

# Machine learning for low-latency inference

Dylan Rankin [UPenn]  
Lepton Photon 2023  
July 20th, 2023

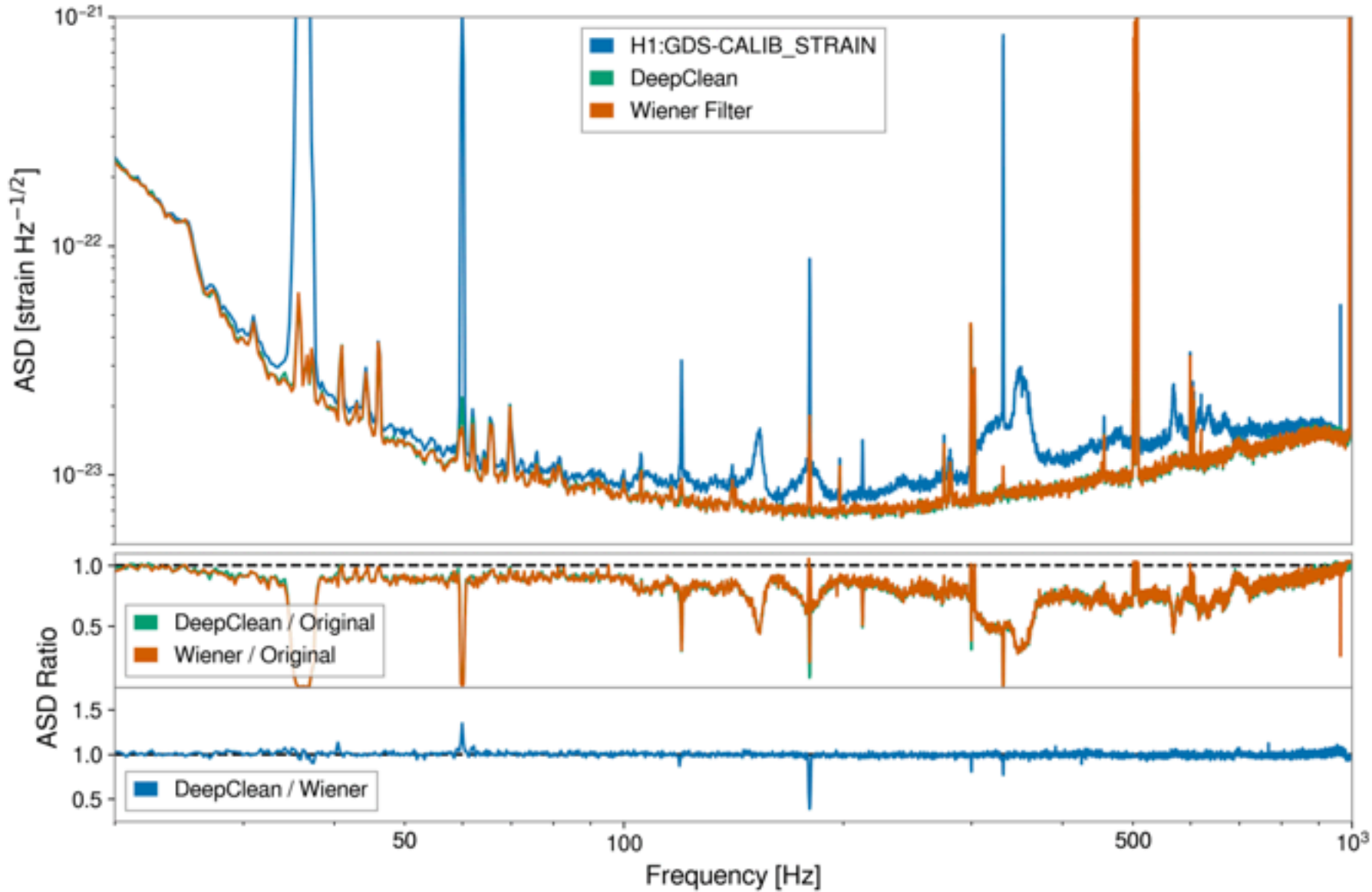
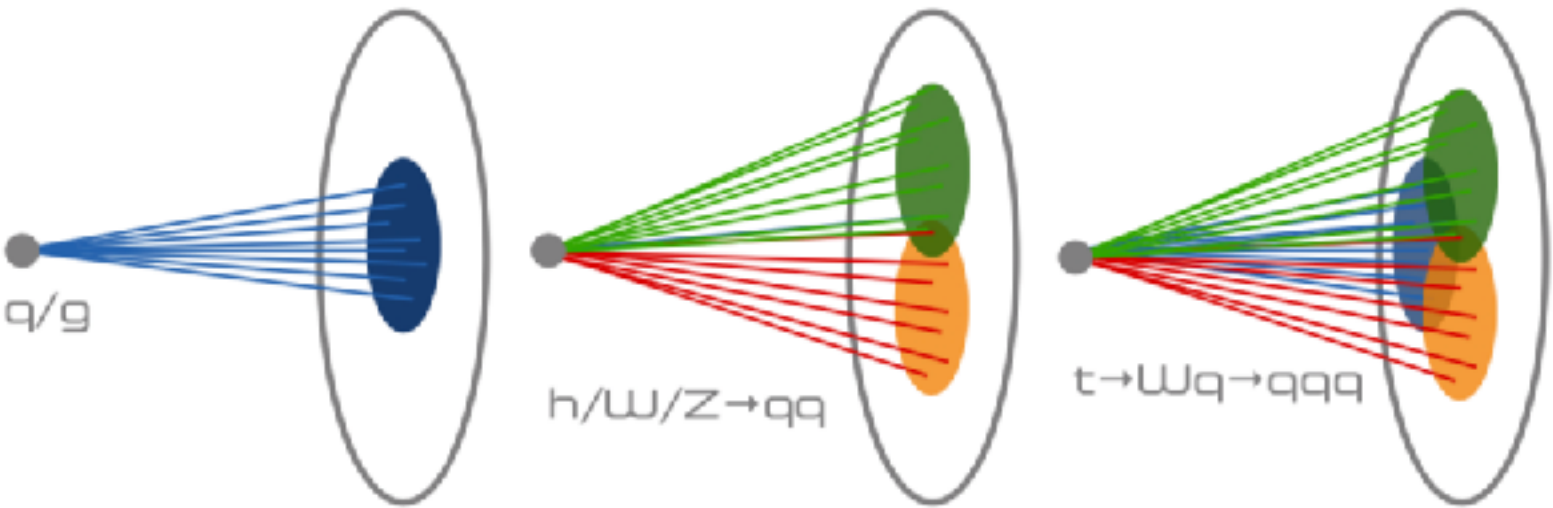
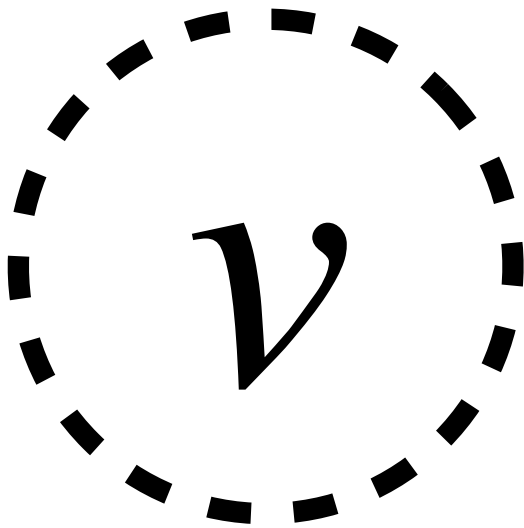
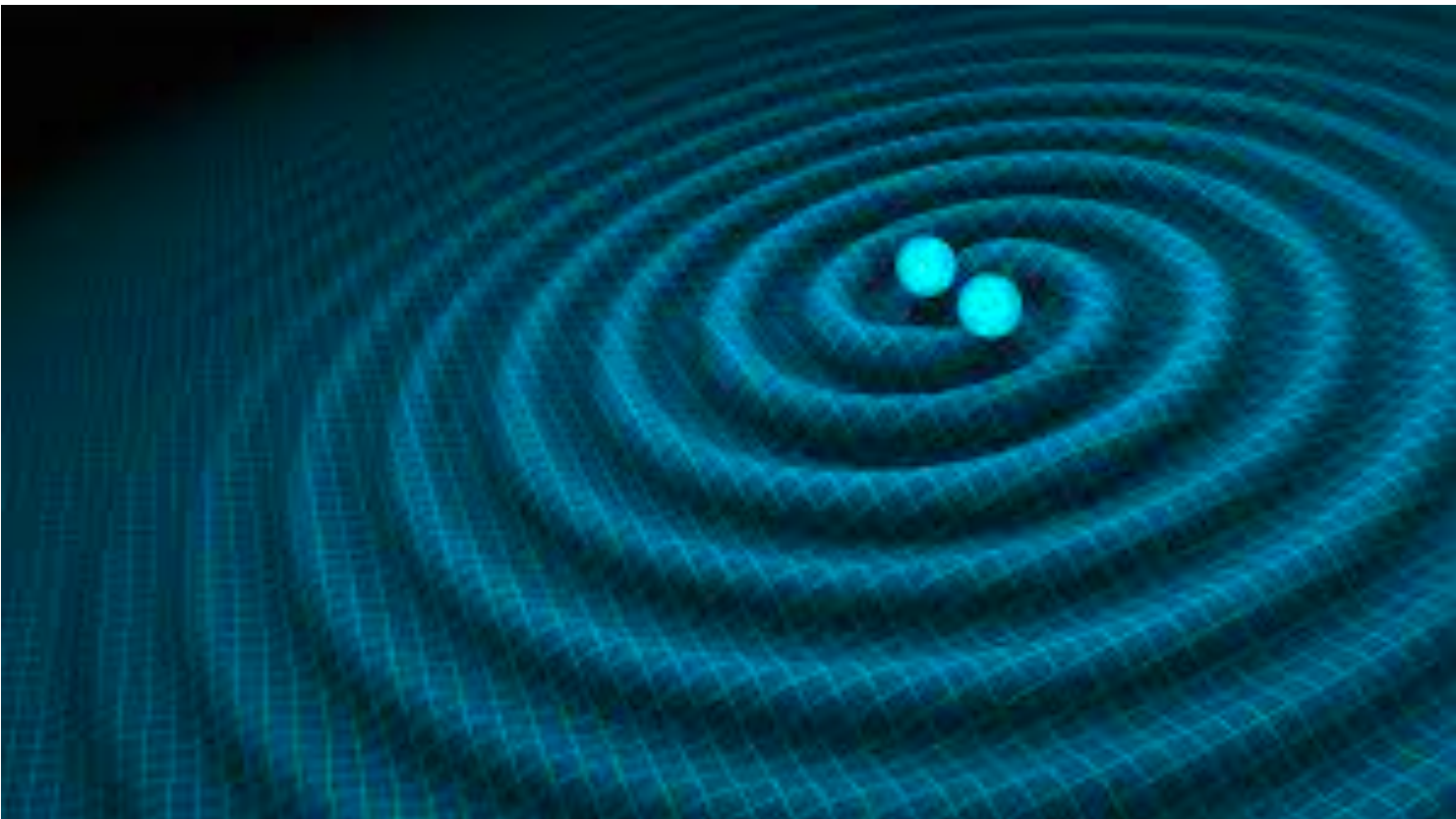




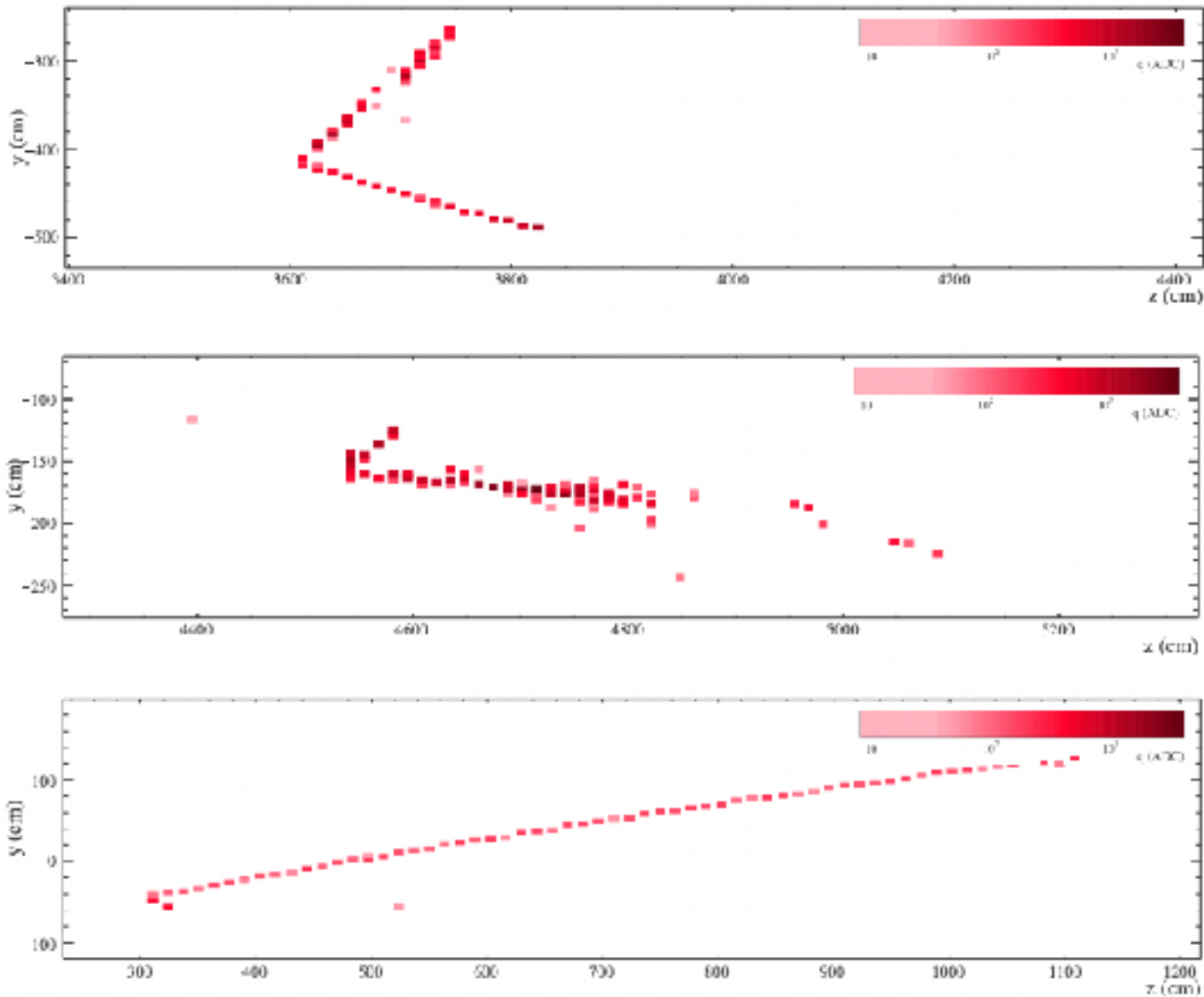
# Introduction

- Machine learning (ML) is becoming more and more popular
  - Availability of CPUs, GPUs, software has accelerated adoption
- What about low-latency (FPGA/ASIC-based) systems?
  - (How) can ML inference be run effectively in  $O(100\text{ ns})$ ?
- Applications
- Future

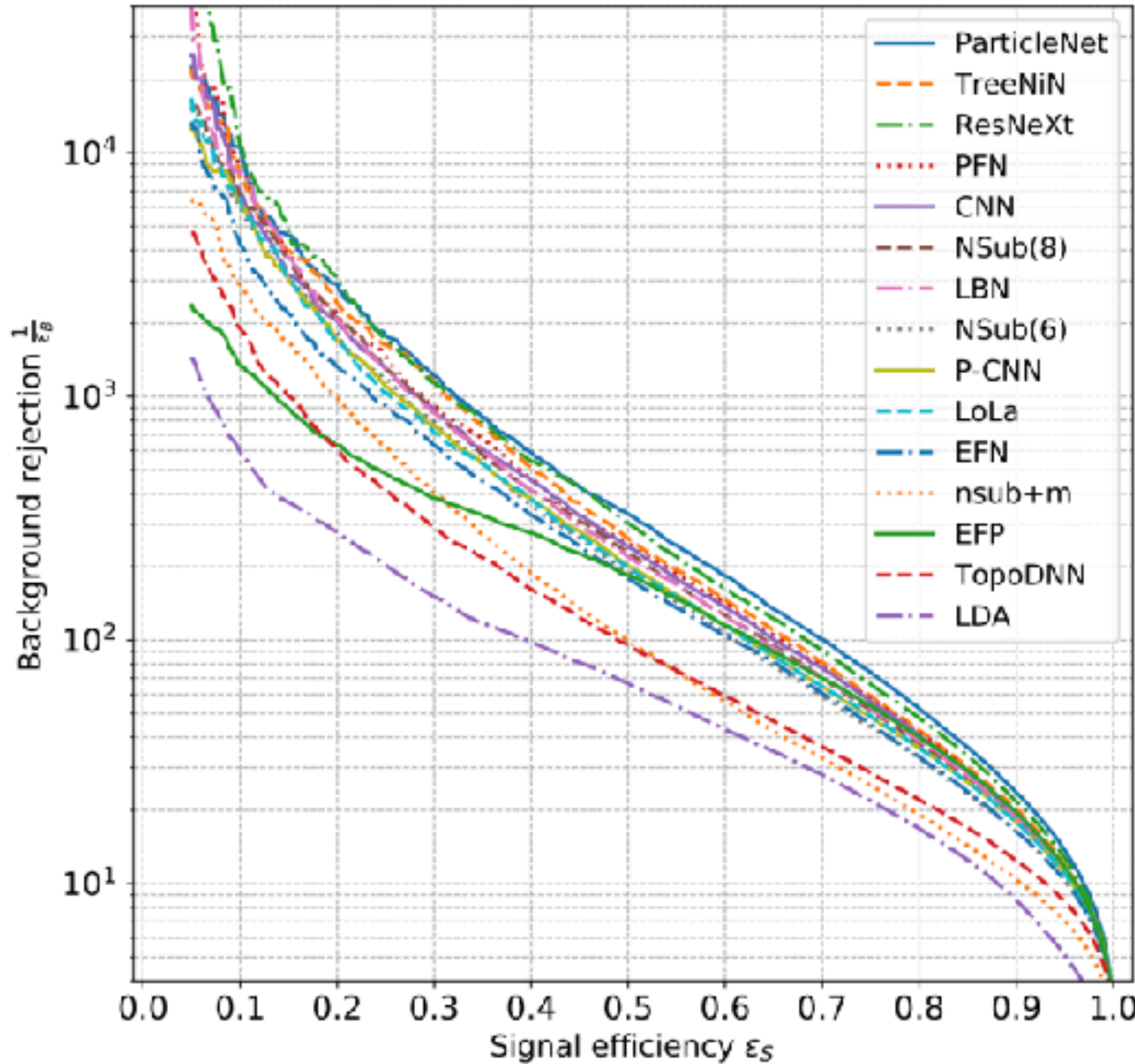
# Machine Learning (ML) is Everywhere



arXiv:2005.06534



arXiv:1604.01444



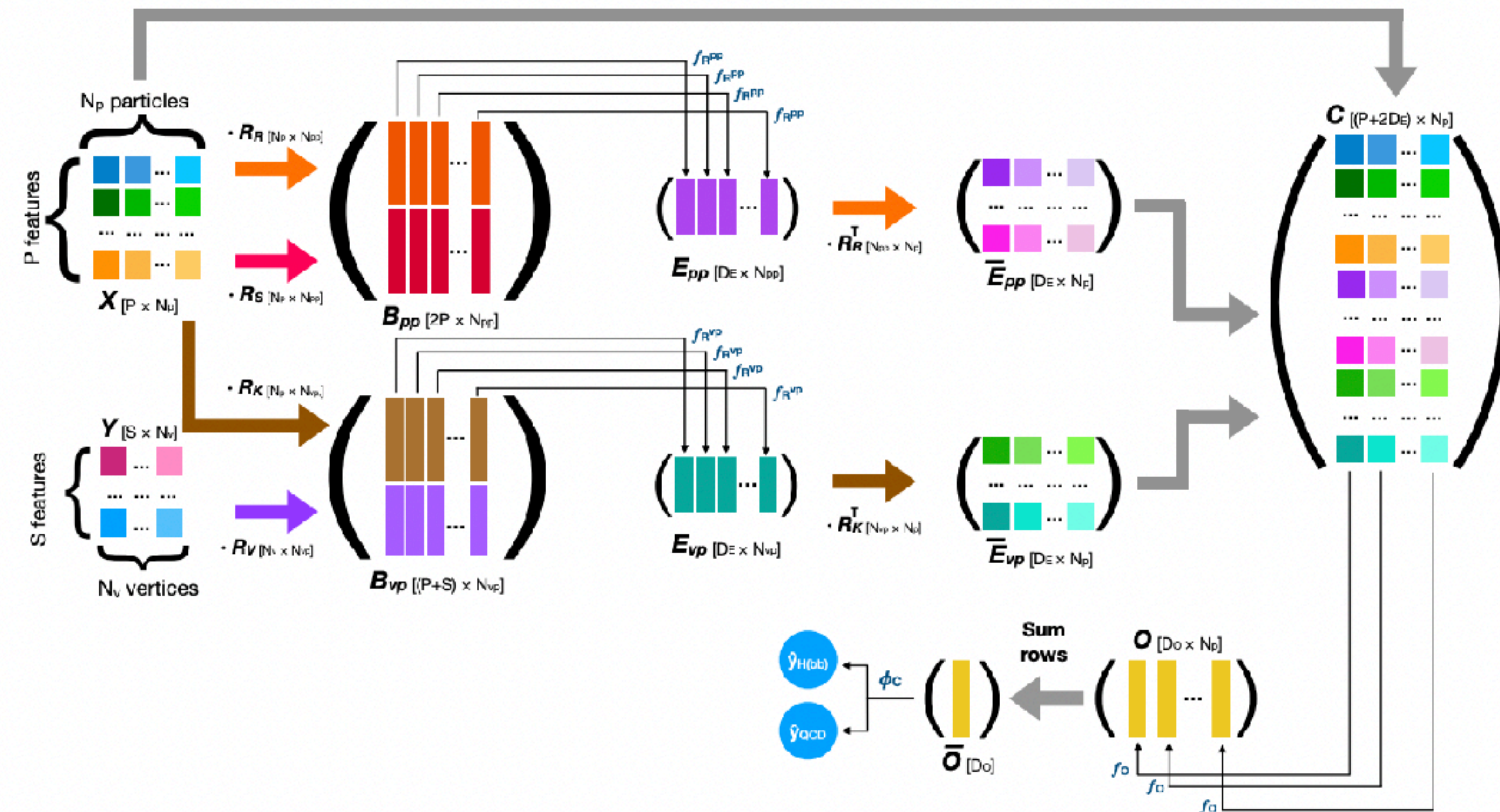
arXiv:1902.09914



# Machine Learning is (almost) Everywhere

- Trends in ML towards bigger and more complicated models, more computing (GPUs)
- → **Majority of ML in physics is “off detector”**

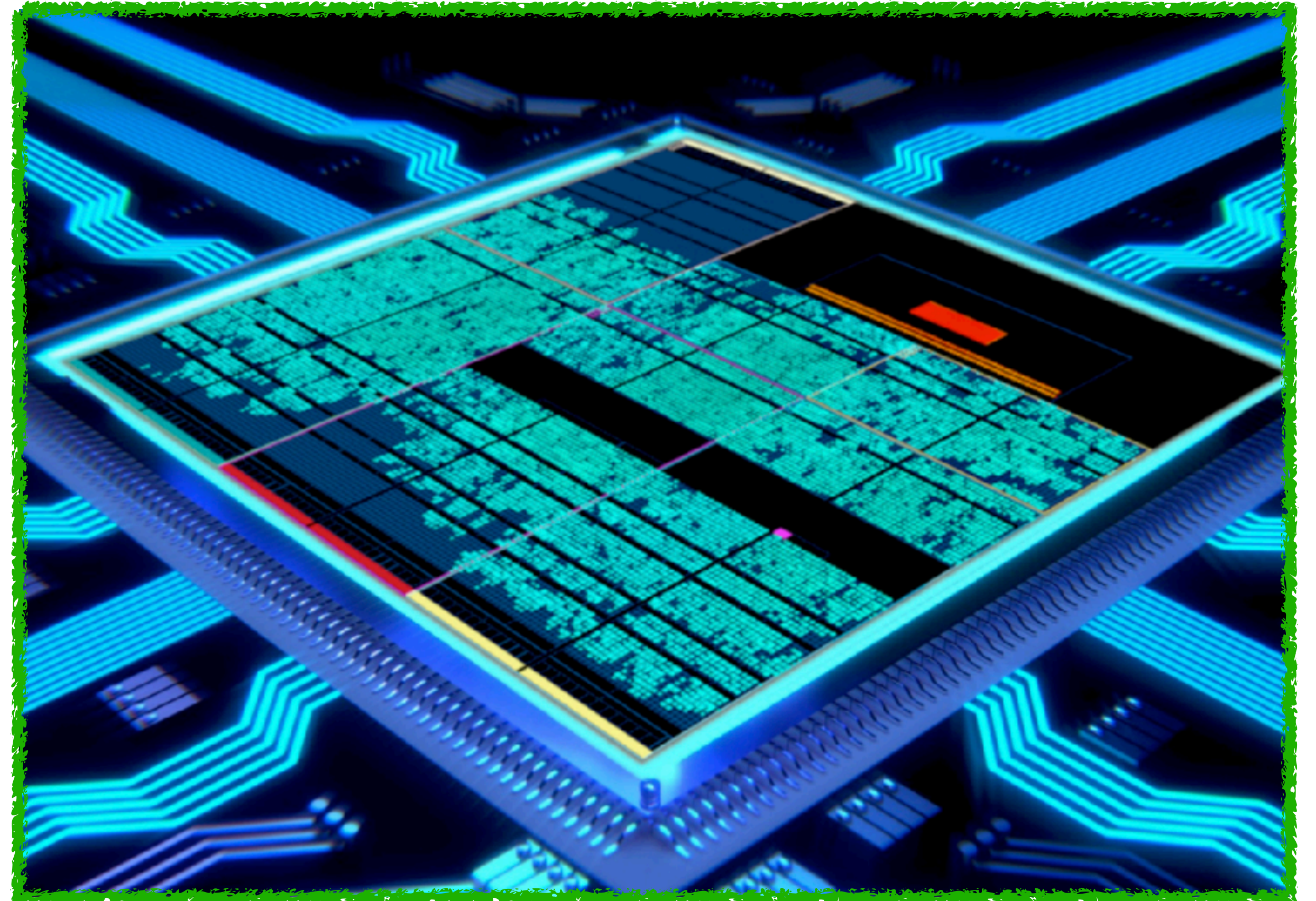
- System latency limits are typically soft (if at all)
- No radiation
- Issues do not impact data collection
- Can re-run algorithms/workflows





# What if...

- What if:
  - **System latency limits are low?**  
( $\lesssim \mu s$ )
  - High radiation?
- Requires dedicated hardware
- **FPGAs** (or ASICs)



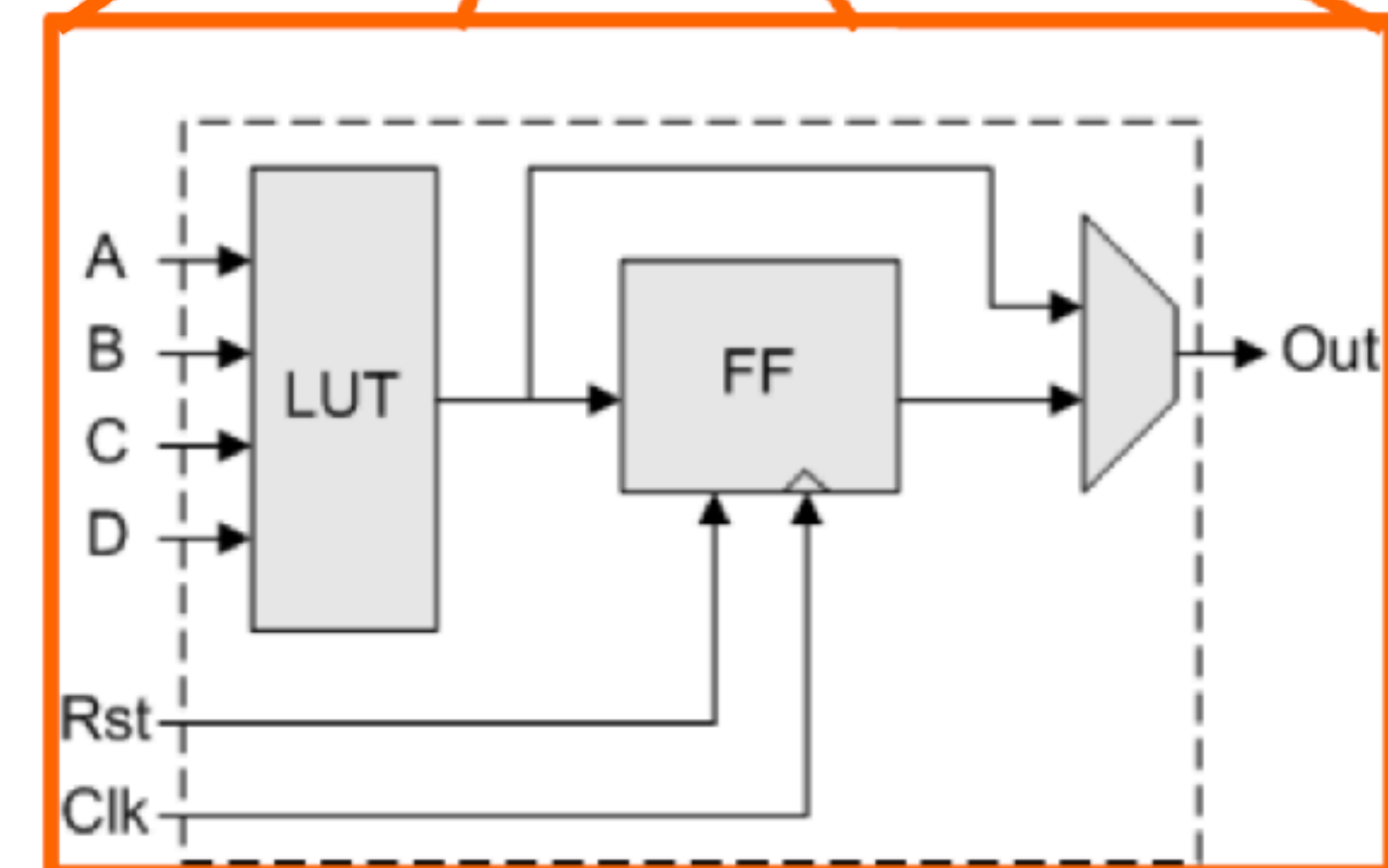
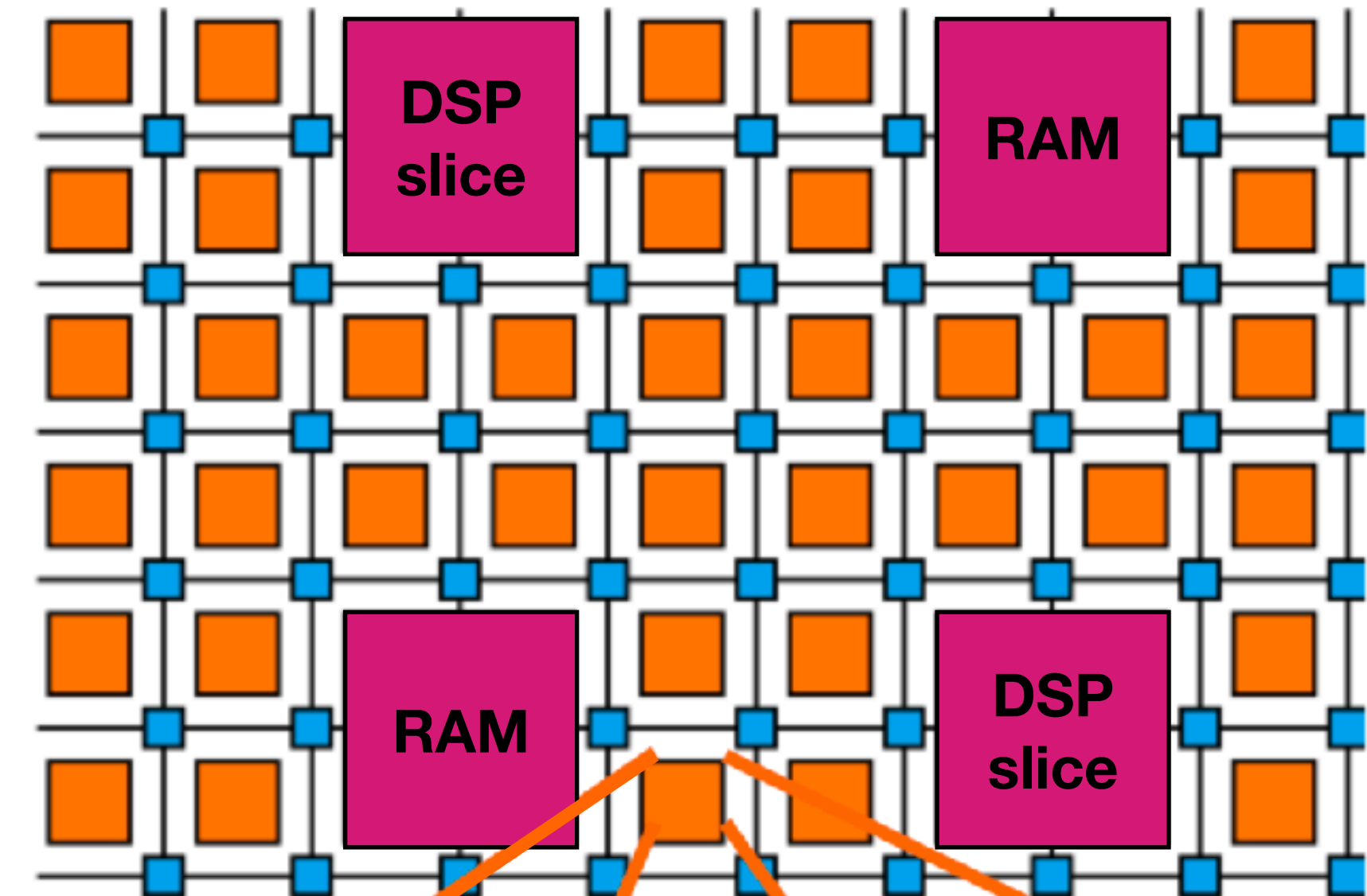


# What is an FPGA?

- Field-programmable gate array
- Building blocks:
  - Multiplier units (DSPs) [arithmetic]
  - Look Up Tables (LUTs) [logic]
  - Flip-flops (FFs) [registers]
  - Block RAMs (BRAMs) [memory]
- Algorithms are wired onto the chip
  - Can only use the resources on the chip
- Run at high frequency: hundreds of MHz, O(ns) runtime

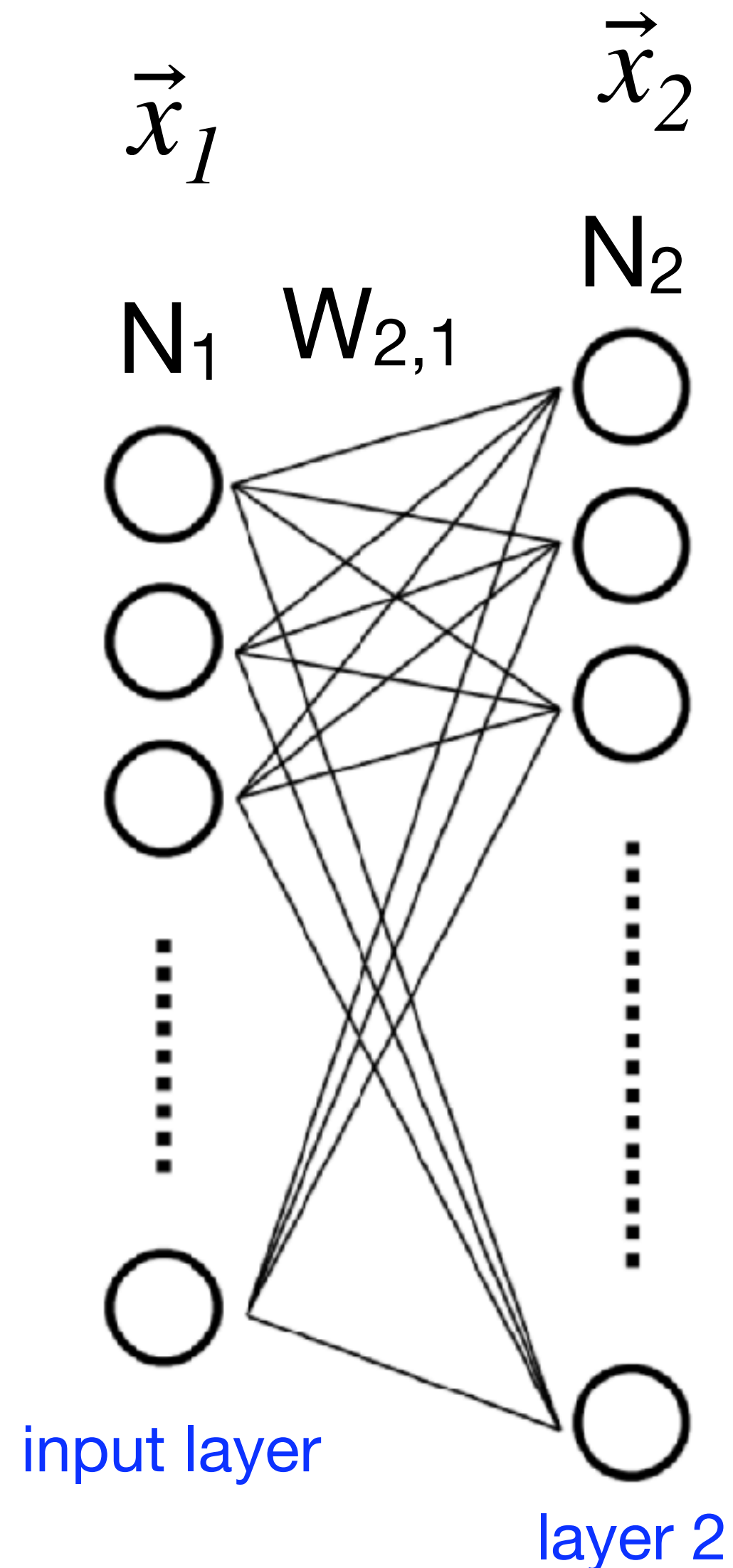
## Xilinx Virtex Ultrascale+ VU13P

12288 Multipliers  
1.7M LUTs  
3.4M FFs  
95 Mb BRAM





# What is a Neural Network?

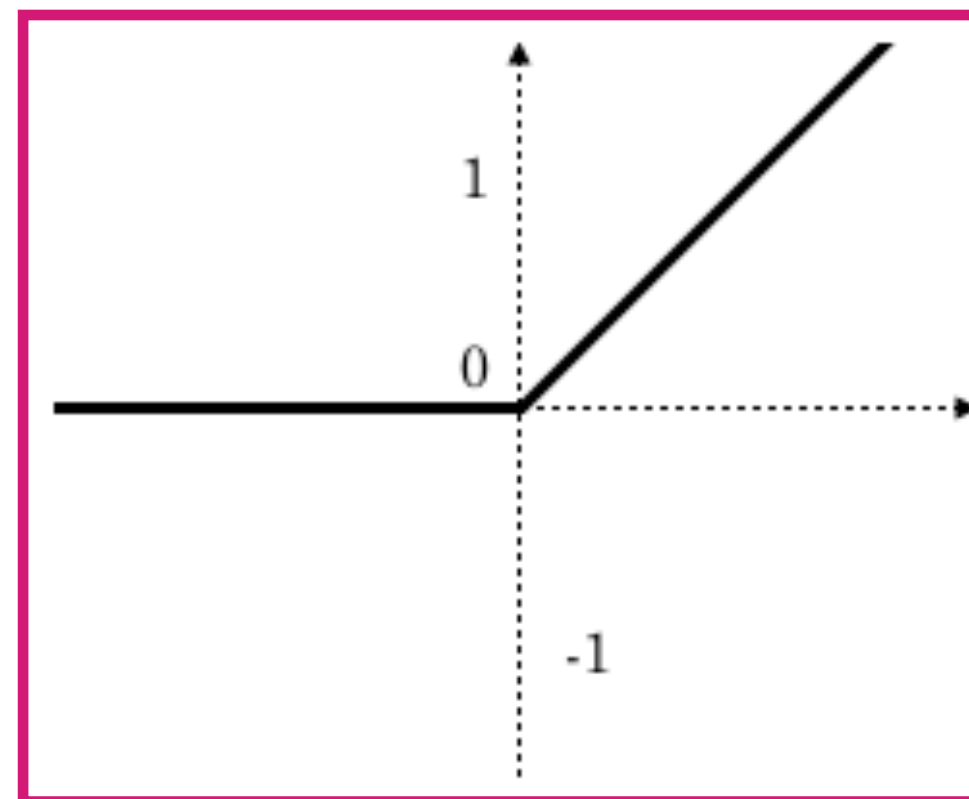


$$\vec{x}_1 = \begin{pmatrix} x_{1,1} \\ x_{1,2} \\ x_{1,3} \\ \vdots \\ x_{1,N_1} \end{pmatrix}$$

$$W_{2,1} = \begin{pmatrix} w_{1,1}^{(2,1)} & w_{2,1}^{(2,1)} & w_{3,1}^{(2,1)} & \dots & \dots & \dots & w_{N_1,1}^{(2,1)} \\ w_{1,2}^{(2,1)} & w_{2,2}^{(2,1)} & w_{3,2}^{(2,1)} & \cdot & \cdot & \cdot & w_{N_1,2}^{(2,1)} \\ w_{1,3}^{(2,1)} & w_{2,3}^{(2,1)} & w_{3,3}^{(2,1)} & & & & w_{N_1,3}^{(2,1)} \\ & \vdots & & \cdot & \cdot & \cdot & \\ & \vdots & & & \cdot & & \\ w_{1,N_2}^{(2,1)} & w_{2,N_2}^{(2,1)} & w_{3,N_2}^{(2,1)} & & & & w_{N_1,N_2}^{(2,1)} \end{pmatrix}$$

$$\vec{b}_2 = \begin{pmatrix} b_{2,1} \\ b_{2,2} \\ b_{2,3} \\ \vdots \\ b_{2,N_2} \end{pmatrix}$$

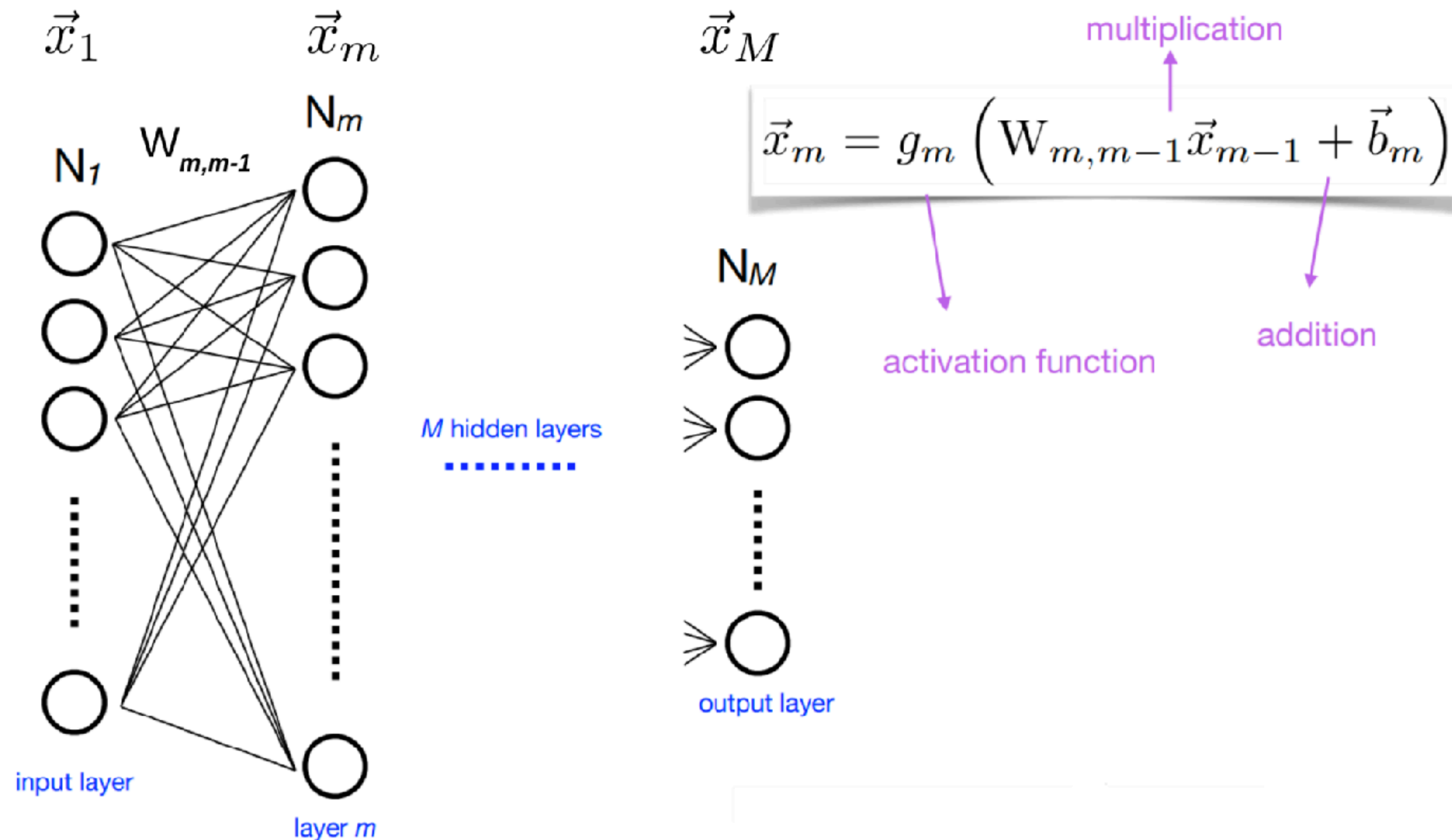
$$g_2(\cdot) =$$



$$\vec{x}_2 = g_2(W_{2,1}\vec{x}_1 + \vec{b}_2)$$

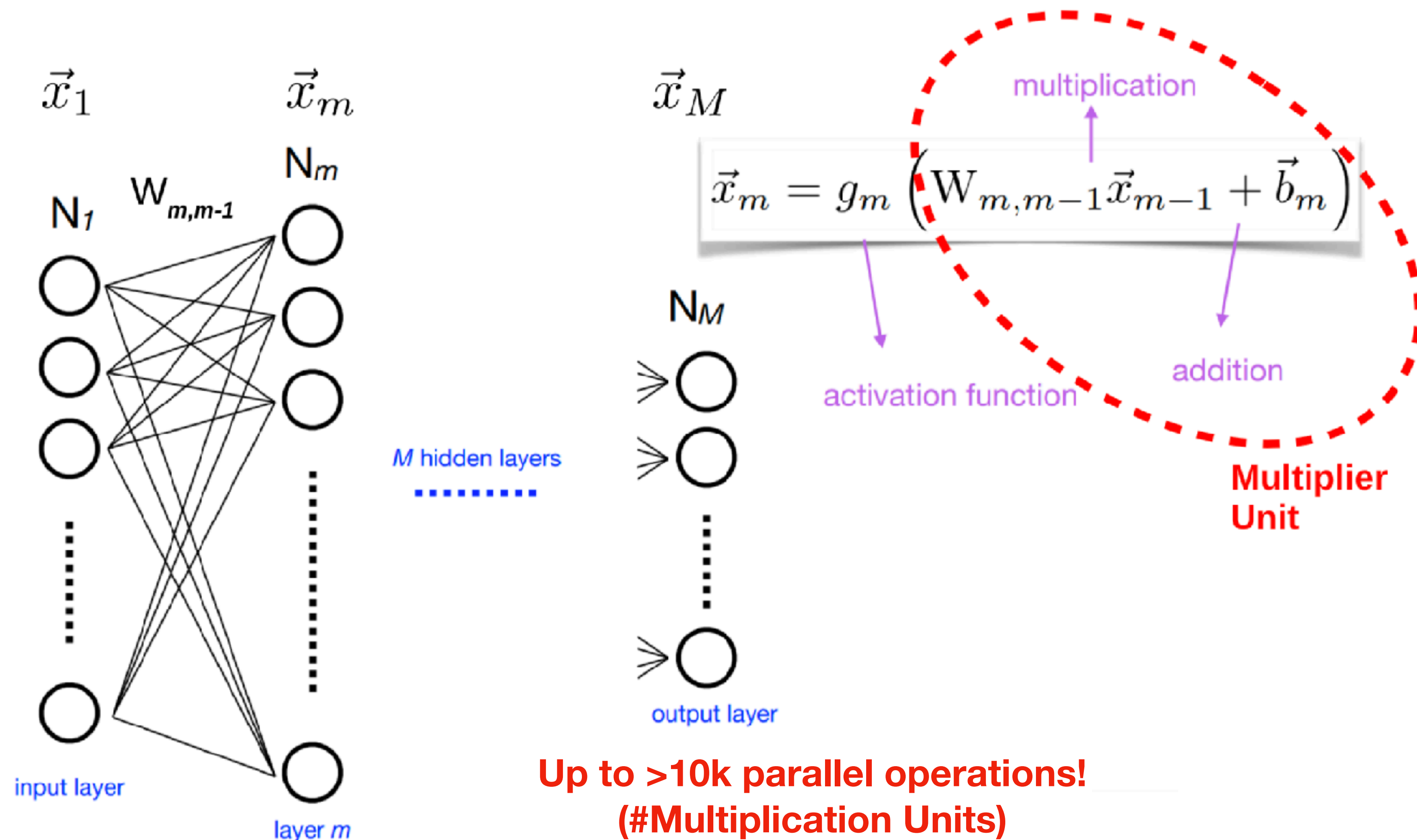


# What is a Neural Network?



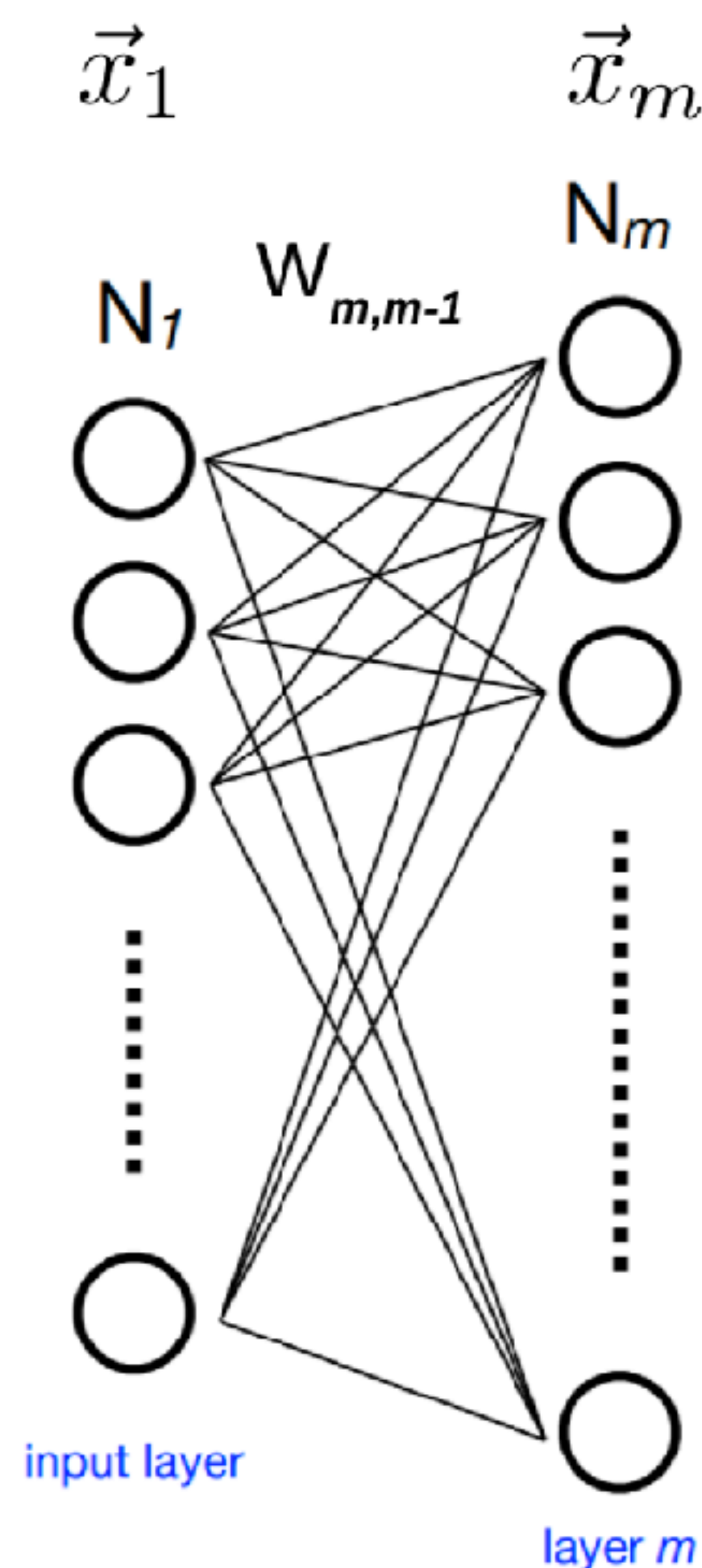


# Inference on FPGAs

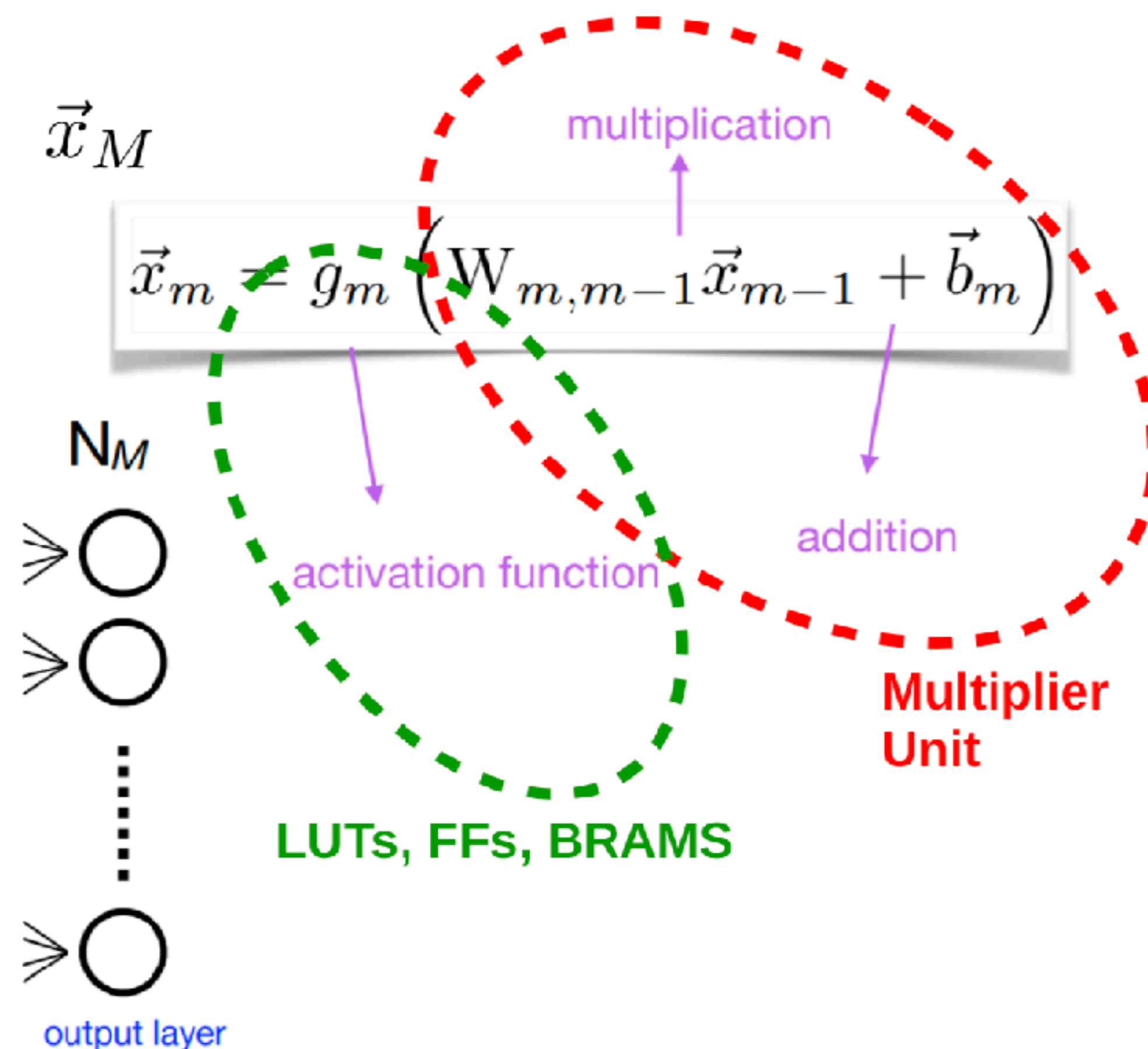




# Inference on FPGAs



$M$  hidden layers  
.....

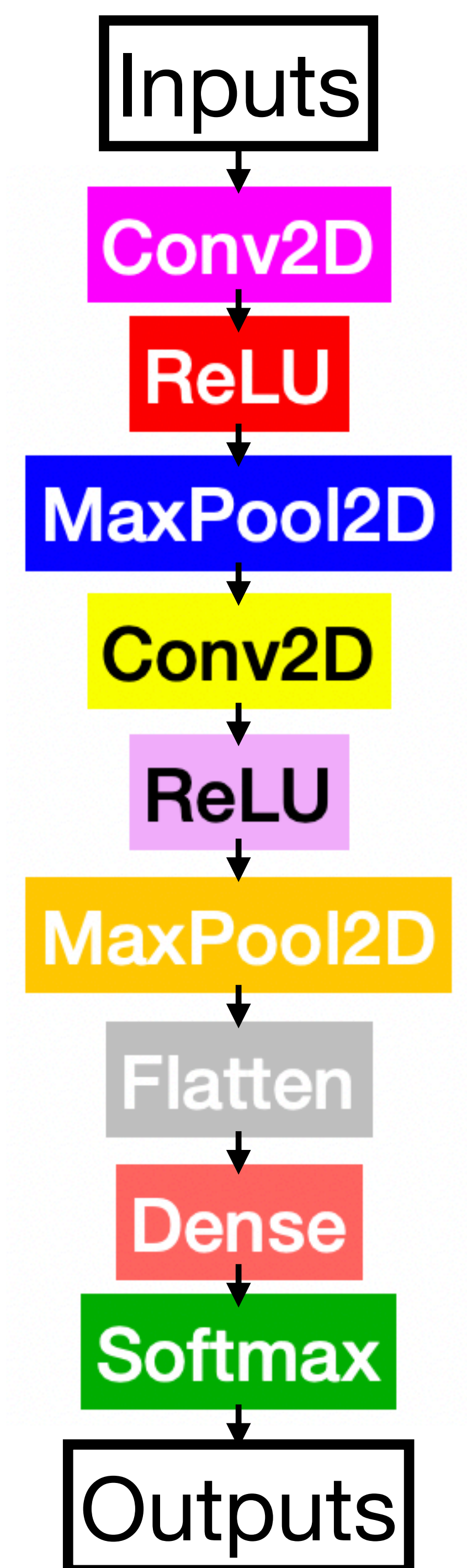
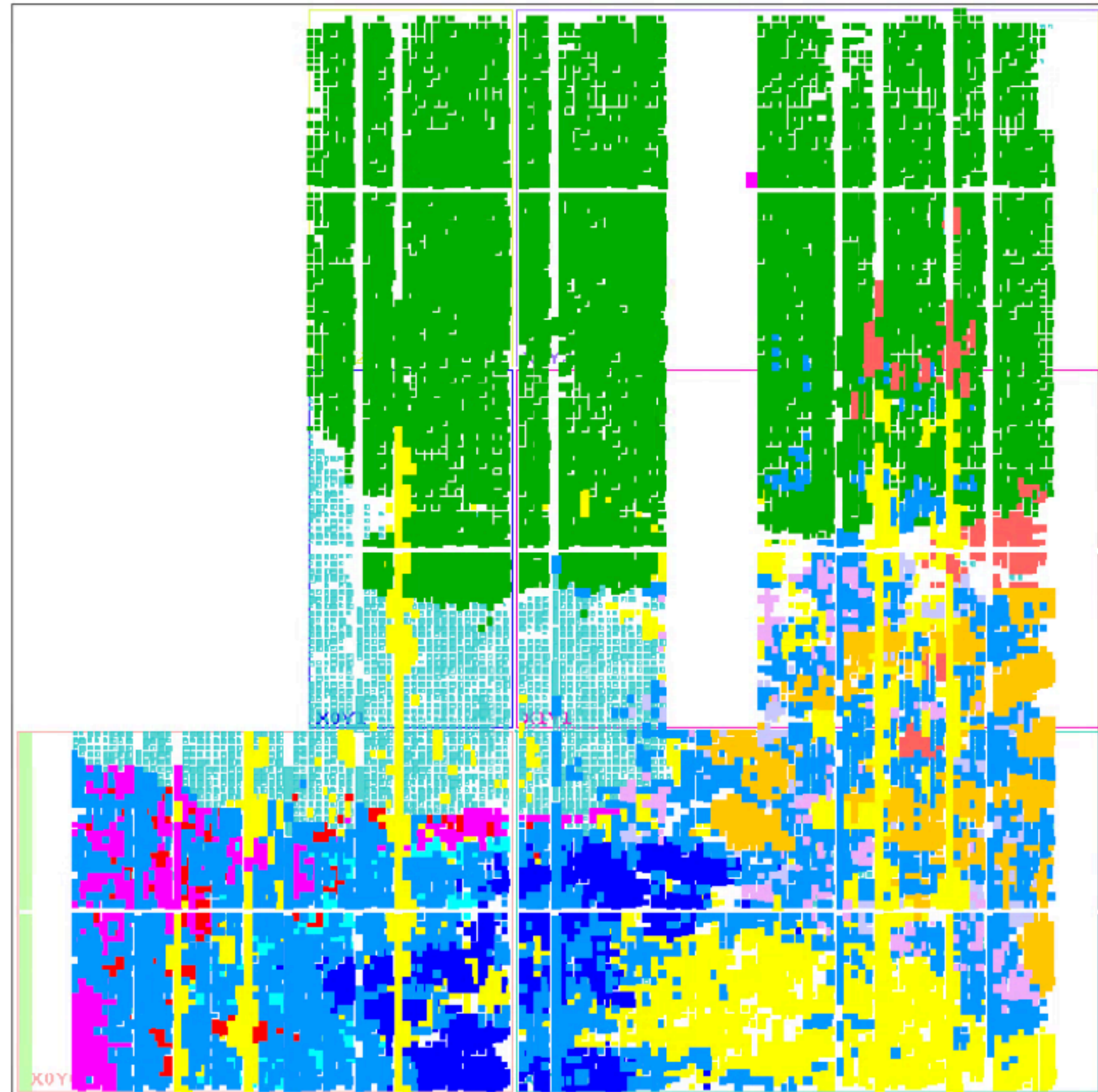


**Up to >10k parallel operations!**  
**(#Multiplication Units)**

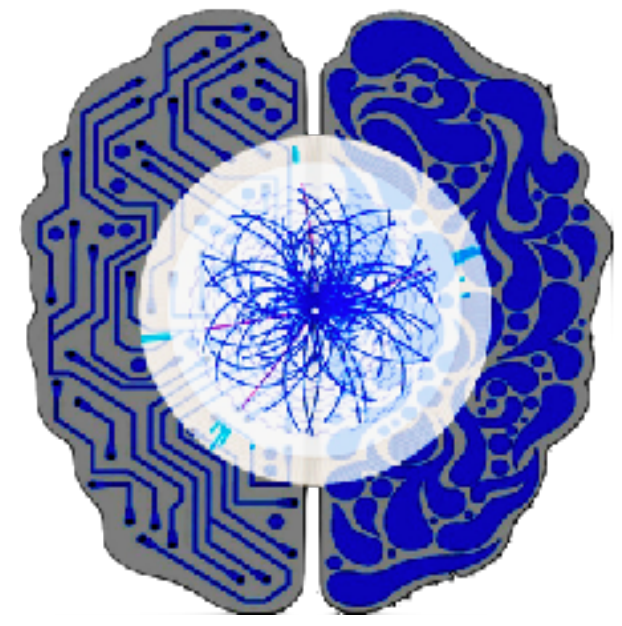


# Inference on FPGAs

- Each part of network must be placed on the FPGA, connected together
- Cannot implement an algorithm if there are no resources left



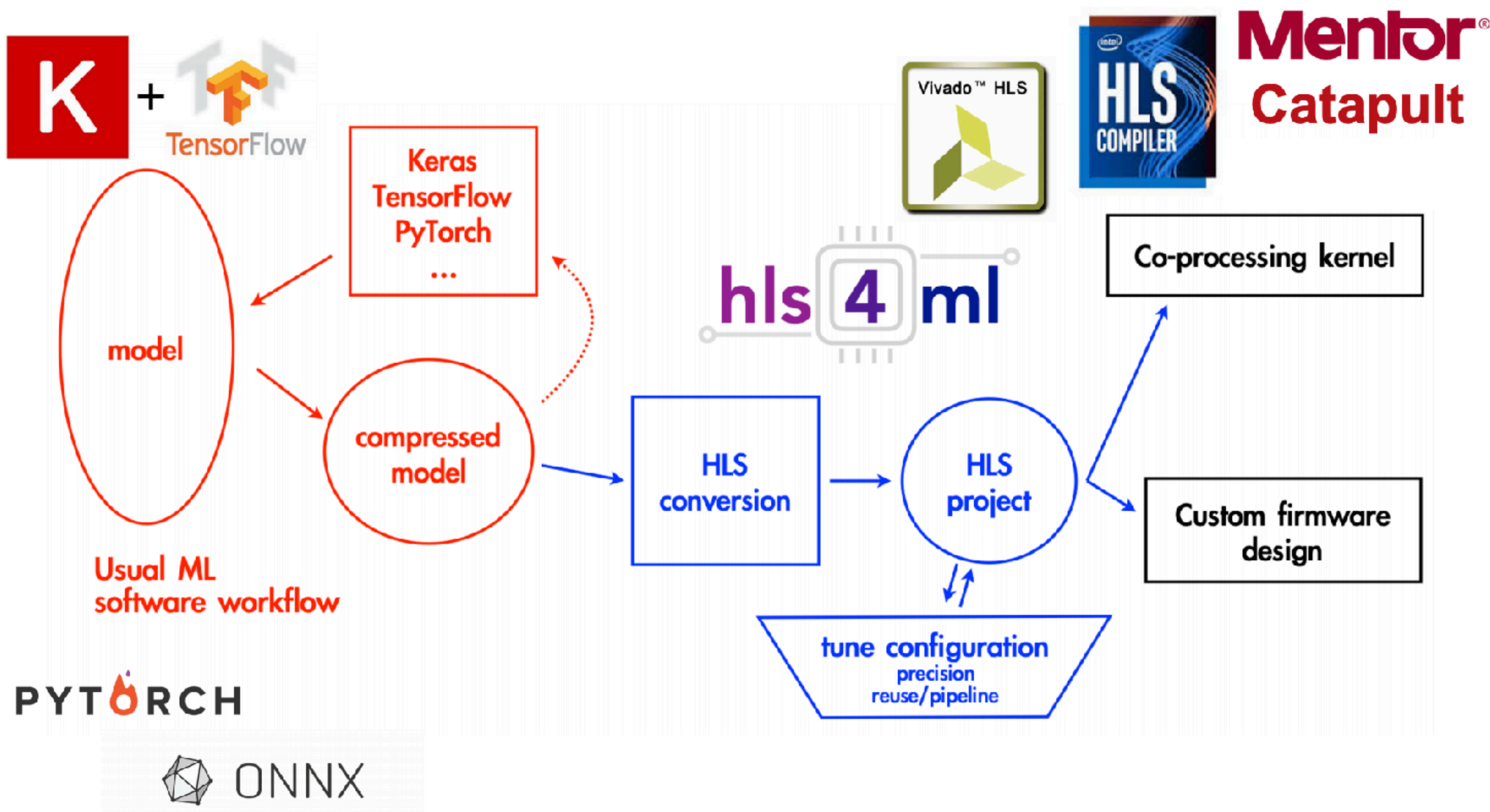




- hls4ml is a software package for automatically creating implementations of neural networks for FPGAs and ASICs
  - <https://fastmachinelearning.org/hls4ml/> [arXiv:1804.06913]
  - pip installable
- Supports common layer architectures and model software (keras, tensorflow, pytorch, ONNX)
- Part of larger Fast Machine Learning collaboration



# hls4ml Workflow





# Many Others

- NNs:

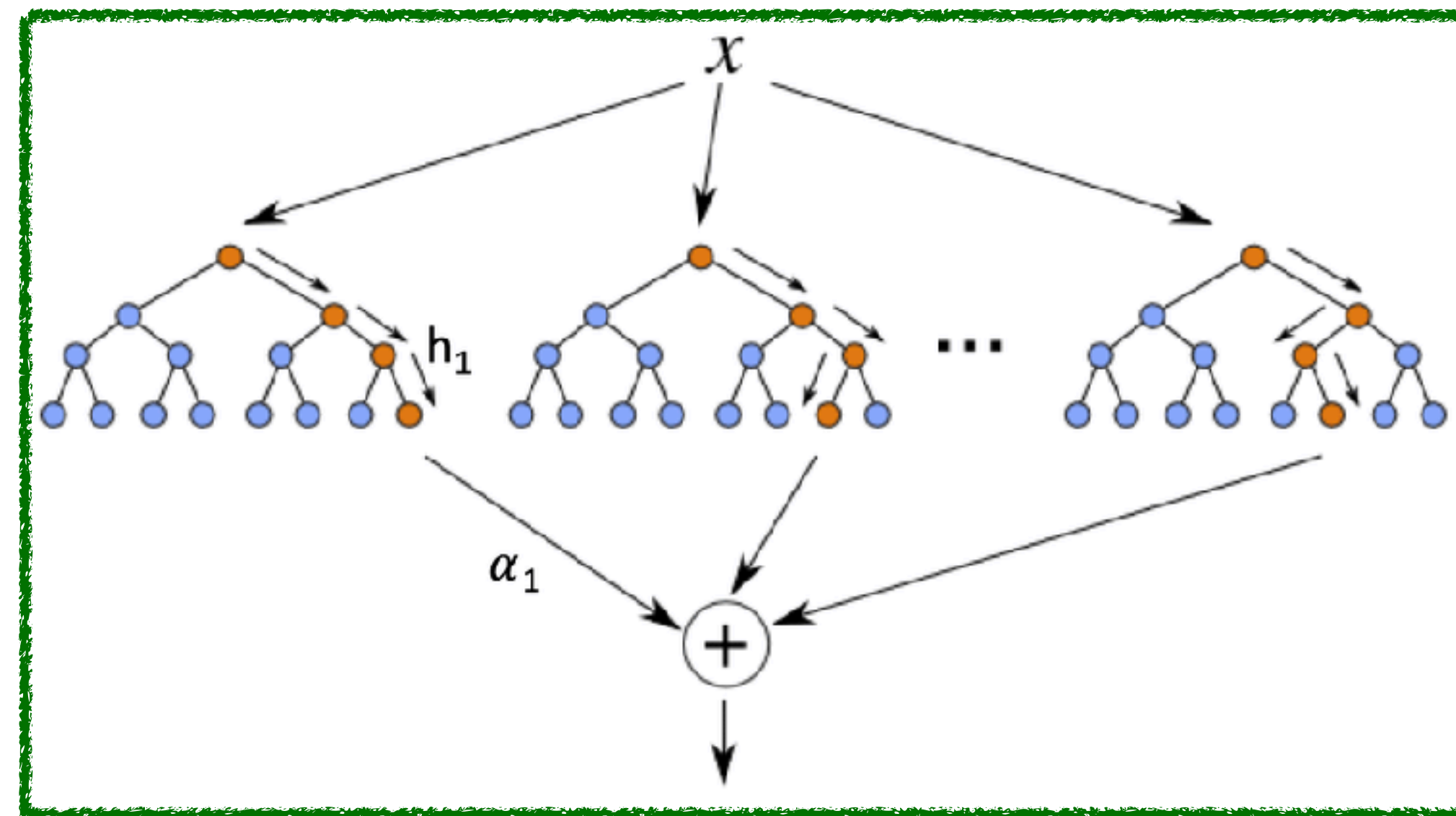
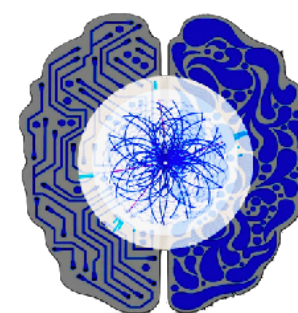


arXiv: 2004.03021

- Boosted Decision Trees (BDTs):



arXiv: 2002.02534



arXiv: 2104.03408

- Entirely non-exhaustive list...



# ML Size / Complexity

- Regardless of toolkit, limitation of doing low latency ML is FPGA size
  - Bigger FPGA → more resources → more computation

## Xilinx Virtex Ultrascale+ VU13P

12288 Multipliers

1.7M LUTs

3.4M FFs

95 Mb BRAM



- *Pruning* and *quantization* are ways to reduce resources
  - Challenge is maintaining performance



# Pruning

- Are all the pieces a given network necessary?

- M

- M

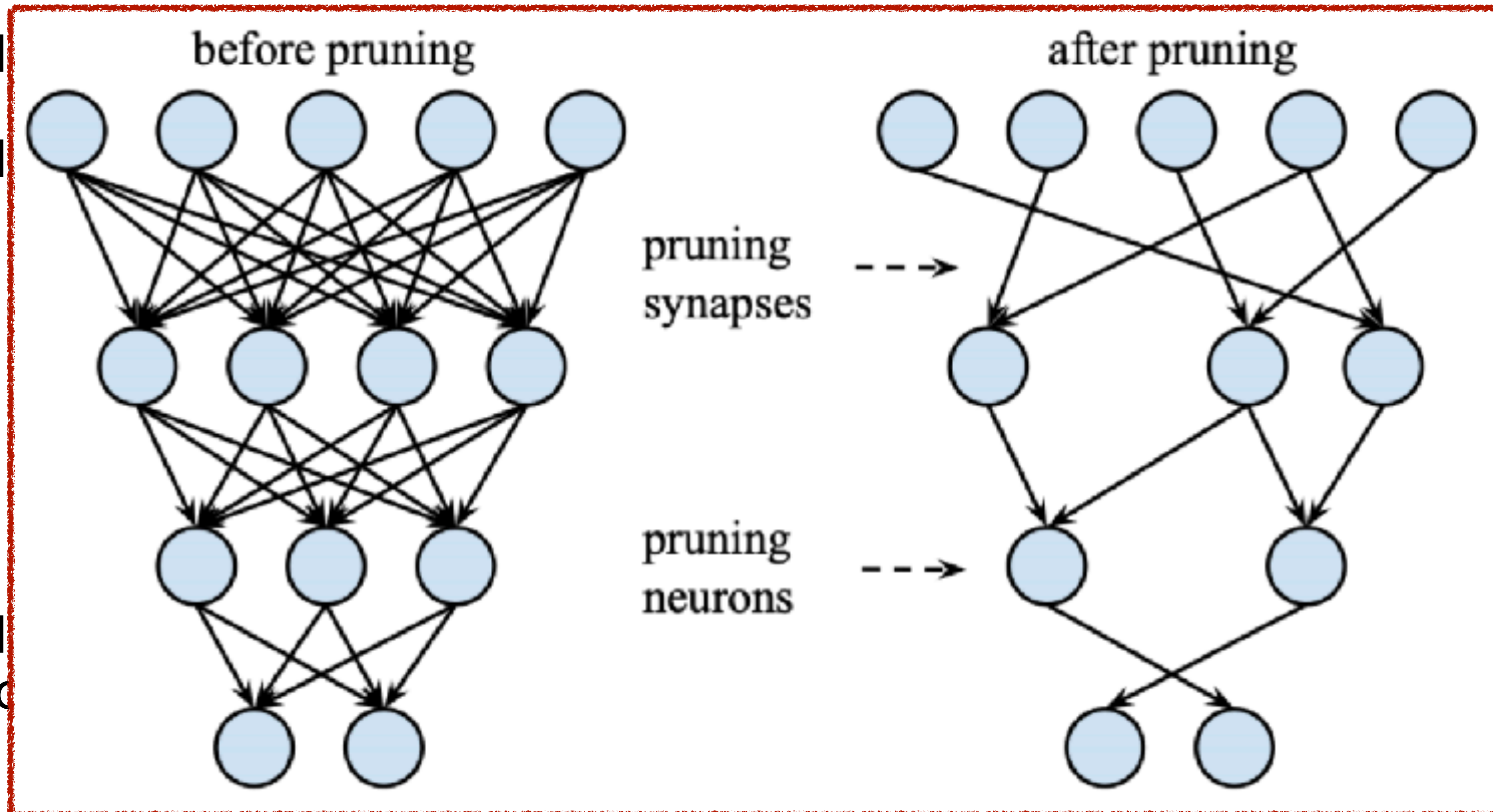
- 

- 

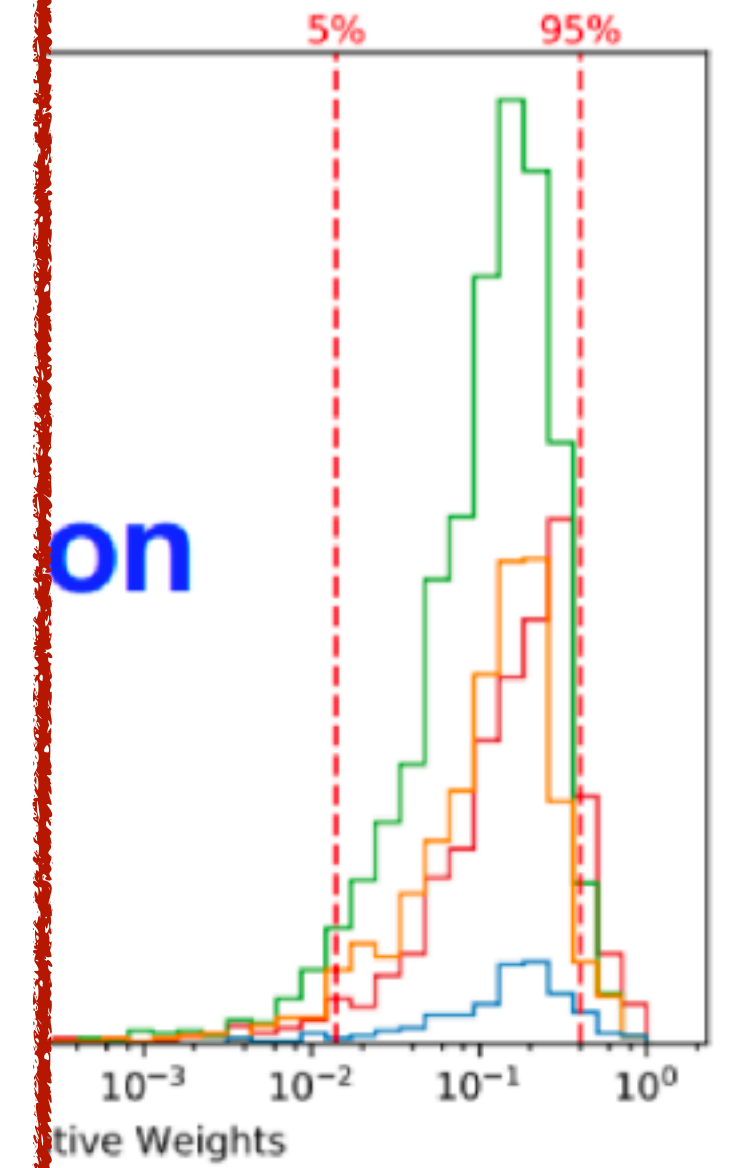
- 

- M

- fro



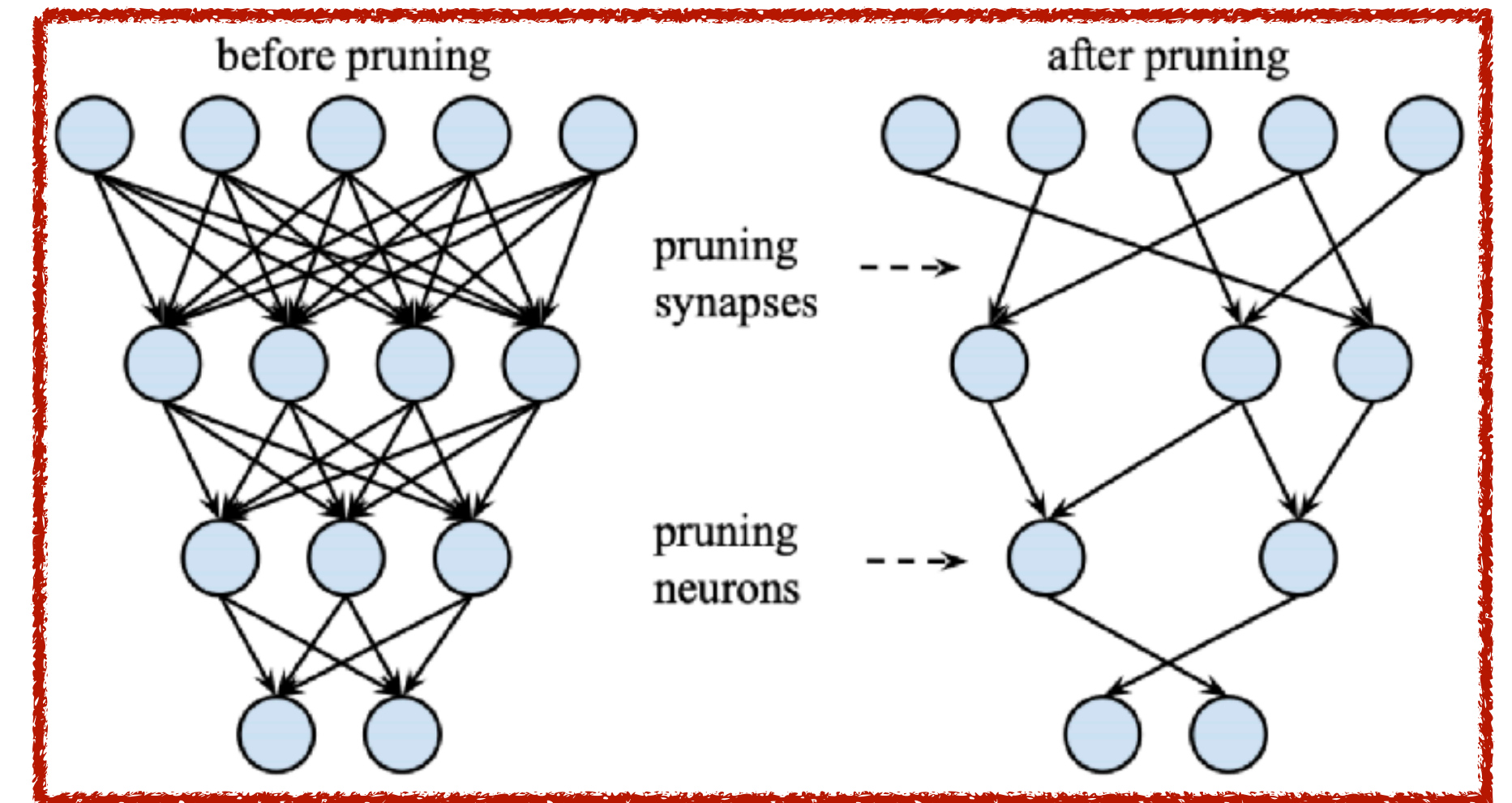
$$v) + \lambda ||\mathbf{w}||$$



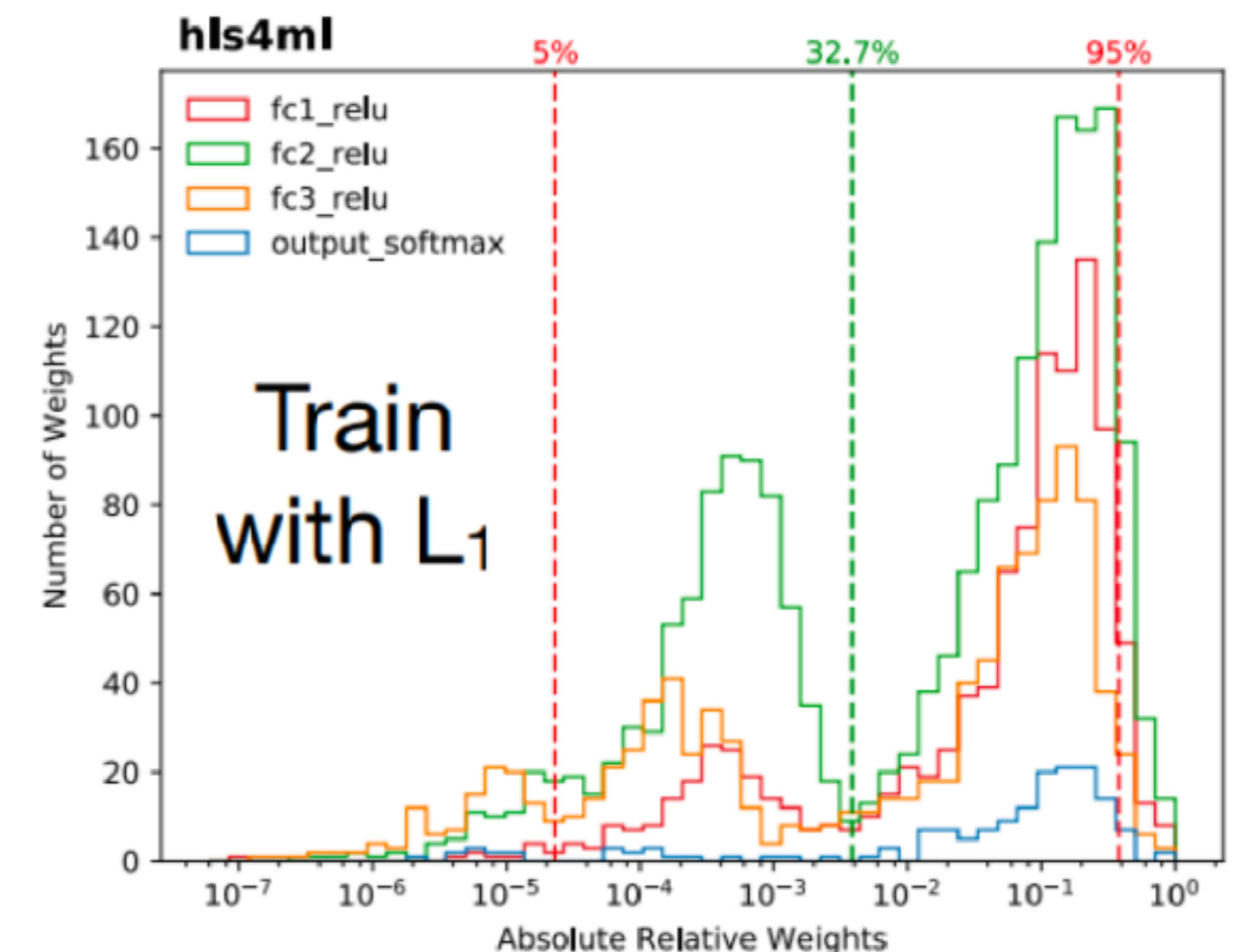


# Pruning

- Are all the pieces a given network necessary?
- Many different types of pruning
- Magnitude-based:
  - Use **regularization** (penalty term in loss function for large weights)
  - Remove smallest weights
  - Repeat
- Multiplications by 0 can be completely removed from FPGA design



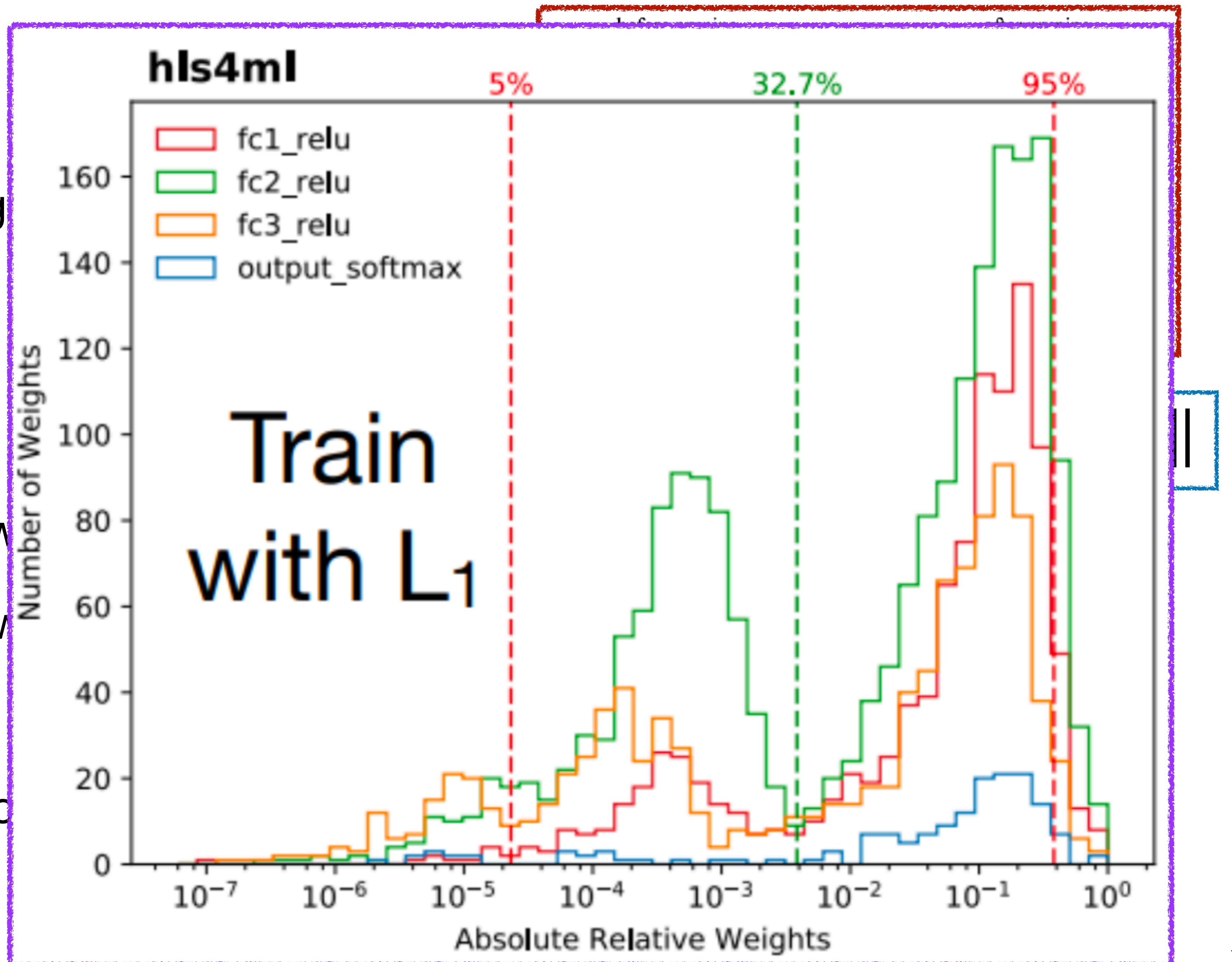
$$L_{\lambda}(\mathbf{w}) = L(\mathbf{w}) + \lambda \|\mathbf{w}\|$$





# Pruning

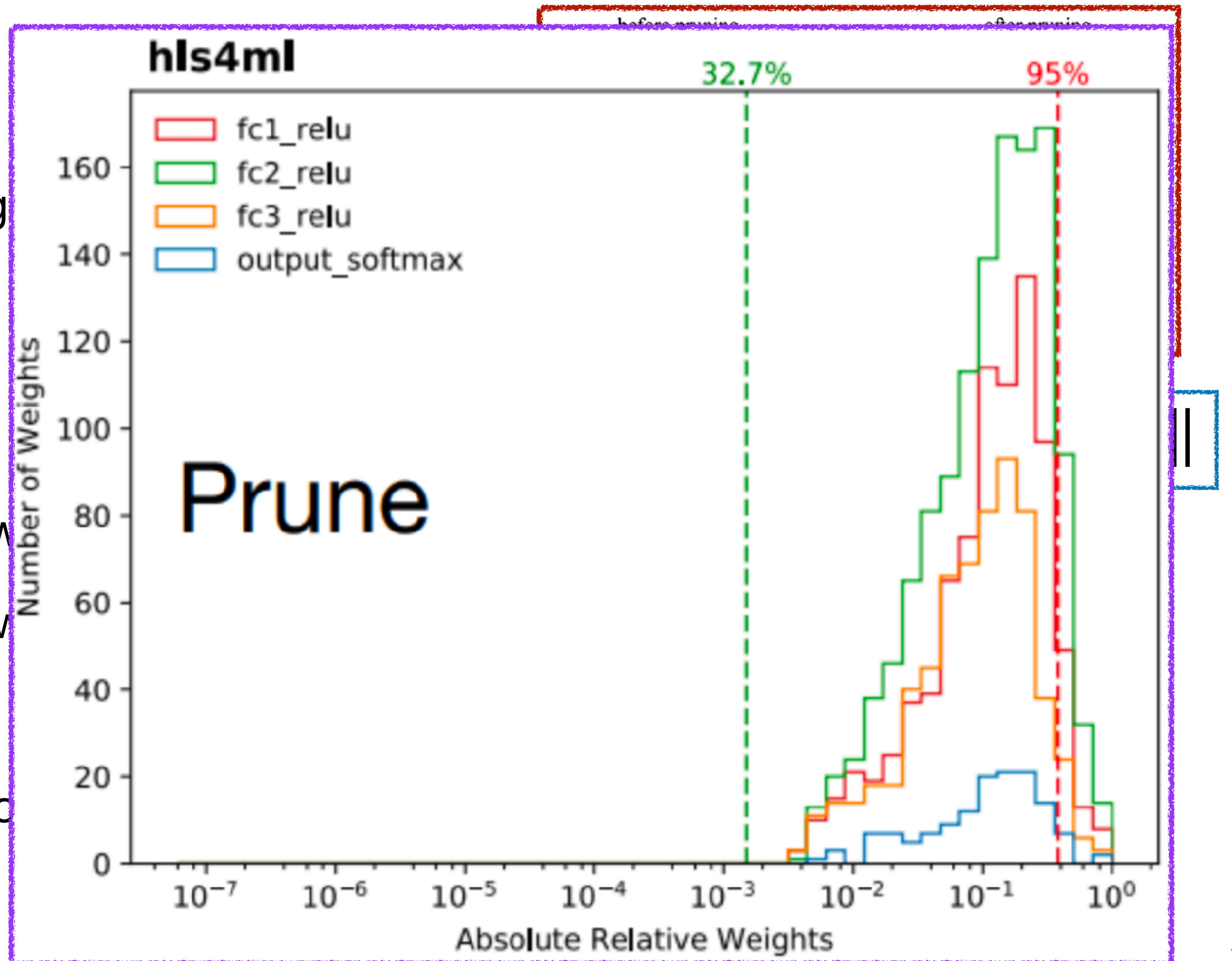
- Are all the pieces a good fit?
- Many different types
- Magnitude-based:
  - Use [regularization](#) function for large weights
  - Remove smallest weights
  - Repeat
- Multiplications by 0 can be removed from FPGA design





# Pruning

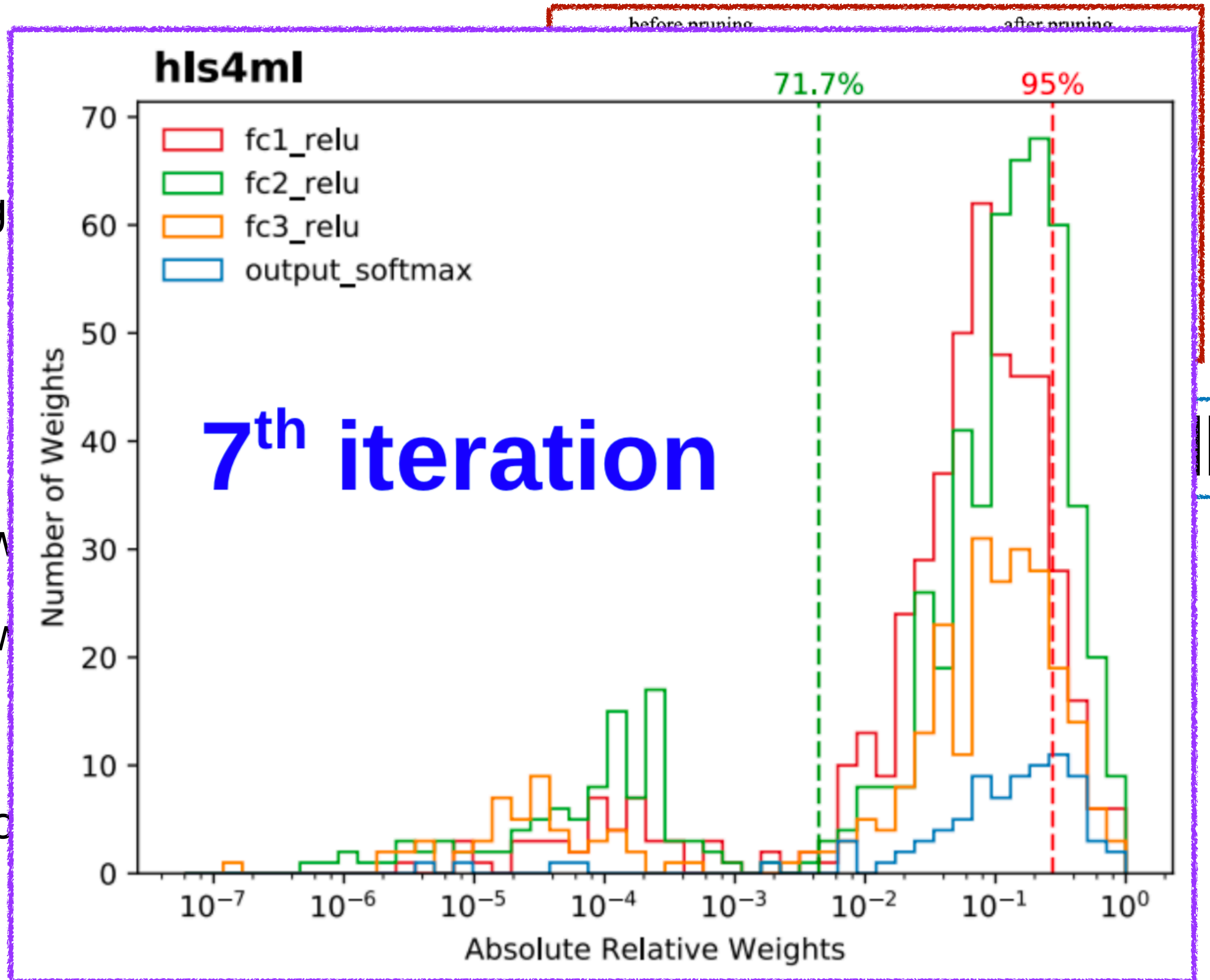
- Are all the pieces a good fit?
- Many different types
- Magnitude-based:
  - Use **regularization** function for large weights
  - Remove smallest weights
  - Repeat
- Multiplications by 0 can be removed from FPGA design





# Pruning

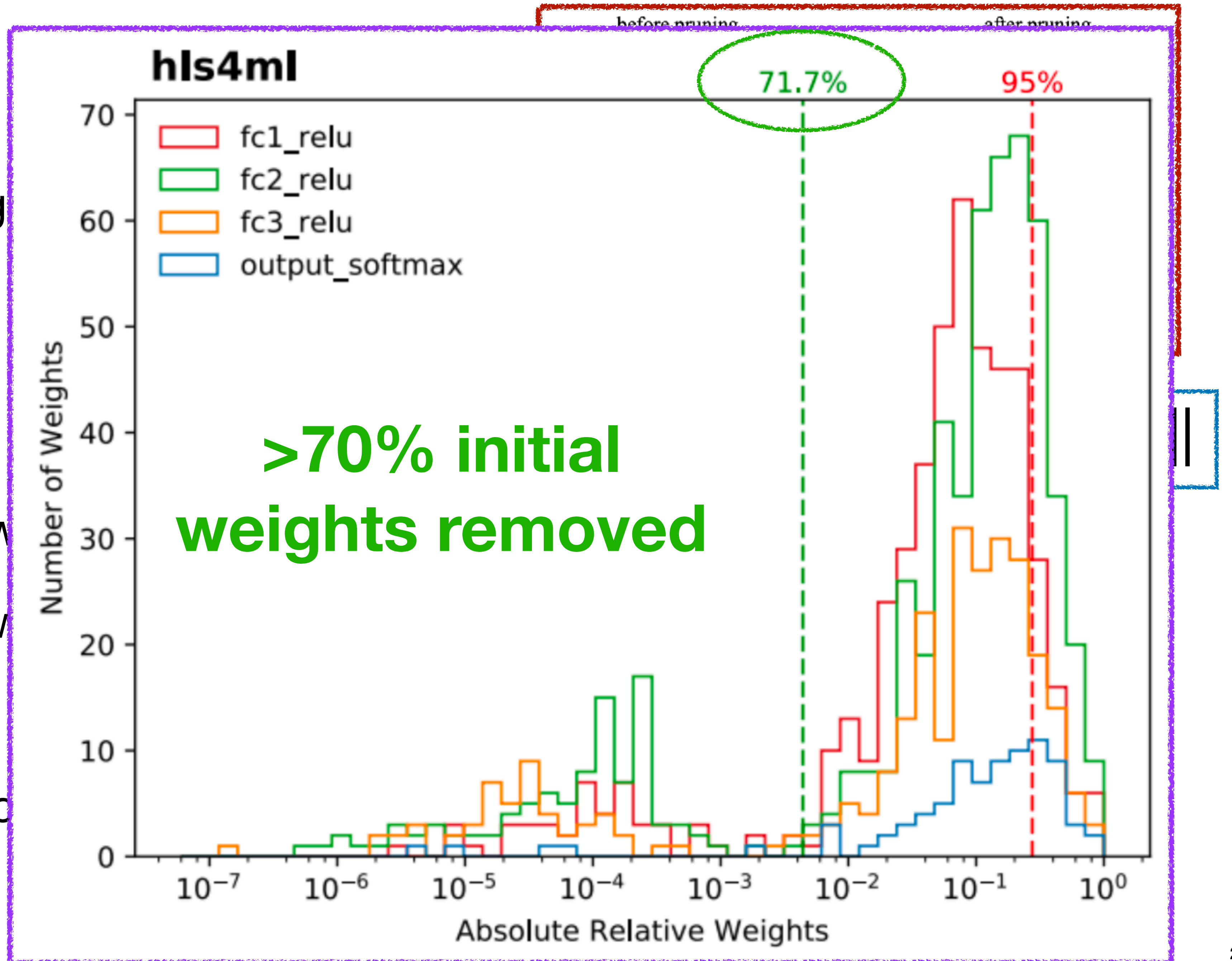
- Are all the pieces a good fit?
- Many different types
- Magnitude-based:
  - Use **regularization** function for large weights
  - Remove smallest weights
  - Repeat
- Multiplications by 0 can be removed from FPGA design





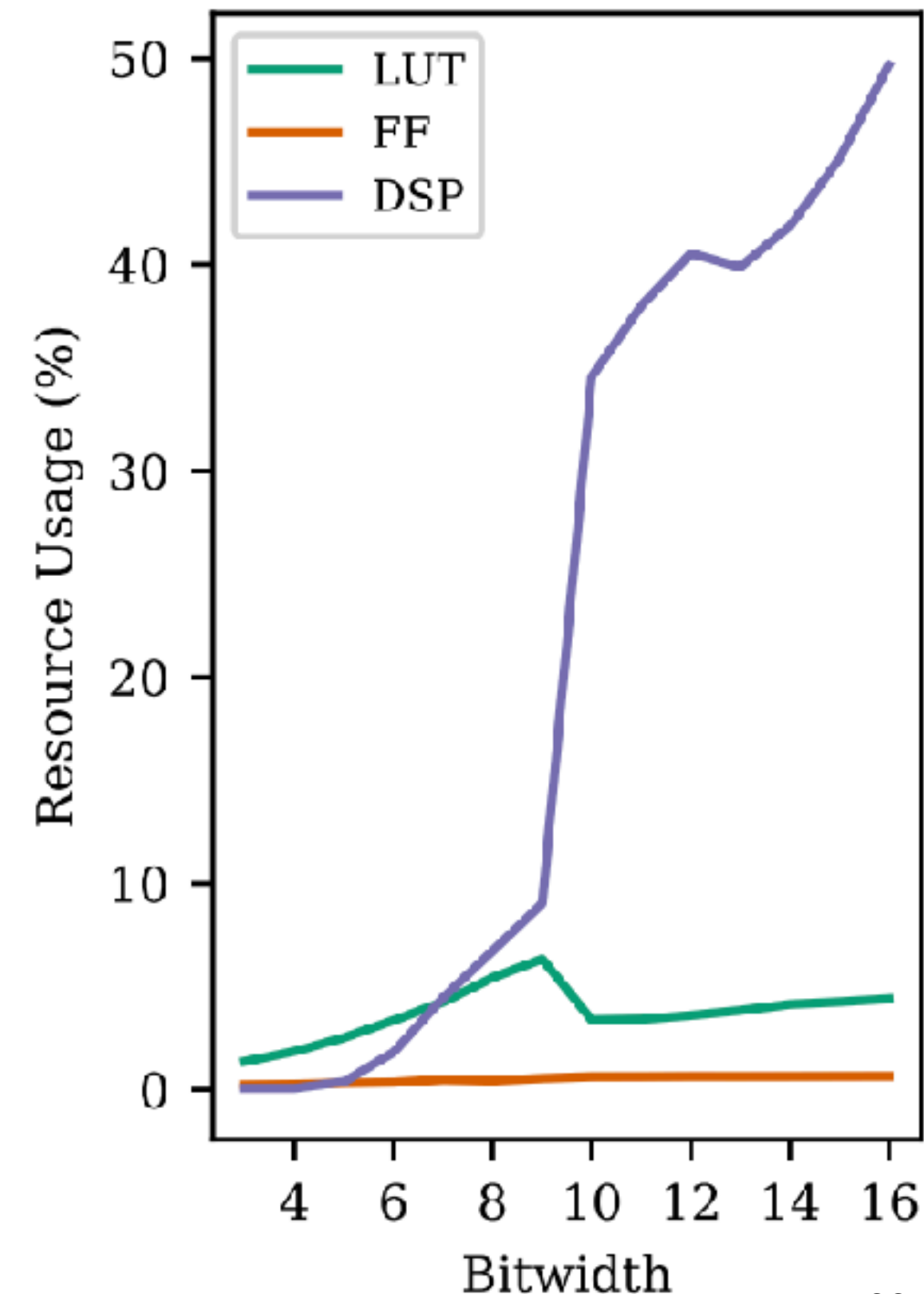
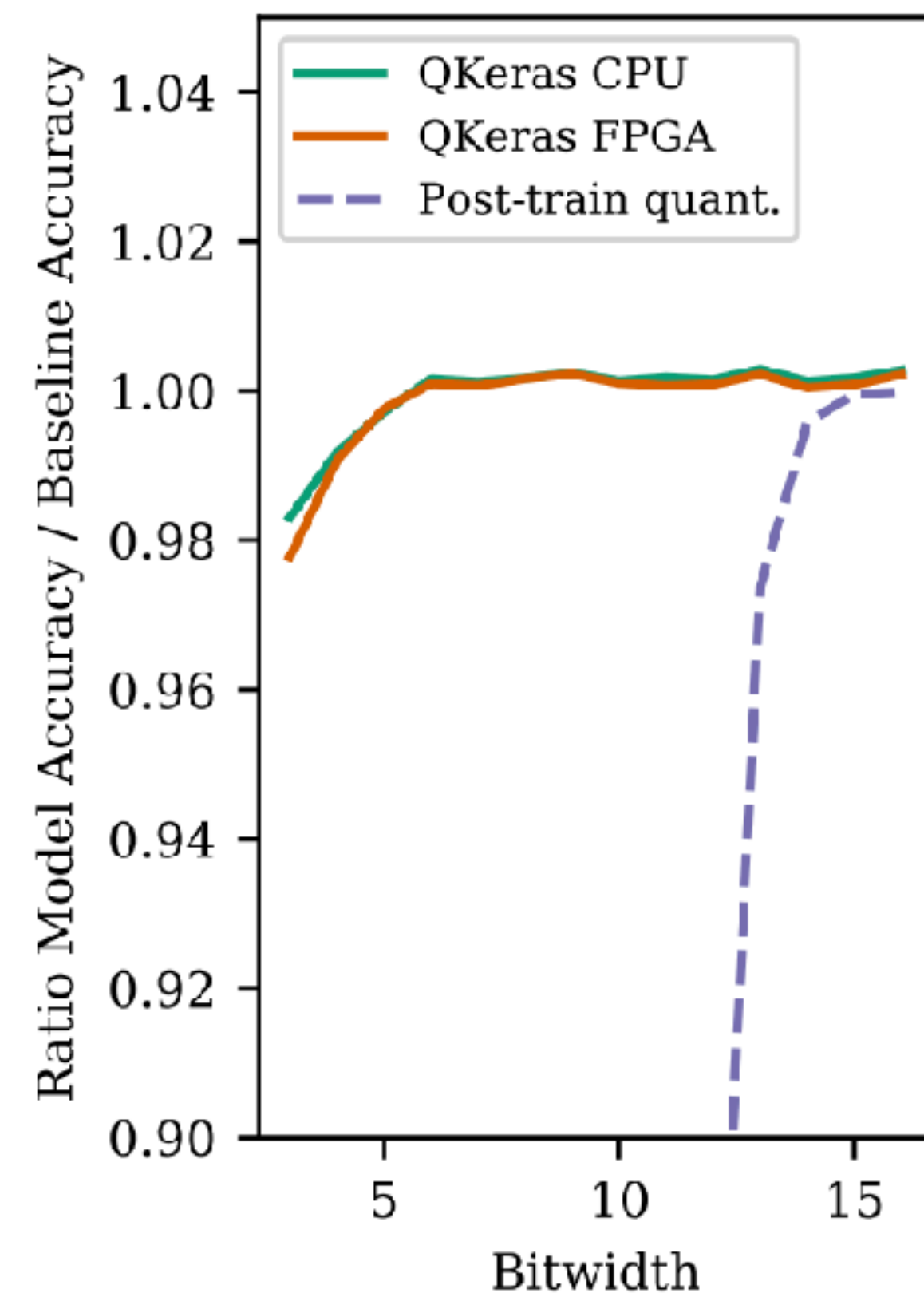
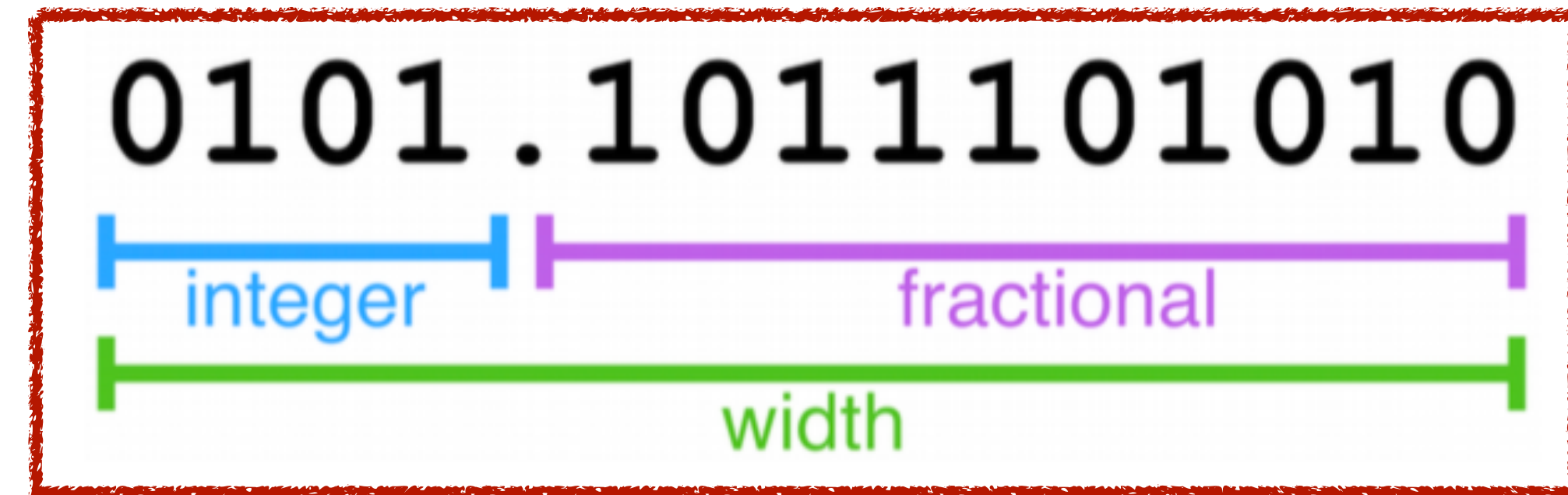
# Pruning

- Are all the pieces a good fit?
- Many different types
- Magnitude-based:
  - Use **regularization** function for large weights
  - Remove smallest weights
  - Repeat
- Multiplications by 0 can be removed from FPGA design



# Quantization

- FPGAs are well suited to fixed-point numbers, not floating point
- Bitwidth can be adjusted as needed (impacts accuracy, performance, resources)
  - Can be combined with other customizations
- Quantization-aware training [[arXiv:2006.10159](#)]
  - Can greatly reduce size of network by training with knowledge of quantization





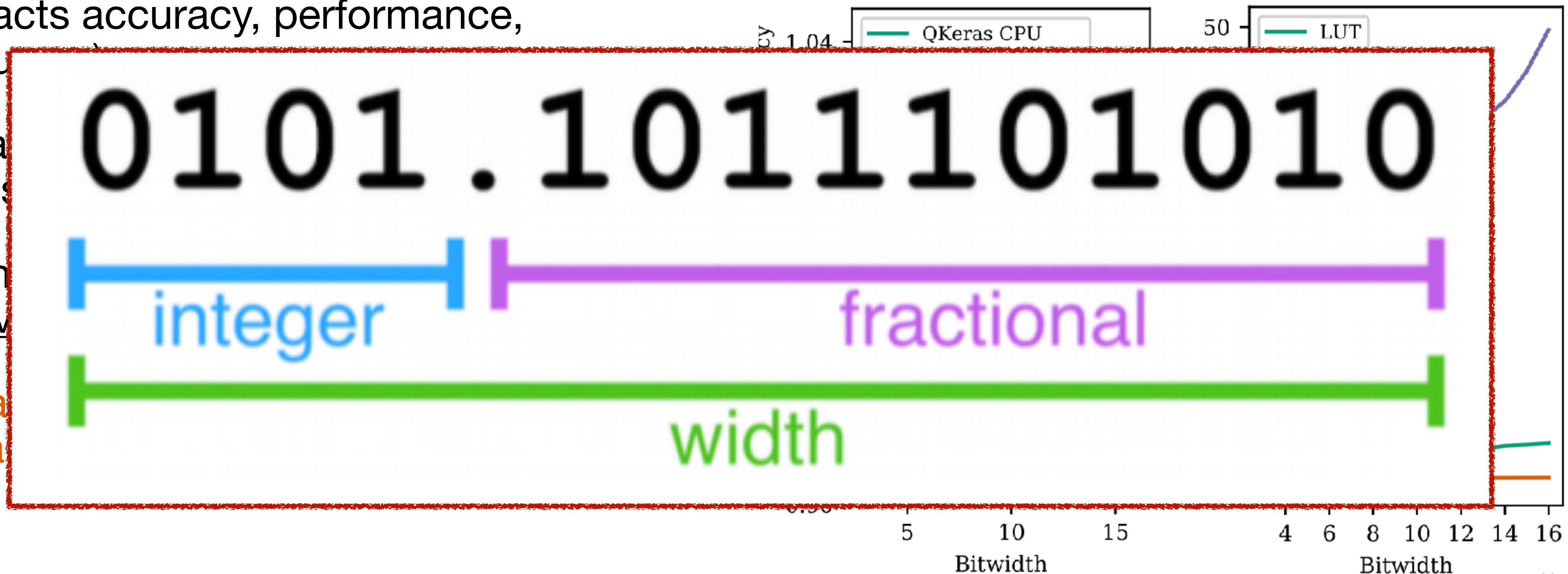
# Quantization

- FPGAs are well suited to fixed-point numbers, not floating point
- Bitwidth can be adjusted as needed (impacts accuracy, performance, resou

- Ca  
cus

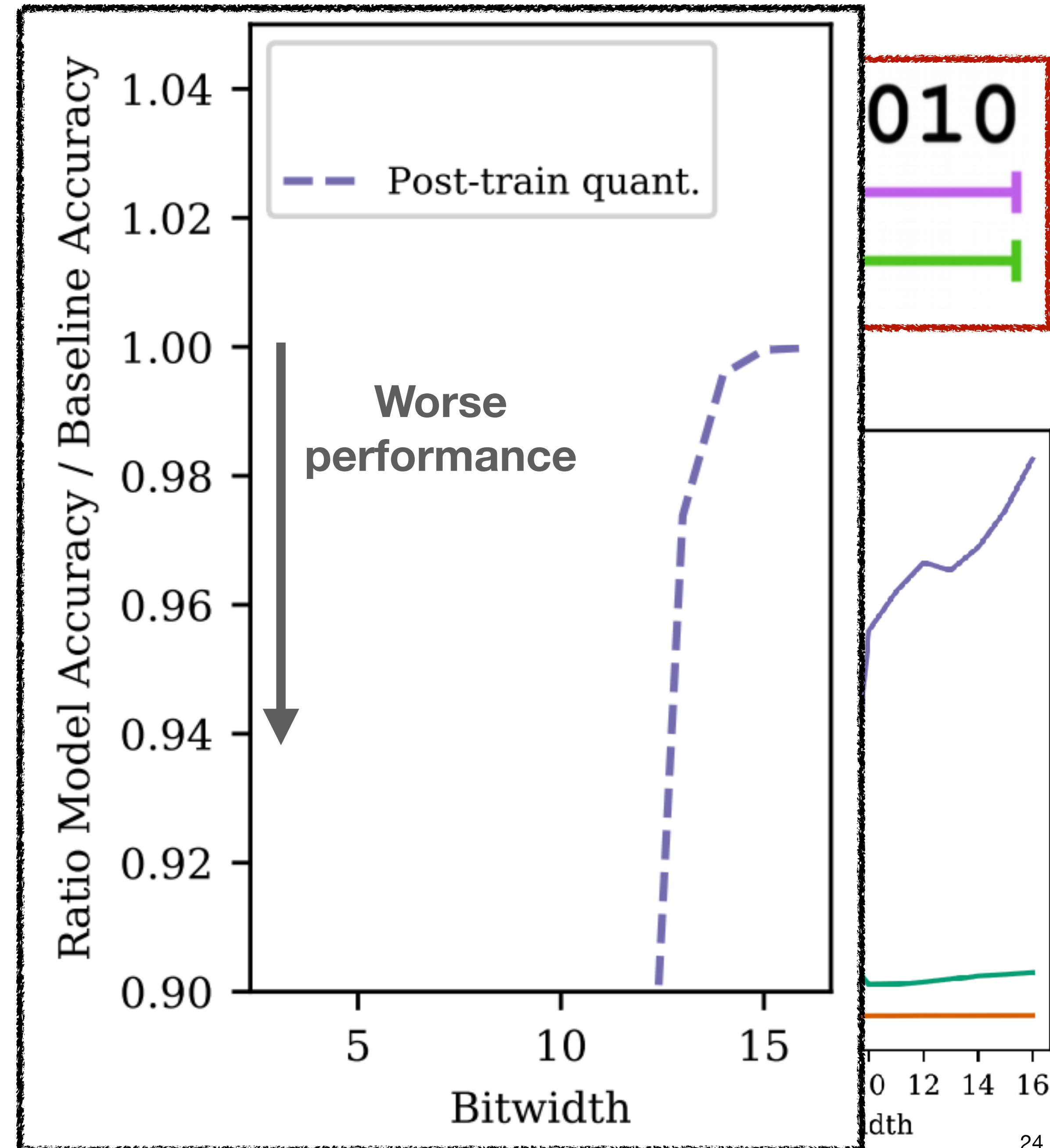
- Quan  
[arXiv

- Ca  
tra



# Quantization

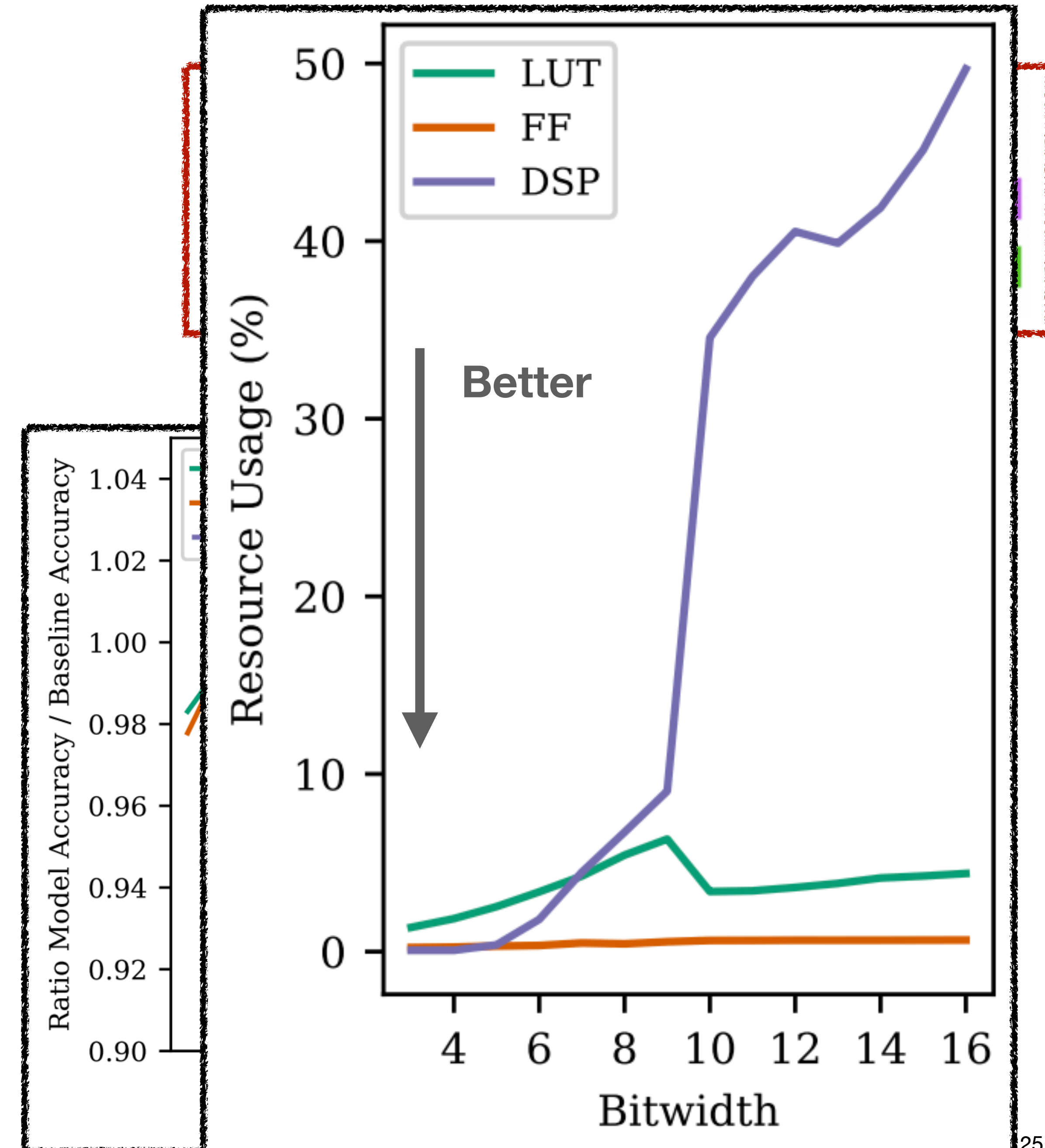
- FPGAs are well suited to fixed-point numbers, not floating point
- Bitwidth can be adjusted as needed (impacts accuracy, **performance**, resources)
  - Can be combined with other customizations
- Quantization-aware training [[arXiv:2006.10159](#)]
  - Can greatly reduce size of network by training with knowledge of quantization





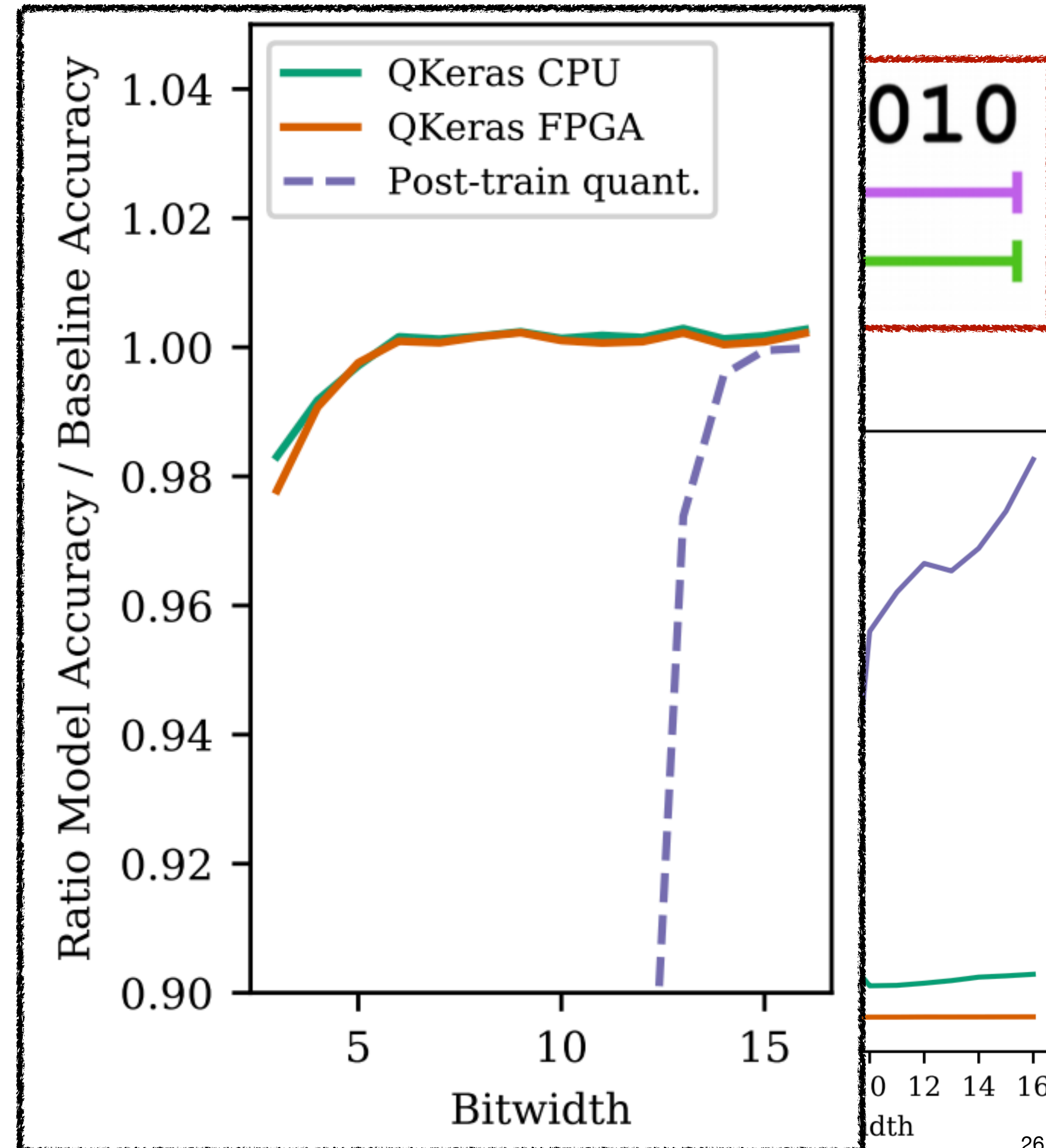
# Quantization

- FPGAs are well suited to fixed-point numbers, not floating point
- Bitwidth can be adjusted as needed (impacts accuracy, performance, **resources**)
  - Can be combined with other customizations
- Quantization-aware training [[arXiv:2006.10159](#)]
  - Can greatly reduce size of network by training with knowledge of quantization



# Quantization

- FPGAs are well suited to fixed-point numbers, not floating point
- Bitwidth can be adjusted as needed (impacts accuracy, performance, resources)
  - Can be combined with other customizations
- **Quantization-aware training** [[arXiv:2006.10159](#)]
  - Can greatly reduce size of network by training with knowledge of quantization

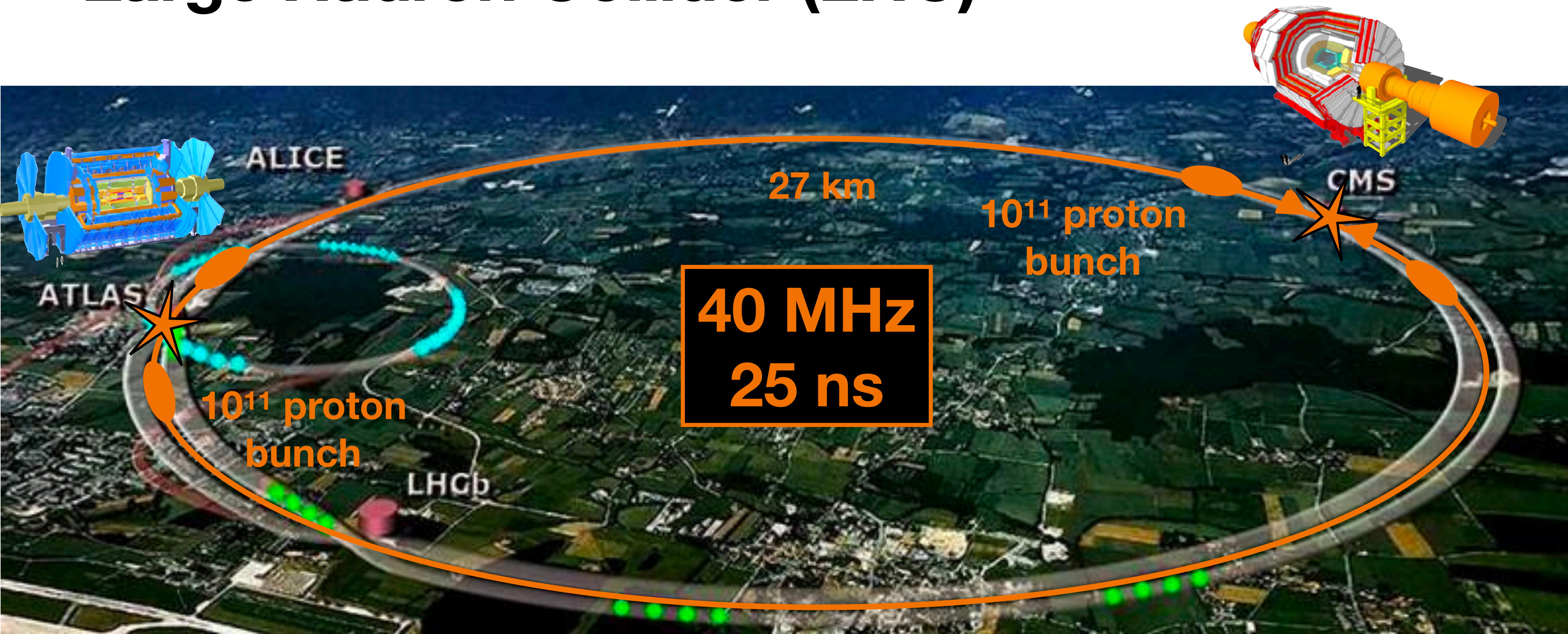




# **An Application**

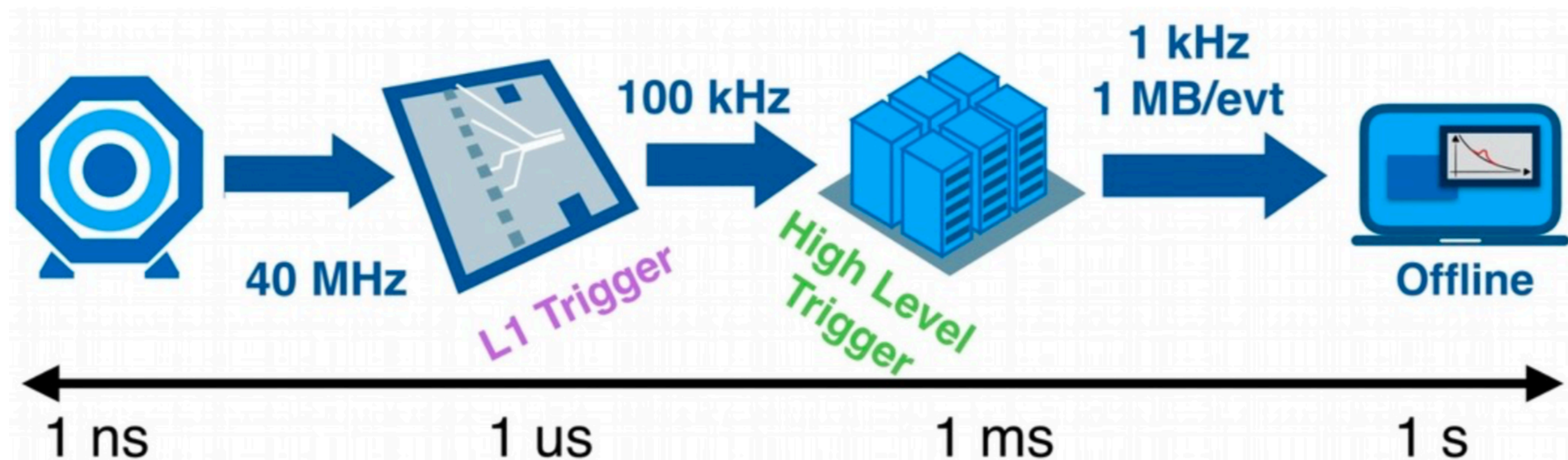


# Large Hadron Collider (LHC)



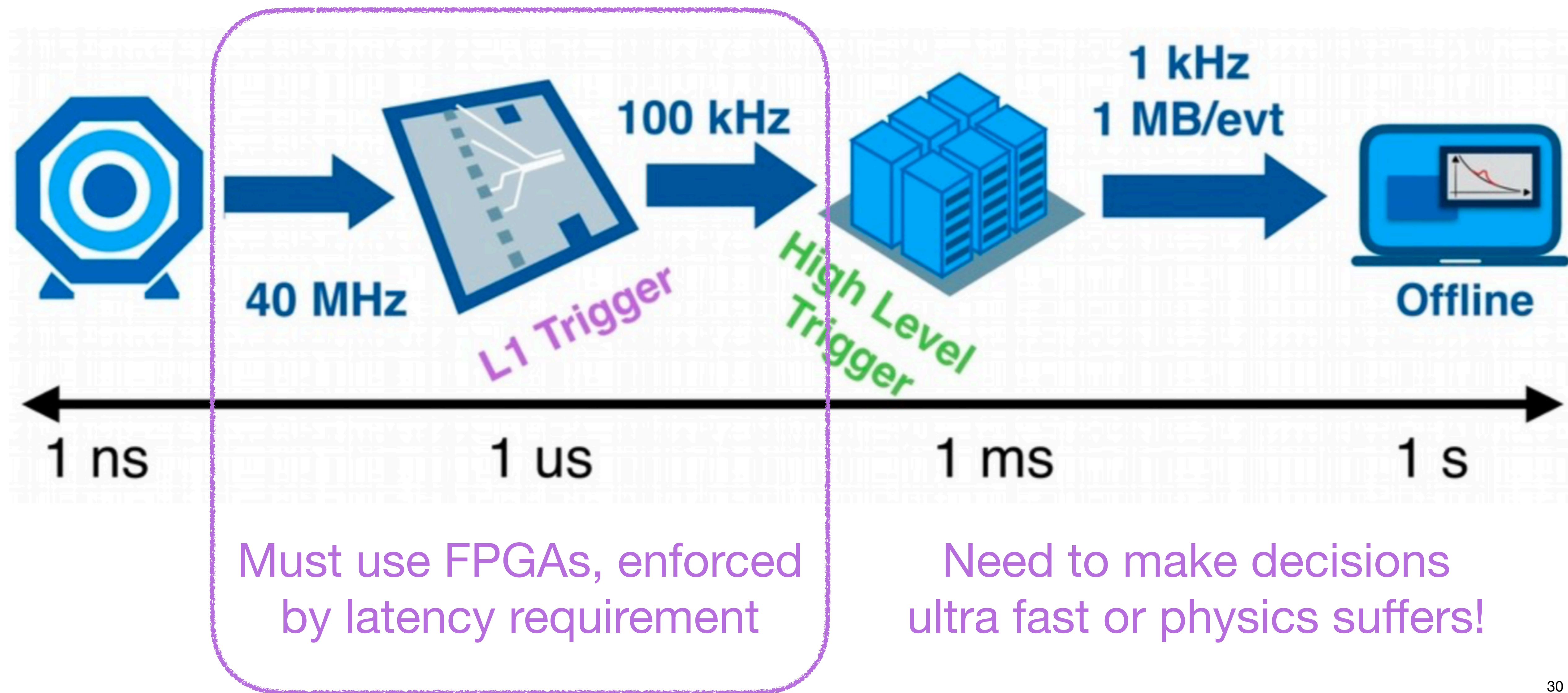


# LHC Data Processing / Readout





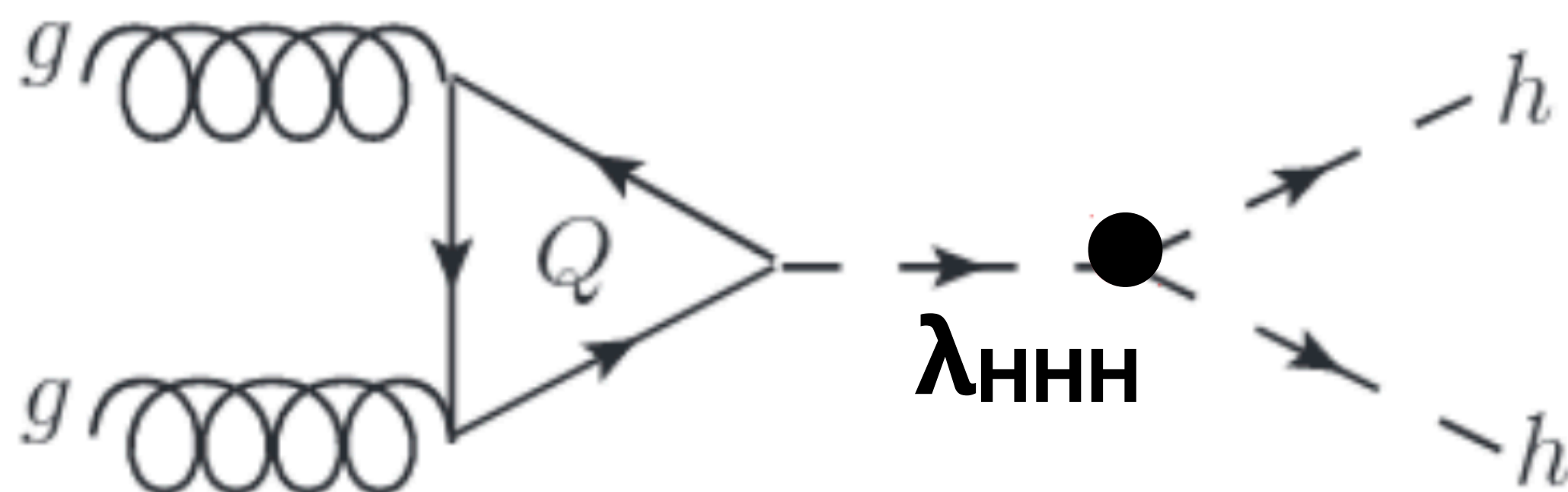
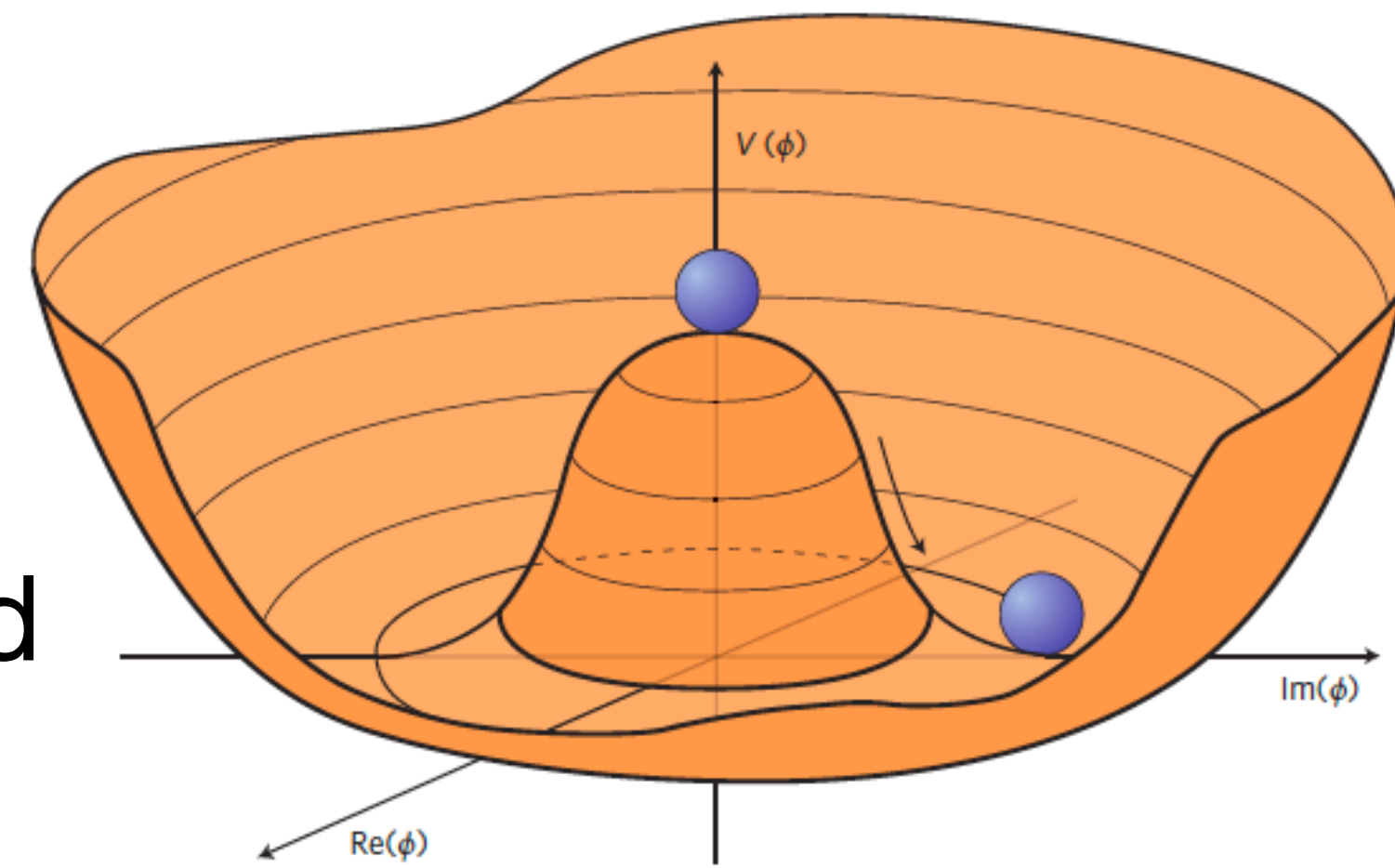
# LHC Data Processing / Readout





# Di-Higgs Production

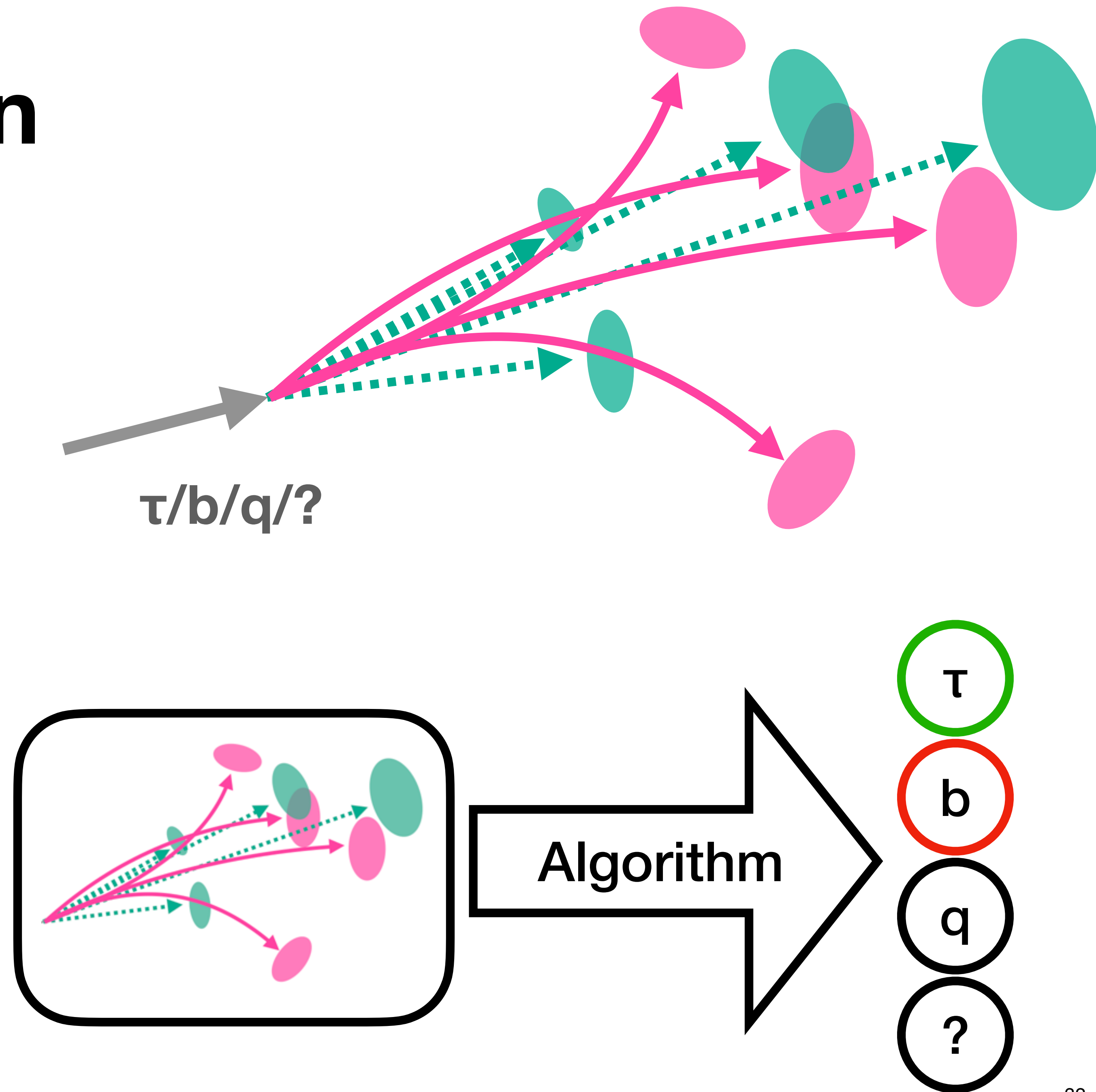
- HH is best way to measure scalar potential of Higgs field
  - Higgs self-coupling:  $\lambda_{HHH}$
- Major target for CMS and ATLAS experiments
- **bbbb**, **bbWW**, **bb $\tau\tau$**  have largest branching ratios
- Lots of background that mimics these signals  $\rightarrow$  very difficult to record



	bb	WW	$\tau\tau$	ZZ	$\gamma\gamma$
bb	33%				
WW	25%	4.6%			
$\tau\tau$	7.4%	2.5%	0.39%		
ZZ	3.1%	1.2%	0.34%	0.076%	
$\gamma\gamma$	0.26%	0.10%	0.029%	0.013%	0.0005%

# Particle Identification

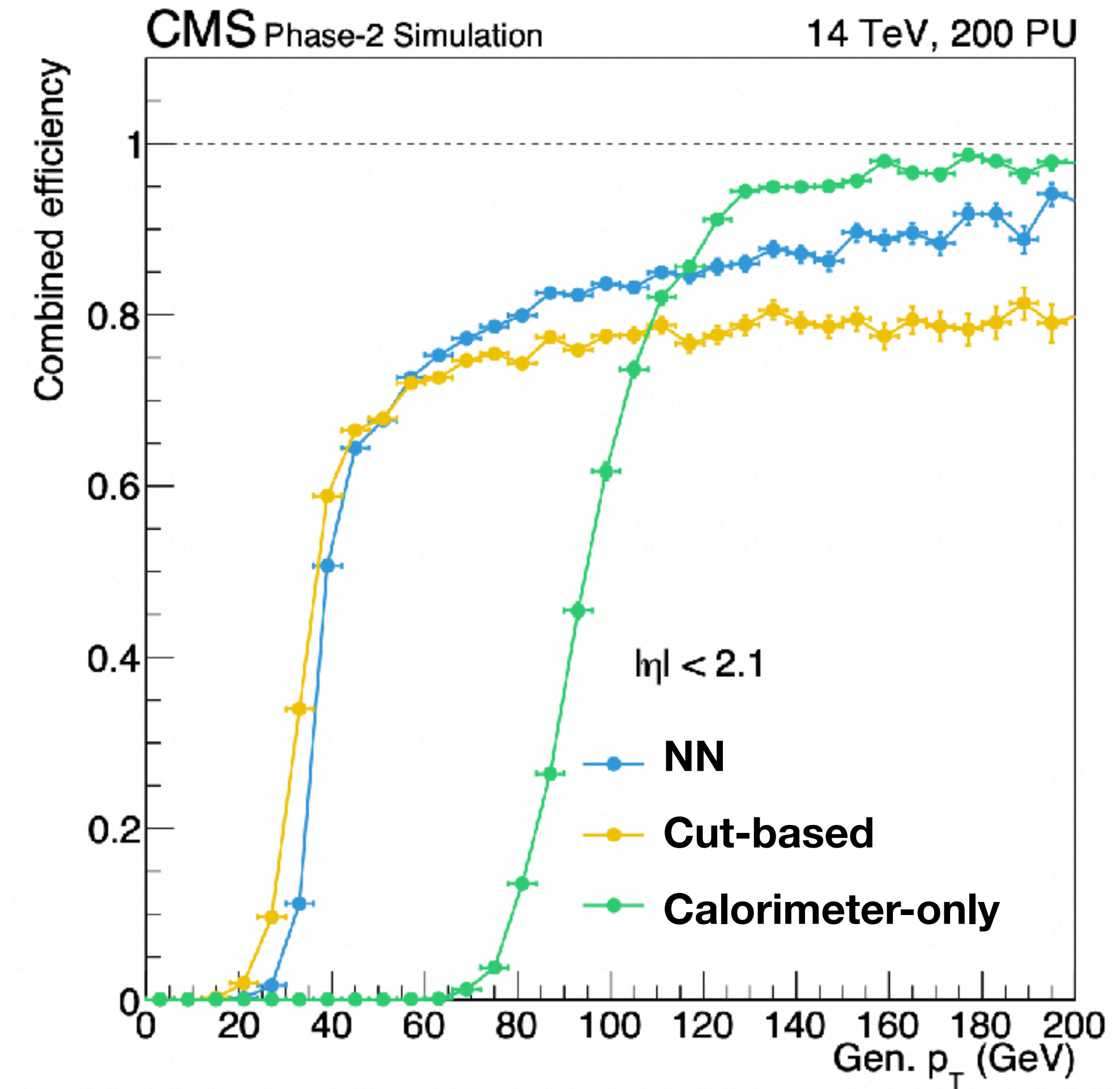
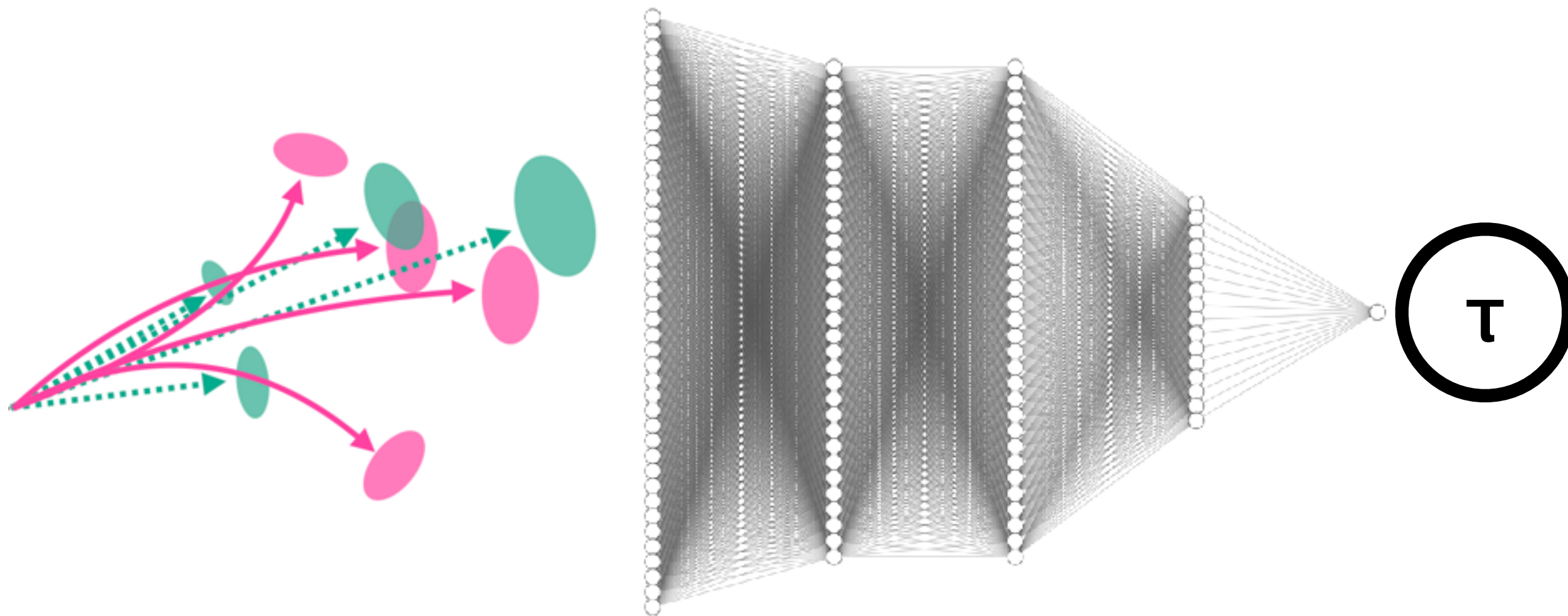
- $HH \rightarrow b\bar{b}b\bar{b}, b\bar{b}\tau\tau$
- Can we design algorithms to differentiate different collections of particles / detector signals
  - $\tau$  lepton,  $b$  bottom quark
  - Light quarks, gluons, noise, combinatorics
- Can we do it every 25 ns on FPGAs?





# L1 $\tau$ Identification

- **NN algorithm** capable of accepting more  $\tau$  leptons than traditional **cut-based method**
- Network is 3 layer dense model, uses information about particle  $p_T$ ,  $\eta$ ,  $\phi$ , and type
- Outputs decision in 38 ns (9 clocks @ 240 MHz)



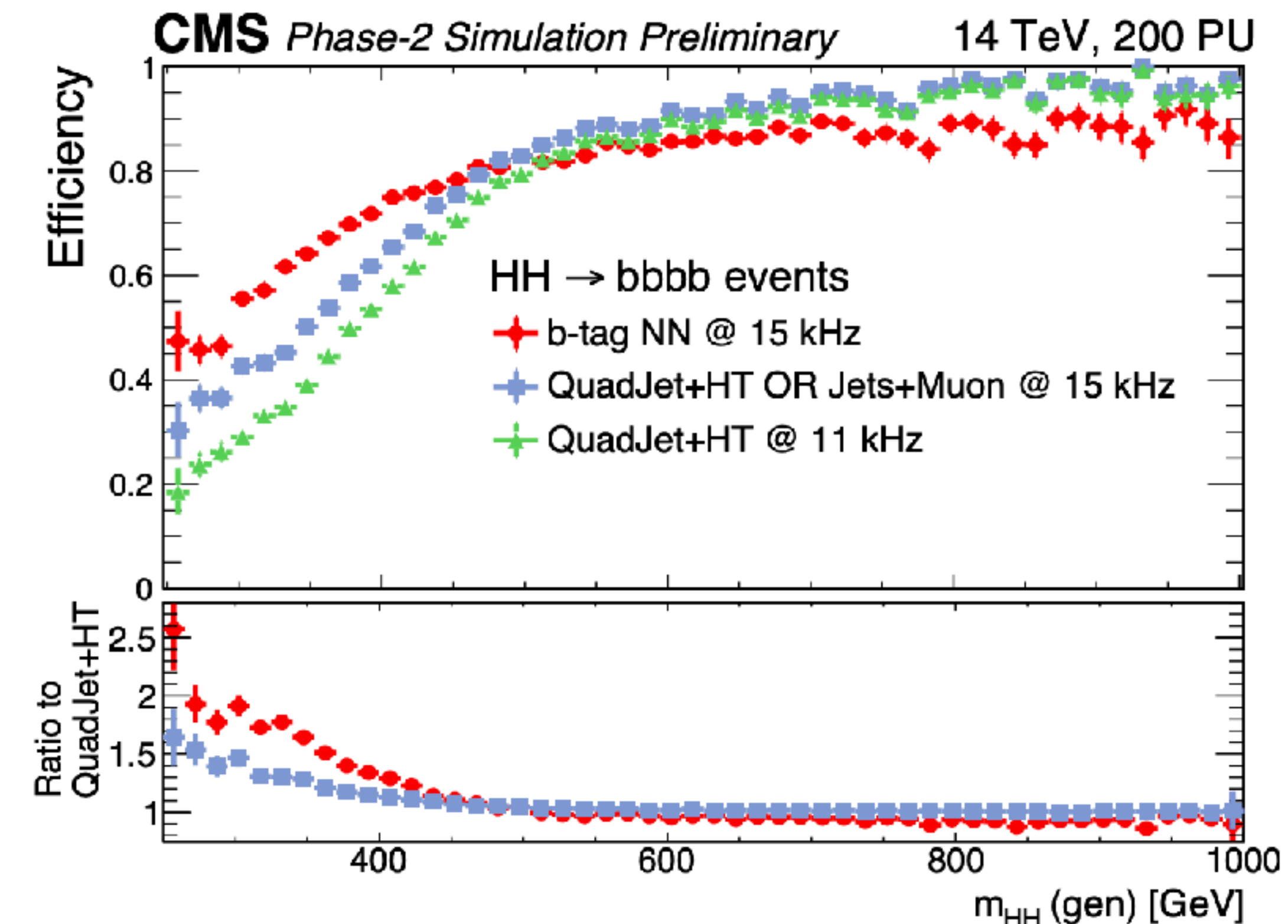
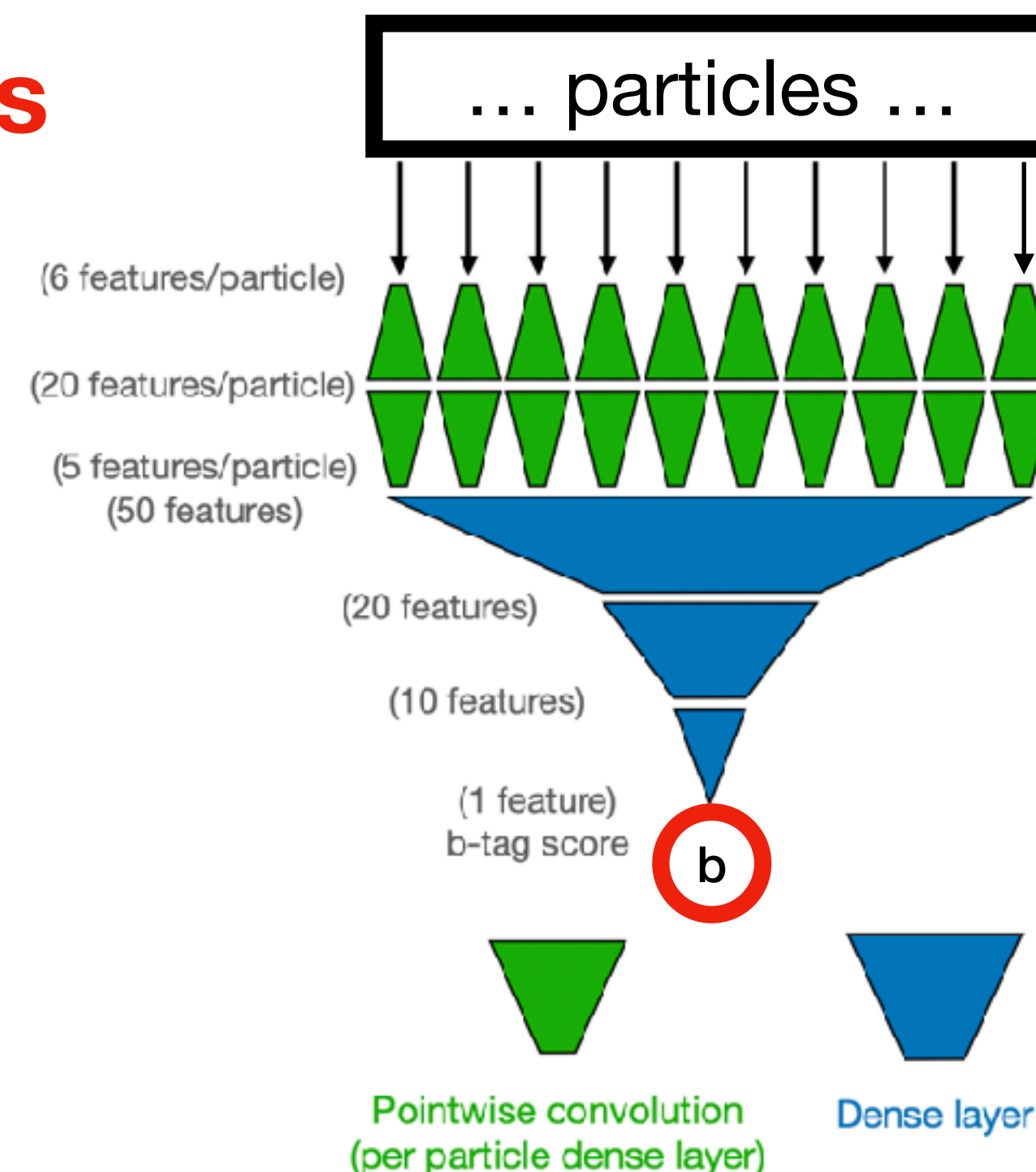
CMS TDR-021



# L1 b-quark Identification

- NN trained to identify b-quarks using collection of particles
- Architecture includes featurizers that act on each particle individual

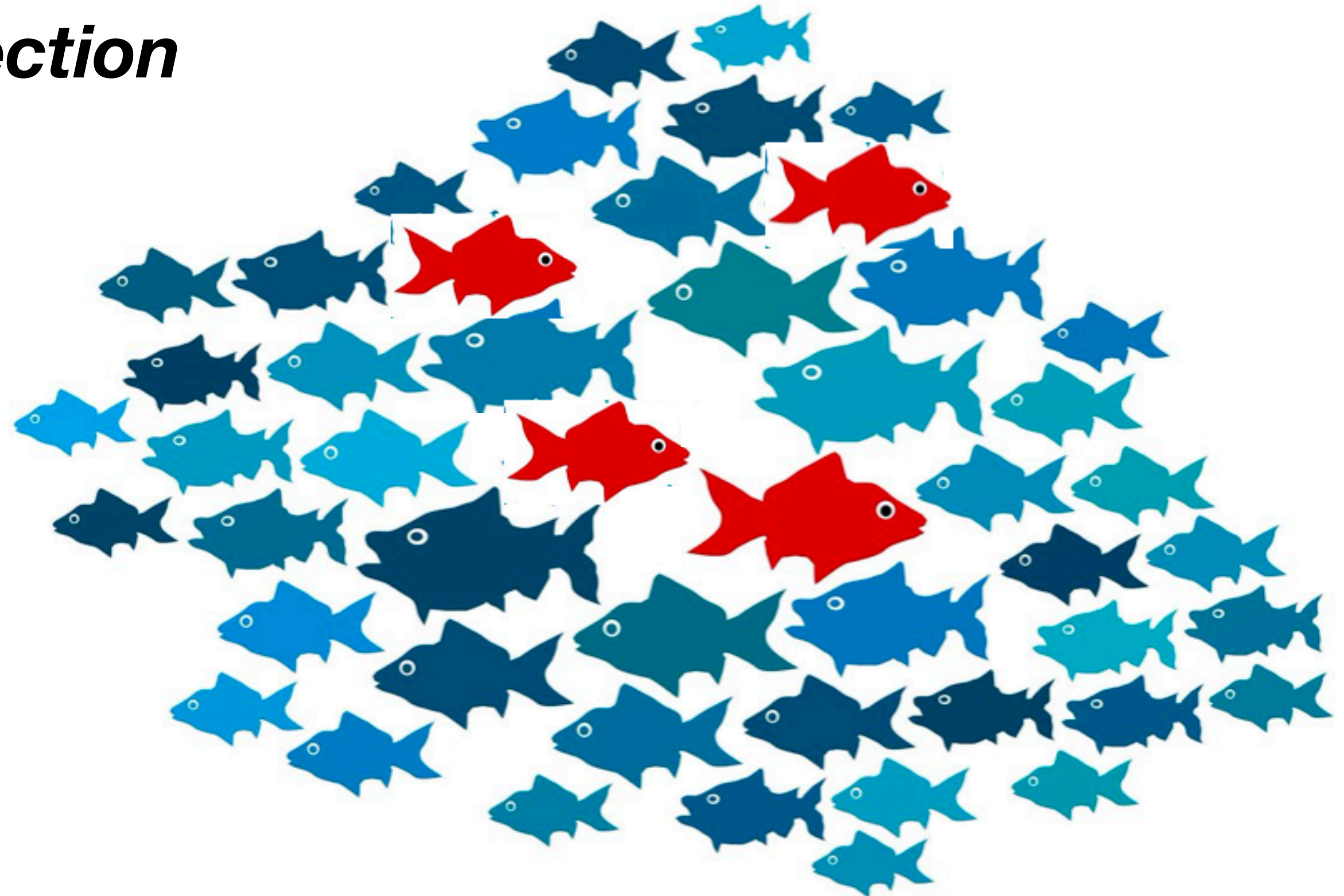
- **Significantly improved acceptance for  $HH \rightarrow bbbb$  events with low  $m_{HH}$**  (compared to traditional cut-based methods)





# Anomaly Detection

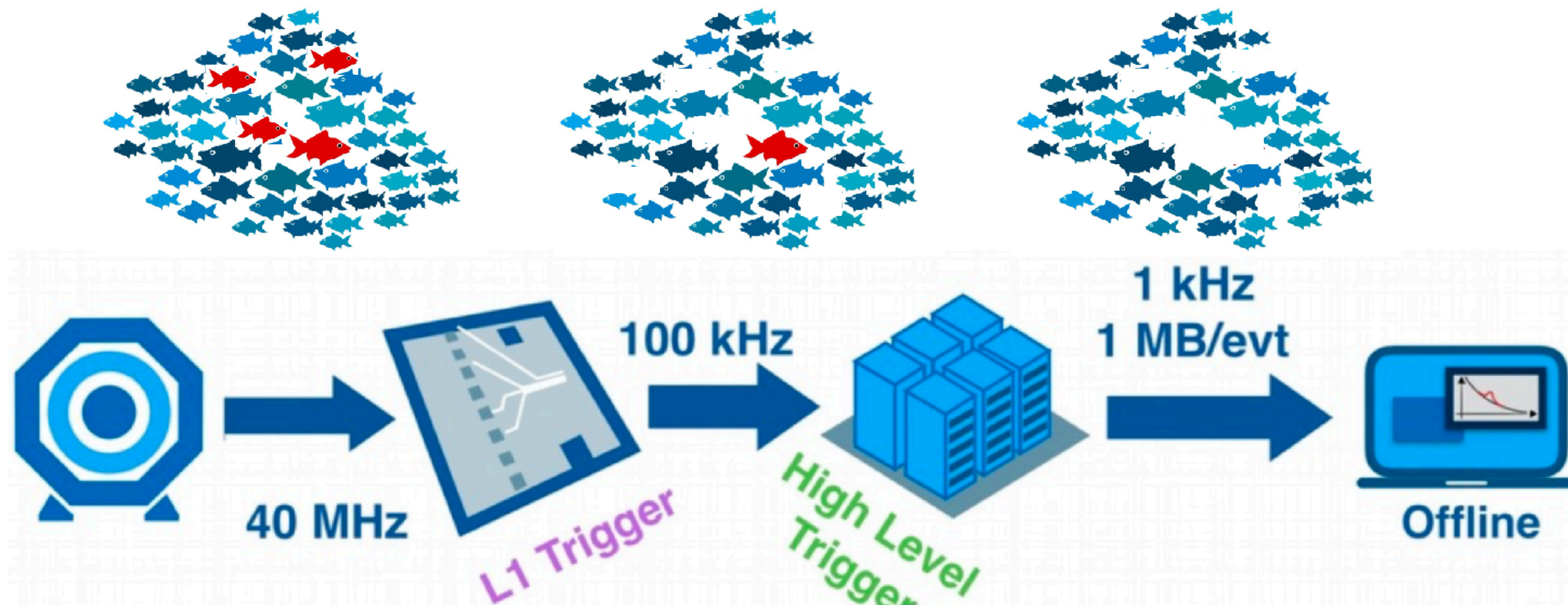
- What if we don't know exactly what we are looking for?
- ML offers unique solution to this challenge (no traditional alternative)
- Broad field of ***anomaly detection***





# Fast Anomaly Detection?

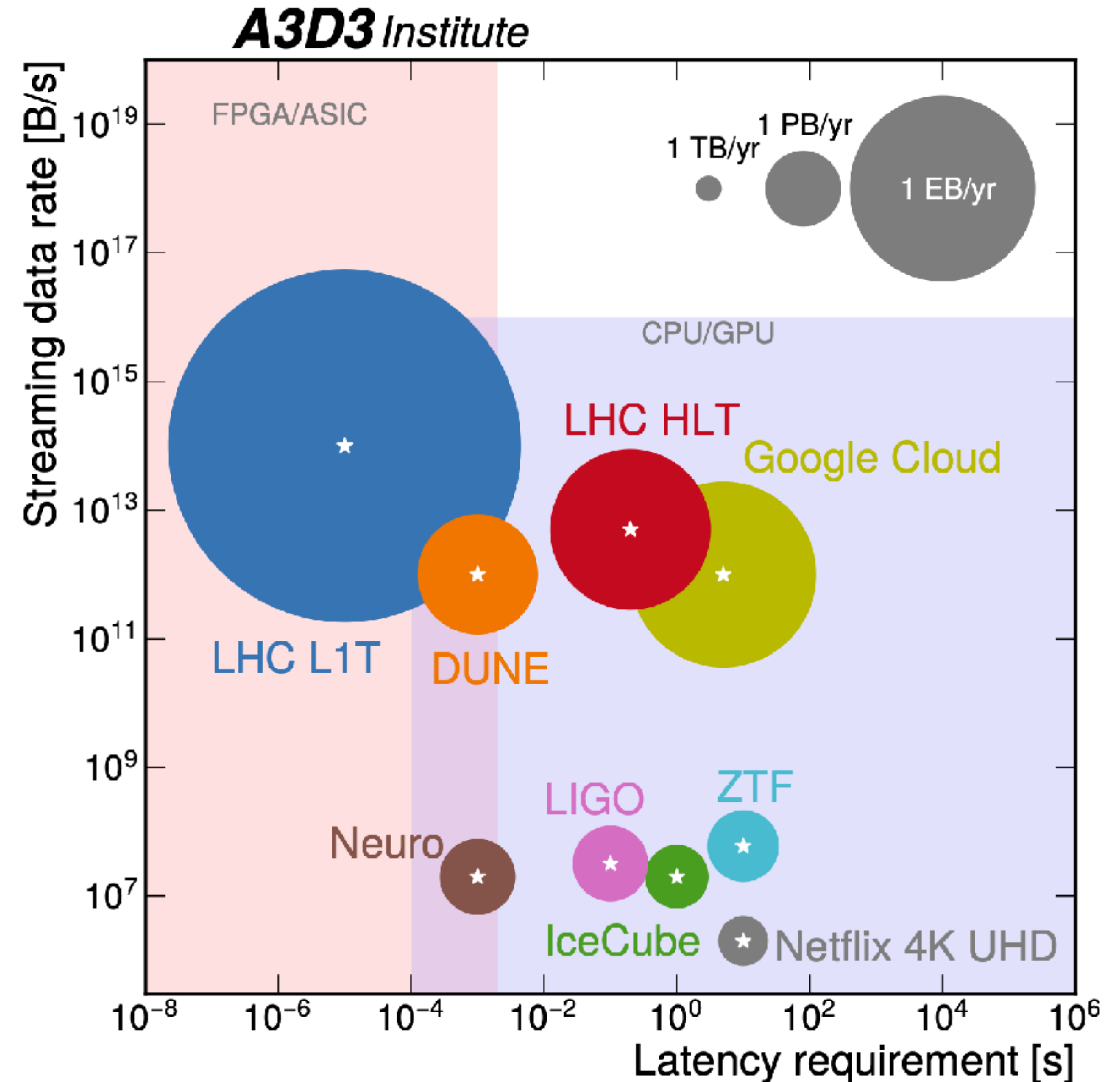
- Depending on anomaly, we could have none left in recorded data
- Low-latency ML is the only option!





# Conclusions

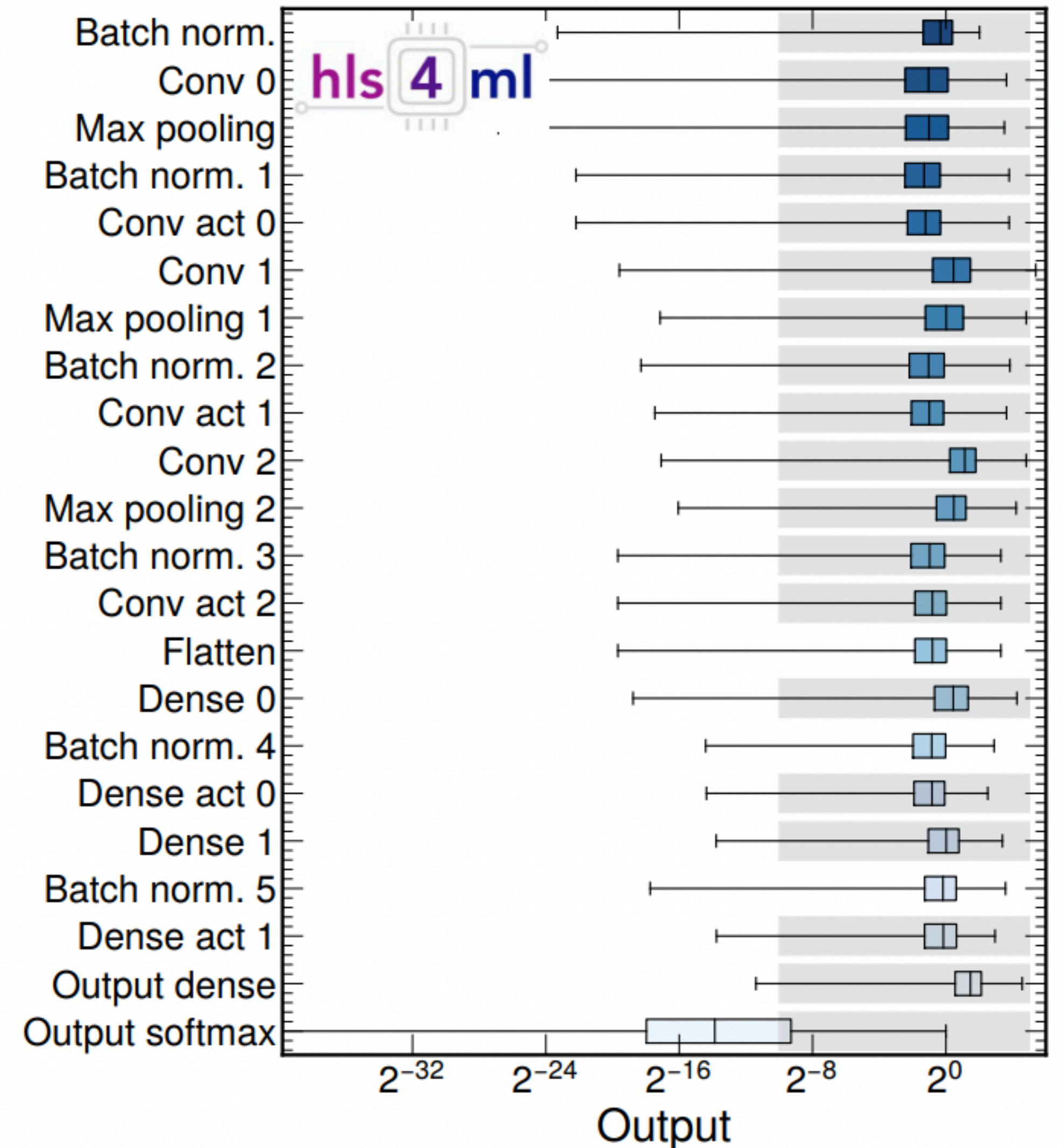
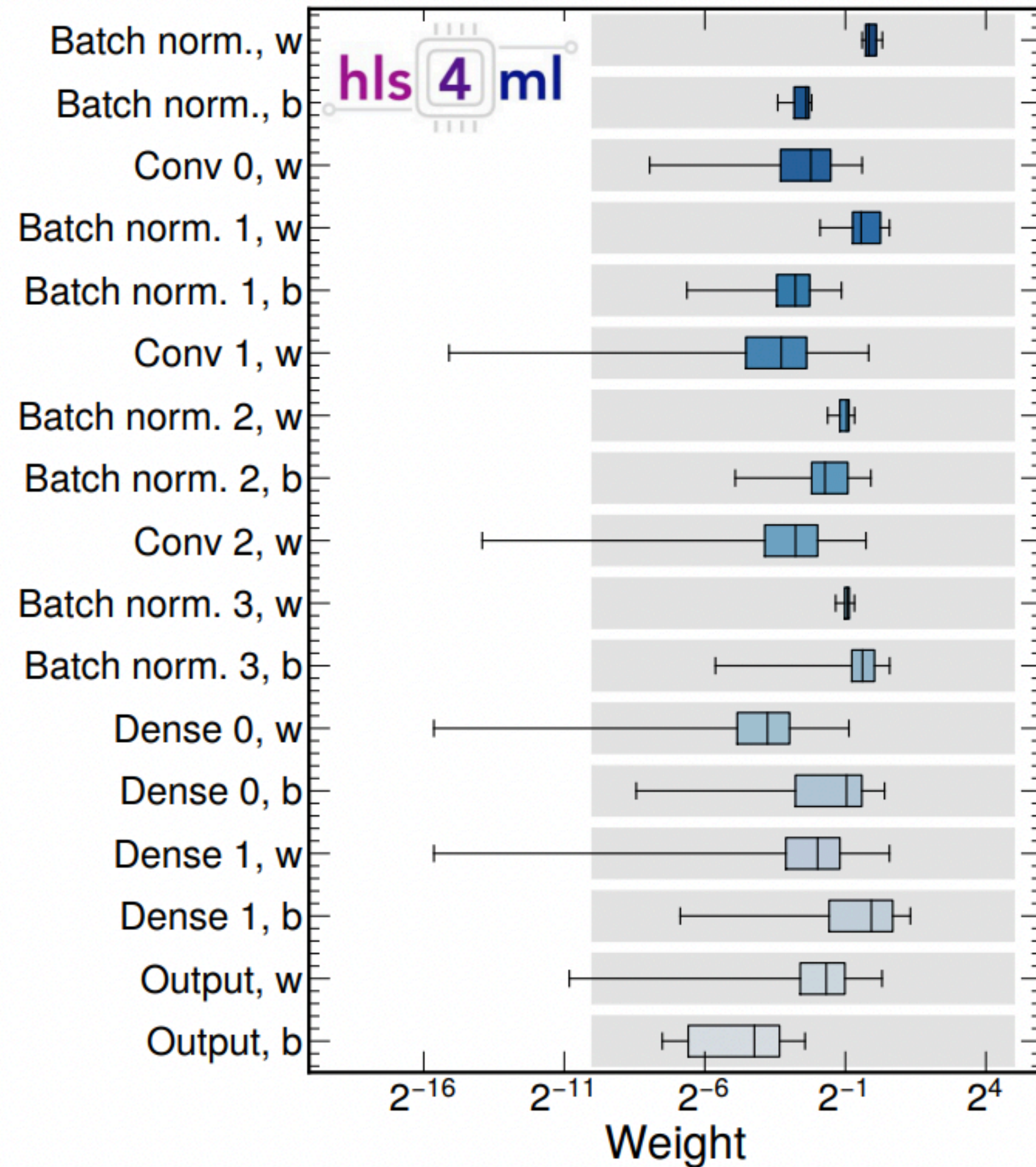
- Increasingly possible to perform low latency inference of ML models
  - Also low-power, high radiation
- ML offers improved performance over traditional algorithms
  - Potential for better alignment of offline and online algorithms
- Applications in many fields, areas



# BACKUP



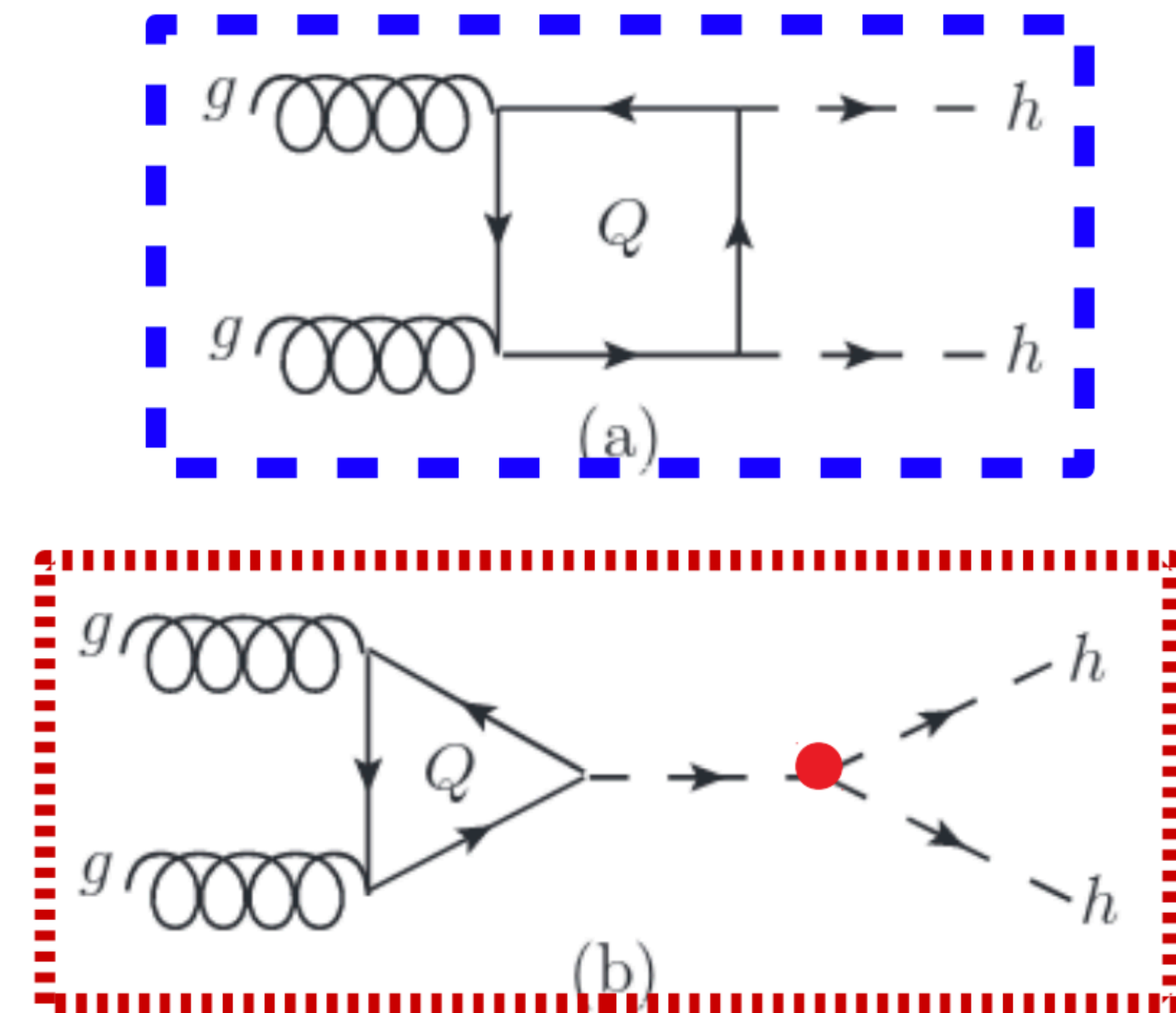
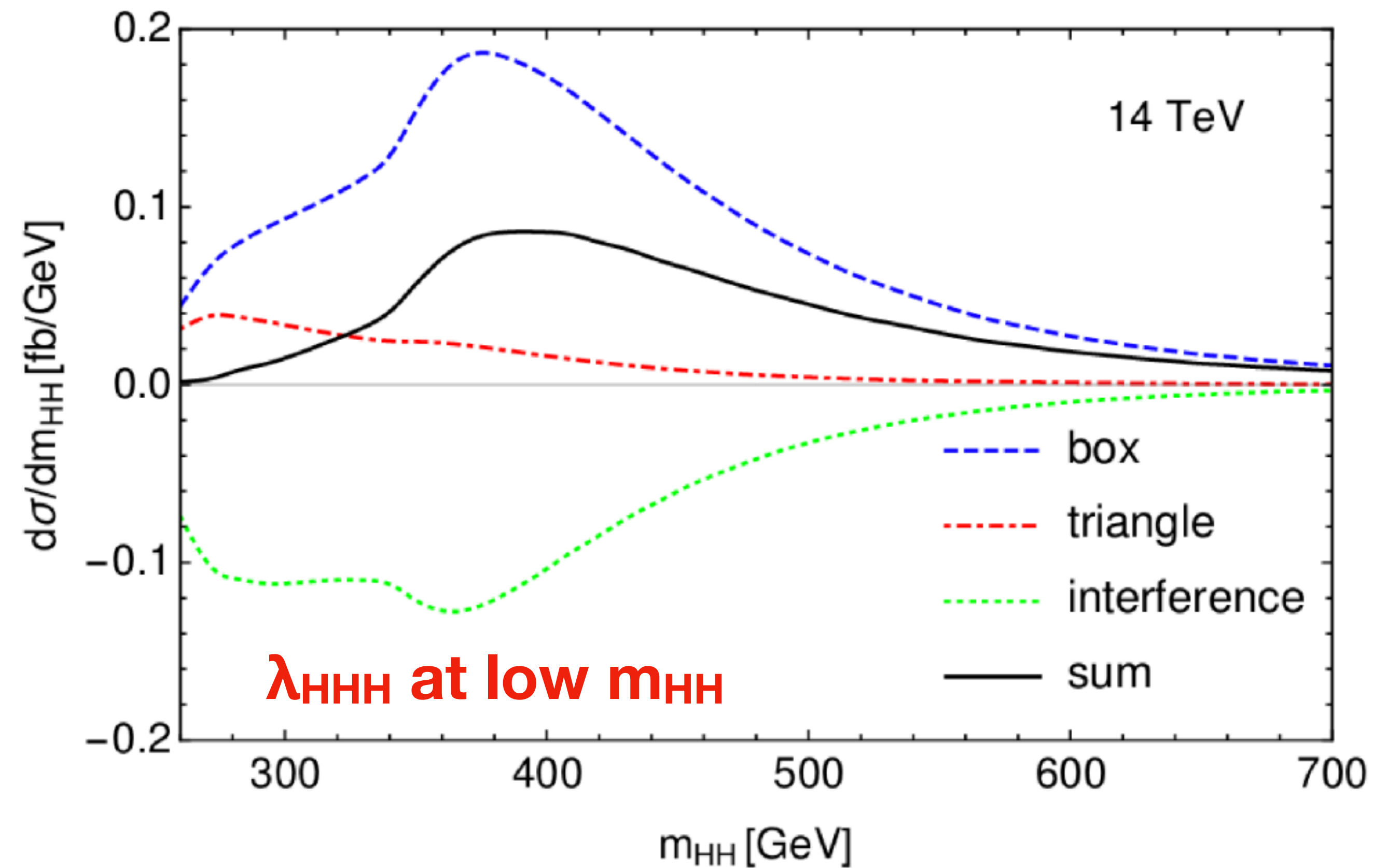
# Precision





# Di-Higgs Decays

- Lots of background that mimics these signals → very difficult to record
- Low  $m_{HH}$  is most critical, but produced object have lower energy

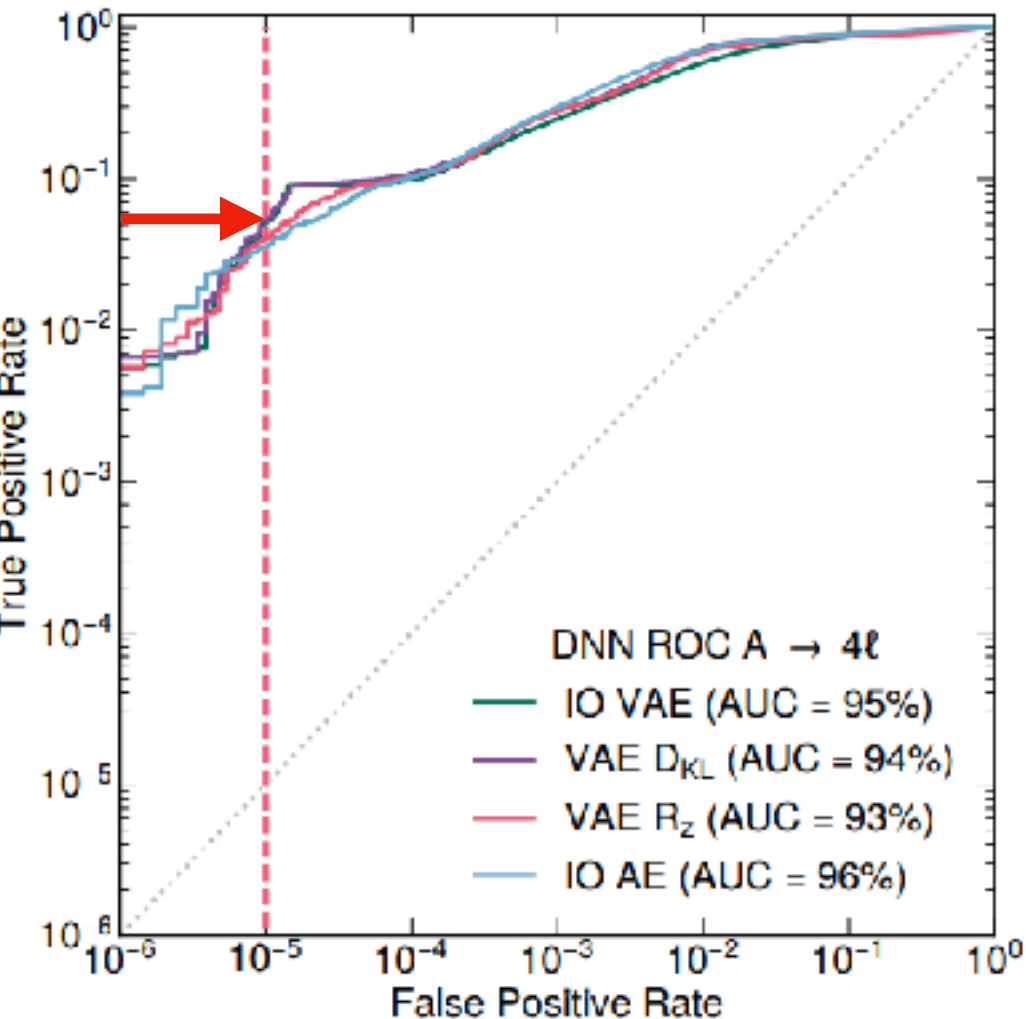
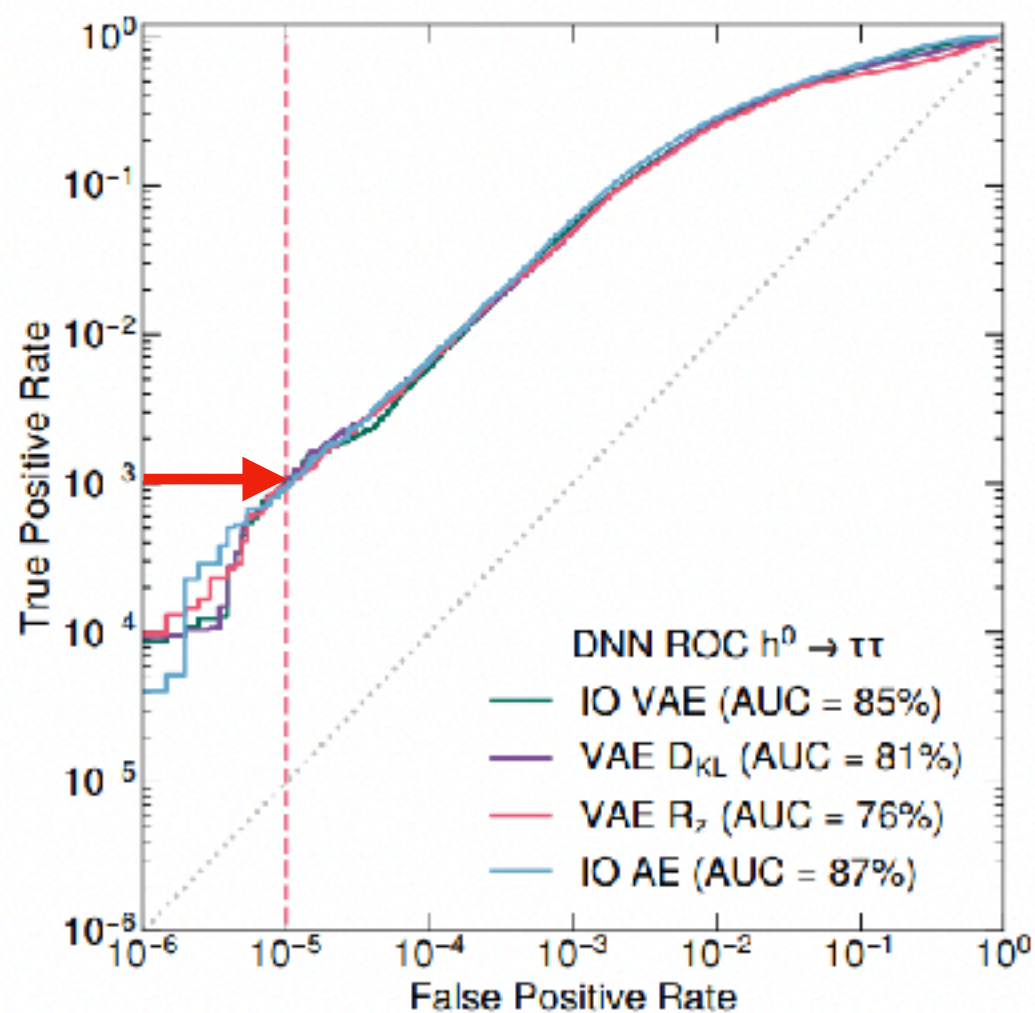
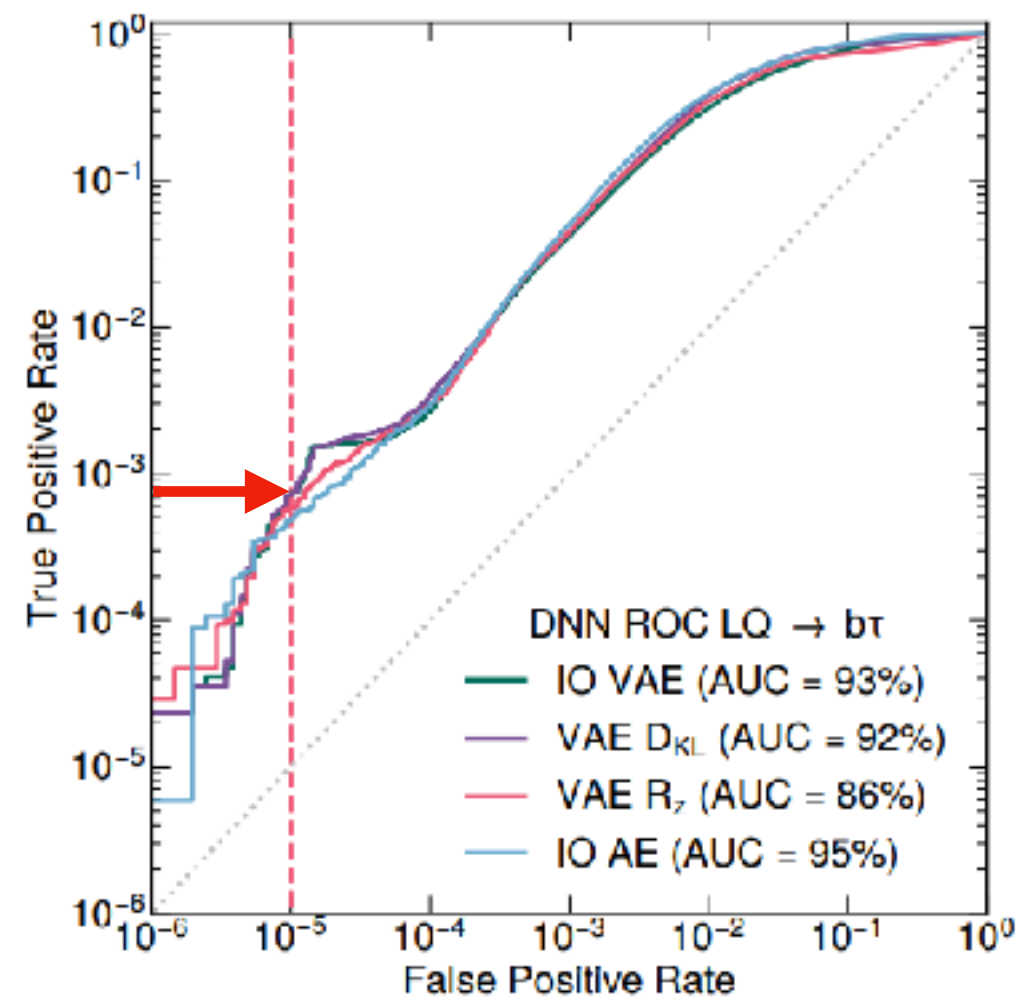
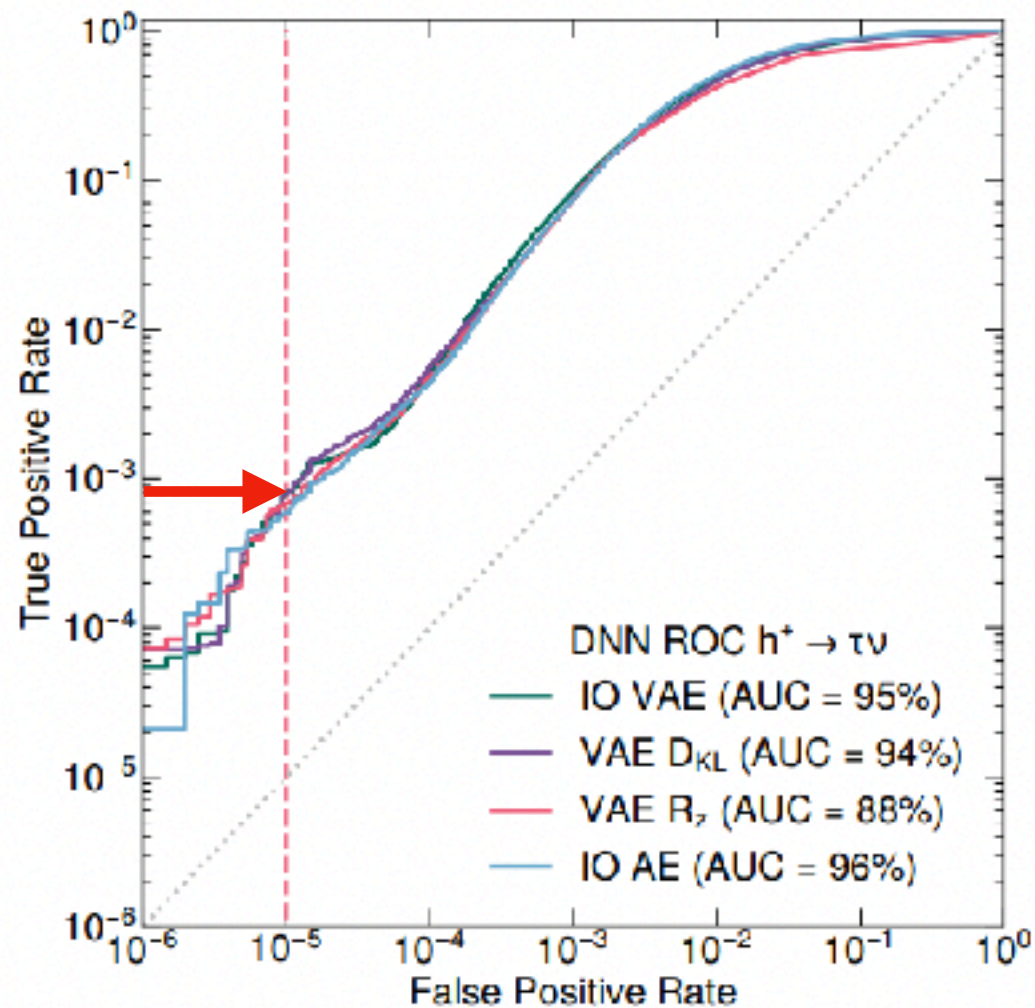




# Fast Anomaly Detection

- Algorithm could take in relevant objects in each event
- Low latency is significant limitation on anomaly detection
- Performance depends on signals

Model	DSP [%]	LUT [%]	FF [%]	BRAM [%]	Latency [ns]	II [ns]
DNN AE QAT 8 bits	2	5	1	0.5	130	5
CNN AE QAT 4 bits	8	47	5	6	1480	895
DNN VAE PTQ 8 bits	1	3	0.5	0.3	80	5
CNN VAE PTQ 8 bits	10	12	4	2	365	115



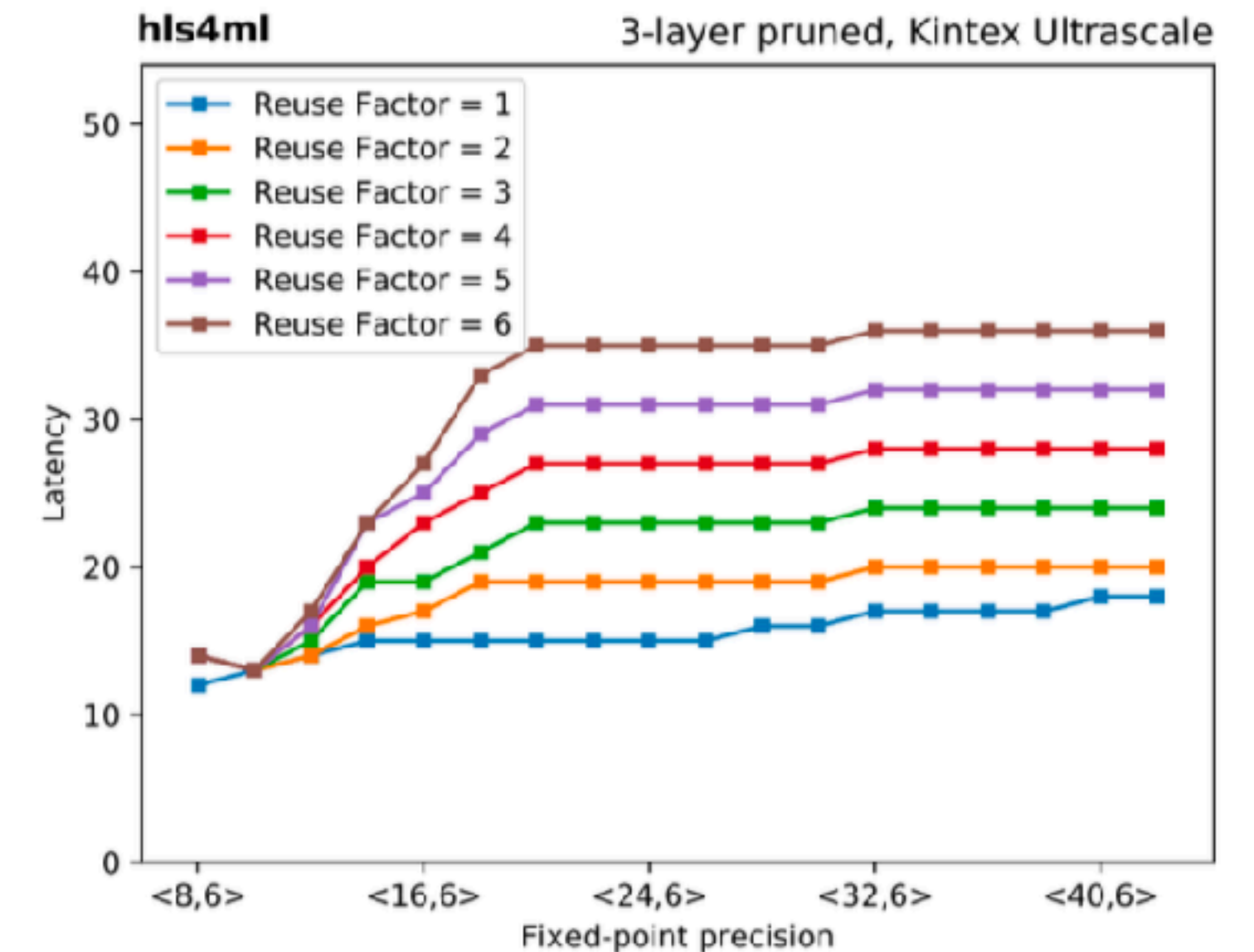
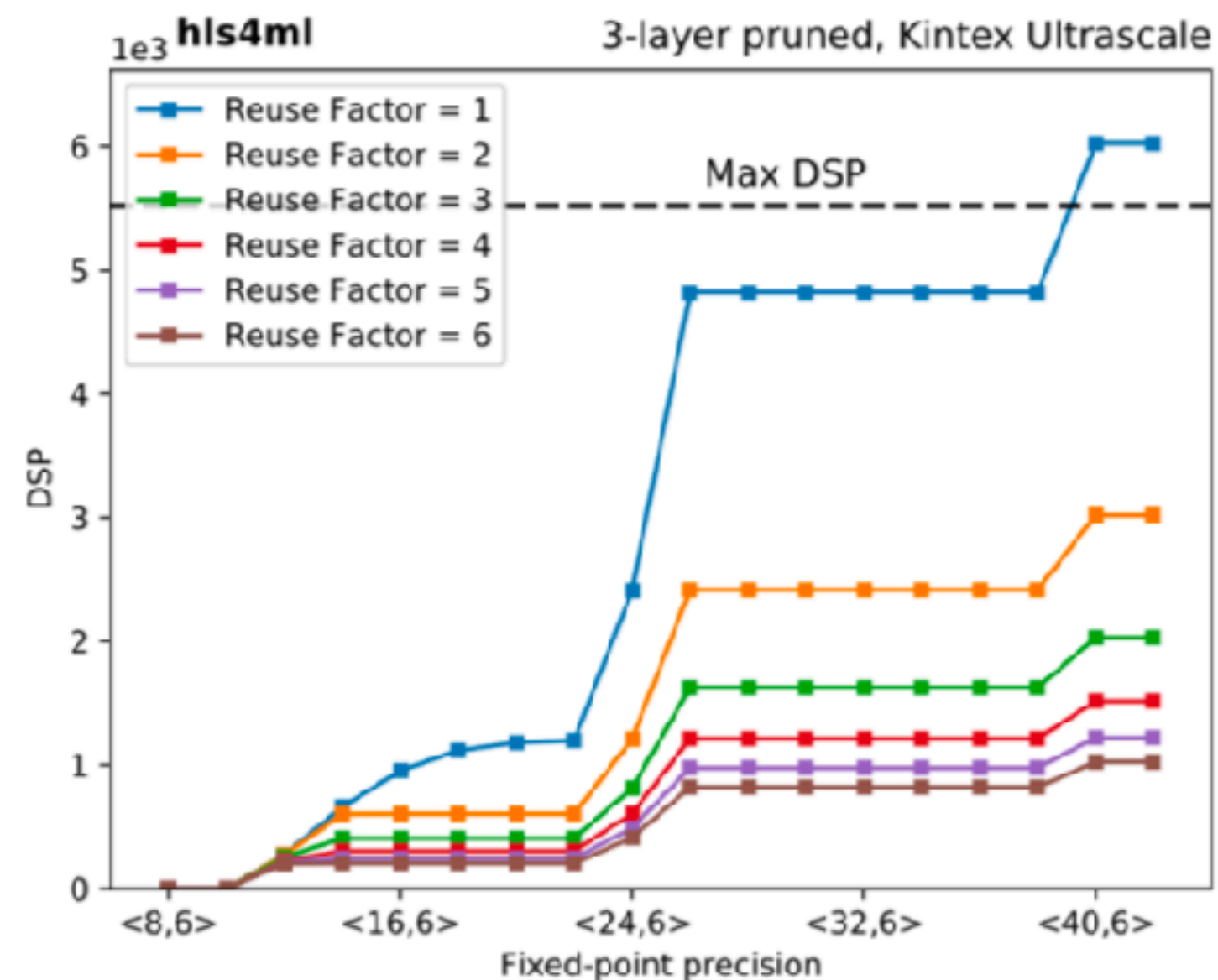
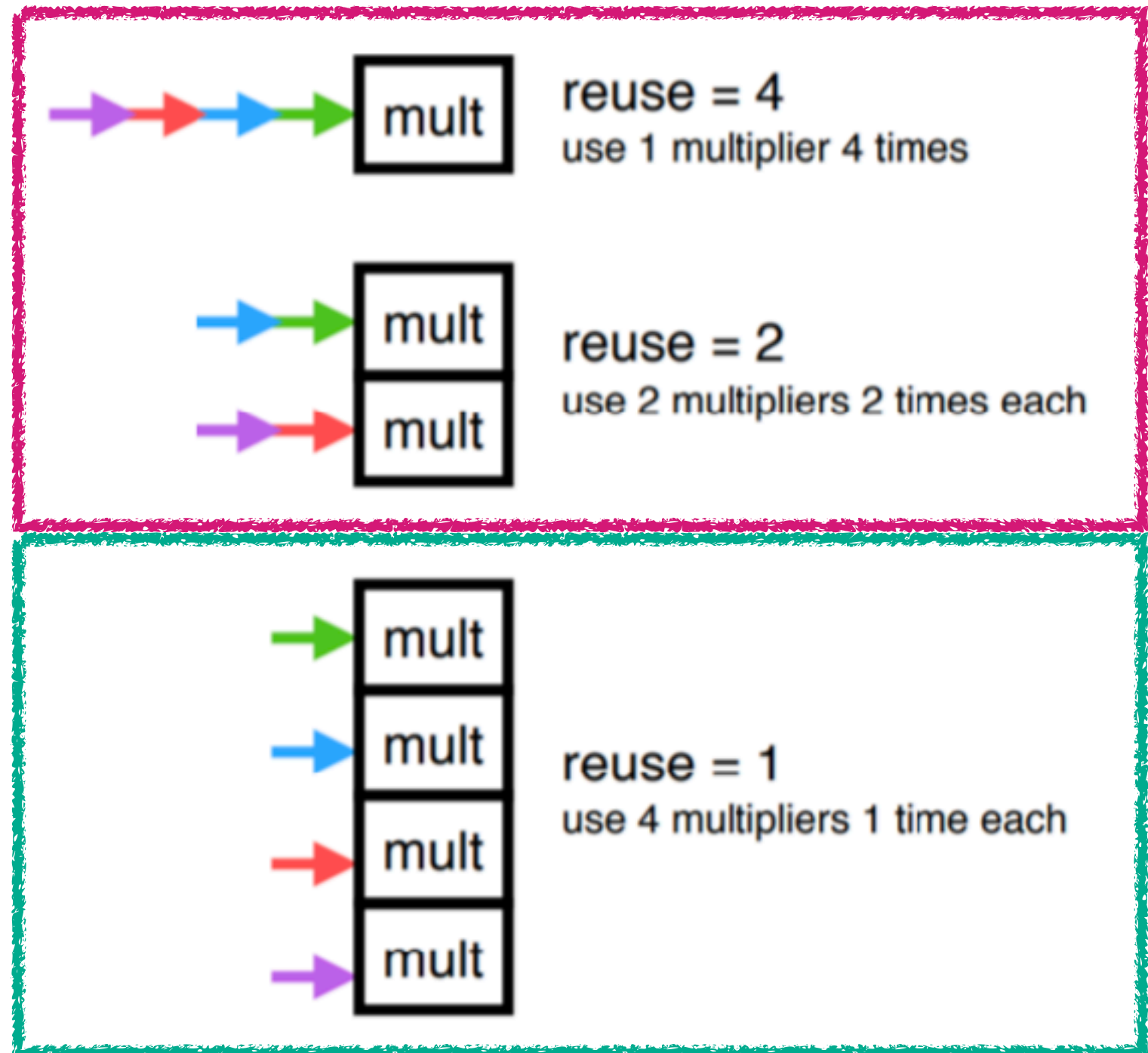
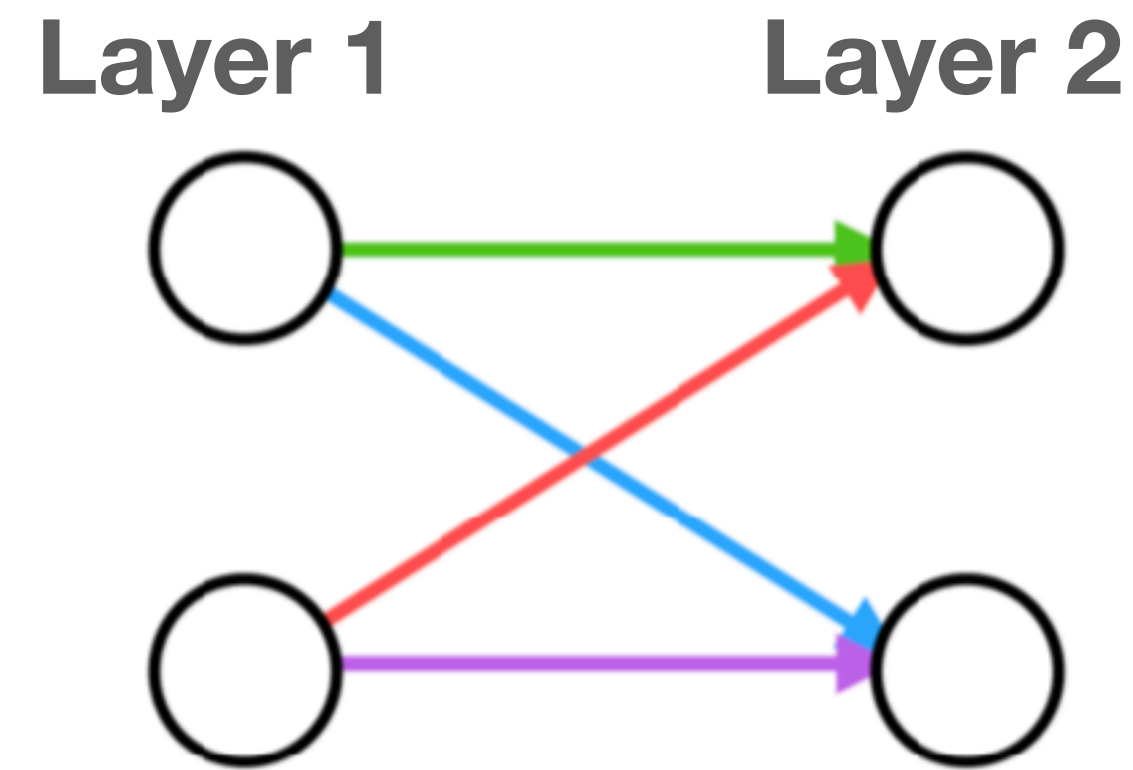
# hls4ml Support

- Support for:
  - **MLPs, BDTs** [[arXiv:2002.02534](#)], **CNNs** [[arXiv:2101.05108](#)], **Binary & Ternary NNs** [[arXiv:2003.06308](#)], **Quantization-aware training (QKeras)** [[arXiv:2006.10159](#)], **Modified GarNet architecture (GraphNN)** [[arXiv:2008.0360](#)], **RNNs/LSTMs/GRUs** [[arXiv:2207.00559](#)]
- Active maintenance and development
  - Many applications for fast ML in physics (low latency, low power)



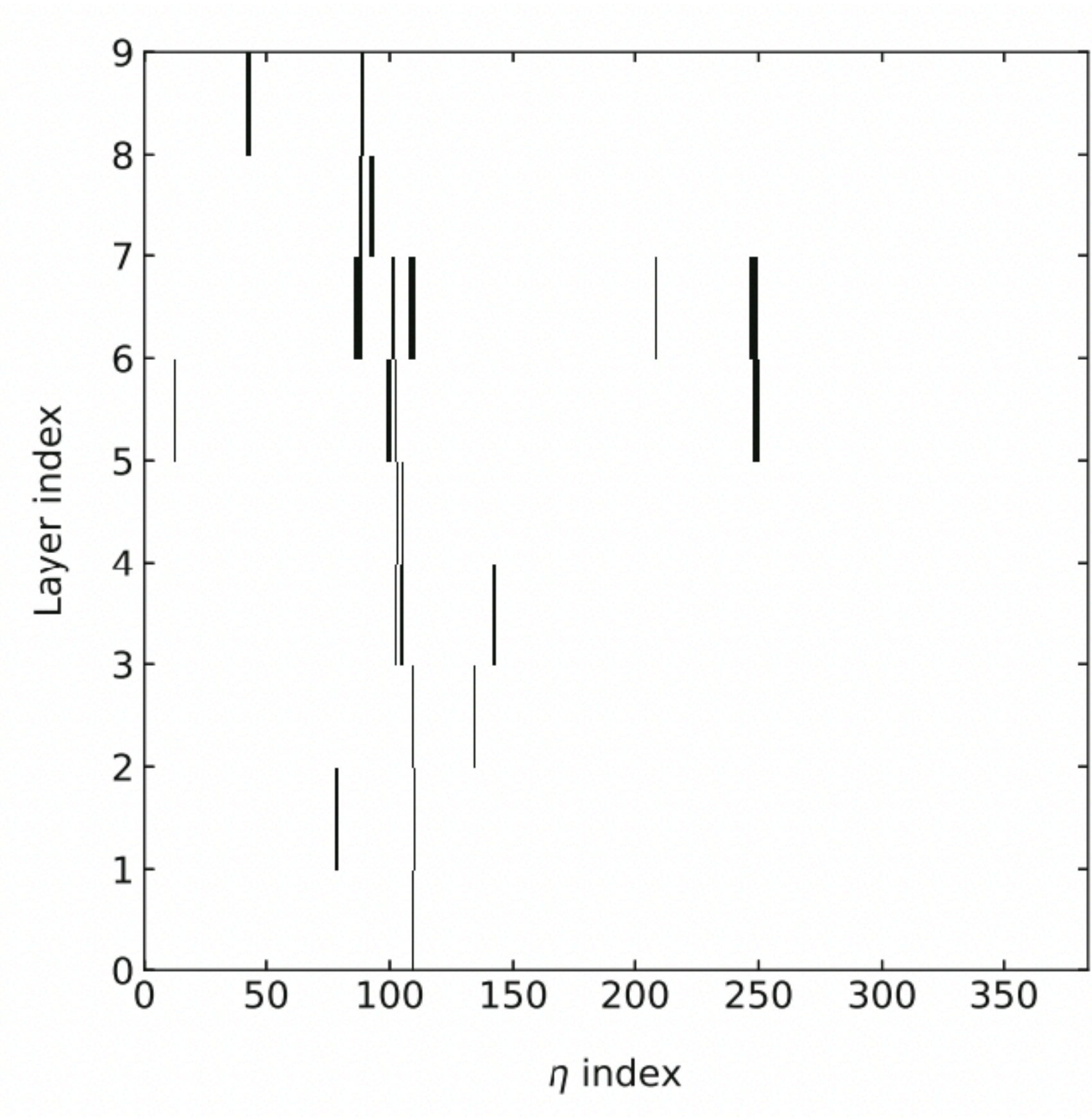
# Reuse

- For lowest latency, compute all multiplications at once
- Reuse = 1 (fully parallel)  
→ latency = # layers
- Larger reuse implies more serialization
- Allows trading higher latency for lower resource usage

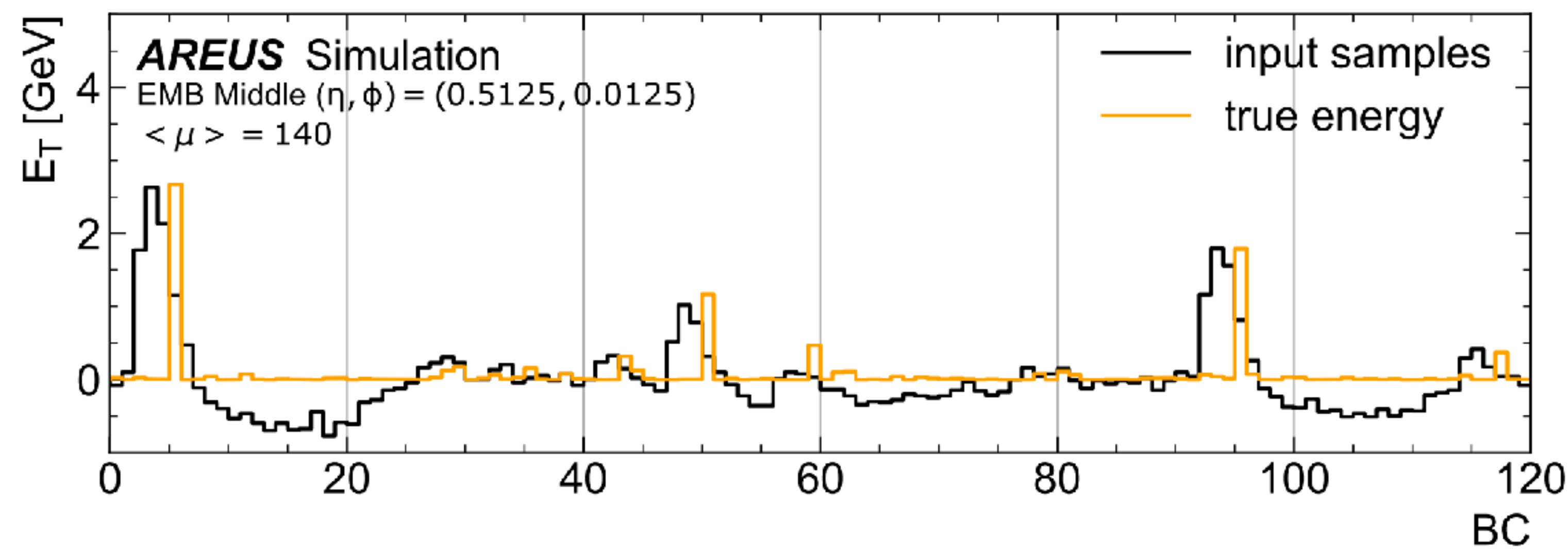




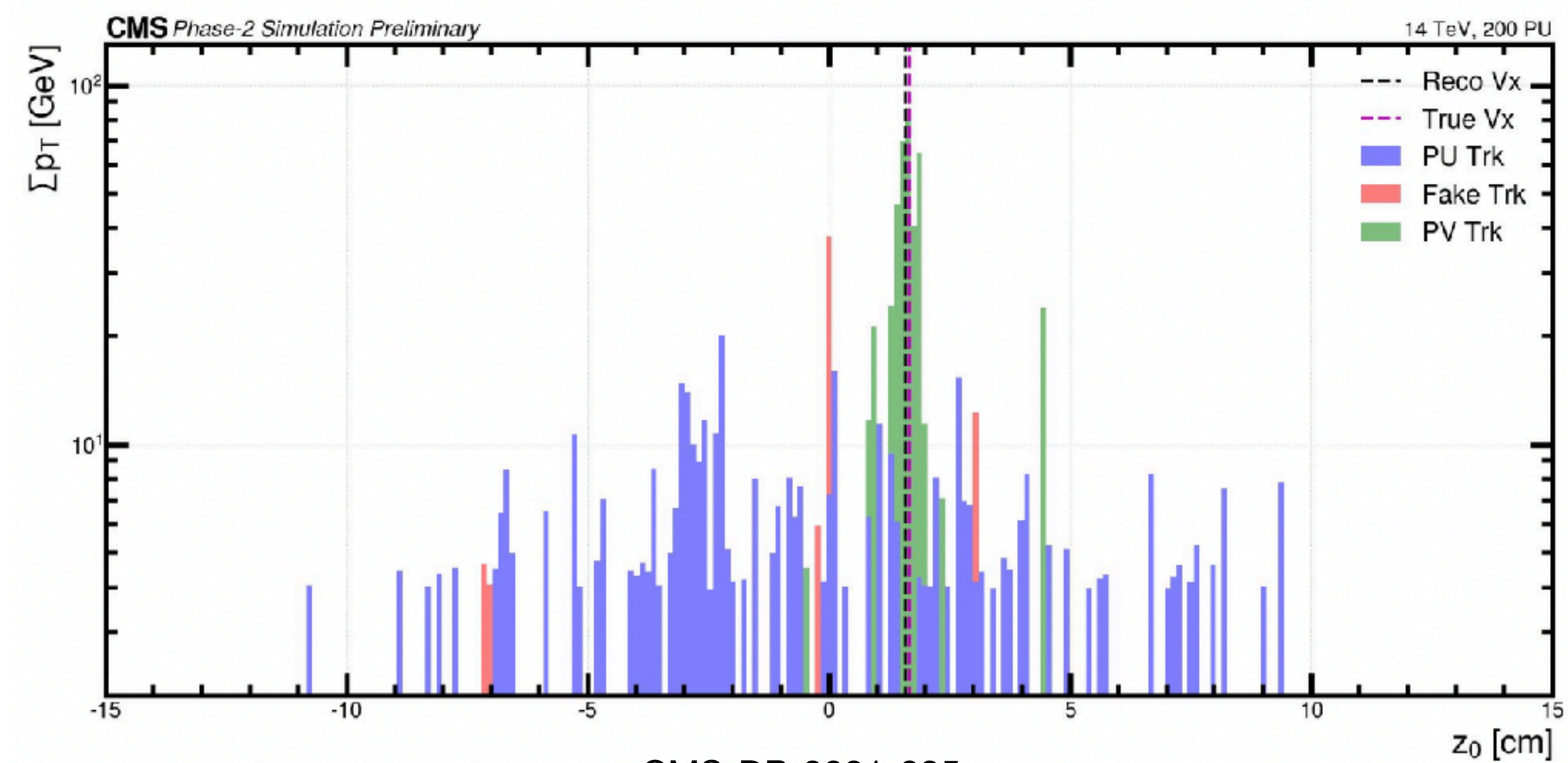
# Applications



Eur. Phys. J. C (2021) 81 :969



arXiv: 2111.08590



CMS-DP-2021-035