The ATLAS Trigger system: the current and future HLT

Francesco Giuli (on behalf of the ATLAS Collaboration)



CEPC workshop Nanjing (China) 26/10/2023



Introduction





Run 3 Multi-Threaded HLT

- > The HLT was redesigned to share the same code with offline reconstruction
 - Support the Multi-Threaded mode
 - Reduce the memory footprint of the code (not an issue for online operation)
- The upgrade benefits:
 - Simplified maintenance of the code
 - General performance improvements
 - Integration of computing accelerators for future running periods
- More details can be found in <u>ATL-DAQ-PROC-2019-004</u>



= Algorithms

3

26/10/23

Run 3 HLT configuration

- The Run 3 HLT Control Flow is generated based on a list of algorithms organised in steps, performing reconstruction and selection
- The steps are combined in chains and organised in a selection menu
- The configuration is stored in JSON format and can be provided transparently to HLT in different ways:
 - From a database
 - > From a file
 - From a configuration in Python
 - From 'in-file meta-data' (mostly used for offline reconstruction)



Online performance

- The configuration of the HLT Processing Unit and its CPU resource utilisation is defined by 3 parameters
- Number of process forks
- Number of threads within the process
- Number of event slots defining how many events can be executed in parallel per node



Figure from <a>TriggerCoreSWPublicResults

Online performance

- In 2022 the online event processing was maximized with a pure Multi-Processing configuration
- A pure Multi-Threaded configuration shows lower throughput
 - It is still used for MC production, where memory savings are necessary
- Hybrid configurations also considered, giving similar gains in memory usage without throughput penalty



Figure from TriggerCoreSWPublicResults

26/10/23

ATLAS HL-HLC Trigger System

- Trigger decisions <u>much more</u> <u>challenging</u> at HL-LHC
 - ▶ Luminosity: $2 \rightarrow 7.5 \cdot 10^{34} \text{ cm}^{-2} \text{ s}^{-1}$
 - ➢ Pile-up: 60 → 200
 - From 100kHz (2kHz) to 1MHz (10 kHz) for L1 (HLT)

> ATLAS detector upgrade

New Tracker, new Timing Detector, additional muon chambers, new Tile electronics, ...





Event Filter

- Commodity hardware:
 - CPU (7.8/11.4 MHS06 for Run4/5)
 - Possibility w/ accelerators: GPU, FPGA
- Preliminary feasibility studies
 - CPU showed x8 speed-up
 - Use of GPU/FPGA looked promising
- First demonstrators started
 - Tracking, calorimeters, muons
- Better software, new algorithms and commercial accelerators to handle the computer challenges
- Technology decision about the use of accelerators in 2025



EF tracking

M and get ing FPGA and GPU
 Sti EF Fracking
 Trought of the state of the

track fitting, ambiguity removal, ...

- Exploring use of High Level Synthesis
- Hough transforms with FPGAs and Graph Neural Network
- Plan to use <u>Acts Common Tracking</u> <u>Software</u>
- Experiment independent toolkit for track reconstruction
- Support for accelerators and heterogenous options



EF Calo and PPES

- EF Calo: demonstrated topological Cell Clustering with GPUs
- Speed-up wrt CPU
 - \succ ~3.5 for di-jets at < μ > ~ 20
 - \blacktriangleright ~5.5 for $t\bar{t}$ at < μ > ~ 80
- Exploring FPGAs alternatives
- Physics, Performance & Event Selection group coordinates simulation, performance and trigger menu development for Level-0 and EF algorithms



- Toy model inspired by the New Small Wheel in ATLAS
- 4 samples produced with different noise rates: 2, 5, 10 and 15 kHz/cm²
- A target is used to emulate effect from correlated background
- Evaluation of clustering and pattern recognition performance on CPU, GPU and FPGAs





Muon cluster reconstruction

- > A cluster is formed from neighbouring hits and typically the weighted centroid of the cluster is used : $x_{cluster} = \frac{\sum_{strips} q_{strips} \cdot x_{strips}}{q_{strips}}$
- The known challenges with the standard approach are:
 - Depending on the incidence angle of the muon, a degradation is expected
 - "Correlated" background that originates from interactions with material prior to the active layers



12

MC's talk

ML is good candidate to improve clustering performance

Deep Neural Network

Input variables:

- Total number of hits belonging to the cluster
- The charge/position of the strip with highest charge
- The charge/position of its two left-right closest neighbours

Muon cluster reconstruction

Model inferred in FPGA using Vitis-AI Flow (Xilinx)



- Quantization converts 32-bit floating point weights and activations to fixed-point INT8
- Many quantization models available
- No re-training was performed, accepting a small degradation of performance



DNN: inference time results

- > CPU is already well within the latency requirements
- GPU with TensorRT improves a bit further
- > No significant gain observed over both architectures



Pattern recognition and tracking

- AIM: we want to identify which hits belong to the track and which are background hits
- Current algorithm for offline reconstruction are not optimized for NSW:
 - Large number of fakes with high occupancy, with consequent time increase

Current pattern finding algorithm: Hough transform (HT)

Implementing new machine learning algorithm to test on FPGA



Recurrent Neural Network

- Designed to deal with sequential data
- At each step, it takes as input at time t the output at time t-1

16

RNN performance results



0.95

0.96

0.97

rej

0.98

0.99

1.00

- In order to test the algorithm with Alveo cards, a CNN was also developed
- Not optimal approach for pattern recognition but useful for testing FPGA performance
- An event display translated into a 3000x16 pixel 2D image
- Convolution/deconvolution operations are used





CNN: inference time results

- CNN model successfully tested on CPU, GPU and several FPGAs
- > Overall CPU already meets the requirement imposed by the HLT latency
- > Largest improvements is seen with TensorRT on GPU
- CPU load to be studied, as well as the power dissipations



Conclusion & outlook

- During the current Run 3, redesign of ATLAS HLT framework to support the Multi-Threaded mode and to share reconstruction modules with offline
 - HLT farm upgrade, increasing the performance to 2.0M HS06 (start of 2023)
- ▶ HLT upgrade based on a mix of commodity and custom solutions for HL-LHC
 - Most projects already passes many reviews
 - Prototypes available for many projects
- ➤ Event Filter
 - Investigating accelerator options, technology decision in 2025
- Toy models have been implemented in order to investigate timing performance on commercial accelerator cards
- Choice will be a balance between:
 - Implementation complexity of novel technologies to the HLT farm and costs
 - Gain in power consumption and CPU load (yet to be studied)
 - Time performance... as well as cost
- > Preliminary studies suggest that CPU are already suitable for this task

Backup Slides



Hardware platforms



- Direct HLC implementation into FPGAs
- Require implementation of a neural network in VHDL or similar
- Significant effort to do so
- Platform developed and maintained for HEP community exists: <u>hls4ml</u>
- Fand suitable for a Level-0 trigger

Use commercial accelerator cards which offer integrated platform for deployment

Vitis AI &

- Commercially available, no ad-hoc maintenance
- Dedicated hardware and related software to translate from high level python codes, into code executable in dedicated hardware
- Not as fast as the other approach, but suitable for a HLT trigger

Vitis-Al overview

- Xilinx offers several accelerator card designed and built to accelerate ML algorithms (mostly CNN)
- The claim is that inference and throughput are improved over standard COU and GPU
- Improvements also expected in terms of power consumption



Quantisation & compilation



Testbed installation

Hardware:

- Supermicro Server, installed in Bat40-5-D05
- Graphics card: NVIDIA
 Quadro RTX A5000 PB 24 GB
- FPGA cards : Xilinx U50, U250 and Alveo VCK5000 Accelerator Card
- Study of new ML algorithms for clusters positions reconstruction and pattern recognition with the NSW (HLT Run-3)
- Performance studies of ML algorithms in FPGA(Run-4)

• SW

- Ubuntu 18.04
 - OS needs to be validated against Xilinx tools
- Docker
 - Docker GPU support
- Nvidia drivers, CUDA tools
- Xilinx Vitis (2021.2) development tool
- Xilinx Vitis-AI (2.0) docker images (Tensorflow, Pytorch, compilers, Realtime support)
- Xilinx Realtime environment (XRT) and platform support.
- Samba
- X2Go server
- DHCP server for local network
- Support
 - many tickets opened to Xilinx Support to reach this stage.
 - only a subset of all potentially interesting NN architectures is supported by the current tool, but a significant improvement is found with new releases.

cost²

25

Inference studies

Xilinx offers several accelerator card designed and built to accelerate ML \succ algorithms (mostly CNN)



The processing units









More about the Alveo VCK5000 accelerator board later in the talk

Muon identification and tracking

- Work started during software development for the New Small Wheel (NSW)
- Study the implementation of Deep Neural Network (DNN) for:
 - Identification of cluster produced by muons in Micromega (MM) and Small TGC (STGC) chambers
 - > Optimization of the single-point (i.e. by layer) position resolution
 - Cancel or mitigate the effect of backgrounds, in particular those correlated to the muon track
- Study Recurrent Neural Network (RNN) or other possible AI-based approaches (i.e. Convoluted Neural Network – CNN) for pattern recognition, muon identification and momentum measurement
 - Combine the layers info in a pattern across all muons stations
 - Compare to other existing approached i.e. Hough Transform
- Models initially thought for the offline, proved to be a good benchmark for studies on Al-accelerators

28

Why ML?

- Current methods performances worsen with correlated background cause of showers produced by muons in the last part of the JD shielding
- > ML algorithms implemented to:
 - \succ Classify clusters \rightarrow assign different errors to the fit in track reconstruction
 - Increase the single point resolution
 - Mitigate effects of correlated background with dedicated NN training



Pattern recognition and tracking



- > Truth on top, black points represent background hits, while orange ones muon signal
- Predictions from the RNN at the bottom (from black to orange, according to the RNN output)
- > The 'x' represents the hits which the output cuts (i.e. the hits used in the fit to determine p_T), while the blue line the reconstructed p_T $p_T = 0.3BR$

Inference studies

- > To run the inference:
 - ONNX: "session.Run(input_names, batch_input_tensors, output_names)"
 - TensorFlow: "Output = model.predict(input)"

Optimized CPU and GPU based model inference - Workflow



Inference studies

- Main point is that tensorRT does not work with dynamic batch sizes
- > We want to study the inference time when varying the batch size
- > Below is report the general workflow to change the batch size



Alveo VCK5000 accelerator board

- We have 1 Alveo VCK5000 accelerator board, which allows us to run Al algorithms on Vector processor Arrays
- To look at new features, like pipelining together different engine Kernels
- Vitis 2022.3.5 have been installed
- Preliminary studies on the DNN with batch size = 8:
 - > 2x (3x) slower than CPU (GPU)
 - 2x faster than U50
- Timing seems to be independent from the complexity of the model (number of layers, weights, etc.)
- <u>Next steps?</u> Same studies for the CNN (not supported currently)



