

JAVIER DUARTE ML4FE WORKSHOP U HAWAII, MAY 19, 2025



HL-LHC EVENT PROCESSING



Challenges:

Each collision produces O(10³) particles The detectors have O(10⁸) sensors Extreme data rates of O(100 TB/s)



SIMPLIFIED HL-LHC L1 TRIGGER MENU

- Single/double/triple muons/electror
- Photons
- Taus
- Hadronic
- Missing transverse energy
- "Cross" triggers (not shown)



J	Triaaer	Threshold [Ge\
hs ds set by ounds, olution @ e budget	1 μ	22
	2 µ	15,7
	3 μ	5, 3, 3
	1 e	36
	2 e	25, 12
	1γ	36
	2γ	22, 12
	1 т	150
	2 т	90, 90
	1 jet	180
	2 jet	112, 112
	H _T	450
	4 jet + H _T	75, 55, 40, 40, 4
	PT ^{miss}	200



WHAT COULD WE BE MISSING?

- How can we trigger on more complex low-energy hadronic signatures? Long-lived/displaced particles?
- What if we don't know exactly what to look for?
- What if our signatures require complex multivariate algorithms (e.g. b tagging)?
- How can we improve on our traditional (often slow) reconstruction algorithms?









WHAT MAKES THIS HARD?

- Reconstruct all events and reject 98% of them in O(10) µs



Algorithms have to be <1 μ s and process new events every 25 ns





CODESIGN

- Codesign: intrinsic development loop between ML design, training, and implementation
- Pruning
 - Maintain high performance while removing redundant operations
- Quantization
 - Reduce precision from 32-bit
 floating point to 8-bit fixed point, ...
- Parallelization
 - Balance parallelization (how fast)
 with resources needed (how costly)







QUANTIZATION-AWARE TRAINING

- Small NN benchmark correctly identifies 5 classes of jets 70-80% of the time
- Full performance with 6 bits instead of 14 bits
- Much smaller fraction of resources





PRUNING + QUANTIZATION-AWARE TRAINING

- Iterative pruning further reduces hardware complexity of a quantized model
 - Can remove 90% of parameters with almost no loss in performance!
- After pruning and quantization-aware training, can achieve 50 × reduction in bit operations compared to original 32-bit model



Bit operations (BOPs) definition: <u>arXiv:1804.10969</u>





HIGH-GRANULARITY QUANTIZATION

- High-granularity quantization (HGQ) further extends this paradigm by per-weight and per-activation granularity
- Demonstration on a CNN image classification (SVHN) benchmark



implementing gradient-based automatic bit width optimization at (up to)





NEURAL ARCHITECTURE CODESIGN

- Combine neural architecture search with hardware-aware codesign to find **Pareto-optimal front**
- [<u>arXiv:2402.01876</u>]







PROGRAMMING HARDWARE (FPGAS)

How do we implement an optimized AI model (e.g., in <u>OONNX</u>) for an FPGA?





High-Level Synthesis







HLS4ML

- <u>hls4ml</u> for scientists or
 ML experts to translate
 ML algorithms into RTL
 firmware
- NN model is parsed through the front ends into the IR, modified by the optimizer, and passed to the back ends, which create the output





LATEST FEATURES IN HLS4ML 1.1.0

- OONNX front end
- Automatic precision inference
- Support for HGQ models
- Symbolic regression, hardware-aware optimization API, and more...

```
config = hls4ml.utils.config_from_keras_model(model, default_precision='ap_fixed<16,6>', granularity='model')
config['LayerName'] = {
   # Infer all types of these layers
    'first_layer': {
        'Precision': 'auto',
   },
hls_model = hls4ml.converters.convert_from_keras_model(
   model, hls_config=config, io_type=io_type, output_dir=odir, backend=backend
```





APPLICATION: SENSOR DATA COMPRESSION ON ASIC

- area/power/radiation tolerance requirements
- new physics signatures)

8" hexagonal silicon module (1 out of ~27,000)





Autoencoder on the detector front-end for data compression: ~100 ns latency and

Reconfigurable ASIC to address: evolving LHC conditions (beam-related), detected performance (noise, dead channels), and updated performance metrics (resolutio





APPLICATION: JET TAGGING @ L1

- Upgraded HL-LHC level-1 track trigger information enables b-tagging with a **neural network** to improve the $HH \rightarrow 4b$ search
 - Input features for 10 particles within each jet: particle type, momentum, and vertex information







APPLICATION: ANOMALY DETECTION AT 40 MHZ

- Challenge: if new physics has an unexpected signature that doesn't align with existing triggers, precious signal events will be discarded
- Can we use unsupervised algorithms to detect non-SM-like anomalies?
 - Autoencoders (AEs): compress input to a smaller dimensional latent space then decompress and calculate difference
 - Variational autoencoders (VAEs): model the latent space as a probability distribution; possible to detect anomalies purely with latent space variables





Key observation: Can build an anomaly score from the latent space of VAE directly! No need to run decoder!

$$R_z = \sum_i \frac{\mu_i^2}{\sigma_i^2}$$



APPLICATION: ANOMALY DETECTION AT 40 MHZ

- AXOL1TL anomaly detection algorithm for the trigger based on a variational autoencoder
- Selects unique events relative to existing triggers
- Preference for high multiplicity events



<u>CMS-DP-2023-079</u> <u>CMS-DP-2024-059</u>

17



APPLICATION: SMART PIXELS

- ~300 µW and latency of 3.9 ns
- hls4ml + Catapult AI to ASIC workflow





Mach. Learn. Sci. Tech. 5 (2024) 3, 035047 arXiv:2406.14860

Locally customized neural network for p_T filtering of charged particles in sensor readout can enable data reduction by 55-75% with low power consumption of



SUMMARY AND OUTLOOK

- ML allows us to better reconstruct our data and save potentially overlooked data
- Codesign principles can enable ML on hardware with stringent constraints
- Community (<u>fastmachinelearning.org</u>, e-group <u>hls-</u> fml@cern.ch), Institute (a3d3.ai) and CERN Project (<u>NextGen Triggers</u>) developing open-source tools and techniques to enable this
 - <u>hls4ml</u>: expanding open-source toolkit for translating ML into hardware aimed at trigger applications and more...
- Applications range from jet tagging to tracking, and more!
 - Enhance future particle physics program









JAVIER DUARTE ML4FE WORKSHOP U HAWAII, MAY 19, 2025





PARALLEL OR STREAMING INPUTS

 hls4ml supports two modes for parallel or streaming inputs



(a) io_parallel

(b) io_stream

