

hls4ml: deploying deep learning on FPGAs for L1 trigger and Data Acquisition





Javier Duarte, Sergo Jindariani, Ben Kreis, Ryan Rivera, Nhan Tran (Fermilab) Jennifer Ngadiuba, Maurizio Pierini, Sioni Summers, **Vladimir Loncar** (CERN)

Edward Kreinar (Hawkeye 360)

Phil Harris, Song Han, Dylan Rankin (MIT)

Zhenbin Wu (University of Illinois at Chicago)

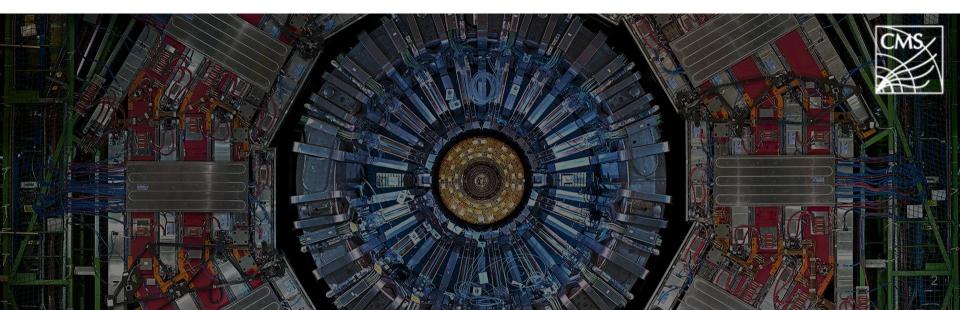
Giuseppe di Guglielmo (Columbia University)

Challenges in LHC

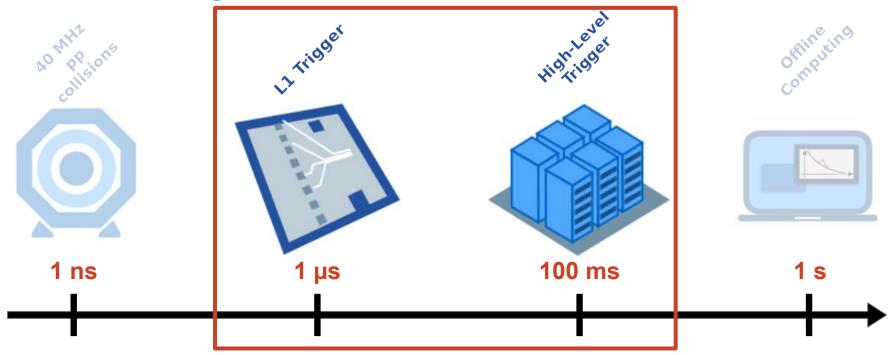
At the LHC proton beams collide at a frequency of 40 MHz

Extreme data rates of O(100 TB/s)

"Triggering" - Filter events to reduce data rates to manageable levels



The LHC big data problem



Deploy ML algorithms very early

Challenge: strict latency constraints!

Field-Programmable Gate Array

Reprogrammable integrated circuits

Configurable logic blocks and embedded components

- Flip-Flops (registers)
- LUTs (logic)
- DSPs (arithmetic)
- Block RAMs (memory)

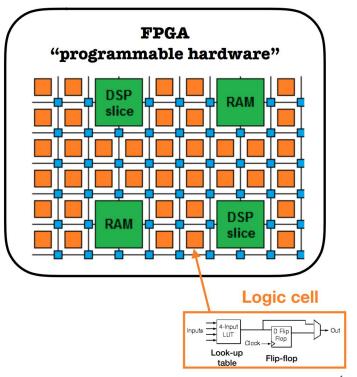
Massively parallel

Low power

Traditionally programmed with VHDL and Verilog

High-Level Synthesis (HLS) tools

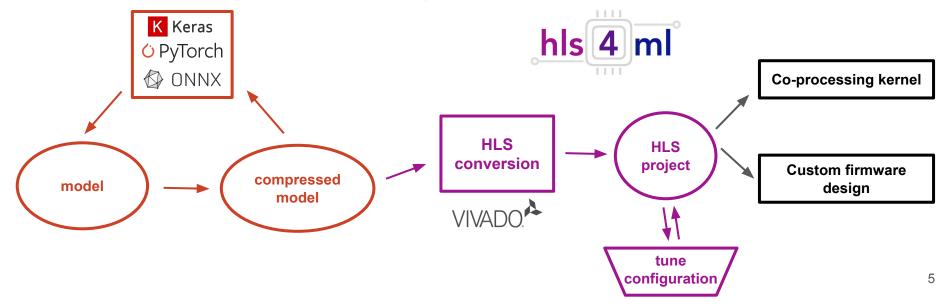
Use C, C++, System C



high level synthesis for machine learning

User-friendly tool to automatically build and optimize DL models for FPGAs:

- Reads as input models trained with standard DL libraries
- Uses Xilinx HLS software
- Comes with implementation of common ingredients (layers, activation functions, binary NN ...)





The main idea: Store the full architecture and weights on chip

- Much faster access times
- For longer latency applications, weights storage in on-chip block memory is possible
- No loading weights from external source (e.g. DDR, PCle)

Limitations:

- Constraints on model size
- Not reconfigurable without reprogramming device

Solution: User controllable trade-off between resource usage and latency/throughput

Tuned via "reuse factor"



hls 4 ml : exploiting FPGA hardware

Parallelization: Use reuse factor to tune the inference latency versus utilization of FPGA resources

Can now be specified per-layer

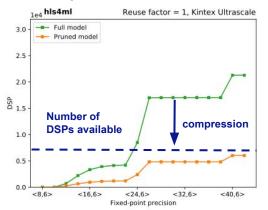


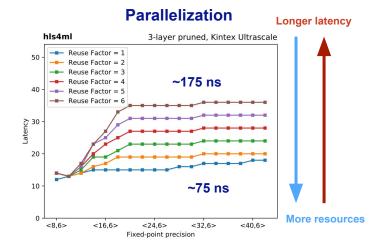
Quantization: Reduce precision of the calculations

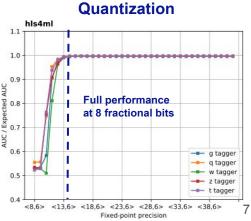
Compression: Drop unnecessary weights (zero or close to zero) to reduce

reuse = 1 use 4 multipliers 1 time each the number of DSPs used

70% compression ~ 70% fewer DSPs







reuse = 4use 1 multiplier 4 times

reuse = 2 use 2 multipliers 2 times each





hls 4 ml : compression by binarization/ternarization

Replace floating/fixed-point with 1/2-bit arithmetics

- Binary: 1-bit (<u>arXiv:1602.02830</u>)
- Ternary: 2-bits (<u>arXiv:1605.04711</u>)

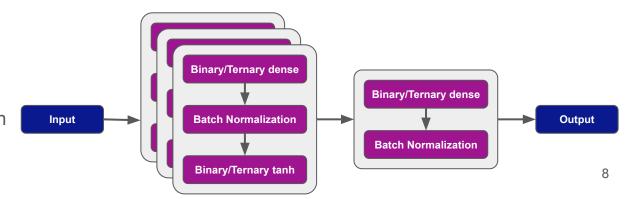
Multiplications (a * w) as bit-flip operations:

- Binary: res = w == 0 ? -d : d;
- Ternary: res = w == 0 ? 0 : w == -1 ? -d : d;

Real-valued Networks sign(□)

Binary/ternary architecture:

- Binary/Ternary Dense
- **Batch Normalization**
- Binary/Ternary tanh activation



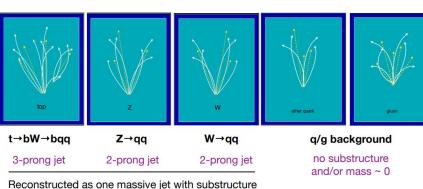


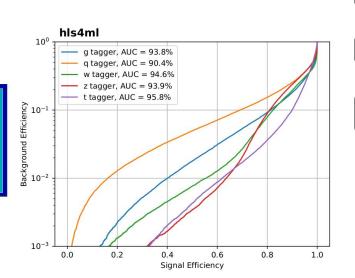
hls 4 ml : Jet tagging benchmark model

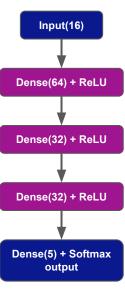
Multi-classification task:

- Discrimination between highly energetic (boosted) q, q, W, Z, t initiated jets
- 16 inputs, 5 outputs

Average accuracy ~ 0.75









hls 4 ml : Jet tagging benchmark model

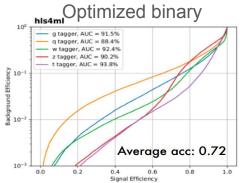
Run hyper-parameter bayesian optimization:

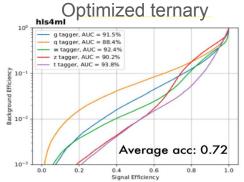
Number of neurons/layers, batch size, learning rate

Recover performance with larger models

- Binary: **16x448x224x224x5** (7x more neurons)
- Ternary: **16x128x64x64x64x5** (2x more neurons + one more layer)

Model	Accuracy	Latency	DSP	BRAM	FF	LUT
Base model	0.75	0.06 µs	60%	0%	1%	7%
Optimized Binary	0.72	0.21 µs	0%	0%	7%	15%
Optimized Ternary	0.72	0.11 <i>µ</i> s	0%	0%	1%	6%







Dense networks trained with the MNIST dataset

- 784 inputs (28x28 grayscale image), 10 outputs (digits)

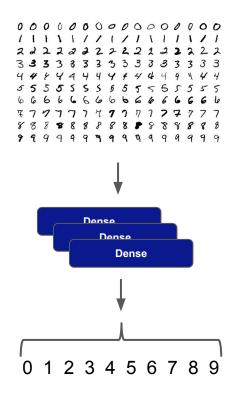
Base model:

- 3 hidden layers with 128 neurons and ReLU activation

Binary/Ternary model:

- 3 hidden layers with batch normalization and binary/ternary tanh Xilinx VU9P FPGA at 200 MHz, reuse factor 128

Model	Accuracy	Latency	DSP	BRAM	FF	LUT
Dense model	0.97	2.6 µs	21%	45%	12%	33%
Binary dense model	0.93	2.6 µs	0%	33%	7%	39%
Ternary dense model	0.95	2.6 µs	0%	33%	7%	40%





Supported architectures:

- DNN
 - Support for very large layers NEW
 - Zero-suppressed weights
- Binary and Ternary DNN NEW
 - 1- or 2-bit precision with limited loss of performance
 - Computation without using DSPs, only LUTs
- Convolutional NNs
 - 1D and 2D with pooling
 - Currently limited to very small layers, working on support for larger layers

Other:

- Batch normalization
- Merge layers (concatenation, addition, subtraction etc)
- Numerous activation functions

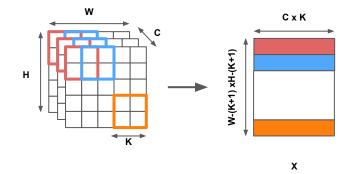


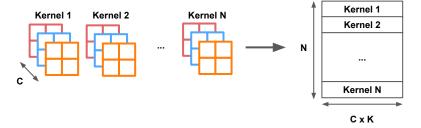
Convolutional layers

Support for "large" convolutional layers



- Express convolution as matrix multiplication
- im2col algorithm
- Reuse "large" matrix multiplication algorithm from MLP
- Quantized (binary and ternary) weights







Convolutional layers

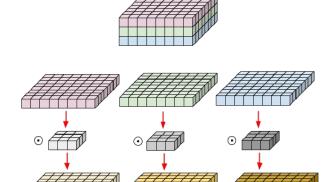
Depthwise separable convolution (arXiv:1610.02357)

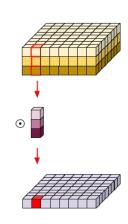
- First step: depthwise convolution
- Second step: pointwise convolution
- For 3x3 kernels this can yield 8-9 times less multiplications

LeanConvNet (arXiv:1904.06952)

- Depth-wise (block diagonal) operator operating on each channel separately and 1×1 convolution
- 5-point convolution kernel





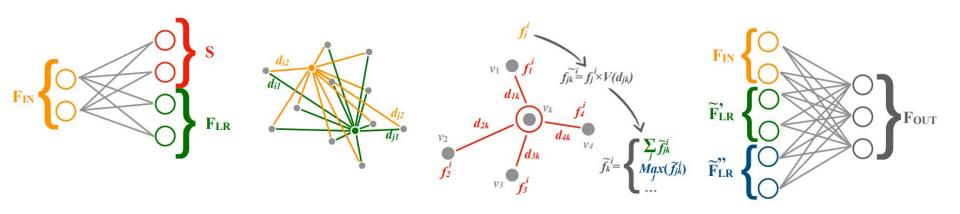




Graph networks (GarNet)



- Distance-weighted GNN capable of learning irregular patterns of sparse data (arXiv:1902.07987)
- Suitable for irregular particle-detector geometries
- Early stage of HLS implementation





Multi-FPGA inference

H1 2020

- Main idea: place layers onto multiple FPGAs and pipeline the execution

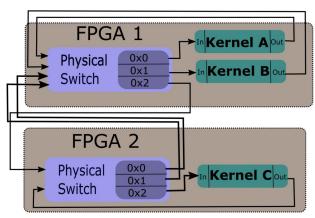
Leverage Galapagos framework (https://github.com/tarafdar/galapagos)

- "...a framework for creating network FPGA clusters in a heterogeneous cloud data center."

- Given a description of how a group of FPGA kernels are to be connected, creates a ready-to-use

network device

- Possible to use MPI programming model



Credit: Naif Tarafdar, Phil Harris



Recurrent Neural Networks (RNNs)

Q4 2019

Boosted decision trees

Q4 2019

Autoencoders

H2 2020

HLS implementations beyond Xilinx/Vivado

H1 2020

- Quartus HLS Compiler for Intel/Altera FPGAs
- Mentor Catapult HLS

Inference engine for CPUs based on hls4ml

H1 2020

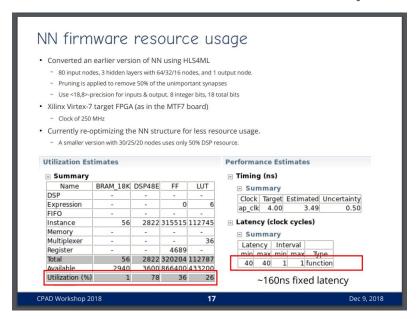
Targeting integration with CMSSW

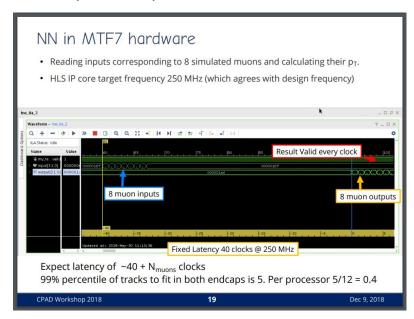
Many more...



CMS designing DL-based triggers for Run III, using hls4ml for deployment

- Reduce muon rate by factor 4 (link)
- Run inference in 160ns on currently used boards (Virtex 7)





Conclusions

hls4ml - software package for translation of trained neural networks into synthesizable FPGA firmware

- Tunable resource usage latency/throughput
- Fast inference times, O(1µs) latency

More information:

- Website: https://hls-fpga-machine-learning.github.io/hls4ml/
- Paper: https://arxiv.org/abs/1804.06913
- Code: https://github.com/hls-fpga-machine-learning/hls4ml



Bonus



Install:

pip install hls4ml SOON (for now: git clone ... && cd hls4ml && pip install .)

Translate to HLS:

hls4ml convert -c my_model.yml

Run synthesys etc.:

hls4ml build -p my project dir -a

Get help:

hls4ml <command> -h

...or visit: https://fastmachinelearning.org/hls4ml/

...or contact us at hls4ml.help@gmail.com

Degree of parallelism

K Keras

♠ ONNX OnnxModel: models/my model.onnx InputData: data/my input features.dat OutputPredictions: data/my predictions.dat OutputDir: my project dir ProjectName: myproject XilinxPart: xcku115-flvb2104-2-i ClockPeriod: 5 IOType: io parallel HLSConfia: Model: Precision: ap fixed<16,6> ReuseFactor: 2 Strategy: Resource Default precision

Support for large models

(weights, biases...)



hls 4 ml : Advanced configuration example

```
KerasJson: models/my model.json
      KerasH5: models/my model weights.h5
      OutputDir: my project dir
      ProjectName: myproject
      XilinxPart: xcku115-flvb2104-2-i
      ClockPeriod: 5
      IOType: io parallel
      HLSConfia:
        Model:
          Precision: ap fixed<16,6>
          ReuseFactor: 8
          Strategy: Resource
        LayerName:
                                       Applies to the
          fc1 relu:
                                        whole model
            Precision:
              weight: ap fixed<18,6>
Specific to this
              bias: ap fixed<16,8>
layer by name
               result: ap fixed<18,8>
            ReuseFactor: 4
```

```
Applies to all other
LayerType:
                           Dense layers
  Dense:
    Precision:
      default: ap fixed<18,8>
      weight: ap fixed<14,6>
    ReuseFactor: 2
  Activation:
    Precision: ap fixed<12,8>
```

Applies to all Activation layers



Boosted decision trees

Q4 2019

- BDTs have been popular for a long time in HEP reconstruction and analysis
- Suitable for highly parallel implementation in FPGAs
- Implementation in hls4ml optimised for low latency
- No 'if/else' statement in FPGAs → evaluate all options and select the right outcome
 - Compare all features against thresholds, chain together outcomes to make the 'tree'

Test for model with 16 inputs, 5 classes, 100 trees, depth 3 on VU9P FPGA:

- 4% LUTs, 1% FFs (0 DSPs, 0 BRAMs)
- 25 ns latency with II=1

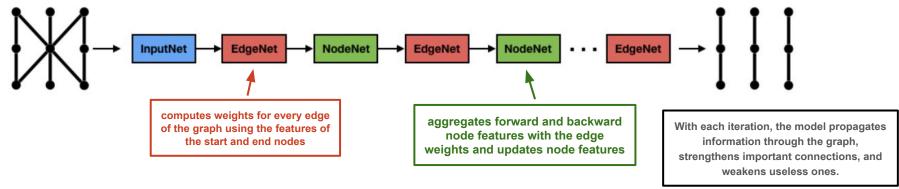
Credit: Sioni Paris Summers



Graph networks

H1 2020

Natural solution for reconstructing the trajectories of charged particles



Preliminary implementation:

- Implemented as an HLS project, not supported in conversion tools
- Successfully tested a small example with 4 tracks, 4 layers
- Major effort required to scale up to larger graphs



Recurrent neural networks

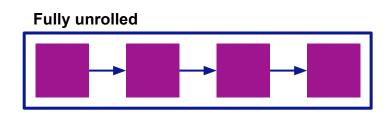
Q4 2019

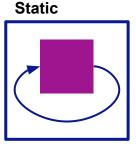
Simple RNN, LSTM, GRU

Two implementations:

- Fully unrolled:
 - Latency optimized with II=1
 - Large resource usage
- **Static:** same resources used for weights and multiplications
 - N (N=latency of layer) copies can go through at the same time
 - Latency is larger and II limited to clock time for each layer

Supports small networks → scale it up using "large" matrix multiplication algorithm







Training on FPGAs

H2 2020

- Build on top of multi-FPGA idea

Use synthetic gradients (SG) to remove the update lock

Individual layers to learn in isolation

Train SGs by another NN

- Each SG generator is only trained using the SGs generated from the next layer
- Only the last layer trains on the data

