# Fast inference of deep neural networks in FPGAS



hls 4 ml° Maurizio Pierini











- (<10 µsec) to reconstruct and filter these events

  - deploying invariant mass trigger was a big deal)





• The LHC delivers ~40M events/sec. This gives us little time

• Coarse reconstruction, with worse resolution than offline

Simple algorithms, e.g., sums and products (at some point,

• Extended use of Look Up Tables to parameterise complex functions (with limitation in terms of dimensionality)











 But the process is slow
 • We can use DL solutions as a shortcut: we teach neural networks how to give us the answer we want directly from the raw data

В



Problem: how do we put a NN on an FPGA?









• HLS4ML aims to be this automatic tool

- reads as input models trained on standard DeepLearning libraries

- Could also be used to create co-processing kernels for HLT environments





• comes with implementation of common ingredients (layers, activation functions, etc) • Uses HLS softwares to provide a firmware implementation of a given network







mPP

• Dataflow architecture: each layer is an independent compute unit

• With tunable parallelism and quantization





- Fully on-chip: NN must fit within available FPGA resources (pynq-z2 floorplan shown)











• You have a jet at LHC: spray of hadrons coming from a "shower" initiated by a fundamental particle of some kind (quark, gluon, W/Z/H bosons, top quark)

• You have a set of jet features whose distribution depends on the nature of the initial particle

• You can train a network to start from the values of these quantities and guess the nature of your jet

• To do this you need a sample for which you know the answer

## Example: fast inference

6





https://arxiv.org/abs/1804.06913













## <u>Example: jet taqqinq</u>



















### Lowier Duarta Lbl c/ml

## The full model





### Making the MM smaller: pruning

really contribute to performances

> possible (regularization)







### Making the NN smaller: pruning

Pruning: remove
 parameters that don't
 really contribute to
 performances

force participation procession of the provided as a set of the

 $L_{\lambda}(\vec{w}) = L(\vec{w}) + \lambda ||\vec{w}_1||$ 

Remove the small





### → 70% reduction of weights



20

Eui Re: Co





Xilinx Vivado 2017.2 Clock frequency: 200 MHz **FPGA: Xilinx Kintex Ultrascale** (XCKU115-FLVB2104)



- Big reduction in DSP usage with pruned model!
- ~15 clocks @ 200 MHz = 75 ns inference



12





● Quantisation: reduce101e.1011101010 number of bits used til represent numbers (i.e., reduce used memory)

• models are usually trained at 64 or 32 bits

• this is not necessarily needed in real feg\_relu ftg\_relu

In our case, <sup>\*</sup>
We could to 16 bits w/o loosi precision

 Beyond that, one would have to
 accept some performance loss





### ReuseFactor: how much to parallelize



### related to the Initiation Interval = when new inputs are introduced to the algo.











reuse = 1 <16, 6> bits	BRAM	DSP	FF	LUT
Total	13	954	<b>53k</b>	36k
% Usage	~0%	17%	3%	5%

15













### Parallelisation



### Quantization-aware Training

 Post-training quantisation
 can affect accuracy

• for a given bit allocation, the loss minimum at floating-point precision might not be the minimum anymore

One could specify quantisation while look for the minimum

Maximize accuracy for minimal FPGA resources



• We teamed up with Google to exploit this strategy in a QKeras+h1s4m1 bund1e



Model	Accuracy [%]	Latency [ns]	Latency [clock cycles]	DSP [%]	LUT [%]	
Baseline	74.4	45	9	56.0 (1826)	5.2 (48321)	
Baseline pruned	74.8	70	14	7.7 (526)	1.5 (17577)	
Baseline heterogeneous	73.2	70	14	1.3 (88)	1.3 (15802)	
QKeras 6-bit	74.8	55	11	1.8 (124)	3.4 (39782)	
QKeras Optimized	72.3	55	11	1.0 (66)	0.8 (9149)	



https://arxiv.org/abs/2006.10159











One can push quantisation to extremes

• binary & ternary networks

Multiplications can be replaced by bit manipulations, saving resources

• Can achieve low latencies at small accuracy cost and minimal resource consumption



### Extreme Quantization



A	В	$\overline{A \oplus B}$
0	0	1
0	1	0
1	0	0
1	1	1



https://arxiv.org/abs/2003.06308





### Fast CNN inference on n' - o







Layer name	Layer type	Input shape	Weights	MFLOPs	Energy [nJ]
Conv O	Conv2D	(32, 32, 3)	432	0.778	1,795
Conv 1	Conv2D	(15, 15, 16)	2,304	0.779	1,802
Conv 2	Conv2D	(6, 6, 16)	3,456	0.110	262
Dense O	Dense	(96)	4,032	0.008	26
Dense 1	Dense	(42)	2,688	0.005	17
Output	Dense	(64)	65	0.001	4
Model total			12,858	1.71	3,918











### Fast CNN inference on FPGAs





# **Execution time reduced to 5 µsec to basically no accuracy loss down to 6 bits**







Slightly simplified version of GarNet

• Double message passing (vertices <-> aggregators) can be condensed in FPGA-friendly math

$$g_{v}^{\prime k} = \sum_{a=1}^{S} W_{av} \left( \sum_{j=1}^{F_{in}} \tilde{w}_{ja}^{k} G_{a}^{j} + \tilde{b}_{a}^{k} L_{a} \right) + c$$

$$\tilde{w}_{ja}^{k} = \sum_{i=1}^{F_{\text{LR}}} u_{ia}^{k} w_{j}^{i}, \quad \tilde{b}_{a}^{k} = \sum_{i=1}^{F_{\text{LR}}} u_{ia}^{k} b^{i}, \quad G_{a}^{j} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \text{ and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{\text{max}}} \sum_{v=1}^{V} W_{av} g_{v}^{j}, \quad \text{and } L_{a} = \frac{1}{V_{max}} \sum_{v=1}^{V} W_{av}$$

Non-linearity in distance weighting

$$V_{av} = \exp(-d_{av}^2)$$



## GraphNets on FPG





## GraphNets on FF



Model	$V_{\max}$	$R_{\mathrm{reuse}}$	$\begin{array}{c} \text{Latency} \\ \text{(cycles)} \end{array}$	$\frac{\text{Interval}}{(\text{cycles})}$	DSP $(10^{3})$	LUT $(10^3)$	$FF (10^3)$	BRAM (Mb)	ROC AUC	Response RMS
Continuous	128	32	155	55	3.1 [56%]	57 [9%]	39  [2.9%]	1.8  [2.3%]	0.98	0.23
Quantized	128	32	148	50	1.6 $[29%]$	$70 \ [11\%]$	$41 \ [3.1\%]$	$1.9 \ [2.4\%]$	0.98	0.24
Quantized	64	16	99	34	1.6 [29%]	$63 \ [9\%]$	$38 \ [2.9\%]$	1.8  [2.3%]	0.96	0.24
Quantized	32	8	75	26	1.4  [25%]	52 $[8%]$	33~[2.5%]	1.8  [2.3%]	0.86	0.37
Quantized	16	4	63	22	1.5~ig[27%ig]	57 $[9%]$	37[2.8%]	$1.8\;ar{2.3\%}$	0.64	0.36





Regression







### Anomaly Detection on FPGA

- Autoencoders are compression-decompression algorithms that learn to describe a given dataset in terms of points in a lower-dimension latent space
- Can be used for anomaly detection: compression less effective when anomalous events are given
- In principle, same compression techniques apply
- In practice, special considerations apply
  - Figure of merit for anomaly detection might be different than loss -> QAT might run in trouble
  - Variational autoencoders require sampling from random numbers (possible, but with some resource utilisation)











### Anomaly Detection on FPGF







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















Sometimes the model is too big to be deployed

(the student model)

• Every problem is turned into a regression

• Easier to compress than unsupervised models





- Still, a network is a function that can be approximated by ... a neural network

  - Does not require resource-consuming output activation for classifications











• Sometimes the model is too big to be deployed

• Still, a network is a function that can be approximated by ... a neural network (the student model)

• Every problem is turned into a regression

• Easier to compress than unsupervised models

 Does not require
 resource-consuming output activation for classifications







### hls4ml: An Open-Source Codesign Workflow to Empower **Scientific Low-Power Machine Learning Devices**

Farah Fahim\* Benjamin Hawks Christian Herwig James Hirschauer Sergo Jindariani Nhan Tran<sup>\*</sup> Fermilab Batavia, IL, USA

Manuel Blanco Valentin Josiah Hester Yingyi Luo John Mamish Seda Orgrenci-Memik Northwestern University Evanston, IL, USA

Thea Aarrestad Hamza Javed Vladimir Loncar Maurizio Pierini Adrian Alan Pol Sioni Summers European Organization for Nuclear Research (CERN) Geneva, Switzerland

Scott Hauck Shih-Chieh Hsu University of Washington Seattle, WA, USA

> Duc Hoang **Rhodes** College Memphis, TN, USA



### https://github.com/hls-fpga-machine-learning/hls4ml

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

https://github.com/thesps/conifer



Luca P. Carloni Giuseppe Di Guglielmo Columbia University New York, NY, USA

Philip Harris Jeffrey Krupa Dylan Rankin MIT Cambridge, MA, USA

Jennifer Ngadiuba Caltech Pasadena, CA, USA

Mia Liu Purdue University

West Lafayette, IN, USA

Edward Kreinar HawkEye360 Herndon, VA, USA

Zhenbin Wu University of Illinois at Chicago Chicago, IL, USA



Javier Duarte UC San Diego La Jolla, CA, USA jduarte@ucsd.edu