

# Neural networks in FPGAs for trigger and DAQ

Jennifer Ngadiuba, Maurizio Pierini [CERN]

Javier Duarte, Sergo Jindariani, Ben Kreis, Ryan Rivera, *Nhan Tran* [Fermilab]

Edward Kreinar [Hawkeye 360]

Phil Harris [MIT]

Song Han [Stanford University/Google]

Zhenbin Wu [University of Illinois at Chicago]

Connecting The Dots

March 20, 2018

## motivation

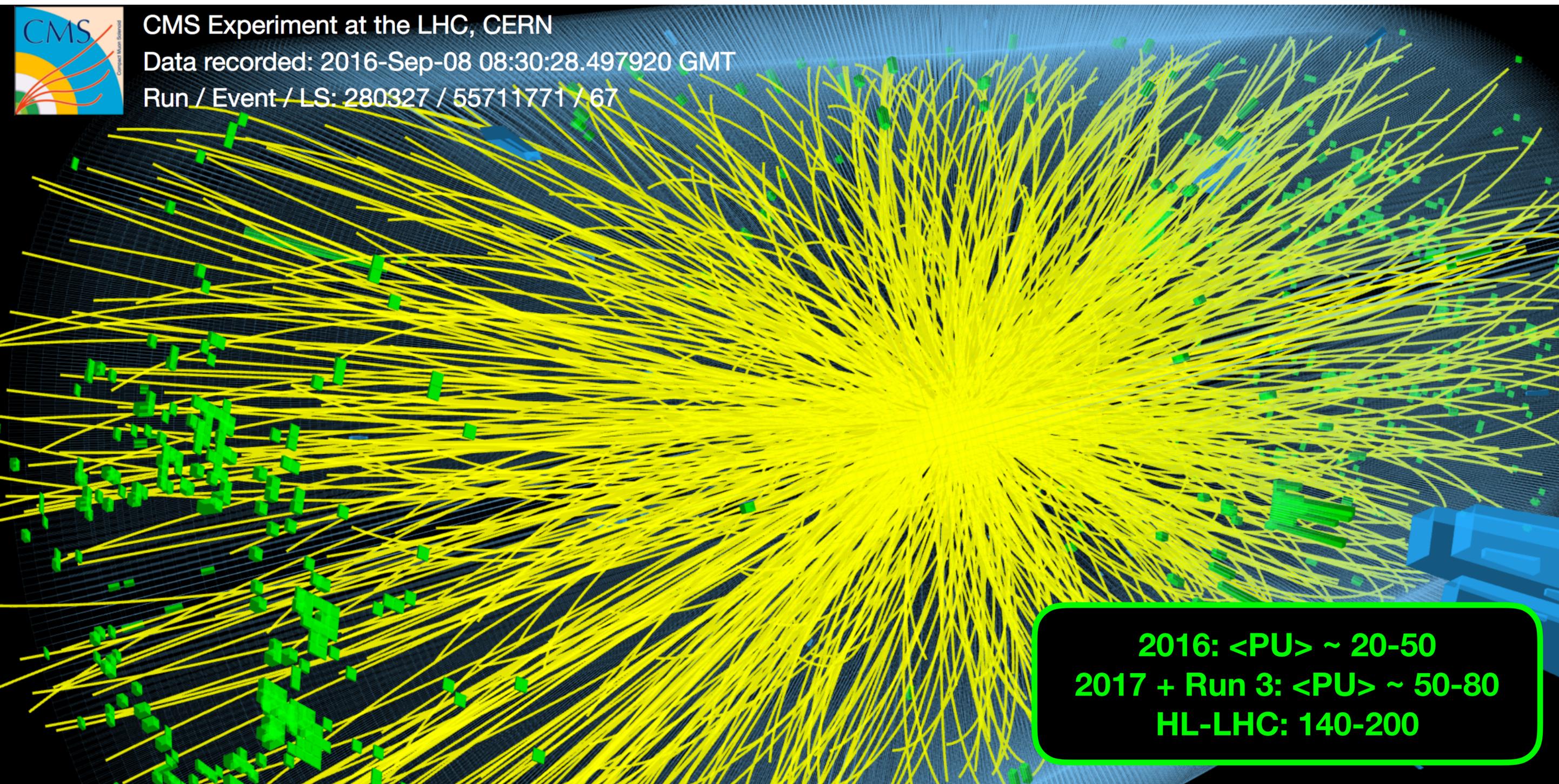
### hls4ml

case study: jet substructure  
neural networks with HLS  
results

## status and outlook

**motivation**

**A CHALLENGE TO MAINTAIN PHYSICS IN  
INCREASINGLY COMPLEX COLLISION ENVIRONMENTS**  
detector design, object performance, physics sensitivity,  
**lose untriggered events forever**



**A CHALLENGE TO MAINTAIN PHYSICS IN  
INCREASINGLY COMPLEX COLLISION ENVIRONMENTS**  
detector design, object performance, physics sensitivity,  
**lose untriggered events forever**

**Need sophisticated techniques to preserve the physics!**

e.g. Dedicated advanced reconstruction algorithms,  
**Particle Flow and PUPPI**, and **machine learning algorithms**

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

**A CHALLENGE TO MAINTAIN PHYSICS IN  
INCREASINGLY COMPLEX COLLISION ENVIRONMENTS**  
detector design, object performance, physics sensitivity,  
**lose untriggered events forever**

**Need sophisticated techniques to preserve the physics!**

e.g. Dedicated advanced reconstruction algorithms,  
**Particle Flow and PUPPI**, and **machine learning algorithms**

**See poster by Ben Kreis on  
PF + PUPPI in the trigger**

**this talk**

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

## Goal:

reduce event rate from 40MHz to 500 Hz

**How:** multi-tier system

### custom hardware (“L1”)

latency,  $O(\mu\text{s})$

rate in/out: 40 MHz / 100 KHz

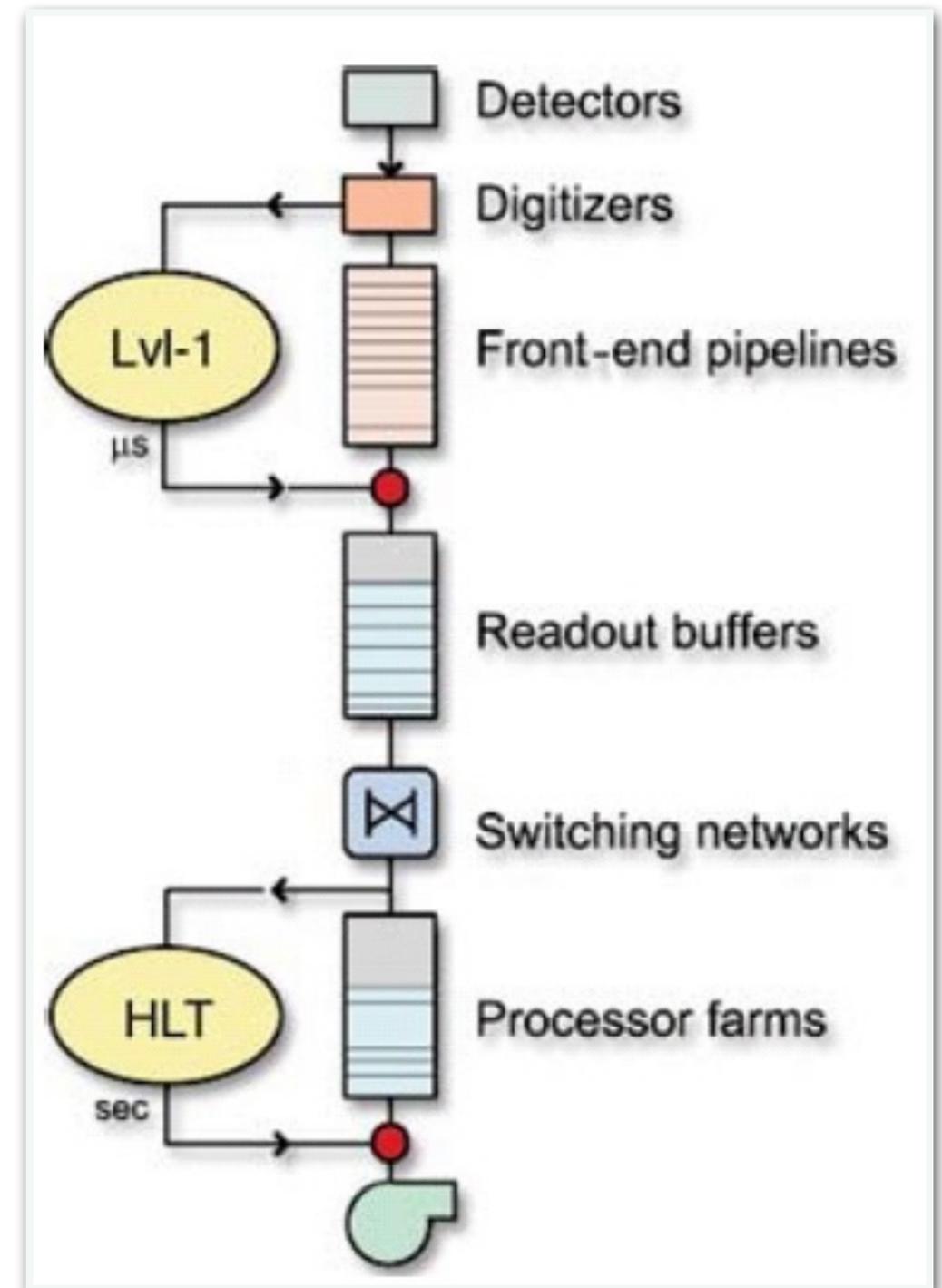
### computing farm (“HLT”)

latency,  $O(100\text{ ms})$

rate in/out: 100 KHz / 500 Hz

n.b. all numbers approximate

For HL-LHC upgrade: latency and output rates go up  $\sim 5$



# THE TRIGGER @ ATLAS AND CMS

## Goal:

reduce event rate from 40MHz to 500 Hz

## How: multi-tier system

### custom hardware (“L1”)

latency,  $O(\mu\text{s})$

rate in/out: 40 MHz / 100 KHz

### computing farm (“HLT”)

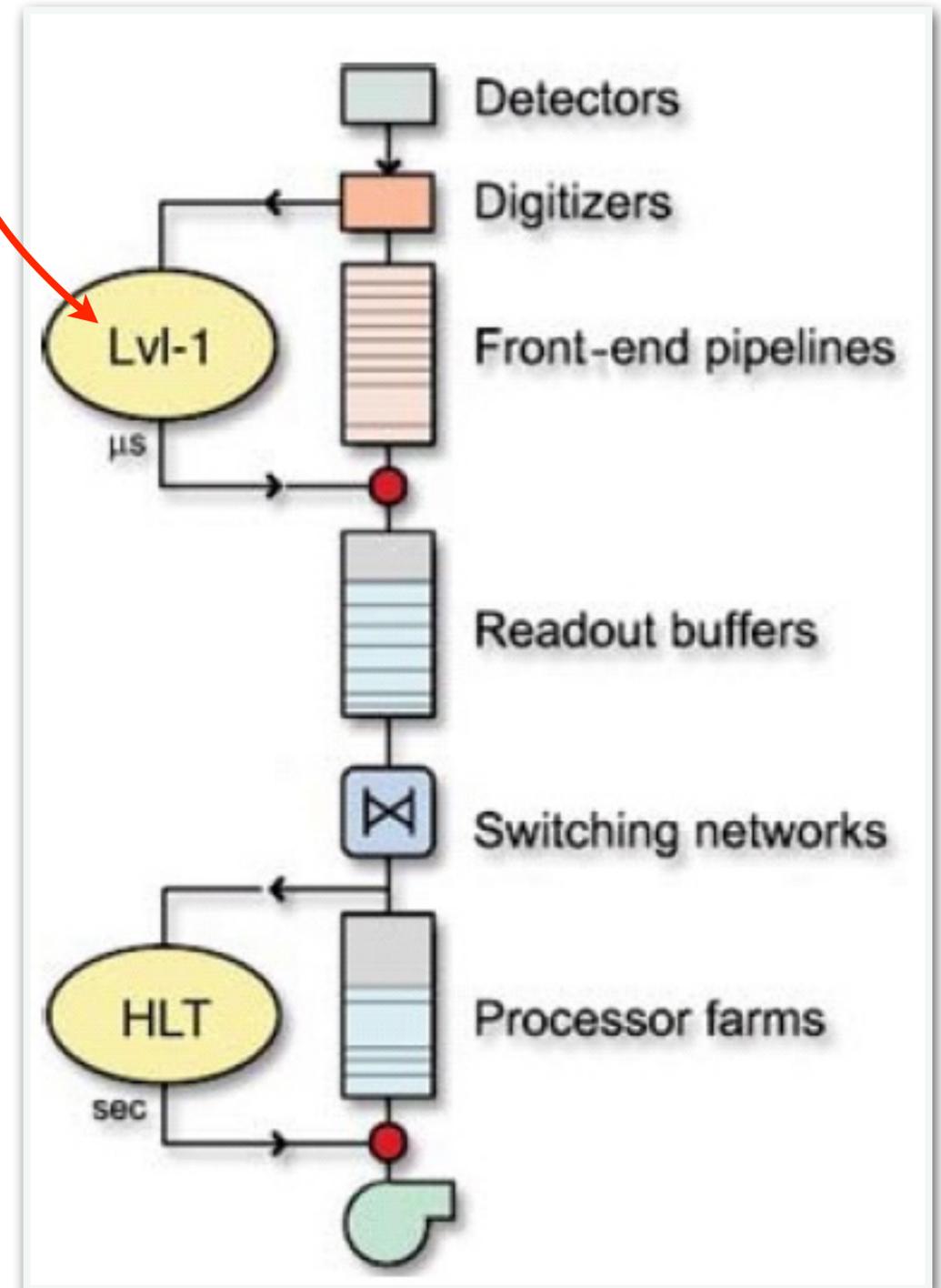
latency,  $O(100 \text{ ms})$

rate in/out: 100 KHz / 500 Hz

n.b. all numbers approximate

For HL-LHC upgrade: latency and output rates go up  $\sim 5$

**Latencies necessitate all-FPGA design!**

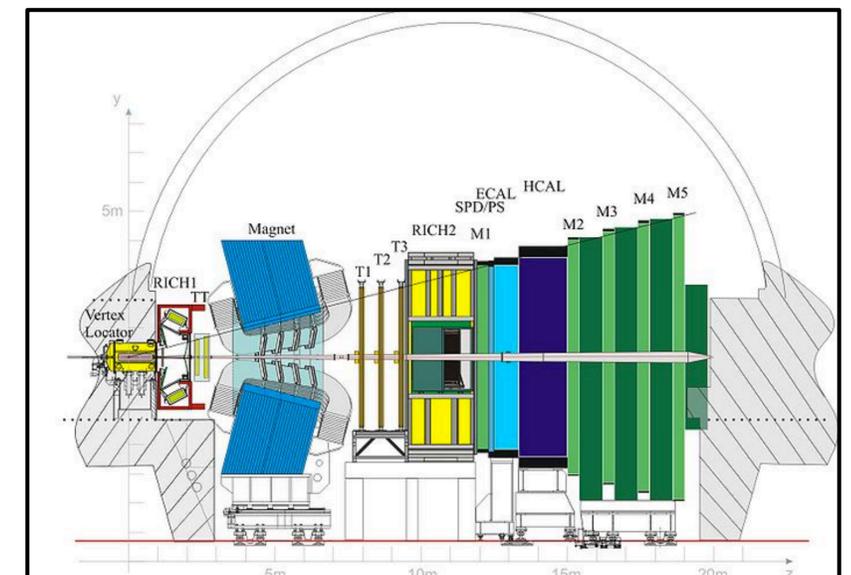
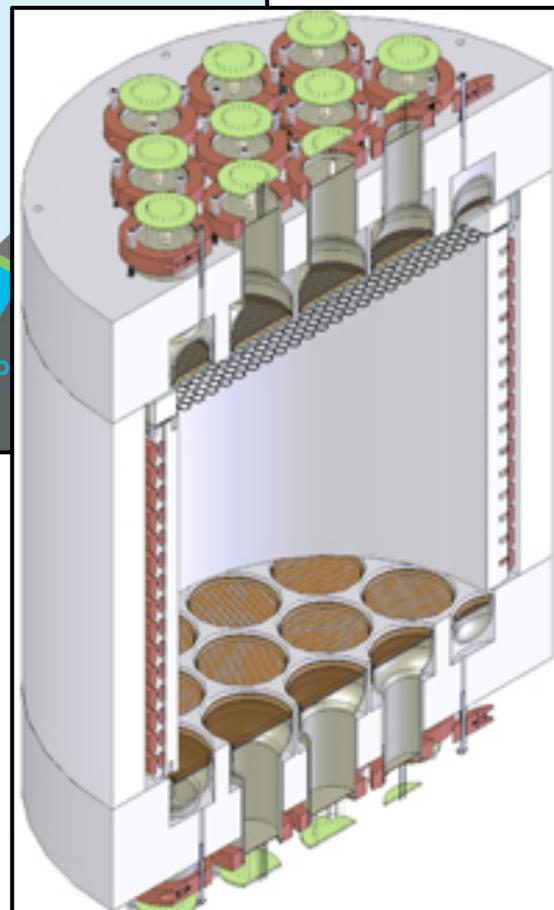
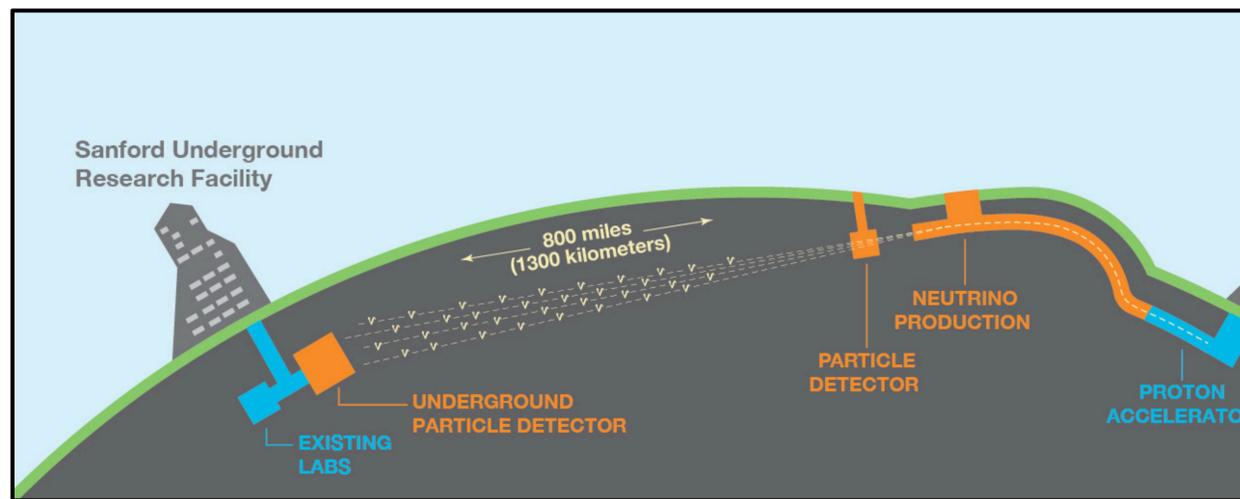


# MORE OPPORTUNITIES

In the era of big science, more sophisticated triggers and DAQ systems are required

Even in traditional “low” rate experiments

Other LHC applications, like LHCb, and ATLAS/CMS HLT and cosmic and intensity frontier experiments



## Machine learning algorithms are ubiquitous in HEP

### FPGA usage broad across HEP experiments

Centered on DAQ and trigger development

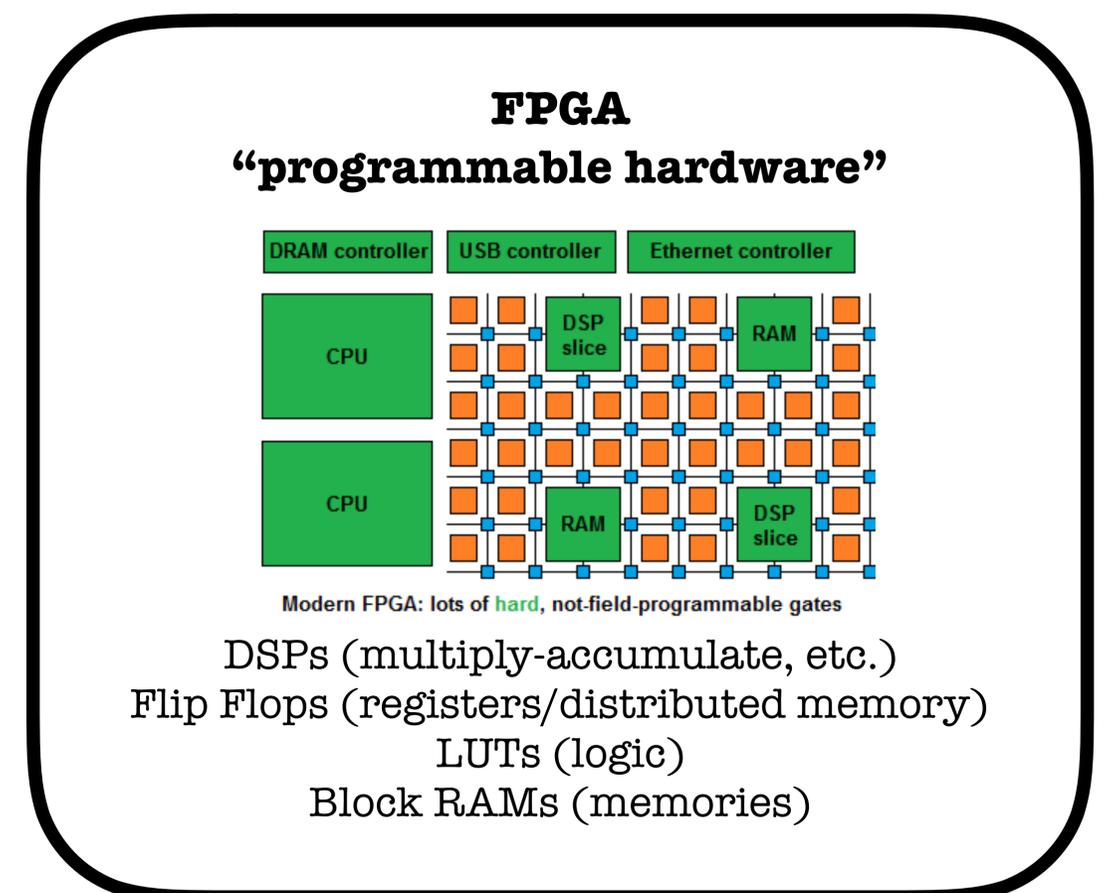
Some early adaptations of ML techniques in trigger [1]

FPGA development becoming more accessible

### High Level Synthesis, OpenCL

### FPGA interest in industry is growing

Programmable hardware with structures that maps nicely onto ML architectures



[1] Carnes et al., <https://indico.cern.ch/event/567550/contributions/2629686/>

high level synthesis for machine learning

~~hlsfml~~

# **hls4ml**

neural networks in HLS

case study: jet substructure

results

1st application:

**Focus on L1 trigger**

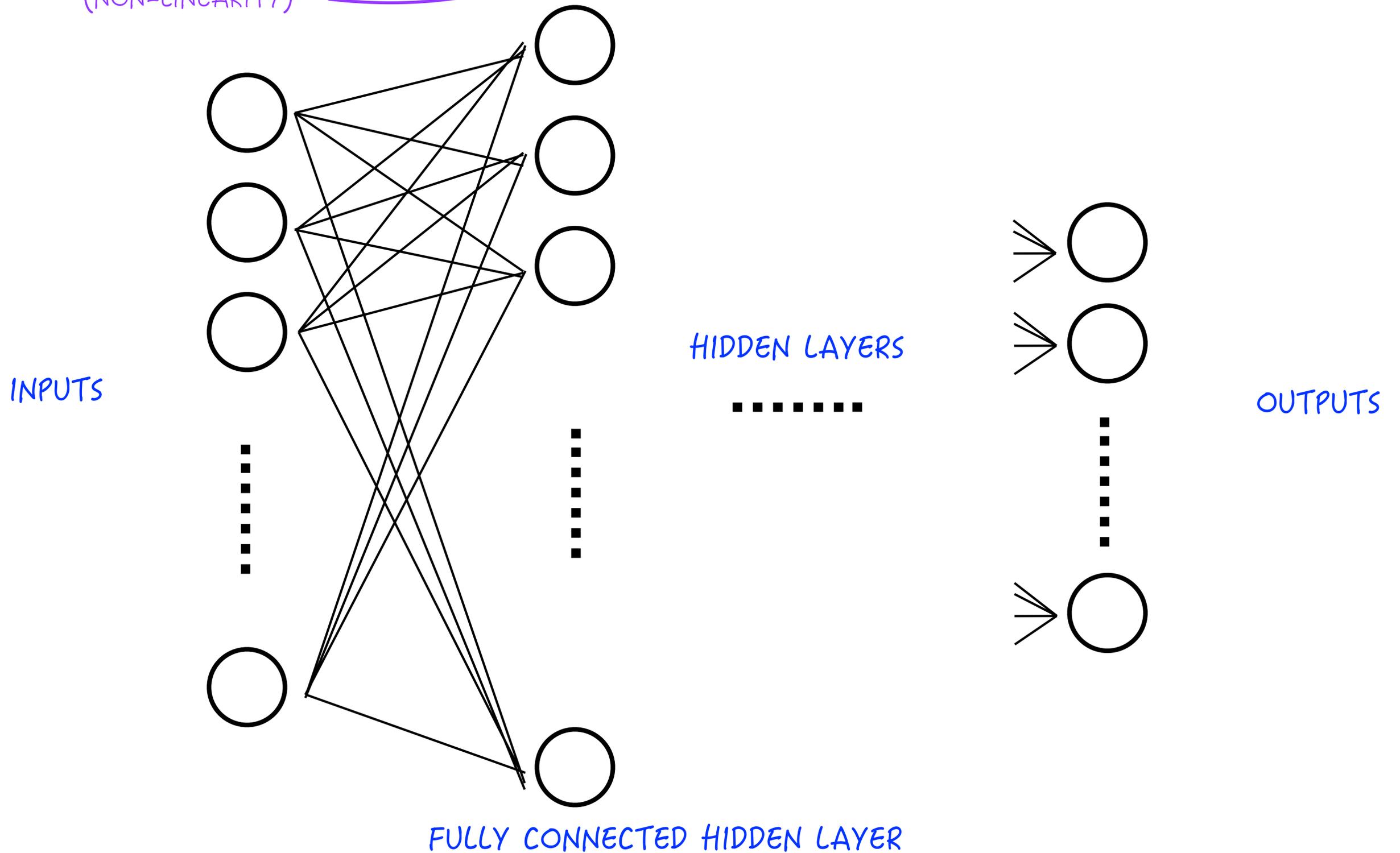
**What can we do in  $< \mu\text{s}$  on one FPGA?**

Discussion on other applications at the end

# NN INFERENCE IN A NUTSHELL

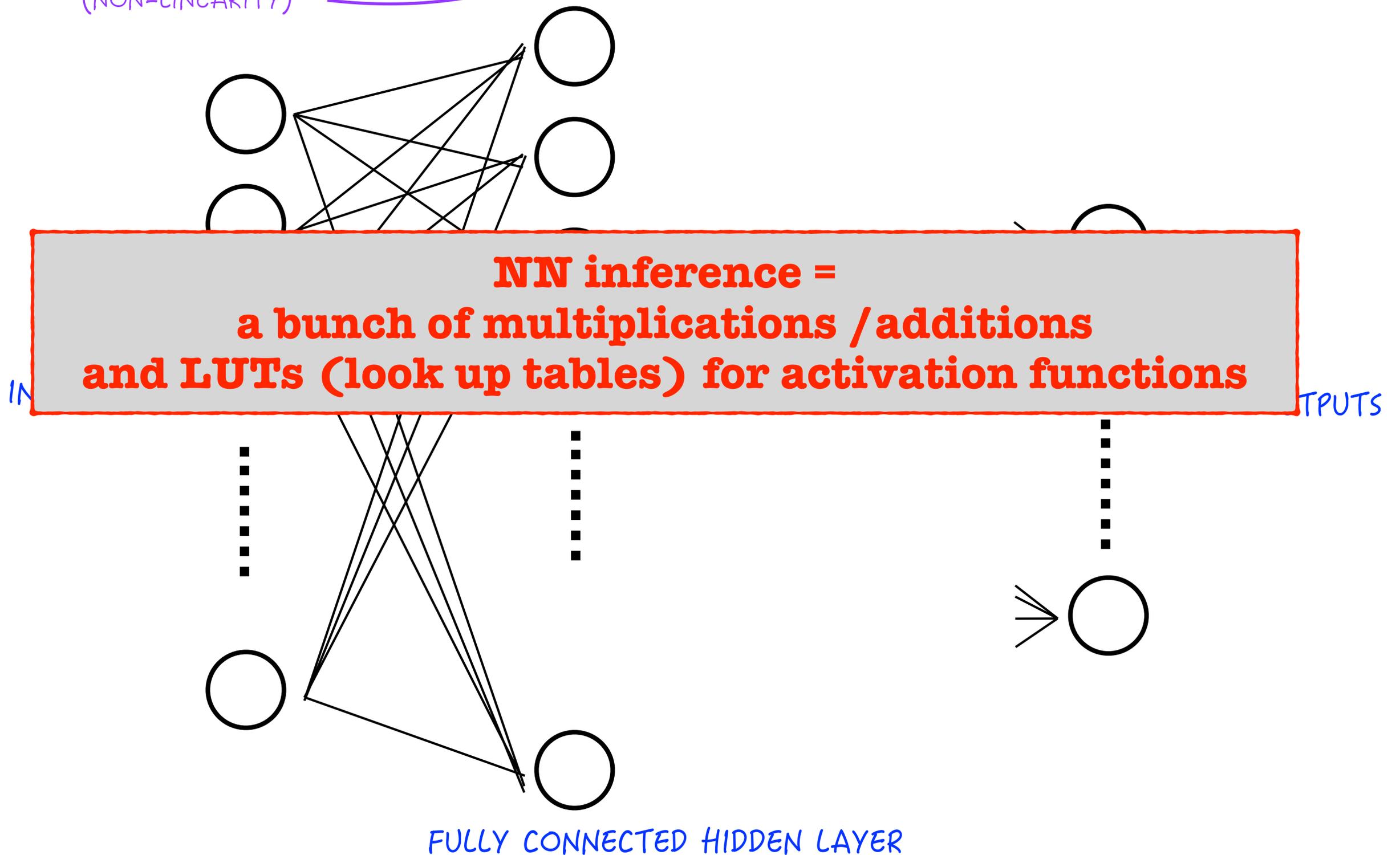
$$\vec{O}_j = \Phi(\vec{l}_i \times \vec{W}_{ij} + \vec{b}_j)$$

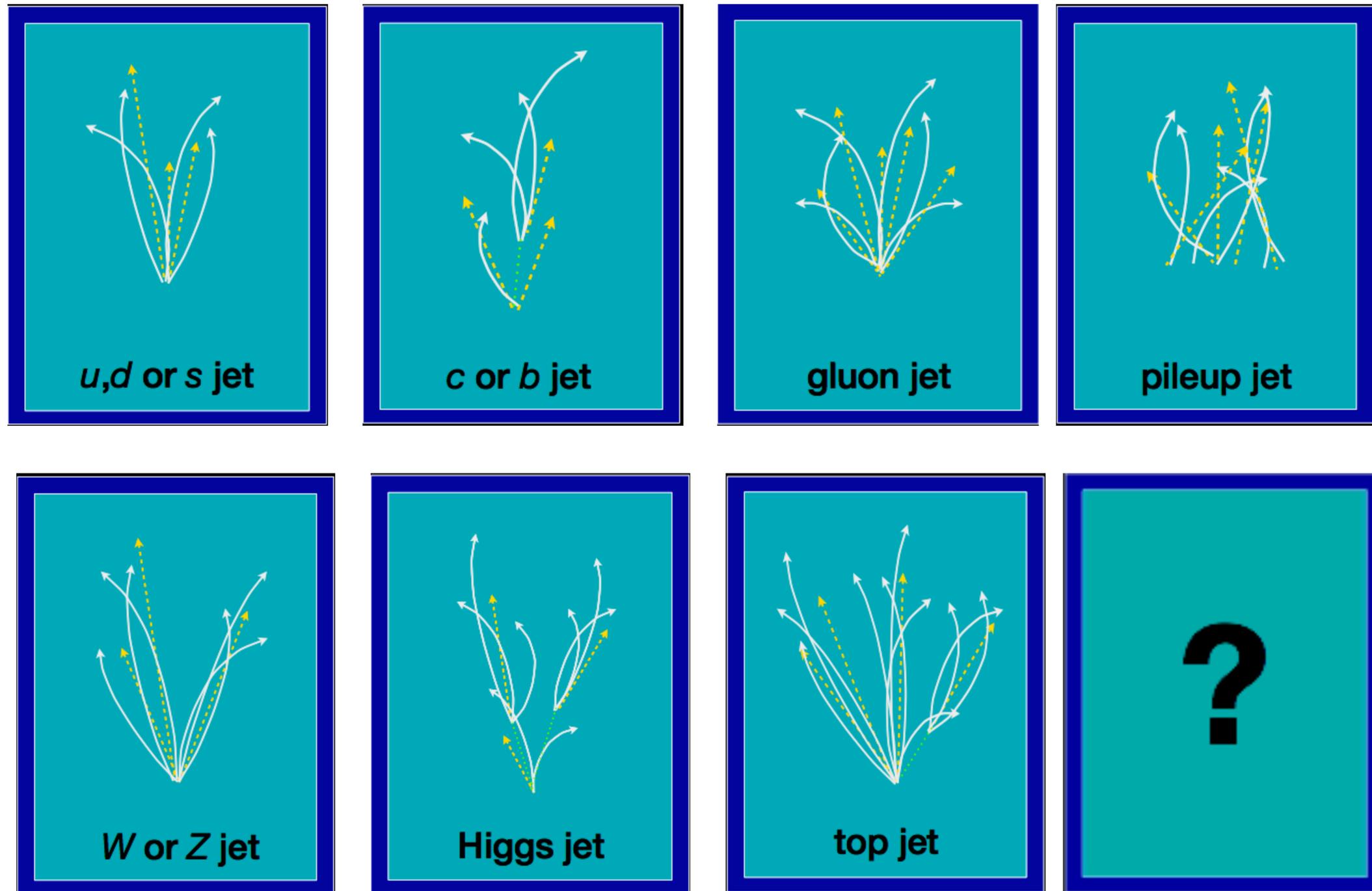
$\Phi$  = ACTIVATION FUNCTION  
(NON-LINEARITY)



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

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





**Just an illustrative example, lessons are generic!**

Might not be the best application, but a familiar one  
ML in substructure is well-studied

5 output multi-classifier:

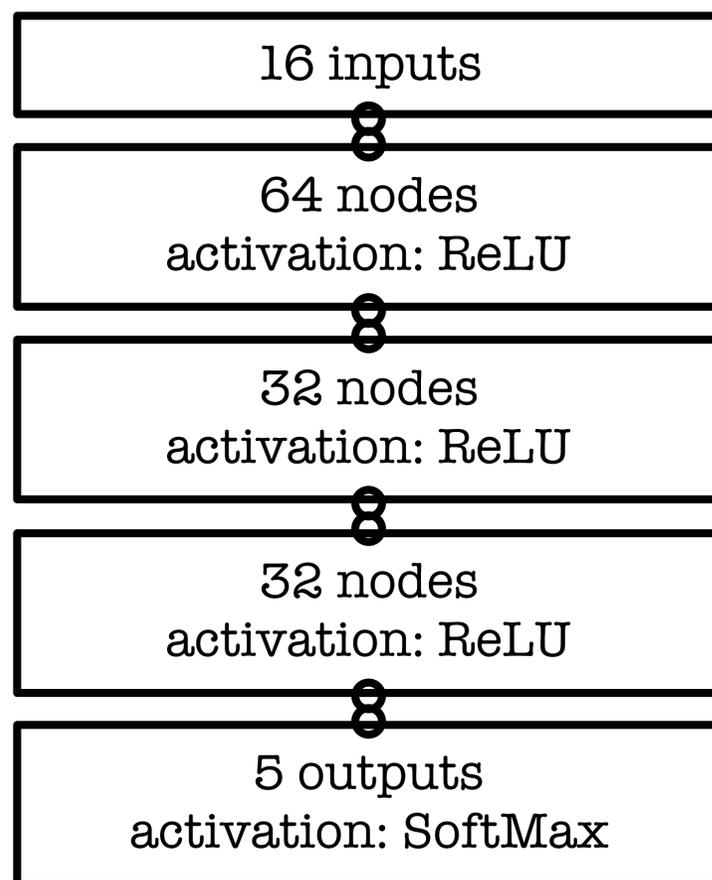
Does a jet originate from a quark, gluon, W/Z boson, top quark?

## Network architecture

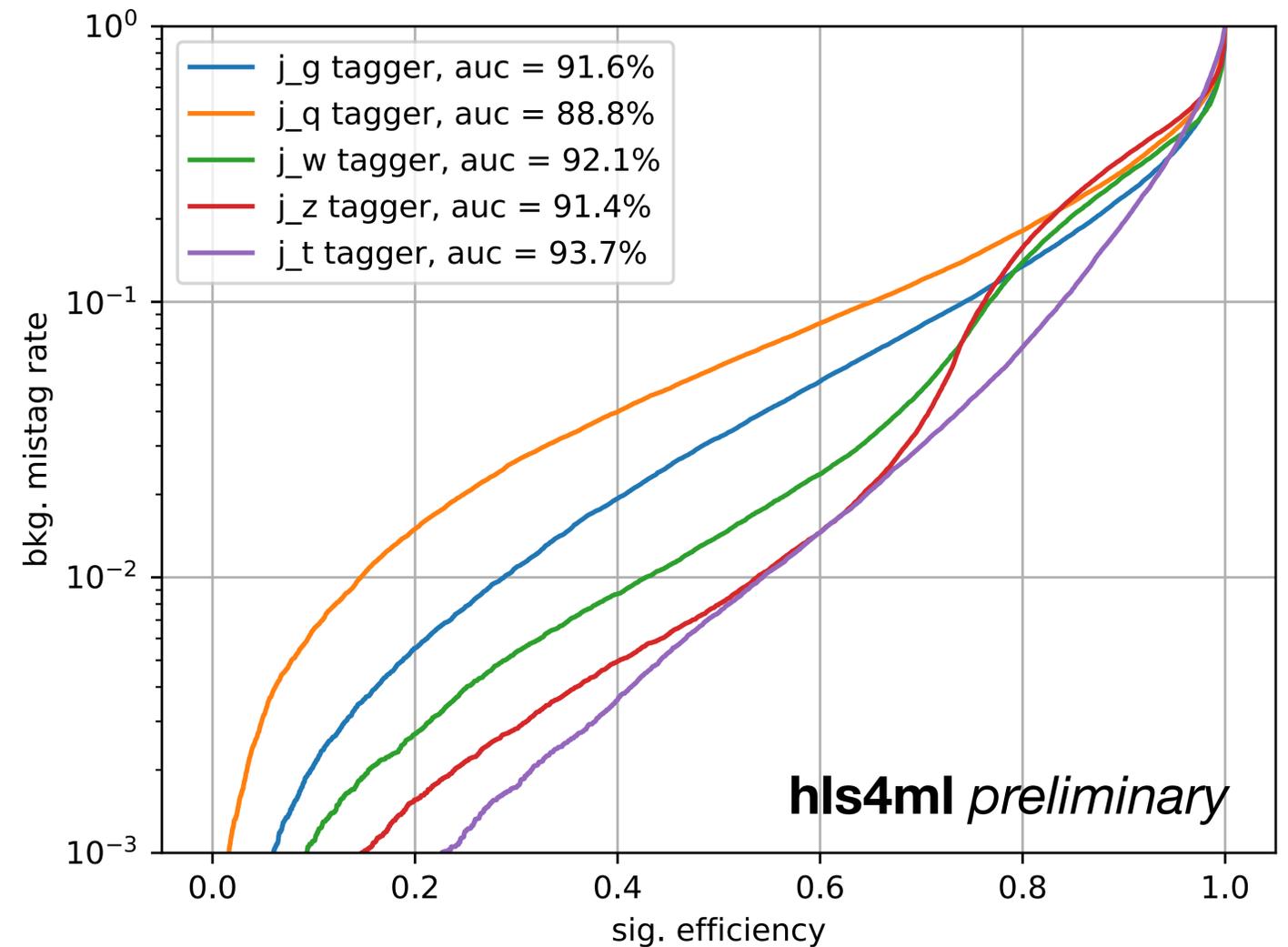
16 expert inputs

jet masses, multiplicity

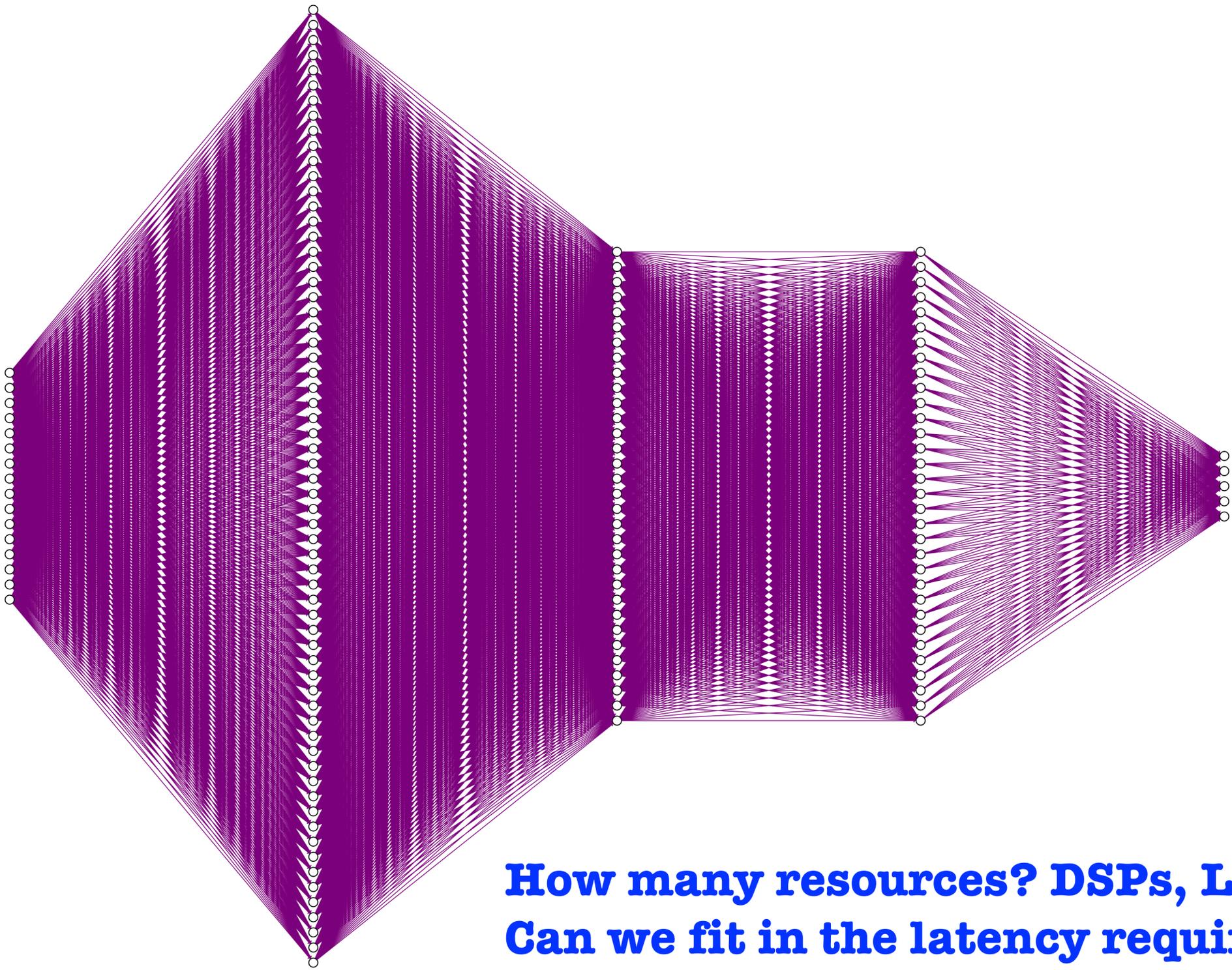
ECFs ( $\beta=0,1,2$ )



*Fully connected deep  
neural network*



## FPGA



**How many resources? DSPs, LUTs, FFs?  
Can we fit in the latency requirements?**

Emergent engineering field, efficient implementation of NN architecture

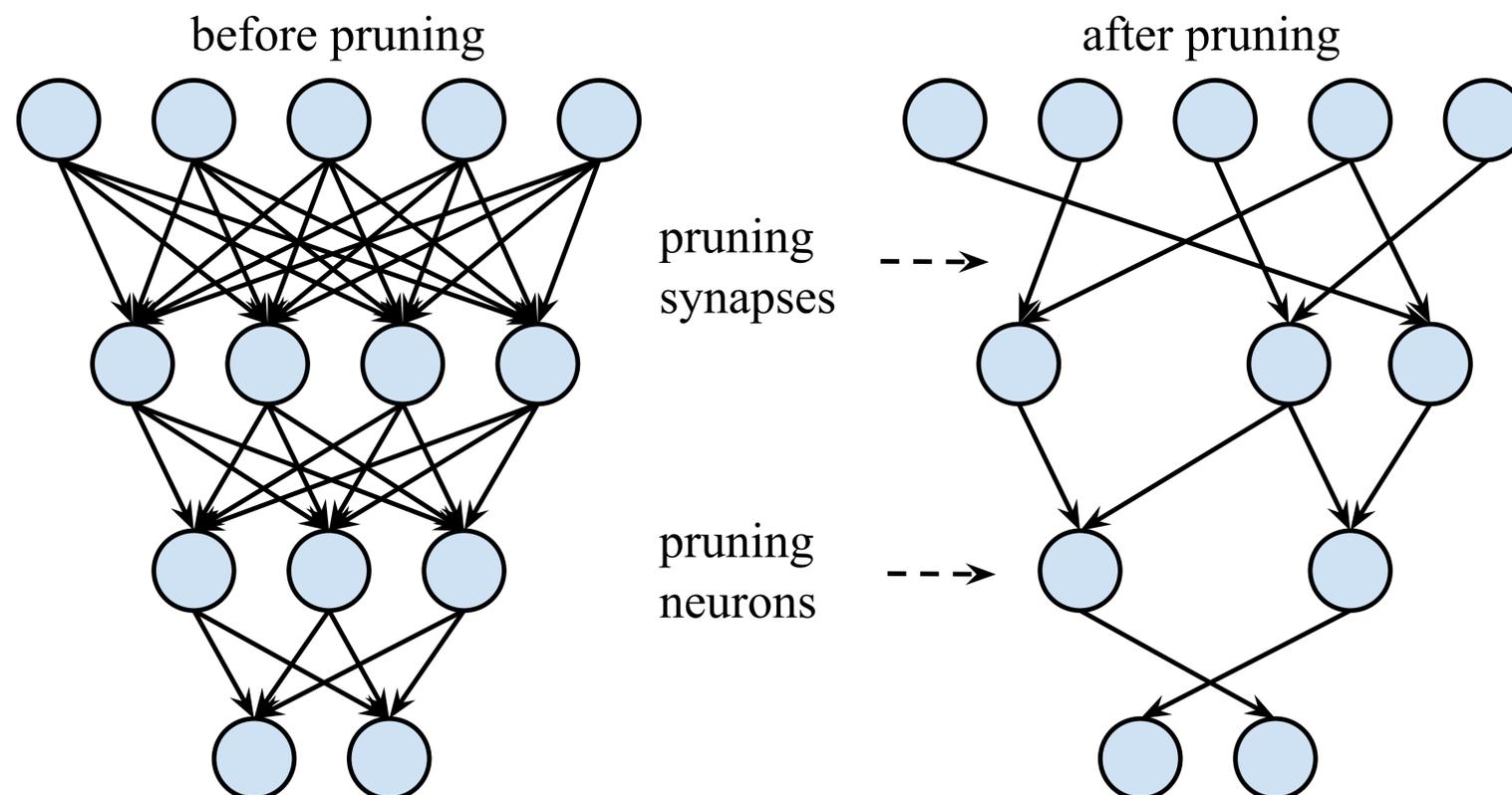
## Compression:

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

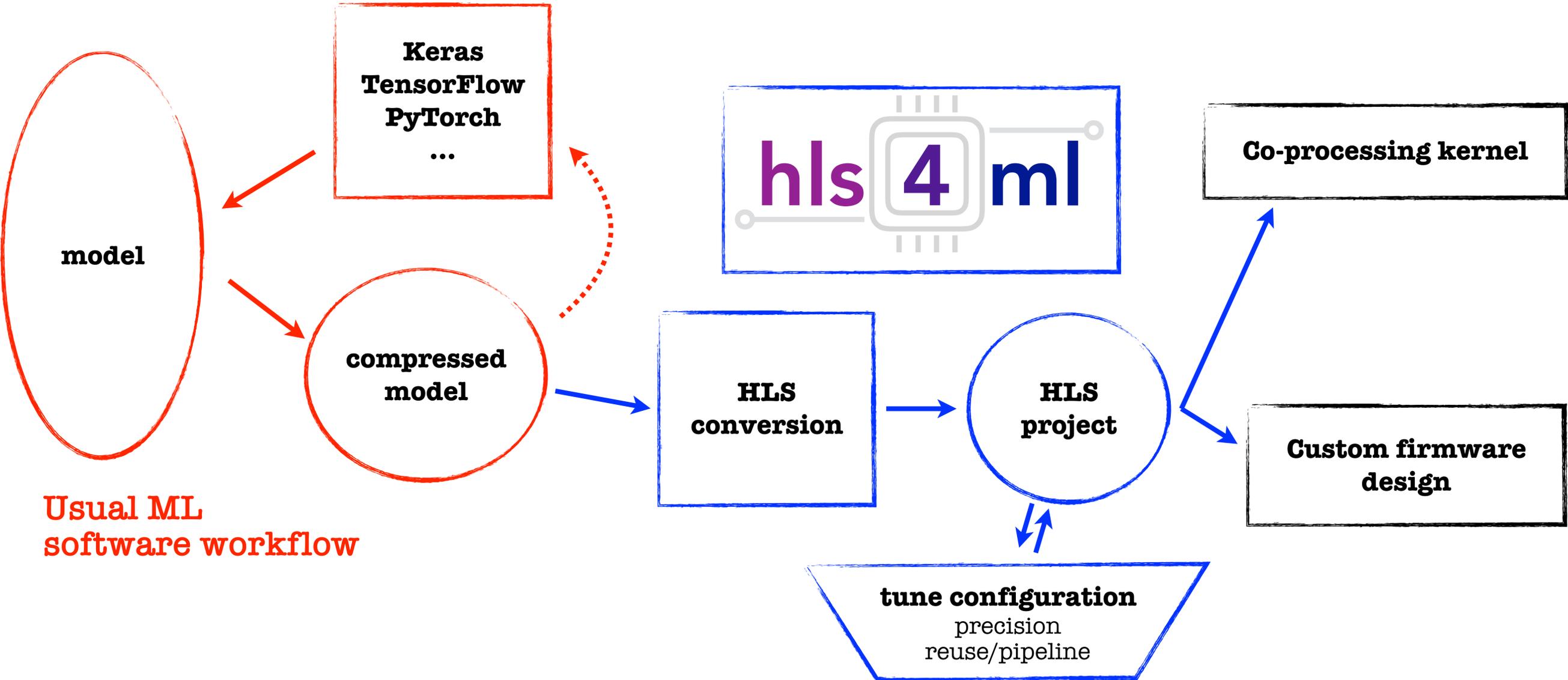
## Quantization:

32-bit floating point math is overkill

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



# PROJECT OVERVIEW



Usual ML software workflow



*Translation in one line!*

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

```
KerasJson: example-keras-model-files/KERAS_1layer.json  
KerasH5:   example-keras-model-files/KERAS_1layer_weights.h5  
OutputDir: my-hls-test  
ProjectName: myproject  
XilinxPart: xc7vx690tffg1927-2  
ClockPeriod: 5  
  
IOType: io_parallel # options: io_serial/io_parallel  
ReuseFactor: 1  
DefaultPrecision: ap_fixed<18,8>
```

} Keras inputs

**IOType:** parallelize or serialize

**ReuseFactor:** how much to parallelize

**DefaultPrecision:** inputs, weights, biases



*Translation in one line!*

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

```
KerasJson: example-keras-model-files/KERAS_1layer.json
KerasH5:   example-keras-model-files/KERAS_1layer_weights.h5
OutputDir: my-hls-test
ProjectName: myproject
XilinxPart: xc7vx690tffg1927-2
ClockPeriod: 5

IOType: io_parallel # options: io_serial/io_parallel
ReuseFactor: 1
DefaultPrecision: ap_fixed<18,8>
```

} Keras inputs

**IOType:** parallelize or serialize

**ReuseFactor:** how much to par

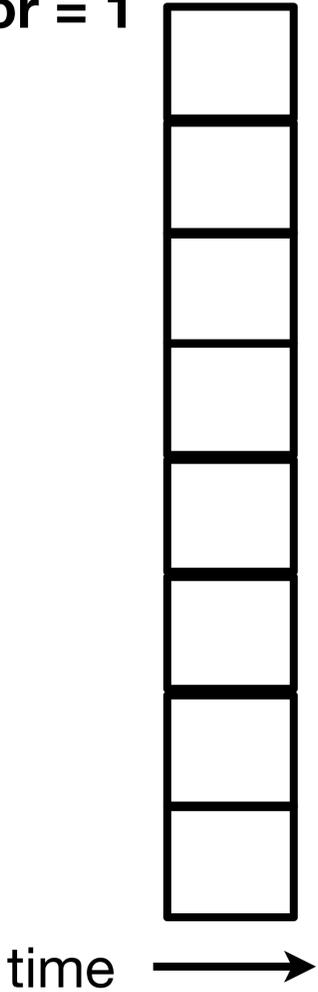
**DefaultPrecision:** inputs, weig

**FIRMWARE IN MINUTES,  
FAST DEVELOPMENT CYCLES!**

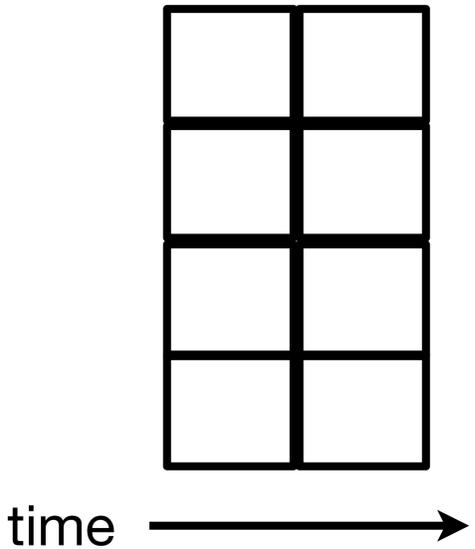
**ReuseFactor:** how much to parallelize operations a hidden layer

# of multiplications per clock (DSPs usage)

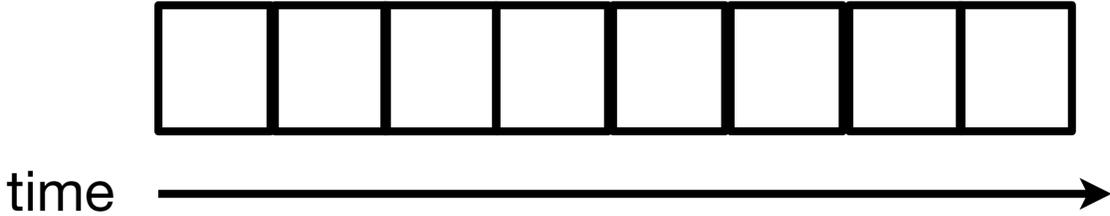
parallel  
8 multipliers in 1 clock  
reuse factor = 1



parallel  
4 multipliers in 2 clocks  
reuse factor = 2



serial  
1 multiplier in 8 clocks



less resources/less throughput



defines the **Pipeline Interval** = when new inputs are introduced to algo

**Xilinx Vivado 2017.2**

**Clock frequency: 200 MHz**

**FPGA: Xilinx Kintex Ultrascale (XCKU115-FLVB2104)**

*[We'll come back to these later]*

## **A note on inputs:**

We assume network inputs have already been computed

*Not a good assumption in the jet substructure case*

We ignore this here, but the choice to use “raw” or “expert” features will depend on the problem

**On the coming plots: resource usage comes from HLS estimates**  
discussion on differences w.r.t. implementation later

There are many schemes for compression

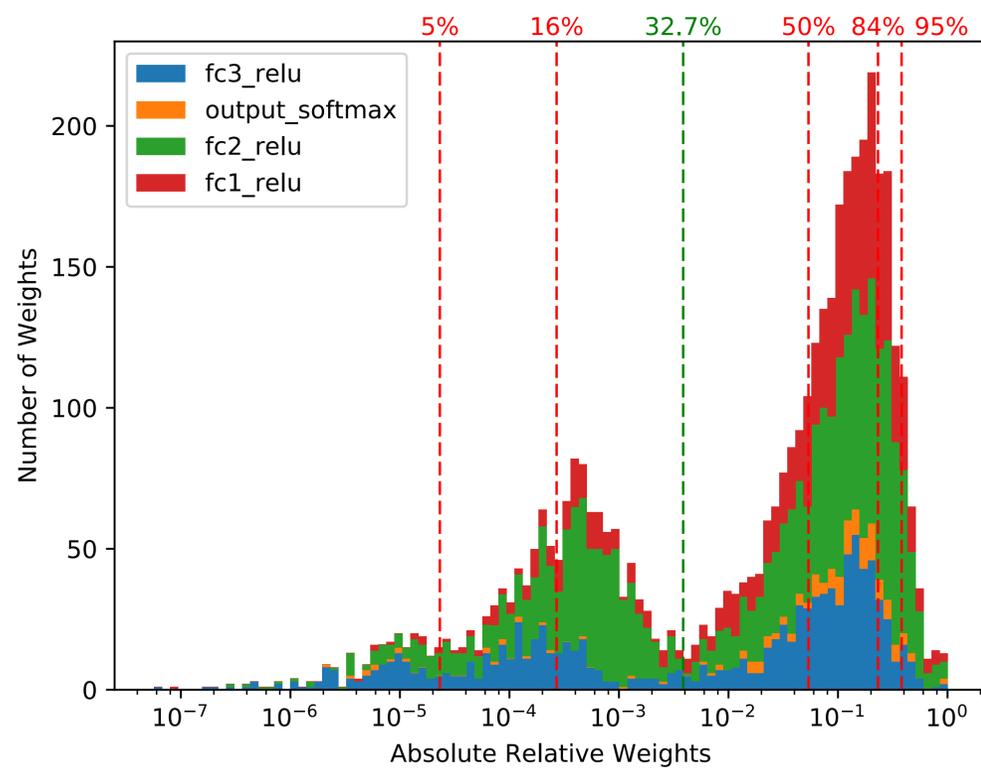
We do a simplistic, iterative version

Training with “L<sub>1</sub>” regularization, down-weight unimportant synapses

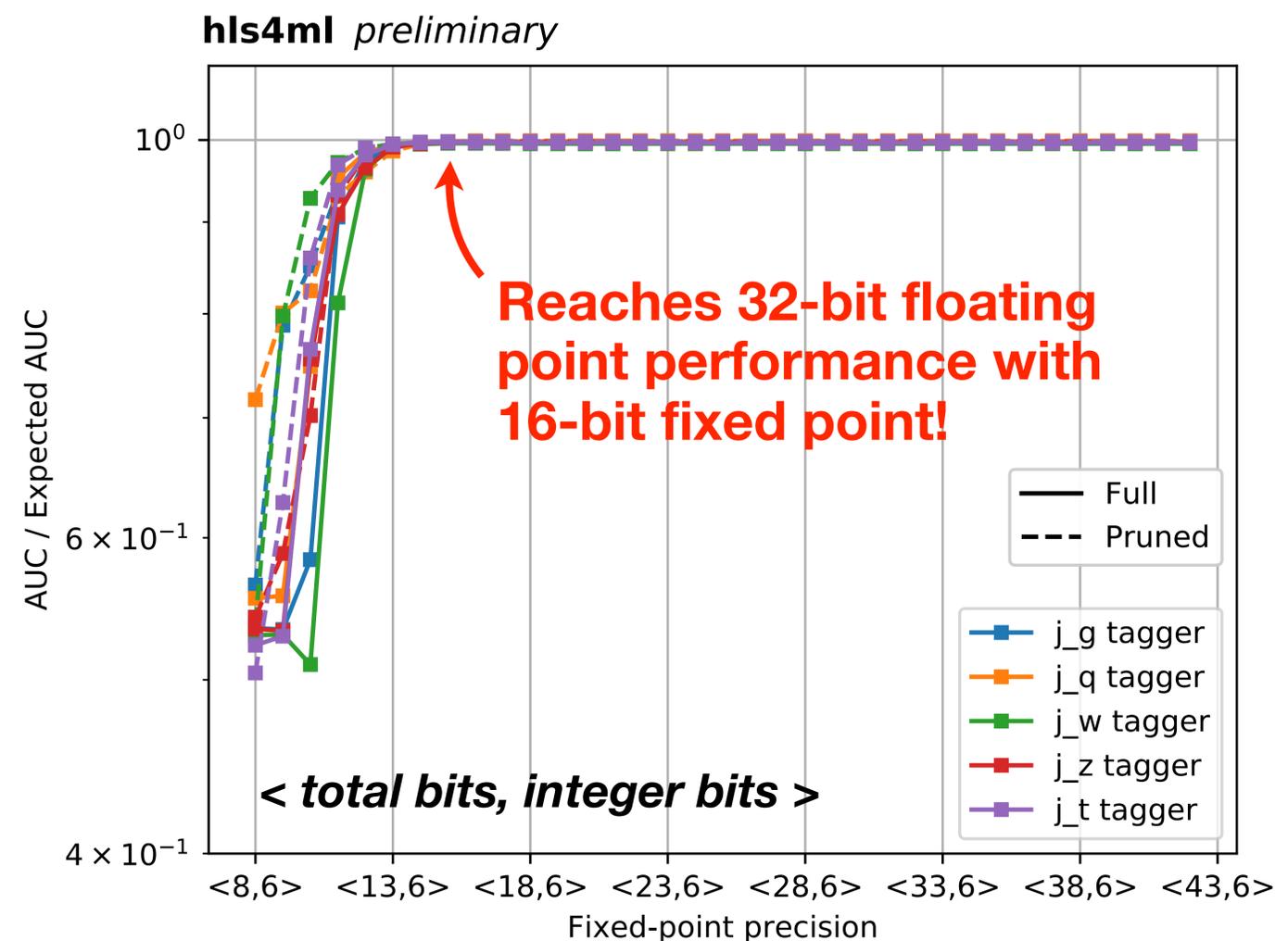
Remove X% of weights and retrain

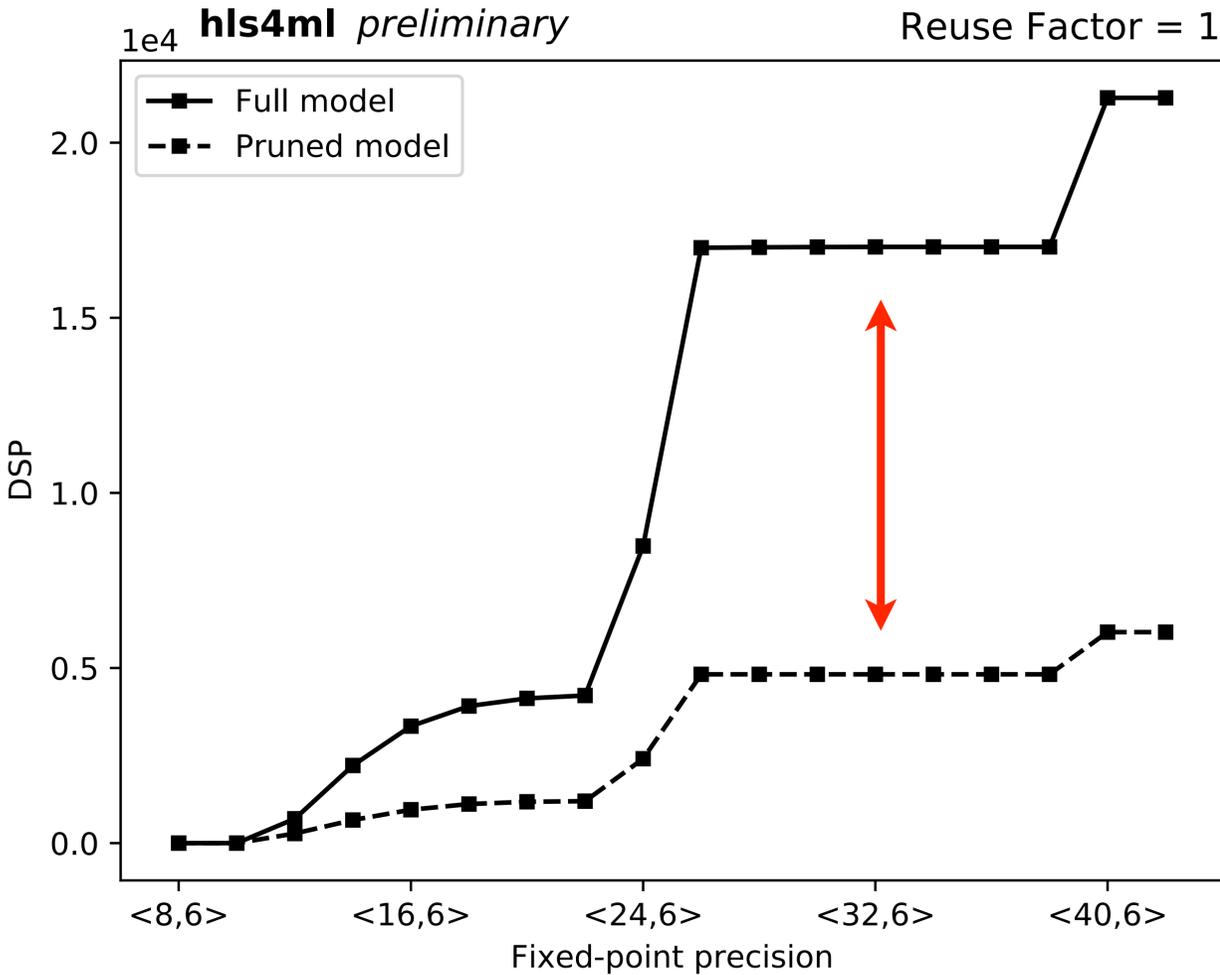
Rinse, repeat

## Our case study: 70% weight reduction with no performance loss

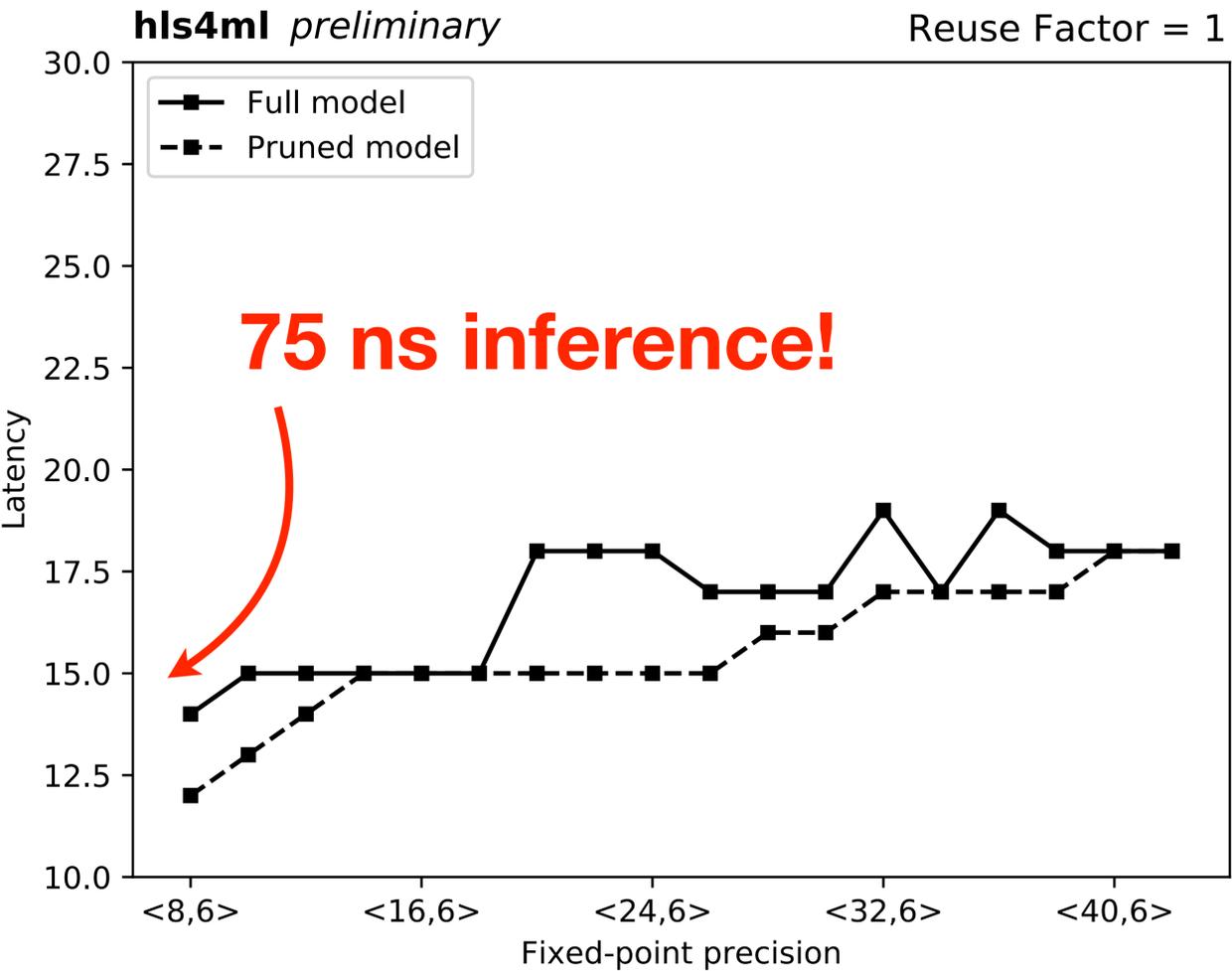


Distribution of weights in NN  
with L<sub>1</sub> regularization



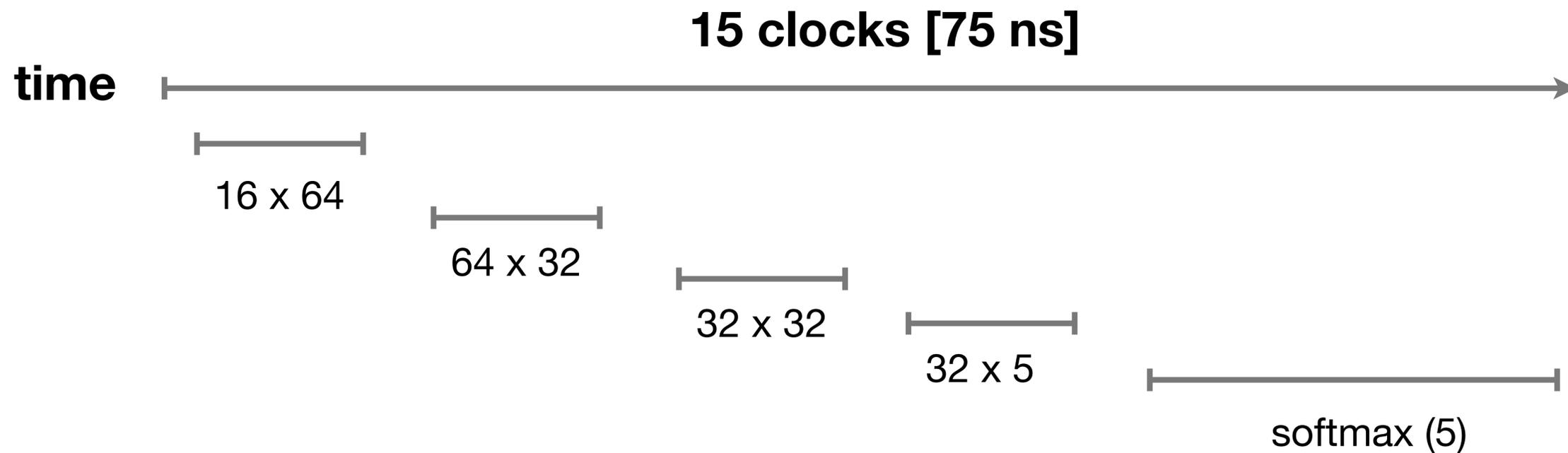


**Big reduction in DSP resources**



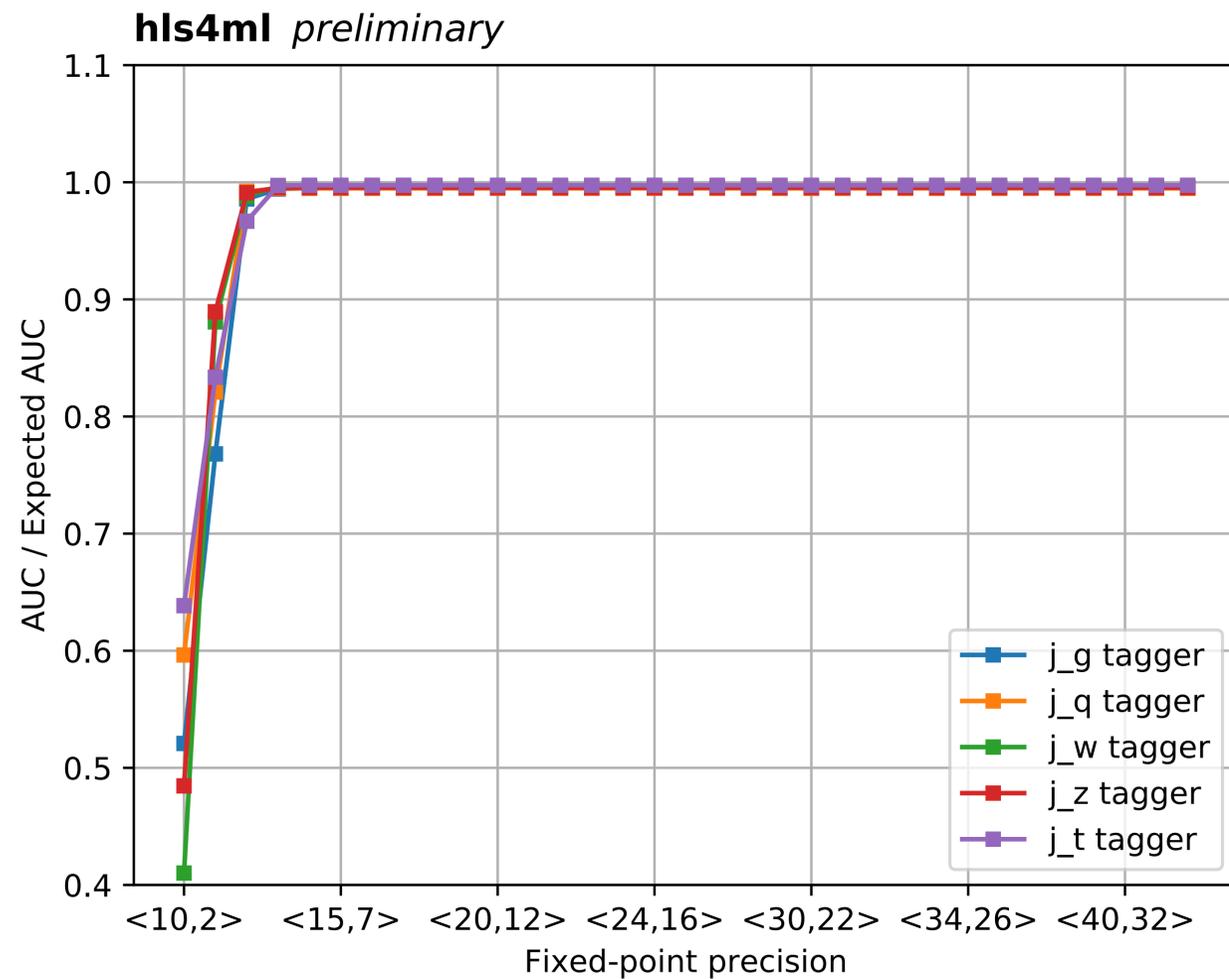
# 3-LAYER PRUNED MODEL - UNDER THE HOOD

Let's take our **3-layer pruned neural network** as a benchmark and study it in more detail

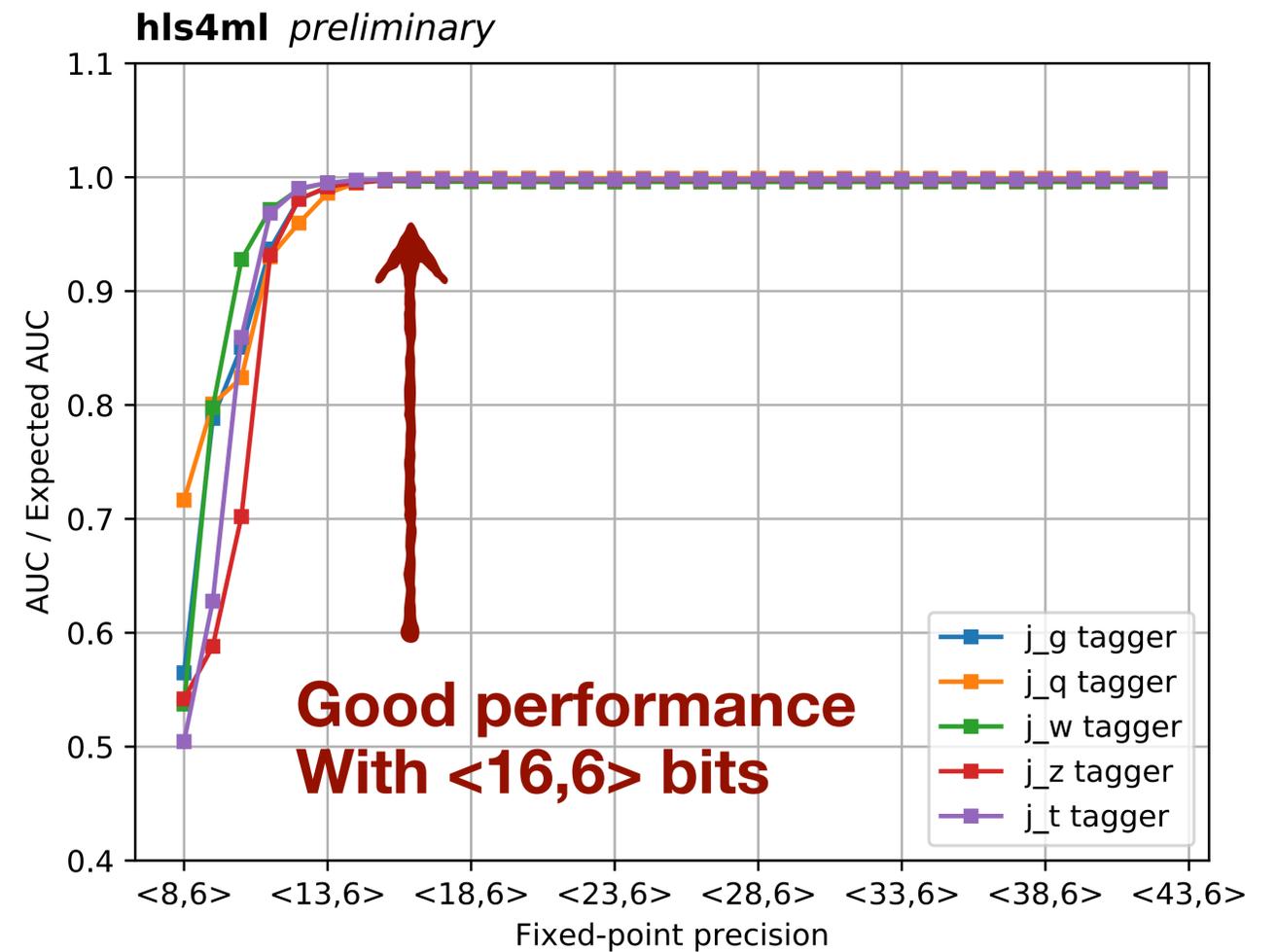


Reuse = 1	BRAM	DSP	FF	LUT
Total	13	1116	47k	35k
% Usage	~0%	20%	3%	5%

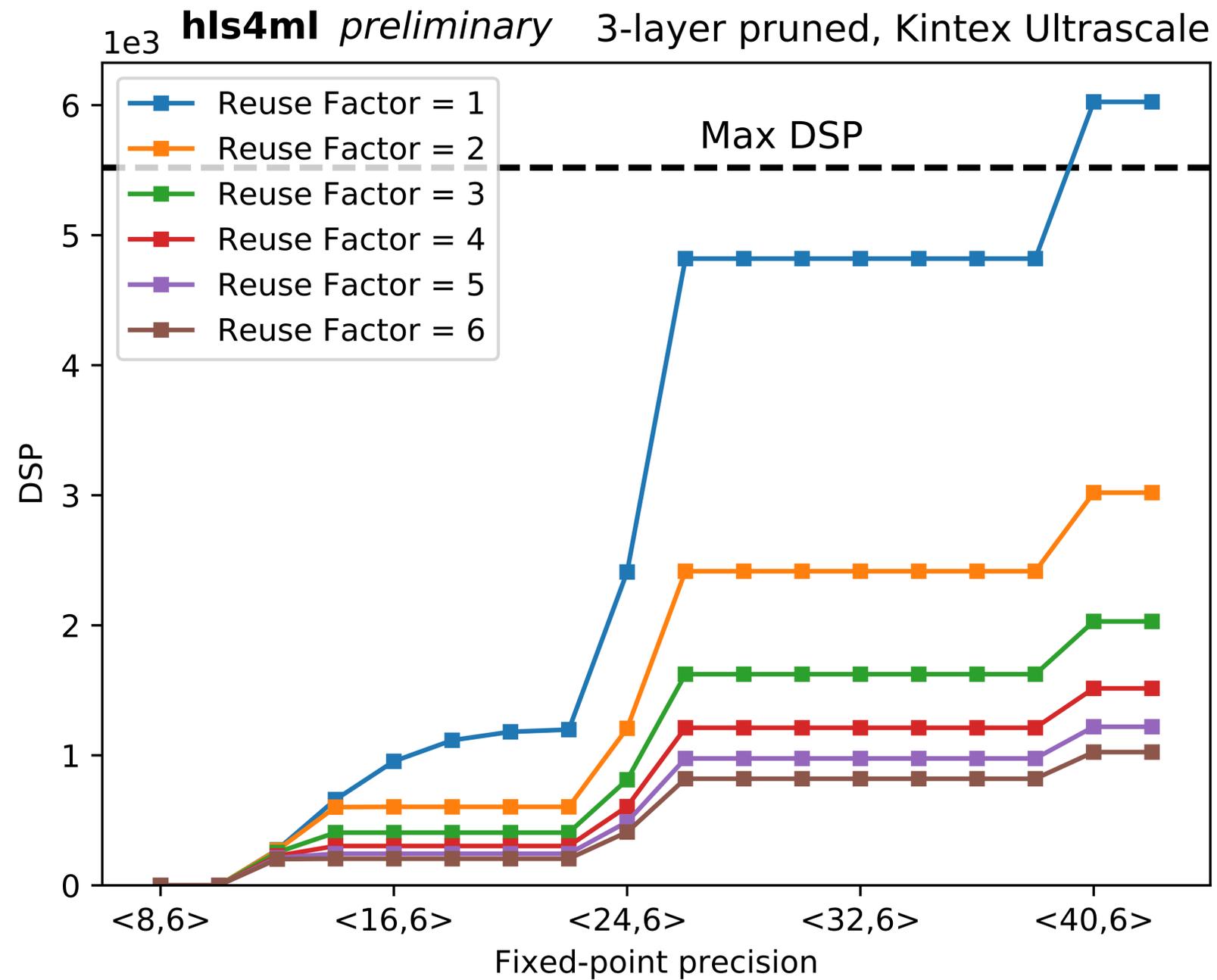
Scan integer bit



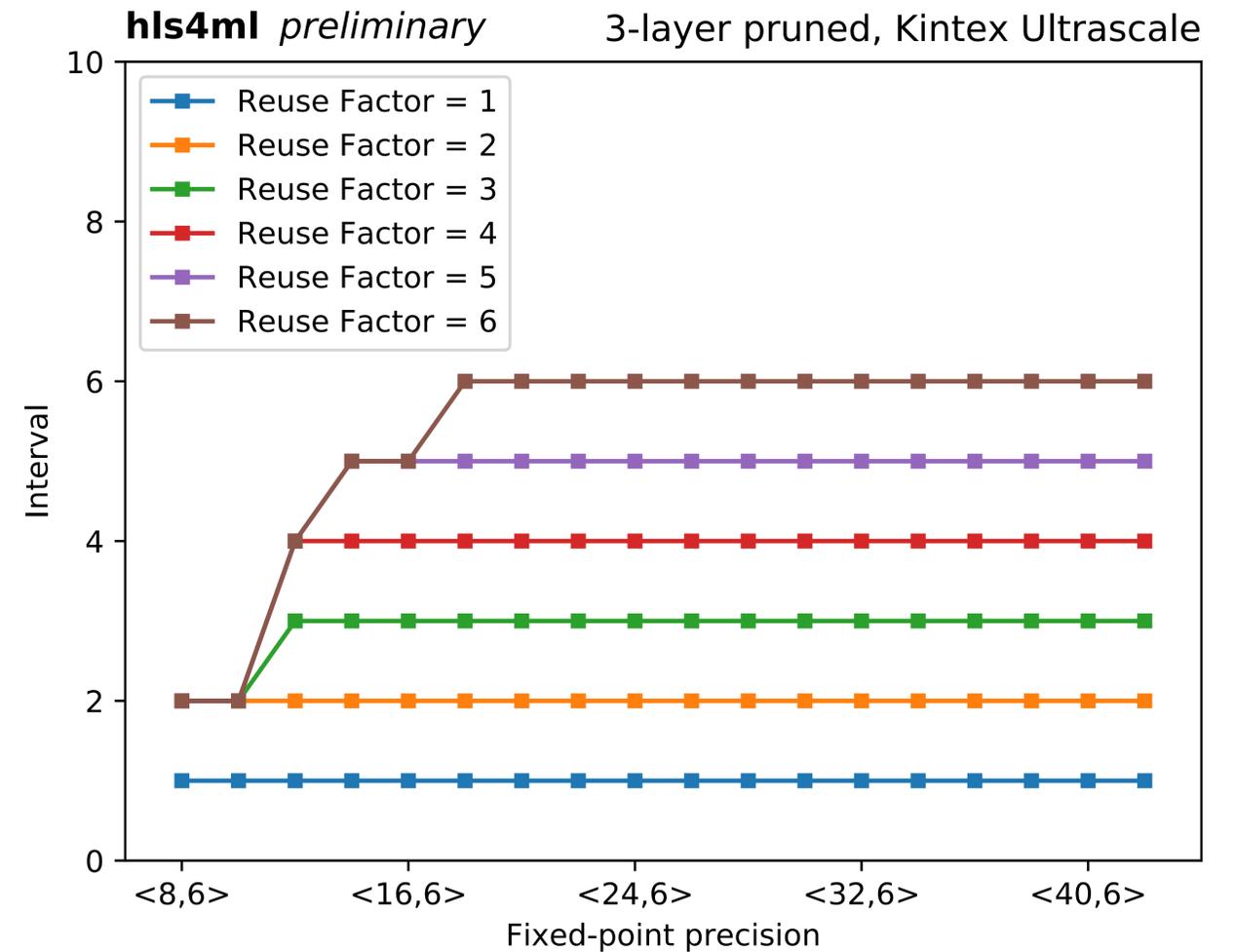
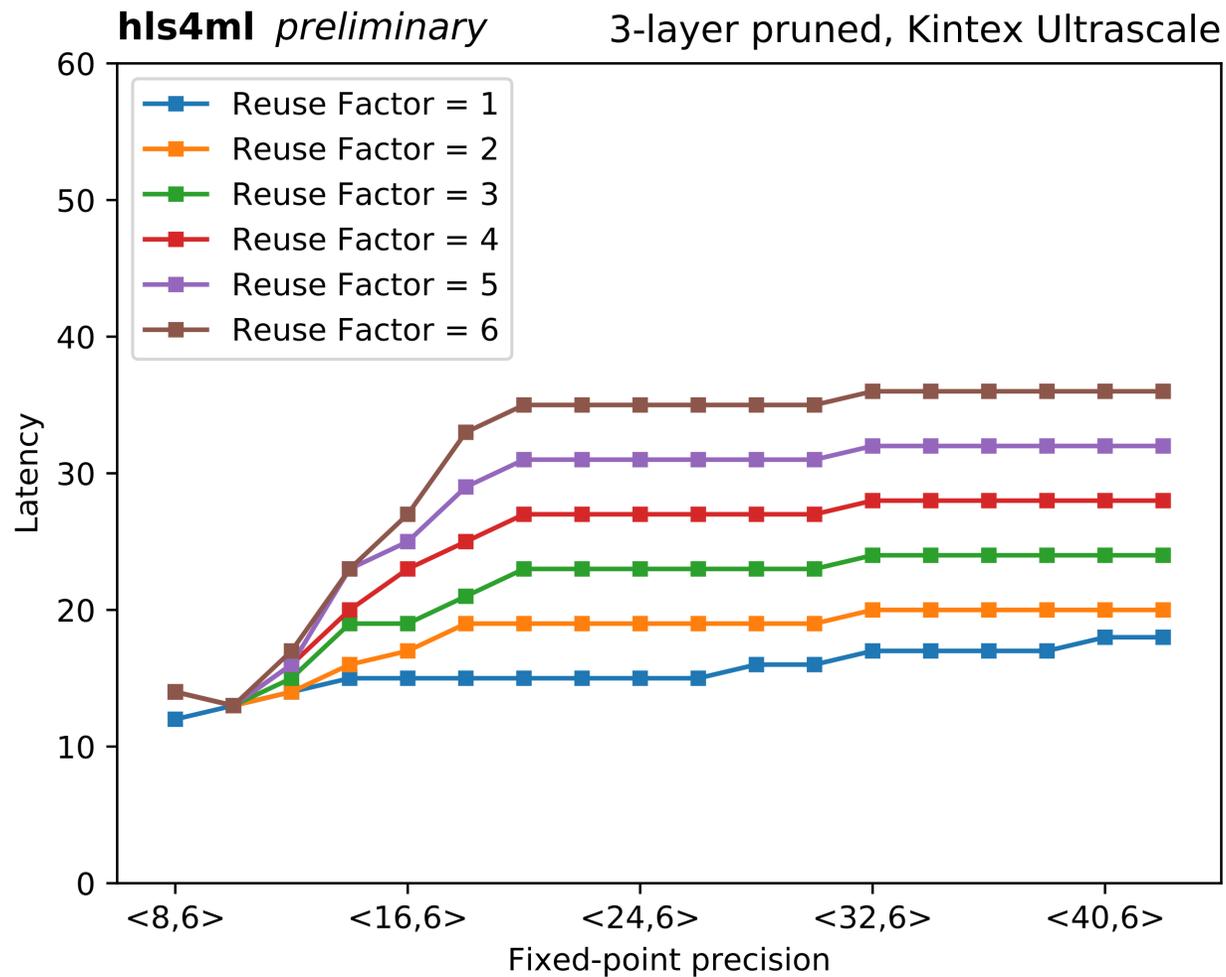
Scan decimal bit



General strategy, avoid overflows in integer bit;  
then scan the decimal bit until reaching optimal performance

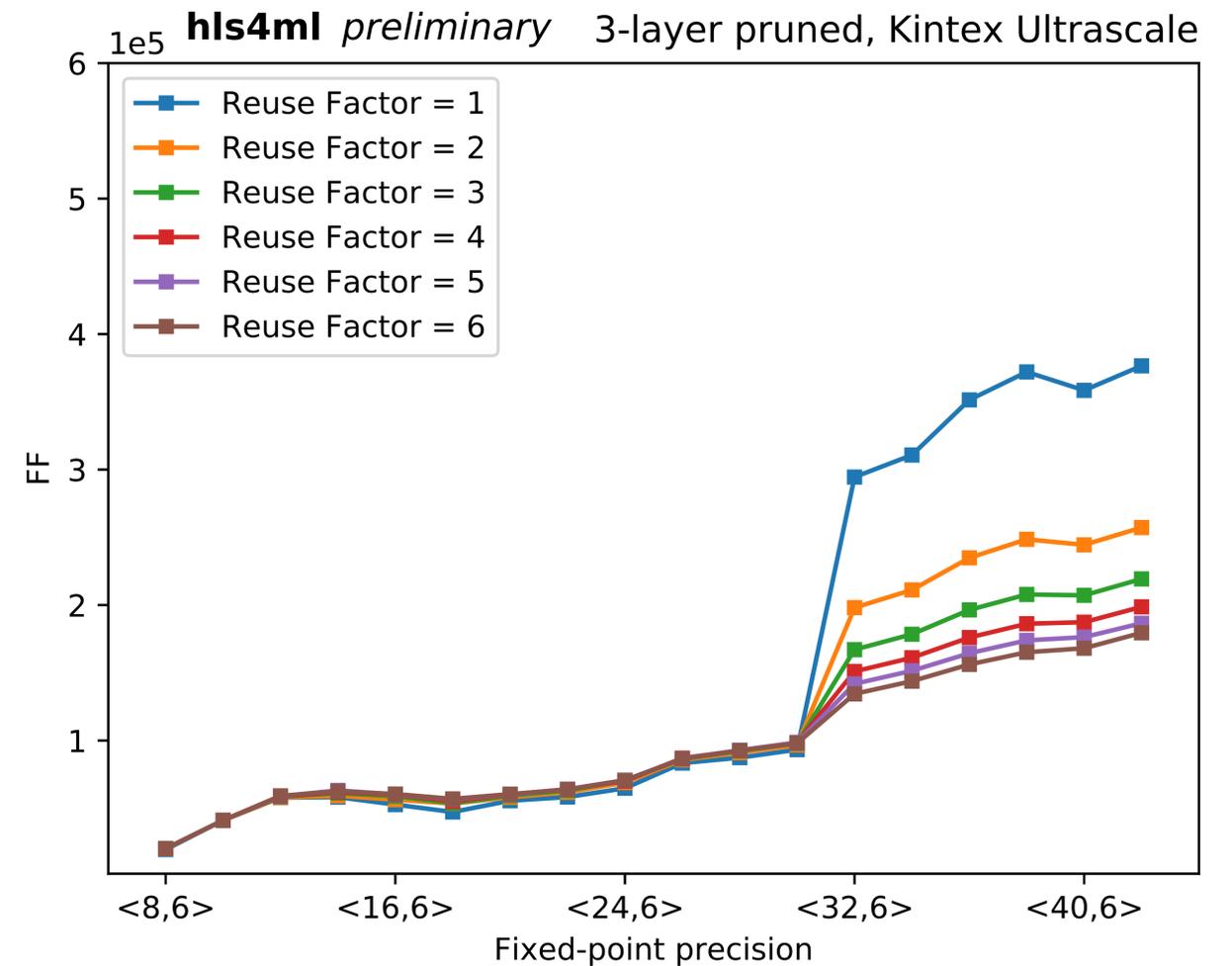
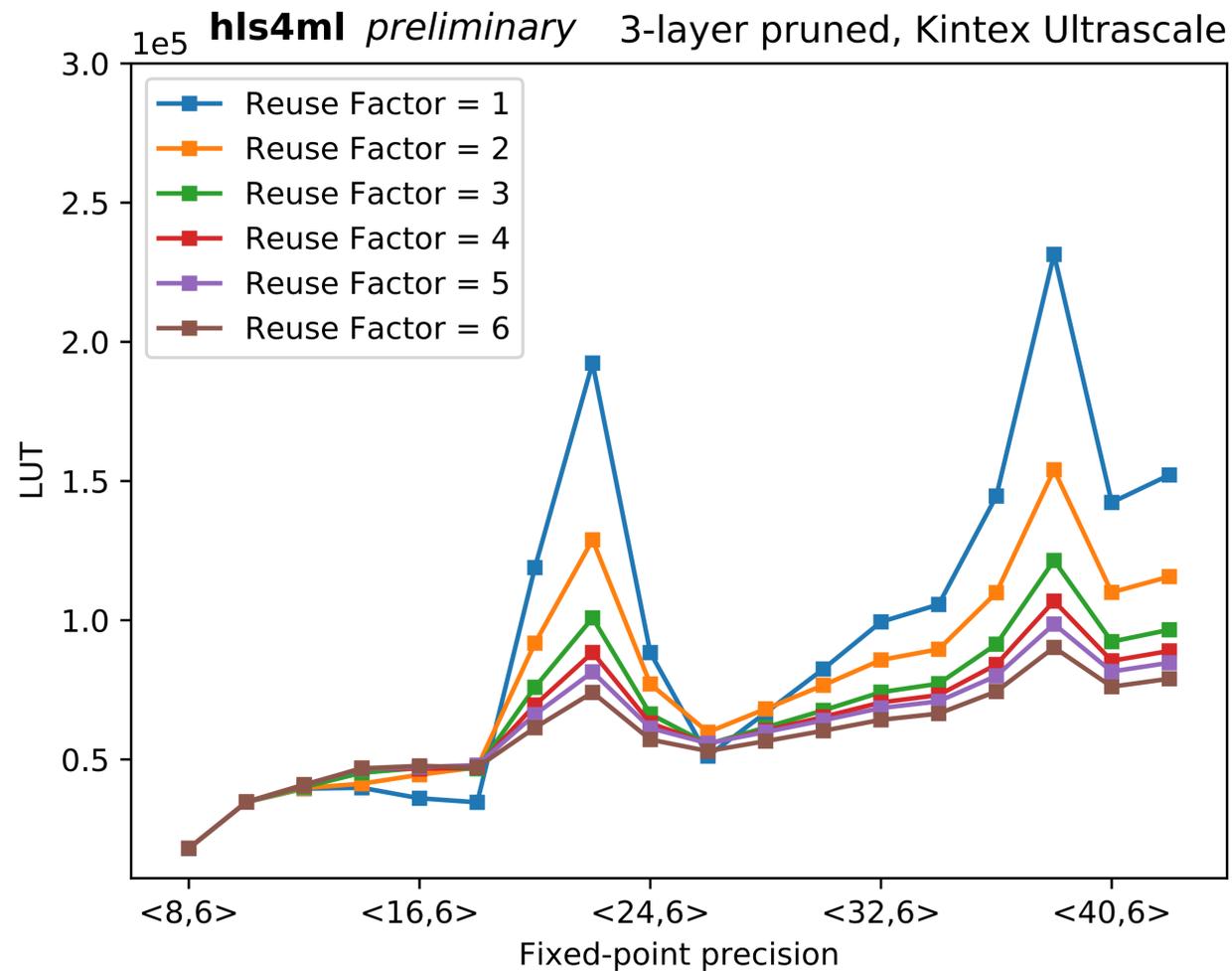


**Tuning the throughput with reuse factor will reduce the DSP usage**



**Additional latency introduced by reusing the multipliers**

**Behavior of pipeline interval controlled well by the reuse factor**



**Steady increase of LUTs and FFs vs. bit precision**

*Spikes in LUTs at the DSP precision transitions*

*Will come back to these estimates later in final implementation*

## **Xilinx Vivado 2017.2**

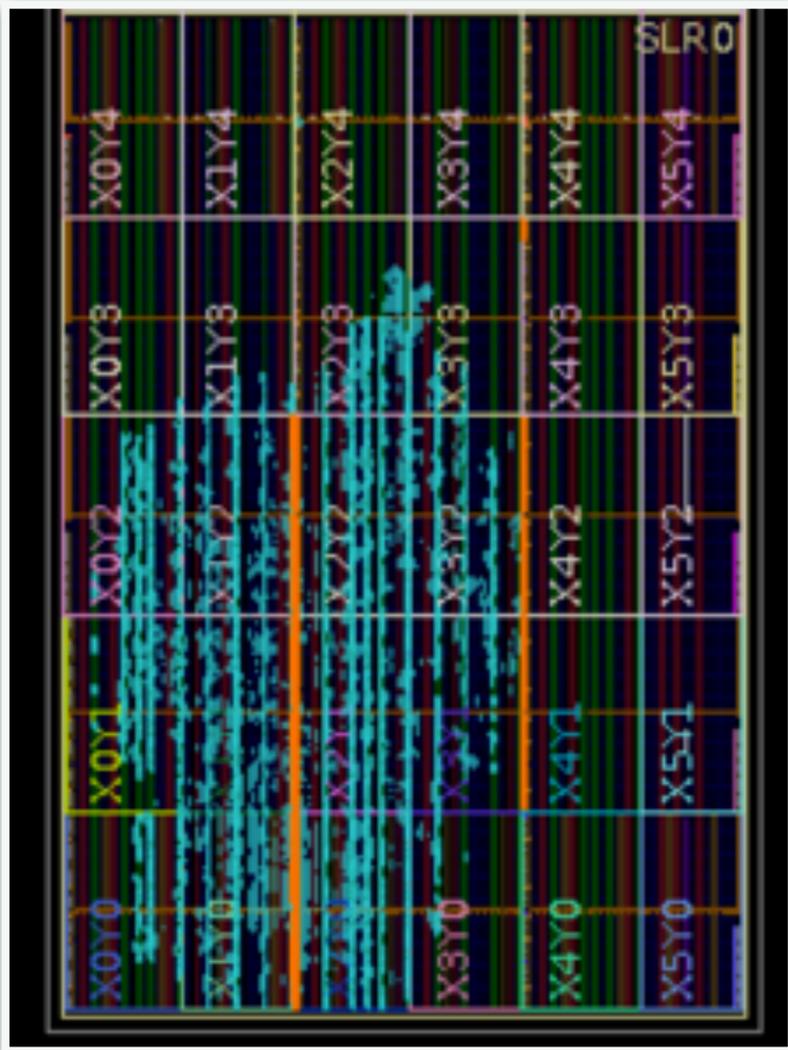
Results are slightly different in other versions of Vivado  
e.g. 2016.4 optimization is less performant for Xilinx ultrascale FPGAs

## **Clock frequency: 200 MHz**

Latency results can vary (~10%) with different clock choices

## **FPGA: Xilinx Kintex Ultrascale (XCKU115-FLVB2104)**

Results are slightly different in other FPGAs  
e.g. Virtex-7 FPGAs are slightly differently optimized



## How optimal is the HLS design?

For DSPs, one can get an optimal “back-of-the-envelope” estimate

DSPs are the dominant resource usage

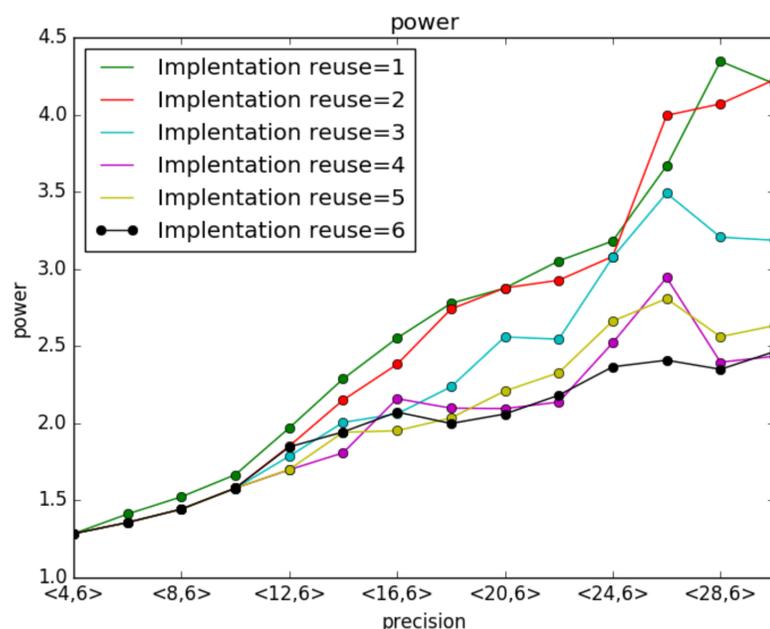
HLS resource estimates are great for quickly getting an estimate of resources

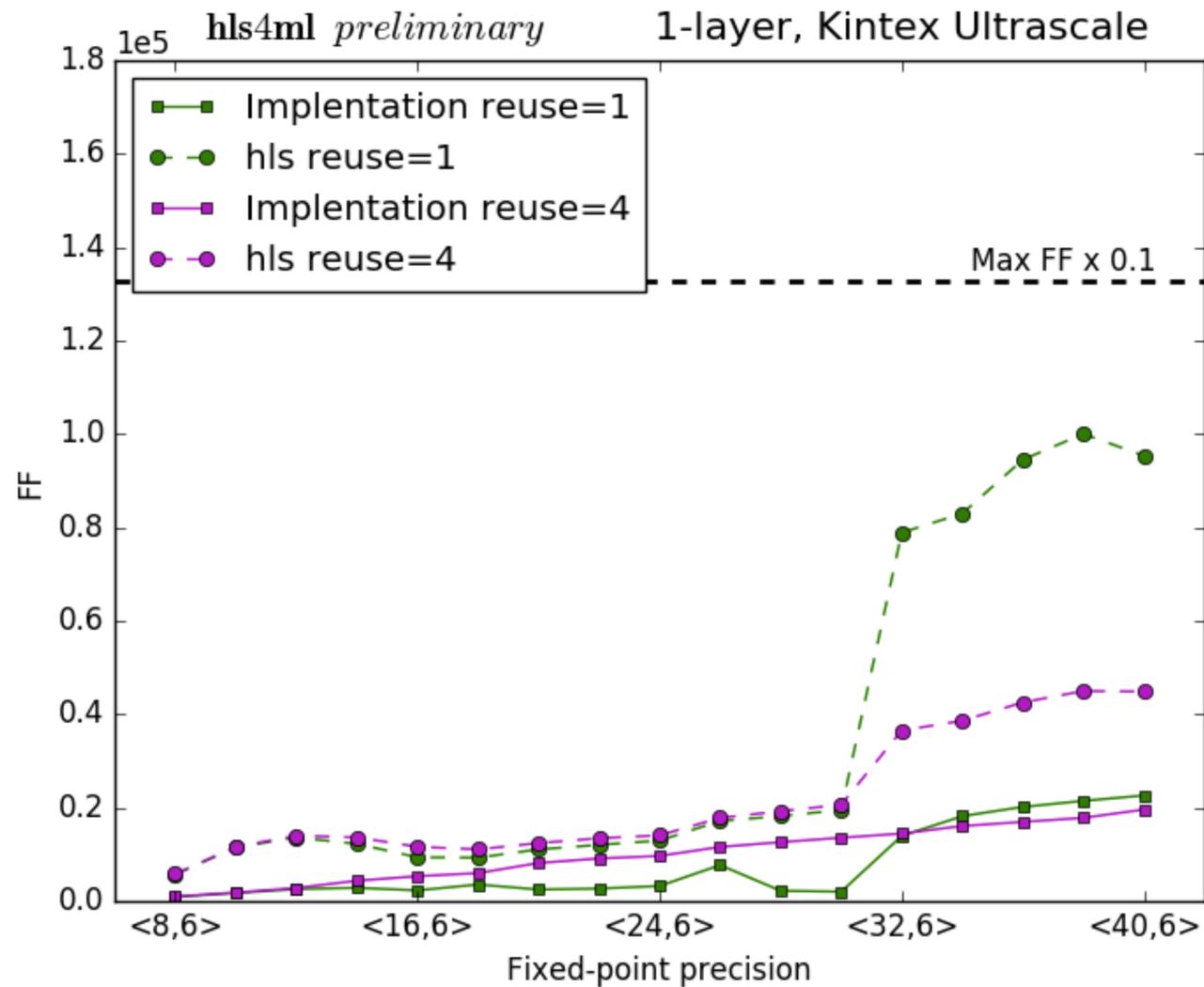
## Final firmware implementation gives actual resource usage

This depends on what else is in the rest of the firmware design

We implement the network in a “barebones” design and compare

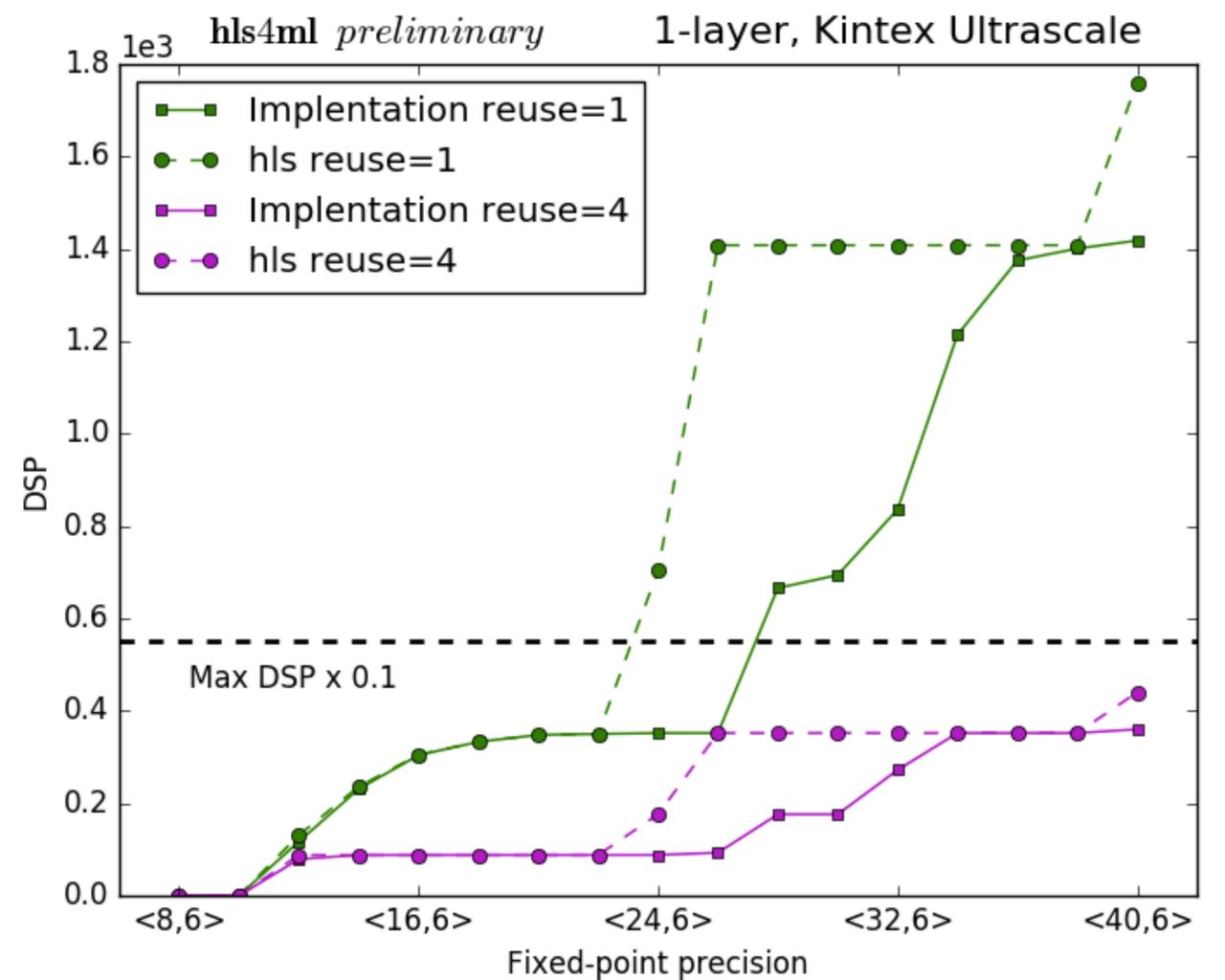
Does not always meet timing (see bonus)





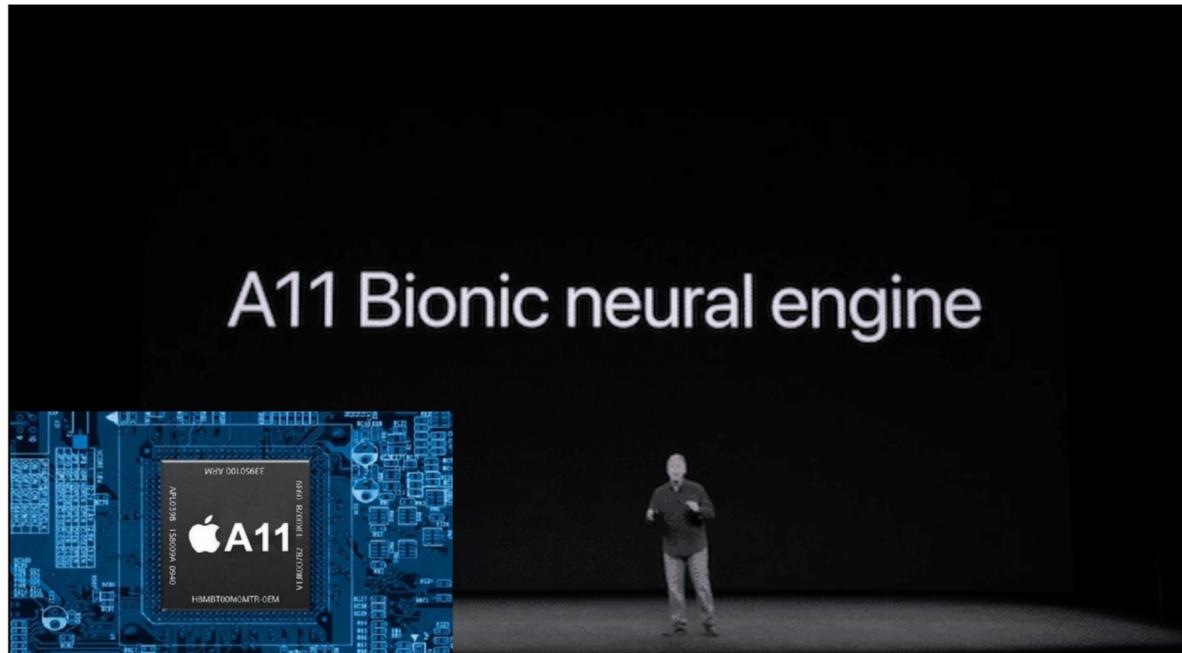
**DSPs are estimated fairly well**  
 Around DSP precision transition, Vivado does further optimizations

**FFs are overestimated by a factor of 2-4**  
**LUTs are overestimated by a factor of ~2**



# **status and outlook**

Specialized co-processor hardware for machine learning inference



## INTEL<sup>®</sup> FPGA ACCELERATION HUB

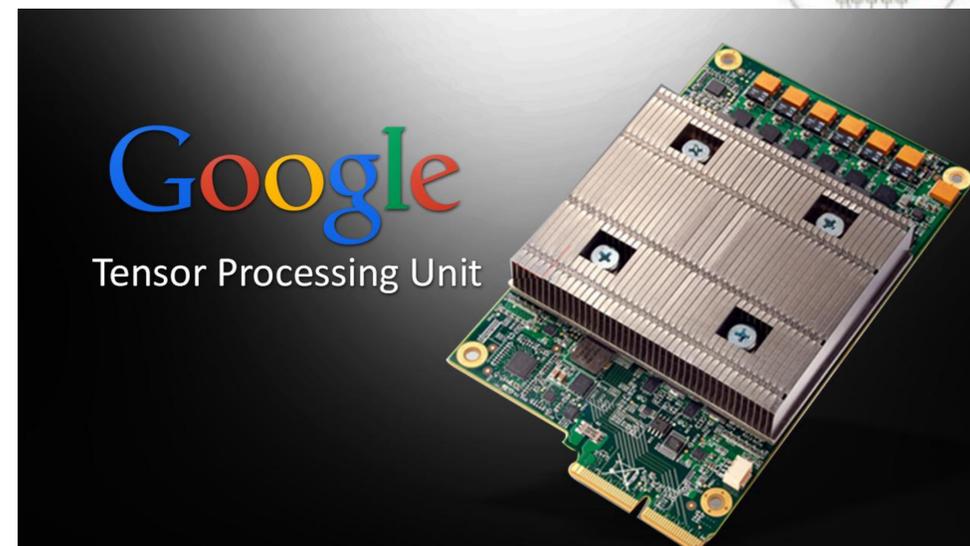
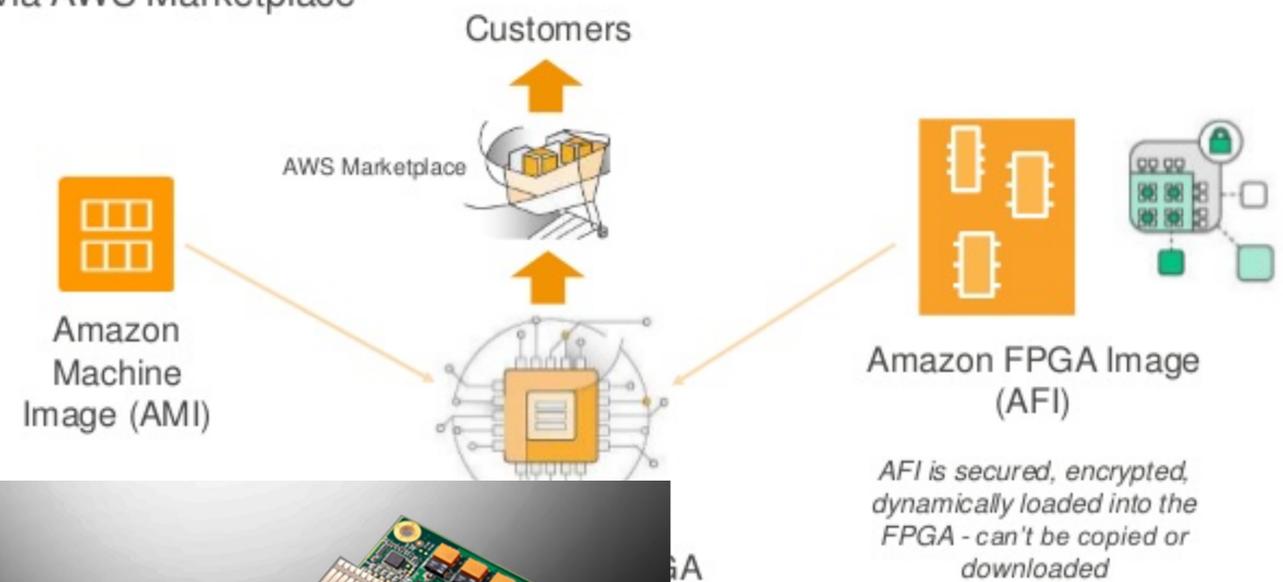
The Intel<sup>®</sup> Xeon<sup>®</sup> Acceleration Stack for FPGAs is a robust framework enabling data center applications to leverage an FPGA's potential to increase



## Catapult/Brainwave



## Delivering FPGA Partner Solutions on AWS via AWS Marketplace



Specialized co-processor hardware for machine learning inference



Microsoft

Catapult/Brainwave

**The same NN libraries can be adapted for other applications**  
From  $\mu$ s to ms timescales and small to large networks



## 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



<https://hls-fpga-machine-learning.github.io/hls4ml/>

hls4ml

Setup

Dependencies

Quick Start

Configuration

Concepts

Release Notes

Reference and Contributors

Code Repository

Published with GitBook



A package for machine learning inference in FPGAs. We create firmware implementations of machine learning algorithms using high level synthesis language (HLS). We translate traditional open-source machine learning package models into HLS that can be configured for your use-case!

The project is currently in development, so please let us know if you are interested, your experiences with the package, and if you would like new features to be added.

contact: [hls4ml.help@gmail.com](mailto:hls4ml.help@gmail.com)

## Project status

*in beta version (March 17, 2018), v0.1.2*

**BETA VERSION IS LIVE!**

Watch for arXiv preprint in a few weeks

## More network architectures

Library for **convolutional layers** are implemented, *in testing*

Similar to deep neural networks, more parameters

Prototype libraries exist for **recurrent layers and LSTM**

Similar to DNN for in multiplication, different data flow paradigm

## More ML codes

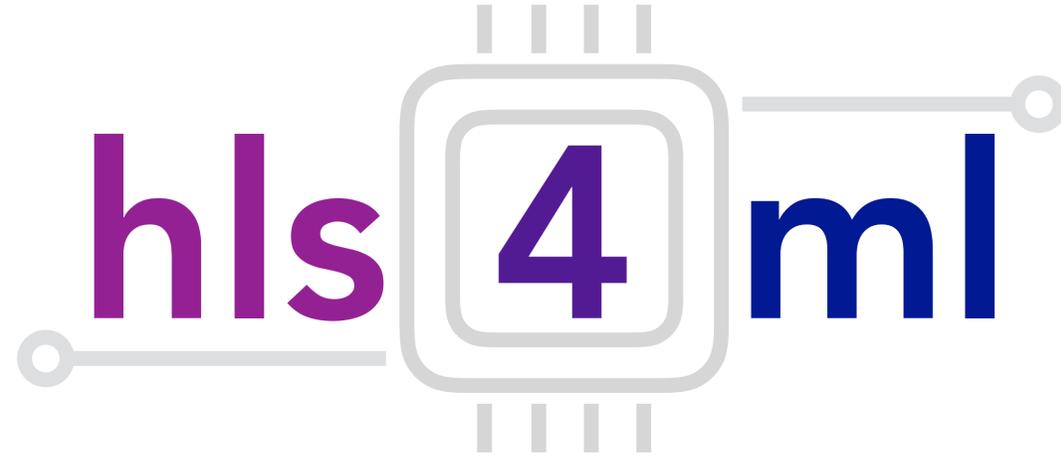
Support currently for **Keras** and **TensorFlow**

**PyTorch** has a first implementation, in testing

Others?

## More FPGAs

Support planned for **Altera Quartus HLS**



We introduce a software/firmware package, **hls4ml**

Automated translation of neural networks into firmware using HLS

Case study present with jet substructure in L1 trigger

Tunable configuration for a broad range use cases

More info here:

<https://hls-fpga-machine-learning.github.io/hls4ml/>

# Bonus

Timing result, 3 hidden pruned, <16,8>, reuse=2

