

Imperial College 25.01.2023

# Ultrafast Machine Learning at the Large Hadron Collider

Thea Klæboe Årrestad (ETH Zürich)





Inspire: ("machine learning" or "deep learning" or neural) and (hep-ex or hep-ph or hep-th)

### Date of pape



Selected Papers: 420 Total Papers: 420 Year: 2022







#### <u>Nature Review</u>



Dataset increase factor for 5o discovery

Google	chat gpt X 🕹 🤇		
	Q All 🗉 News 🖾 Images 🕞 Videos 🖺 Books : More	Tools	
	About 10,700 results (0.31 seconds)		
	👅 The New York Times	J. CO	
	Opinion   ChatGPT Has a Devastating Sense of Humor	TTO AK	
	The chat bot makes a lot of mistakes. But it's fun to talk to, and it knows its limitations.	TSA	
	5 weeks ago		
	The New York Times	1 4 <sup>-1</sup>	
	Can ChatGPT Make This Podcast?		
	It's writing podcast scripts, finishing students' homework and correcting mistakes in computer code: ChatGPT, the A.I. chatbot from OpenAI,		
	4 weeks ago		
	The New York Times		
	How to Use ChatGPT and Still Be a Good Person		
	It's a turning point for artificial intelligence, and we need to take advantage of these tools without causing harm to ourselves or others.		
	2 weeks ago		
	The New York Times	The Daily is a sh	
	Did Artificial Intelligence Just Get Too Smart?	Hosted by Michael E	
	The power and potential of a technology called ChatGPT have led some to claim it heralds a new era in computing.	Their team of journa They cover the news	
	3 weeks ago		
	The New York Times		
	ChatGPT is Social Media's Newest Star		
	Social media's newest star is a robot: a program called ChatGPT that tries to answer questions like a person. Since its debut last week		











<u>GPT-3: 175 billion parameters (0.16% of the human brain)</u>





<u>GPT-3: 175 billion parameters (0.16% of the human brain)</u>





T1037 / 6vr4 90.7 GDT (RNA polymerase domain) **T1049 / 6y4f** 93.3 GDT (adhesin tip)

Experimental result
Computational prediction

#### AlphaFold nature cover











T1037 / 6vr4 90.7 GDT (RNA polymerase domain) **T1049 / 6y4f** 93.3 GDT (adhesin tip)

Experimental result
Computational prediction

#### AlphaFold nature cover









#### <u>100 million jets for training</u>

"Particle Transformer For Jet Tagging" H. Qu, C. Li, S. Qian







#### 0(1) ms



#### ASIC/GPU





#### 0(1) ms





#### ASIC/GPU



ASIC



#### 0(1) ms





#### ASIC/GPU



0(1) ns

ASIC

FPGA







#### CMS Experiment at the LHC, CERN

Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST) Run / Event: 151076/1405388







#### CMS Experiment at the LHC, CERN

Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST) Run / Event: 151076/1405388

### billion collisions per second -TPB of data per second









~0(1) billion collisions per second ~0(1) PB of data per second





### CMS CT

Geneva Lake



#### High Level Trigger: Latency 0(100) ms

CMS

A DECEMBER OF A DECEMBER OF

DATA ~99.75% of events rejected! O(100) kHz ~Tb/s

Geneva



#### High Level Trigger: Latency 0(100) ms

### CMS

#### DATA ~99.75% of events rejected! O(100) kHz ~Tb/s

Geneva



### ATLAS ~0.02% of collision events remaining ALICE

DATA ~99.98% of events rejected O(1) kHz ~Gb/s

Geneva

LHC

LHCh



#### High Level Trigger: Latency 0(100) ms

(intel)

Xeon\* 7500

Genev

### CMS

#### DATA 100 kHz ~Tb/s

**Detector:** 40 MHz ~Pb/s

Level-1 trigger: Latency O(1) µs

LHCh



### High Luminosity LHC

#### New Physics is produced 1 in 10<sup>12</sup>

• Need <u>more collisions</u> to observe rare processes

### High Luminosity LHC

#### New Physics is produced 1 in 10<sup>12</sup>

• Need <u>more collisions</u> to observe rare processes

#### High Luminosity HLC

- $\bullet$   $\times 10$  increase in data size
- ×3 collisions per second

#### How

- ×2 protons per bunch
- Squeeze beam at interaction point (B\*)





/m

#### ructure $\rightarrow$ pile-up of ~ 60 events/x-ing (s/x-ing)





### High Luminosity LHC

200 vertices (average 140)



### High Luminosity LHC

Must maintain physics acceptance  $\rightarrow$  better detectors

CMS High Granularity (endcap) calorimeter

• 85K (today)  $\rightarrow$  6M (HL-LHC) readout channels

More collisions More readout channels





CMS HGCAL TDR

### High Luminosity LHC



Need innovation and new techniques to maintain physics reach while staying within throughout requirements!

#### **CMSOfflineComputingResults**

... flat computing budget



# Todays algorithms will not be sustainable in HL-LHC!

## → Utilise modern Machine Learning to become

faster better and do more

#### Fast ML on FPGA

Faster and better decisions



#### Fast ML on ASIC

• Better, smaller data







#### Fast ML on GPU

• Faster and better decisions



#### ML offline

Improve analysis sensitivity



#### Nanosecond inference on specialised hardware





#### **ASIC inference**





#### Low latency

 Strictly limited by collisions occurring every 25 ns



#### Low latency

 Strictly limited by collisions occurring every 25 ns

#### Low resource usage Several algorithms in parallel on single device





#### Low latency • Strictly limited by collisions occurring every 25 ns 10<sup>11</sup> protons • 25 ns • 25 ns

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

#### Power efficient

- On detector: limited to mW
- L1: Cooling key challenge













11

-

ASIC inference






### CMS High Granularity calorimeter

• 6.5 million readout channels, 50 layers





VDE LI (...)

es











• 1.5 µs latency

Thorben Quast | Edinburgh PPE Seminar, 11 June 2021







### All silicon cassette





# Variational Autoencoder





See more here





See more here



## Al for compression



<u>ECON-T, D. Noonan</u>



#### **16 ReLU activated nodes**











### **FPGA** inference











# The level-1 trigger

# To cope with increased data complexity, new at L1:

- Tracking
- Particle Flow
- O(1)M channel HGCal

### Input data

• 2 Tb/s  $\rightarrow$  63 Tb/s

### Latency

• 4  $\mu$ s  $\rightarrow$  12  $\mu$ s













High parallelism = Low latency

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







High parallelism = Low latency

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

**Power efficient** 

• FPGAS ~x10 more power efficient than GPUs (for Phase-2, FPGAs dissipate heat of ~7W/cm2 while processing 5% of total internet traffic!)







High parallelism = Low latency

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

**Power efficient** 

• FPGAS ~x10 more power efficient than GPUs (for Phase-2, FPGAs dissipate heat of ~7W/cm2 while processing 5% of total internet traffic!)

Latency deterministic

• CPU/GPU has processing randomness, FPGAs repeatable and predictable latency







High parallelism = Low latency

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

**Power efficient** 

• FPGAS ~x10 more power efficient than GPUs (for Phase-2, FPGAs dissipate heat of ~7W/cm2 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?



## What are FPGAs?

### Digital signal processors (DSPs):

specialised for multiplication



Memory (BRAM)



# Programming an FPGA



- 3. HLS translates to hardware-description



# Programming an FPGA



Not well served by industry tools!

## FPGA trigger code

```
library ieee;
use ieee.std logic 1164.all;
use ieee.std logic unsigned.all;
use ieee.std logic arith.all;
use work.gtl pkg.all;
entity invariant mass is
    generic (
        upper_limit: real := 15.0;
        lower limit: real := 10.0;
       pt1_width: positive := 12;
        pt2 width: positive := 12;
        cosh cos width: positive := 28;
       INV MASS PRECISION : positive := 1;
        INV MASS COSH COS PRECISION : positive := 3
   );
    port(
        pt1 : in std logic vector(pt1 width-1 downto 0);
        pt2 : in std_logic_vector(pt2_width-1 downto 0);
        cosh_deta : in std_logic_vector(cosh_cos_width-1 downto 0); -- cosh of etal - eta2
        cos_dphi : in std_logic_vector(cosh_cos_width-1 downto 0); -- cos of phi1 - phi2
        inv mass comp : out std logic;
       sim inv mass sq div2 : out std logic vector(pt1 width+pt2 width+cosh cos width-1 downto 0)
    );
end invariant mass;
```

architecture rtl of invariant mass is

constant INV\_MASS\_VECTOR\_WIDTH : positive := pt1\_width+pt2\_width+cosh\_cos\_width; constant INV MASS PRECISION FACTOR : real := real(10\*\*INV MASS PRECISION);.pkg.

```
constant FACTOR 4 VECTOR : std logic vector((INV MASS COSH COS PRECISION+1)*4-1 downto 0) := conv std logic vector(10**(INV MASS COSH COS PRECISION+1),(INV MASS
signal inv_mass_sq_div2 : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);
signal upper limit vector : std logic vector(INV MASS VECTOR WIDTH-1 downto 0);
signal lower limit vector : std logic vector(INV MASS VECTOR WIDTH-1 downto 0);
```

#### begin

```
-- Converting the boundary value for the comparison
```

-- Calculation of invariant mass with the formula: M\*\*2/2 = pt1\*pt2 \* (cosh(etal - eta2) - cos(phi1 - phi2)) inv mass sq div2 <= pt1 \* pt2 \* (cosh deta - cos dphi); sim\_inv\_mass\_sq\_div2 <= inv\_mass\_sq\_div2;</pre>

-- Comparison with boundary values inv mass comp <= '1' when (inv mass sq div2 >= lower limit vector and inv mass sq div2 <= upper limit vector) else '0';

#### end architecture rtl;

#### <u>M. Jeitler</u>

upper limit vector <= conv std logic vector((integer(upper limit\*INV MASS PRECISION FACTOR)), INV MASS VECTOR WIDTH-FACTOR 4 VECTOR'length)\*FACTOR 4 VECTOR; lower limit vector <= conv std logic vector((integer(lower limit\*INV MASS PRECISION FACTOR)), INV MASS VECTOR WIDTH-FACTOR 4 VECTOR'length)\*FACTOR 4 VECTOR;



## FPGA trigger code

```
library ieee;
use ieee.std logic 1164.all;
use ieee.std logic unsigned.all;
use ieee.std logic arith.all;
use work.gtl pkg.all;
entity invariant mass is
    generic (
        upper_limit: real := 15.0;
        lower limit: real := 10.0;
       pt1 width: positive := 12;
       pt2 width: positive := 12;
       cosh cos width: positive := 28;
       INV MASS PRECISION : positive := 1;
        INV MASS COSH COS PRECISION : positive := 3
   );
    port(
       pt1 : in std logic vector(pt1 width-1 downto 0);
       pt2 : in std_logic_vector(pt2_width-1 downto 0);
       cosh_deta : in std_logic_vector(cosh_cos_width-1 downto 0); -- cosh of etal - eta2
        cos_dphi : in std_logic_vector(cosh_cos_width-1 downto 0); -- cos of phi1 - phi2
        inv mass comp : out std logic;
        sim inv mass sq div2 : out std logic vector(pt1 width+pt2 width+cosh cos width-1 downto 0)
end invariant mass;
```

architecture rtl of invariant mass is

constant INV MASS VECTOR WIDTH : positive := pt1 width+pt2 width+cosh cos width; constant INV MASS PRECISION FACTOR : real := real(10\*\*INV MASS PRECISION);.pkg. constant FACTOR 4 VECTOR : std logic vector((INV MASS COSH COS PRECISION+1)\*4-1 downto 0) := conv std logic vector(10\*\*(INV MASS COSH COS PRECISION+1),(INV MASS signal inv\_mass\_sq\_div2 : std\_logic\_vector(INV\_MASS\_VECTOR\_WIDTH-1 downto 0); signal upper limit vector : std logic vector(INV MASS VECTOR WIDTH-1 downto 0);

```
signal lower limit vector : std logic vector(INV MASS_VECTOR_WIDTH-1 downto 0);
```

#### begin

```
-- Converting the boundary value for the comparison
```

-- Calculation of invariant mass with the formula: M\*\*2/2 = pt1\*pt2 \* (cosh(etal - eta2) - cos(phi1 - phi2)) inv\_mass\_sq\_div2 <= pt1 \* pt2 \* (cosh\_deta - cos\_dphi);</pre> sim\_inv\_mass\_sq\_div2 <= inv\_mass\_sq\_div2;</pre>

-- Comparison with boundary values inv mass comp <= '1' when (inv mass sq div2 >= lower limit vector and inv mass sq div2 <= upper limit vector) else '0';

#### end architecture rtl;

# $\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$ Generic HLS implementations for DNN inference

upper limit vector <= conv std logic vector((integer(upper limit\*INV MASS PRECISION FACTOR)), INV MASS VECTOR WIDTH-FACTOR 4 VECTOR'length)\*FACTOR 4 VECTOR; lower limit vector <= conv std logic vector((integer(lower limit\*INV MASS PRECISION FACTOR)), INV MASS VECTOR WIDTH-FACTOR 4 VECTOR 'length)\*FACTOR 4 VECTOR;

#### M. Jeitler





#### TensorFlow / TF Keras / PyTorch / ONNX









pip install hls4ml pip install conifer





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











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

# write and compile the HLS hmodel.compile()

# run bit accurate emulation y tf = model.predict(x) y\_hls = hmodel.predict(x)

# do some validation

# run HLS synthesis hmodel.build()



**Prediction** 

```
hmodel = convert_from_keras_model(model, cfg)
```

```
np.testing.assert_allclose(y_tf, y_hls)
```

#### pynq-z2 floorplan





# Ideally



# Reality

# Efficient NN design for edge compute

### **During training**

- Quantization: do you really need 32-bit FP precision?
- Pruning: removal insignificant synapses
- Knowledge distillation

#### Post-training

Parallelisation (lower latency ↔ more resources)

From 8 GPU server to tiny FPGA!







## Quantization







### Nature Machine Intelligence 3 (2021)

www.nature.com/natmachintell/August 2021 Vol. 3 No. 8

# nature machine intelligence

### Quantized neural networks on the edge





## Google AI

**QKeras** model



<u>hls4ml</u> Fixed-point translation Parallelisation Firmware generation







from tensorflow.keras.layers import Input, Activatio
from qkeras import quantized\_bits
from qkeras import QDense, QActivation
from qkeras import QBatchNormalization

 $\mathbf{x} = \text{Input}((16))$ x = QDense(64,kernel\_quantizer =  $quantized_bits(6,0,alpha=1),$  $bias_quantizer = quantized_bits(6,0,alpha=1))$ 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 = QBatchNormalization()(x) $x = QActivation('quantized_relu(6,0)')(x)$ x = QDense(32,kernel\_quantizer = quantized\_bits (6, 0, alpha=1),  $bias_quantizer = quantized_bits(6,0,alpha=1))e$ 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)



on	from hls4ml import import tensorflow as tf
	<pre># train or load a model model = tf.keras.models.load_model()</pre>
$(\mathbf{x})$	<pre># make a config cfg = config_from_keras_model(model, granularity=`name')</pre>
	<pre># do the conversion hmodel = convert_from_keras_model(model, cfg)</pre>
$(\mathbf{x})$	<pre># write and compile the HLS hmodel.compile()</pre>
$(\mathbf{x})$	<pre># run HLS synthesis hmodel.build()</pre>



## FPGA performance



### Nature Machine Intelligence 3 (2021)



### **OpenReview**.net

Go to NeurIPS 2022 Track Datasets and Benchmarks h...

## Why do tree-based models still outperform deep learning on typical tabular data? PDF

### Leo Grinsztajn, Edouard Oyallon, Gael Varoquaux

06 Jun 2022 (modified: 16 Jan 2023) NeurIPS 2022 Datasets and Benchmarks Readers: 🚱 Everyone Show Bibtex Show Revisions

**Abstract:** While deep learning has enabled tremendous progress on text and image datasets, its superiority on tabular data is not clear. We contribute extensive benchmarks of standard and novel deep learning methods as well as tree-based models such as XGBoost and Random Forests, across a large number of datasets and hyperparameter combinations. We define a standard set of 45 datasets from varied domains with clear characteristics of tabular data and a benchmarking methodology accounting for both fitting models and finding good hyperparameters. Results show that tree-based models remain state-of-



Often the best way to compress is: Just use BDTs!

Conifer = hls4ml for BDTs

If resource/latency constrained, BDT might be solution

- Depending on data, can 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%

**Conifer** 





Often the best way to compress is: Just use BDTs!

Conifer = hls4ml for BDTs

If resource/latency constrained, BDT might be solution

- Depending on data, can be as accurate as a DNN
- Usually significantly faster and more resource efficient



Model	Python AUC	HLS AUC	Latency (clk)	LUT $\%$	FF %
NN	0.985	0.982	8	0.104	0.029
GBDT	0.986	0.981	3	0.140	0.027
### ML for reconstruction



On FPGA: 3.5 µs to cluster energy deposits

### ML for reconstruction



On FPGA: 3.5 µs to cluster energy deposits



### ML for reconstruction



On FPGA: 3.5 µs to cluster energy deposits

• Graph Neural Networks for fast clustering of irregular geometry detectors

Work done on speeding this up from <u>Imperial</u>!

## ML for tracking

In HL-LHC, will need to do track finding at L1

• O(1000) hits, O(100) tracks, 40 MHz rate, ~5 µs latency

Graph Neural Networks for fast charged particle tracking

"Throughput-optimised" for L1 applications, "resource-optimised" for co-processing

	Design	(n <sub>nodes</sub> , n <sub>edges</sub> )	RF	Precision	Latency [cycles]	ll [cycles]	DSP [%]	LUT [%]	FF [%]	BRAN
	Throughput-opt.	(28, 56)	1	ap_fixed<14,7>	59	1	99.9	66.0	11.7	0.
	Resource-opt.	(28, 56)	1	ap_fixed<14,7>	79	28	56.6	17.6	3.9	13.

#### DOI:10.3389/fdata.2022.828666







# Which ML algorithms are we currently exploring to do things completely different

?





### Bias in particle physics



Some variable of interest



Need to exploit the full capabilities of the LHC and be more generic!



## Limitations of current trigger



Trigger threshold

Energy (GeV)

### Level-1 rejects >99% of events! Is there a smarter way to select?



Trigger threshold

Energy (GeV)

### Look at data rather than defining signal hypothesis a priori

Can we "classify" objects/events?



## ML for anomaly detection

### Autoencoders: Learns from data

- Trains unsupervised
- Learns to compress, then reconstruct data
- Often used for financial fraud detection
  - Low rate of anomalous events versus high rate "background"



#### Real data X





## ML for anomaly detection

### Autoencoders: Learns from data

- Trains unsupervised
- Learns to compress, then reconstruct data
- Often used for financial fraud detection
  - Low rate of anomalous events versus high rate "background"



### • Difference $\mathbf{X}$ - $\hat{\mathbf{X}}$ defines "degree of abnormality"







### ML for anomaly detection



Nature Machine Intelligence 4, 154 (2022)

### Select based on degree of abnormality!

















- 300 ms latency
- Thousands of "modules" on many collision events in parallel

ATLAS ALICE

## HLT: More Particle Flow

#### Particle Flow: Best reconstruction at HLT

• Slow, cannot run on all events (currently 17%)









#### CPU node (16/20 cores)



### Events from L1 @ 750 kHz







To handle HL-LHC data rates

• Offload resource-intensive computations to GPU





#### Graph Deep Neural Networks: "fast" approximations of ParticleFlow

#### **Classical Particle Flow**



CMS Simulation Preliminary  $t\bar{t} + PU, \sqrt{s} = 14 \text{ TeV}$ Particle Flow reconstruction



Electror

Muons



#### Graph Neural Network **Graph neural network**



CMS Simulation Preliminary  $t\bar{t} + PU, \sqrt{s} = 14 \text{ TeV}$ Machine-Learned Particle Flow reconstruction







Electrons Muons











### Real-time ML in other experiments





<u>F. Capel et al.</u>





### Real-time ML in other experiments





<u>F. Capel et al.</u>





## Real-time ML in other experiments

### Signals and backgrounds



<u>F. Capel et al.</u>



## hls4ml in other CERN experiments





#### <u>DOI:10.1007/s41781-021-00066-y</u>

### Muon segment finding and reconstruction

- Regression of muon position and angle
- 400 ns budget



<u>R. Teixeira de Lima, R Rojas Caballero et al.</u>



## ... and outside of HEP

### Semantic segmentation for autonomous vehicles



N. Ghielmetti et al.

#### Other examples

- For fusion science phase/mode monitoring
- <u>Crystal structure detection</u>
- <u>Triggering in DUNE</u>
- <u>Accelerator control</u>
- Magnet Quench Detection
- <u>MLPerf tinyML benchmarking</u>
- Food contamination detection
- etc....







# Join the community: fastmachinelearning.org Sign up to the hls-fml group







le ponte



## Backup

#### Data challenge on real-time anomaly detection

- Dataset: Nature Scientific Data (2022) 9:118
- Deadline: November 2022

#### Tutorial: Anomaly detection on FPGA with hls4ml

github:thaarres/quantumUniverse\_pynqZ2

Help us find new physics!

mpp-hep.github.io/ADC2021/

Welcome to the Anomaly Detection Data Challenge 2021!





Model (quantized/pruned)



Convert model to internal representation

Write HLS project targeting specified backend (configurable parallelization/ quantization)

Quantized:









Run emulation

Run synthesis



### Co-processing kernel (Xilinx accelerators/SoCs)

# FPGA custom designs



### **ASICs**















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

# write and compile the HLS hmodel.compile()

# run bit accurate emulation y tf = model.predict(x) y\_hls = hmodel.predict(x)

# do some validation

# run HLS synthesis hmodel.build()



**Prediction** 

```
hmodel = convert_from_keras_model(model, cfg)
```

```
np.testing.assert_allclose(y_tf, y_hls)
```

#### pynq-z2 floorplan



## Compression







### 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!

• Big engineering field in its own right

#### **During training**

- Quantization: do you really need 32-bit FP precision?
- Pruning: removal insignificant synapses

#### Post-training

Parallelisation (lower latency ↔ more resources)

From 8 GPU server to tiny FPGA!







Fixed point post-training quantization

FP 32 arithmetic ~ x3-5 more resources,
x2 higher latency than fixed-point → convert to fixed-point



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

Can we do better?


## Quantization-aware training



# Quantization-aware training

#### Lossless quantization for deep neural networks!



<u>arxiv:2103.13630</u>



#### Nature Machine Intelligence 3 (2021)

www.nature.com/natmachintell/August 2021 Vol. 3 No. 8

# nature machine intelligence

#### Quantized neural networks on the edge





## Google AI

**QKeras** model



<u>hls4ml</u> Fixed-point translation Parallelisation Firmware generation







from tensorflow.keras.layers import Input, Activatio
from qkeras import quantized\_bits
from qkeras import QDense, QActivation
from qkeras import QBatchNormalization

 $\mathbf{x} = \text{Input}((16))$ x = QDense(64,kernel\_quantizer =  $quantized_bits(6,0,alpha=1),$  $bias_quantizer = quantized_bits(6,0,alpha=1))$ 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 = QBatchNormalization()(x) $x = QActivation('quantized_relu(6,0)')(x)$ x = QDense(32,kernel\_quantizer = quantized\_bits (6, 0, alpha=1),  $bias_quantizer = quantized_bits(6,0,alpha=1))e$ 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)



on	from hls4ml import … import tensorflow as tf
	<pre># train or load a model model = tf.keras.models.load_model()</pre>
$(\mathbf{x})$	<pre># make a config cfg = config_from_keras_model(model, granularity=`name')</pre>
	<pre># do the conversion hmodel = convert_from_keras_model(model, cfg)</pre>
$(\mathbf{x})$	<pre># write and compile the HLS hmodel.compile()</pre>
$(\mathbf{x})$	<pre># run HLS synthesis hmodel.build()</pre>



# QKeras quantisers



# FPGA performance



#### Nature Machine Intelligence 3 (2021)



# FPGA performance



#### Nature Machine Intelligence 3 (2021)



# Ideally



# Reality

# Estimating energy and size

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

heterogeneous quantization

# Estimating energy and size

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

heterogeneous quantization

For edge inference, need best possible quantization configuration for

- Highest accuracy  $\uparrow$ ...

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

# Estimating energy and size

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

heterogeneous quantization

For edge inference, need best possible quantization configuration for

- Highest accuracy  $\uparrow$ ...
- ... and lowest resource consumption  $\downarrow$

 $\rightarrow$  hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

QTools: Estimate QKeras model bit and energy consumption, assuming 45 nm Horowitz process

- Relative model size in bits
- Relative energy consumption in Watts

Model A	ccuracy [%	ő]	Р	er-layer	energy
		Dense	ReLU	Dense	ReLU
BF	74.4	1735	53	3240	27
$\mathbf{Q6}$	74.8	794	23	1120	11
					4



**AutoQKeras** 

# AutoQKeras

#### AutoQ Bayesian optimization at work!

• Simultaneously scan quantizers and N filters/neurons (often less/more filters/neurons needed when quantizing)



#### DOI 10.1088/2632-2153/ac0ea1







Model	Accuracy (%)		Precision							$\frac{E}{E_{Q6}}$	Bits Bits <sub>Q6</sub>
		Dense	ReLU	Dense	ReLU	Dense	ReLU	Dense	Softmax		
QE	72.3	(4, 0)	<b>(</b> 4, 2 <b>)</b>	Ternary	(3, 1)	(2,1)	<b>(4, 2</b> )	w: Stoc. bin. b: (8, 3)	(16, 6)	0.27	0.18*

#### Nature Machine Intelligence 3 (2021)

Example with target: **Energy reduction x4** Accuracy degradation max 5%

#### \*w.r.t homogeneously quantized 6 bit model





# QPYTÖRCH ?

#### Brevitas like QKeras, but for PyTorch

- QAT library
- Support most common layers and activation functions

Other quantization techniques:

- HAWQ: Hessian AWare Quantization
- Quantization Aware Pruning (B. Hawks et al.)



#### https://github.com/Xilinx/brevitas

import brevitas.nn as qnn	
qnn.	
🖧 quant_bn (brevitas.nn)	
🕻 😋 QuantCat	brevitas.nn.quant_eltwise
📀 QuantTanh	brevitas.nn.quant_activation
📀 ScaleBias	brevitas.nn.quant_scale_bias
🖧 quant_conv (brevitas.nn)	
👆 🖧 hadamard_classifier (brevit	as.nn)
🚽 🖧 quant_accumulator (brevitas	.nn)
🖞 🖧 quant_activation (brevitas.	nn)
- 🖧 quant_avg_pool (brevitas.nn	
🗛 quant_convtranspose (brevit	as.nn)
🗛 quant_dropout (brevitas.nn)	
🕴 💑 quant_eltwise (brevitas.nn)	
🖧 quant_linear (brevitas.nn)	
💑 quant_max_pool (brevitas.nn	
🗛 quant_scale_bias (brevitas.	nn)
🗛 quant_upsample (brevitas.nn	
BatchNorm1dToQuantScaleBias	brevitas.nn.quant_bn
BatchNorm2dToQuantScaleBias	brevitas.nn.quant_bn
🕒 🕒 ClampQuantAccumulator	brevitas.nn.quant_accumulator









# QPYTÖRCH ?

#### Brevitas like QKeras, but for PyTorch

- QAT library
- Support most common layers and activation functions

Other quantization techniques:

- HAWQ: Hessian AWare Quantization
- Quantization Aware Pruning (B. Hawks et al.)

hls4ml collaborate with Xilinx Research Labs to develop QOONX

- Introducing 'Quant' node to ONNX graph
- Brevitas (PyTorch) and QKeras (Keras) can export QONNX (HAWQ export in progress): then hls4ml and FINN can import QONNX

#### <u>Quantized ONNX (QONNX), J. Mitrevski et. al</u>





Often the best way to compress is: Just use BDTs!

Conifer is hls4ml for Boosted decision trees (scikit-learn, XGBoost)

If resource/latency constrained, 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%

<u>Conifer</u>

# AConifer





Often the best way to compress is: Just use BDTs!

Conifer is hls4ml for Boosted decision trees (scikit-learn, XGBoost)

If resource/latency constrained, 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



Model	Python AUC	HLS AUC	Latency (CIK)	LUI %	$\mathbf{F}\mathbf{F}$ $\gamma_0$	
NN	0.985	0.982	8	0.104	0.029	
GBDT	0.986	0.981	3	0.140	0.027	
						<u> </u>

#### In HL-LHC, will need to do track finding at L1

• O(1000) hits, O(100) tracks, 40 MHz rate, ~5 µs latency

Graph Neural Networks for fast charged particle tracking

• Custom converter for PyTorch Geometric integrated in hls4ml

Throughput-optimized for L1 applications, resource-optimised for co-processing

Design	(n <sub>nodes</sub> , n <sub>edges</sub> )	RF	Precision	Latency [cycles]	ll [cycles]	DSP [%]	LUT [%]	FF [%]	BRAM [%]
Throughput-opt.	(28, 56)	1	ap_fixed<14,7>	59	1	99.9	66.0	11.7	0.7
Throughput-opt.	(28, 56)	8	ap_fixed<14,7>	75	8	21.9	23.8	4.7	0.7
Resource-opt.	(28, 56)	1	ap_fixed<14,7>	79	28	56.6	17.6	3.9	13.1
Resource-opt.	(448, 896)	1	ap_fixed<14,7>	470	174	56.6	25.0	7.4	16.5
Resource-opt.	(448, 896)	8	ap fixed<14,7>	1590	520	5.6	25.0	7.4	16.3

#### DOI:10.3389/fdata.2022.828666













#### More and more dedicated AI processors on the market

• Can we utilise highly specialised ML hardware at CERN?



#### More and more dedicated AI processors on the market

• Can we utilise highly specialised ML hardware at CERN?

#### Xilinx Versal AI processors

- Programmed in C/C++
- Running at 1 GHz
- Example Xilinx ACAP board: 400 AI processors, ~2M logic cells (FPGA), 2k DSPs, Arm CPU, Arm RPU
- Data can move back and forth between AI Engines and FPGA

#### Currently explored for real-time tracking in trigger application

- Interaction Network for pattern recognition (similar to **DeZoort et al**)
- Deployed on Xilinx Versal VC1902 ACAP





#### **CMS***Public* Total CPU HL-LHC (2031/No R&D Improvements) fractions 2022 Estimates



#### **CMS Offline Computing Results**

# **FPGAs as accelerators**

Our DAQ FPGAs are idle  $\sim$ 50% of the time (no collisions)

• Could these be utilised for co-processing?

#### Heterogeneous compute on-site with FPGA co-processors

- E.g LHCb: No hardware trigger, but ~200 FPGA read-out boards receiving data from sub detectors
- Repurpose for data processing when LHC is off?





<u>C. Beteta, I. Bezshyiko, N. Serra</u>



# FPGAs as accelerators

Our DAQ FPGAs are idle ~50% of the time (no collisions)

• Could these be utilised for co-processing?

#### Heterogeneous compute on-site with FPGA co-processors

- E.g LHCb: No hardware trigger, but ~200 FPGA read-out boards receiving data from sub detectors
- Repurpose for data processing when LHC is off?

#### Alternative: FPGA-as-a-Service toolkit for Cloud inference

• Using hls4ml to deploy large models on FPGA, run inference in the cloud

Algorithm	Platform	Number of Devices	Batch Size	Inf./ [Hz]
 FACILE	AWS EC2 F1	1	16,000	36 N
FACILE	Alveo U250	1	16,000	86 N
FACILE	T4 GPU	1	16.000	8 M



#### Data center/ experimental site



FaaST, D. Rankin et. al

# hls4ml in other CERN experiments



# hls4ml in other CERN experiments



# hls4ml in other CERN experiments



### Muon segment finding and reconstruction

- Regression of muon position and angle
- Feasible within 400 ns budget



#### <u>R. Teixeira de Lima, R Rojas Caballero et al.</u>



# ...and outside of HEP



# ...and outside of HEP



# ...and outside of HEP





# ... and outside of HEP

#### Semantic segmentation for autonomous vehicles



N. Ghielmetti et al.

#### Other examples

- For fusion science phase/mode monitoring
- <u>Crystal structure detection</u>
- <u>Triggering in DUNE</u>
- <u>Accelerator control</u>
- Magnet Quench Detection
- <u>MLPerf tinyML benchmarking</u>
- Food contamination detection
- etc....





# Benchmarking

#### Datasets: Common FastML Science Benchmarking datasets

• guide design of edge ML hardware and software for sub-microsecond inference!

#### Algorithms: hls4ml-FINN benchmarked in MLPerf<sup>™</sup>

- how fast systems can process inputs and produce results
- demonstrate efficient and low-latency solutions on FPGAs with hls4ml and FINN

#### Consistently competitive (QKeras+hls4ml, semantic segmentation, MLPerf)

Model	LU	JT	LUT	<b>RAM</b>	FF		BRAM [36 kb]		
Pynq							-Z2		
Available	53 200		17 400		106	400	140		
IC (hls4ml)	28 544	53.7%	3 7 5 6	21.6%	49 215	46.3%	42	30.0%	
IC (FINN)	24 502	46.1%	2 0 8 6	12.0%	34 354	32.3%	100	71.4%	
AD	40 658	76.4%	3 6 5 9	21.0%	51 879	48.8%	14.5	10.4%	
KWS	33 7 32	63.4%	1 0 3 3	5.9%	34 405	32.3%	37	26.4%	

#### https://mlcommons.org/en/inference-tiny-07/



arxiv:2103.05579





#### Detector

- 100% of events
- Latency: 25 ns

CMS Experiment at the LHC, CERN Data recorded; 2010-Nov-14-18:37:44.420271 GMF(19:37:44 CEST) Bun / Event: 151076 11005388

## ∼Pb/s





Level-1 hardware trigger Reject 99.75% of events

• Latency: 0(1) **µs** 



#### Level-1 hardware trigger Reject 99.75% of events

• Latency: 0(1) **µs** 

#### High Level Trigger

- Reject 99.982% of events
- Latency: 0(100) ms





#### Level-1 hardware trigger Reject 99.75% of events

• Latency: O(1) µs

#### High Level Trigger • Reject 99.982% of events Offline reconstruction and storage • Latency: 0(100) ms 19.7 fb<sup>-1</sup> (8 TeV) + 5.1 fb<sup>-1</sup> (7 TeV) CMS G ~1 kHz $H \rightarrow \gamma \gamma$ $m_H = \sqrt{2E_{\gamma_1}E_{\gamma_2}(1-\cos\theta_{\gamma_1\gamma_2})}$ ~Gb/s event Ś S)

(intel) Xeon\* 7500

110


# Do physics with 0.018% of collision events, the rest is discarded! Throughput is main limitation

#### Detector

CMS Experiment at the LHC, CERN

- 100% of events
- Latency: 25 ns

40 MHz

~Tb/s



750 kHz ~Tb/s

### accept/reject



### Level-1 hardware trigger • Reject 99.75% of events

• Latency: O(1) µs

## High Level Trigger

- Reject 99.982% of events
- Latency: 0(100) ms



#### Offline reconstruction and storage



