Track reconstruction for the ATLAS Phase-II Event Filter using GNNs on FPGAs

Sebastian Dittmeier (Physikalisches Institut - Heidelberg University)

on behalf of the ATLAS TDAQ collaboration CHEP 2023 May 08 – 12, 2023@ Norfolk, Virginia, USA



Event Filter Tracking for ATLAS @ HL-LHC

- Event Filter forms trigger decision after L0 accept
- Track reconstruction computationally intensive
 - Track finding biggest challenge
- Ongoing R&D on acceleration with new algorithms and GPUs or FPGAs

Table 2.4: The CPU required in HS06 × sec to reconstruct $t\bar{t}$ MC events with $\langle \mu \rangle = 140$ and 200 in the ITk, using full-scan and regional (5% η - ϕ coverage) track selection cuts.

$\langle \mu \rangle$	Tracking	Byte Stream	Cluster	Space	Si Track	Total	
		Decoding	Finding	Points	Finding	ITk	
140	full-scan	2.2	6.1	1.0	13.4	22.7	
	regional	0.33	0.90	0.15	1.11	2.49	
200	full-scan	3.2	8.1	1.2	23.2	35.7	
	regional	0.48	1.23	0.18	1.92	3.81	[1]



Finding Track Candidates with Graph Neural Networks



- Exa.TrkX / GNN4ITK pipeline with different methods per step (computational constraints)
- More talks at CHEP on <u>Exa.TrkX</u> and developments in GNN4ITK (<u>1</u>) and (<u>2</u>)
- FPGAs can potentially, compared to CPUs / GPUs:
 - Speed up inference considerably
 - **Reduce power** consumption

Challenge: retaining performance with smaller models due to resource constraints on the FPGA

Considerations for GNNs on FPGAs for EF Tracking

Model architecture

 ↔ Translation using available tools (e.g. <u>HLS4ML</u>, <u>FINN</u>), dedicated hardware implementations

Hardware implementations under development



Model architecture

↔ Translation using available tools (e.g. <u>HLS4ML</u>, <u>FINN</u>),
 <u>dedicated hardware implementations</u>

VHDL implementations under development:

- Module Map for graph construction
- Walkthrough for graph segmentation

Sebastian Dittmeier - Heidelberg University

Hardware implementations under development

Model architecture

 ↔ Translation using available tools (e.g. <u>HLS4ML</u>, <u>FINN</u>), dedicated hardware implementations



VHDL implementations under development:

- Module Map for graph construction
- Walkthrough for graph segmentation

Resource estimate from synthesis (Stratix 10 GX versus number of graph edges

Considerations for GNNs on FPGAs for EF Tracking

- Model architecture
- ↔ Translation using available tools (e.g. <u>HLS4ML</u>, <u>FINN</u>), dedicated hardware implementations
- Model optimization
- ↔ Smaller models, but also:
 Quantization + Pruning to meet resource constraints

Model optimization studies

Graph construction using Metric Learning



- Goal: Limit amount of false edges while keeping true edges
- Embed hit points into latent space using a **multi-layer perceptron (MLP)**
- Connect hits that are close to each other in latent space
- Trade-off efficiency versus purity by cut on radius



Sebastian Dittmeier - Heidelberg University

Quantization and Pruning studies with TrackML dataset

- Simplified detector geometry, adapted from early ATLAS ITk designs
- Pile-up 200 conditions like @ HL-LHC
- The following studies make use of a small pre-processed dataset of particles with $p_T > 1$ GeV
- <u>Repository on Github</u>: fork of Exa.TrkX pipeline



dimensions



Quantization Aware Training

Using **Brevitas** (Xilinx), quantizing weights and activations to 8 bits



Hit input data using signed fixed point representation <2,10>

Model size evaluation in Bit Operations per cluster: BOPs $\propto (1-f_p)b_a b_w$ f_p : pruning fraction, $b_a(b_w)$: bit width of activations (weights)





Sebastian Dittmeier - Heidelberg University







Further Considerations

- Model architecture
- Model optimization
- Event / graph size

Throughput

Deployment and integration

- ↔ Translation using available tools (e.g. <u>HLS4ML</u>, <u>FINN</u>),
 dedicated hardware implementations
- ↔ Smaller models, but also:
 Quantization + Pruning to meet resource constraints
- ← Segmentation of graphs , restrict on sub-detectors, or external memory / HBM to meet memory constraints
- ↔ Online track reconstruction for the trigger
 EF system rates: 1 MHz (regional), 150 kHz (global)
- ↔ Accelerator cards and tools (Xilinx, Intel), heterogeneous software platforms (e.g. oneAPI)
 Sebastian Dittmeier - Heidelberg University



Conclusions and Outlook

- Ongoing studies to implement GNNs on FPGAs for the ATLAS Event Filter @ HL-LHC
- Graph construction methods under investigation
 - Metric learning MLP quantization and pruning studies look promising
 - Pending next steps:
 - Translation for FPGA and complete algorithmic implementation
 - Application to ATLAS ITk simulation
 - Planned comparison with VHDL Module Map implementation
- More ongoing studies:
 - Performance of smaller interaction networks for ATLAS ITk
 - Graph segmentation into detector regions
 - MLP and Interaction Network to FPGA translation



Conclusions and Outlook

- Ongoing studies to implement GNNs on FPGAs for the ATLAS Event Filter @ HL-LHC
- Graph construction methods under investigation
 - Metric learning MLP quantization and pruning studies look promising Ο
 - Pending next steps: Ο
 - Translation for FPGA and complete algorithmic implementation
 - Application to ATLAS ITk simulation
 - Planned comparison with VHDL Module Map implementation Ο
- More ongoing studies:
- Thank you for your attention! Performance of smaller interaction networks for ATLAS ITk Ο
 - Graph segmentation into detector regions Ο
 - MLP and Interaction Network to FPGA translation Ο



Sebastian Dittmeier - Heidelberg University

ATLAS upgrade for the High-Luminosity LHC

- Peak instantaneous luminosity of 7.5 × 10^{34} cm⁻²s⁻¹ leads to average pile-up $\langle \mu \rangle \approx 200$
- Upgrades include:
 - New all-silicon tracking detector: Inner Tracker (ITk)
 - Upgrade of Trigger and Data
 Acquisition System (TDAQ) ^[3]
 - And more (HGTD, muon system, calorimeter)



The approach: Graph Neural Networks on FPGAs

- Finding track candidates with **Graph Neural Networks**
- **Demonstrated good performance** with ATLAS ITk simulation ^[5](GNN4ITK)
- Reference to GNN4ITK talks
- FPGAs can potentially, compared to CPUs / GPUs:
 - Speed up inference considerably
 - **Reduce power** consumption



Collision events as graphs

- Graphs consist of a set of nodes and edges
- Take an event in a detector:
 - Represent each **hit** as a **node**
 - Connect nodes by edges
 - Edges suggest two hits belong to the same track
- 3 levels of information:
 - Node-level (position, energy deposited, etc.)
 - Edge-level (belongs to track or not, geometric info, etc.)
 - Graph-level (event, detector region, etc.)
- Similarly, predictions can be made with a GNN on node-level, edge level, and graph-level
 - In this case, make edge-level predictions to construct tracks (i.e. are edges true or false?)



Graph construction methods

Goal: Limit amount of false edges while keeping true edges

- Metric learning
 - Embed hit points into latent space using a **multi-layer perceptron (MLP)**
 - Connect hits that are close to each other in latent space

• Module map

- For each module pair, find max and min of geometric values
- Apply geometric cuts for each pair, and construct map of possible connections/edges
- Store in permutation invariant matrix



Edge classification (GNN)

Goal: classify edges as "True" or "False" - could they belong to a track?

GNNs rely on **Message Passing**:

- Node vectors (properties of node) are updated through an update function (MLP)
- 2. Information is passed along edges to neighbouring nodes
- 3. All messages are aggregated (summed) at each node
- 4. After one or more message passing steps, use classifier to make edge-level predictions



Edge classification (GNN)

Use type of GNN called "Interaction Network"

- Interaction Network adds an extra step to the message passing algorithm
- "Edge network" updates edge features, allowing two nodes to form unique relationships
- Improves quality of edge-level predictions

Finally, edge scores are assigned to all edges in an event

Use threshold value (e.g. ~0.5) to discard false edges



Track reconstruction methods

Method 1: Walkthrough

- Identify starting node
- Traverse edges with high scores
- Longest path found \rightarrow track candidate



Method 2: Connected components

- For edges above threshold score, identify connected paths
- Assign component index to nodes



Method 3: Connected components followed by walkthrough

A closer look at the data set

- **TrackML** Detector Example
- Input to MLP for metric learning (graph construction)
- Pre-processed dataset of particles $p_T > 1$ GeV
 - Input data: cluster data with 12 input features
 - cell_count: the number of cells in the cluster
 - cell_val: the total amount of energy deposited in the cluster



- geta, gphi: the angle of the shape vector in global coordinates
- r, phi, z:cluster coordinates
- Detector information:
 - Pixels: 50 μm × 50 μm
 - ο Short strips: 80 μm × 1200 μm
 - Long strips: 120 μm × 10800 μm
- $t \overline{t}$ production overlaid with 200 soft QCD interactions
- More details about the TrackML data set can be found here ^[6]



Top view

v[mm]~ch1

Described in https://arxiv.org/abs/2012.04533

D(100) µm

Side view

w[mm]

0(100)

dimensions



Sebastian Dittmeier - Heidelberg University

dimensions



 f_p : pruning fraction, $b_q(b_w)$: bit width of activations (weights)

Sebastian Dittmeier - Heidelberg University





Exploration of HLS4ML

- Translation from offline GNN tools to FPGA implementation required
- Graph construction Metric Learning MLP as test bed
 - Translation of metric learning MLP successful PyTorch → ONNX → HLS4ML
 - Comparable performance, evaluated for 10 events
 - Efficiency and purity on edge level, connecting 2 hits belonging to the same track

Model Implementation	Efficiency	Purity	
PyTorch	0.986	0.221	
HLS converted from PyTorch	0.974	0.222	

Exploration of tools for MLP and GNN translation

- OpenHLS: study for metric learning MLP
 - <u>https://arxiv.org/abs/2302.06751</u>
 - <u>https://github.com/makslevental/openhls</u>
- FINN: Fast, Scalable Quantized Neural Network Inference on FPGAs
 - <u>https://github.com/Xilinx/finn</u>
- FlowGNN: A Dataflow Architecture for Universal Graph Neural Network Inference via Multi-Queue Streaming
 - <u>https://arxiv.org/abs/2204.13103</u>
 - <u>https://github.com/sharc-lab/FlowGNN</u>