

# Al at the Extreme Edge: Nanosecond inference for New **Physics Discovery at the LHC**

Thea Klaeboe Aarrestad (ETH Zürich)

# FIRZURICH





Improved Precision in  $Vh(\rightarrow bb)$  via Boosted Decision Trees

Philipp Englert (Jul 30, 2024) e-Print: 2407.21239 [hep-ph]
Def Cite 🗟 claim
Accelerating template generation in resonant anomaly detection searches with optimal transport Matthew Leigh (Geneva U.), Debajyoti Sengupta (Geneva U.), Benjamin Nachman (LBL, Berkeley), Tobias Golling (Geneva U.) (Jul 29, 2024) e-Print: 2407 11: The Company of the Storner Of pdf Cite Cite Calim
Accuracy version and imposted person of the ATLAS detector ATLAS Collaboration - Georges Add (Marseille, CPPM) et al. (Jul 29, 2024) e-Print: 2407.20127 [hep-ex] pdf [attention]] and [hep-
Probing Charm Yukawa through <i>ch</i> Associated Production at the Hadron Collider Nuoyu Dong Construction at Collider e-Print: 2407.19797 [hep-ph]
Universal New Physics Latent Space Anna Hallin (Hamburg U.), Gregor Kasieczka (Hamburg U.), Sabine Kraml (LPSC, Grenoble), André Lessa (ABC Federal U.), Louis Moureaux (Ha e-Print: 2407.20315 [hep-ph] pdf 🖃 cite 🐻 claim
Constraints on the fuzzy dark matter mass window from high-redshift observables Hovav Lazare (Ben Gurion U. of Negev), Jordan Flitter (Ben Gurion U. of Negev), Ely D. Kovetz (Ben Gurion U. of Negev) (Jul 28, 2024) e-Print: 2407.19549 [astro-ph.CO]
Comparison of Geometrical Layouts for Next-Generation Large-volume Cherenkov Neutrino Telescopes Tong Zhu (UC, Berkeley), Miaochen Jin (Harvard U.), Carlos A. Argüelles (Harvard U.) (Jul 26, 2024) e-Print: 2407.19010 [physics.ins-det]
Generative Flow Networks in Covariant Loop Quantum Gravity Joseph Bunao (Algoma U. and Western Ontario U.), Pietropaolo Frisoni (Algoma U. and Western Ontario U.), Athanasios Kogios (Perimeter Inst. e-Print: 2407.19036 [gr-qc]
Can we learn from matter creation to solve the H <sub>0</sub> tension problem? Emilio Elizalde (Barcelona, IEEC), Martiros Khurshudyan (Barcelona, IEEC), Sergei D. Odintsov (Barcelona, IEEC) (Jul 26, 2024) Published in: <i>Eur.Phys.J.C</i> 84 (2024) 8, 782 · e-Print: 2407.20285 [gr-qc]
Anomaly Detection Based on Machine Learning for the CMS Electromagnetic Calorimeter Online Data Quality Monit CMS Collaboration • Abhirami Harilal (Carnegie Mellon U.) et al. (Jul 25, 2024) Contribution to: CALOR2024 • e-Print: 2407.20278 [physics.ins-det]
The Observation of a 95 GeV Scalar at Future Electron-Positron Colliders Karabo Mosala (U. Witwatersrand, Johannesburg, Sch. Phys. and iThemba LABS), Anza-Tshilidzi Mulaudzi (U. Witwatersrand, Johannesburg, S Kumar (IISER, Mohali) et al. (Jul 23, 2024)

Enhancing High-Energy Particle Physics Collision Analysis through Graph Data Attribution Techniques

A. Verdone (Rome U.), A. Devoto (Rome U.), C. Sebastiani (Liverpool U.), J. Carmignani (Liverpool U.), M. D'Onofrio et al. (Jul 20, 2024)

Benno Käch (DESY), Isabell Melzer-Pellmann (DESY), Dirk Krücker (DESY) (Aug 9, 2024) e-Print: 2407.14298 [astro-ph.CO] e-Print: 2408.04997 [hep-ex] 🗟 reference search 🔁 Qeitations claim A pdf → cite cite Unruh-DeWitt Quantum Computing: Realizing Quantum Shannon Theory With Quantum Fields Extracting Signal Electron Trajectories in the COMET Phase-I Cylindrical Drift Chamber Using Deep Learning Fumihiro Kaneko, Yoshitaka Kuno (Osaka U.), Joe Sato (Yokohama Natl. U.), Ikuya Sato (Saitama U.), Dorian Pieters (Osaka U.) et al. (Aug 8, 2024) e-Print: 2408.04795 [hep-ex] C reference\_search ① 0 citations 🚯 pdf 🛛 🖃 cite 📑 clair

### Date of

AI for Nuclear Physics: the EXCLAIM project

Simonetta Liuti, Douglas Adams, Marie Boër, Gia-Wei Chern, Marija Cuic et al. (Jul 31, 2024) e-Print: 2408.00163 [hep-ph]

999

e-Print: 2407.16806 [hep-ph]

🔓 pdf 🛛 🖃 cite 🛛 🔂 claim

Deep learning-driven likelihood-free parameter inference for 21-cm forest observations Tian-Yang Sun, Yue Shao, Yichao Li, Yidong Xu, Xin Zhang (Jul 19, 2024)

### Selected Papers: 30 Total Papers: 30 Year: 2024

based Graph Attention Network for Particle Track Reconstruction y Chan (LBL, Berkeley), Loic Delabrouille (LBL, Berkeley and Unlisted, FR), Brandon Wang (LBL, I -an]

ed Track Finding in the LHCb Vertex Detector otis L Giasemis (LPNHE, Paris and LIP6, Paris), Nabil Garroum (LPNHE, Paris), Vladimir Vava nst. Part. Phys.); Xiao-Hua Li (U. South China, Hengyang and Jishou U.), Kai Zhou (Shenzhen

tures of Heavy Flavor-Violating Scalars at the LHC with Parametrized Neural Ne Almeida, Alex G. Dias, Diego S.V. Gonçalves (Jul 16, 2024)

L, Berkeley and Sao Paulo, IFT), Benjamin Nachman (LBNL, Berkeley and Sao Paulo, IFT and UC,

s-enriched semi-visible jets

essandro Russo (Stanford U., Phys. Dept.), Emre Sitti (Zurich, ETH), Annapaola de Cosa (Zurich

raction Particle Transformer

Kun Wang (Shanghai U. Sci. Tech.), Jingya Zhu (Henan U.) (Jul 11, 2024)

Imaging Cherenkov Detectors with Swin Transformers and Normalizing Flow M II.), James Giroux (William-Mary Coll.), Justin Stevens (William-Mary Coll.) (Jul 10, 2024)

orgi-Machacek scenario at the Large Hadron Collider Rituparna Ghosh (IISER, Kolkata), Biswarup Mukhopadhyaya (Saha Inst.) (Jul 10, 2024)

dard-Model Event Discovery at the LHC with a Novel Quantum Autoencoder

Callum Duffy (University Coll. London), Mohammad Hassanshah (University Coll. London), Marcin Jastrzebski (University Coll. London)

2024

; EMMI and Heidelberg, Max Planck Inst.), S.R. Stroberg (Notre Dame U.), J.D. Holt (TRIUMF and









T1037 / 6vr4 90.7 GDT (RNA polymerase domain) **T1049 / 6y4f** 93.3 GDT (adhesin tip)

Experimental result

Computational prediction

sequence—the structure prediction component of the 'protein folding problem'<sup>8</sup>—has been an important open research problem for more than 50 years<sup>9</sup>. Despite recent









T1037 / 6vr4 90.7 GDT (RNA polymerase domain) **T1049 / 6y4f** 93.3 GDT (adhesin tip)

Experimental result

Computational prediction

sequence—the structure prediction component of the 'protein folding problem'<sup>8</sup>—has been an important open research problem for more than 50 years<sup>9</sup>. Despite recent





### PFN P-CNN ParticleNet

ParT

### 100 million jets for training

ParT (plain)



Accuracy	# params
0.772	86.1 k
0.809	354 k
0.844	370 k
0.861	2.14 M
0.849	2.13 M

"Particle Transformer For Jet Tagging" H. Qu, C. Li, S. Qian



## GPT-3



### 175,000,000,000 (0.16% of neurons in your brain)

# GPT-4 (MoE)



## **1,800,000,000,000** (1.6% of neurons in your brain)

### semianalysis 2023



### <u> Train (GPT-4):</u>

- 2.15<sup>25</sup> floating point operations
- ~25,000 A100 GPUs
- 90-100 days
- \$63 million
- Trained on 13 trillion tokens

### semianalysis 2023



### <u>Train (GPT-4):</u>

- 2.15<sup>25</sup> floating point operations
- ~25,000 A100 GPUs
- 90-100 days
- \$63 million
- Trained on 13 trillion tokens



I'm looking for a really good Svensken, Dansken and Nordmannen joke. Could you help me?



Certainly! Here's a classic "Svensken, Dansken, och Nordmannen" joke:

Svensken, Dansken, and Nordmannen were stranded on a desert island. They found a bottle, rubbed it, and out popped a genie.

The genie said, "I'll grant each of you one wish."

Svensken said, "I wish I was back in Stockholm, sitting at my favorite café with a cold beer."

Poof-he was gone.

Dansken said, "I wish I was back in Copenhagen, enjoying a hot dog and a Carlsberg."

Poof-he was gone.

Nordmannen looked around and aid, "I'm feeling lonely. I wish Svensken and Dansken were back here with me!"

# Inference (GPT-4): Multiple clusters of 128 GPUs Model carefully mapped onto hardware



### FP16 vs FP32





GPT-3

175,000,000,000

Range 10 <sup>-38</sup> - 10 <sup>38</sup>	Accuracy 0.000006%	~700 GB $\rightarrow$ (175B × 4B/par) $\rightarrow$ O(10 <sup>1</sup> )> single GPU me
5x10 <sup>-5</sup> - 6x10 <sup>4</sup>	0.05%	→ $\sim 350 \text{ GB}$ (175B × 2B/par) → 11 NVIDIA V100 (\$10 0







AI & Memory Wall

### **AI and Memory Wall**



### Kaplan et al. (2020)



### Kaplan et al. (2020)

Test loss

### 0(1) ms





### ASIC/GPU



0(1) ns

ASIC

FPGA

High-throughput

(Data-efficient)

### Low power

On-device

# EFFICIENT ML ALGORITHMS

Low-latency

### **Resources:** 128 interconnected GPUs Latency: 10 seconds

ChatGPT 40  $\, \smallsetminus \,$ 

\$

A JSJ



I'm looking for a really good Svensken, Dansken and Nordmannen joke. Could you help me?



### **<u>Resources:</u>** 128 interconnected GPUs <u>Latency:</u> 10 seconds

ChatGPT 40  $\vee$ 

\$



I'm looking for a really good Svensken, Dansken and Nordmannen joke. Could you help me?

# Resources:1 single chipsLatency:1 millionth of a secondThroughput:5% of internet traffic



### CMS Experiment at the LHC, CERN

Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST) Run / Event: 1510767 1405388









ATLAS ALICE

2,500 bunches of 100 billion protons, 11,000 rotations per second











### CMS Experiment at the LHC, CERN

Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST) Run / Event: 15107671405388





### CMS Experiment at the LHC, CERN

Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST) Run / Event: 15107671405388

## 1 billion collisions /s MB of data / collision PB of data / s.



## Quarks u C



t

### Leptons





### Masses span 9 orders of magnitude!











### Leptons













	L⁺
	Ľ
	• +
	L`
	L

### cmsexperiment.web.cern.ch



We had to collide billions of protons, only around 10 signal events were needed to claim discovery!

	L⁺
	Ľ
	• +
	L`
	L

### cmsexperiment.web.cern.ch



We had to collide billions of protons, only around 10 signal events were needed to claim discovery!

	L⁺
	Ľ
	• +
	L`
	L

# The Standard Model

$$\begin{split} & -\frac{1}{2}\partial_{\nu}g_{\mu}^{\mu}\partial_{\nu}g_{\mu}^{\mu} - g_{\lambda}f^{abc}\partial_{\mu}g_{\nu}^{\mu}g_{\mu}^{\mu}g_{\nu}^{\nu} - \frac{1}{4}g_{s}^{\mu}f^{abc}f^{abc}f^{abc}g_{\mu}^{\mu}G_{\nu}^{\mu} + \frac{1}{2}ig_{s}^{\mu}(\tilde{q}^{\sigma}\gamma^{\mu}q_{j}^{\sigma})g_{\mu}^{\mu} + \tilde{G}^{\alpha}\partial^{2}G^{\alpha} + g_{s}f^{abc}\partial_{\mu}G^{\alpha}G^{b}g_{\mu}^{\mu} - \partial_{\nu}W_{\mu}^{+}\partial_{\nu}W_{\mu}^{-} - M^{2}W_{\mu}^{+}W_{\mu}^{-} - \frac{1}{2}\partial_{\mu}Q^{0}\partial_{\mu}Q^{0} - \frac{1}{2}e_{w}^{-}MQ^{0}Q^{0} - \partial_{h}[\frac{2M^{2}}{2} + \frac{2M^{2}}{2}M^{2}H^{2} - \frac{1}{2}\partial_{\mu}A_{\nu}\partial_{\mu}A_{\nu} - \frac{1}{2}\partial_{\mu}H\partial_{\mu}H - \frac{1}{2}m_{h}^{h}H^{2} - \partial_{\mu}\phi^{+}\partial_{\mu}\phi^{-} - M^{2}\phi^{+}\phi^{-} - \frac{1}{2}\partial_{\mu}\phi^{0}\partial_{\mu}\phi^{0} - \frac{1}{2}e_{w}^{-}M\phi^{0}Q^{0} - \partial_{h}[\frac{2M^{2}}{2} + \frac{2M^{4}}{2}H^{2} - \partial_{\mu}\phi^{+}\partial_{\mu}\phi^{+} - M^{2}\phi^{+}\phi^{-} - \frac{1}{2}\partial_{\mu}\phi^{0}\partial_{\mu}\phi^{0} - \frac{1}{2}e_{w}^{-}M\phi^{0}Q^{0} - \partial_{h}[\frac{2M^{2}}{2} + \frac{2M^{4}}{2}\partial_{\mu}H^{-} - W_{\nu}^{+}W_{\nu}^{-} - Z_{\nu}^{\mu}(W_{\mu}^{+}W_{\nu}^{-} - W_{\nu}^{-}\partial_{\nu}W_{\mu}^{+}) - Z_{\nu}^{\mu}(W_{\mu}^{+}\partial_{\nu}W_{\mu}^{-} - W_{\nu}^{-}\partial_{\nu}W_{\mu}^{+}) - A_{\nu}(W_{\mu}^{+}\partial_{\nu}W_{\mu}^{-} - W_{\nu}^{-}\partial_{\nu}W_{\mu}^{+}) + A_{\mu}(W_{\nu}^{+}\partial_{\nu}W_{\mu}^{-} - W_{\nu}^{-}\partial_{\nu}W_{\mu}^{+}) - A_{\nu}(W_{\mu}^{+}\partial_{\nu}W_{\mu}^{-} - W_{\nu}^{-}\partial_{\nu}W_{\mu}^{+}) + A_{\mu}(W_{\nu}^{+}\partial_{\nu}W_{\mu}^{-} - W_{\nu}^{-}\partial_{\nu}W_{\mu}^{+}) - \frac{1}{2}g^{2}W_{\mu}^{+}W_{\nu}^{-}W_{\nu}^{-}W_{\nu}^{-} + \frac{1}{2}g^{2}W_{\mu}^{+}W_{\nu}^{-}W_{\nu}^{-}W_{\nu}^{-} + \frac{1}{2}g^{2}W_{\mu}^{+}W_{\nu}^{-}W_{\nu}^{-} + M^{2}\partial_{\nu}W_{\nu}^{-} - W_{\nu}^{-}W_{\nu}^{-}W_{\nu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}(A^{\mu})^{\mu}W_{\nu}^{-} - A_{\mu}(W_{\mu}^{+}\partial_{\nu}W_{\nu}^{-} - W_{\nu}^{+}W_{\nu}^{-}) + g^{2}S_{\mu}^{\mu}(A^{\mu})^{\mu}W_{\nu}^{-} - A_{\mu}(W_{\mu}^{+}W_{\nu}^{-}) + g^{2}S_{\mu}^{\mu}(A^{\mu})^{\mu}W_{\nu}^{-} + 2g^{\mu}(M_{\nu}^{+}W_{\nu}^{-}) + g^{2}S_{\mu}^{\mu}(A^{\mu})^{\mu}W_{\nu}^{-} - A_{\mu}A^{\mu}W_{\nu}^{+}W_{\nu}^{-}) + g^{2}S_{\mu}^{\mu}(M_{\mu}^{+}W_{\nu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}(M_{\mu}^{+}W_{\nu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}(M_{\mu}^{+}W_{\nu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}(M_{\mu}^{+}W_{\nu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}(W_{\mu}^{+}W_{\nu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}(W_{\mu}^{+}W_{\mu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}D_{\mu}^{\mu}(W_{\mu}^{+}W_{\mu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{\mu}D_{\mu}^{\mu}(W_{\mu}^{+}W_{\mu}^{-}) + \frac{1}{2}g^{2}S_{\mu}^{$$


#### CMS Experiment at the LHC, CERN

Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST) Run / Event: 15107671405388

## 1 billion collisions /s MB of data / collision PB of data / s.





Higgs produced ~1 in a billion collisions!

Saving all collisions not useful (even if we could)!



#### "Probability" of producing "anything"

"Probability" of producing a Higgs

### Geneva Lake

### CMS

#### Software rate reduction (GPU+CPU)

LHC

LHCb

2 step rate reduction (hardware+software)

# Geneva ATLAS ALICE

#### 2 step rate reduction (hardware+software)

#### Continous read-out (CPU+GPU)







# Geneva A Data temporarily stored and a local state of in detector electronics for 4 $\mu$ s

LHC







#### 5% internet traffic to L1 [63 Tb/s]

A

L1 trigger: ~1000 Xilinx/AMD Ultrascale FPGAs

Decide which event to keep within ~4 µs latency

Discard >99% of collisions!







#### DATA 99.72% of events rejected! 110 thousand events/

CMS

6 13 Mg 1 1 1 1 1 1

LHCb

Geneva Lake



#### DATA 99.72% of events rejected! 110 thousand events/

CMS

A THE T

LHCb

Geneva Lake



High Level Trigger: 25'600 CPUs / 400 GPUs Latency: 3-400 ms

TATAT

9 M Z 8 8 8 4 9 **0** 29 € € 5

•

.

**Reject further 99%!** 

-HCb







## ATLAS

DATA 99.9975% of events rejected! 1000 events/second ~5 GB/s





## ATLAS

TIFR 0:∞

DATA 99.9975% of events rejected! 1000 events/second ~5 GB/s







### ATLAS

TIFR 0:∞

DATA 99.9975% of events rejected! 1000 events/second ~5 GB/s





Geneva

#### To make sure we select "the right" 0.0025%, algorithms must be • Fast (get more data through) • Accurate (select the right data)



New Physics is produced less than 1 in a trillion (if at all)

Need <u>more</u> data!



#### "Probability" of producing "anything"

#### **New Physics?**



### High Luminosity LHC

#### New Physics is produced 1 in a trillion

Need more collisions to observe rare processes

#### High Luminosity LHC

- ×10 data size
- ×3 collisions/s









MAJOR UPGRADE

2026 - 2028

2029 - 2038

HL-LHC



#### $\frac{1}{100} ructure \rightarrow pile-up of \sim 60 events/x-ing$ ts/x-ing)



### High Luminosity LHC

200 vertices (average 140)



#### Maintain physics acceptance $\rightarrow$ better detectors

#### CMS High Granularity (endcap) calorimeter • X20 times more readout channels (6.5 million!!)

More collisions More readout channels





### Our current algorithms won't be able to cope with HL-LHC data rates! Need innovation!

Can we use modern Machine Learning to be faster more accurate and do more?

#### **GPU** inference

Ge

#### HLT trigger: Latency 0(100) ms

#### Fast inference on specialised hardware

#### **ASIC inference**

**Detector:** 40 MHz ~Pb/s

#### **FPGA** inference

LHCb

Level-1 trigger: Latency O(1) µs

ATLAS

Geneva

VIRTEX"5 VIRTEX"5 XC5VLX30" FTEG76E0005 D1030908A

ALICE



Complete re-design of Level-1



#### Complete re-design of Level-1

Charged particle tracks (6.4 Tb/s, 200 FPGAs)



ics y, no tracking information





#### Complete re-design of Level-1

- Charged particle tracks
- Particle Flow (40 FPGAs)





#### Complete re-design of Level-1

- Charged particle tracks
- Particle Flow
- HGCal (4 Tb/s, 200 FPGAs)





#### Complete re-design of Level-1

- Charged particle tracks
- Particle Flow
- HGCal

## Input data • 2 Tb/s → 63 Tb/s

#### Latency

•4 µs → **12 µs** 

Extremely high data complexity,

#### Extremely little time









#### Nanosecond ML inference on FPGAs! ~40 billion inferences/s during HL-LHC L1 trigger Hardware-based, implemented in sustom-built electronics My, no tracking information **Journey to HL-LHC** run: $7 \times 10^{33}$ , PU = 30, E = 7 TeV, 50 nsec bunch spacing Detectors 40 MHz Detectors **TLAS, CMS operating:** 40 MHz Front end Front end pipelines ccept ≤ 100 kHz, pipelines L1 output: 75 kHz 76 Comput: 100 kHz ; 100 kHz Readout Readout MU $1cy \leq 2.5$ (AT), 4 µsec (CM) L1 trigger decision Lvl-2 buffers buffers ~3 kHz in ~2.5 (4) µs for Switching Switching **I**rigger Accept $\leq 1 \text{ kHz}$ network network **ATLAS (CMS)** ept/reject LAS & CMS will be: \* Processor Processor HLT Lvl-3 farms

HLT output:

~1 kHz

**5** v 1034

Simulated event display with average pileup of 140





## Nanosecond ML inference on FPGAs! ~40 billion inferences/s during HL-LHC

#### L1 trigger

### **Journey to HL-LHC**

Hardware-based, implemented in sustom-built electronics My, no tracking information  $\sigma_{in}^{pp} = 69.2 \ mb$ Conifer <u>hls</u>4ml Mean number of interactions per crossing Mean number of interactions per crossing  $7 \times 10^{33}$ , PU = 30, E = 7 TeV, 50 nsec bunch spacing Detectors 40 MHz TLAS, Civis operating: 40 MHz Detectors Front end Front end pipelines ccept  $\leq$  100 kHz, pipelines L1 output: 75 kHz 75 Comput: 100 kHz ; 100 kHz Readout Readout MU  $1Cy \leq 2.5$  (AT), 4 µsec (CM) L1 trigger decision Lvl-2 buffers buffers ~3 kHz in ~2.5 (4) µs for Switching Switching **I**rigger Accept  $\leq 1 \text{ kHz}$ network network **ATLAS (CMS)** ept/reject LAS & CMS will be: \* Processor Processor 17 HLT Lvl-3 farms farms 7100 HZ 12.5 µs HLT output: ~1 kHz **5 v 1034** 

Simulated event display with average pileup of 140

<µ> = 32





## Why FPGAs?

## Why FPGAs?

Latency (resource parallelism)



## Why FPGAs?

Throughput (pipeline parallelism)



### pipeline parallelism



Latency (resource parallelism) Can work on different parts of problem, different data simultaneously Latency strictly limited by detector frontend buffer

High bandwidth (pipeline parallelism) Phase 2 L1T processes 5% of total internet traffic

Latency deterministic CPU/GPU processing randomness, FPGAs repeatable and predictable latency

**TRACK FINDER TMUX=18** RS = 9 (phi) **FPGAs = 162** 




#### Digital signal processors (DSPs) 0(5,000) units



16\*64 +64\*32 +32\*32 +32\*5 = 4,256 synapses 

stored in BRAMs

#### Digital signal processors (DSPs) 0(5,000) units





Firmware block



-0

#### **KERAS / PyTorch / ONNX**









pip install hls4ml pip install conifer https://github.com/fastmachinelearning/hls4ml https://fastmachinelearning.org/hls4ml/

VIRTEX









#### Model (quantized/pruned)

#### Quantized:













#### **ASICs**













Prediction









Data flow architecture Tailored hardware for a model Each layer is separate compute unit 

- Stay on-chip
- "Decisions are design time"

Prediction

# 





## Ideally



## Reality



### Ideally



- Quantization
- Pruning
- Parallelisation
- Knowledge distillation

### Reality

# Quantization



## Floating point 32: 4B numbers in [-3.4e38, +3.4e38]



# Quantization



### **Quantising:** int8 2<sup>8</sup>=256 numbers in [-128,127]

 $x_q = Clip(Round(\frac{x_f}{scale}))$ 



Precision	Approx. Peak
1b	64 000
<b>4b</b>	16 000
<b>8b</b>	4 000
32b	300

**Trillions** of quantized operations per second

AMD UltraScale+ MPSoC ZU19EG (conservative estimates)



Weights can stay entirely on-chip

## Weights Layer 1



### Weights Layer 2



### Fixed point

## Weights Layer 1



### Weights Layer 2



## **Fixed point** 0101.1011101010

width

integer

fractional

## Weights Layer 1



### Weights Layer 2





## hls4 + Google Quantization-aware training



#### Nature Machine Intelligence 3 (2021)

#### Forward pass →





## Back propagation







## hls4 + Google Quantization-aware training



#### Nature Machine Intelligence 3 (2021)

#### Forward pass →





## Back propagation







## Jogle Quantization-aware training



#### Nature Machine Intelligence 3 (2021)

#### Forward pass



```
x = QDense(32, kernel_initializer=tf.keras.initializers.HeUniform(),
               kernel_quantizer='quantized_bits(8,3,1, alpha=1.0)',
               bias_quantizer='quantized_bits(8,3,1, alpha=1.0)')(inputArray)
x = QDense(16, kernel_initializer=tf.keras.initializers.HeUniform(),
               kernel_quantizer='quantized_bits(8,3,1, alpha=1.0)',
               bias_quantizer='quantized_bits(8,3,1, alpha=1.0)')(x)
encoder = QDense(latent_dim, kernel_initializer=tf.keras.initializers.HeUniform(),
               kernel_quantizer='quantized_bits(16,6,1, alpha=1.0)',
               bias_quantizer='quantized_bits(16,6,1, alpha=1.0)', name='bottleneck')(x)
```





1.0

## Estimating energy and size

Some layers more accommodating for aggressive quantization, others require expensive arithmetic

heterogeneous quantization

## Estimating energy and size

Some layers more accommodating for aggressive quantization, others require expensive arithmetic

heterogeneous quantization

For edge inference, need best possible quantization configuration for

- Highest accuracy  $\uparrow$ ...
- ullet ... and lowest resource consumption igslash

ightarrow hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

## Estimating energy and size

Some layers more accommodating for aggressive quantization, others require expensive arithmetic

heterogeneous quantization

For edge inference, need best possible quantization configuration for

- Highest accuracy  $\uparrow$ ...
- ... and lowest resource consumption  $\downarrow$

 $\rightarrow$  hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

QTools: Estimate QKeras model bit and energy consumption, assuming 45 nm Horowitz process

Model Accuracy [%]			Per-layer energy consumption [pJ]							Total en	ergy $[\mu J]$	Total bits
		Dense	ReLU	Dense	ReLU	Dense	ReLU	Dense	Softmax			
BF	74.4	1735	53	3240	27	1630	27	281	11		0.00700	61446
$\mathbf{Q6}$	74.8	794	23	1120	11	562	11	99	11		0.00263	26334
Forgiving Factor = $1 + \Delta_{accuracy} \times \log_{rate}(S \times \frac{Cost_{ref}}{Cost_{trial}})$												
<u>Maximize accuracy + minimizing cost in hyper parameter scan over quantizers:</u>												

#### **AutoQKeras**



As optimization progresses, best model accuracy/size trade-off is found!



#### Example: One convolutional layer





#### Nature Machine Intelligence 3 (2021)







# Pruning

before pruning





# Pruning

before pruning







Original image

# Shot Noise Gaussian Noise

#### Less accurate



Hooker et al. (2021)

From Brian Bartoldson

Less robust to noise













#### There exists a optimal network WITHIN each network (lottery ticket) Uncover it through pruning!



Diffenderfer, Bartoldson, et al. (2021)







Can we have the best of both worlds?

## Knowledge Distillation



#### Inference








is cat
is dog













### Soft labels contain information!!





### Train student to learn both true and predicted (teacher) labels!

 $L_{total} = \beta \times L_{Distillation} + \alpha \times L_{student}$ 

Student learns subtle learned features from teacher!



Go to NeurIPS 2022 Track Datasets and Benchmarks h...

### Why do tree-based models still outperform deep learning on typical tabular data? PDF

### Leo Grinsztajn, Edouard Oyallon, Gael

### Varoquaux

06 Jun 2022 (modified: 16 Jan 2023) NeurIPS 2022 Datasets and Readers: 🚱 Everyone Benchmarks Show Bibtex Show Revisions

**Abstract:** While deep learning has enabled tremendous progress on text and image datasets, its superiority on tabular data is not clear. We contribute extensive benchmarks of standard and novel deep learning methods as well as tree-based models such as XGBoost and Random Forests, across a large number of datasets and hyperparameter combinations. We define a standard set of 45 datasets from varied domains with clear characteristics of tabular data and a benchmarking methodology accounting for both fitting models and finding good hyperparameters. Results show that tree-based models remain state-ofthe art on medium-sized data (~ 10K samples) even without accounting





#### Computer Science > Machine Learning

[Submitted on 11 Oct 2022 (v1), last revised 25 Oct 2022 (this version, v3)]

### Neural Networks are Decision Trees

#### Caglar Aytekin

In this manuscript, we show that any neural network with any activation function can be represented as a decision tree. The representation is equivalence and not an approximation, thus keeping the accuracy of the neural network exactly as is. We believe that this work provides better understanding of neural networks and paves the way to tackle their black-box nature. We share equivalent trees of some neural networks and show that besides providing interpretability, tree representation can also achieve some computational advantages for small networks. The analysis holds both for fully connected and convolutional networks, which may or may not also include skip connections and/or normalizations.

Subjects: Machine Learning (cs.LG)

arXiv:2210.05189 [cs.LG] Cite as: (or arXiv:2210.05189v3 [cs.LG] for this version) https://doi.org/10.48550/arXiv.2210.05189 🚯

#### Submission history

From: Çağlar Aytekin [view email] [v1] Tue, 11 Oct 2022 06:49:51 UTC (216 KB) [v2] Mon, 17 Oct 2022 15:18:14 UTC (224 KB) [v3] Tue, 25 Oct 2022 17:32:33 UTC (240 KB)









hyperparameters. Results show that tree-based models remain state-ofthe art on modium-sized data (~ 10K samples) even without accounting

<u>https://arxiv.org/abs/2210.05189</u>



%VU9P	Accuracy	Latency	DSP	LUT
qDNN	75.6%	40 ns	22 (~0%)	1%
BDT	74.9%	5 ns	-	0.5%





### Quantised input data

Floating point model

Compressed model (Quantised + Pruned)



### Firmware design





es





















### Variational Autoencoder

<u>ECON-T, D. Noonan</u>







### <u>ECON-T, D. Noonan</u>

#### AEs for compression also at LHCb!







<u>ECON-T, D. Noonan</u>







ECON-T, D. Noonan







ECON-T, D. Noonan





# Nanosecond ML inference on FPGAs! 40 billion inferences/s during HL-LHC





### HEP developed libraries for fast ML on FPGAs

### ML for reconstruction



### ML for reconstruction





### ML for reconstruction



On FPGA: 3.5 µs to cluster energy deposits into disentangled showers from individual particles

#### <u>EPJC Vol 79 608 (2019)</u>



## ML for tracking

In HL-LHC, will need to do track finding at L1

• O(1000) hits, O(100) tracks, 40 MHz rate, ~5 µs latency

Graph Neural Networks for fast charged particle tracking

Throughput-optimized for L1 applications, resource-optimised for co-processing



#### DOI:10.3389/fdata.2022.828666



Latency [cycles]	ll [cycles]	DSP [%]	LUT [%]	FF [%]	BR
59 295 ns	1	99.9	66.0	11.7	
79 395 ns	28	56.6	17.6	3.9	





## Fast jet tagging



Sets: Information is only assigned to individual nodes.



Graphs: Information is assigned to edges, i.e., pairs of nodes.

#### cds.cern.ch/record/2814728/



## Fast jet tagging



Sets: Information is only assigned to individual nodes.

**Graphs:** Information is assigned to edges, i.e., pairs of nodes.



14

12

10

8

6

1/ FPR

#### cds.cern.ch/record/2814728/









### The scientific method



# Form and test hypothesis



### Searches at LHC



### Searches at LHC (almost) always start with by

- assuming Standard Model
- and some signal hypothesis

#### No longer learn from observation

• Blind analysis only way we perform searches

## **Anomaly detection for New Physics searches**

LEARN THIS FROM DATA



### LOOK FOR ANYTING THAT DOESNT LOOK **LIKE THIS**



# **Anomaly Detection triggers**



Trigger threshold

Energy (GeV)

### Level-1 rejects >99% of events! Is there a smarter way to select?

# **Anomaly Detection triggers**



Trigger threshold

Energy (GeV)

- - LOST DATA SELECTED DATA - - POSSIBLE NP SIGNAL

### **Everything here** is normal

**Everything here** is abnormal

AD threshold



## Types of anomaly detection

### Outlier detection

Find (non-resonant) out-of-distribution datapoints



## Detecting overdensities

Find (resonant) overdensities in distributions



# Types of anomaly detection

### Outlier detection



#### Non-resonant, tail of distributions

- Often (variational) auto-encoders
- Useful for triggering!

#### Caveats

- What's a good metric for optimisation?
- How to use selected events in analysis?

## Detecting overdensities



#### Resonant, similar to a bump hunt

- Density estimation methods
- Useful for offline analysis

#### Caveats

Relies on a definition of "sideband" and a sizeable signal




Compressed representation of x. Latent space  $\Re^k$ , k < m×n prevents memorisation of input, must learn





 $\mathscr{L}(\mathbf{x}, \hat{\mathbf{x}})$  is Mean Squared Error $(\mathbf{x}, \hat{\mathbf{x}})$ , "high error events" proxy for "degree of abnormality"





 $n \times m$ 



#### SciPost Physics

other component contain ugg the invariant mass  $(k_{\mu,i}) =$  $k_{2,1}$ *g*̃ (400 GeV)  $\mathcal{K}_{0,jk}$ the left panel of Fig<sub>i</sub> 1 we use N = 40 constituents, after checking t 20 does not maks for stable difference." For jets with fewer of ill the entries remaining in  $\tilde{t}$  he soft regime with zeros. Y ove all information from the jet-level kinematics we boost all 4-mo ' of the fat jet. This also improves the performance of our net -ization jet algorithans avocatorad dulinear compitations of the sout  $\circ$ nRteixASij39efiningbaneembinationERSBRILLarge[26].odfoarchitec ver immediately after the LOLA contains of maked at a er after LoLA and the last layer, the autoencoder netw  $\overline{k_{\mu,i}} \xrightarrow{\sim} k_{\mu,j} = k_{\mu,i} C_{ij}$ with C = $C_{N,N+1}$ 

Ve allow  $f_{OT} \stackrel{CoLa}{M} = 10$  trainable linear combinations. These combined 4-vectors of on on the hadronically decaying massive particles. In the original LOLA appr the momenta  $\tilde{k}_j$  onto observable Lorentz scalars and related observables [13] napping is not easily invertible we do not use it for the autoencoder. Instead, vectors by another component containing the invariant mass,









AD threshold



#### ....in 50 nanoseconds! Currently taking data in CMS!

## **Continual learning**



Many reasons for changing conditions

- Detector position slightly changes
- Radiation damage

#### Radiation damage of silicon detector



## Continual learning



Many reasons for changing conditions

- Detector position slightly changes
- Radiation damage

Continual learning to the aid for self-supervised training?

• Avoid re-training on TBs of data, adapt to gradual changes!





.....

## Real-time ML in other experiments



Taking plasma accelerators to market



<u>F. Capel et al.</u>



#### Real-time ML in other experiments





<u>F. Capel et al.</u>



## Real-time ML in other experiments

#### Signals and backgrounds



F. Capel et al.













#### MULTI-MESSENGER ASTRONOMY (MMA)

Neutinos





## Real-time gravitational wave detection



120

- source location
- properties

Crucial to facilitate multi-messenger follow-up of sources !

Using gravitational waves, we can detect neutron star merger (kilo novae) before it occurs! • Use real time ML for event cleaning, parameter estimation and anomaly detection!



Rapid communication of gravitational-wave date on



### LSST TAKES 20TB OF IMAGES PER NIGHT **Rubin Observatory 10 million transient**

**Transient searches have** relied on human eyes for alerts → real-time ML!

alerts per night!

# such high volumes of data?











#### Semantic segmentation for autonomous vehicles



N. Ghielmetti et al.

#### Other examples

- For fusion science phase/mode monitoring
- <u>Crystal structure detection</u>
- <u>Triggering in DUNE</u>
- <u>Accelerator control</u>
- Magnet Quench Detection
- MLPerf tinyML benchmarking
- Food contamination detection
- etc....







#### The New York Times

A.I. and Chatbots > How the A.I. Race Began One Year of ChatGPT Key Figures in A.I. How A.I. Could Be Regulated

#### THE SHIFT

#### Maybe We Will Finally Learn More **About How A.I. Works**

Stanford researchers have ranked 10 major A.I. models on how openly they operate.

## Al Explainer: Foundation models and the <sup>2</sup> next era of Al

Published March 23, 2023



Figure 5: Representative sample of companies that have publicly stated that they are using, building, or enabling

BigScience	BLOOM	176B	July 2022
	Т0рр	11B	October 2021
	GPT-J	6B	July 2021
	GPT-NeoX	20B	February 2022
(读) 消華大学 Tsinghua University	GLM	130B	August 2022
Google Research	UL2	20B	October 2022
	T5	11B	February 2020
	OPT	175B	June 2022







## Foundation Models



#### Heterogeneous detector Multi-modal input!









 $x = (x_1, x_2, \dots, )$ 









# One model, learn



 $x = (x_1, x_2, \dots, )$ 

Some new space



CMS Simulation I Temmary  $t\bar{t} + PU, \sqrt{s} = 14 \text{ TeV}$ Machine-Learned Particle Flow reconstruction/



HFEM Electrons Muons



# 

CMS

CMS Simulation  $t\bar{t} + PU, \sqrt{s} = 14 \text{ TeV}$ Machine-Learned Particle Flow reconstruction,

> Charged hadrons eutral hadrons HFHAD

HFEM Electrons



# Backup

# Types of anomaly detection

#### Outlier detection

Find (non-resonant) out-of-distribution datapoints



## Detecting overdensities

Find (resonant) overdensities in distributions



# Types of anomaly detection

#### Outlier detection



#### Non-resonant, tail of distributions

- Often (variational) auto-encoders
- Useful for triggering!

#### Caveats

- What's a good metric for optimisation?
- How to use selected events in analysis?

## Detecting overdensities



#### Resonant, similar to a bump hunt

- Density estimation methods
- Useful for offline analysis

#### Caveats

Relies on a definition of "sideband" and a sizeable signal





Compressed representation of x. Latent space  $\Re^k$ , k < m×n prevents memorisation of input, must learn



 $\mathscr{L}(\mathbf{x}, \hat{\mathbf{x}})$  is Mean Squared Error $(\mathbf{x}, \hat{\mathbf{x}})$ , "high error events" proxy for "degree of abnormality"











#### SciPost Physics

Ι



Ve allow for M = 10 trainable linear combinations. These combined 4-vectors of on on the hadronically decaying massive particles. In the original LOLA app the momenta  $\widetilde{k}_j$  onto observable Lorentz scalars and related observables [13] mapping is not easily invertible we do not use it for the autoencoder. Instead, vectors by another component containing the invariant mass,





# Outlier detection in analysis E.g <u>CASE</u>







# Outlier detection in ana E.g <u>CASE</u>









# **Outlier detection in analysis**



# **Example for semi-visible jets**

## F. Eble: Normalized autoencoders





#### R. Seidita: Lund Graph autoencoders






# Finding overdensities



# Finding overdensities - CWoLa bumphunt

#### S enriched sample in data





B enriched sample in data















#### E.g <u>CASE</u>











 $Z(\ell \ell)$ 



q/g jet



#### **JETS FROM MET+JET TOPOLOGY** $\rightarrow$ SIGNAL REGION

#### MIXED SAMPLE 1



#### **JETS FROM** $\ell\ell$ +JET TOPOLOGY $\rightarrow$ SIGNAL NOT EXPECTED HERE

#### MIXED SAMPLE 2



# **Density estimation** Various methods

<u>ML-based interpolation from sidebands to signal region:</u>

**ANODE:** interpolates densities from sidebands to the signal-region & constructs likelihood ratio

**CATHODE:** samples from the background model in signal region after interpolating and estimates likelihood ratio with classifier

LaCATHODE: Use a in flow to perform CATHODE in latent space

**CURTAINS:** Train invertible NN conditioned on mass to map between sidebands

ML-based MC reweighting:

SALAD: Reweight simulation to match sideband, interpolate into the signal region and use a second classifier to get the likelihood

**FETA:** Map simulation to data in sidebands, then compare to SR data







# Why these methods are good for DM searches







We could cast a huge net to catch a broad range of signals in a single search!



#### Do physics with 0.0025% of collision events, the rest is discarded!

# Level-1 hardware trigger0.3% of events left

40 MHz



100% of events left



# High Level Trigger CPU farm0.0025% of events left



Probing smaller and smaller couplings, lower and lower masses

Need <u>more</u> statistics!



# **Anomaly Detection triggers**



Trigger threshold

Energy (GeV)

#### Level-1 rejects >99% of events! Is there a smarter way to select?

# **Anomaly Detection triggers**



Trigger threshold

Energy (GeV)

- - LOST DATA SELECTED DATA - - POSSIBLE NP SIGNAL

#### **Everything here** is normal

**Everything here** is abnormal

AD threshold



**AXOLITL** 

#### Anomaly Detection in the CMS Level 1 µGT for Run3!

#### Input from Run 3 µGT quantities:

•( $p_T$ , n,  $\phi$ ) hardware integer inputs from: 1 MET, 4 e/ $\gamma$ , 4  $\mu$ , and 10 jet objects





**AXOLITL** 





loss =  $|| \mathbf{x} - \mathbf{x}^{\prime} ||^{2} + KL[N(\mu_{x}, \sigma_{x}), N(0, I)]$ 

**AXOLITL** 

#### Only deploy encoder, compute degree of abnormality from patent space only

- Do not need to keep input around for MSE
- Half network size and latency!









# CNN in Level-1 Calorimeter Trigger! Represent calorimeter tower as image and use CNN auto encoder







#### E.g Higgs $\rightarrow$ A(15 GeV) A(15 GeV) $\rightarrow$ 4b

AXOL1TL Rate	1 kHz	5 kHz	10 kHz
Signal Efficiency Gain	46%	100%	133%

#### E.g Higgs $\rightarrow$ A(15 GeV) A(15 GeV) $\rightarrow$ 4b

#### We can do both of these efficiently, model-agnostic and datadriven!





# End-to-end-approach: NPLM

Alternative approach: End-to-end DNN search

- How do we get around defining a signal hypothesis?
- What is alternate hypothesis to test reference?

Idea: Assume alternate model n(x|w) can be parametrised in terms of reference model n(x|R)

$$n(x \mid \overrightarrow{w}) = n(x \mid R)e^{f(x; \overrightarrow{w})}$$
 - Set of real functions

• Let DNN parametrise alternative model

$$f(x; \overrightarrow{w}) = NN$$

unctions



# One model, learn



 $x = (x_1, x_2, \dots, )$ 

Some new space



# One model, learn neural embedding?











#### Something New

# Metric Learning

#### What if we really try to focus on this space



NN





#### Something New

# Neural embedding

#### What if we really try to focus on this space



NN







Something New

# Neural embedding

# Learning the space





# Learning the space

By looking at data, we can learn a lot

- Go over input piece by piece
- Analyze every aspect
- Compare every feature
- Find distinctive style of the input
  - can be done e.g by looking for a deviation





## Cat A



# Dog A



## Cat A



## Dog A



#### Augmented Cat A



## Augmented Dog A





### Cat A



# Dog A



#### Augmented Cat A



#### Augmented Dog A









## Dog A



#### Augmented Cat A



## Jugmented Dog A



# Physically motivated augmentations?



• Minimizing and maximizing distances learns a space

#### Augmented Cat A



#### Cat A



#### Cat B 🔮





# Physically motivated augmentations?



#### No class labels used in training! How do we augment detector data?

# Physically motivated augmentations?



#### No class labels used in training! How do we augment detector data?


Embedded Space can use any NN to embed

## QM foundation models



→ embedding quantum mechanics into AI algorithm

gluon

quark

Н





 $x = (x_1, x_2, \dots, )$ 

### **Training 2: Fine tune for specific task** (fast, small dataset, simulation)

### Training number one: learn embedding







### Something New



### NN

## Training number two: fine tuning



### Something New

## Foundation model of the Level-1 trigger



Charged Fadron (e.g. Pion) Photon

Photo

### 63 Tb/s



# Do I really think this will be possible?



## Careful software-hardware co-design

## O(1M)parameter model on **1000 FPGAs** and do inference in **O(1)µS?**



Accept / Reject



## Careful software-hardware co-design

**Designed our** own protocol to make boards talk to each other fast enough

(25 Gbs to transfer data LHCsynchronously between boards)





Accept / Reject



## TheoristsN-D Space



## Capture Physics

# We can replace the QCD theorist with a NN (And it works better)

(Graph) NN

NN



## NN N-D Space

## Capture Physics

## **Towards end-to-end reconstruction**



# Masked language modelling

#### Next-token-prediction

The model is given a sequence of words with the goal of predicting the next word.

Example: Hannah is a \_\_\_\_

Hannah is a *sister* Hannah is a *friend* Hannah is a *marketer* Hannah is a *comedian* 

# Self-supervised pre-training

#### Masked-languagemodeling

The model is given a sequence of words with the goal of predicting a 'masked' word in the middle.

Example Jacob [mask] reading

Jacob *fears* reading Jacob *loves* reading Jacob *enjoys* reading Jacob *hates* reading