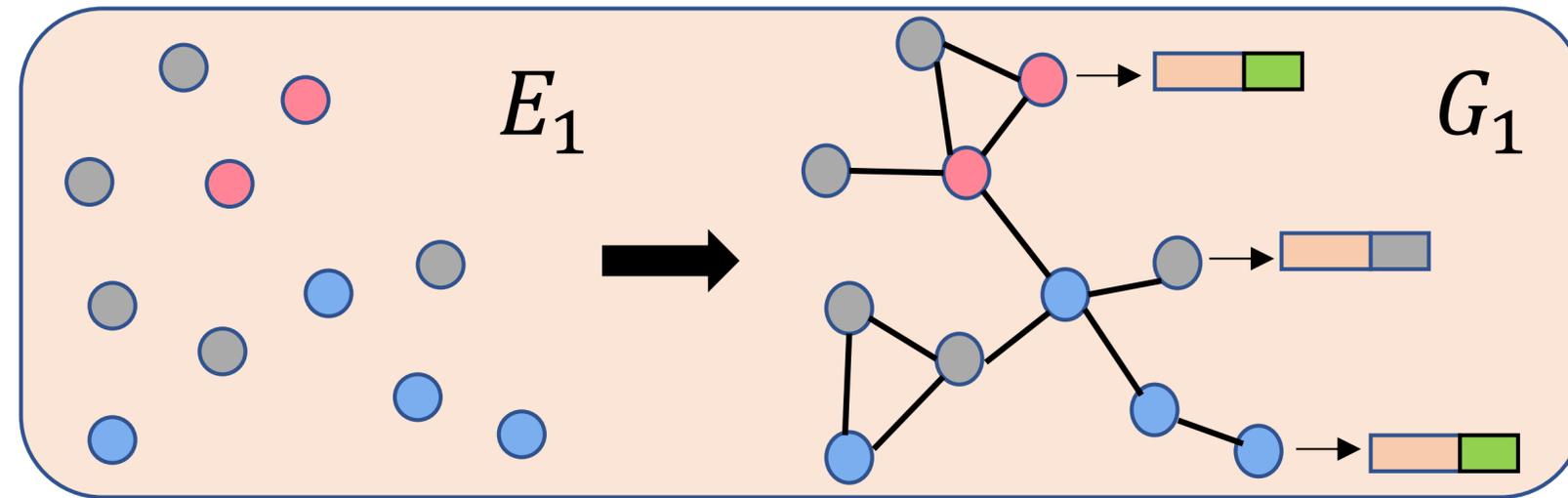
**Machine learning offers opportunities to significantly boost the discovery potential at the LHC.**

- Need accelerated machine learning inference in online & offline processing.
  - Multidisciplinary teams to realize optimal ML on targeted (e.g. L1 trigger) or heterogeneous systems (e.g. CMS HLT & offline).
  - User-friendly prototype tools for domain experts to streamline ideas to implementation.
  - Discussions focused not limited to reconstructions. e.g. fast simulation, experiment control/real-time calibration..etc

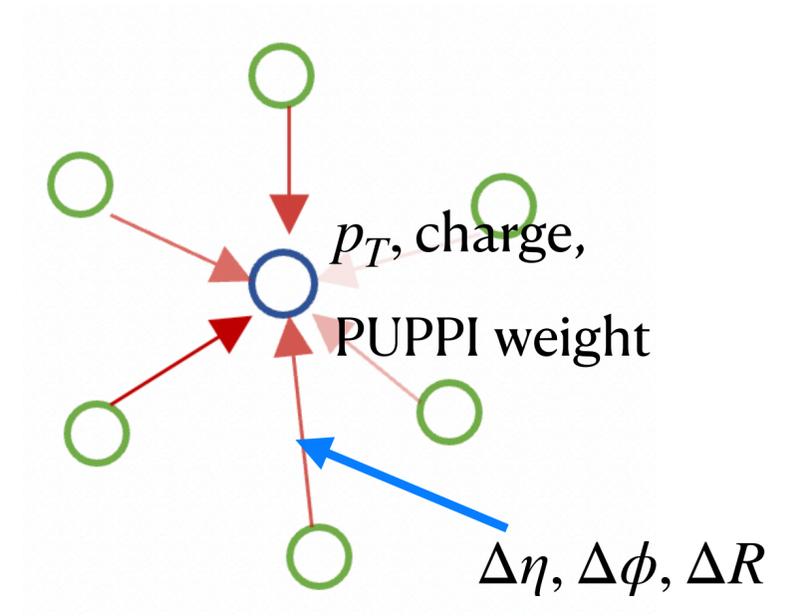
**Looking forward to continuing with tackling these challenges in A3D3 community**

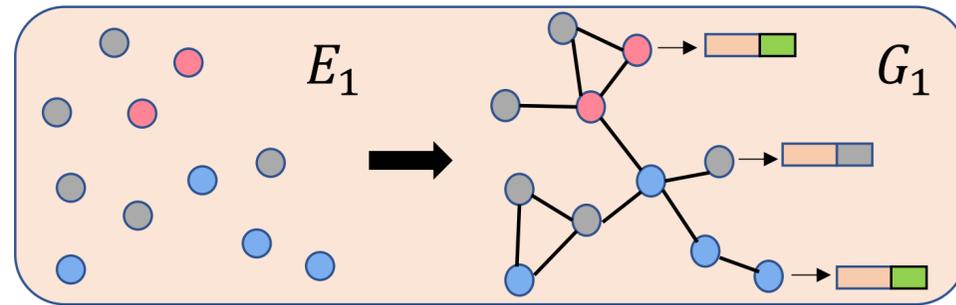
- Cross-domain synergies in ML mathematical formulations
- Advanced optimal ML methods in targeted&heterogeneous systems

# Semi-supervised Graph Puppi

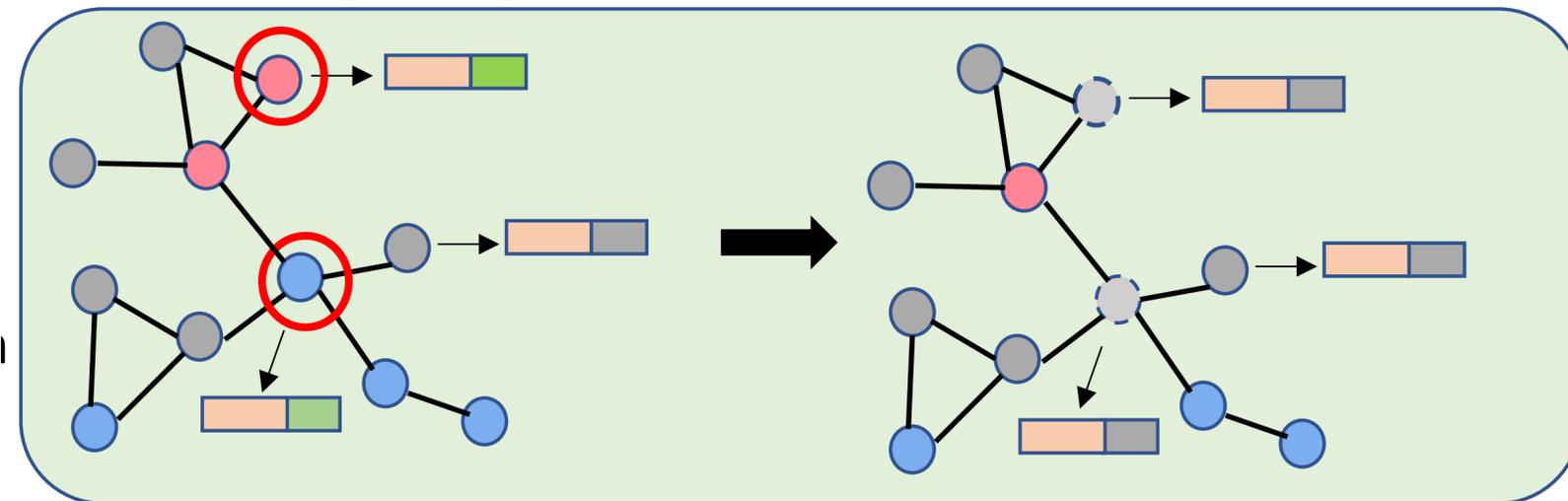


- Charged LV particles
- Charged PU particles
- Neutral particles
- Common feature domain
- Charged-specific feature domain
- Neutral-specific feature domain

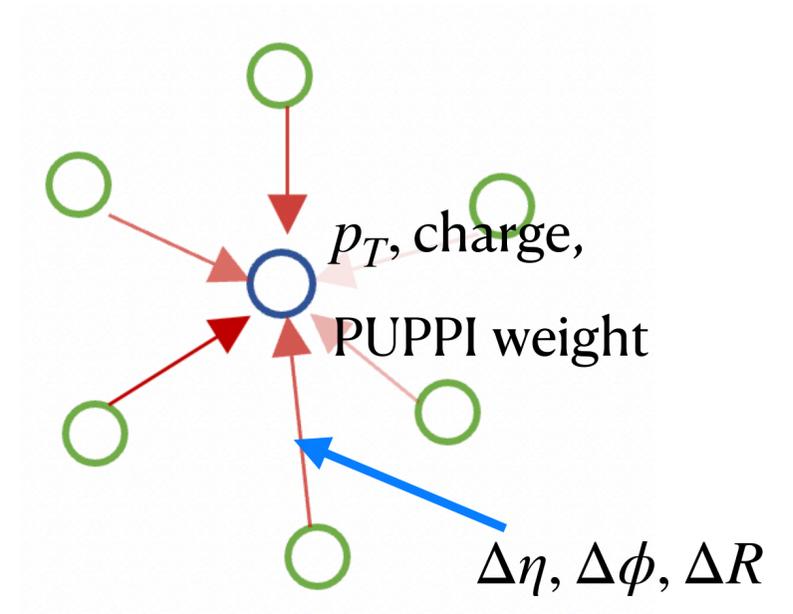




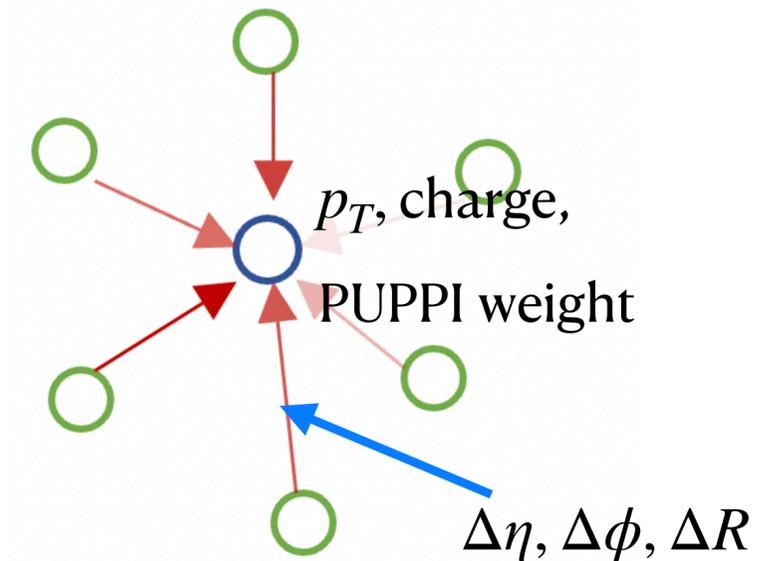
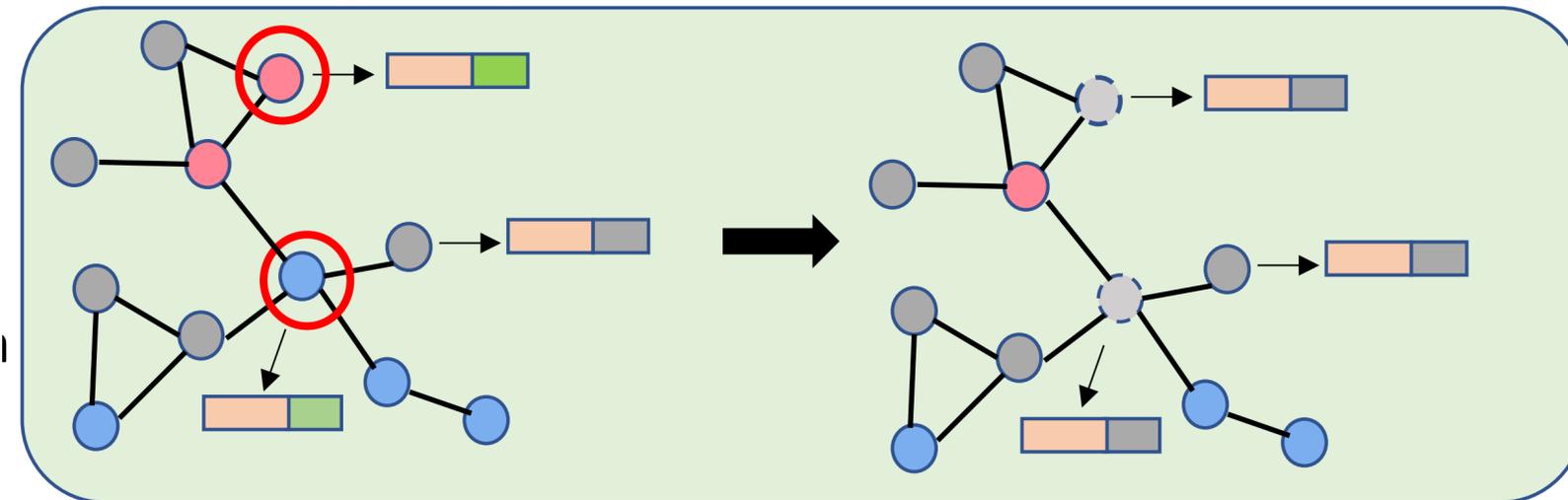
(b). Randomly select charged LV/PU particles, and mask charged-specific feature domain for training



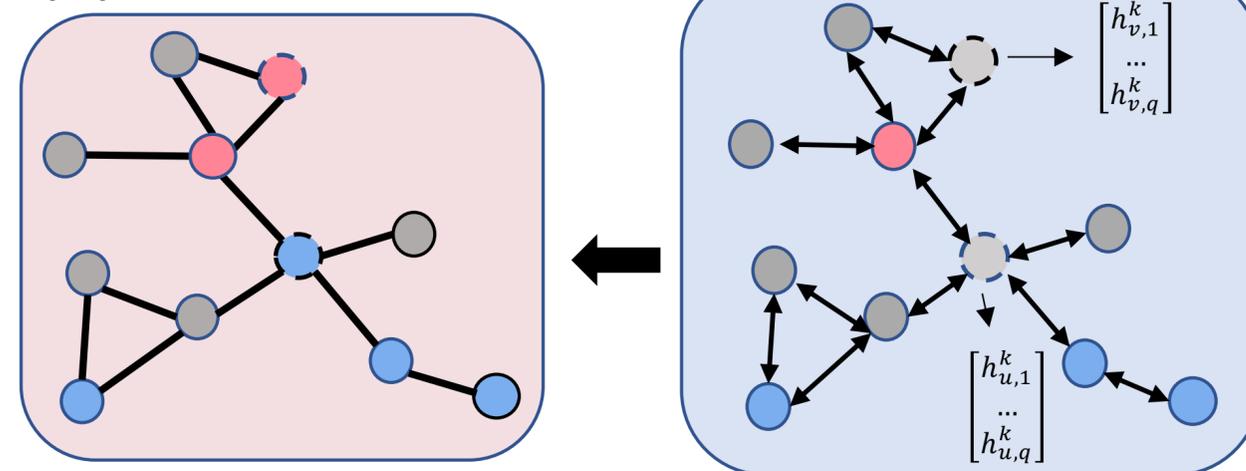
- Charged LV particles
- Charged PU particles
- Neutral particles
- Common feature domain
- Charged-specific feature domain
- Neutral-specific feature domain



(b). Randomly select charged LV/PU particles, and mask charged-specific feature domain for training



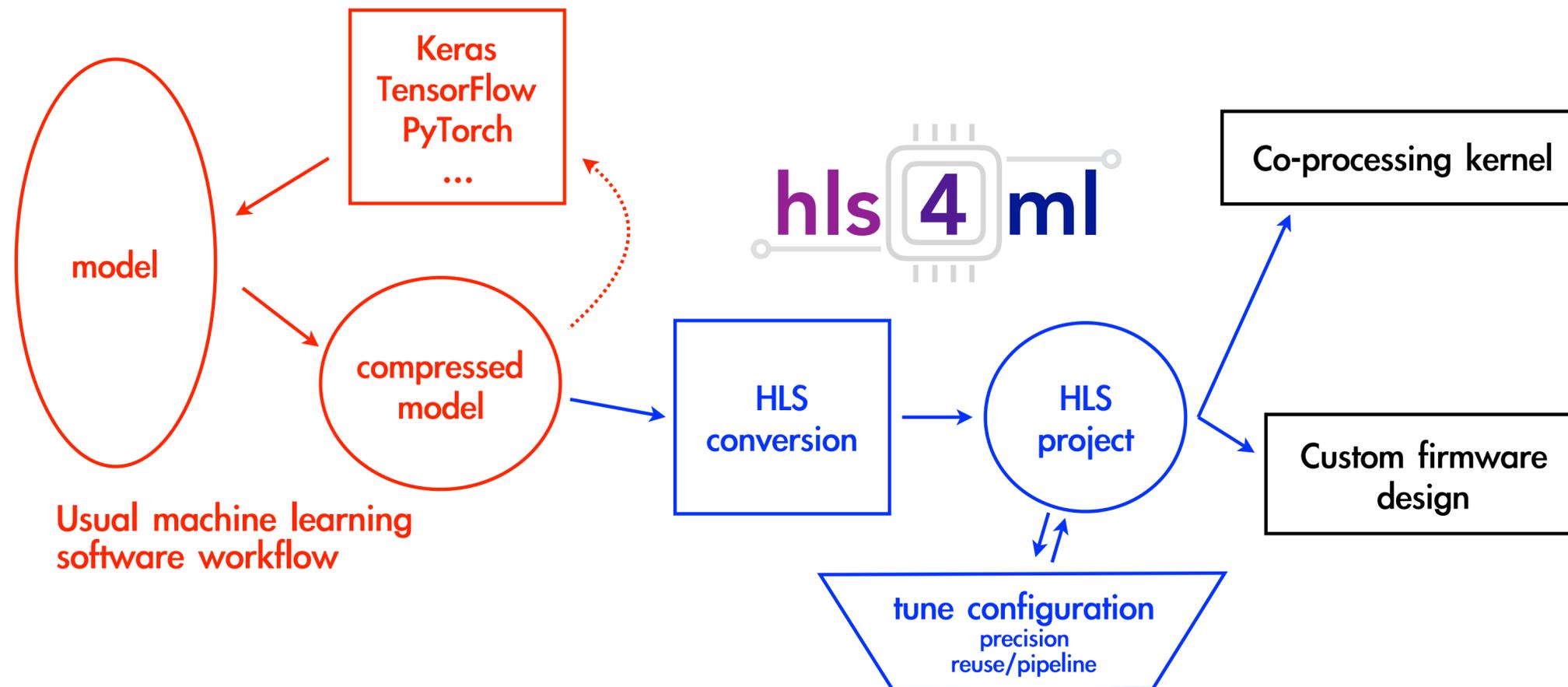
(d). Predict LV/PU



- Algorithm outperforms Puppi, comparable to fully supervised method. Can be adapted to different pile up conditions. No need for tuning as the particle itself is represented as a node.
- Presented at [BOOST 2021](#), Short version of the paper submitted to [NeurIPS 2021 AI for Science Workshop](#). Long version paper targeting PRD in preparation.
- Next: Apply to CMS simulation & data. Neutral particle vertex association in for the forward region.

# High-Level Synthesis 4 Machine learning

33



**HLS** - High Level Synthesis - compiler for C, C++, SystemC into FPGA IP cores

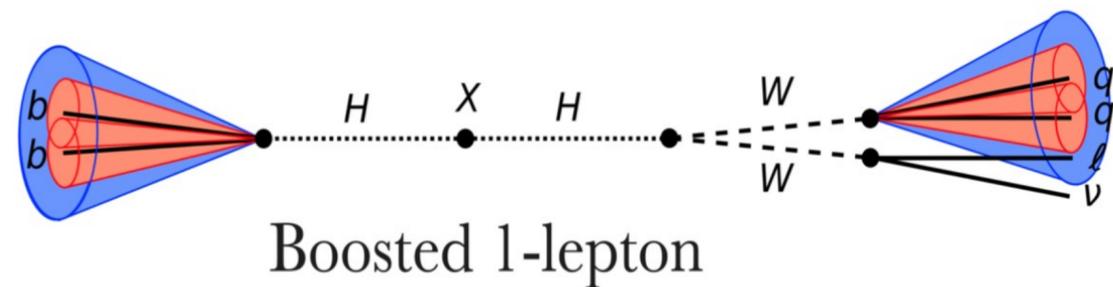
**HLS 4 Machine learning** :Prototype ML algorithms for FPGA WITHOUT Verilog/  
VHDL: firmware in a few hours

# Boosted Higgs Boson Tagging using ML

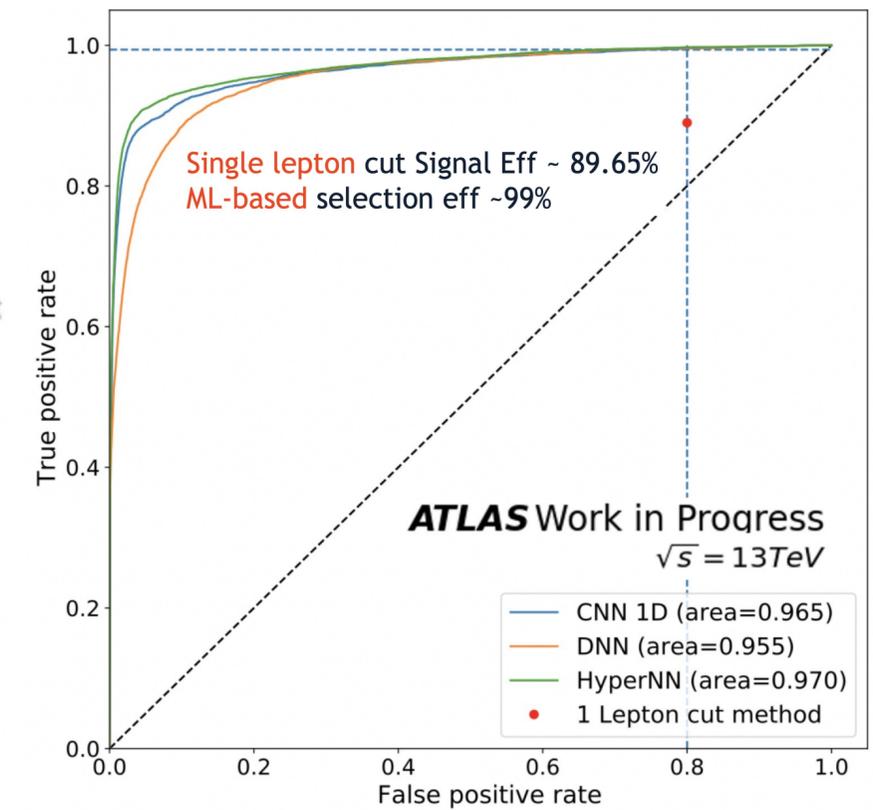
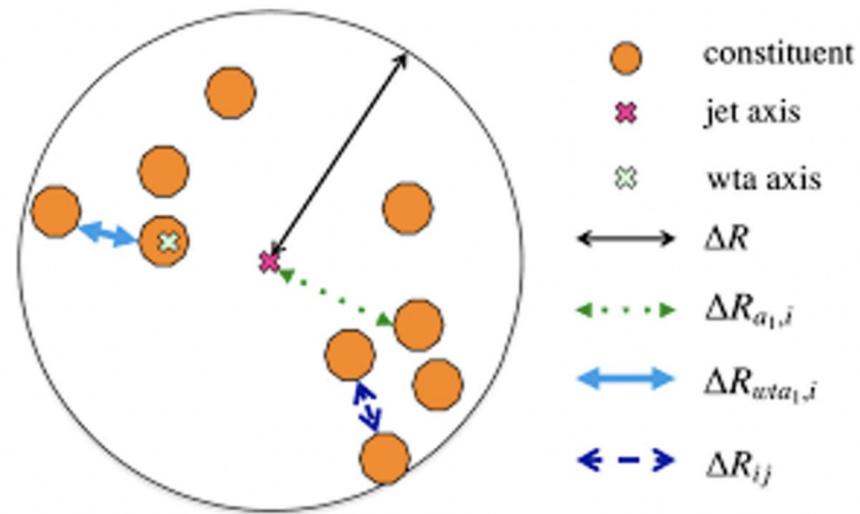
## Boosted $H \rightarrow bb$

## Boosted $H \rightarrow WW^{(*)}$

High  $p_T$   $H \rightarrow WW^{(*)}$  semileptonic decays, e.g. from a heavy particle  $X \rightarrow HH$  decay, are difficult to identify due to merging of lepton and jets



Train Neural Network taggers on jet 4-vectors and jet constituents



# As-a-service Computing Model

How many CPU can this GPU serve?

CPU-to-GPU ratio:

$$t_{total\_cpu} - t_{othercpu} = t_{ml}$$

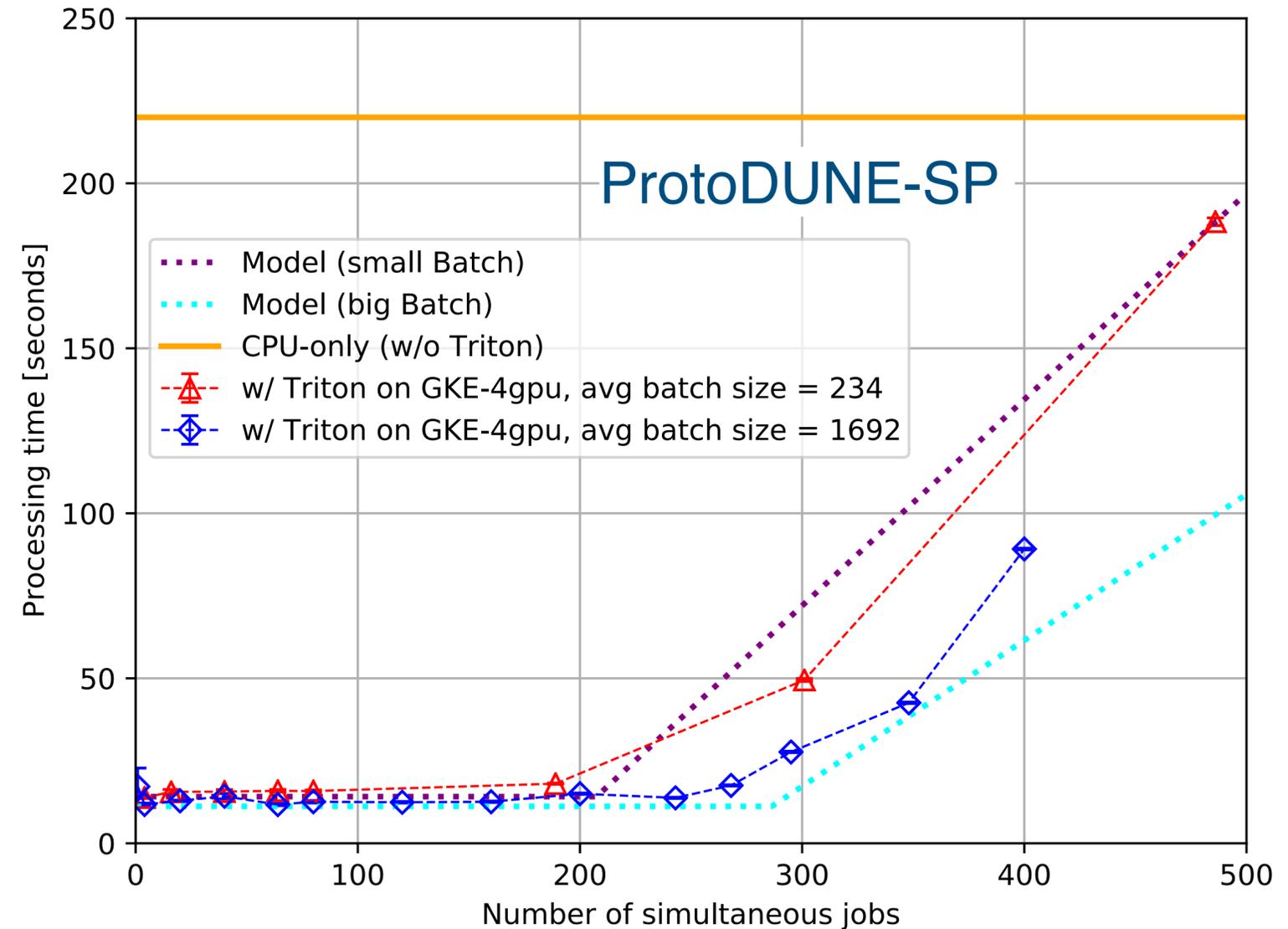
$$t_{total\_sonic} - t_{othercpu} - t_{sonic\_cpu} = t_{transfer} + t_{scheduling} + t_{sonic\_gpu}$$

$$t_{ml} / (t_{sonic\_cpu} + t_{transfer} + t_{scheduling} + t_{sonic\_gpu}) =$$

$$t_{sonic\_cpu\_part} = t_{sonic} - t_{transfer} - t_{scheduling}$$

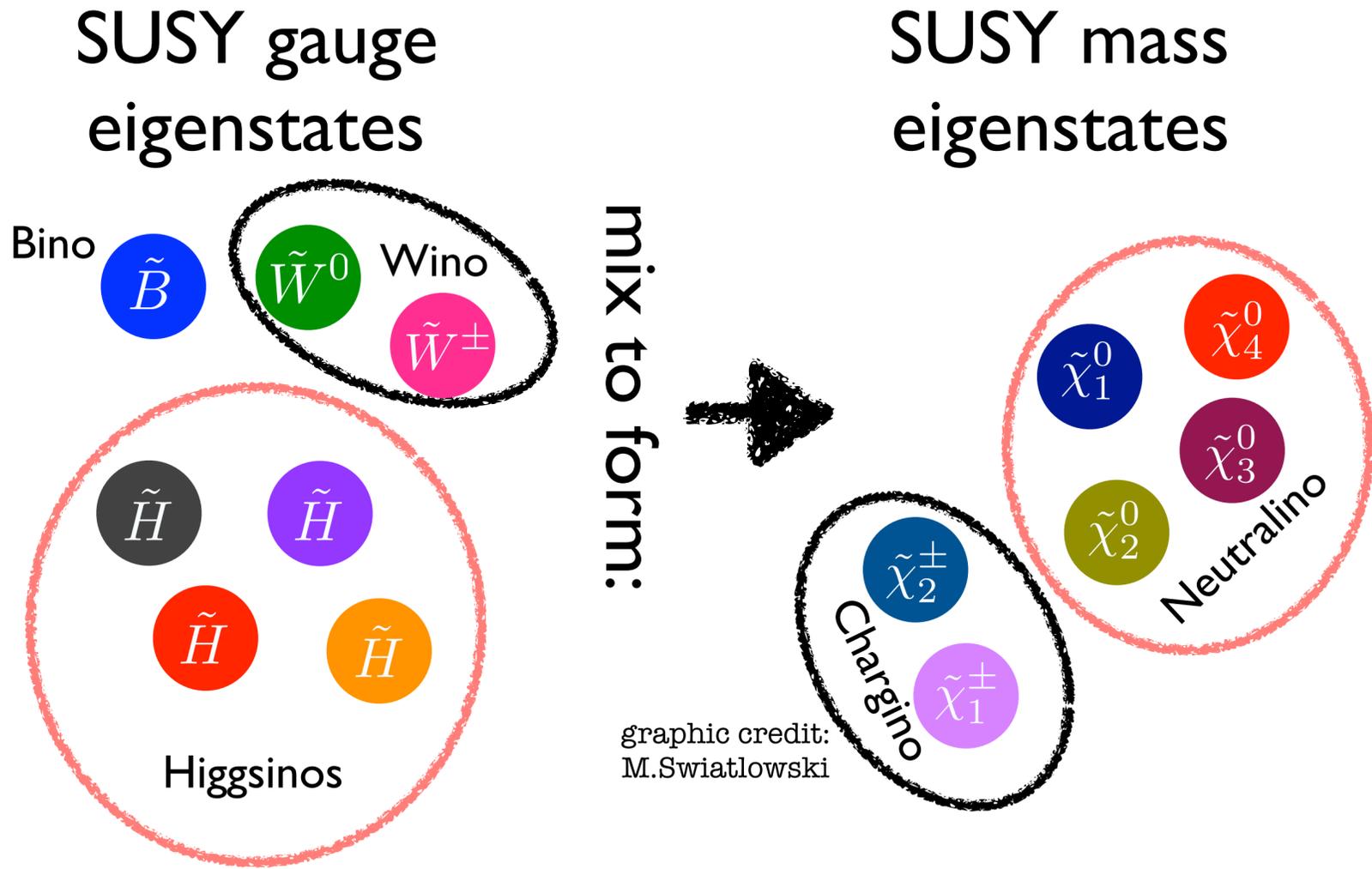
$$ratio = t_{sonic\_cpu\_part} / t_{gpu}$$

Passes this ratio, GPU saturates and average processing time increase.

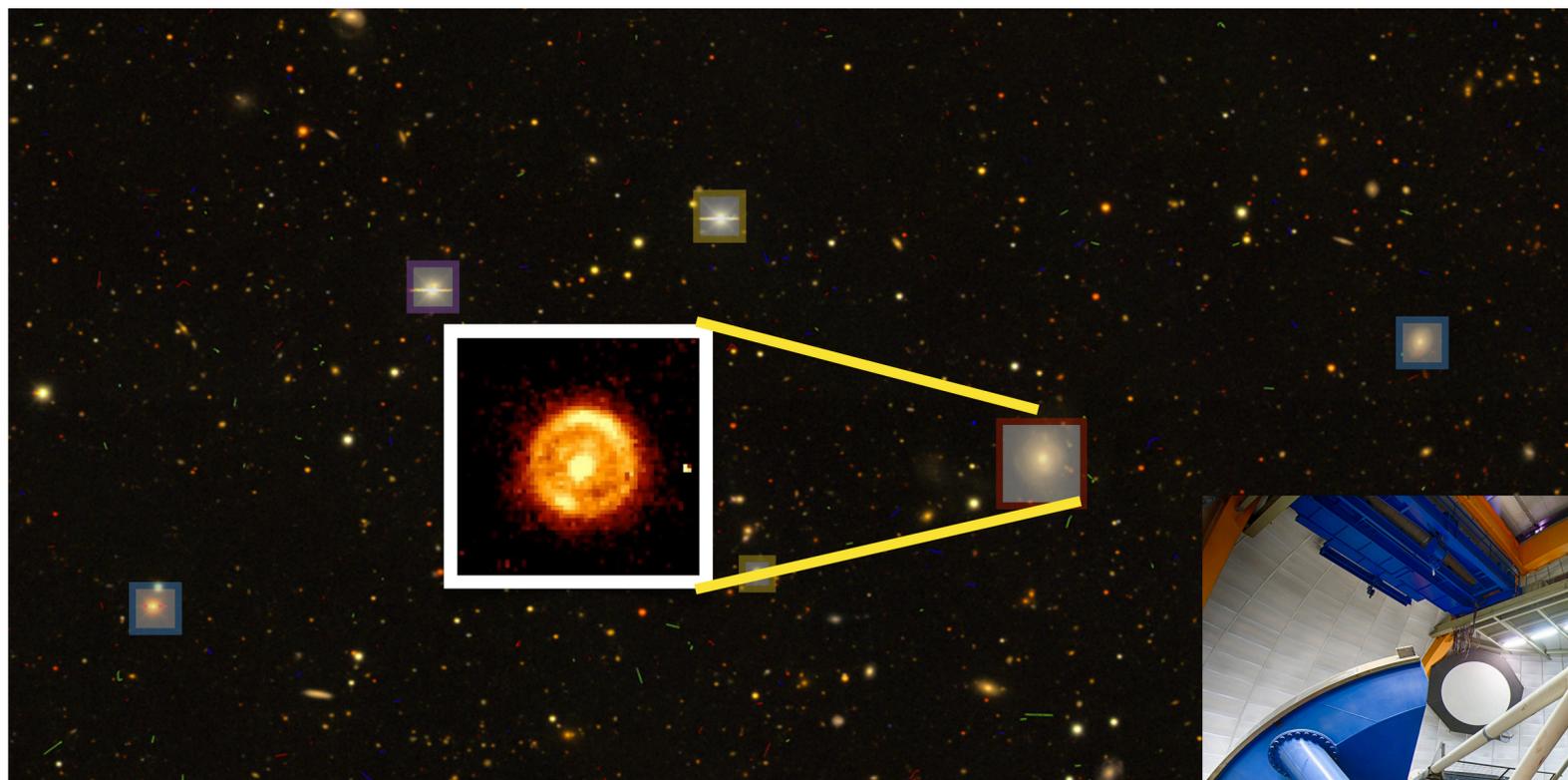


- SUSY partners of the SM electroweak sector:
  - $U(1) \rightarrow$  Bino,  $SU(2) \rightarrow$  Winos
  - Higgs  $\rightarrow$  Higgsinos
  - Leptons  $\rightarrow$  sleptons
- Could be light and accessible at the LHC

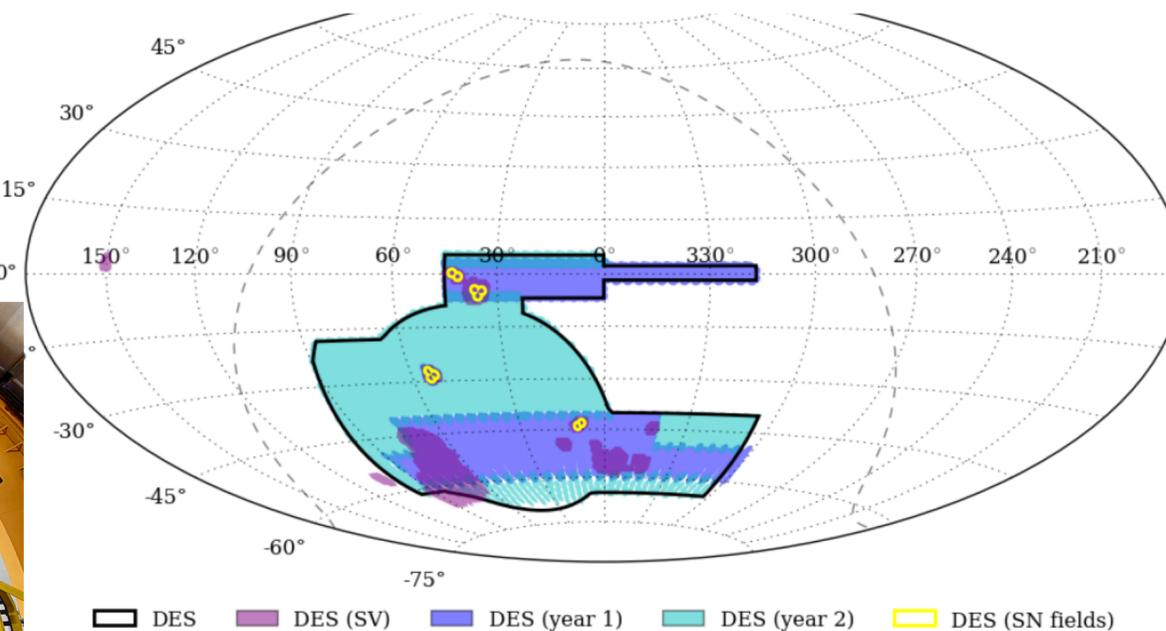
Important to search for them at the LHC!



# ML in the Sky: it's full of stars



## Dark Energy Survey: Sky Footprint of Observations



Populations of objects show **dark matter, dark energy**

**Region-based CNNs** on heterogeneous compute devices

- LSST: 20 Tb / night
- 1 Billion transient alerts /night

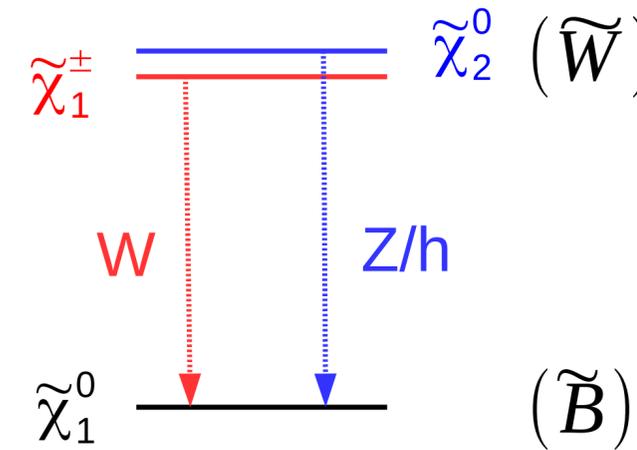
### Challenge of scheduling on multiple time

**Long:** competition between **faint galaxies**, **transient objects**

**Short:** Weather, annual modulation of sky positions

**Smart telescopes:** reinforcement learning for optimal scheduling and control

- 3000 fb-1 data expected at the HL-LHC
- e.g. Higgsinos: Low cross section, challenging signatures
  - $\Delta m \sim \text{tens of GeV}$  : Soft decay products
  - $\Delta m \sim \text{hundreds of MeV}$  : Long-lived signatures



Wino-like  
45 fb\*

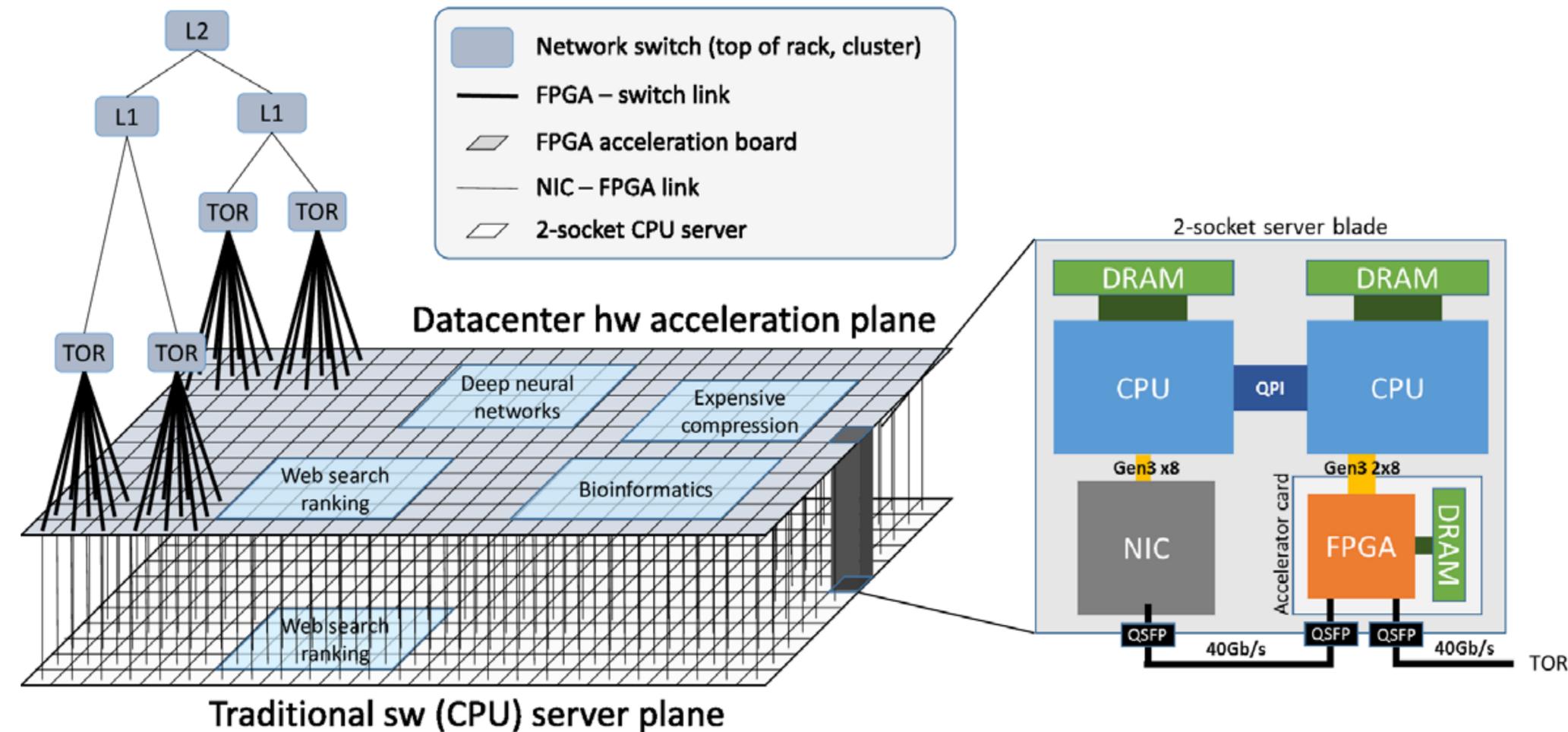
$\Delta m \sim \text{hundreds of MeV to tens of GeV}$

Higgsino-like  
11 fb\*

\* Cross-sections for 500 GeV sparticles @ 13 TeV ( $\tilde{\chi}_2^0 \tilde{\chi}_1^\pm$  only)

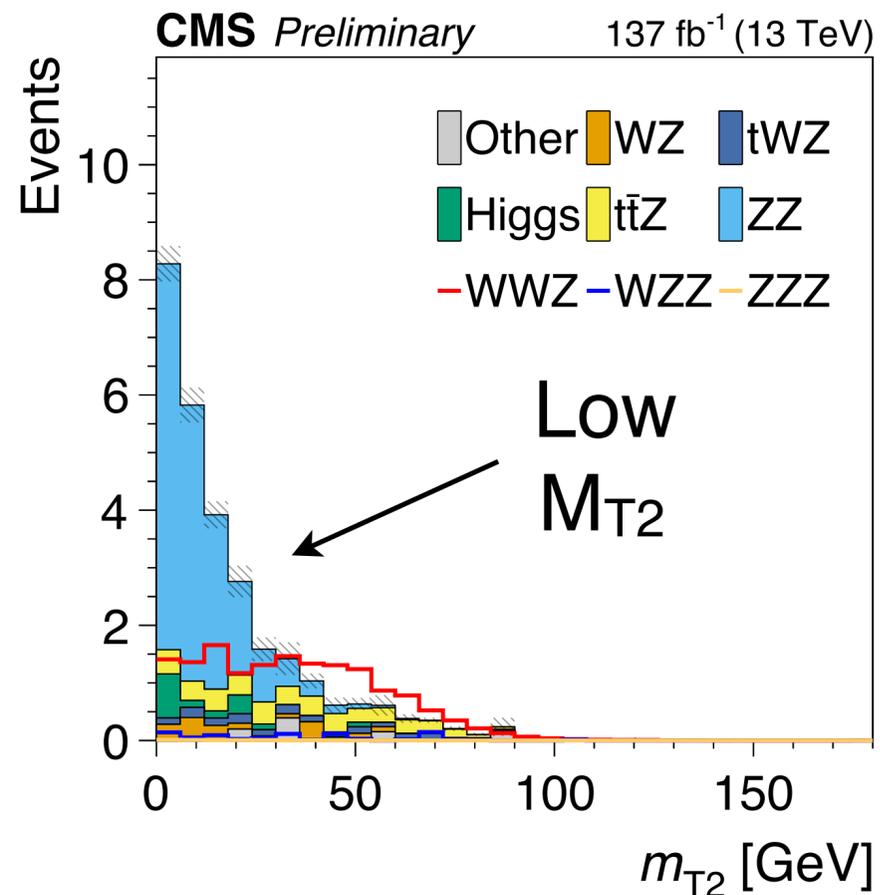
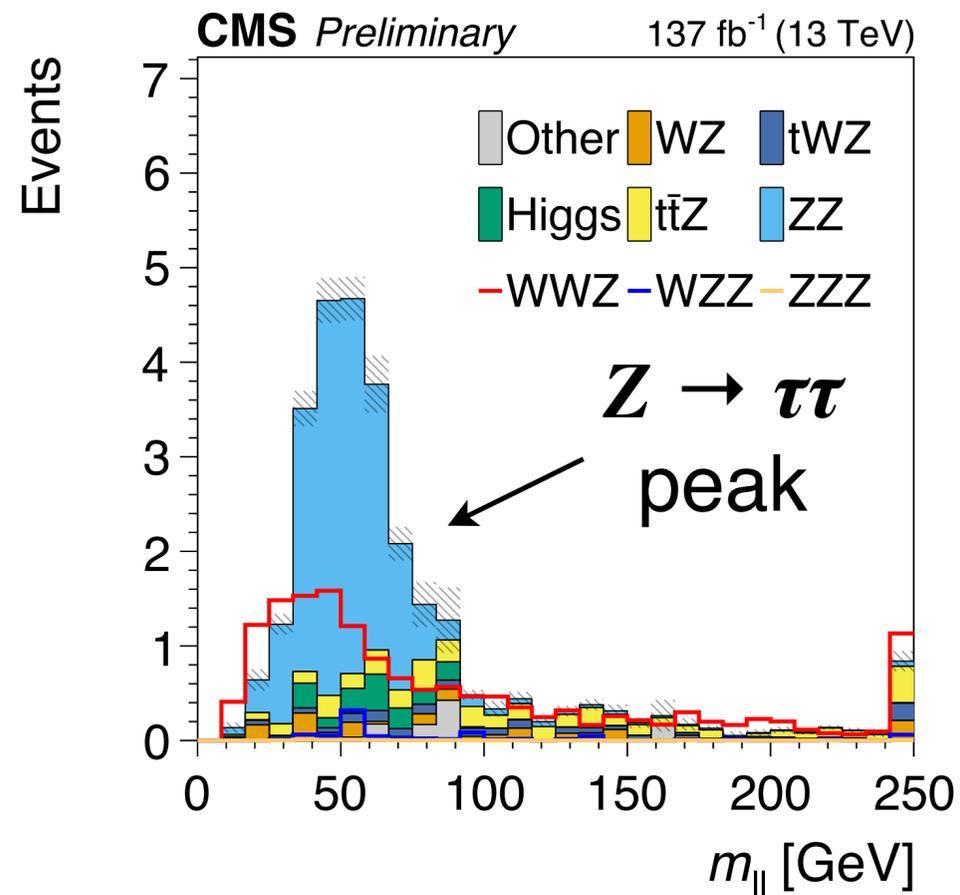
# Microsoft Brainwave

39



- **Mature service** at scale (more than just a single co-processor)
- Multi-FPGA/CPU fabric accelerates *both* **computing** and **network**
- Models supported:
  - ResNet50, ResNet152, DenseNet121, VGGNet16...
  - Partially fixed neural network architecture. weights can be retuned.

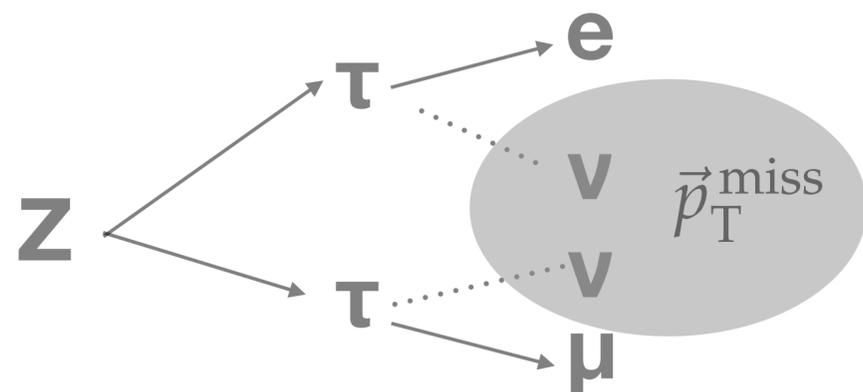
# WWZ: smaller rate but clean



40% smaller than WWW ,  
leptonic decays

$4\ell$ : Tag  $Z \rightarrow \ell\ell$ ,  $WW \rightarrow (e\mu/ee/\mu\mu)$

- $e\mu$  shown as example:  
kinematic selections against  
ZZ



$$m_{T2} = \min_{\vec{p}_T^{\nu(1)} + \vec{p}_T^{\nu(2)} = \vec{p}_T^{\text{miss}}} \left[ \max \left( m_T^{(1)}(\vec{p}_T^{\nu(1)}, \vec{p}_T^e), m_T^{(2)}(\vec{p}_T^{\nu(2)}, \vec{p}_T^\mu) \right) \right]$$

Edge data box at Feymann computing center



Docker container on server  
(PCIe connection):

$14 \pm 25$  ms

Fermilab computing  
cluster:

$20 \pm 30$  ms

Local laptop:

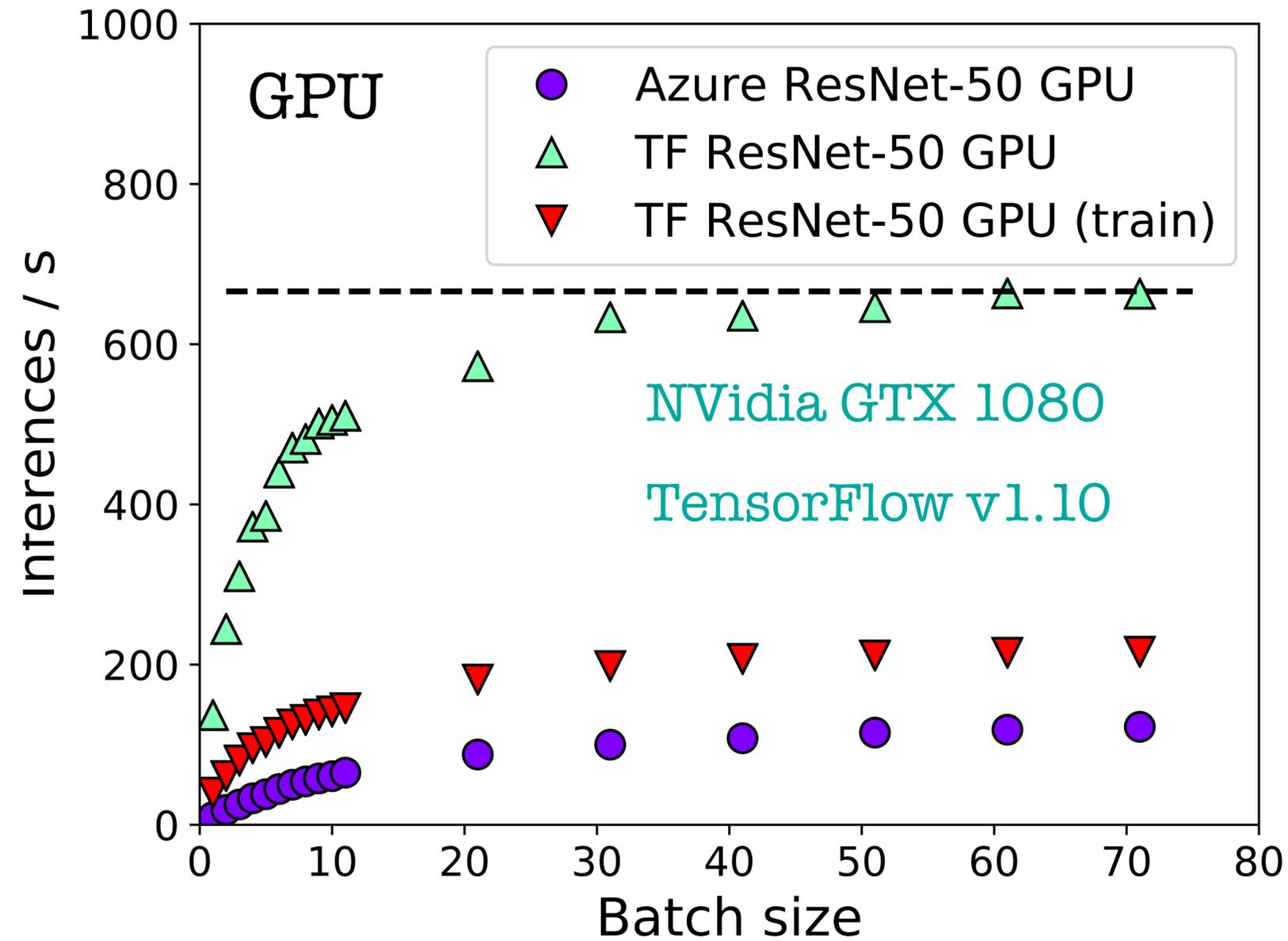
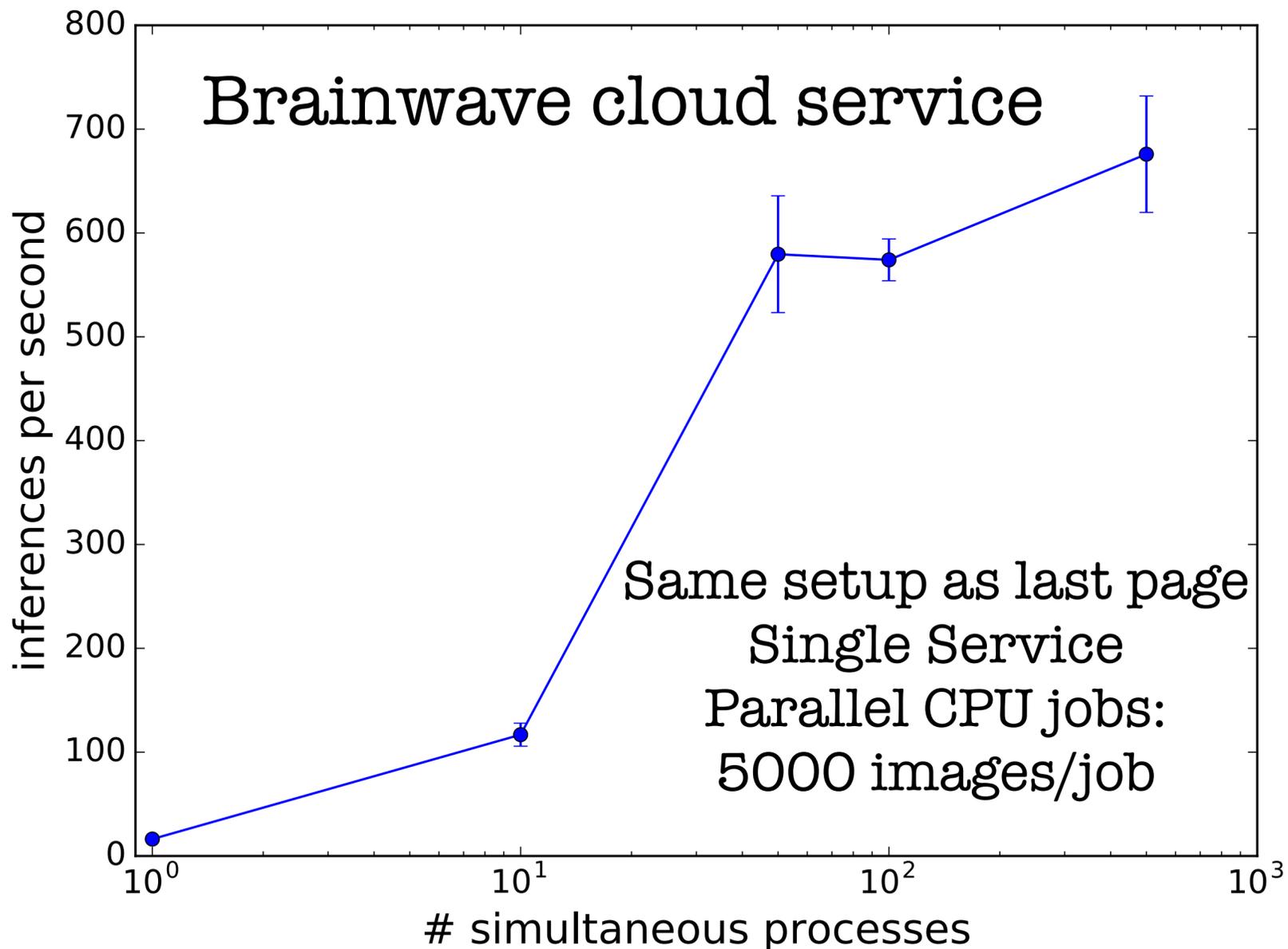
$68 \pm 27$  ms

CERN (Geneva):

$168 \pm 62$  ms

- Gain experience in deploying co-processors in local clusters with cloud native tools: docker image, kubernetes
- Benchmark latency and scaling performance, compare with previous studies
- Can be used for neutrino and cosmology experiments as of ~today—> next slide

# data throughput compared to



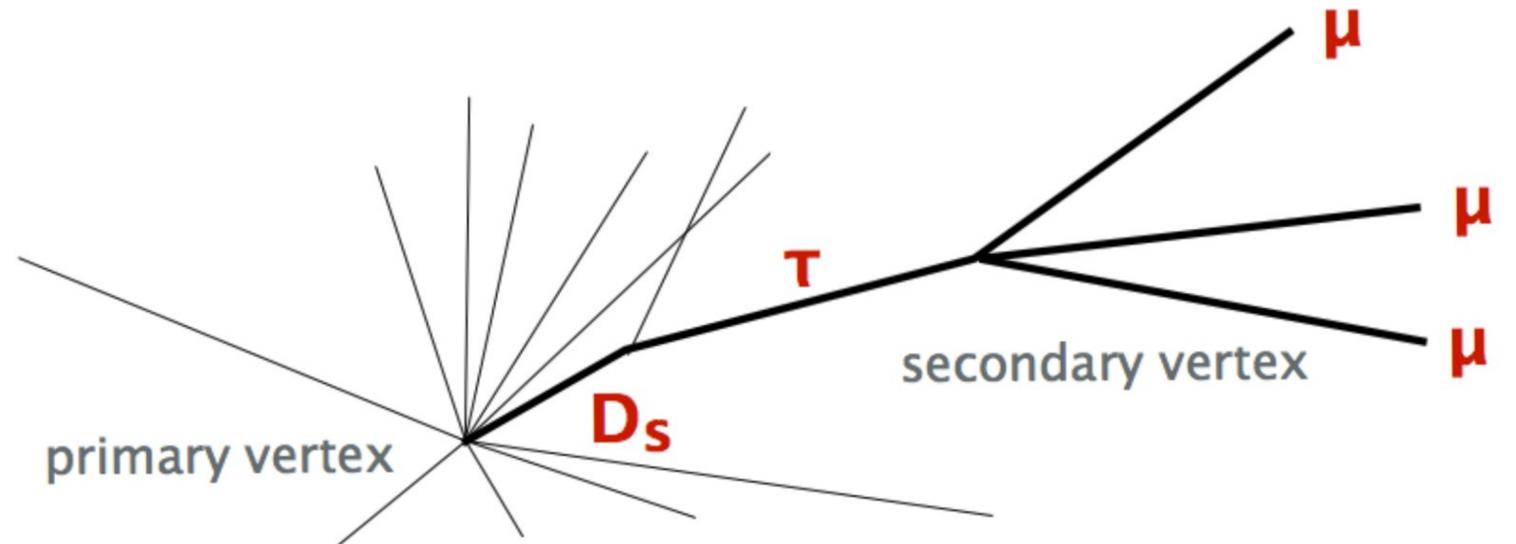
Comparable max data throughput: 600-700 images/sec

**Neutrinos oscillate:**  
Lepton number not conserved



**At the LHC**

**What about in charged leptons?  $\tau \rightarrow 3\mu$**



SM:  $10^{-25}$

Currently:  $< 10^{-8}$



**Need to cope with more challenging LHC environment in Run 2 & Run 3 ( $300 \text{ fb}^{-1}$ ) until HL-LHC upgrade (2023).**

**Module designed to reduce dynamic inefficiency**

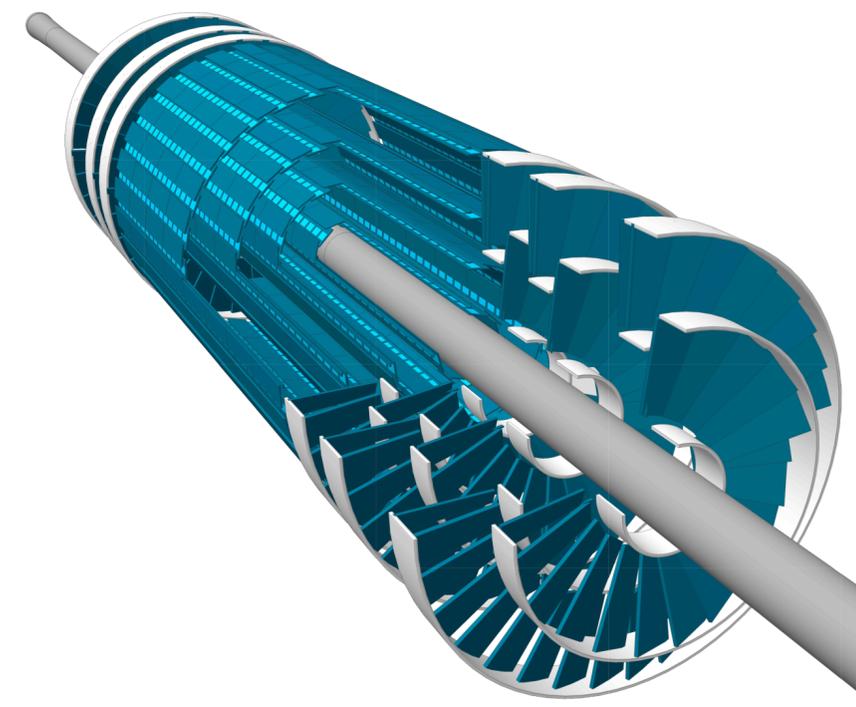
Digital readout chip (ROC). Faster readout.

**Geometry design: ensure tracking and vertex quality**

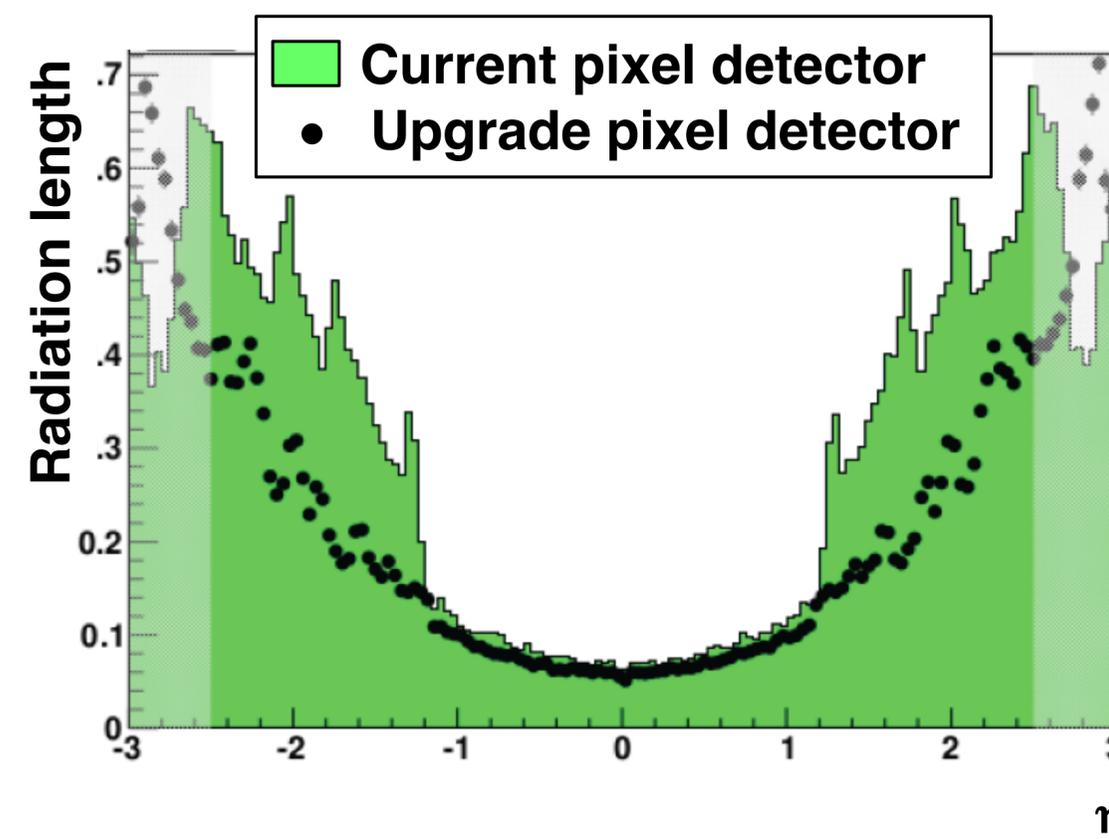
Added layers, channels doubled

**Services: reduce material budget**

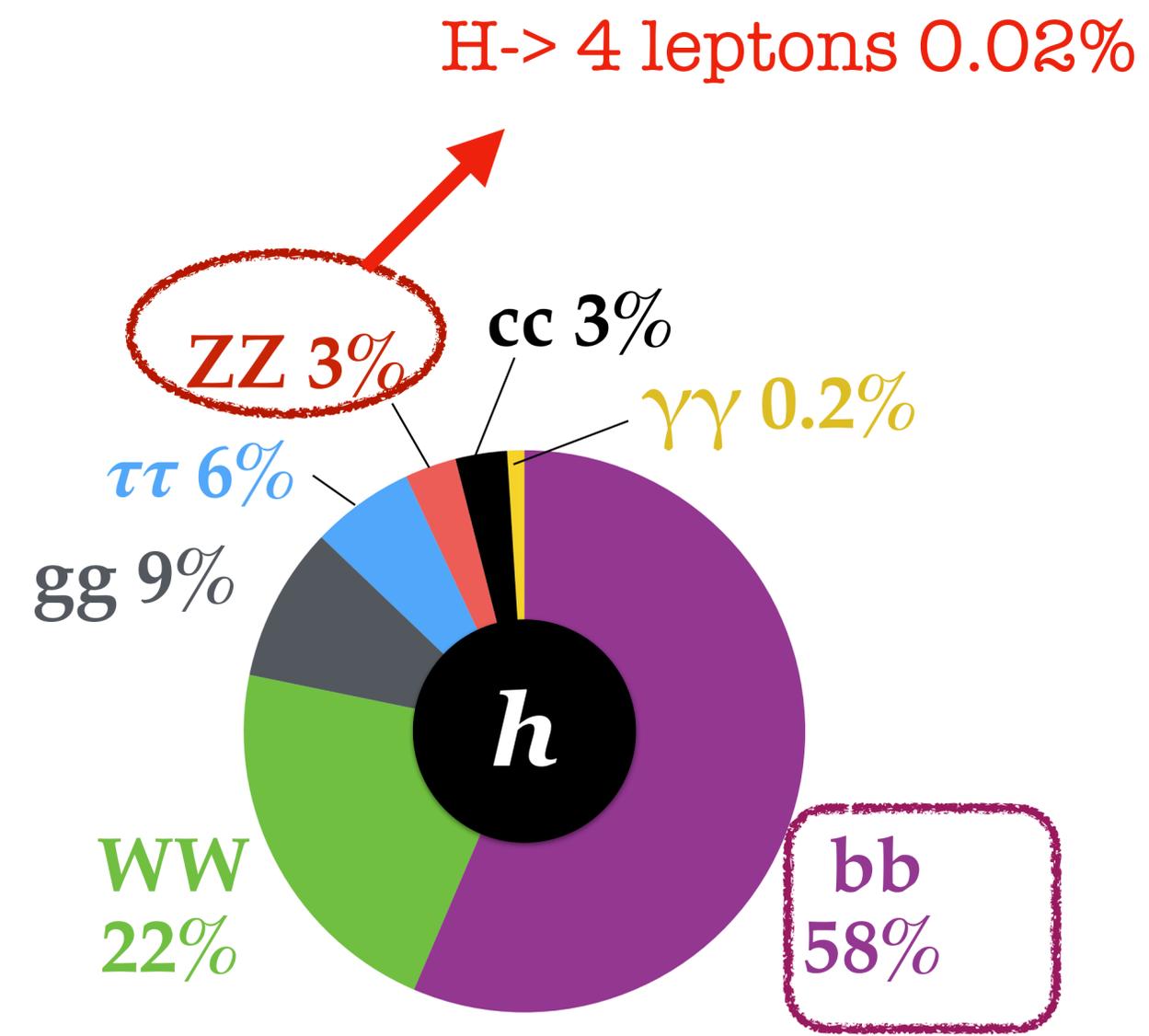
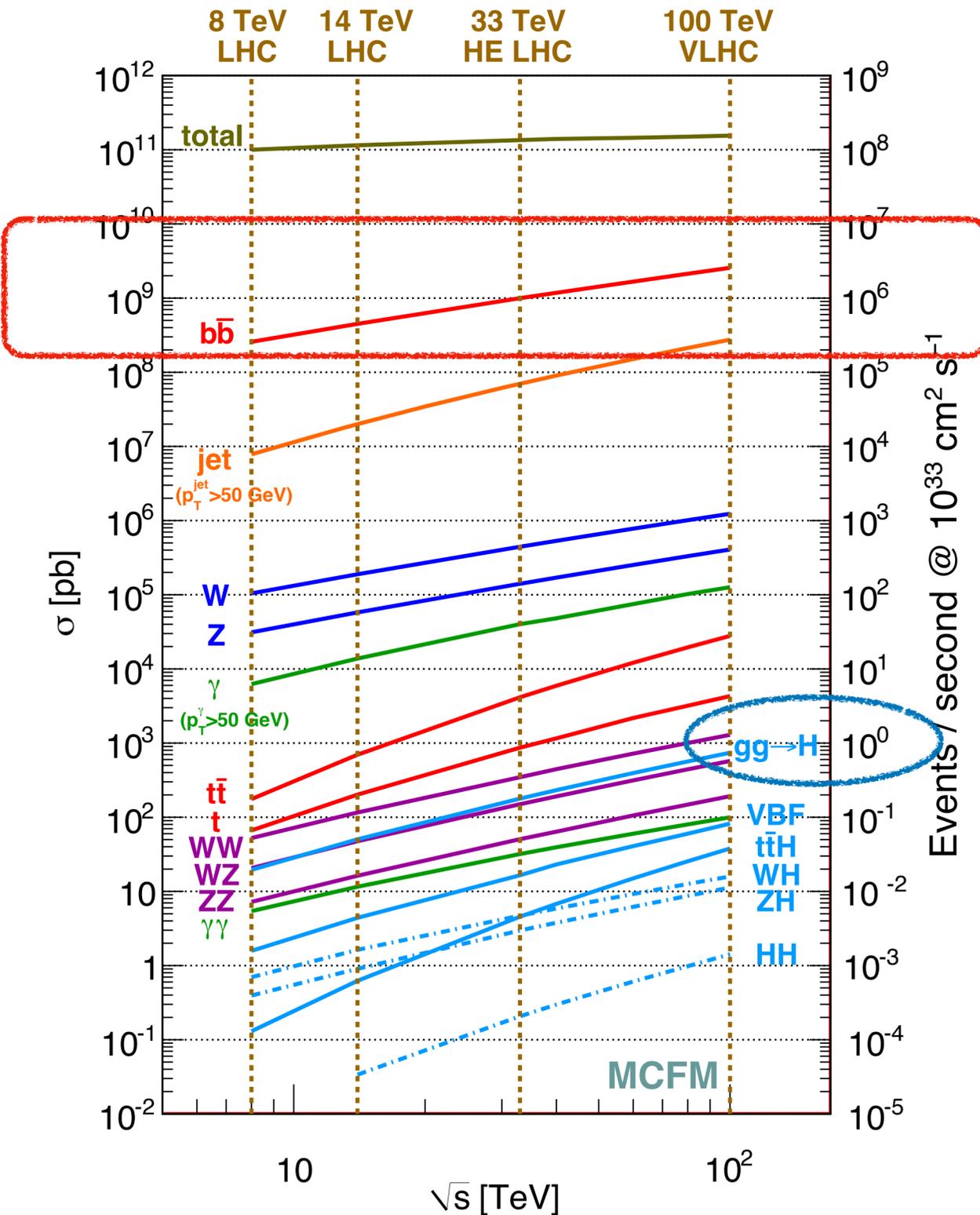
CO<sub>2</sub> cooling, DCDC powering, Service electronics out of tracker volume.



**Material budget**

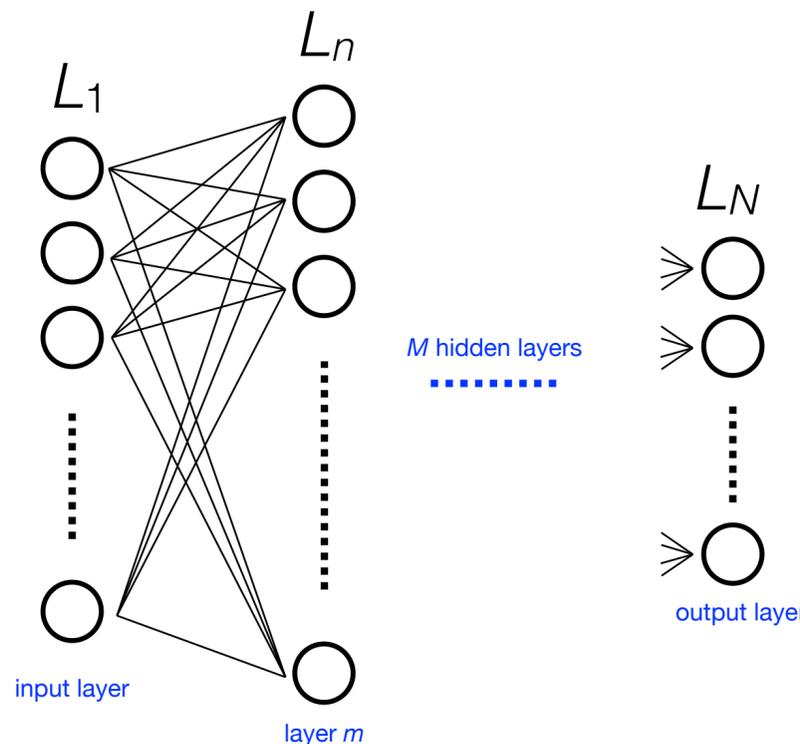


# Signal vs background

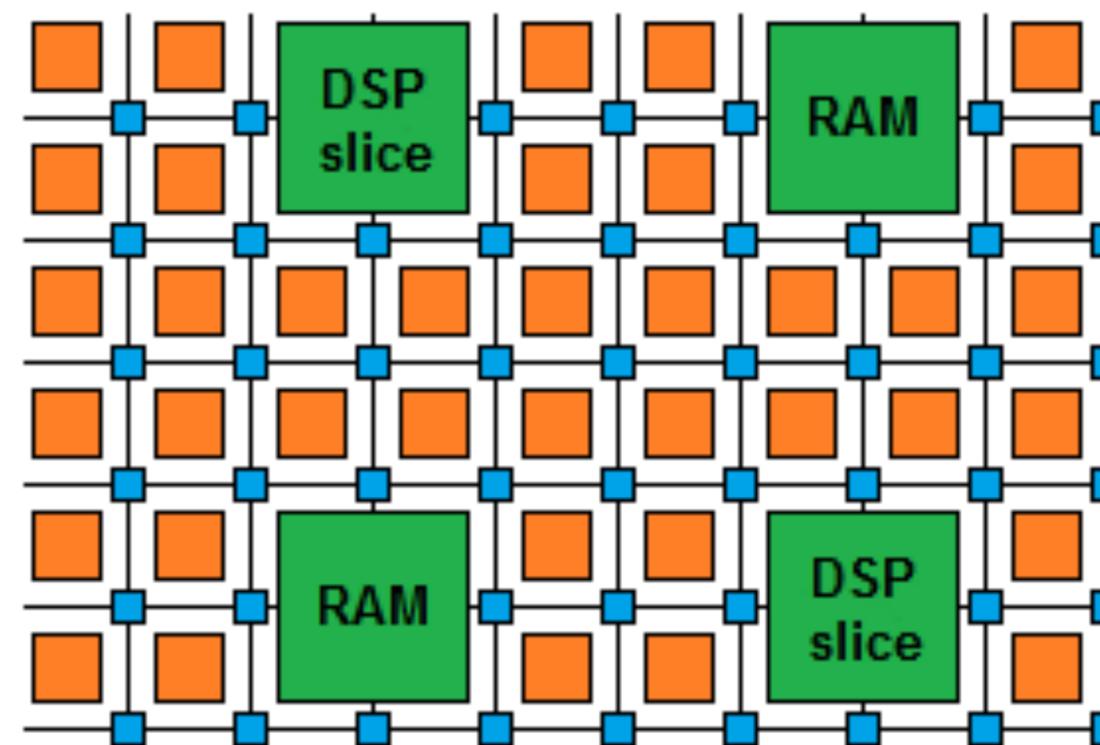


Higgs boson discovery: decay modes of lower backgrounds ( $WW/ZZ/\gamma\gamma$ ).

# NN on FPGAs



## FPGA diagram



$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

Activation functions  
Precomputed, and  
stored in BRAMs

Multiplications  
**Digital Signal Processing  
DSPs**

Addition  
Logic cells

# Natural fit for FPGAs... limited resources 47

$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

Activation functions  
Precomputed, and  
stored in BRAMs

Multiplications  
**DSPs**

Addition  
**Logic cells**

$$N_{\text{multiplications}} = \sum_{n=2}^N L_{n-1} \times L_n$$

Small network: thousands of connections

**Limitation: Number of DSPs**



**Virtex Ultrascale+ VU9P**

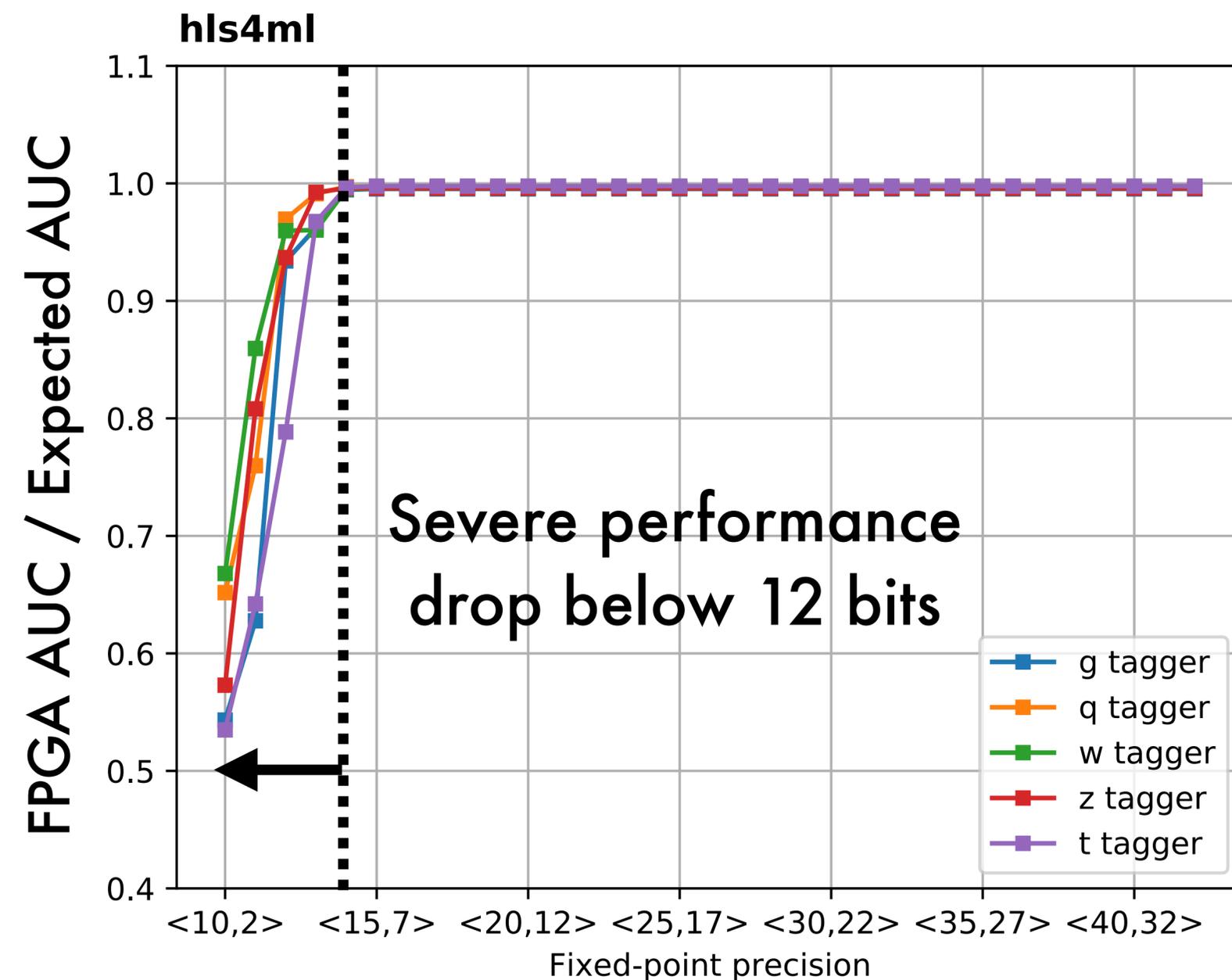
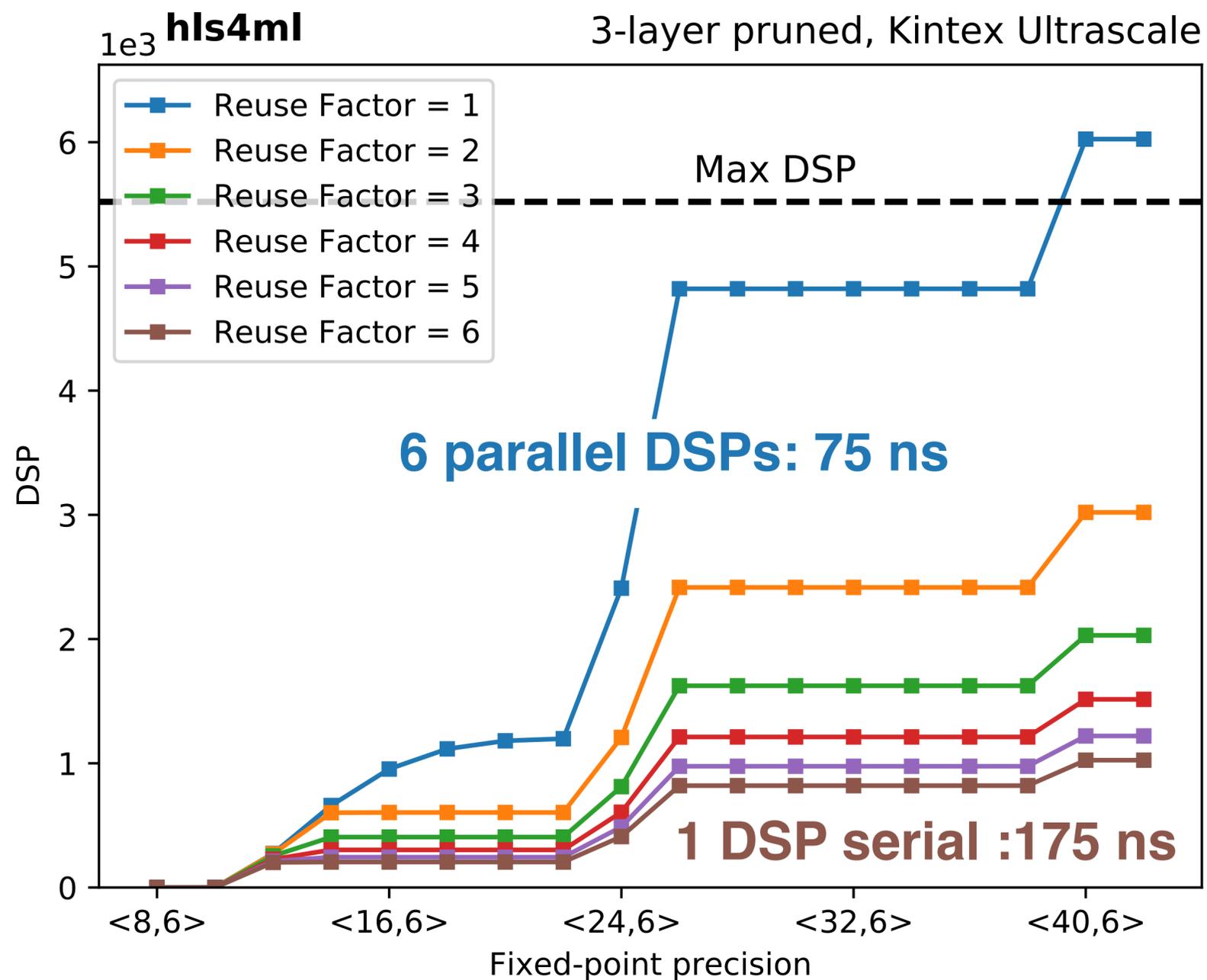
6800 DSPs

1M LUTs

2M FFs

75 Mb BRAM

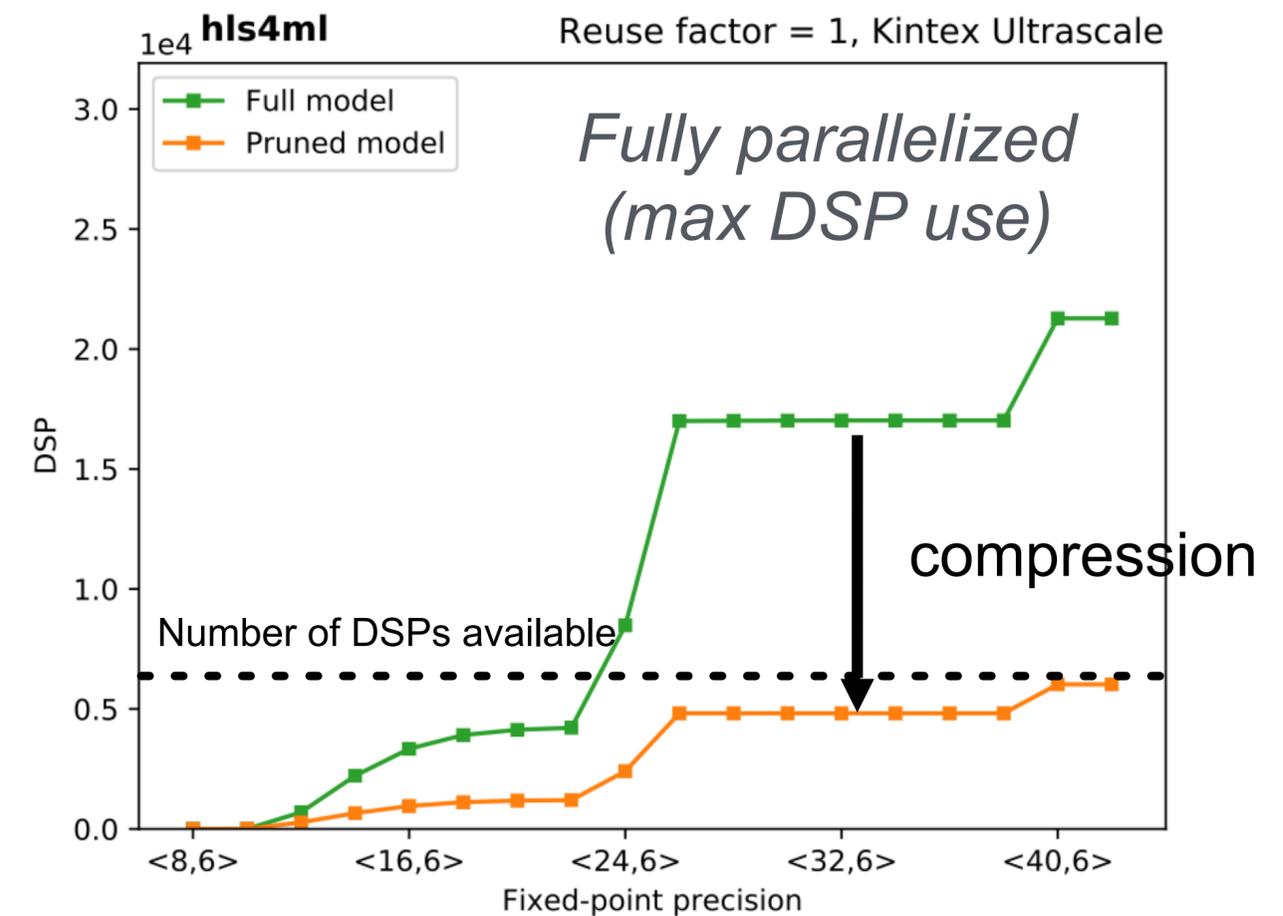
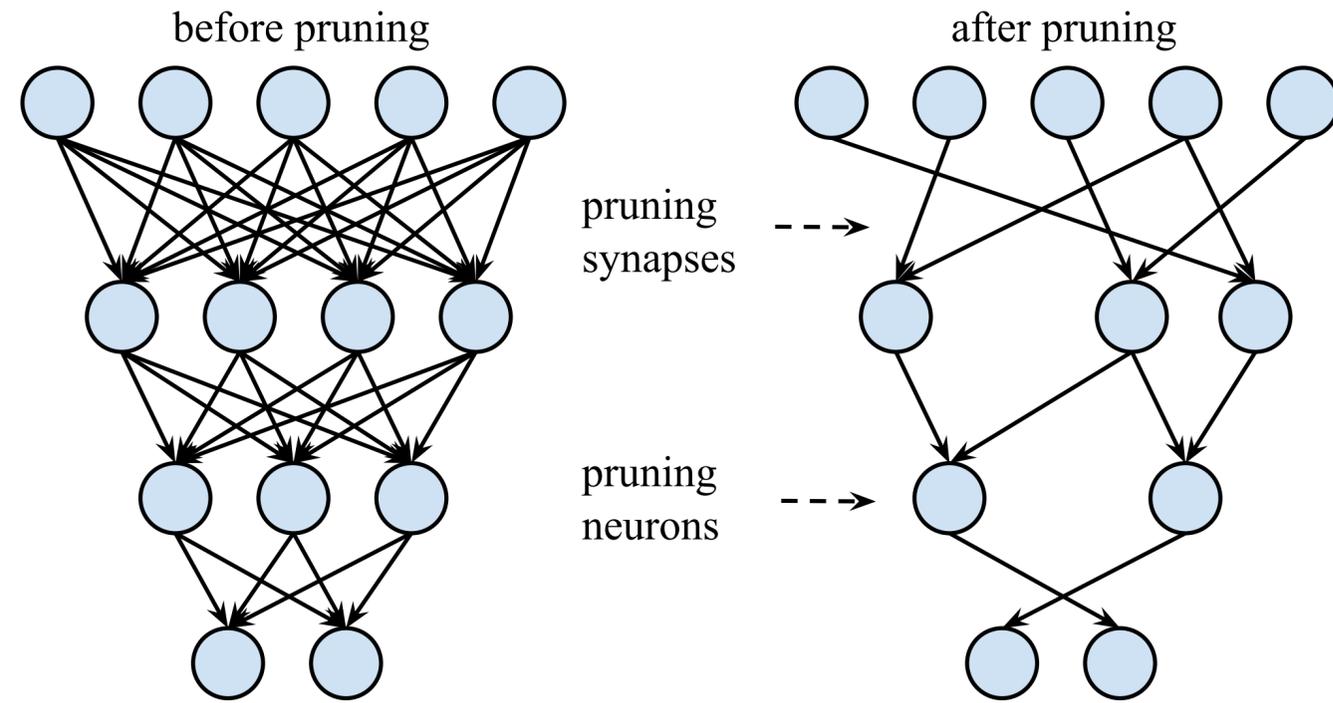
# Fit NN on FPGA: Quantization & Reuse



# Efficient NN design: compression

49

[<https://arxiv.org/abs/1804.06913>]



**Neural Network compression is a widespread technique to reduce the size, energy consumption, and overtraining of deep neural networks**

Several approaches in literature

arxiv.1510.00149, arxiv.1712.01312, arxiv.1405.3866, arxiv.1602.07576,  
doi:10.1145/1150402.1150464

**First paper demonstrated a fully connected NN in 100 ns.**

**HLS4ML in CMS**

Run 3: muon momentum regression in CMS

More models demonstrated for Phase-2 trigger upgrade TDR

**Advanced models:**

binary/ternary, CNNs, RNNs, auto-encoders. Support for Graph Neural Network Models

**Advanced Pruning/quantization:**

Quantization-aware training with QKeras/Quantization-aware pruning

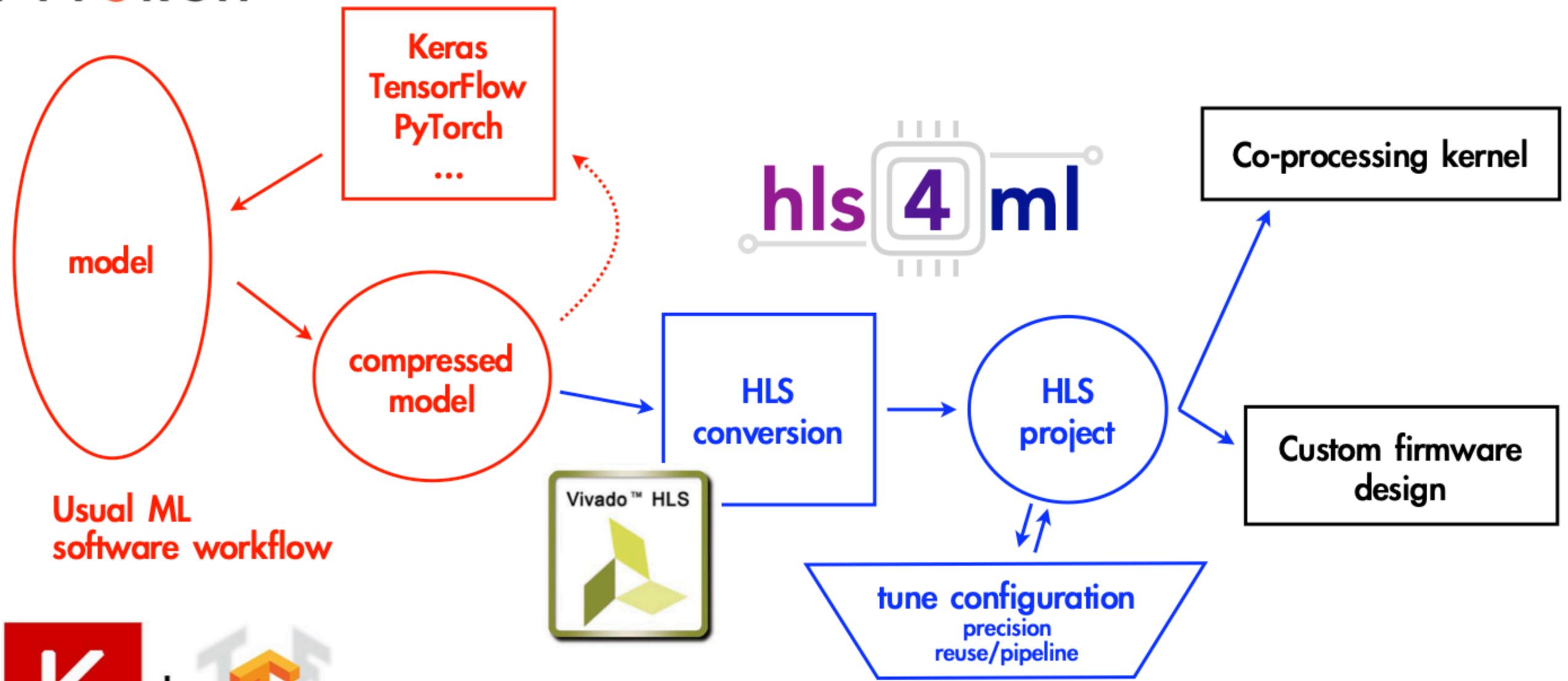
**On ASICs and Low power devices.**

For latest status: please check [hls4ml website](#), [CPAD 2021 talk](#),

Try it out: [hls4ml tutorials](#)

# High Level Synthesis 4 Machine Learning

PYTORCH



Usual ML software workflow



<https://fastmachinelearning.org/hls4ml/>