# Fast inference of Deep Learning Applications with FPGAs

Maurizio Pierini

CERN

hls 4 ml

erc European Research Council

# hls4ml

Jennifer Ngadiuba Vladimir Loncar, Maurizio Pierini

Javier Duarte, Burt Holzman, Sergo Jindariani, Ben Kreis, Mia Liu, Kevin Pedro, Ryan Rivera, Nhan Tran, Aris Tsaris

Edward Kreinar

Sioni Summer

Song Han, Phil Harris, Dylan Rankin

Zhenbin Wu

**https://github.com/hls-fpga-machine-learning/hls4ml**

**https://github.com/hls-fpga-machine-learning/SonicCMS**

**https://hls-fpga-machine-learning.github.io/hls4ml/**

**https://arxiv.org/pdf/1804.06913.pdf**

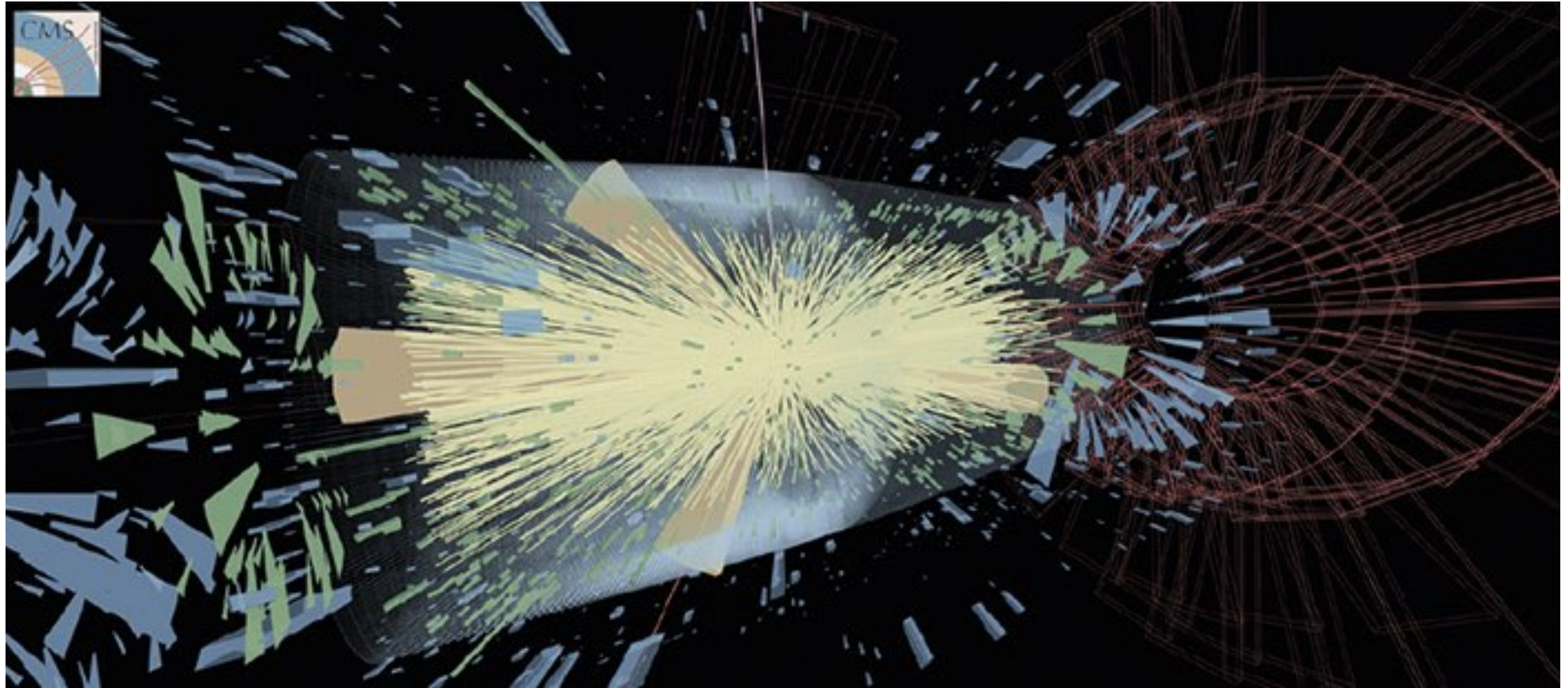**Recent talks at conferences:**

**CERN Data Science Seminar**
**FNAL Seminar**
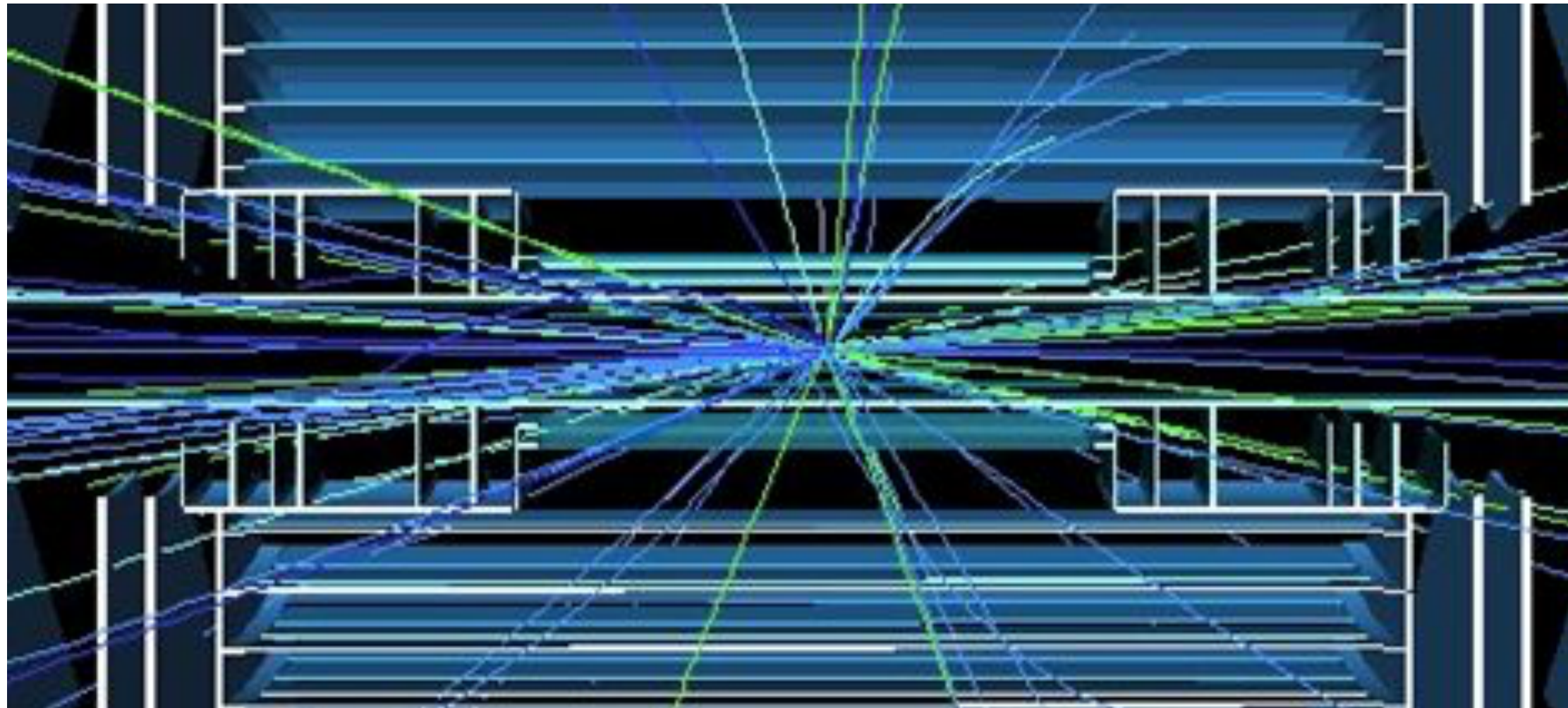**Connecting The Dots 2018**
**TWEPP 2018**

# Outline

- *The problem: High-Luminosity LHC*

- *Why Deep Learning: a few application examples*

- *Deploying Deep NNs online*

  - *HLT accelerated inference*

  - *L1 NN on FPGAs with HLS*

- *Conclusions*

# ML and HEP future challenges

# HL-LHC: elephant in the room

5 interactions/beam cross

400 interactions/beam cross



| 6 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 | 2025 | | 2035 |

This is when the R&D has to happen

LHC  **Today**

HL_LHC
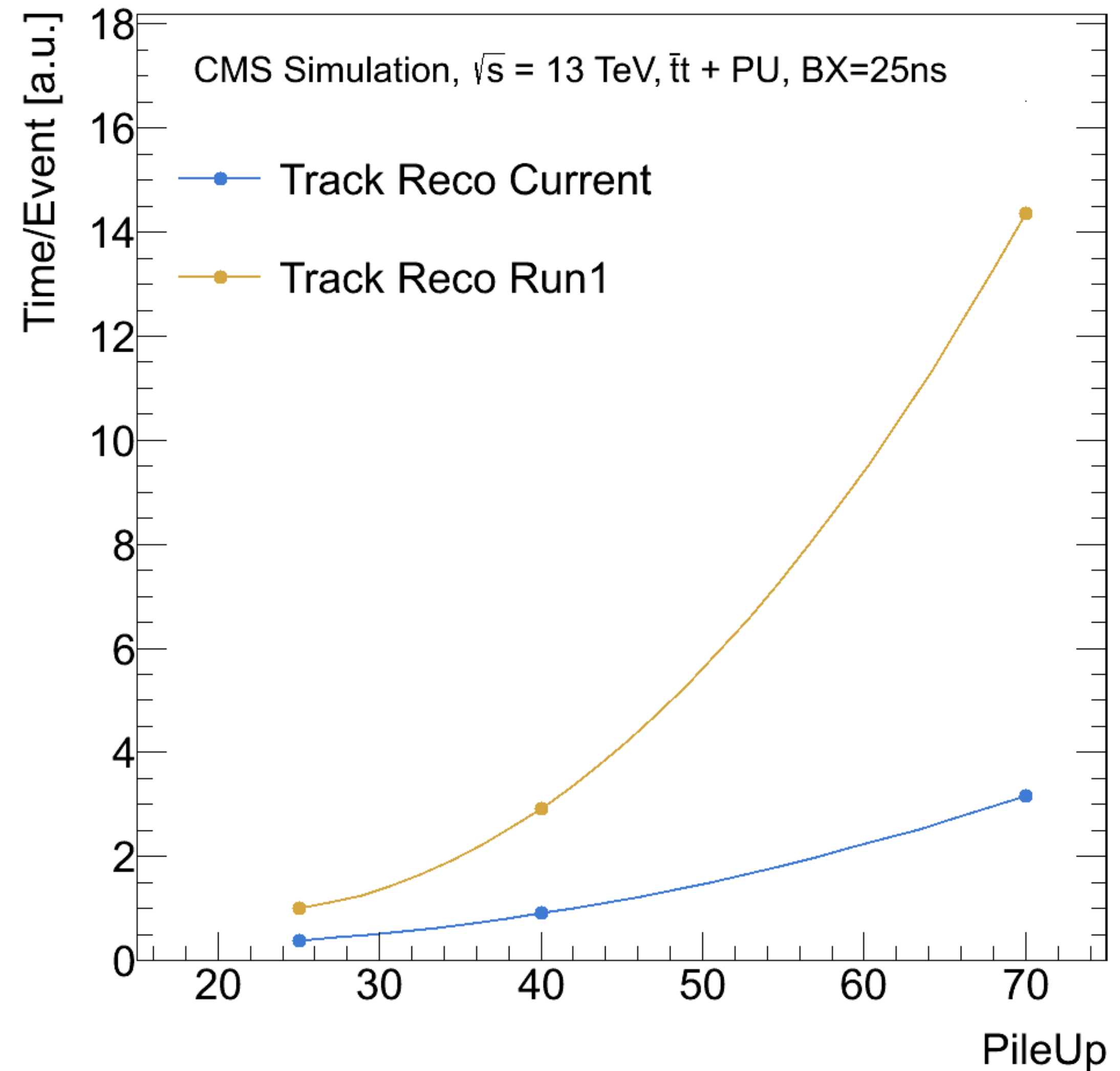
- ~40 collisions/event
- ~10 sec/event processing time
- (at best)Same computing resources as today

- ~200 collisions/event
- ~minute/event processing time
- (at best)Same computing resources as today

# HL-LHC: elephant in the room

- *Flat budget vs. more needs = **current rule-based reconstruction algorithms will not be sustainable***

- ***Adopted solution:** more granular and complex detectors → more computing resources needed → more problems*

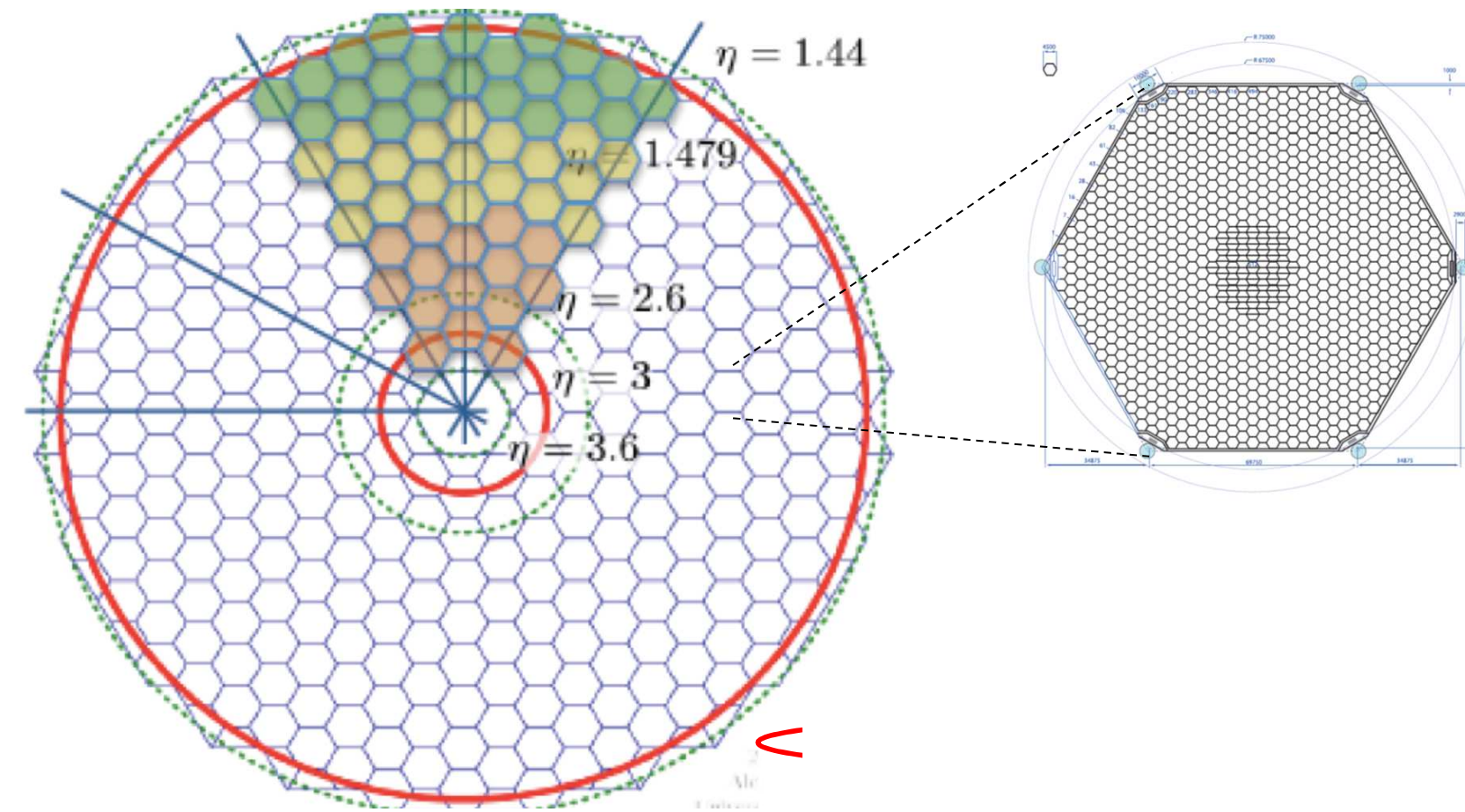- ***Modern Machine Learning might be the way out***

CMS Simulation, $\sqrt{s}$ = 13 TeV, $\bar{t}t$ + PU, BX=25ns

Track Reco Current

Track Reco Run1

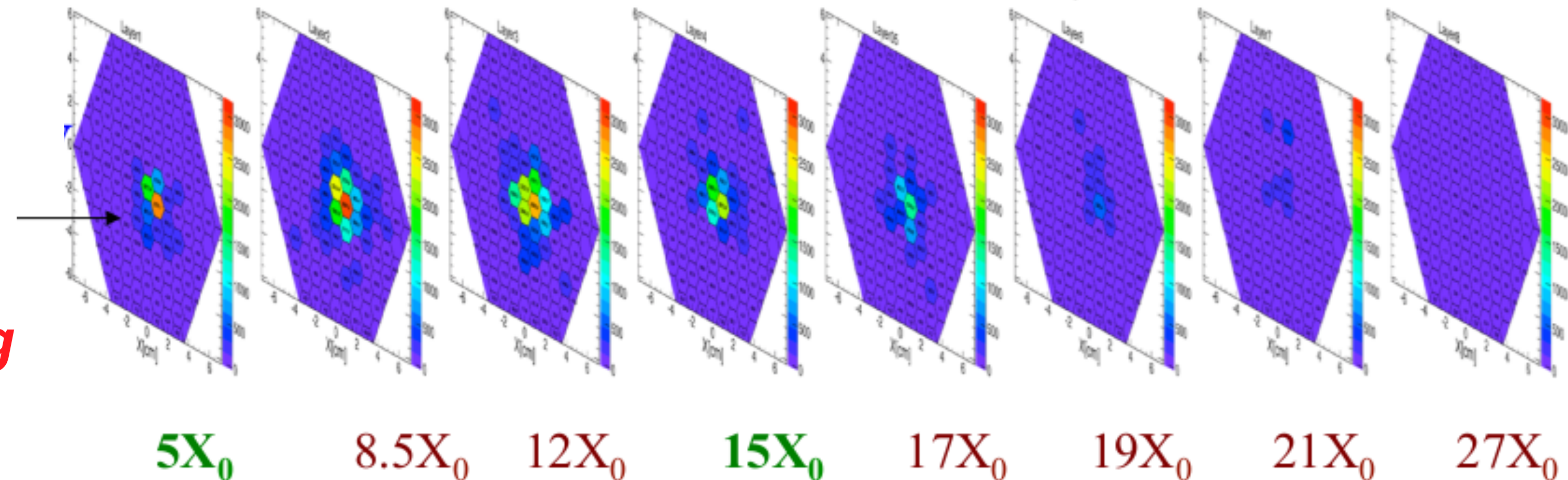Time/Event [a.u.]

PileUp

# HL-LHC: elephant in the room

- *Flat budget vs. more needs = **current rule-based reconstruction algorithms will not be sustainable***

- ***Adopted solution:** more granular and complex detectors → more computing resources needed → more problems*
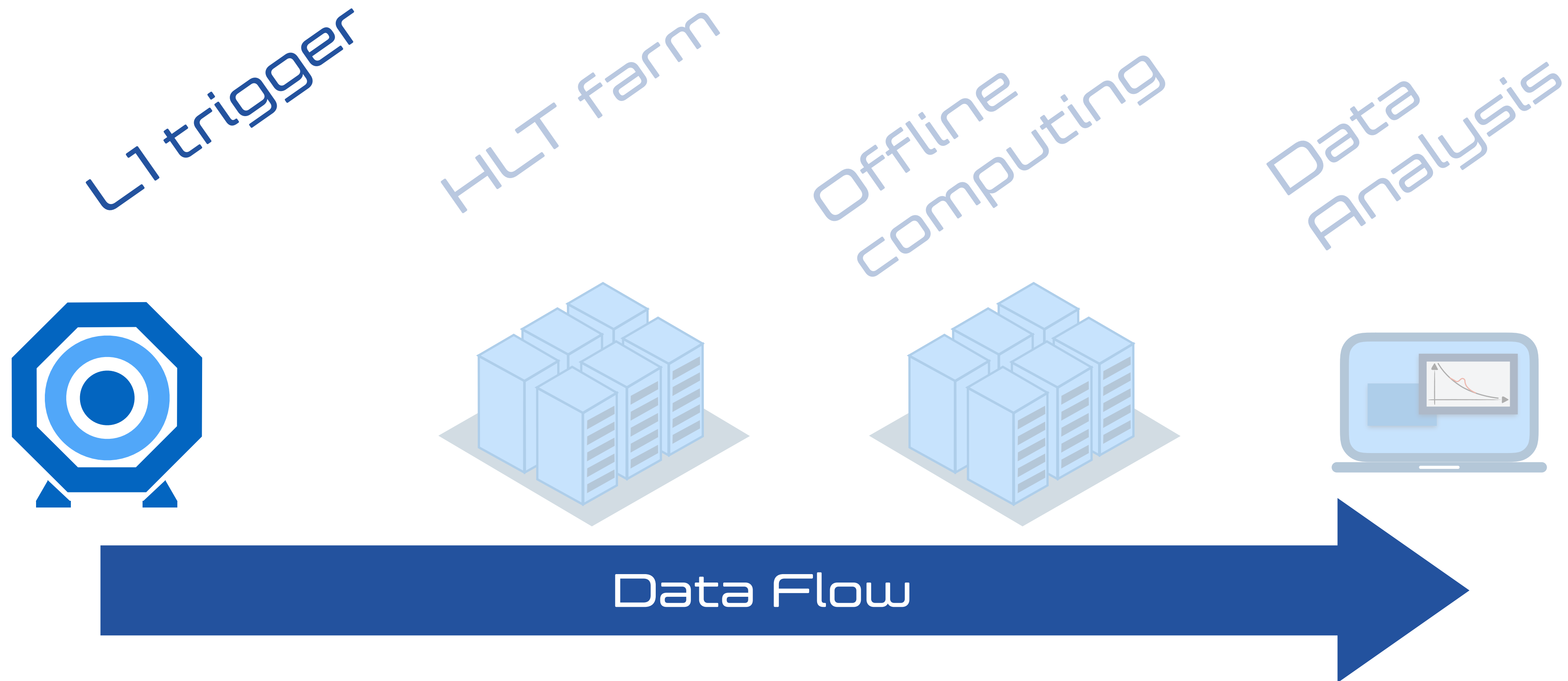
- ***Modern Machine Learning might be the way out***

$\eta = 1.44$

$\eta = 1.479$

$\eta = 2.6$

$\eta = 3$

$\eta = 3.6$

250 GeV electron passing through 8 layers (27 X$_0$)

$5X_0$   $8.5X_0$   $12X_0$   $15X_0$   $17X_0$   $19X_0$   $21X_0$   $27X_0$

# The LHC Big Data problem

L1 trigger

HLT farm

Offline computing

Data Analysis

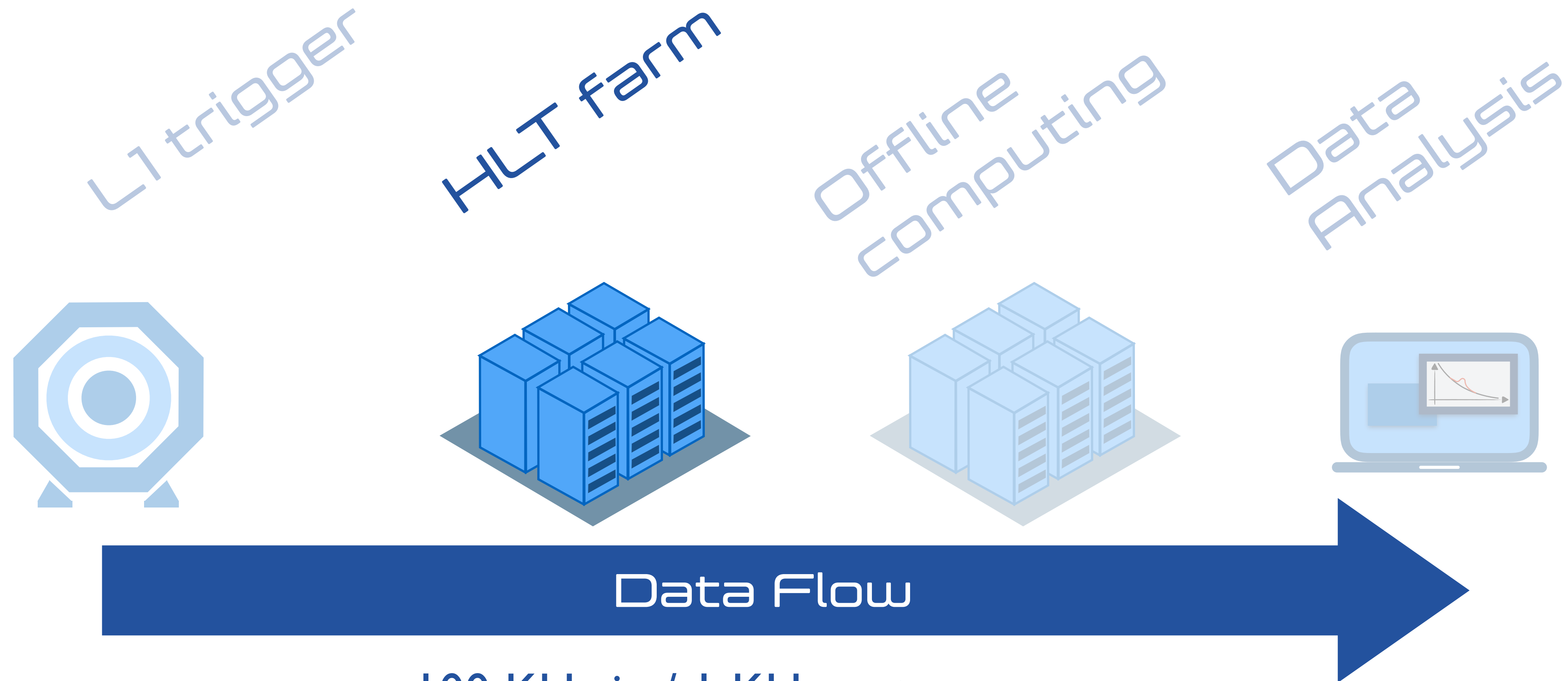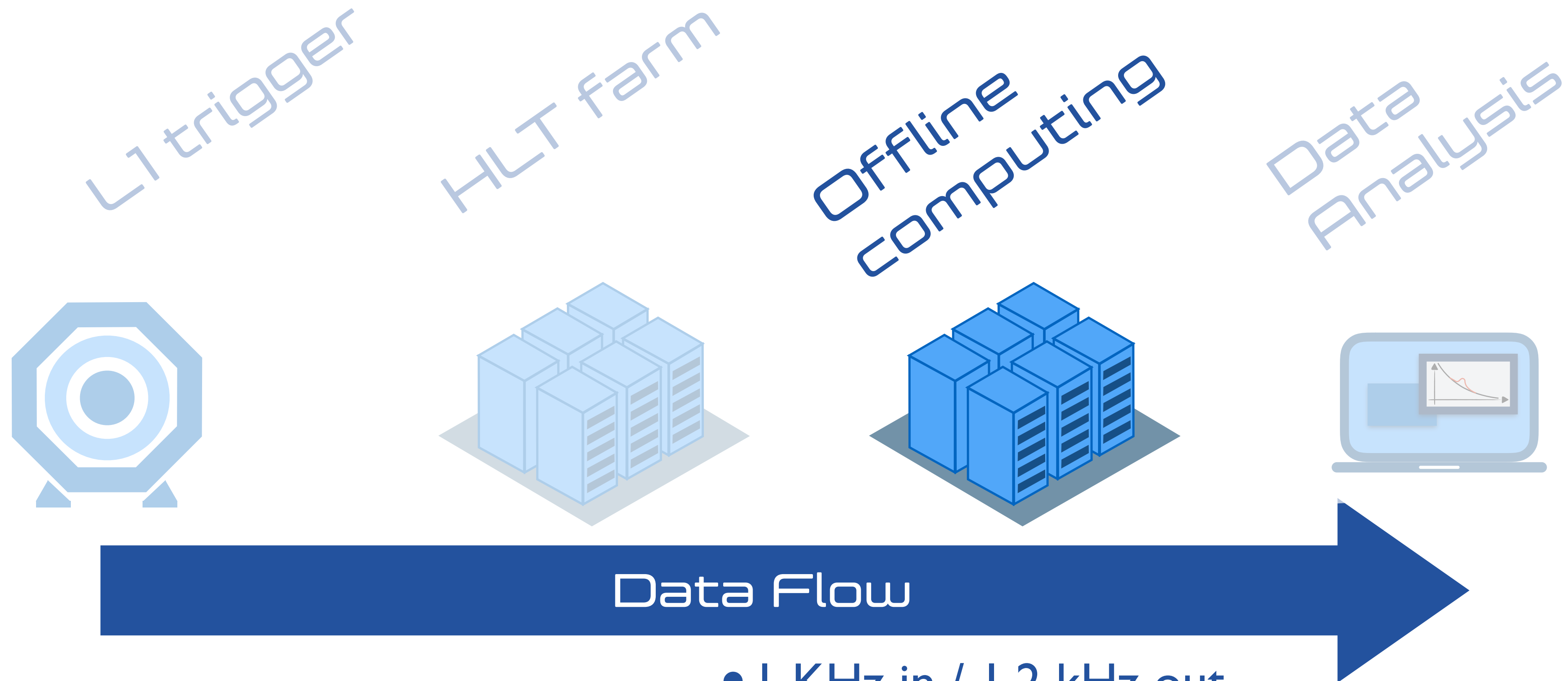**Data Flow**

- 40 MHz in / 100 KHz out
- ~ 500 KB / event
- Processing time: ~10 μs
- Based on coarse local reconstructions
- FPGAs / Hardware implemented

# The LHC Big Data problem

L1 trigger

HLT farm

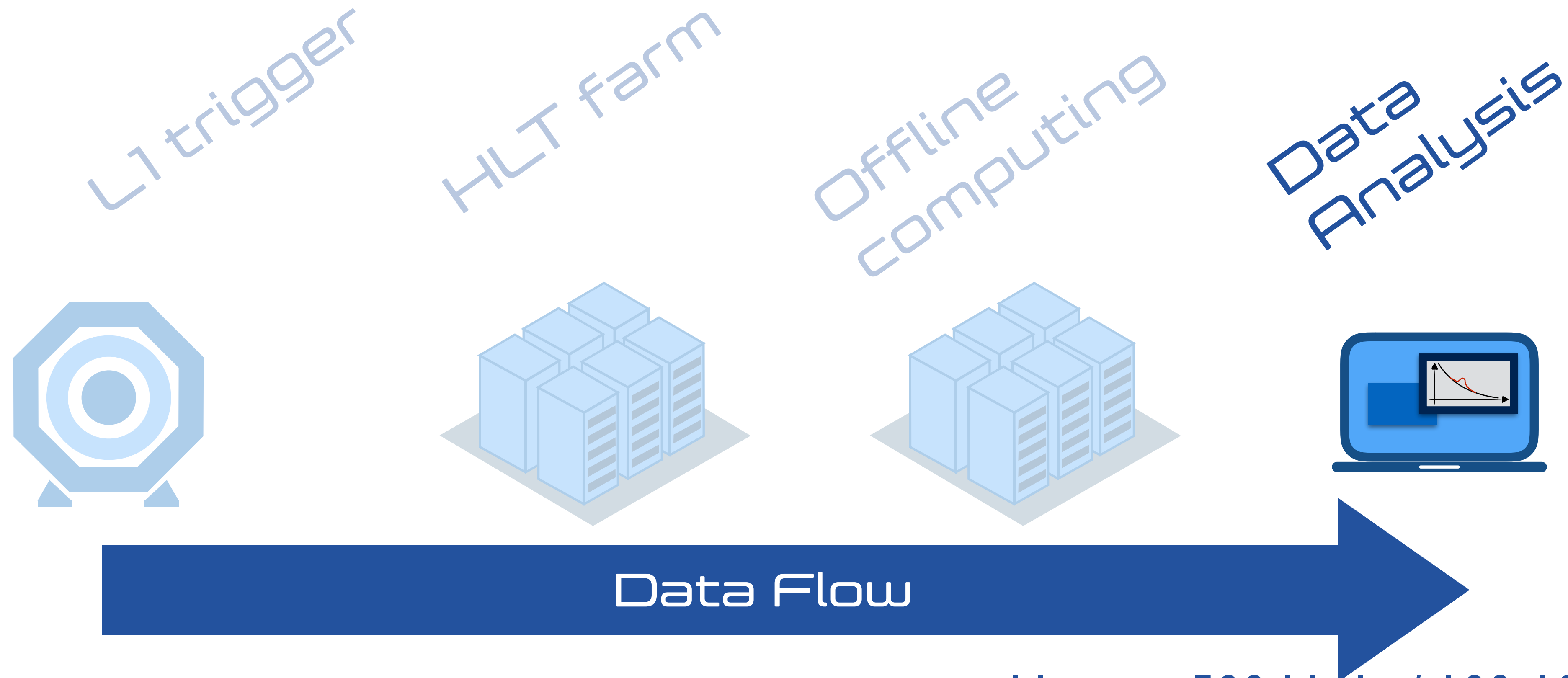Offline computing

Data Analysis

**Data Flow**

- 100 KHz in / 1 KHz out

- ~ 500 KB / event

- Processing time: ~30 ms

- Based on simplified global reconstructions

- Software implemented on CPUs

# The LHC Big Data problem

L1 trigger

HLT farm

Offline computing

Data Analysis

**Data Flow**

- 1 KHz in / 1.2 kHz out
- ~ 1 MB / 200 kB / 30 kB per event
- Processing time: ~20 s
- Based on accurate global reconstructions
- Software implemented on CPUs

# The LHC Big Data problem

L1 trigger

HLT farm

Offline computing

Data Analysis

**Data Flow**

- Up to ~ 500 Hz In / 100-1000 events out
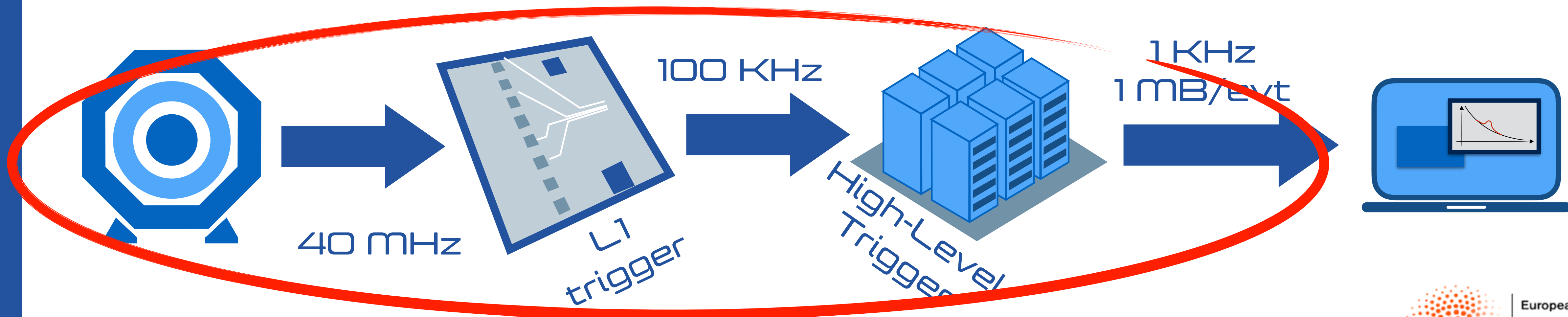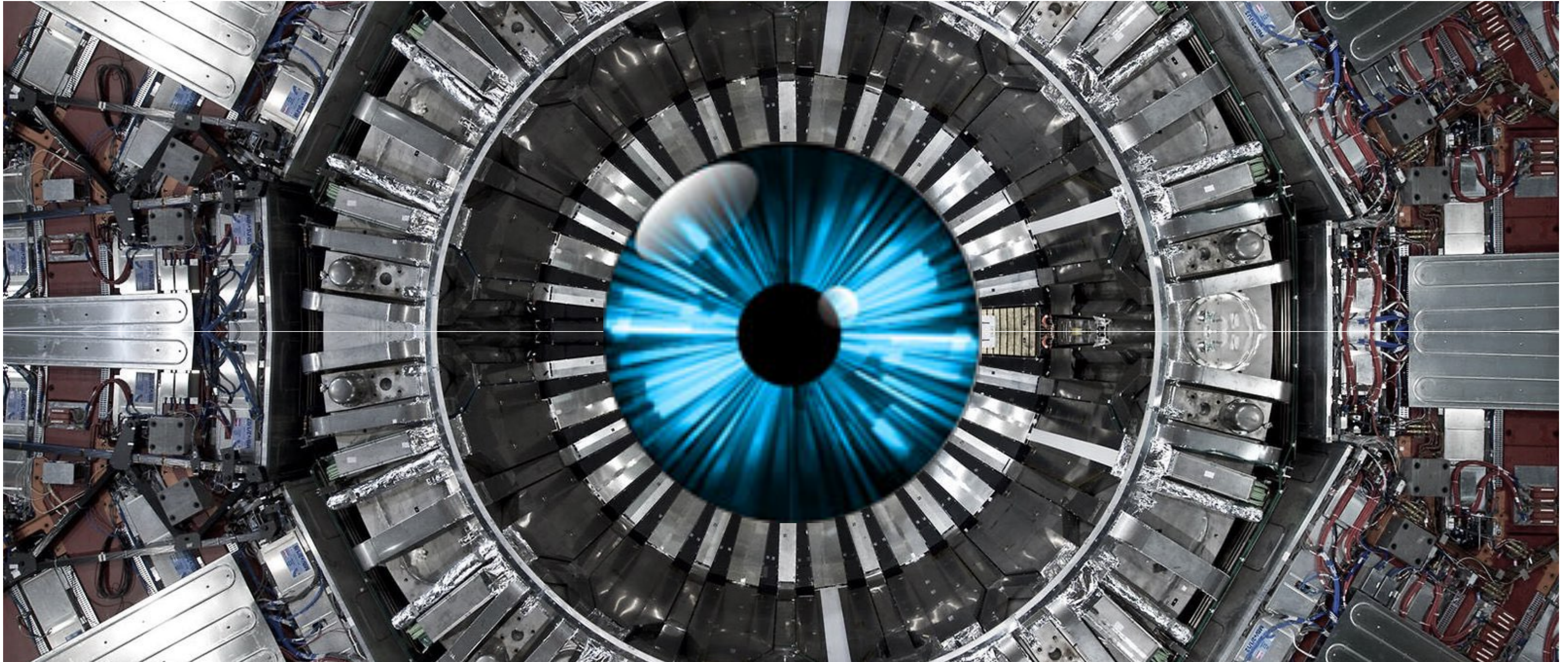- <30 KB per event
- Processing time irrelevant
- User-written code + centrally produced selection algorithms

# Deep Learning and LHC Big Data

◉ *Possible solution to the HL-LHC problem: modern Machine Learning to be <u>faster</u> and <u>better</u> in what we do today, freeing resources for new ideas*

◉ *This ML deployment need to happen in between collisions and data analysis (trigger, reconstruction, …), where freeing resources will make a difference*



100 KHz

1 KHz
1 MB/evt

40 MHz

L1 trigger

High-Level Trigger
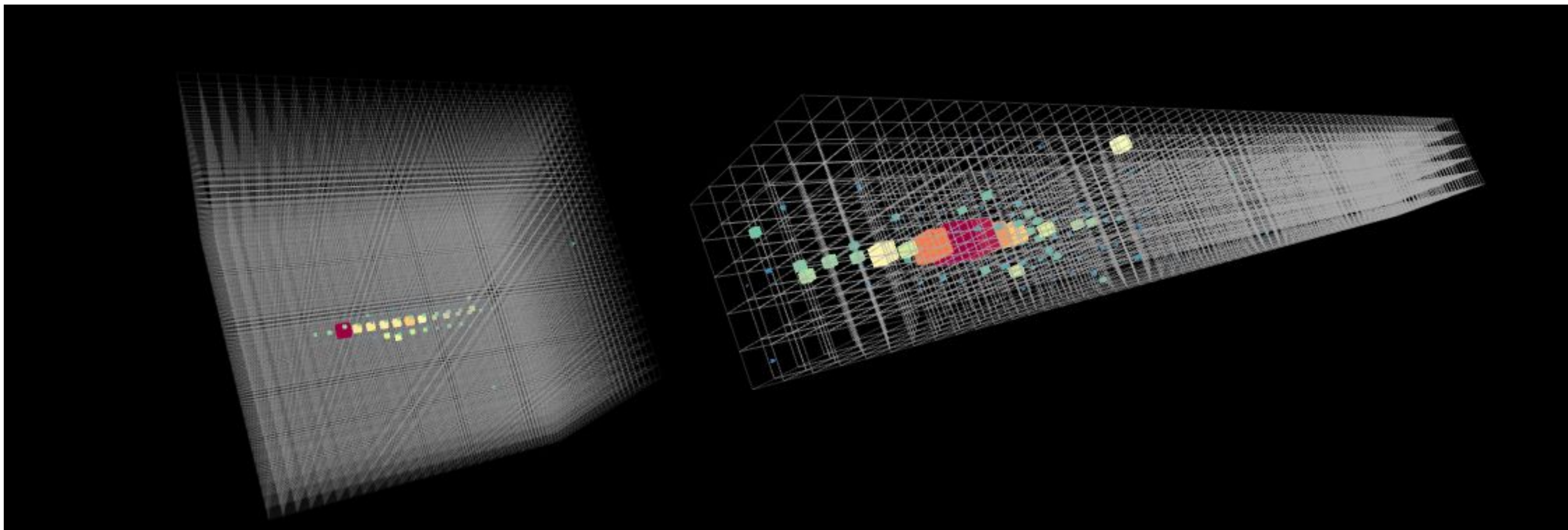
European Research Council
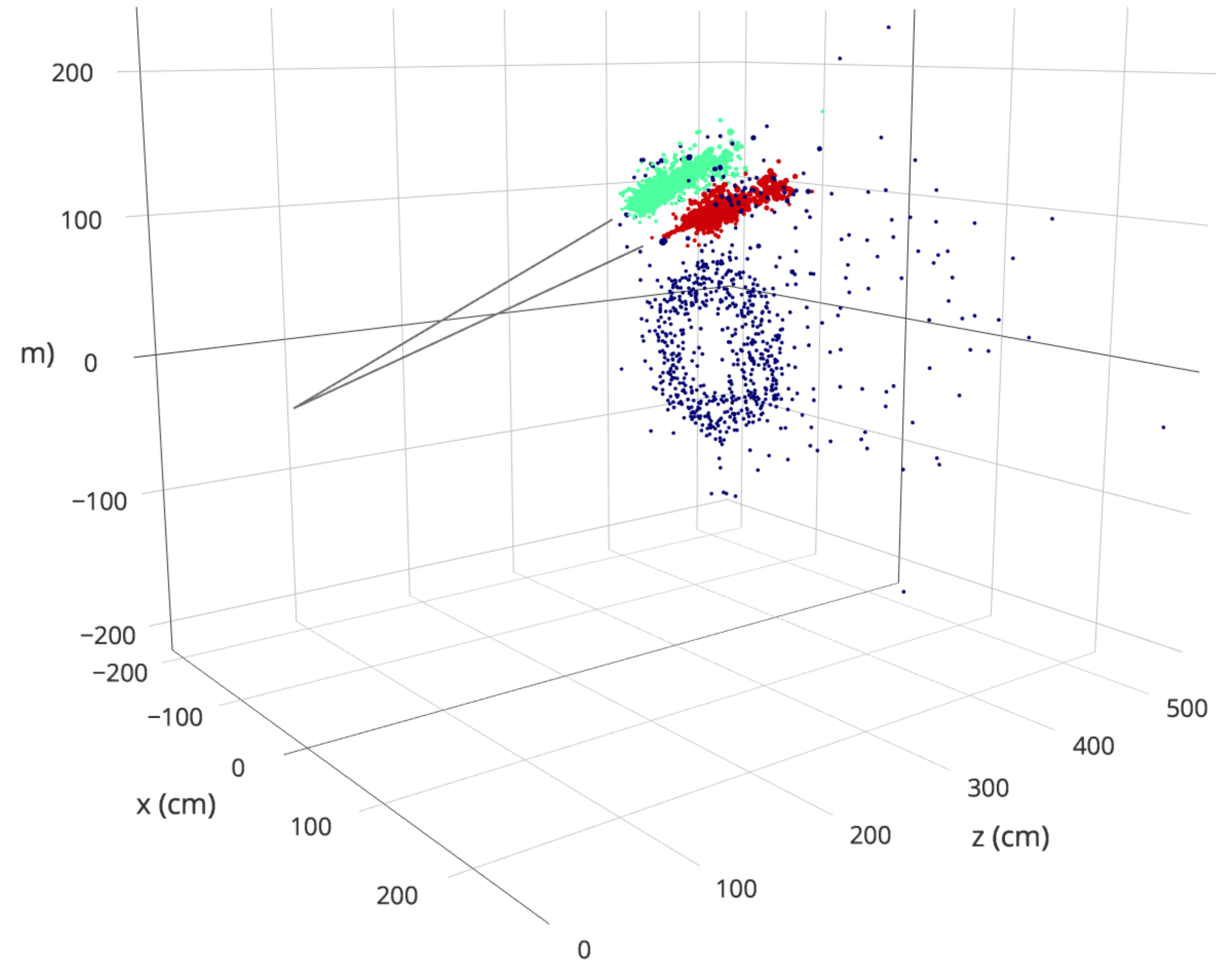
# Faster Particle Reconstruction With Computer Vision

# Calorimetry & Computer Vision

◉ *(next generation) digital calorimeters: 3D arrays of sensors with more regular geometry*

◉ *Ideal configuration to apply Convolutional Neural Network*

   ◉ *speed up reconstruction at similar performances*

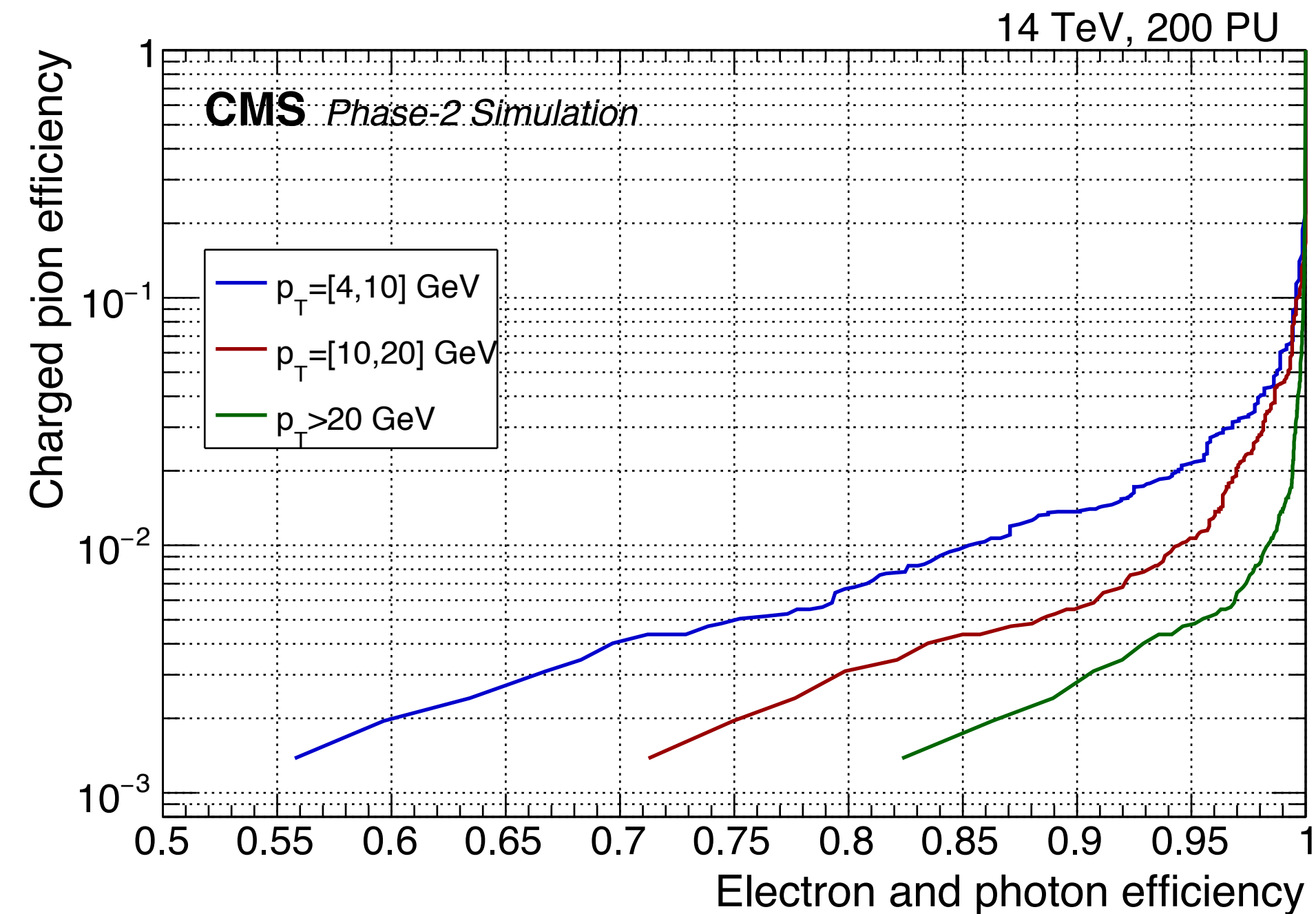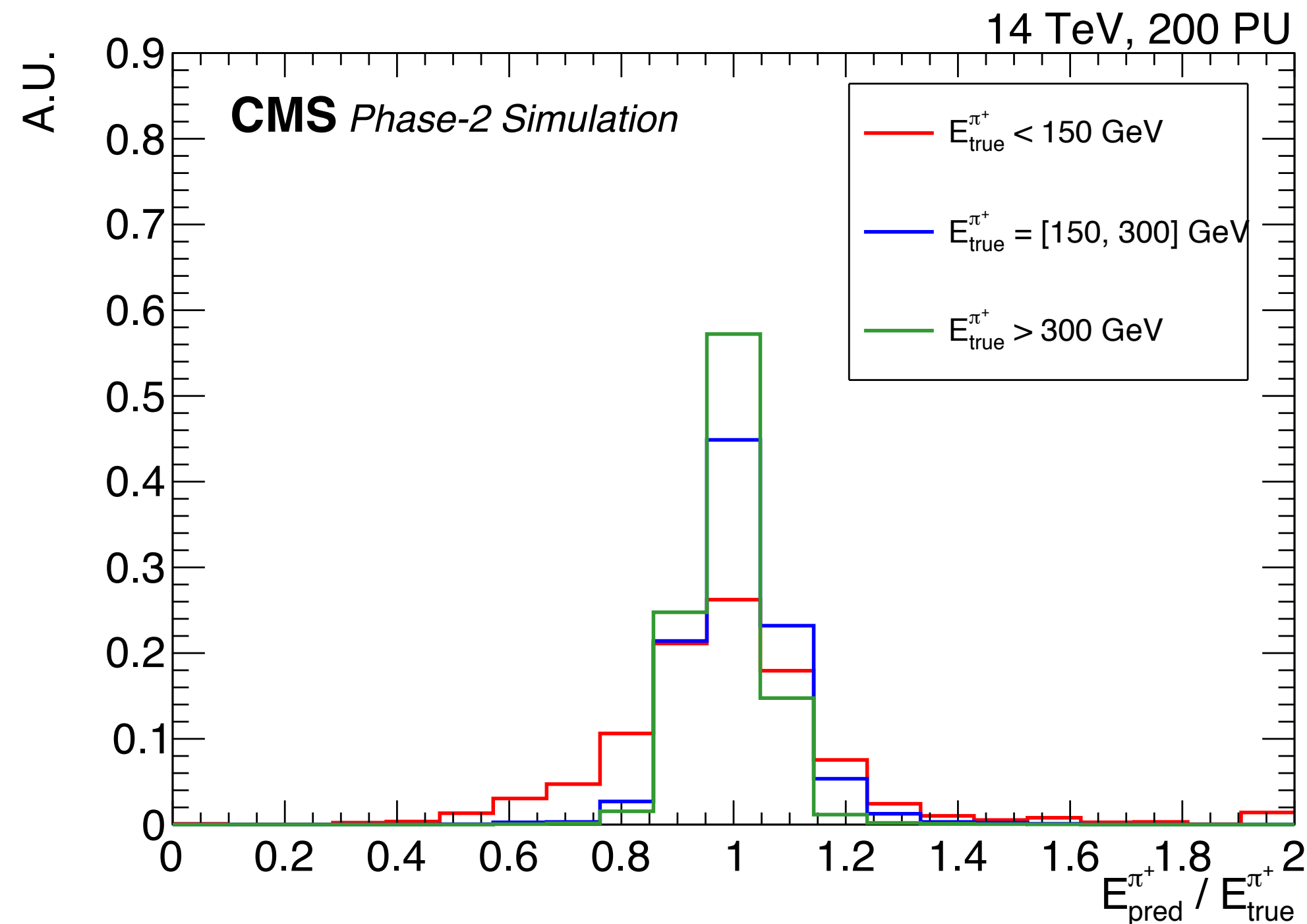   ◉ *and possibly improve performances*

# HGCAL: Why Deep Learning

- *High granularity to distinguish individual particles even with many simultaneous collisions*

- *Standard algorithms slowed down by combinatorial*

- *3D Convolutional Neural Networks much faster in going from raw data to answer*

- *Need to develop models to guarantee same performances, possibly better*
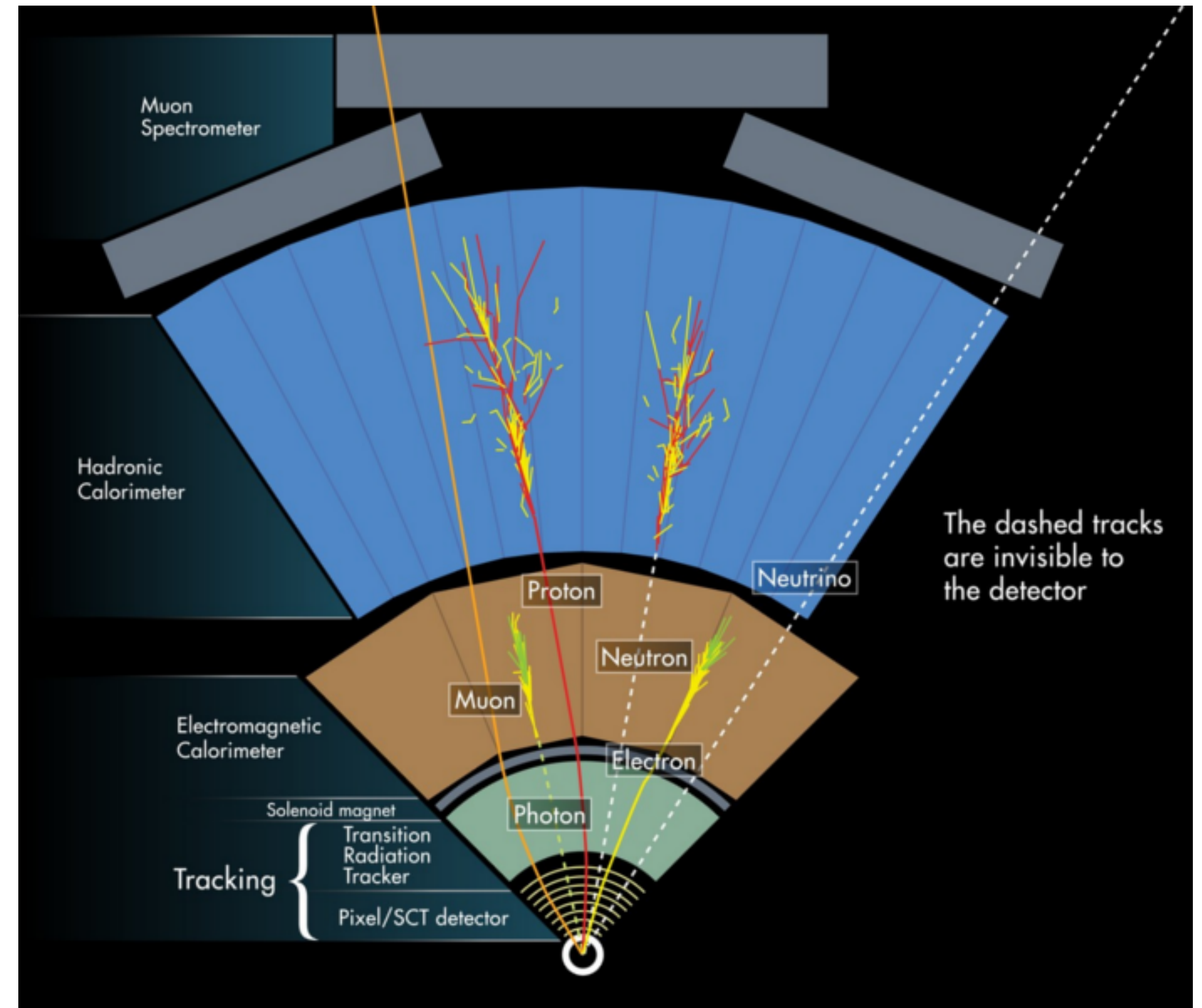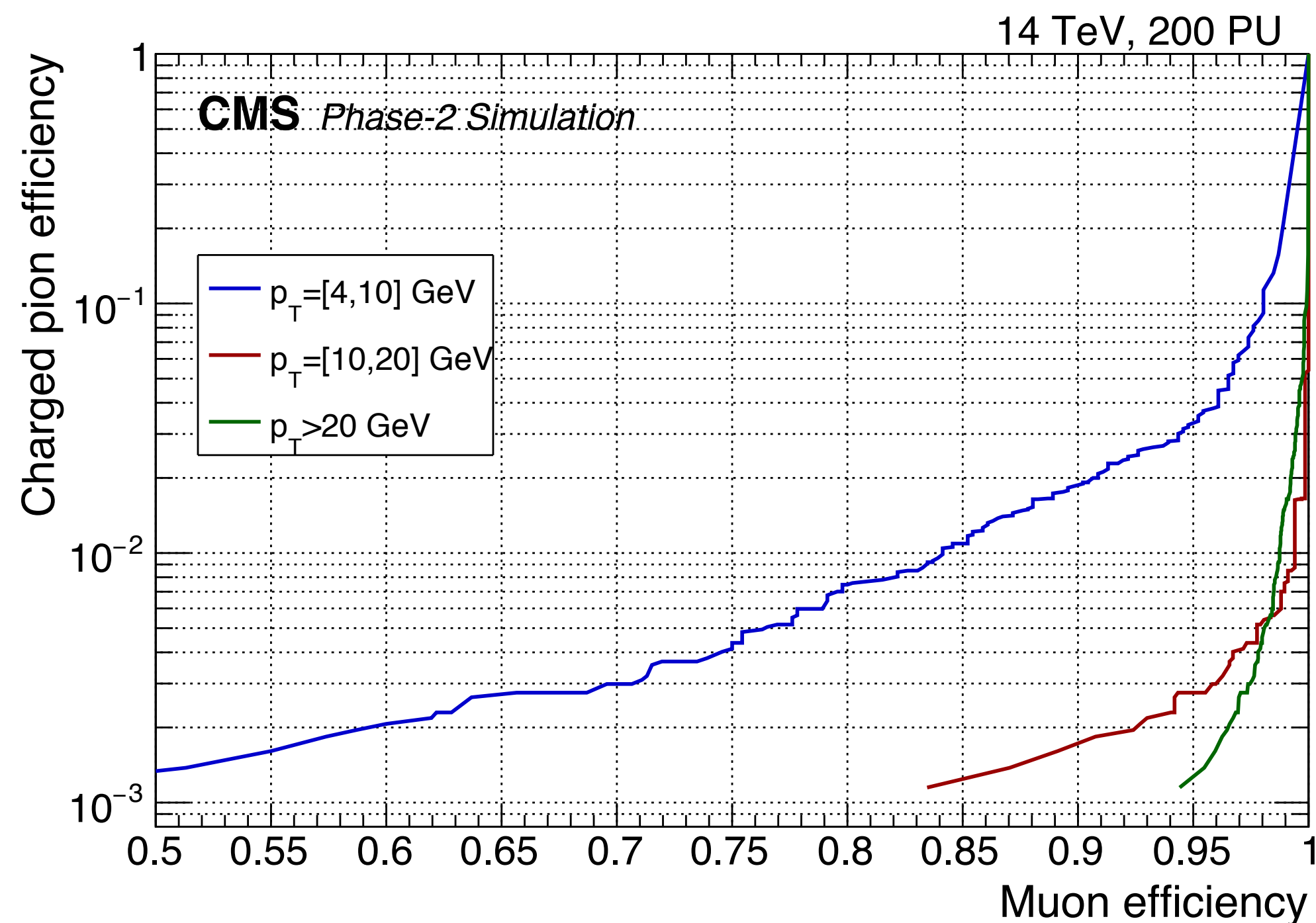


16

# HGCAL: DL reconstruction



- ◉ *State-of-the-art performances in terms of particle identification & energy measurement*

- ◉ *Sizeable speed-up at reconstruction time*

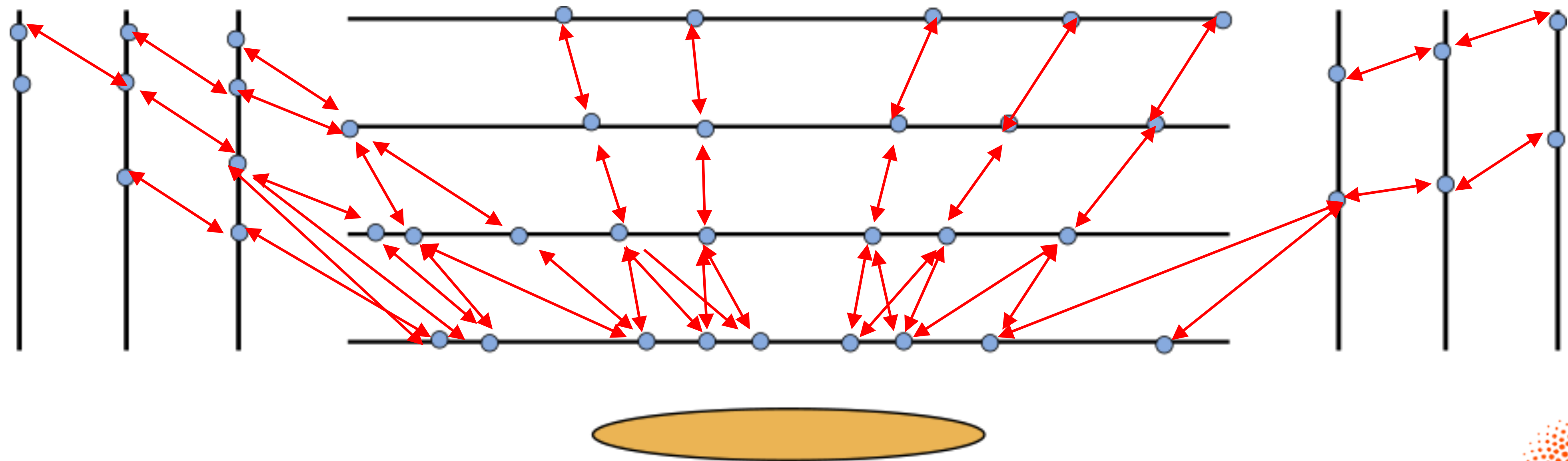- ◉ *Can get even better performances with model optimization*

# HGCAL: Opportunities

- New hardware + new techniques = new opportunities & paradigm breaking

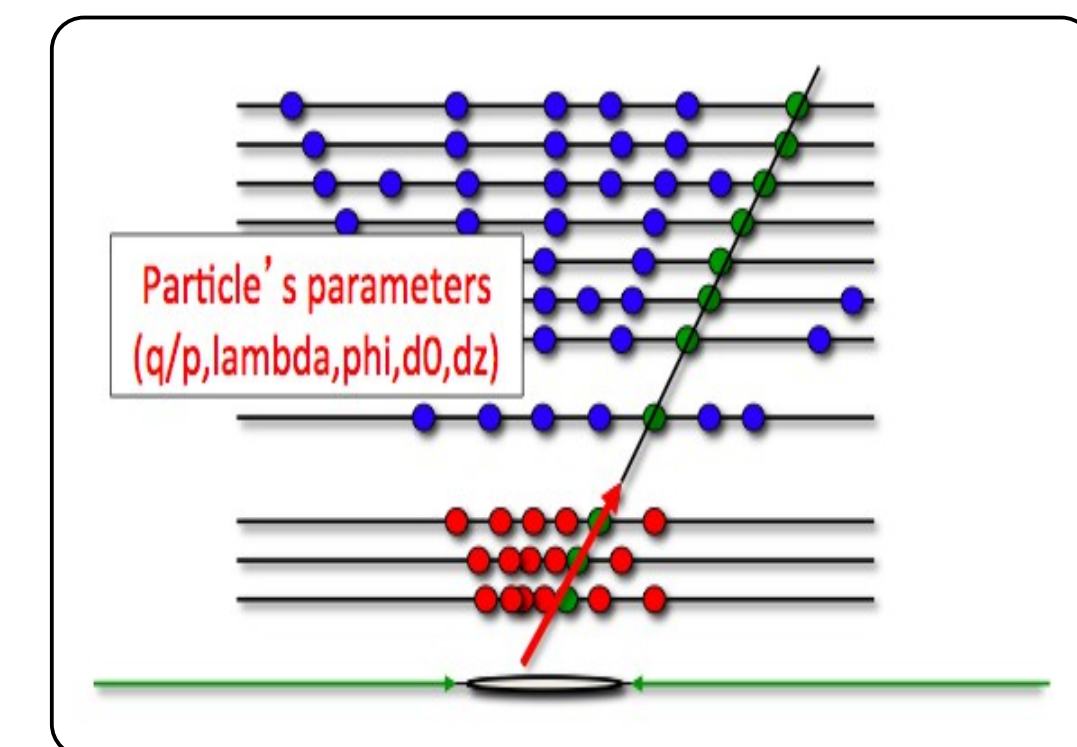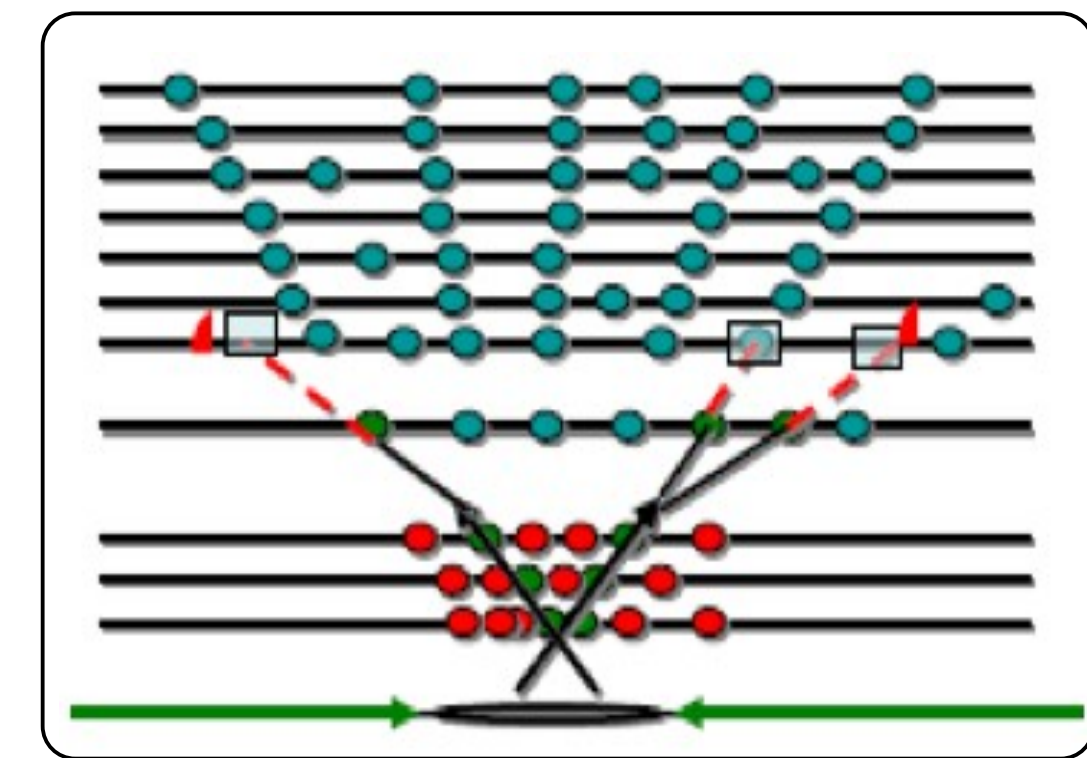- Muon reconstruction with calorimeters

# Tracking

- *Tracking is the pattern-recognition task that builds particle trajectories from a set of recorded hits*

- *One of the slowest tasks we perform to reconstruct particles in LHC collisions*

- *Non-linear slow-down with number of simultaneous collisions, due to combinatoric effects*

# Tracking

- *Works in three steps*

  - *seeding: start from pair of hits in the inner detector*

  - *hit-to-track association: propagate the seed and look for hits close ton the predicted trajectory*

  - *Track fitting: measure the track parameters (particle energy) from a fit of the points to an helix trajectory*

# Deep Learning to the Rescue

◉ *The detector sees the charge deposited by the crossing*

◉ *A hit is a window of sensors (16x16 here) with its deposited charge. This can be seen as a sparse digital image.*

◉ *Given two images, one can train a network to decide if a pair of hits is a good or bad match*

Barrel 1 - Inner Cluster

Fwd_Endcap 1 - Outer Cluster

# PixelSeed ConvNN

- The final model uses two sets of inputs:

  - the hit images

  - a set of expert features (e.g., position of the hits in the detector) to help the learning process
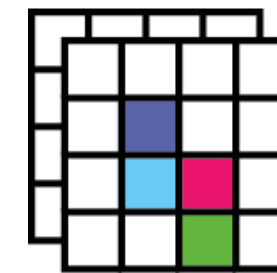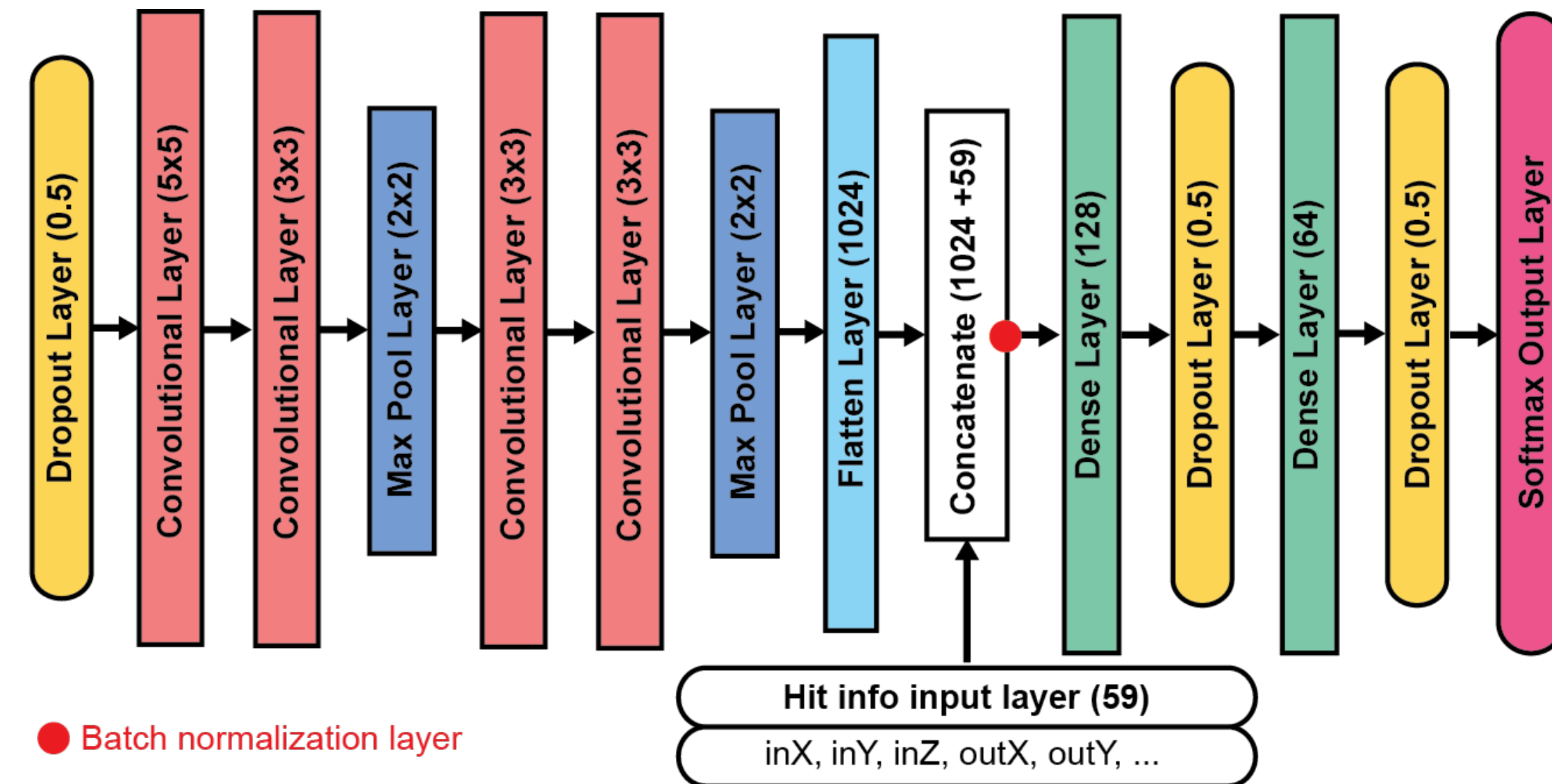
```
Dropout Layer (0.5) → Convolutional Layer (5x5) → Convolutional Layer (3x3) → Max Pool Layer (2x2) → Convolutional Layer (3x3) → Convolutional Layer (3x3) → Max Pool Layer (2x2) → Flatten Layer (1024) → Concatenate (1024 +59) → Dense Layer (128) → Dropout Layer (0.5) → Dense Layer (64) → Dropout Layer (0.5) → Softmax Output Layer
```

● Batch normalization layer

**Hit info input layer (59)**

inX, inY, inZ, outX, outY, ...

- **One can reduce the fake rate by one order of magnitude with a few % loss in efficiency**

**Efficiency (tpr) @ fake rejection**

```
tpr @ rej 50%: 0.998996700259
tpr @ rej 75%: 0.990524391331
tpr @ rej 90%: 0.922210826719
tpr @ rej 99%: 0.338669401587
```

## Layer Map Output Score

Legend:
- Test - True
- Train - True
- Test - Fake
- Train - Fake

counts/0.02 vs $p_{true}$

# HEP & Language processing networks

# Particle (language) processing

- CMS uses particle flow for event reconstruction:
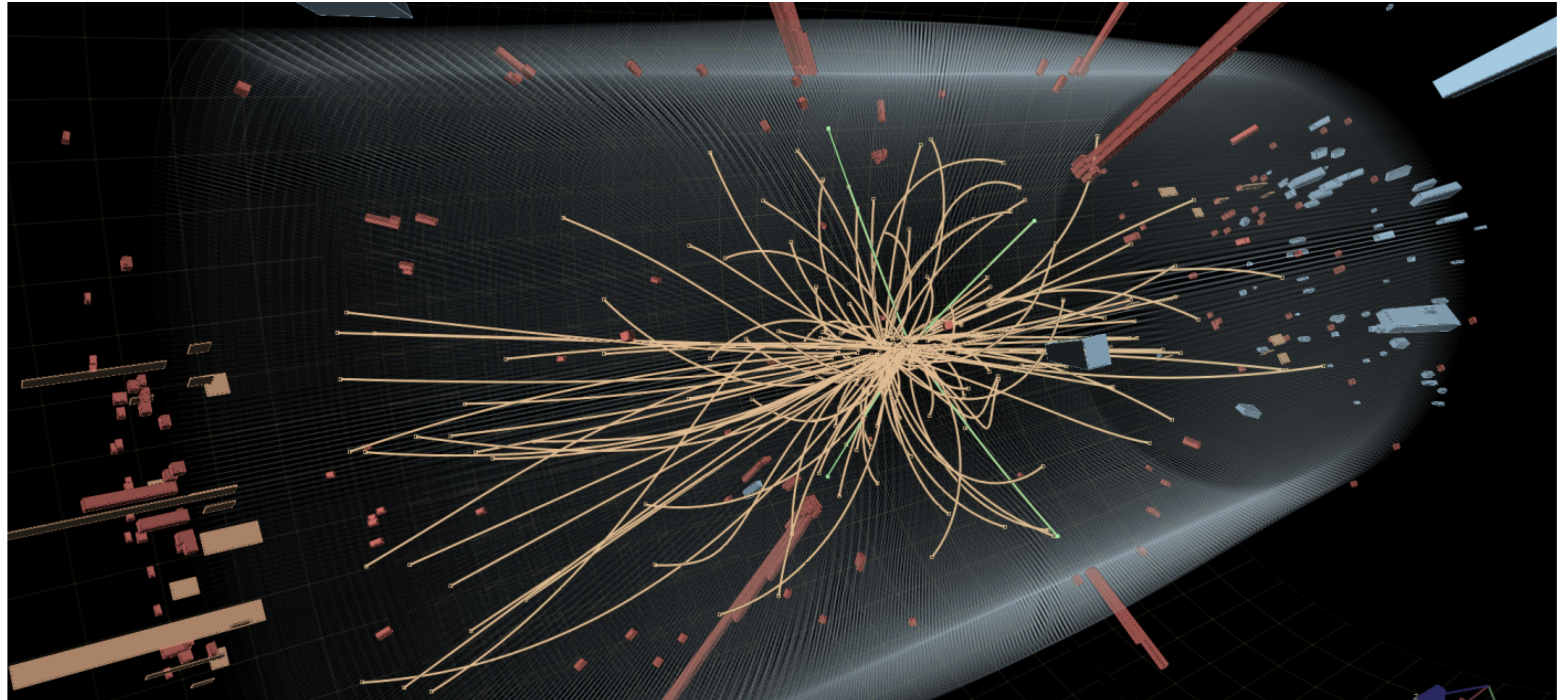
  - At some point in the central processing, collision images are turned into a list of particles.

  - From these particles, complex objects (e.g., jets) are formed

- In this framework, Computing vision approaches are not necessarily ideal

- One can instead use language-processing approaches (e.g., recurrent neural networks

  - particles are words in a sentence

  - QCD is the grammar

Fermilab has a herd of bisons

Fermilab has a herd of bisons

$k_t$

# Recurrent Neural Networks

- *A network architecture suitable to process an ordered sequence of inputs*

  - *words in text processing*

  - *a time series*

  - *particles in a list*

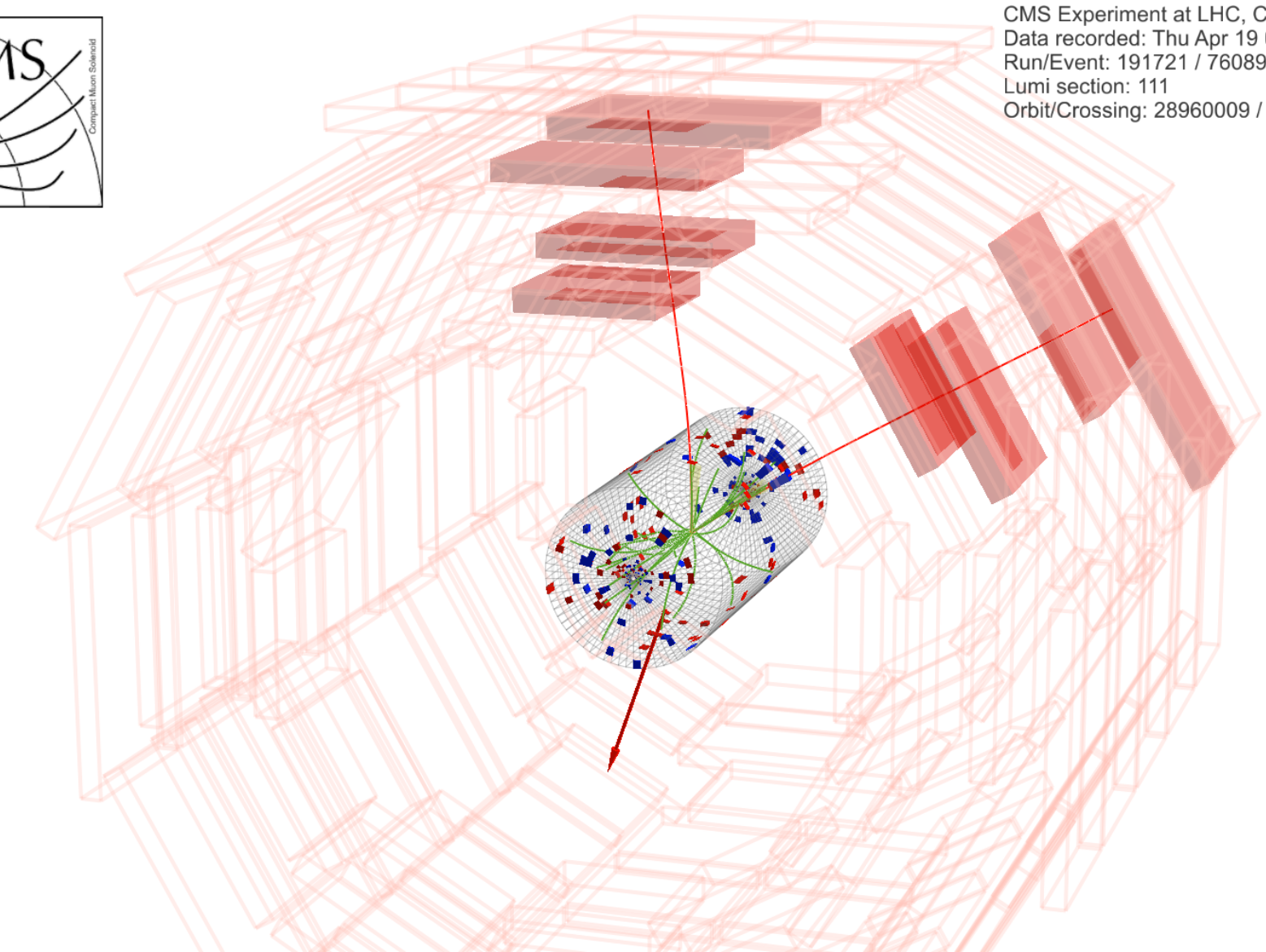- *Could be used for a single jet or the full event*

- *Next step: graph networks (active research direction)*



Jet

# A Topology Classifier

*A typical example: leptonic triggers*

- at the LHC, producing an isolated electron or muon is very rare. Typical smoking gun that something interesting happened (Z,W,top,H production)-> *TAKE THEM!*

- Triggers like those are very central to ATLAS/CMS physics

- *The sample selected is enriched in interesting events, but still contaminated by non-interesting ones*

- Can we clean this up w/o biasing the physics? yes, with ML

CMS Experiment at LHC, CERN
Data recorded: Thu Apr 19 09:14:14 2012 CEST
Run/Event: 191721 / 76089774
Lumi section: 111
Orbit/Crossing: 28960009 / 815

lepton Isolation + pT threshold

- tt
- W+jets
- QCD

7%

92%

*See contribution to NIPS workshop*

# A Topology Classifier



$W + \text{jets}$ (45.6%)

QCD (53.9%)

$t\bar{t}$

Deep Topology Classifier

**Calo Image Classifier**

Raw images of the calorimetry hits fed to a convolutional NN.

**Particle Sequence Classifier**

A sequence of particles taken as input to a recurrent NN.

**Abstract Image Classifier**

Based on an abstract representation of the reconstructed particles as an image to feed to a convolutional NN.

**High-level Feature Classifier**

Use high-level features as inputs to a fully connected NN.

$W + \text{jets}$

Photons — Charged Tracks — Neutral Hadrons

Forward Endcap Barrel Endcap Forward Forward Endcap Barrel Endcap Forward Forward Endcap Barrel Endcap Forward

$W + \text{jets}$ QCD $t\bar{t}$

Photons — Charged Tracks — Neutral Hadrons

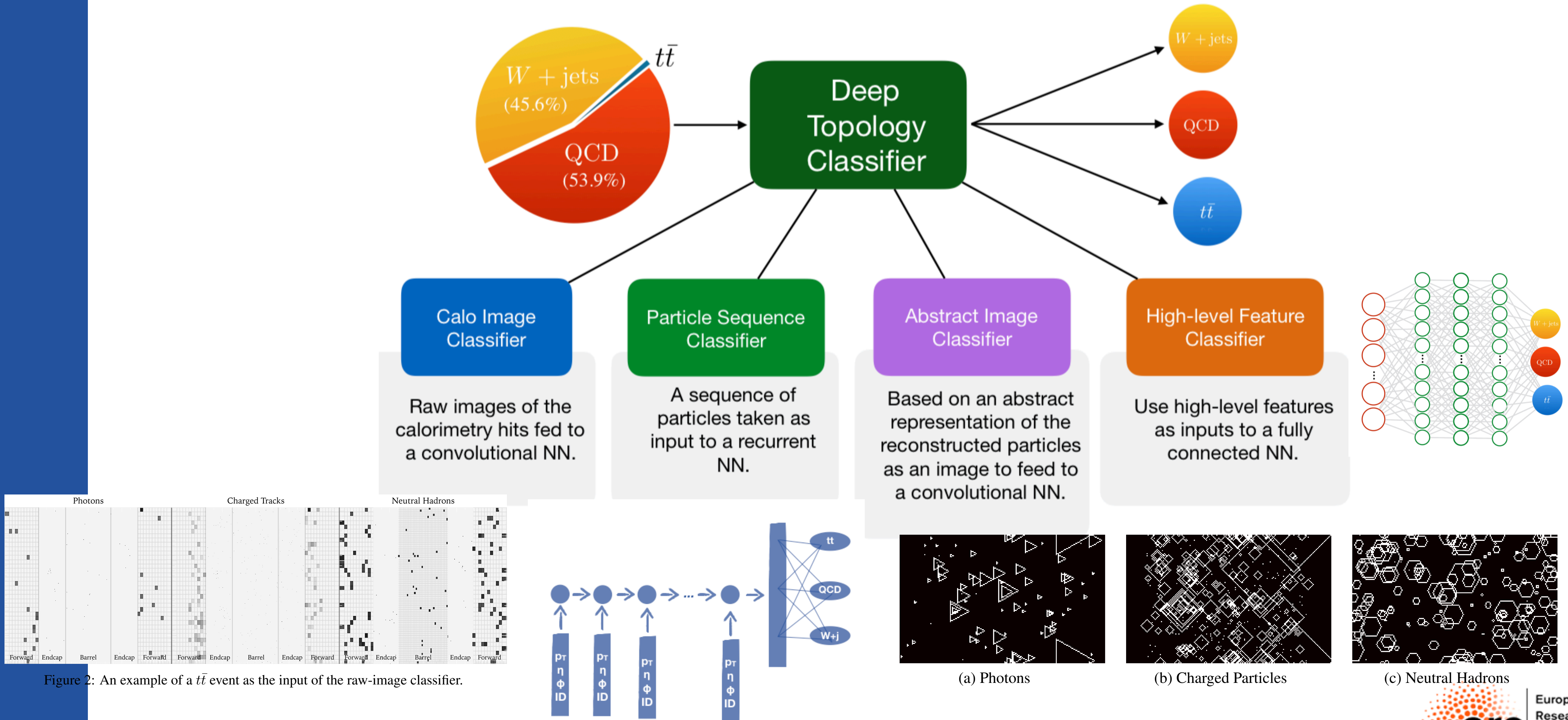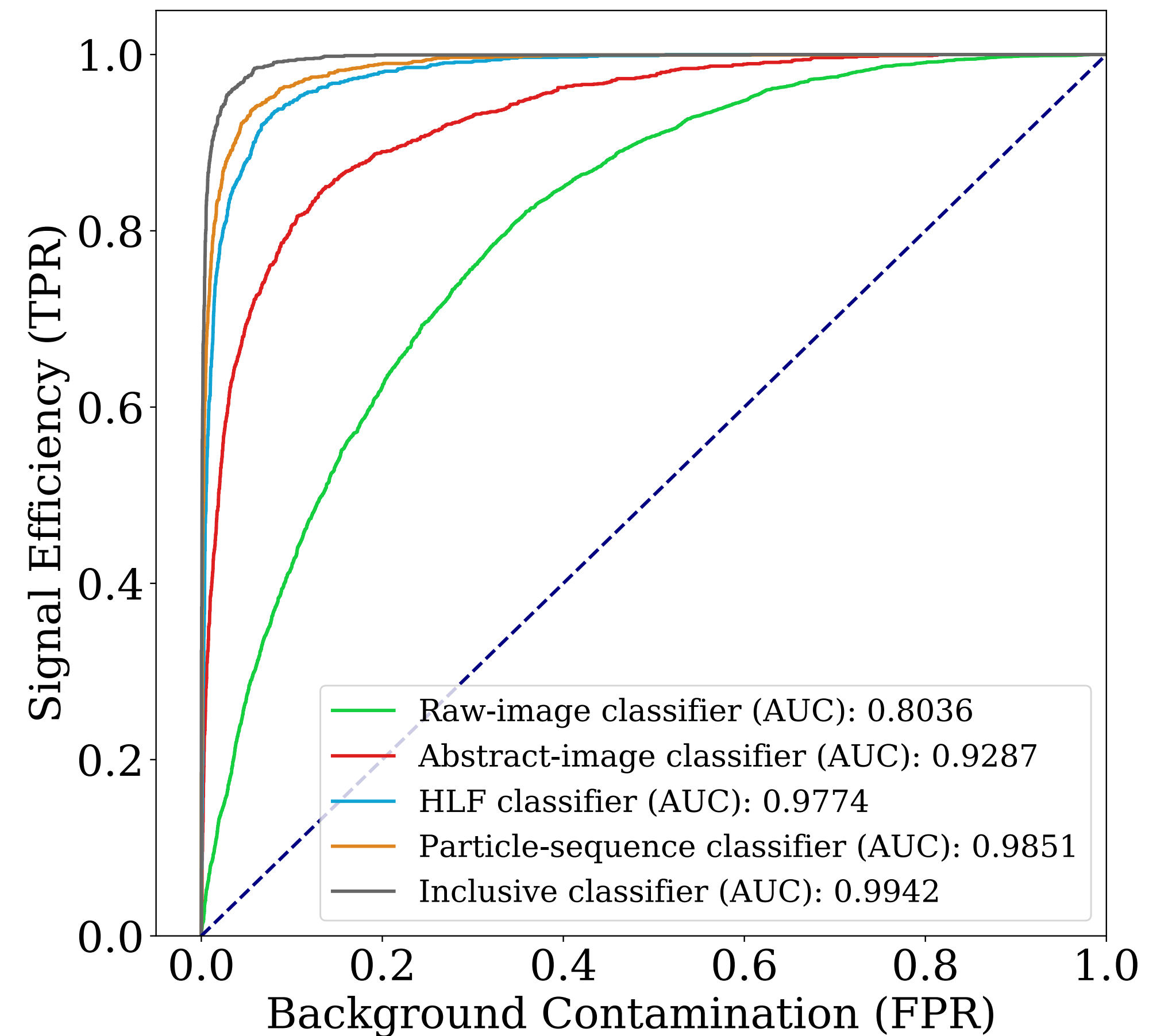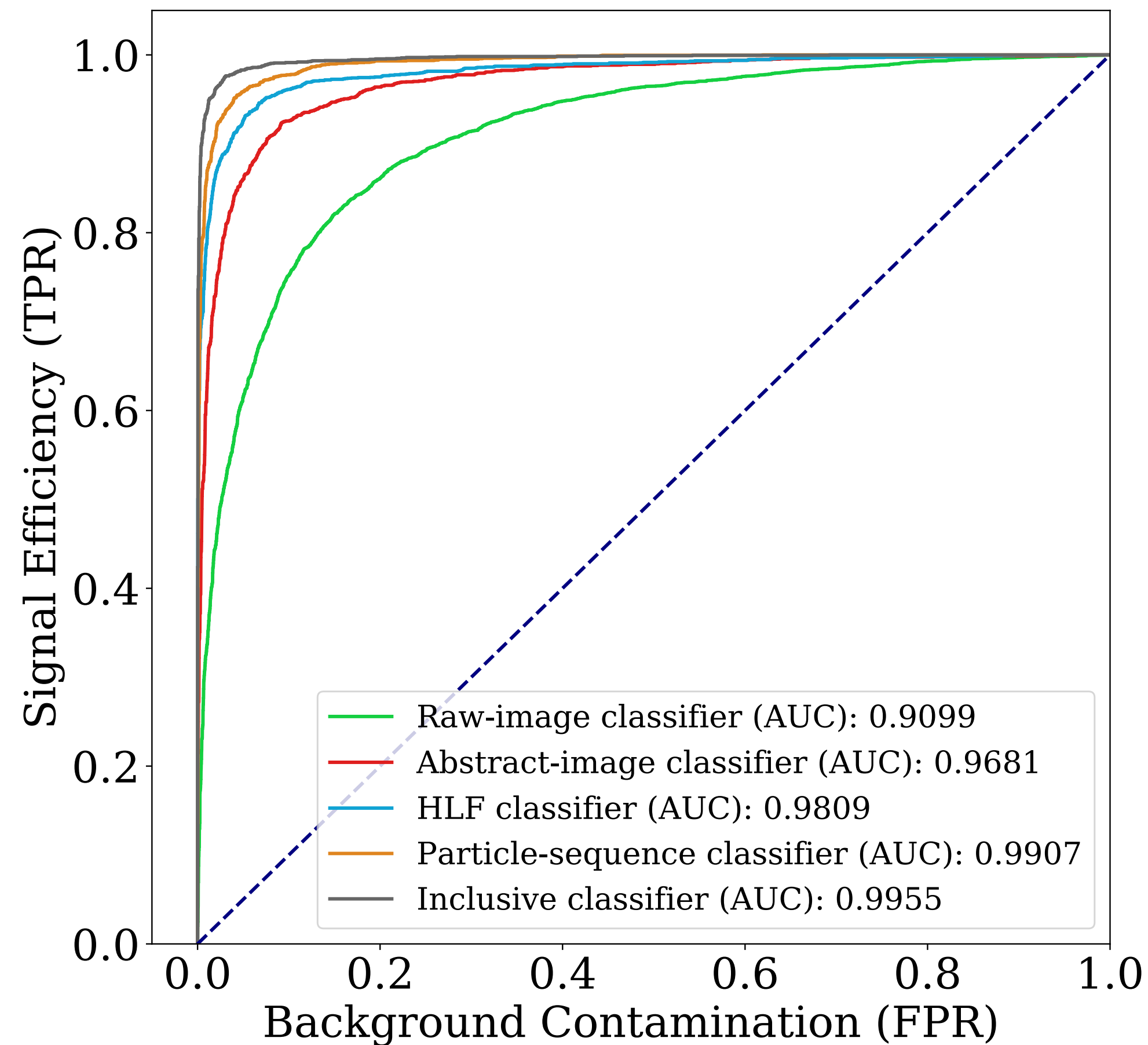Forward Endcap Barrel Endcap Forward Forward Endcap Barrel Endcap Forward Forward Endcap Barrel Endcap Forward

Figure 2: An example of a $t\bar{t}$ event as the input of the raw-image classifier.

$p_T$ η φ ID

tt QCD W+j

(a) Photons   (b) Charged Particles   (c) Neutral Hadrons

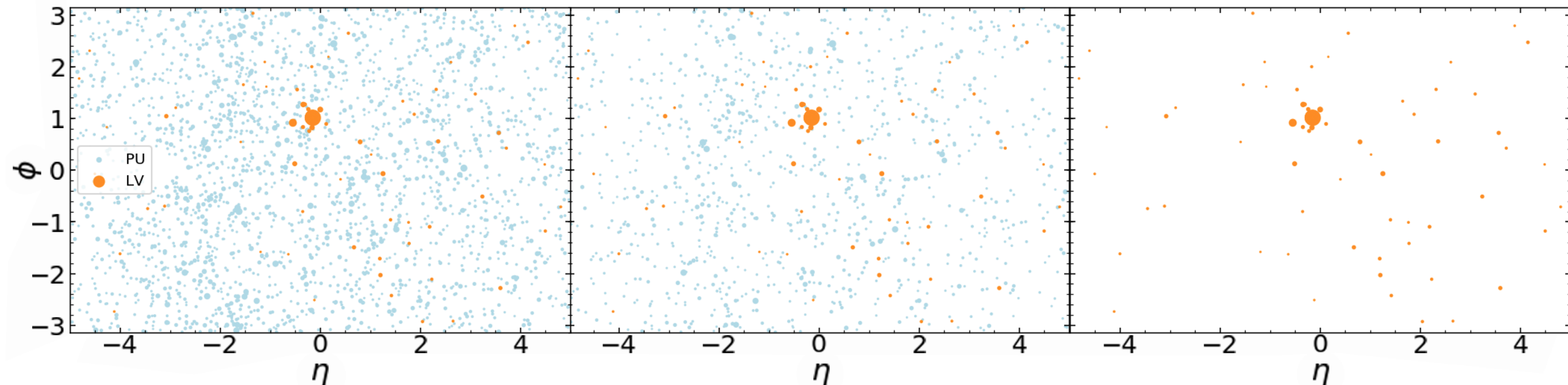European Research Council

erc

# Selection performances



**Can select 99% of the top events and reduce the fraction of written events by a factor ~ 7**

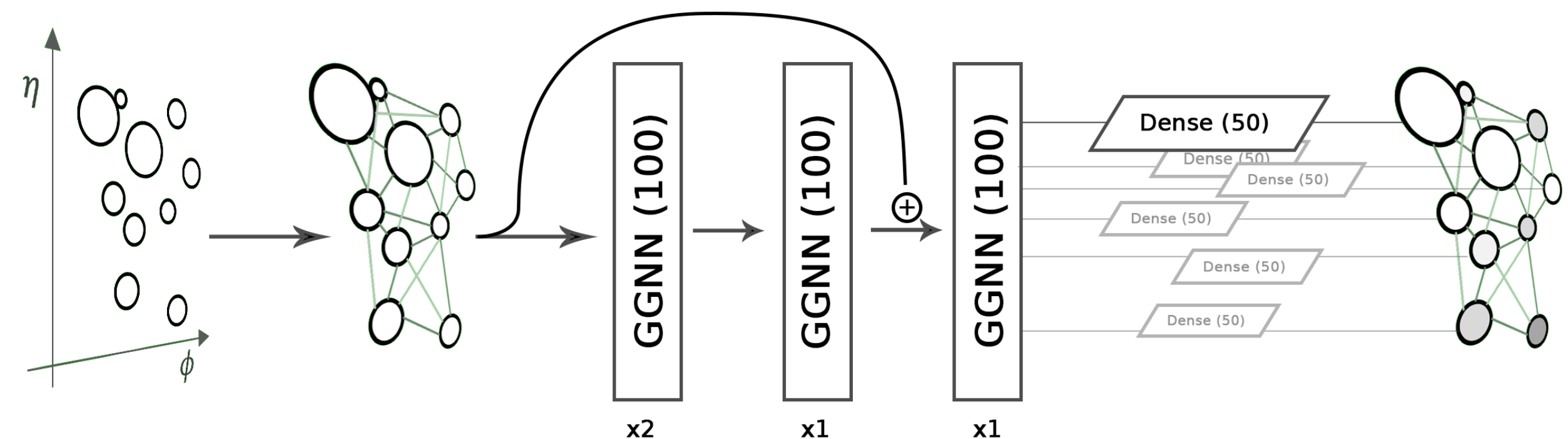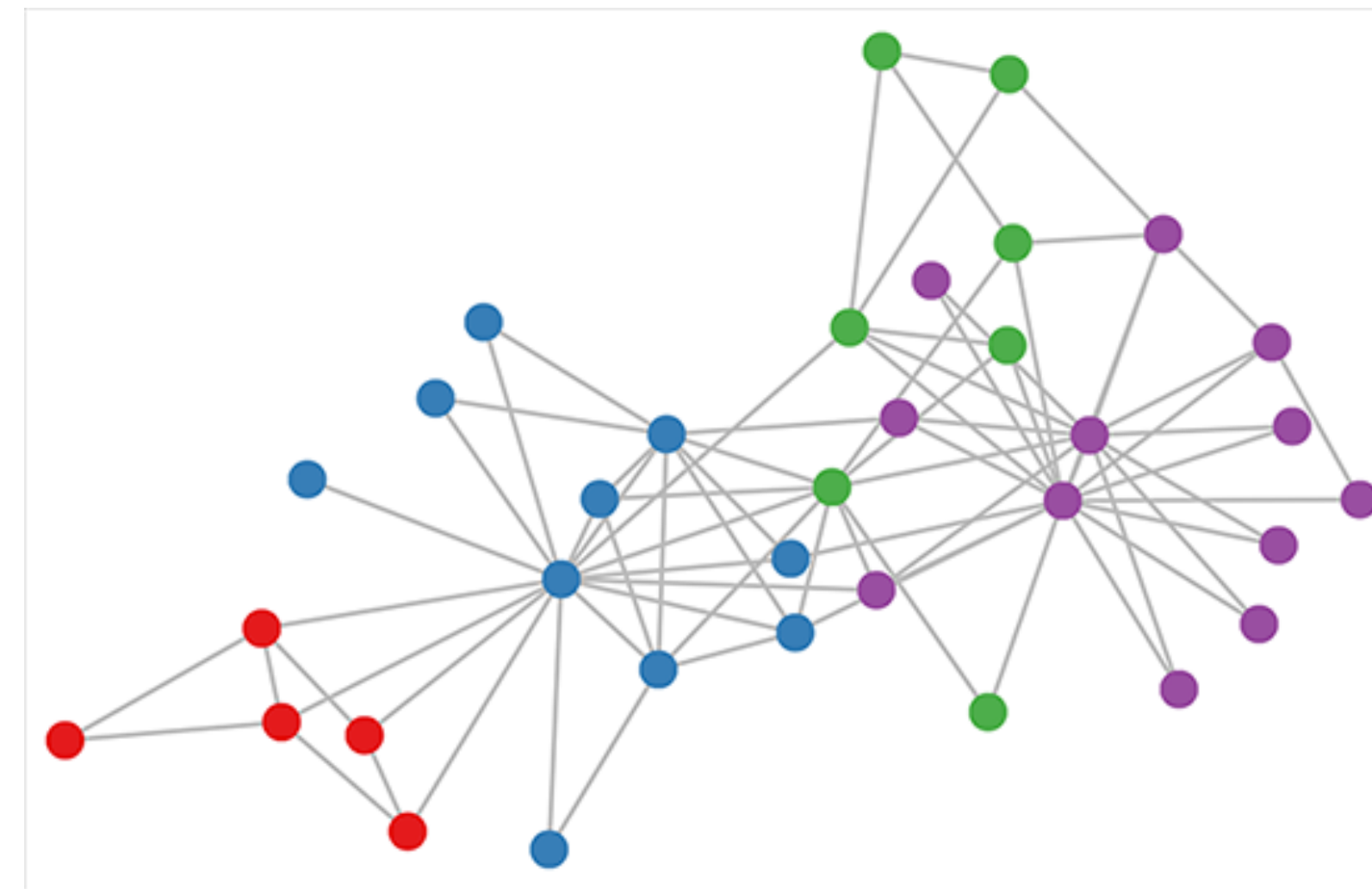# Graph Networks for particle physics

# The pileup problem



- ◉ *Pileup introduces noise to any reconstruction algorithm*

- ◉ *Usually, one runs a PU subtraction algorithms first*

  - ◉ *Usually based on global information of the event (average occupancy vs observed local occupancy)*

  - ◉ *OK offline, sort of OK @HLT, complicated @ L1*

- ◉ *State-of-the-art algorithms (Softkiller, PUPPI) improved situation dramatically wrt Run I*

- ◉ *Now ATLAS and CMS deal with ~40 collisions/bunch crossing with no problems*
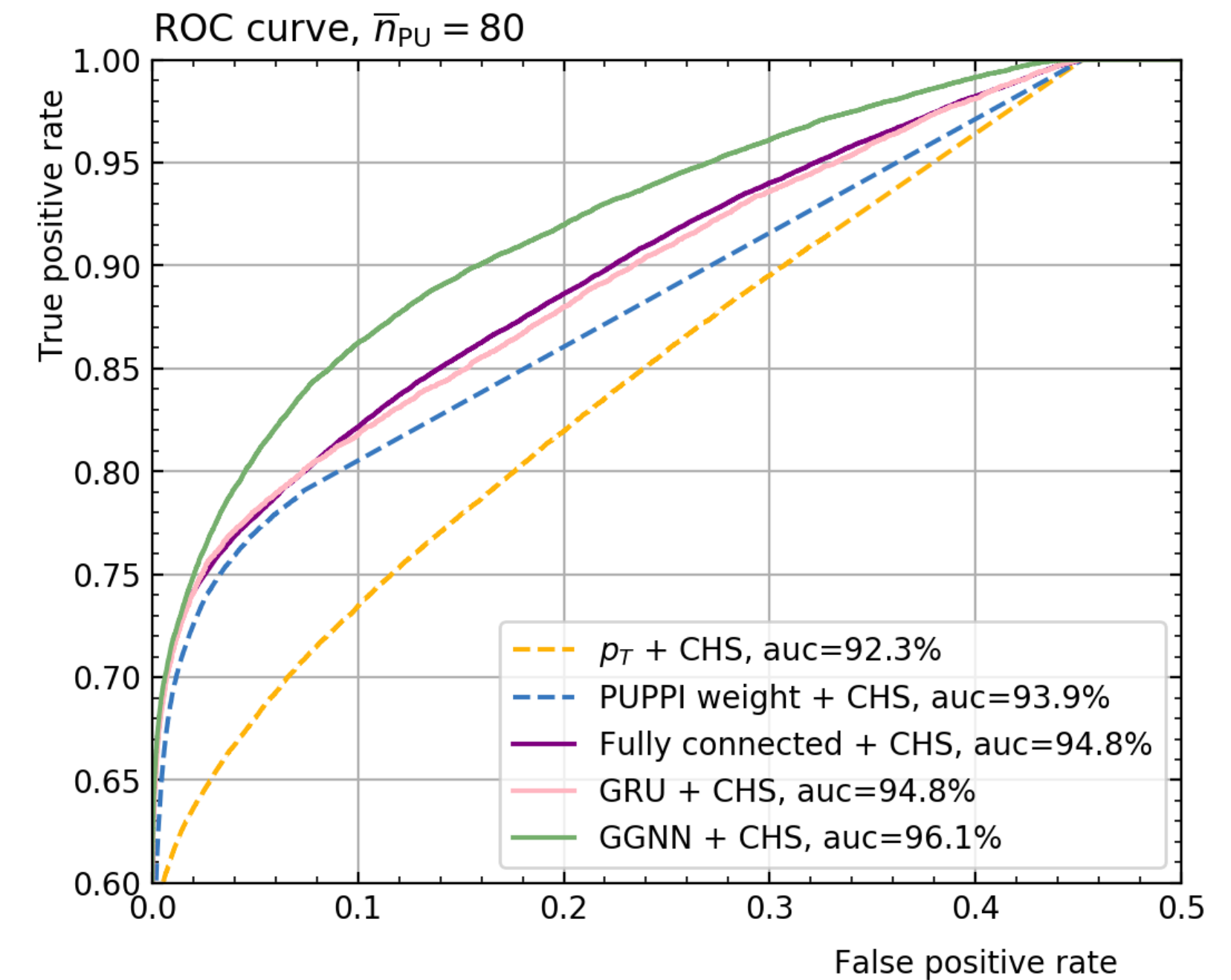
# Graph Networks

- Graph networks can be seen as generalization of Conv NN

  - Network learns from single "pixel" (graph node) and its neighbours

  - The concept of neighbour is not driven by geometrical proximity

  - Instead, what is "close" and what is not depends on connections (graphs) which are learned in the training

- We used a Gated Graph NN to decide if a given particle is from PU or not, based on its neighbours charged particles (which can be tracked to a vtx) are pileup or not

# PUPPIML: Graph Nets for PU subtraction

- ⊚ *Improve state-of-the-art algorithms substantially*

| $\overline{n}_{\mathrm{PU}}$ | 20 (CHS) | 80 (CHS) | 140 (CHS) | 80 (No CHS) |
|---|---|---|---|---|
| $p_T$ | 92.3% | 92.3% | 92.5% | 64.9% |
| PUPPI weight | 94.1% | 93.9% | 94.4% | 65.1% |
| Fully-connected | 95.0% | 94.8% | 94.8% | 68.5% |
| GRU | 94.8% | 94.8% | 94.7% | 68.8% |
| GGNN | **96.1%** | **96.1%** | **96.0%** | **70.1%** |

- ⊚ *Little dependence of algorithm tuning on pileup conditions*

- ⊚ *Small/No performance loss with average number of PU collisions*



Pileup mitigation at the Large Hadron Collider with Graph Neural Networks

**J. Arjona Martínez,**[a,b,c] **O. Cerri,**[b] **M. Pierini,**[c] **M. Spiropulu**[b] **and JR. Vlimant**[b]

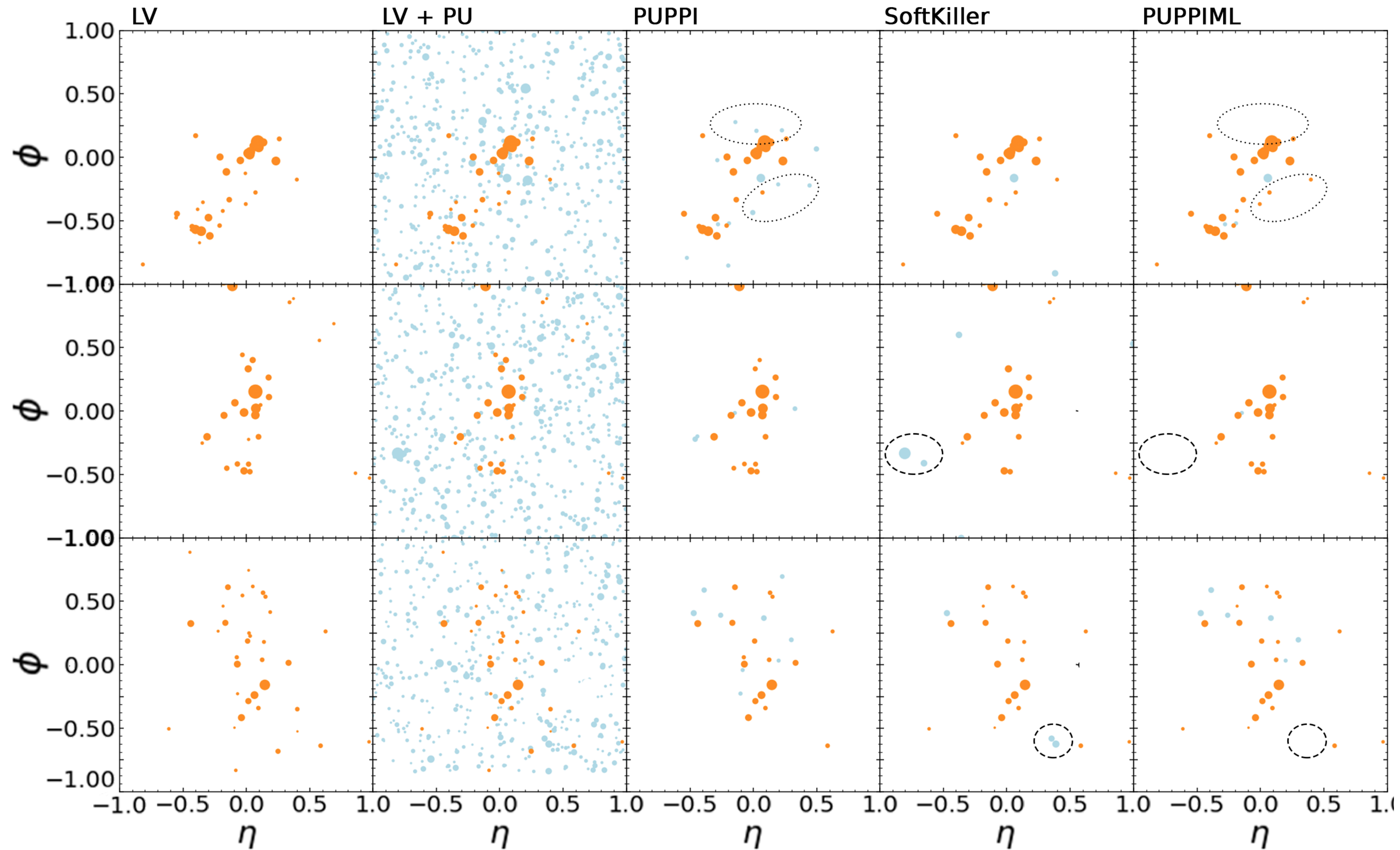[a] *University of Cambridge, Trinity Ln, Cambridge CB2 1TN, UK*
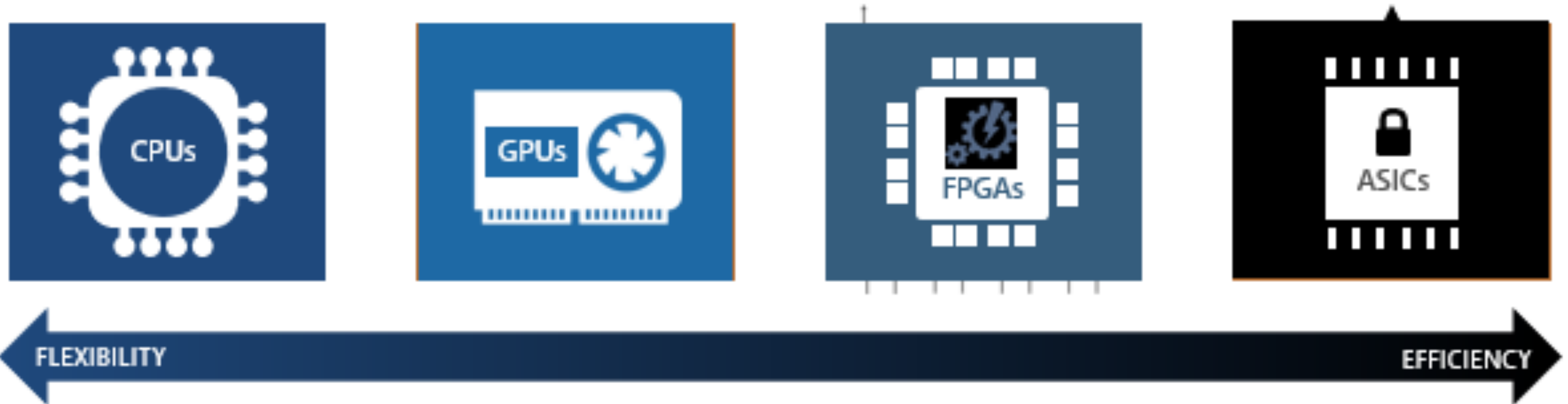
[b] *California Institute of Technology, 1200 E. California Blvd, Pasadena, CA 91125*

[c] *CERN, CH-1211 Geneva, Switzerland*

*E-mail:* ja618@cam.ac.uk, olmo.cerri@cern.ch, maurizio.pierini@cern.ch, smaria@caltech.edu, jvlimant@caltech.edu
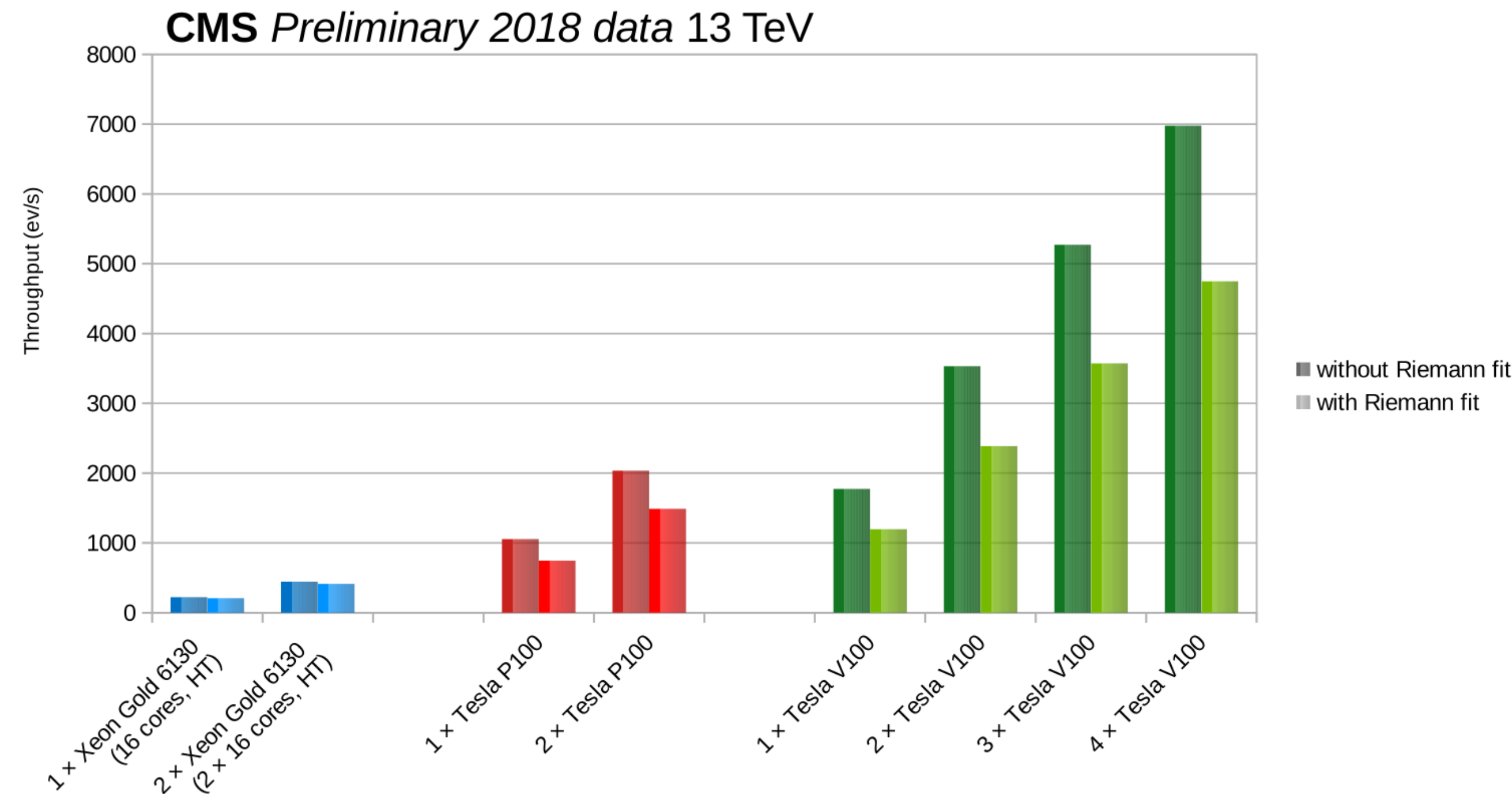
FLEXIBILITY ← → EFFICIENCY
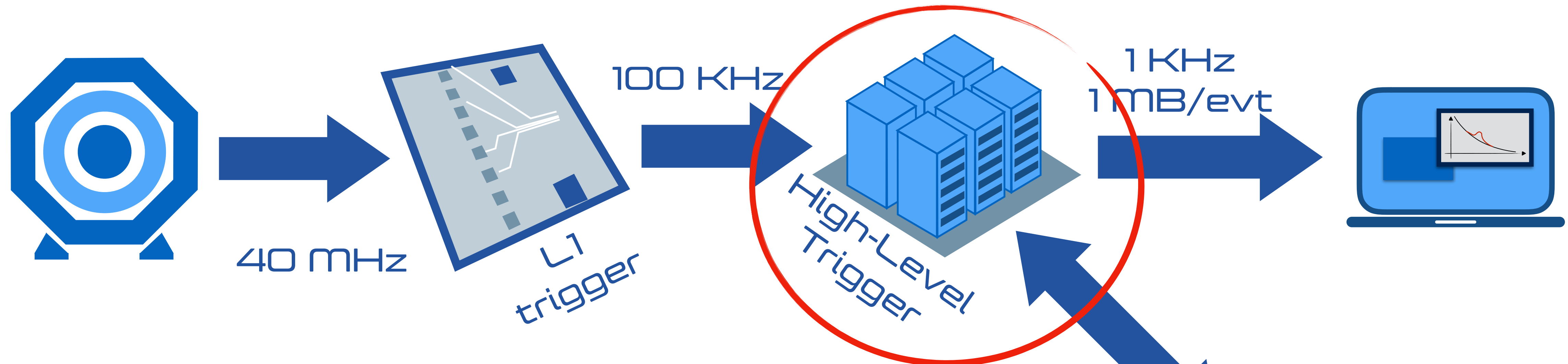
# Porting Deep Learning to Trigger/DAQ system

# The next HLT

⦿ *Looking at current tendency, we expect the next trigger system to be based on heterogenous computing, with GPUs & FPGAs used as accelerators to compensate saturation of Moore's law*

  ⦿ *for tracking, clustering, etc*

⦿ *In such a system, Deep Learning inference could be made very fast*

  ⦿ *On GPUs, as long as batching can be exploited*

    ⦿ *No big gain running one inference at once*

    ⦿ *Gain if many "samples" are sent at once. Example: 1K tracks per event*

    ⦿ *If objects are made on GPUs, no need to move them back and forth*

  ⦿ *On FPGAs, without need of batching, as long as the model can fit the available resources (including resource recycle with fast access to memory)*

**Patatrack project for CMS HLT on GPUs**



CMS *Preliminary 2018 data* 13 TeV

# Heterogeneous HLT

40 MHz

L1 trigger

100 KHz

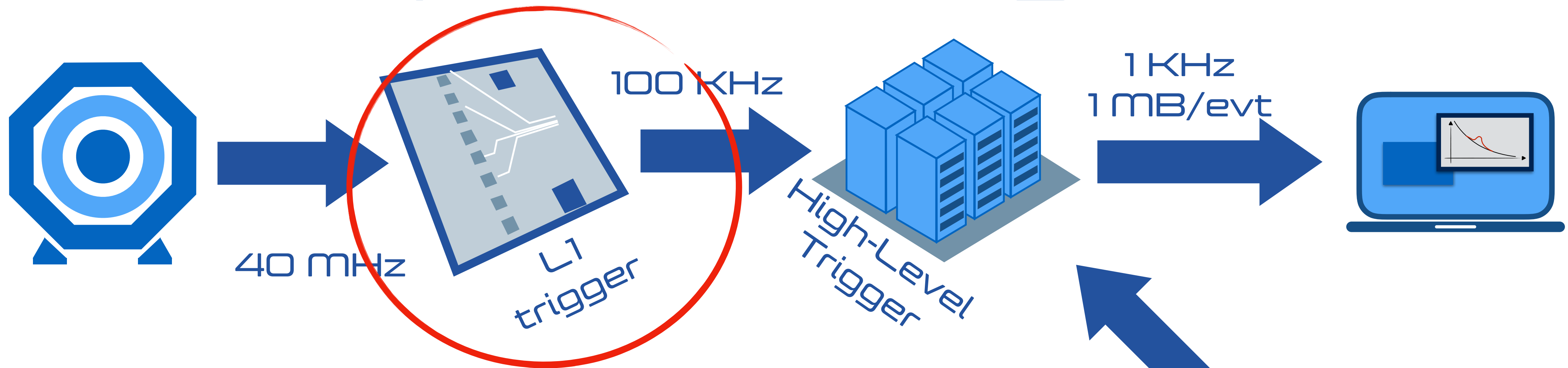High-Level Trigger

1 KHz
1 MB/evt

- *With heterogenous hardware in place (for other reasons) Deep Learning inference @HLT quite easy*

- ***Example:*** *the seed-selection for tracking I showed you before*

  - *1 μsec to know if a seed is good or not*

  - *1M seeds/event -> 1sec to process an event serially*
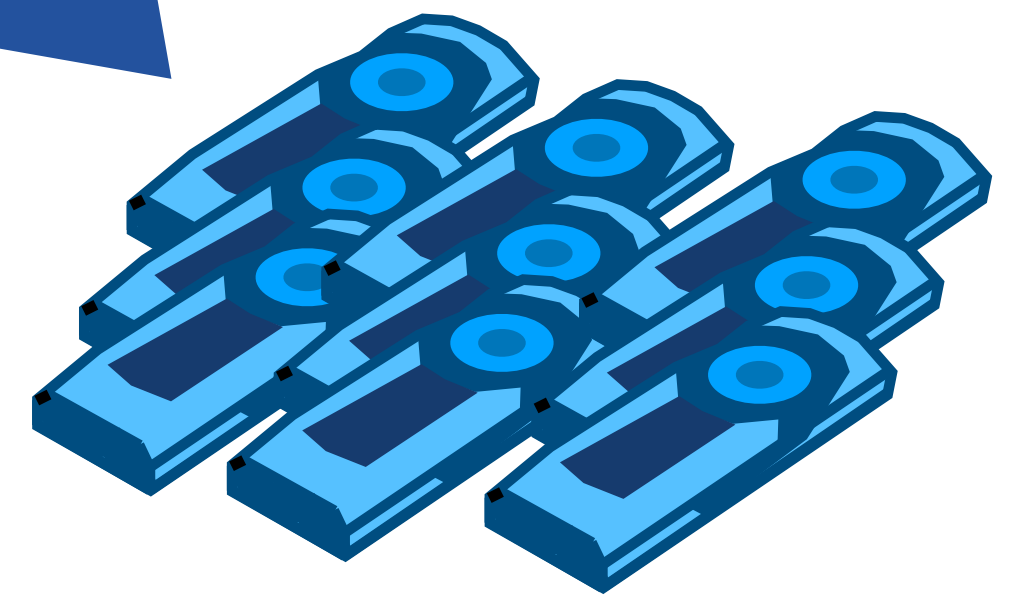
European Research Council

# Deep Learning at L1

◉ *Situation at L1 is different, mainly due to the typical latency (<10 μsec)*

◉ *Custom cards connected to detector electronics by optic links*

◉ *Data flow in the cards one by one*

◉ *Networks need to be implemented in FPGA firmware*

  ◉ *advanced design by expert engineers (not common resource in HEP)*

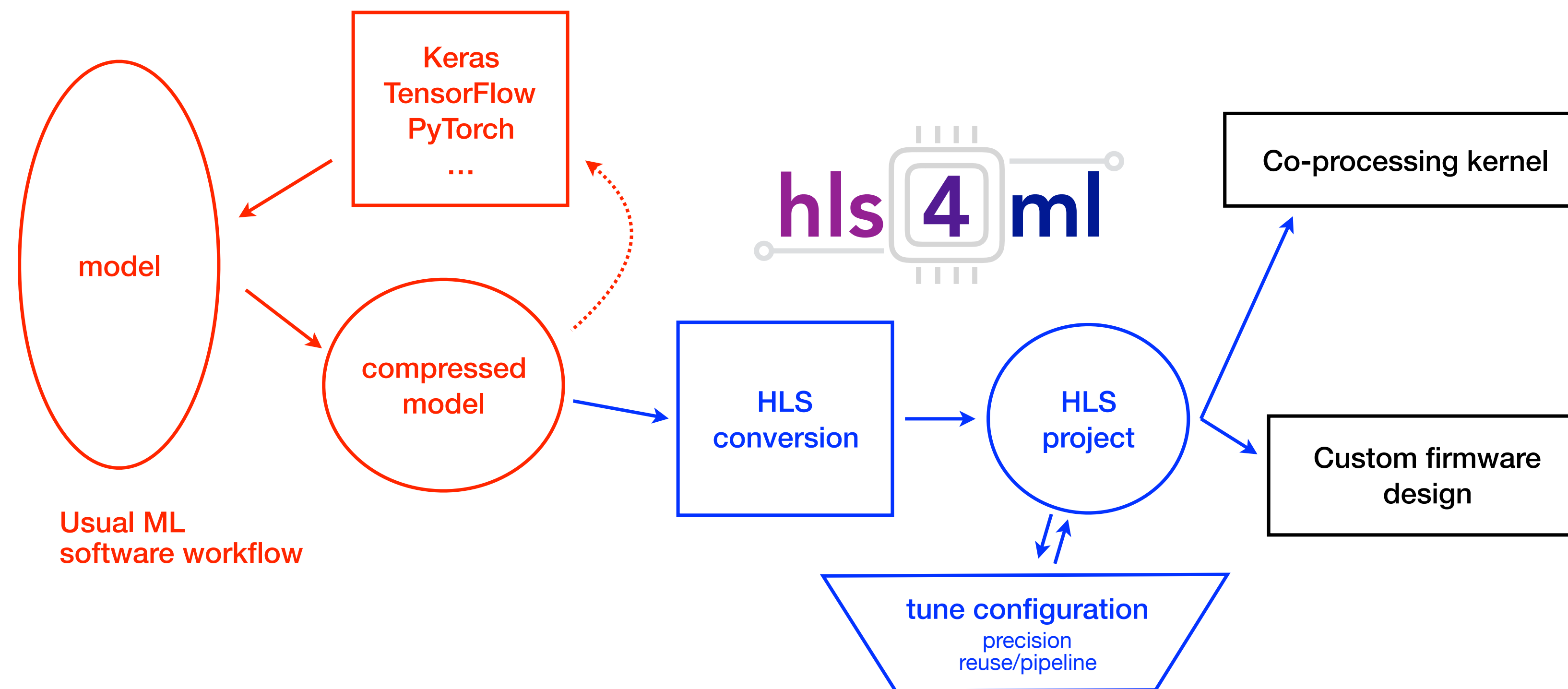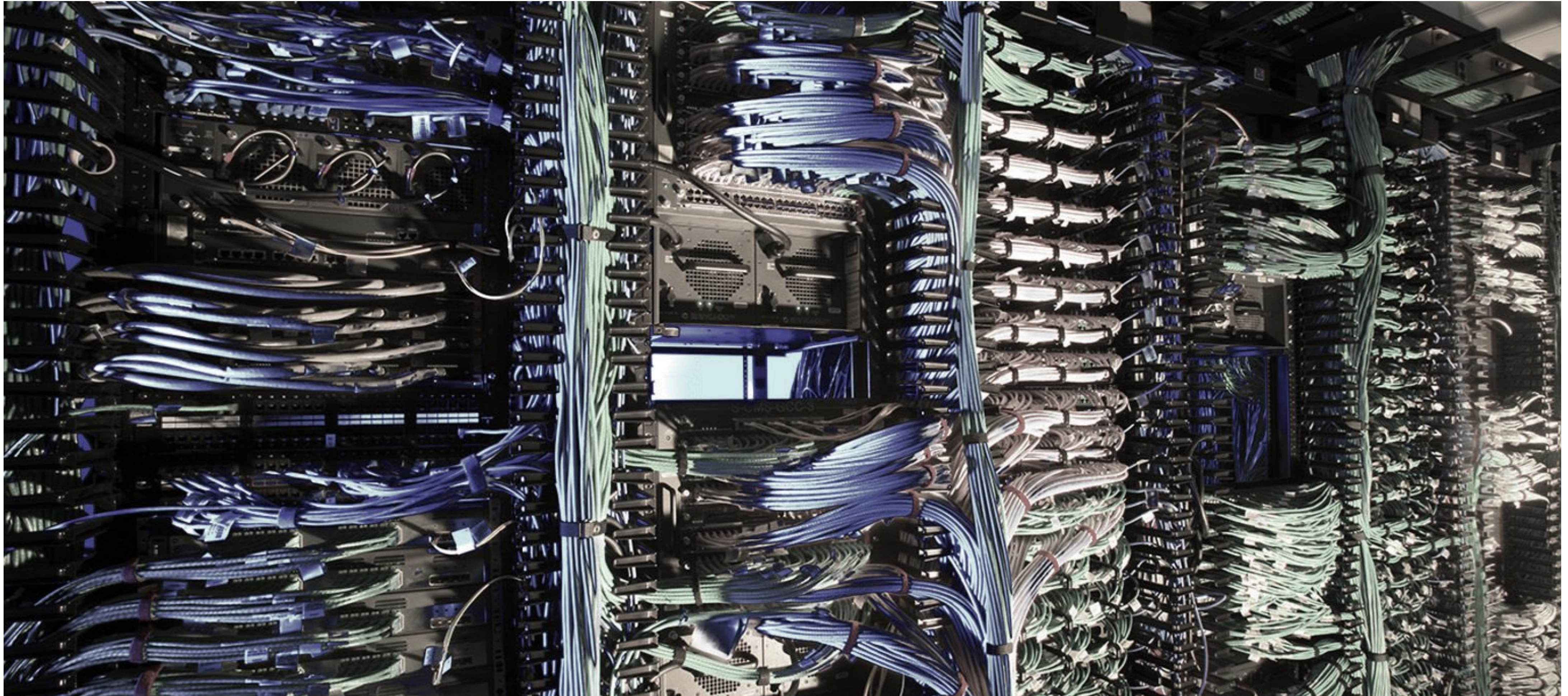  ◉ *automatic translation tools doing the job*

# Deep Learning at L1

100 KHz

1 KHz
1 MB/evt

40 MHz

L1 trigger

High-Level Trigger

◉ *Situation at L1 is different, mainly due to the typical latency (<10 μsec)*

◉ *Custom cards connected to detector electronics by optic links*

◉ *Data flow in the cards one by one*

◉ *Networks need to be implemented in FPGA firmare*

   ◉ *advanced design by expert engineers (not common resource in HEP)*

   ◉ *automatic translation tools doing the job*

European Research Council

# HLS4ML

◉ *HLS4ML aims to be this automatic tool*

  ◉ *reads as input models trained on standard DeepLearning libraries*

  ◉ *comes with implementation of common ingredients (layers, activation functions, etc)*

  ◉ *Uses HLS softwares to provide a firmware implementation of a given network*

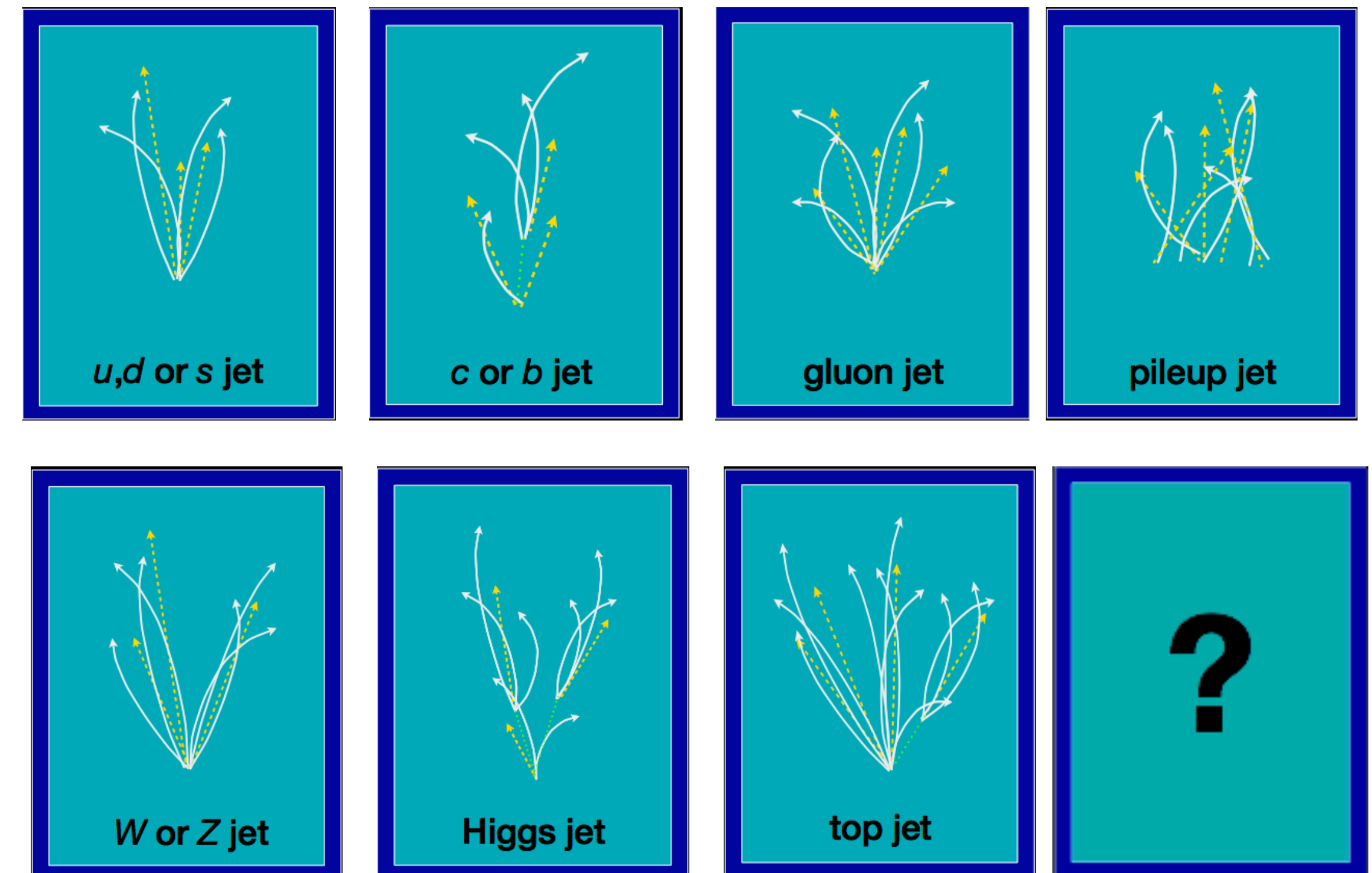  ◉ *Could also be used to create co-processing kernels for HLT environments*

# Fast Decision Taking

# Example: fast inference

◉ *You have a jet at LHC: spray of hadrons coming from a "shower" initiated by a fundamental particle of some kind (quark, gluon, W/Z/H bosons, top quark)*

◉ *You have a set of jet features whose distribution depends on the nature of the initial particle*

◉ *You can train a network to start from the values of these quantities and guess the nature of your jet*

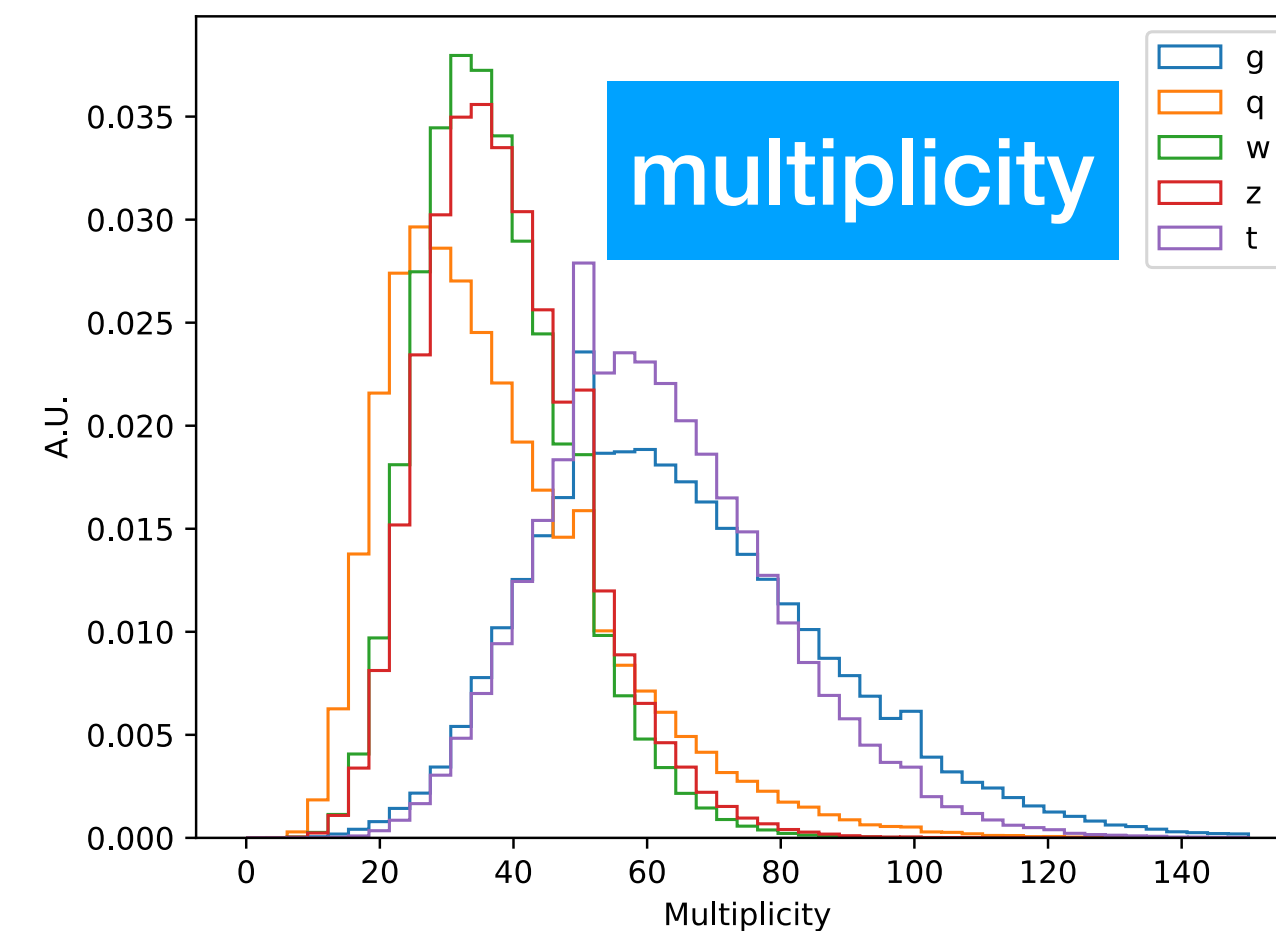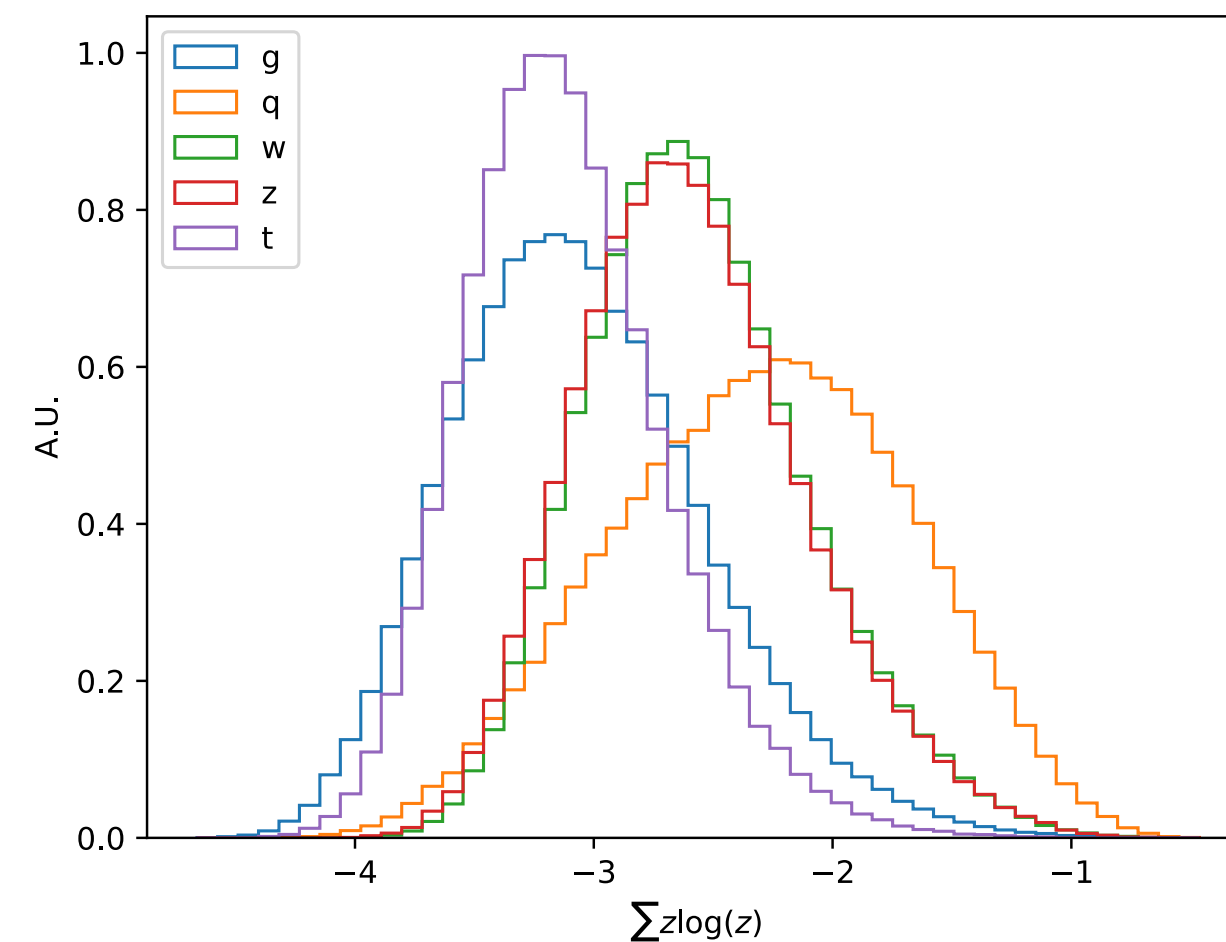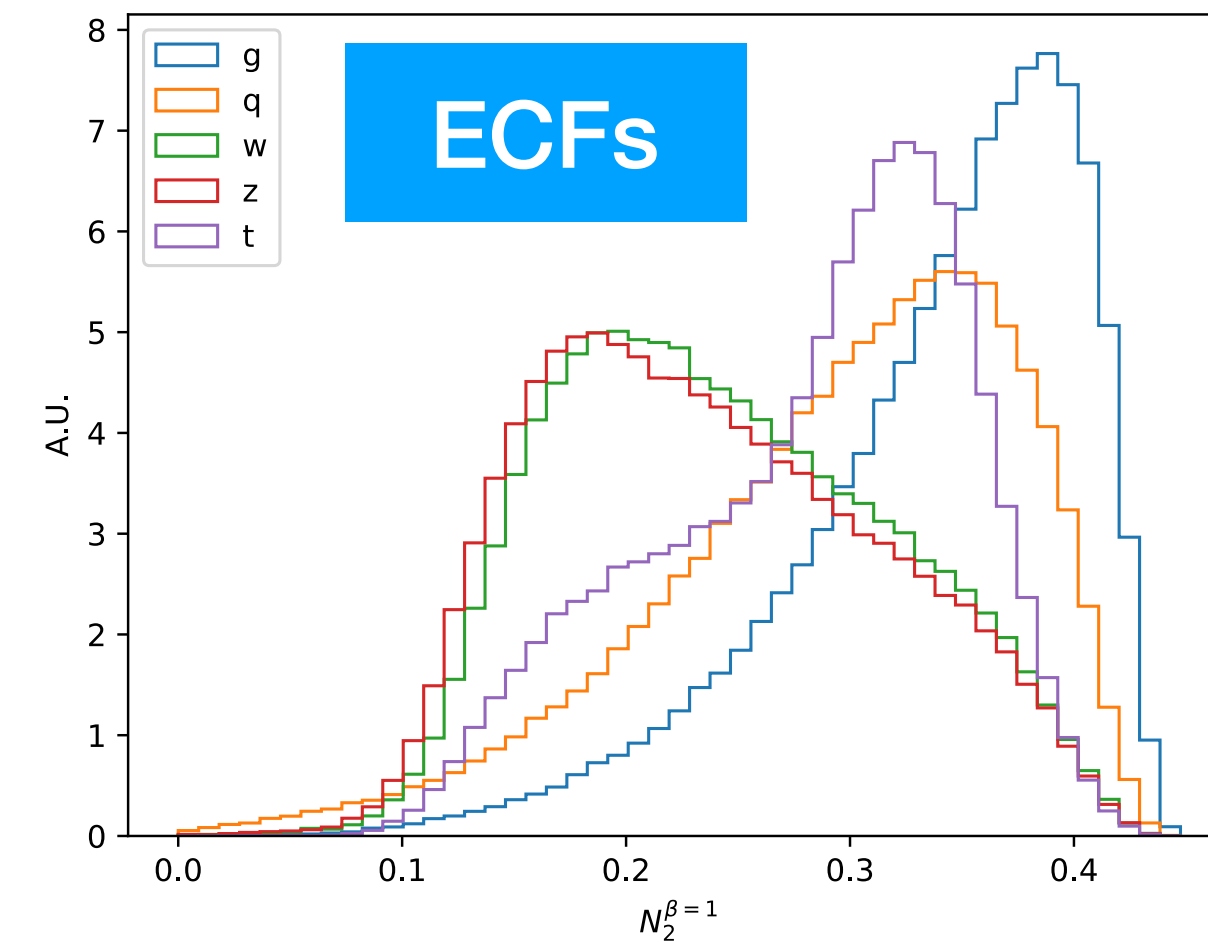◉ *To do this you need a sample for which you know the answer*

| | |
|---|---|
| *u,d* or *s* jet | *c* or *b* jet |
| gluon jet | pileup jet |
| *W* or *Z* jet | Higgs jet |
| top jet | ? |

- *Simple DNN bas
  multiplicities*

| **Observables** |
|:---:|

$$m_{\text{mMDT}}$$
$$N_2^{\beta=1,2}$$
$$M_2^{\beta=1,2}$$
$$C_1^{\beta=0,1,2}$$
$$C_2^{\beta=1,2}$$
$$D_2^{\beta=1,2}$$
$$D_2^{(\alpha,\beta)=(1,1),(1,2)}$$
$$\sum z \log z$$
Multiplicity

- 3-layer model trained without regularization

- No pruning applied

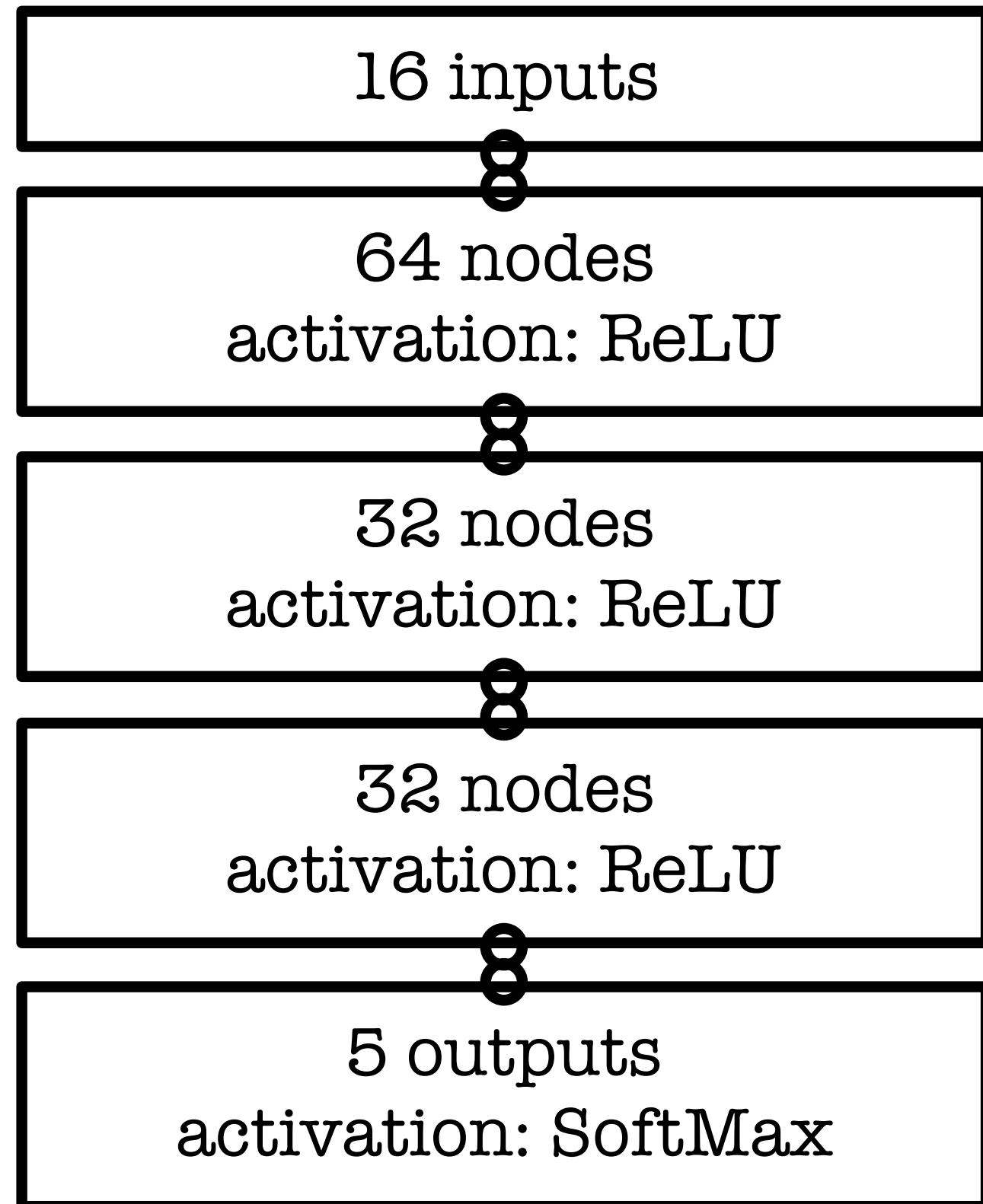- Resulting distribution of weights

*Simp...*
*high...*
*(je...*
*mul...*
*cor...*

16 inputs

64 nodes
activation: ReLU

32 nodes
activation: ReLU

32 nodes
activation: ReLU

5 outputs
activation: SoftMax

16 in

64 (r

32 (relu)



j_g tagger, auc = 91.0%
j_q tagger, auc = 88.7%
j_w tagger, auc = 92.1%
j_z tagger, auc = 91.4%
j_t tagger, auc = 93.7%

bkg. mistag rate

sig. efficiency

**Better**

5%  16%  50% 84% 95%

Number of Weights

Absolute Relative Weights

4

# Network Operations
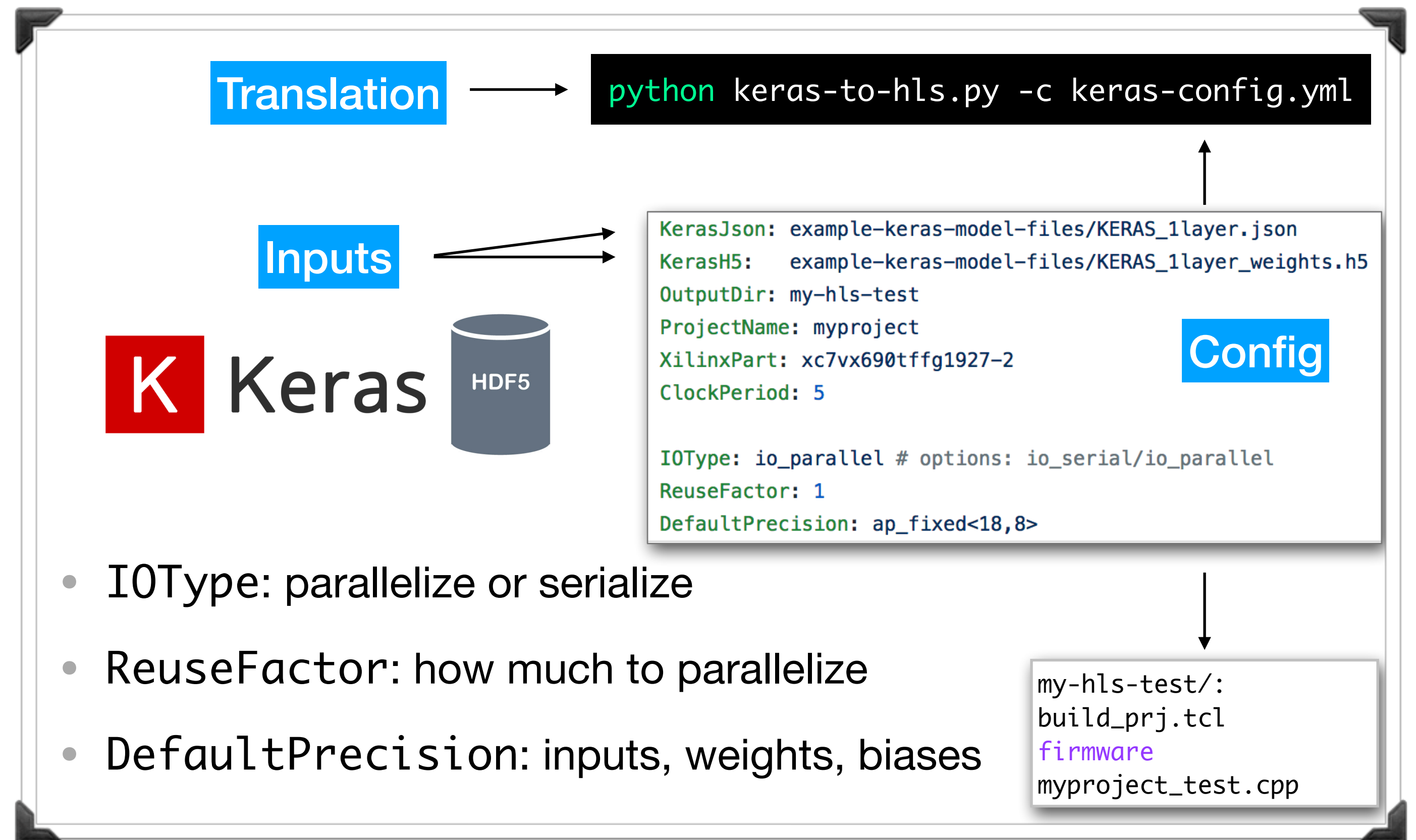
- A classic Dense NN manipulate. the inputs in three ways

  - multiplying by weights

  - adding biases

  - applying activation functions

- All these operations map nicely into an FPGA

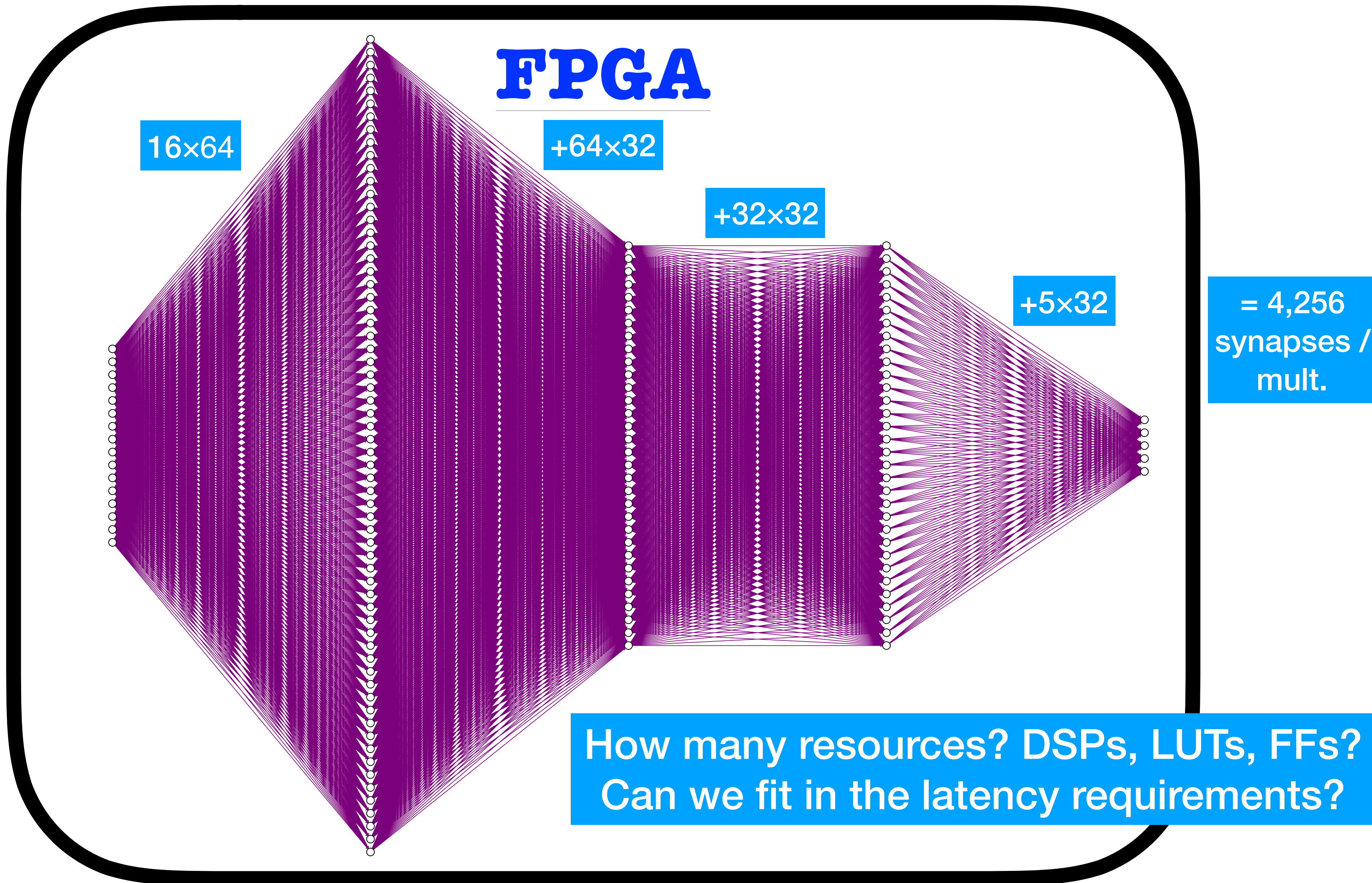  - high IO, DSPs, LUTs, tunable precision

**multiplication**

$$\ell_j^k = \phi(W_{ij}\ell_i^{k-1} + b_j)$$

**activation function**

**addition**

$\ell_i^{k-1}$

$W_{ij}$   $\ell_j^k$

# Bring the model to FPGA

● *How this works in practice*

  ● *A python based library that takes inputs via a yam file*

  ● *Model architecture with supported format*

  ● *FPGA configuration parameters (reuse factor, FPGA model, Clock period, etc)*

● *The library provides inputs for Vivado HLS*



```
Translation ──────▶ python keras-to-hls.py -c keras-config.yml

Inputs ──────▶

KerasJson: example-keras-model-files/KERAS_1layer.json
KerasH5:   example-keras-model-files/KERAS_1layer_weights.h5
OutputDir: my-hls-test
ProjectName: myproject                          Config
XilinxPart: xc7vx690tffg1927-2
ClockPeriod: 5

IOType: io_parallel # options: io_serial/io_parallel
ReuseFactor: 1
DefaultPrecision: ap_fixed<18,8>
```

- IOType: parallelize or serialize

- ReuseFactor: how much to parallelize

- DefaultPrecision: inputs, weights, biases

```
my-hls-test/:
build_prj.tcl
firmware
myproject_test.cpp
```

Javier Duarte | hls4ml

**FPGA**

16×64

+64×32

+32×32

+5×32

= 4,256 synapses / mult.

How many resources? DSPs, LUTs, FFs?
Can we fit in the latency requirements?

Javier Duarte I hls4ml
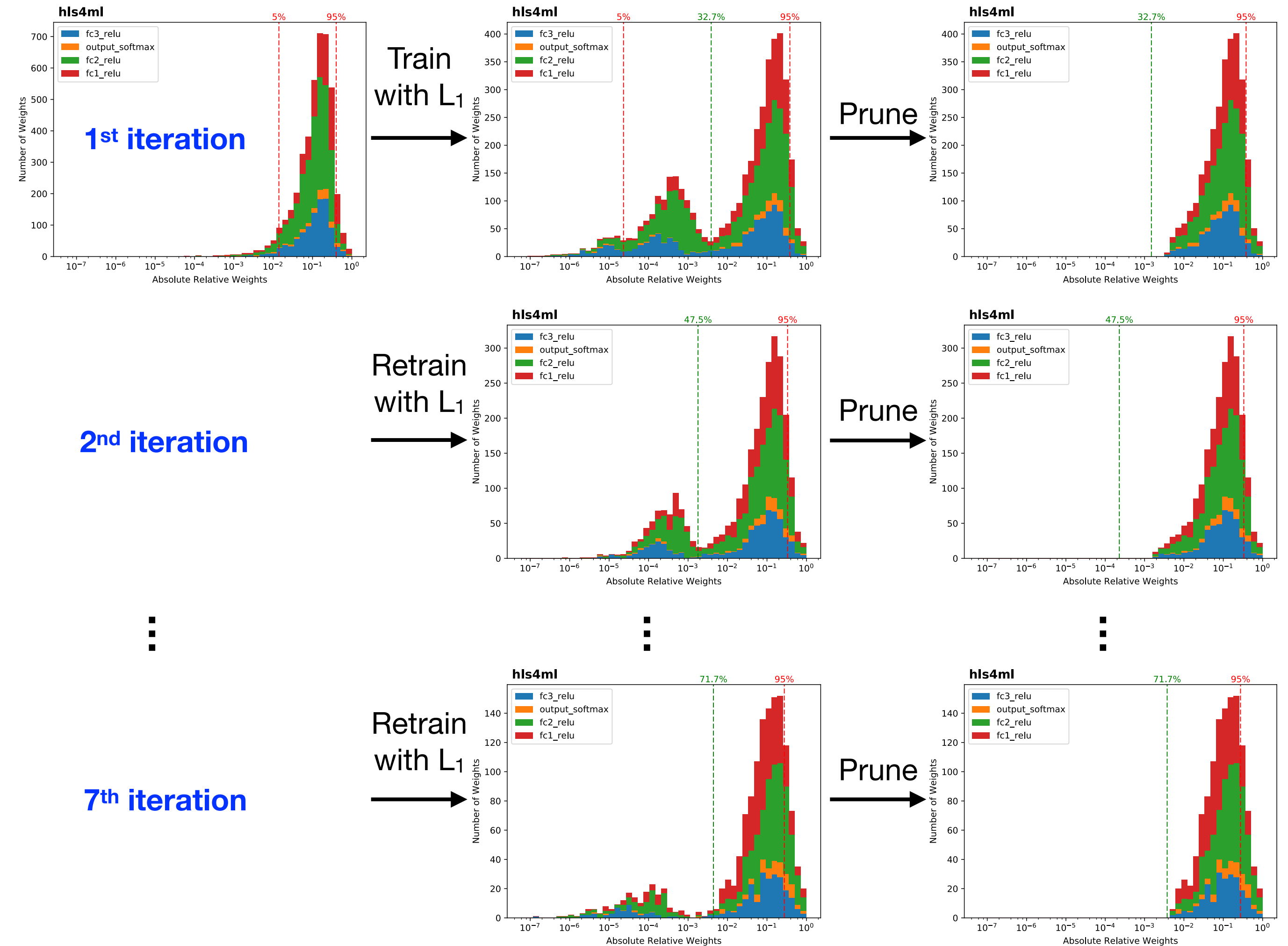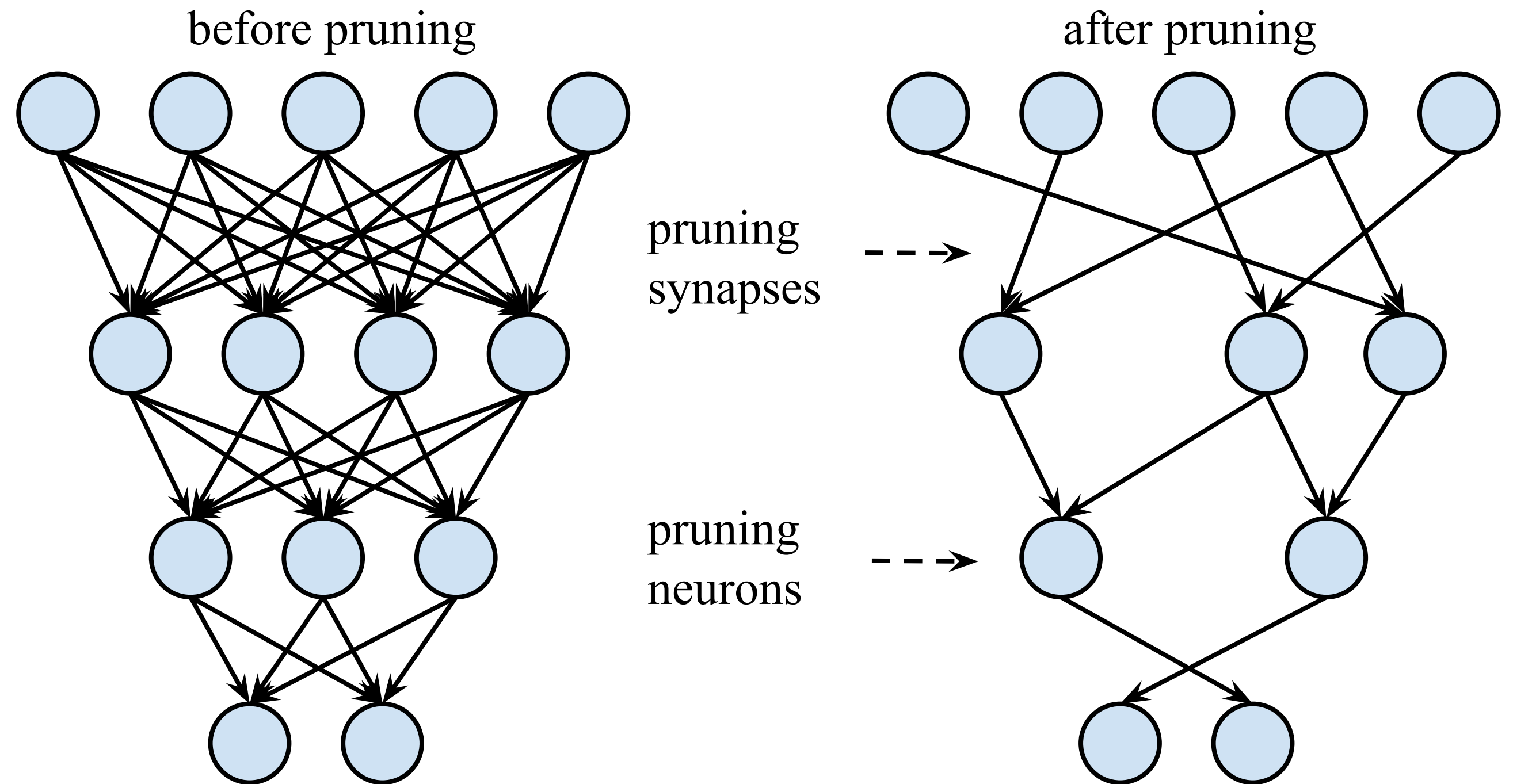
# Compression

- **Pruning:** *remove parameters that don't really contribute to performances*

  - *force parameters to be as small as possible (regularization)*

$$L_\lambda(\vec{w}) = L(\vec{w}) + \lambda||\vec{w}_1||$$

  - *Remove the small*

# Compression

● **Pruning:** *remove parameters that don't really contribute to performances*

　　● *force parameters to be as small as possible (regularization)*

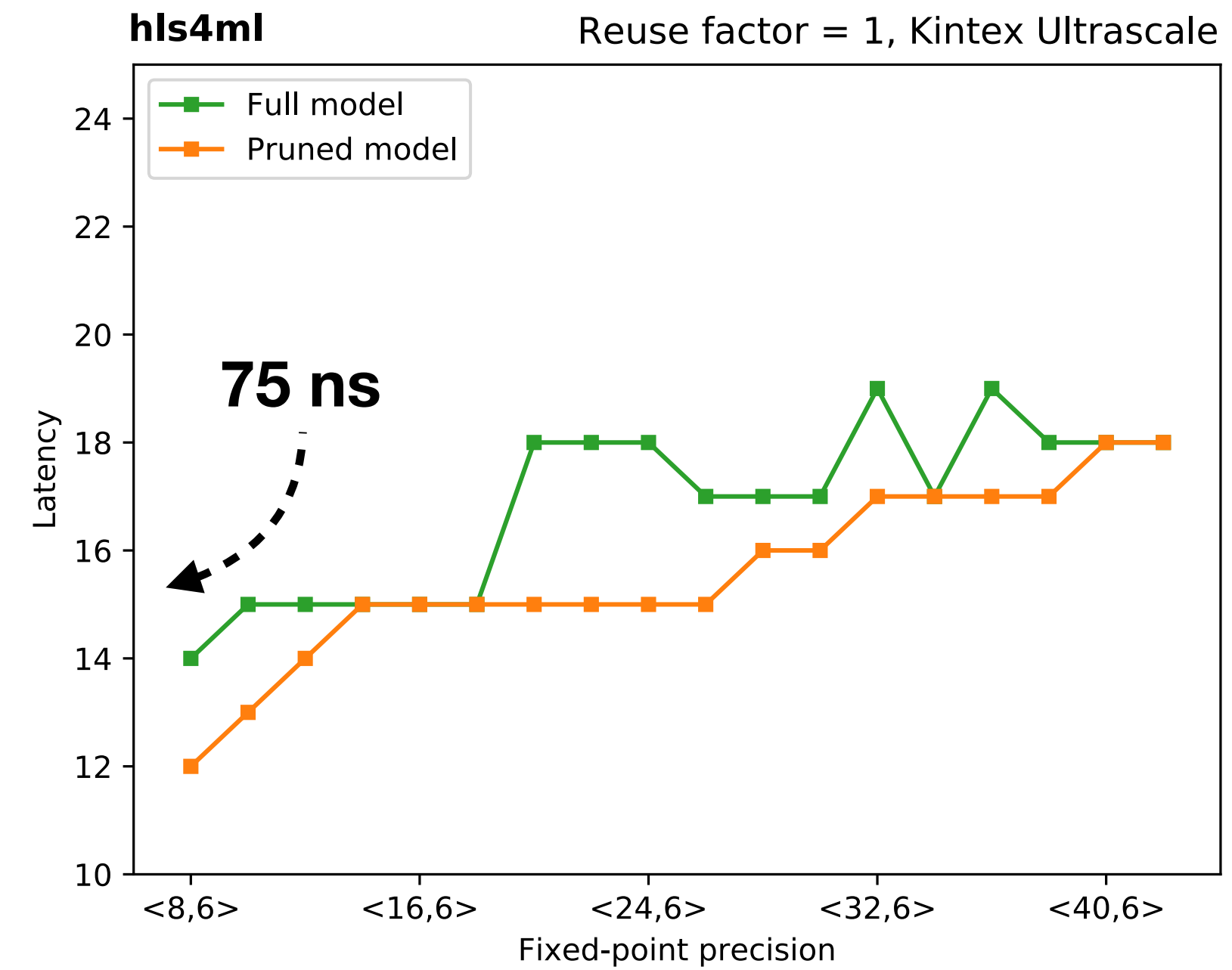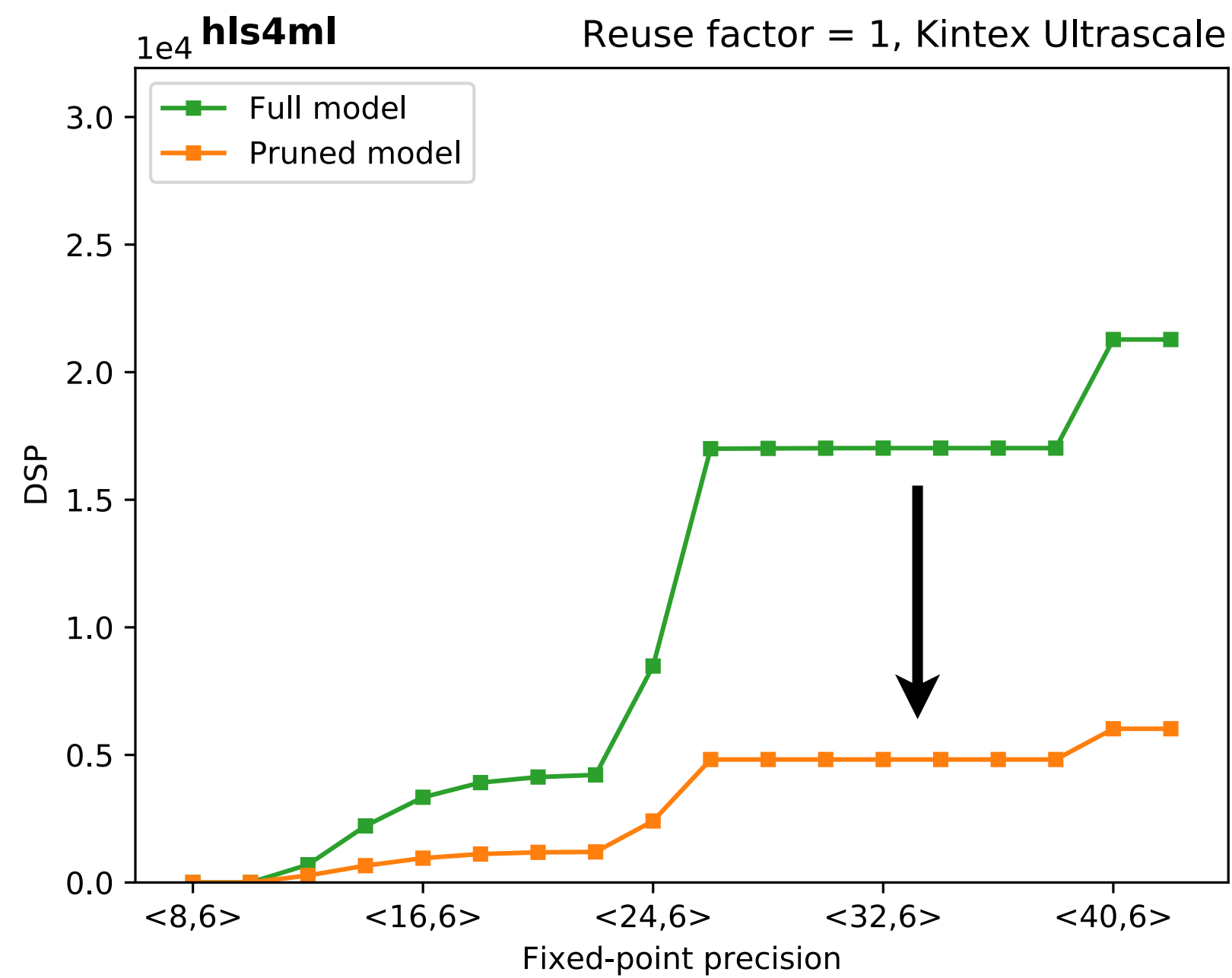$$L_\lambda(\vec{w}) = L(\vec{w}) + \lambda||\vec{w_1}||$$

　　● *Remove the small*

before pruning

after pruning

pruning synapses - - →

pruning neurons - - →

→ 70% reduction of weights
and multiplications w/o

hls4ml

| | |
|---|---|
| ■ fc3_relu | |
| ■ output_softmax | |
| ■ fc2_relu | |
| ■ fc1_relu | |

5%　　95%

700

600

500

400

Train with L$_1$

hls4ml

| | |
|---|---|
| ■ fc3_relu | |
| ■ output_softmax | |
| ■ fc2_relu | |
| ■ fc1_relu | |

5%　　32.7%　　95%

400

350

300

250

hls4ml

| | |
|---|---|
| ■ fc3_relu | |
| ■ output_softmax | |
| ■ fc2_relu | |
| ■ fc1_relu | |

32.7%　　95%

400

350

300

250

Prune

European
Research
Council

# Compression

- Big reduction in DSP usage with pruned model!

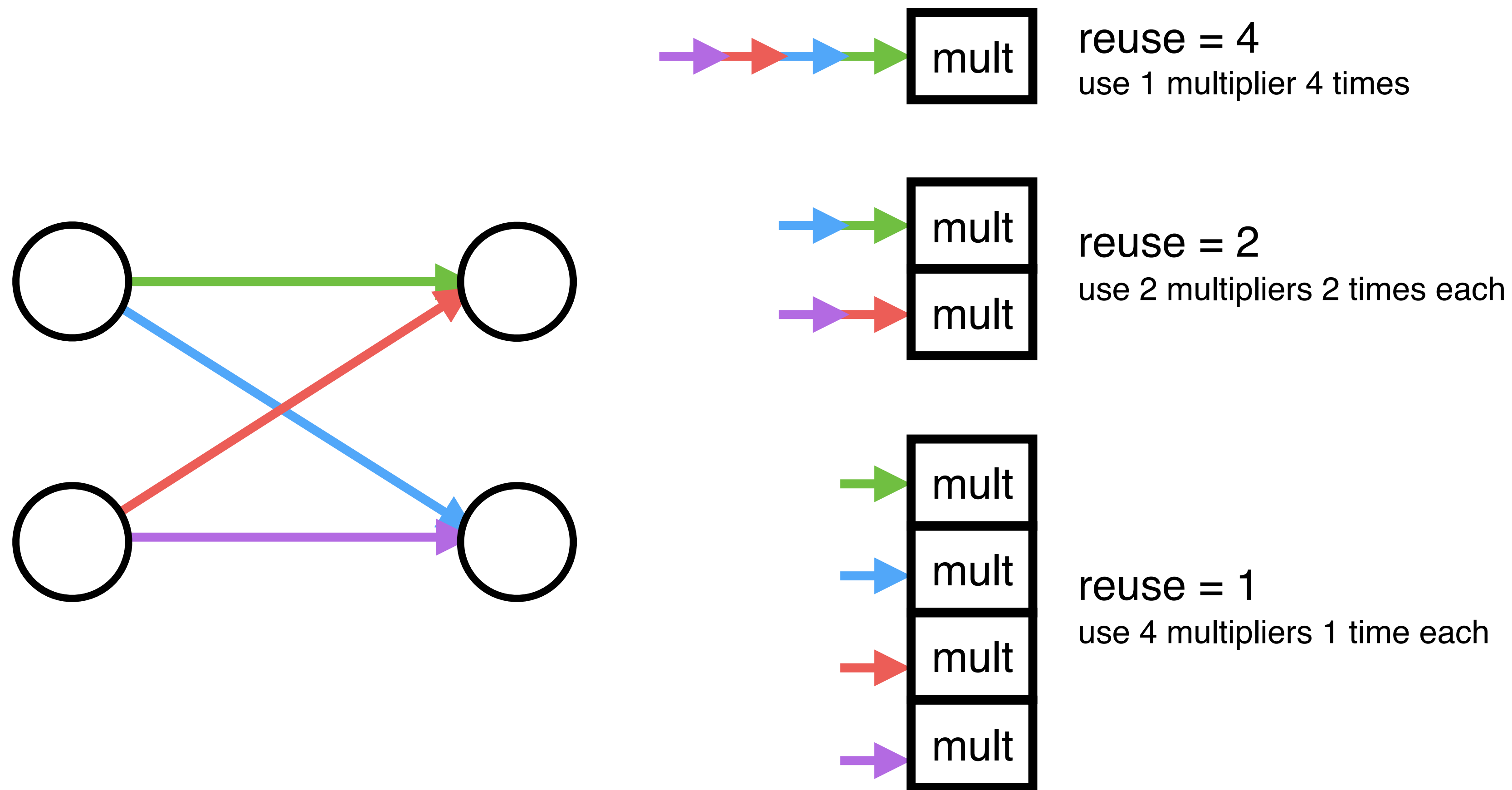- ~15 clocks @ 200 MHz = 75 ns inference

# Quantisation

- **Quantisation:** *reduce the number of bits used to represent numbers (i.e., reduce used memory)*

  - *models are usually trained at 64 or 32 bits*

  - *this is not necessarily needed in real life*

- *In our case, we could reduce to 16 bits w/o loosing precision*

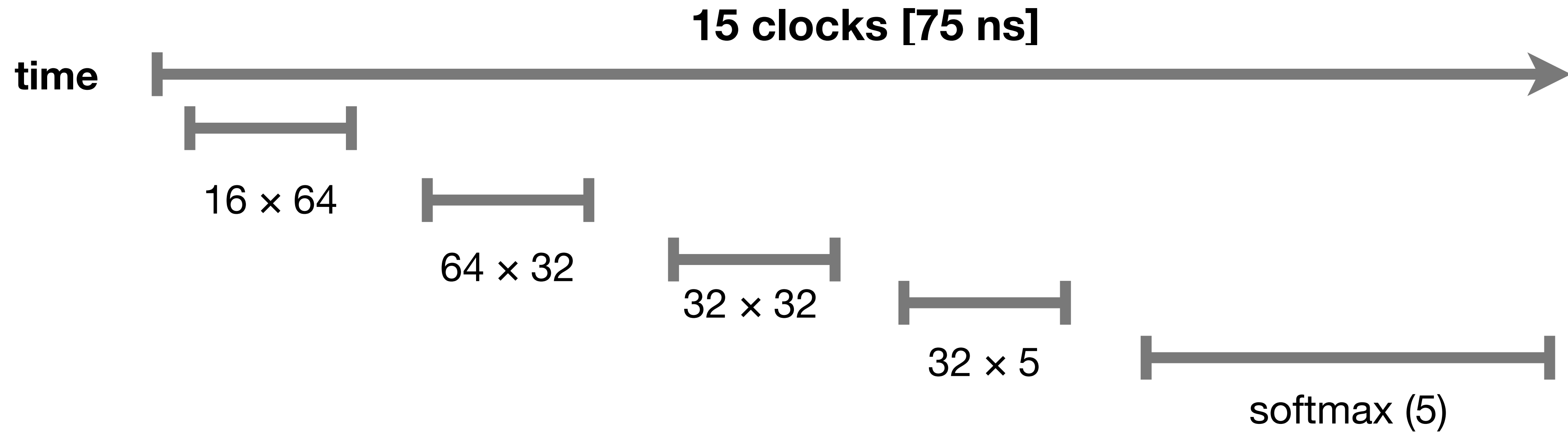- *Beyond that, one would have to accept some performance loss*

loss in performance [XXX], but this must be done with care. In Fig. 7
the absolute value of the weights after the compression described in S
overflow in the weights, at least three bits should be assigned above the b
the largest absolute value and one for the sign. The neuron values, $x_m$, a
FPGA used to compute them, require more bits, given the form of Equ
number of bits to assign *below* the binary point by scanning physics p
these bits.

**weights**

```
0101.1011101010
```
integer | fractional | width

`ap_fixed<14,4>`

**Figure 7**: Distribution of the absolute value of the weights a

**Reaches 32-bit floating point performance with 16-bit fixed point!**

In addition to saving on resources used for signal routing, reducing
and latency used for mathematical operations. For many applications t
the DSP resources of the FPGA used for multiplication. The number
depends on the precision of the numbers being multiplied and can chan
Xilinx DSP48 block [XXX] can multiply a 25-bit number with an 18-bi
to multiply a 25-bit number with a 19-bit number. Similarly, the latency
precision, though they can remain pipelined. Detailed exploration of the
is presented in Sec. 3.

hls4ml

AUC / Expected AUC

Fixed-point precision

*< total bits, integer bits >*

# Parallelisation

- ReuseFactor: how much to parallelize



**reuse = 4**
use 1 multiplier 4 times

**reuse = 2**
use 2 multipliers 2 times each

**reuse = 1**
use 4 multipliers 1 time each

**related to the Initiation Interval = when new inputs are introduced to the algo.**

Javier Duarte I hls4ml

# Parallelisation

**15 clocks [75 ns]**

time

16 × 64

64 × 32

32 × 32

32 × 5

softmax (5)

| reuse = 1 <16, 6> bits | BRAM | DSP | FF | LUT |
|---|---|---|---|---|
| Total | 13 | 954 | 53k | 36k |
| % Usage | ~0% | 17% | 3% | 5% |

# Parallelisation



**hls4ml** *preliminary*  3-layer pruned, Kintex Ultrascale

Latency vs Fixed-point precision. Reuse Factor = 1, 2, 3, 4, 5, 6.

*15-40 clock cycles (75-200 ns)*

**hls4ml** *preliminary*  3-layer pruned, Kintex Ultrascale

DSP vs Fixed-point precision. Reuse Factor = 1, 2, 3, 4, 5, 6. Max DSP.

*Foreseen architecture (FPGAs) will handle these networks*
*Inference-optimized GPUs could break the current paradigm*
*Looking forward to R&D projects with nVidia & E4 on this*

Tuning the throughput with reuse factor
will reduce the DSP usage

53

# Implementing HLS Design

◉ *HLS gives us a **conservative** estimate of the resources needed*

◉ *It actually seems to give estimates close to the back-of-the-envelope optimal estimate*

◉ *Real life much more "smooth" than emulation: no spikes observed for LUTs*

# Altera Support

- *HLS4ML support Xilinx FPGAs from the beginning*

- *Working to extend the package to work with Altera*

  - *Work in progress*

  - *Technical complications slowed us down (software licences, Quartus HLS ... on CERN, etc)*

- *First results encouraging based on emulation. To be confirmed with actual deployment on*

**Virtex V vs. Stratix V
(based on simulation)**

# Deep Learning on the Cloud

# Inference on the cloud

- *In the (near) future, DAQ/HLT farms will be based on heterogenous computing*

  - *CPU+GPU / CPU+FPGA*

  - *Mainly to accelerate slow algorithms (e.g., tracking) through parallelisation*

  - *Also useful for ML inference*

- *R&D on heterogeneous environments on commercial clouds*

  - *provides easy-to-use CPU+FPGA (or GPU) ecosystem*

  - *allows further R&D: inference on demand from the CPU-based HLT farm to the FPGAs/GPUs on the cloud*

# Microsoft Brainwave

## A Reconfigurable Fabric for Accelerating Large-Scale Datacenter Services

Andrew Putnam     Adrian M. Caulfield     Eric S. Chung     Derek Chiou[1]
Kypros Constantinides[2]  John Demme[3]  Hadi Esmaeilzadeh[4]  Jeremy Fowers
Gopi Prashanth Gopal     Jan Gray     Michael Haselman     Scott Hauck[5]  Stephen Heil
Amir Hormati[6]  Joo-Young Kim     Sitaram Lanka     James Larus[7]  Eric Peterson
Simon Pope     Aaron Smith     Jason Thong     Phillip Yi Xiao     Doug Burger

https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/Catapult_ISCA_2014.pdf

**Heterogeneous Edge Resource**

gRPC protocol

Trigger

CPU ≡ FPGA

- Cloud service has latency

- Run CMSSW on Azure cloud machine → simulate local installation of FPGAs ("on-prem" or "edge")

- Provides test of "HLT-like" performance

2xCPU

European Research Council

# Brainwave at scale

- Provides a full service at scale
  (more than just a single co-processor)

- Multi-FPGA/CPU fabric, accelerating
  both compute and network

- ➤ Demonstrated large improvements in
    processing time for Bing searches

- Caveat: only selected pre-trained
  DNN models currently available

# Pros & Cons

◉ *Commercial clouds focus on what is sellable*

  ◉ *supports computing-vision off-the-shelf networks (ResNet50, ResNet512, DenseNet121, VGGNet16*

  ◉ *(for now) reduced flexibility: doesn't allow customised architectures. Input has to be an image*

  ◉ *(longer term) more architectures will become available*

◉ *As long as one of these networks is good for the problem at hand, implementation is optimized (beyond what HLS might do)*



Zoom ×10

DeepAK8

arXiv:1605.07678

- ResNet50: 25M parameters, 7B operations
- Examples of large networks used in CMS:
  - DeepAK8, 500K parameters, 15M operations
  - DeepDoubleB, 40K parameters, 700K operations

# SONIC

- SONIC (a **S**ervices for **O**ptimized **N**etwork **I**nference on **C**oprocessors) is a framework to exploit cloud resources for on-demand inference

- CPU runs "locally" (for us at FNAL) and sends data to the cloud system

- FPGAs there set to run our inference problems

- answer communicated back via gRPC protocol (driven by Microsoft infrastructure boundaries)



**CPU comparison:**

- Intel i7 3.6 GHz (8 core, TF v1.10) ~ 180 ms
- Intel i7 3.6 GHz (1 core, TF v1.10) ~ 500 ms
- Intel i7 3.6 GHz (1 core, TF v1.06) ~ 1.2 s
- Intel Xeon 2.6 GHz (1 core, TF v1.06) ~ 1.75 s [what we are running on cmslpc]

# Testing SONIC



- Good performance in initial tests
  - ○ "remote": cmslpc @ FNAL to Azure (VA),   ‹time› = 56 ms
  - ○ "onprem": run CMSSW on Azure VM,      ‹time› = 10 ms
    (~2 ms on FPGA, rest is classifying and I/O)
  - ○ CPU (cmslpc): 1.75 sec
    (6 min to load ResNet50 session)
- ➤ More than order of magnitude improvement!

# Using GPUs instead



- Benchmark: Nvidia GTX 1080 (GPU), Intel i7 3.6 GHz (CPU)
- All tests use .pb file with Brainwave version of ResNet50
- Using classic ResNet50 implementation w/ CuDNN: faster on GPU by 5–10×

**Brainwave w/ SONIC**

- Featurizer: 1.8 ms (FPGA)

- Classifier: 2 ms (CPU)

- Infrastructure: 6 ms (trying to improve)

- Transit time: 10 ms (speed of light, Chicago to Virginia) + network switching, etc.

- Total: ~10 ms (onprem), ~56 ms (remote)

**CPU**
- 1.75 s (cmslpc CPU)
- 500 ms (new CPU, TF version)

**GPU**
- 100 ms (batch 1)
- 15 ms (batch 10)

# Jet Tagging with ResNet on Cloud

- Use feature set generated by ResNet50, train new fully-connected classifier

  o Brainwave accelerates training
    (heaviest component is evaluating ResNet50 to produce feature set)

- CNNs have been used in jet image classification: arXiv:1709.04464

- Proposed "realistic" test:

  <span style="color:red">Preliminary result w/ small dataset</span>

  o Compute jet discriminator values
    (q, g, W, Z, t)

  o Run module in miniAOD sequence

- Can also be used for other experiments

  o e.g. NOvA: identify neutrino events

### Confusion matrix

| True label | quark | gluon | W | Z | top |
|---|---|---|---|---|---|
| quark | 0.474 | 0.224 | 0.149 | 0.067 | 0.086 |
| gluon | 0.152 | 0.604 | 0.166 | 0.047 | 0.031 |
| W | 0.054 | 0.101 | 0.735 | 0.108 | 0.002 |
| Z | 0.054 | 0.106 | 0.496 | 0.336 | 0.007 |
| top | 0.091 | 0.017 | 0.125 | 0.174 | 0.592 |

European Research Council

# Outlook

◉ *HLS4ML aims to be a flexible tool to implement your home-made NN in a trigger/DAQ system where low latency matters*

◉ *Now works with TensorFlow and PyTorch for Dense Neural networks*

  ◉ *Working to support ONNX format*

◉ *Working on new architecture support*

  ◉ *Boosted Decision Trees*
  ◉ *Convolutional NNs (1D & 2D)*
  ◉ *Recurrent NNs (GRUs, LSTMs, etc)*
  ◉ *Graph Networks*

◉ *Extra functionalities added*

  ◉ *New activation functions*
  ◉ *Batch Normalization*
  ◉ *Layer concatenate*
  ◉ *Max Pooling*
  ◉ *…*

# Backup

# Detector Monitoring

# Data Quality Monitoring

- *When taking data, >1 person watches for anomalies in the detector 24/7*

- *At this stage no global processing of the event*

- *Instead, local information from detector components available (e.g., detector occupancy in a certain time window)*



A

Raw Occupancy (Run: 272011, W: 1.0, St: 1.0, Sec: 6.0)

B

Raw Occupancy (Run: 273158, W: 0.0, St: 2.0, Sec: 12.0)

C

Raw Occupancy (Run: 275310, W: 1.0, St: 2.0, Sec: 7.0)

# Two approaches

- *Given the nature of these data, ConvNN are a natural analysis tool. Two approaches pursued*

  - *Classify good vs bad data. Works if failure mode is known*

  - *Use autoencoders to assess data "typicality". Generalises to unknown failure modes*

**A. Pol et al., to appear soon**

# Two approaches

● *Given the nature of these data, ConvNN are a natural analysis tool. Two approaches pursued*

  ● *Classify good vs bad data. Works if failure mode is known*

  ● *Use autoencoders to assess data "typicality". Generalises to unknown failure modes*

  **A. Pol et al., to appear soon**

# Data Quality Certification

◉ *Autoencoder-based 1-class approach generalises to later stages of quality assessment*

  ◉ *after reconstruction of the events, event reconstruction allows a global assessment (w.g., looking at electrons, muons, etc rather than hits in the detector)*

  ◉ *A global autoencoder can spot all these features*

  ◉ *Monitoring individual contributions to loss function (e.g., MSE) one can track the problem back to a specific physics object/detector component*

**F. Široký  et al., to appear sooner or later**



Reconstruction error for different classes

# HL4ML: FPGA details

## Xilinx Vivado 2017.2

Results are slightly different in other versions of Vivado

e.g. 2016.4 optimization is less performant for Xilinx ultrascale FPGAs

## Clock frequency: 200 MHz

Latency results can vary (~10%) with different clock choices

## FPGA: Xilinx Kintex Ultrascale (XCKU115-FLVB2104)

Results are slightly different in other FPGAs

e.g. Virtex-7 FPGAs are slightly differently optimized

# Why Deep Learning

◉ *Neural network can model non linear functions*

  ◉ *the more complex is the network, the more functions it can approximate*

◉ *Neural network are faster to evaluate (inference) than typical reco algorithm.*

  ◉ *This is the speed up we need*

◉ *Neural Networks (unlike other kind of ML algorithms) are very good with raw (non-preprocessed) data (the recorded hits in the event)*

$$(pT, \eta, \phi, E)_{OFFLINE} = f( (pT, \eta, \phi, E)_{ONLINE} )$$

  ◉ *could use them directly on the detector inputs*

$$(pT, \eta, \phi, E)_{OFFLINE} = g( \text{ Event hits } )$$

**One would have to learn $f$ and $g$ to evaluate them at trigger. Online processing is replaced by offline training**

# Beyond the toy-model

- ◉ *Approach works in principle*

    - ◉ *Can identity easily 2 of the 3 models*

    - ◉ *With enough statistics, could see the third*

- ◉ *Might not work in absolute*

    - ◉ *encoder based on physics motivate quantities which are not model-agnostic*

    - ◉ *Use deep:learning: train on raw data directly.* **To be done next**
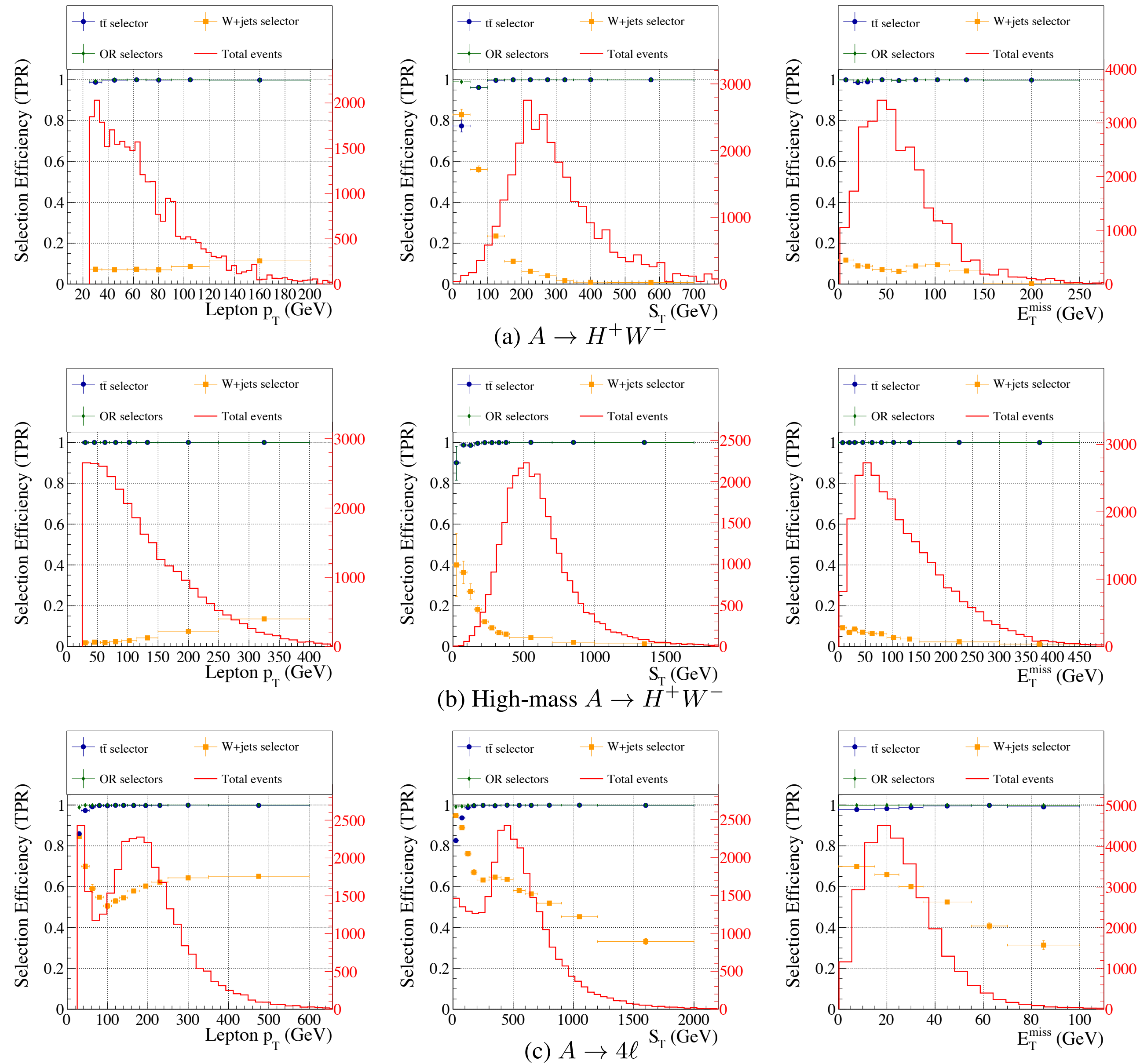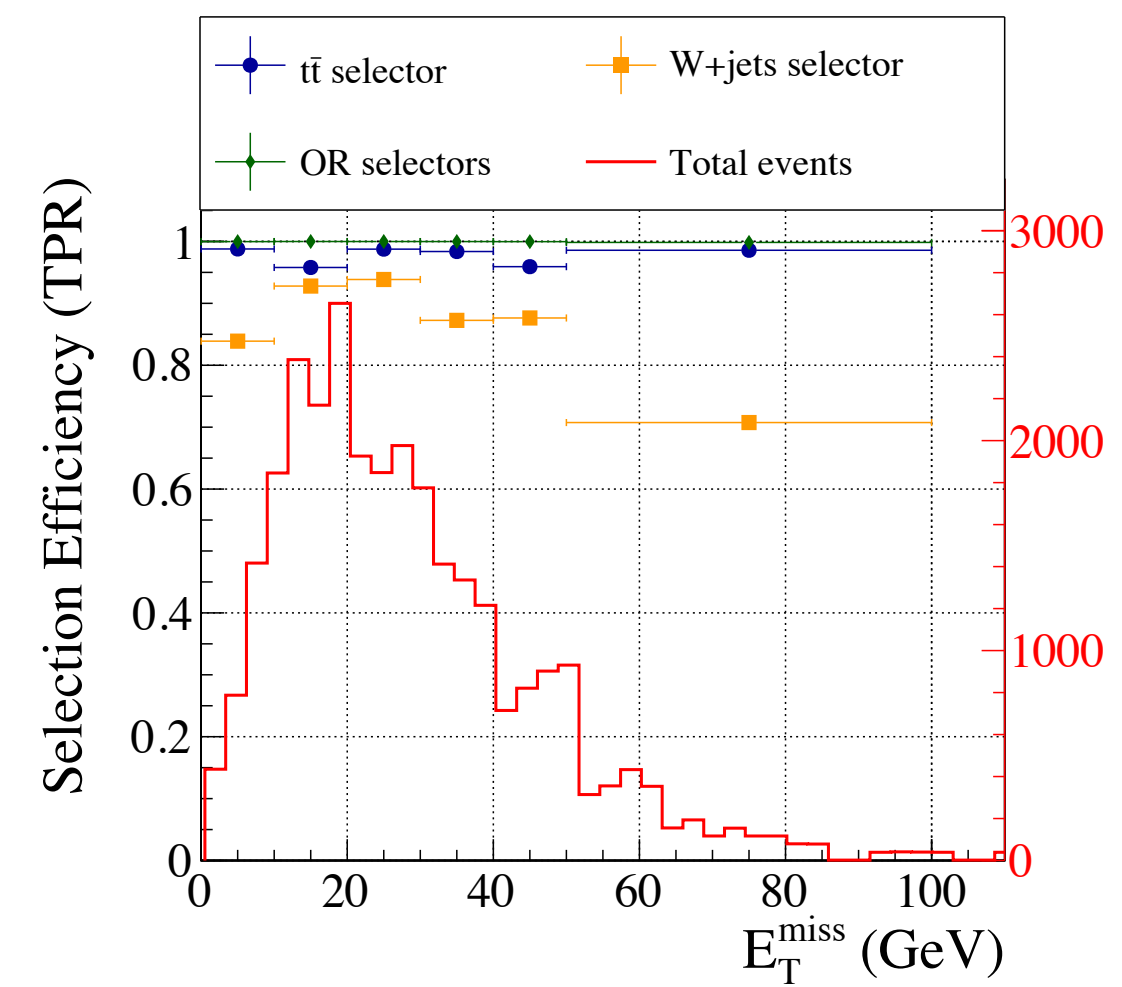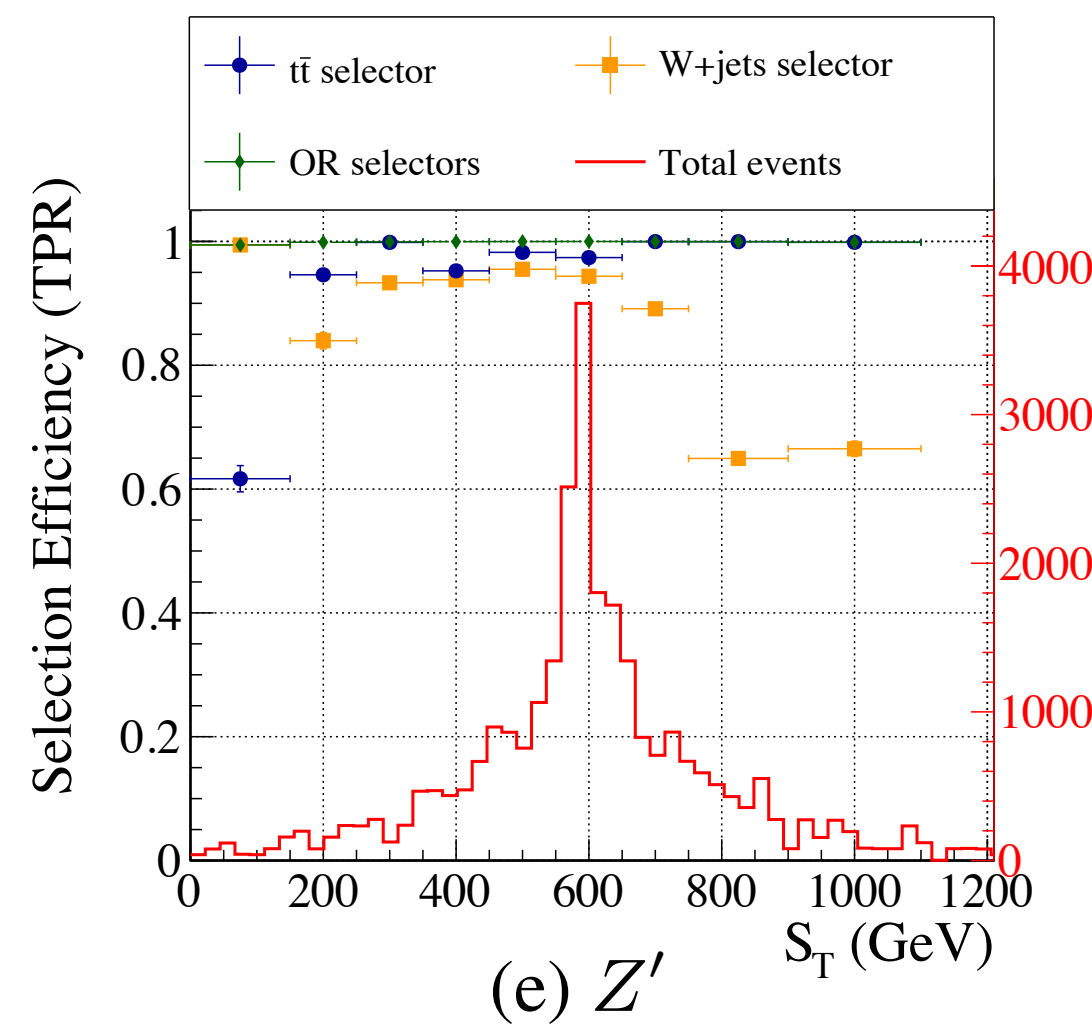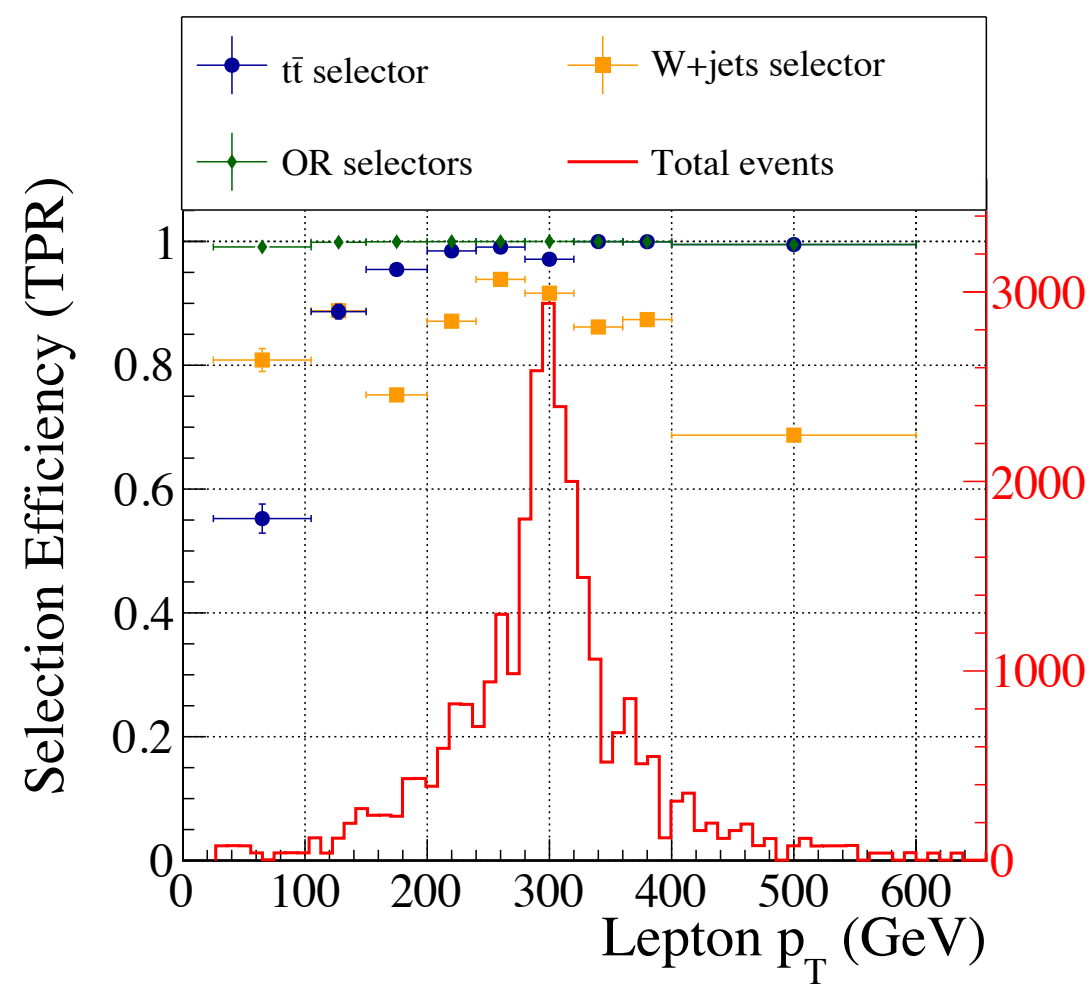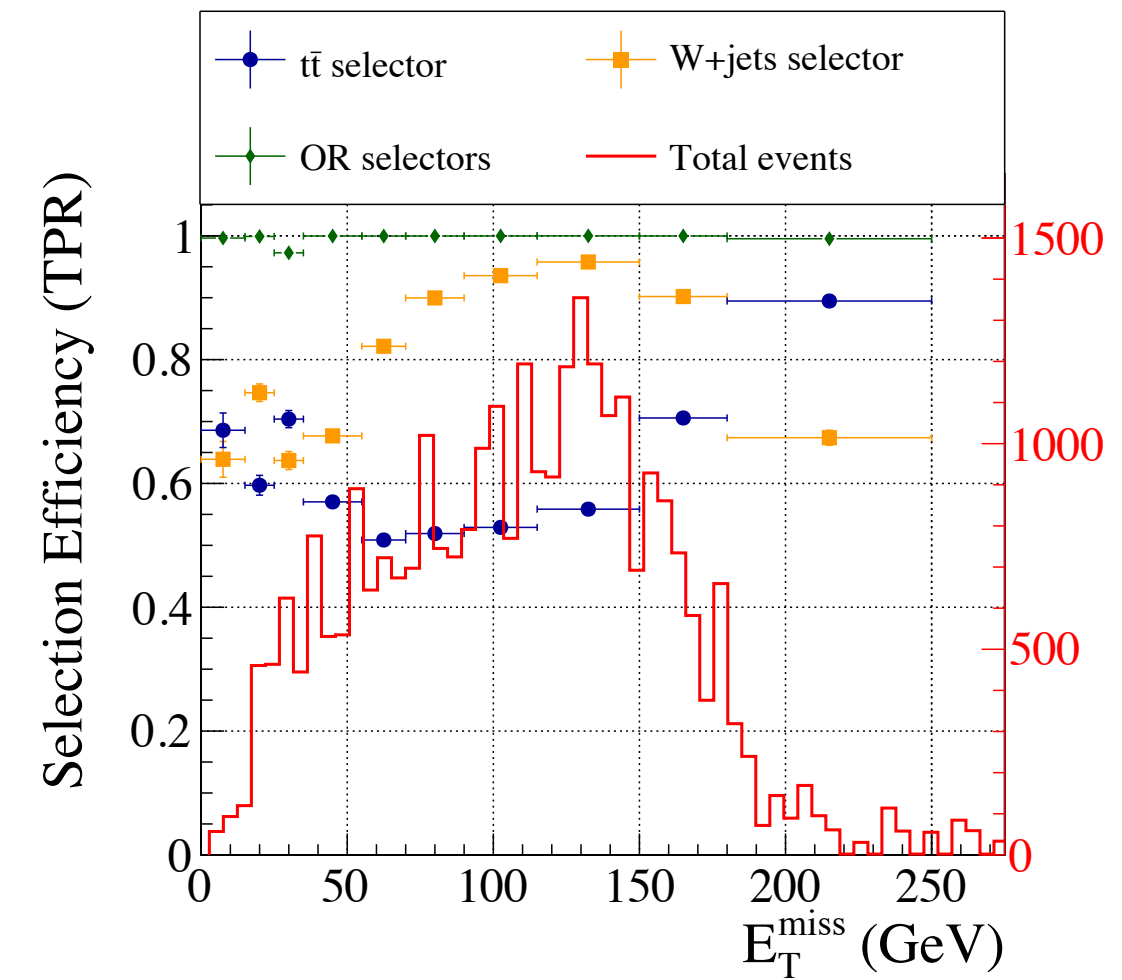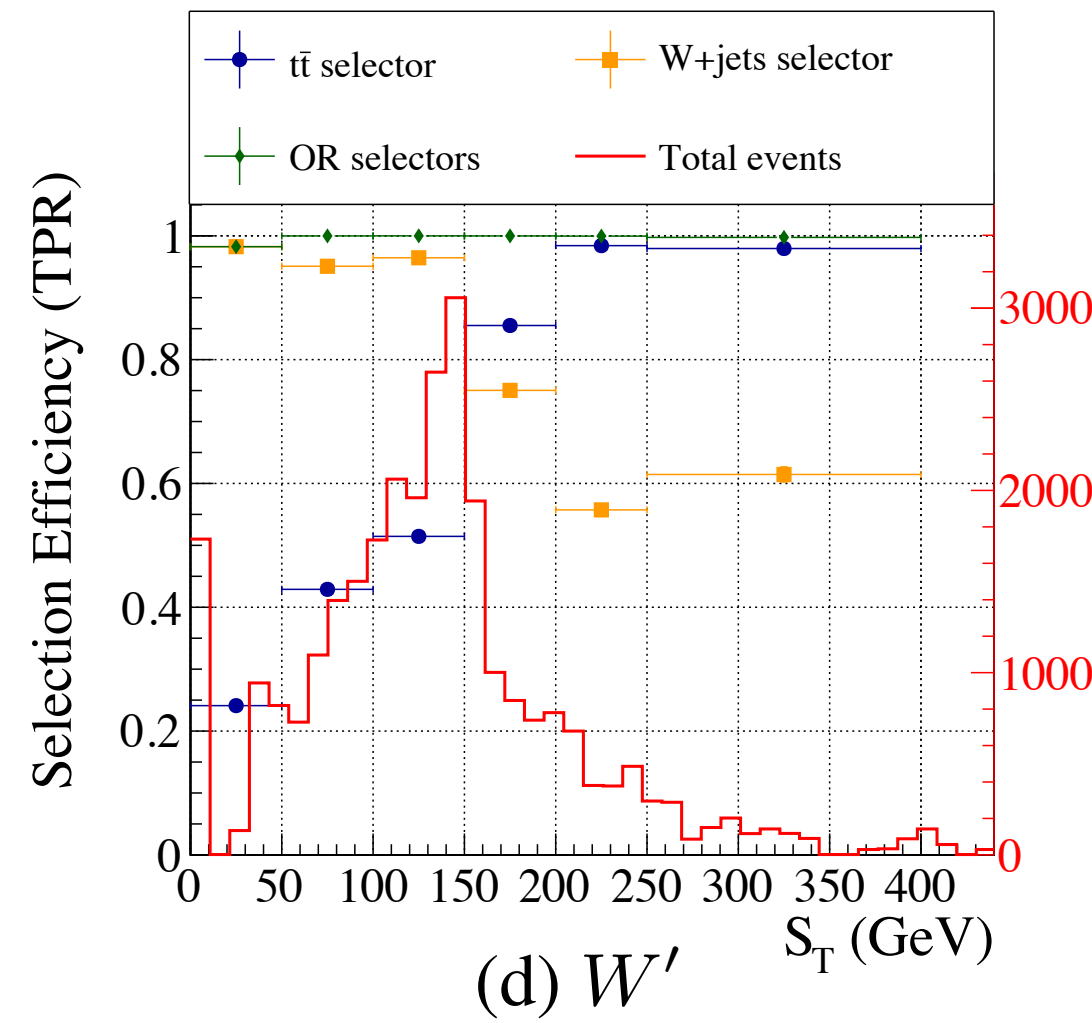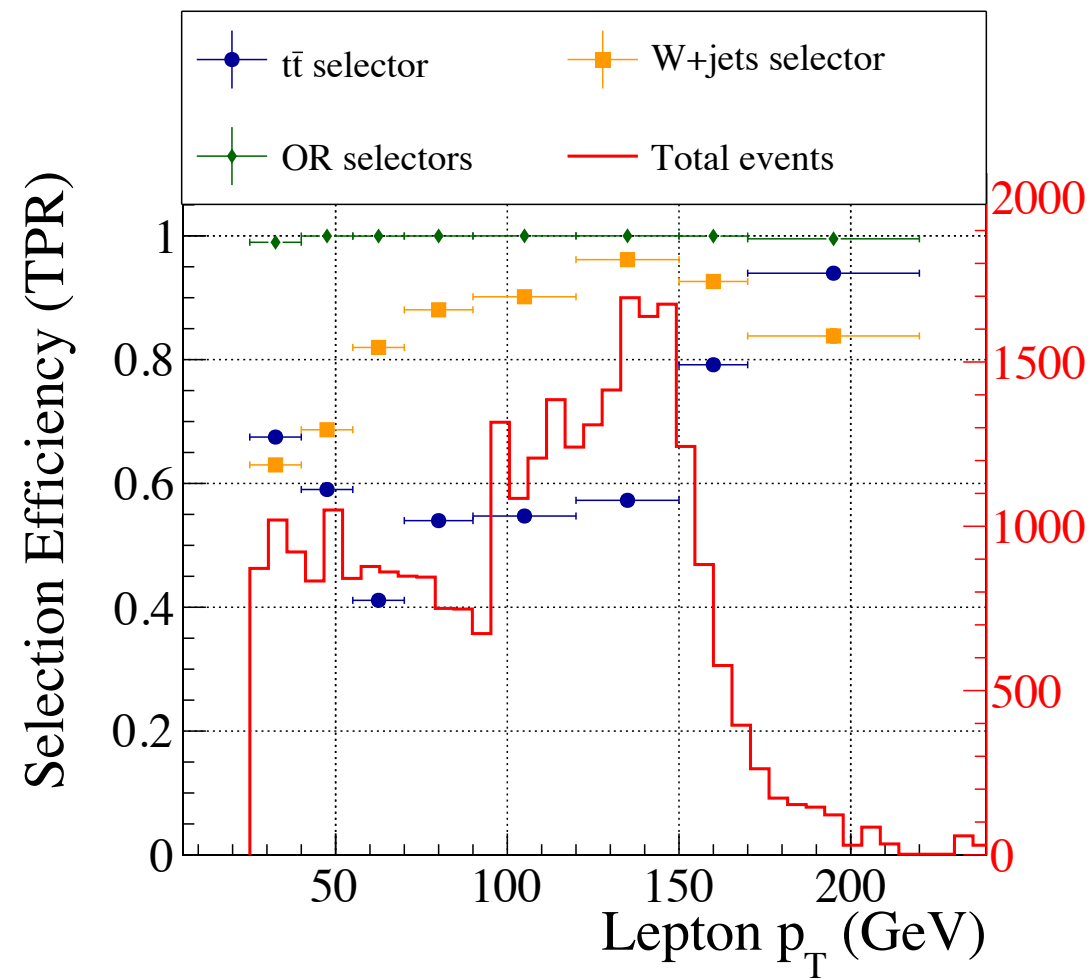
# Kinematic Bias?

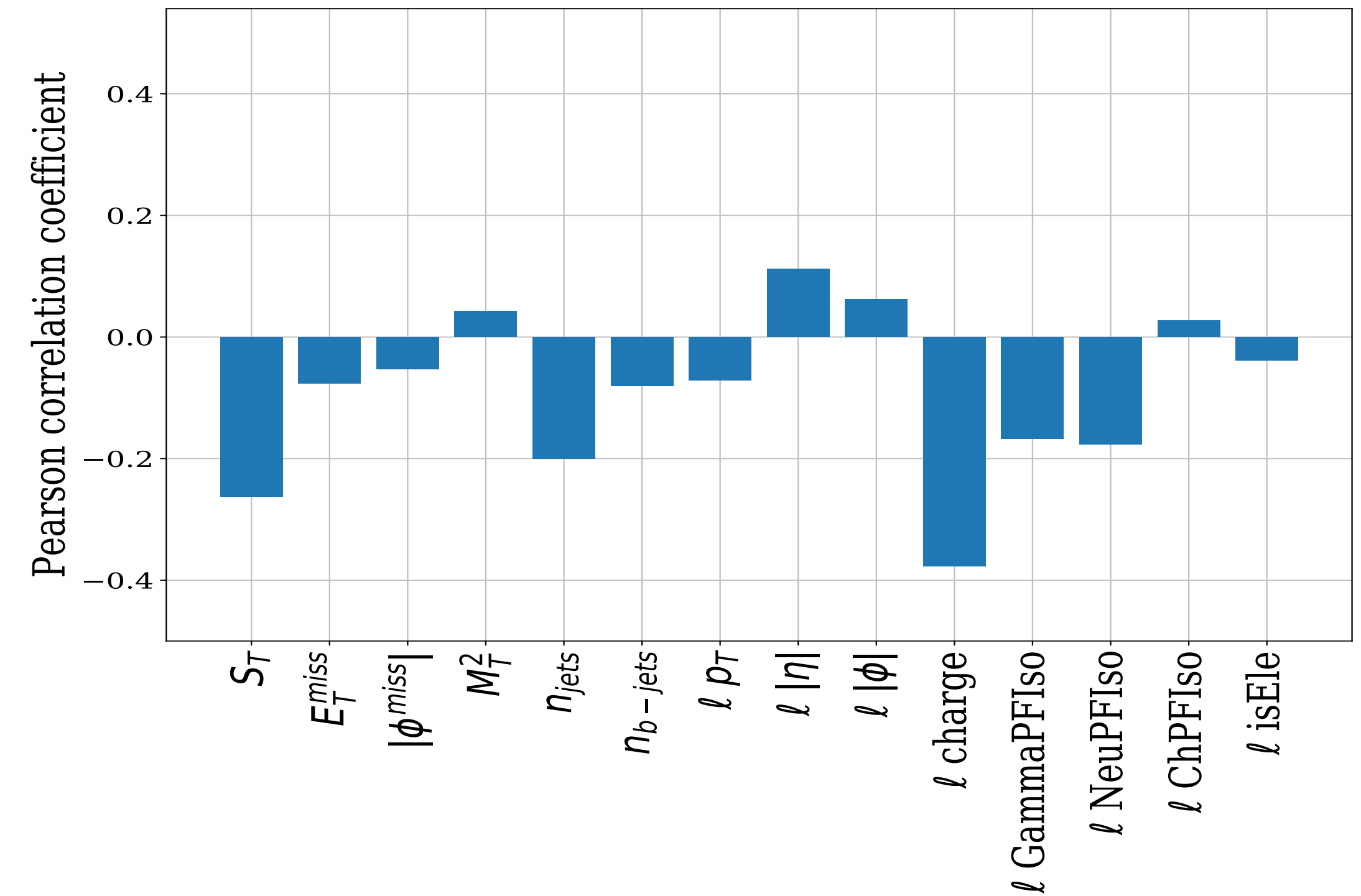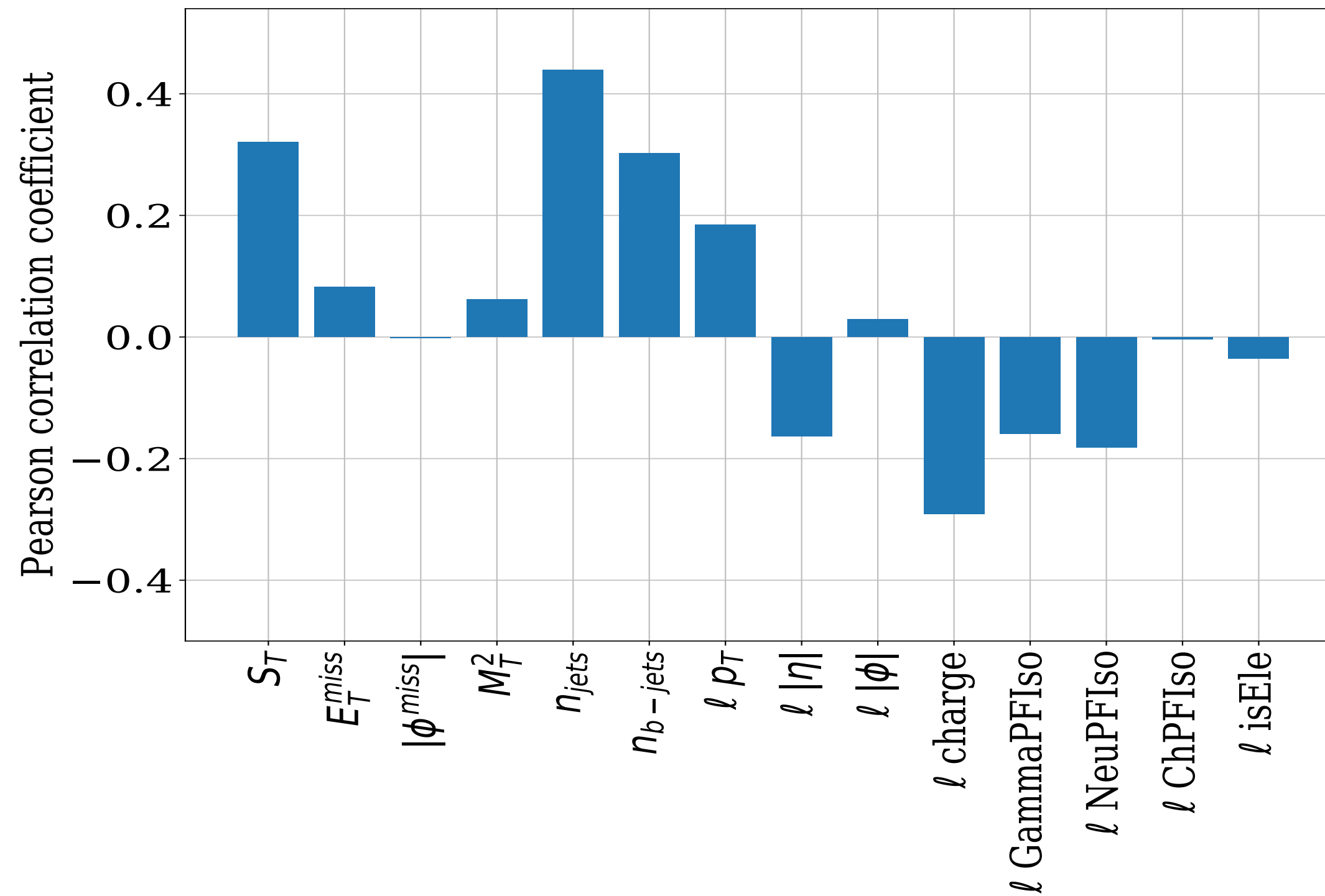- *With 99% signal efficiency, bias on kinematic variables within the uncertainty of a trigger-efficiency measurement*

(a) $A \to H^+ W^-$

(b) High-mass $A \to H^+ W^-$

(c) $A \to 4\ell$

(d) $W'$

(e) $Z'$

# Selection performances



**The network is learning some physics…**
- tt events are more crowded that W events
- leptons in W and tt events are isolated from other particles