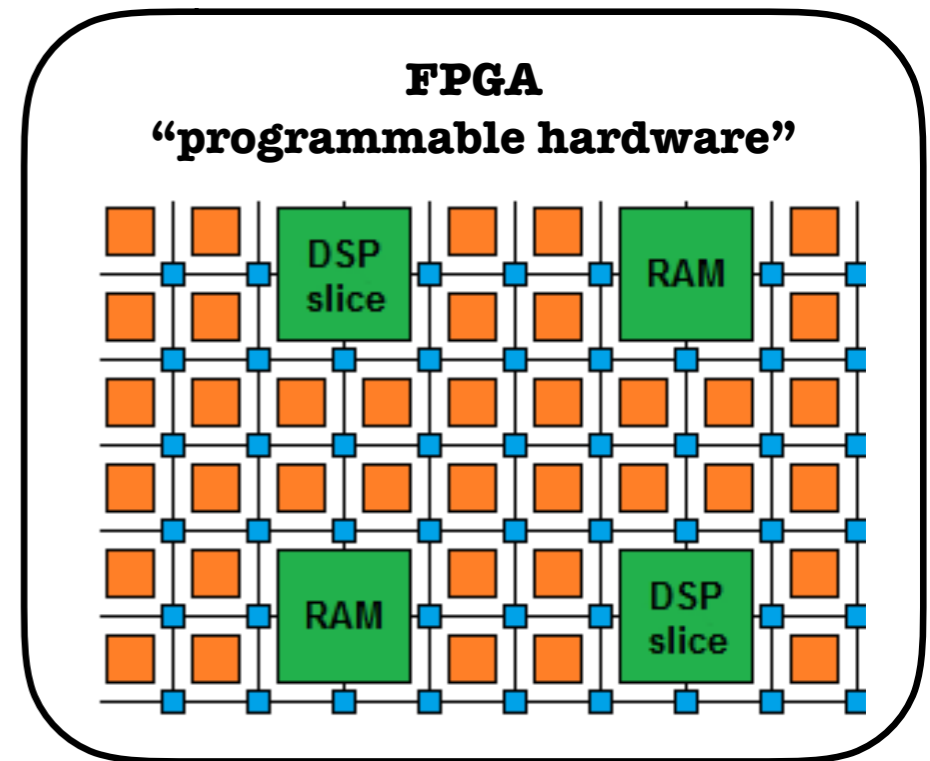JAVIER DUARTE (UCSD)
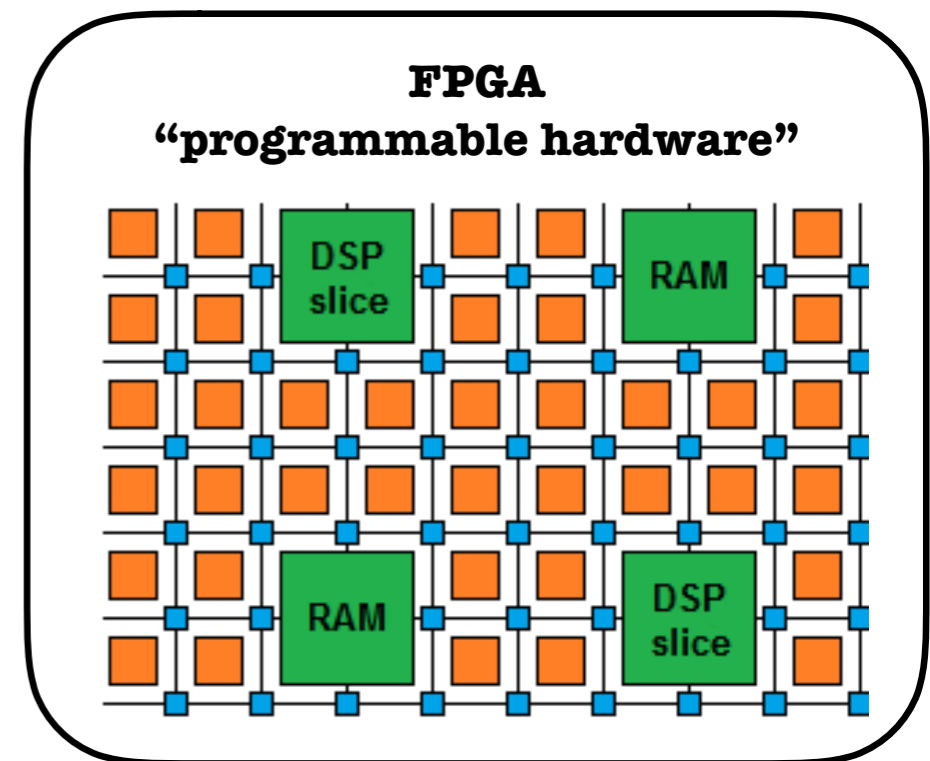OCTOBER 23, 2019
WEST COAST LHC JAMBOREE, SLAC
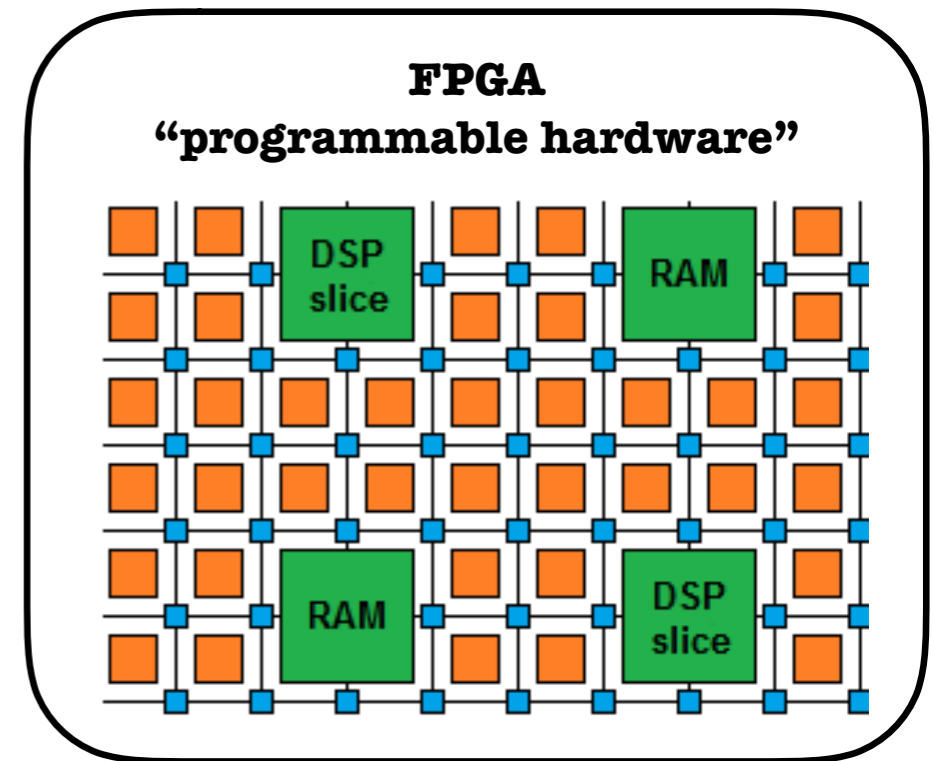
# NEW IDEAS & TECHNOLOGIES IN TRIGGER FOR HL–LHC (AND RUN 3)

**FPGA**
**"programmable hardware"**

‣ **Modern FPGAs** with large amounts of embedded components that perform multiplication (DSPs), apply logical functions (LUTs), or store memory (BRAM)



FPGA
"programmable hardware"

FPGA
"programmable hardware"

▸ **Modern FPGAs** with large amounts of embedded components that perform multiplication (DSPs), apply logical functions (LUTs), or store memory (BRAM)

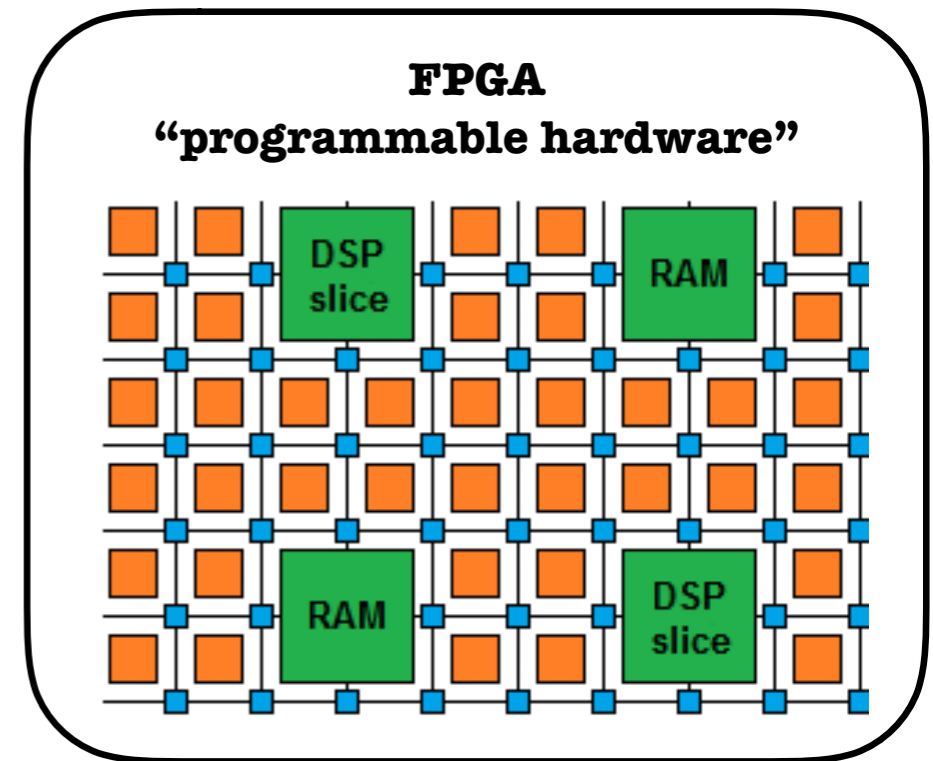▸ **High level synthesis** to more easily program FPGAs

FPGA
"programmable hardware"

▸ **Modern FPGAs** with large amounts of embedded components that perform multiplication (DSPs), apply logical functions (LUTs), or store memory (BRAM)
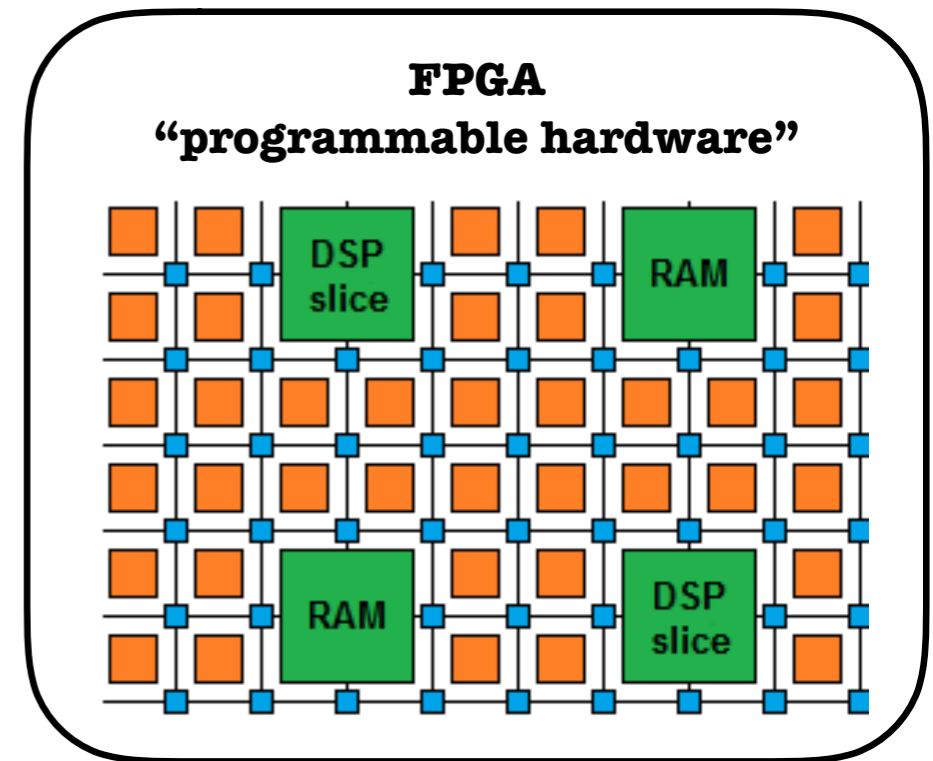
▸ **High level synthesis** to more easily program FPGAs

▸ **Sophisticated algorithms**

FPGA
"programmable hardware"

▸ **Modern FPGAs** with large amounts of embedded components that perform multiplication (DSPs), apply logical functions (LUTs), or store memory (BRAM)
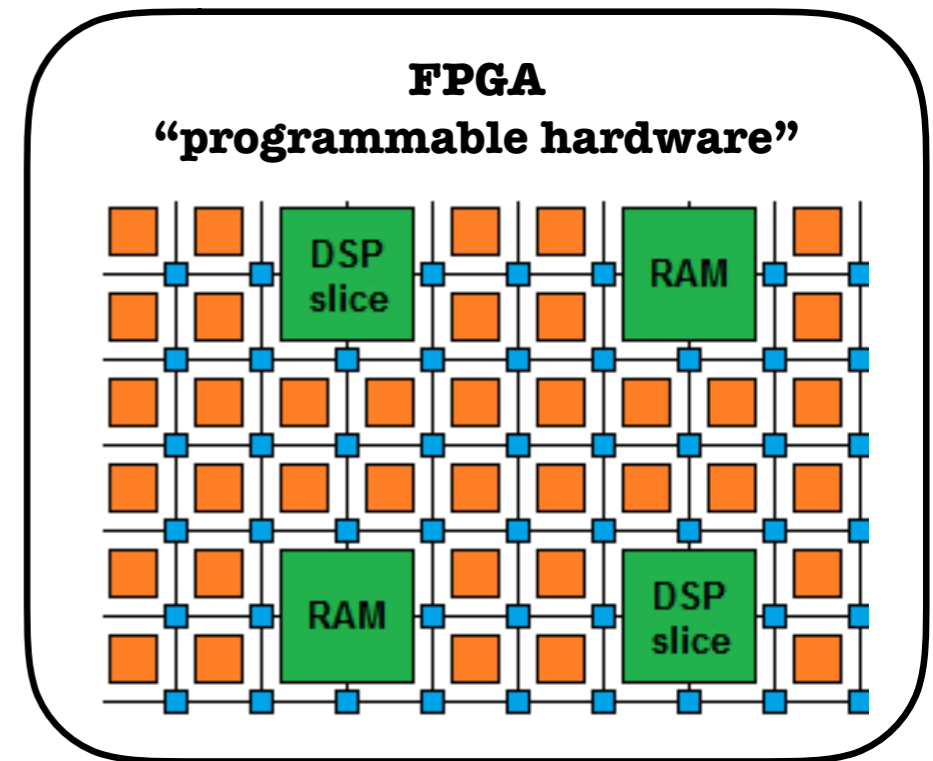
▸ **High level synthesis** to more easily program FPGAs

▸ **Sophisticated algorithms**

▸ **Machine learning**

FPGA
"programmable hardware"



▸ **Modern FPGAs** with large amounts of embedded components that perform multiplication (DSPs), apply logical functions (LUTs), or store memory (BRAM)
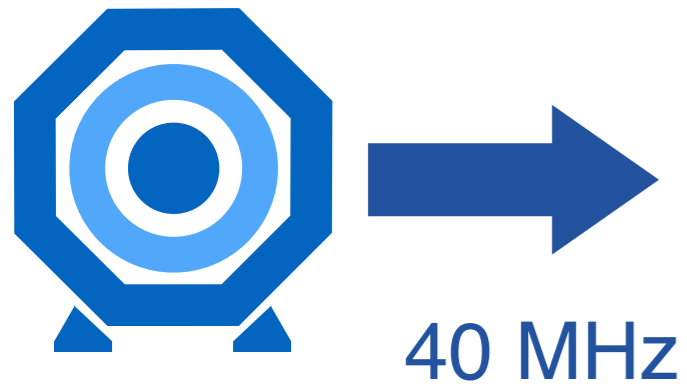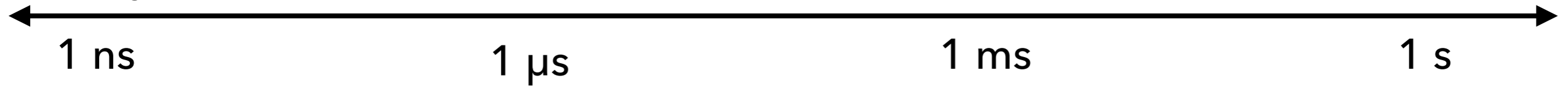
▸ **High level synthesis** to more easily program FPGAs

▸ **Sophisticated algorithms**

▸ **Machine learning**

▸ GPUs or FPGAs or ASICs as **co-processors** for software trigger

# LHC EVENT PROCESSING

Compute
Latency

1 ns             1 μs             1 ms             1 s
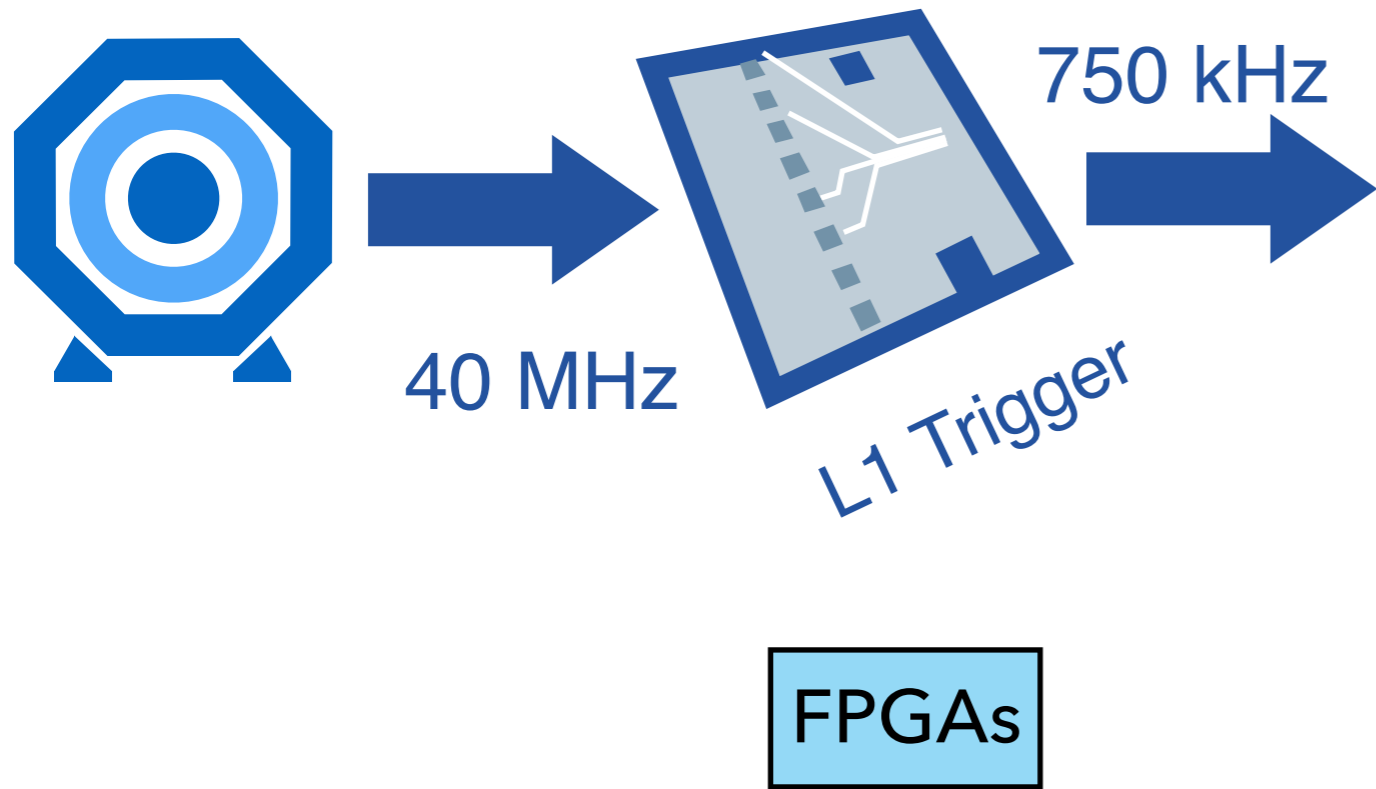
40 MHz

**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



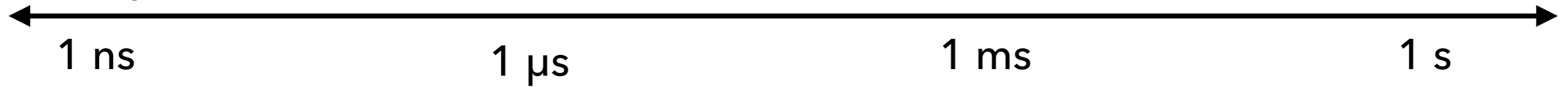1 ns  1 μs  1 ms  1 s

40 MHz  L1 Trigger  750 kHz

FPGAs

**Challenges:**

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

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

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

Compute Latency
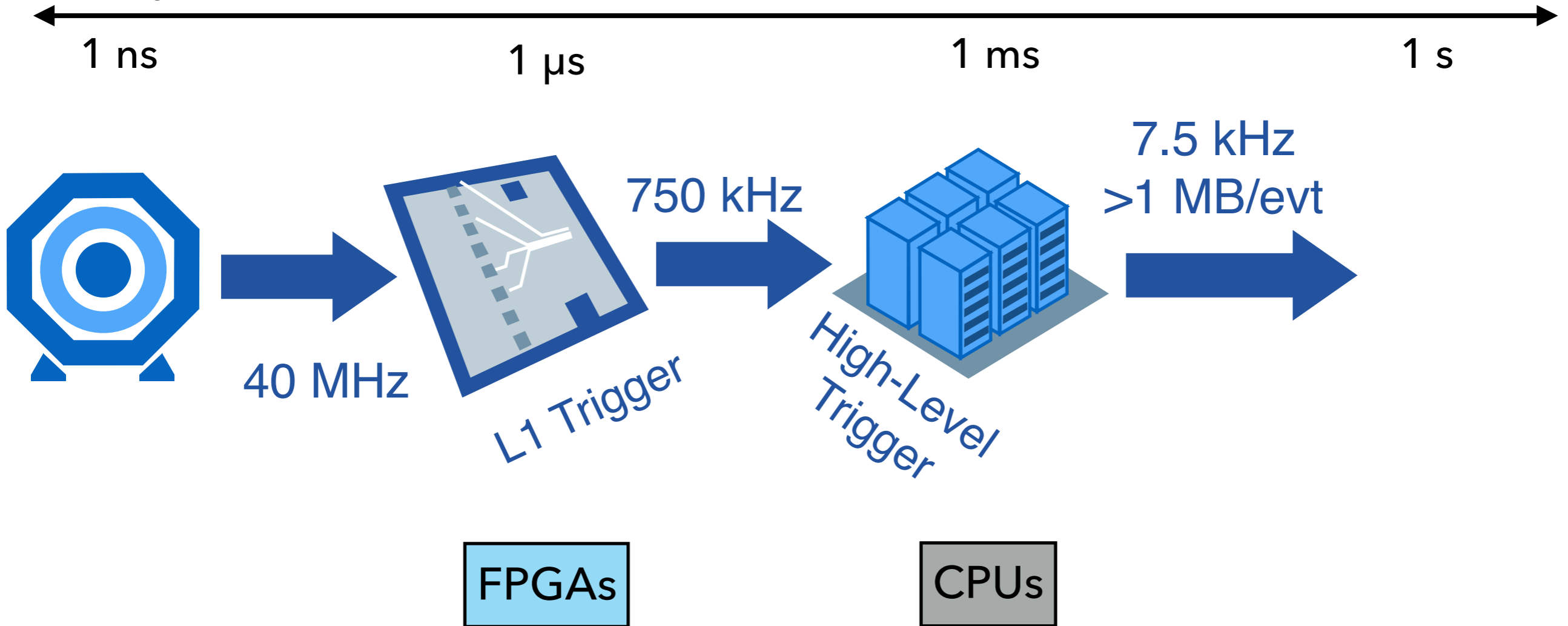
1 ns — 1 μs — 1 ms — 1 s

40 MHz

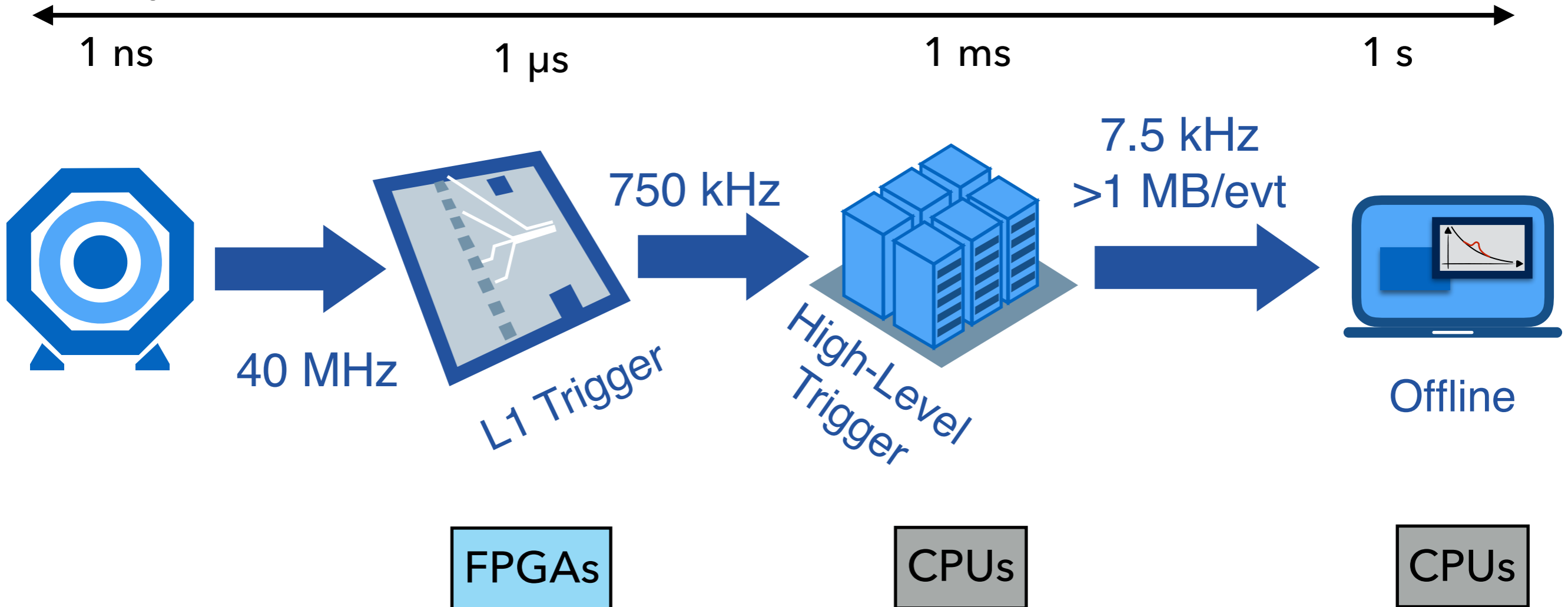L1 Trigger — 750 kHz — High-Level Trigger — 7.5 kHz >1 MB/evt

FPGAs

CPUs

ML

**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

1 ns          1 μs          1 ms          1 s

40 MHz

L1 Trigger

750 kHz

High-Level
Trigger

7.5 kHz
>1 MB/evt

Offline

FPGAs          CPUs          CPUs

ML          ML

**Challenges:**
Each collision produces $O(10^3)$ particles
The detectors have $O(10^8)$ sensors
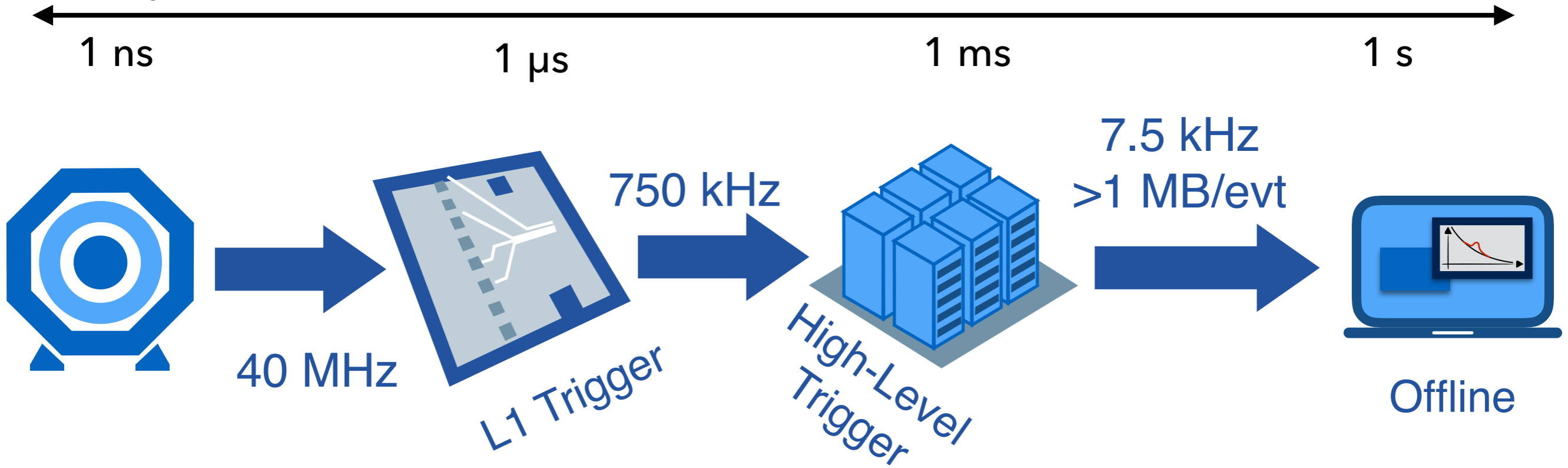Extreme data rates of $O(100\ TB/s)$

Compute Latency

1 ns · 1 μs · 1 ms · 1 s

40 MHz · 750 kHz · 7.5 kHz >1 MB/evt

L1 Trigger · High-Level Trigger · Offline
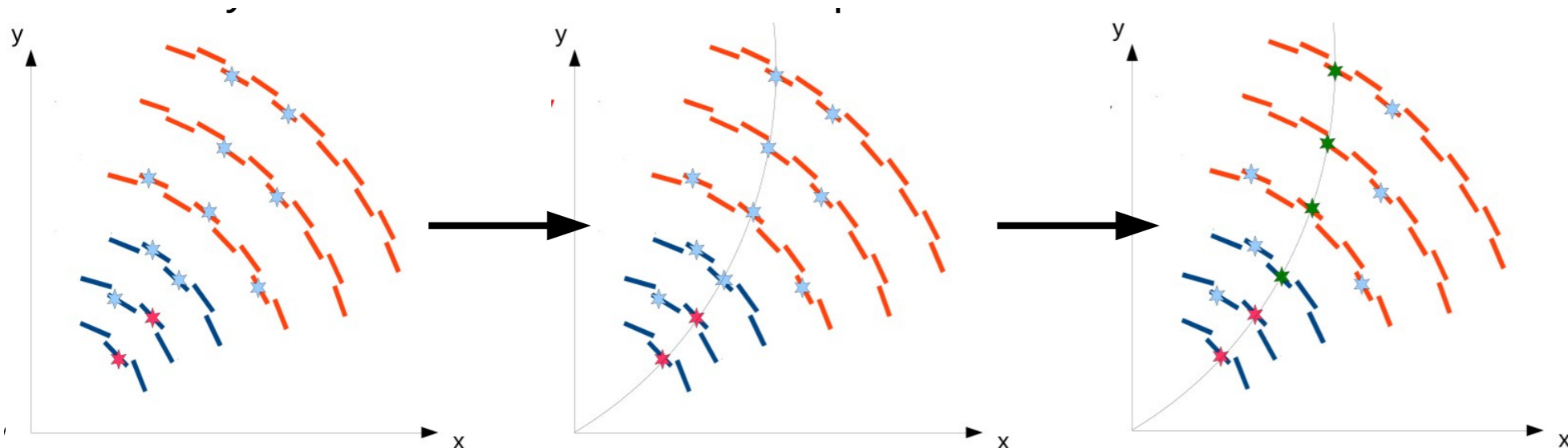
FPGAs
ML

CPUs
GPUs
FPGAs
ML

CPUs
GPUs
FPGAs
ML

**Challenges:**
Each collision produces $O(10^3)$ particles
The detectors have $O(10^8)$ sensors
Extreme data rates of $O(100\ TB/s)$

- Algorithm approach: tracklet and Kalman filter hybrid algorithm written in Vivado HLS to expedite development
    - Tracks are seeded with pairs of stubs in adjacent layers
    - Projections to other layers are calculated (assuming beamline constraint)
    - Full tracks after duplicate removal are inputs to the final track fit (Kalman filter)
- R&D efforts: displaced tracking for long-lived particles, etc.

InputRouter

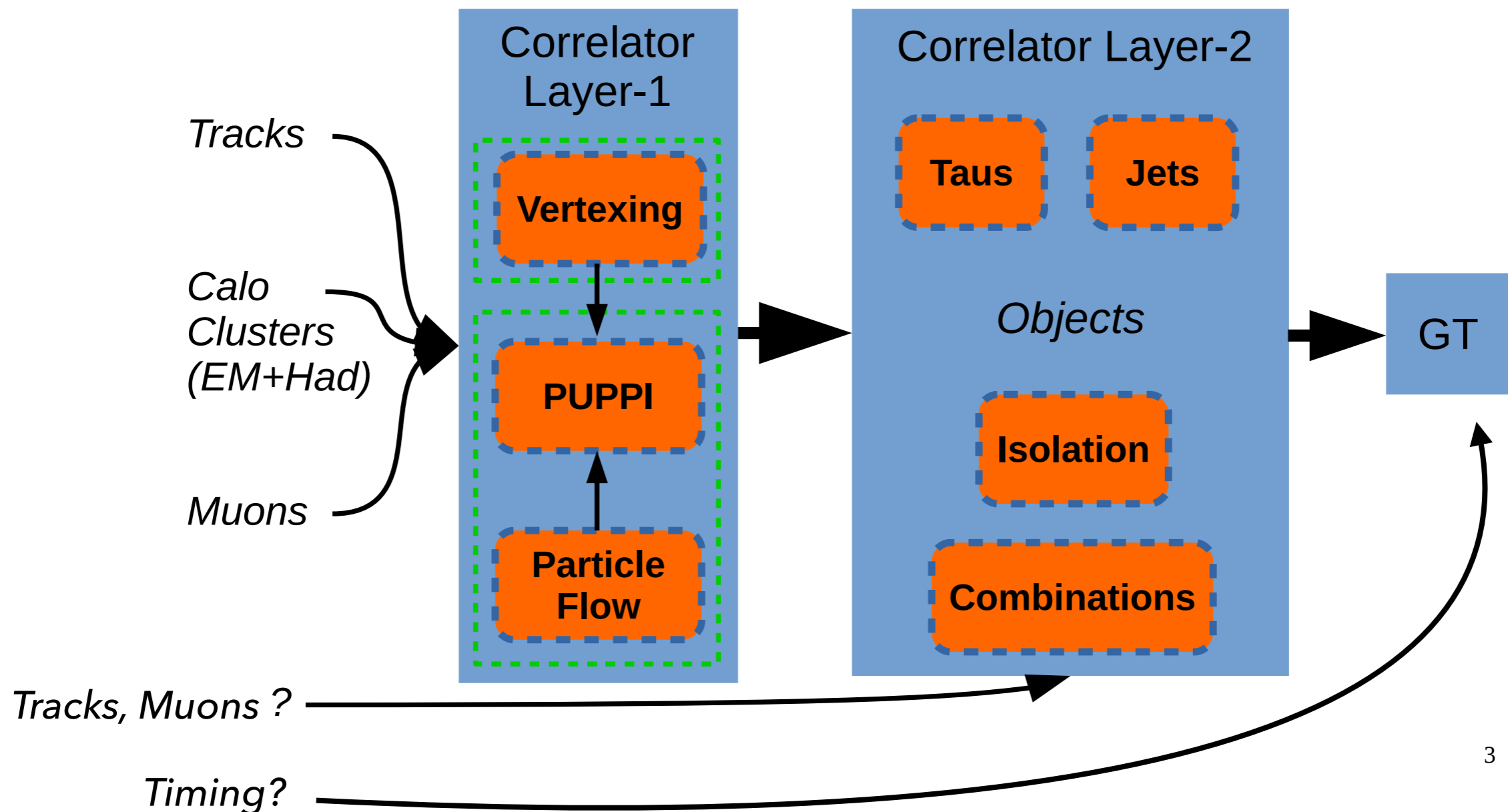VMRouter

TrackletEngine

TrackletCalculator
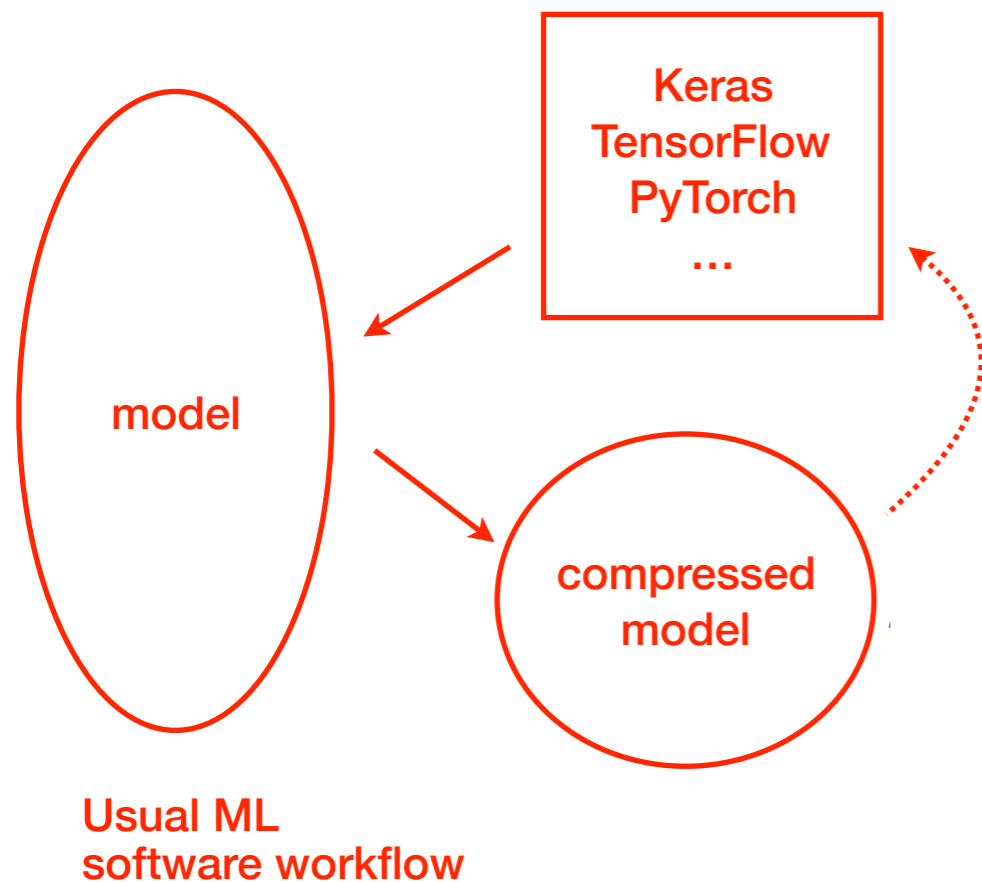
ProjectionRouter

MatchEngine
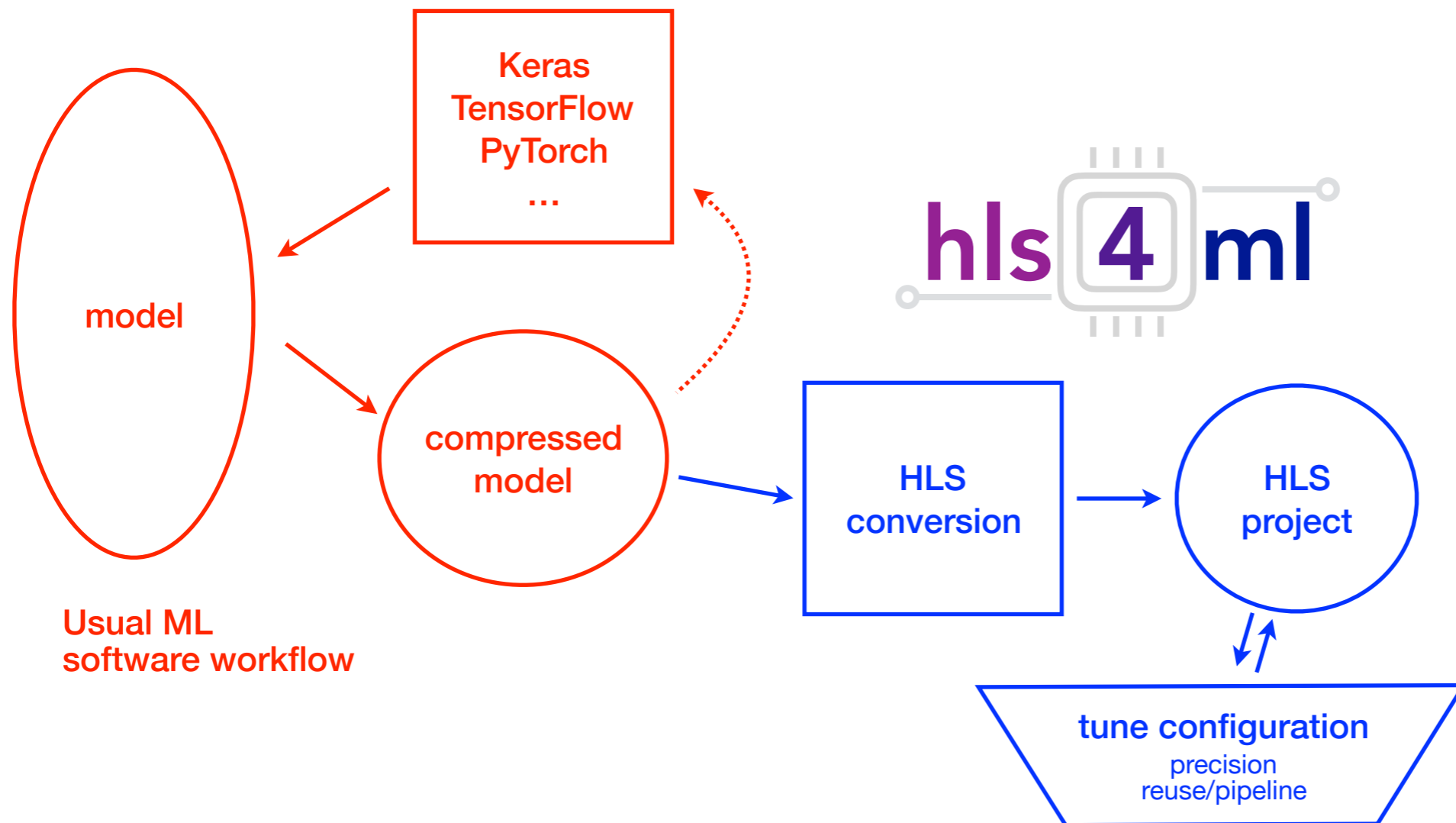
MatchCalculator

DuplicateRemoval

KalmanFilter

▸ Correlator layer 1 will process pileup mitigated candidates $\{\mu, e, \gamma, h^{\pm}, h^0, vtx\}$

▸ Full correlator trigger must complete all processing & transmit trigger objects $\{\mu, e, \gamma, \tau, j, MET, etc.\}$ to the GT within 2.5 μs



3

‣ **hls4ml** for physicists or ML experts to translate **ML algorithms** into **FPGA firmware**



Keras
TensorFlow
PyTorch
…
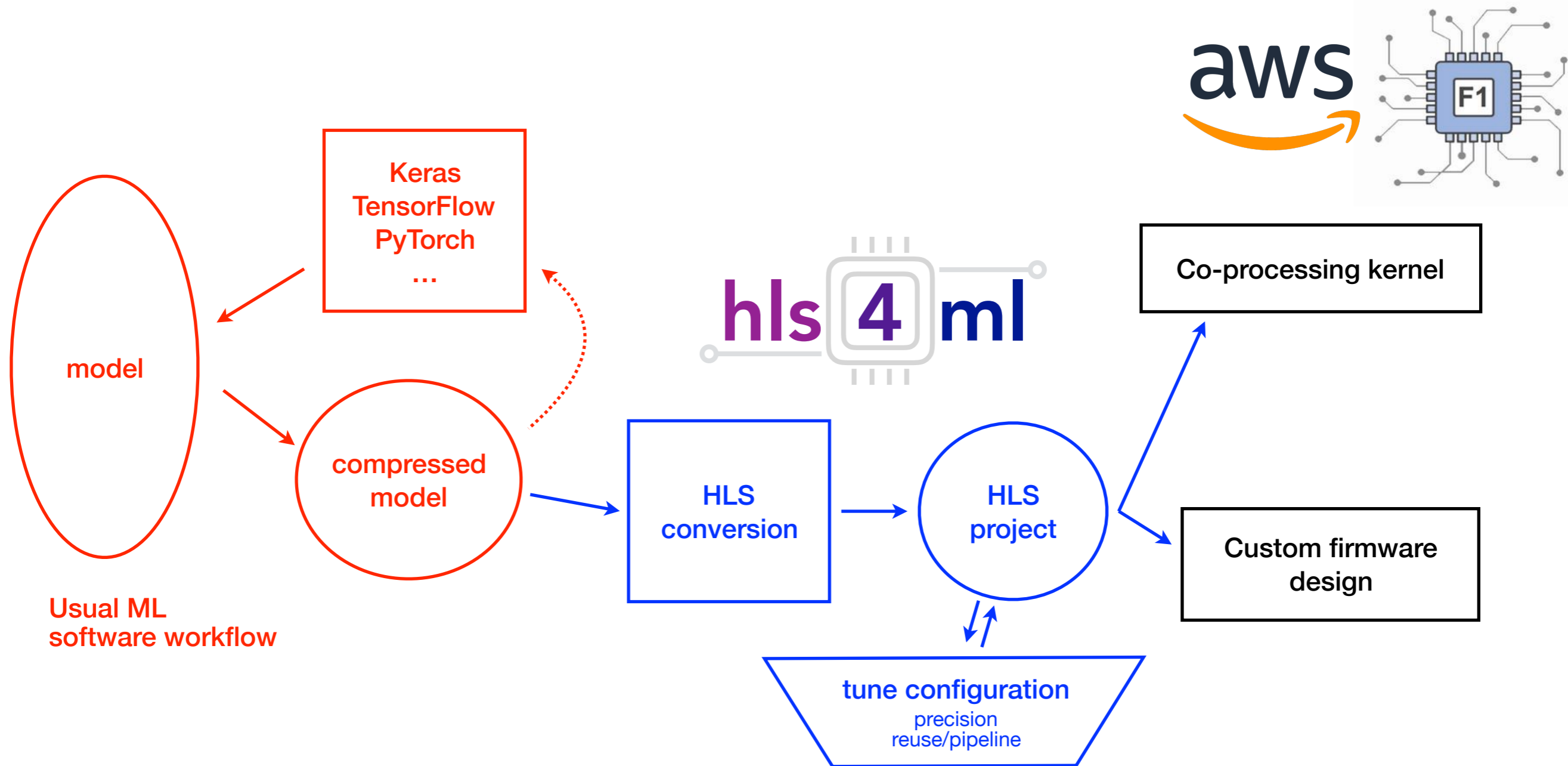
model

compressed
model

Usual ML
software workflow

hls 4 ml

▸ **hls4ml** for physicists or ML experts to translate **ML algorithms** into **FPGA firmware**

‣ **hls4ml** for physicists or ML experts to translate **ML algorithms** into **FPGA firmware**

Translation $\longrightarrow$

```
hls4ml convert -c keras-config.yml
```
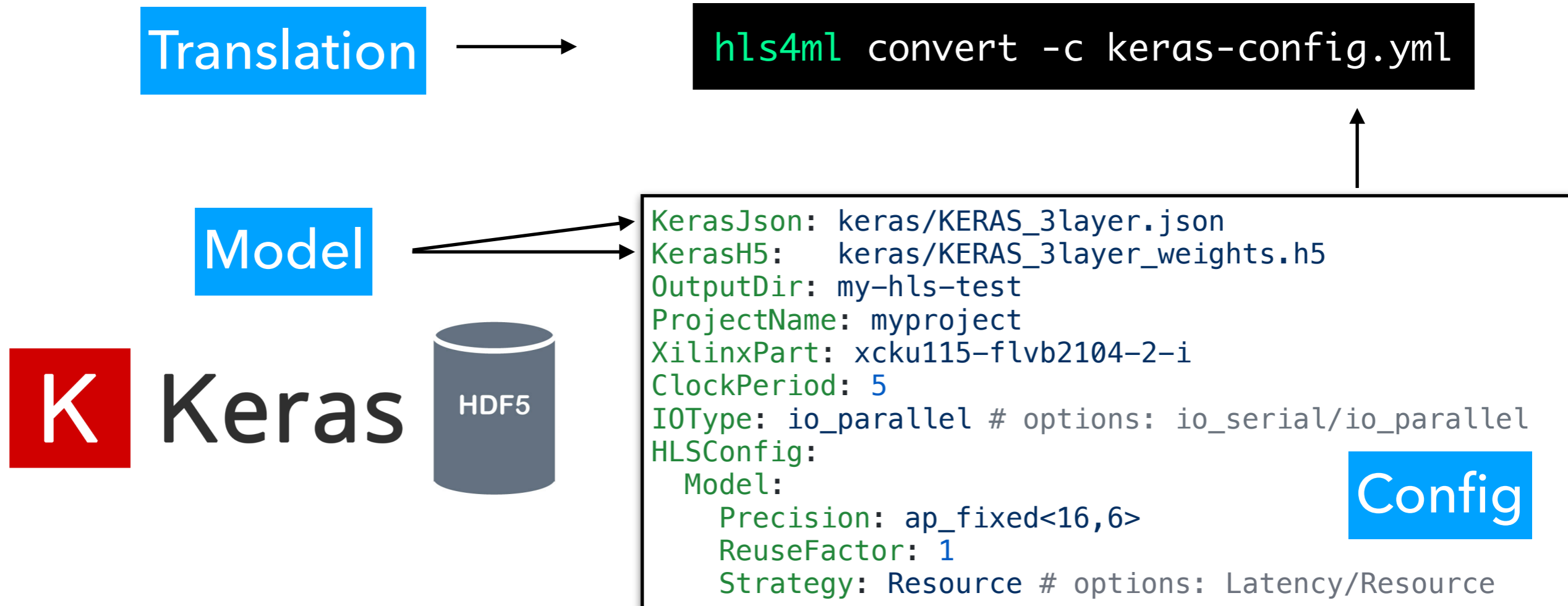
Translation →

```
hls4ml convert -c keras-config.yml
```

```
KerasJson: keras/KERAS_3layer.json
KerasH5:   keras/KERAS_3layer_weights.h5
OutputDir: my-hls-test
ProjectName: myproject
XilinxPart: xcku115-flvb2104-2-i
ClockPeriod: 5
IOType: io_parallel # options: io_serial/io_parallel
HLSConfig:
  Model:
    Precision: ap_fixed<16,6>
    ReuseFactor: 1
    Strategy: Resource # options: Latency/Resource
```

Translation → `hls4ml convert -c keras-config.yml`

Model

**K** Keras  HDF5

Config

```
KerasJson: keras/KERAS_3layer.json
KerasH5:   keras/KERAS_3layer_weights.h5
OutputDir: my-hls-test
ProjectName: myproject
XilinxPart: xcku115-flvb2104-2-i
ClockPeriod: 5
IOType: io_parallel # options: io_serial/io_parallel
HLSConfig:
  Model:
    Precision: ap_fixed<16,6>
    ReuseFactor: 1
    Strategy: Resource # options: Latency/Resource
```
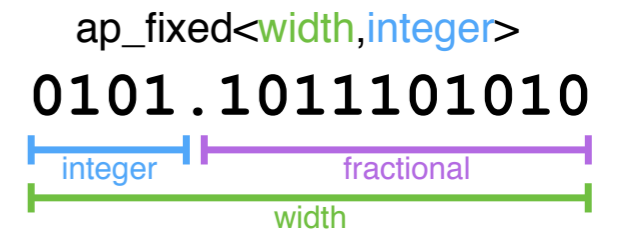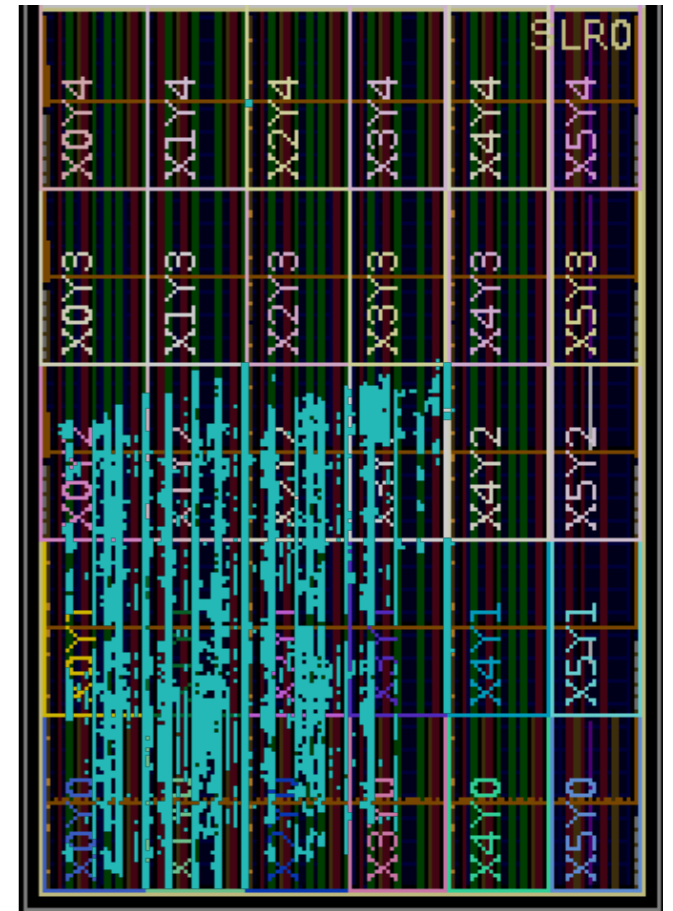
Translation ⟶ `hls4ml convert -c keras-config.yml`

Model

Keras   HDF5

Config

```
KerasJson: keras/KERAS_3layer.json
KerasH5:   keras/KERAS_3layer_weights.h5
OutputDir: my-hls-test
ProjectName: myproject
XilinxPart: xcku115-flvb2104-2-i
ClockPeriod: 5
IOType: io_parallel # options: io_serial/io_parallel
HLSConfig:
  Model:
    Precision: ap_fixed<16,6>
    ReuseFactor: 1
    Strategy: Resource # options: Latency/Resource
```
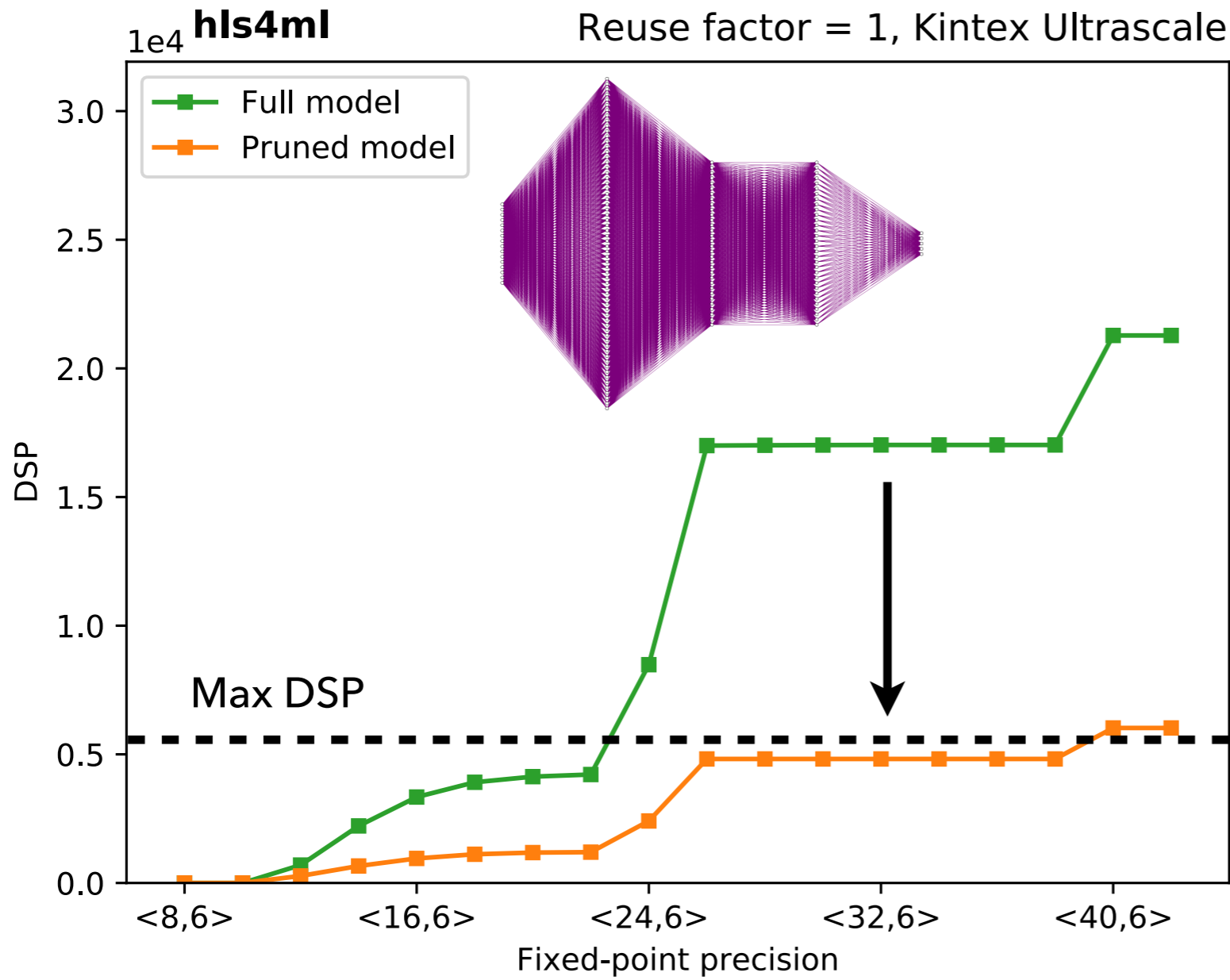
‣ **IOType**: parallel or serial

‣ **ReuseFactor**: how much to parallelize

‣ **Precision**: inputs, weights, biases

‣ **Strategy**:

  ‣ Resource for large NN

  ‣ Latency for small NN (fully pipelined)

Translation → `hls4ml convert -c keras-config.yml`

Model →

```
KerasJson: keras/KERAS_3layer.json
KerasH5:   keras/KERAS_3layer_weights.h5
OutputDir: my-hls-test
ProjectName: myproject
XilinxPart: xcku115-flvb2104-2-i
ClockPeriod: 5
IOType: io_parallel # options: io_serial/io_parallel
HLSConfig:
  Model:
    Precision: ap_fixed<16,6>
    ReuseFactor: 1
    Strategy: Resource # options: Latency/Resource
```
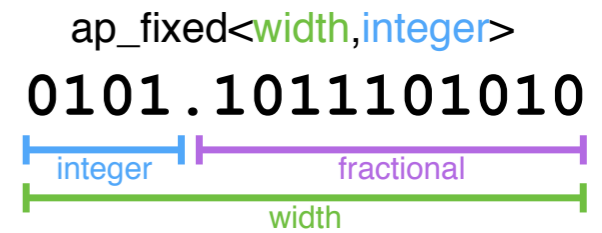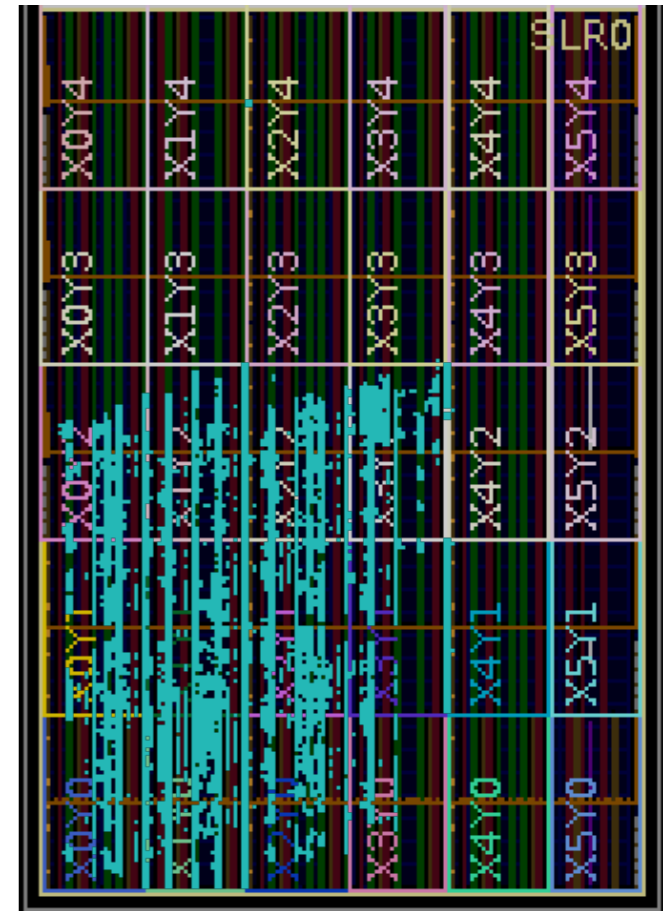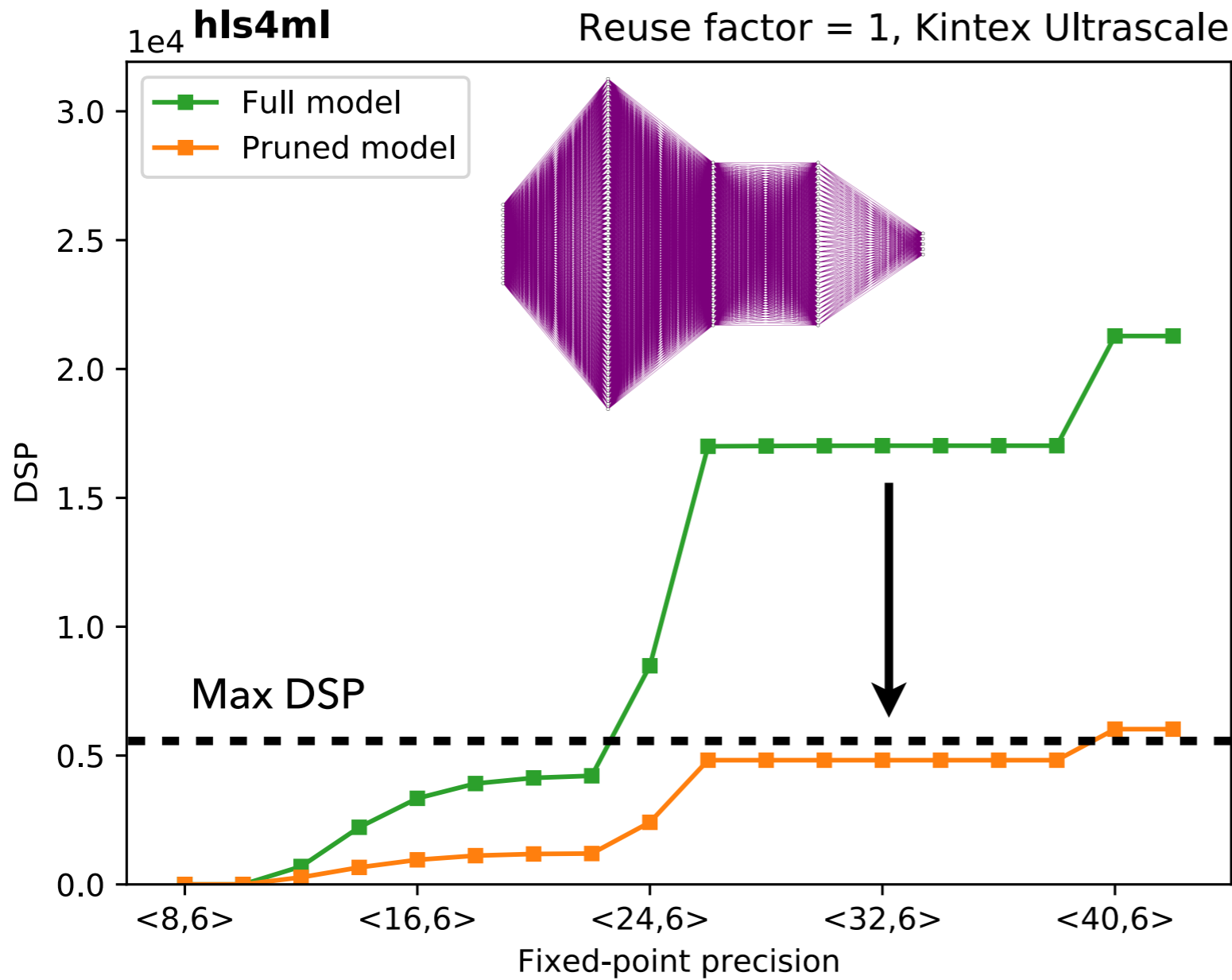
Config

**K** Keras · HDF5

‣ **IOType**: parallel or serial

‣ **ReuseFactor**: how much to parallelize

‣ **Precision**: inputs, weights, biases

‣ **Strategy**:

   ‣ Resource for large NN

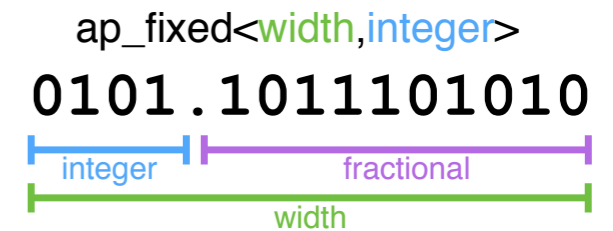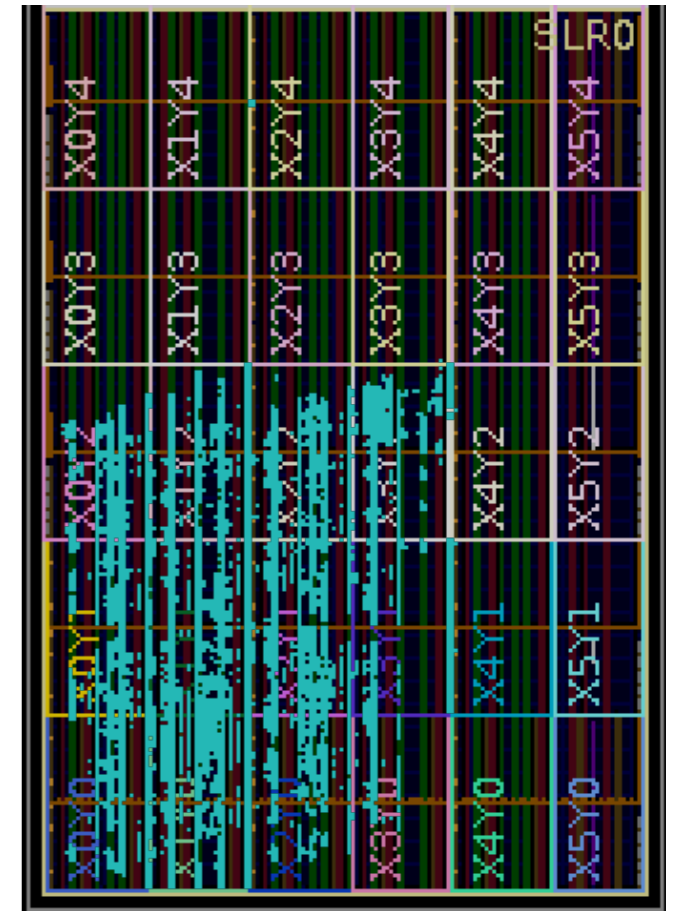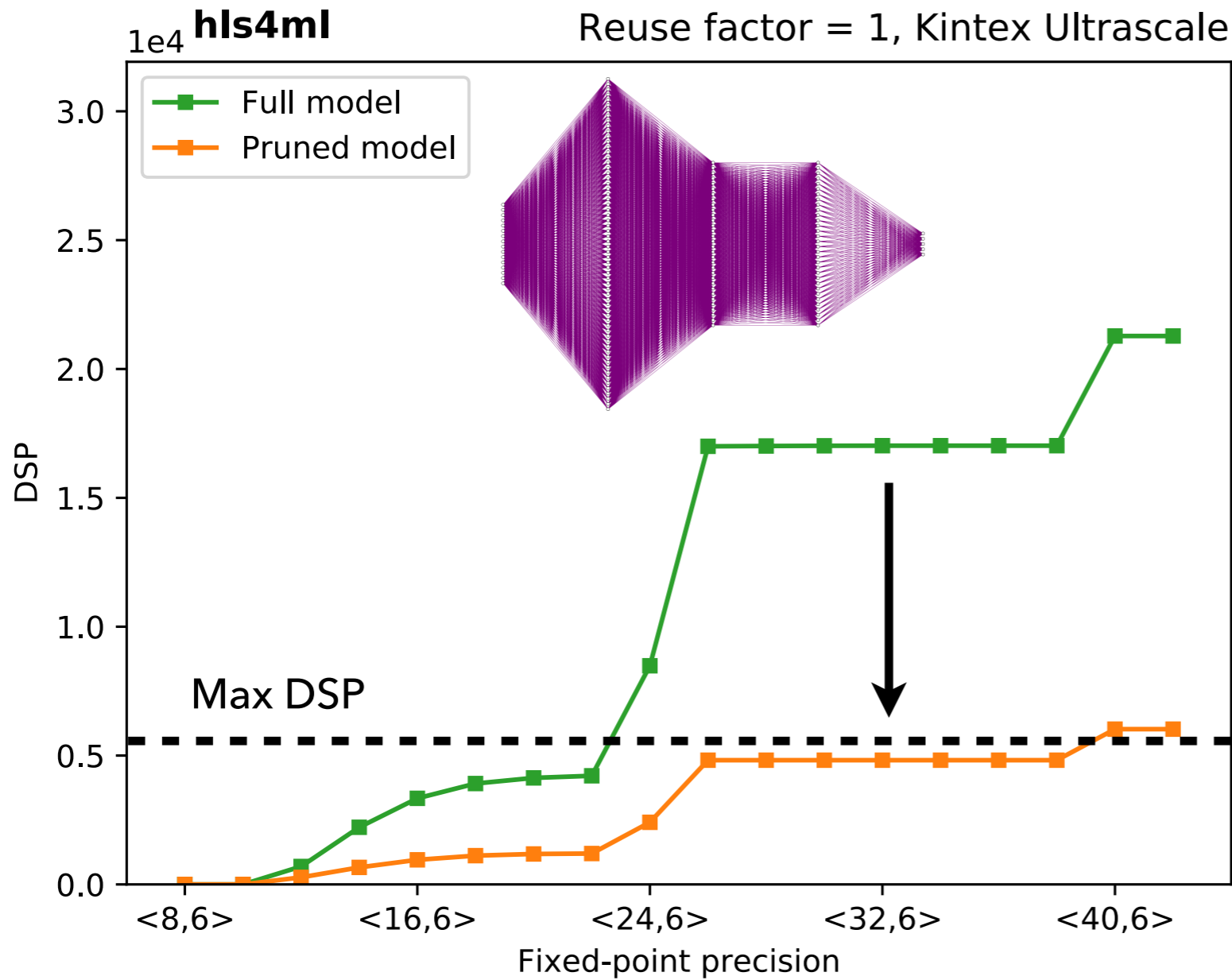   ‣ Latency for small NN (fully pipelined)

Build HLS project

↓

`hls4ml build -p my-hls-test -a`

ap_fixed<width,integer>

0101.1011101010

integer | fractional

width

Reuse factor = 1, Kintex Ultrascale

**hls4ml**

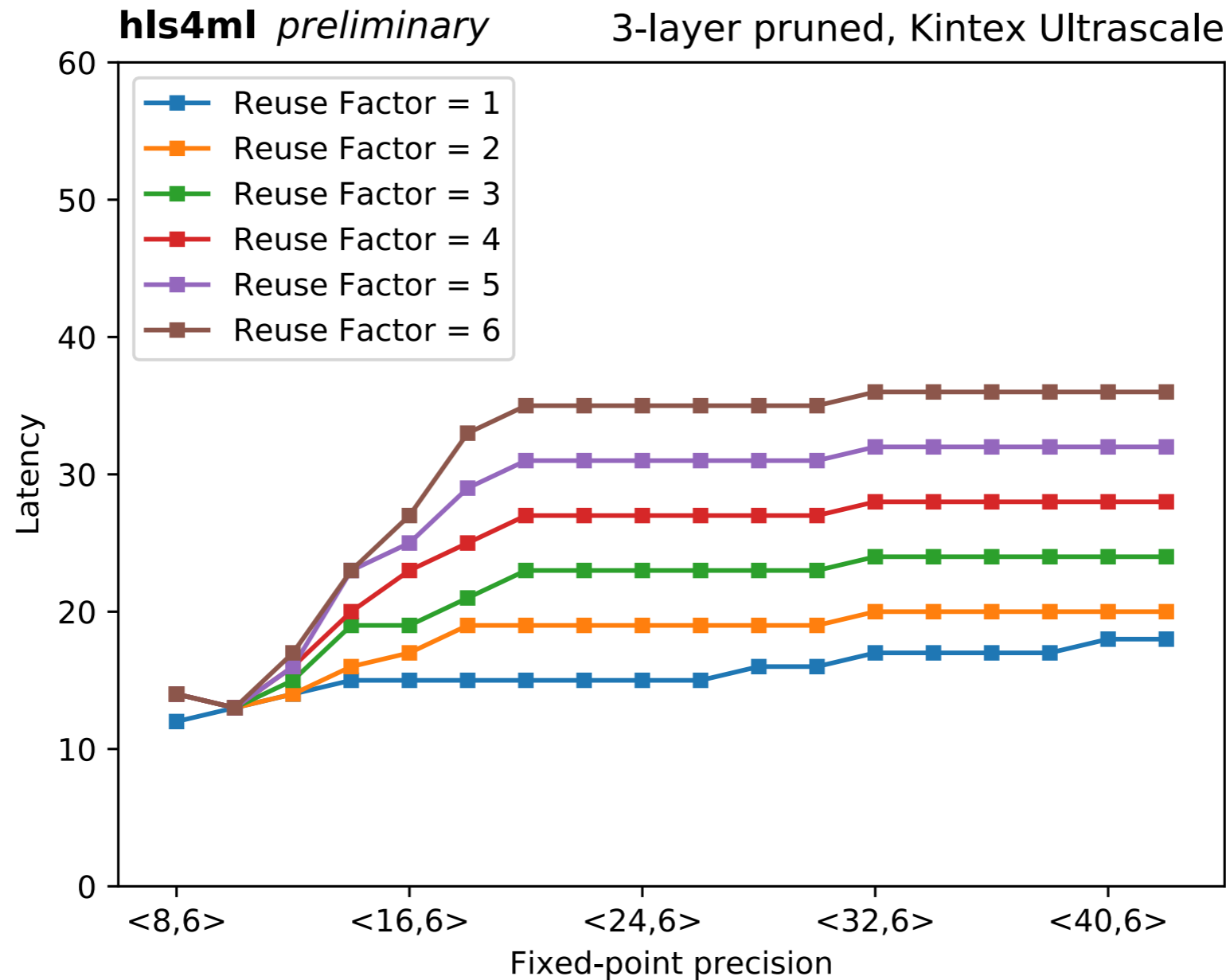ap_fixed<width,integer>

0101.1011101010

integer | fractional
width

▸ Big reduction in DSPs (multipliers) with compression

- ▶ Big reduction in DSPs (multipliers) with compression
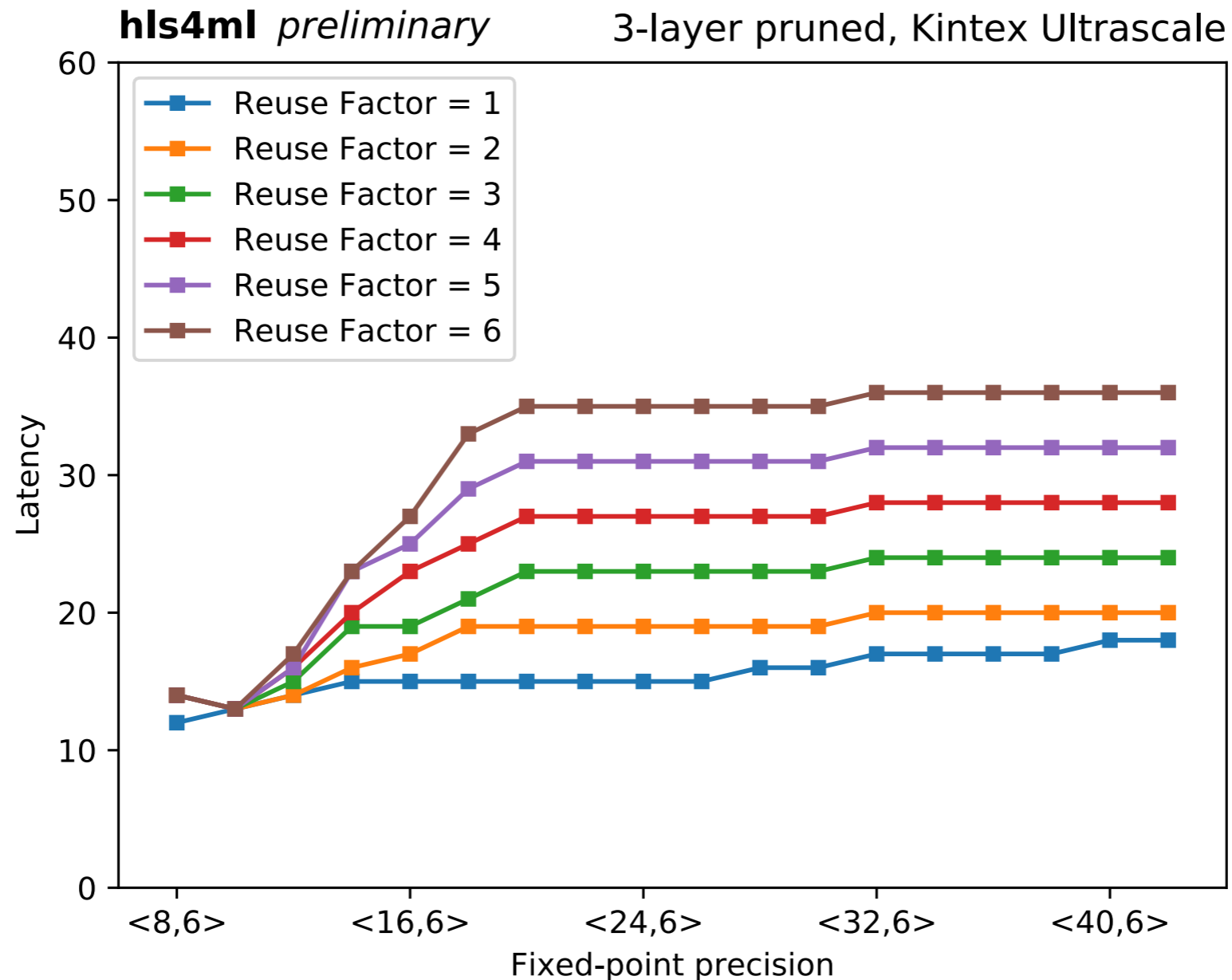- ▶ Easily fits on 1 FPGA **after compression**

▸ **Increasing** reuse factor, **increases** latency



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

~35 clocks
@ 200 MHz
= 175 ns

~15 clocks
@ 200 MHz
= 75 ns

▶ **Increasing** reuse factor, **increases** latency



~35 clocks
@ 200 MHz
= 175 ns

~15 clocks
@ 200 MHz
= 75 ns

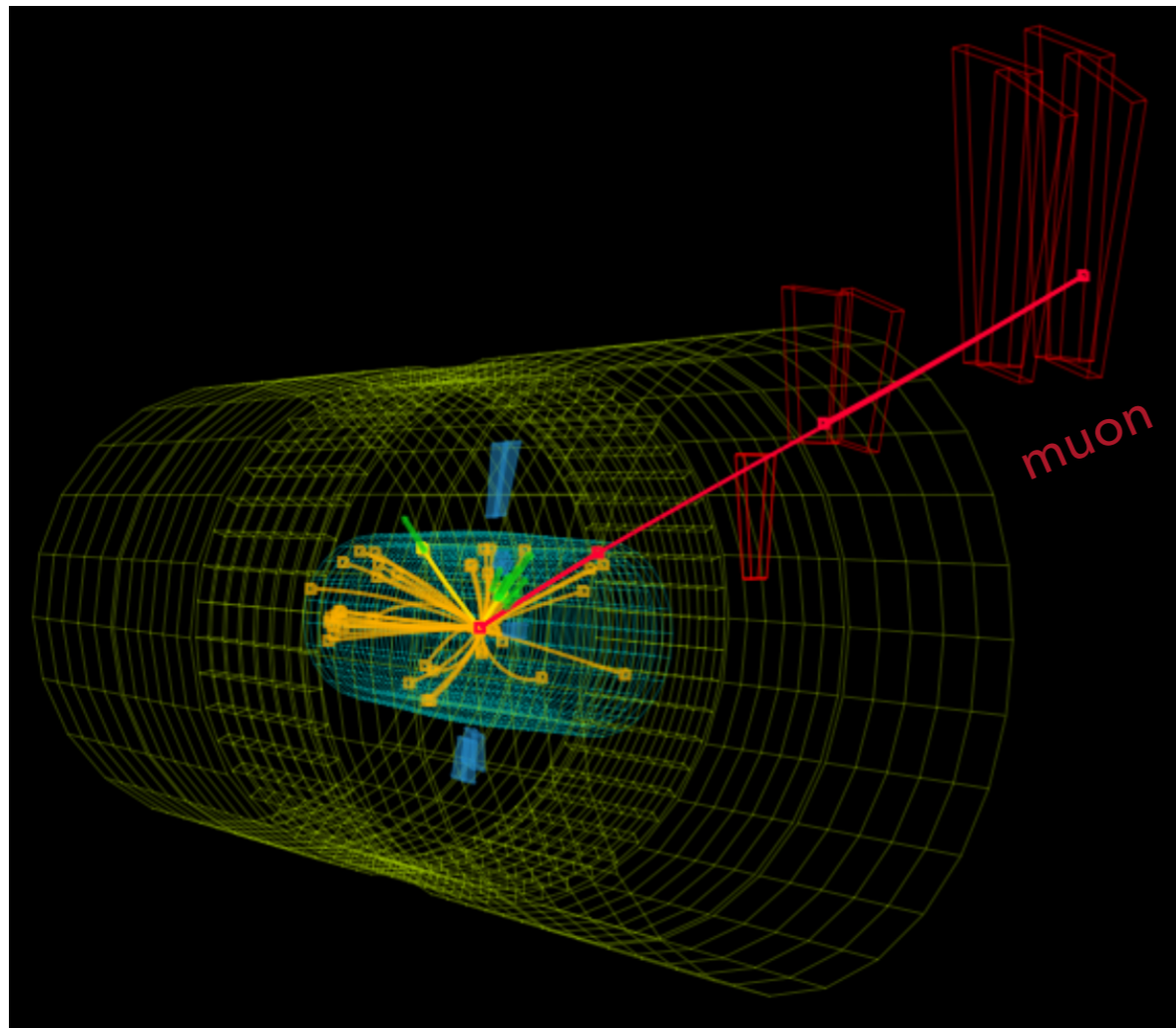For low-latency, small reuse factor, inference in O(100 ns)!
What if we have O(ms)? Can go to **bigger networks!**

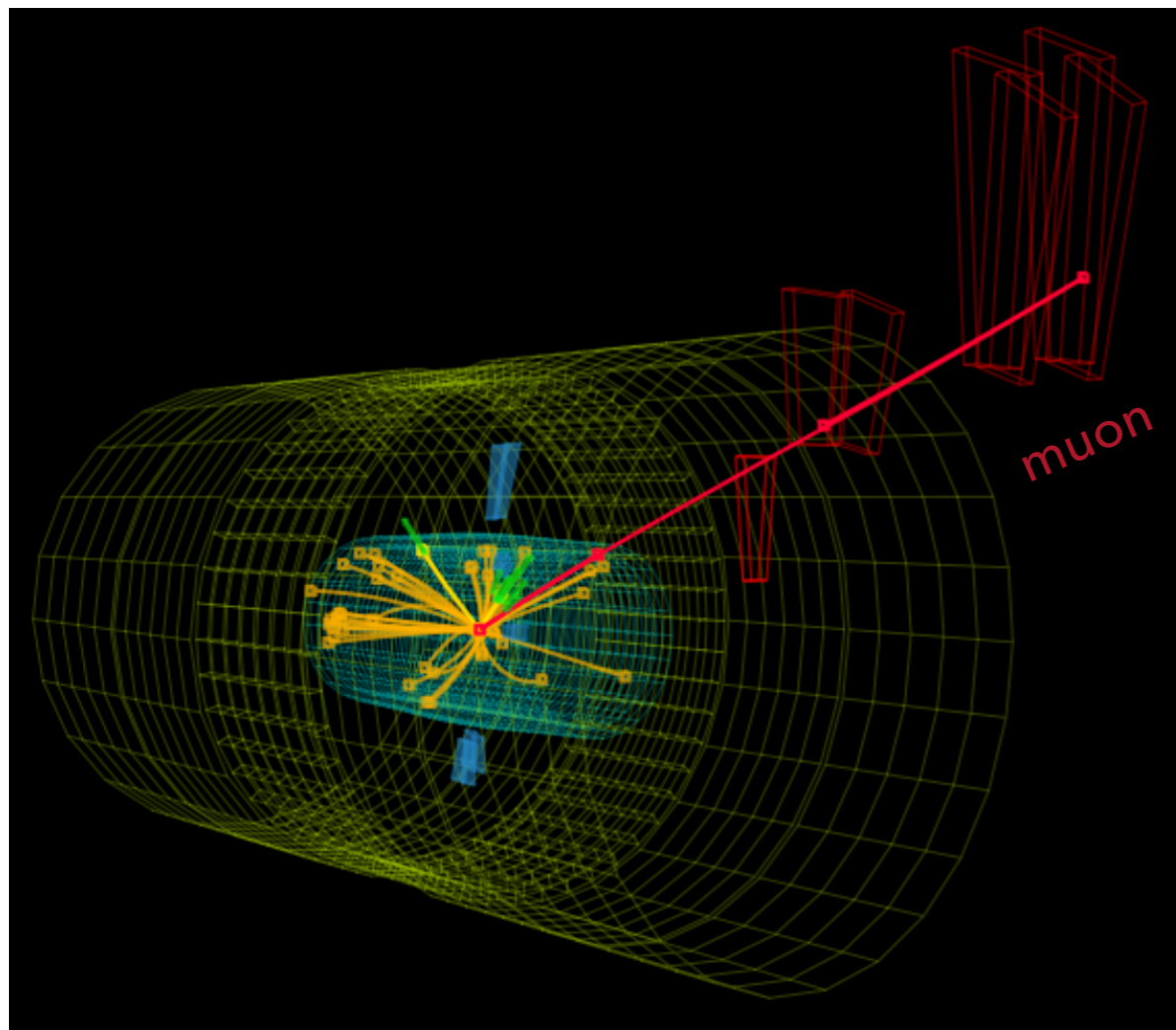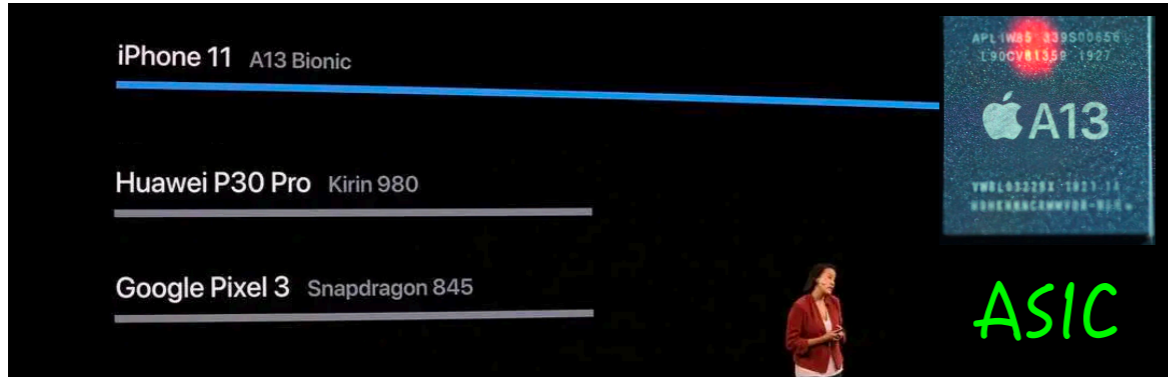▸ Inference of ML algorithms possible in **O(100 ns)** on **1 FPGA** with **hls4ml**!

▸ Inference of ML algorithms possible in **O(100 ns)** on **1 FPGA** with **hls4ml**!

  ▸ Applications across CMS, ATLAS, DUNE, and accelerator controls:

‣ Inference of ML algorithms possible in **O(100 ns)** on **1 FPGA** with **hls4ml**!

  ‣ Applications across CMS, ATLAS, DUNE, and accelerator controls:

    ‣ Muon $p_T$ determination in the CMS endcap with a DNN:
    **runs in 160 ns** on an FPGA and **reduces** the **fake muon rate** by
    **up to 80%**

▸ Inference of ML algorithms possible in **O(100 ns)** on **1 FPGA** with **hls4ml**!

  ▸ Applications across CMS, ATLAS, DUNE, and accelerator controls:

    ▸ Muon $p_T$ determination in the CMS endcap with a DNN:
      **runs in 160 ns** on an FPGA and **reduces** the **fake muon rate** by
      **up to 80%**

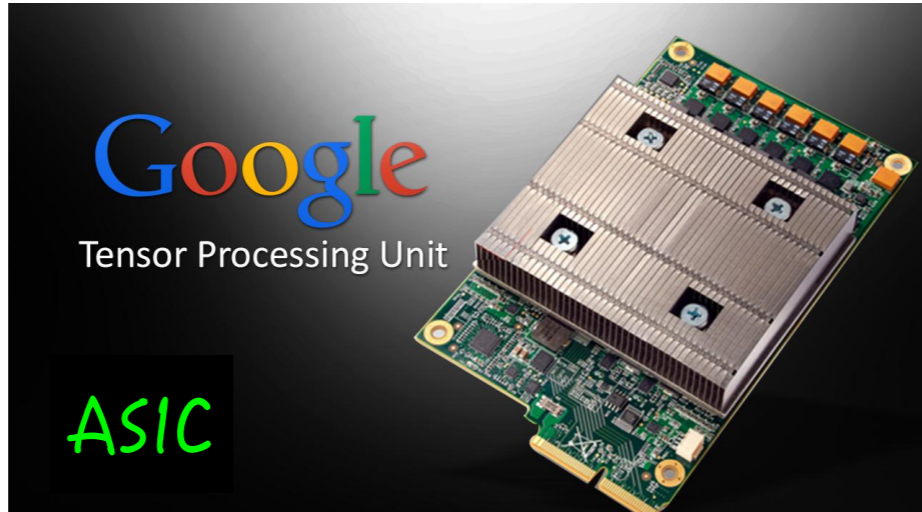    ▸ **Variational autoencoder** for anomaly detection

▸ Inference of ML algorithms possible in **O(100 ns)** on **1 FPGA** with **hls4ml**!

  ▸ Applications across CMS, ATLAS, DUNE, and accelerator controls:

    ▸ Muon $p_T$ determination in the CMS endcap with a DNN:
      **runs in 160 ns** on an FPGA and **reduces** the **fake muon rate** by
      **up to 80%**

    ▸ **Variational autoencoder** for anomaly detection



▸ Currently supported:

  ▸ Small and large dense NNs

  ▸ Binary and ternary NNs

  ▸ Small 1D/2D CNNs

▸ Planned support

  ▸ Big 1D/2D CNNs

  ▸ Graph NNs

  ▸ Other HLS/RTL backends

Specialized co-processor hardware for machine learning inference

FPGA

**Catapult/Brainwave**

Microsoft

ASIC

**Delivering FPGA Partner Solutions on AWS**
via AWS Marketplace

Customers

FPGA

FPGA

AWS Marketplace

Amazon Machine Image (AMI)

Amazon FPGA Image (AFI)

Amazon EC2 FPGA Deployment via Marketplace

AFI is secured, encrypted, dynamically loaded into the FPGA - can't be copied or downloaded

iPhone 11  A13 Bionic

Huawei P30 Pro  Kirin 980

Google Pixel 3  Snapdragon 845

A13

intel
Stratix 10
1SG280LN3F43E3VGS1

INTEL® FPGA ACCELERATION HUB

Google
Tensor Processing Unit

ASIC

XILINX
VERSAL

FPGA

XILINX
VERSAL

Specialized co-processor hardware
for machine learning inference

**FPGA**

Microsoft **Catapult/Brainwave**

Datacenter (CPU farm)

Experimental Software

Network input

Prediction

gRPC protocol

Heterogeneous Cloud Resource

CPU

FPGA

▸ Services for Optimized Network Inference on Coprocessors (SONIC)

  ▸ Send jet images from CMSSW to Microsoft Brainwave FPGA
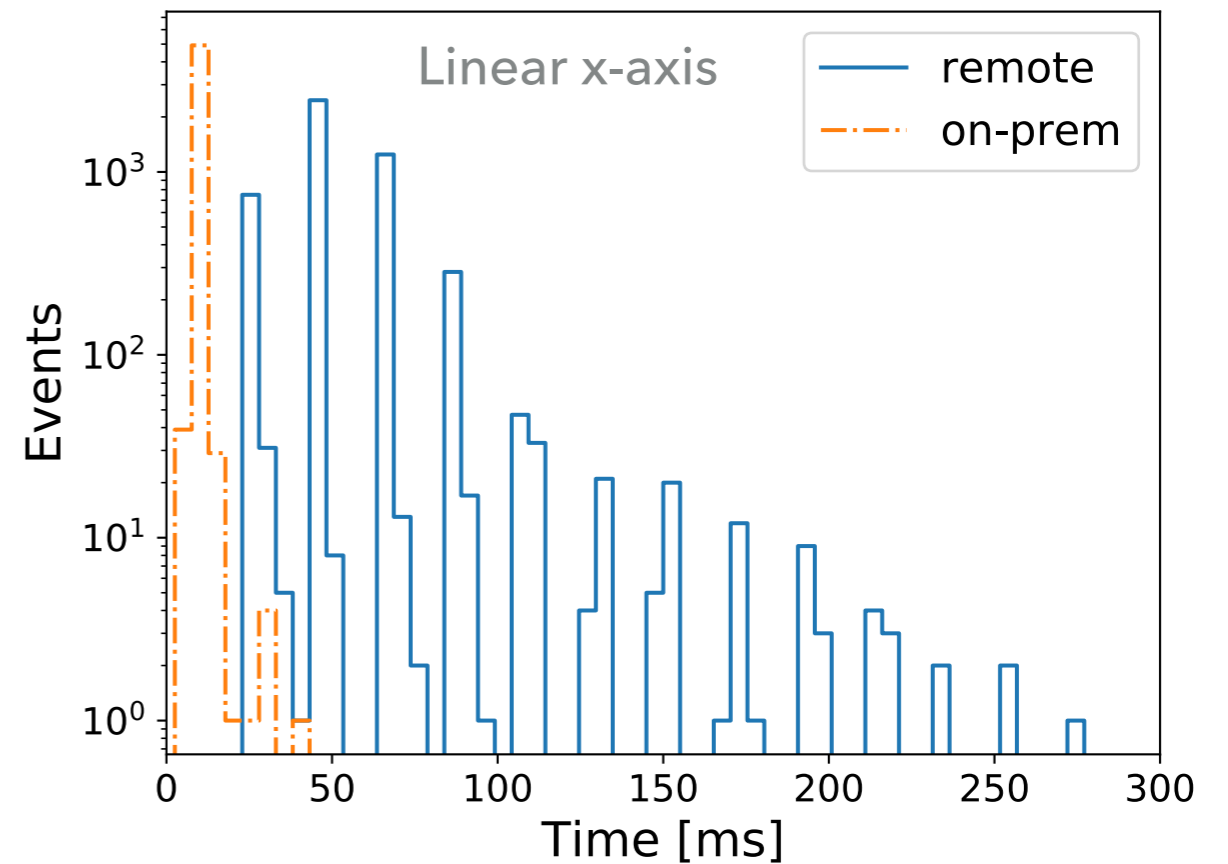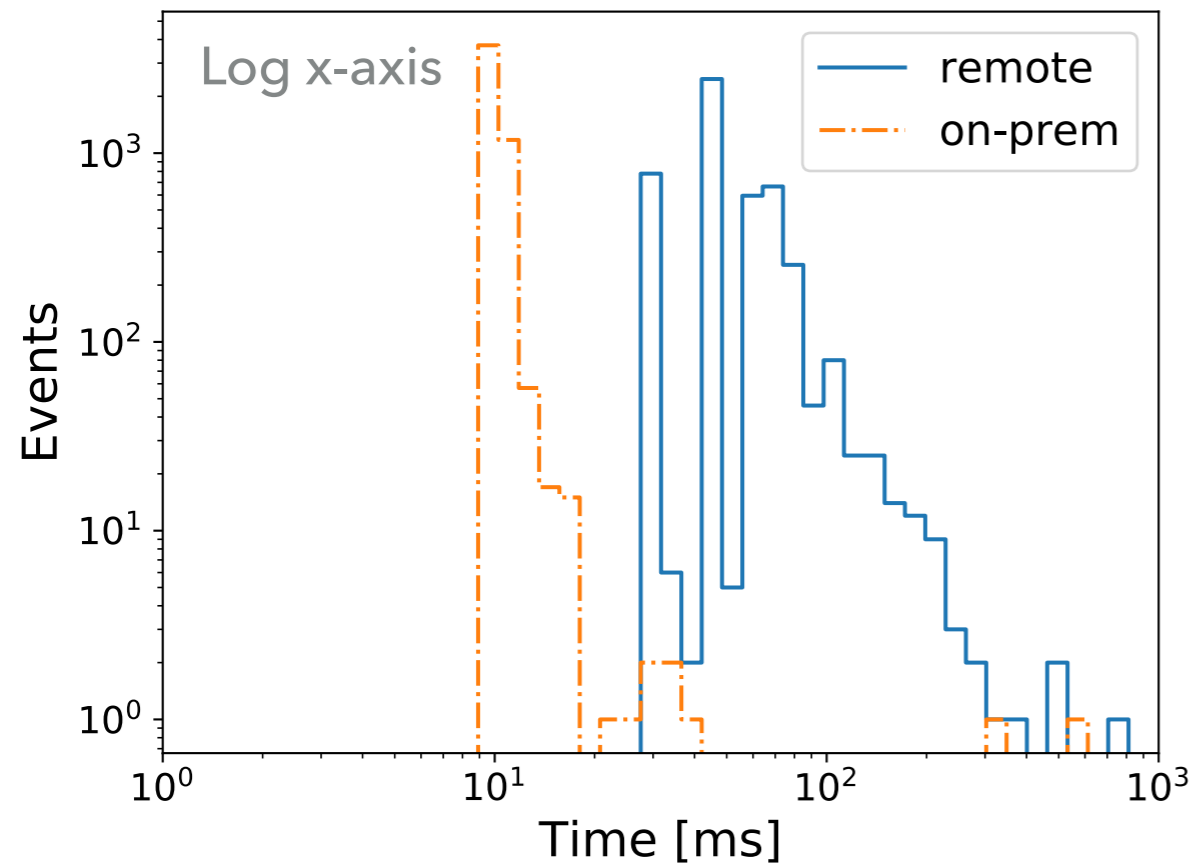
**Datacenter (CPU farm)**



- [Services for Optimized Network Inference on Coprocessors (SONIC)](#)

  - Send jet images from CMSSW to Microsoft Brainwave FPGA
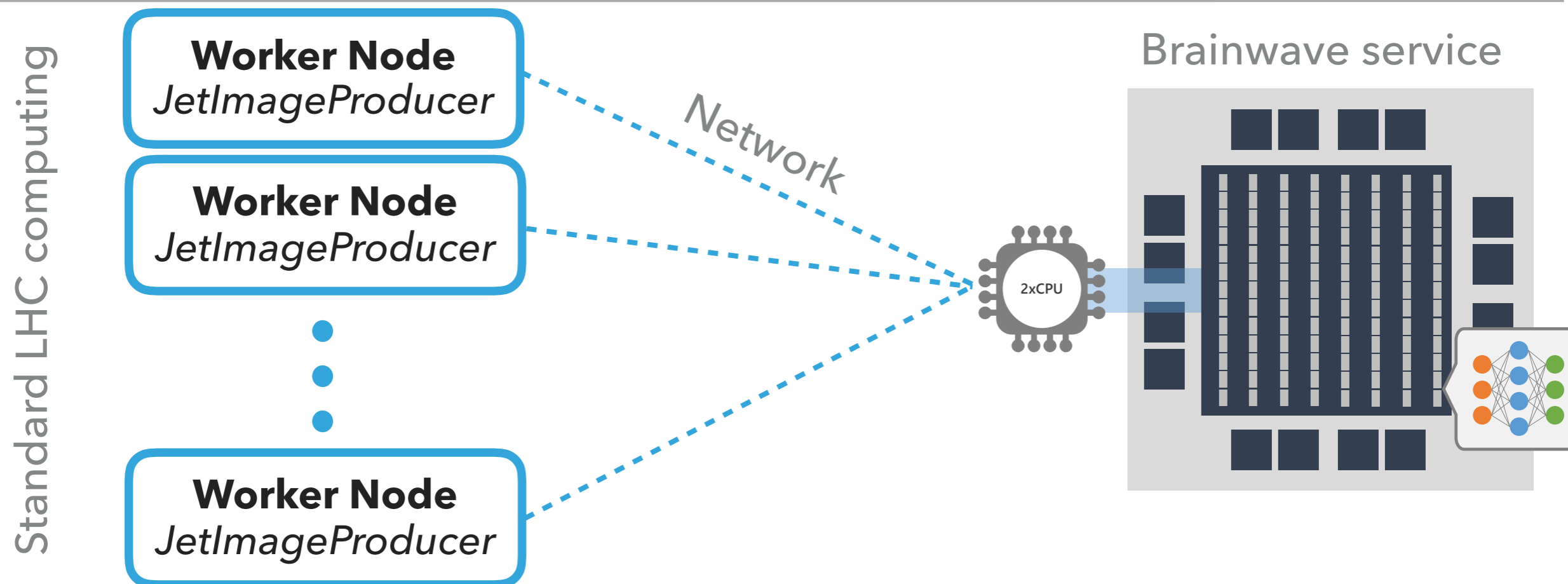
- Two modes: cloud service and on premises

▸ Remote: FNAL (IL) to Azure (VA)    ‹time› = 60 ms

  ▸ Highly dependent on network conditions

- ▸ Remote: FNAL (IL) to Azure (VA)    ‹time› = 60 ms
  - ▸ Highly dependent on network conditions
- ▸ On-prem: run CMSSW on Azure    ‹time› = 10 ms
  - ▸ on FPGA: 1.8 ms for inference
  - ▸ Remaining time used for classifying and I/O

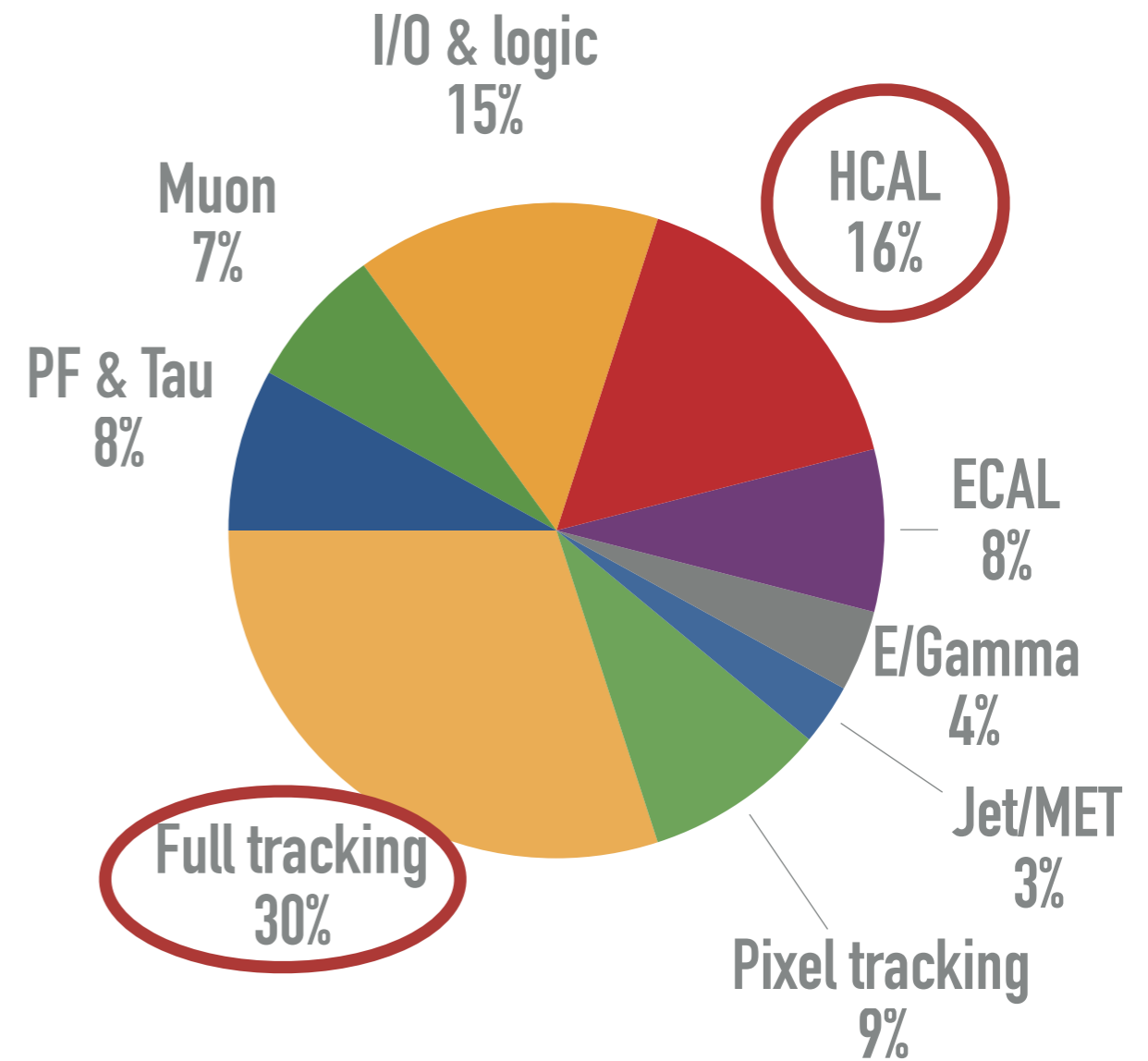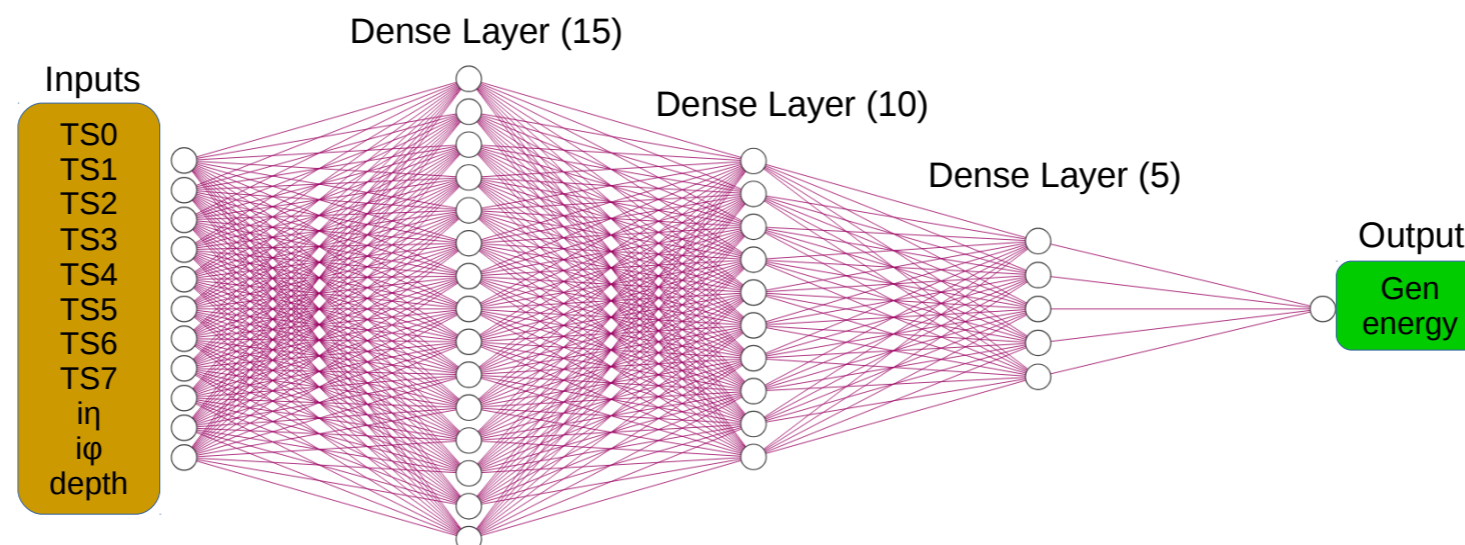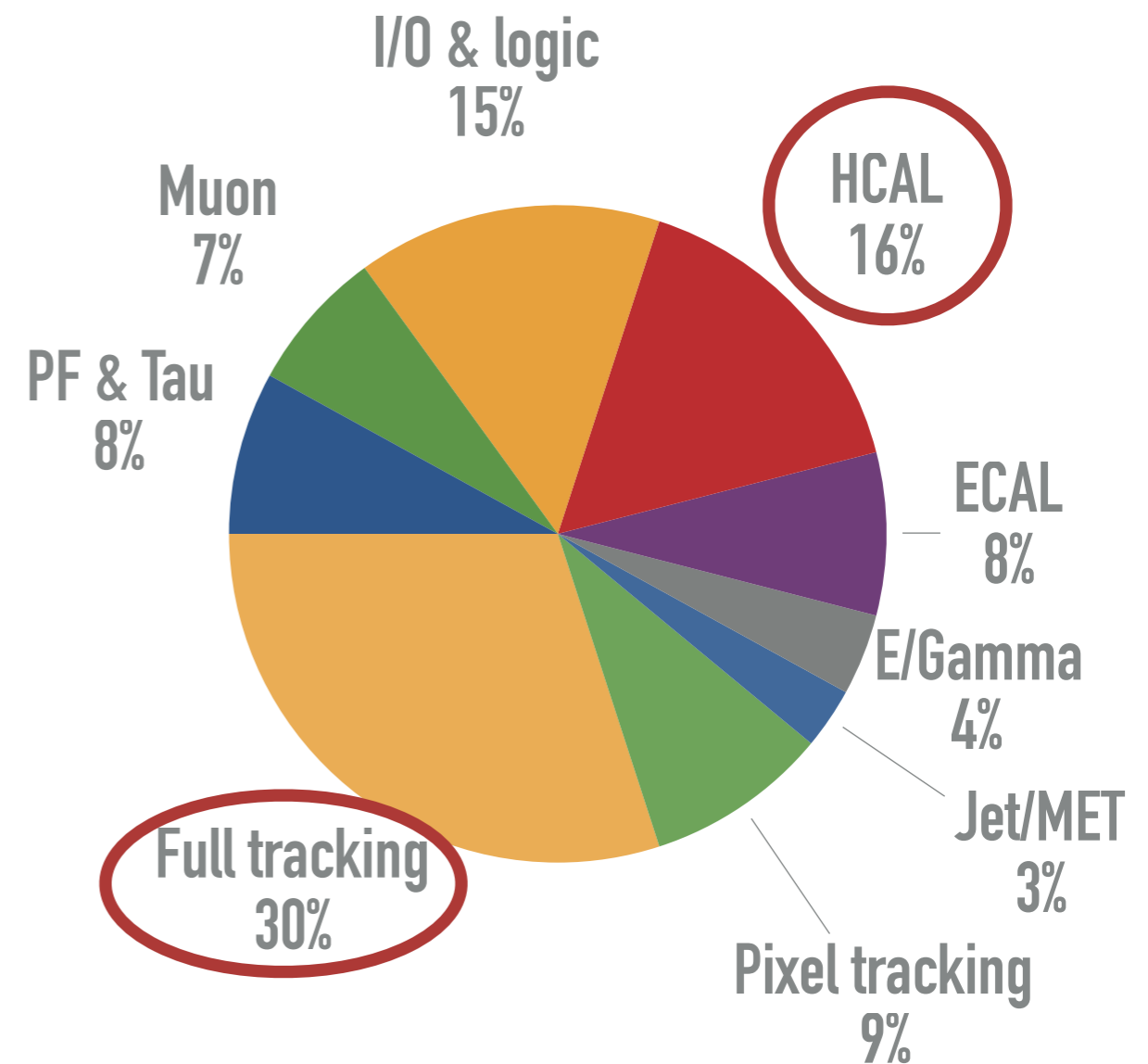▸ **Brainwave + SONIC achieves**

   ▸ **175× (30×)** on-prem (remote) **better latency**
     vs. CMS CPU

   ▸ 1 FPGA service can serve 100s of CPU worker nodes
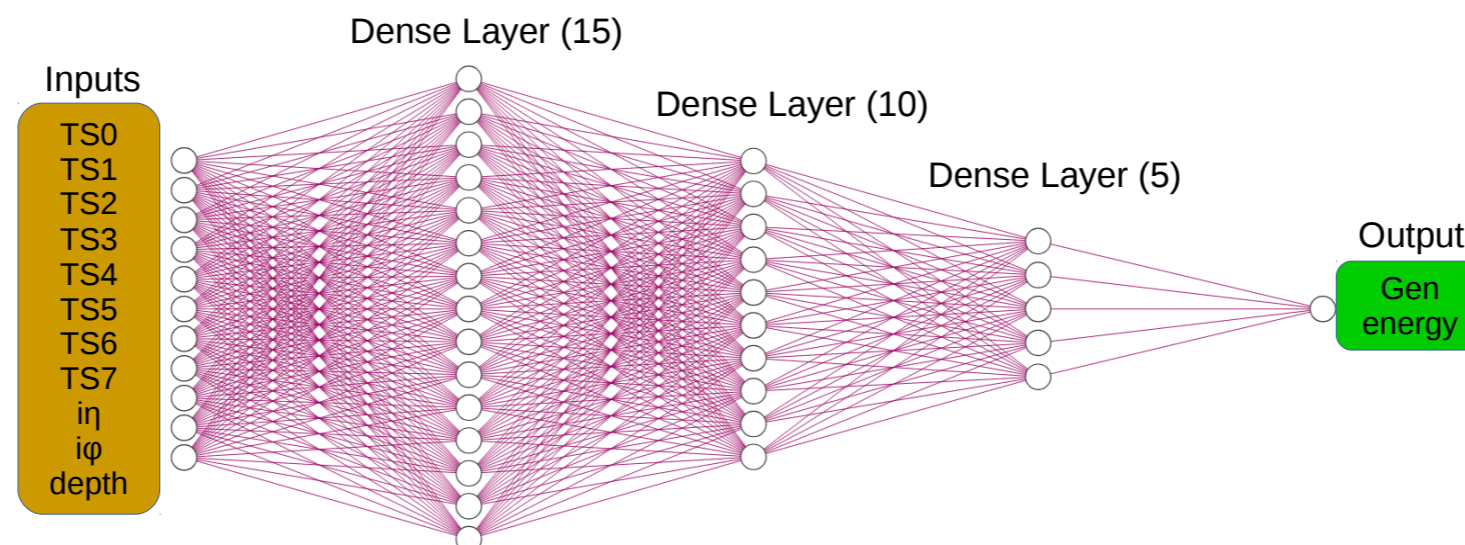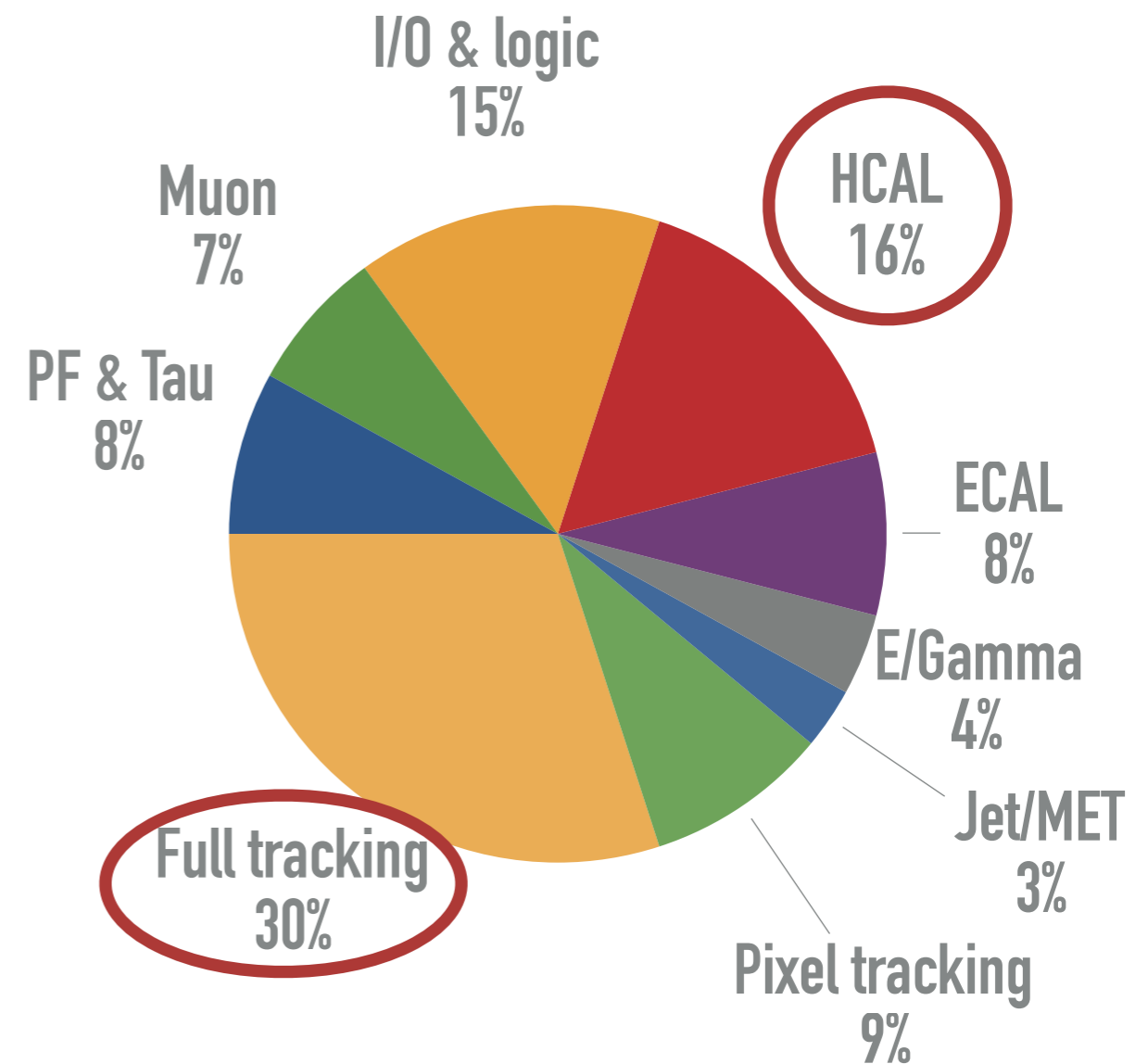
     ▸ **Competitive throughput** vs. GPU as a service

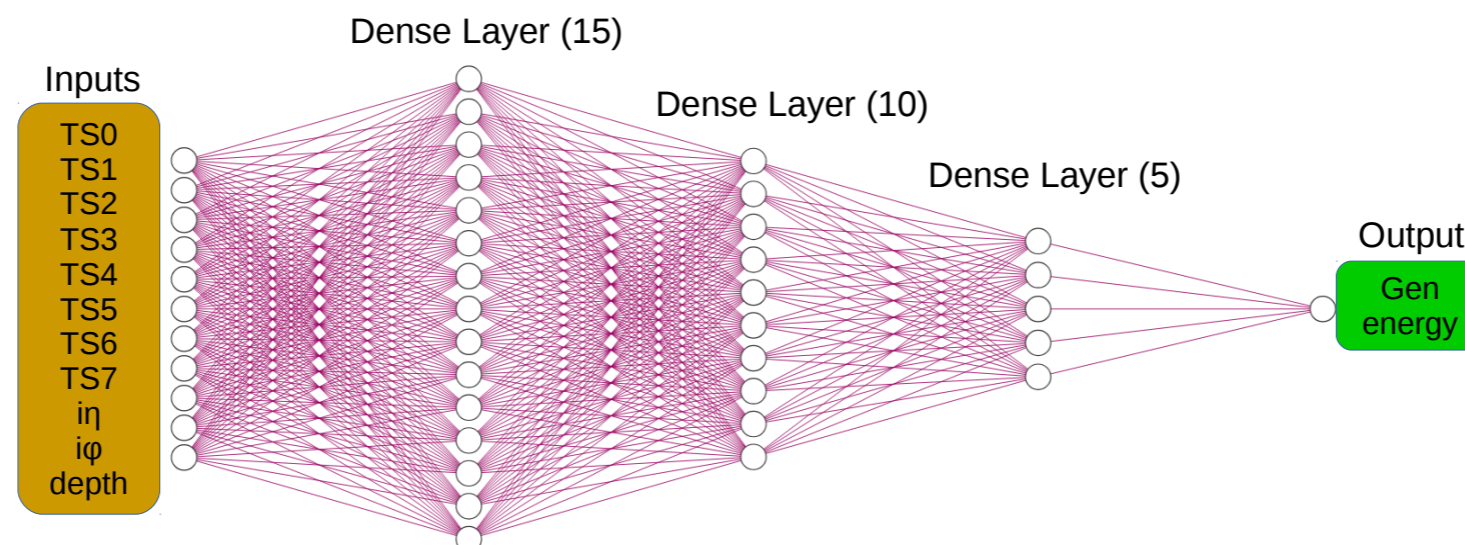▸ HCAL reconstruction and tracking contribute significantly to HLT compute time

- HCAL reconstruction and tracking contribute significantly to HLT compute time

- GPU/FPGA as co-processor can reduce compute time

▸ HCAL reconstruction and tracking contribute significantly to HLT compute time

▸ GPU/FPGA as co-processor can reduce compute time

    ▸ Patatrack pixel reconstruction on GPUs

▸ HCAL reconstruction and tracking contribute significantly to HLT compute time
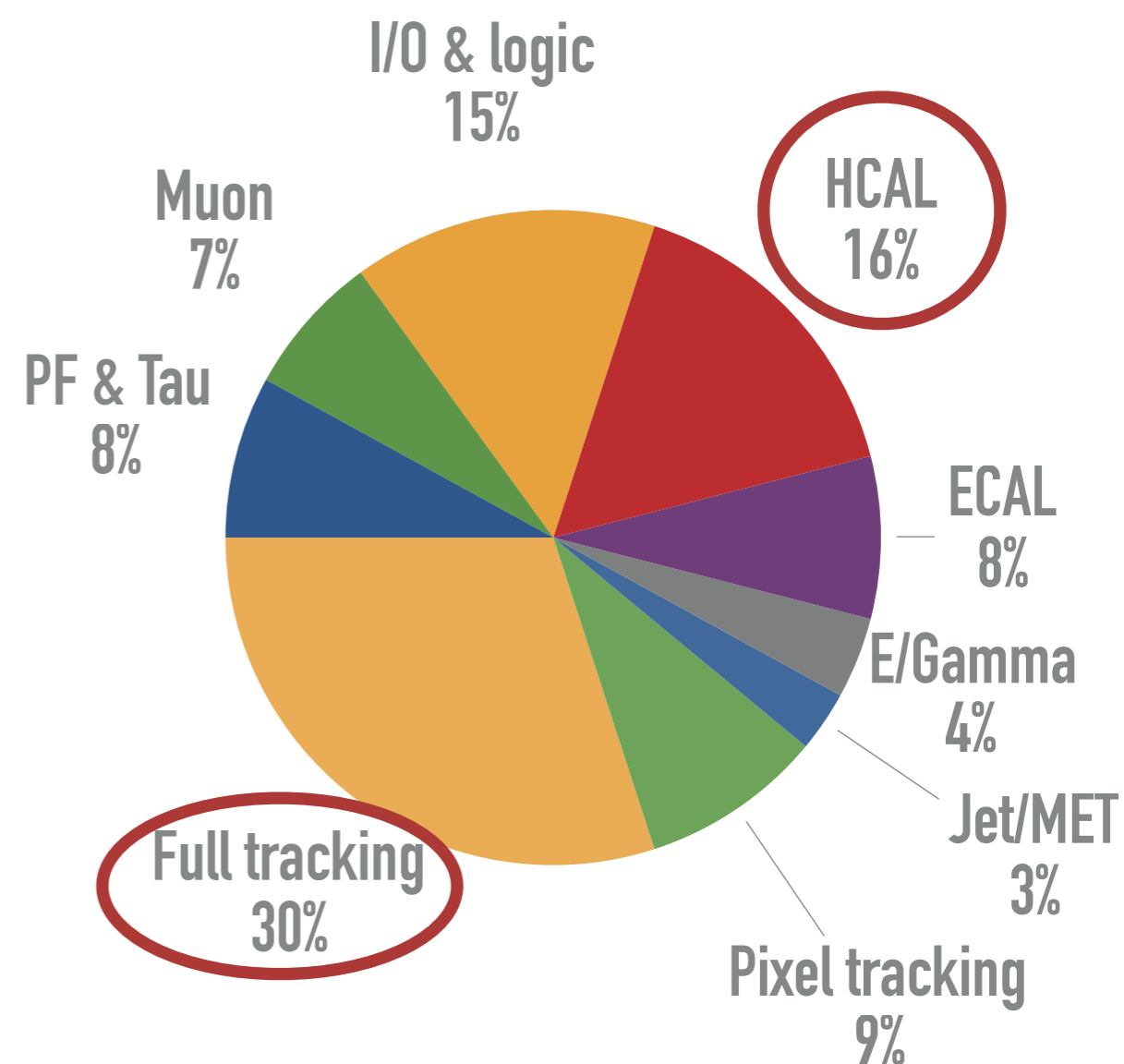
▸ GPU/FPGA as co-processor can reduce compute time

  ▸ Patatrack pixel reconstruction on GPUs
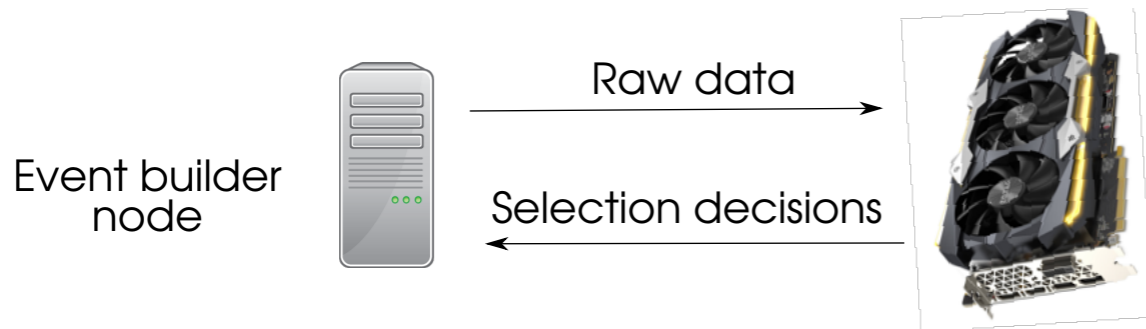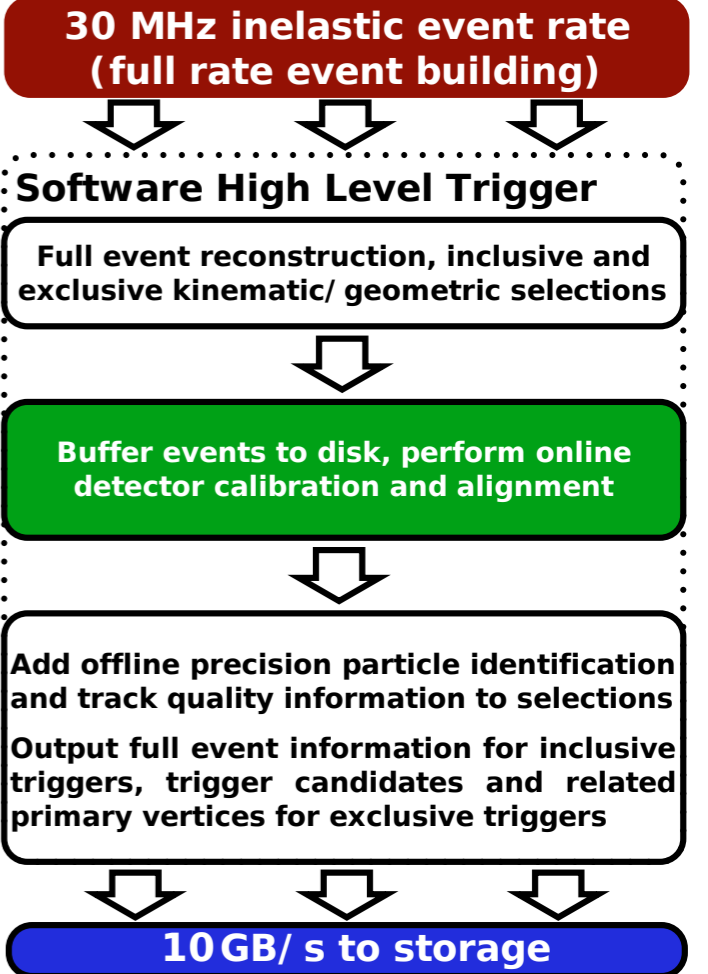
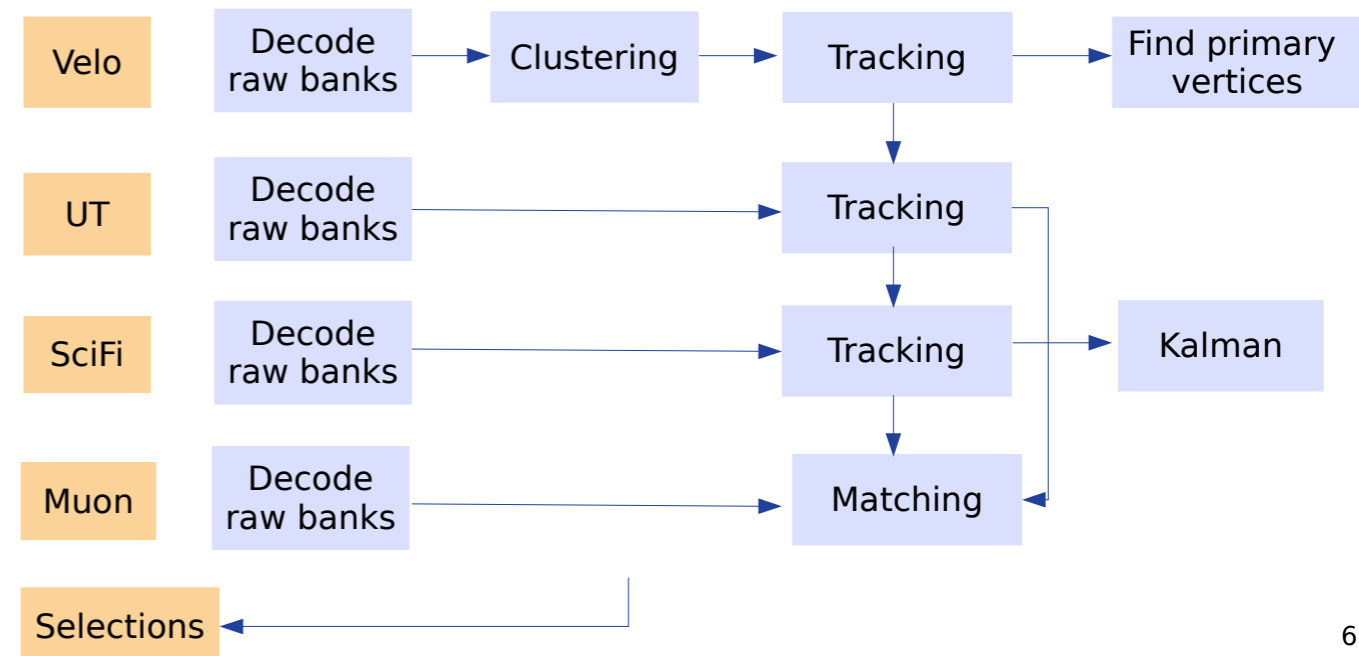  ▸ HCAL reconstruction with ML on GPUs/FPGAs (as a service)

- By 2021, full LHCb trigger chain in software (HLT)

- Run full first stage of HLT (HLT1) on GPUs

- One GPU has to process 30/60 k events/s

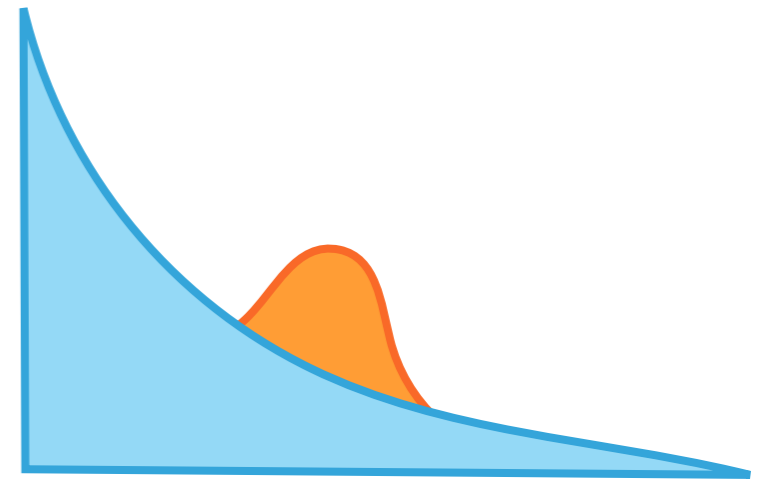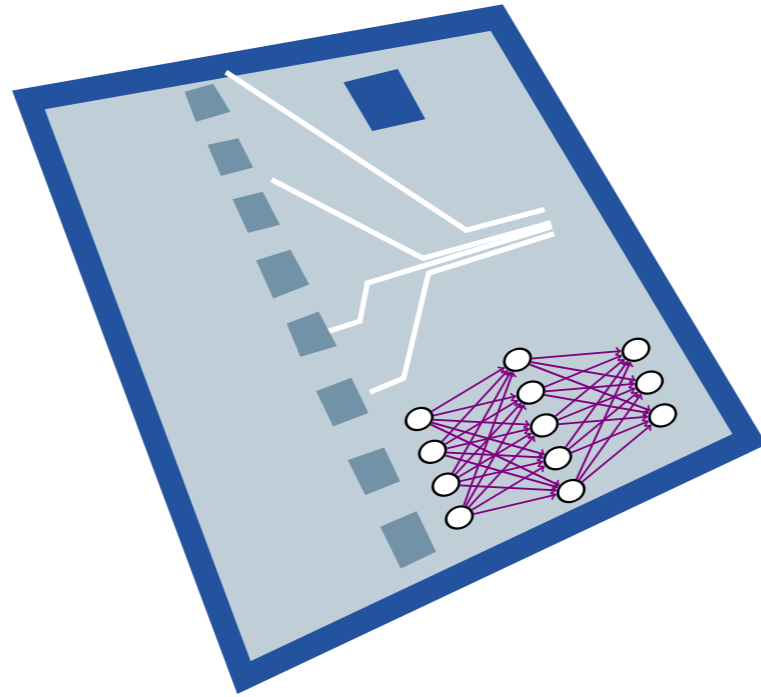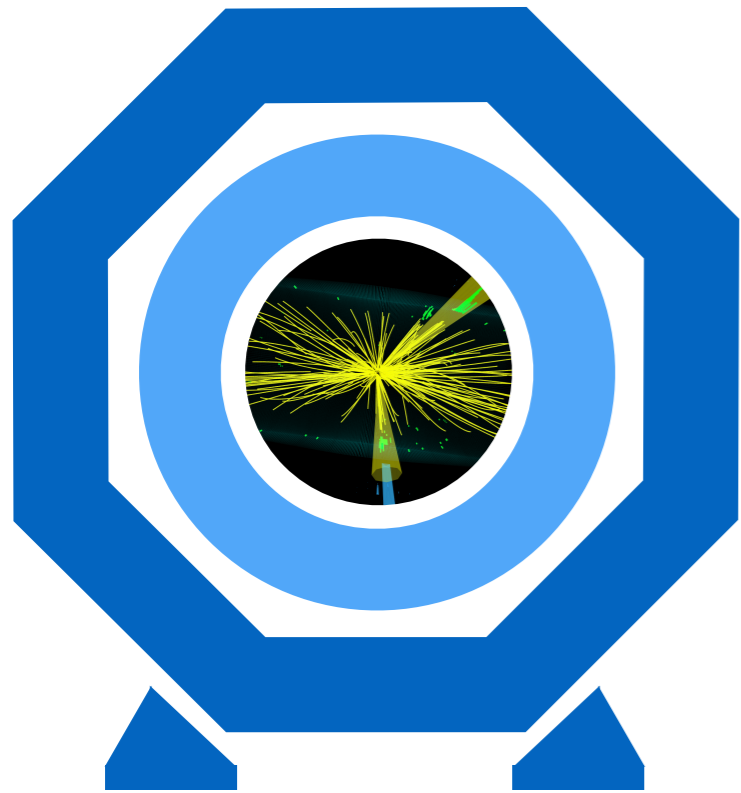- The current sequence of full Velo, primary vertices, full UT, and SciFi decoding runs on an NVIDIA V100 at 112 kHz

**LHCb Upgrade Trigger Diagram**

**30 MHz inelastic event rate (full rate event building)**

**Software High Level Trigger**

Full event reconstruction, inclusive and exclusive kinematic/ geometric selections

**Buffer events to disk, perform online detector calibration and alignment**

Add offline precision particle identification and track quality information to selections

Output full event information for inclusive triggers, trigger candidates and related primary vertices for exclusive triggers

**10 GB/s to storage**

Event builder node

Raw data

Selection decisions

Raw data is decoded on the GPU

| Velo | Decode raw banks | Clustering | Tracking | Find primary vertices |
| UT | Decode raw banks | | Tracking | |
| SciFi | Decode raw banks | | Tracking | Kalman |
| Muon | Decode raw banks | | Matching | |
| Selections | | | | |

6

▸ Particle physics experiments face **extreme trigger challenges** in the coming years

▸ Exploiting **new algorithms**, **new hardware**, and **machine learning** will be key to the success of next-gen experiments

▸ Open questions:

    ▸ With more sophisticated algorithms at earlier trigger, how do we ensure performance/safety? backup triggers?

    ▸ What community tools do we need to deploy ML at the trigger?

    ▸ Which co-processors are best suited to which tasks for the high-level trigger?

    ▸ How do we incorporate timing information at the trigger level?

    ▸ What are the physics use-cases for L1 scouting at 40 MHz?

    ▸ What can we do with the new trigger hardware capabilities which we aren't thinking about?

    ▸ L1 gives us a fundamental limitation but is there more we can exploit at the HLT?

    ▸ Can we make a (realistic) wish-list for triggerable characteristics of events?

JAVIER DUARTE (UCSD)

OCTOBER 23, 2019

WEST COAST LHC JAMBOREE, SLAC

# BACKUP

- ‣ Phase-2 : improved reconstruction using a Kalman filter
  - ‣ Iterative outer-inner tracking to reconstruct tracks and assign track $p_T$ (as offline)
- ‣ Both PV constrained and unconstrained tracks: displaced standalone muons

More and better information available in the Level-1 trigger!

What can we do with it?