

### Al for fast inference in real-time systems: status and future perspectives 2026 Update: European Strategy for Particle Physics

Jun 23–27 2025, Venice Lido

## EnHzürich

#### Thea Klæboe Årrestad (ETH Zürich)







#### **Time to process**





# This readout technology doesn't (yet) exist!

## Ultimate precision at future colliders could be bottlenecked by our data acquisition systems



#### Example vertex detector inner layer from before (#241 FCC-ee):



#### Current technology 3.2 - 6.4 Gbit/s

32 bits pixel data in inner layer @ 200 MHz/cm<sup>2</sup>

#### 24.4 Gbit/s

### We have two options:



### $\rightarrow$ ML on ASIC/FPGA

### <u>#67 #124</u> ML inference at low latency **HEP Tools and Communities**





**HEP quantization libraries:** 







(Collaboration with AMD: Brevitas for PyTorch)



#### Co-processing kernel (Xilinx accelerators/SoCs)

#### **HEP hardware ML libraries:**



## **\Conifer**





#### **ASICs**









### Frontend (QAT)

Q K Keras

## **PYTÖRCH Brevitas**



HQG





### Intelligent frontends **Reduce on-detector**

- Subtle data patterns within single (pixel) layer! ML on ASIC on-detector
  - filtering, and/or
  - featurizing

#### Challenges

Metric	Simulation	Target
Power	48 mW	<100 mW
Energy / inference	1.2 nJ	N/A
Area	2.88 mm <sup>2</sup>	<4 mm²
Gates	780k	N/A
Latency	50 ns	<100 ns

Radiation hardness

#### Signal

#### Background





### #<u>11</u> #<u>95</u> #<u>247</u> #<u>233</u> #<u>93</u> #<u>211</u> #<u>272</u> **Smart Pixels** Pixel readout ASIC with ML

- <u>frontend filtering</u>: discard low-p<sub>T</sub> tracks
- <u>feature extraction</u>: particle position+angle Mixture Density Model



### #<u>11</u> #<u>95</u> #<u>247</u> #<u>233</u> #<u>93</u> #<u>211</u> #<u>272</u> **Smart Pixels** Pixel readout ASIC with ML

- frontend filtering: discard low-p<sub>T</sub> tracks
- <u>feature extraction</u>: particle position+angle Mixture Density Model

![](_page_11_Picture_3.jpeg)

### **eFPGAs** Fully **reconfigurable logic** in ASIC design

• The pathway to put ML on-detector!

![](_page_11_Picture_6.jpeg)

#### BDT classifier in 28nm CMOS ASIC

![](_page_11_Figure_8.jpeg)

### #<u>11 #95 #247 #233 #93 #211 #272</u> **Smart Pixels Pixel readout ASIC with ML**

- frontend filtering: discard low-p<sub>T</sub> tracks
- feature extraction: particle position+angle Mixture Density Model

![](_page_12_Picture_3.jpeg)

### **Calorimeter data concentrator**

![](_page_12_Figure_5.jpeg)

![](_page_12_Picture_6.jpeg)

**On FPGA** 

**Transmit encoded data!** 

![](_page_12_Figure_12.jpeg)

### Intelligent back-ends

![](_page_13_Figure_1.jpeg)

![](_page_13_Picture_2.jpeg)

## Intelligent backends

- Trigger-less for future e<sup>+</sup>e<sup>-</sup> not guaranteed
- FPGA-based inference for improved triggering
  - in L1 trigger
  - as accelerators in HLT
- Possible through HEP experiment-agnostic tools!

#### 2024: Neural hardware triggers making decisions in LHC experiments!

#### CMS:

![](_page_14_Picture_10.jpeg)

### **ML** inference on **FPGA**

![](_page_15_Figure_1.jpeg)

Lower-level information like tracks and particles in hardware triggers has lead to increased usage of set-based and graph-based architectures, like 250 ns DeepSet flavour tagging!

#### <u>Object tagging for Phase-2 CMS</u>

![](_page_15_Picture_4.jpeg)

![](_page_15_Picture_5.jpeg)

![](_page_15_Figure_6.jpeg)

### **Real-time tracking BELLE-II**

![](_page_16_Figure_1.jpeg)

- polar angles of single particle track

#### arxiv:411.13596

•Already use neural track hardware trigger for vertex-reconstruction, and azimuthal and

![](_page_16_Picture_9.jpeg)

#### **Well-established:** MLPs, <u>CNNs</u>, <u>DeepSets</u>, GNNS, RNNS, SYNDOLICAL regression, (variational) autoencoders, BDTs, IsolationTrees

**Everything here Everything here** Experimental: is abnormal Transformers, large **Reconstruction error** distributed CNNs (ResNet, VGG)

Aultilayer Perceptron MLP ep Sets DS nteraction Network IN

NP?

#### **Quantised Interaction** Networks and Deep Sets in <160 ns

![](_page_17_Figure_4.jpeg)

#### P. Odagiu et al. 2024

![](_page_17_Figure_6.jpeg)

### **ESPP** proposal: AI R&D collaborations **EuCAIF** proposal for scalable, robust AI through cross-domain collaboration

- Al for Detector and Accelerator Control: Accelerator performance, calibration, system monitoring.
- Al for Detector Optimization: Differentiable programming, reinforcement learning to maximize detector performance
- Al for Event Reconstruction: Tracking, calorimetry, end-to-end foundation models.
- Al for Data Processing: Front-end electronics, trigger

DRD-7 and EuCAIF AI-RDs

![](_page_18_Picture_6.jpeg)

#### in close collaboration with the Fast Machine Learning Lab, CERN NGT,

![](_page_18_Picture_10.jpeg)

### Conclusion

- ML is essential to address unique data and processing challenges in HEP
  - Custom workflows and tools developed for extreme constraints
- AI/ML is ubiquitous in upcoming and future high-luminosity experiments
  - Intelligent processing near sensors needed to manage data from granular detectors
- Cross-experiment collaboration (e.g. DRDs, AI-RDs) is key to future success!