

CLUSTER OF EXCELLENCE QUANTUM UNIVERSE



Bundesministerium für Bildung und Forschung



# Real-time ML event classification with FPGAs at the LHC

Finn Labe on behalf of the UHH ML@L1 team

Deep Learning Roundtable at DESY | 08.12.2023



#### Why do we need ML in the trigger?

 Trigger performance decides what data is available for offline analysis!

- Sensitivity limited by trigger thresholds
  - Example: HH → bbWW (single lepton)
  - Possible solution: ML@L1







 Trigger performance decides what data is available for offline analysis!

- Sensitivity limited by trigger thresholds
  - Example: HH → bbWW (single lepton)
  - Possible solution: ML@L1







08.12.23

 Trigger performance decides what data is available for offline analysis!

- Sensitivity limited by trigger thresholds
  - Example: HH → bbWW (single lepton)
  - Possible solution: ML@L1







CMS Run 2/3 L1 trigger system



#### Total L1 rate 100 kHz (Run 2/3) 750 kHz (HL-LHC)



### **Event classification in the L1 trigger**

Event classification in global trigger

Uses objects defined further

Binary network, cut on output

upstream in L1 trigger system

node defines L1 trigger condition



### **Event classification in the L1 trigger**



- Event classification in global trigger
  - Uses objects defined further upstream in L1 trigger system

 Binary network, cut on output node defines L1 trigger condition

L1 trigger built from Field-Programmable Gate Arrays



ML on FPGAs

750 kHz (HL-LHC)

#### -

#### 7

Supervised training: simple NN

**Trigger NN development** 

- Discriminating signal from MinBias\*
- Using  $p_T$ ,  $\eta$ ,  $\varphi$  of L1 objects

- Comparing to cut-based triggers:
  - NN outperforms at any given rate

#### ML@L1 is promising!





optimizatior

#### 8

### **Neural Network on FPGA**

Conversion for FPGA using hls 4 ml library

- Compressing network to fit FPGA
  - Single (hidden) layer network
  - Pruning: removing connections
  - Bit precision: less bits per weight
- **75 ns** latency (target < 100 ns)
- 2% of an FPGAs resources (6 available)







00:00 01:00 02:00

### **Integration tests**

Testing NN trigger in CMS operation\*

Probe stability over wide threshold range

\*in so-called test crate  $\rightarrow$  no detector readout

23:00

07-Jun

Rate [Hz]

10<sup>6</sup>

**CMS** *Private Work* 



Ы

0.467 fb<sup>-1</sup>, 2023 (13 TeV)

Data

100

06:00

Time

03:00 04:00 05:00

#### H<sub>T</sub> > 320 GeV L1 topo score > 25



0.7

45

50

55



### **Integration tests**

Testing NN trigger in CMS operation\*

PU dependency comparable

to traditional cut-based triggers

Probe stability over wide threshold range

\*in so-called test crate  $\rightarrow$  no detector readout

**CMS** *Private Work* 

MET > 90 GeV



<PU>

Data

60

### Integration tests





Probe stability over wide threshold range

 PU dependency comparable to traditional cut-based triggers

#### Integration results very important to us: this kind of trigger is feasible!





### Increasing analysis sensitivity

- Studying expected sensitivity gain from NN-based L1 algorithms
  - Simulation-based integration in run 2 HH → **bbWW** analysis
- L1 NN targeting 10 kHz L1 rate
  - Resulting pure rate much lower



### Increasing analysis sensitivity

- Studying expected sensitivity gain from NN-based L1 algorithms
  - Simulation-based integration in run 2 HH  $\rightarrow$  bbWW analysis
- L1 NN targeting 10 kHz L1 rate
  - Resulting pure rate much lower

- How to handle HLT strategy?
  - Use cut-based trigger
  - Train another neural network





Events [a.u.]

### Increasing analysis sensitivity

- Studying expected sensitivity gain from NN-based L1 algorithms
  - Simulation-based integration in run 2 HH  $\rightarrow$  bbWW analysis
- L1 NN targeting 10 kHz L1 rate
  - Resulting pure rate much lower

- How to handle HLT strategy?
  - Use cut-based trigger
  - Train another neural network





#### Summary & outlook

- Supervised ML@L1 promising for otherwise trigger-limited signals
- Integration of NNs in the CMS level 1 trigger demonstrated!
- Interesting open questions
  - Can we generalize supervised ML in the trigger system?
  - Efficiency measurement



Another approach is anomaly detection in the L1T: UHH with <u>AXOL1TL</u> and <u>CICADA</u>





- Project lead Artur Lobanov ( = AL, postdoc)
  - PIs Johannes Haller, Gregor Kasieczka (Profs)
- Higgs Expert Matthias Schroeder (Staff)
- Supervised ML triggers Finn Labe (PhD), Shahin Sepanlou (MSc), Ihor Komarov (MSc), Salome Fresenbet (BSc), Karla Kleinbölting (BSc, MSc), AL
- Anomaly detection at L1 Sven Bollweg (PhD), Lars Emmerich (BSc), Susan Sefidrawan (BSc), Karim El Morabit (postdoc), AL
  - ML jet identification Philipp Rincke (MSc), Karim El Morabit (postdoc), AL
    - L1T Menu Run3: Sven Bollweg (PhD), AL Phase2: Daniel Hundhausen (PhD), Matteo Bonnanomi (postdoc), AL
    - L1 Run3 DQM Mathis Frahm (PhD), AL



## Backup

### L1 neural network details



Binary fully connected network



#### ReLu activation, BCE loss

- 26 input variables
  - MET, MET  $\varphi$
  - 4 jets: p<sub>T</sub>, η, φ
  - 2 muons: *p<sub>T</sub>*, *η*, *φ*
  - 2 "egammas": *p<sub>T</sub>*, η, φ

Electrons and photons indistinguishable without tracker

- **Compression details**
- Lowering **bit precision** to (6,1) for inputs & weights
  - Compression performed using **qkeras** package











The following offline cuts were used to study the trigger:

- Exactly one muon with CutBasedIdTight, PFIsoTight and  $|\eta| < 2.4$
- At least three AK4 jets fulfilling tight PuPPi criteria
- At least one of them fulfilling loose DeepJet b-tagging criteria
- Distance between lepton and jets  $\Delta R > 0.2$
- Transverse mass  $M_T(\ell + E_T^{\text{miss}} + b_{\text{lep}}) > 60 \text{ GeV}$

In the offline analysis, an event classification DNN is next.



- Training & evaluating using Run 3 "ZeroBias" sample
  - Signal efficiency at 5 kHz of rate is below 500 Hz pure rate



