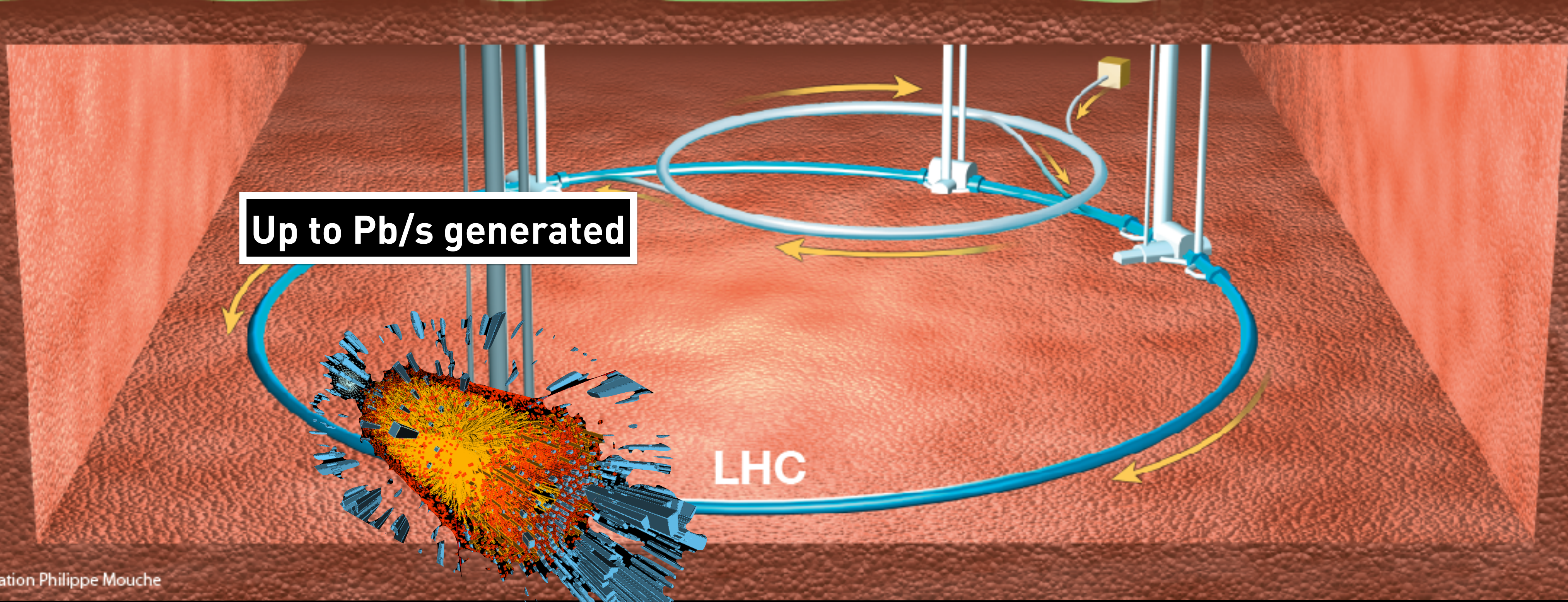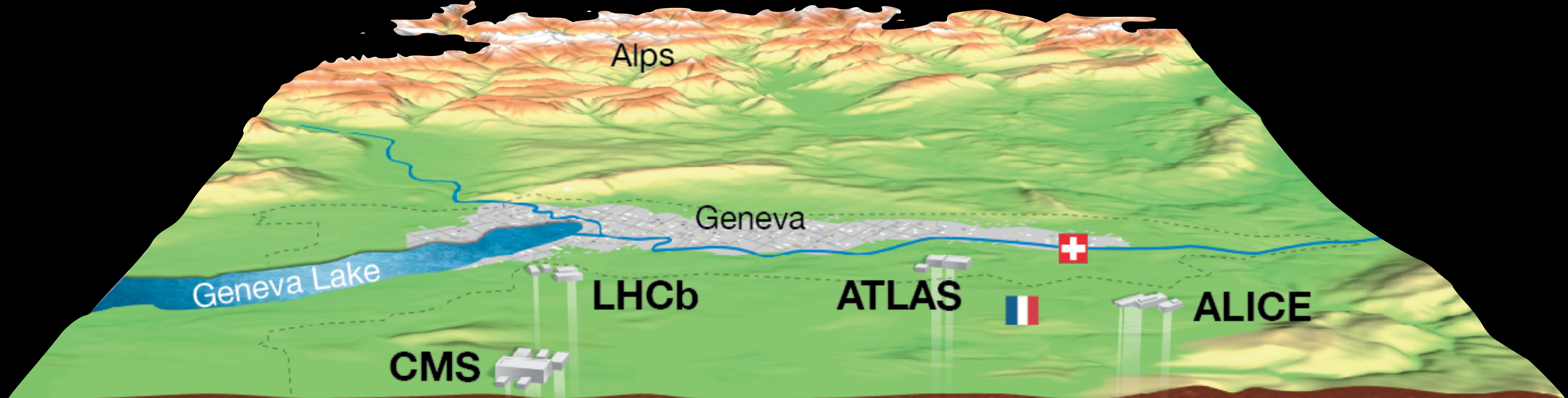# Fast inference on FPGAs at the Large Hadron Collider

Thea Klæboe Årrestad for the hls4ml team
ETH Zürich

Bologna
November 2nd 2022

Alps

Geneva

Geneva Lake

LHCb

ATLAS

ALICE

CMS

Up to Pb/s generated

LHC

Illustration Philippe Mouche

Alps

Geneva

Geneva Lake

LHCb

ATLAS

CMS

**63 Tb/s to L1**

**Data buffered on detector for O(1) μs**

**L1 trigger: O(1) μs latency**

LHC

Illustration Philippe Mouche

Illustration Philippe Mouche
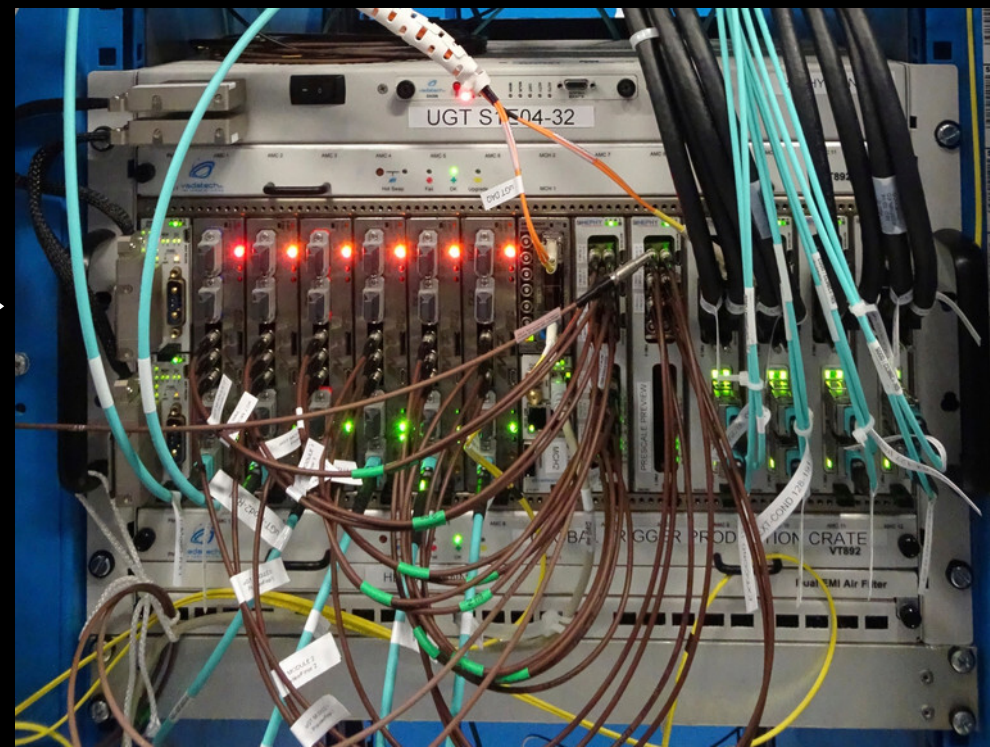
**Level-1 hardware trigger**
- **12.5 µs** to make decision
- Input data bandwidth **63 Tb/s**
- **1000 FPGAs** running thousands of algos

**40 MHz**

**750 kHz**

**7.5 kHz**

**Detector**
- Collisions every 25 ns
- Detector front-end **ASICs**

**FPGA inference**

**ASIC inference**



$$m_H = \sqrt{2E_{\gamma_1}E_{\gamma_2}(1-\cos\theta_{\gamma_1\gamma_2})}$$

19.7 fb⁻¹ (8 TeV) + 5.1 fb⁻¹ (7 TeV)

CMS
H → γγ

S/(S+B) weighted events / GeV

5.7σ

B component subtracted

m_γγ (GeV)

**Nanosecond inference on specialised hardware**

FPGA inference*

ASIC inference*

$$m_H = \sqrt{2E_{\gamma_1}E_{\gamma_2}(1-\cos\theta_{\gamma_1\gamma_2})}$$

5.7σ

*examples in Jennifers talk

# Low latency

- Strictly limited by collisions
  occurring every 25 ns

## Low latency

- Strictly limited by collisions occurring every 25 ns



10¹¹ protons

25 ns / 7.5 m

~60 pp collisions per crossing!

## Low resource usage

- Several algos in parallel on single device

# Low latency
- Strictly limited by collisions occurring every 25 ns



$10^{11}$ protons

25 ns / 7.5 m

~60 pp collisions per crossing!

# Low resource usage
- Several algos in parallel on single device
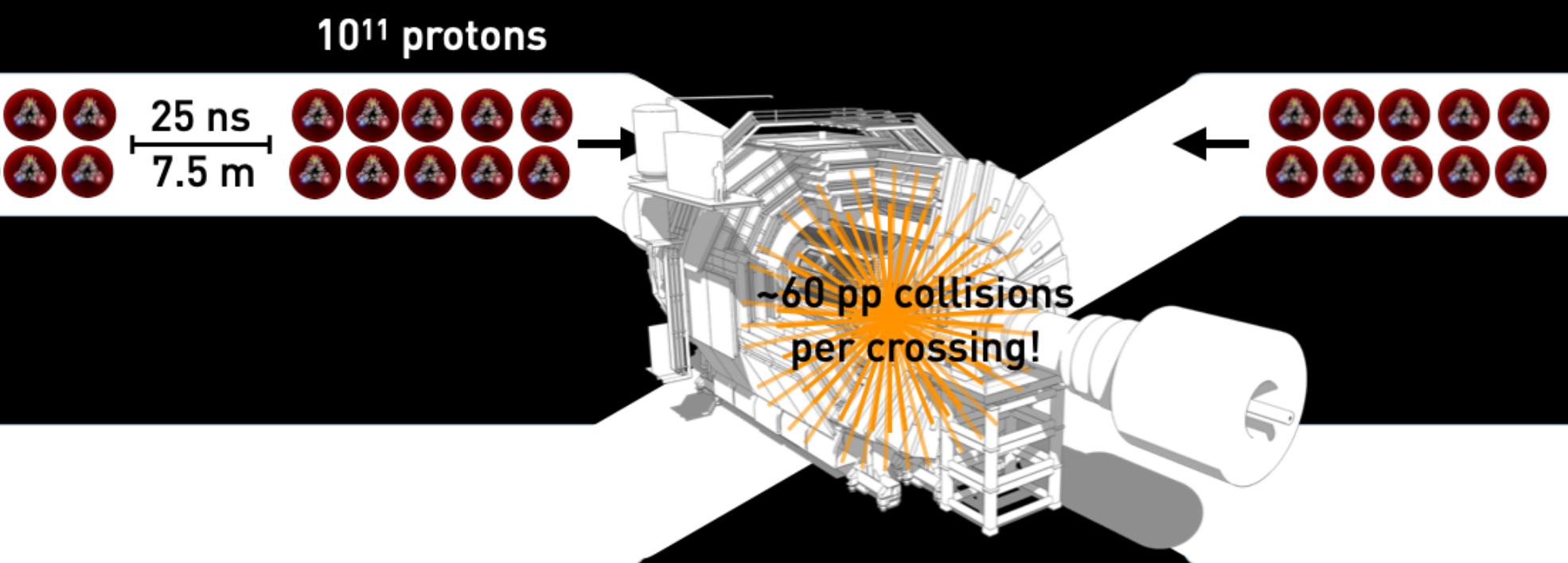


# Low power
- On detector: cooled to -30°C
- L1: Cooling, cooling, cooling

To L1



ASIC

**Extreme combination of low power, low latency, low resource!**

# The level-1 trigger



CALORIMETRY:
370 FPGAs

*54 for HGCAL only!

TRACK FINDER:
174 FPGAs

MUONS:
96 FPGAs

5 μs

COMBINE:
66 FPGAs

CALORIMETRY

PARTICLE FLOW

MUONS

EXTERNAL TRIGGERS

GLOBAL TRIGGER:
24 FPGAs

Trigger accept/reject

12.5 μs

USC55

UXC55

~1000 algorithms in parallel on ~10 FPGAs
~100 ns latency per algorithm

# Why FPGAs at LHC?

# Why FPGAs at LHC?



High parallelism ⬆ = Low latency⬇

- Can work on different data simultaneously (pipelining)! **High bandwidth**

# Why FPGAs at LHC?



High parallelism ⬆ = Low latency⬇
- Can work on different data simultaneously (pipelining)! **High bandwidth**

Power efficient
- FPGAS ~x10 more power efficient than GPUs
  (our L1T FPGA processors pull currents of O(200)A at ~1V, dissipate **heat** of ~7W/cm$^2$
  while processing **5% of total internet traffic**!

# Why FPGAs at LHC?



High parallelism ⬆ = Low latency⬇
- Can work on different data simultaneously (pipelining)! **High bandwidth**

Power efficient
- FPGAS ~x10 more power efficient than GPUs
  (our L1T FPGA processors pull currents of O(200)A at ~1V, dissipate **heat** of ~7W/cm$^2$
  while processing **5% of total internet traffic**!

Latency deterministic
- CPU/GPU has processing randomness, FPGAs **repeatable and predictable latency**

# Why FPGAs at LHC?



High parallelism ⬆ = Low latency⬇
- Can work on different data simultaneously (pipelining)! **High bandwidth**

Power efficient
- FPGAS ~x10 more power efficient than GPUs
  (our L1T FPGA processors pull currents of O(200)A at ~1V, dissipate **heat** of ~7W/cm$^2$
  while processing **5% of total internet traffic**!

Latency deterministic
- CPU/GPU has processing randomness, FPGAs **repeatable and predictable latency**

**Latency is fixed by proton collisions occurring at 40 MHz, cannot tolerate slack**

# What are FPGAs?

**Field Programmable Gate Arrays**
- Different resources with programmable interconnects (reprogrammable)
- Originally ASIC prototyping, now also for high performance computing
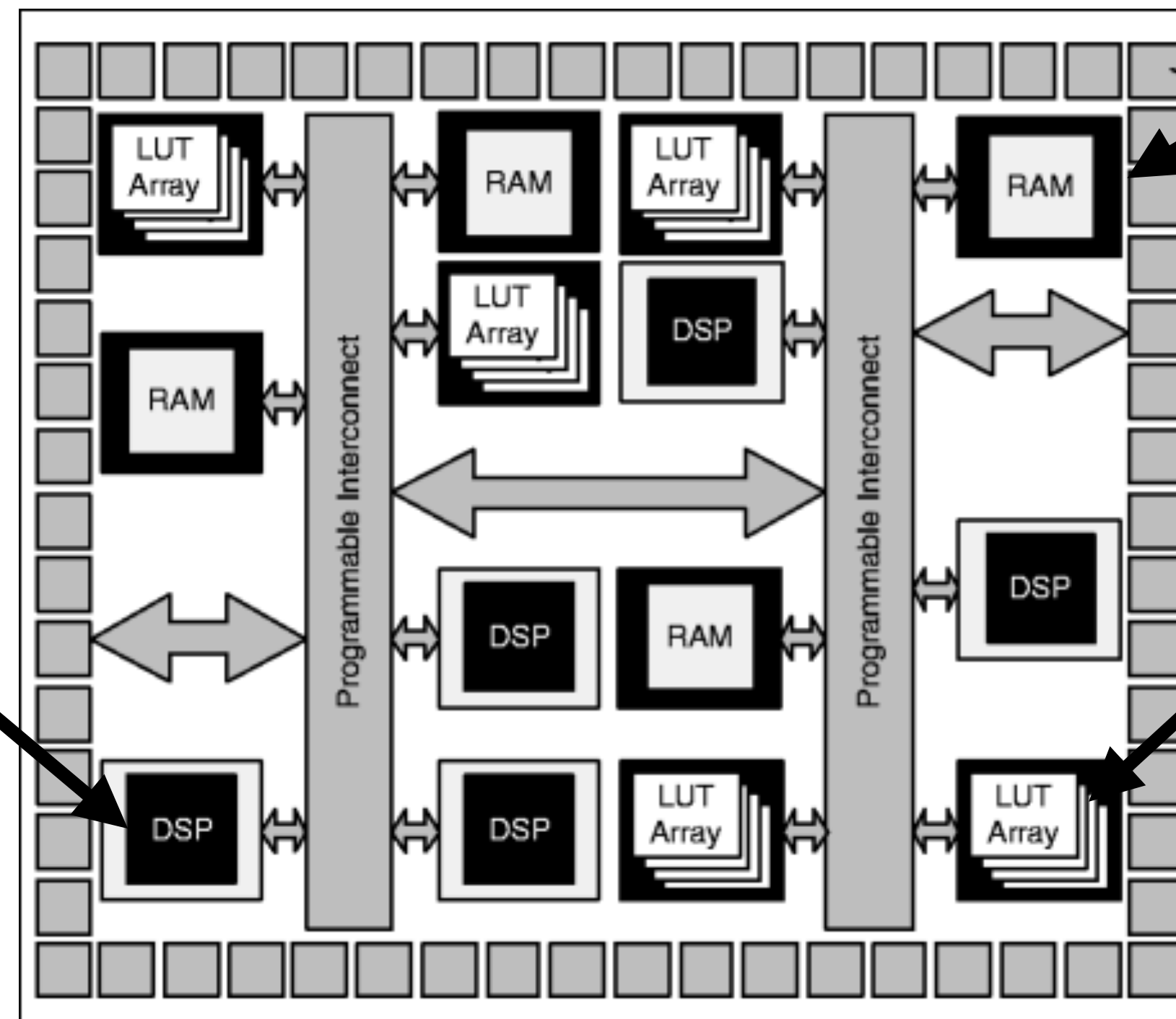
**See Riccardo's talk!**

# What are FPGAs?

Field Programmable Gate Arrays
- Different resources with programmable interconnects (reprogrammable)
- Originally ASIC prototyping, now also for high performance computing

Memory (BRAM)

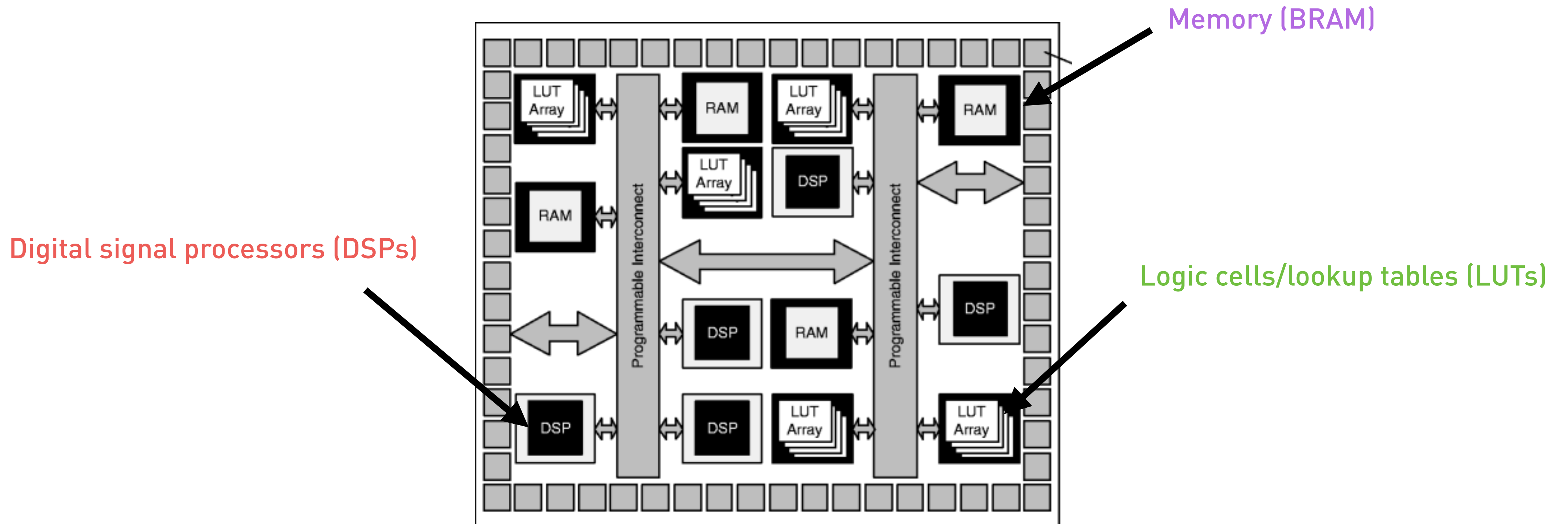Digital signal processors (DSPs):
specialised for multiplication

Logic cells/lookup tables (LUTs):
perform arbitrary functions

flip-flops (FF):
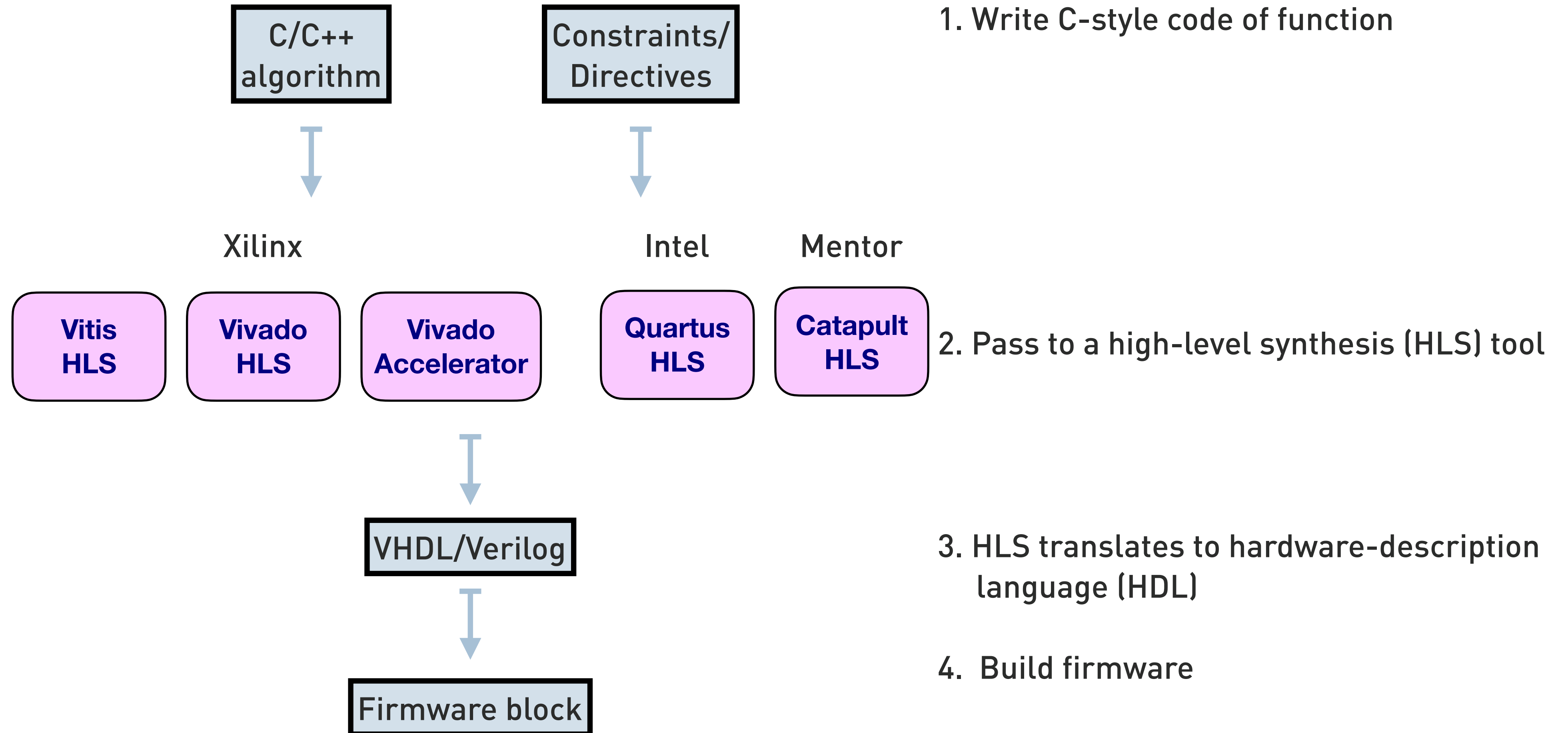registers data in time with clock pulse

See Riccardo's talk!



**13**

$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

activation function       multiplication       addition

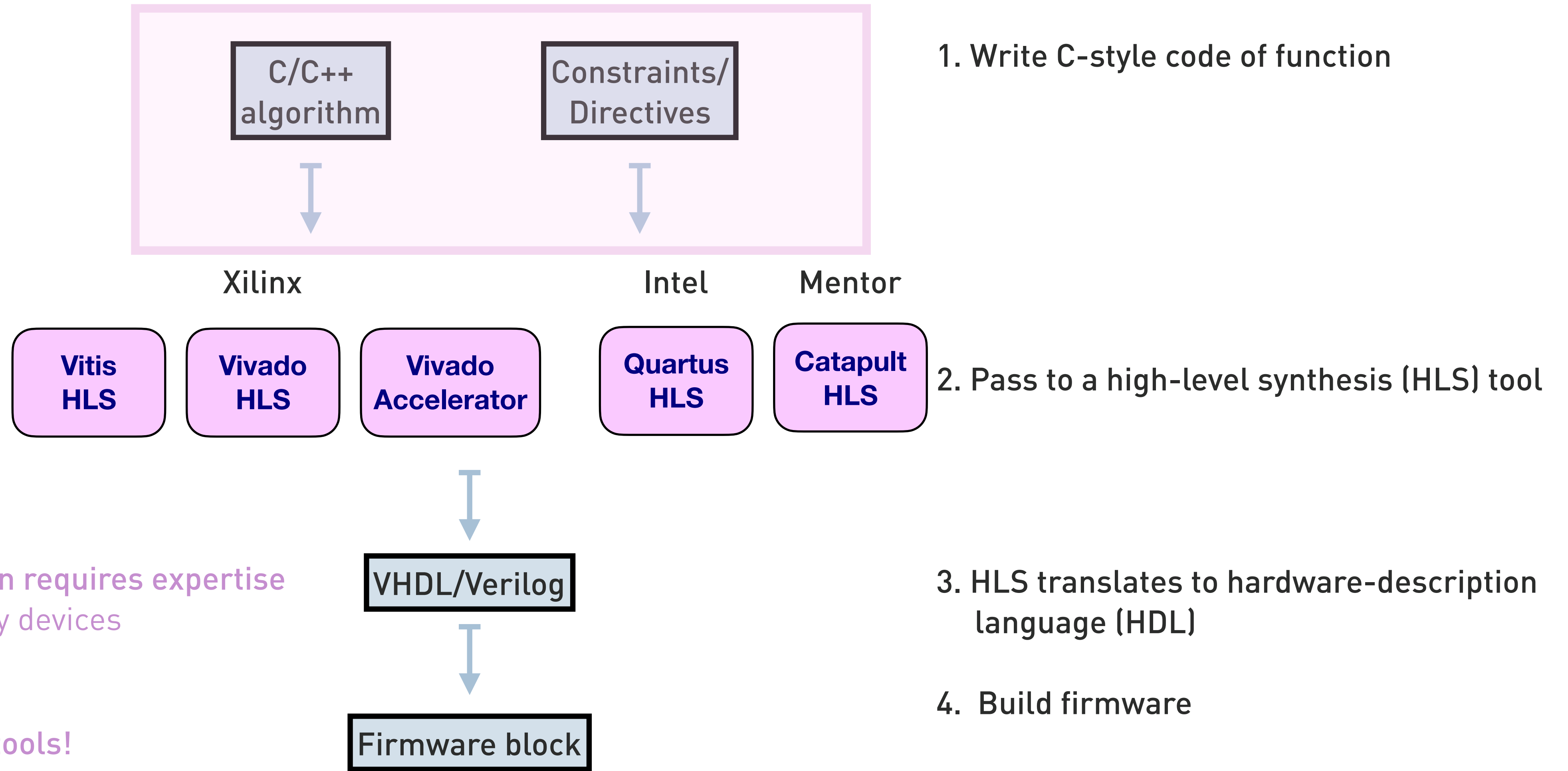precomputed and                DSPs              logic cells
stored in BRAMs

Memory (BRAM)

Digital signal processors (DSPs)

Logic cells/lookup tables (LUTs)

# Programming an FPGA



| | | |
|---|---|---|
| C/C++ algorithm | | Constraints/ Directives |

Xilinx

Intel    Mentor

| Vitis HLS | Vivado HLS | Vivado Accelerator | | Quartus HLS | Catapult HLS |
|---|---|---|---|---|---|

VHDL/Verilog

Firmware block

1. Write C-style code of function

2. Pass to a high-level synthesis (HLS) tool

3. HLS translates to hardware-description language (HDL)

4. Build firmware

# Programming an FPGA

C/C++ algorithm

Constraints/ Directives

Xilinx

Intel    Mentor

**Vitis HLS**

**Vivado HLS**

**Vivado Accelerator**

**Quartus HLS**

**Catapult HLS**

VHDL/Verilog

Firmware block

1. Write C-style code of function

2. Pass to a high-level synthesis (HLS) tool

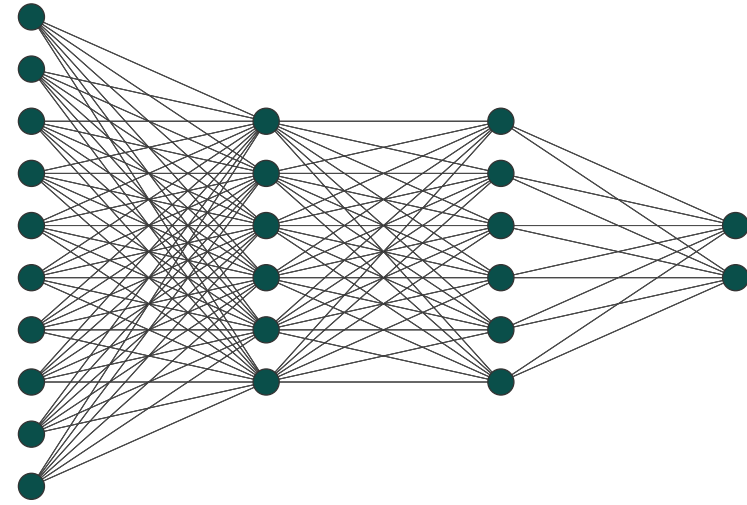3. HLS translates to hardware-description language (HDL)

4. Build firmware

Efficient L1T firmware design requires expertise
- FPGA deployment in busy devices
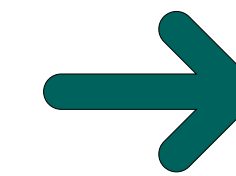- « 1µs latency target
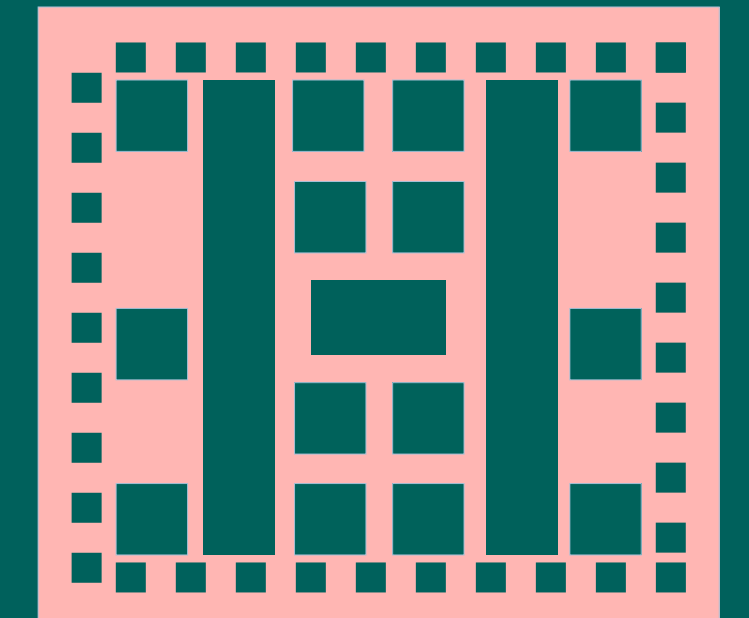
Not well served by industry tools!

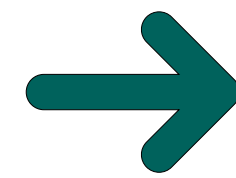**See Riccardo's talk!**

16

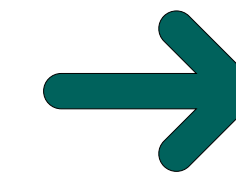TensorFlow / TF Keras / PyTorch / ONNX

scikit-learn / XGBoost / TMVA

hls4ml

Conifer

HLS project:
Xilinx Vitis HLS, Intel Quartus HLS,
Mentor Catapult HLS

```
pip install hls4ml
pip install conifer
```

Keras
TensorFlow
PyTorch
...

model
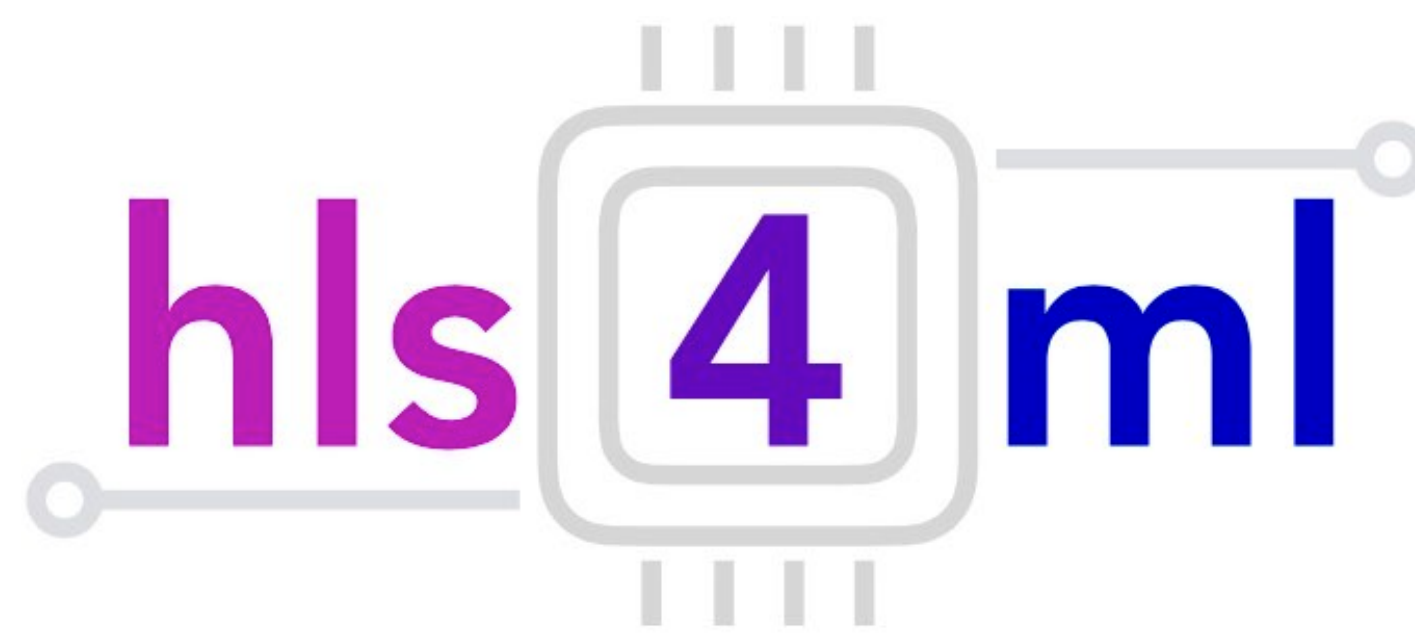
Usual ML
software workflow

compressed
model

1. Pruning
2. Quantization-aware training

hls**4**ml

HLS
conversion

HLS
project

Co-processing kernel

Custom firmware
design

tune configuration
precision
reuse/pipeline

1. Parallelisation
2. Post-training quantization

XILINX®
VIRTEX®
UltraScale+™

**Prediction**

```
from hls4ml import …
import tensorflow as tf

# train or load a model
model = … # e.g. tf.keras.models.load_model(…)

# make a config template
cfg = config_from_keras_model(model,
granularity='name')

# tune the config
cfg['LayerName']['layer2']['ReuseFactor'] = 4

# do the conversion
hmodel = convert_from_keras_model(model, cfg)

# write and compile the HLS
hmodel.compile()

# run bit accurate emulation
y_tf = model.predict(x)
y_hls = hmodel.predict(x)

# do some validation
np.testing.assert_allclose(y_tf, y_hls)

# run HLS synthesis
hmodel.build()
```
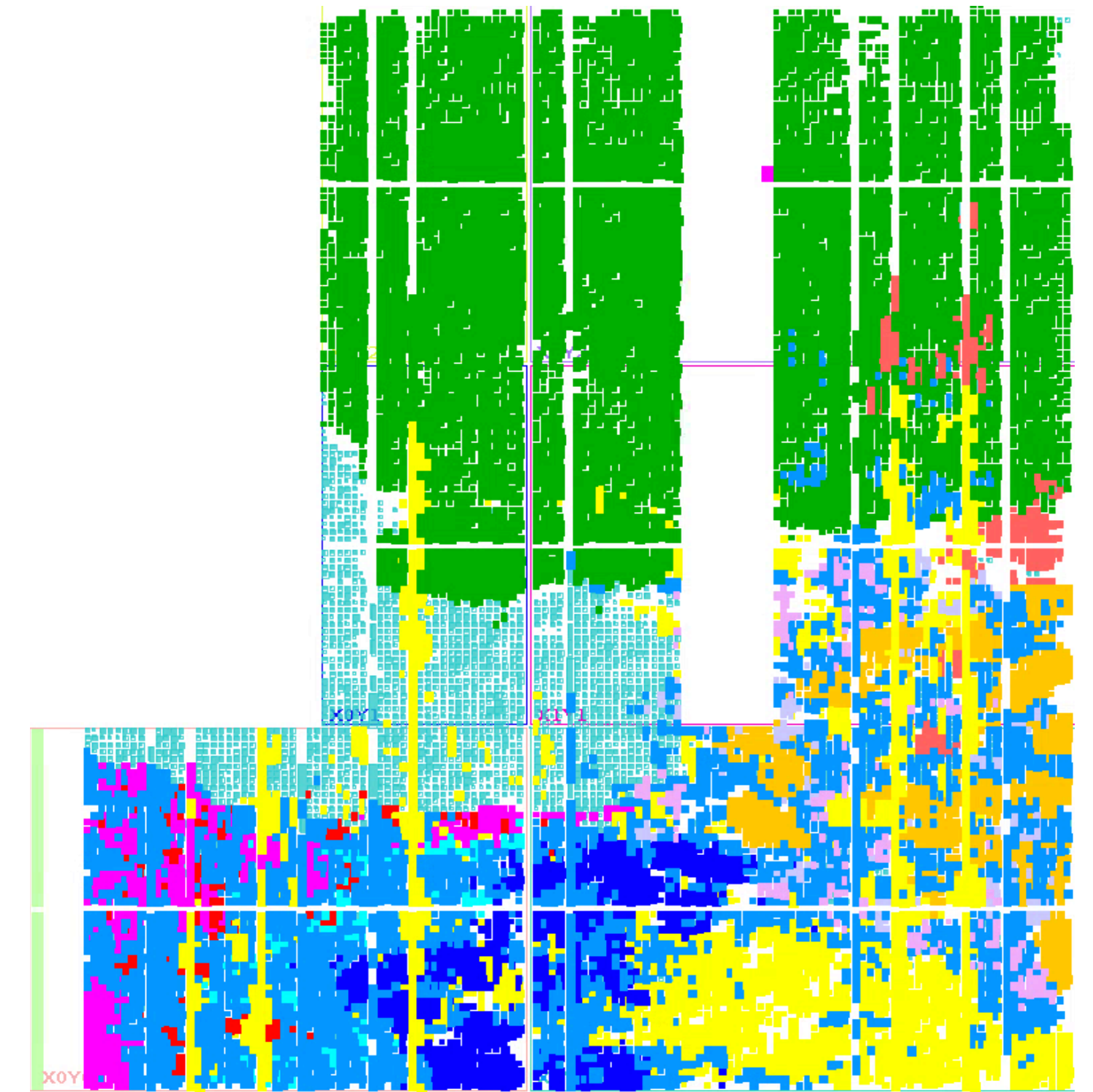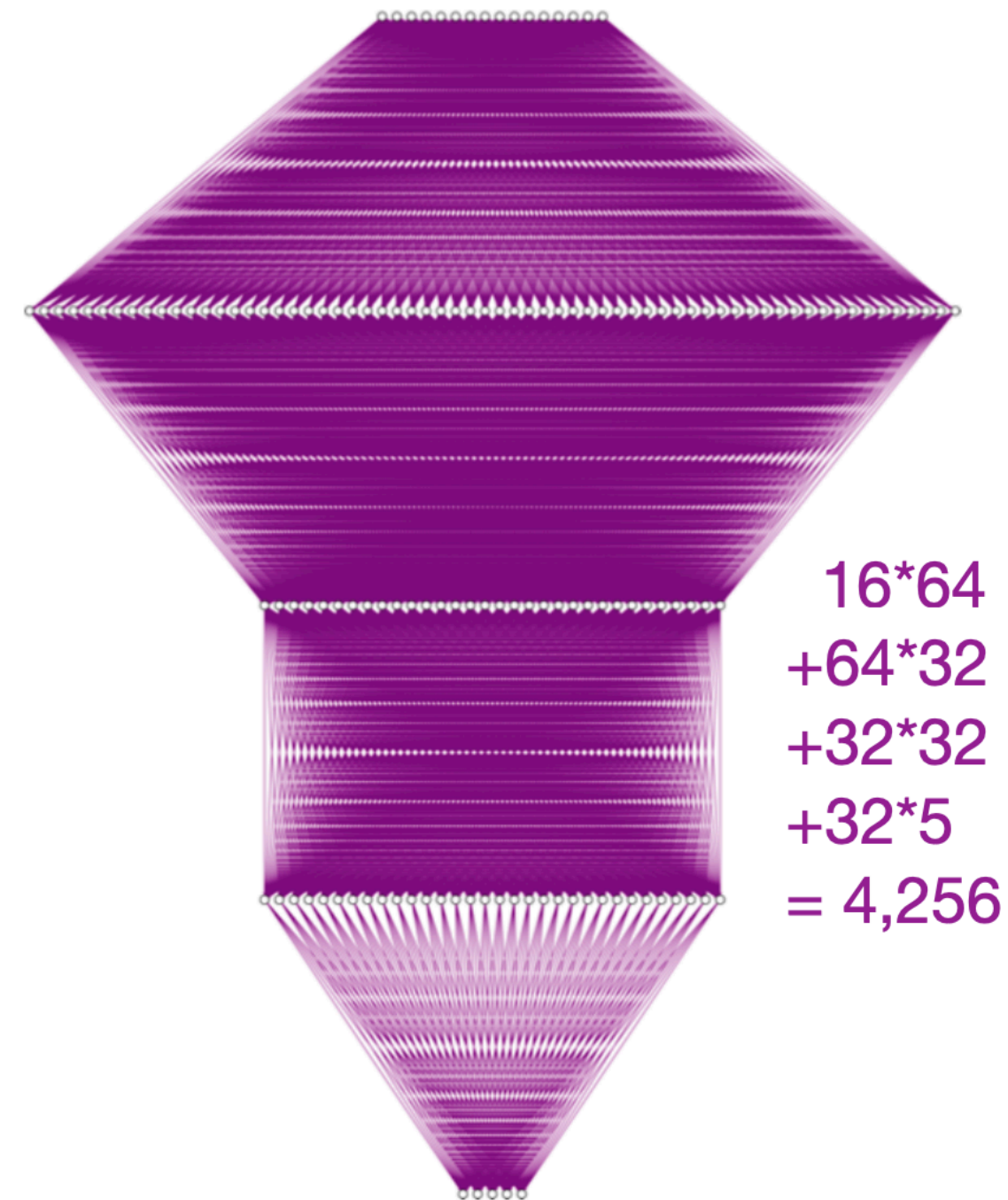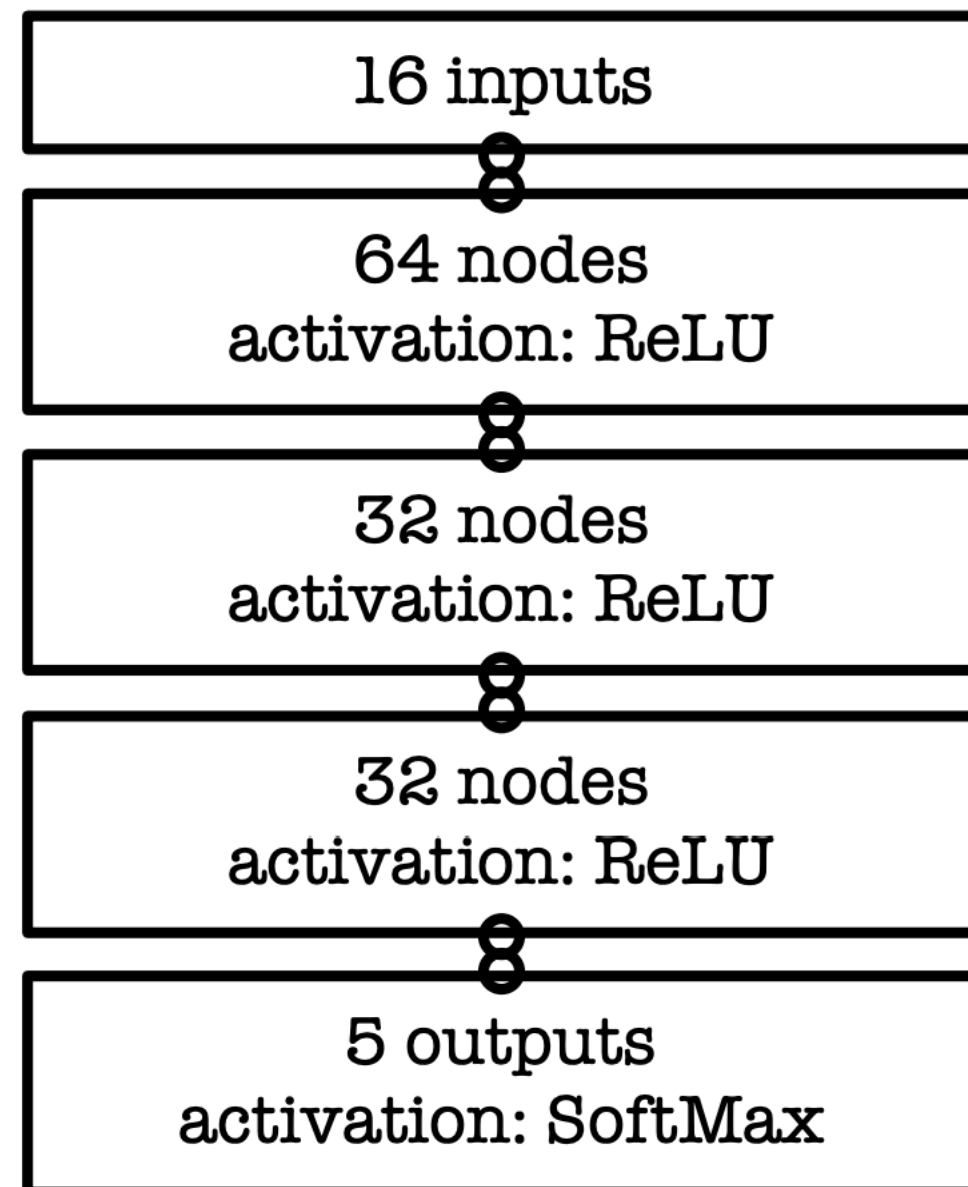
pynq-z2 floorplan

**Learn how to use hls4ml in tomorrows tutorial by Sioni!**

(from Sioni S Summers)

# Compression

| 16 inputs |
|---|

| 64 nodes<br>activation: ReLU |
|---|

| 32 nodes<br>activation: ReLU |
|---|

| 32 nodes<br>activation: ReLU |
|---|

| 5 outputs<br>activation: SoftMax |
|---|

16*64
+64*32
+32*32
+32*5
= 4,256

$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

activation function    multiplication    addition

precomputed and
stored in BRAMs     DSPs     logic cells

**Network size limited by N multiplications**
- E.g, simple dense network, **total multiplications: 4256!**
- A typical FPGA at LHC usually has **4-6000** DSPs
- Can your network fit within the resources?

# Efficient NN design for FPGAs (and other edge compute)



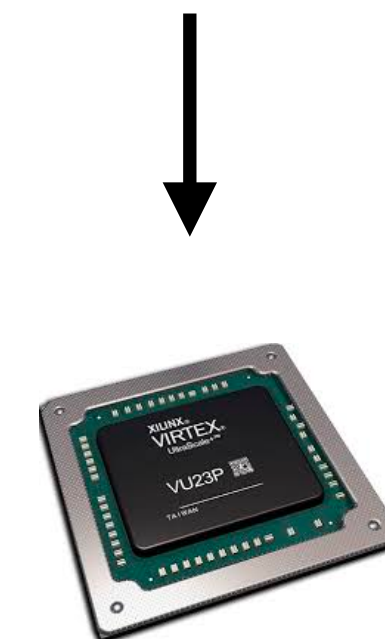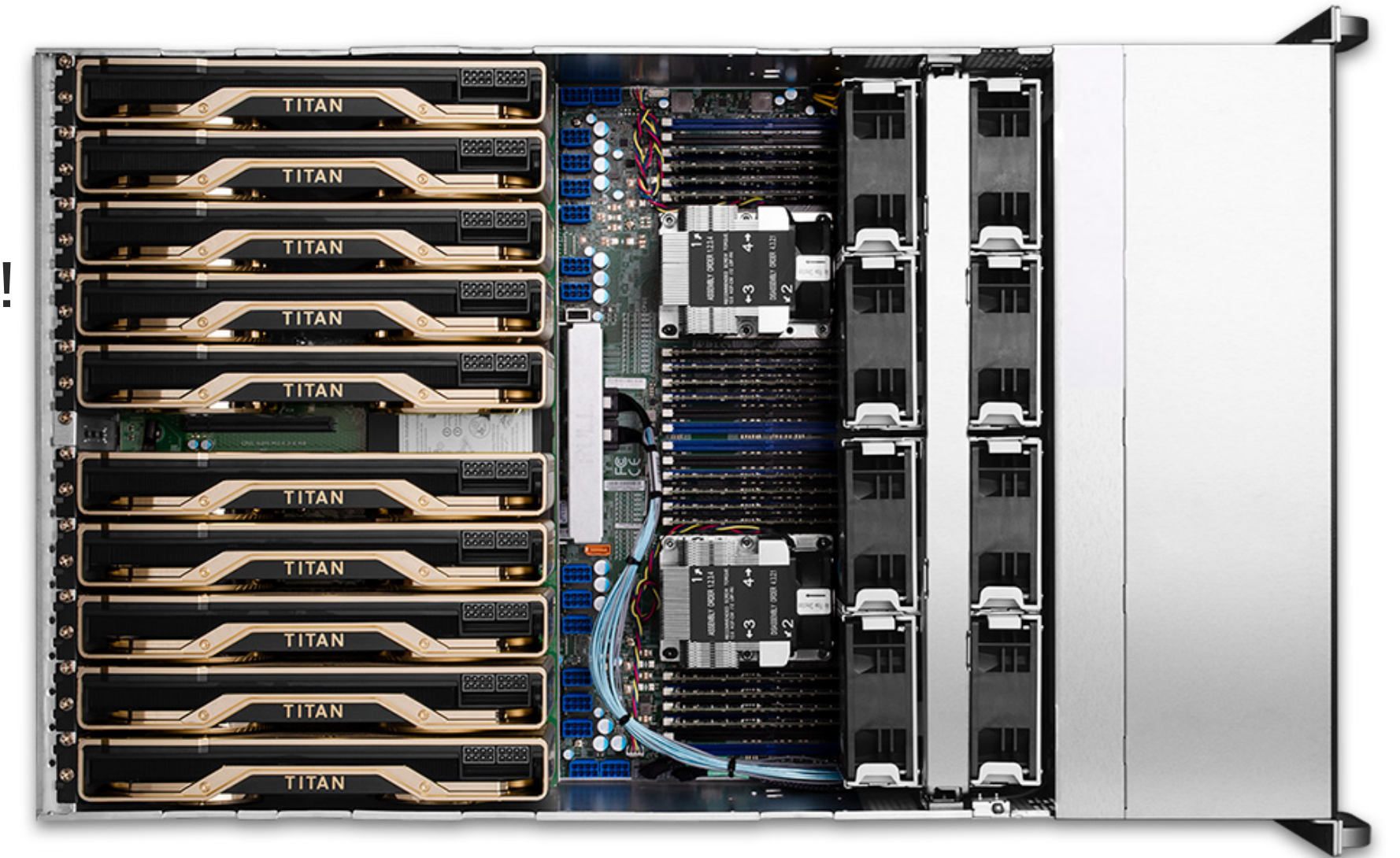Before deploying any DNN on chip (CMS trigger, iPhone), must make it efficient!

During training
- **Quantization**
- Pruning

Post-training
- Parallelisation (lower latency ↔ more resources)

From 8 GPU server to tiny FPGA!



**See RIccardos talk and learn more in tomorrows tutorial by Sioni!**

# Quantization

## Fixed point post-training quantization
- Floating point 32 arithmetic use **x3-5** more resources, **x2** higher latency than fixed-point → convert to fixed-point

Decimal: 3.25

During training: $-1^S \cdot 2^{E-127} \cdot (1.M)$

0100000001010000000000000000000000
S Exponent            Mantissa

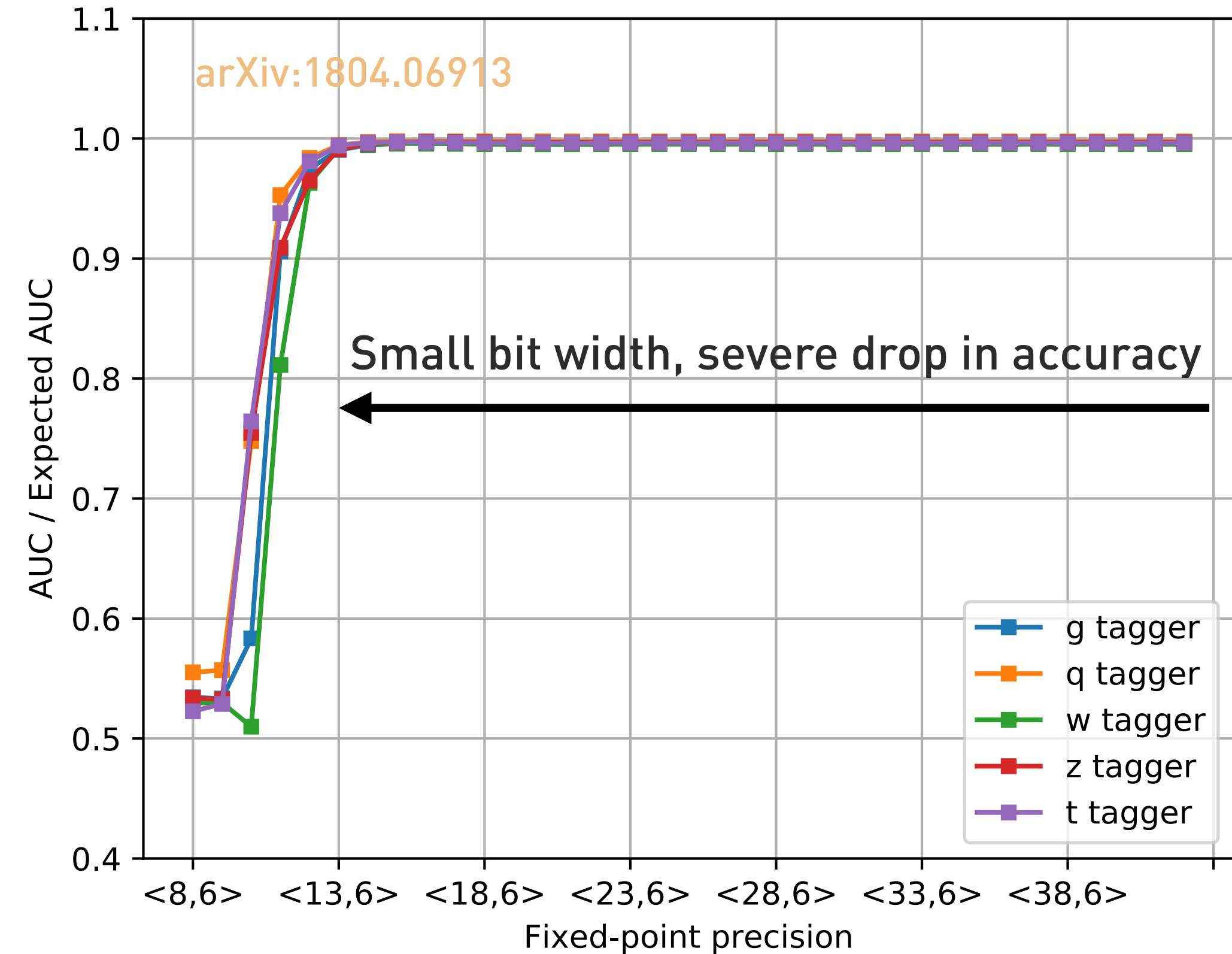On hardware: ap_fixed (W,I)
⟨00011.01⟩

By definition lossy, precision must be tuned carefully (weights usually don't need large dynamic range. But, worse 'resolution')

## Can we do better? Yes!
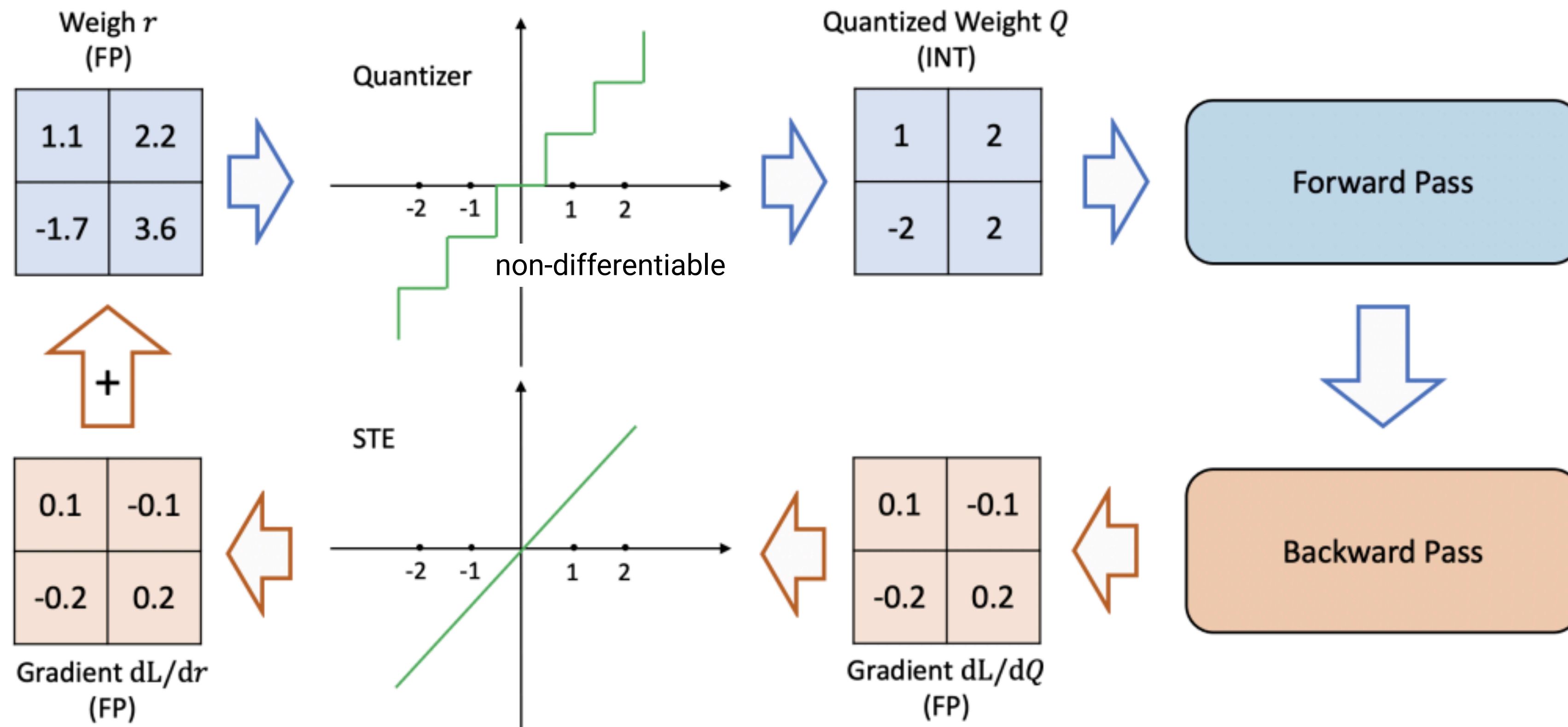- Quantization-aware training (QAT)

**hls4ml**

arXiv:1804.06913

Small bit width, severe drop in accuracy

AUC / Expected AUC

- g tagger
- q tagger
- w tagger
- z tagger
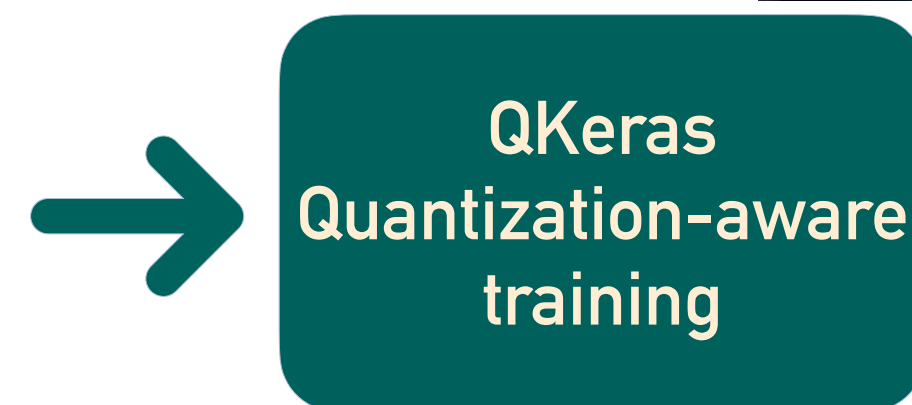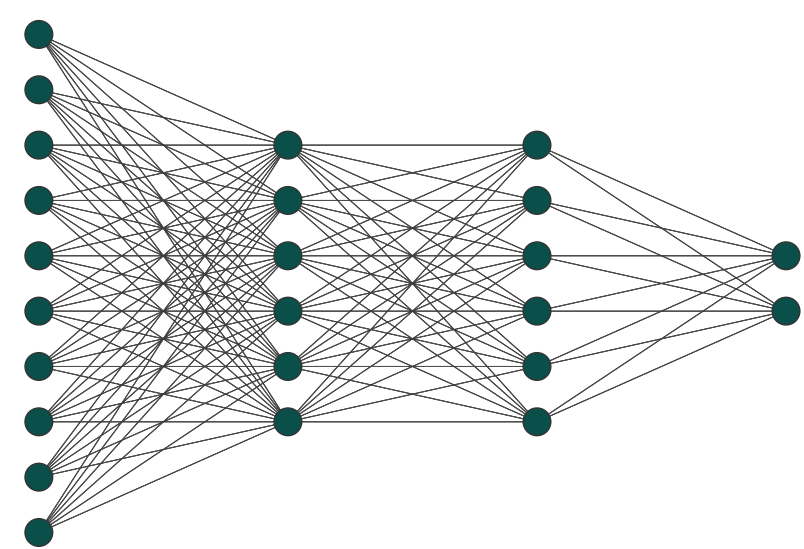- t tagger

Fixed-point precision

**See Riccardo's talk!**

# Quantization-aware training

Lossless quantization for deep neural networks!


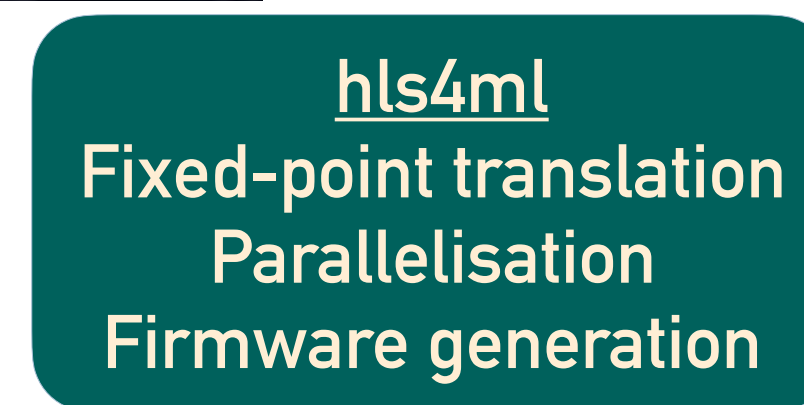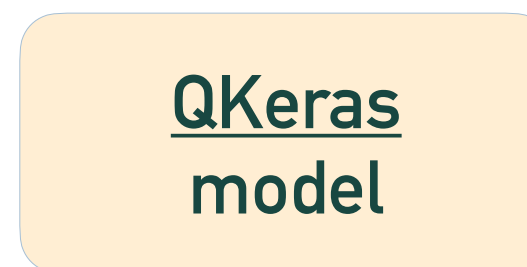
arxiv:2103.13630

**Nature Machine Intelligence 3 (2021)**



QKeras
Quantization-aware
training

QKeras
model

hls4ml
Fixed-point translation
Parallelisation
Firmware generation

Google AI

hls4ml

24

```python
from tensorflow.keras.layers import Input, Activation
from qkeras import quantized_bits
from qkeras import QDense, QActivation
from qkeras import QBatchNormalization

x = Input((16))
x = QDense(64,
    kernel_quantizer = quantized_bits(6,0,alpha=1),
    bias_quantizer   = quantized_bits(6,0,alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(32,
    kernel_quantizer = quantized_bits(6,0,alpha=1),
    bias_quantizer   = quantized_bits(6,0,alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(32,
    kernel_quantizer = quantized_bits(6,0,alpha=1),
    bias_quantizer   = quantized_bits(6,0,alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(5,
    kernel_quantizer = quantized_bits(6,0,alpha=1),
    bias_quantizer   = quantized_bits(6,0,alpha=1))(x)
x = Activation('softmax')(x)
```

```python
from hls4ml import …
import tensorflow as tf

# train or load a model
model = tf.keras.models.load_model(…)

# make a config
cfg = config_from_keras_model(model,
granularity='name')

# do the conversion
hmodel = convert_from_keras_model(model, cfg)

# write and compile the HLS
hmodel.compile()

# run HLS synthesis
hmodel.build()
```
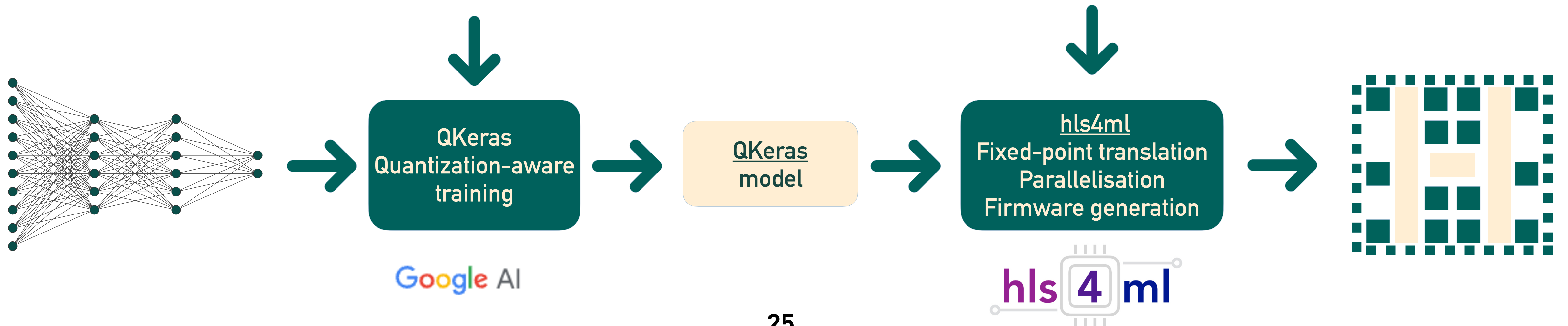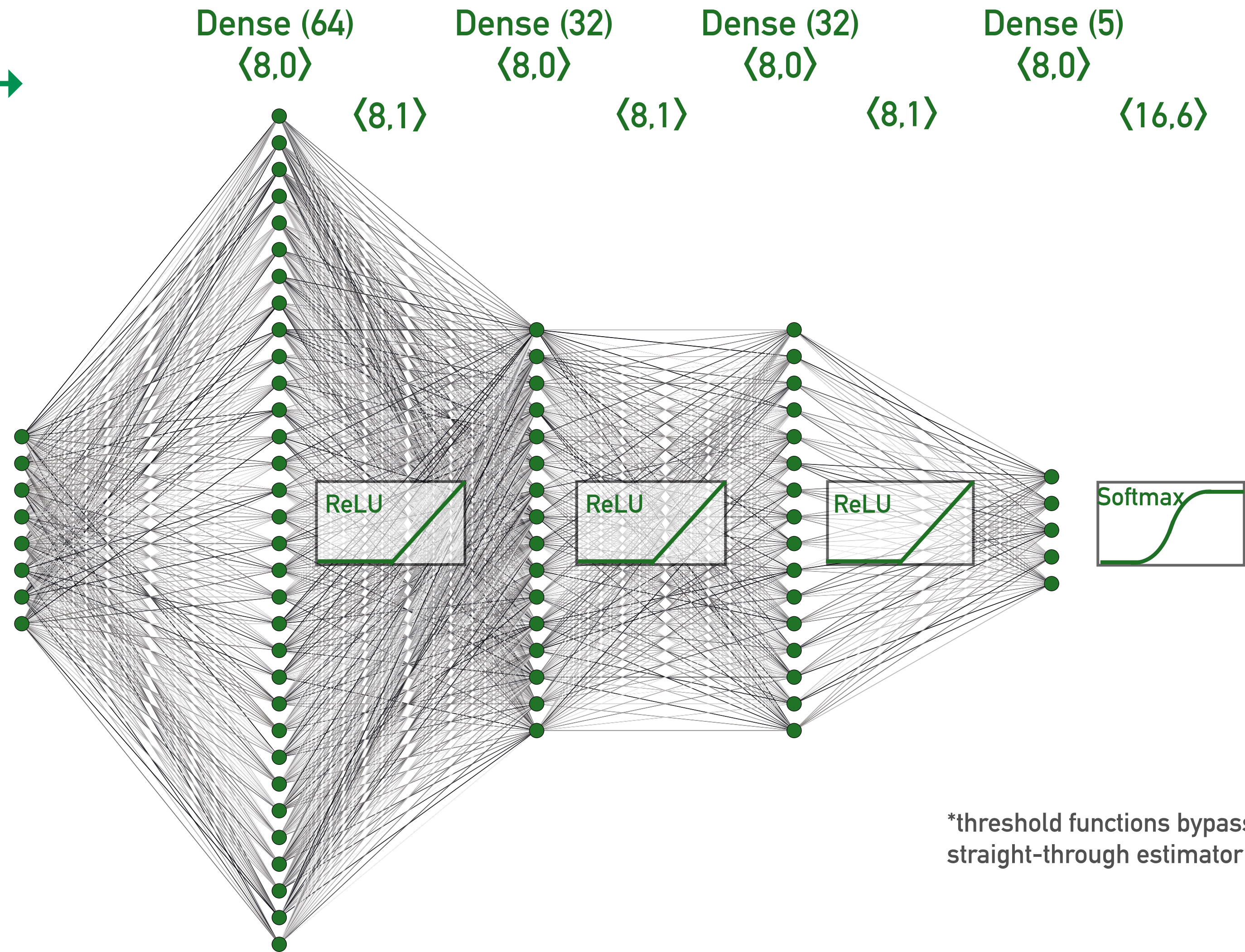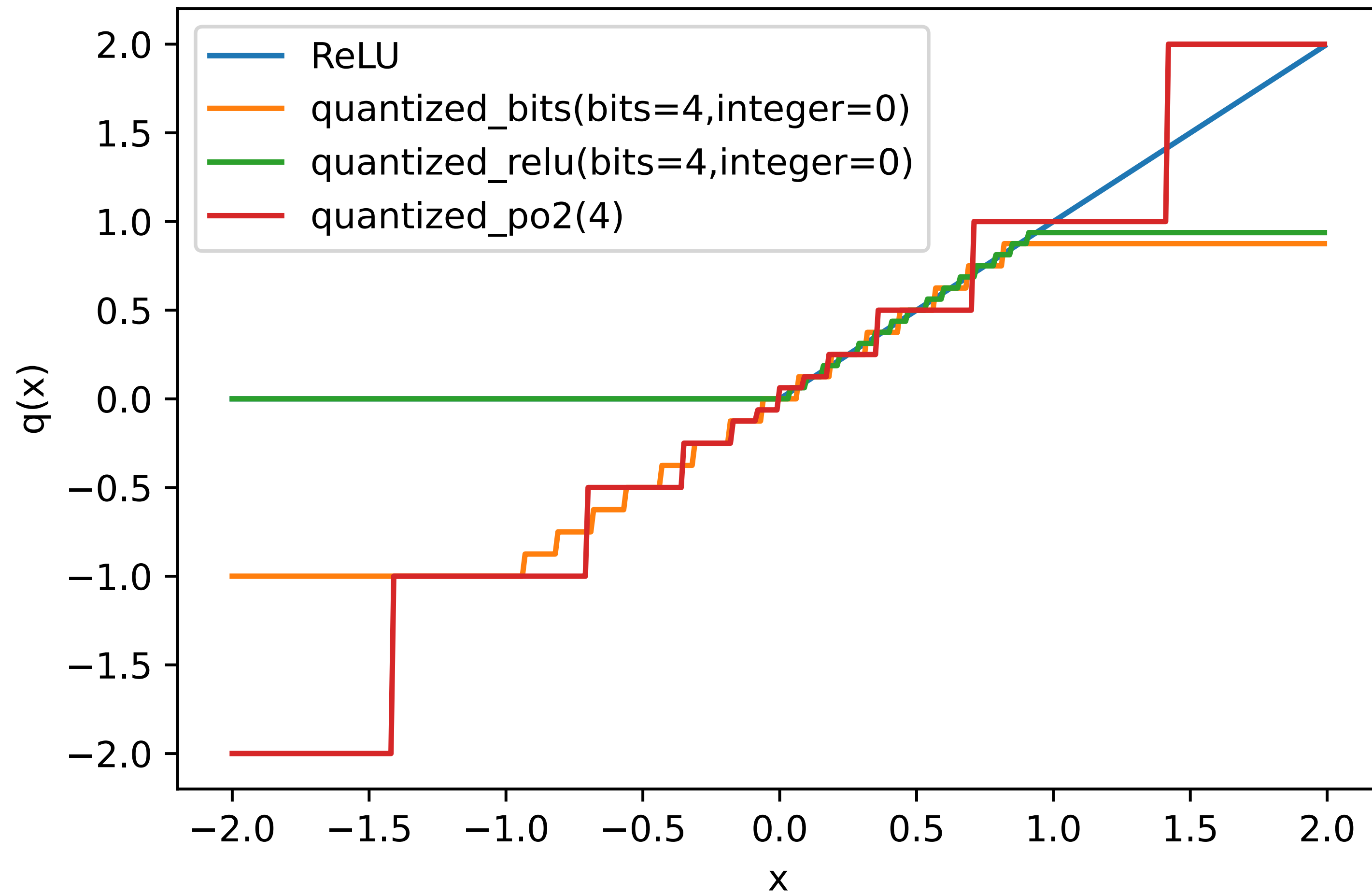


QKeras
Quantization-aware
training

QKeras
model

hls4ml
Fixed-point translation
Parallelisation
Firmware generation

Google AI

hls4ml

Forward pass →

Dense (64)
⟨8,0⟩

⟨8,1⟩

Dense (32)
⟨8,0⟩

⟨8,1⟩

Dense (32)
⟨8,0⟩

⟨8,1⟩

Dense (5)
⟨8,0⟩

⟨16,6⟩

ReLU

ReLU

ReLU

Softmax

*threshold functions bypassed in backward pass,
straight-through estimator

← Back propagation

FP 32

FP 32

FP 32

FP 32

FP 32

FP 32

FP 32

FP 32

# Quantization-aware training

# FPGA performance

28

**Ideally**

**Reality**

# QTools energy estimate

Some layers more accommodating for aggressive quantization, others require expensive arithmetic
- heterogeneous quantization

# QTools energy estimate

Some layers more accommodating for aggressive quantization, others require expensive arithmetic
- heterogeneous quantization

For edge inference, need best possible quantization configuration for
- Highest accuracy ↑...
- ... and lowest resource consumption ↓

→ hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

# QTools energy estimate

Some layers more accommodating for aggressive quantization, others require expensive arithmetic
- heterogeneous quantization

For edge inference, need best possible quantization configuration for
- Highest accuracy ↑…
- … and lowest resource consumption ↓

→ hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

QTools: Estimate QKeras model bit and energy consumption, assuming 45 nm Horowitz process
- Model size in bits
- Energy consumption in Watts

| Model | Accuracy [%] | Per-layer energy consumption [pJ] | | | | | | | | Total energy [$\mu$J] | Total bits |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Dense | ReLU | Dense | ReLU | Dense | ReLU | Dense | Softmax | | |
| **BF** | 74.4 | 1735 | 53 | 3240 | 27 | 1630 | 27 | 281 | 11 | 0.00700 | 61446 |
| **Q6** | 74.8 | 794 | 23 | 1120 | 11 | 562 | 11 | 99 | 11 | 0.00263 | 26334 |

$$\text{Forgiving Factor} = 1 + \Delta_{accuracy} \times \log_{rate}(S \times \frac{Cost_{ref}}{Cost_{trial}})$$

Maximize accuracy + minimizing cost in hyper parameter scan over quantizers:
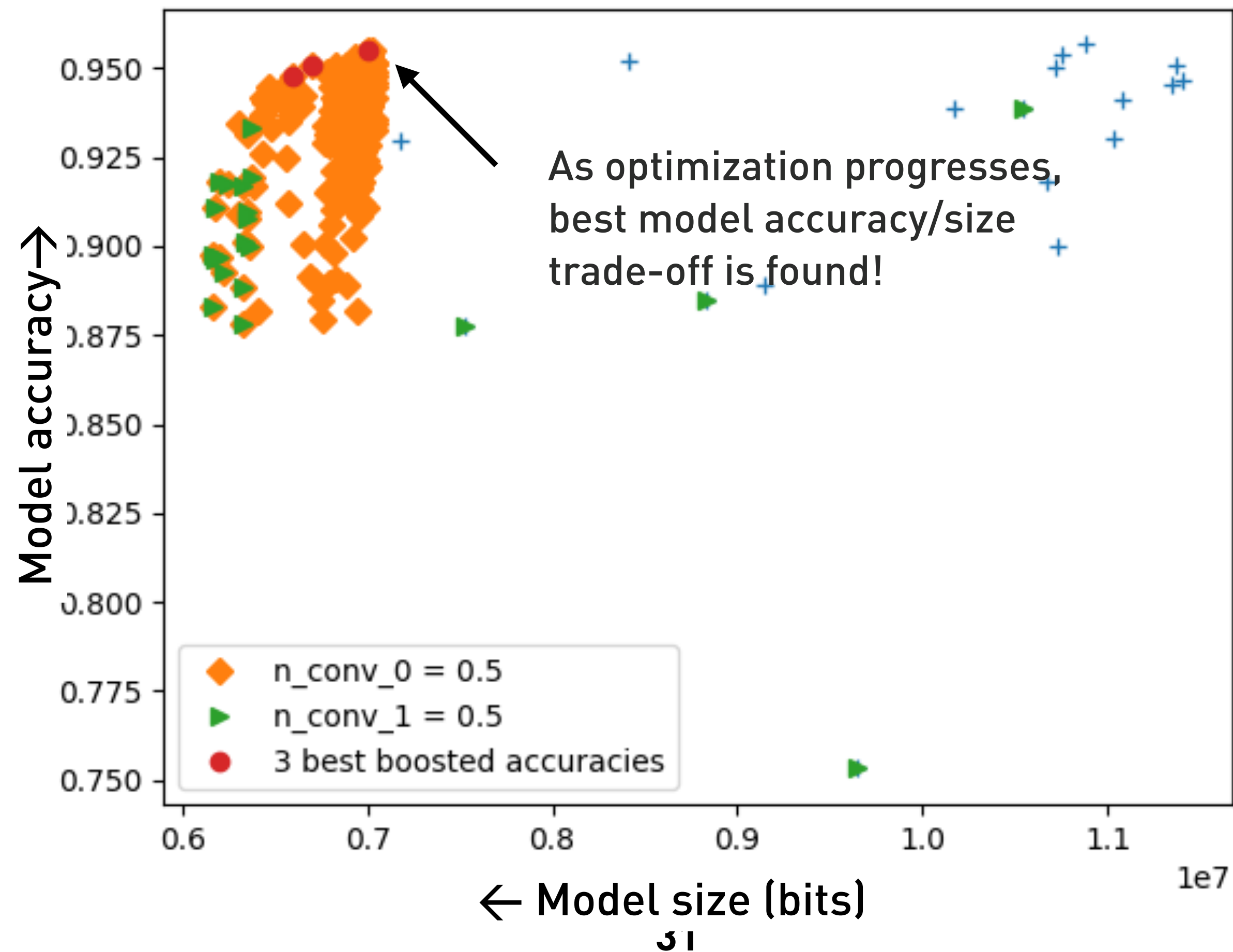
AutoQKeras

# AutoQKeras

**AutoQ Bayesian optimization at work!**
- Simultaneously scan quantizers and N filters/neurons  (often less/more filters/neurons needed when quantizing)
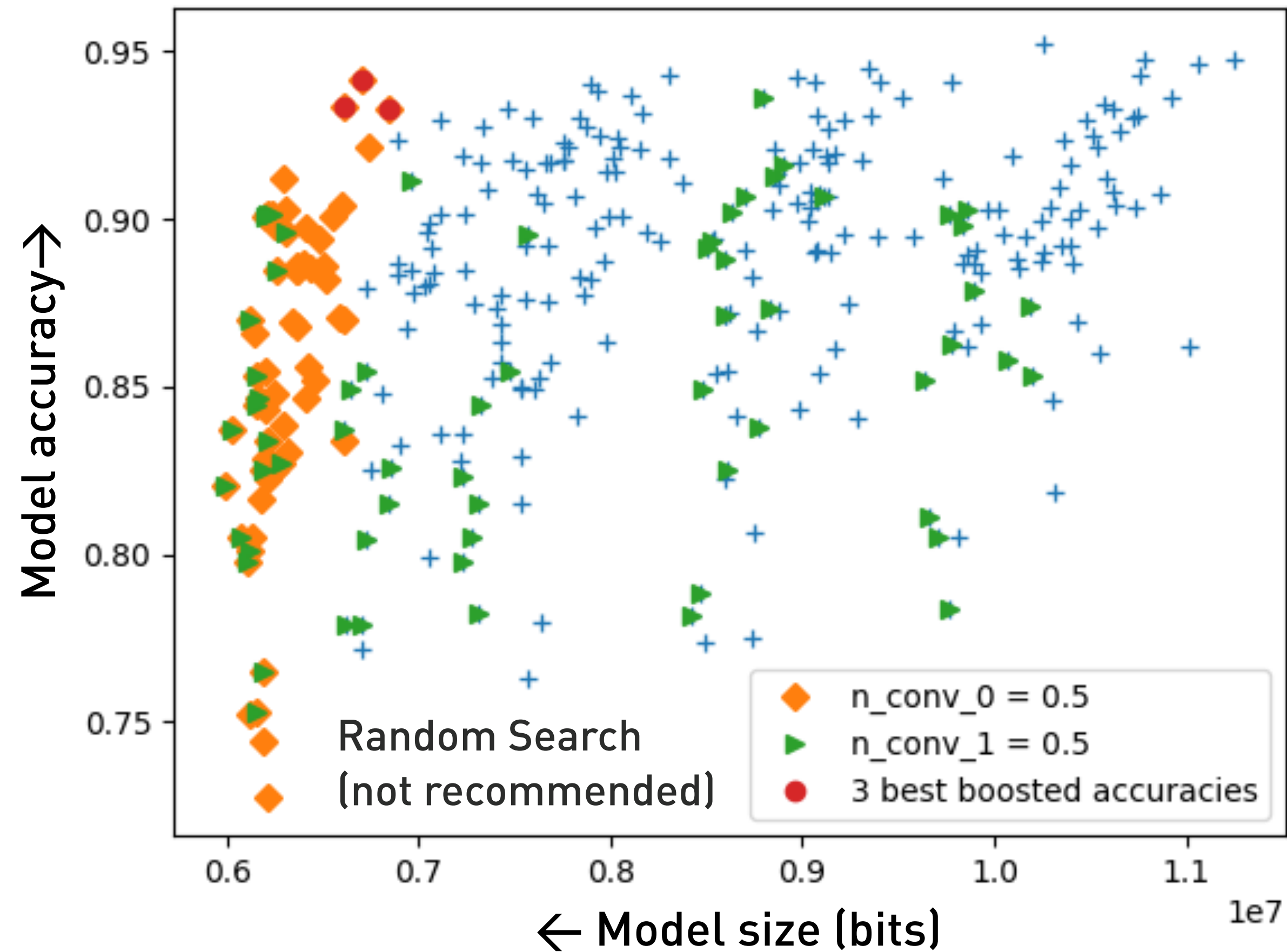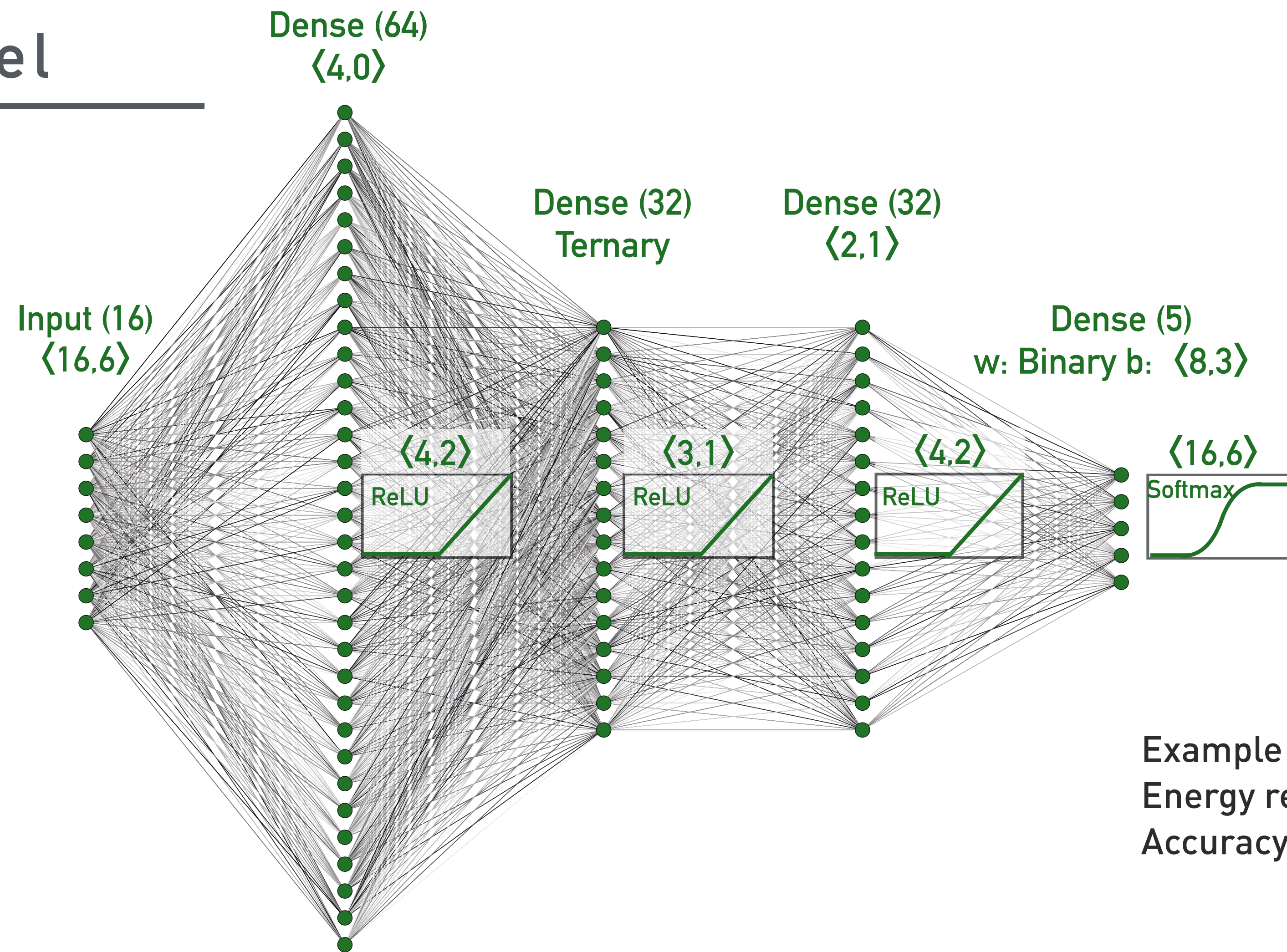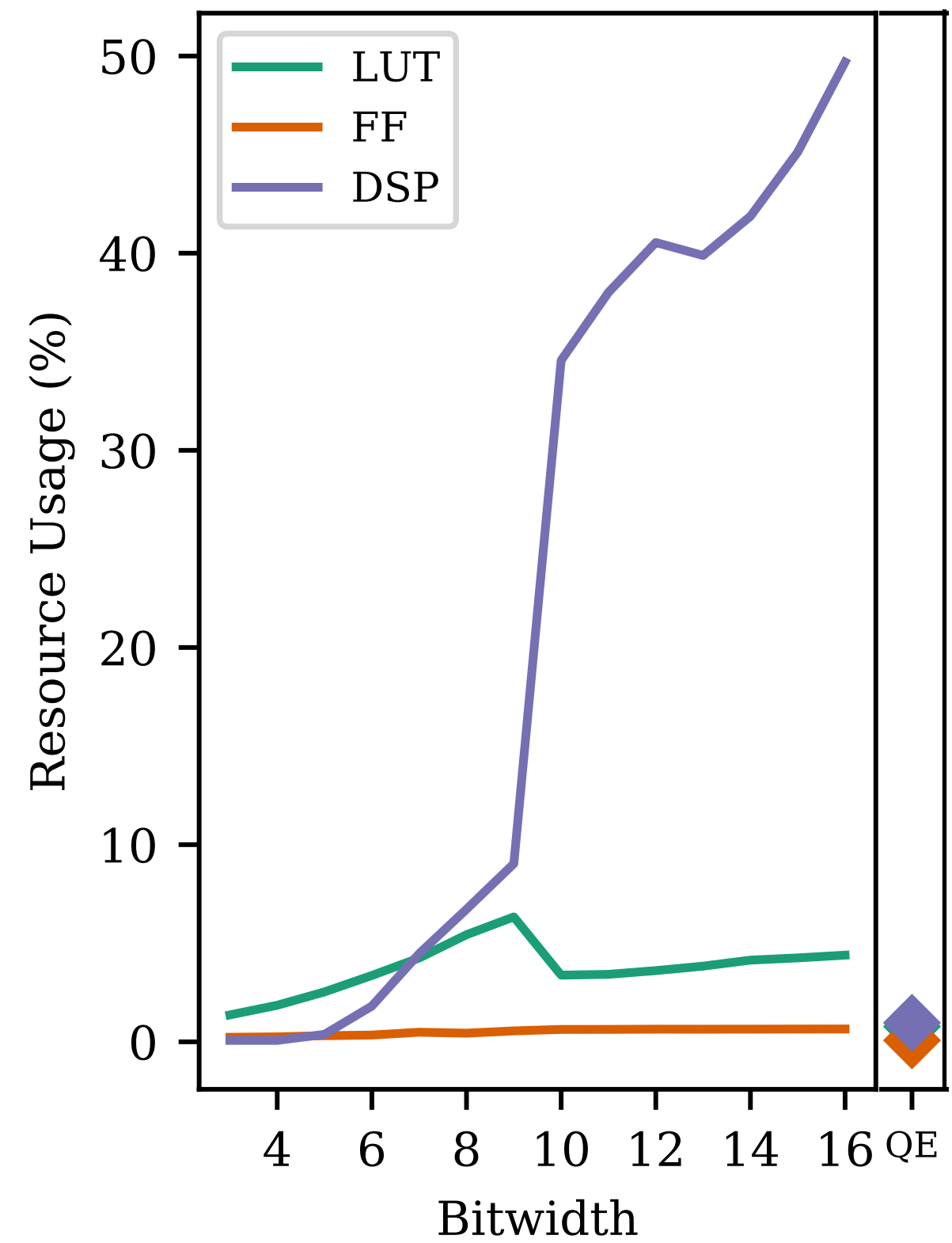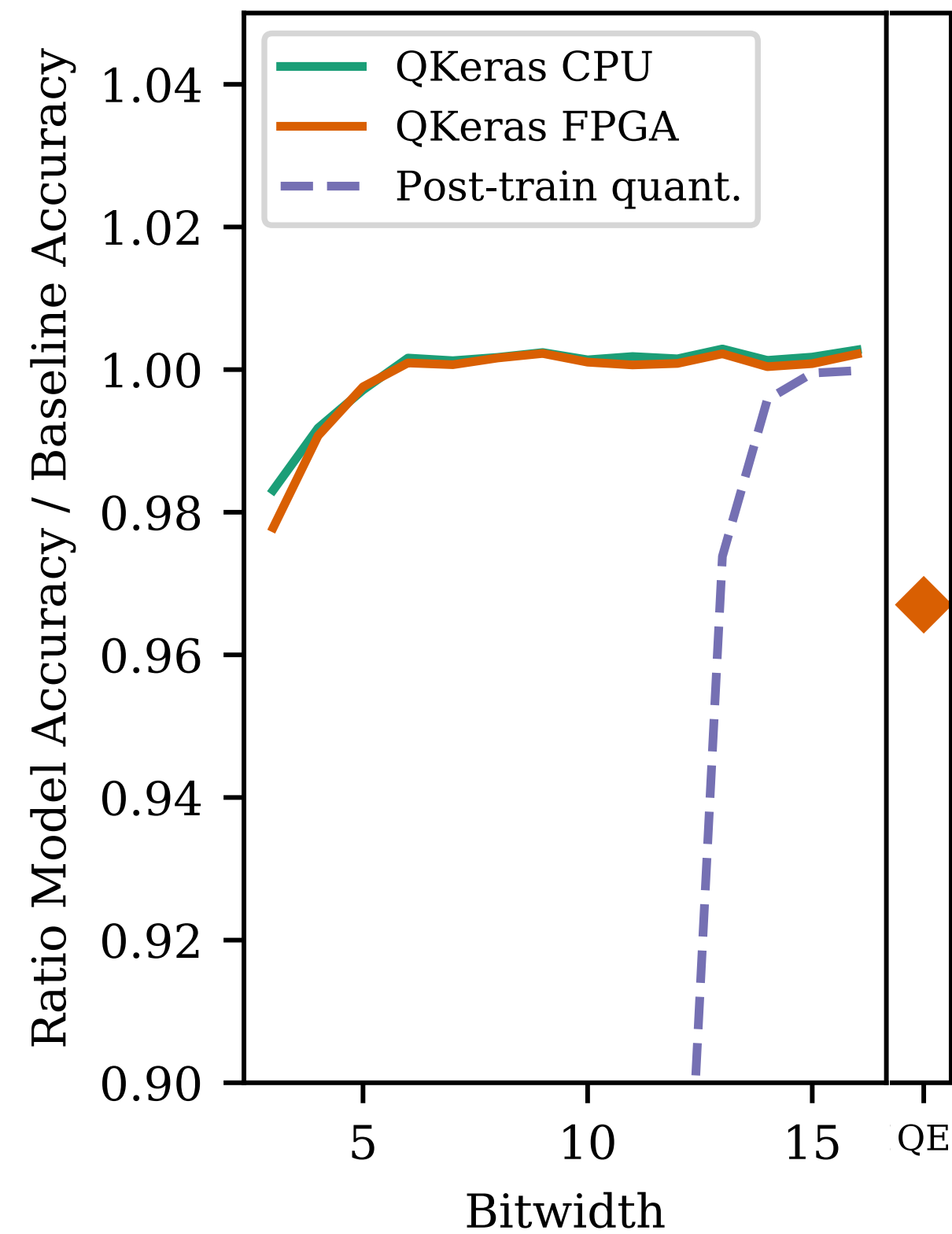


As optimization progresses, best model accuracy/size trade-off is found!

# AutoQKeras

**AutoQ Bayesian optimization at work!**
- Simultaneously scan quantizers and N filters/neurons  (often less/more filters/neurons needed when quantizing)



Random Search
(not recommended)

Legend:
- n_conv_0 = 0.5
- n_conv_1 = 0.5
- 3 best boosted accuracies

Model accuracy →

← Model size (bits)

max 5%

| Model | Acc. [%] | Precision | | | | | | | | | Tot. energy [$\mu$J] | Tot. bits |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | Dense | ReLU | Dense | ReLU | Dense | ReLU | Dense | | Softmax | | |
| **QE** | 72.3 | $\langle 4, 0 \rangle$ | $\langle 4, 2 \rangle$ | Ternary | $\langle 3, 1 \rangle$ | $\langle 2, 1 \rangle$ | $\langle 4, 2 \rangle$ | w: Stoc. Bin. b: $\langle 8, 3 \rangle$ | | $\langle 16, 6 \rangle$ | 0.00095 | 4728 |

# FPGA performance



| Model | Accuracy [%] | Latency [ns] | Latency [clock cycles] | DSP [%] | LUT [%] | FF [%] |
|---|---|---|---|---|---|---|
| **BF** | 74.4 | 45 | 9 | 56.0 (1826) | 5.2 (48321) | 0.8 (20132) |
| **Q6** | 74.8 | 55 | 11 | 1.8 (124) | 3.4 (39782) | 0.3 (8128) |
| **QE** | 72.3 | 55 | 11 | **1.0 (66)** | **0.8 (9149)** | 0.1 (1781) |

Conifer

Same as hls4ml but for Boosted decision trees (scikit-learn, XGBoost)

If resource/latency constainted, BDT might be the way to go
- Depending on your data, might be as accurate as a DNN
- Usually significantly faster and more resource efficient

| %VU9P | Accuracy | Latency | DSP | LUT |
|---|---|---|---|---|
| QKeras 6b | 75.6% | 40 ns | 22 (~0%) | 1% |
| sklearn + conifer | 74.9% | 5 ns | - | 0.5% |

hls 4 ml

Conifer

Join the community:
fastmachinelearning.org

# Backup

We have infinite time and resources to train huge networks

We have infinite time and resources to train huge networks

but very little for inference

Can we have the best of both worlds?

→ Knowledge Distillation

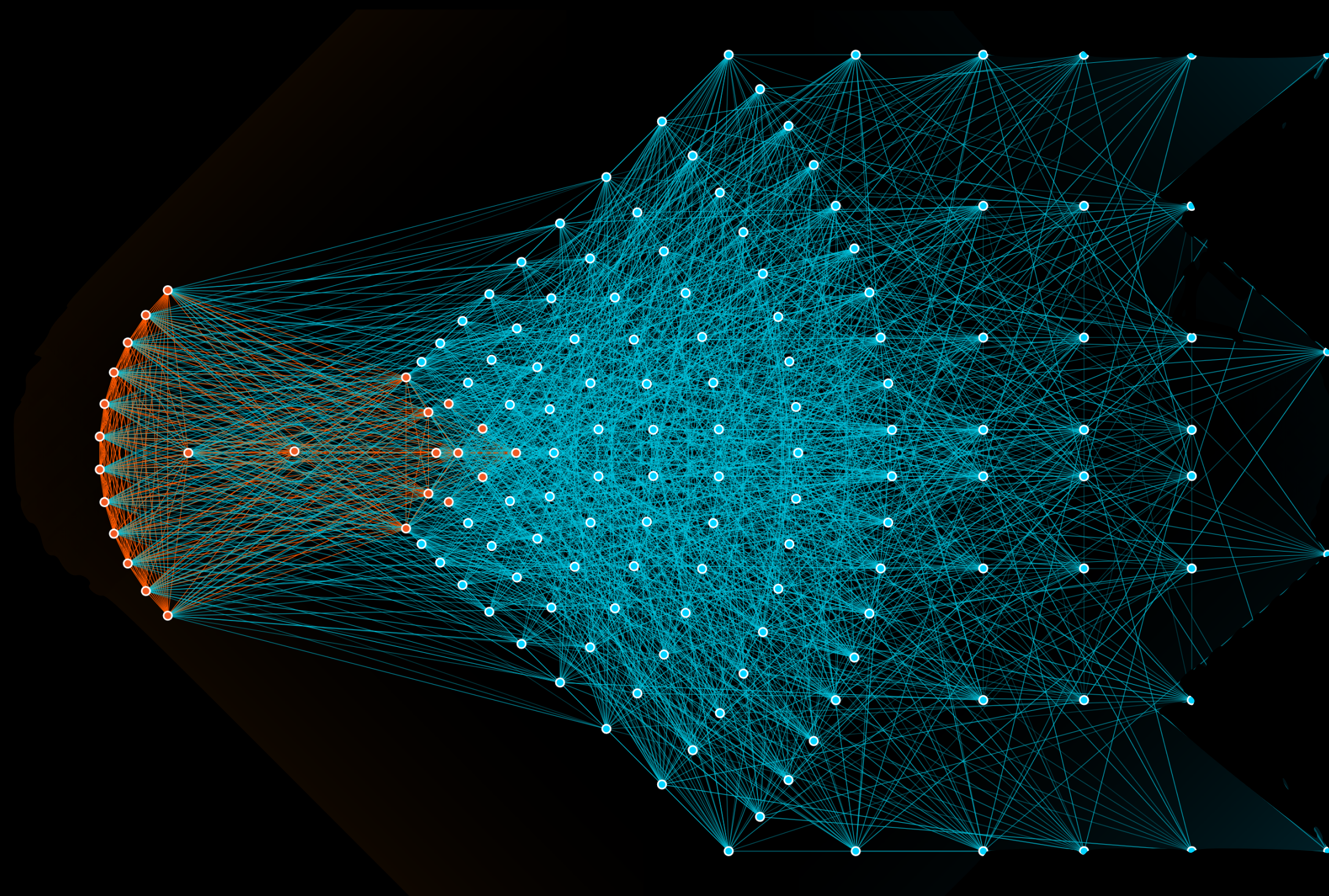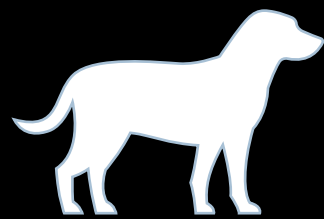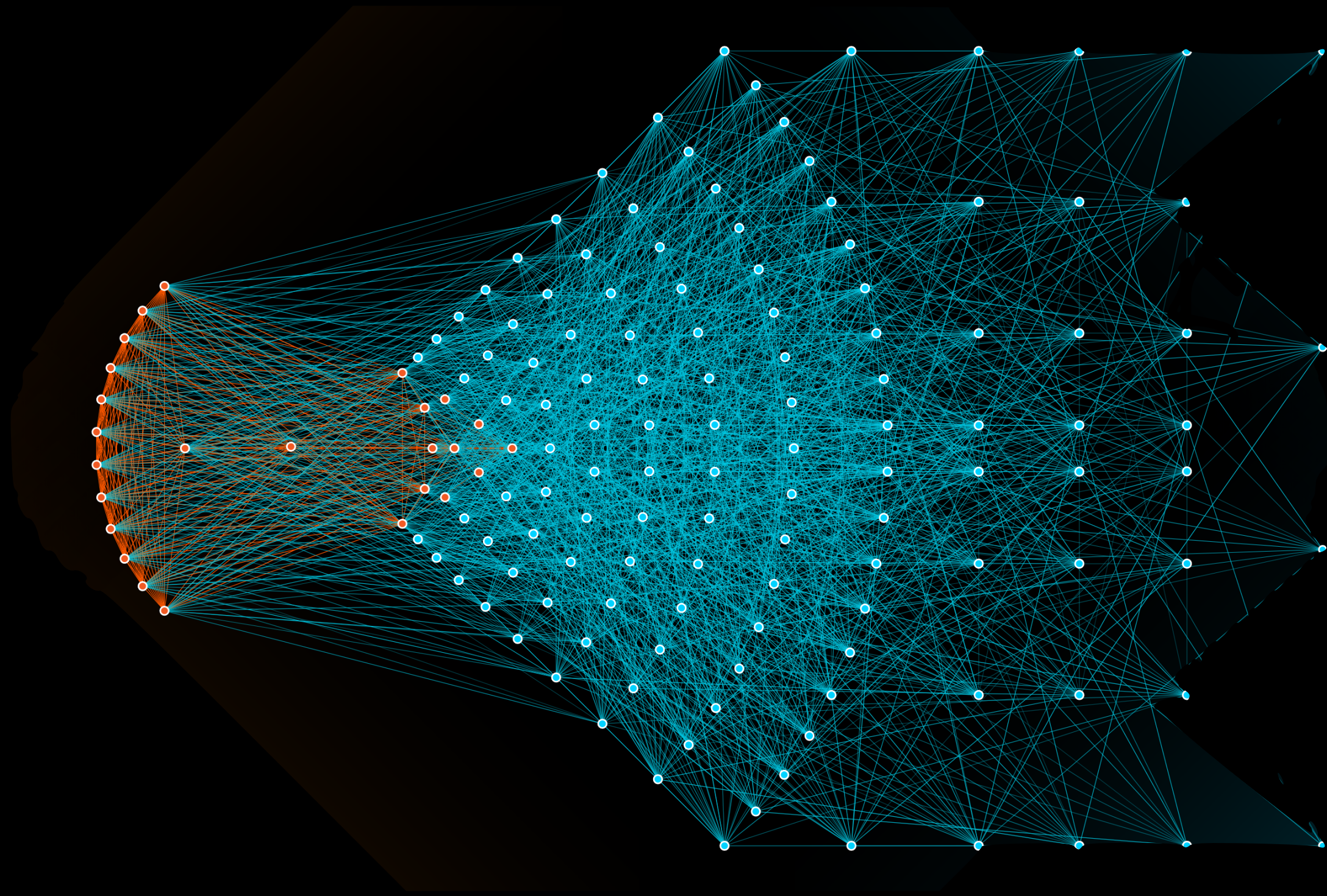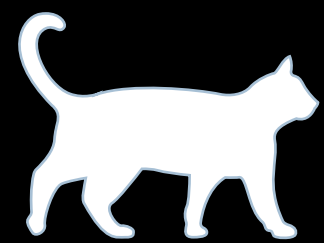Train

Inference

Cat

Dog

Cat

Predicted labels

Teacher
(already trained)

is cat = 0.89

is dog = 0.11

Cat

Teacher (already trained)

Predicted labels

is cat = 0.89

is dog = 0.11

True labels

is cat = 1

is dog = 0

Cat

Teacher
(already trained)

Predicted labels

is cat = 0.89

is dog = 0.11

Distilled knowledge

True labels

is cat = 1

is dog = 0

Train student to learn both
true and predicted (teacher) labels!

$$L_{total} = \beta \times L_{Distillation} + \alpha \times L_{student}$$

Student learns subtle learned features from teacher!