# Status report of FPGA-related activities in UNIPD

*J.Pazzini*
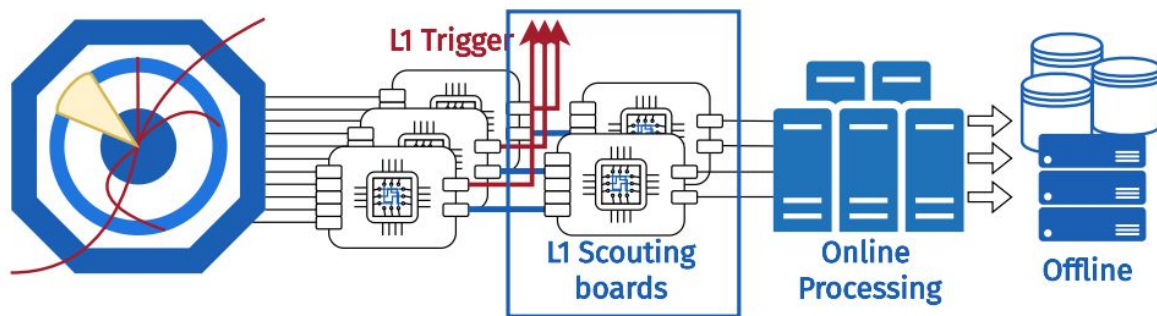*Padova University & INFN*

CN1 - WP2 Meeting
21 Nov 2023

# Introduction

- Group of a few units in UNIPD Physics and Astronomy Department, mainly active in the CMS experiment
- Close collaboration with the CMS DAQ and Level1 Trigger groups, to which we have shared students with periods as Trainee, Doctoral and Technical students
- Activities involving FPGAs at various levels
  - mostly related to algorithms for Trigger and DAQ
  - also targeting the development of data transfer protocols on FPGA
  - not devoted to the usage of FPGAs as computing accelerators
- Person-power
  - A couple of senior members, including Andrea Triossi our in-house expert in FPGA programming at the descriptive level (VHDL)
  - 1 PhD Student (Sabrina Giorgetti) hired directly with PNRR funds
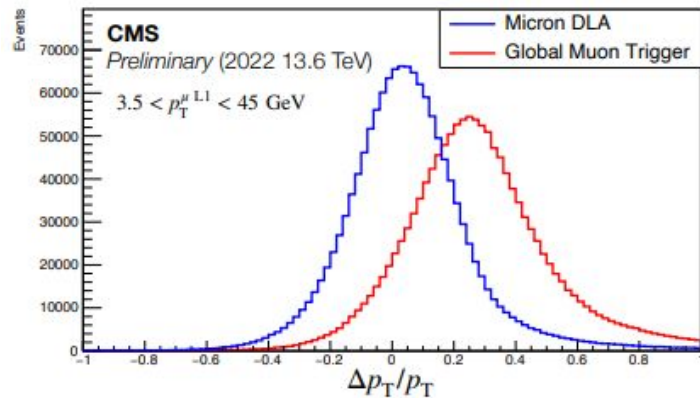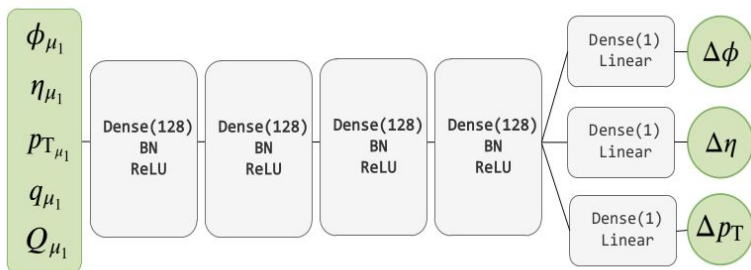  - Some PhDs student working on the topic, at different levels

# The CMS L1/data scouting project

- Project aimed at capturing and processing online the CMS trigger/data streams @ full bx-rate
- Perform on-the-fly analyses on the objects available in the trigger/daq chain before the L1 filter decision
    - Started as early as a few years ago (~2018)
    - Long-term goal is the a running in HL-LHC, leveraging on the L1 trigger upgrade
    - Currently in its first level of deployment at CMS
- Several ongoing activities:
    - Mainly related to data acquisition and utilization from current L1 trigger chain
    - SW and FW studies in anticipation of the Phase-2 upgrade of the CMS L1 trigger, with new, better-quality objects
    - Computing framework for online analysis and anomaly detection techniques
- In addition to the activity on the L1, it is possible to "do scouting" already at the FE level of the detectors, *iff* throughput allows
    - A demonstrator of the 40MHz data scouting of the data from the FE of the CMS's Drift Tubes is in place
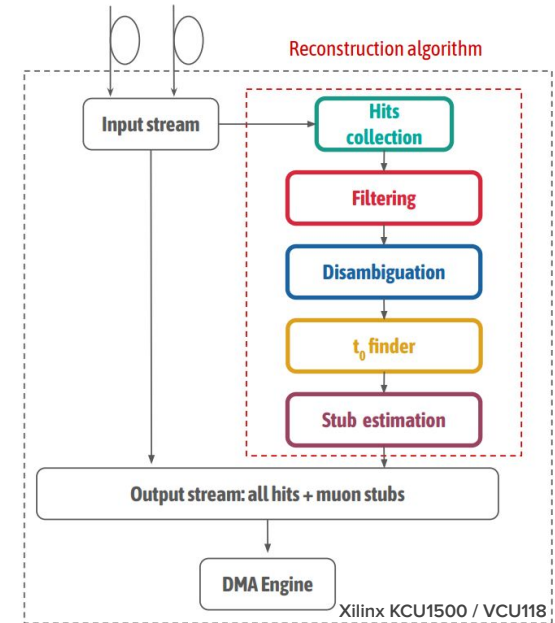
# L1 scouting and Online calibration algorithms

- Scouting activities underway at CMS within a sub-group of DAQ dedicated to the purpose
  - Our PhD students are working on both the FW part (receiving and pre-processing the streams from the L1 boards) and the SW part (online and offline processing)
- Data analysis from L1 scouting is subject to the characteristics of the trigger object, and their mis-calibration:
  - Trigger objects calibrated for a given efficiency at a threshold
  - And missing the offline-level calibration "knowledge"
- Activity on calibration methods in FW was carried out based on ML (primarily, NNs):
  - Use the offline objects as target to correct and re-calibrate the trigger level objects
- Initial development based on Micron SB-852 boards and the MDLA ML-framework, a dedicated framework from Micron
  - Now switching to Xilinx VCU128 boards and HLS4ML
- First deployment targeting the L1 Muon objects:
  - Fake-muon detection
  - Momentum calibration

# Online reconstruction algorithms for data scouting

- Upgrade of the CMS' Drift Tubes Front-End electronics, moving from a ASIC-based to a FPGA-based boards
  - 1 sector equipped with the new electronics at LHC P5
  - Local mock-up of the detector and electronics at Legnaro National Laboratories
- Data scouting of the new DT electronics
  - Trigger-less readout of the detectors' hits
  - Their pre-processing (forming stubs/tracklets from detector hits)
- Developed and tested an hybrid model for the local online reconstruction of the stubs
  - incorporating 2 NNs for noise reduction and hits clustering, combined with an analytical algorithm for the estimation of the stubs' parameters
  - Use of HLS4ML on Xilinx platform (KCU1500 and VCU118 eval boards)
- Recently, new ideas for the development of a similar algorithm based on Spiking Neural Networks
  - Biologically inspired, as the network fires only if neurons receives a series of stimuli within a relatively short time window
  - Work started as a MSc thesis, with a simplified model running on simulated dataset in SW
  - The goal is porting the algorithms to FPGA, by writing the low-level IP cores describing nodes and their activation



[*] https://doi.org/10.48550/arXiv.2105.04428
[**] https://doi.org/10.1016/j.nima.2022.166869

CHEP23 - M.Migliorini et al *"Triggerless data acquisition pipeline for Machine Learning based statistical anomaly detection"* [1]
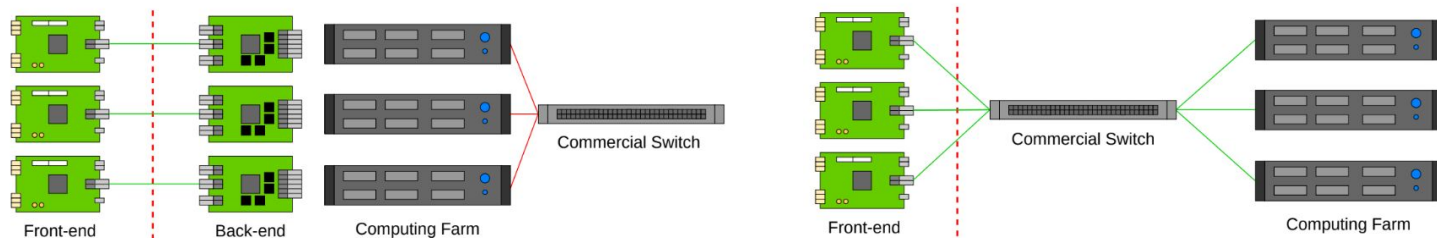TWEPP23 - M.Migliorini et al *"40MHz Triggerless Readout of the CMS Drift Tube Muon Detector"* [2]
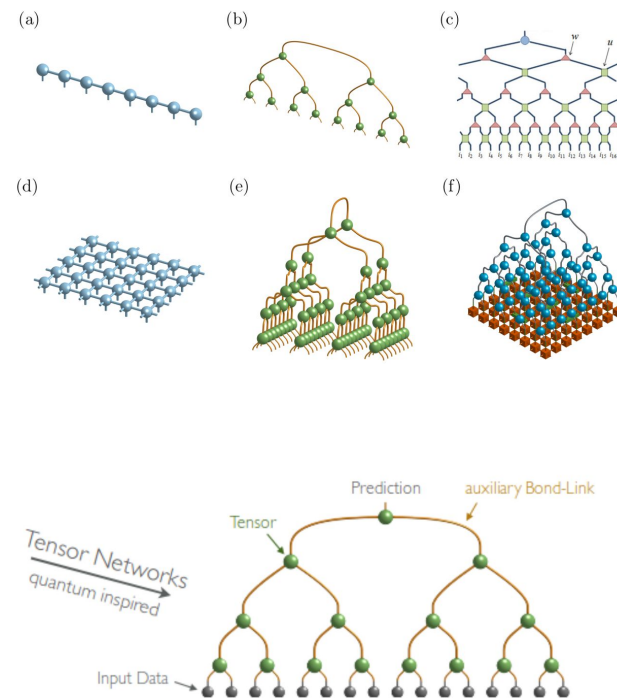
# Remote DMA from detectors' FE

- Ongoing work begun as a 3 year project started beginning 2023, and financed by INFN CSN5, with PI Andrea Triossi
- Remote DMA from FE of detectors to servers for acquisition and processing
- Use of commercial standards, with a multi-vendor ecosystem, based on Ethernet networks
  - ➜ RDMA Over Converged Ethernet (ROCE)
- Move the adoption of the network protocol to the data producer (devices on the detectors' FE):
  - Front-end initiates the RDMA transfer
  - No point-to-point connection between front-end and back-end
  - Dynamical switching routing according to node availability
- Work aimed at creating a very lightweight FW IP for transmission initiation from FPGAs suitable for FE
  - Implementation of a Dynamic Simulation
  - Development starting from the network stack implemented by ETH
  - Tested UDP and TCP in Xilinx VCU118
  - Now moving towards ROCEv2, working on re-writing the FW blocks ➜ aiming at deploying and testing ROCEv2 on VCU118



TIPP23 - G.Bortolato et al "*Front-End RDMA Over Converged Ethernet, real-time firmware simulation*" [3]
TWEPP23 - G.Bortolato et al "*Front-End RDMA Over Converged Ethernet, real-time firmware simulation*" [4]

# Tensor networks

- Activity just very recently started, aimed at implementing Tree Tensor Networks as an alternative to Neural Networks for online classification problems
  - Tensor Networks provide representations of quantum many-body states with a purely algebraic description, based on contraction and factorizations of very large tensors into networks
  - TTNs are a subcategory of TNs with a specific topology that can be used for classification
- The contraction operations are purely algebraic and do not require nonlinearities such as activation functions in NNs
- To each node and connection can be associated a notion of "entropy", which can be used to simplify the network while retaining its expressivity
  - Similar to NN pruning, but informed base on the actual "information-content" rather than purely on the wight connection
- Colleagues from Padova already released studies on a TTN-based flavor tagger for LHCb for offline analysis
- The TTN features make them potentially very attractive for online use:
  - "Simple" algebraic operations are expected to be easily deployable with manageable use of resources such as DSPs
  - The reduction of the network complexity can be performed while retaining a given
- Development has begun on two fronts:
  - On the SW-side, model development and training, devoted to solving typical L1 classification problems (identification/tagging)
  - On the HW-side, development of prototypes of the TTN FW nodes and the related contraction operations
- The activity is carried out in close collaboration and synergy with colleagues from CN1 Spoke10





7

# System on Chip

- New Hybrid platforms are increasingly emerging, with chips integrating FPGAs and processors
- AMD/Xilinx is pushing in this direction with the Versal
  - As of right now, a very broad and diverse product lineup
- Potentially a very interesting platform for various applications, including the development of the projects described in this short report:
  - AI Engines (units suitable for matrix/tensor computation)
    ➜ ML for L1 Scouting, Spiking NNs and Tensor Network applications
  - High-Bandwidth Memories (HBM)
    ➜ potentially, but not necessarily, suitable for high throughput applications, if also equipped with large IO
  - Co-existence of the custom logic and processors (including real-time)
    ➜ suitable for the FEROCE application
- One evaluation board mainly devoted to IO and AI Engines (but not including HBMs) is being purchased with PNRR funds