

ATLAS





Fast Deep Learning on FPGAs for the Phase-II LO Muon Barrel Trigger of the ATLAS Experiment

IML2019 April 17th 2019

Luigi Sabetta

#### **Motivations**

The ATLAS Level-0 muon trigger will face a complete upgrade HiLumi-LHC

#### Detector parameters:

- Pile up:  $30-40 \rightarrow 200$
- Luminosity:  $(2 \rightarrow 7.5) \times 10^{34} \text{ cm}^{-2} \text{s}^{-1}$

#### Greater muon hit rate in the spectrometer: Up to 600 Hz/cm<sup>2</sup>



#### ATLAS improvements: New trigger processor

- New FPGA based system
  - Virtex UltraScale+ XCVU13P
  - Logic cells (K): 3780
  - Memory (Mb): 455
  - GTY Transreceivers (32.75 GB/s): 128
  - I/O Pins: 832

#### New trigger station

- New RPC layer

#### Trigger algorithms improvement

- Need to be very fast and flexible

#### Neural networks→ valid candidate

ATLAS-TDR-026



Luigi Sabetta – IML2019

Resistive Plate Chamber(RPC):

- Fast gas detector
- Trigger sources in ATLAS

#### Sector:

\_

- Divide ATLAS in 2 sides ( $z \ge 0$ ) \_
- Divide  $\phi$  coordinate in 8 sectors, each one divided in \_ 2 (Small and Large)
- Take all the strips from the obtained sector



# Atlas Trigger Requirements

# Design Trigger Requirements $\rightarrow$ up to three candidates

#### **Candidate** = muon with $p_T$ > threshold, + extimate of the position ( $\eta^{muon}$ )

Processing time < 6μs

#### Standard trigger Algorithm



Standard algorithm:

 Check the presence of coincidences inside windows sequentially opened based over the previous layers's hits

台 Stable and reliable 合 Good performances

The coincidence windows need to be tuned "by hand"

- Windows dimensions depend over the  $\mathbf{p}_{\mathrm{T}}$  trigger threshold
- Strong dependency over the trigger choices and local detector geometry
- P Assumes pointing tracks
- Decay vertex far from the IP?



Luigi Sabetta – IML2019

ATLAS-TDR-026

# From strip map to images



Example of a muon with  $p_T$ =19 GeV + noise

# **Convolutional Neural Networks**





Convolutional Neural Network (CNN): Deep Neural networks optimized for image recognition.

Highly effective for problems with rotational or translational symmetries

In this way the number of the parameters can be highly reduced

# **CNN Floating Point structure**



Architecture:

- (Conv2D + Batch Norm. + Max Pooling) x 3
- (Dense) x 2 layers to get to the output

Total number of parameters: 500k

Decaying Learning rate:  $10^{-2} \rightarrow 10^{-5}$ 

Activation= ReLu

**BatchNormalization**:

- $\epsilon = 10^{-6}$
- Momentum= 0.9

#### Input DataSet

Sample:

- 1 muon + Background (random hits ~uniformely distributed)
- 2 muons + Background
- 3 muons + Background
- Just noise (10% of the total number of images)

$$\label{eq:pt} \begin{split} 0 < p_T < 20 \; GeV \\ 0 < |\eta| < 1.05 \end{split}$$

Parametrization of full phase-2 events in the ATLAS detector

Grand total of  $\sim$  900k images





#### **CNN FP performances – Physics quantities**



Interesting physics quantities are well represented

Events are well classified in  $\,n^{muons}$ 



# **CNN FP performances- Trigger**

CNN

**Standard Algorithm** 



Efficiency curves are comparable to the ones obtained with the standard algorithm

#### **Ternary CNN – implementation**

Level-0 trigger sector logic will be implemented in an FPGA

- Standard NNs work with weight described by 32 bit floating Point precison numbers
- FP weights aren't the optimal choice Great logic resources consumption

Ternary CNN:-101Weights =(only two bits)

A Ternary CNN may imply a loss in performances

 Few percentage points in respect to an FP32 NN with the same structure

#### Smaller size:

- Smaller logic resources consumption
- Up to 16 times smaller

#### In principle it can be made deeper

More layer recover the loss in precision

# **Ternary CNN-Performances**



Luigi Sabetta – IML2019

p<sup>ML</sup>(GeV)

MAE=1.9 GeV

# Ternary CNN performances – Number of muons



## NN on FPGAs – How to Implement



## NN on FPGAs – How to Implement





# Next Steps



# Thank you for the attention