# Generalized Machine Learning Quantization Implementation for High Level Synthesis Targeting FPGAs

Matthew Trahms - 15th March, 2022

# Outline

- The Large Hadron Collider
- Particle Tracking
- hls4ml
- FINN and Brevitas
- Particle Tracking GNN Quantization Aware Training
- FINN Collaboration and QONNX
- QONNX Ingest into hls4ml
- Future Work
- Acknowledgements

# The Large Hadron Collider (LHC)

- Large particle accelerator in Europe
  - Accelerate particles near the speed of light
  - Collisions split atoms into subatomic particles
  - Sensors track particle interactions
- Generates large quantities of data
  - 1 petabyte of data / second while operating
    - Data management challenge
    - Set to increase in the future with the high luminosity upgrade
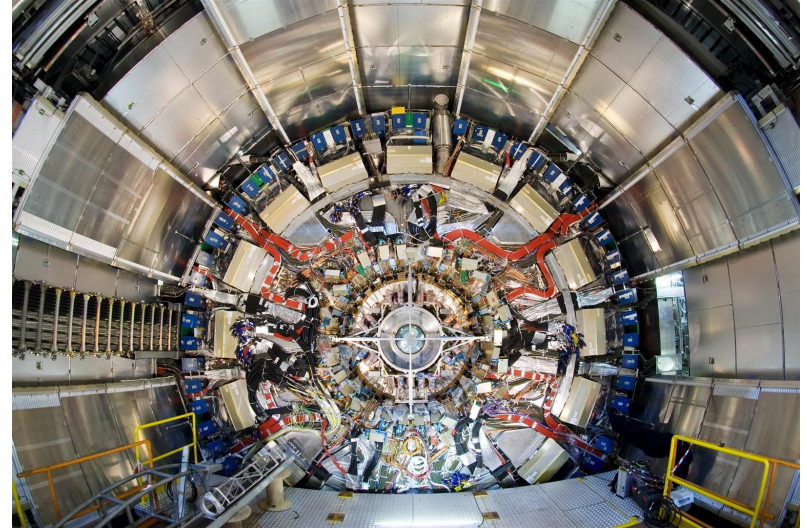  - Multiple different experiments
    - ATLAS
    - CMS



Fig 1. The ATLAS Detector
(atlas.cern)

# Particle Tracking

- Need to track individual particles from collisions in order to make observations
  - Large number of collisions happening in a small space
    - Need to accurately separate particle paths
  - Non-machine learning algorithm algorithm scales poorly with increasing number of particles
    - Higher luminosity = more particles



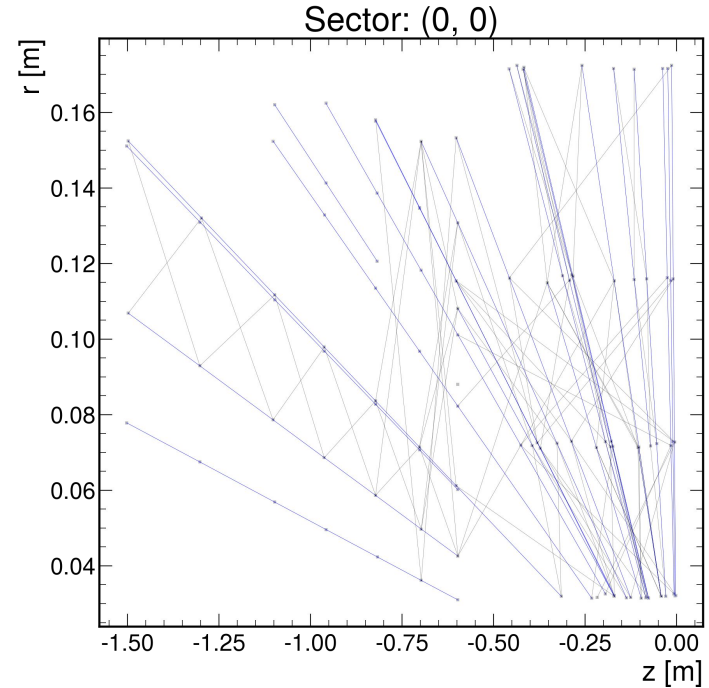Fig. 2: Possible particle tracks. Blue edges are correct edges while black edges were potential particle paths based on detections. (A. Elabd et al)

# Particle Tracking Graph Neural Network (GNN)

- Possible particle tracks represented as directional acyclic graphs
  - Nodes = locations where particles were detected
  - Edges = possible particle paths, determined by distance
- Graph Neural Nets (GNNs) transform graphs
  - Interaction networks determine true edges, or the actual path of the particle
  - Able to match performance of non-machine learning algorithm and scales with luminosity
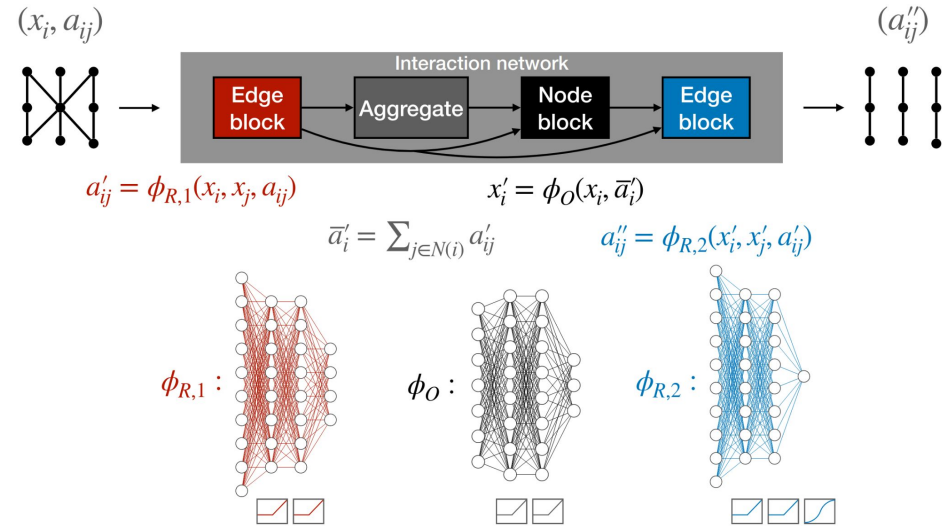
$(x_i, a_{ij})$             $(a_{ij}'')$

Interaction network

Edge block   Aggregate   Node block   Edge block

$a_{ij}' = \phi_{R,1}(x_i, x_j, a_{ij})$      $x_i' = \phi_O(x_i, \bar{a}_i')$

$\bar{a}_i' = \sum_{j \in N(i)} a_{ij}'$     $a_{ij}'' = \phi_{R,2}(x_i', x_j', a_{ij}')$

$\phi_{R,1}:$     $\phi_O:$     $\phi_{R,2}:$

Fig. 3: The structure of an Interaction Network Notice how the graph is transformed at the output
(A. Elabd et al)

# hls4ml

- Machine learning algorithms work as alternatives for high energy physics applications
  - Need high throughput
  - Field Programmable Gate Array (FPGA): Reprogrammable logic that can be used to implement digital algorithms
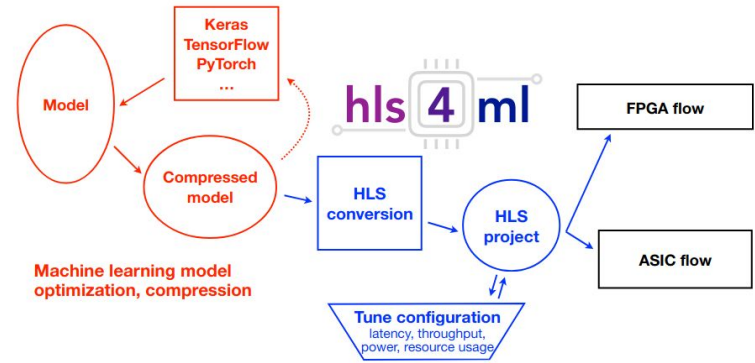    - Generally able to achieve higher throughput than acceleration on CPUs



Fig. 4: The hls4ml flow to generate a hardware implementation of a machine learning model (Luca Carloni et al)

# FINN and Brevitas

- FINN: hls4ml alternative developed by Xilinx Research
    - Targets extremely low bit width deployment
        - Low latency, high throughput
    - Brevitas: Quantization Aware Training (QAT) library developed for FINN
        - QAT allows for high accuracy at low bit widths compared to post training quantization (PTQ)
        - Based on Pytorch
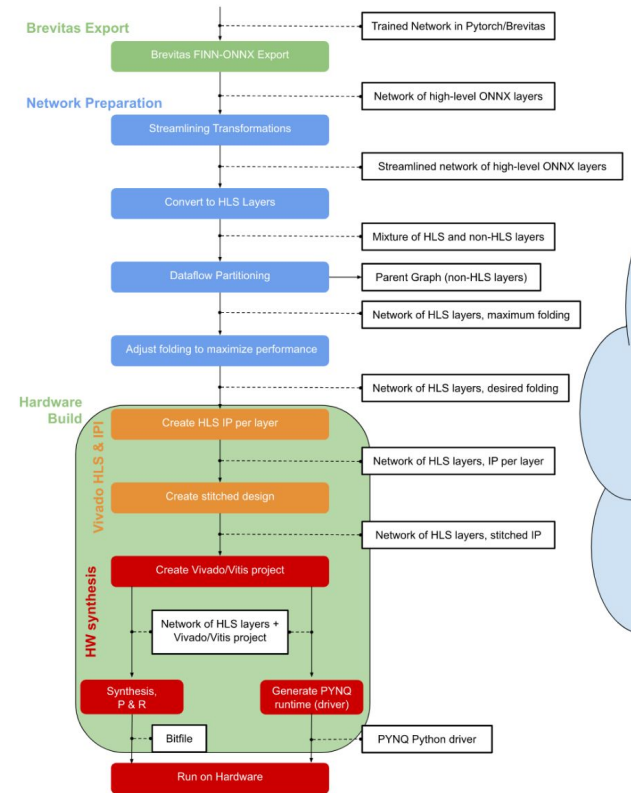    - Uses ONNX graphs with custom nodes for internal representation



Fig. 5: The FINN Flow, from Brevitas export to deployment
(xilinx.github.io)

# Tracking GNN Quantization Aware Training

- ● FPGA acceleration of particle tracking GNN
  - ○ Need for high throughput
  - ○ Network originally implemented in Pytorch
    - ■ Only option for hls4ml is PTQ
    - ■ Loss in accuracy
  - ○ Re-implemented and trained network in Brevitas
    - ■ Layer by layer replacement
    - ■ Retrained on same dataset
  - ○ Achieved near equivalent performance
    - ■ AUC: Area under ROC curve
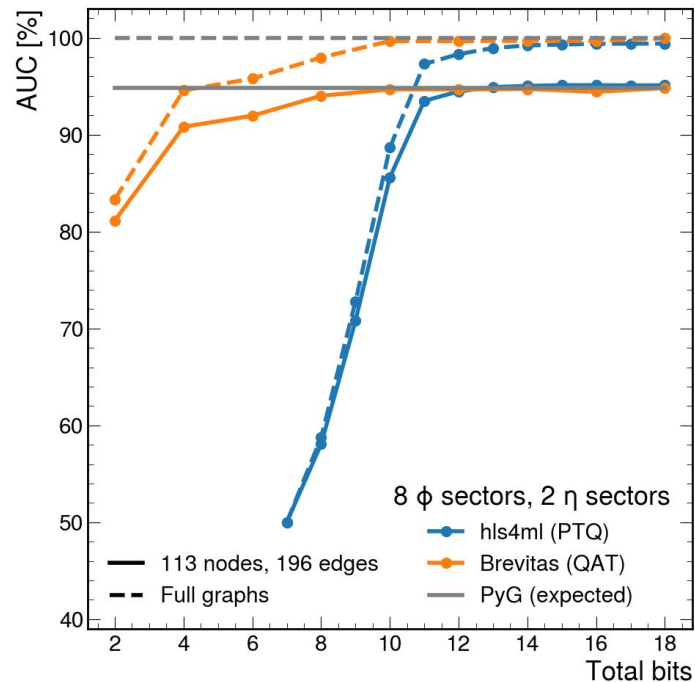      - ● Compares true positive rate and false positive rate



Fig. 6: Tracking GNN QAT Results
(A. Elabd et al)

# FINN Collaboration and QONNX

- FINN and hls4ml accomplish similar tasks
  - Cross organizational collaboration - develop a shared model format that can be used by hls4ml and FINN
    - Generalized version of FINN ONNX
    - Extends the ONNX framework to add quant nodes
      - Represents either weight or input quantization
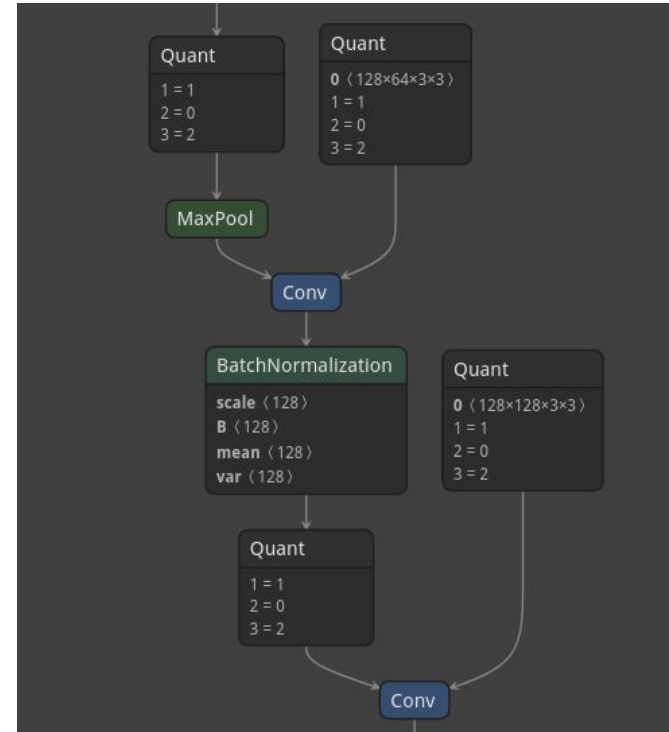  - Interoperability means that users can choose the solution that works better for their purposes



Fig. 7: Visual representation of a QONNX network (netron.app)

# QONNX Ingest Into hls4ml

- ## Need to convert QONNX to HLSModel to synthesize in hls4ml
  - QONNX quant nodes specify quantization
  - HLSModel layers have quantization attributes built into layers
  - Set of transformations
    - Ingest complete structure into HLSModel
    - Incorporate quantizations from Quant nodes into layers
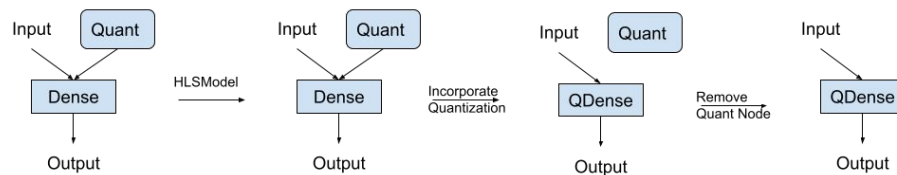    - Remove Quant nodes



Fig. 8: Process to incorporate QONNX quantizations into HLSModel

# Future Work

- QONNX ingest into hls4ml needs further work
  - Bugs with convolutional models and different model architectures
  - Needs to be pulled into the master branch of hls4ml
    - Needs to be converted to the new workflow
  - Need to test latency synthesis implementation
    - Currently testing with resource
  - Need to test with more/different model architectures
- Take a QAT particle tracking GNN through synthesis
  - Only a numerical study currently, need to validate actual performance on an FPGA
- More collaboration between FINN and hls4ml
  - Streamlining in hls4ml
  - MLIR

# Acknowledgements

# Questions?