REALTIME MACHINE LEARNING IN THE CMS LEVEL-1 TRIGGER

CLUSTER OF EXCELLENCE

QUANTUM UNIVERSE



Artur Lobanov Universität Hamburg Institut für Experimentalphysik

Workshop on Realtime Machine Learning | Gießen | 10.4.20224





WHAT WE DO TODAY @ THE LARGE HADRON COLLIDER (LHC)

How collisions help us

What we want to study





Production of a Higgs boson (H) through Vector Boson Fusion (W/Z)

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24



What actually happens

Η



Partons and hadronization





THE CMS EXPERIMENT AT THE LHC



The CMS experiment: LHC camera with 100 Mpixel



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24





HOW CMS SEES PARTICLES

Different particle types can be measured with different detectors



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24





MACHINE LEARNING IN CMS / HEP EXPERIMENTS



ML applications you have seen at this mers, Jennifer Ngadiuba et al.





Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





ML IN HEP EXPERIMENTS



ML applications you have seen at this

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





EVENT SELECTION Trigger





SEARCHING FOR THE NEEDLE IN THE LHC HAYSTACK





Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





1000 W/Z bosons produced / second

1 Higgs boson is produced / second

New physics (= Anomalies) hiding here?



THE CMS TRIGGER SYSTEM

- CMS exploits a two-level trigger (filter):
 - **1. Level-1 Trigger** (L1T)
 - Implemented in hardware on FPGAs
 - Receives coarse detector data
 - **Decision within O(µs)**
 - 2. High-Level Trigger (HLT)
 - Uses CPU/GPUs in a computing farm
 - Full resolution of detector data
 - **Decision within < 1 second**



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24









L1 vs HLT resolution











THE CURRENT CMS LEVEL-1 TRIGGER



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24



















ANOMAL OCMS



DETECTION L1 TRIGGER



ANOMALY DETECTION IN CMS

- Searching for new physics at the LHC multiple fronts:
 - **Direct**: e.g. looking for exotic particles (peak or excess searches)
 - **Indirect**: precision measurements of particle parameters (e.g. H couplings)
 - **Anomaly detection** using recorded data (examples at this conference)
- All rely on existing selection (trigger) algorithms -> Model dependent or high energy thresholds

What if anomalous collisions are NOT RECORDED? \bigcirc -> Anomaly detection at trigger level!

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24









ANOMALY DETECTION WITH AUTO-ENCODERS

- Autoencoders train unsupervised on data
 - Learn to compress and to reconstruct the data
 - Difference $\hat{x} x =$ "degree of abnormality"

Real data X



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24





m/ R







ANOMALY DETECTION WITH AUTO-ENCODERS

- Autoencoders train unsupervised on data
 - Learn to compress and to reconstruct the data
 - Difference $\hat{x} x =$ "degree of abnormality"

If trained on "background" -> "signal" is anomalous!

Real data X



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24















ANOMALY DETECTION @ CMS LEVEL-1 TRIGGER

Raw detector data "in"

Raw detector images: CICADA

Reconstructed objects: AXOL1TL





Artur Lobanov | Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





HIGH-LEVEL INPUTS: AXOL1TL





AXOL1TL: ANOMALY DETECTION WITH OBJECT TOPOLOGY

- AXOL1TL (Anomaly eXtraction Online Level-1 Trigger aLgorithm) is a variational auto-encoder: Encodes input as a distribution over the latent space
 - Add regularisation term in loss: KL divergence, how different is distribution from Gaussian
 - **Inputs: L1 trigger objects 4-vectors** (pT, η , ϕ)

hls 4 ml

Most energetic 4 electron/photons, 4 muons, 10 jets and missing transverse energy (MET)























AXOL1TL: ARCHITECTURE OPTIMISATION

• Full NN architecture does not fit the L1/FPGA constraints

-> only use encoder half of the network

- Compute degree of abnormality from latent space directly •
- No need to use inputs for anomaly score computation •
- Half network size and latency! •



Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24



CMS-DP-2023-079





AXOL1TL: COMPRESSION

- - Narrow, shallow model, aggressively quantised
- Output is one vector [13,1], corresponding to μ part of [μ , σ] KL loss (dropping σ as it is small -> reduces processing time)
- Anomaly score: sum squared of the µ vector \bigcirc



Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24



Quantization-aware training with <u>QKeras</u> and FPGA adaptation with <u>HLS4ML</u>



AXOL1TL: FPGA IMPLEMENTATION



<u>CMS-DP-2023-079</u>

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





Implemented on Xilinx Virtex-7 XCVU9P FPGA Met requirements on latency and resources

50 ns latency & ~1% resources

Resource utilization of Virtex-7 FPGA chip on Imperial College MP7 µGT board

	Latency	LUTs	FFs	DSPs	BRAM
XOLITL	2 ticks 50 ns	2.1%	~0	0	0

AXOL1TL: COMMISSIONING



<u>CMS-DP-2023-079</u>

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24



• AXOL1TL is trained with unbiased data collected testing)

during lard triggers



AXOL1TL: PIPELINE



Development

Artur Lobanov | Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24 |



Conversion

Implementation / Validation



AXOL1TL: FIRMWARE VALIDATION WITH TEST BENCH

- Trigger bits for the L1 menu including 4 anomaly detection thresholds: scores >1250, >250, >25, and >5 from top to bottom
- Test vector column: generated from inference results of a standalone C++ emulator
- HW count: comes from standard global trigger firmware simulation workflow using ModelSim. Perfect bit agreement is observed

ldx	L1 Menu Algorithm Name	Test Vector Count	HW Count	Agreement
94	L1_ADT_20000	0	0	\checkmark
95	L1_ADT_4000	29	29	\checkmark
103	L1_ADT_400	2618	2618	\checkmark
108	L1_ADT_80	3331	3331	\checkmark



Test vectors generated from Run 3 data

















AXOL1TL: FW VALIDATION

Test Crate Validation

L1 Menu Algorithm Name	Test Crate Count	Standalone Emulator Count
L1_ADT_20000	1	1
L1_ADT_4000	742	741
L1_ADT_400	21236	21229
L1_ADT_80	25468	25481

Anomaly Detection hardware vs. emulation trigger mismatches. Events from promptly reconstructed 2023 Ephemeral ZeroBias data where hardware bits are recorded from configured μ GT test crate. In table (left), Test Crate Count shows events triggered in hardware and read out into data and Standalone Emulator Count is evaluated via offline inference with L1 objects. Anomaly score distribution of all events (right): red segments represent mismatches between hardware and emulation. Clustering near decision boundaries implies issue is due to precision/rounding problem. Minimal mismatches in hardware vs. emulation ($\leq 1\%$) observed.





AXOL1TL: EVENT DISPLAY



CMS Experiment at the LHC, CERN Data recorded: 2023-May-24 01:42:17.826112 GMT Run / Event / LS: 367883 / 374187302 / 159





Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24





- Example of an anomalous event during 2023 pp collisions (from random trigger dataset)
 - **Highest anomaly score** event not triggered by L1
- L1 objects:
 - 11 jets with pT > 20 Gev
- Offline objects:
 - 7 jets with pT > 15 GeV from the same vertex
 - 75 identified vertices

















AXOL1TL: PHYSICS PERFORMANCE

- Use simulated hypothetical exotic signal as a anomaly candidate
- Significant performance improvement on various SM and be by adding AXOL1TL to the 2023 trigger menu

L1 Efficiency w/ AXOL1TL@freq L1 Efficiency w/o AXOL1TL Improvement =

• Example performance improvement for H->aa[15 GeV]->4b s

A	AXOLITL Rate	1 kHz	
	Signal Efficiency Gain	46%	
Sig	nal Efficiency Gain	46%	

Starting data-taking with ~O(100) Hz L1 rate in 2024 pp collisions soon!

<u>CMS-DP-2023-079</u>

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24



- 5 kHz 10 kHz 100% 133%
- 133% 100%























EATURES:



CICADA: ANOMALY TRIGGER ON RAW INPUTS



- CICADA (CMS DP-2023/086): Anomaly Detection Algorithm
- Using raw inputs of calorimeter: Image of 18 x 14 energy deposits

 - Independent of domain knowledge (standard trigger algorithms)
- Convolutional auto-encoder trained on background dataset: signal -> anomaly!



CMS DP-2023/086



Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24











CICADA: KNOWLEDGE DISTILLATION

- Full CICADA model is too complex for FPGA resources / L1 Trigger requirements -> use Student-Teacher Knowledge Distillation
 - **Teacher model**: complete encoding and decoding of the original input data
 - **Anomaly score (reconstruction error)**: average of the squared error (predicted input) in reconstruction for each of the 252 individual energy deposits (Mean Squared Error)
 - Student model: regresses the anomaly score of the teacher model



CMS DP-2023/086

Smaller convolutional layer with only 4 filters - his 4 mi Ke layers -> 10x faster & less resources -> fits FPGA/L1T requirements













CICADA: COMMISSIONING

• CICADA currently being commissioned in the L1 Trigger test system

- Software-based emulation based on Firmware (HLS4ML) and validated
- Preliminary performance estimates promising + operational stability tested

• This is the first anomaly detection on low-level inputs in a LHC trigger system!



Artur Lobanov | Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24 |







TOWARDS THE High-Luminosity LHC















CMS L1 TRIGGER FOR THE HIGH-LUMINOSITY LHC

- High-Luminosity phase of the LHC (HL-LHC) will start in 2029: 3x higher instantaneous luminosity and pileup wrt current conditions
 - CMS will upgrade most of its detectors, including all (trigger) electronics
- L1 Trigger for the HL-LHC: \bigcirc
 - Bandwidth: 2 -> 63 TB/s
 - Output 100 -> 750 kHz
 - Latency: 4 -> 12 us
- Tracking @ L1T + new processing systems will enable "offline-like" reconstruction



Latency

5 us

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24









FPGAS: WORKHORSE OF THE CMS LEVEL-1 TRIGGER





Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24







L1 EVENT CLASSIFICATION @ HL-LHC

- ML-based triggers proposed in the <u>L1T "TDR</u>" for the High-Luminosity LHC
- **Classifier approach**: binary classifier for known signals trained on simulation (<u>Note</u>)
- **Anomaly detection**: auto-encoder based on L1 trigger objects (as AXOL1TL) \bigcirc
 - Sensitivity at the ~same order as of the classifier approach (e.g. VBF H>inv)





• Tests of AXOL1TL and CICADA pave the way for anomaly triggering at the HL-LHC in CMS!







VERTEX RECONSTRUCTION





CONTINUAL LEARNING FOR VERTEXING

- Need to deal with changing detector conditions: ageing, noise, LHC conditions, etc
 - Normally do dedicated training/algo optimisation with full dataset
- **Continual Learning: train a model with** a continuous stream of data
 - Learns from a sequence of partial experiences rather than all the data at once
 - Update model to changing conditions without large MC production
- Method tested on Vertex reconstruction:
 - data



Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





CL outperforms a simple retrained model when detector defects are applied to the training







JET FLAVOUR CLASSIFICATION <u>CMS Note DP2022 021</u>

Vertex Finding: b-tagging:



Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24











MUON RECONSTRUCTION

- **Regressing the muon momentum** based on the hits in the muon detectors
- Based on features extracted from previous track finding





	ME1/1	ME1/2	ME2	ME3	N
ф	1	1	1	1	
θ	1	1	1	1	
bend	1	1	1	1	
quality	1	1	1	1	
time					



Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





HADRONIC TAU RECONSTRUCTION

- Hadronic Tau reconstruction very challenging in the L1 Trigger environment
- New CNN approach for tau reconstruction with calorimeter-only information \bigcirc followed by two NNs for pT regression and ID
 - Outperforms baseline algorithm: better efficiency and lower fake rate





Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





Single-τ_h Rate [kHz

- - Triggerless readout <> lower quality objects
- Run-3 demonstrator system reads L1 objects (muons, calo jets, EG, taus) with very heterogenous system:
 - 3 boards: KCU1500, SB-852, VCU128)
 - Output technologies: DMA, TCP/IP
- Studying **ML methods for** realtime calibrations with HLS4ML & proprietary SW (Micron Deep Learning Accelerator)
- **Detailed overview here and here** \bigcirc

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24















ML IN DETECTOR READOUT ASICs



ML applications you have seen at this

Artur Lobanov | Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24 |







AUTOENCODER IN CMS CALORIMETER TRIGGER ASIC

- CMS will be using an (auto)encoder for compressing the data intelligently in the High-Granularity Calorimeter concentrator ASIC (HL-LHC)



FPGAs were designed for ASIC prototyping —> used HLS4ML for design!



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24



Data volume needs compression at detector readout: done by ASICs (fast/efficient!)

Decoding done on the FPGA side of the Level-1 Trigger (or direct use in NN@FPGA?)













MARY



FAST ML IN CMS



ML advancing from OFFLINE (Hz) to ONLINE (MHz) applications





Artur Lobanov | Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24 |



lavior Duarta I hl c/ml





ANOMALY DETECTION WITH THE CMS LEVEL-1 TRIGGER

- Various anomaly searches for new physics performed at the LHC \bigcirc
- Opening **a new direction**: anomaly detection in the CMS Level-1 Trigger
 - **Challenging environment for L1T**:
 - Hardware/FPGAs: restricted resources and latency (ns!)
 - Physics: <60> simultaneous collisions, only calorimeter and muon detector data

Two auto-encoder approaches being commissioned in CMS: \bigcirc

- **AXOL1TL**: using high-level physics objects [CMS-DP-2023-079]
- **CICADA:** using raw detector data [CMS DP-2023/086]
- **Promising prospects for anomaly triggering in CMS!** [HL-LHC L1T]

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24













ML IN CMS L1 TRIGGER FOR HL-LHC

- - All based on HLS: <u>HLS4ML</u> (NN) and <u>Conifer</u> (BDT)



Image by Sioni Summers

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24





• Current Run-3 algorithms (AXOL1TL/CICADA) are an important test bed for the future













xkcd "Machine Learning"



REFERENCES

- CMS Collaboration. "Anomaly Detection in the CMS Global Trigger Test Crate for Run 3". CMS-DP-2023-079, CERN-CMS-DP-2023-079 (2023): <u>https://cds.cern.ch/record/2876546</u>
- CMS Collaboration. "Level-1 Trigger Calorimeter Image Convolutional Anomaly Detection Algorithm", CMS-DP-2023-086;
 CERN-CMS-DP-2023-086 <u>https://cds.cern.ch/record/2879816</u>
- J. Pearkes, "Realtime Anomaly Detection in the CMS Experiment Global Trigger Test Crate", ML4Jets 2023, <u>https://indico.cern.ch/event/1253794/timetable/?view=standard#57-realtime-anomaly-detection</u>
- N. Zipper, "Testing a Neural Network for Anomaly Detection in the CMS Global Trigger test crate during Run 3", TWEPP 2023 <u>https://indico.cern.ch/event/1255624/contributions/5444028/</u>
- C. Sun, "Realtime Anomaly Detection in the CMS Experiment Global Trigger Test Crate", FastML 2023, <u>https://indico.cern.ch/event/1283970/contributions/5554350/</u>
- CMS Collaboration. "CMS Technical Design Report for the Level-1 Trigger Upgrade", CERN-LHCC-2013-011; CMS-TDR-12 <u>https://cds.cern.ch/record/1556311</u>
- E. Govorkova, et al. "Autoencoders on field-programmable gate arrays for real-time, unsupervised new physics detection at 40 MHz at the Large Hadron Collider". Nat. Mach Intell. 4, 154 (2022). <u>https://doi.org/10.1038/s42256-022-00441-3</u>
- FastML Team. hls4ml (Version v0.7.1) [Computer software]. <u>https://doi.org/10.5281/zenodo.1201549</u>
- J. Duarte, et al. "Fast inference of deep neural networks in FPGAs for particle physics". JINST 13, P07027 (2018). <u>https://doi.org/10.1088/1748-0221/13/07/P07027</u>





,

าท

REFERENCES

- More results on Fast ML in CMS and beyond were shown at the Fast ML for Science Workshop series, e.g. <u>https://indico.cern.ch/event/1283970/</u>:
 - Fast ML inference in FPGAs for the Level-1 Scouting system at CMS
 - Realtime Anomaly Detection in the CMS Experiment Global Trigger Test Crate
 - Harnessing charged particle tracks in the Phase-2 CMS Level-1 Trigger with ultrafast Machine Learning
 - B-tagging and Tau reconstruction in the Level-1 Trigger with real-time Machine Learning
 - A Convolutional Neural Network for topological fast selection algorithms in FPGAs for the HL-LHC upgrade of the CMS experiment











original input from the latent space.

Loss =
$$(1 - \beta) \|x - \hat{x}\|^2 + \beta \frac{1}{2} (\mu^2 + \sigma^2 - 1 - \log \sigma^2)$$

Reconstruction term Full regularization term

Equation: VAE loss function. The reconstruction term is computed from the difference between the input (x) and output (\hat{x}) of the VAE. The second, full regularization term, is the Kullback–Leibler divergence (KL-divergence) between the latent space distribution and a standard normal distribution with mean μ and standard deviation σ . The parameter β can be tuned to balance the reconstruction performance with more efficient latent space encoding. At inference time, the loss is approximated by the mean-squared term $\Sigma \mu_{i^2}$ of the KL-divergence for latency considerations. This approximation has no impact on performance.



The AXOL1TL anomaly detection uses a Variational Autoencoder (VAE). A dense feed-forward neural network reads in (p_{τ} , η , ϕ) hardware inputs of 19 L1 objects. The encoder network computes a latent space vector of Gaussian probability distributions, $N(\mu_8, \sigma_8)$. The decoder network reconstructs the







CICADA: Anomaly detection on Raw inputs



Shown here is a comparison of the teacher model ability to reconstruct a Zero Bias (ZB) beam event (original: far left, reconstructed: center left) versus a signal sample, Soft Unclustered Energy Patterns (SUEP) on the right (original: center right, reconstructed: far right). In general, the teacher model is better able to reconstruct the Zero Bias beam event as evidenced by a far lower loss (0.81) compared to the SUEP loss (14.21). This example shows how the CICADA anomaly detection mechanism works to find anomalies. From [CMS DP-2023/086]

Artur Lobanov Machine Learning @ CMS L1 Trigger Workshop on Realtime ML, 10.4.24













• <u>hls4ml</u>: package for translating NN to FPGA firmware



Artur Lobanov Machine Learning @ CMS L1 Trigger | Workshop on Realtime ML, 10.4.24



