# Real-time Artificial Intelligence
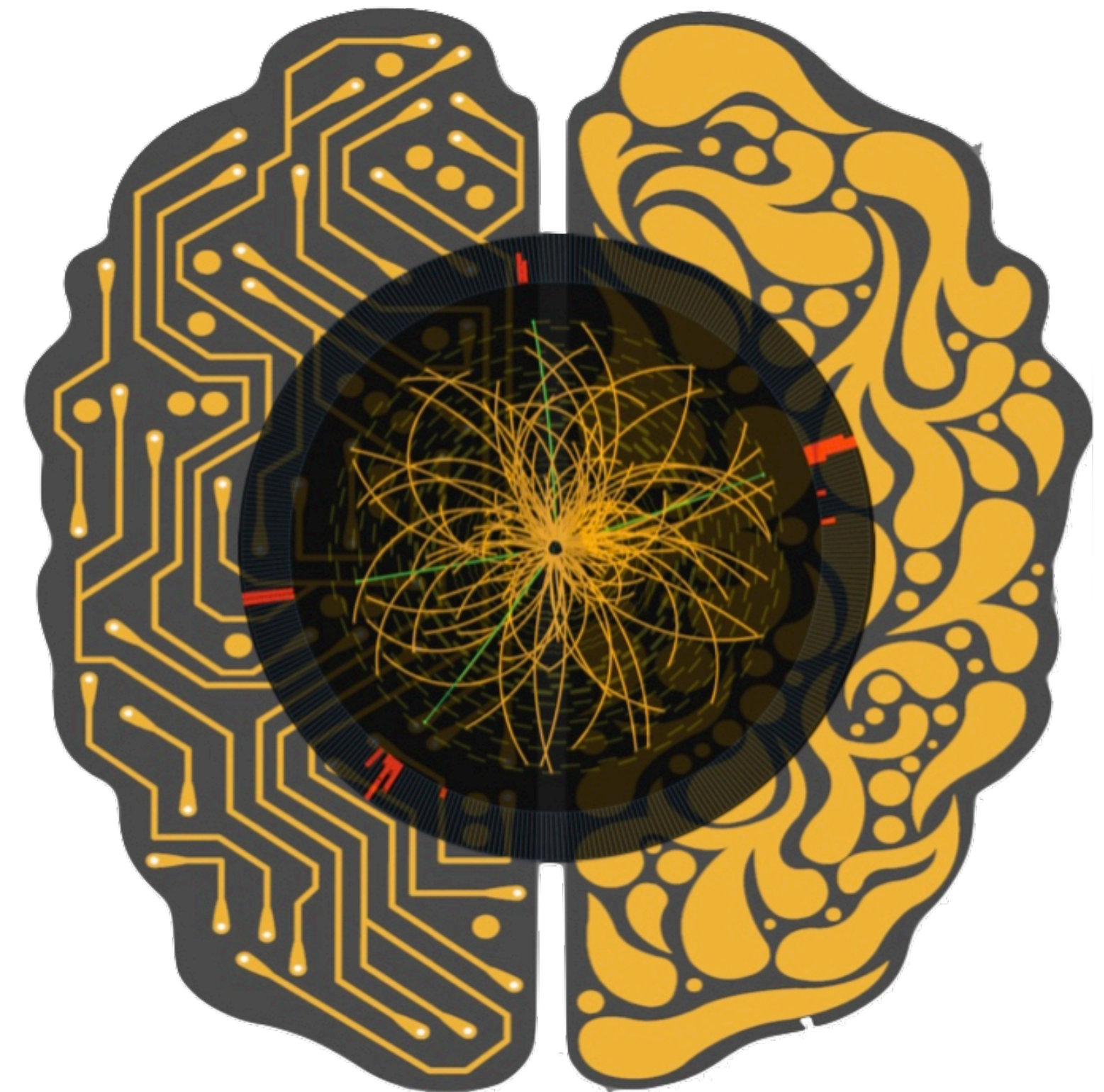## — *with heterogeneous compute*

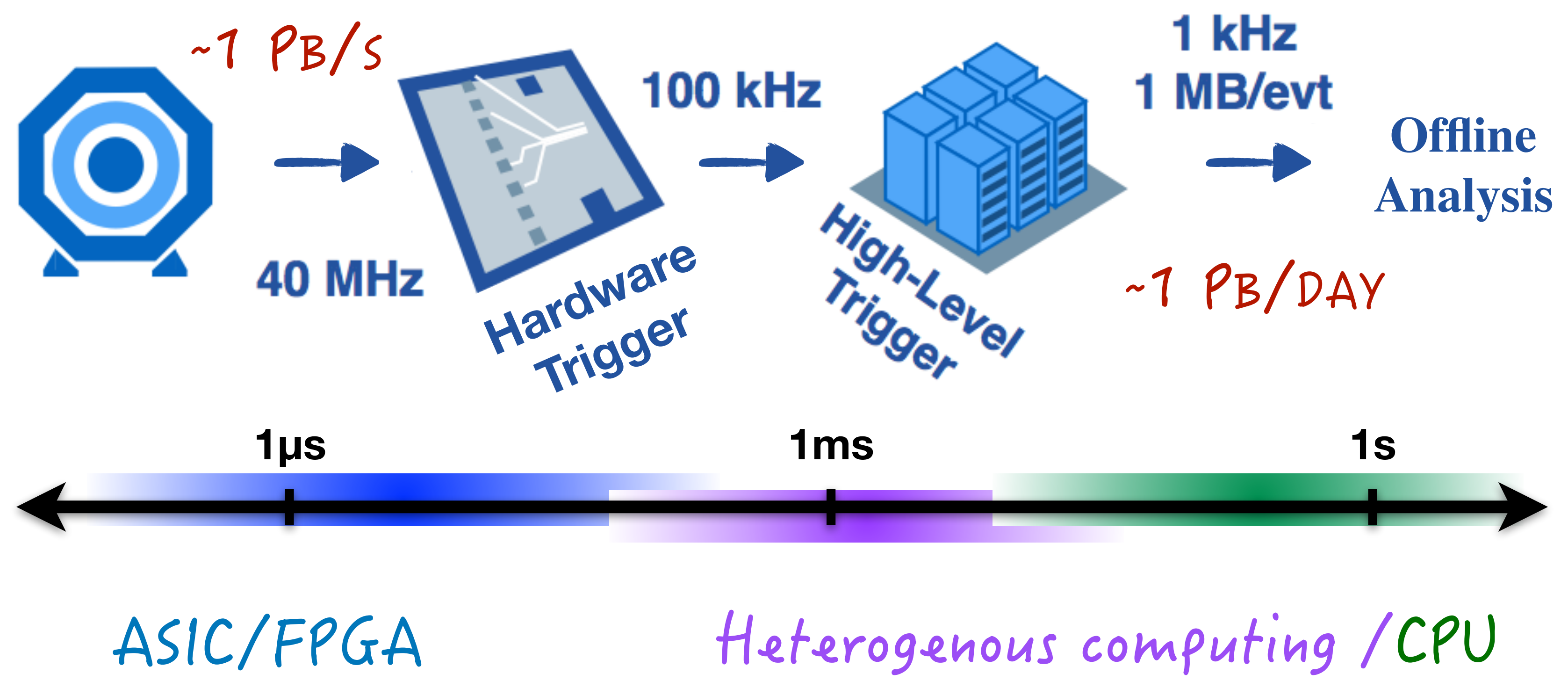**Mia Liu**
**Purdue University**
**Oct. 13. 2020**
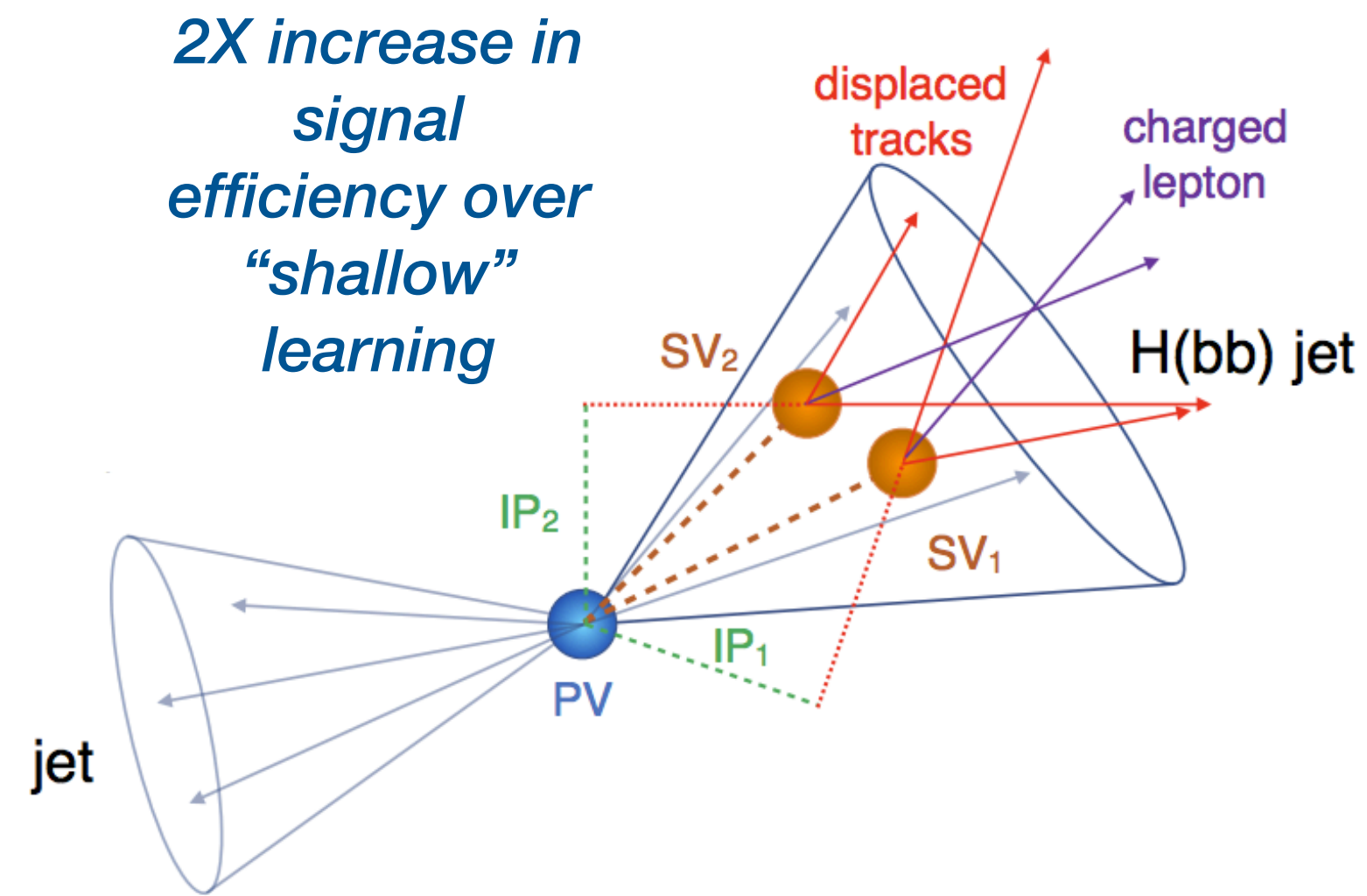**IEEE Real-Time Conference**
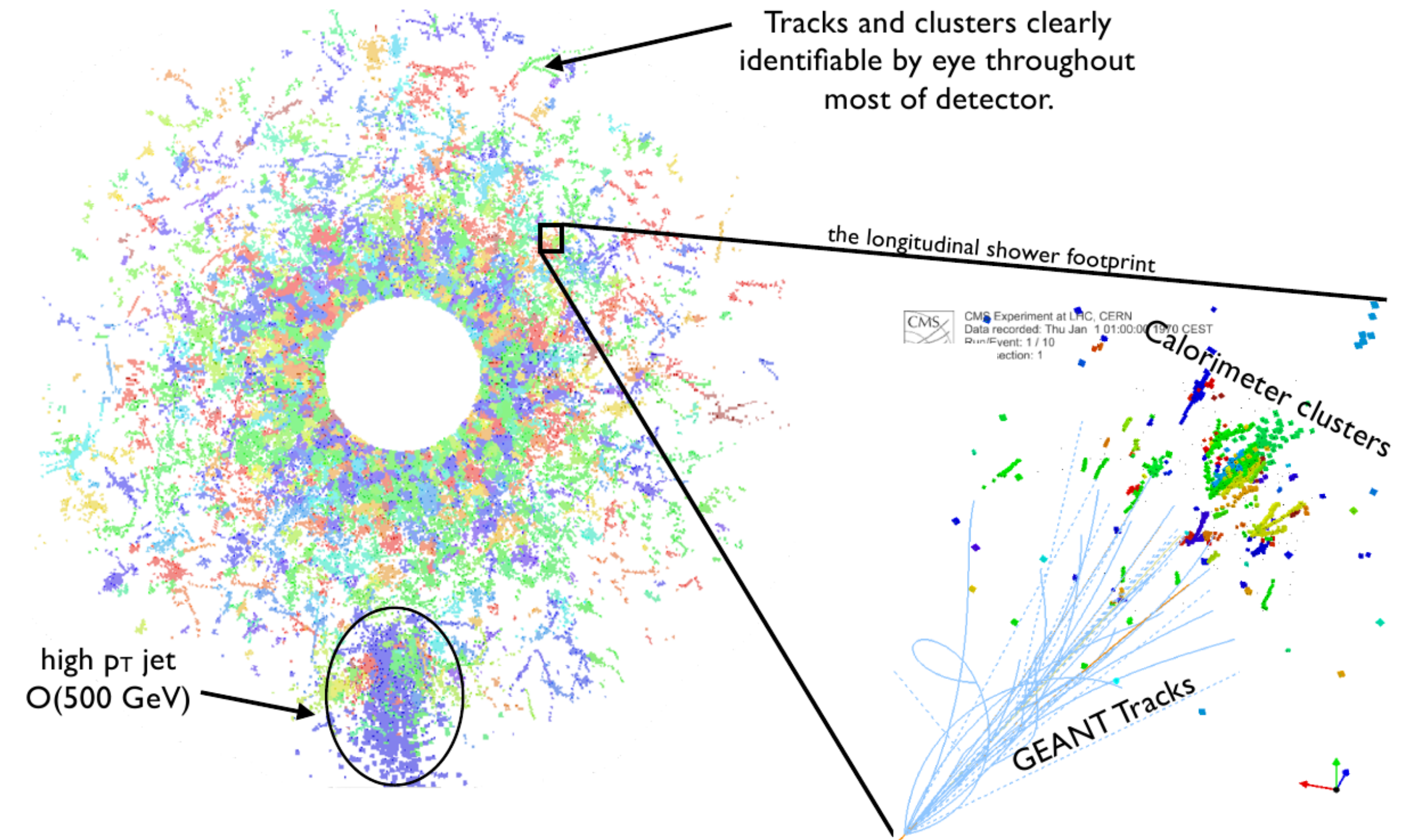
# Data processing in Particle Physics



CMS as an example

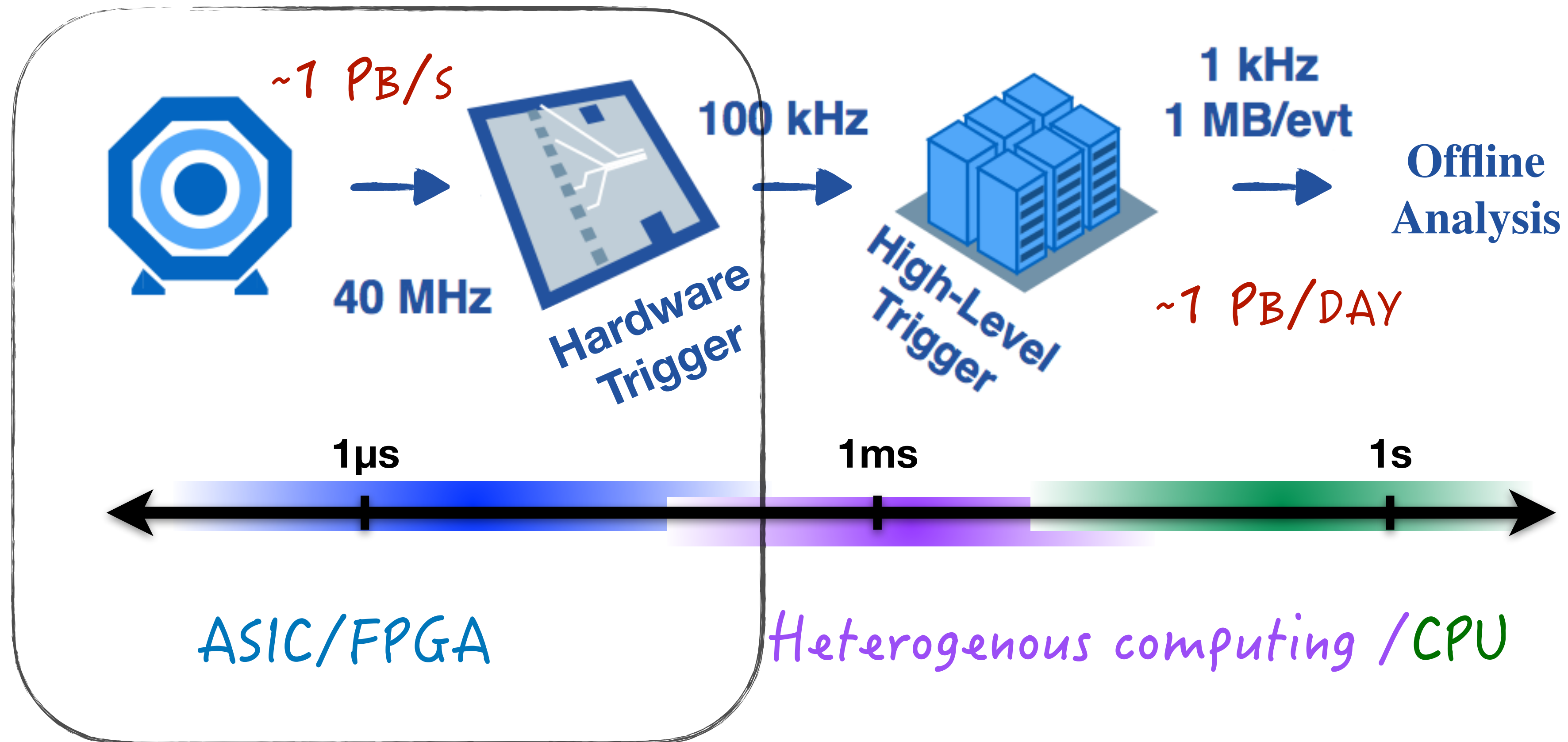# challenges

A quest for accelerated Machine Learning inference

136PU event (2018)

*2X increase in signal efficiency over "shallow" learning*

displaced tracks

charged lepton

$SV_2$

$SV_1$

H(bb) jet

$IP_2$

$IP_1$

PV

jet

$\nu_\mu$ CC

$\nu_e$ CC

NC

$\nu_\tau$ CC

*Deep learning improvement:*

*Effective 40% increase in NOvA active volume*

Tracks and clusters clearly identified by eye throughout

Observed (Q, U)

Reconstructed

ResUNet

O(500 GeV)

GEANT Tracks

*Heavy flavor jet tagging*

*CMS High Granularity Calorimeter Reconstruction*

y frontier: DUNE

est liquid argon detector designed

Accelerated machine learning opens up AI application domain in real-time system and offers novel solutions to computing challenges. See our white paper.
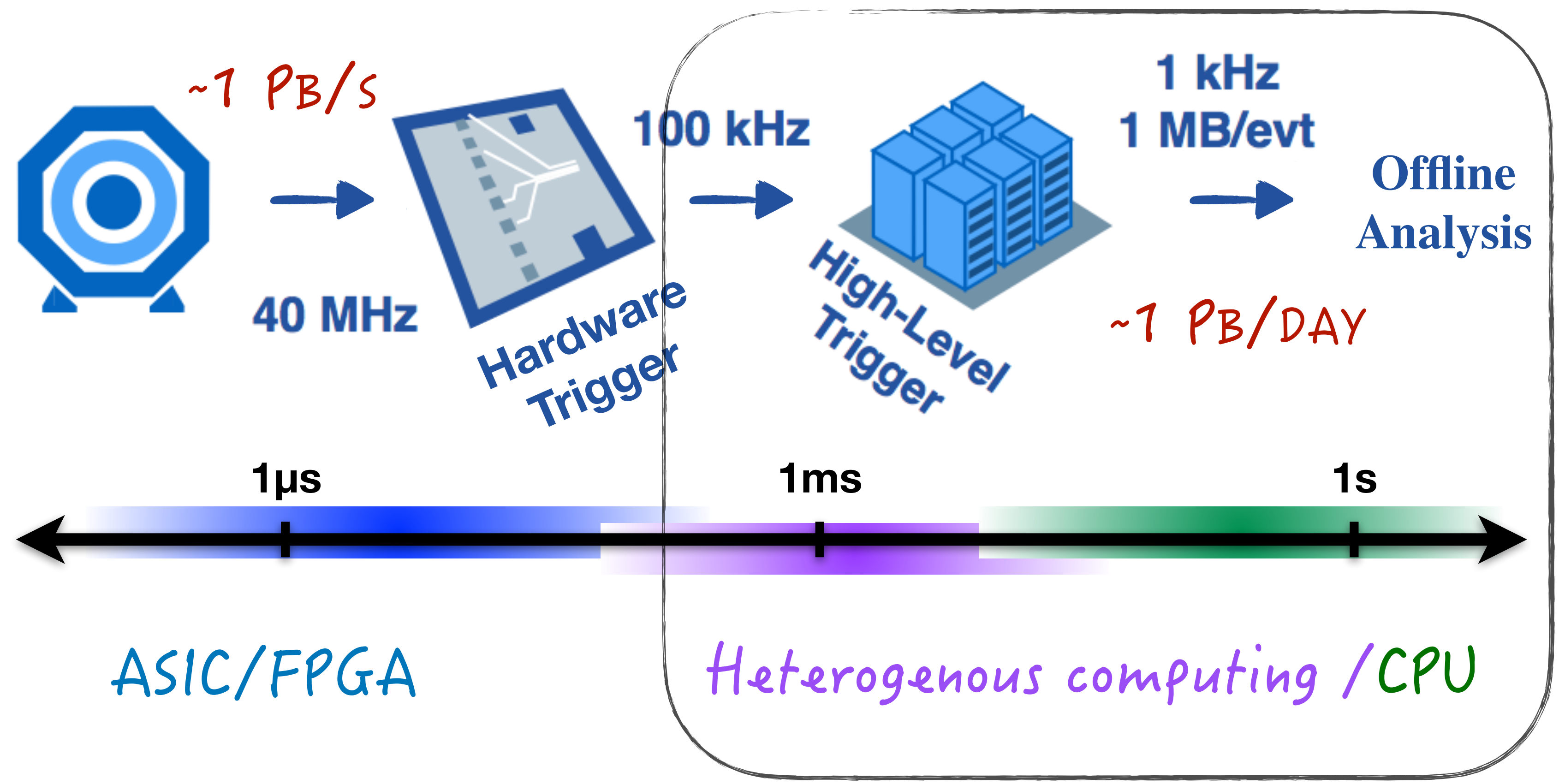
channels, 1 ms integration

# Accelerated ML in embedded systems



~1 PB/S

40 MHz

Hardware Trigger

100 kHz

High-Level Trigger

1 kHz
1 MB/evt

Offline Analysis

~1 PB/DAY

1μs      1ms      1s

ASIC/FPGA

Heterogenous computing /CPU
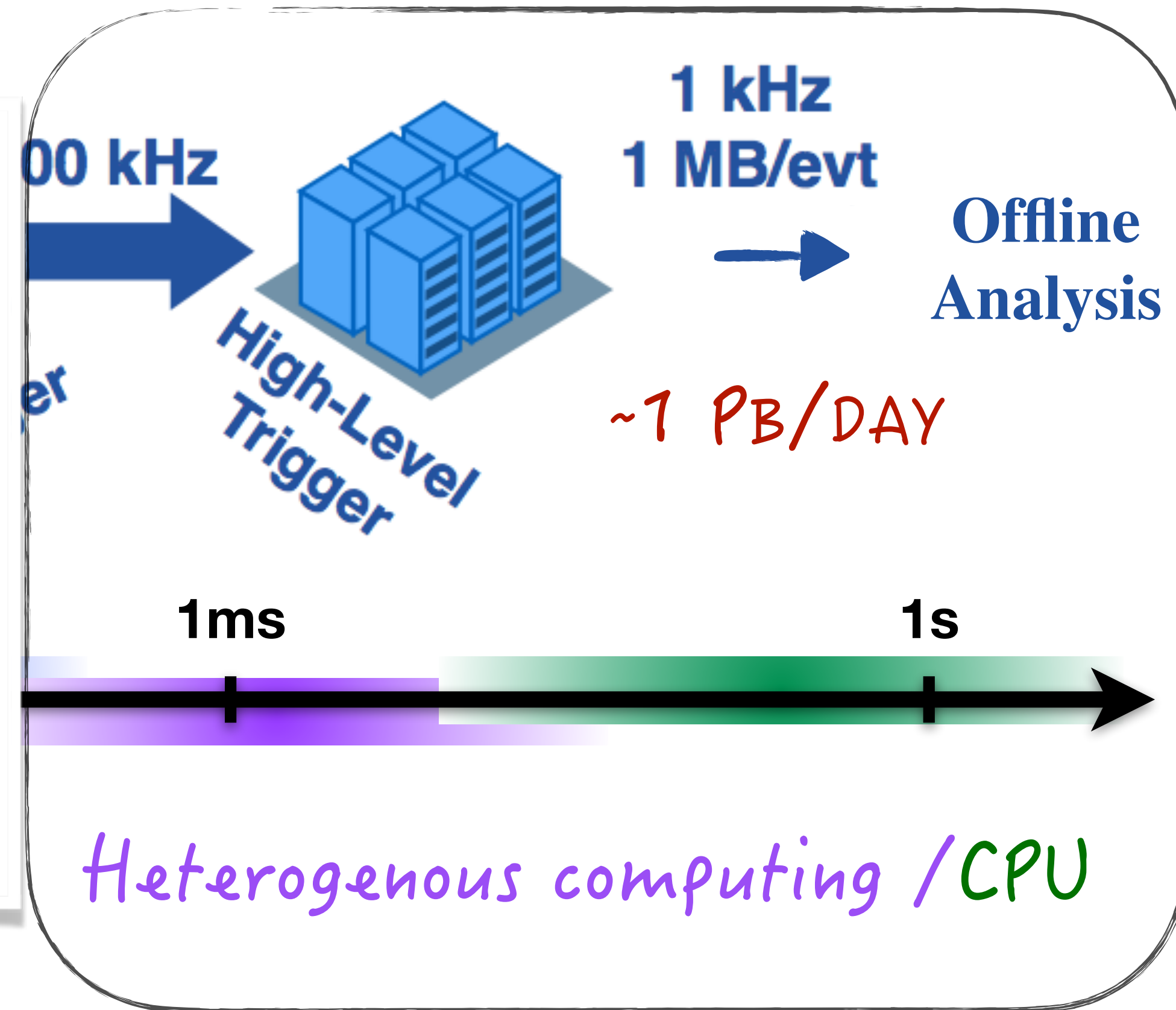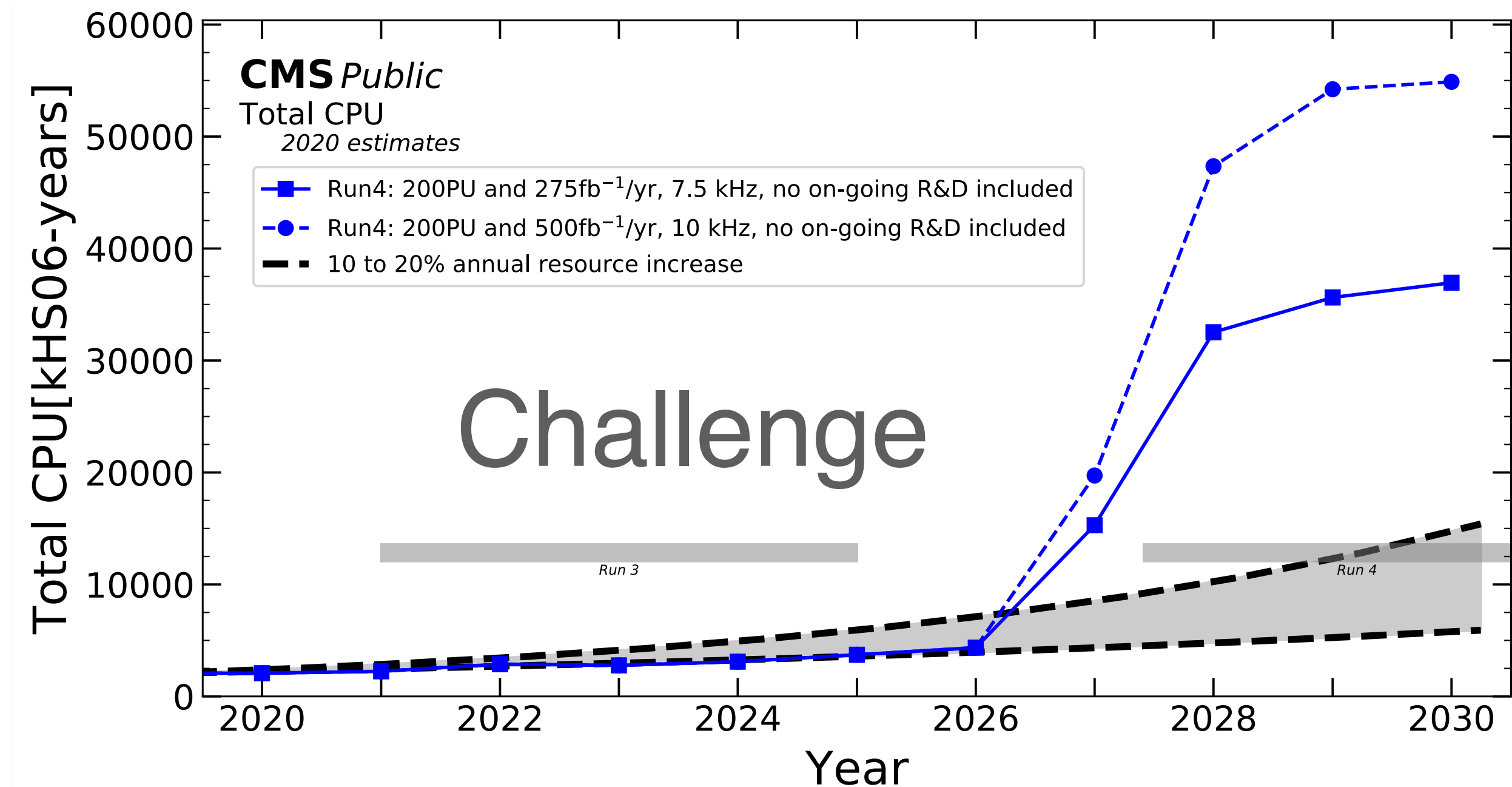
See Nhan's talk on Monday:
Real-time machine learning in embedded systems for
particle physics

# Accelerated ML for HLT/offline



~1 PB/S

40 MHz
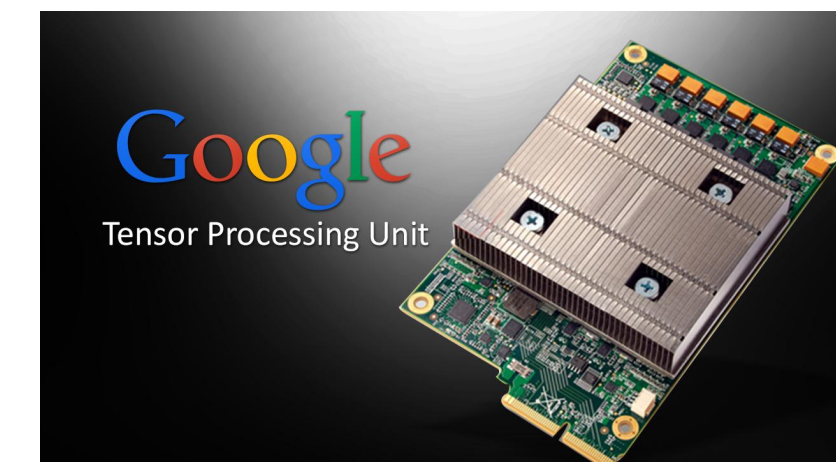
100 kHz

**Hardware Trigger**
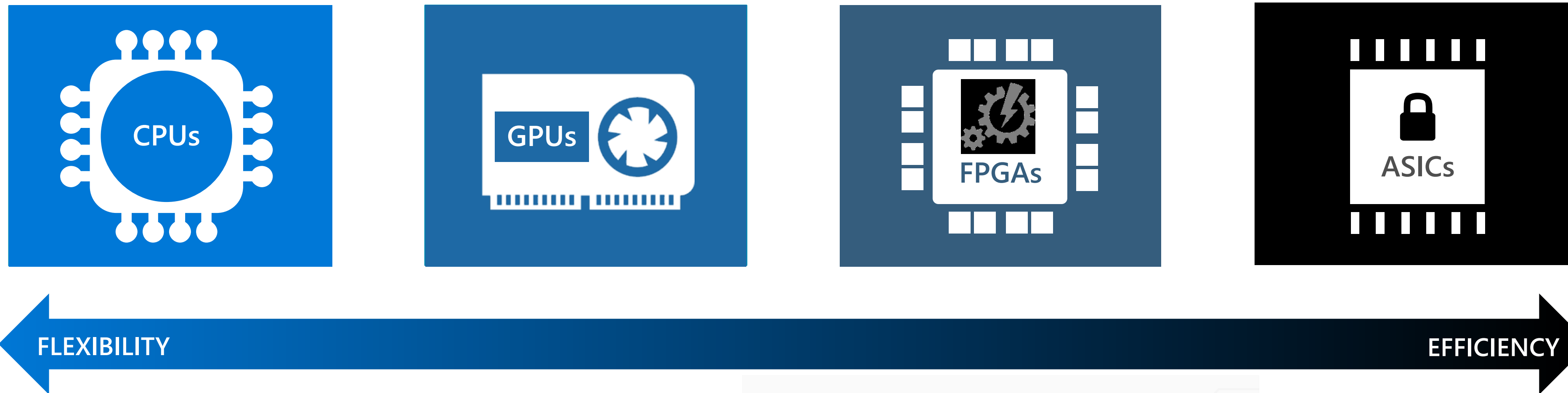
**High-Level Trigger**

1 kHz
1 MB/evt

**Offline Analysis**

~1 PB/DAY

1μs

1ms

1s

ASIC/FPGA

Heterogenous computing /CPU

I will focus on this

# No faster CPUs for free

**CMS** *Public*
Total CPU
*2020 estimates*

- ■ — Run4: 200PU and 275fb$^{-1}$/yr, 7.5 kHz, no on-going R&D included
- ● -- Run4: 200PU and 500fb$^{-1}$/yr, 10 kHz, no on-going R&D included
- - - 10 to 20% annual resource increase

Total CPU[kHS06-years]

60000
50000
40000
30000
20000
10000
0

Challenge

Run 3

Run 4

2020  2022  2024  2026  2028  2030

Year

1 kHz
1 MB/evt

00 kHz

High-Level Trigger

Offline Analysis

~1 PB/DAY

1ms                1s

Heterogenous computing /CPU

**ATLAS** Preliminary
2020 Computing Model - CPU

Run 3 (μ=55)     Run 4 (μ=88-140)     Run 5 (μ=165-200)

[MHS06-years]

80
70
60

○ Baseline
▲ Conservative R&D
▼ Aggressive R&D

e opportunities and challenges.
ent developments… with a bias

# *#Trending in Industry*: Heterogeneous Computing
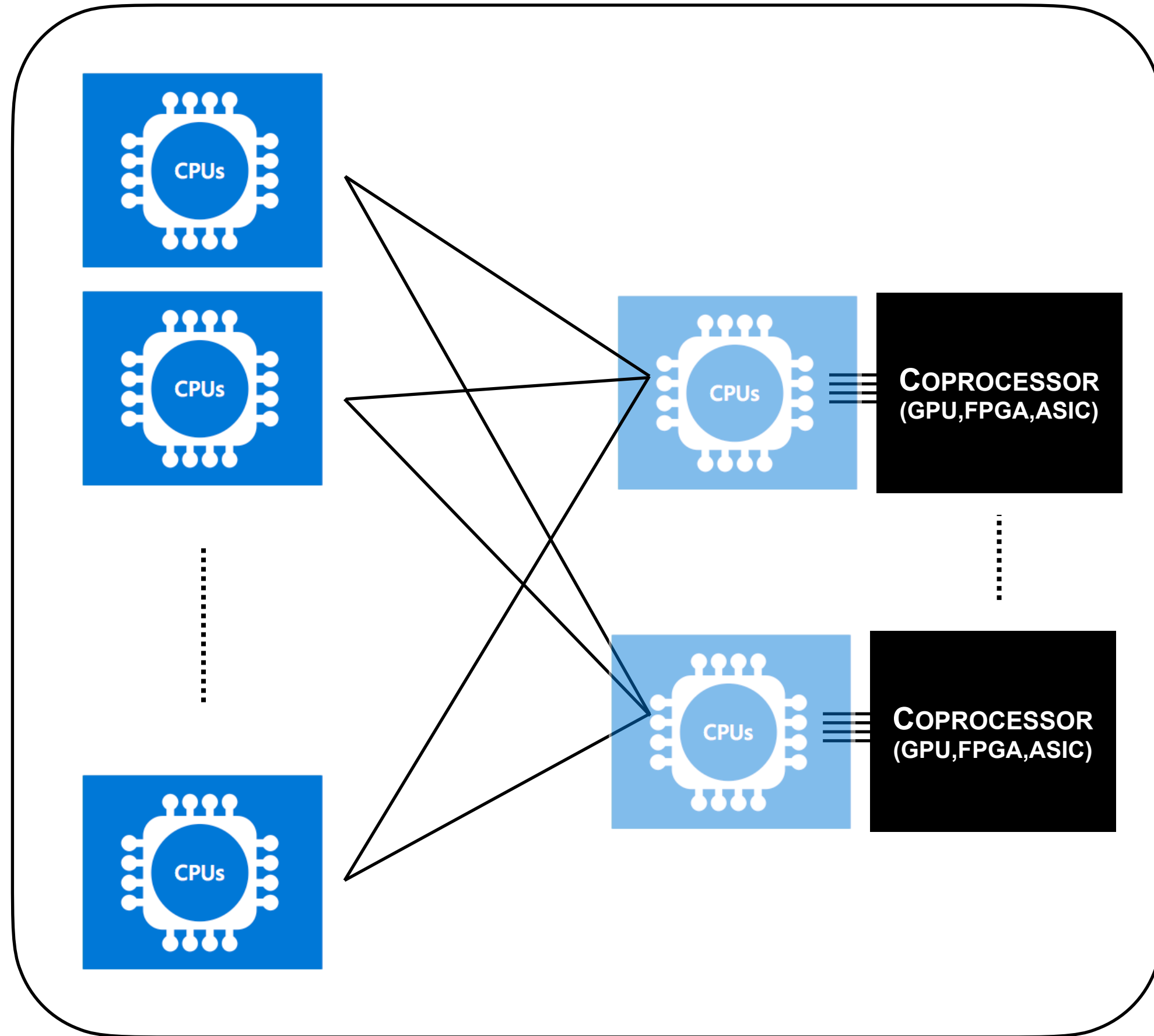


FLEXIBILITY ← → EFFICIENCY

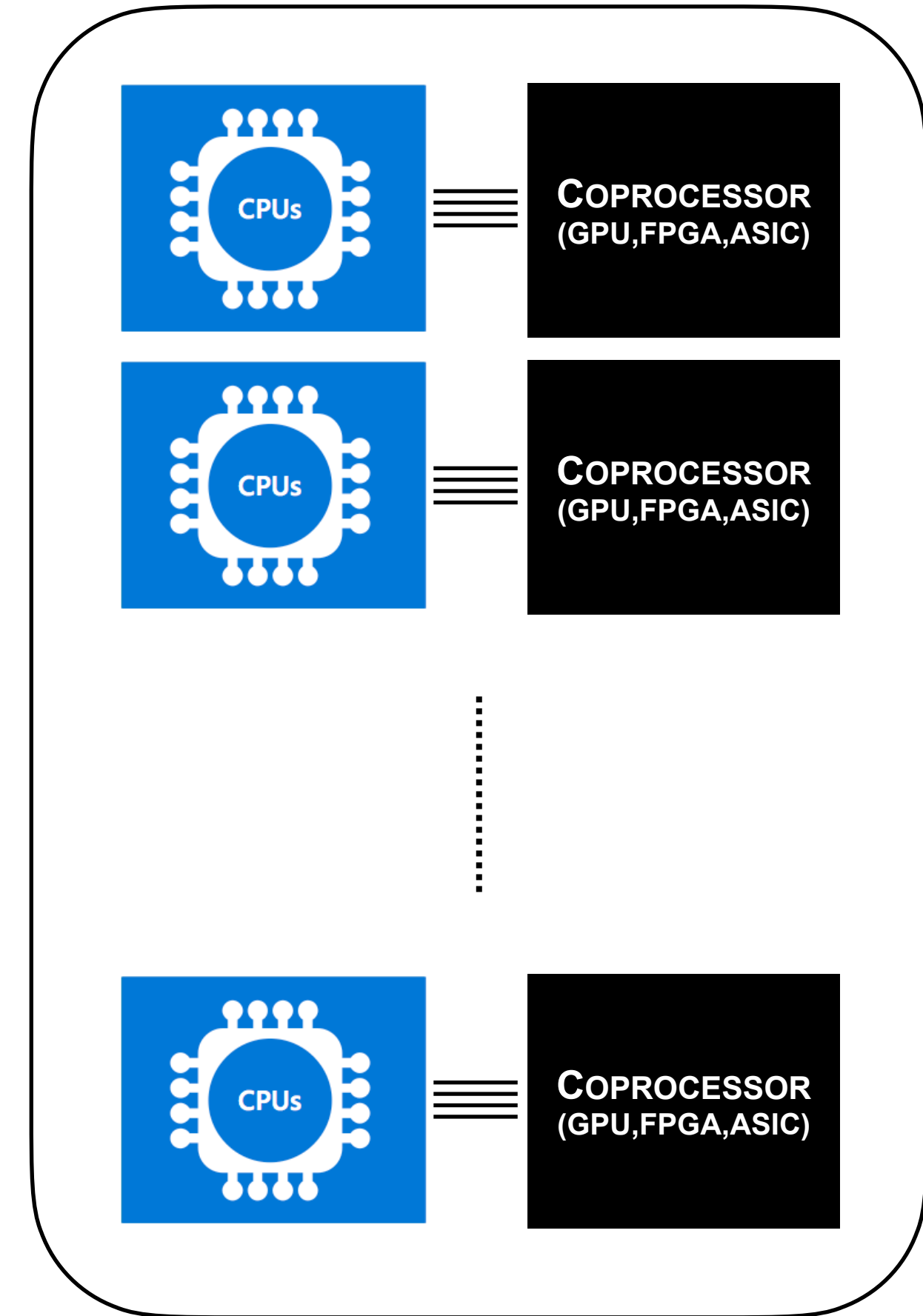Advances driven by
**big data explosion
& machine learning**

# #Heterogeneous Computing Paradigm

**Pros:**
scalable algorithms
scalable to the grid/cloud
Heterogeneous heterogeneity (mixed hardwares)

**Pros:**
less system complexity
no network latency

full story

Fermilab-led team tests Azure AI
for particle physics data challenge
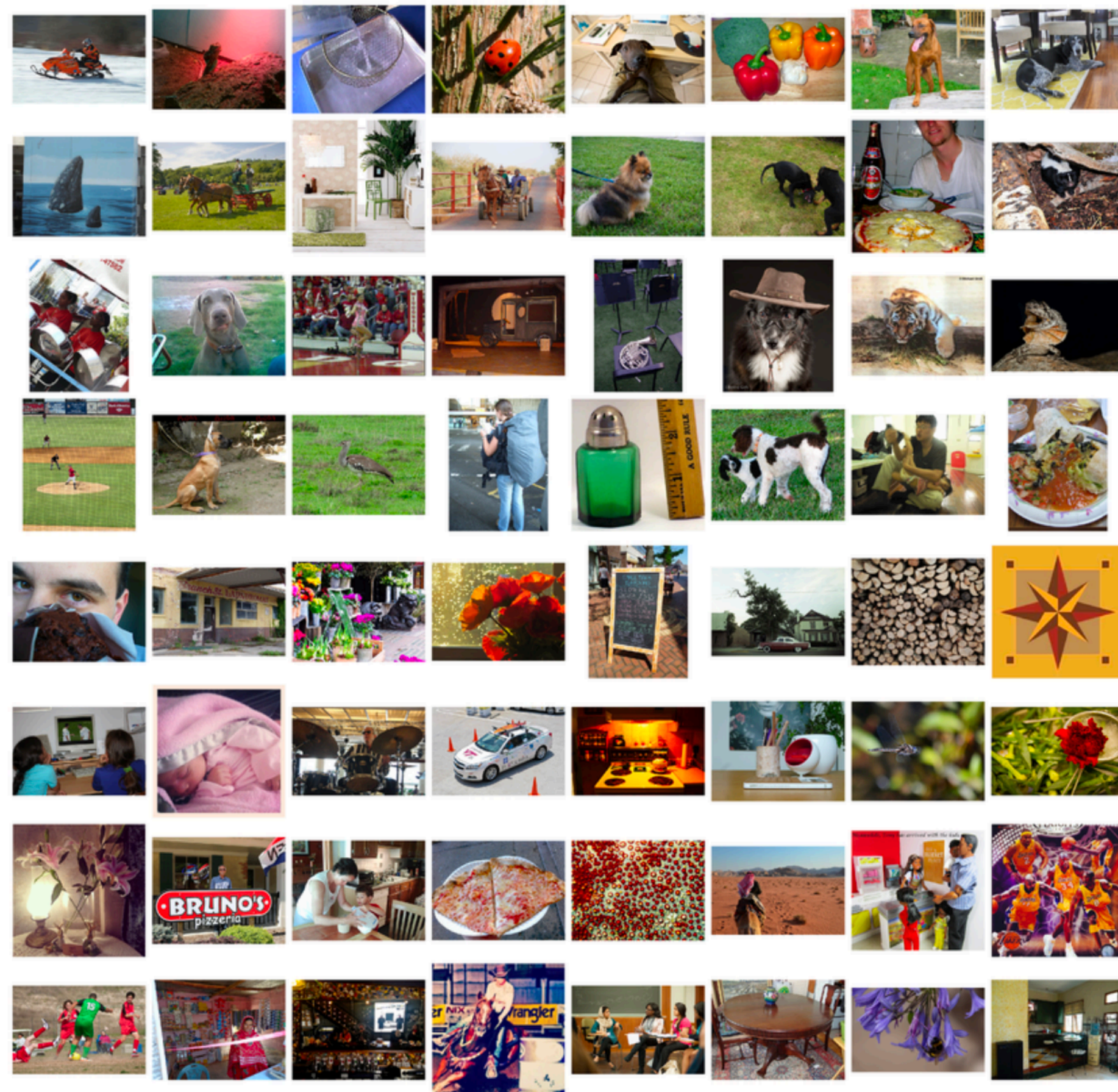
**Brainwave**

Microsoft

FPGA

**Question:**
Can we/How can we take advantage of
**heterogenous computing as-a-service** for our big
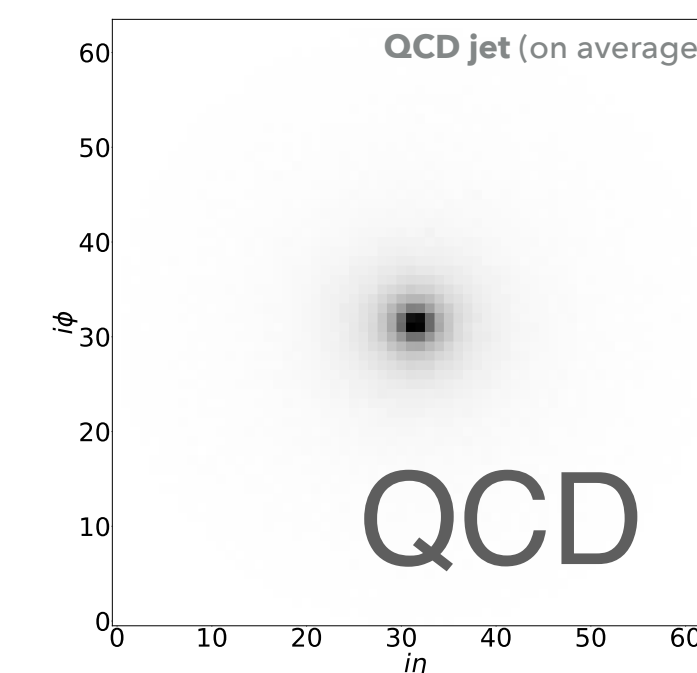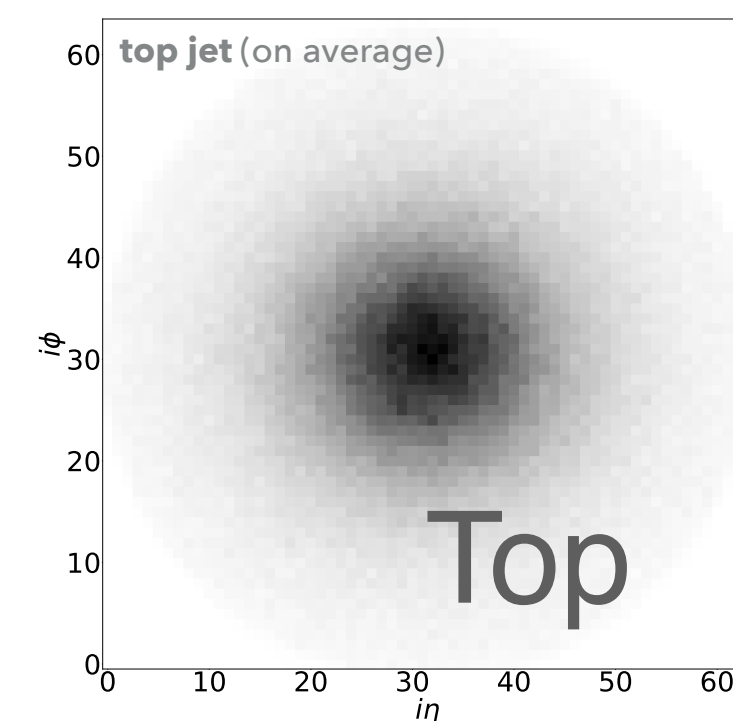data problems?

SONIC

Published in CSBS

# 'Teach' Res-Net 50 about particle physics

Res-Net 50 (25M parameters)

Featurizer

FPGA

Classifier
**CPU**

rain with top
CD labels

Top

QCD

1000 classes
(cats, dogs…)

# Quantized Res-Net 50 performance

**Quantization matters:**

- Floating point—> Quantized model brainwave's implementation of ResNet50 on FPGA
  - Loss in performance

- Re-train the model with fixed precision regains the performance



Legend:
- Floating point: AUC = 98.0%, acc. = 90.1%, $1/\varepsilon_B$ = 671
- Quant.: AUC = 97.5%, acc. = 84.1%, $1/\varepsilon_B$ = 415
- Quant., f.t.: AUC = 98.2%, acc. = 93.0%, $1/\varepsilon_B$ = 971
- Brainwave: AUC = 98.2%, acc. = 92.6%, $1/\varepsilon_B$ = 935
- Brainwave, f.t.: AUC = 98.3%, acc. = 93.5%, $1/\varepsilon_B$ = 1000

Better

Axes: Background efficiency (y) vs Signal efficiency (x)

# Is it faster? Inference speed

Test integrated in CMS software stack

External processing

*FPGA, GPU, etc.*

EVENT DATA

CALLBACK

CMSSW module

*acquire(*

*produce(*

Speed of light→10 ms

| | Inference time |
|---|---|
| **local** | **10 ms** ~2 ms on FPGA classifying, I/O HLT |
| **remote** | **60 ms** (includes travel latency) (10/100) faster than CPU-only computations |

# Computing: data throughput

**Fermilab**

Single Brainwave service

**JetImageProducer**

**JetImageProducer**

Network

5000 images

**JetImageProducer**



**Max data throughout: 600-700 images/sec**, Comparable with V100 GPU (with large batch sizes).

# SONIC: latest explorations

**GPU-as-a-service at the LHC**

https://arxiv.org/abs/2007.10359

*Hardware platforms*

**GPU-as-a-service for DUNE**

https://arxiv.org/pdf/2009.04509.pdf

# GPU as-a-service with Triton

Standard HEP computing

**Client CPU**

**Client CPU**

**Client CPU**

Network

Load Balancer

AI Inference Cluster (CPU | GPU)

**Triton Inference Server**

**Triton Inference Server**

**Triton Inference Server**

**Triton Inference Server**

AI Model Repository

**Example in neutrino: speedup, saturate GPUs**

https://arxiv.org/pdf/2009.04509.pdf

**Proto-DUNE**

**Largest LArTPC ever built**

**Busy environment: cosmic ray muons & beam**

**Reconstruction chain:**

Noise mitigation,

hit finding,

pandora pattern recognition,

**-> CNN EmTrkMichelId**



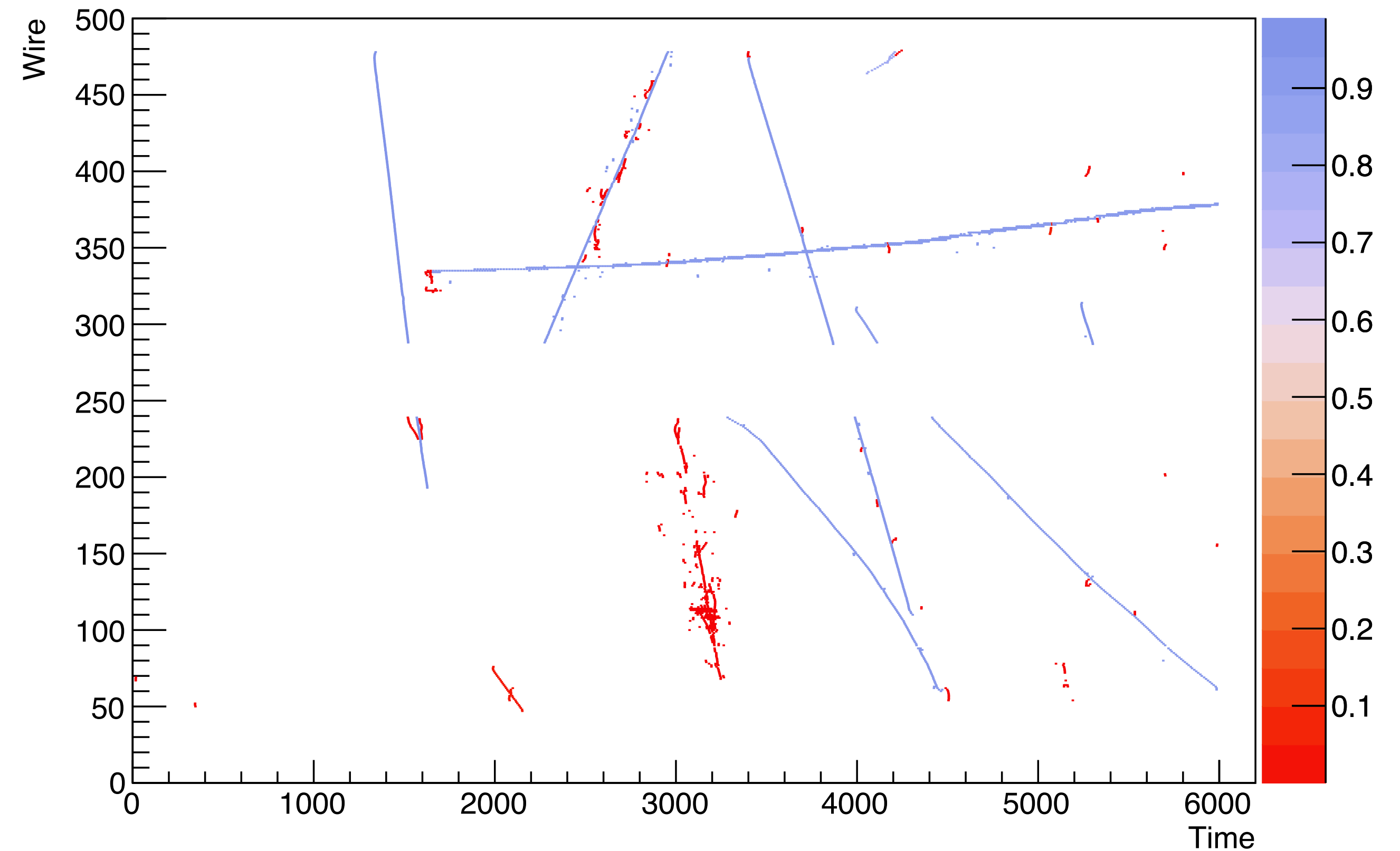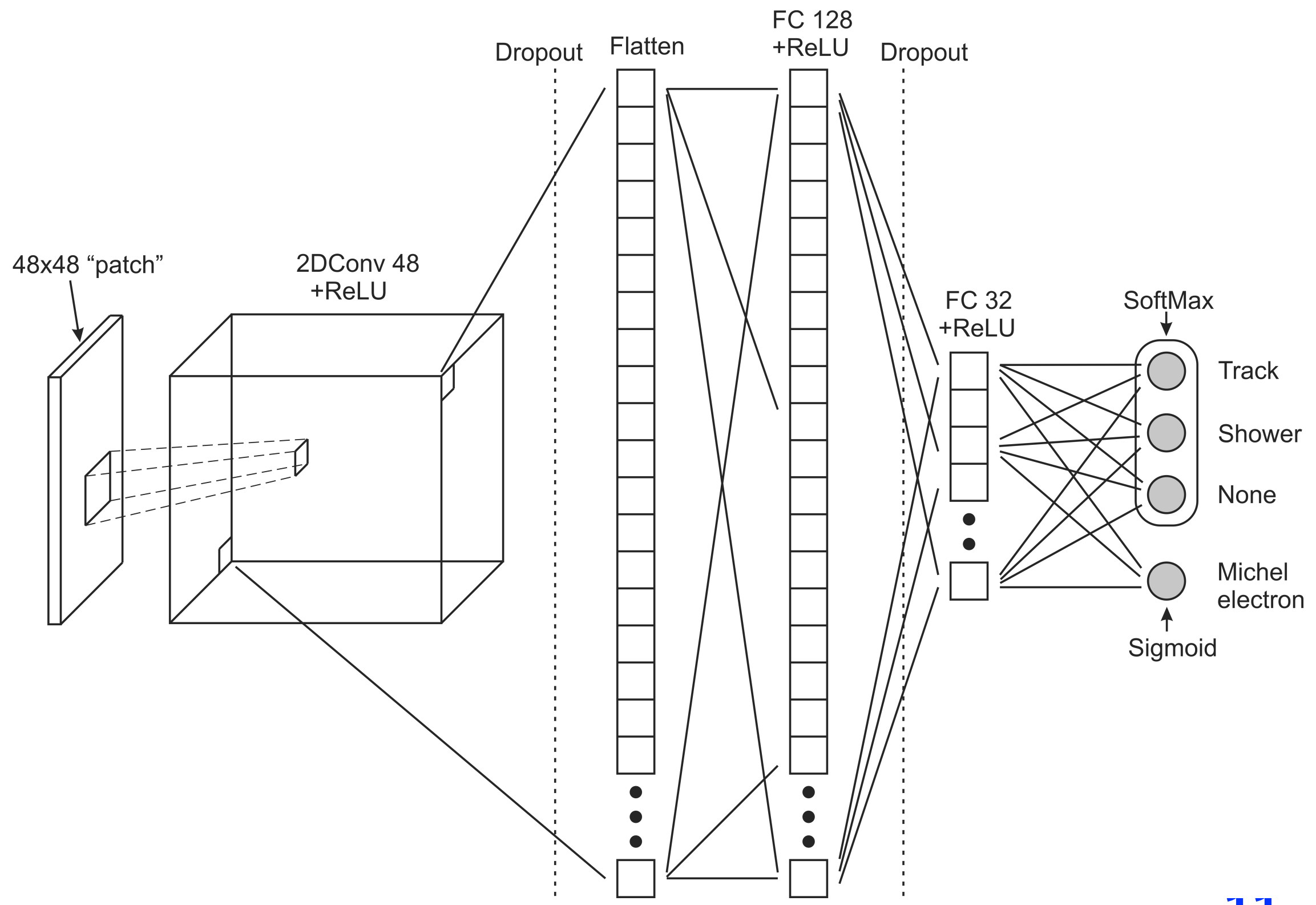Reconstructed ProtoDUNE-SP Event Labelled with CNN Track Score. Run: 5387, Event: 128178, TPC: 1.

**CPU**

| | Wall time (s) | | |
|---|---|---|---|
| ML module | non-ML modules | | Total |
| 220 | 110 | | 330 |

*Wang et al.* **GPU-accelerated ML inference aaS for computing in neutrino experiments**

| CPU type | fraction (%) |
|---|---|
| AMD EPYC 7502 @ 2.5 GHz | 11.7 |
| AMD Opteron 6134 @ 2.3 GHz | 0.6 |
| AMD Opteron 6376 @ 2.3 GHz | 4.6 |
| Intel Xeon E5-2650 v2 @ 2.6 GHz | 30.8 |
| Intel Xeon E5-2650 v3 @ 2.3 GHz | 5.2 |
| Intel Xeon E5-2670 v3 @ 2.3 GHz | 7.3 |
| Intel Xeon E5-2680 v4 @ 2.4 GHz | 17.3 |
| Intel Xeon Gold 6140 @ 2.3 GHz | 22.6 |

**Single GPU server (NVIDIA T4)**

**Table 1.** CPU types and distribution for the grid worker nodes used for the "big-batch" clients (see text for more details).

| | Wall time (s) | | |
|---|---|---|---|
| ML module | non-ML modules | Total | |
| 220 | 110 | 330 | |

**Table 2.** The average CPU-only wall time per job for the different module categories.

**~20-11s**

### 3.2 Server-side performance

To get a standardized measure of the performance, we first use standard tools for benchmarking the GPU performance. Then we perform a stress test on our GPUaaS instance to understand the server-side performance under high load.

| | Wall time (s) | | |
|---|---|---|---|
| ML module | non-ML modules | Total | |
| 220 | 110 | 330 | |

$$t_S \sim t_{preprocess} + t_{transmit} + t_{travel} + t_{GPU}$$

*Server standalone performance*

The baseline performance of the GPU server running the *EmTrackMichelId* model is measured using the *perf_client* tool included in the Nvidia Triton inference server package. The tool emulates a simple client by generating requests over a defined time period. It then returns the latency and throughput, repeating the test until the results are stable. We define the baseline performance as the throughput obtained under the saturation of the model on the GPU. We attain this by increasing the client-side request concurrency—the maximum number of unanswered requests by the client—until the throughput saturates. We find that the model reaches this limit quickly at a client-side concurrency of only 2 requests. At this point, the throughput is determined to be $20,000 \pm 2,000$ inferences per second. This corresponds to an event processing time of $2.7 \pm 0.3$ s. This is the baseline expectation of the performance of the GPU server.

**7s**
On CPU preparing NN inputs

**2s**
Based on 2Gbps ethernet bandwidth

**0.4s**
Ping latency between Iowa and FNAL

**1.8s**
Time on the GPU



**CNN EmTrkMichelId**

~20x speedup of EMMichelTrackID module

48x48 "patch" — 2DConv 48 +ReLU — Dropout — Flatten — FC 128 +ReLU — Dropout — FC 32 +ReLU — SoftMax — Track, Shower, None — Michel electron — Sigmoid

# Saturating the GPUs

Left plot legend:
- Model (small Batch)
- Model (big Batch)
- CPU-only (w/o Triton)
- w/ Triton on GKE-4gpu, avg batch size = 234
- w/ Triton on GKE-4gpu, avg batch size = 1692

Right plot legend:
- Model (big batch)
- Model (dynamic batching, big batch)
- CPU-only (w/o Triton)
- w/ Triton on GKE-4gpu, dyn bat Off, avg bat sz = 1692
- w/ Triton on GKE-4gpu, dyn bat On, avg bat sz = 1692

Both plots: x-axis "Number of simultaneous jobs", y-axis "Processing time [seconds]"

$$t_{\mathrm{SONIC}} = (1-p) \times t_{\mathrm{CPU}} + t_{\mathrm{GPU}} \left[ 1 + \max \left( 0, \frac{N_{\mathrm{CPU}}}{N_{\mathrm{GPU}}} - \frac{t_{\mathrm{ideal}}}{t_{\mathrm{GPU}}} \right) \right] + t_{\mathrm{latency}}$$

**GPU saturates**

2.7x speed up of the full ProtoDUNE-SP processing chain

1 GPU can handle 68 CPU

processes simultaneously

# SONIC: latest explorations

**GPU-as-a-service**

https://arxiv.org/abs/2007.10359

*Hardware platforms*

**GPU-as-a-service for DUNE**

https://arxiv.org/pdf/2009.04509.pdf

*Algorithm complexity*

**More benchmarks driven by use cases to test scaling for HLT/offline**

**FACILE**

**DeepCalo***

**ResNet**

**CMS Hadronic Calorimeter channel regression**

**ECAL cluster regression**

**top quark image classification**

**2k parameters**

**10 M parameters**

| | FACILE | DeepCalo* | ResNet |
|---|---|---|---|
| | **HCAL channel regression** | **ECAL cluster regression** | **top quark image classification** |
| CPU | 16 ms | 75 ms | ~1 s◆ |
| GPU as-a-service | 2 ms (GPU) | 0.1 ms | 1-2 ms (GPU/FPGA) |
| **Gain** | **8x (GPU)** | **750x** | **500x** |

**1–24 GPU server**

## Significant gain from dynamic batching

GPU usage: 45%

# Where do we gain?

# SONIC: latest explorations

**FPGA-as-a-service Toolkit**

**GPU-as-a-service**
https://arxiv.org/abs/2007.10359

**hls 4 ml**

Hardware platforms

Open source tools: flexibility

**GPU-as-a-service for DUNE**
https://arxiv.org/pdf/2009.04509.pdf

Algorithm complexity

**More benchmarks driven by use cases to test scaling for HLT/offline**

# hls4ml: accelerating ML on hardware  **fastmachinelearning.org/hls4ml**

Originally designed for LHC triggers applications but broad and growing user base

**VITIS +** hls 4 ml

**FACILE**

**ResNet**

**Xilinx Machine Learning Suite**

| Algorithm | Platform | Number of Devices | Batch Size | Inf./s [Hz] | Bandwidth [Gbps] |
|-----------|----------|-------------------|------------|-------------|------------------|
| FACILE | AWS EC2 F1 | 1 | 16,000 | 36 M | 23 |
| FACILE | Alveo U250 | 1 | 16,000 | 86 M | 55 |
| FACILE | T4 GPU | 1 | 16,000 | 8 M | 5.1 |
| ResNet-50 | AWS EC2 F1 | 8 | 10 | 1400 | 6.7 |
| ResNet-50 | V100 GPU | 8 | 10 | 1,700 | 8.1 |
| ResNet-50 | ASE | 1 | 1 | 460 | 2.2 |
| ResNet-50 | T4 GPU | 1 | 10 | 250 | 1.2 |

**FPGA-as-a-service Toolkit**

# Future prospects

**Integration in full scale production in experiments**

*: processing for full-scale protoDune-SP reconstruction, FACILE@HLT in CMS…scaling with multiple models.*

**AI algorithms suitable for physics data/with domain knowledge embedded**

Graph neural networks, Energy flow networks, see Nhan's talk.

**Hardware awareness in training for co-processor workflow**

*e.g. Brainwave studies explore re-training with quantized version to achieve the best performance in precision.*

*Most studies performed using Cloud services/on-premises clusters*
**SONIC in High Performance Computers (HPCs)**

*: accelerate ML-based simulation/reconstruction etc…*

Data

AI Algorithms

# SONIC: Not limited to ML

*Given a heterogenous computing hardware:*

**re-write physics algorithms for new hardware**

Language: OpenCL, OpenMP,TBB, HLS, …?
Hardware: FPGA, GPU

**re-cast physics problem as a machine learning problem**
Language: C++, Python (TensorFlow, PyTorch,…)
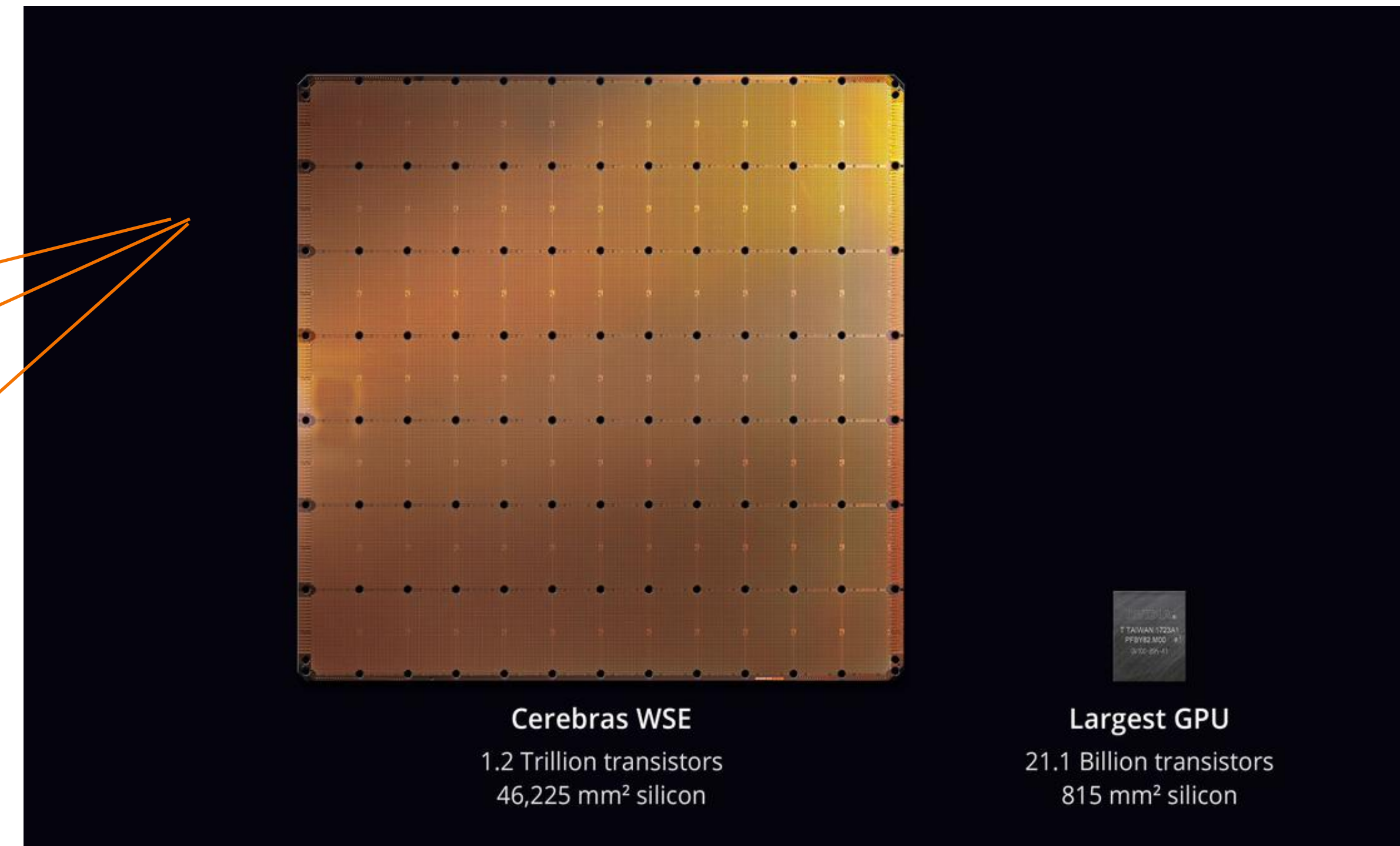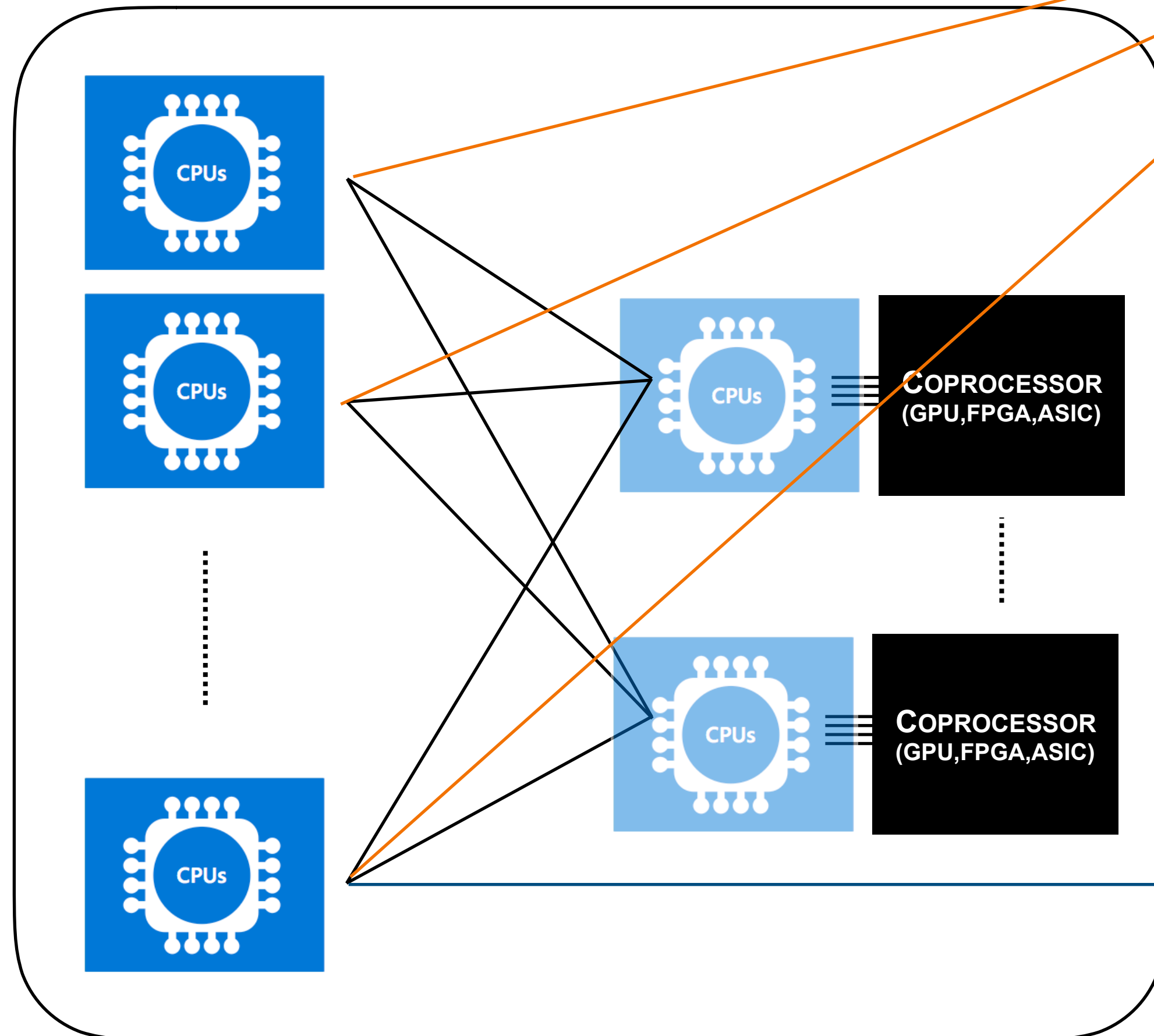
Hardware: FPGA, GPU,ASIC



**Parallelized and Vectorized Tracking Using Kalman Filters**
  • e.g.*On GPUs*

**Tracking with ML**

  • Algorithms parallelizable

  • Solutions with ML e.g., HEP.TrkX.

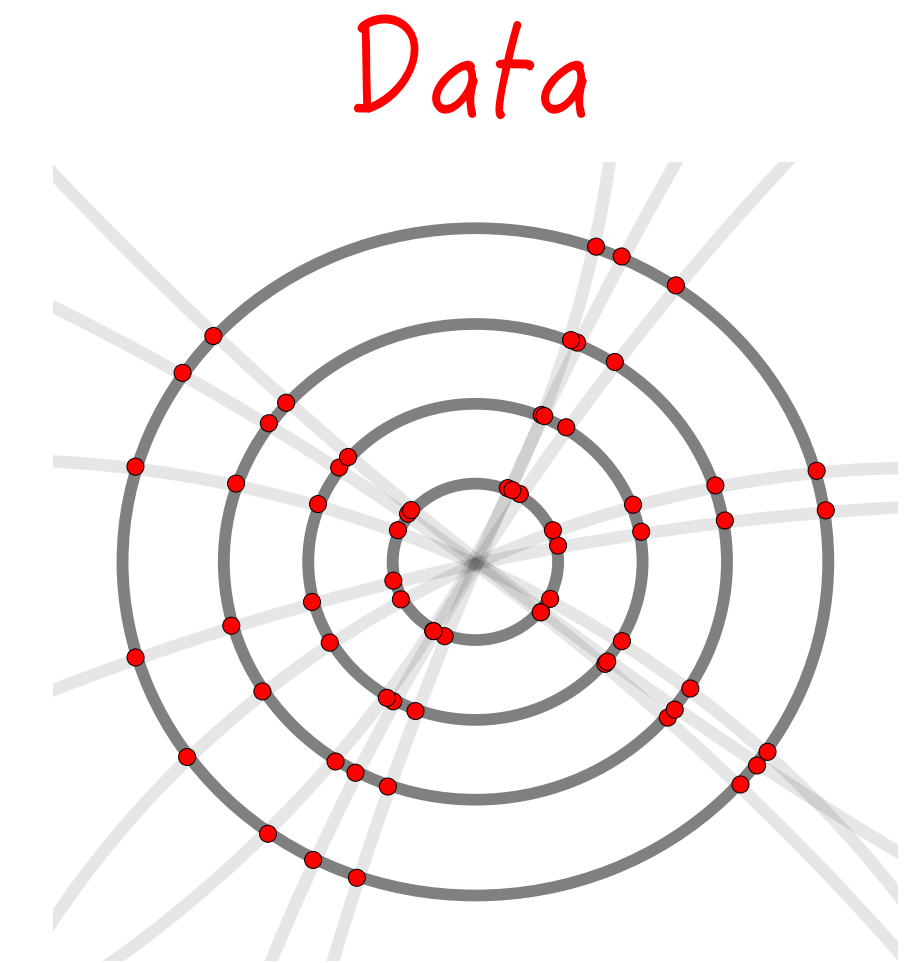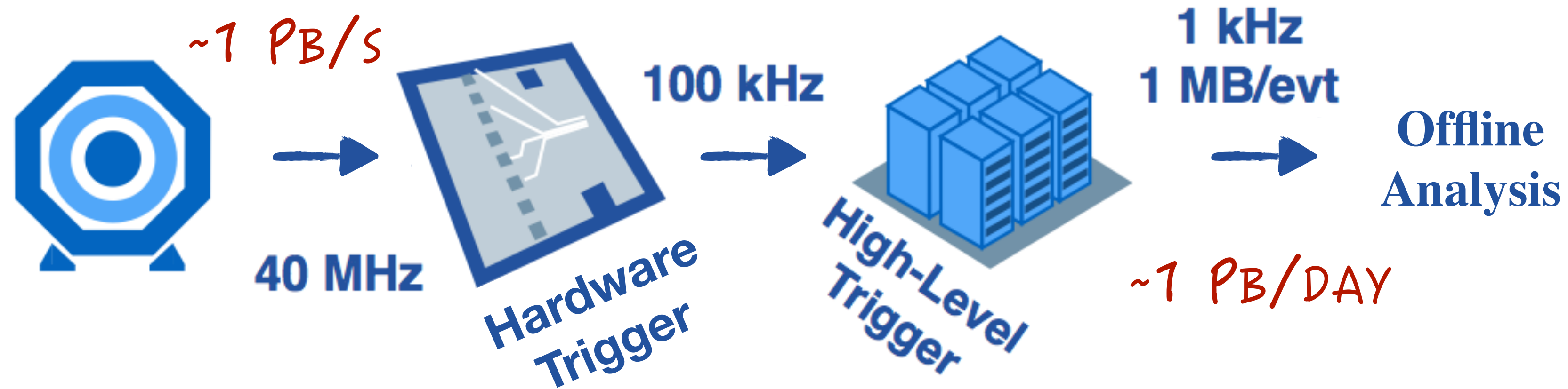**Charged particle Tracking With graph neural networks**

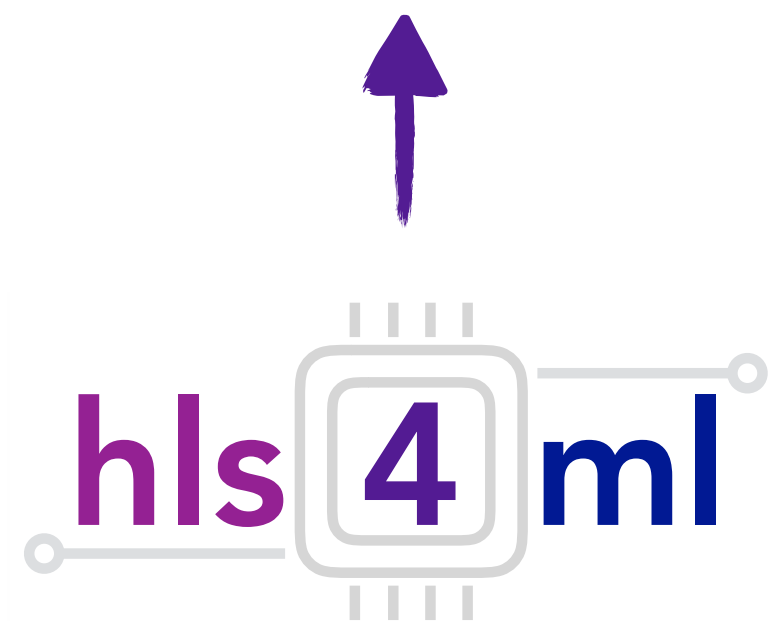# The now/future computing paradigm?

Heterogeneous
computing as-a-service



Emerging technologies ...
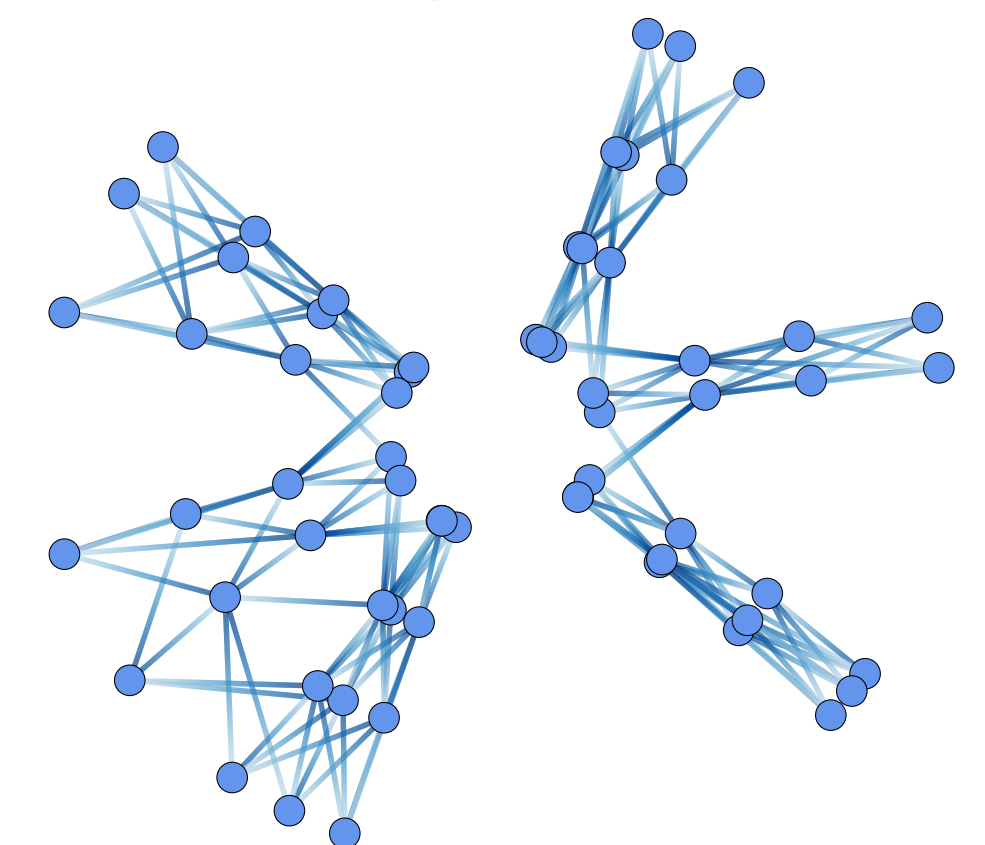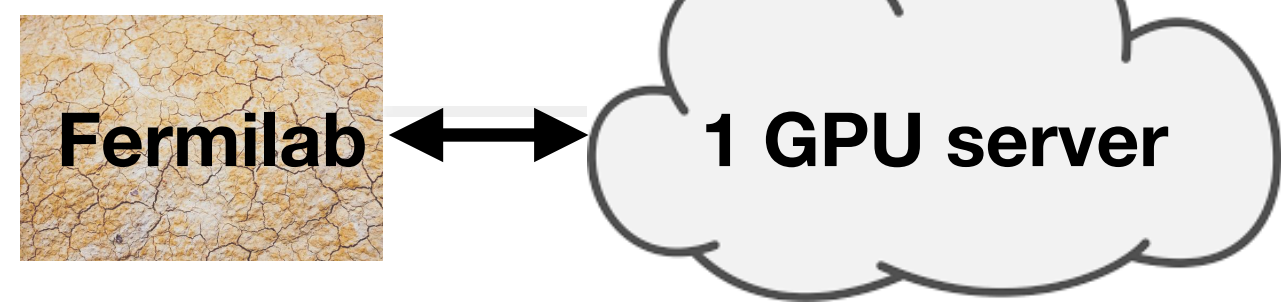
# Accelerated discoveries with Real-Time AI

# Throughput

Fermilab ↔ 1 GPU server

**FACILE**

**DeepCalo**

**ResNet**



*T4:V100*
*16000 batch*

Throughput [events/s]
Simultaneous processes
T4
V100

*10 batch*

Throughput [events/s]
Simultaneous processes

*10 batch*

Throughput [events/s]
Simultaneous processes