



Jets/MET with pileup and machine learning

Nhan Tran
Fermilab

Princeton HL-LHC Trigger Workshop
January 16, 2018

Last talk, Giovanni:

Particle Flow - efficient combination of detector information to extract best physics performance

Building of the technology presented by Giovanni...

This talk: more advanced algorithms

Dealing with pileup

PUPPI proof-of-concept: jets, MET, jet substructure (?),...

More sophistication with machine learning and **HLS4ML**

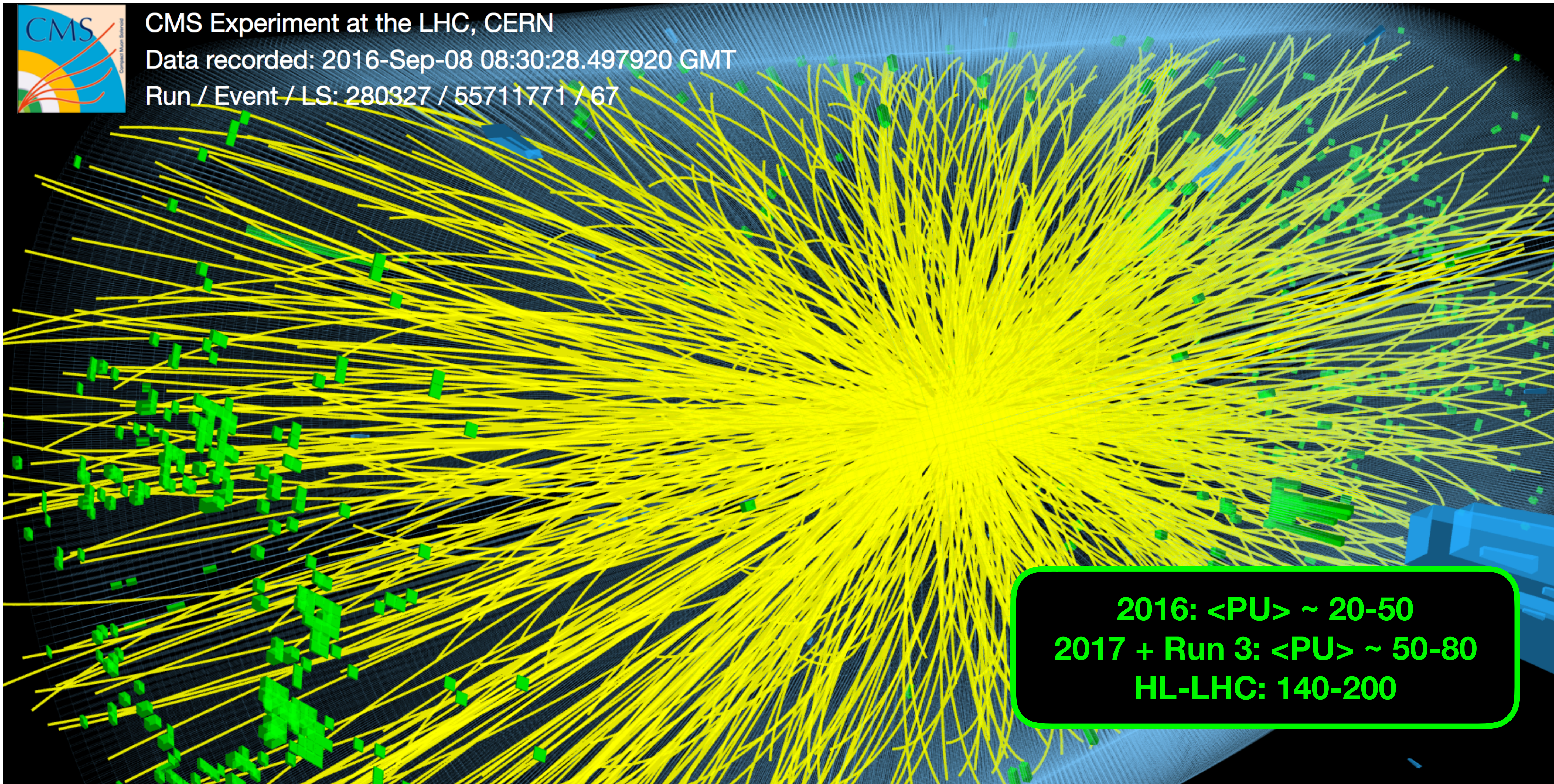
Multiple pp collisions in the same beam crossing
To increase data rate, squeeze beams as much as possible



CMS Experiment at the LHC, CERN

Data recorded: 2016-Sep-08 08:30:28.497920 GMT

Run / Event / LS: 280327 / 55711771 / 67



2016: $\langle \text{PU} \rangle \sim 20\text{-}50$
2017 + Run 3: $\langle \text{PU} \rangle \sim 50\text{-}80$
HL-LHC: 140-200

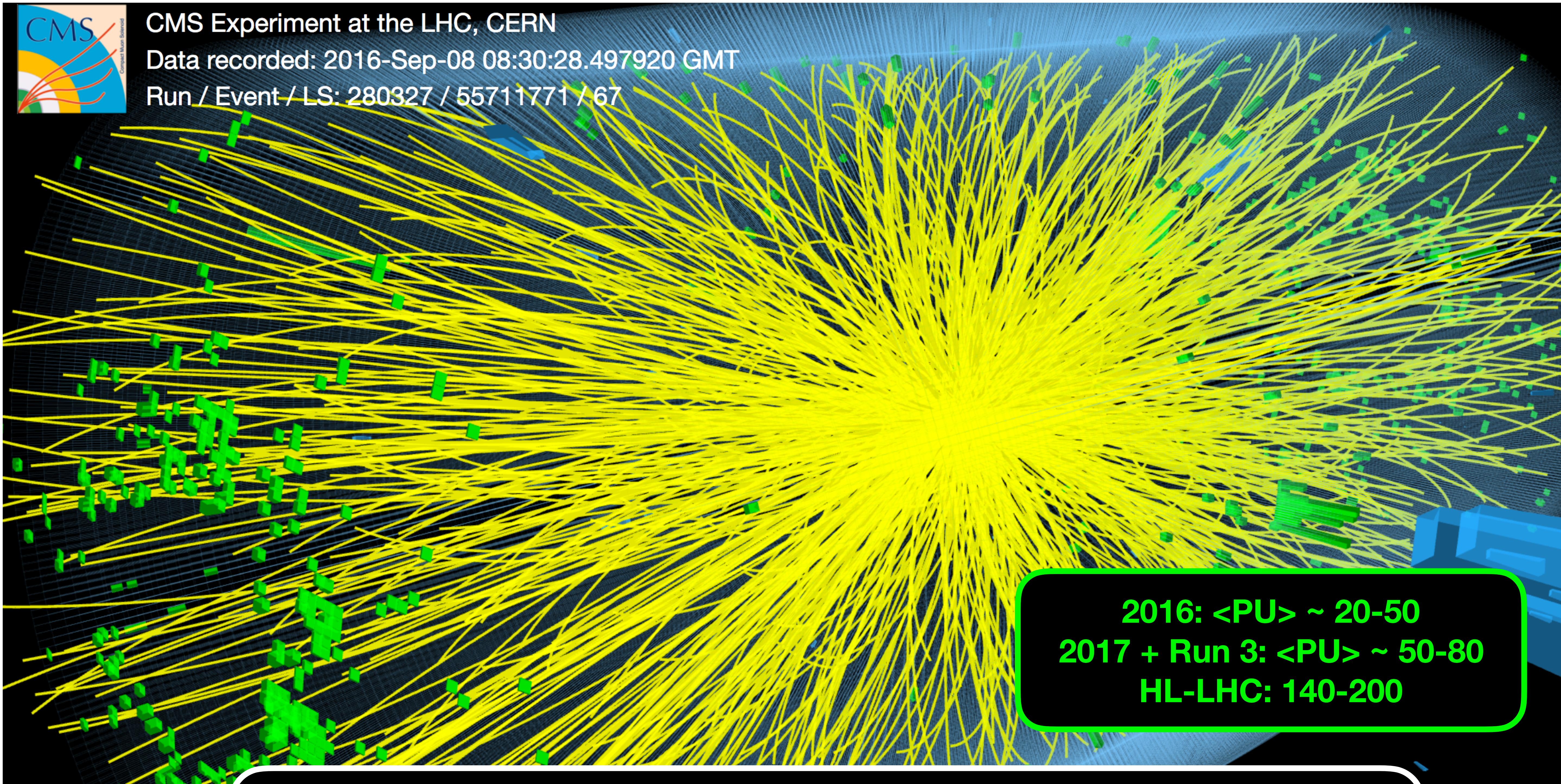
Multiple pp collisions in the same beam crossing
To increase data rate, squeeze beams as much as possible



CMS Experiment at the LHC, CERN

Data recorded: 2016-Sep-08 08:30:28.497920 GMT

Run / Event / LS: 280327 / 55711771 / 67



2016: $\langle \text{PU} \rangle \sim 20\text{-}50$
2017 + Run 3: $\langle \text{PU} \rangle \sim 50\text{-}80$
HL-LHC: 140-200

Need sophisticated techniques to preserve the physics!

PUPPI (PileUp Per Particle Id): based on PF paradigm

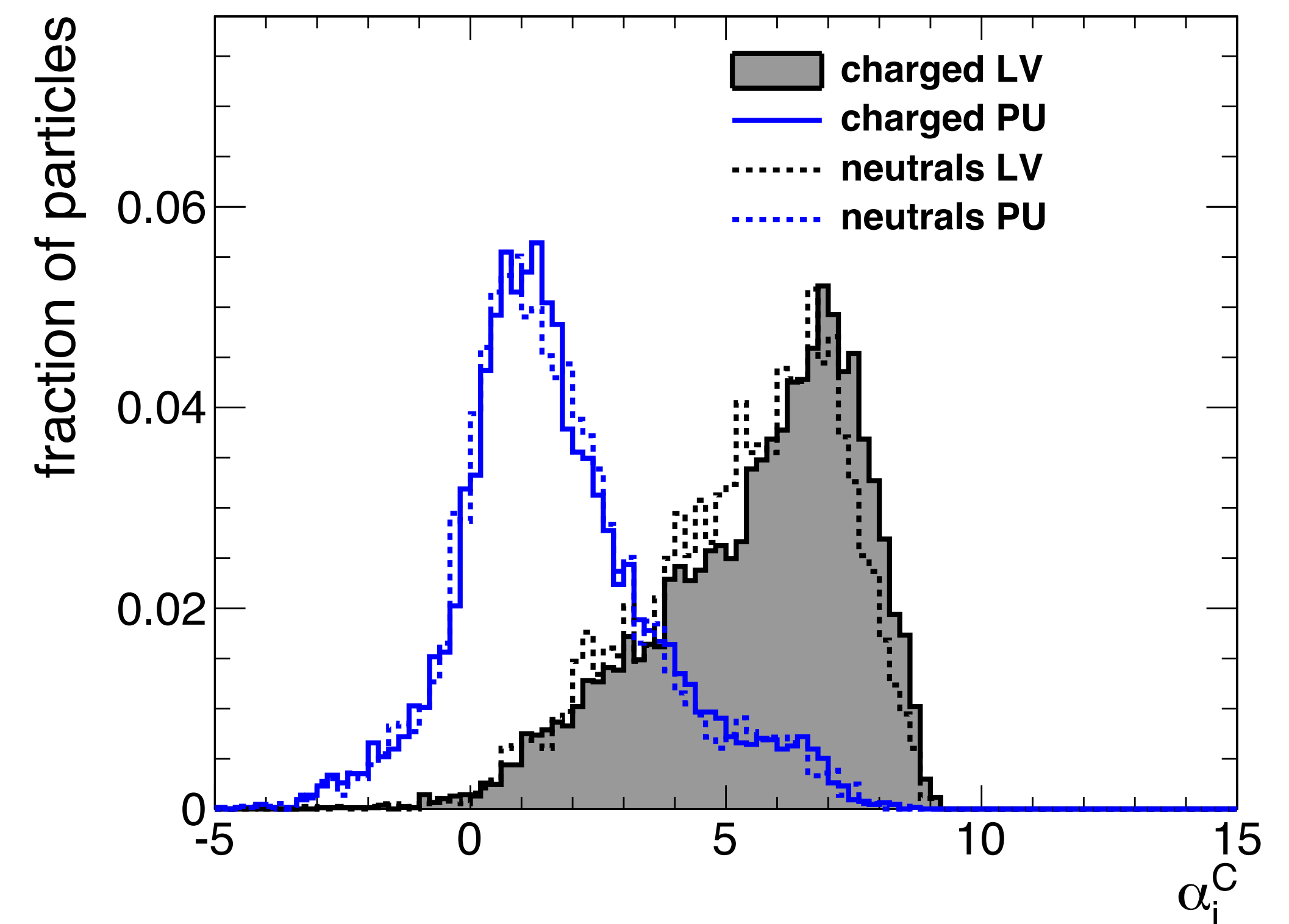
a general framework that determines, per particle, weight for **how likely** a particle is from PU

key insight: using QCD ansatz to infer neutral pileup contribution

[1] define a local discriminant, α , between pileup (PU) and leading vertex (LV)

$$\alpha_i^C = \log \left[\sum_{j \in \text{Ch, LV}} \frac{p_{T,j}}{\Delta R_{ij}} \Theta(R_0 - \Delta R_{ij}) \right]$$

[2] get data-driven α distribution for PU using charged PU tracks



PUPPI (PileUp Per Particle Id): based on PF paradigm

a general framework that determines, per particle, weight for **how likely** a particle is from PU

key insight: using QCD ansatz to infer neutral pileup contribution

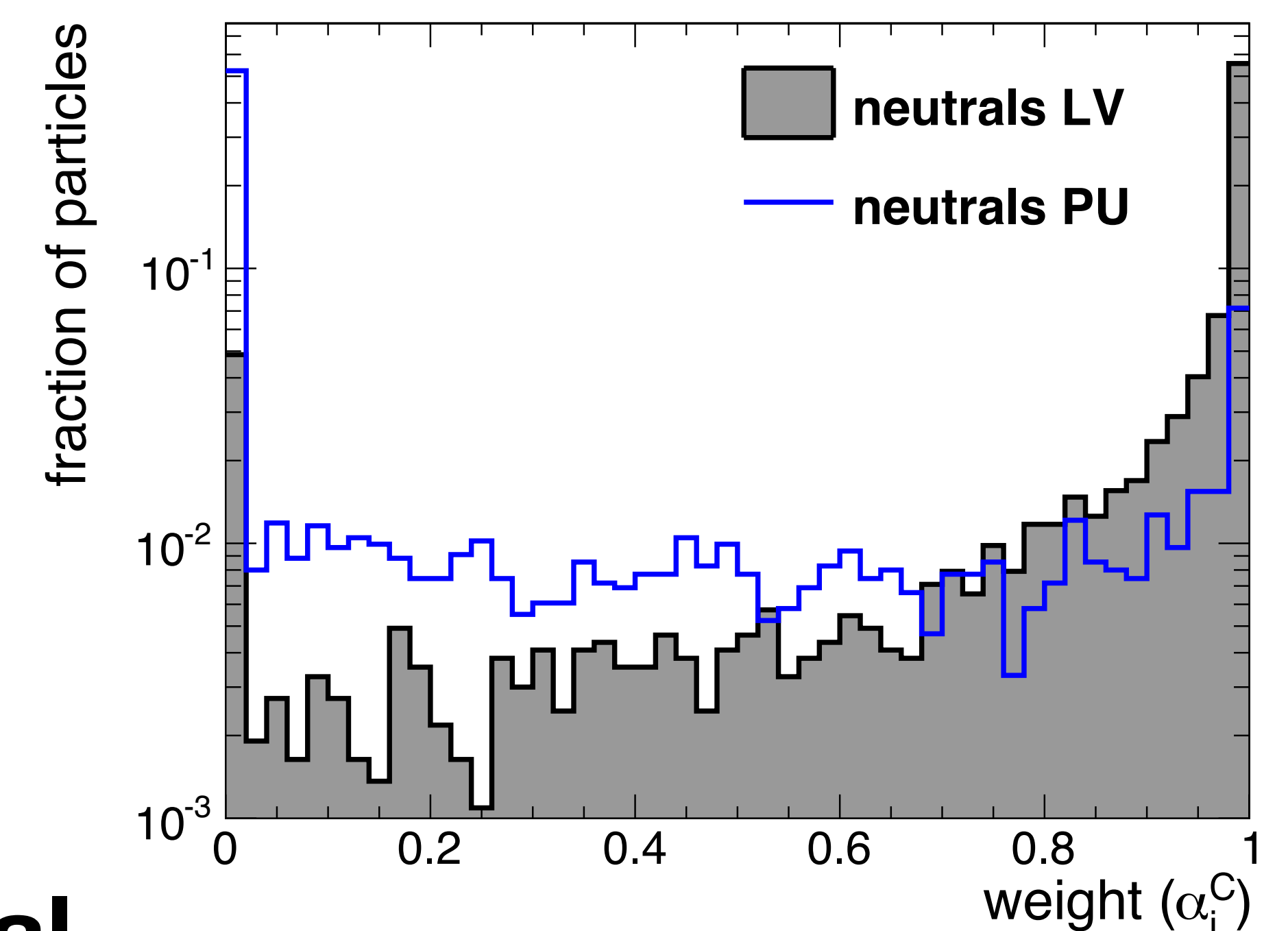
[1] define a local discriminant, α , between pileup (PU) and leading vertex (LV)

$$\alpha_i^C = \log \left[\sum_{j \in \text{Ch, LV}} \frac{p_{T,j}}{\Delta R_{ij}} \Theta(R_0 - \Delta R_{ij}) \right]$$

[2] get data-driven α distribution for PU using charged PU tracks

[3] for the neutrals, ask “how un-PU-like is α for this particle?”, compute a weight

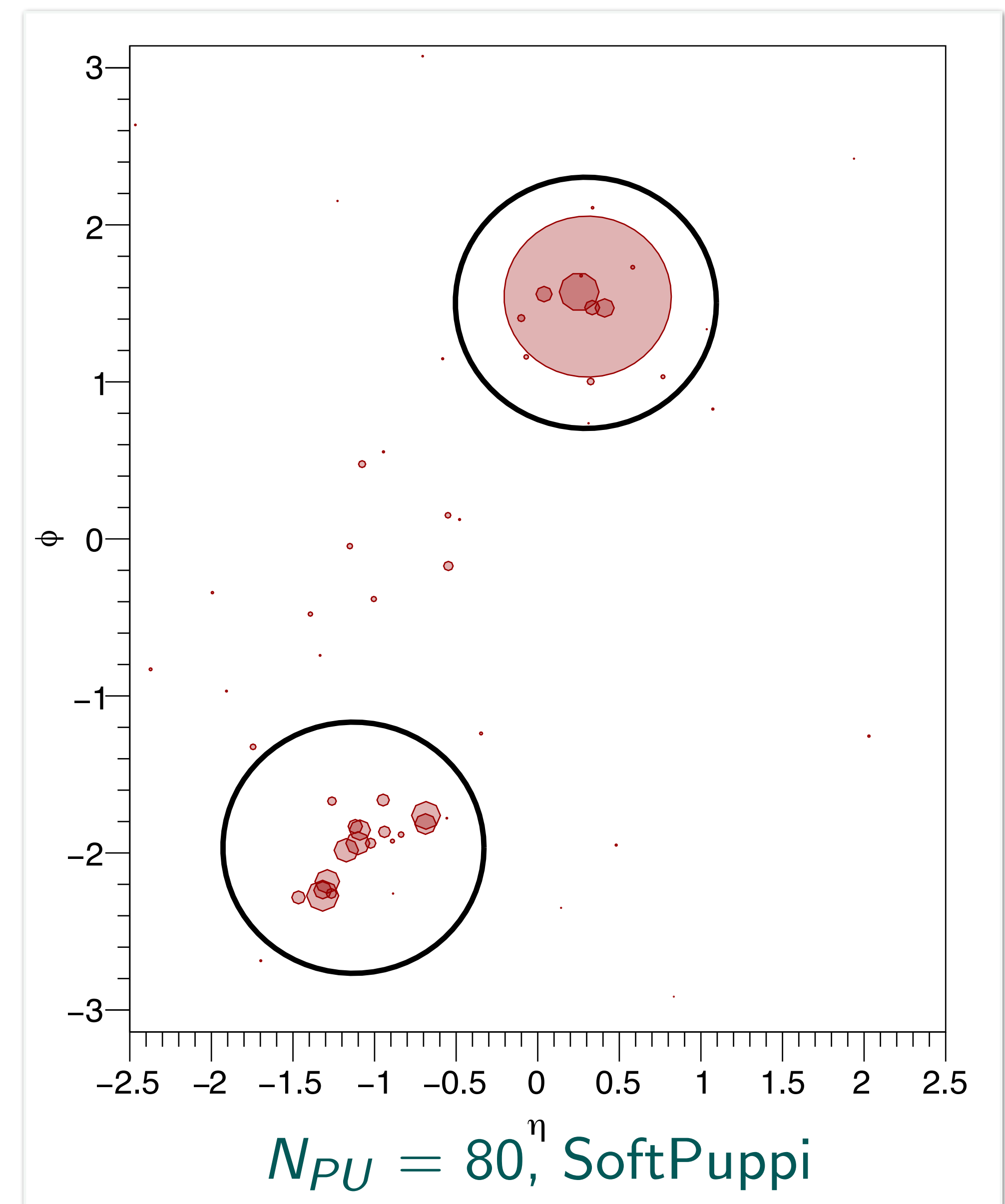
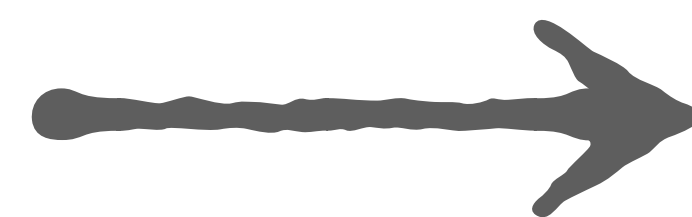
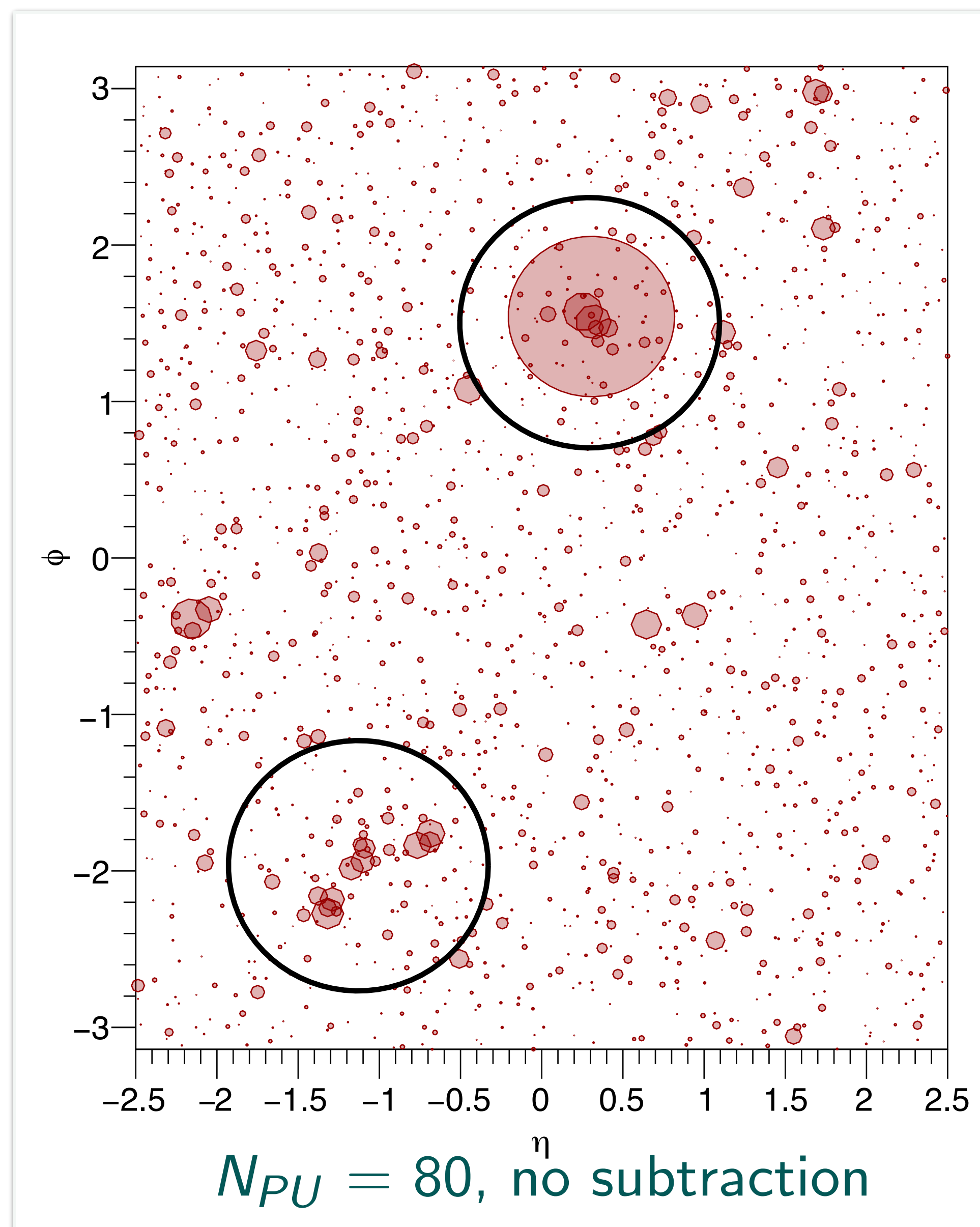
[4] reweight the four-vector of the particle by this weight, then proceed to interpret the event as usual



PUPPI (PileUp Per Particle Id): based on PF paradigm

a general framework that determines, per particle, weight for **how likely** a particle is from PU

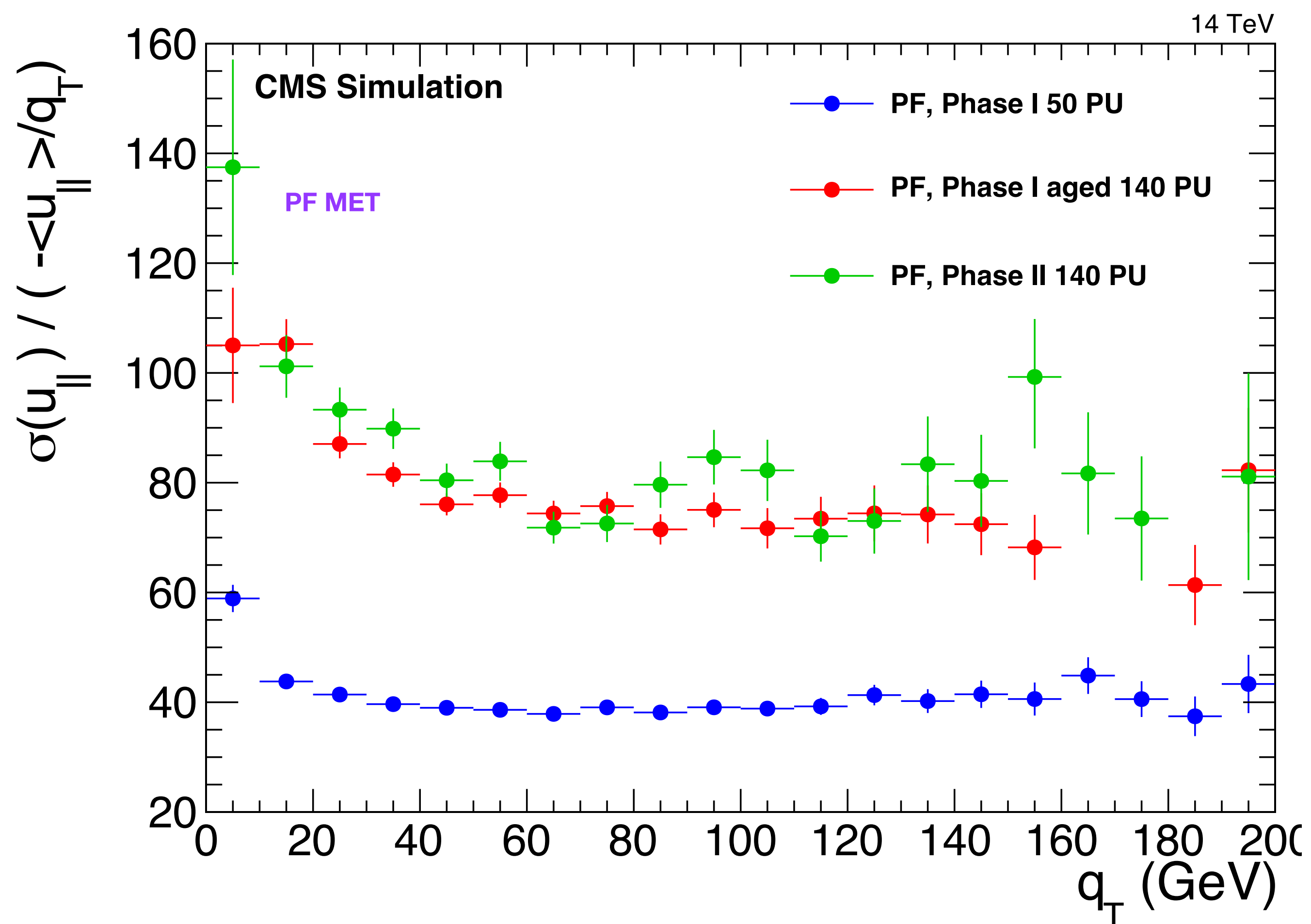
key insight: using QCD ansatz to infer neutral pileup contribution



Large gains from PUPPI, especially at high PU

arXiv:1407.6013 (original), LHCC-P-008 (CMS TP), JME-14-001, CMS analyses,...

MET resolution, 140 PU
[LHCC-P-008]



many gains at high PU

jet pT resolution

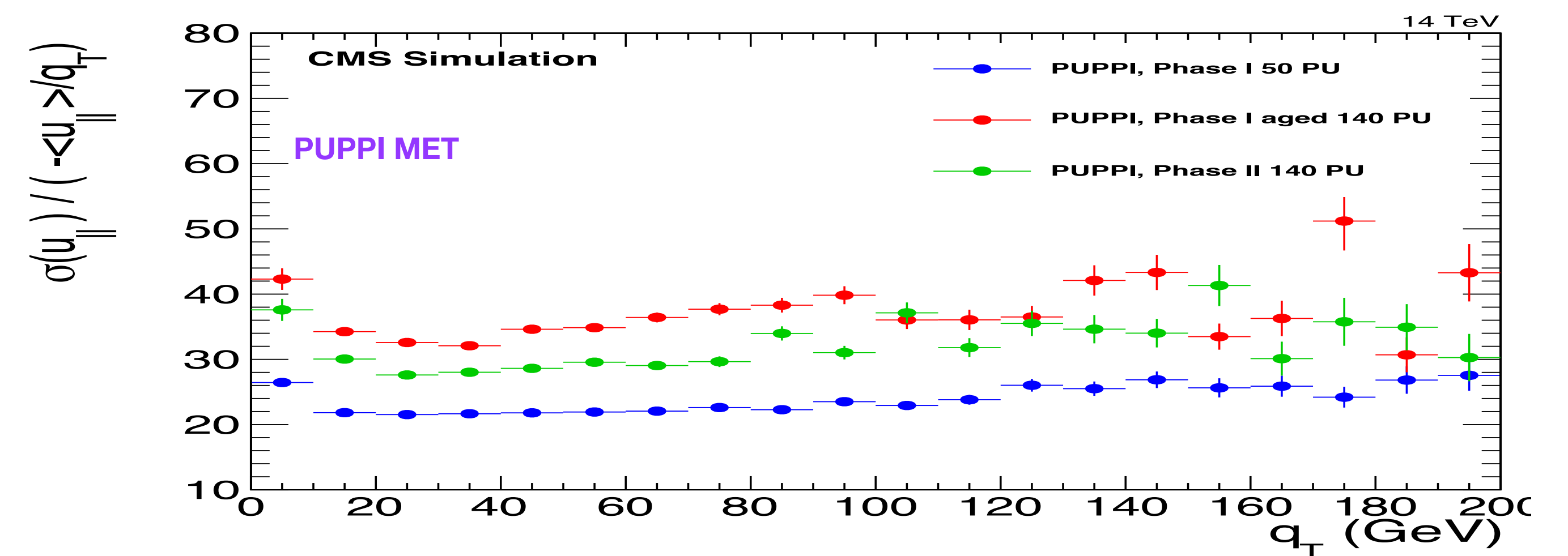
fake jet rate

MET resolution

jet substructure

lepton isolation

...

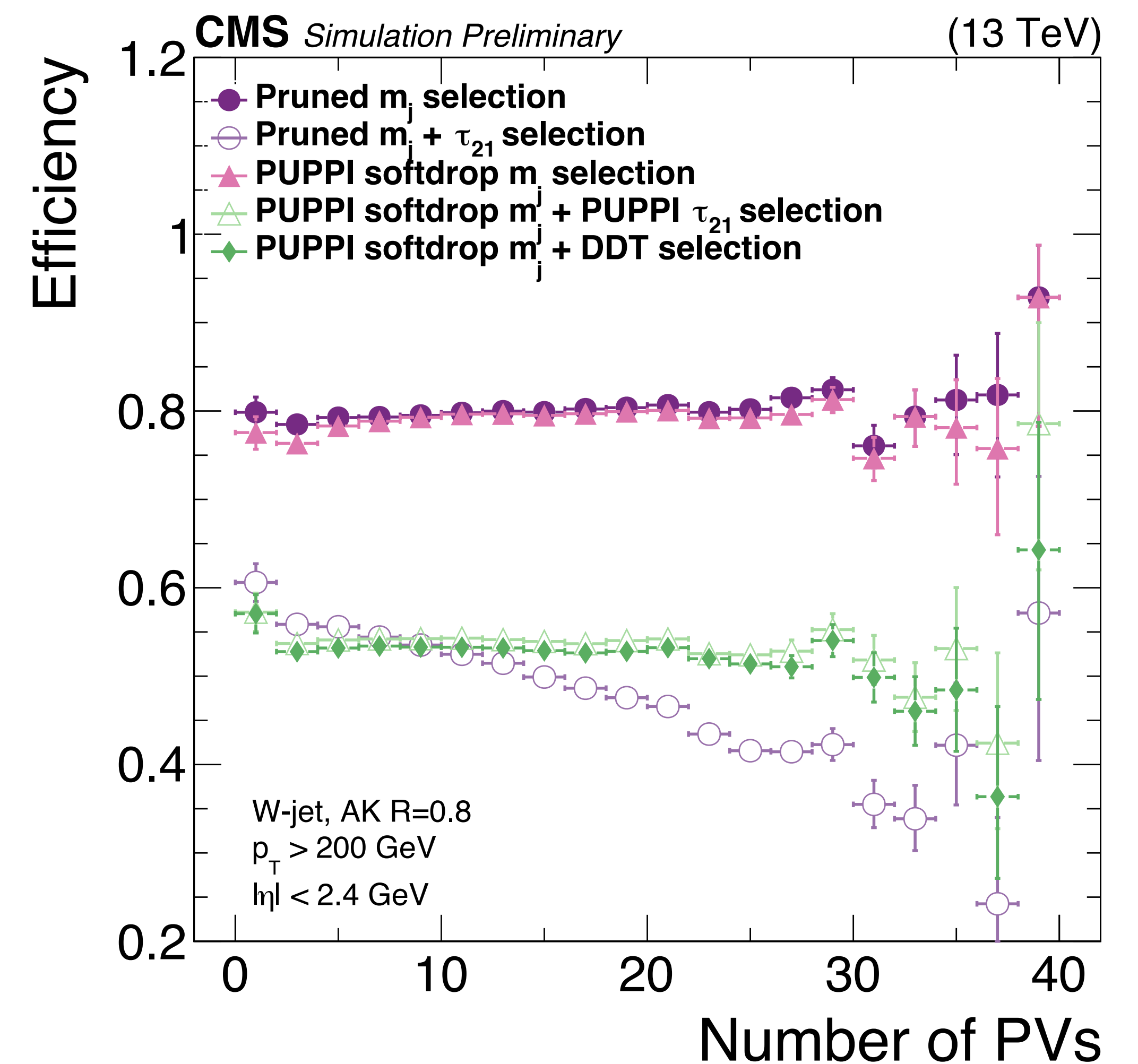
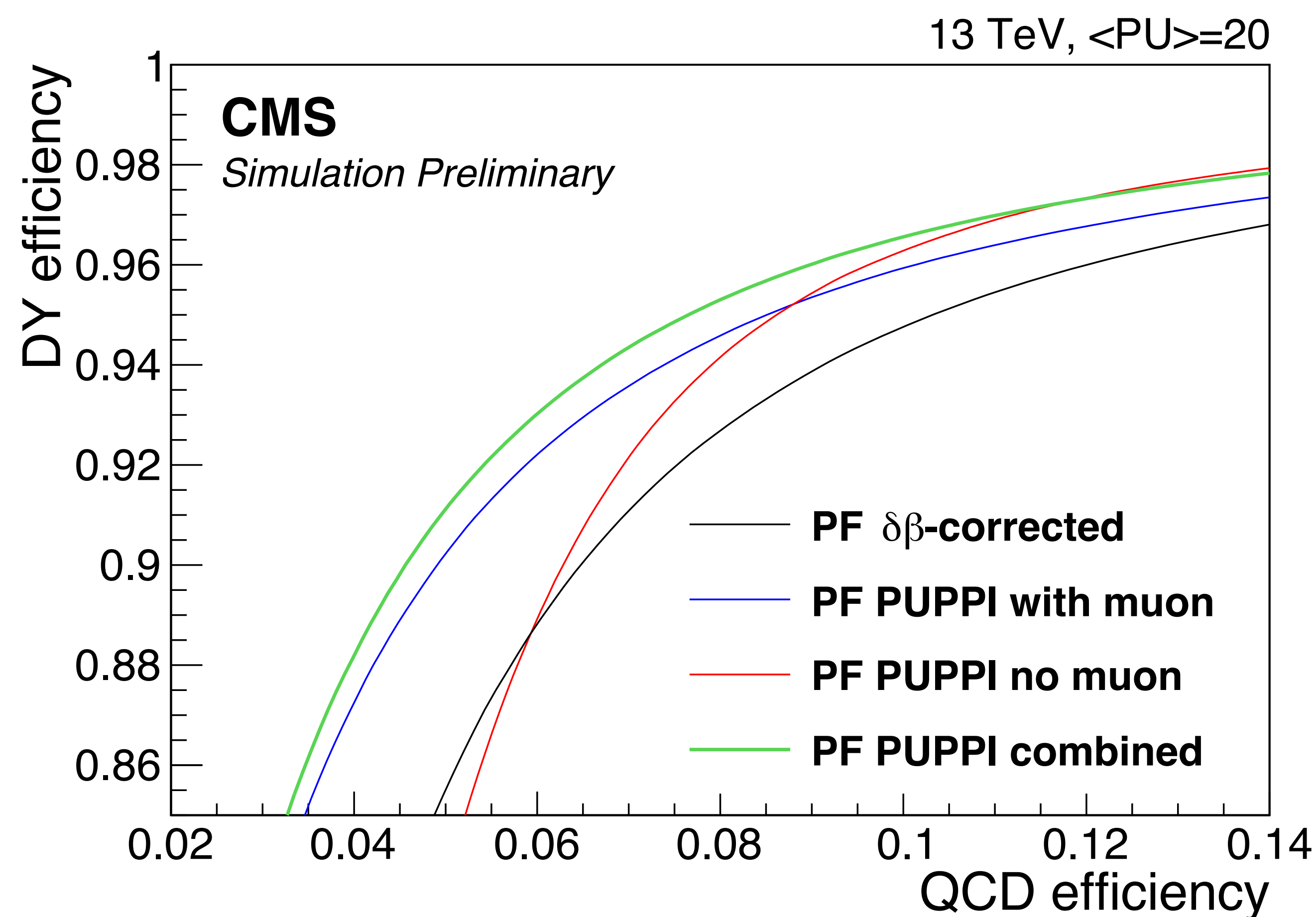


Trying to preserve soft, hidden physics

Things hidden in jets and jet substructure

Isolated, soft leptons in high PU environments

* *Examples plots from offline studies*



Large gains is soft muon backgrounds and jet substructure

Implementation of Puppi proof-of-concept using High level synthesis (HLS) as well

COMPUTE FOR EACH NEUTRAL

[1] define a local discriminant, α , between pileup (PU) and leading vertex (LV)

$$\alpha_i^C = \log \left[\sum_{j \in \text{Ch, LV}} \frac{p_{T,j}}{\Delta R_{ij}} \Theta(R_0 - \Delta R_{ij}) \right]$$

Implementation of Puppi proof-of-concept using High level synthesis (HLS) as well

COMPUTE FOR EACH NEUTRAL

[1] define a local discriminant, α , between pileup (PU) and leading vertex (LV)

$$\alpha_i^C = \log \left[\sum_{j \in \text{Ch, LV}} \frac{p_{T,j}}{\Delta R_{ij}} \Theta(R_0 - \Delta R_{ij}) \right]$$

[2] get data-driven α distribution for PU using charged PU tracks

[3] for the neutrals, ask “how un-PU-like is α for this particle?”, compute a weight

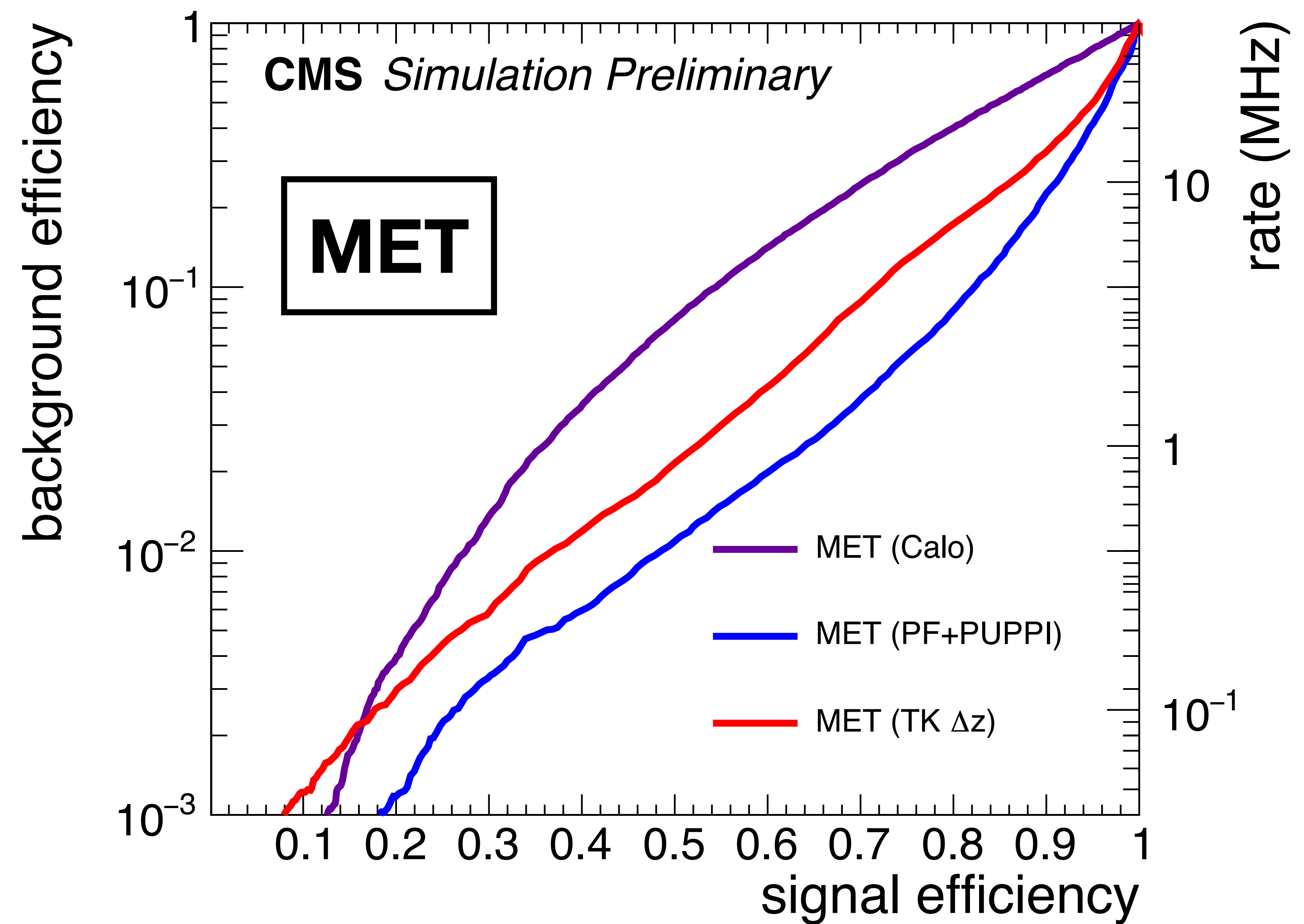
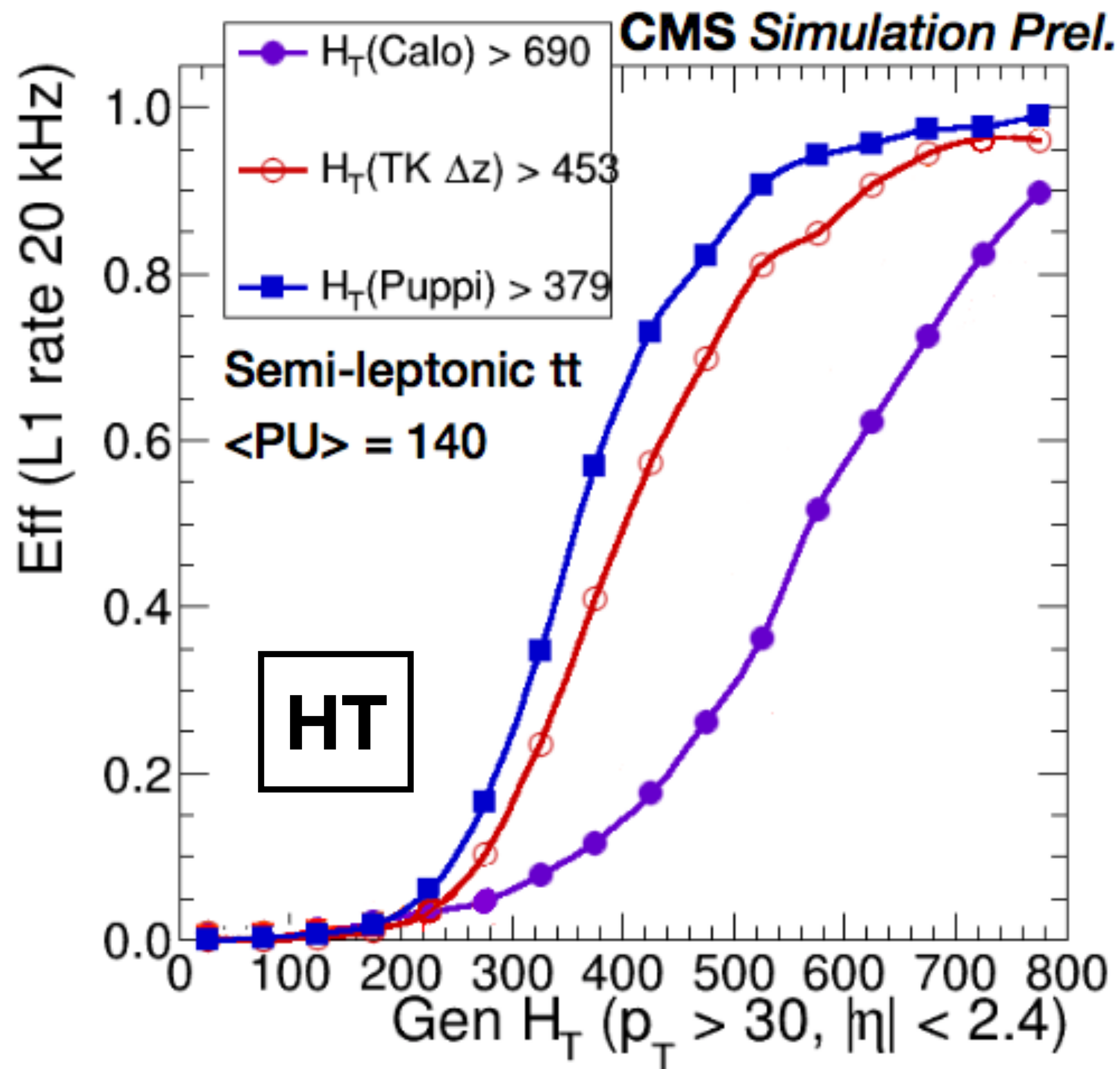
[4] reweight the four-vector of the particle by this weight, then proceed to interpret the event as usual

PRECOMPUTE STEP 2 OFFLINE
WITH CONSTANTS
(FOR GIVEN PILEUP LEVEL)

DO STEP 3/4 WITH A
LOOK-UP TABLE

RESOURCE USAGE ONLY FEW %
OF FPGA AND 100S OF NS
LATENCY WITH LITTLE
DEGRADATION IN PERFORMANCE

First physics results on HT and MET triggers for CMS phase-2 trigger interim document

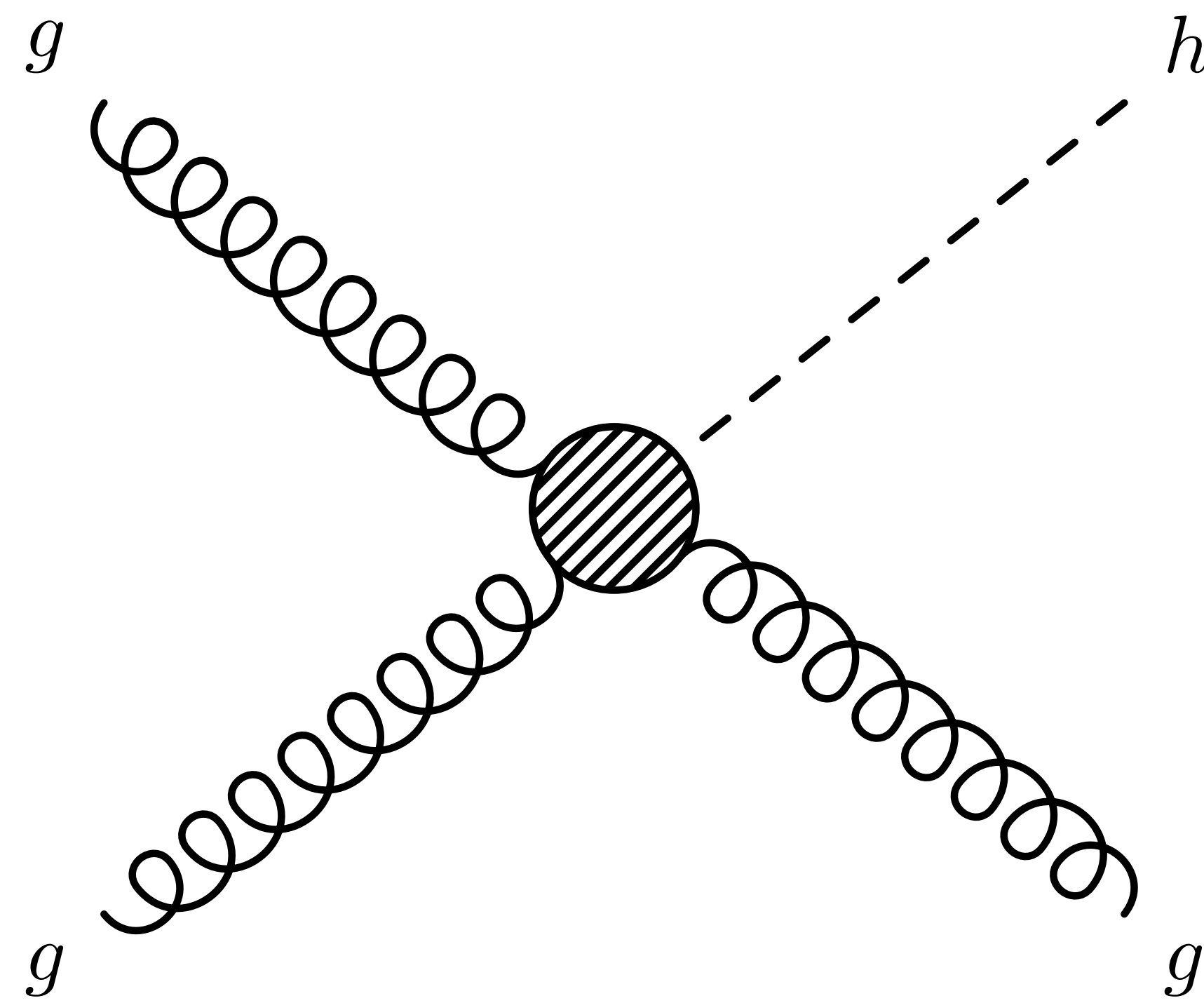
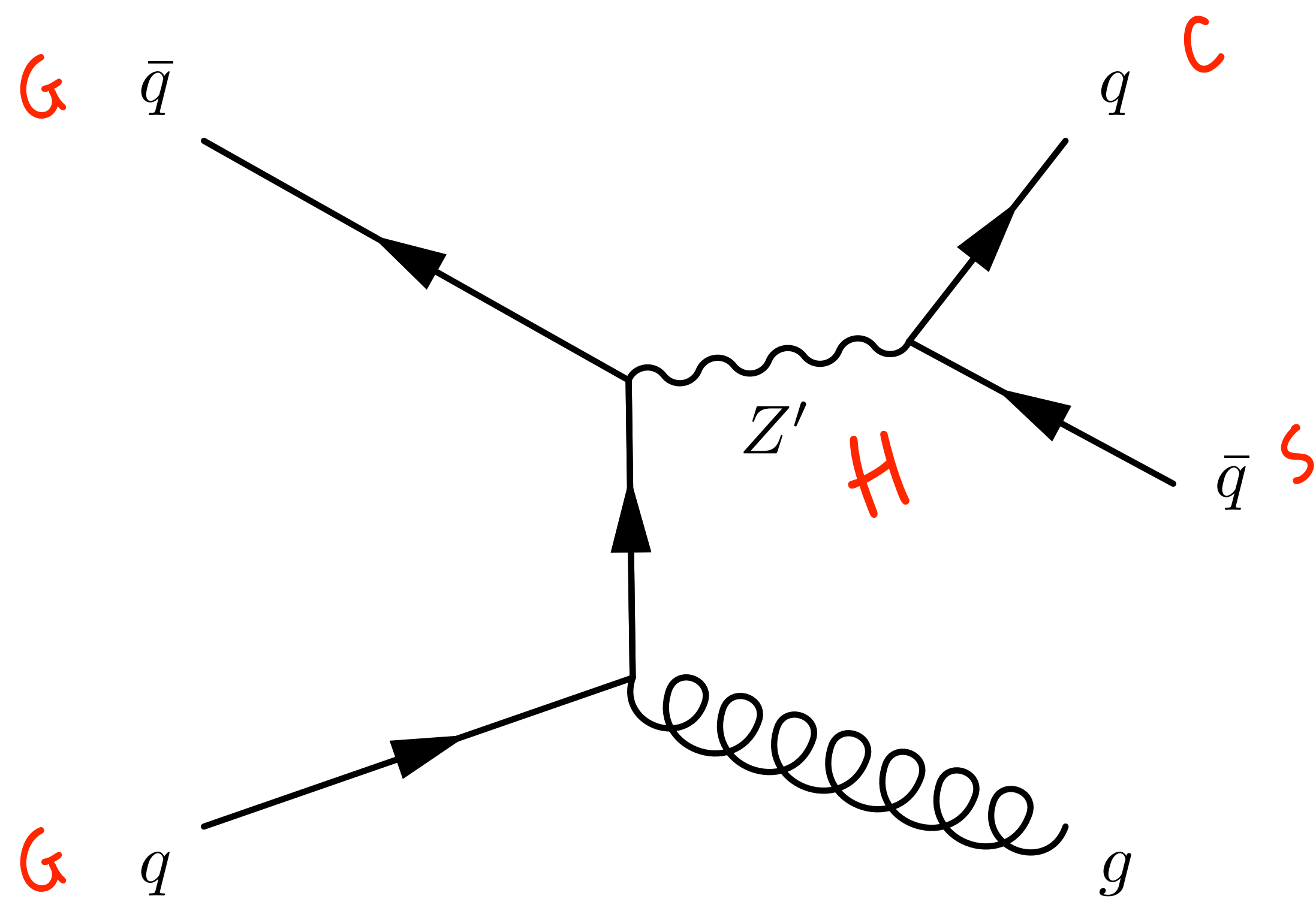


Gains in rate reduction, signal efficiency, lower thresholds

Bringing advanced physics algorithms to the hardware trigger!

Proof-of-concept PF+PUPPI running at L1

large physics gains: HT, MET, jet (substructure), lepton isolation



Other advanced algorithms...

how about **machine learning**?

TRIGGER ON
TOPOLOGIES LIKE THIS
AT L1!

high level synthesis for machine learning

~~HLS4ML~~

HLS4ML

JENNIFER NGADIUBA, MAURIZIO PIERINI (CERN)
JAVIER DUARTE, SERGO JINDARIANI, BEN KREIS, NHAN TRAN (FNAL)
PHIL HARRIS (MIT)
ZHENBIN WU (UIC)

+ EJ KREINAR (HAWKEYE 360) AND SONG HAN (GOOGLE/STANFORD)

Many parts of the trigger could benefit machine learning

clustering, fitting (regression), classification, anomaly detection

Not just LHC physics or triggering

DAQ, neutrino physics, intensity frontier, ...

No industry solutions: LHC latency constraints are unheard of

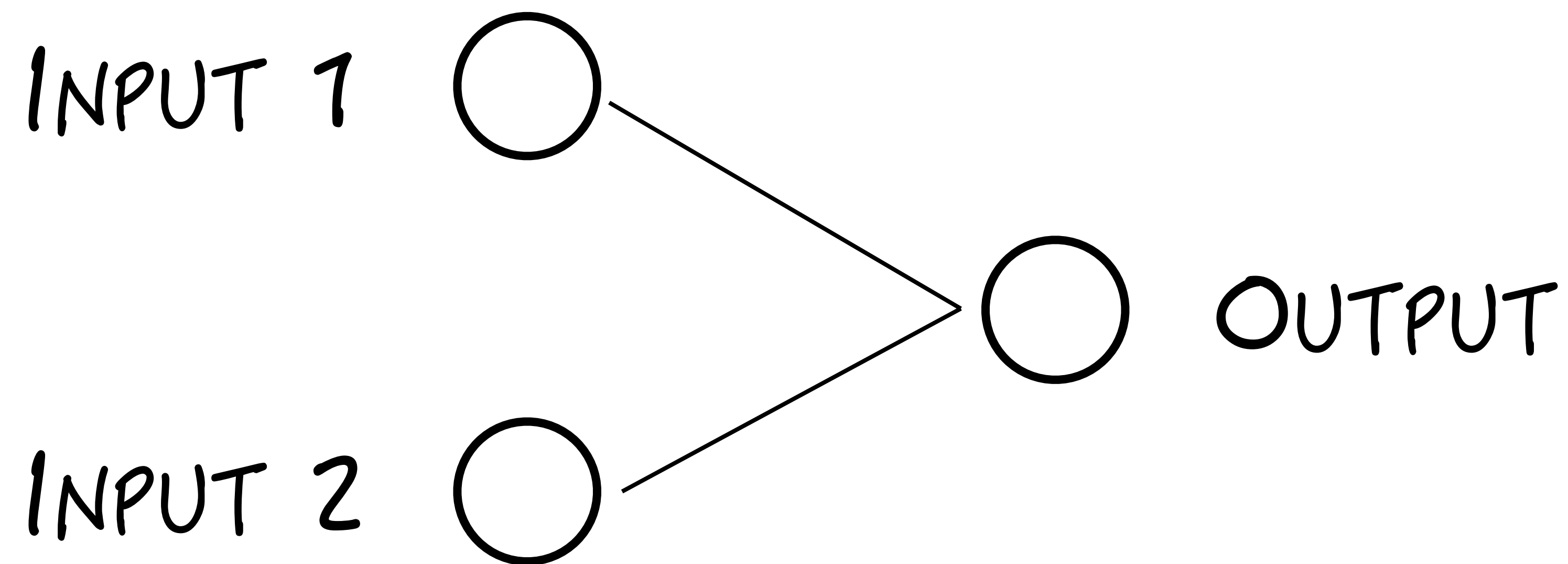
Why HLS?

HLS allows (super)-fast algorithm development

Write a tool for machine learning *inference** at low latencies:

HLS4ML

*for training, GPUs remain top dog

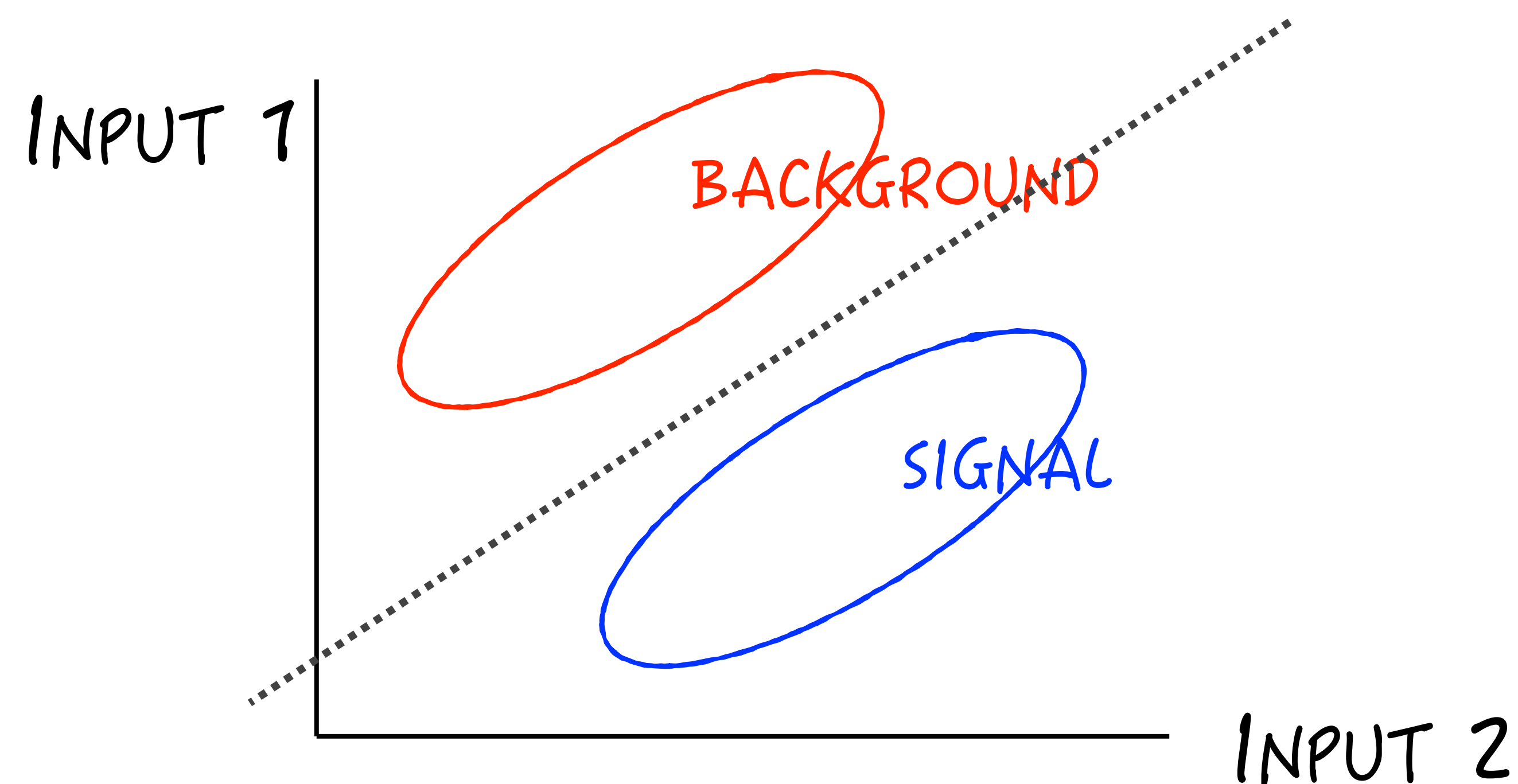


$$\vec{O}_j = \vec{l}_i \times \vec{W}_{ij} + \vec{b}_j$$

Simple 2 input example

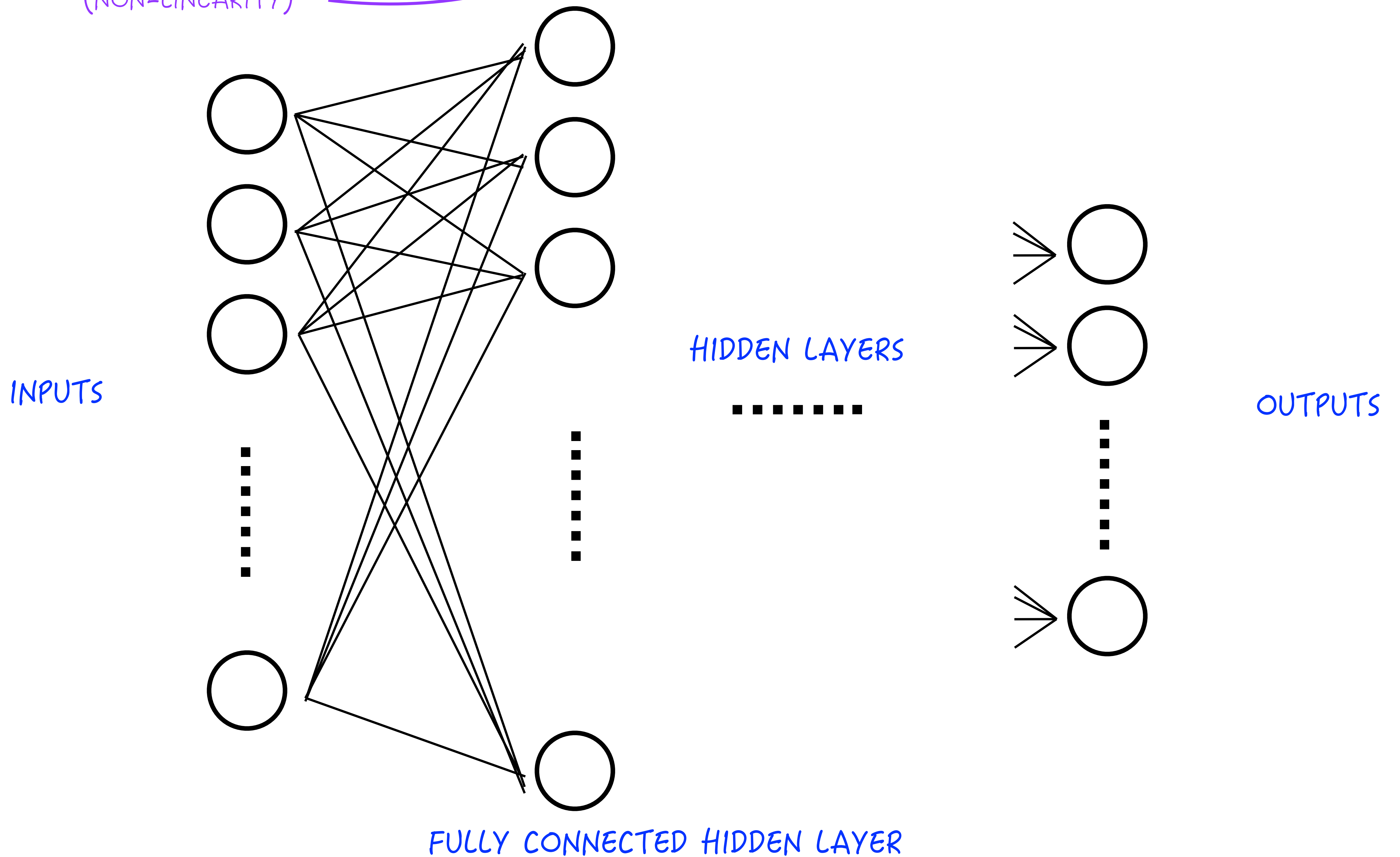
(Fisher linear discriminant, linear support vector machine,...)

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



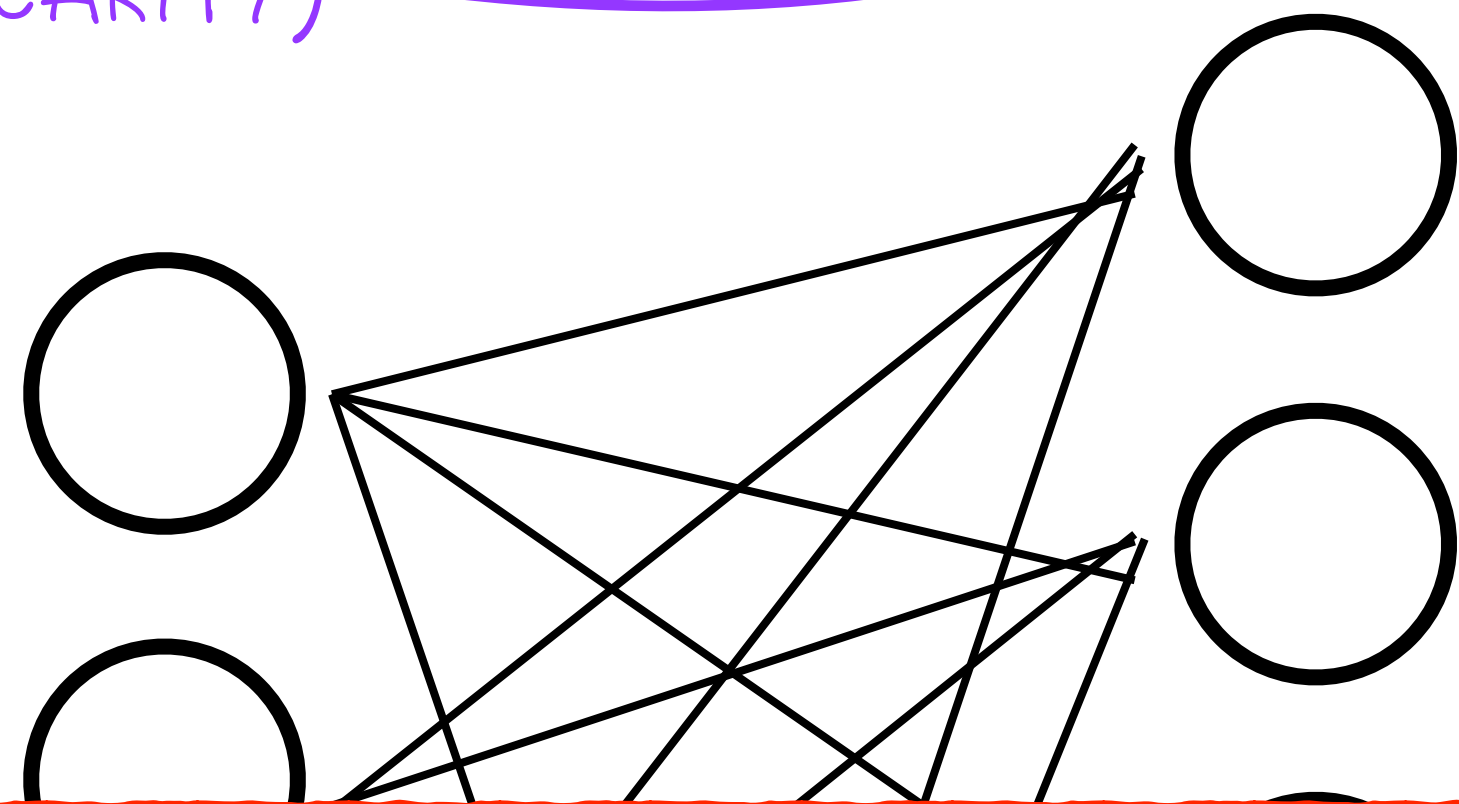
$$\vec{O}_j = \Phi(\vec{l}_i \times \vec{W}_{ij} + \vec{b}_j)$$

Φ = ACTIVATION FUNCTION
(NON-LINEARITY)



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

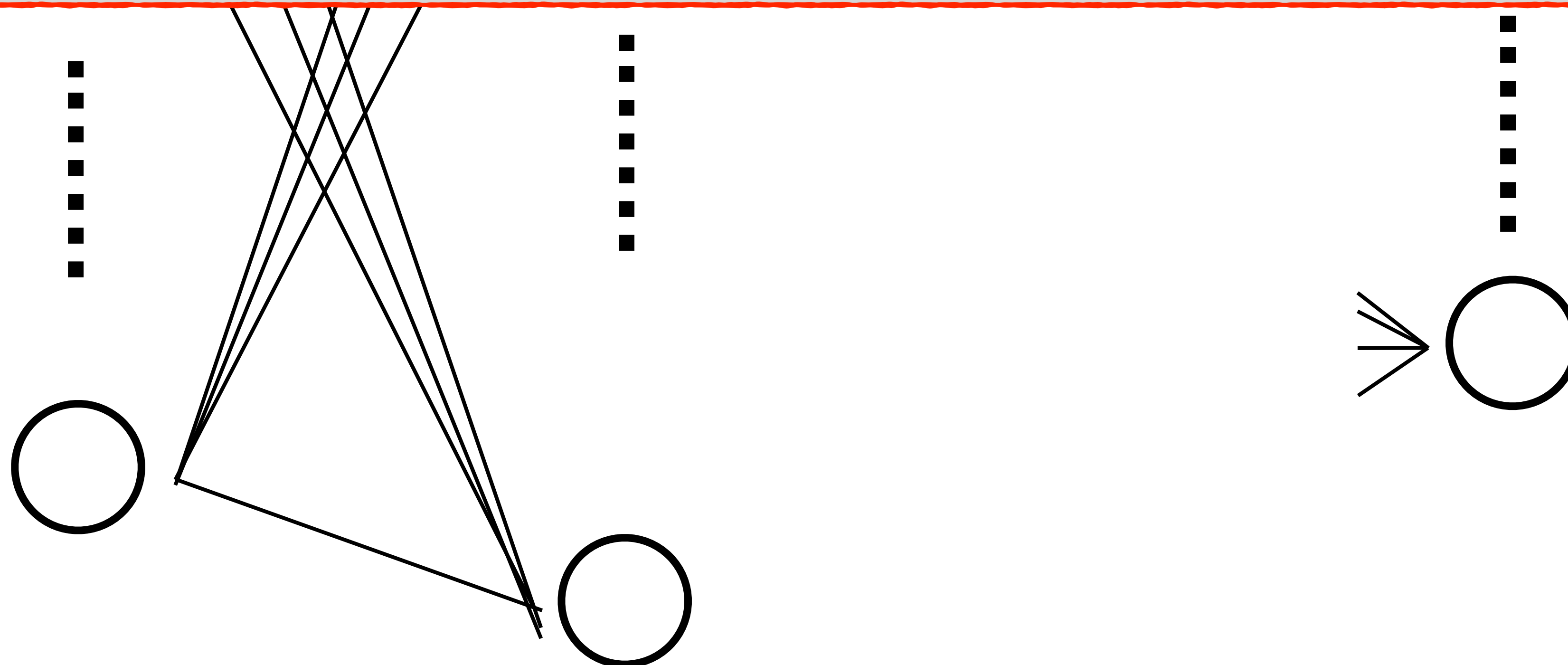
Φ = ACTIVATION FUNCTION
(NON-LINEARITY)



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

IN

TPUTS



FULLY CONNECTED HIDDEN LAYER

Emergent engineering field, efficient implementation of NN architecture

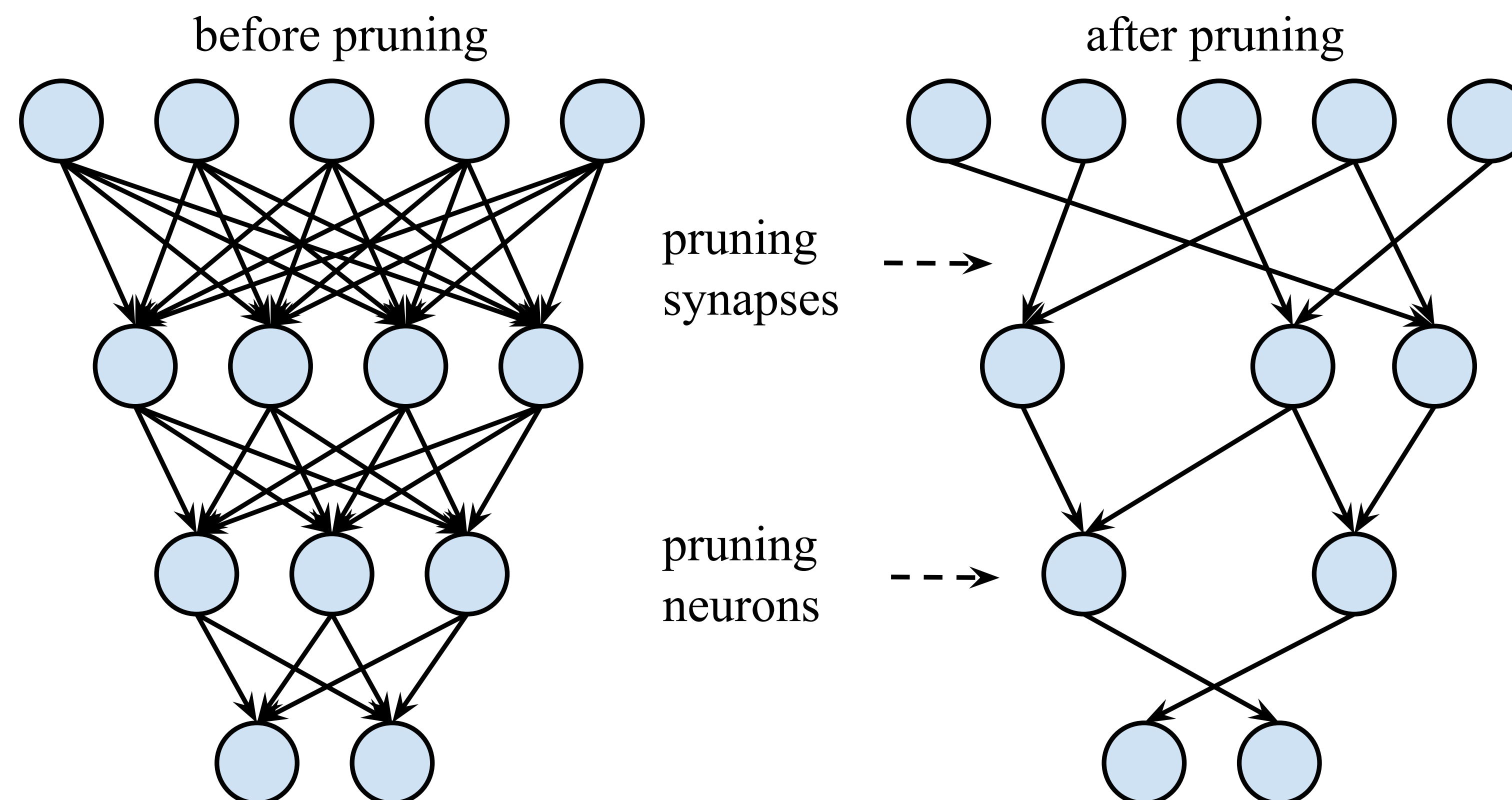
Compression/Pruning:

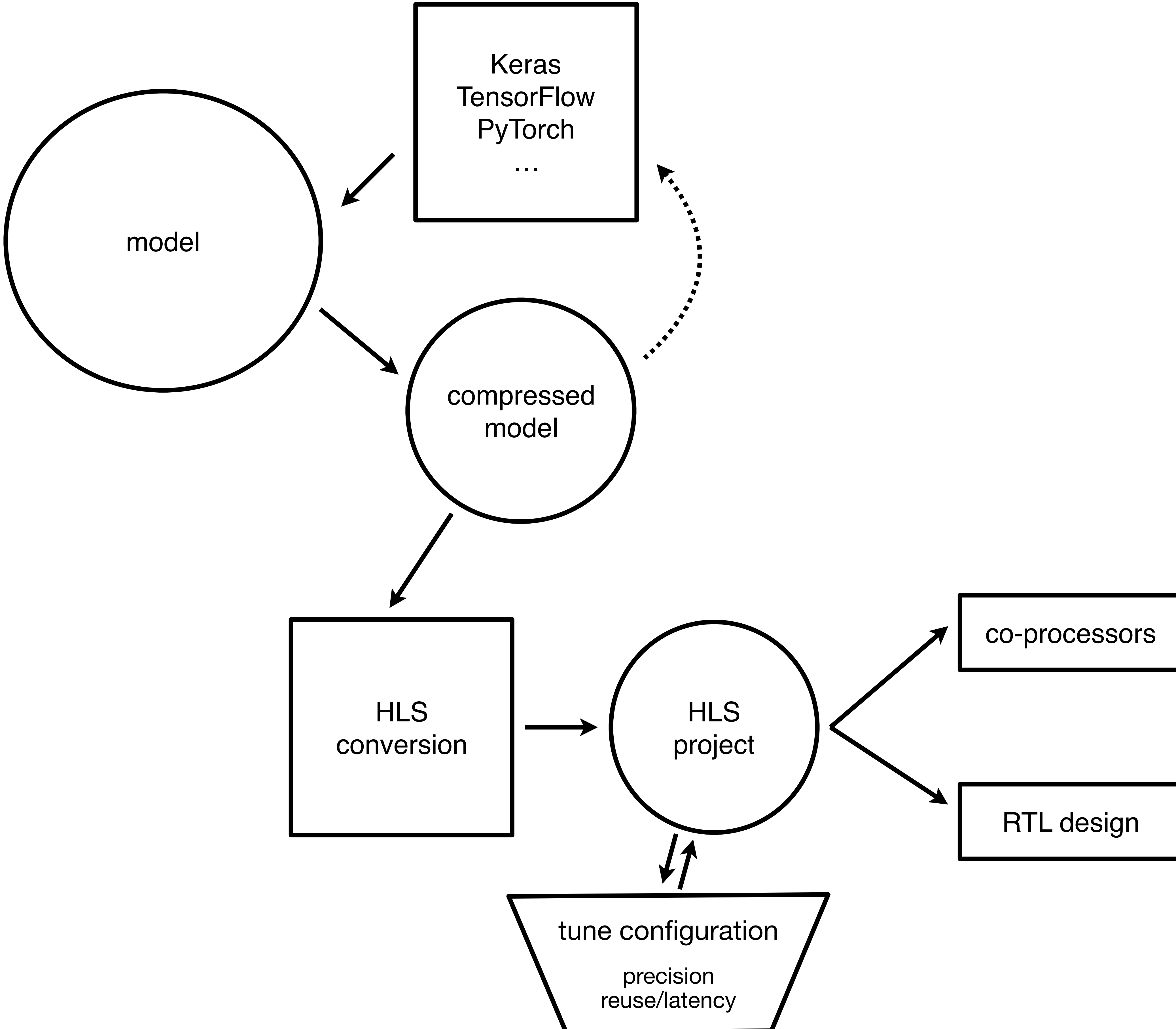
maintain the same performance while removing low weight synapses and neurons (many schemes)

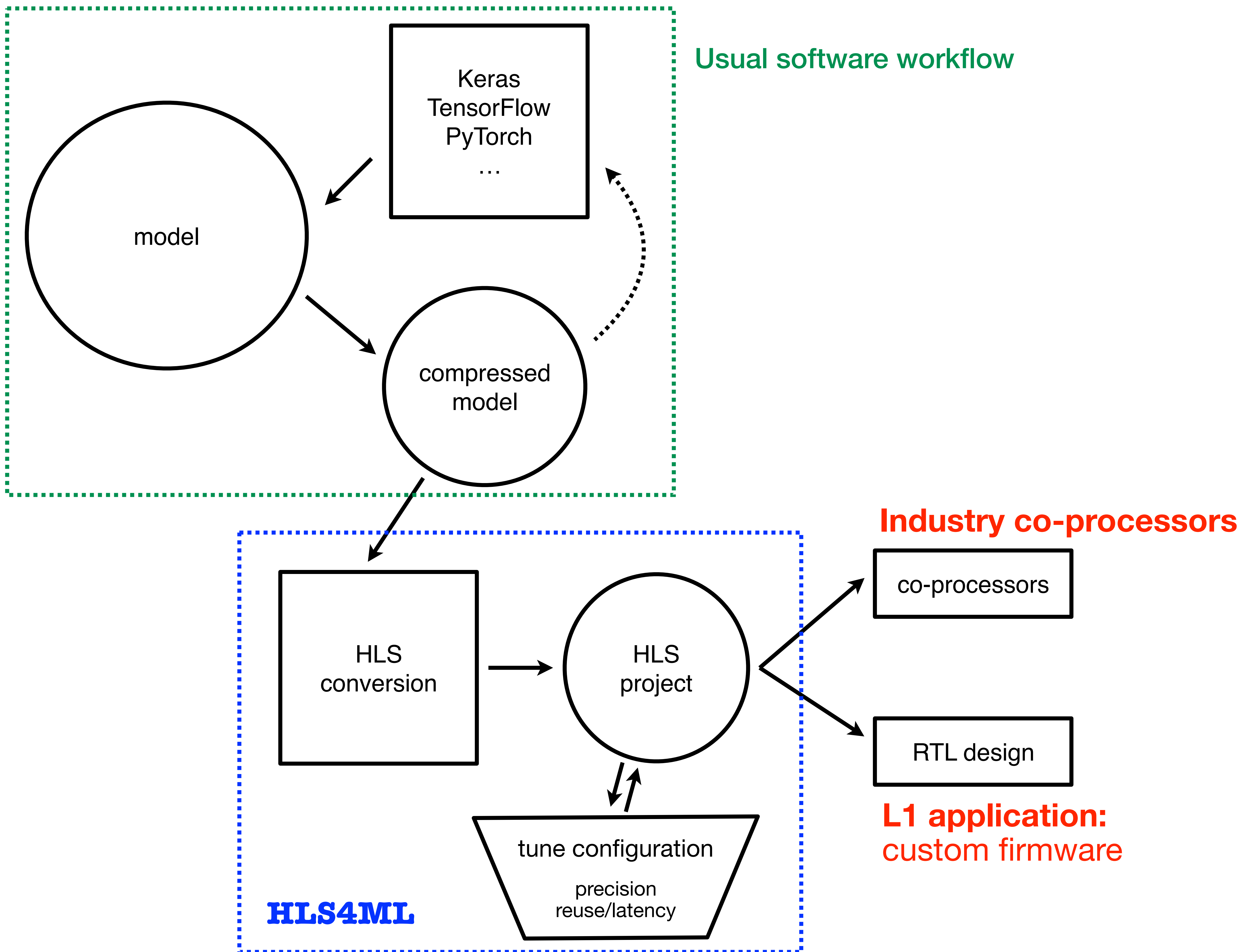
Quantization/Approximate math:

32-bit floating point math is overkill

20-bit, 18-bit, ...? fixed point, integers? binarized NNs?







```
python keras-to-hls.py -c keras-config.yml
```

```
1 KerasJson: example-keras-model-files/KERAS_1layer.json
2 KerasH5:   example-keras-model-files/KERAS_1layer_weights.h5
3 OutputDir: my-hls-dir-test
4
5 IOType: io_parallel # options: io_serial/io_parallel
6 ReuseFactor: 1
7 DefaultPrecision: ap_fixed<18,8>
```

} Keras/TF inputs

IOType: parallelize or serialize

ReuseFactor: how much to parallelize

DefaultPrecision: self-explanatory :)

EXAMPLE: JET SUBSTRUCTURE

5 output multi-classifier:

Does a jet originate from a quark, gluon, W/Z boson, top quark?

Network architecture

16 expert inputs

jet masses, multiplicity

ECFs ($\beta=0,1,2$)

16 inputs

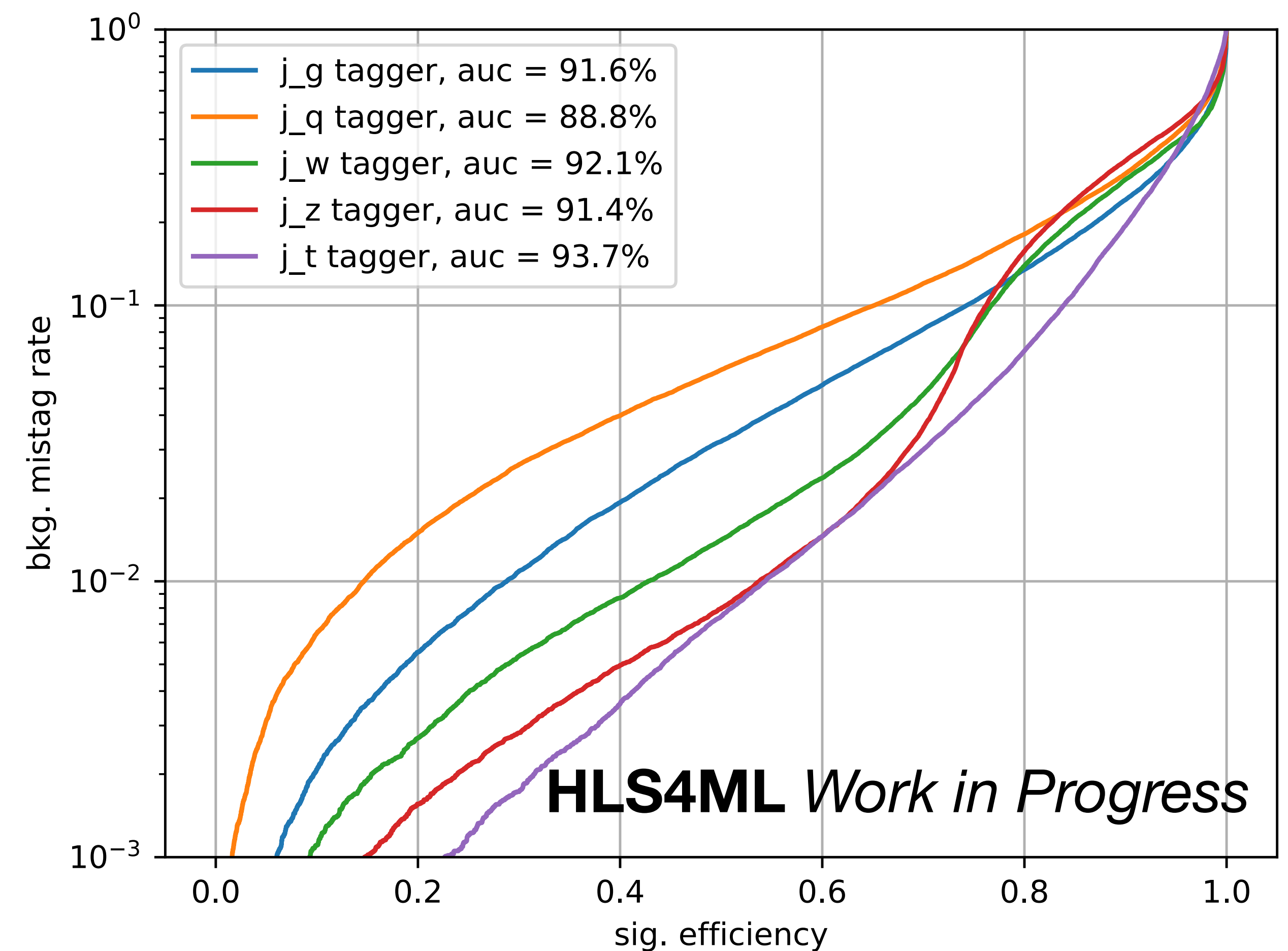
64 (relu)

32 (relu)

32 (relu)

5 outputs (softmax)

*Fully connected deep
neural network*



EXAMPLE: NETWORK (NOT JET) PRUNING

Resource usage:

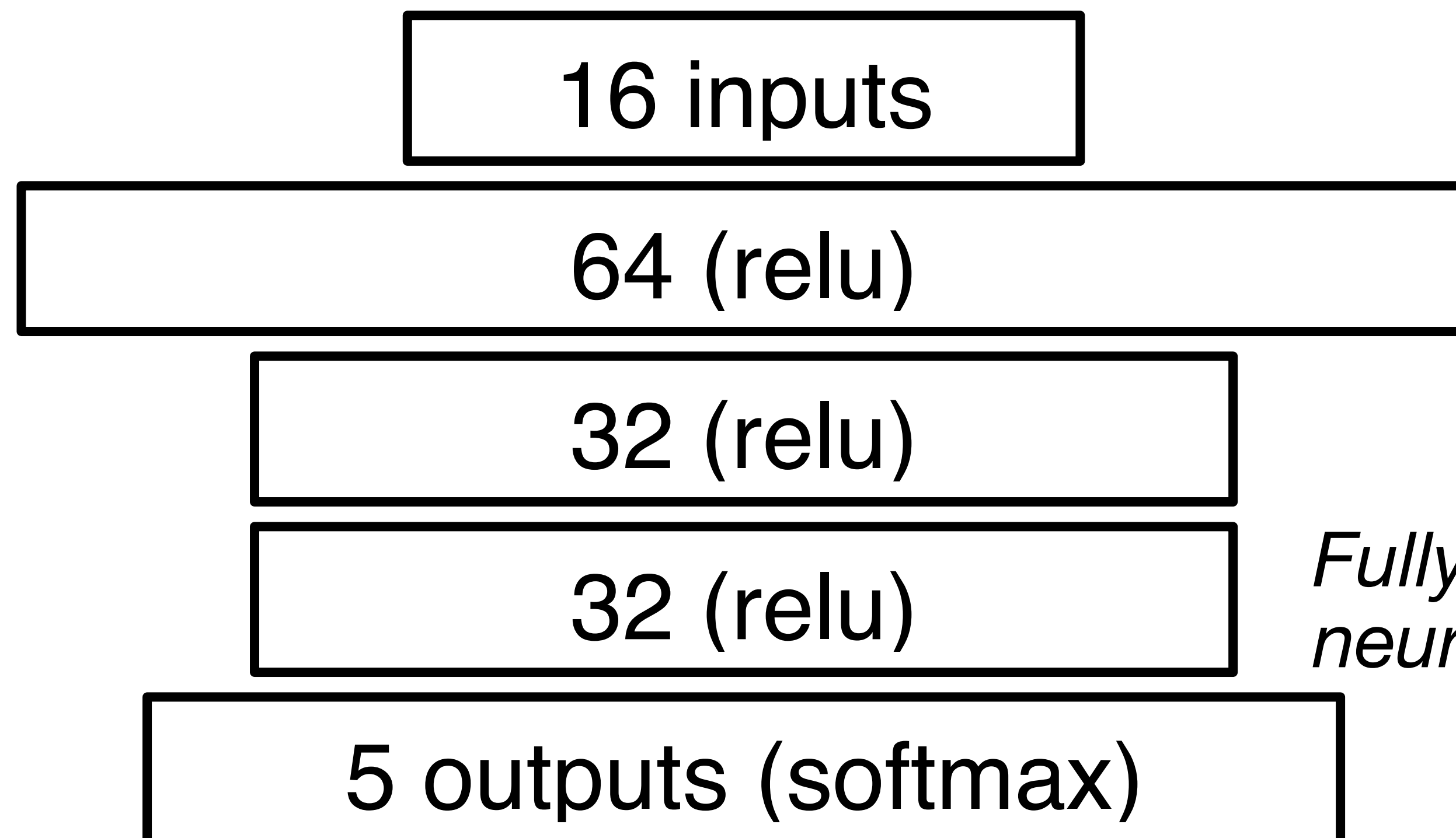
92% DSP usage for Virtex 7

61 clocks (305 ns), Pipeline = 1

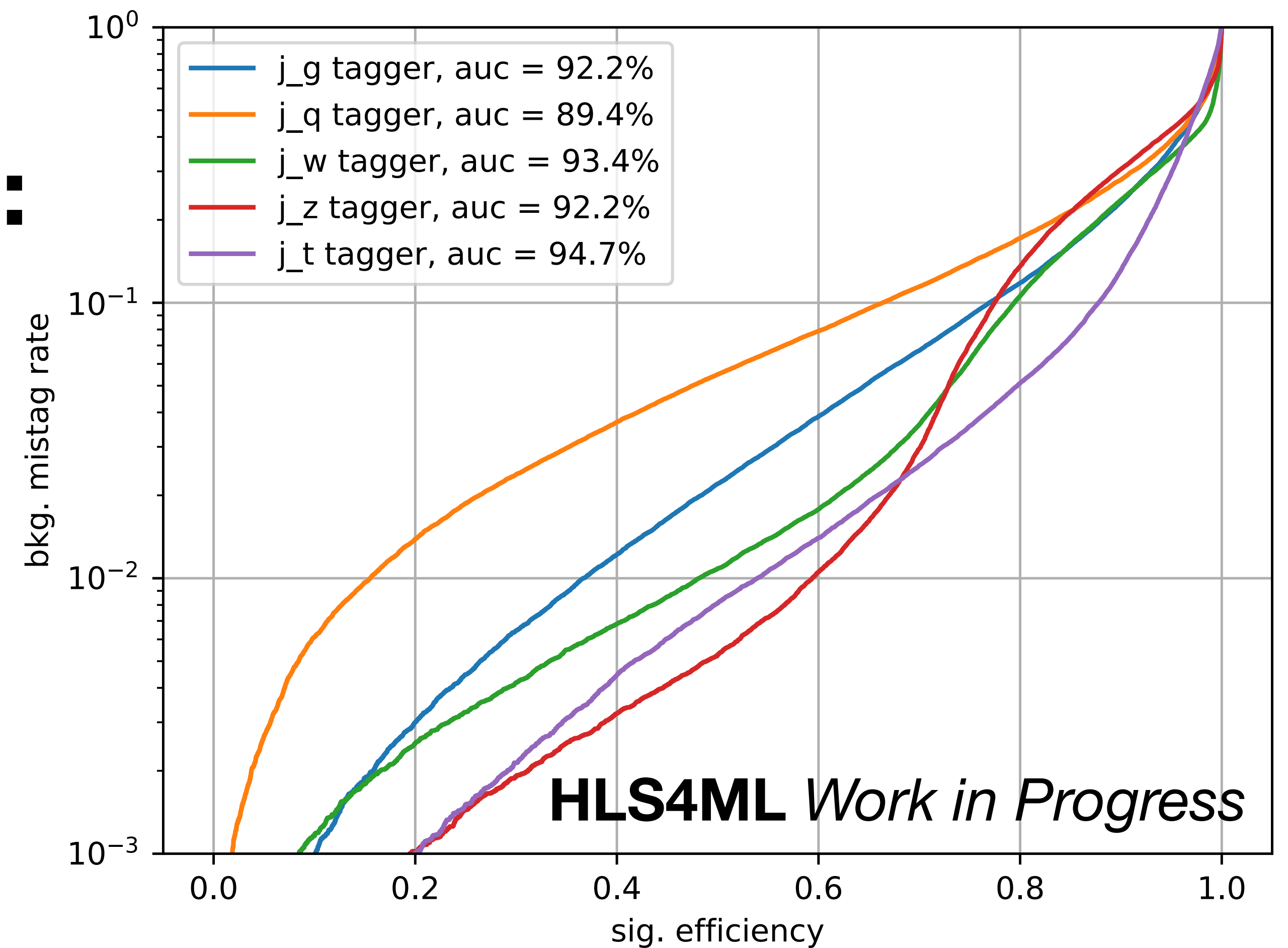
Compression (50%) + reuse = 2:

29% DSP usage for Virtex 7

60 clocks (300 ns), Pipeline = 2



*Fully connected deep
neural network*



HLS4ML

a tool to translate ML algorithms for FPGAs in minutes
highly parallelizable with user controls for resource usage and latency
tunable precision, resource reuse
very efficient network design with model compression

Work in progress

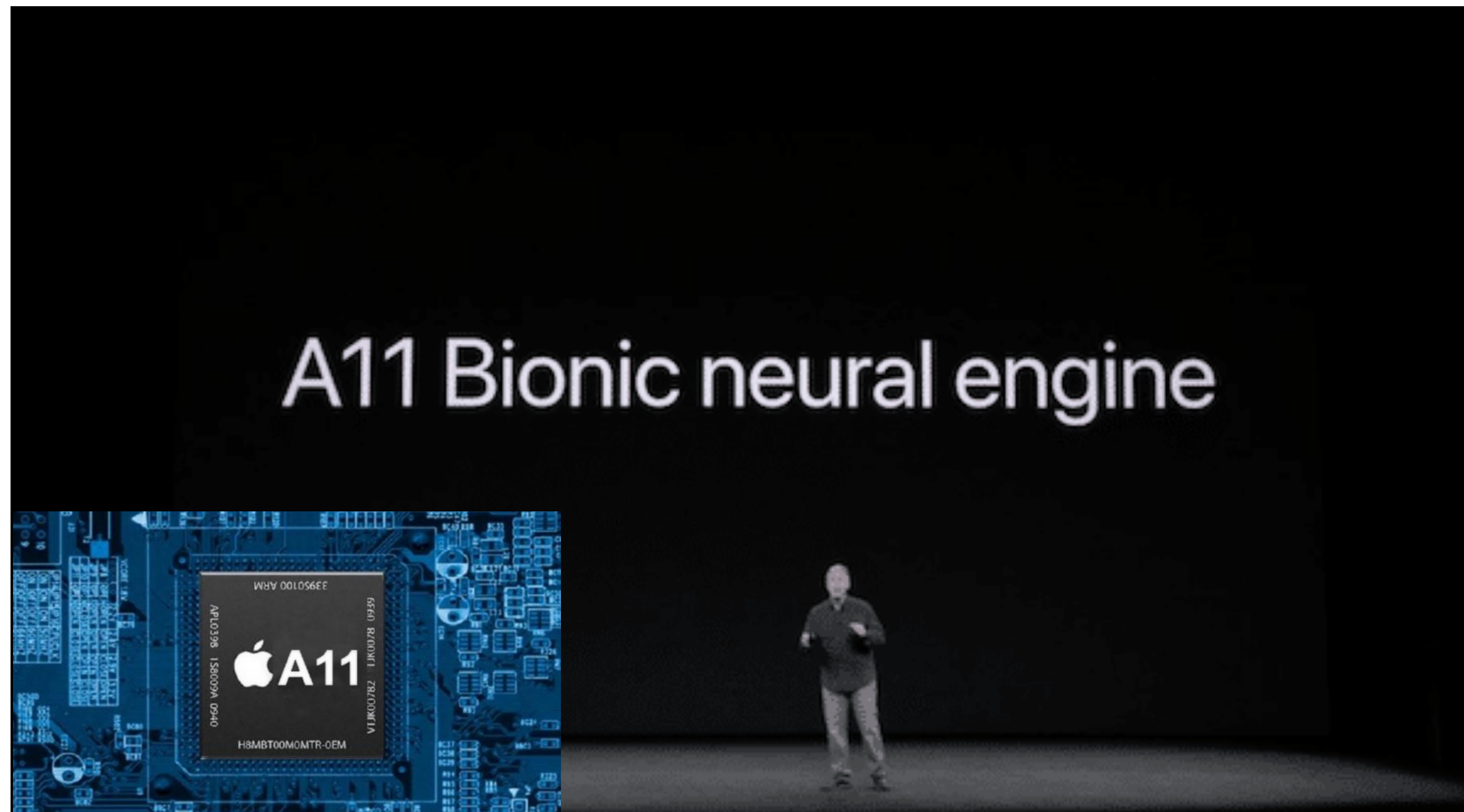
Mapping out resource usage and latency as a function of neural network hyper parameters
More network architectures: CNN (in progress), RNN/LSTM (tricky!), TMVA BDT (efficient?)

Status

Alpha version - few weeks; Targeting March-April release of Beta version
Please contact us if you are interested! hls4ml.help@gmail.com

**one more fun thing to think about
for the high level trigger (and beyond?)**

Specialized co-processor hardware for machine learning inference



INTEL[®] FPGA ACCELERATION HUB

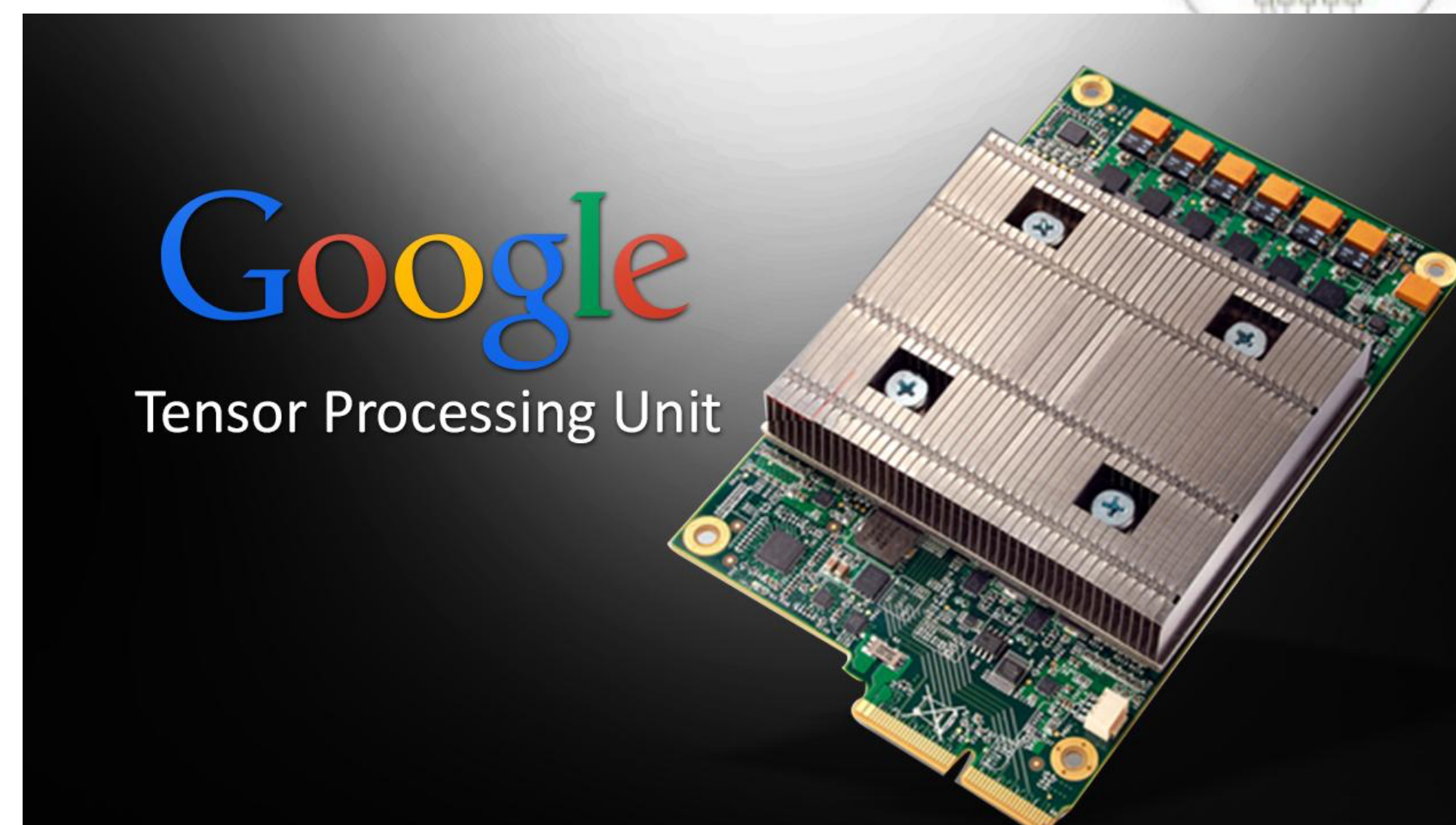
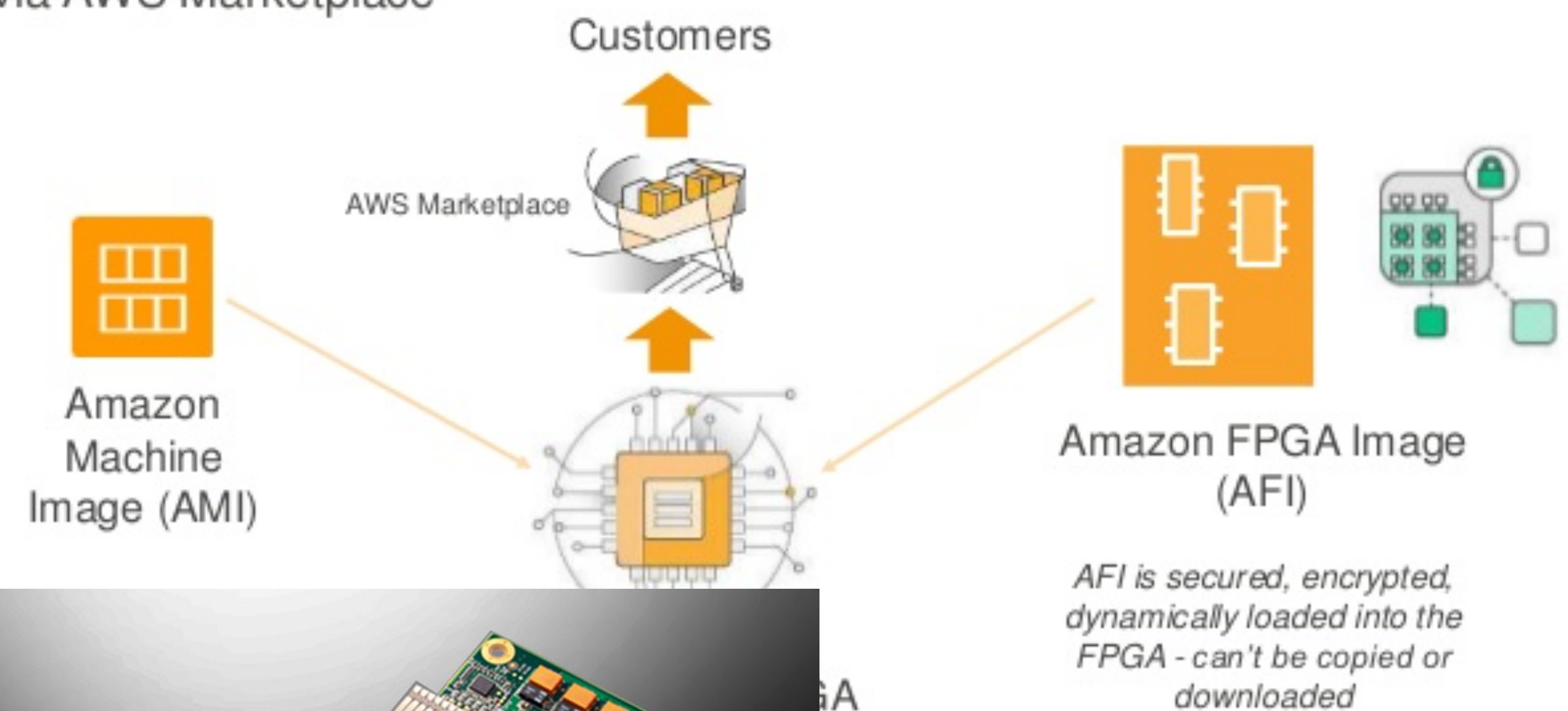
The Intel[®] Xeon[®] Acceleration Stack for FPGAs is a robust framework enabling data center applications to leverage an FPGA's potential to increase



Catapult/Brainwave



Delivering FPGA Partner Solutions on AWS via AWS Marketplace



Specialized co-processor hardware for machine learning inference



It already exists!
One example: Microsoft catapult

The screenshot shows a web browser window with a translation tool. The tool is set to translate from Wikipedia to Spanish. It displays statistics for Wikipedia: over 5.2 million articles and approximately 3.1 billion words. A processor selection window is open, showing the "Azure FPGA Server - SV4-D5-1U" with a "Stratix V D5-accelerator" and a peak power of 240 Watts. A horizontal bar chart shows the compute capacity of the FPGA server, with a slider set to 1E (1 Exa-op). Below the chart, the tool reports a compute capacity of 1,000,000 Tera-ops, an estimated time of 0.098 seconds, and a rate of 78,120,000 pages per second. A "TRANSLATE" button is visible at the bottom right of the tool window.

Processor Type	Azure FPGA Server - SV4-D5-1U	Type:	10 CPU cores + 4 FPGAs
		Model:	Stratix V D5-accelerator
		Peak Power/Unit:	240 Watts
Compute Capacity	10T 100T 1P 10P 100P 1E		
Compute Capacity:	1 Exa-op		1,000,000 Tera-ops
Estimated Time:			0.098 seconds
Pages Per Second:			78,120,000

Translation of all of wikipedia in 0.1 seconds!
~O(100) times faster than CPU

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



Large gains from hardware accelerating co-processors
Industry trending towards specialized computing paradigms

Option 1

re-write physics algorithms for new hardware

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

Hardware: FPGA, GPU

Option 2

re-cast physics problem as a machine learning problem

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

Hardware: FPGA, GPU, ASIC

PF & PUPPI @ L1
TRIGGER

HLS4ML/
INDUSTRY?

Why (Deep) Machine Learning?

a common *language* for solving problems
which can universally be expressed on
optimized computing hardware and follow industry trends

Recent advances in hardware and compilation/synthesis allow for sophisticated techniques at low latency

Big improvements in performance,
preserve soft and hidden signatures

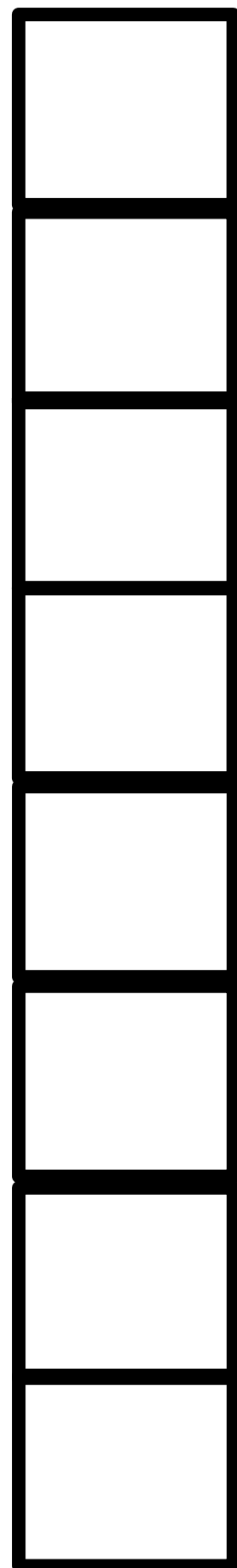
Proof-of-concept holistic pileup mitigation techniques such as PUPPI
Efficient machine learning at Level-1 Trigger
New paradigms for HLT and offline?

BONUS

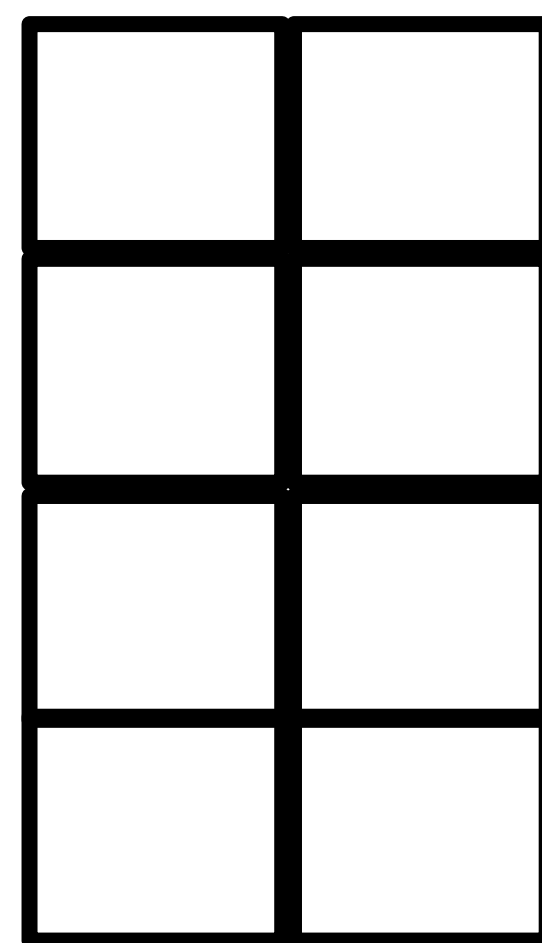
ReuseFactor: how much to parallelize operations a hidden layer

of multiplications per clock (DSPs usage)

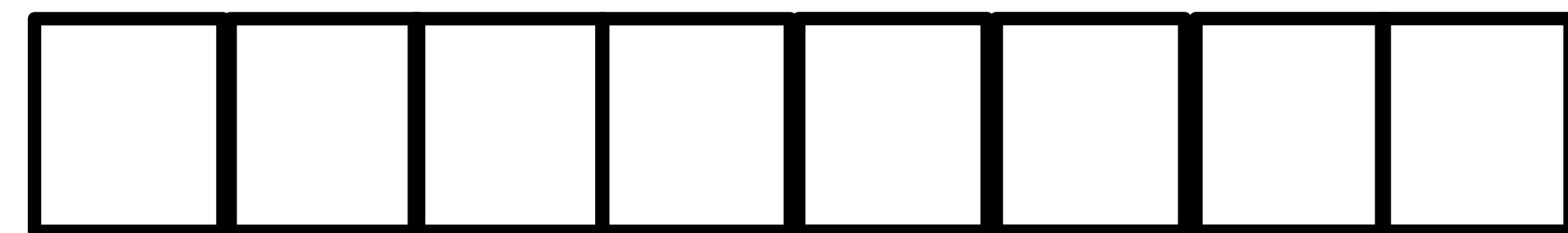
parallel
8DSPs in 1 clock
reuse = 1



parallel
4DSPs in 2 clocks
reuse = 2



serial
1 DSP in 8 clocks



time



(decreasing throughput)

EXAMPLE: NETWORK (NOT JET) PRUNING

Resource usage:

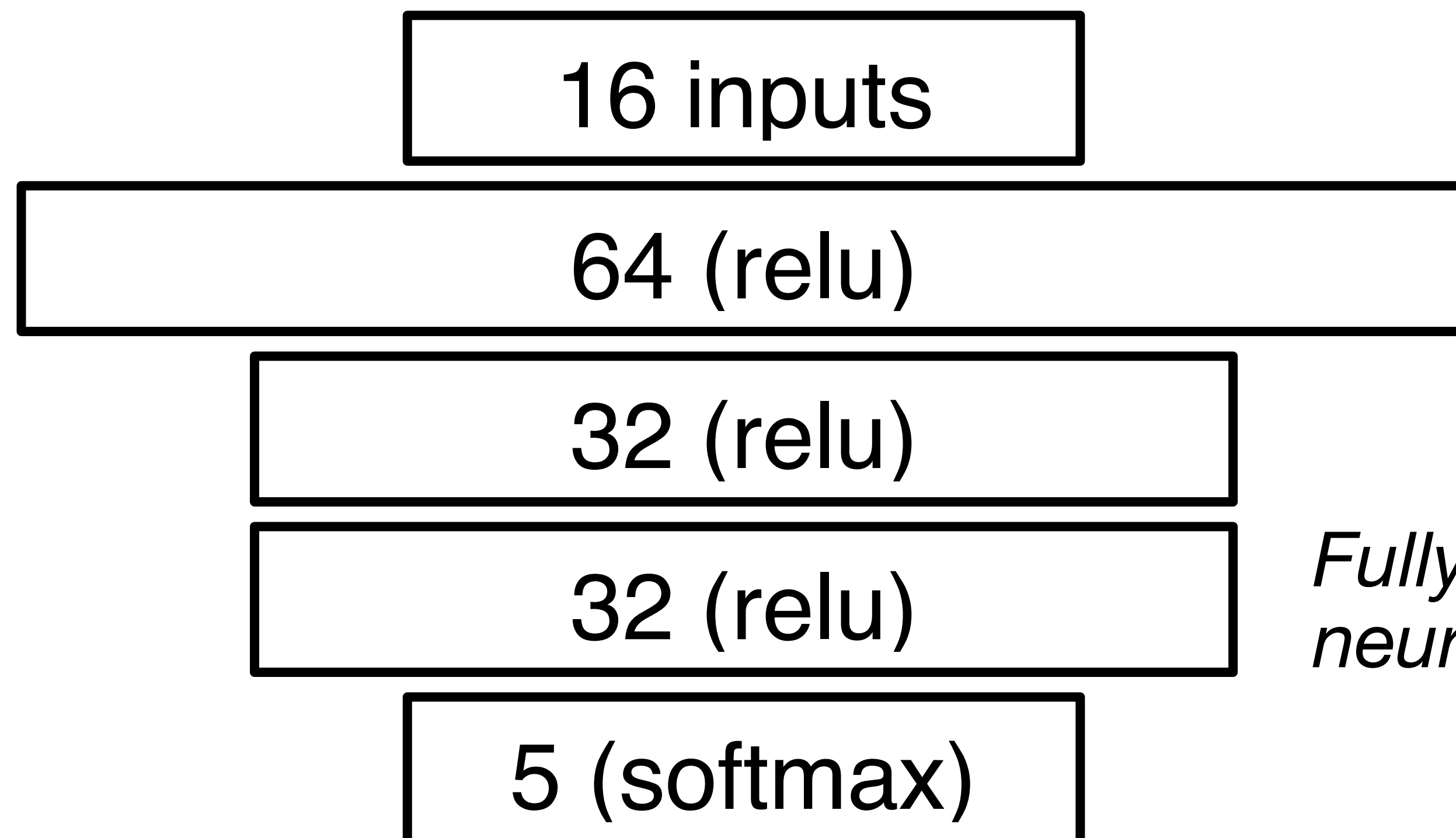
92% DSP usage for Virtex 7

61 clocks (305 ns), Pipeline = 1

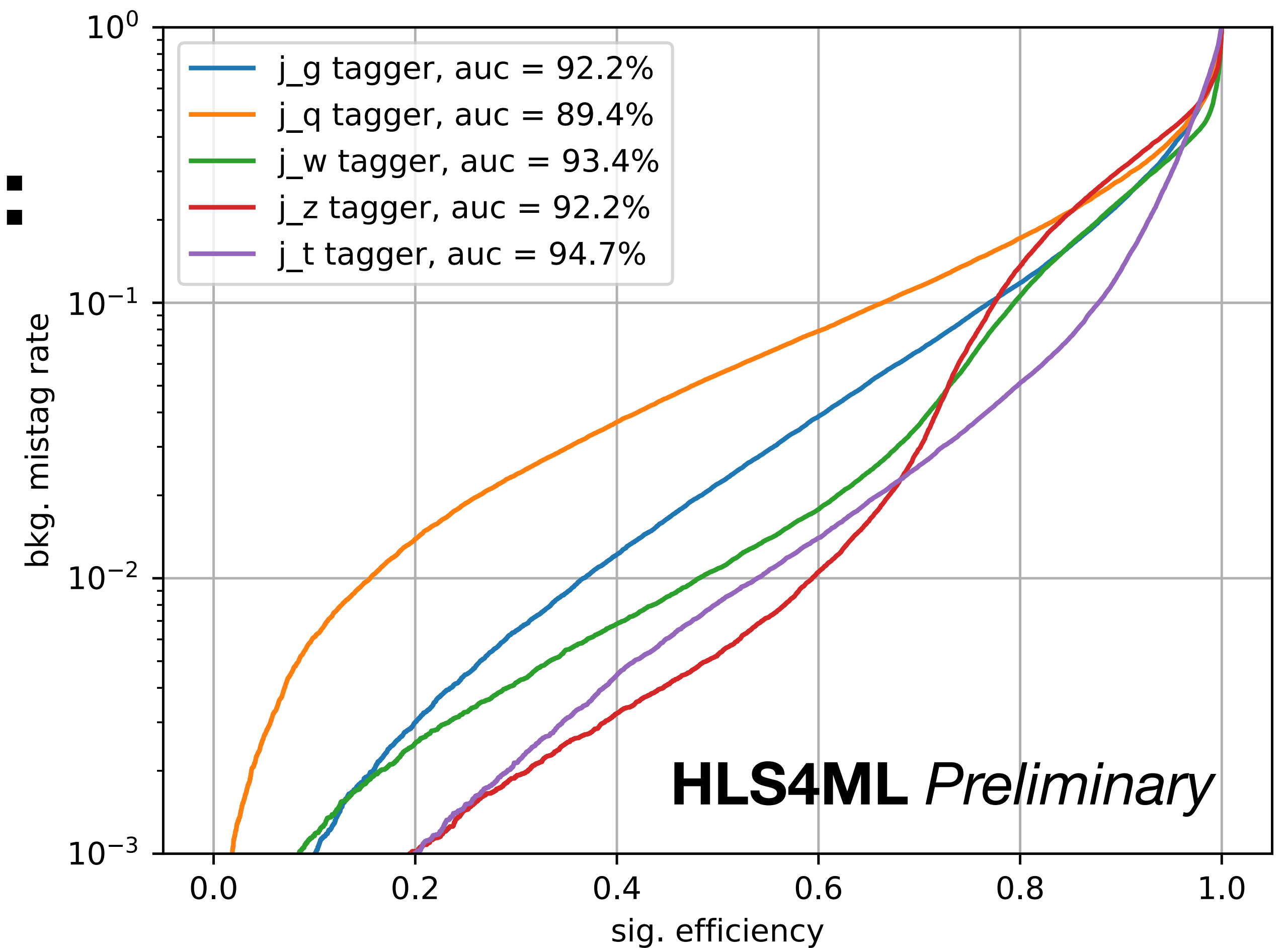
Compression (50%) + reuse = 2:

29% DSP usage for Virtex 7

60 clocks (300 ns), Pipeline = 2



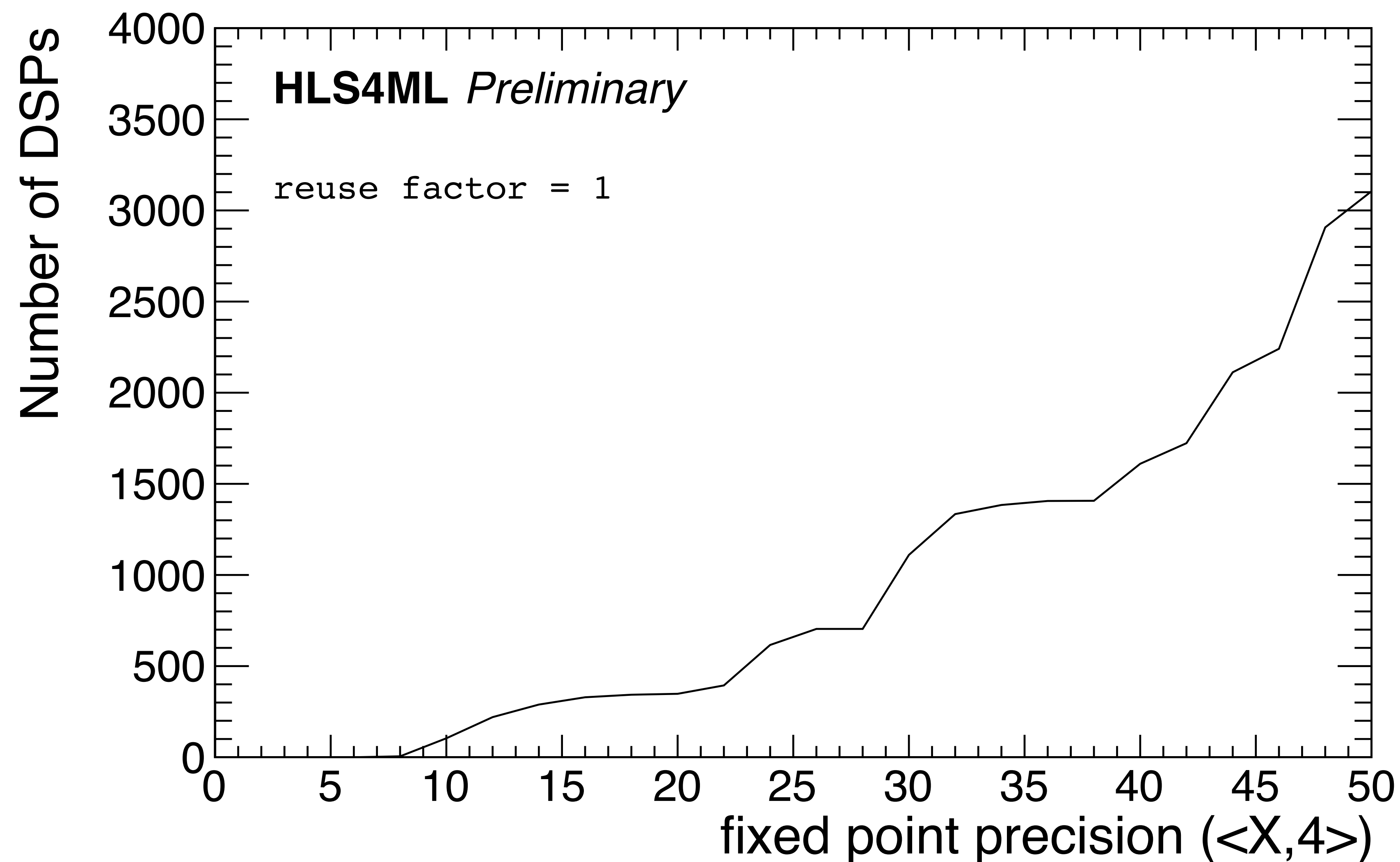
Fully connected deep neural network

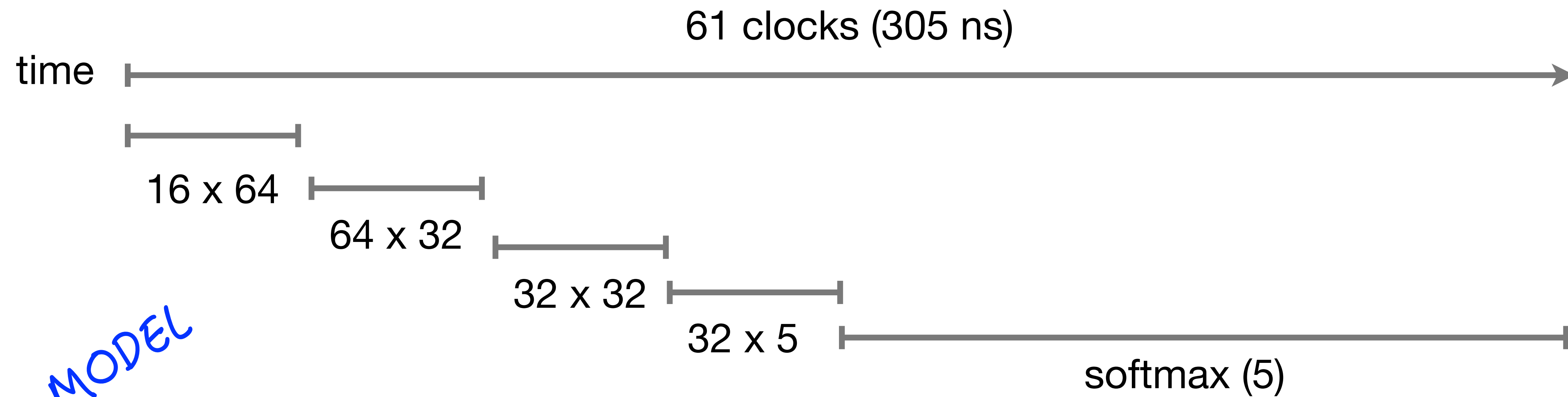


Take a simple 1-layer network and scan in input/weight precision

Reduced precision can greatly reduce resource usage

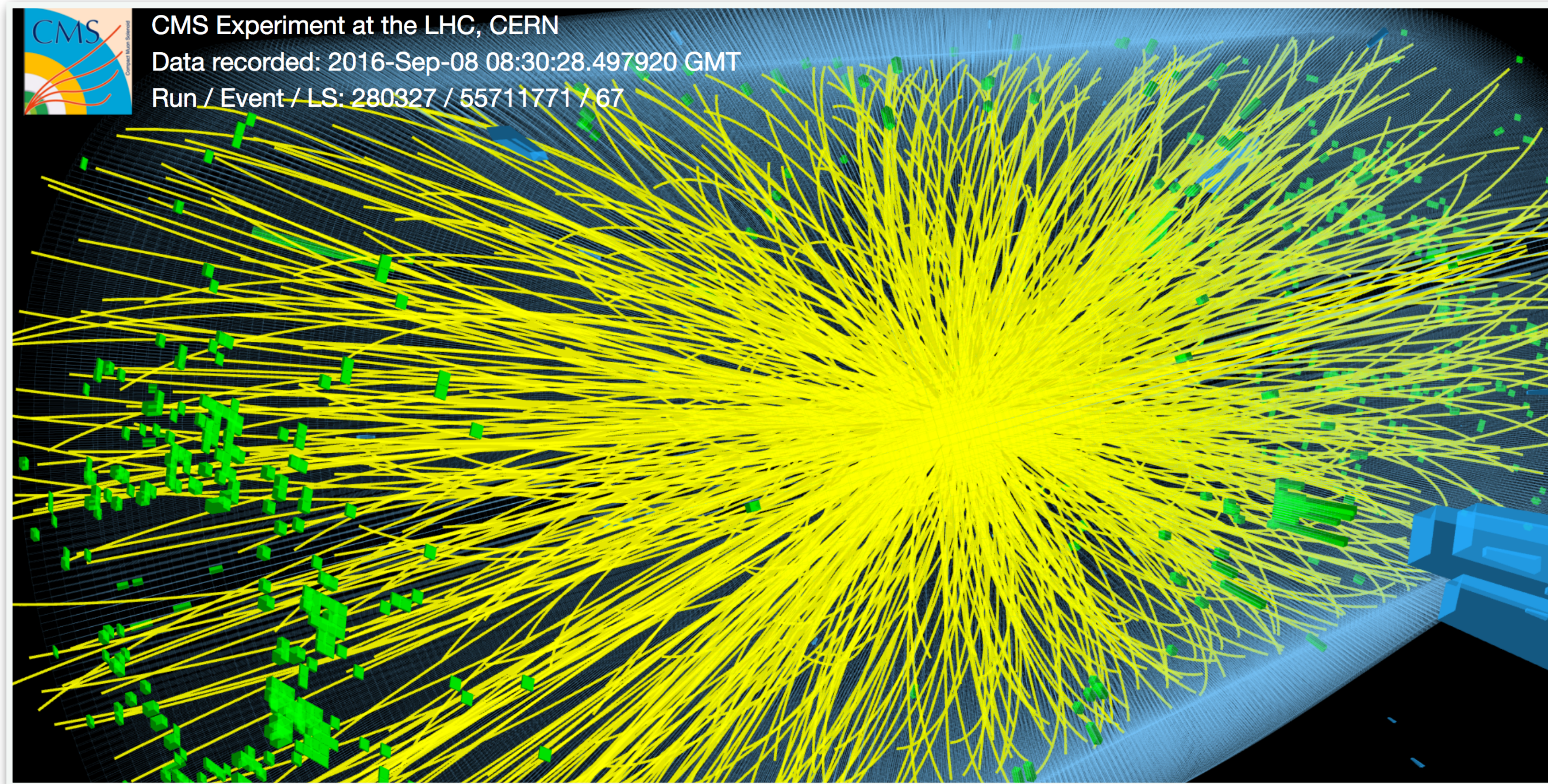
e.g. factor of 4 reduction with 18 instead of 32 bits with minimal loss in performance





ORIGINAL MODEL

	BRAM	DSP	FF	LUT
Total	3	3329	95924	8127
% Usage	~0	92	11	18

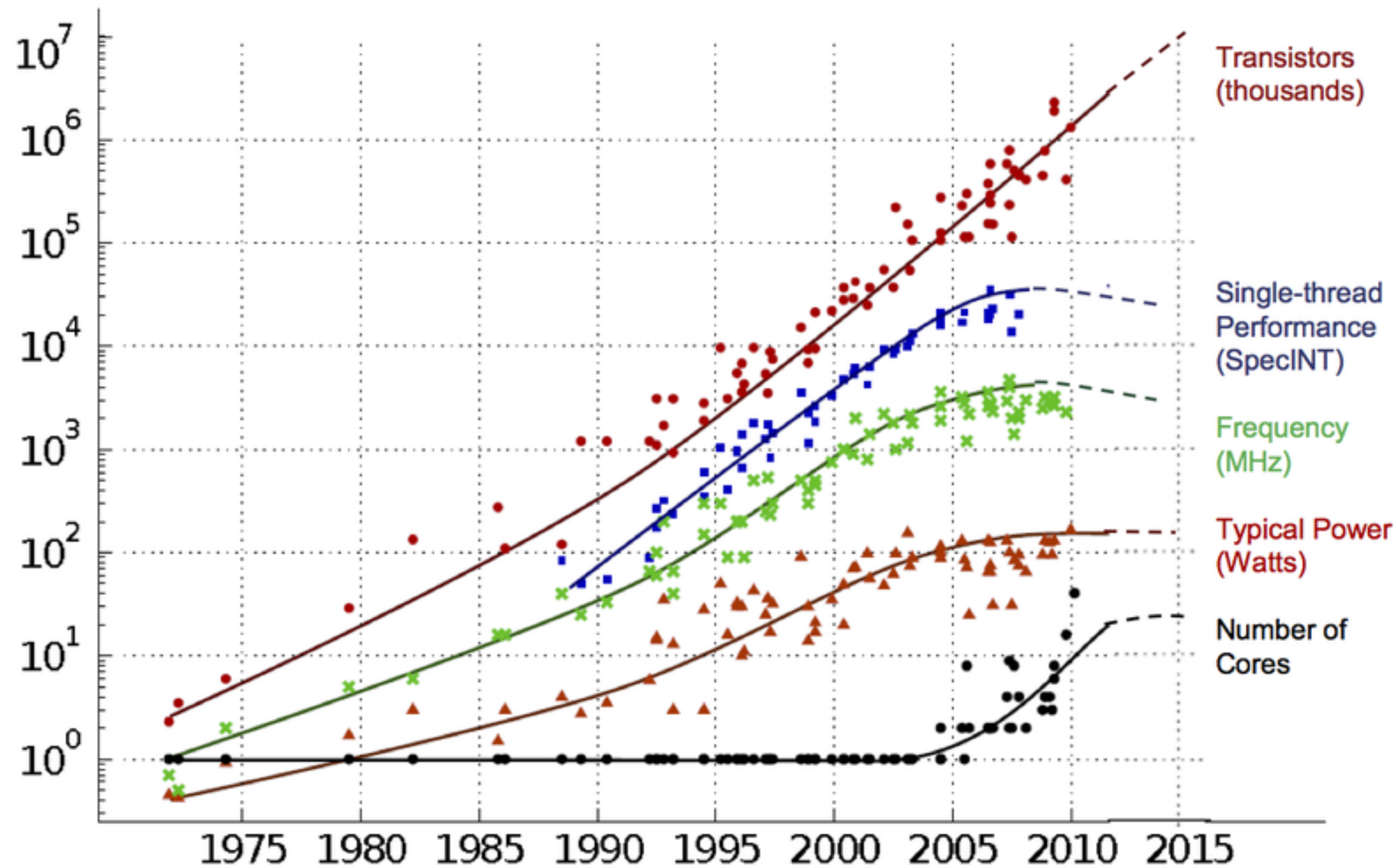


**Current:
~5 minutes per
HL-LHC event**

**100 times the
data...
exabytes!**

Major HLT and computing challenges going forward!

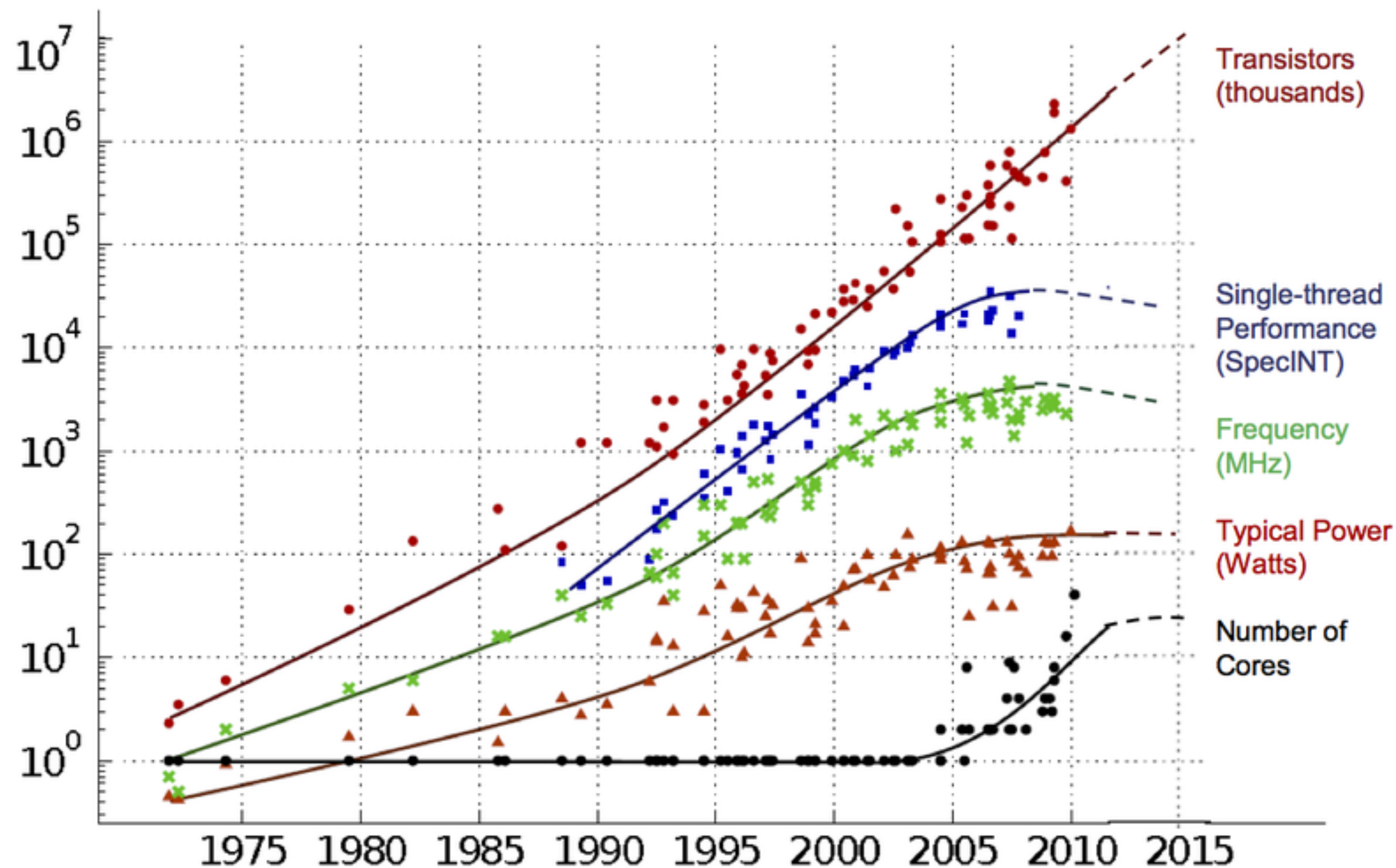
MOORE'S LAW AND DENNARD SCALING



Moore's Law continues

Dennard Scaling fails

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



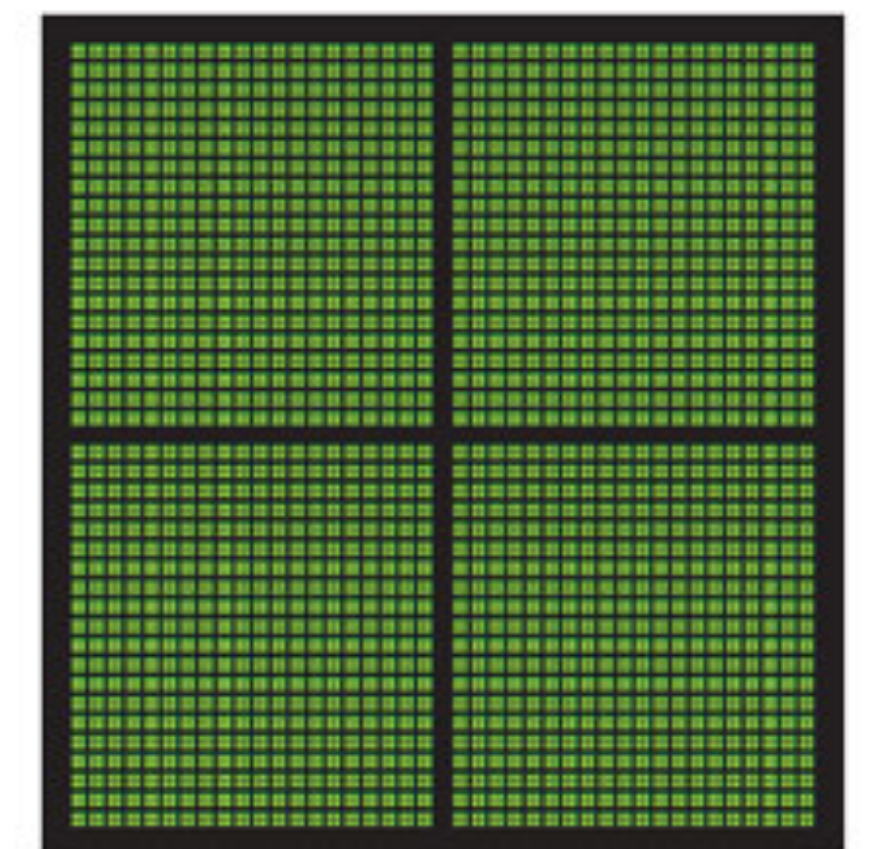
Moore's Law continues

Dennard Scaling fails

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



CPU
MULTIPLE CORES

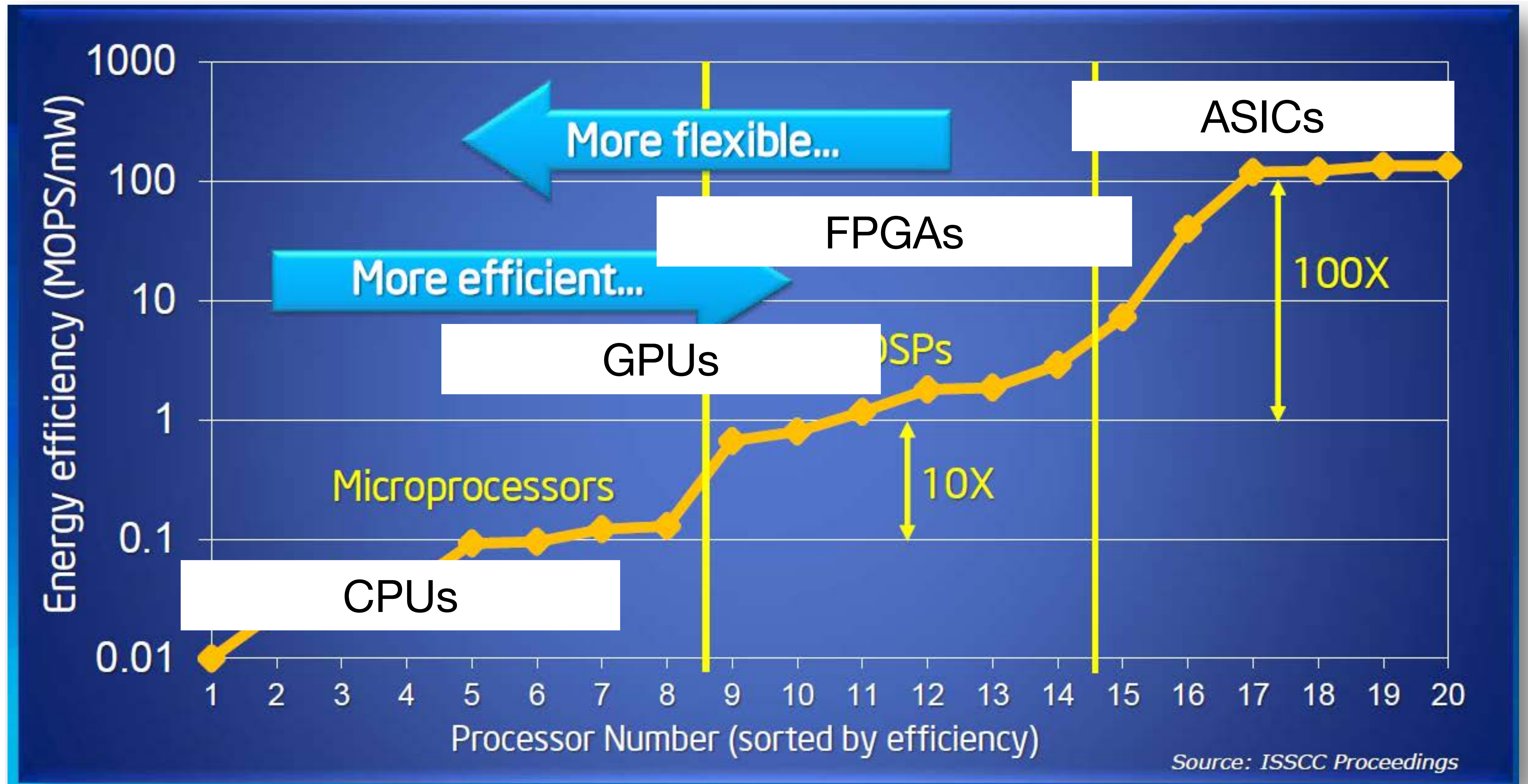


GPU
THOUSANDS OF CORES

Single threaded performance not improving

Circa ~2005: "The Era of Multicore"

→ Today: Transition to the "Era of Specialization"? (c.f. Doug Burger)



Source: Bob Broderson, Berkeley Wireless group

- * GPUs still best option for training
- * FPGAs generally much more power efficient