

![](_page_0_Figure_1.jpeg)

![](_page_0_Picture_2.jpeg)

Workshop // on // Fast Realtime Systems // and // Realtime Machine Learning April 9, 2024 https://indico.belle2.org/event/10782/contributions/75182/

![](_page_1_Picture_0.jpeg)

#### Classification

Jinst Published by I	OP PUBLISHING FOR SISSA MEDIALAB Received: April 9, 2021 Accepter: June 29, 2021 Published: August 4, 2021	Jinst
Nanosecond machine learning event of boosted decision trees in FPGA for high	classification with gh energy physics	Nanosecond boosted dec
T.M. Hong,* B.T. Carlson, B.R. Eubanks, S.T. Racz, S.T. Racz Department of Physics and Astronomy, University of Pittsburgh, 100 Allen Hall, 3941 O'Hara St., Pittsburgh, PA 15260, U.S.A. <i>E-mail</i> : tmhong@pitt.edu Anstrakct: We present a novel implementation of classification up intelligence method called boosted decision trees (BDT) on field The firmware implementation of binary classification requiring depth of 4 using four input variables gives a latency value of abs speed from 100 to 320 MHz in our setup. The low timing value BDT layout and reconfiguring its parameters. The FPGA resou a range from 0.01% to 0.2% in our setup. A software package implementation. Our intended user is an expert in custom electror energy physics experiments or anyone that needs decisions at the event classification. Two problems from high energy physics at electrons vs. photons and in the selection of vector boson fusion rejection of the multijet processes. KEYWORDS: Digital electronic circuits; Trigger algorithms; Trigg and software); Data reduction methods ArXIV EPRINT: 2104.03408 *Corresponding author.	the, J. Stelzer and D.C. Stumpp sing the machine learning/artificial programmable gate arrays (FPGA). 100 training trees with a maximum but 10 ns, independent of the clock s are achieved by restructuring the arce utilization is also kept low at called FWXMACHINA achieves this nics-based trigger systems in high lowest latency values for real-time re considered, in the separation of on-produced Higgs bosons vs. the er concepts and systems (hardware	<ul> <li>B.T. Carlson, <sup>a,b</sup> Q. H</li> <li><sup>a</sup> Department of Physic 955 La Paz Road, Sa</li> <li><sup>b</sup> Department of Physic 100 Allen Hall, 3941</li> <li>E-mail: tmhong@pi</li> <li>ABSTRACT: We prese called boosted decisis (FPGA). The softwar paths that allows for new optimization sch timal physics results: of proton collisions : transverse momentum experiments, with a s the firmware perform eight input variables speed, and O(0.1)%</li> <li>KEYWORDS: Data rec cepts and systems (ha ARXIV EPRINT: 2207</li> </ul>
© 2021 IOP Publishing Ltd and Sissa Medialab https://doi.c	rg/10.1088/1748-0221/16/08/P08016	© 2022 IOP Publishing Ltd

PUBLISHED BY IOP PUBLISHING FOR SISSA MEDIALAB RECEIVED: July 13, 2022 ACCEPTED: August 23, 2022 PUBLISHED: September 27, 2022

N

 $\bigcirc$ N

N

**N** SNL

Н

-

Р

 $\bigcirc$ 

6

 $\bigcirc$ 

ω

6

2023

Apr

[ -

[hep-ex]

arXiv:2304.03836v1

nosecond machine learning regression with deep oosted decision trees in FPGA for high energy physics				
T Carlson $a,b$ O Bayer $b$ TM Hond <sup>b,*</sup> and ST Roche $b$				
Department of Physics and Engineering, Westmont College, 055 La Par Road Santa Barbara CA 03108 U.S.A.				
Department of Physics and Astronomy, University of Pittsburgh, 100 Allen Hall, 3941 O'Hara St., Pittsburgh, PA 15260, U.S.A.				
E-mail: tmhong@pitt.edu				
BSTRACT: We present a novel application of the machine learning / artificial intelligence method				

ion trees to estimate physical quantities on field programmable gate arrays re package FWXMACHINA features a new architecture called parallel decision deep decision trees with arbitrary number of input variables. It also features a neme to use different numbers of bits for each input variable, which produces opand ultraefficient FPGA resource utilization. Problems in high energy physics at the Large Hadron Collider (LHC) are considered. Estimation of missing  $m(E_T^{miss})$  at the first level trigger system at the High Luminosity LHC (HL-LHC) simplified detector modeled by Delphes, is used to benchmark and characterize nance. The firmware implementation with a maximum depth of up to 10 using of 16-bit precision gives a latency value of O(10) ns, independent of the clock of the available FPGA resources without using digital signal processors.

duction methods; Digital electronic circuits; Trigger algorithms; Trigger conardware and software)

#### 05602

and Sissa Medialab

https://doi.org/10.1088/1748-0221/17/09/P09039

#### Regression + deep Anomaly detection

PITT-PACC-2311

#### Nanosecond anomaly detection with decision trees for high energy physics and real-time application to exotic Higgs decays

S.T. Roche<sup>a,b</sup>, Q. Bayer<sup>b</sup>, B.T. Carlson<sup>b,c</sup>, W.C. Ouligian<sup>b</sup>, P. Serhiayenka<sup>b</sup>, J. Stelzer<sup>b</sup>, and T.M. Hong\*<sup>b</sup>

<sup>a</sup>School of Medicine, Saint Louis University <sup>b</sup>Department of Physics and Astronomy, University of Pittsburgh <sup>c</sup>Department of Physics and Engineering, Westmont College

#### April 11, 2023

#### Abstract

We present a novel implementation of the artificial intelligence autoencoding algorithm, used as an ultrafast and ultraefficient anomaly detector, built with a forest of deep decision trees on FPGA, field programmable gate arrays. Scenarios at the Large Hadron Collider at CERN are considered, for which the autoencoder is trained using known physical processes of the Standard Model. The design is then deployed in real-time trigger systems for anomaly detection of new unknown physical processes, such as the detection of exotic Higgs decays, on events that fail conventional threshold-based algorithms. The inference is made within a latency value of 25 ns, the time between successive collisions at the Large Hadron Collider, at percent-level resource usage. Our method offers anomaly detection at the lowest latency values for edge AI users with tight resource constraints.

Keywords: Data processing methods, Data reduction methods, Digital electronic circuits, Trigger algorithms, and Trigger concepts and systems (hardware and software).

\*Corresponding author, tmhong@pitt.edu

Hong et al., JINST 16, P08016 (2021) http://doi.org/10.1088/1748-0221/16/08/P08016

Carlson et al., JINST 17, P09039 (2022) http://doi.org/10.1088/1748-0221/17/09/P09039

Roche et al., accepted for publication https://arxiv.org/abs/2304.03836

#### Outline

TM Hong

![](_page_2_Picture_2.jpeg)

#### Introduction

- Autoencoders for anomaly detection
- Machine learning at L1
- Decision tree autoencoder
  - Novel training method
- Firmware design
  - Novel latent-spaceless design for FPGA

# Physics & FPGA results

- Exotic decay of Higgs to pseudoscalars to 2e 2µ
- "LHC anomaly detection" dataset

Save How to find BSM without models at L1

# **Prior work**

![](_page_3_Picture_2.jpeg)

#### Autoencoder

- Typically constructed using neural networks
- Challenge to implement in pure digital logic on FPGA
- NN example shown on right \_\_\_\_

#### Decision tree?

- Used in our work
- Has certain advantages: technical (no multiplication) & philosophical (interpretable)

Govorkova et al., Autoencoders on field-programmable gate arrays for real-time, unsupervised new physics detection at 40 MHz at the Large Hadron Collider, Nature Mach. Intell. **4** (2022) 154–161 https://doi.org/10.1038/s42256-022-00441-3

![](_page_3_Figure_11.jpeg)

**Extended Data Fig. 1** | Network architectures. Network architecture for the DNN AE (top) and CNN AE (bottom) models. The corresponding VAE models are derived introducing the Gaussian sampling in the latent space, for the same encoder and decoder architectures (see text).

#### Anomaly detection in HEP

![](_page_4_Picture_2.jpeg)

#### Model-agnostic detection of BSM signals

- Many anomaly detection methods have been devised and tested on a variety of different HEP problems [https://iml-wg.github.io/HEPML-LivingReview]
- Anomaly detection in ATLAS analysis [ATLAS-CONF-2022-045]

#### Can't analyze data that's not saved

- L1 triggers at ATLAS & CMS use custom electronics such as FPGAs to discard 99.8%
- Implementing anomaly detection at the L1 is challenging and possible (this talk)

![](_page_4_Figure_9.jpeg)

# Decision trees

![](_page_5_Picture_1.jpeg)

![](_page_6_Picture_0.jpeg)

<sup>1</sup> Denby, <u>Comp. Phys. Comm. 49-3, 429 (1988)</u>
 <sup>2</sup> Duarte et al., <u>J. Instrum. 13, P07027 (2018)</u>
 <sup>3</sup> CMS Collaboration, <u>Phys. Lett. B 716, 31 (2012)</u>
 <sup>4</sup> Summers et al., <u>J. Instrum 15, P02056 (2020)</u>
 <sup>5</sup> Hong et al., <u>J. Instrum. 16, P08016 (2021)</u>
 <sup>6</sup> Carlson et al., <u>J. Instrum. 17, P09039 (2022)</u>

#### Neural Network

Popular Depth Score Been around HEP since the 80s<sup>1</sup> Challenging, so ~3 on FPGA<sup>2</sup>  $y = \Theta(M \cdot x + b)$ 

Activation Multiplication

![](_page_6_Picture_6.jpeg)

#### Decision Tree

Popular Depth Score Discovered the Higgs!<sup>3</sup> Challenging, so 4 to 8 on FPGA<sup>4,5,6</sup>  $y = \Theta(x < \text{threshold})$ 

![](_page_6_Figure_10.jpeg)

• FWX Decision Tree

PhysicsComparable results vs. NN on FPGAFloat / fixedBit integer  $\rightarrow$  bit shifts  $\rightarrow$  efficientOptimizedParallelize  $\rightarrow$  one step  $\rightarrow$  low latency

![](_page_7_Picture_0.jpeg)

## Regression

#### Look-up table

- Toy problem in 1-d
- Train / test on f(x) = sin(x) + Gaussian(x)
- For sample of x: y = f(x) in 16 bits

![](_page_8_Figure_5.jpeg)

![](_page_8_Picture_6.jpeg)

![](_page_8_Picture_7.jpeg)

![](_page_9_Figure_0.jpeg)

![](_page_10_Figure_0.jpeg)

# Autoencoder

**Tree-based** 

![](_page_11_Picture_2.jpeg)

#### Autoencoder intro

![](_page_12_Picture_1.jpeg)

#### Example: handwritten numbers

• Teach it about the number 4

![](_page_12_Picture_4.jpeg)

#### Corresponding data set

Image	Pixel I	Pixel 2	 Pixel 300	 Pixel 783	Pixel 784
I	0	0	 240	 0	0
2	0	I	 255	 0	0
i			 	 	
500k	0	0	 231	 0	0

#### Details

• Each pixel in the data set are unrelated to each other

=

#### Autoencoder intro

![](_page_13_Picture_1.jpeg)

#### Example: handwritten numbers

• Teach it 0, 1, 2, 3, 4 with a sample

![](_page_13_Figure_4.jpeg)

#### Details

• Input-output distance is relatively small = good compression

#### **Autoencoder intro**

![](_page_14_Picture_1.jpeg)

#### Example: handwritten numbers

Teach it 0, 1, 2, 3, 4 with a sample (doesn't know about 9!)

![](_page_14_Figure_4.jpeg)

Output looks bad!

#### Details

Input-output distance is relatively large = bad compression

![](_page_15_Picture_1.jpeg)

#### Training philosophy (novel method described in paper)

- Place small "bins" around locations of high event density
- Example
  - 2d toy dataset, say  $x = p_T$  and y = eta for some SM sample

![](_page_15_Figure_6.jpeg)

![](_page_16_Picture_1.jpeg)

Training philosophy (novel method described in paper)

- Place small "bins" around locations of high event density
- Choose variable by sampling the max of the distributions

![](_page_16_Figure_5.jpeg)

Х

![](_page_16_Picture_6.jpeg)

![](_page_17_Picture_1.jpeg)

Training philosophy (novel method described in paper)

- Place small "bins" around locations of high event density
- Sample the variable for a cut, then repeat

![](_page_17_Figure_5.jpeg)

![](_page_17_Picture_6.jpeg)

![](_page_18_Picture_1.jpeg)

Training philosophy (novel method described in paper)

- Place small "bins" around locations of high event density
- Iteratively repeat for subsamples

![](_page_18_Figure_5.jpeg)

![](_page_18_Picture_6.jpeg)

![](_page_19_Picture_1.jpeg)

#### Latent space is bin number

- Encoding: Event  $\rightarrow$  which bin it's in
- Decode by returning a "reconstruction point"
  - Decoding: Bin  $\rightarrow$  median of the training data in bin

![](_page_19_Figure_6.jpeg)

![](_page_20_Picture_1.jpeg)

#### How does this detect anomalies?

• Define: Distance between input – output = anomaly score

![](_page_20_Figure_4.jpeg)

![](_page_21_Picture_1.jpeg)

#### How does this detect anomalies?

- Define: Distance between input output = anomaly score
- Non-anomaly
  - Input is similar to training data
  - Will likely land in a small bin → close to reconstruction point

![](_page_21_Figure_7.jpeg)

![](_page_22_Picture_1.jpeg)

#### How does this detect anomalies?

- Define: Distance between input output = anomaly score
- Non-anomaly
  - Input is similar to training data
  - Will likely land in a small bin → close to the reconstruction point
- Anomaly
  - Input is not similar to training data
  - Will likely land in a large bin → far from the reconstruction point

![](_page_22_Figure_10.jpeg)

## Toy dataset (2 input variables) %

![](_page_23_Picture_1.jpeg)

![](_page_23_Figure_2.jpeg)

#### FWXMACHINA

![](_page_24_Picture_1.jpeg)

#### Latent spaceless implementation

• Closer look at what it means to encode

![](_page_24_Figure_4.jpeg)

• Skip the encoding & decoding

![](_page_24_Figure_6.jpeg)

### FWXMACHINA

# Logic flow

- Left-to-right data flow (see right)
- Realized that we can bypass the latent space!

![](_page_25_Figure_4.jpeg)

![](_page_25_Picture_6.jpeg)

![](_page_26_Figure_0.jpeg)

![](_page_26_Figure_1.jpeg)

![](_page_26_Figure_2.jpeg)

# SM 2e 2µ vs. ?

![](_page_27_Picture_1.jpeg)

#### Proof of concept problem

- Background: we generate all SM with 2e 2µ (predominantly ZZ\*)
- Signal: ggF H $\rightarrow$  a<sub>1</sub> a<sub>2</sub>  $\rightarrow$  e<sup>+</sup> e<sup>-</sup>  $\mu$ <sup>+</sup>  $\mu$ <sup>-</sup> (different m<sub>H</sub> & m<sub>a</sub>)

![](_page_27_Figure_5.jpeg)

#### Veto events with lepton $p_T > 23$ GeV

• Consider only events that won't be already captured by L1 trigger

# SM 2e 2µ vs. ?

![](_page_28_Picture_1.jpeg)

#### Proof of concept problem

- Design
  - 40 decision trees with maximum depth of 5
  - 3 variables:  $m_{ee}$ ,  $m_{\mu\mu}$ ,  $m_{4l}$
- Physics results (see figure)
  - Great separation for H<sub>125</sub>
  - May need a "window selection" for H<sub>70</sub>
- FPGA results (see table)
  - Latency within 25 ns = 1 BC
  - Percent-level (or smaller) resource usage
  - No multiplications!

![](_page_28_Figure_13.jpeg)

Parameter	Value		
Clock speed	320 MHz		
Latency	8 ticks (25 ns)		
Interval	1 tick (3.125 ns)		
FF	10k (0.4 %)		
LUT	31k (2.6%)		
DSP	3 (0.04%)		
BRAM	0		

# **Compare with hls4ml**

![](_page_29_Picture_1.jpeg)

#### LHC anomaly detection ds [Sci Data 9, 118]

- Background
  - W  $\rightarrow$  Iv, Z  $\rightarrow$  II, multijet, ttbar
- Signal
  - 4 BSM scenarios
- Input variables
  - 54 variables
  - p<sub>T</sub>, η, φ of the 4 leading μ, 4 leading
     e, 10 leading jets, MET
  - See distributions on the right
- Sample selection
  - Require  $\geq$ 1 lepton w/ p<sub>T</sub> > 23 GeV
  - (L1 will already save these...)

![](_page_29_Figure_14.jpeg)

# **Cross-check with public data**

TM Hong

![](_page_30_Figure_2.jpeg)

**BRAM** 

0.3%

0

# What I presented

![](_page_31_Picture_1.jpeg)

#### Decision tree-based autoencoder

- New training method by sampling, it's density estimation
- More transparent (to me) than neural network-based designs
- Can do problems in high energy physics (3 50 variables)
- Competitive performance vs. hls4ml

#### Efficient implementation

- Latent space-less design where encoding = decoding
- Performance on Xilinx Virtex Ultrascale+ VU9P
  - O(1)% level resource usage
  - ► Fast at 30 ns latency
  - Try it yourself with the provided testbench & IP available online

# What I think about

#### Then what

![](_page_32_Picture_2.jpeg)

- What are we going to do with the events that we save?
  - Everyone is saving rare events that are uncategorized. Who's going to categorize them? CMS recently showed an event display of the most anomalous event. Will we go through one-by-one to try to guess at the physics?
  - There are ideas, but more needed

#### What about benchmarks?

- By construction, it's supposed to pick up events that we don't know about. But to benchmark it, we choose models that we know about. Is this a contradiction? How do we avoid it? Who gets to choose?
- How much trigger bandwidth do we devote to it if we don't know what may be in it?

#### Backup slides

![](_page_33_Picture_1.jpeg)

# **Neural networks basics**

From Bruce Denby, *Tutorial on Neural Network Applications in High Energy Physics: A 1992 Perspective*, FERMILAB-CONF-92 / 121-E

![](_page_34_Figure_2.jpeg)

![](_page_35_Figure_0.jpeg)

Sum of step functions can approximate the desired contour

![](_page_36_Figure_0.jpeg)

![](_page_36_Picture_1.jpeg)

The contour is converted to the final step function

# **Activation function**

Fuzzy boundary using a function

![](_page_37_Figure_2.jpeg)

![](_page_37_Picture_3.jpeg)

Activation fn gives users a handle to control true / false positive rates

# **Decision tree basics**

And how it achieves the same result as NN

![](_page_38_Figure_2.jpeg)

#### Step function for 2d

![](_page_38_Figure_4.jpeg)

![](_page_38_Picture_5.jpeg)

# Flip book

![](_page_39_Figure_1.jpeg)

![](_page_39_Picture_2.jpeg)

#### Unit gaussians of two variables

![](_page_40_Picture_1.jpeg)

![](_page_40_Figure_2.jpeg)

![](_page_40_Figure_3.jpeg)

![](_page_40_Picture_4.jpeg)

#### **Binary classification**

41

![](_page_41_Figure_1.jpeg)

![](_page_41_Picture_2.jpeg)

S

#### **Binary classification**

#### 42

tree1 depth1

# tree1 depth2

![](_page_42_Figure_2.jpeg)

![](_page_42_Figure_3.jpeg)

![](_page_42_Picture_4.jpeg)

#### **Binary classification**

![](_page_43_Figure_1.jpeg)

![](_page_43_Figure_2.jpeg)

![](_page_43_Figure_3.jpeg)

![](_page_43_Picture_4.jpeg)

#### **Binary classification**

#### tree1 depth4

# tree1 depth4

![](_page_44_Figure_2.jpeg)

tree1 depth8

![](_page_44_Picture_4.jpeg)

#### Draws diagonal

# Depth 2

vary trees

![](_page_45_Picture_2.jpeg)

![](_page_46_Figure_0.jpeg)

![](_page_46_Picture_1.jpeg)

![](_page_47_Figure_0.jpeg)

![](_page_47_Picture_1.jpeg)

![](_page_48_Figure_0.jpeg)

![](_page_48_Picture_1.jpeg)

![](_page_49_Figure_0.jpeg)

![](_page_49_Picture_1.jpeg)

![](_page_50_Figure_0.jpeg)

![](_page_50_Picture_1.jpeg)

![](_page_51_Figure_0.jpeg)

![](_page_51_Picture_1.jpeg)

![](_page_52_Figure_0.jpeg)

![](_page_52_Picture_1.jpeg)

![](_page_53_Figure_0.jpeg)

![](_page_53_Picture_1.jpeg)

#### becomes very blurry

## Put it together on one slide

![](_page_54_Figure_1.jpeg)

![](_page_54_Picture_2.jpeg)

Sweet spot depends on the physics problem

# Forest of decision trees

Fuzzy boundary by averaging step functions

![](_page_55_Figure_2.jpeg)

![](_page_55_Picture_3.jpeg)

Forest of decision trees provides the gradient

# **Activation function**

Fuzzy boundary using a function

![](_page_56_Figure_2.jpeg)

![](_page_56_Picture_3.jpeg)

#### Different approach, but same result