

# Real-time machine learning in embedded systems for particle physics

Nhan Tran

Fermilab/Northwestern ECE

October 12, 2020



# Outline

**Disclaimers & Motivations**

**Real-time opportunities**

**Efficient machine learning**

**Example applications**

**LHC Trigger, Accelerator Controls, Reconfigurable AI ASICs**

# Disclaimers & Motivation

# Disclaimers

- **Disclaimer 1:**

- This is a **very exciting and rapidly moving** field! I apologize for any **biases** I have based on my research focus. I am excited to hear about other work in this space (I am member of CMS experiment at LHC and based at Fermilab)

- **Disclaimer 2:**

- Though I will talk about tools, methods, and example applications - it is still early days - **much is still unexplored**

- **Disclaimer 3:**

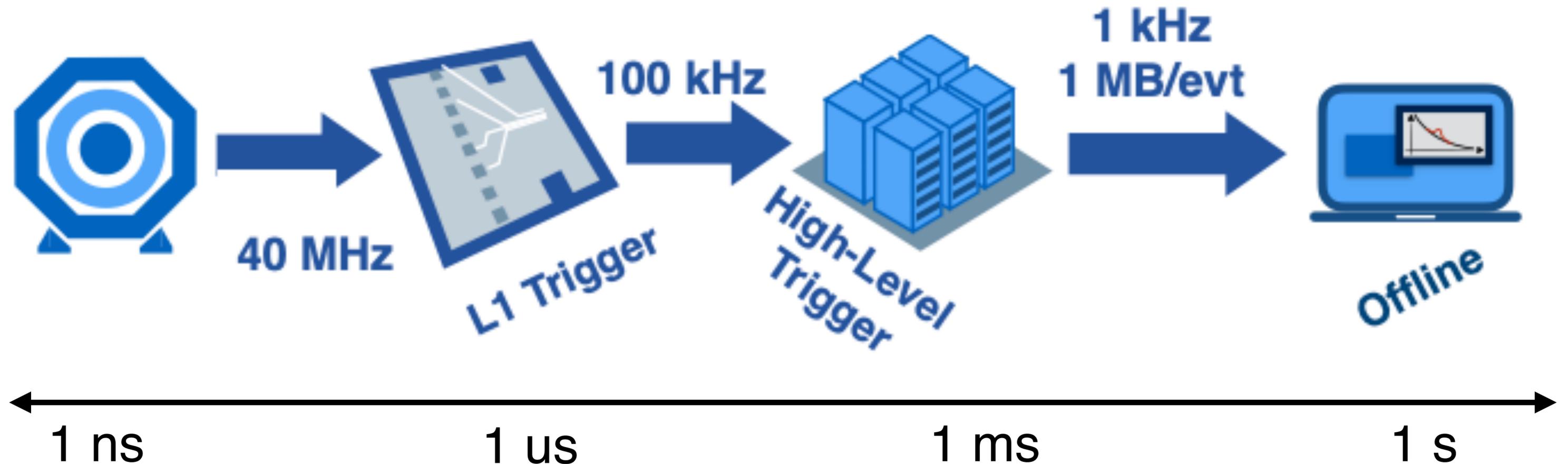
- I have cut off the scope of this talk at embedded systems — but FPGAs and ASICs are also used for heterogeneous computing and have common techniques and lessons, **see Mia Liu's talk on Tuesday**

# Motivation

(a biased view that I hope you share!)

- Cutting edge science requires...  
**Faster, more precise, bigger, more granular,...**
- Science at the ~MHz scale and beyond creates massive amounts of data  
**Online compression, featurization, filtering, feedback, etc. is vital - arguably more so than offline machine learning**
- Real-time intelligence can greatly accelerate science allowing us to test hypotheses faster, enhance performance of detectors/accelerators, and save potentially lost data

# CMS real-time processing



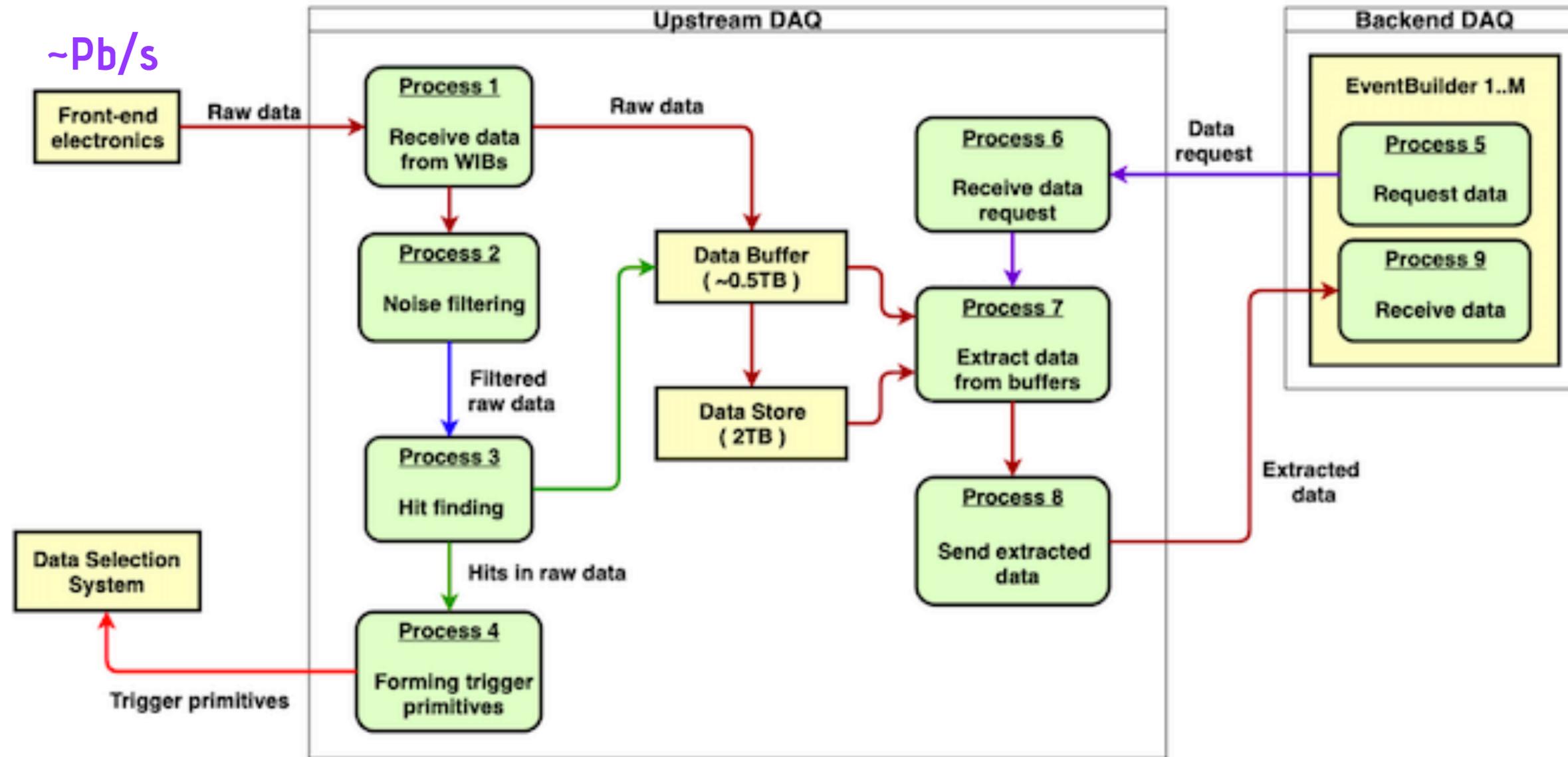
Custom electronics  
Latency  $\sim 25\text{ns} - 10\ \mu\text{s}$

FPGAs/ASICs - high bandwidth low latency specialized compute hardware

Off-the-shelf computing  
Latency  $\sim 0(1+ \text{ms})$

“standard” CPU computing, coprocessors

# DUNE upstream DAQ



# Accelerator controls

Fermilab accelerator complex

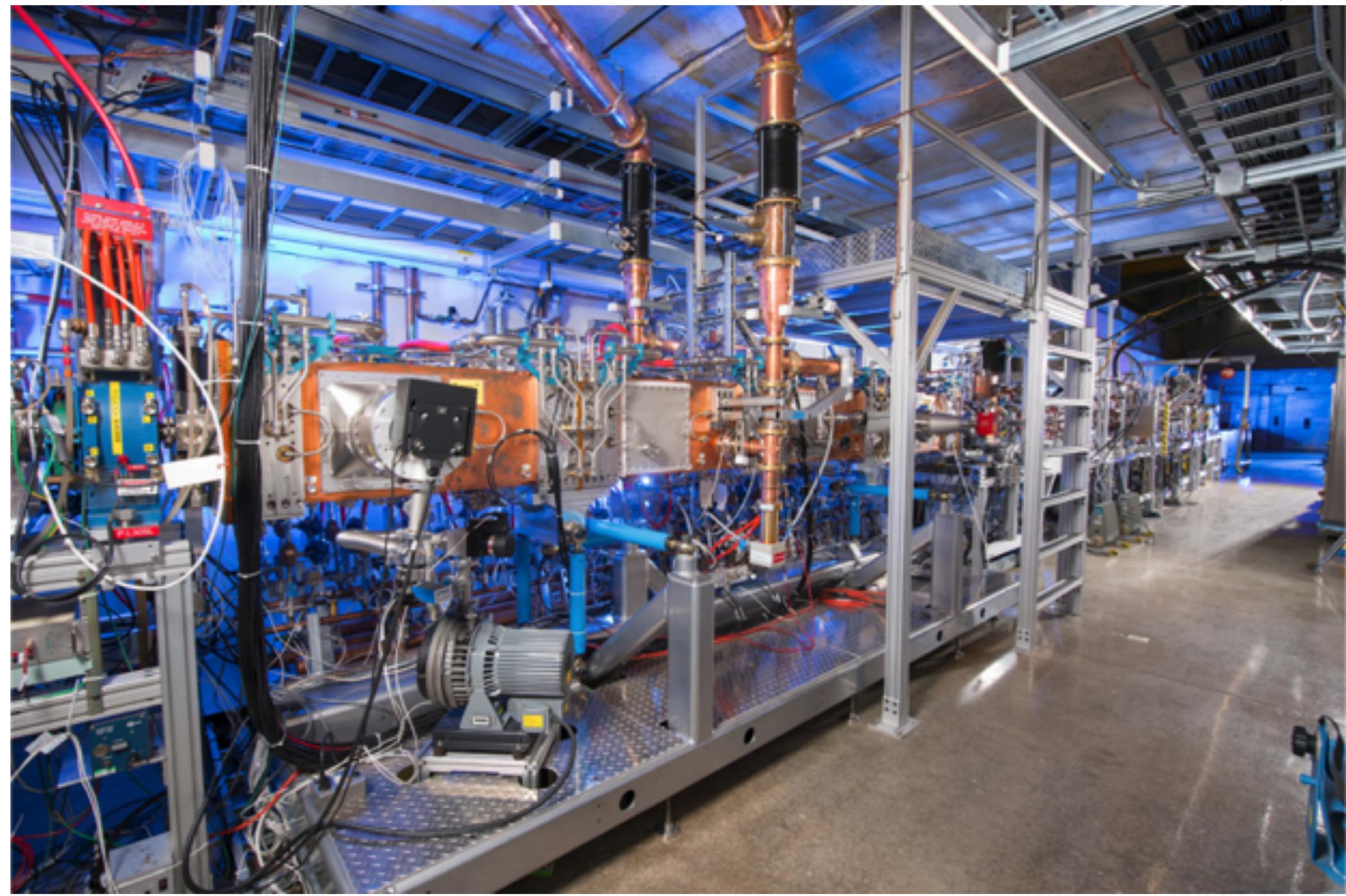


300,000 devices monitored

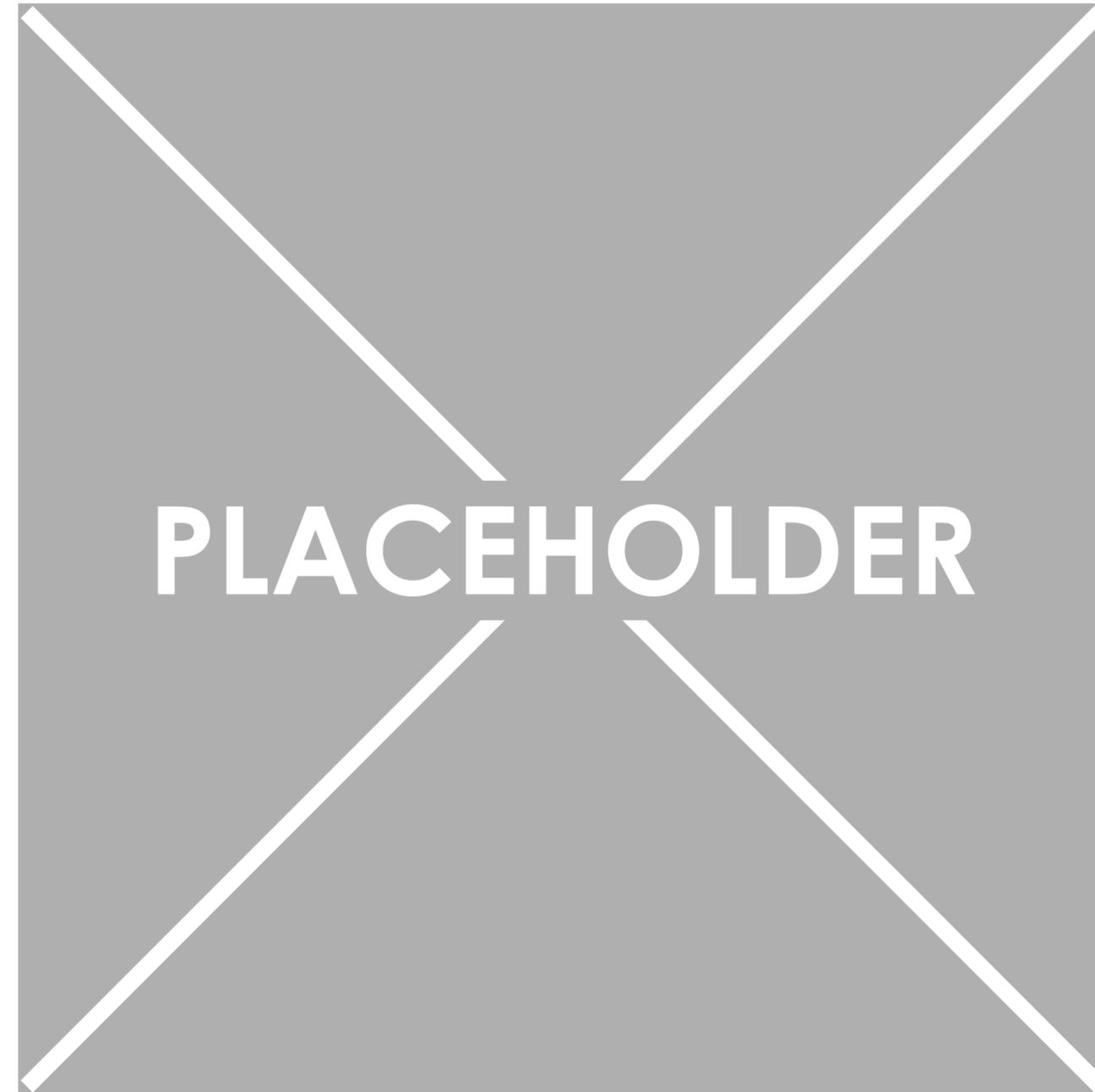
18 miles of accelerator

25,000 daily alarms

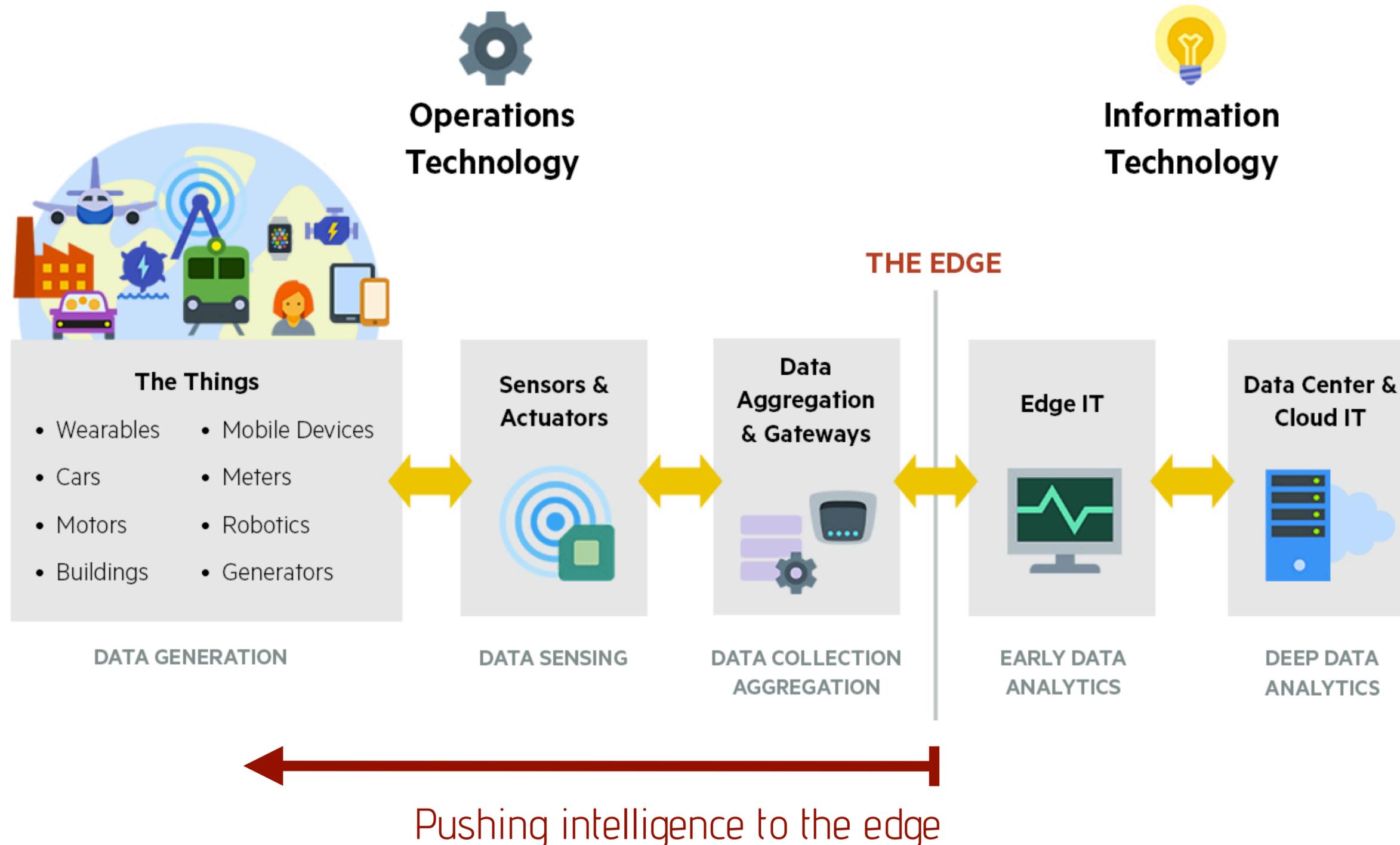
The left column contains three vertically stacked images. The top image shows a server room with rows of racks and glowing green lights. The middle image is an aerial view of a large industrial or research facility with a circular pond. The bottom image shows a woman in a red shirt sitting at a control desk with multiple computer monitors.



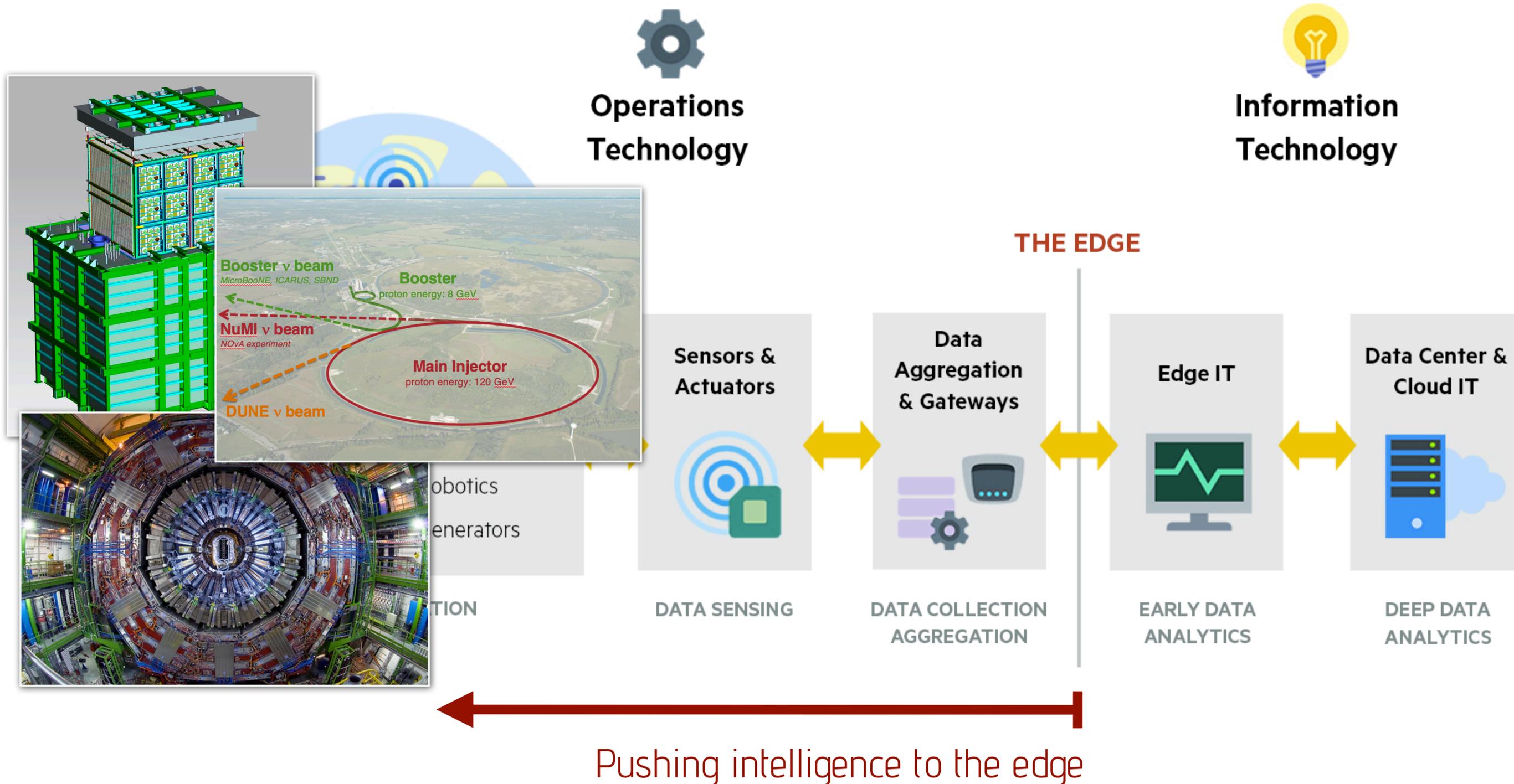
# Insert your application here



# Internet of things...particle physics



# Internet of things...particle physics



# Efficient machine learning

{implementations, optimizations, tool flows for}

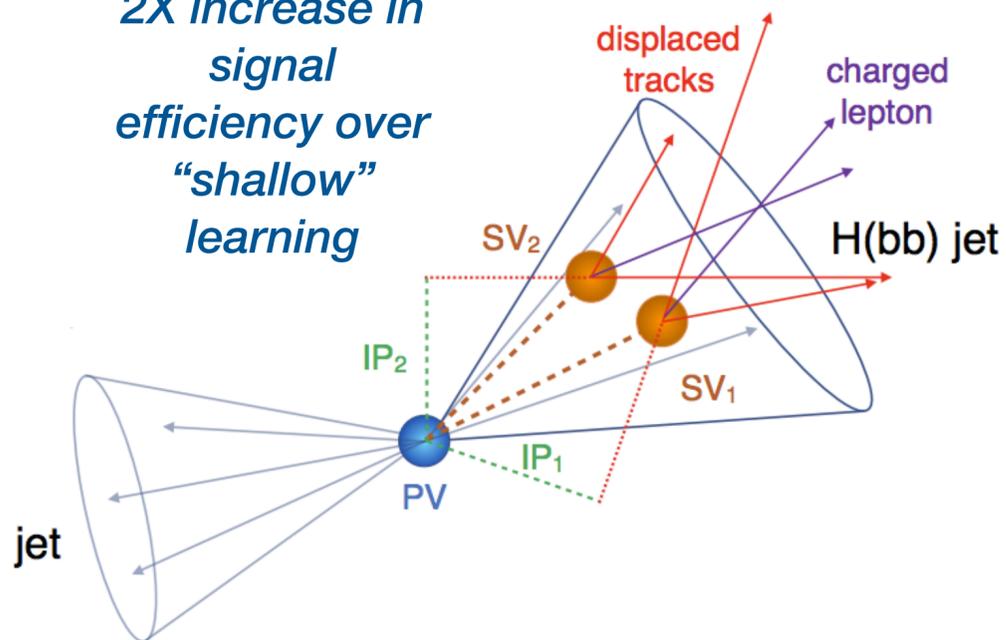
**Efficient<sup>^</sup> machine learning**

# Machine learning

## Energy

Identification of boosted Higgs jet decay to two bottom quarks

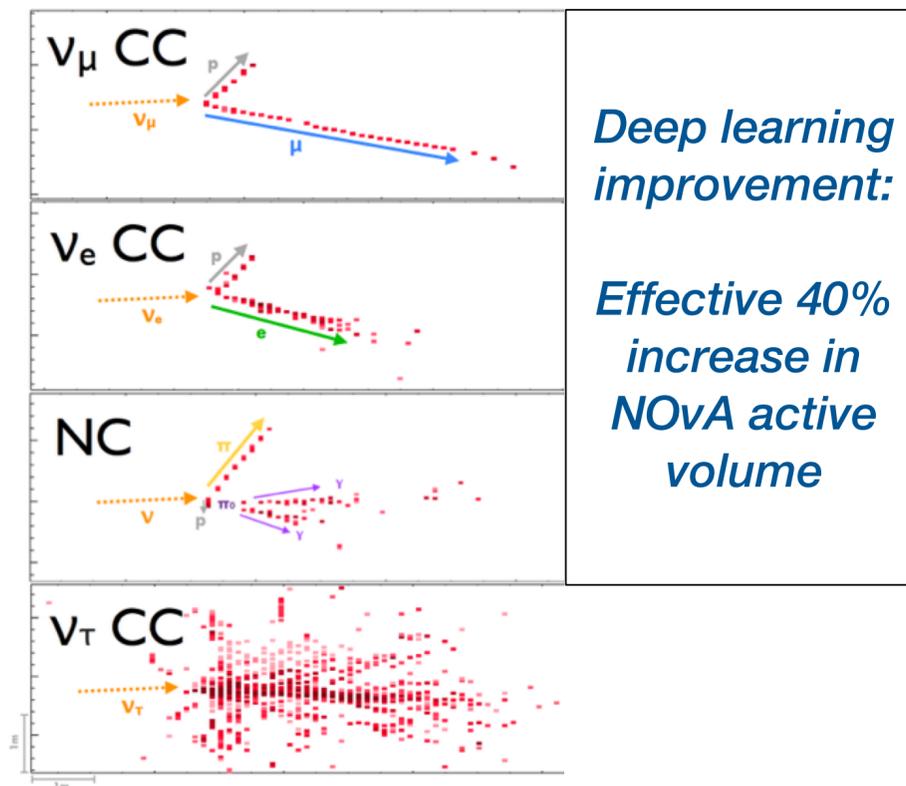
*2X increase in signal efficiency over "shallow" learning*



J. Duarte et al., CMS DP-2018/046

## Intensity

NOvA event classification



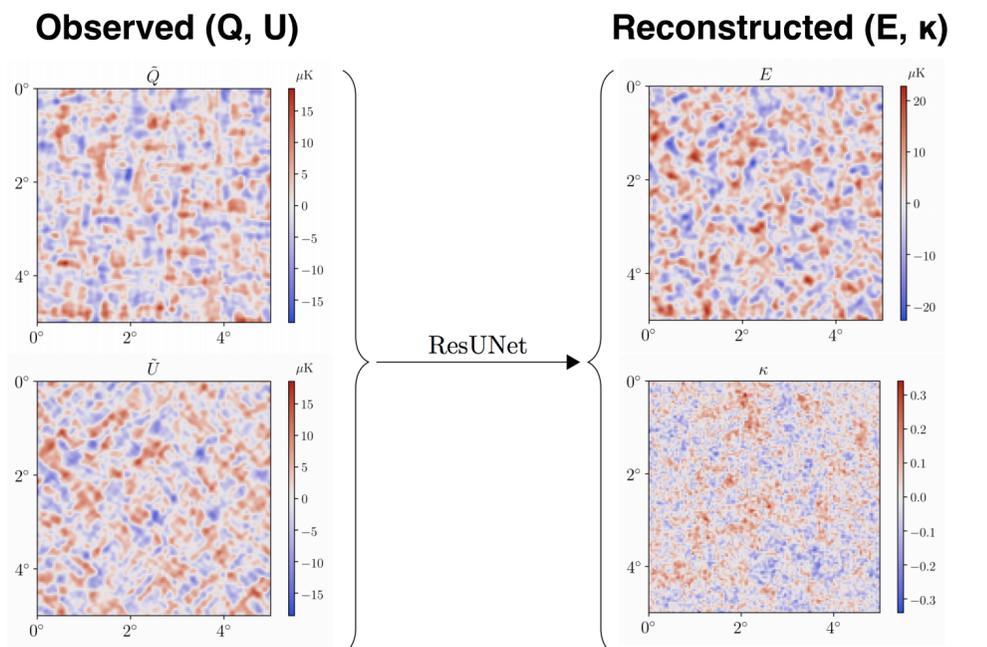
*Deep learning improvement:*

*Effective 40% increase in NOvA active volume*

A. Himmel, E. Niner, F. Pshihias et al.  
<https://arxiv.org/abs/1604.01444>  
 1st deployed in oscillation analysis  
<https://arxiv.org/abs/1703.03328>

## Cosmic

Reconstruction of CMB polarization map from Stokes parameters

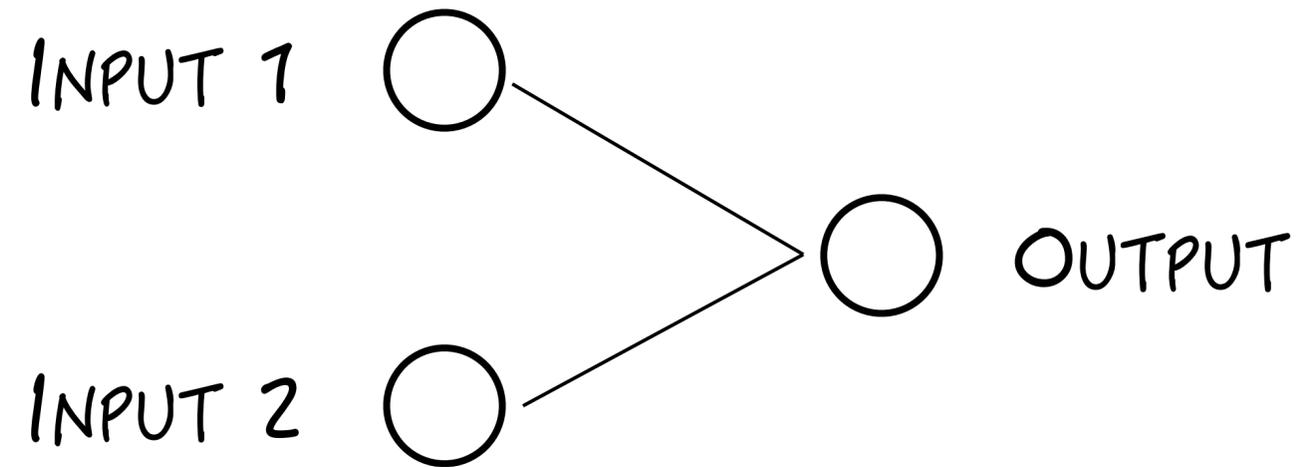


*50% less noise vs. traditional methods across large range of scales*

J. Caldeira, B. Nord, et al.,  
<https://arxiv.org/abs/1810.01483>

**ML applications is huge field by itself - see talk by Christof Buehler on Thursday**  
**Many open opportunities for new applications as well!**

# NN inference in a nutshell

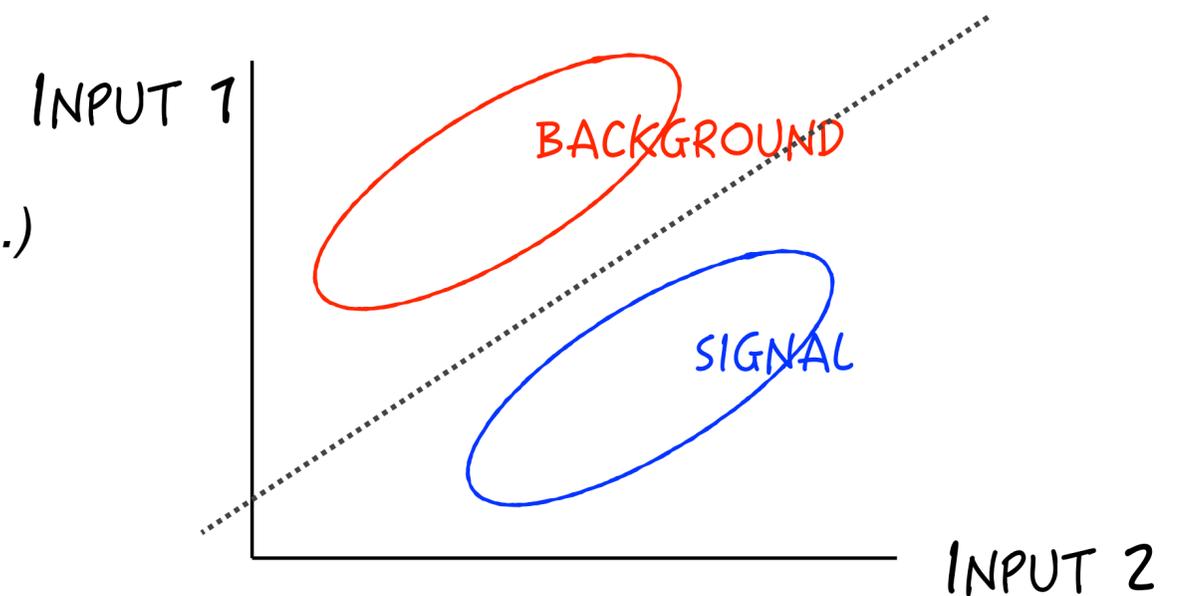


$$\vec{O}_j = \vec{l}_i \times \vec{W}_{ij} + \vec{b}_j$$

Simple 2 input example

(Fisher linear discriminant, linear support vector machine,...)

$$O_1 = l_1 \times W_{11} + l_2 \times W_{21} + b_1$$

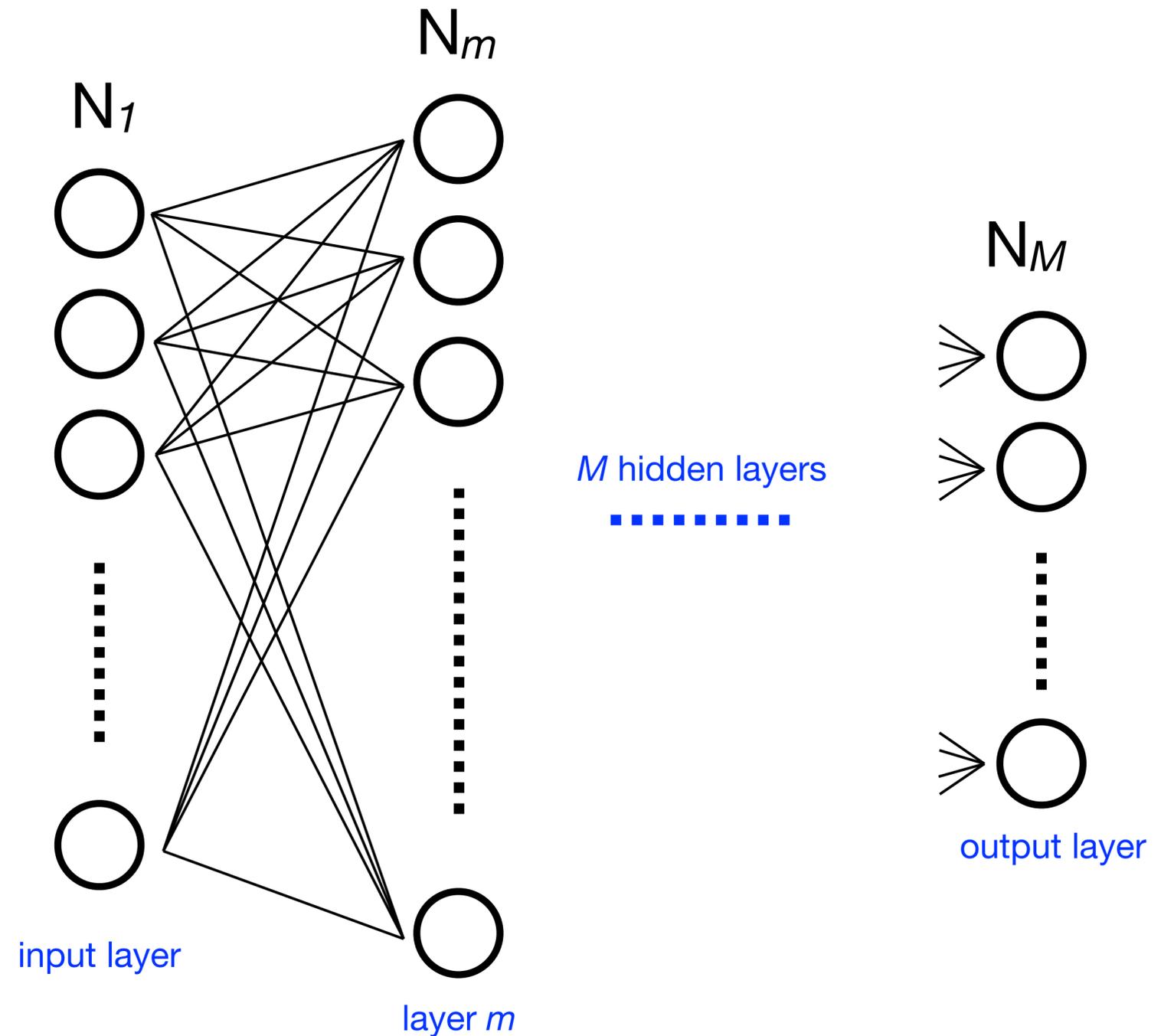


# NN inference in a nutshell

$$\vec{O}_j = \Phi(\vec{l}_i \times \vec{W}_{ij} + b_j)$$

$\Phi$  = ACTIVATION FUNCTION  
(NON-LINEARITY)

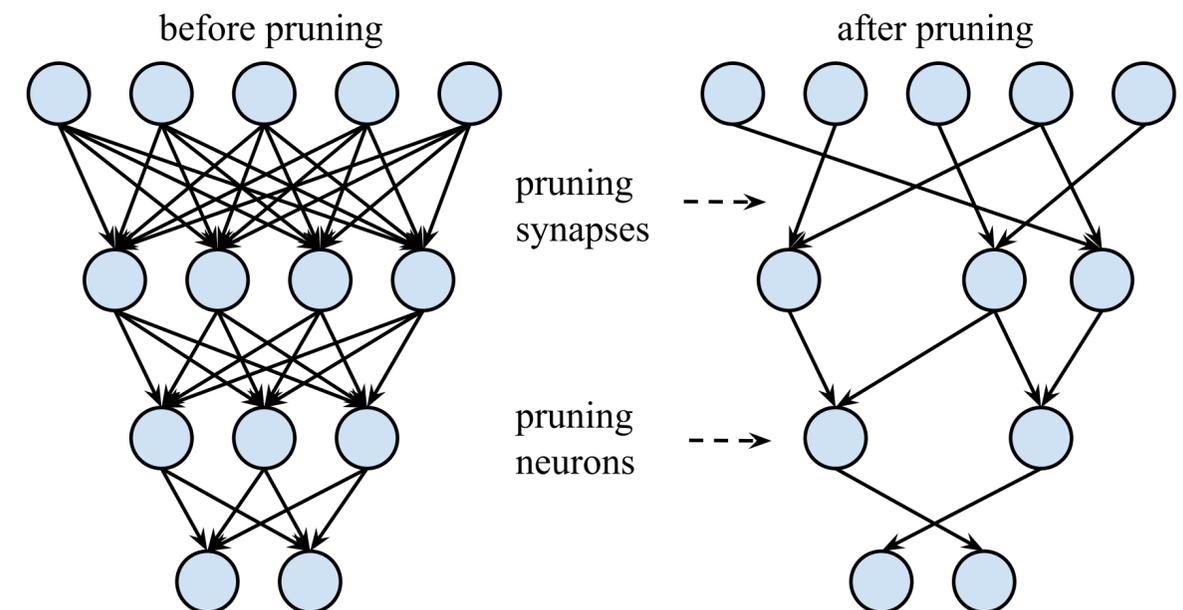
**NN inference =  
a bunch of multiplications /additions  
and LUTs (look up tables) for activation  
functions**



FULLY CONNECTED HIDDEN LAYER

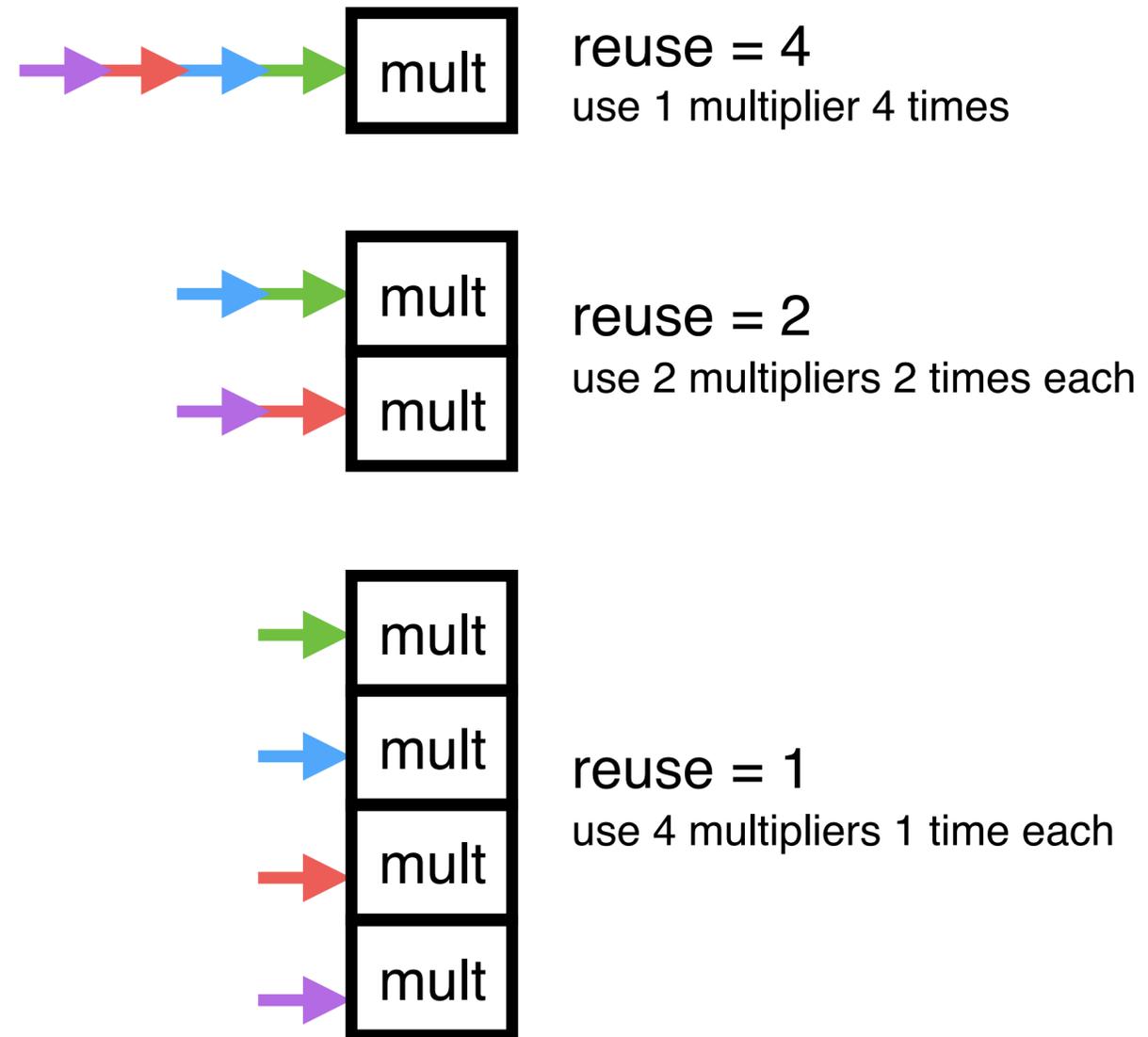
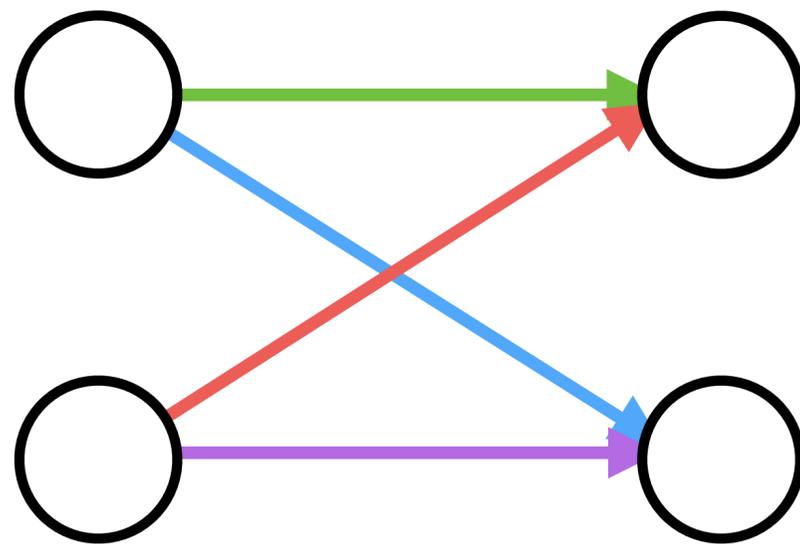
# (Energy) Efficient Neural Networks

- Important engineering field, efficient implementation of NN architecture
- **Parallelization**: performing operations simultaneously
- **Compression/Pruning**:
  - maintain the same performance while removing low weight synapses and neurons (many schemes)
- **Quantization/Approximate math**:
  - 32-bit floating point math is overkill
  - 20-bit, 18-bit, ...? fixed point, integers? binarized NNs?



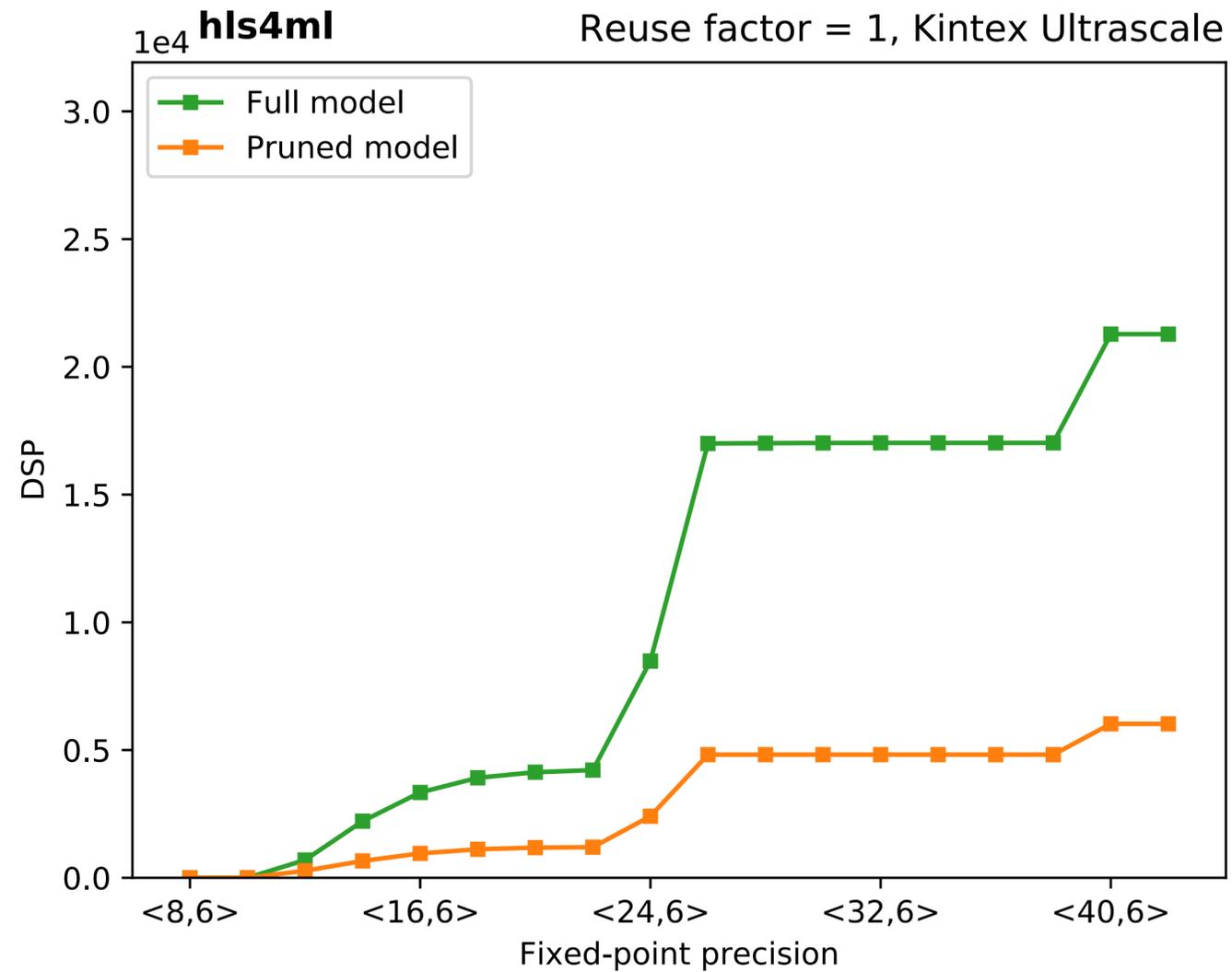
# Example: Parallelization

**ReuseFactor:** how much to parallelize operations a hidden layer

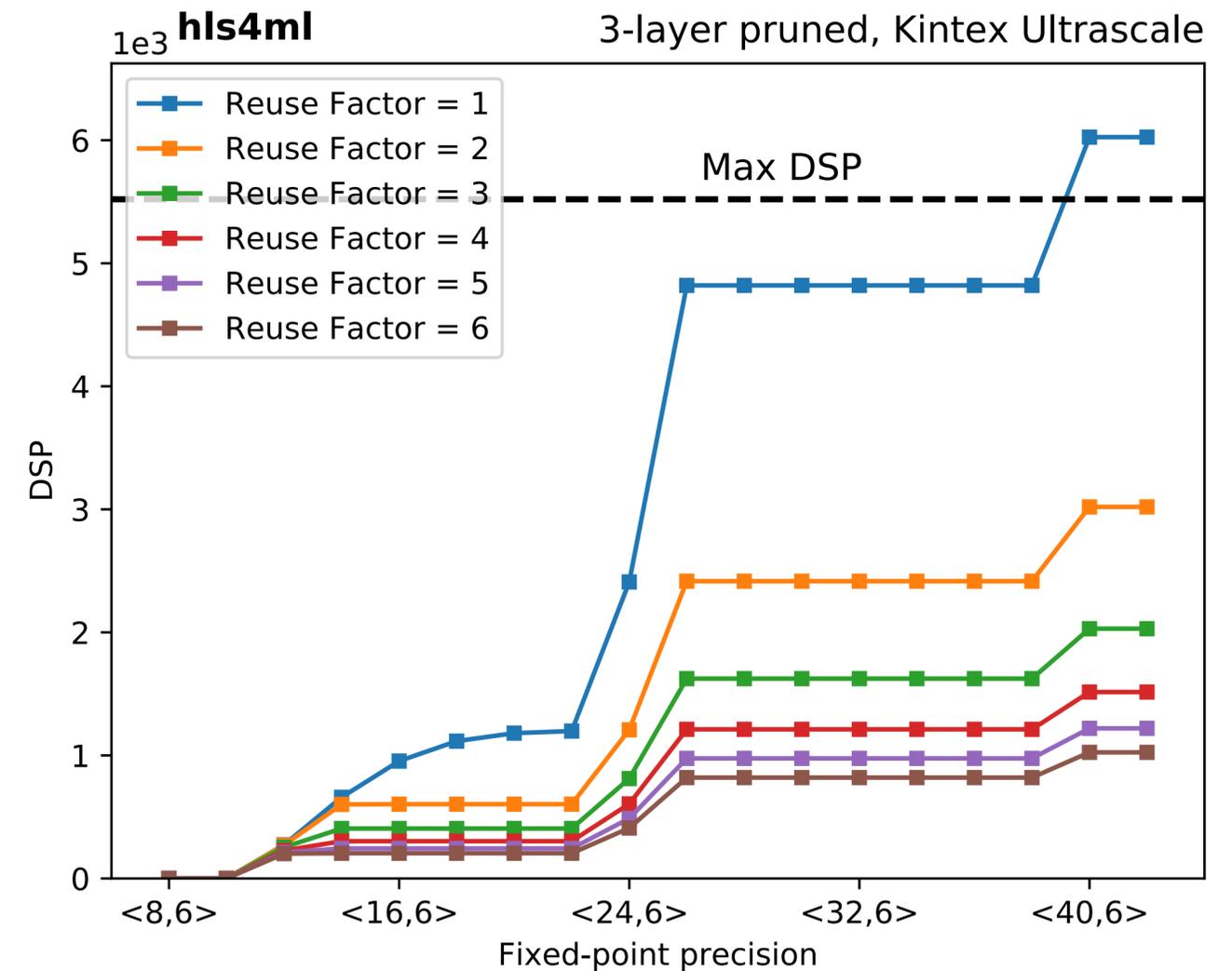


# Parallelization, pruning, quantization

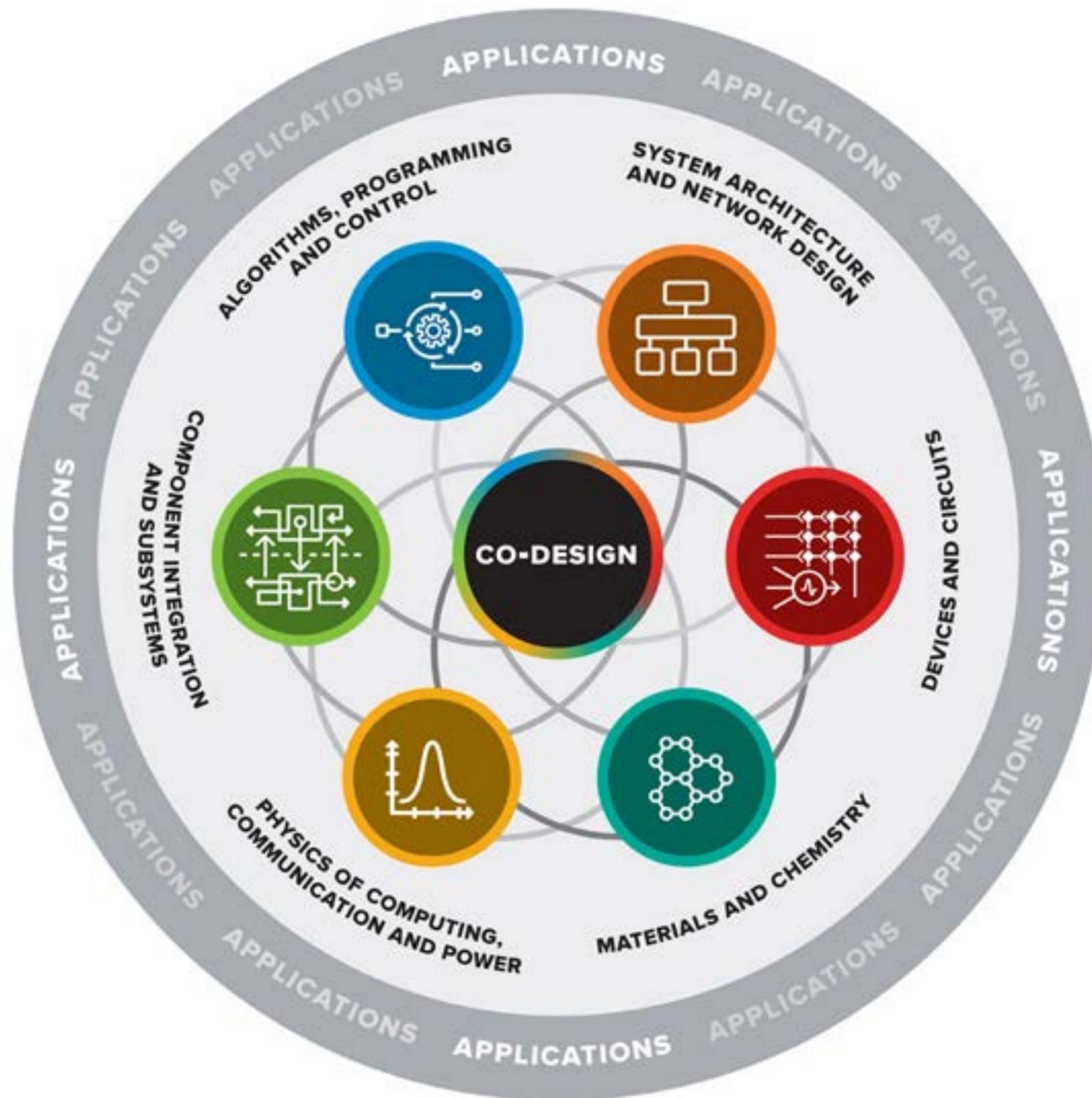
Increased unroll factor to  
reduce DSP usage  
(at latency cost)



Unstructured iterative pruning  
with L1 regularization

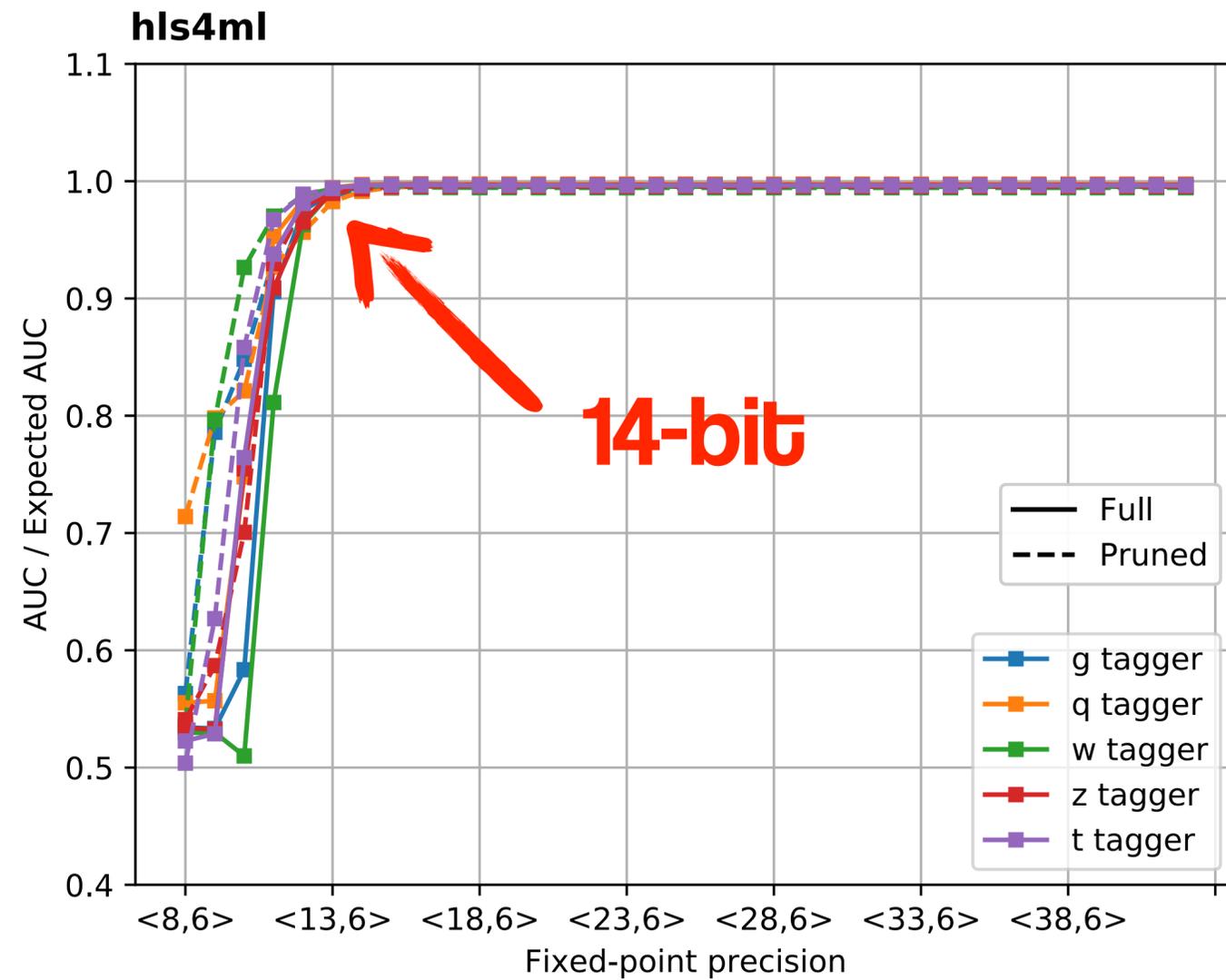


# Codesign



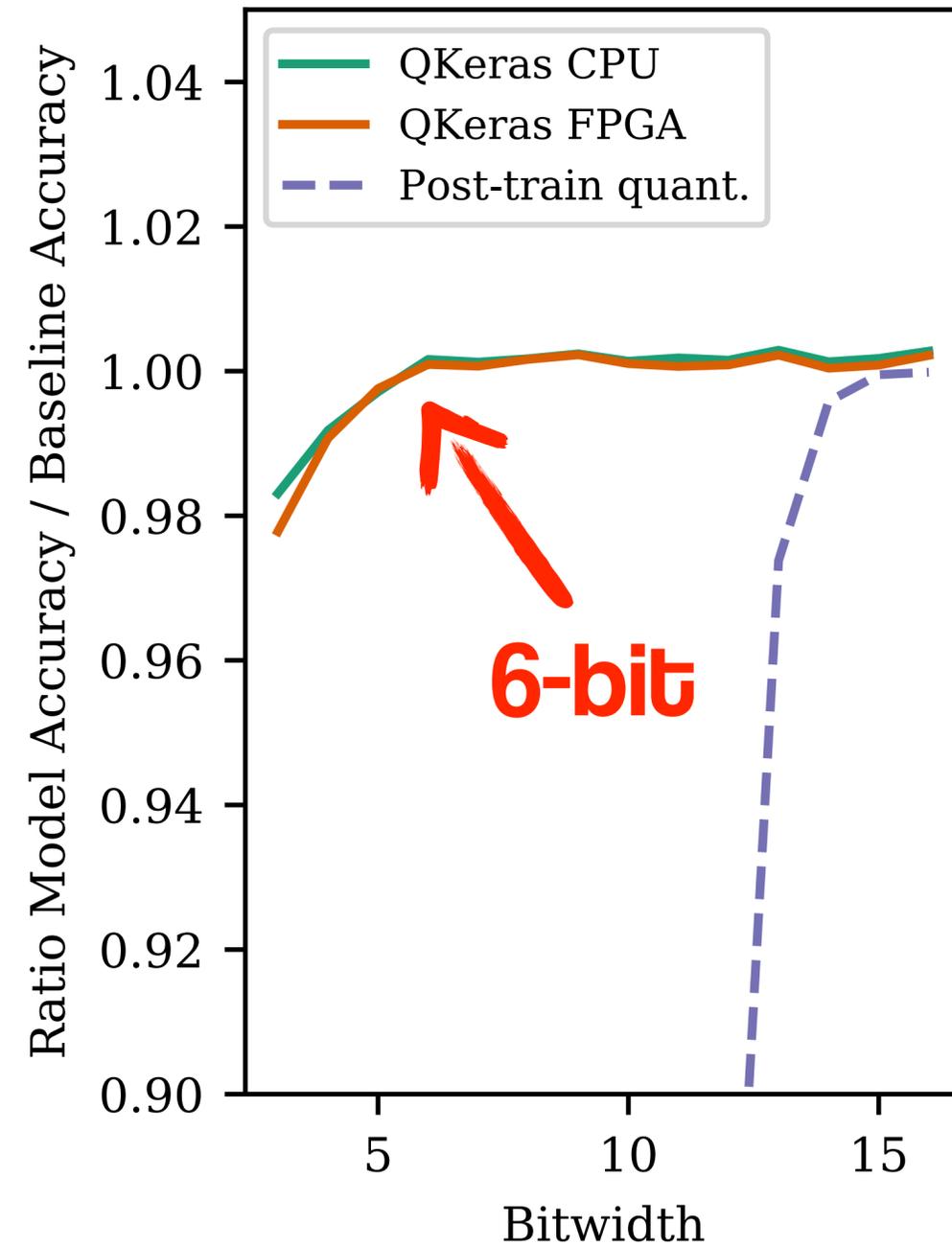
**Intrinsic development loop  
between algorithm design,  
training, and  
implementation**

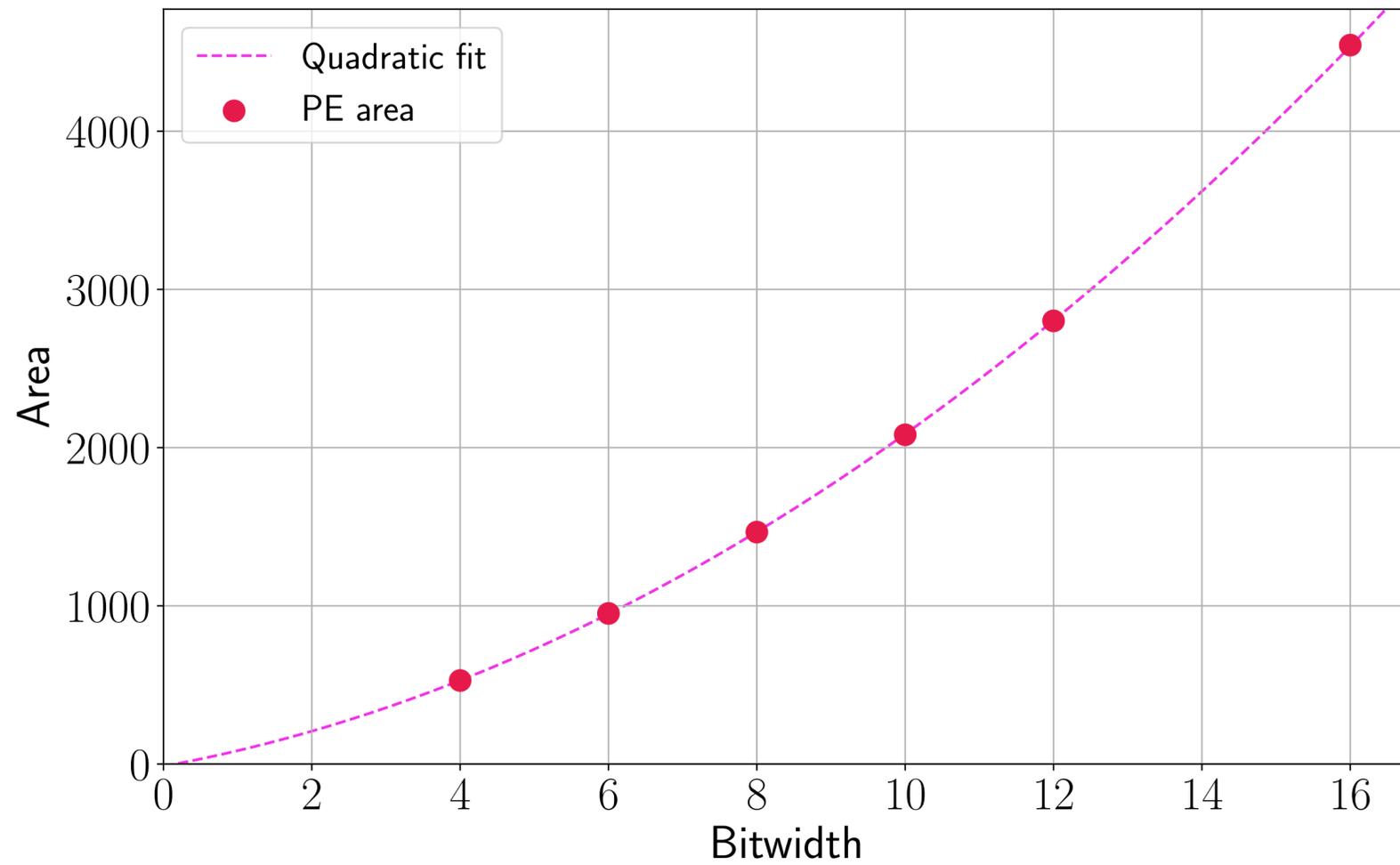
**Acceleration of design  
space exploration is  
important**



Post-training quantization

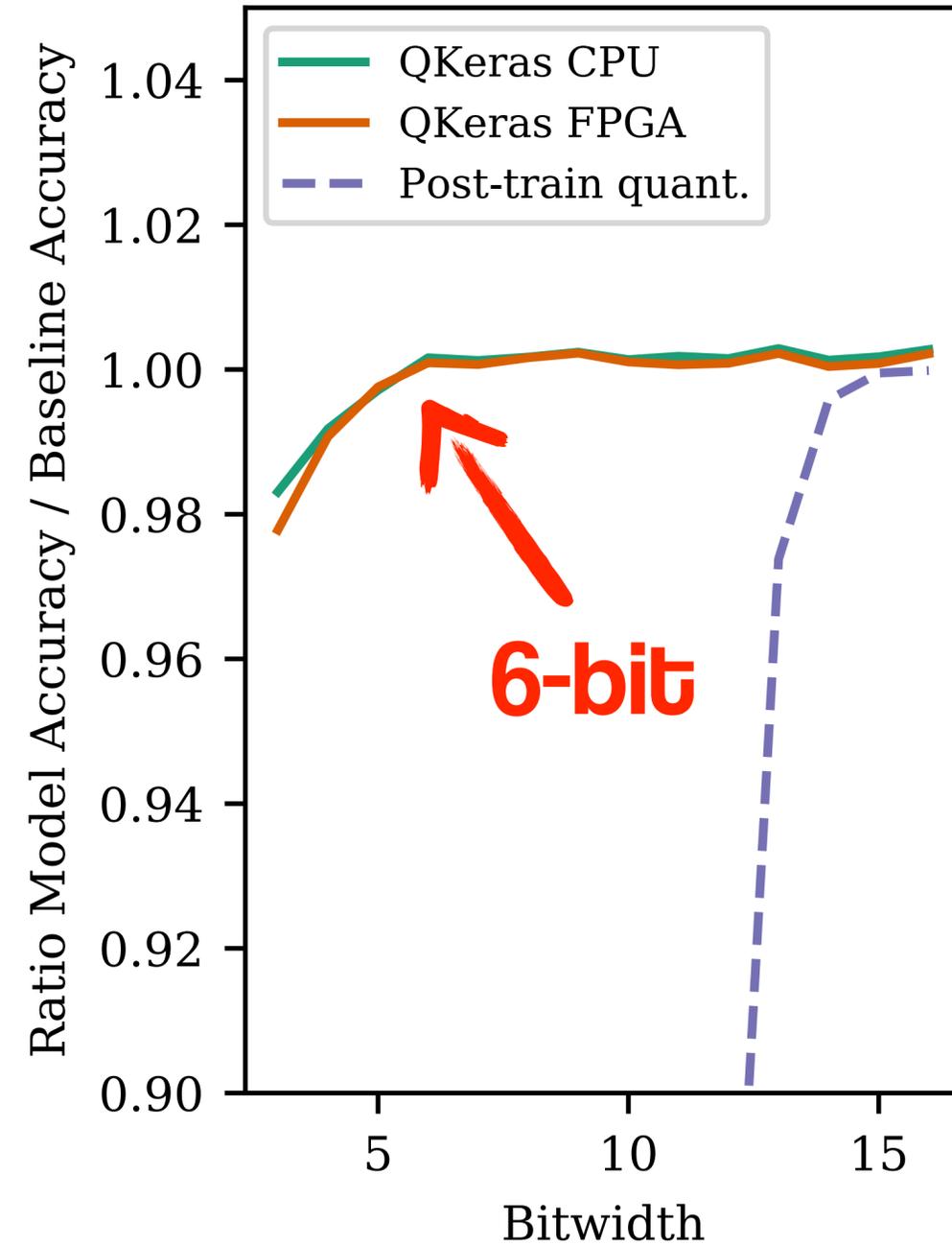
## Quantization-aware training





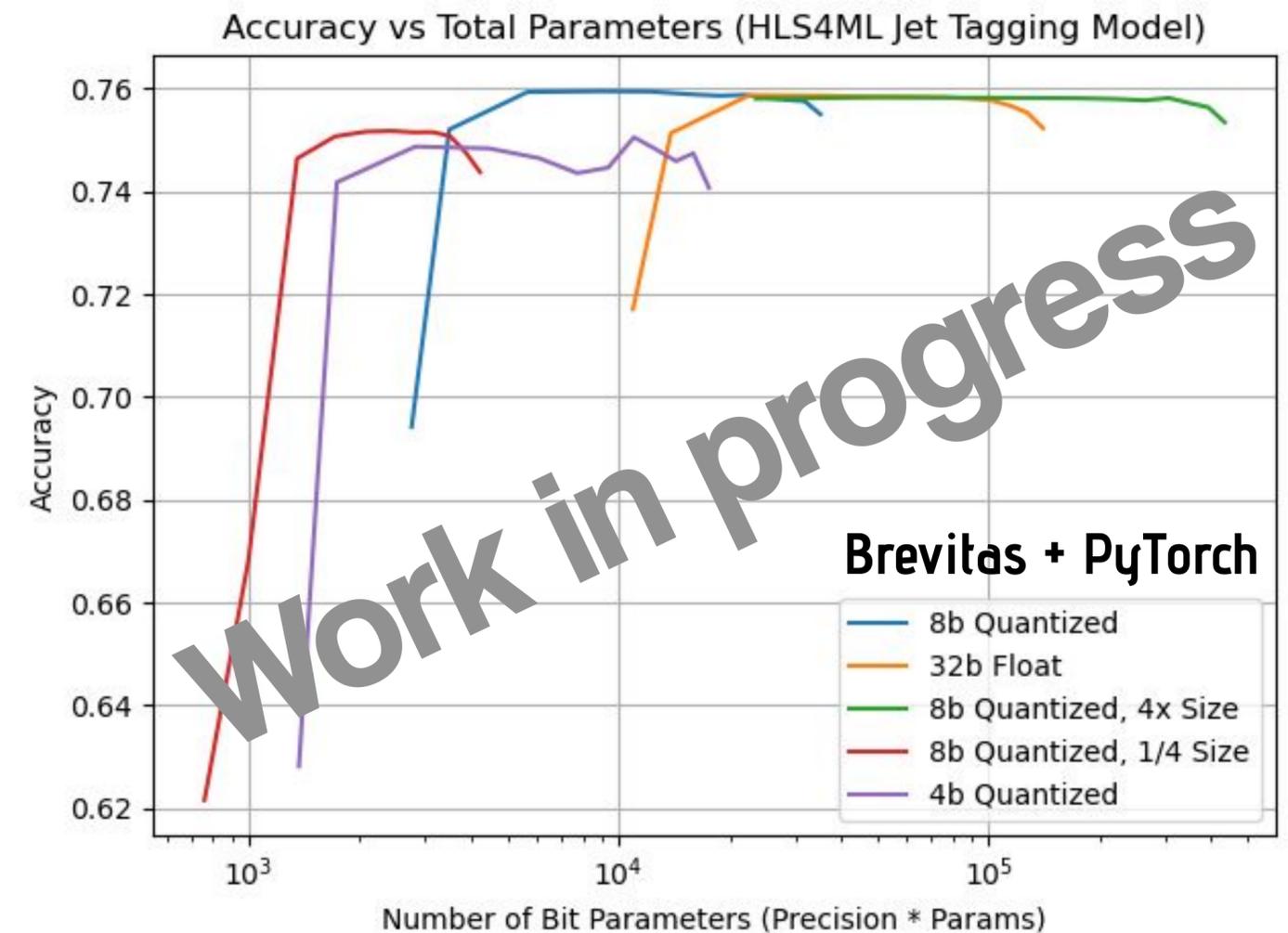
Area & power go quadratically with bit width

## Quantization-aware training



# E.g. deeper ML questions (a personal view)

- What about combinations of techniques?  
Unstructured (synaptic) pruning + QAT
  - How does that compare to hyper-parameter optimization and Bayesian Optimization?
- What did the network learn? How comparable are they? Active research on “neural efficiency” of networks and the “distance” between neural networks
- How does the learning (training) process affect efficiency? 2nd order optimization vs. stochastic gradient descent?



# Domain-inspired ML

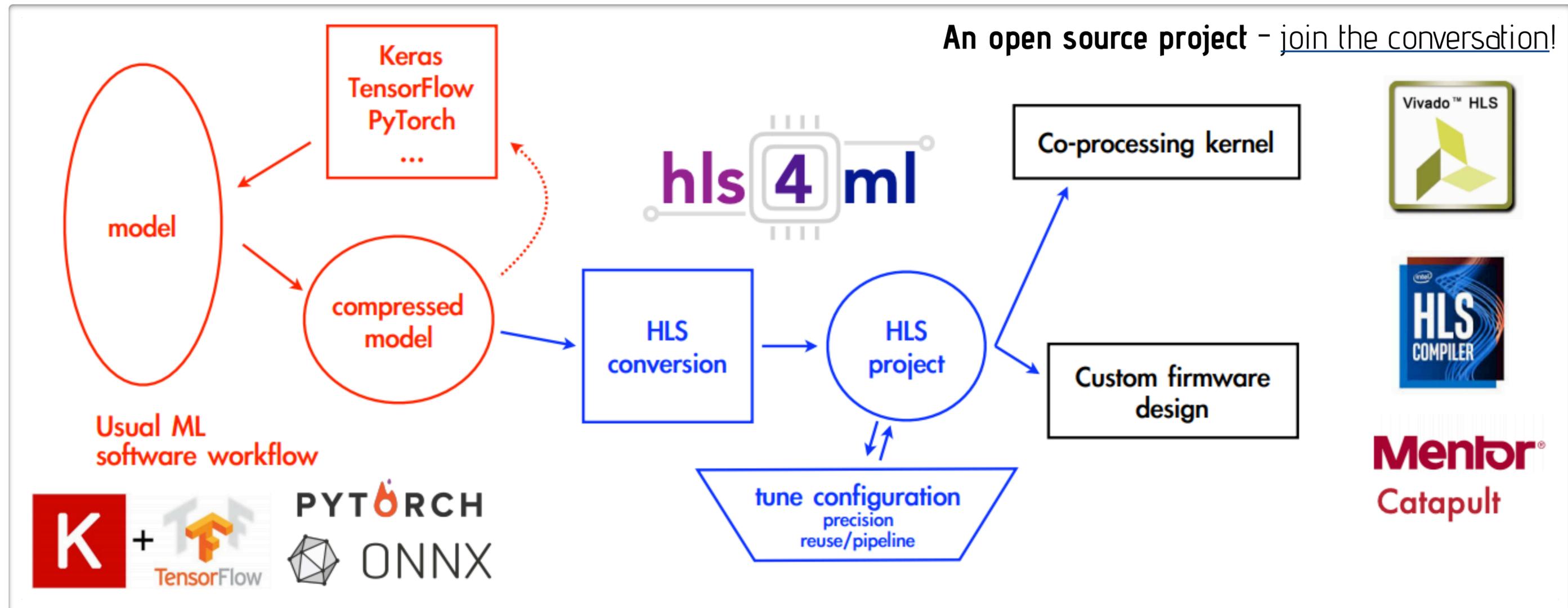
Architecture	Accuracy	AUC	$1/\epsilon_B$	#Param
ParticleNet	0.938	0.985	$1298 \pm 46$	498k
P-CNN	0.930	0.980	$732 \pm 24$	348k
ResNeXt	0.936	0.984	$1122 \pm 47$	1.46M
EFP	0.932	0.980	384	1k
EFN	0.927	0.979	$633 \pm 31$	82k
PFN	0.932	0.982	$891 \pm 18$	82k
TopoDNN	0.916	0.972	$295 \pm 5$	59k
LGN	$0.929 \pm .001$	$0.964 \pm 0.018$	$435 \pm 95$	4.5k

Potentially HUGE gains in efficiency from more intelligent network design incorporating domain knowledge

Energy flow polynomial neural network

Lorentz group equivariant neural network

# hls4ml: accelerating ML on hardware



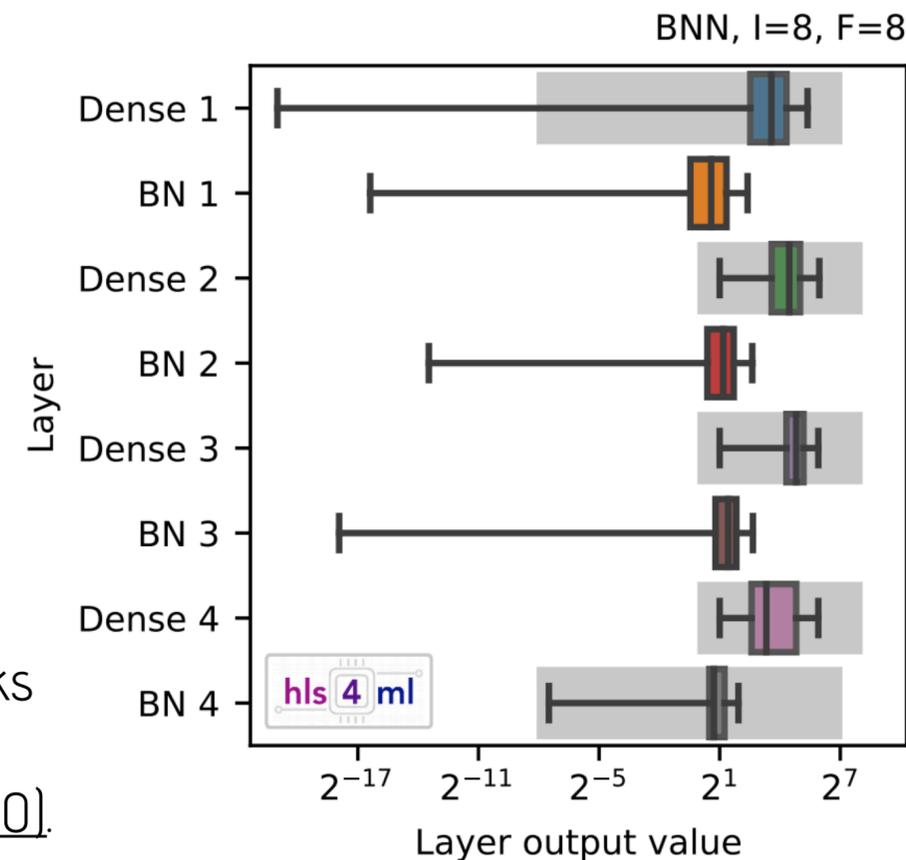
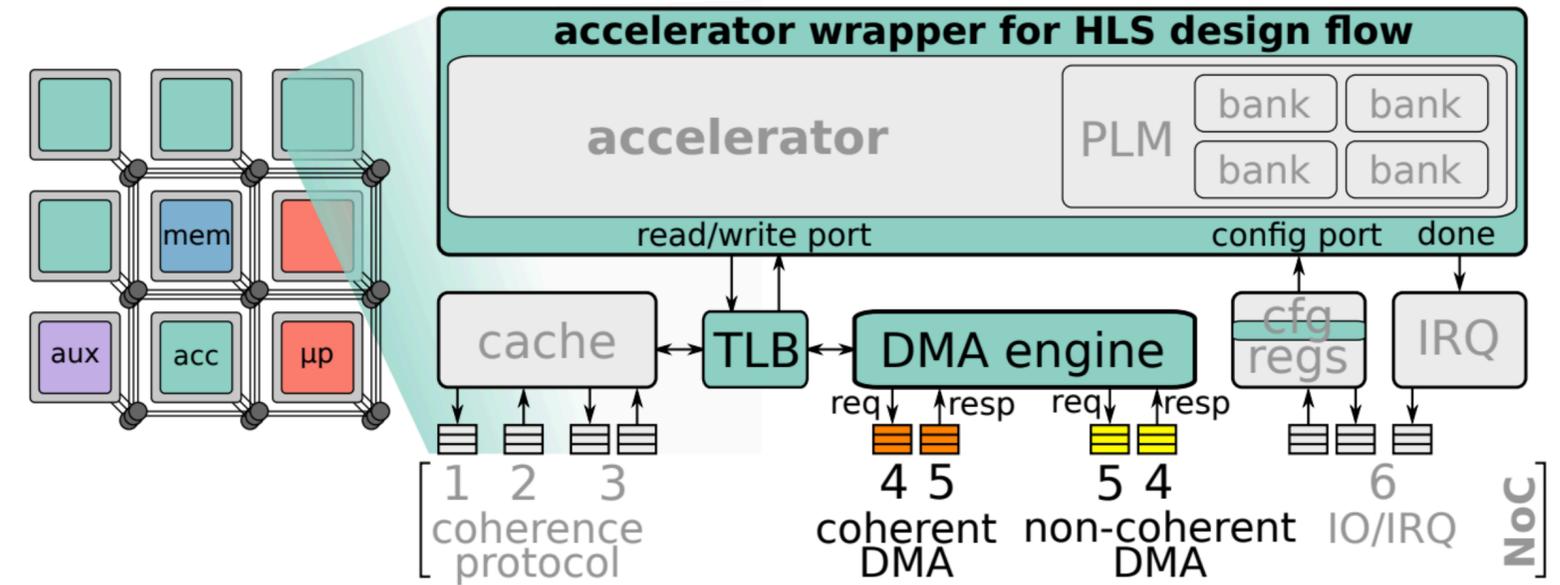
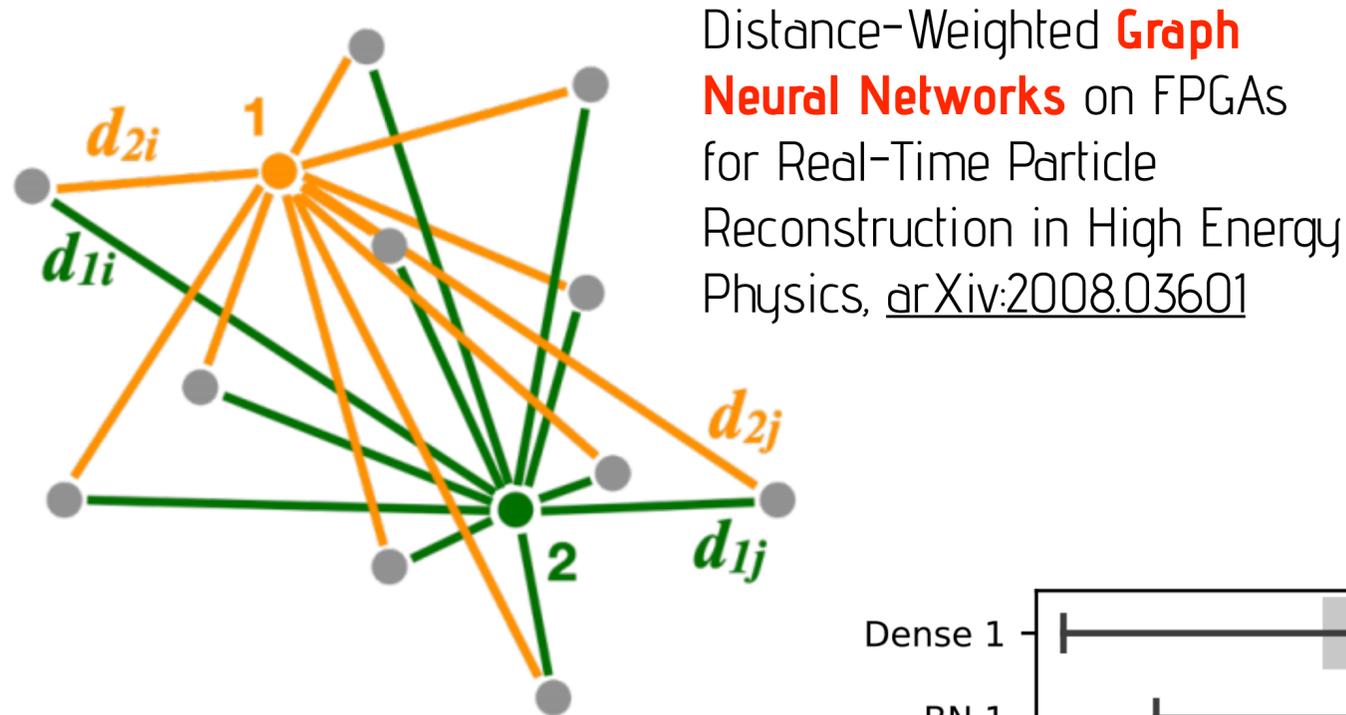
Originally designed for LHC triggers applications but broad and growing user base  
Example: ML training (qkeras/brevitas), backend devices (Xilinx, Intel, Microsemi, ASIC),  
network architectures (Conv, RNN/LSTM, Graph), and applications (...)

**Join us for the tutorial on Thursday!**

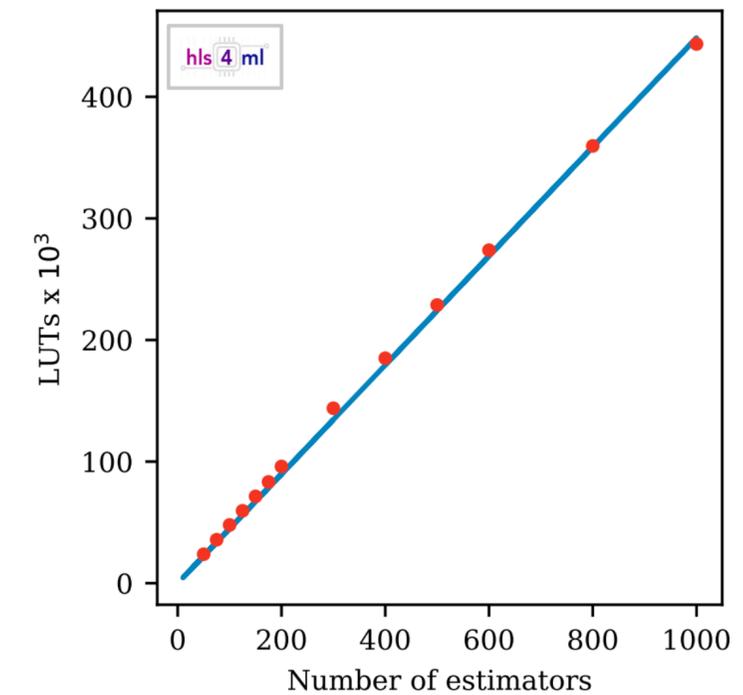
Originally designed for LHC triggers applications but broad and growing user base  
Example: ML training (qkeras/brevitas), backend devices (Xilinx, Intel, Microsemi, ASIC),  
network architectures (Conv, RNN/LSTM, Graph), and applications (...)

# More hls4ml

**ESP4ML**: Platform-Based Design of Systems-on-Chip for Embedded Machine Learning, DATE Conference 2020.



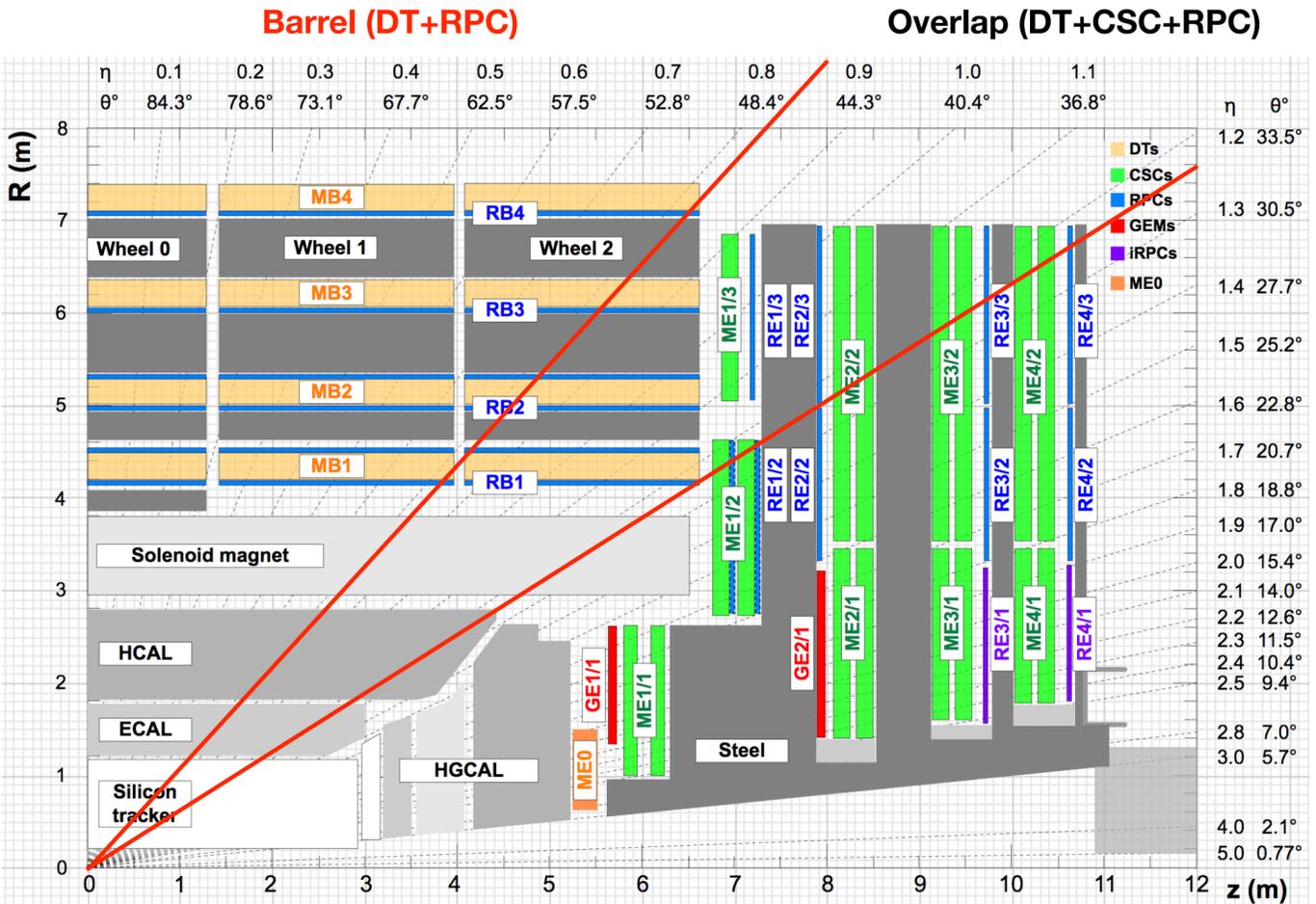
Compressing deep neural networks on FPGAs to **binary and ternary** precision with hls4ml, MLST (2020).



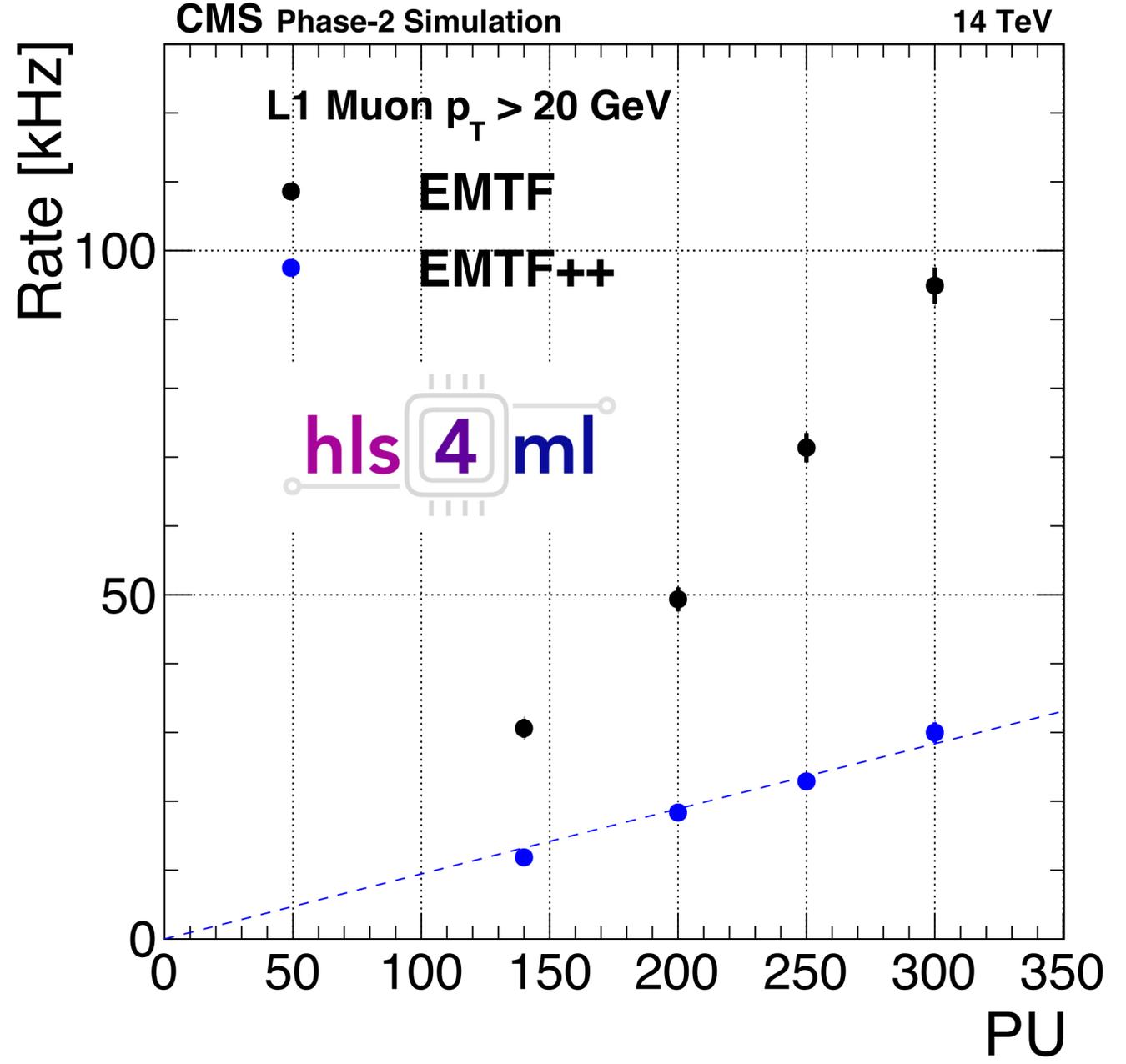
Fast inference of **Boosted Decision Trees** in FPGAs for particle physics, JINST 15, P05026 (2020).

# Example applications

# Case study: muon trigger upgrade



Endcap (GEM+CSC+RPC)



EMTF = BDT (external memory)

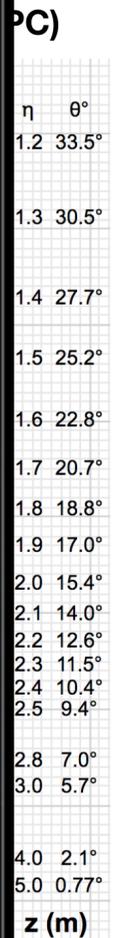
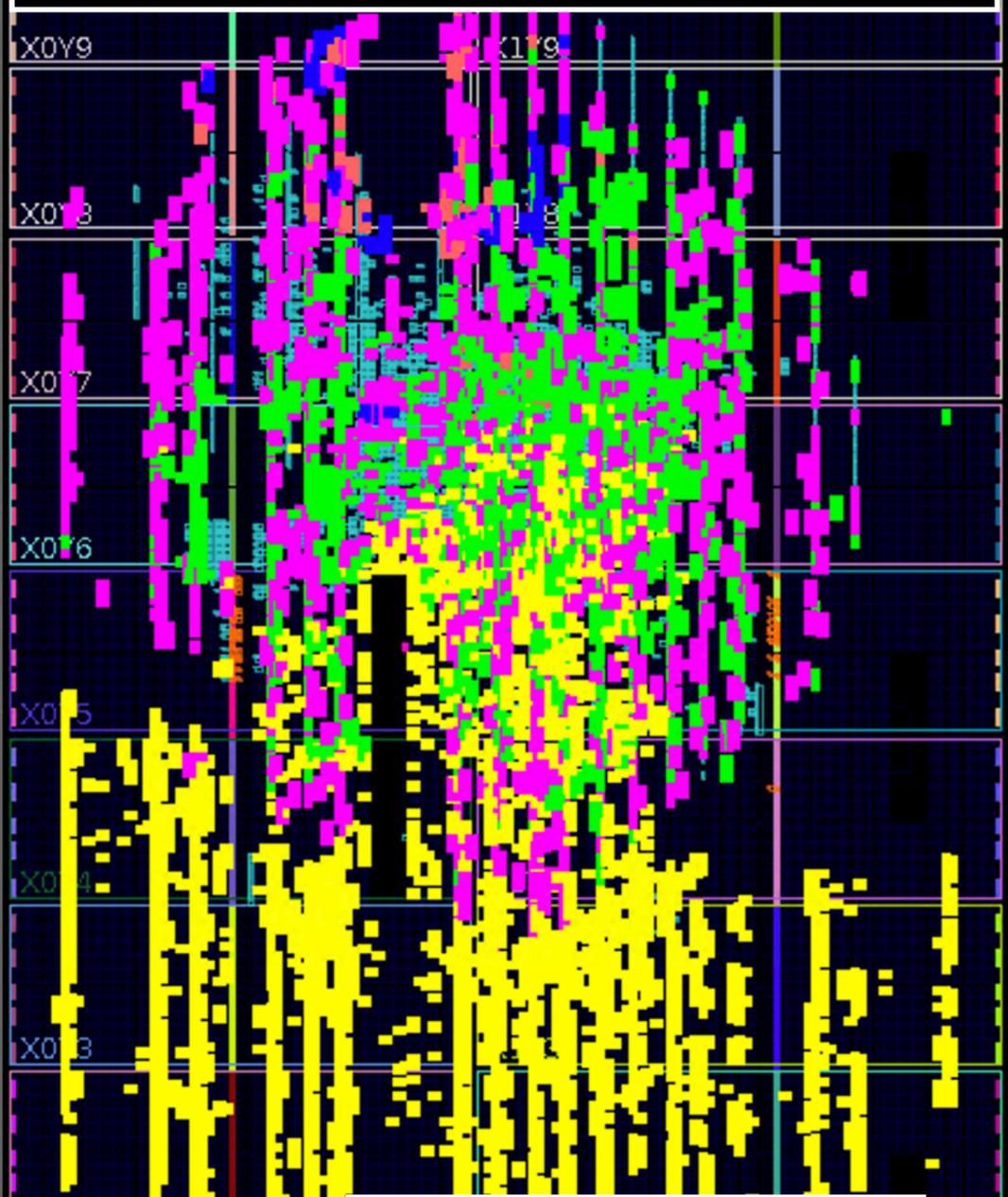
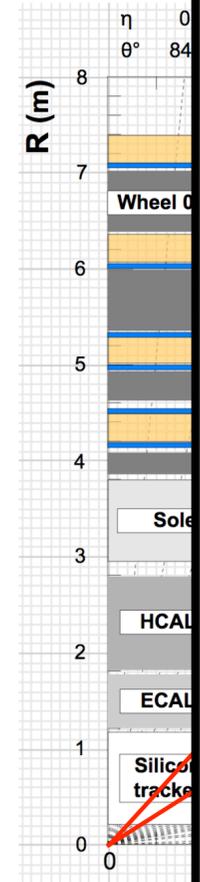
EMTF++ = NN

~3x reduction in the trigger rate for neural network!

# Case

# trigger upgrade

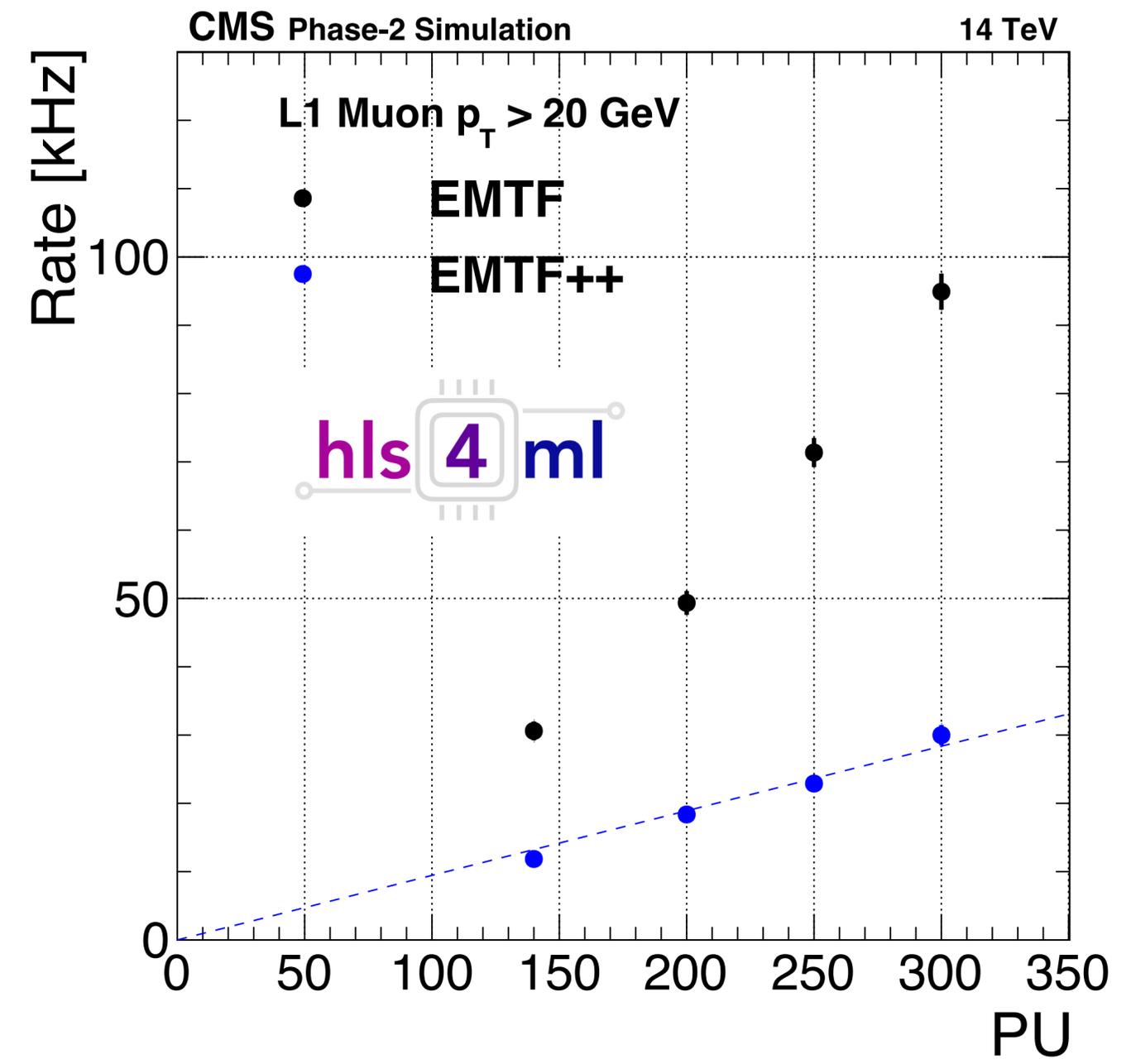
## AI circuit for ultrafast inference on FPGA



Inference time: 280 ns  
Throughput: 104 Gb/s

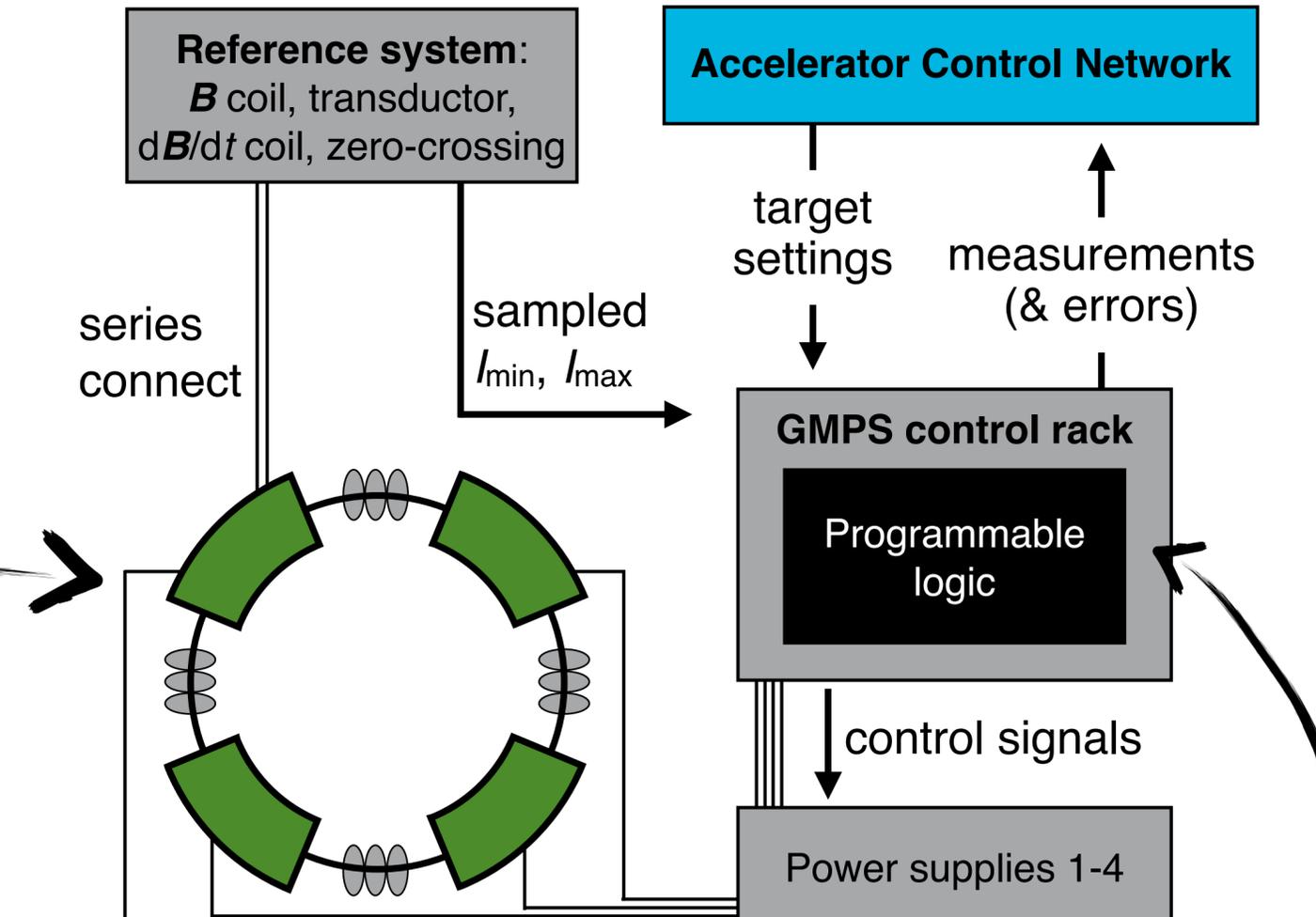
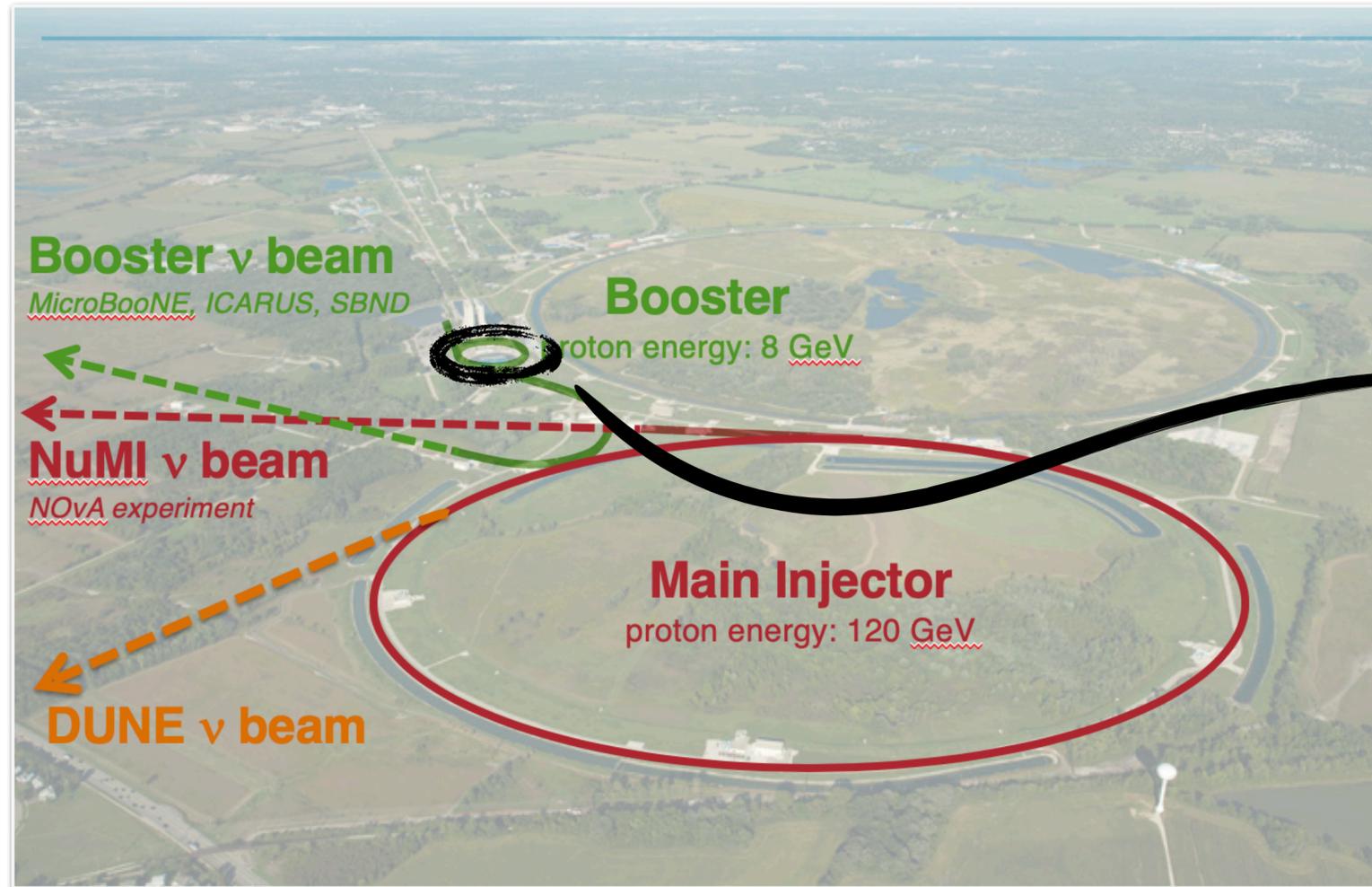
Dense Network  
23 → 30 → 25 → 20  
→ momentum & classifier

EMTF  
EMTF  
~3x



neural network!

# Accelerator controls



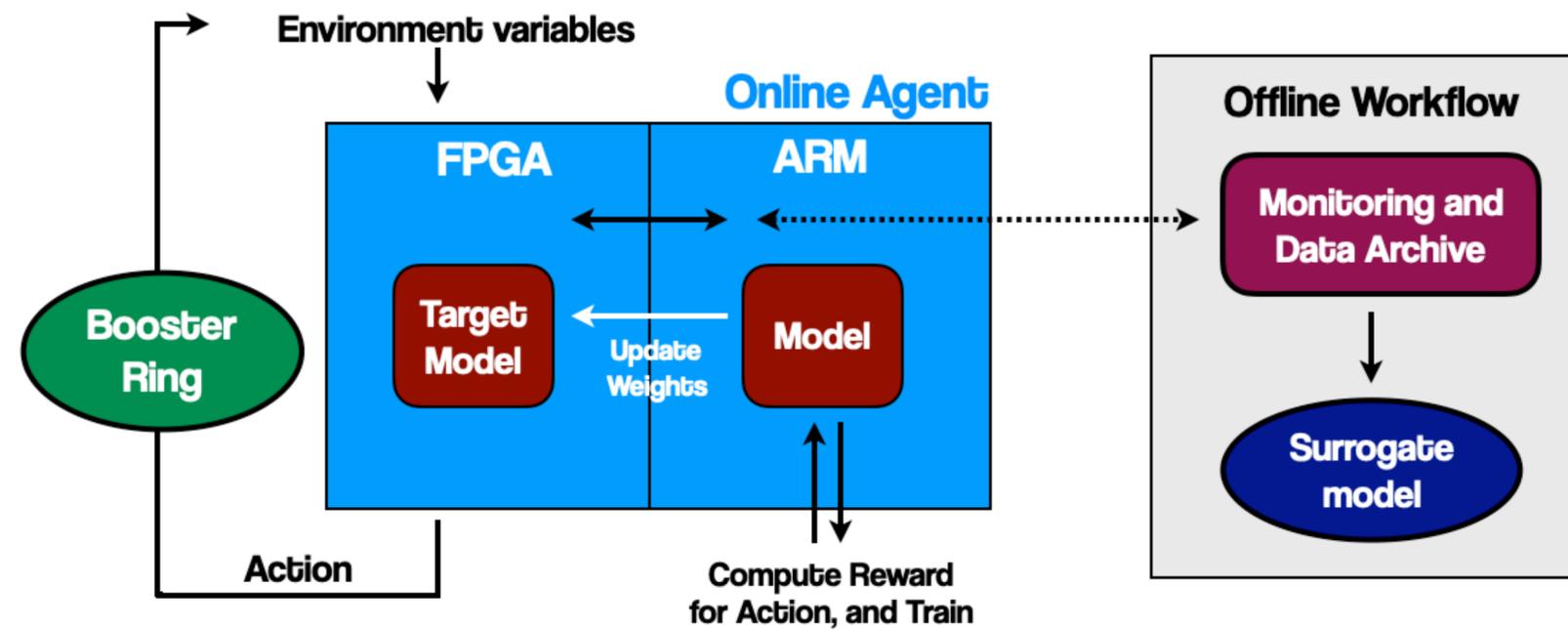
## The Task:

Control Gradient Magnet Power Supplies for Booster Ring at 15 Hz to reduce beam losses using reinforcement learning setup

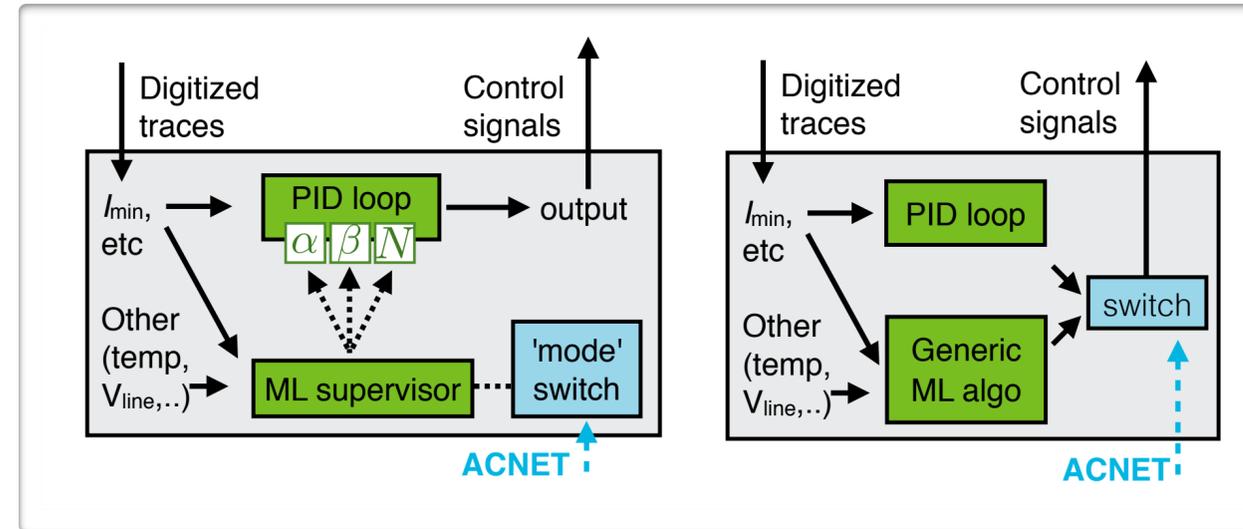


Arria10 SoM

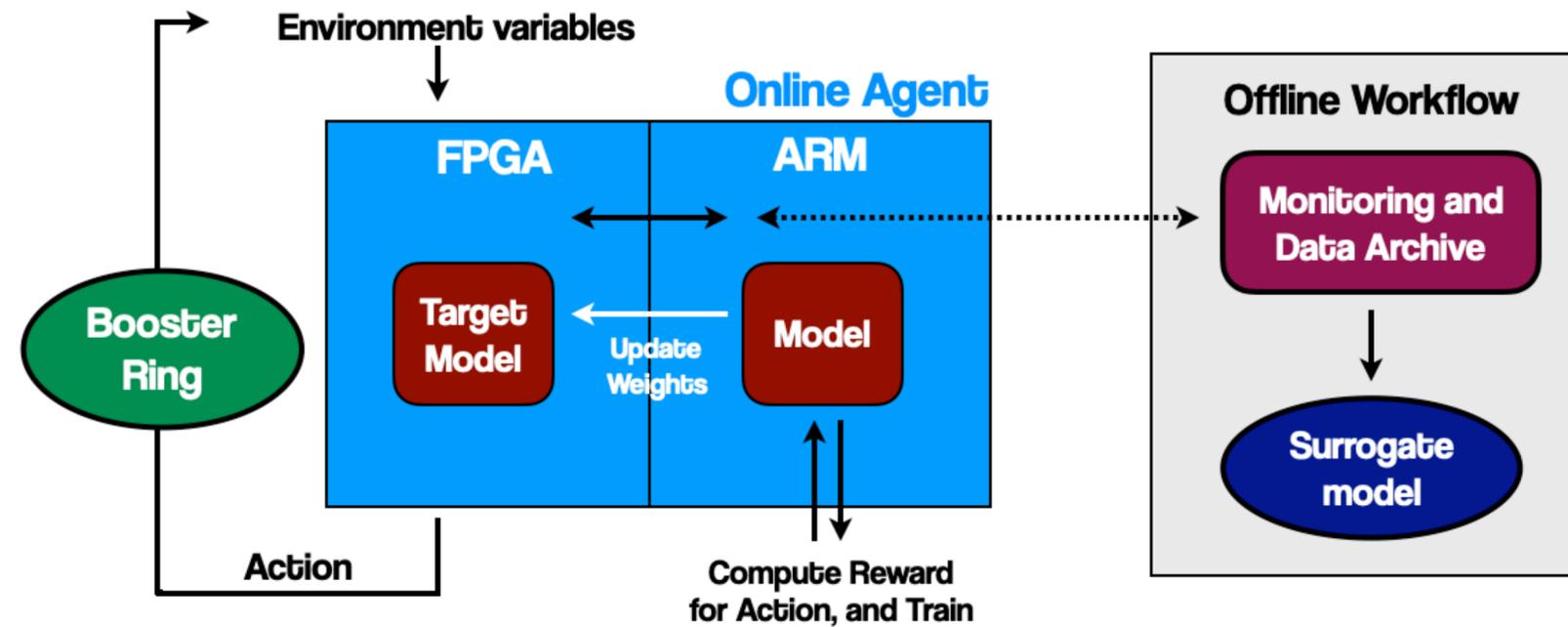
# RL setup



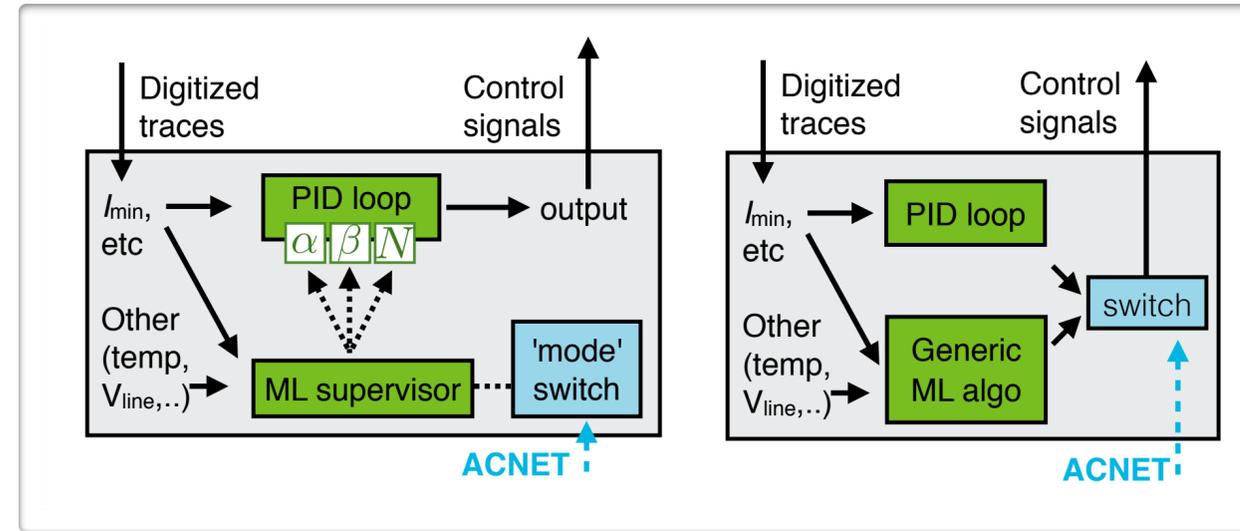
# RL setup



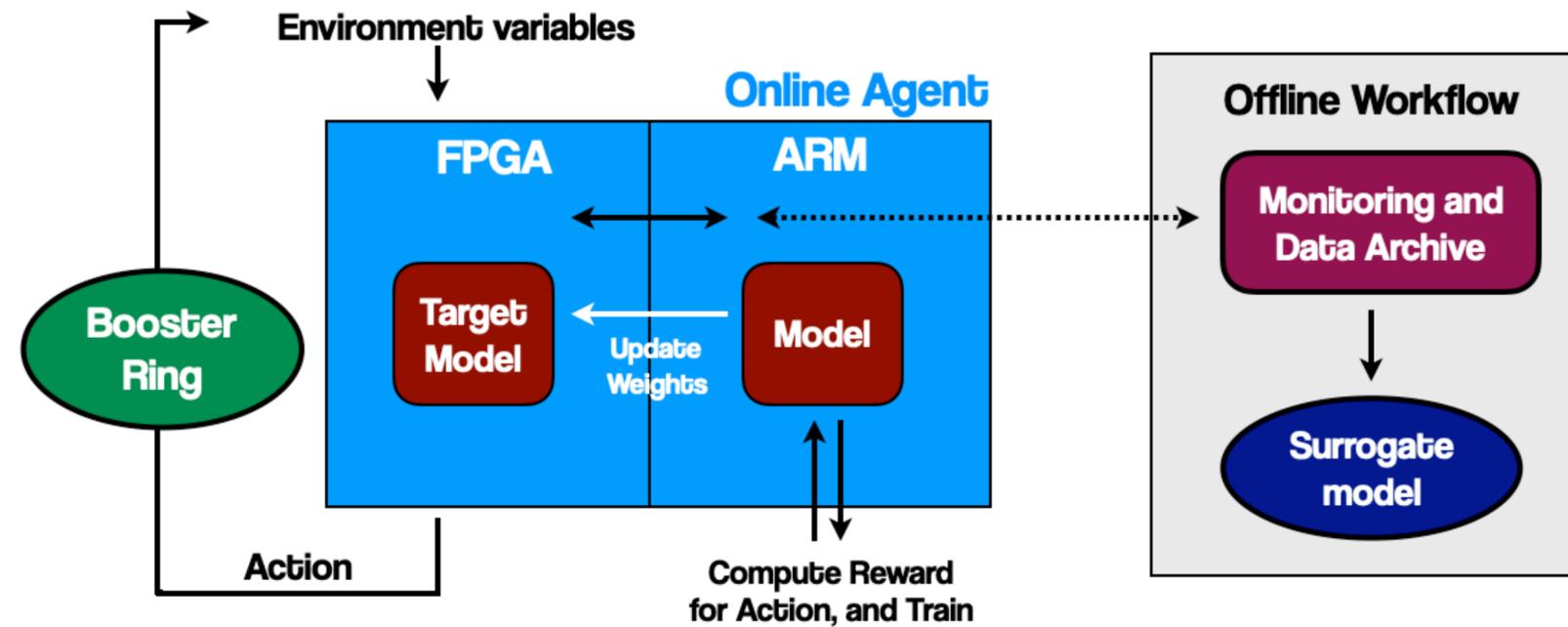
For safe-guarding online algorithm, start with model-based RL; plan to compare against model-free online agent



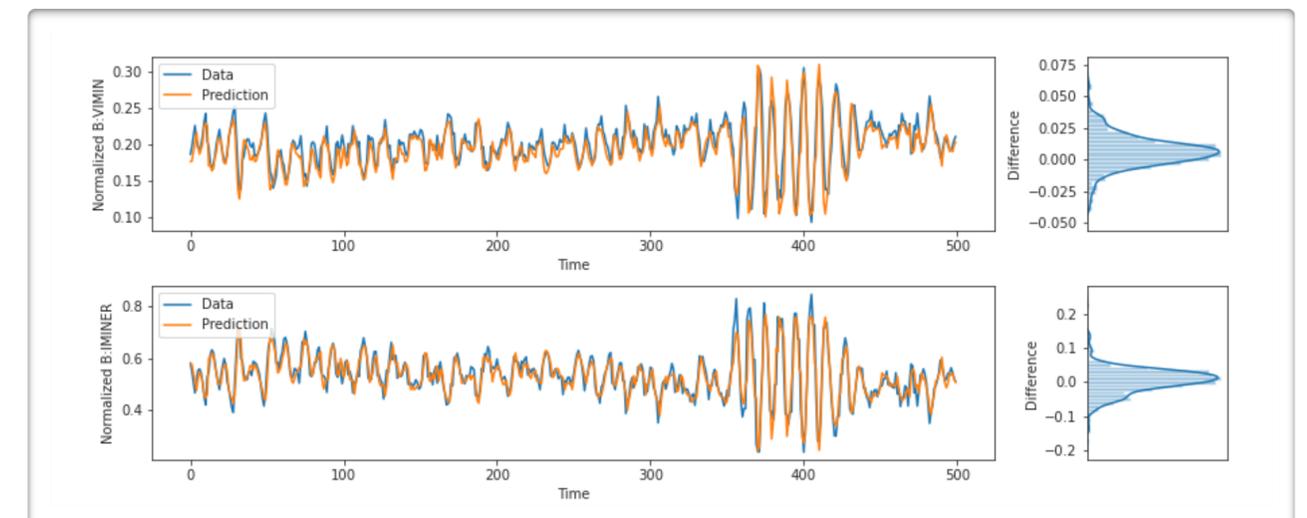
# RL setup



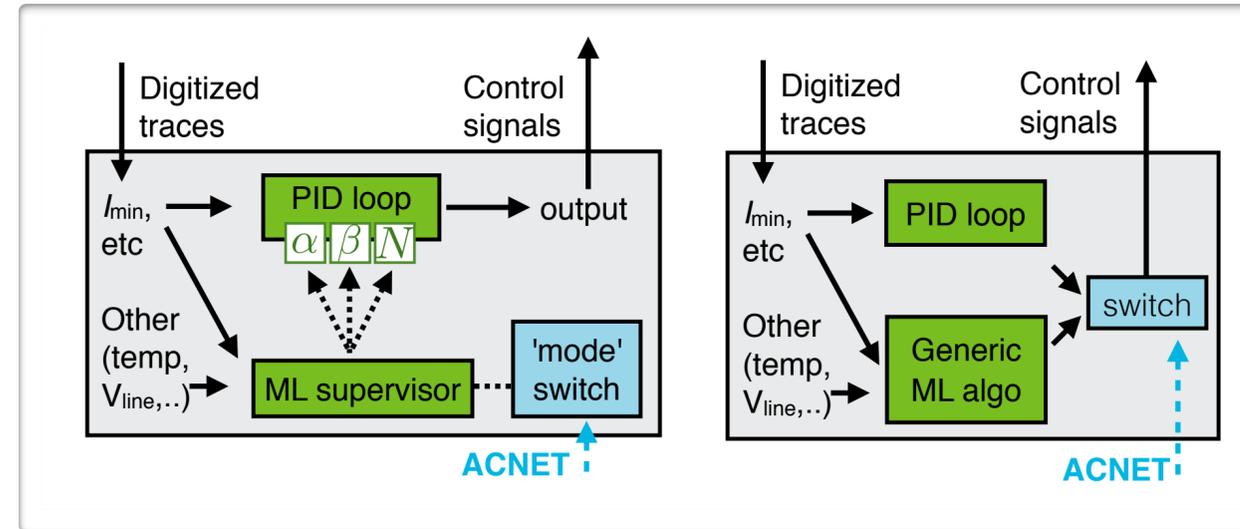
For safe-guarding online algorithm, start with model-based RL; plan to compare against model-free online agent



LSTM-based network using Deep Q-learning framework for surrogate model to mimic behavior of Booster

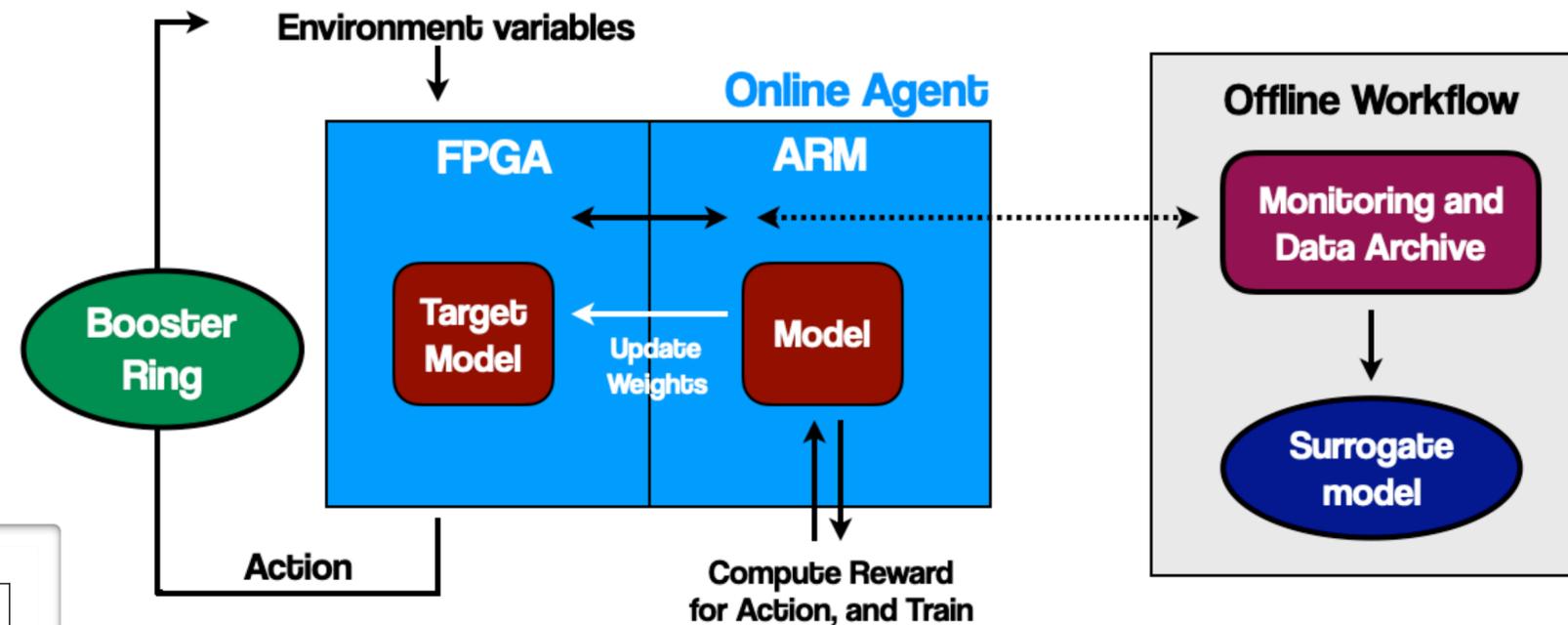
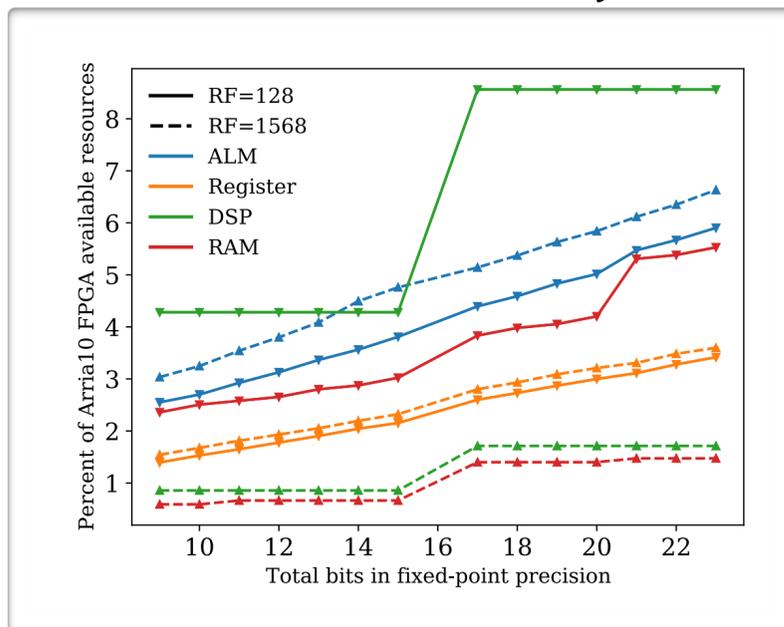


# RL setup



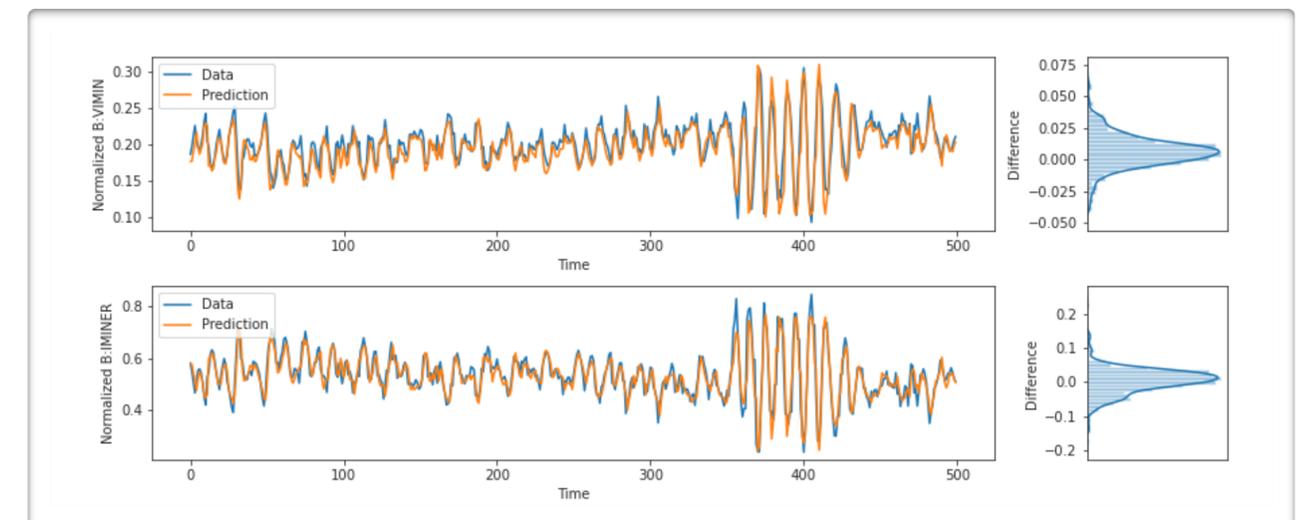
For safe-guarding online algorithm, start with model-based RL; plan to compare against model-free online agent

DNN-based ensemble model for RL online agent



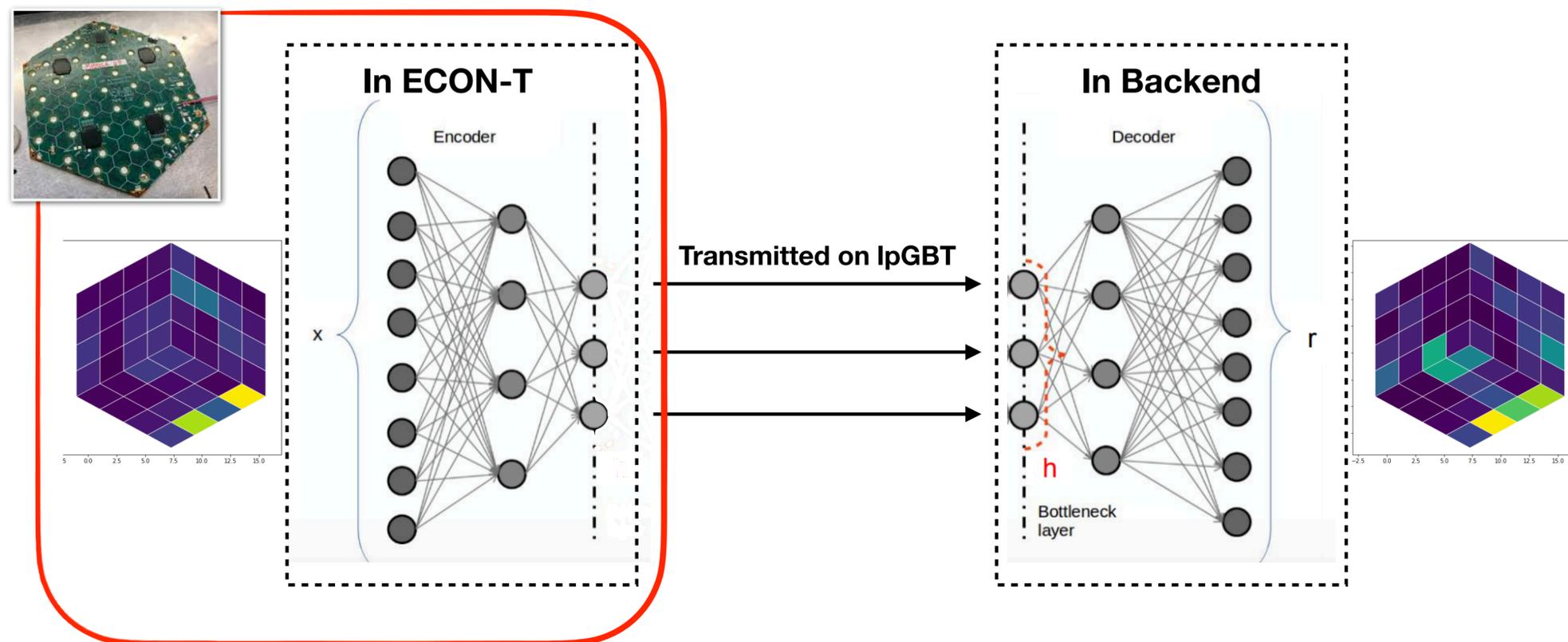
LSTM-based network using Deep Q-learning framework for surrogate model to mimic behavior of Booster

Requires development of PL-PS interface for weight updating



# What about ASICs?

- Putting a neural network on the detector front-ends for data compression
  - Application: **CMS HGCaI ECON (concentrator) ASIC**
  - ASIC required due to radiation tolerance and power requirements
- **Fully reconfigurable ASIC** to address: future unknown unknowns' including **evolving LHC conditions** (pileup, beam bkgs), **detector performance** (noise, dead channels), **performance metrics** (resolution, substructure, new physics signatures)



## Auto-encoder design: Conv2D + Dense

Performance metrics defined based on Energy Mover's Distance and physics resolutions/rates

Quantization aware training very important!

# ASIC workflow

See poster here by Manuel Blanco Valentin  
Talks at IEEE NSS (Christian Herwig, Farah Fahim)

ALGORITHM  
DEVELOPMENT



ML Model  
Training



- **hls4ml** simplifies the design of ML accelerators
  - | hls4ml directives | << | HLS directives |
  - C++ library of ML functionalities optimized for HLS

Part: ...  
ReuseFactor: ...  
Precision: ...  
IOType: ...  
Backend: ...

**hls4ml**  
Directives



**hls4ml**



```
inline bar 'ON'  
unroll 11 factor=10  
bundle b1=A,B, b2=C
```

HLS  
Directives

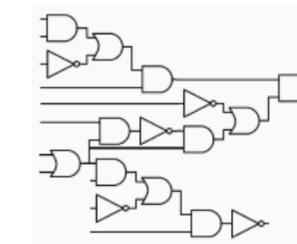
```
void foo(int A[10],  
         int B[10], int C[10])  
{  
    int i = 0;  
    11: for (; i < 10; i++) {  
        A[i] = B[i] * i;  
    }  
    i = 0;  
    12: for (; i < 10; i++) {  
        B[i] = A[i] * B[i];  
        C[i] = B[i] / 10;  
    }  
    bar(A, B);  
}
```

C++  
Specification

HLS

Technology Library

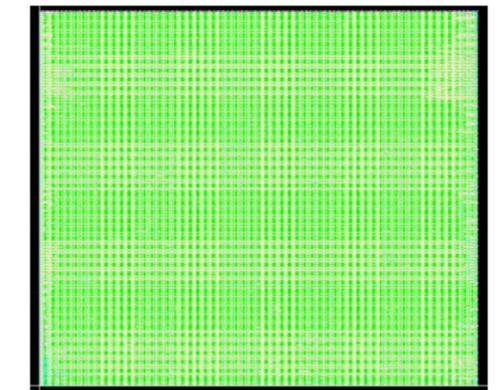
**TMR4sv\_hls**



RTL  
Hardware

Implementation(s)

ECON  
AUTOENCODER



GDSII

**First implementation of AI in an ASIC for particle physics**

TSMC 65nm process: ~95 mW, 3.6mm<sup>2</sup>, 25ns latency, 40M img/s, ~4000 ops, 200 MRad spec

# Summary

# Summary and outlook

Fast ML for Science Workshop  
<https://indico.cern.ch/e/fml2020>  
Join us! Nov 30-Dec 3

- Real-time intelligence can greatly accelerate science allowing us to test hypotheses faster, enhance performance of detectors/accelerators, and save potentially lost data
- I look forward to seeing exciting tools, techniques, and applications of real-time machine learning at this conference
- Early days towards smart sensors, autonomous detectors, and accelerated experimental loops