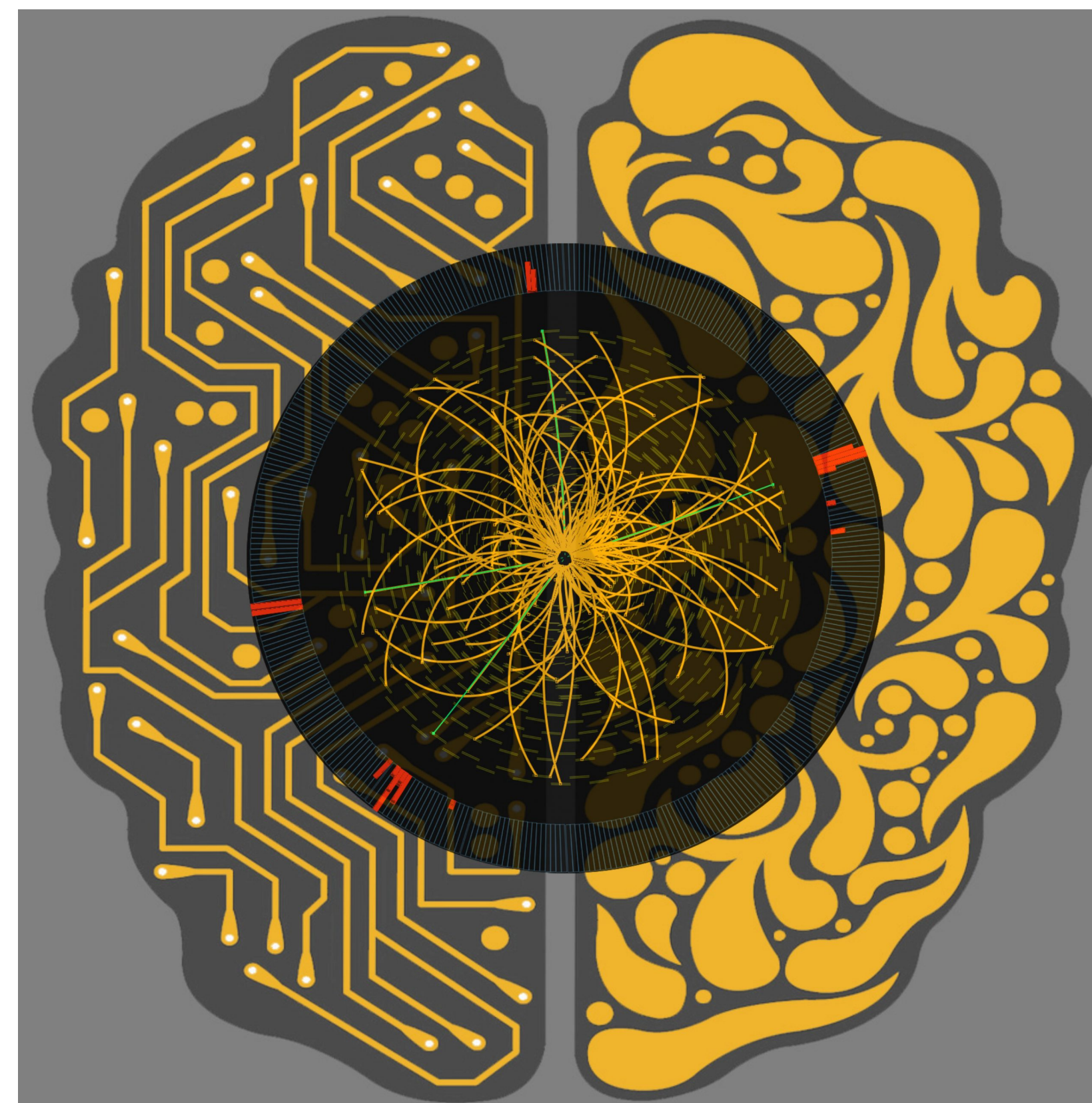


FPGA-accelerated machine learning inference as trigger and computing solutions in particle physics

Mia Liu
FNAL

Fast Machine Learning Workshop

September.11.2019



Big data challenge in particle physics

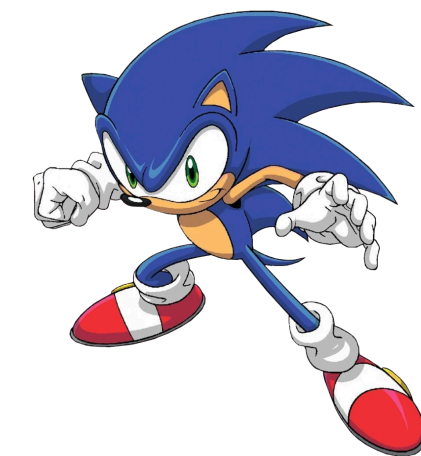
an example:

Trigger and computing challenges@LHC: speed, volume, complexity

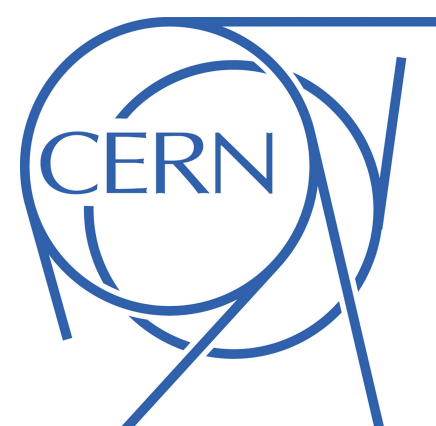
Potential machine learning solutions

SONIC paper

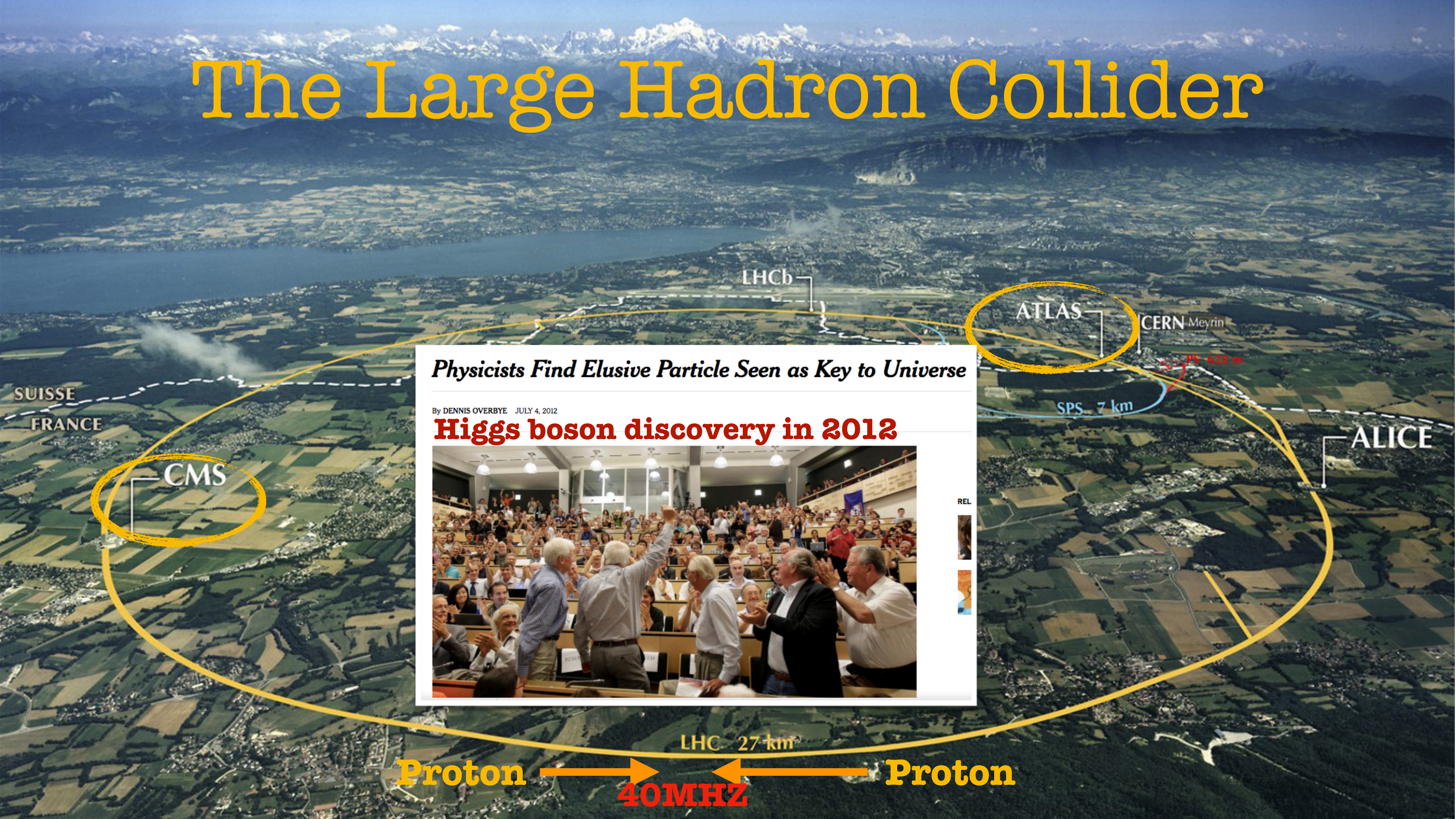
Proof of concept study with Brainwave:
Heterogenous computing for particle physics



Outlook& takeaways



The Large Hadron Collider



Physicists Find Elusive Particle Seen as Key to Universe

By DENNIS OVERBYE JULY 4, 2012

Higgs boson discovery in 2012



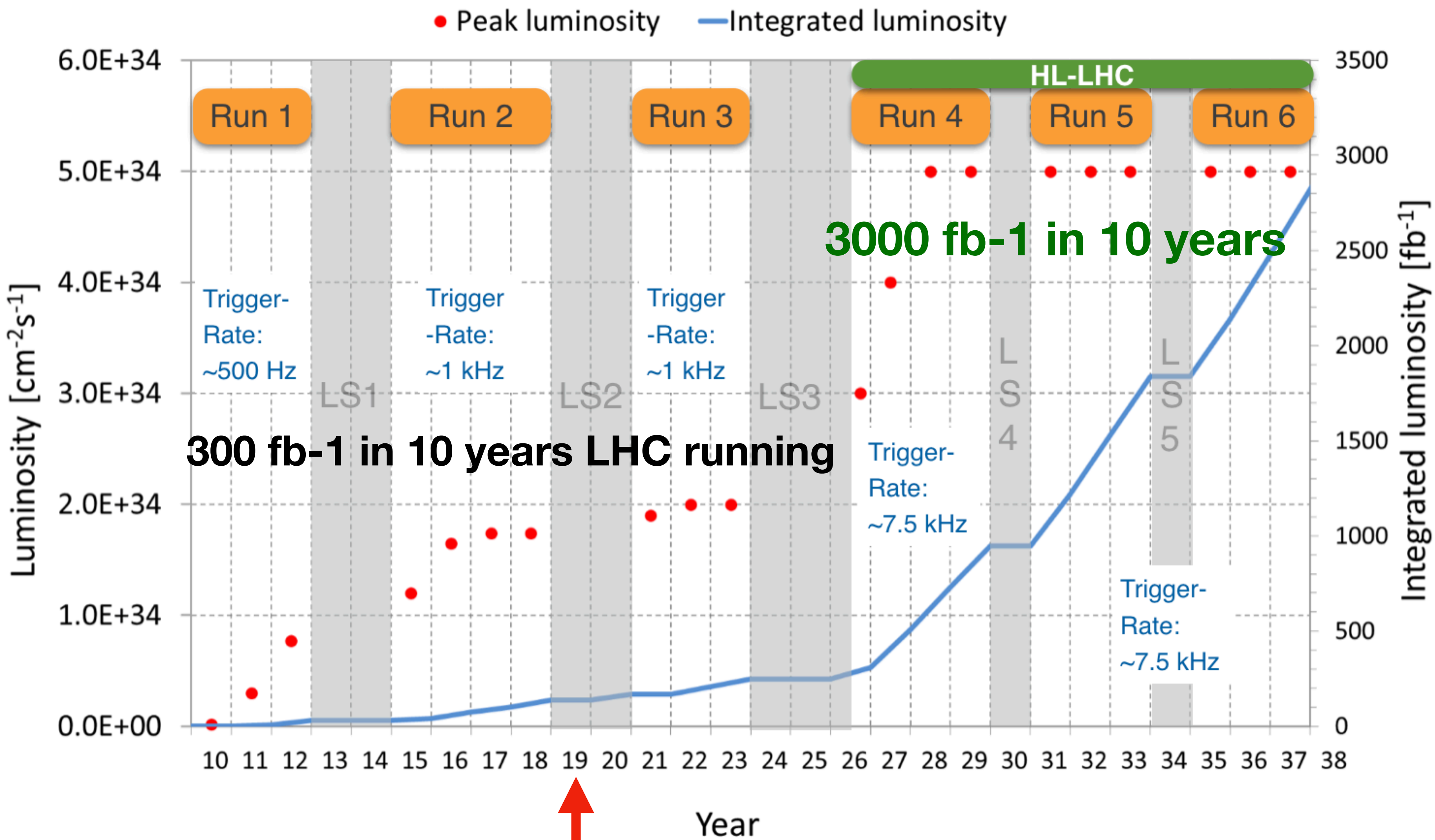
Proton

40MHZ

Proton

LHC 27 km

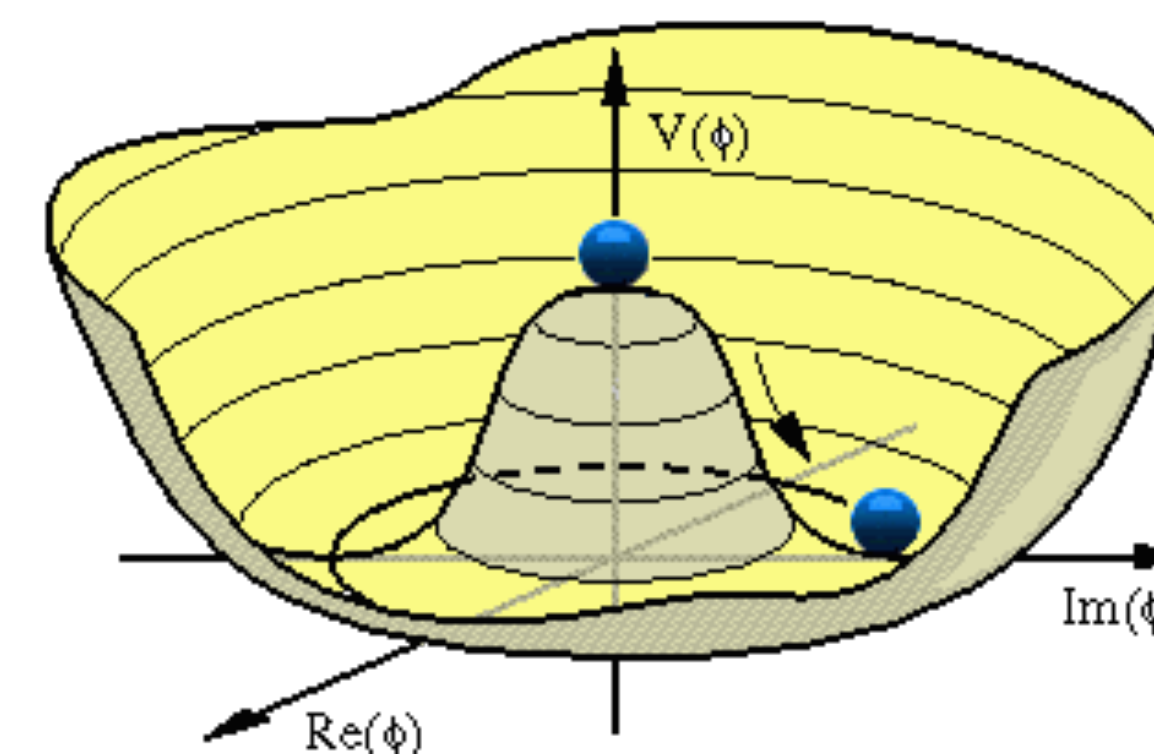
LHC running plan



we are here

We need the full dataset

example: Study the Higgs boson potential



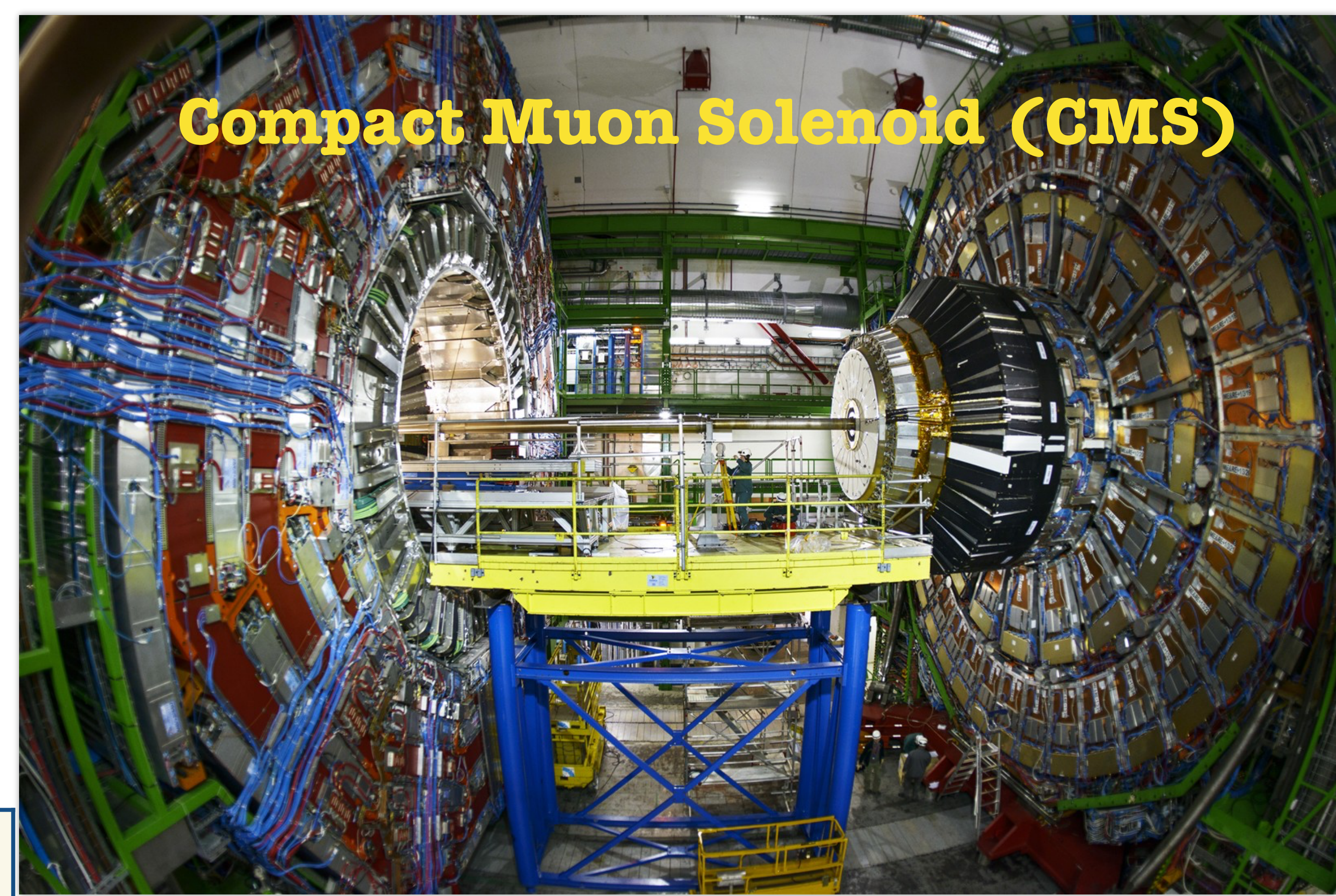
Multiple pp collisions in the same beam crossing
To increase data rate, squeeze beams as much as possible

Run 2: $\langle \text{PileUp} \rangle \sim 20\text{-}50$
Run 3: $\langle \text{PileUp} \rangle \sim 50\text{-}80$
HL-LHC: 140-200

FASTER AND MORE COMPLEX DETECTORS



Phase 1 forward pixel detector @ Fermilab



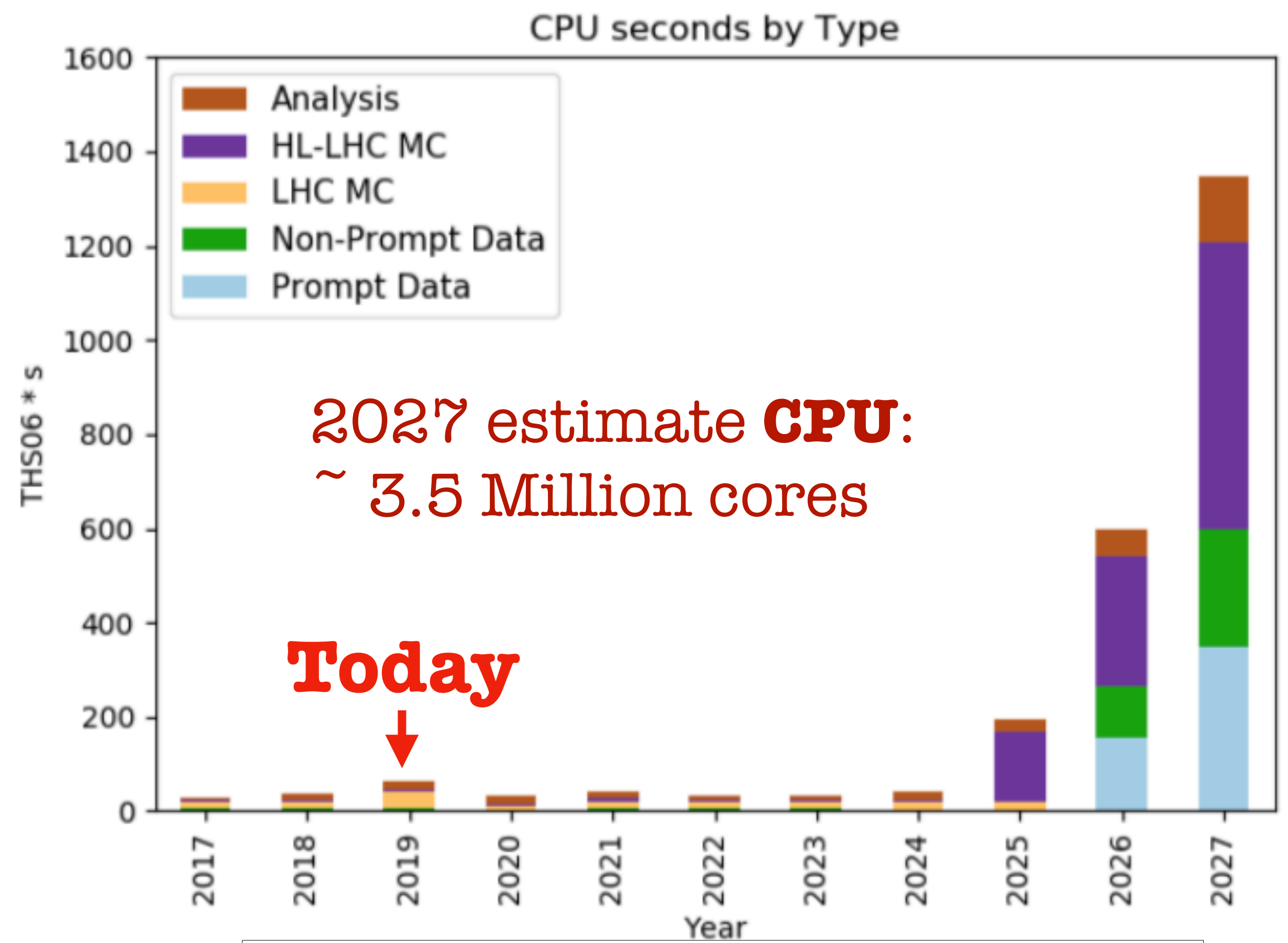
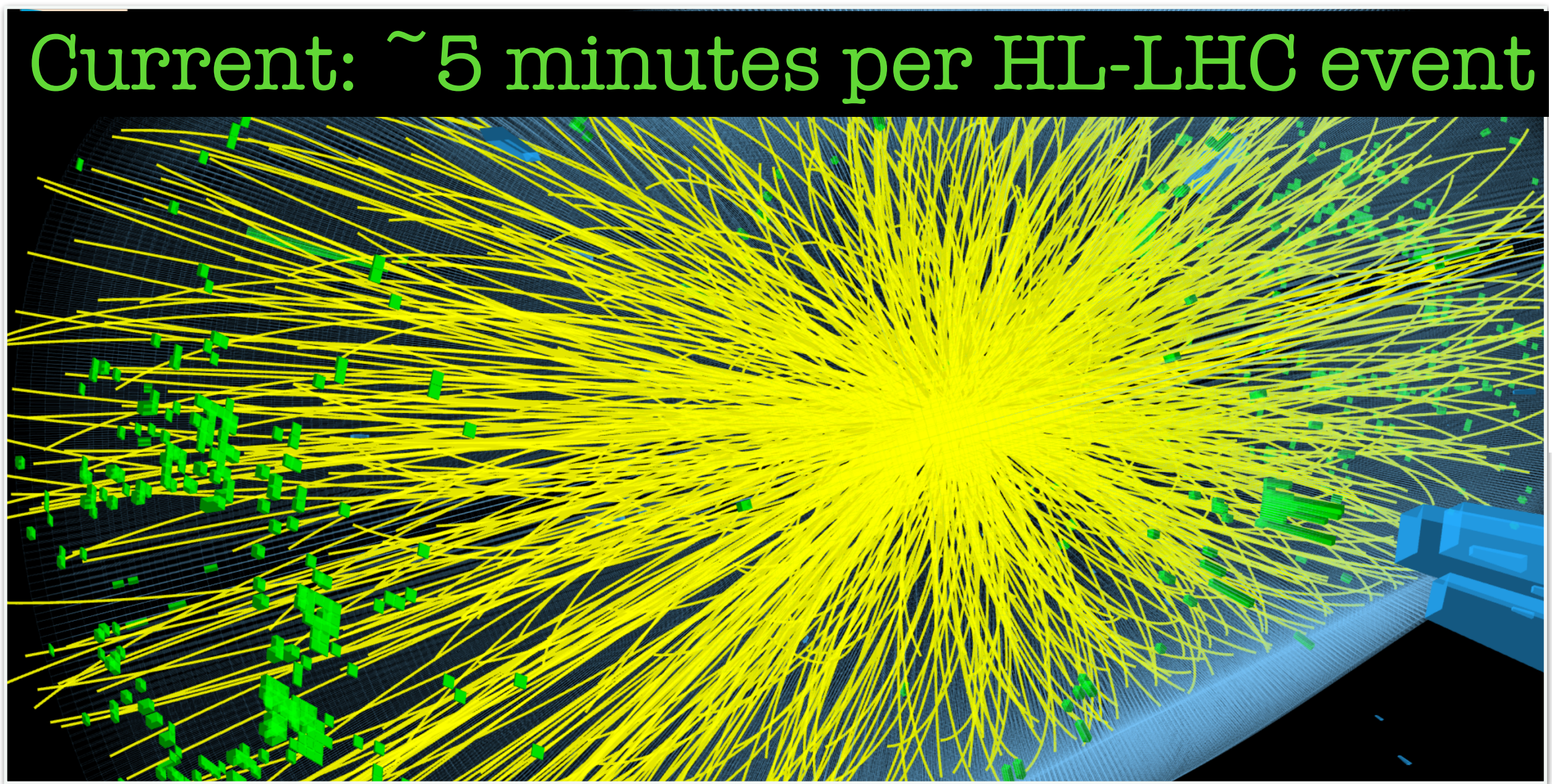
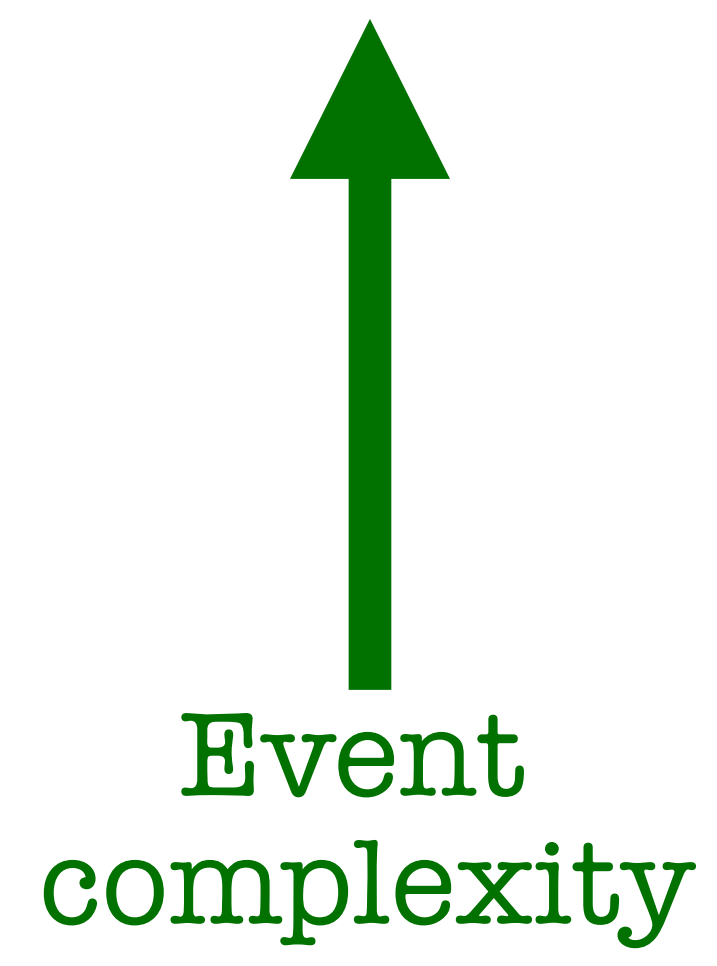
Compact Muon Solenoid (CMS)

CMS pixel ('ultra high speed/resolution camera')	Number of channels
'Phase-0'	66 M
'Phase-1'	123 M
'Phase-2'	2B

- Building faster detectors with better resolution
- Trigger & computing challenges @ HL-LHC: machine learning solutions?

INCREASED DATA VOLUME AND COMPLEXITY

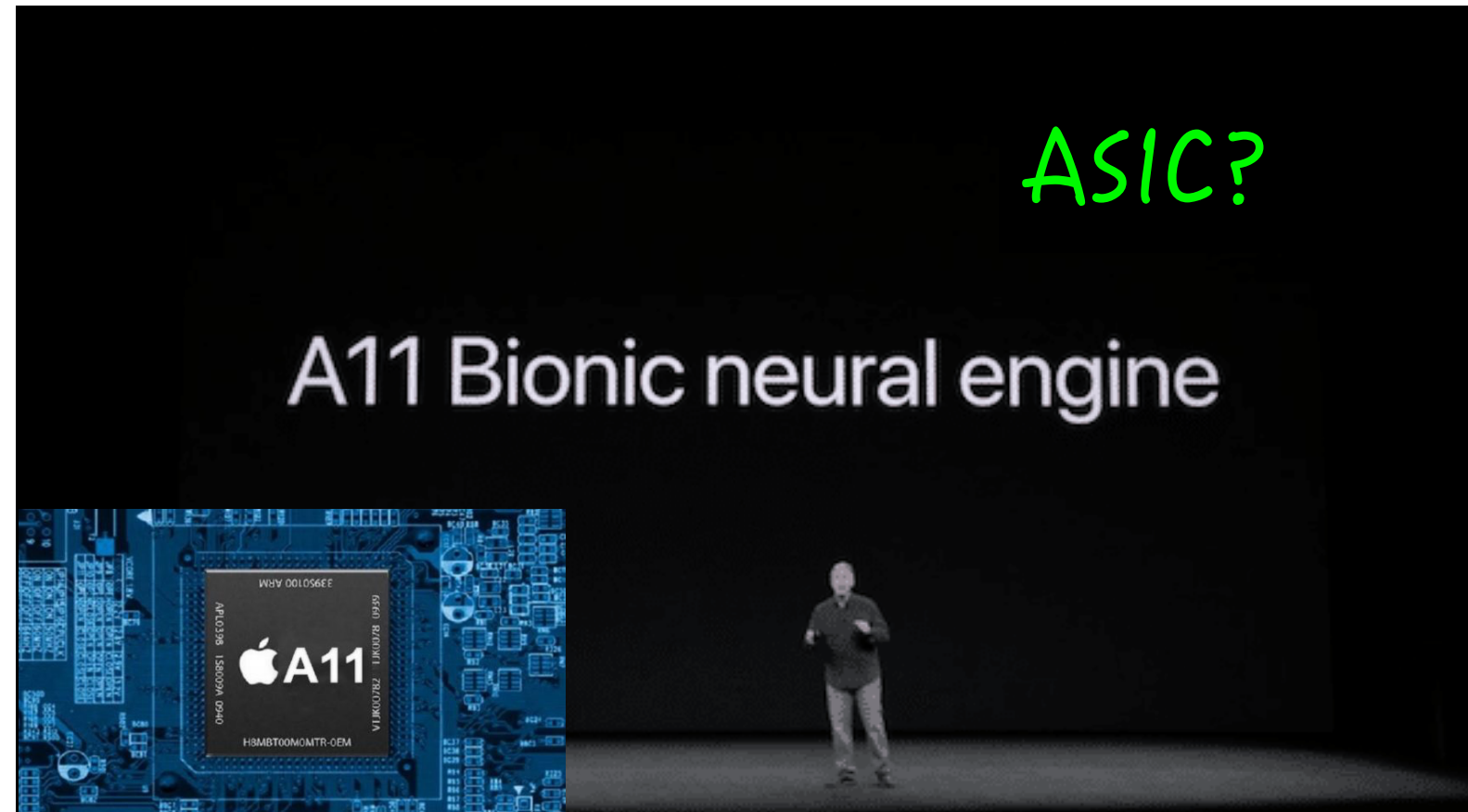
10X more data
at the High-
luminosity LHC



Problem won't get solved by
itself: **Moore's Law**
continues
...but **Dennard Scaling** fails

ASIC?

A11 Bionic neural engine



Specialized co-processor hardware for machine learning:
flexibility vs speed (efficiency)

Catapult/Brainwave
Microsoft



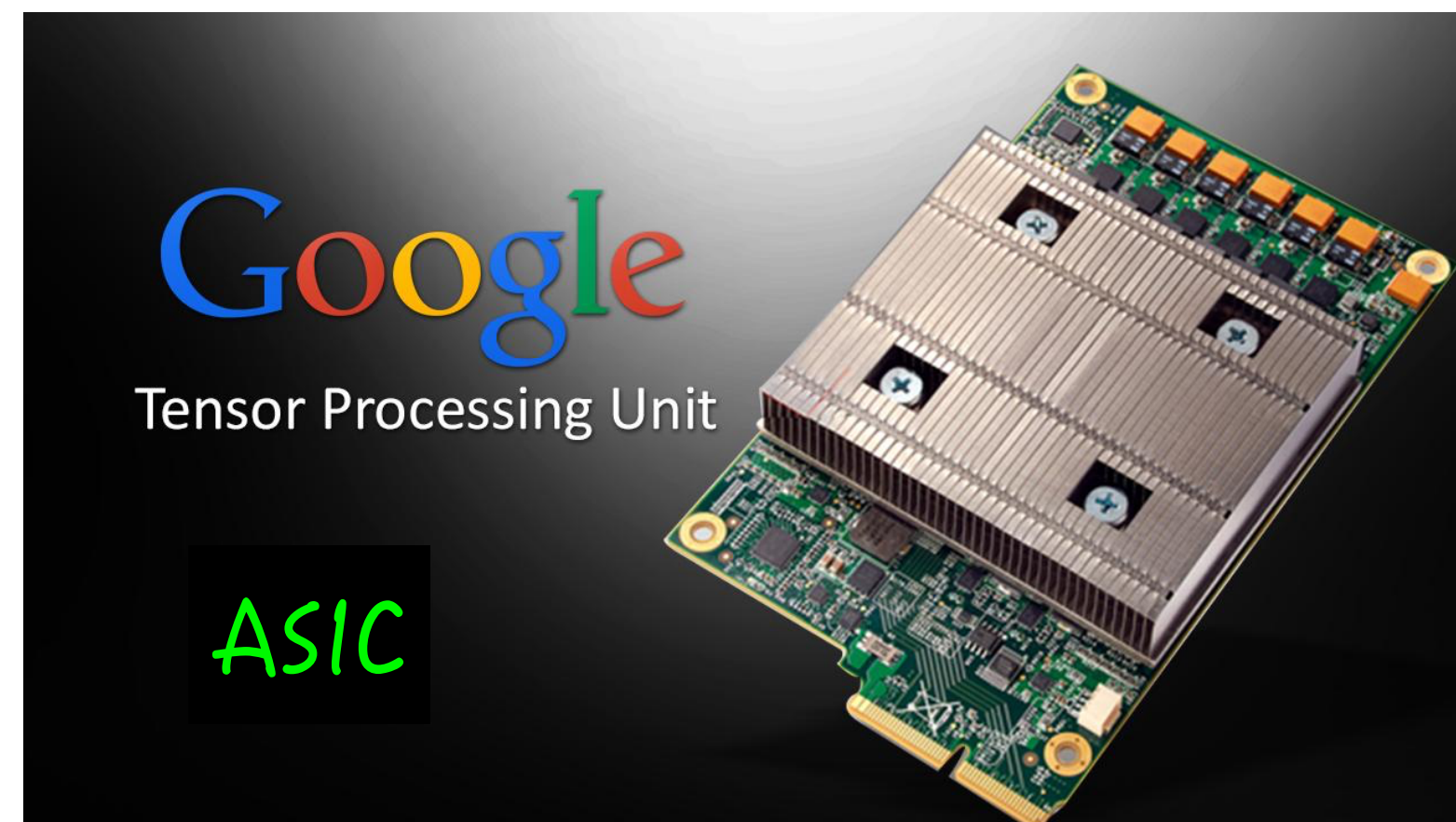
FPGA

FPGA



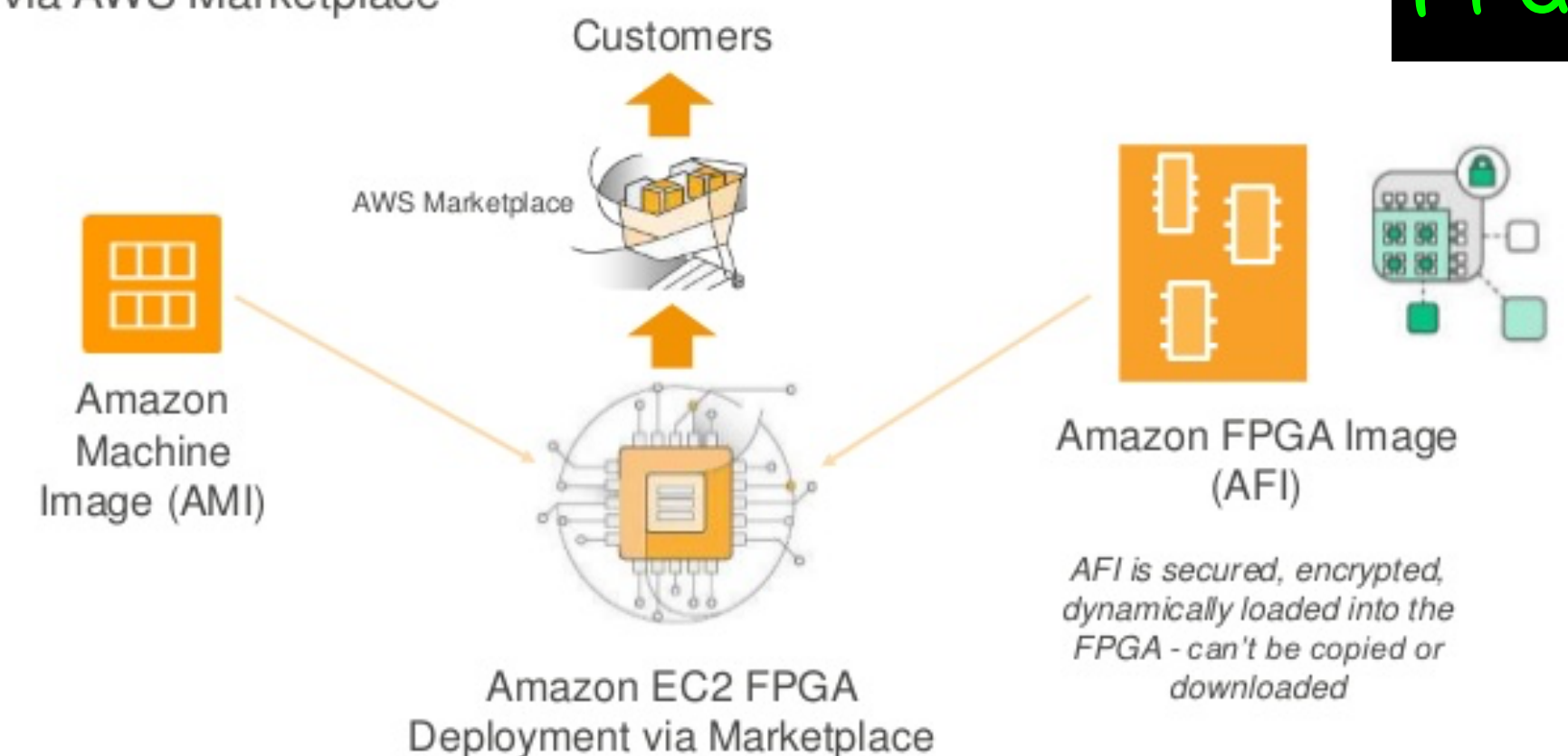
Google
Tensor Processing Unit

ASIC



Delivering FPGA Partner Solutions on AWS
via AWS Marketplace

FPGA



INTEL® FPGA ACCELERATION HUB

The Intel® Xeon® Acceleration Stack for FPGAs is a robust framework enabling data center applications to leverage an FPGA's potential to increase

Option 1

re-write physics algorithms for new hardware

Language: OpenCL, OpenMP, TBB, HLS, ...?

Hardware: FPGA, GPU

Option 2

re-cast physics problem as a machine learning problem

Language: C++, Python (TensorFlow, PyTorch,...)

Hardware: FPGA, GPU, ASIC

e.g. track reconstruction

Option 1: Parallelized and Vectorized Tracking Using Kalman Filters

Option 2: Tracking using CNN and Graph Networks: e.g. HEP.TrkX project.

Advantages:

- Algorithms expressed as matrix multiplications: intrinsically parallelizable
- Take advantage of co-processors optimized for ML fast inference

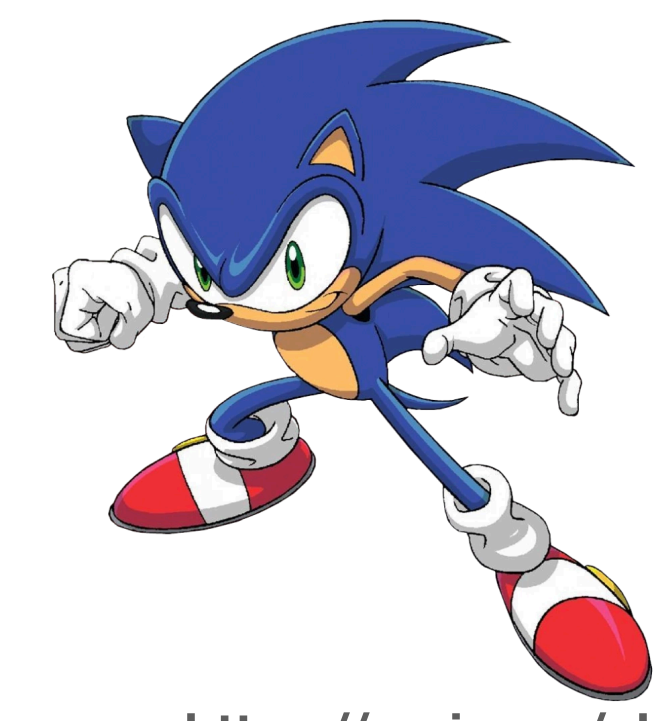
Caveat: challenges of solving our reconstruction problems with NNs.

PROOF OF CONCEPT: SONIC

Services for **O**ptimized **N**etwork **I**nference on **C**o-processors

FPGA-accelerated machine learning inference as a service for particle physics computing

Javier Duarte · Philip Harris · Scott Hauck · Burt Holzman · Shih-Chieh Hsu · Sergo Jindariani · Suffian Khan · Benjamin Kreis · Brian Lee · Mia Liu · Vladimir Lončar · Jennifer Ngadiuba · Kevin Pedro · Brandon Perez · Maurizio Pierini · Dylan Rankin · Nhan Tran · Matthew Trahms · Aristeidis Tsaris · Colin Versteeg · Ted W. Way · Dustin Werran · Zhenbin Wu



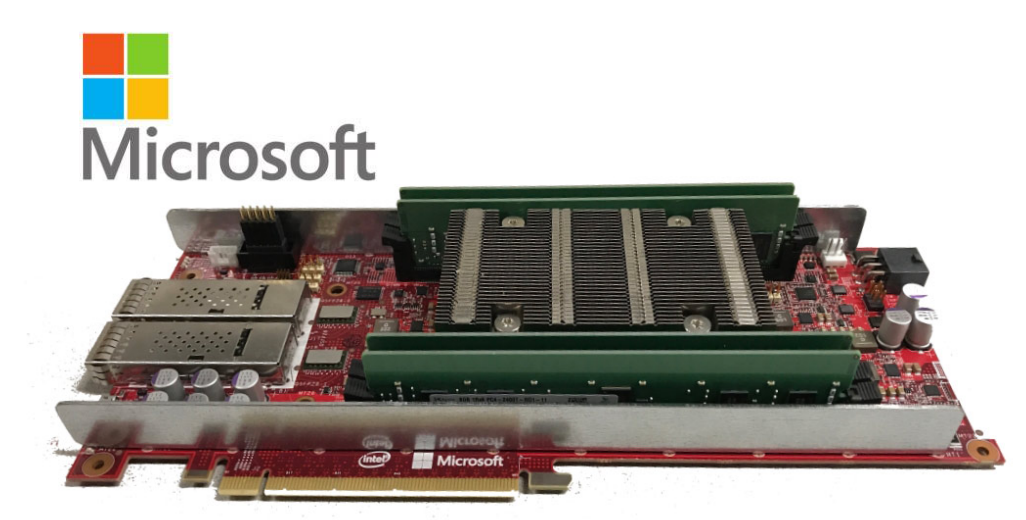
<https://arxiv.org/abs/1904.08986>

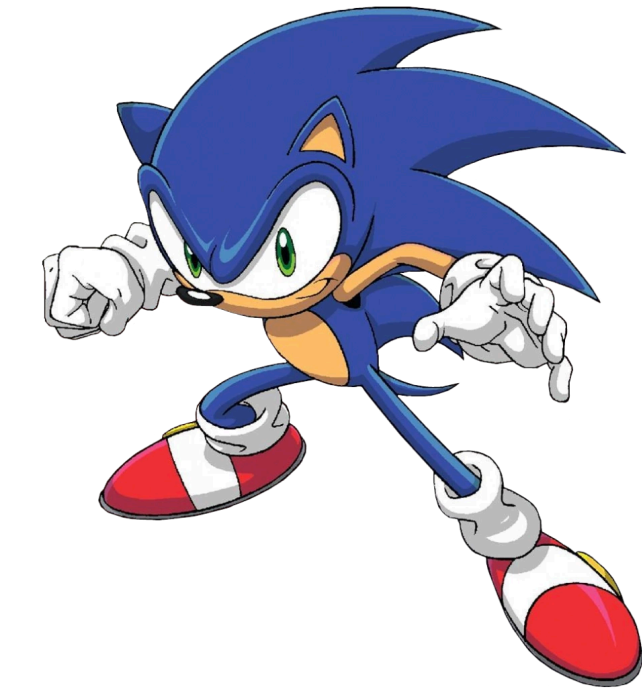
Question:

How do we help with physics event data processing model with industry developments in co-processors?

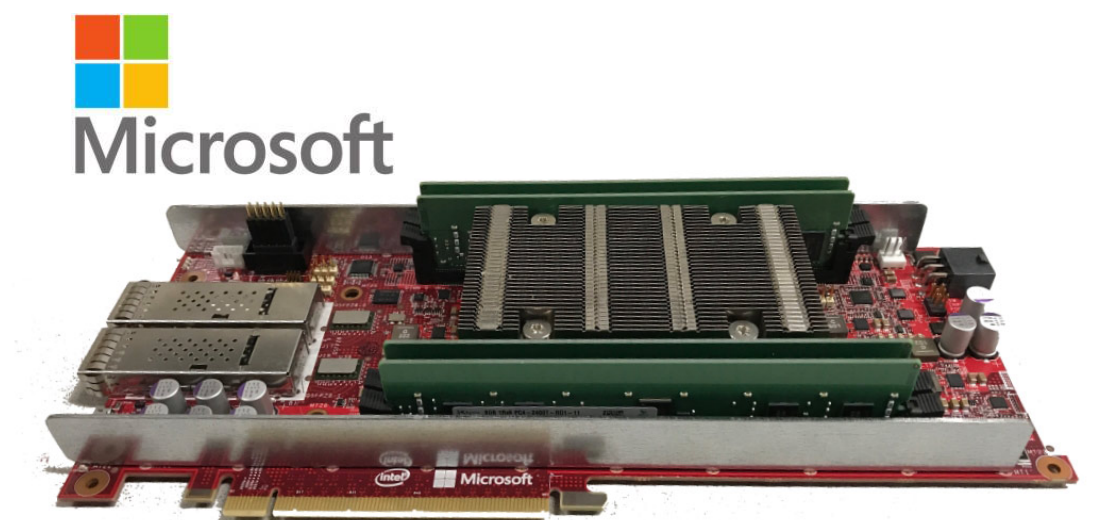
Focus on speeding up the **inference**.

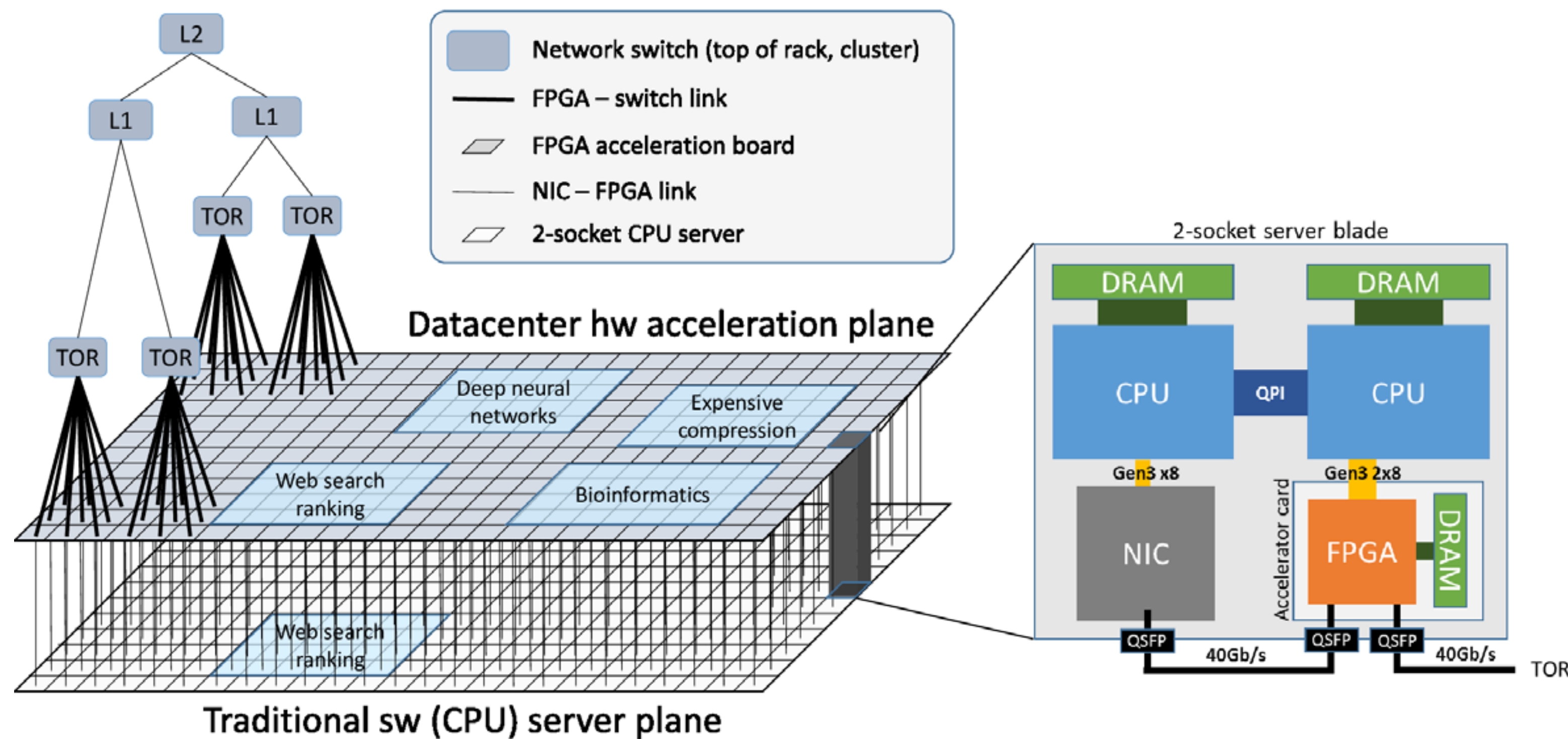
Catapult/Brainwave





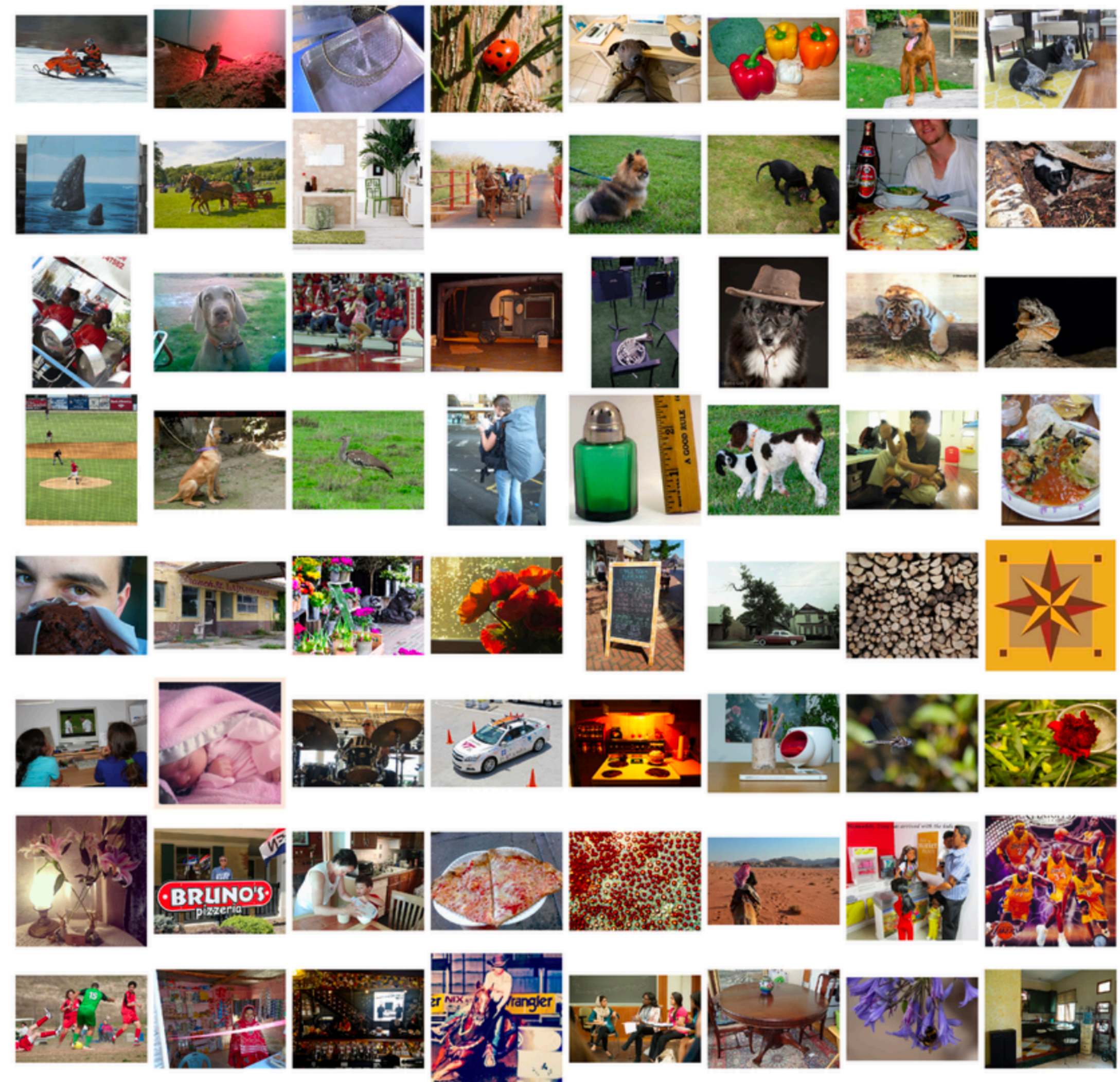
- Proof of concept: Top tagging/neutrino image classification on Brainwave
- Implementation as service in CMSSW in non-disruptive way.
- Speed and data throughput performance



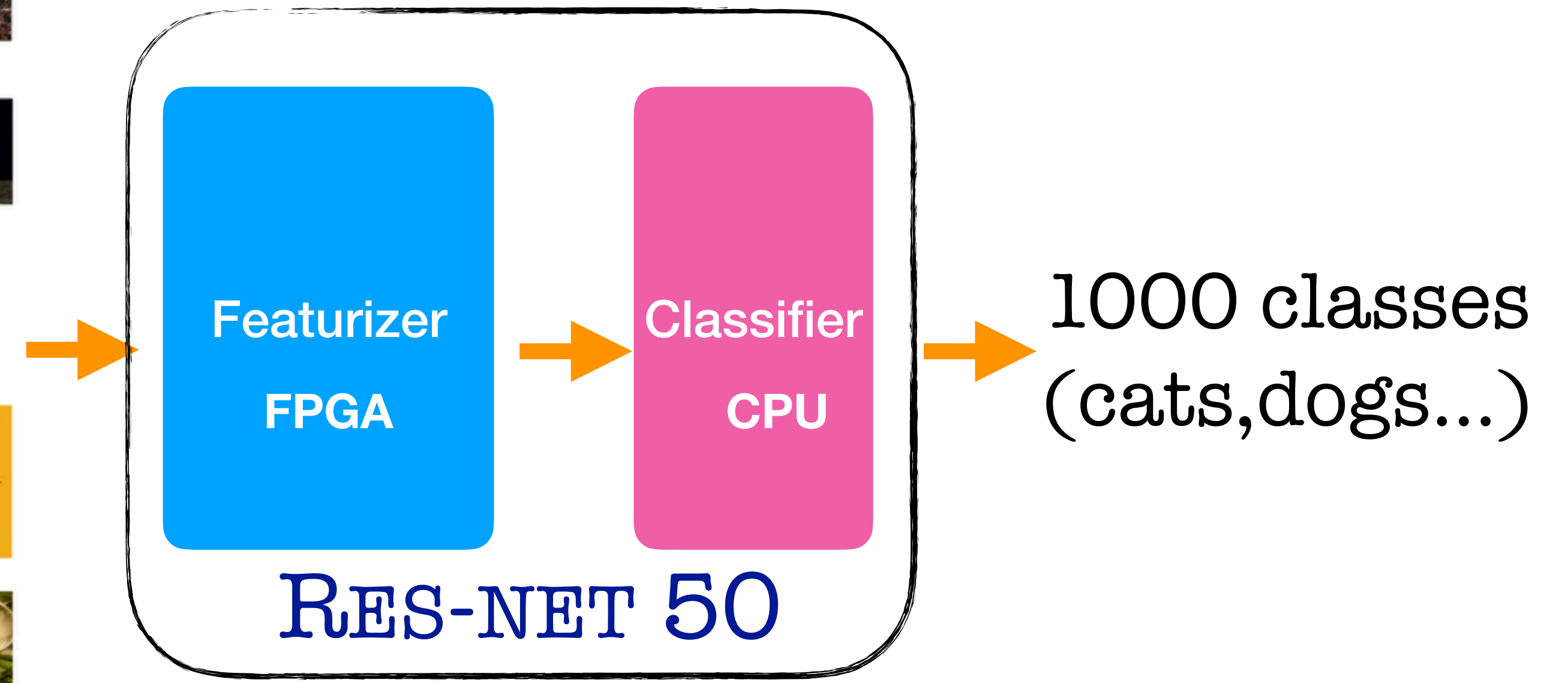


- Provides a full service at scale
(more than just a single co-processor)
- Multi-FPGA/CPU fabric accelerates both **computing** and **network**
- Models supported:
 - ResNet50, ResNet152, DenseNet121, VGGNet 16
 - Weight retuning available

A PHYSICS CASE: JET TAGGING WITH RES-NET 50¹³



ImageNet



RES-NET 50

25M parameters

4 G-ops/inference

PASSING JET IMAGES TO RE-TRAIN RES-NET 50

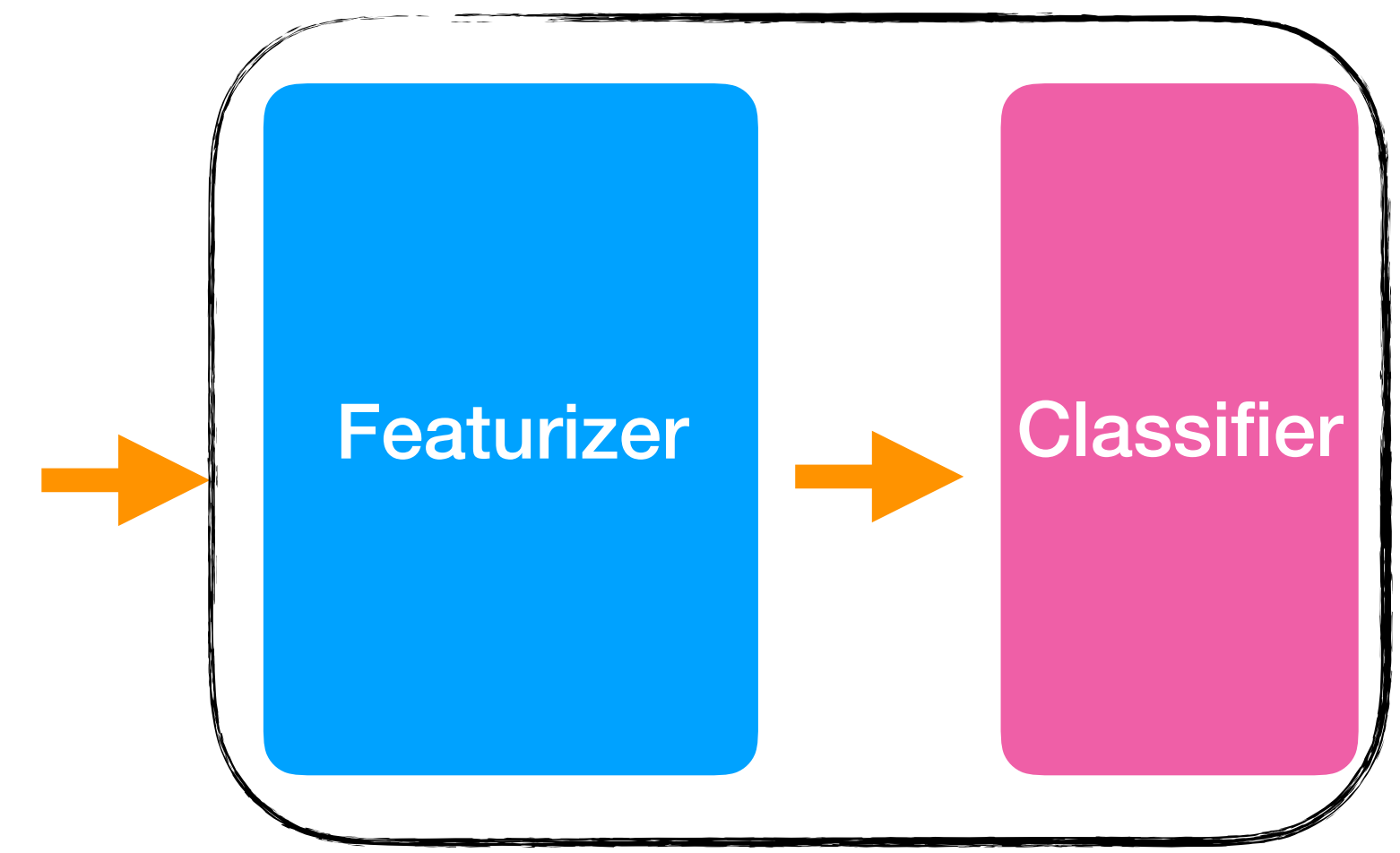
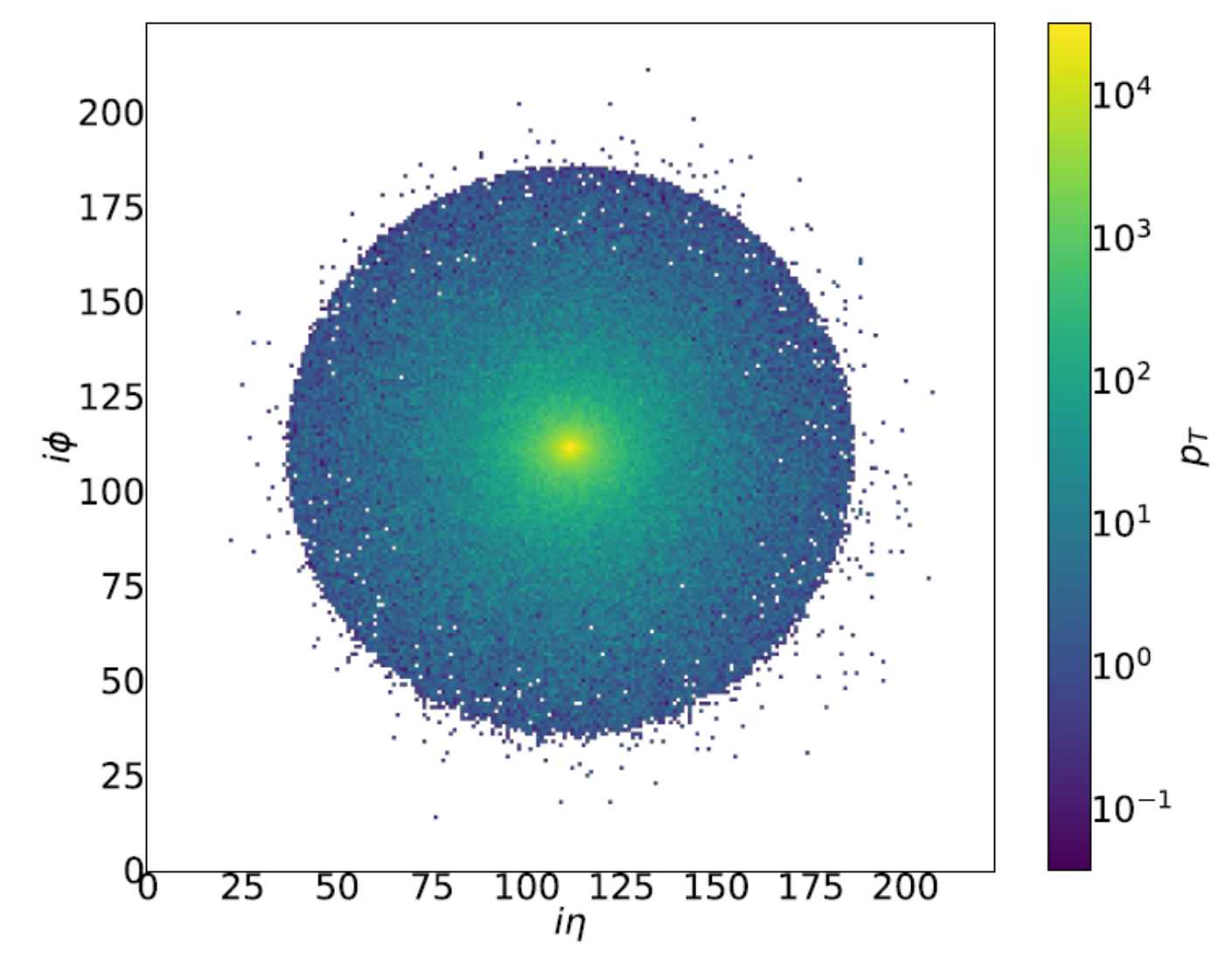
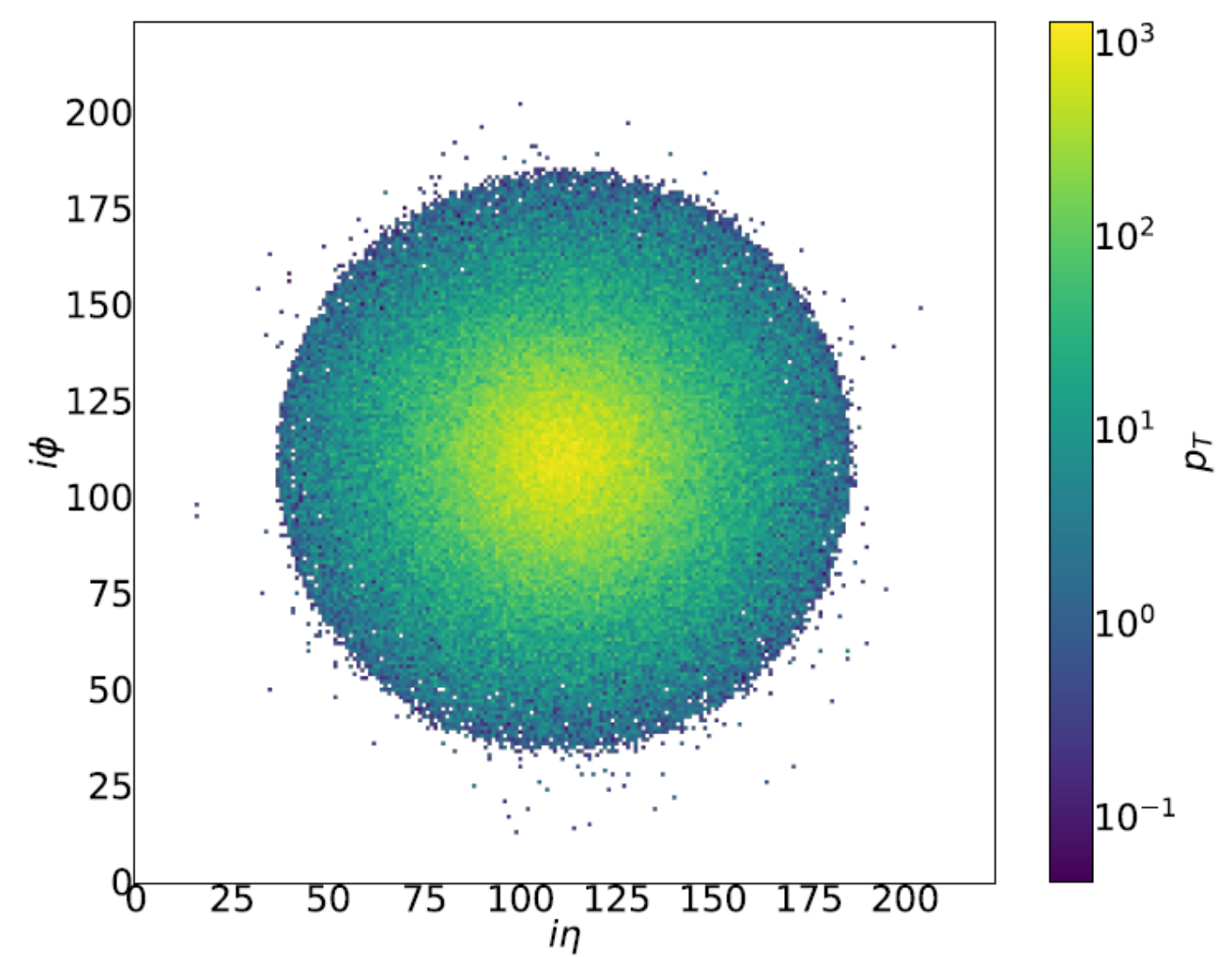


Top tagging benchmark dataset

- Images made from density map of the pt of jet constituents in $\eta * \phi$ space.
- Grey image, duplicated to RGB.

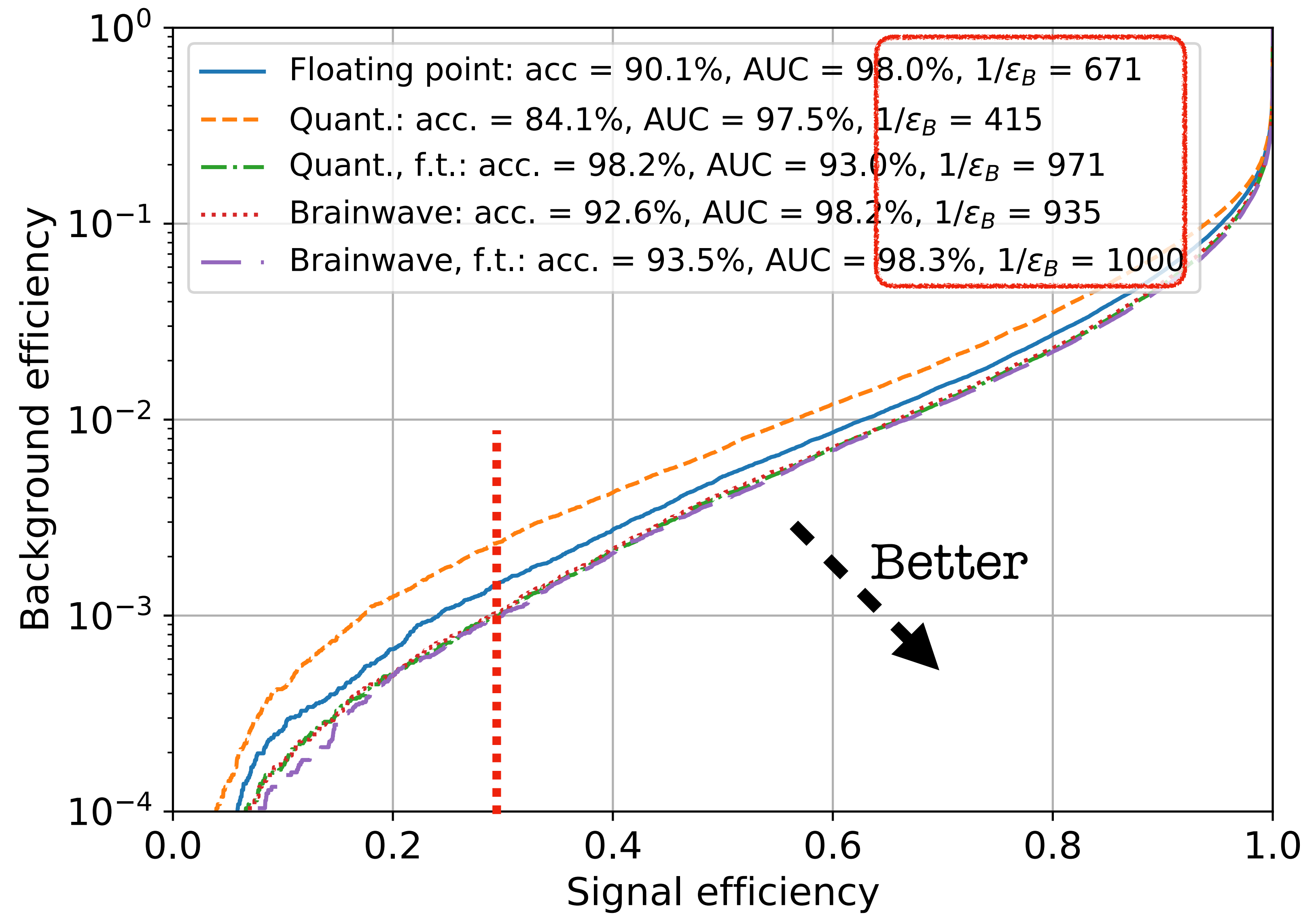
top, averaged over 5k jets

QCD, averaged over 5k jets



RETRAIN RES-NET 50

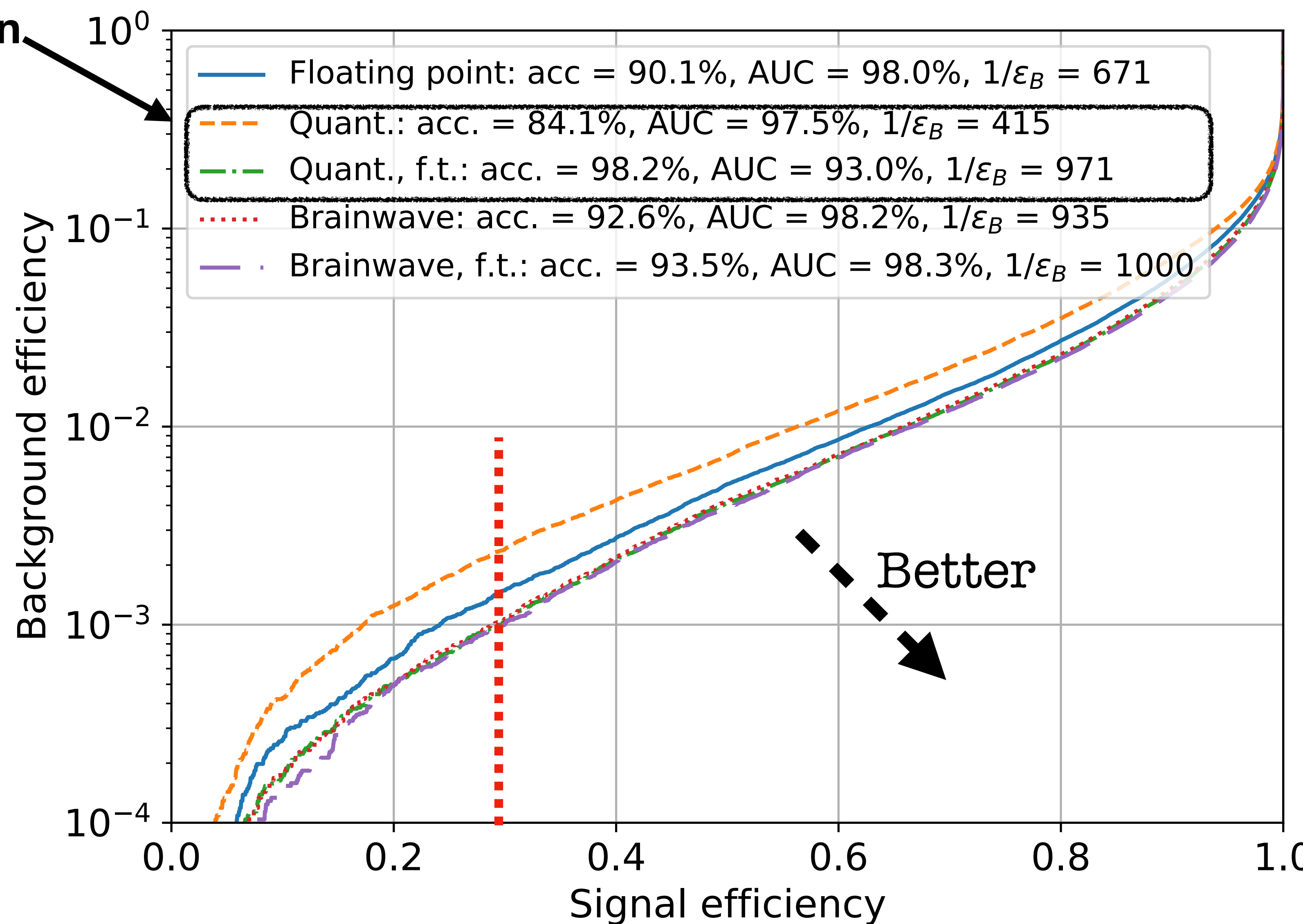
RE-TRAIN RESNET-50 FOR TOP TAGGING

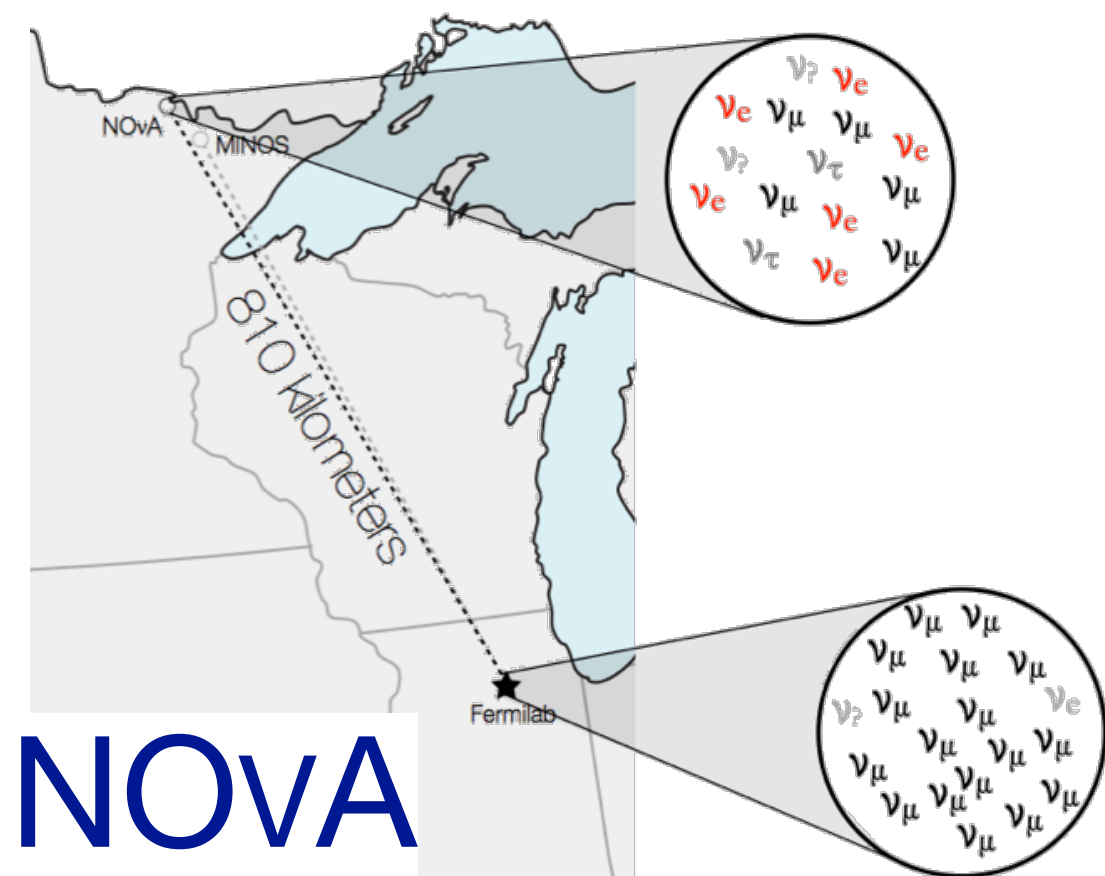


Quantized model:

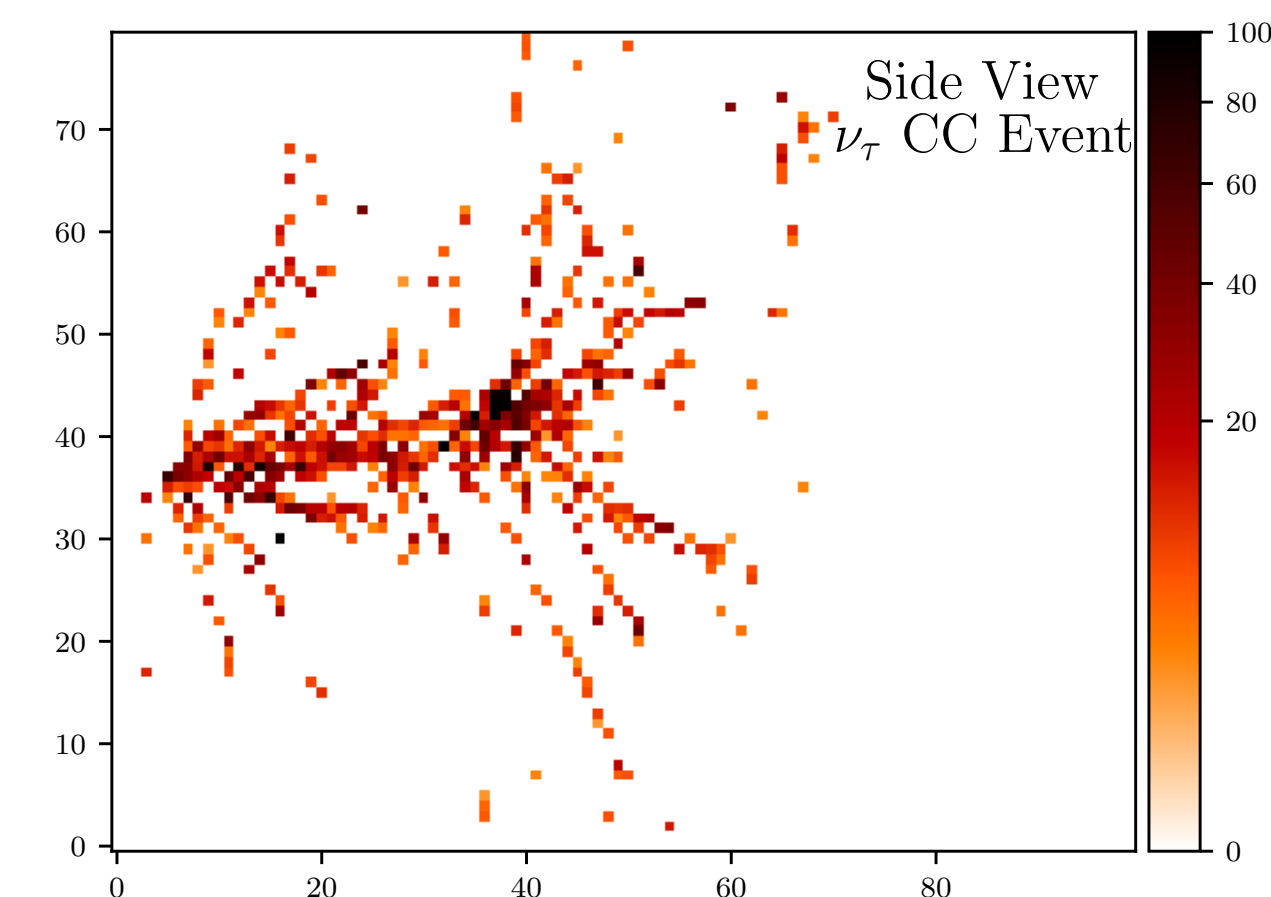
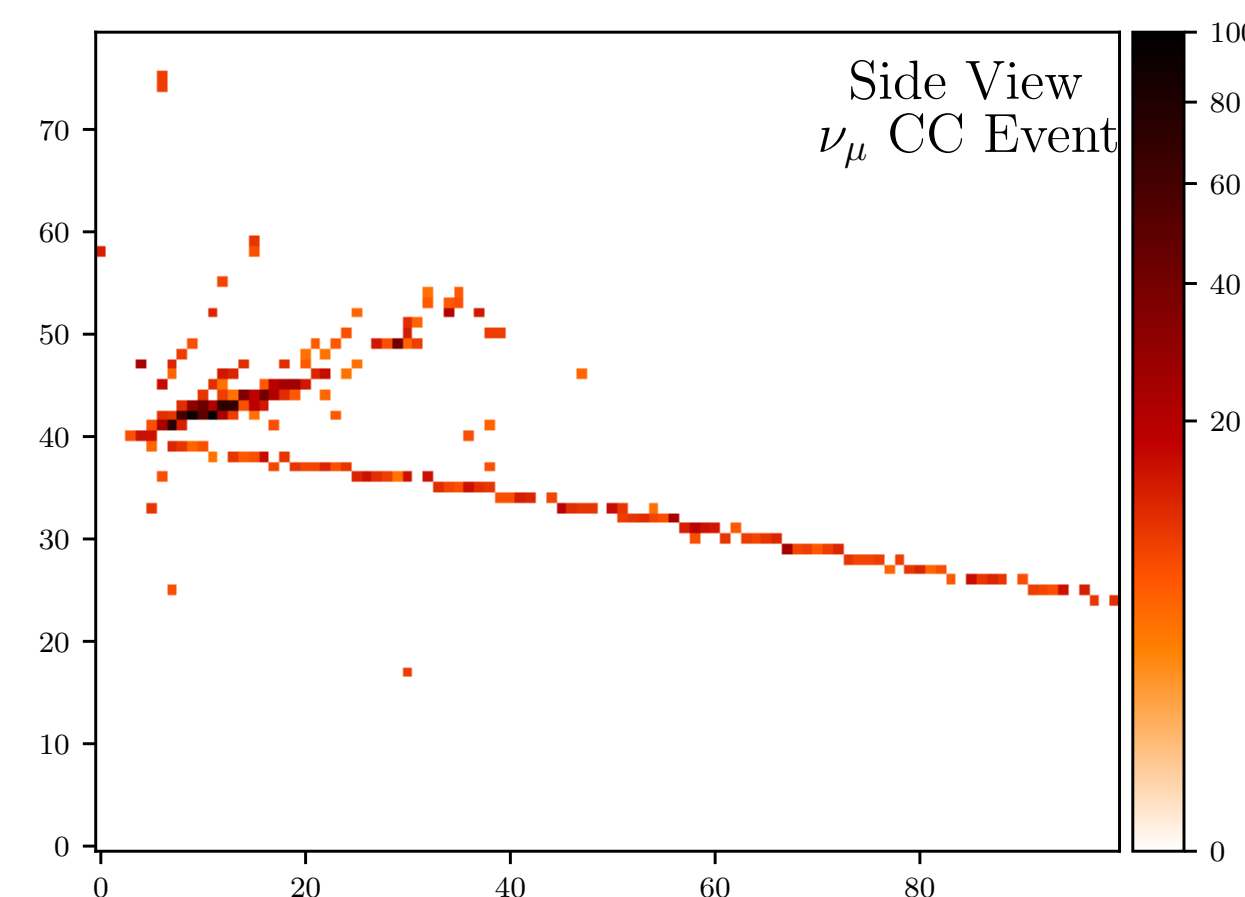
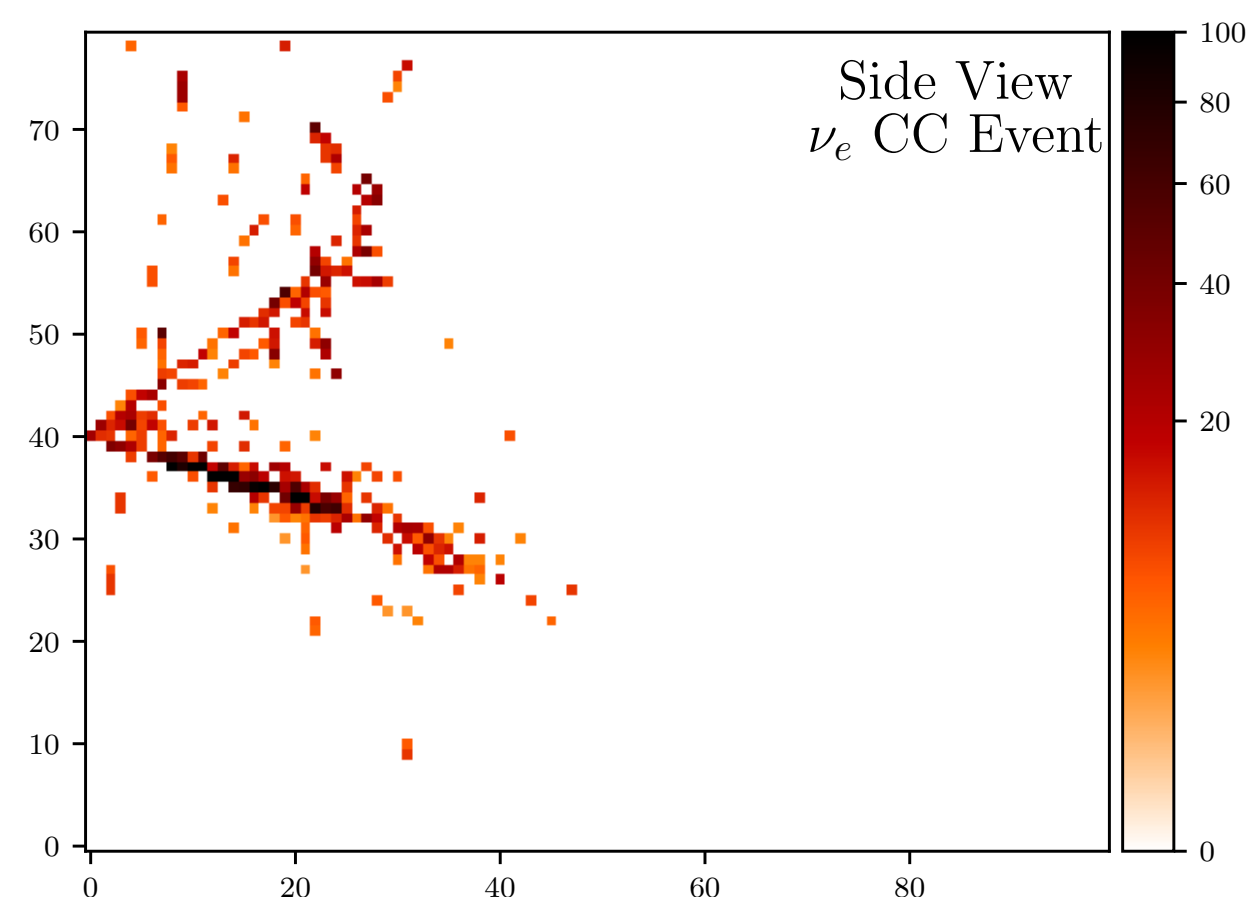
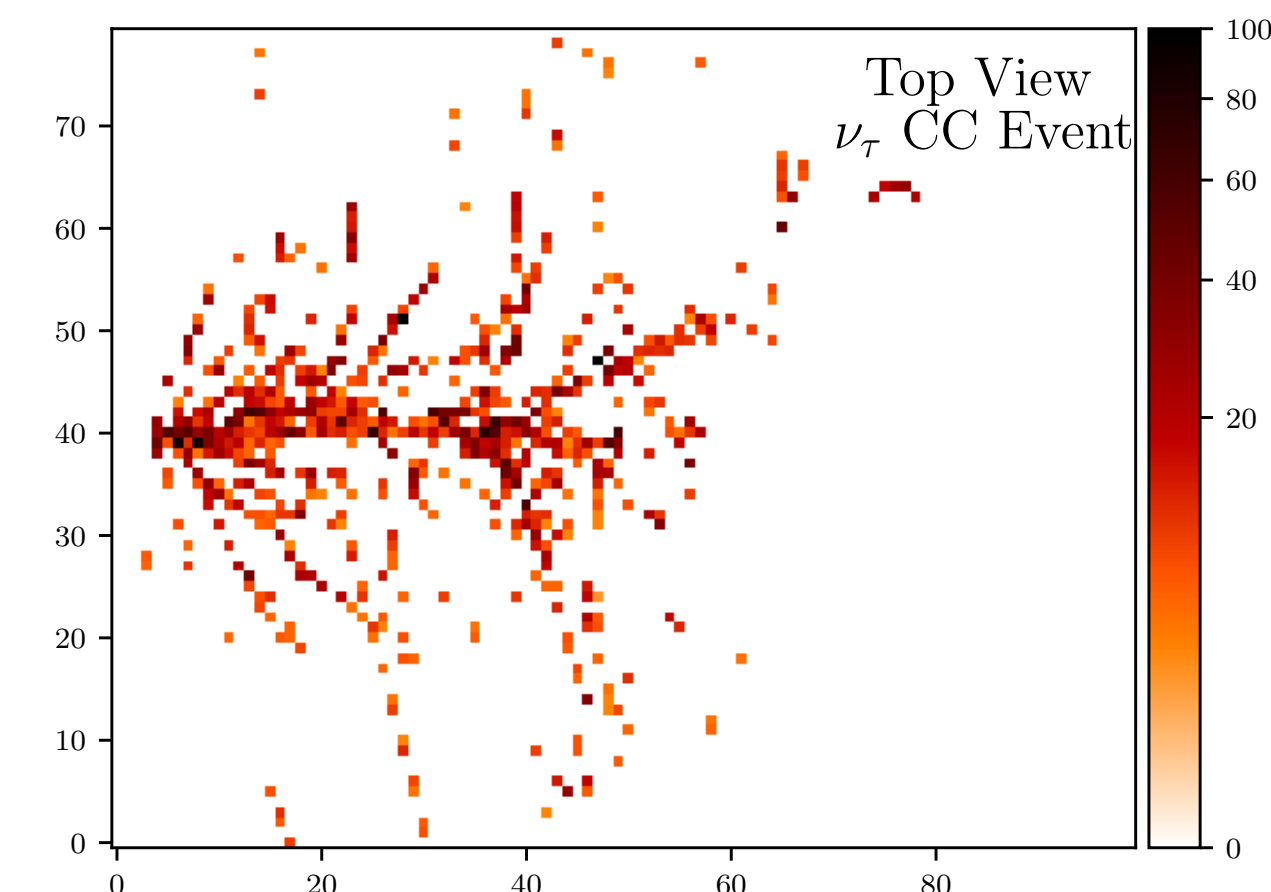
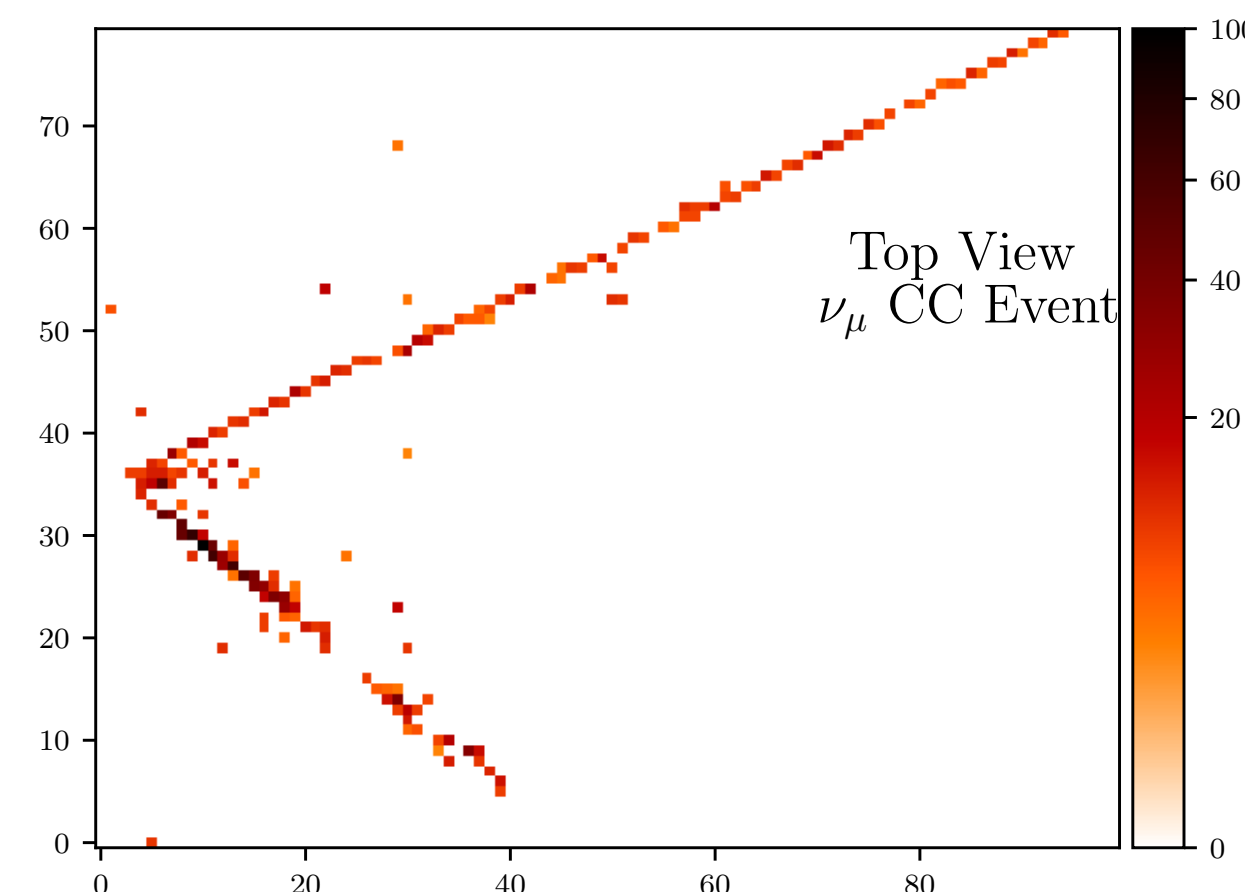
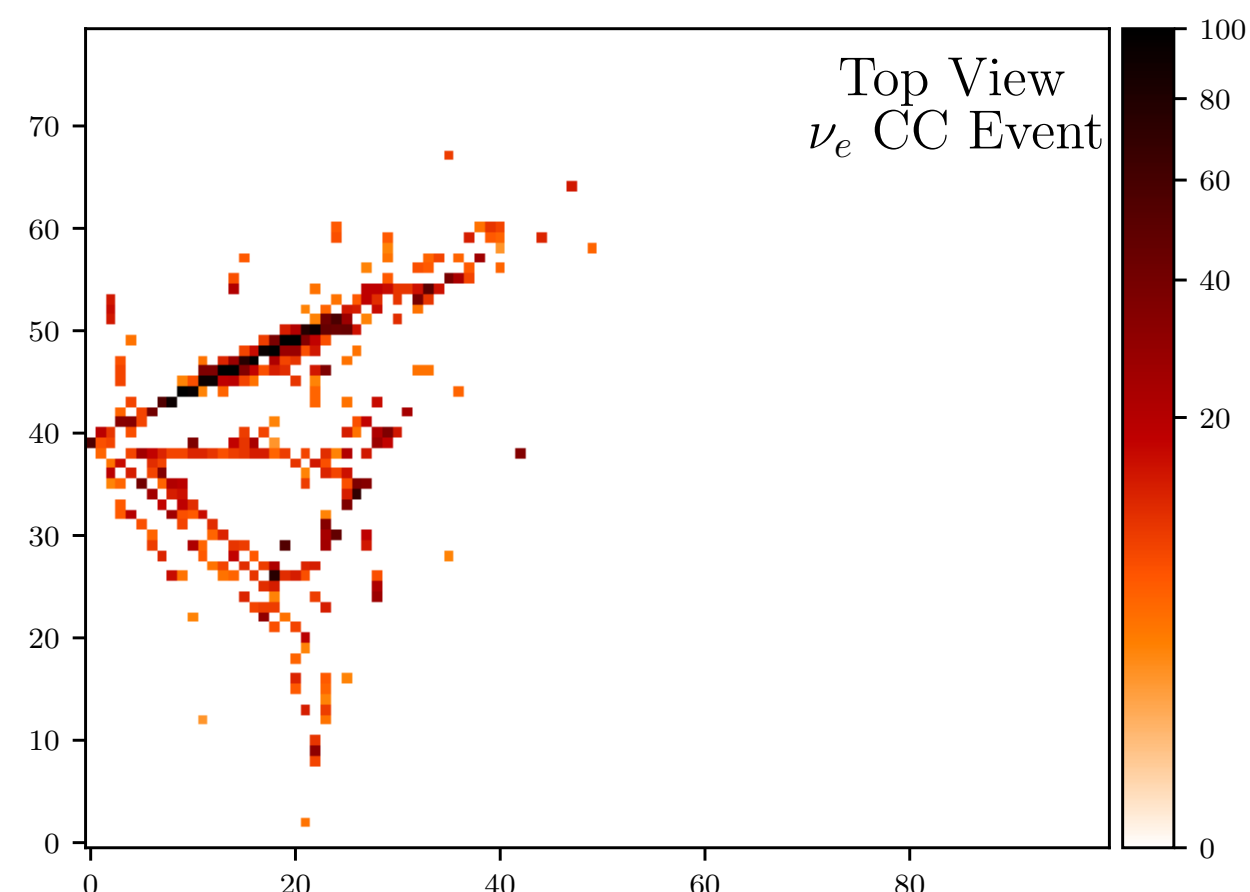
- Brainwave's implementation of ResNet50 on FPGA
- Can tune weights
- State of art performance achieved with quantized ResNet 50 on BrainWave service

Emulation





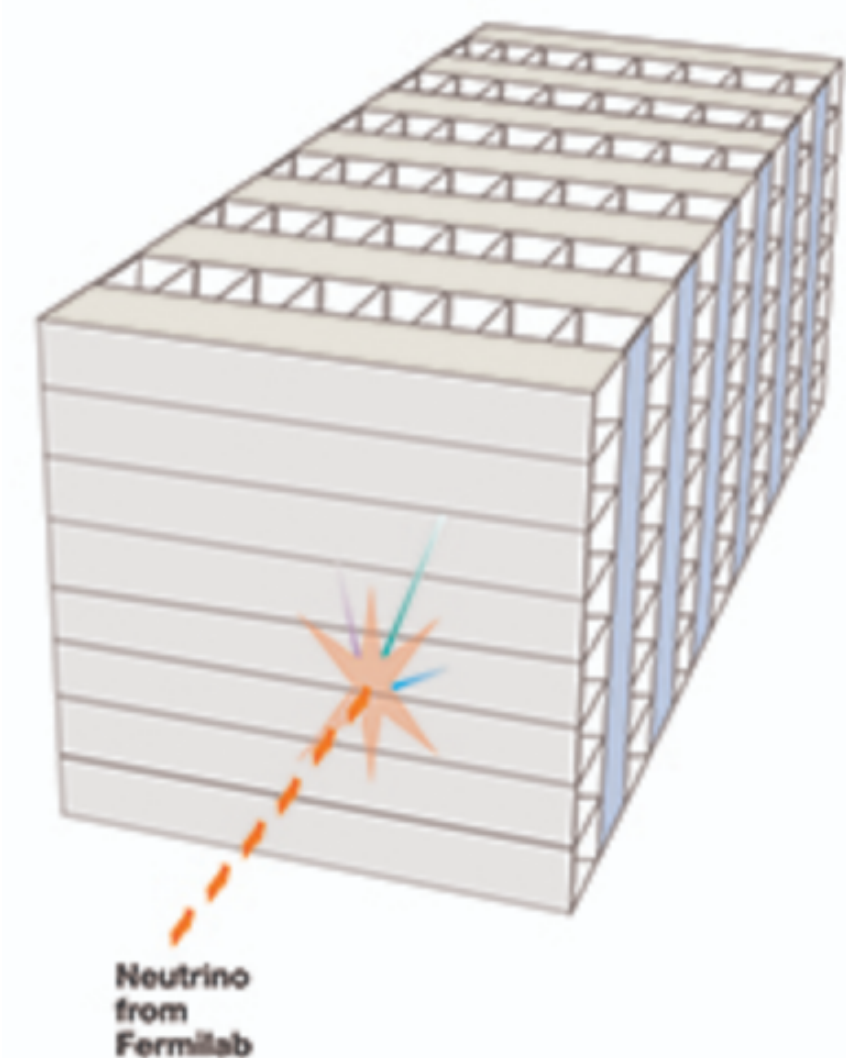
NOvA



Electron neutrino

Muon neutrino

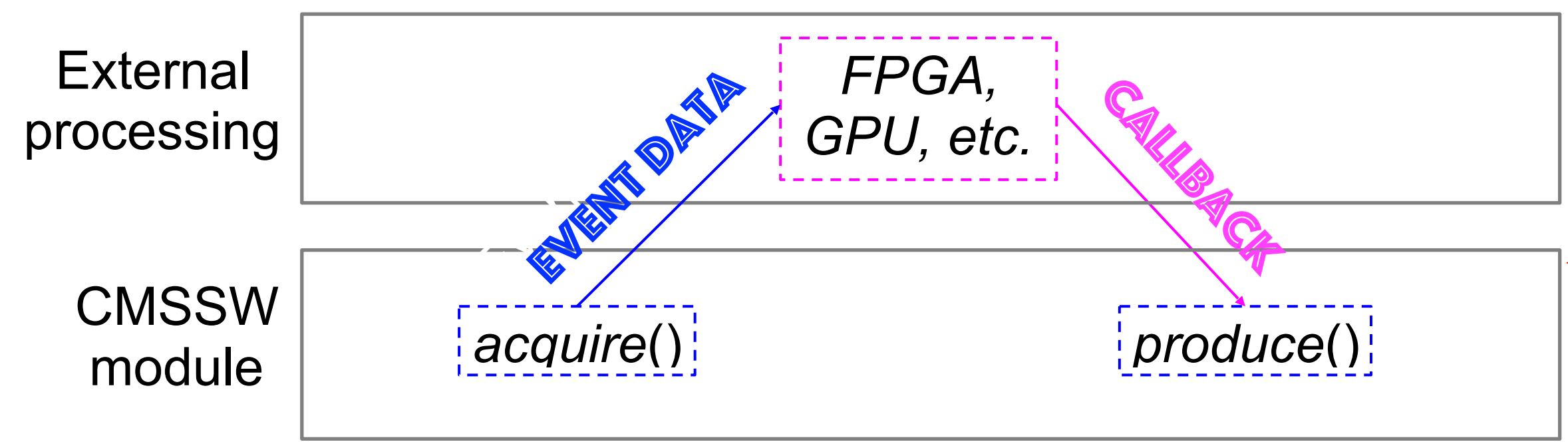
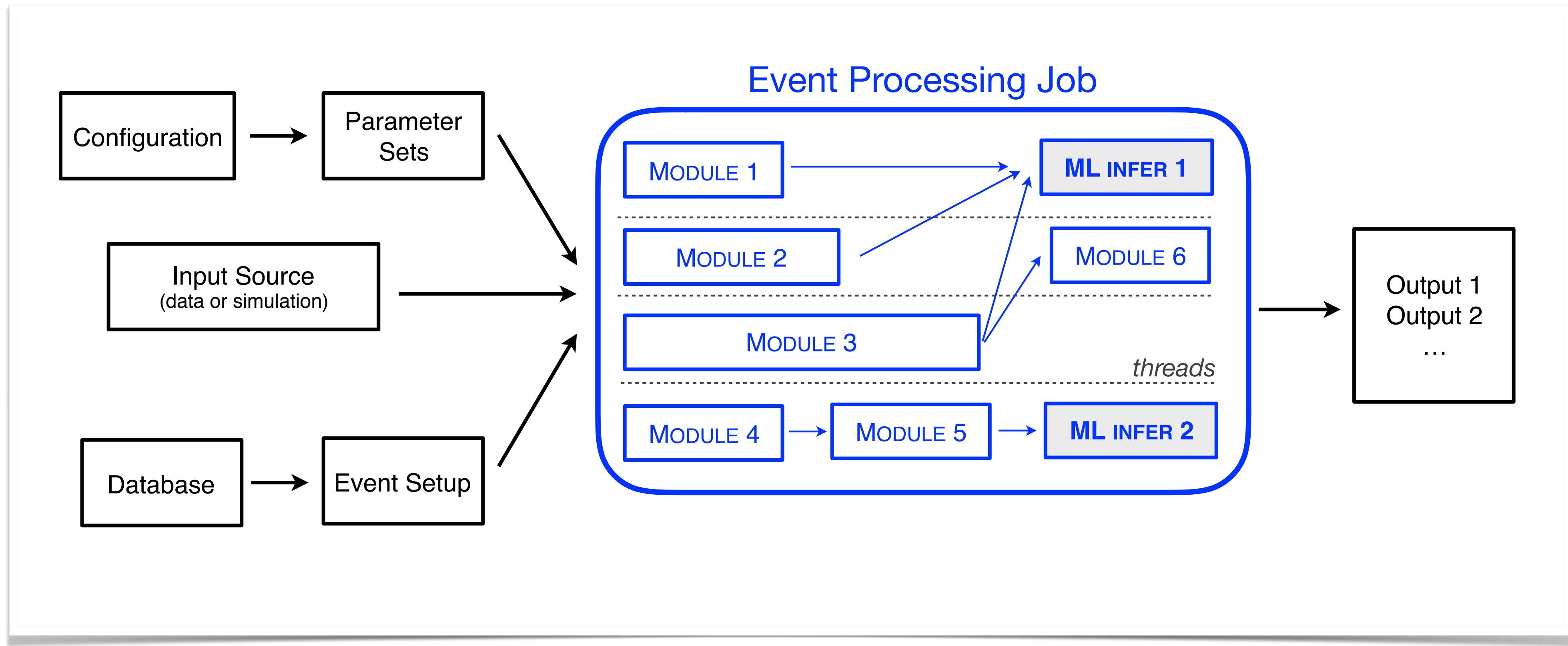
Tau neutrino



- **Primary goal of NOvA:** measurement of neutrino oscillations via $\nu_\mu \rightarrow \nu_e$: Classifying neutrinos with ResNet50 (transfer learning).

- Can be used in NOvA event processing as of today: see Thomas's lightning round talk.

CMS COMPUTING MODEL

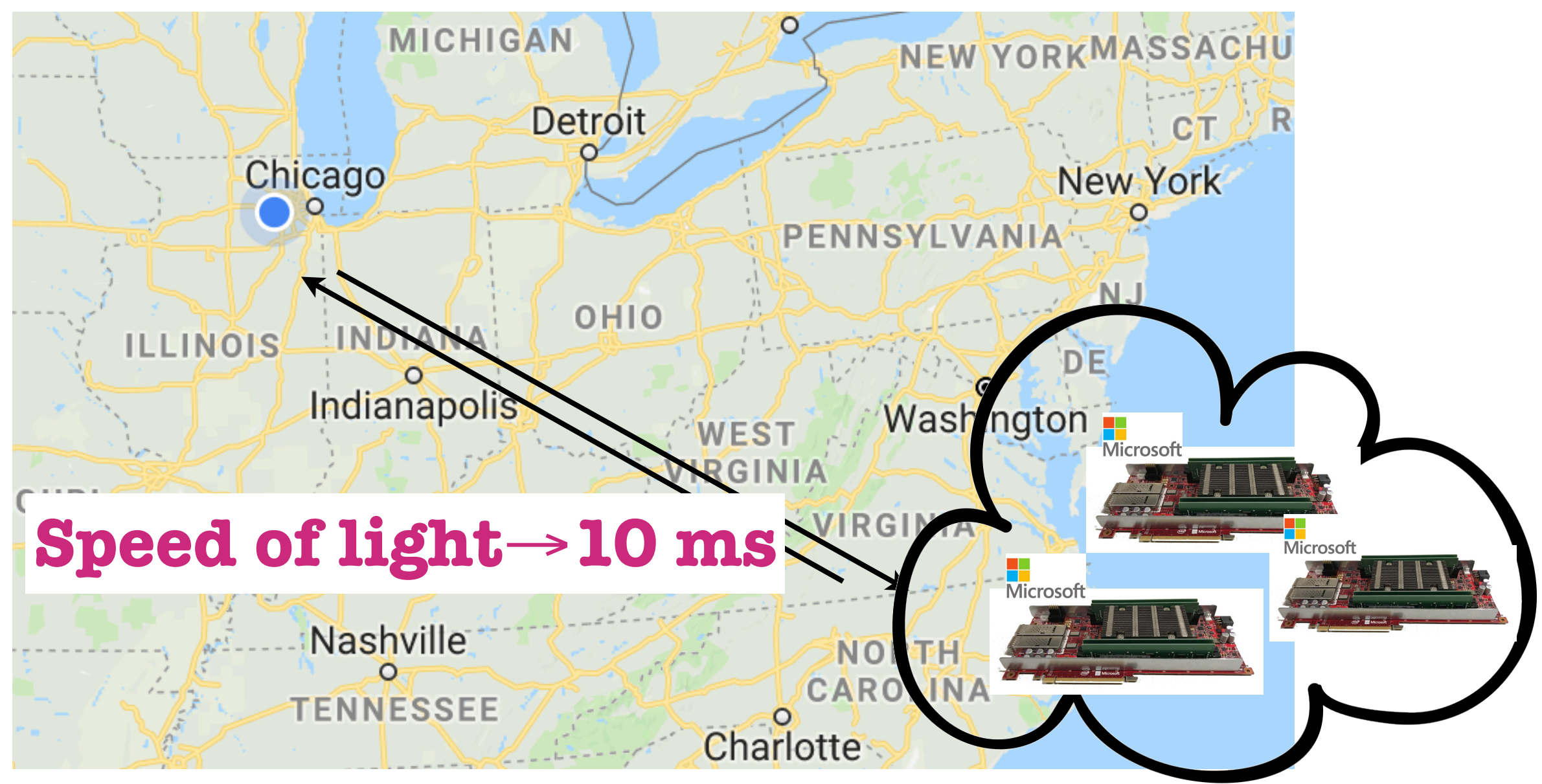
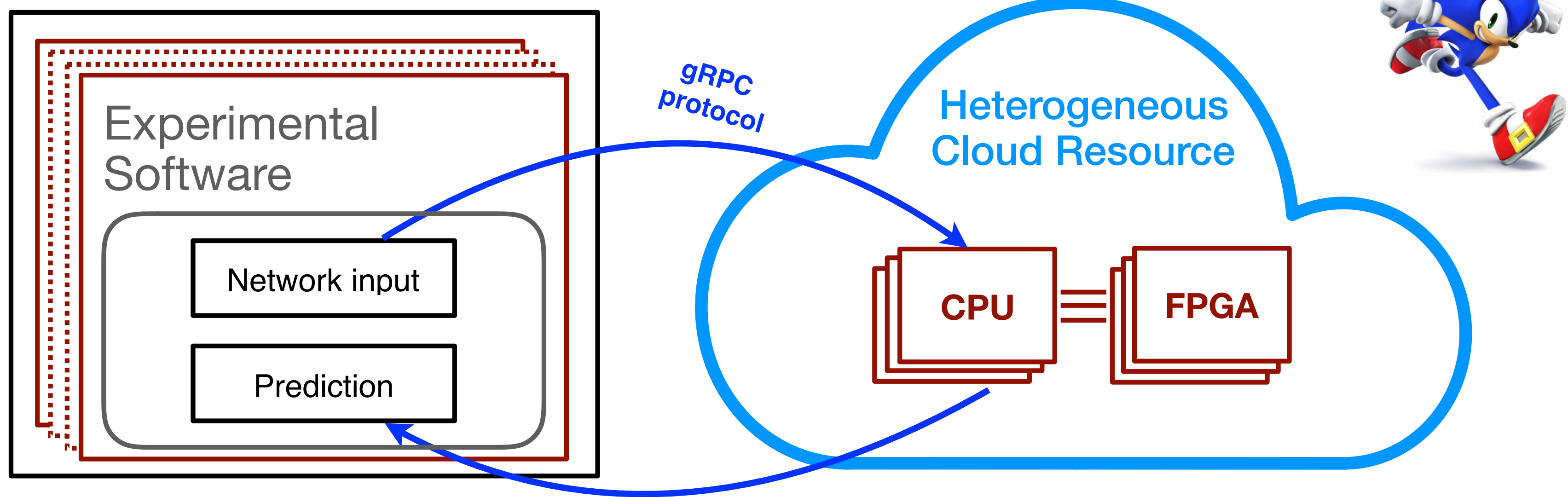


Deploy MS Brainwave as a service:

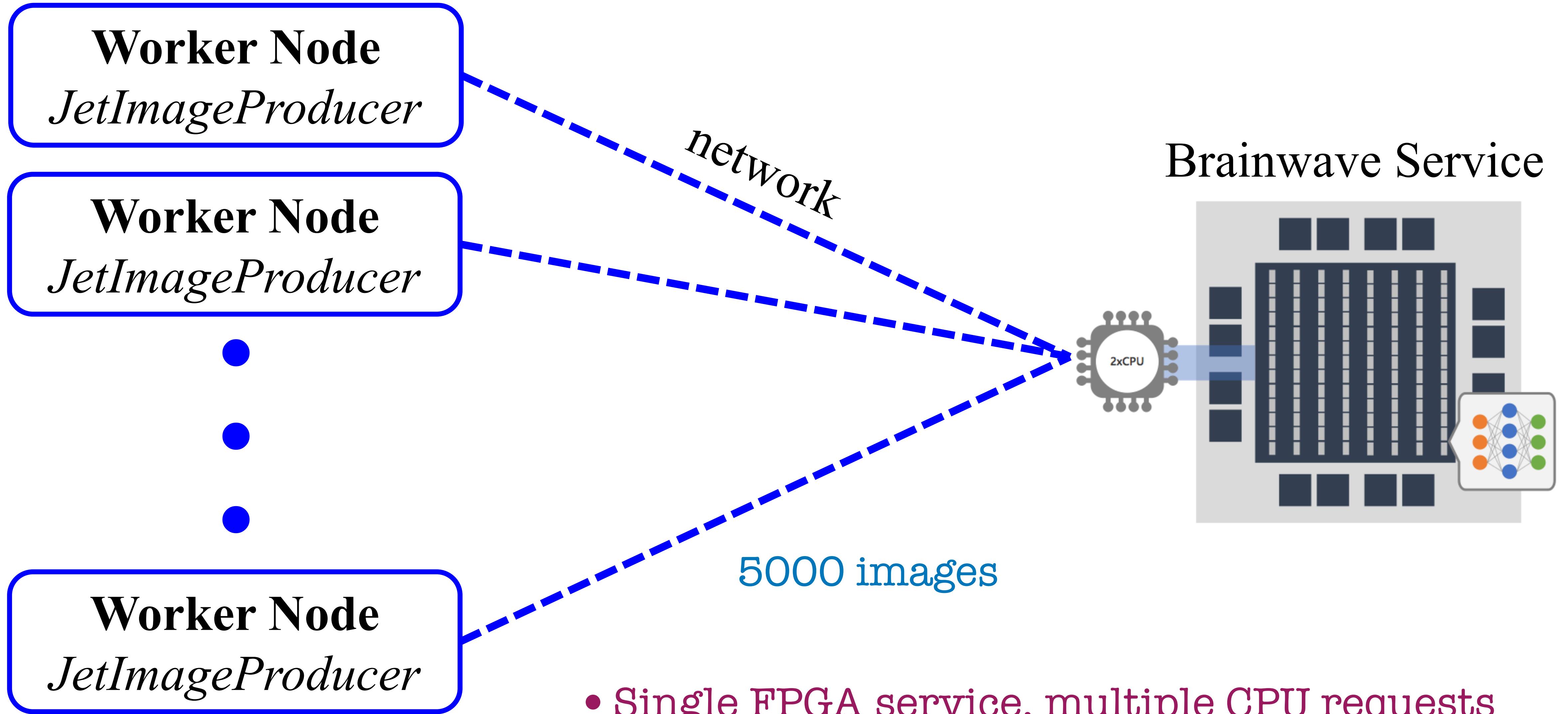
- Implemented with CMSSW ExternalWork module
- Fits CMS computing model in a non-disruptive way

SINGLE INFERENCE SPEED TESTS

Datacenter (CPU farm)

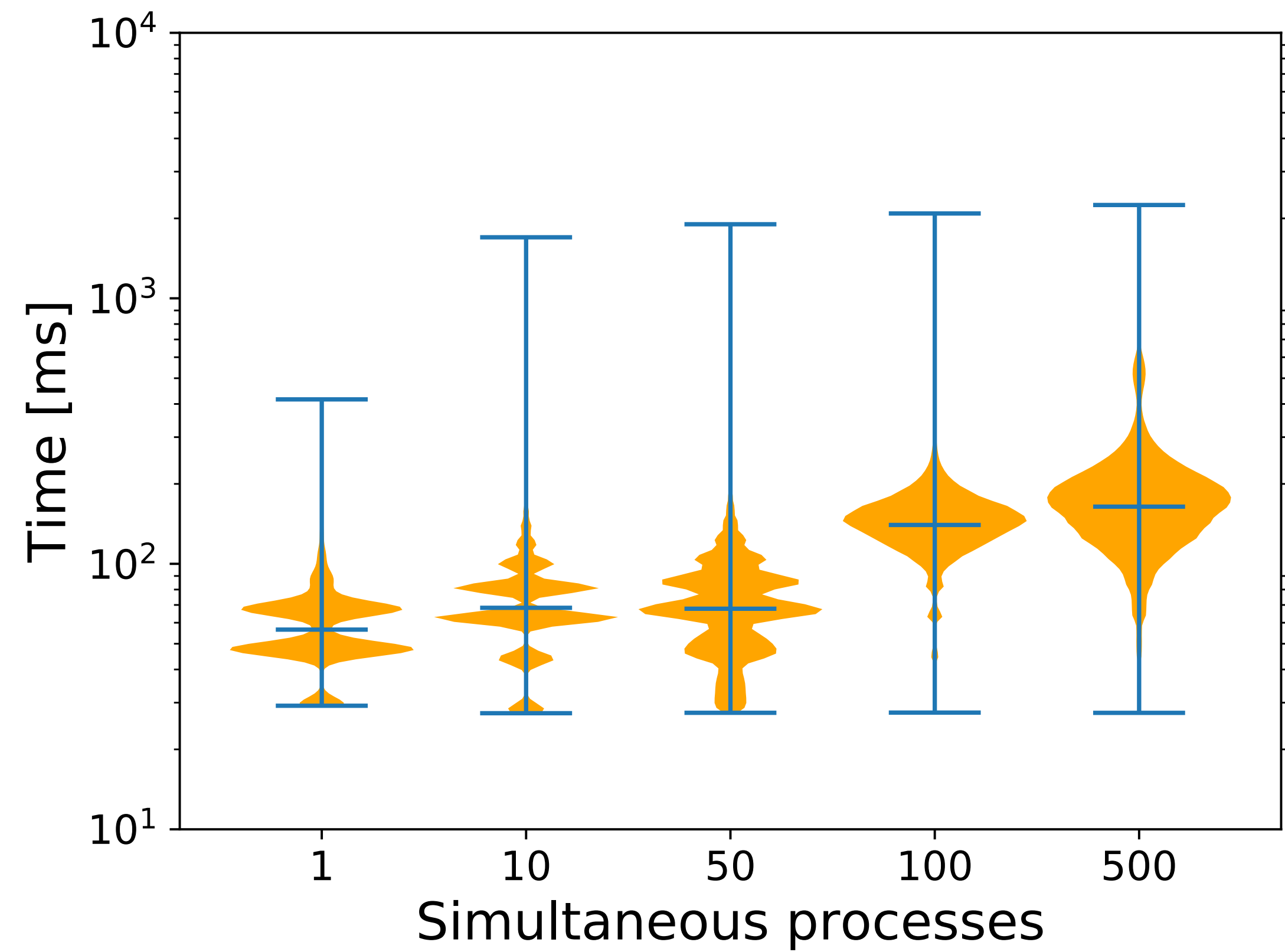
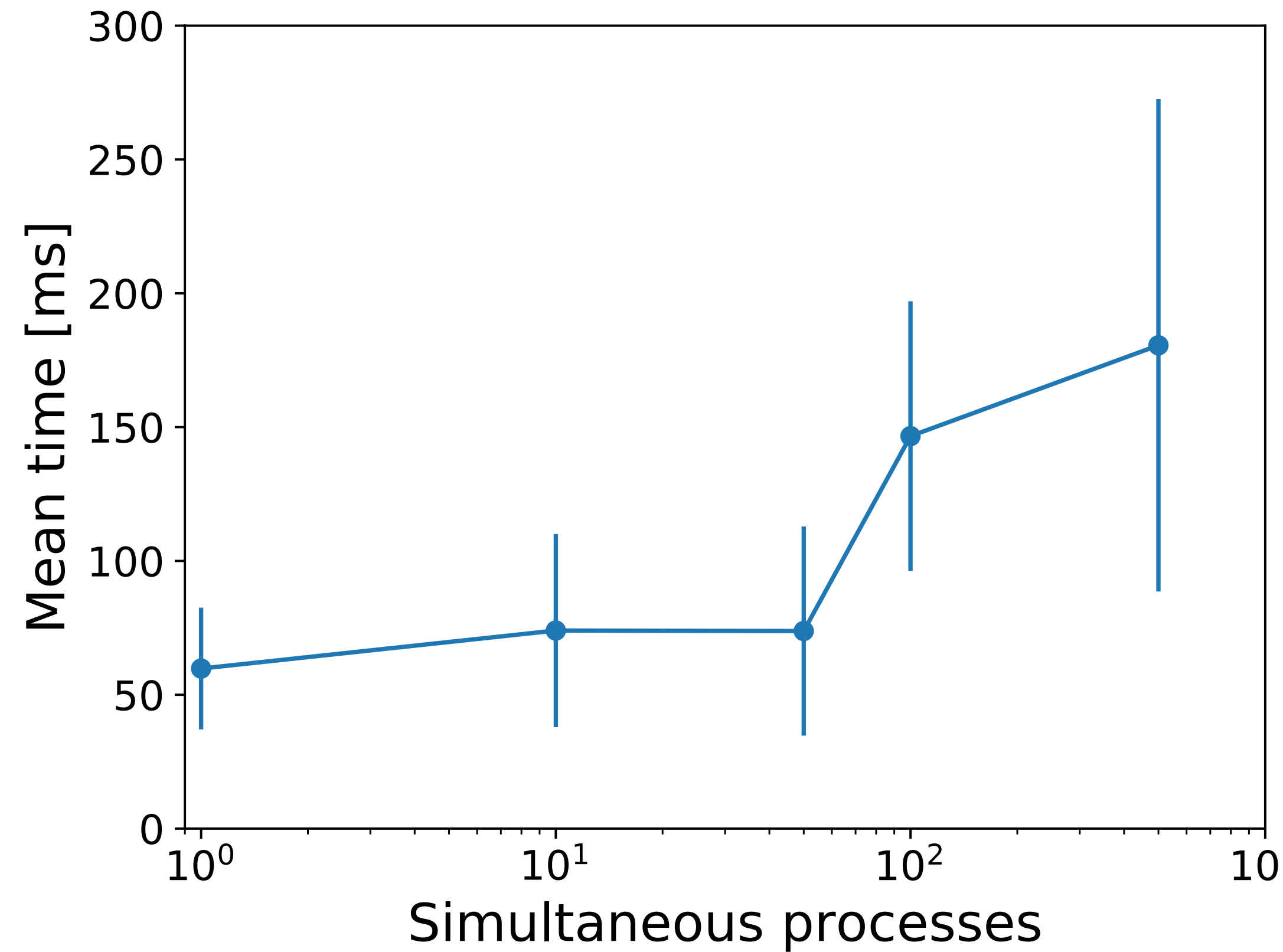


Test	Inference time
local	<ul style="list-style-type: none"> • 10 ms (~ 2 ms on FPGA + classifying, I/O) • Meets HLT latency requirement
remote	<ul style="list-style-type: none"> • 60 ms (includes travel latency) • (4/10/100) faster than CPU-only computations



- Single FPGA service, multiple CPU requests
- Each request sends 5000 images

N: simultaneous processes

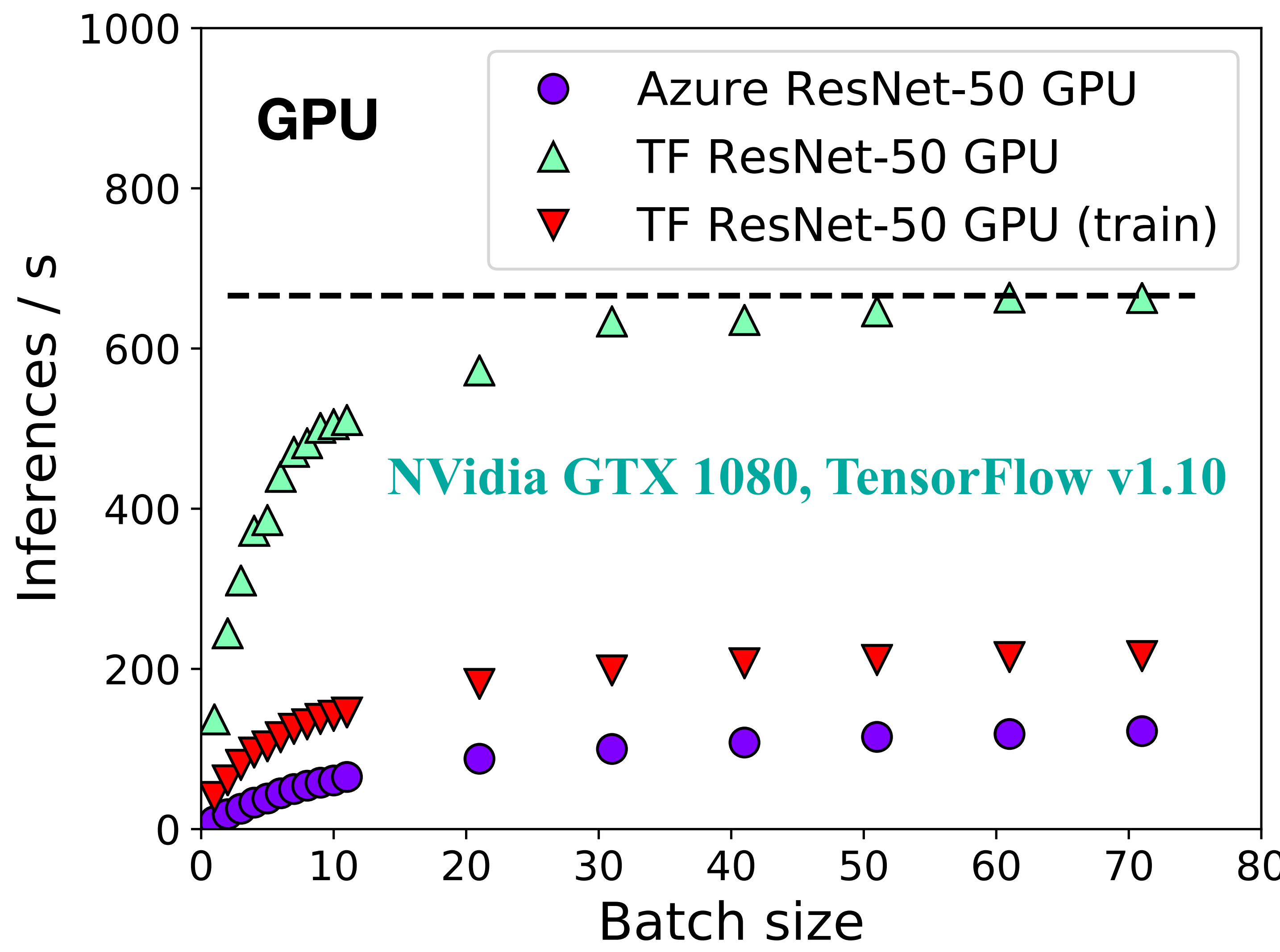
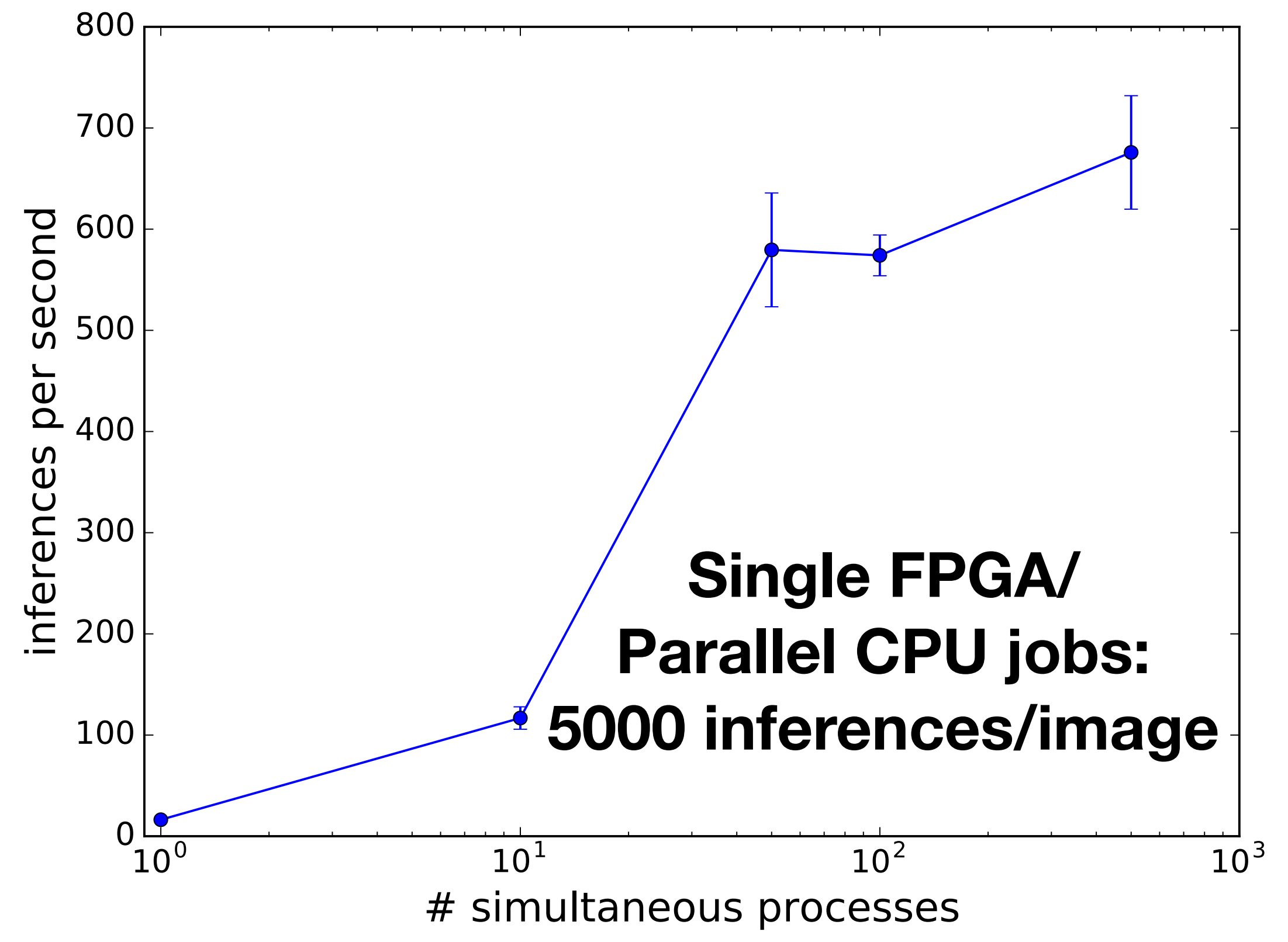


Tests: N=1,10,50,100,500

- Moderate increases in mean, standard deviation, and long tail for latency
- o Fairly stable up to N = 50

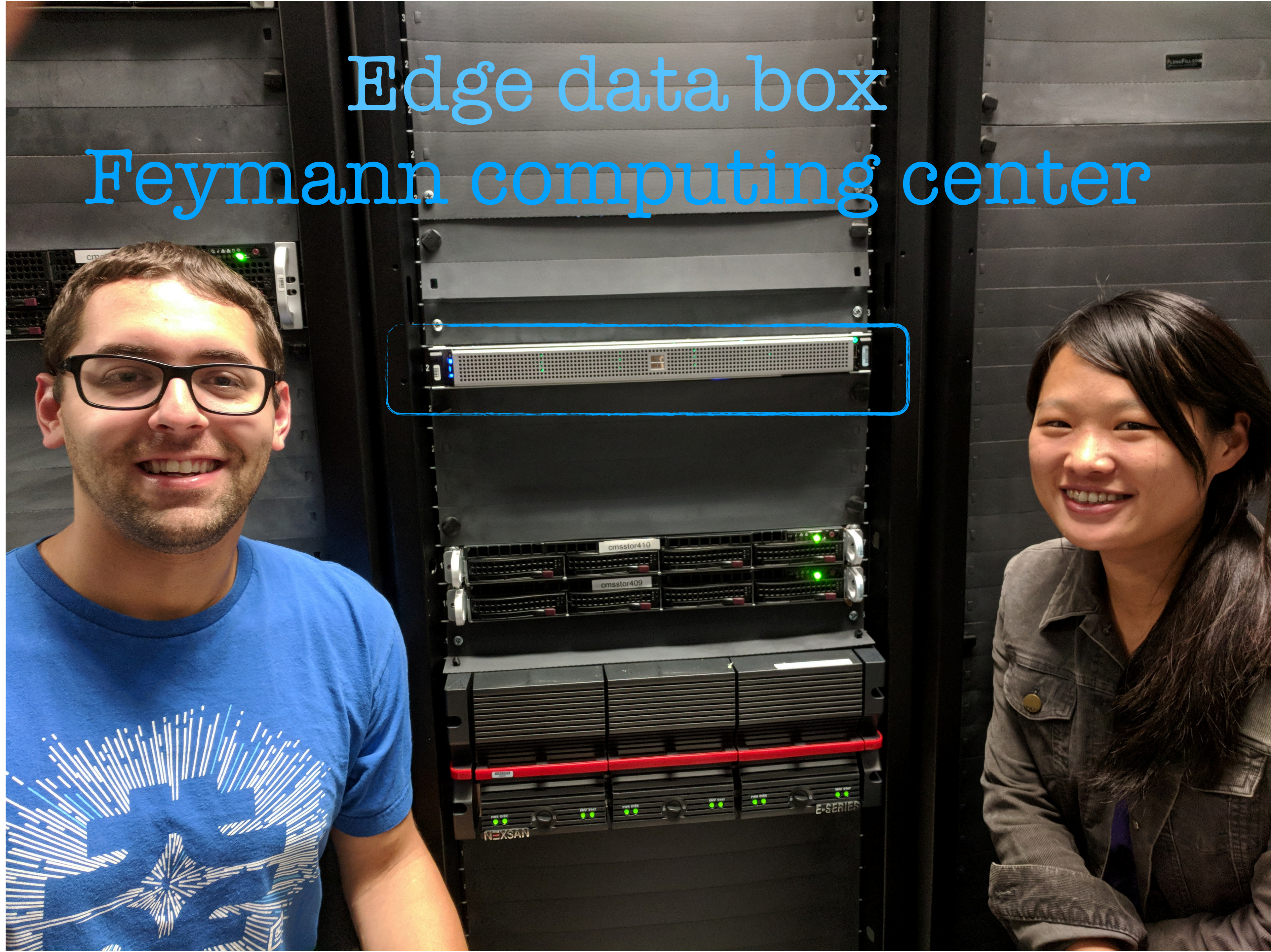
DATA THROUGHOUT COMPARED TO GPUS

Brainwave cloud service



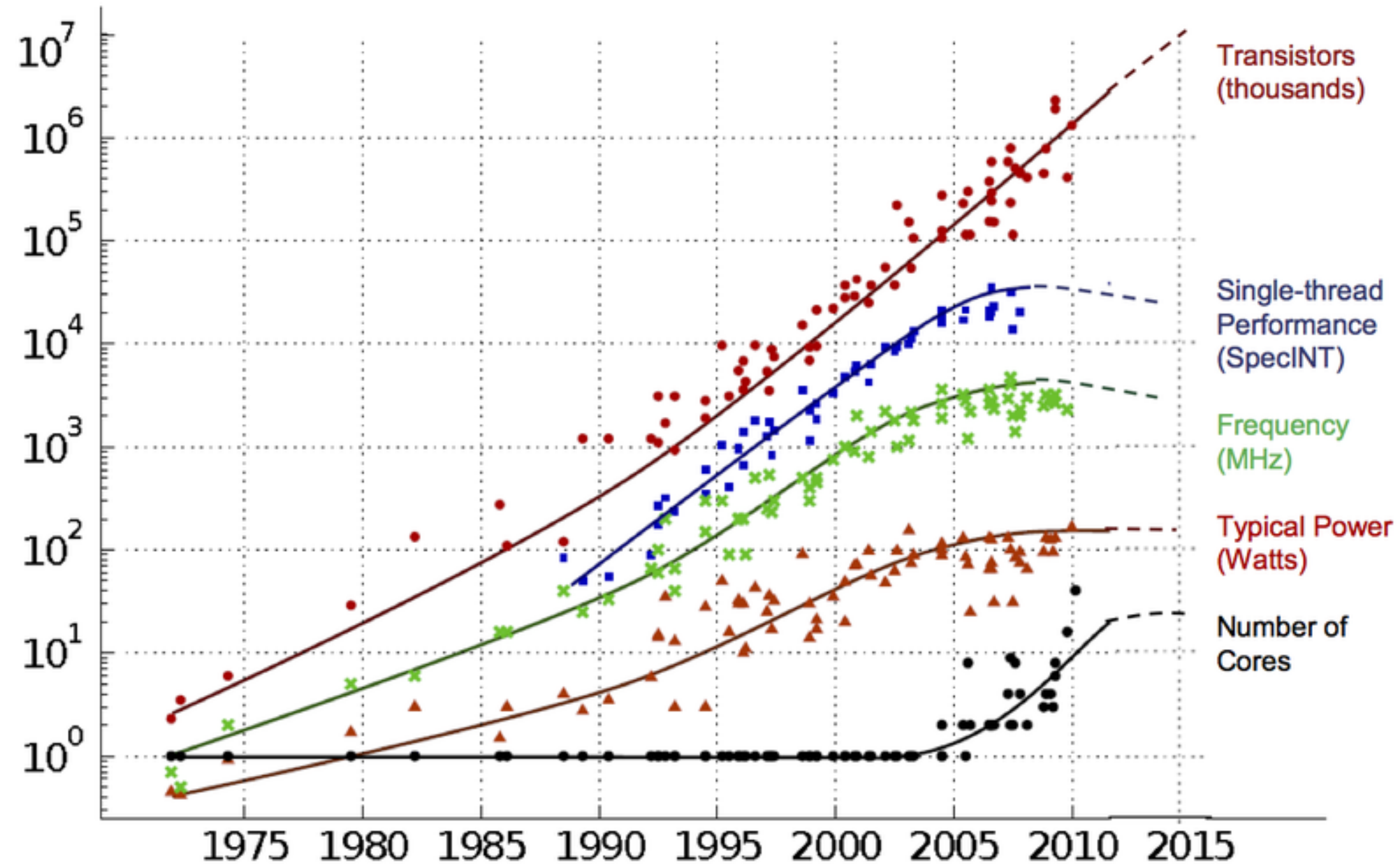
Comparable max data throughout: 600-700 images/

Edge data box Feymann computing center



Summary and outlooks

- **Proof of concept study with Brainwave:**
 - Integrate heterogeneous computing in our software framework
- **We are doing studies to benchmark other options (speed and scaling):**
 - Intel Open VINO, AWS, Google TPU...
- **Other (top) user's considerations for a dream Heterogeneous Computing Platform:**
 - Flexibility: Model support
 - e.g. Support for Graph Neural Networks
 - Support for ML framework
 - As a service: Cloud and Edge. Cost model.
- **Only works if we can solve our problems with ML! See showcases in the Lightning round talks.**



Moore's Law continues
...but Dennard Scaling fails

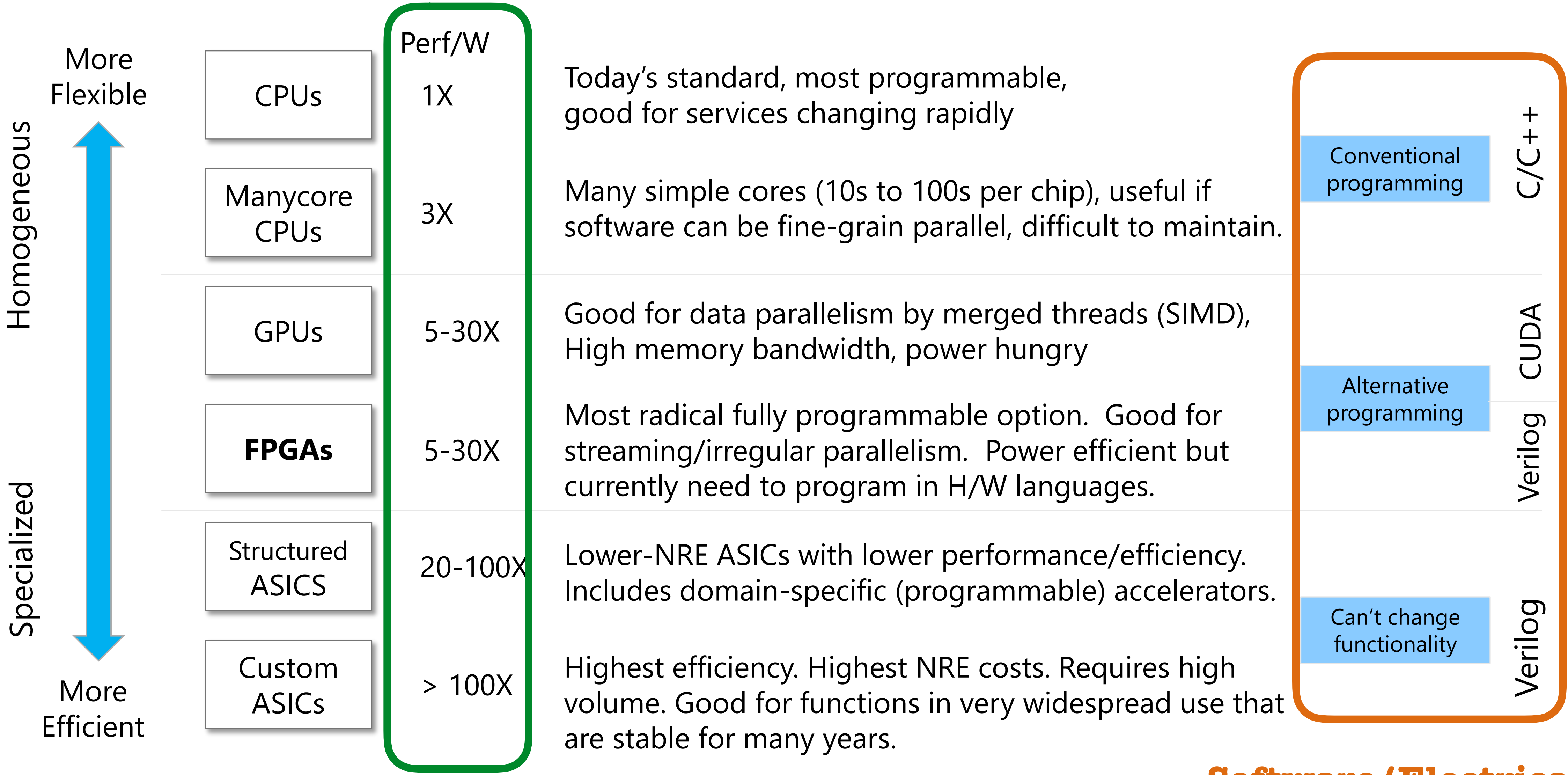


CPU

Original data collected and plotted by M. Horowitz, F. Labonte, O. Shacham, K. Olukotun, L. Hammond and C. Batten
Dotted line extrapolations by C. Moore

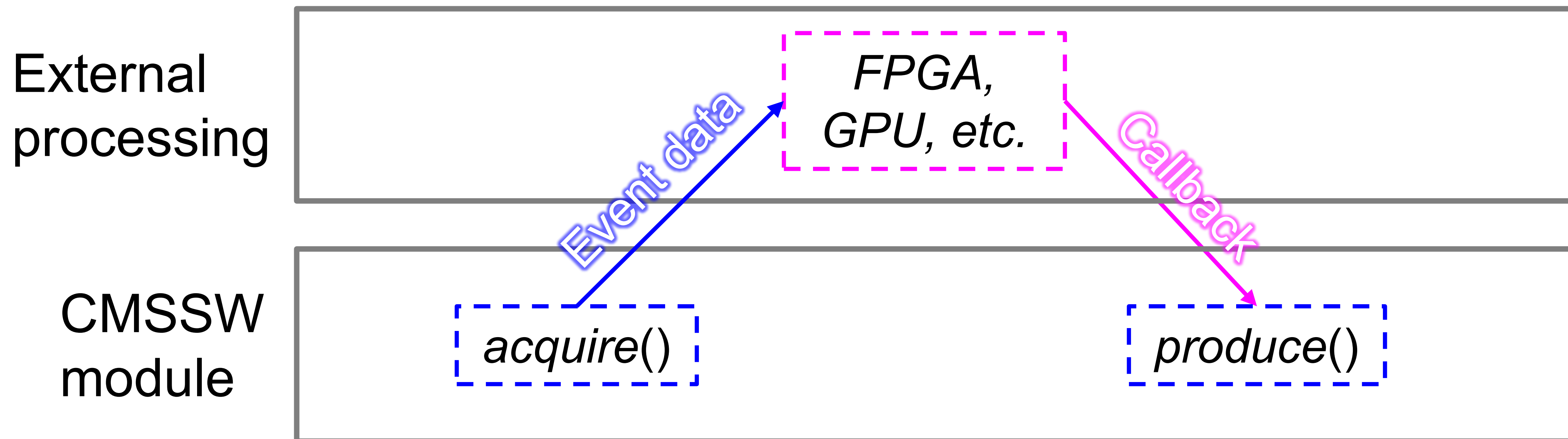
Single threaded performance not improving
~2005: "The Era of Multicore"

COMPUTING PLATFORM PROS & CONS:



Electricity bill

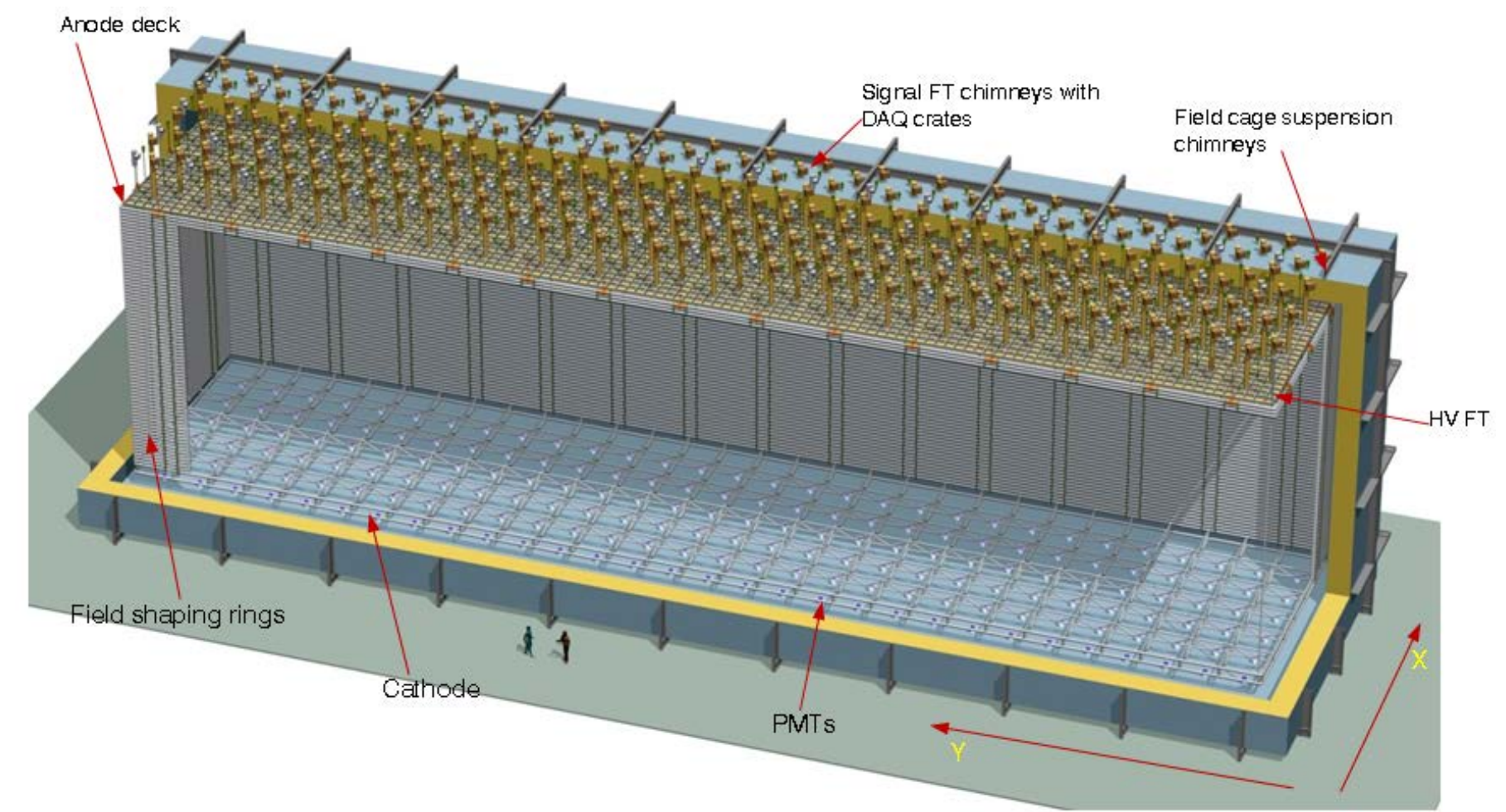
Software/Electrical engineer hours



Deploy MS Brainwave as a service:

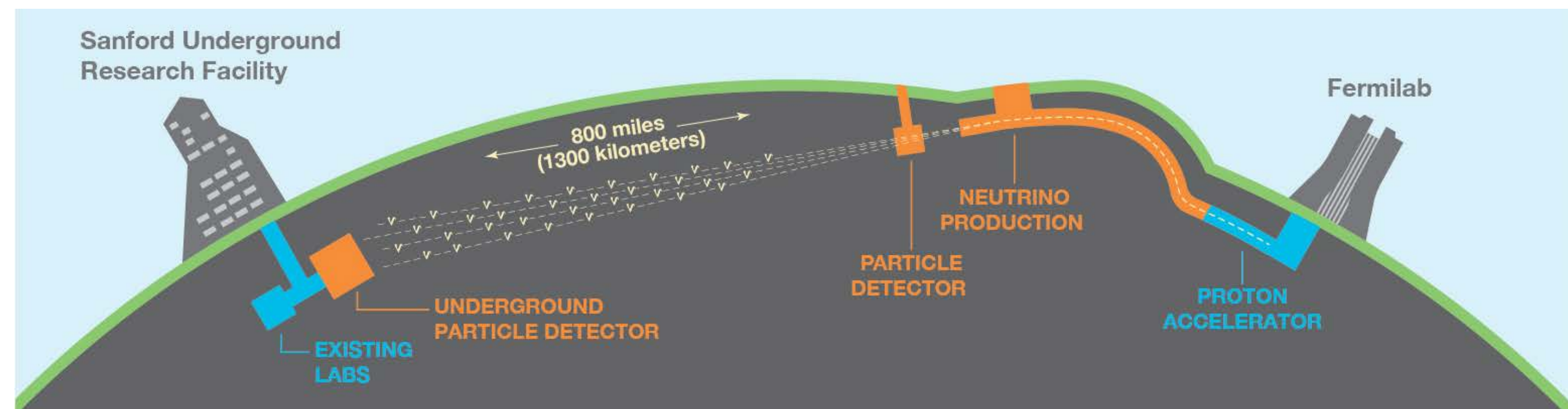
- Implemented with CMSSW ExternalWork module
- Fits CMS computing model in a non-disruptive way

Neutrino Computing Challenges



Intensity frontier: DUNE

- Largest liquid argon detector ever designed
- ~1M channels, 1 ms integration time w/ MHz sampling
→ 30+ petabytes/year



➤ CPU needs for particle physics will increase by *more than an order of magnitude* in the next decade

Example: First observation of $t\bar{t}H$ using particle Run 2 data.

$t\bar{t}H$ production

