

Outline

Introduction

Who we are

FPGA design

- Parallelizing decision trees
- HLS trees → VHDL trees
- Autoencoder ⊕ model distillation

Thoughts on ML4FE

More info

- Code structure & git
- Slides & video tutorials









(Preliminary) new results

1.Classification

• parallel cuts using HLS

inst Published by IOP Publishing for Sissa Medialab	Inst Published by IOP Publishing For Sissa Medialab
Recurves: April 9, 2021 Accestric: June 29, 2021 Pendissheit: August 4, 2021	Ruseuros July 13, 202 Accestras: August 23, 202 Punissino: September 27, 202
Nanosecond machine learning event classification with	Nanosecond machine learning regression with deep boosted decision trees in FPGA for high energy physics
boosted decision nees in FFGA for high energy physics	
II. Hong,* B.T. Carlson, B.R. Eubanks, S.T. Racz, S.T. Roche, J. Stelzer and D.C. Stumpp Department of Physics and Astronomy, University of Pittsburgh,	all Solid Department of Physics and Engineering, Westmant College, OSS Is the Department of Physics and Engineering, Westmant College, OSS Is the Department of Physics and South Reviews (A 92/18 JL S A
100 Allen Hall, 3941 O'Hara St., Pittsburgh, PA 15260, U.S.A. E-mail: tmhono@pitt.edu	b Department of Physics and Astronomy, University of Pittsburgh, 100 Allen Hall, 3941 O'Hara St., Pittsburgh, A 15260, USA.
BSTRACT: We present a novel implementation of classification using the machine learning/artificial	E-mail: tmhong@pitt.edu
elligence method called boosted decision trees (BDT) on field programmable gate arrays (FPGA). e firmware implementation of binary classification requiring 100 training trees with a maximum ph of 4 using four input variables gives a latency value of about 10ns, independent of the clock zed from 100 to 320 MHz in our setup. The low timing values are achieved by restructuring the JT layout and reconfiguring its parameters. The FPAA resource utilization is also kept low at ange from 0.01% fo 0.2% in our setup. A software package called vs/Xnsctrurs achieves this plementation. Our intended user is an expert in custom electronics-based trigger systems in high ergy physics experiments or anyone that needs decisions at the lowest latency values for real-time and classification. Two problems from high energy physics are considered, in the separation of xctrons vs. photons and in the selection of vector boson fusion-produced Higgs bosons vs. the lection of the multijet processes. INVWORDS: Digital electronic circuits; Trigger algorithms; Trigger concepts and systems (hardware d software); Data reduction methods AXIV EPRINT: 2104.03408	ABSTRACT: We present a novel application of the machine learning / artificial intelligence method called boosted decision trees to estimate physical quantities on field programmable gate arrays (FPGA). The software package re%tractine features a new architecture called parallel decision paths that allows for deep decision trees with arbitrary number of input variables. It also features a new optimization scheme to use different numbers of bits for each input variables. It also features a new optimization scheme to use different numbers of bits for each input variable, which produces optimal physics results and ultraefficient FPGA resource utilization. Problems in high energy physics of proton collisions at the Large Hadron Collider (LHC) are considered. Estimation of missing transverse momentum (<i>E</i> ₁ ^{mix}) at the first level trigger system at the High Luminosity LHC (HL-LHC) (Experiments, with a simplified detector modeled by Delphes, is used to benchmark and characterize the firmware performance. The firmware implementation with a maximum depth of up to 10 using eight input variables of 16-bit precision gives a latency value of <i>O</i> (10) ns, independent of the clock speed, and <i>O</i> (0.1)% of the available FPGA resources without using digital signal processors. Kevworks: Data reduction methods; Digital electronic circuits; Trigger algorithms; Trigger concepts and systems (hardware and software) ArXiv EPRINT: 2207.05602
*Corresponding author.	*Corresponding author.
© 2021 IOP Publishing Lad and Sissa Medialab https://doi.org/10.1088/1748-0221/16/08/P08016	© 2022 IOP Publishing Ltd and Sissa Medialab https://doi.org/10.1088/1748-0221/17/09/P09039

2.Regression

• parallel paths using HLS

202

N

JINS н

 \vdash \sim

Ы

 \bigcirc

6

 \bigcirc

- in-house training
- bypassing latent space

nature communications

Nanosecond anomaly detection with decision trees and real-time application to exotic Higgs decays

Received: 23 May 2023	S. T. Roche (1 ² , Q. Bayer (1 ² , B. T. Carlson (1 ²³ , W. C. Ouligian ² , P. Serhiayenka
Accepted: 9 April 2024	J. Stelzer Φ ² & T. M. Hong Φ ² ⊠
Published online: 25 April 2024	We present an interpretable implementation of the autoencoding algorithm
Published online: 25 April 2024 Check for updates	used as an anomaly detector, built with a forest of deep decision trees on
	FPCA, field programmable gate arrays. Scenarios at the Large Hadron Collid at CEBN are considered, for which the autoencoder is trained using known physical processes of the Standard Model. The design is then deployed in re- time trigger systems for anomaly detection of unknown physical processes such as the detection of rare exotic decays of the Higgs boson. The inference made with a latency value of 30 ns at percent-level resource usage using the Xilim Virtue UltraScale+ VVDP FFCA. Our method offers anomaly detection.
	low latency values for edge AI users with resource constraints.

agnostic searches beyond the Standard Model (RSM physics at the Large Hadron Collider (LHC) at CRN The LHC is the highest energy proton and heavy ion collider that is designed to discover the Higgs board" and study its properties" as well as to probe the mixtown and discover the study of the study of the study of the study of the ESM in the collected data despite the piebfors of searches conducted at the LHC, dedicated studie look for arce RSM eversity hard well and dedig processors". An acceler area of A research in high energy phy- sics is in using autoencoders for anomaly detection, much of which provides methods find mare and unanticipated RSM physics. Muchol	of the trigger system accepts between 100 kHz to 11 kHz of collisions, discarding the remaining ~9% of the collisions. Therefore, it is essential to discovery that the PTA-based trigger system is capable of tiggering potential BSA events. A provision study almost at LHC data triggering potential BSA events. A provision study almost at LHC data triggering potential BSA events. A provision study almost at LHC data triggering potential BSA events. A provision study almost popending on the design ²⁷ . In this paper, we present an interpretable implementation of an aconcorder using des decision trest truth and enflorent PCA. imple- tisebald comparison resulting in fast and efficient PCA. Imple-
provides methods to find rare and unanticipated ISM physics. Mucho for the existing iterature, mostly using earnal network-based approaches, focuses on identifying ISM physics in already collected data . Such data: collected at the Life . A feated but separate endown, which is the subject of this paper, is enabling the identification of rare and anomalous data on the real-time trigger path for more detailed investigation offline. The trigger of the subject of the Si as in the proof but but that collected and not real-time to the Si as in the proof but that collected at the corresponding to the Si as in the proof but that collected in the corresponding to the Si as in the proof but but secretise collisions. The real-time trigger path for the ALSA and CMS experiments , e.g., processes data using catcom detectorics using field programmable gate arrays (IPCA) followed by software trigger	threshold comparisons resulting in fast and efficient FPGA imple- mentation with minimal relance on digital signal processes. We train the autoencoder on innown Standard Model RoM processes to help in scenarios for which a specific RSM model is trajected and its dynamics are known, dedicated supervised training against the SM sample, Le, ESM-SM classification, would likely outperform an unsupervised approach of SM-only training. The physics scenarios encoder is abite toringer on ESM-centrols as anomalies without this prior knowledge of the ISM specifics. Nevertheless, we consider a benchmark where our autoencoder outperforms the existing con- ventional cut-based algorithms.

- 3.Autoencoder 4.Hardware trees faster & more efficient
 - in VHDL, no more HLS



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. Nat. Comm. 15 (2024) 3527 http://doi.org/10.1038/s41467-024-47704-8

Serhiayenka et al. NIM A 1072 (2025) 170209 http://doi.org/10.1016/j.nima.2025.170209

Parallel paths

bin



Implementation of decision tree on FPGA

• B. Carlson et al., J. Instrumentation 17 (2022) P09039



Classification \oplus **estimation**



Example: Hard-scatter jet vs. Pileup jet ⊕ E_T regression

• S. Roche et al., Pheno May 19, 2025: https://indico.global/event/812/contributions/126530/ (paper in preparation)



FPGA results

• Parallel paths implementation in VHDL, Serhiayenka et al., NIM A 1072 (2025) 170209





Using Xilinx Ultrascale+ VU9P (vcu118) at 200 MHz

Feature	Value
Latency	2 clock ticks (50 ns)
Interval	1 clock tick (25 ns)
Flip-flops (FF)	10399 (0.44%)
Look-up tables (LUT)	13274 (1.1%)
Digital signal processors (DSPs)	0
Block-RAM (BRAM)	9 (0.36%)
Ultra-RAM (URAM)	0





Autoencoder intro



Example: handwritten numbers

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



Details

- Input-output distance is relatively small = good compression
- Input-output distance is relatively large = bad compression

Tree autoencoder, What?!



NN AE

- Training is a black box, done offline
- Latent space is complex



From CMS Machine Learning Group https://cms-ml.github.io/documentation/training/autoencoders.html

"Starcoder" tree AE

- Training is sampling of 1d pdfs
- Latent space is simple / interpretable



Distance

FPGA version simplified for anomaly at CMS FPGA version can optionally skip latent sp.





https://cds.cern.ch/record/2876546/files/DP2023_079.pdf

From CMS Public Note, DP-2023/079

Image from

https://medium.com/@rushikesh.shende/autoencoders-variational-

Training developed my us



Train by sampling 1d projections

- Encoding: Event → which bin it's in
- Decoding returns "reconstruction point"
 - Decoding: Bin \rightarrow median of the training data in bin



AE to anomaly detector



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



Realized we can skip latent space

Decode?

- Encode: input var \rightarrow bin #
- Decode: bin $\# \rightarrow$ coord.

No need to encode

Incoming

• Starcode: input var \rightarrow coord.







Decode bin 3: return (5,4)

Encode is Decode:

return (5,4)

Encode:

return bin 3



Block diagram





Starcoder vs. hls4ml

VHDL trees projected to be smaller

by 2-5x (preliminary).





DSP

BRAM

1%

0.3%

0

13

0.8%

VAE model distillation



Convert NN model with DT

• ATLAS work presented by R. Gupta today at Pheno 2025, https://indico.global/event/812/contributions/126571



Compression



Jet images in multiple calorimeter layers

• Study by R. Gupta, paper in preparation



Python-based code

Availability

- <u>gitlab.com/PittHongGroup/fwX</u> parallel cuts (paper 1)
- Shared by email request
 parallel paths (paper 2)
 autoencoder (paper 3)
 hardware tree (paper 4)

Collaborators welcome

			Pittsbu	rgh	•	9
••• \	PittHongGroup /	fwX · GitLab ×	+		``	<i>_</i>
\leftrightarrow C	🔿 🔒 http	s://gitlab.com/Pitt	HongGre 50%	☆	± ₹	ງ ≡
₩ Why GitLab Pricing Explo	pre				Sign in Ge	et free trial
Q Search or go to	PittHongGroup / 🔀 fwX					
Project						
X fwX	X, fwX ⊕				☆ Star	0
Anage >	° master ∽ fwX	History	Find file Code ~	Project information	n	
code >	Herge branch 'e Tae Min Hong auth	dev-rajat' into 'master' ••• ored 1 day ago	58c21dc0 [0	29 Commits		
Operate >				P 3 Branches		
Monitor >	Name	first commit	2 vests ago	🖉 1 Release		
💾 Analyze >		removed small error	1 day ago	README		
		udpate	1 week ago	CHANGELOG		
		undate stuff	3 years and	Created on		
	♦ .gitignore	first commit	3 years ago	May 11, 2021		
	CHANGELOG	update stuff	3 years ago			
	H EULA.md	first commit	3 years ago			
	M README.md	Update README.md	3 years ago			
	ne fwX.py	udpate	1 week ago			
	netup.py	debug setup	5 months ago			
	_					
	 Doxygen is avail Doxygen is avail WWW Machine Impact #Dependencies Wirado HLS D Navigate to http Click the icon of Navigate back ti Select the desire that supports you devices) Scroll down a lit for example, Wi Once that is down through the inst. 	Able at https://fwx.pitt.edu/ Muther of the second	In a to delate an account o select a version ons support all stallation method. re rd and progress do' and "Vivado			

Hond

Git structure





- Xconfig creates model configuration tutorial - part 1
 Xfirmware writes HLS or VHDL tutorial - part 2
- Vivado
 synthesize & testbench
 tutorial part 3

FW testbench w/ IP available



http://d-scholarship.pitt.edu/45784/

Screenshots in the document



Please download Vivado 2019.2 at the following link, if you do not currently have it: <u>https://www.xilinx.com/support/download/index.html/content/xilinx/en/downloadNav/vivado-design-tools/archive.html</u>

Before Beginning

Before beginning, please make sure that you have (and know the location of) the autoencoder IP folder, and the VHDL testbench files:

Name	Date modified	іуре	Size
autoencoder8var_ip	2/7/2024 1:30 PM	File folder	
tb_vhd_files	2/8/2024 11:50 AM	File folder	

Creating New Project in Vivado

Open Vivado 2019.2 and select "create new Project." On the following pop-up, select "next," and you will be prompted to name the project. Name the project as you wish and choose a location to store it. Keep clicking next until you reach a page that prompts you to select the part/ board. For this tutorial, we will be using the Virtex UltraScale+ VCU118 board. After you have selected your part or board, keeping clicking "next" until you have reached the end of the setup page.

New Project												
fault Part pose a default Xilinx part or board for your project.												
Parts Boards												
Reset All Filters Update Board Repositor												
Vendor: All 🗸 Name: All			✓ Boar	d Rev: Latest 🔍								
Search: Q-vcu118 💿 🗸 (1 match)											
Display Name	Preview	Vendor	File Version	ersion Part								
Virtex UltraScale+ VCU118 Evaluation Platform		xilinx.com	2.3	xcvu9p-flga2104-2L-e								
<												
2		r Back	Most >	Einich Canc								







More info

Start page

• <u>fwx.pitt.edu</u>





Information regarding the fwX project will be available on this page. This project is developed by members of the Hong Group in the Department of v and collaborators Physics and As

What is fwX

• Its full name is "firmware ex machina," a play of the phrase in Latin / Greek deus ex machina / θεὸς ἐκ μηχανῆς. Since it's a mouthful to say, we refer to it as fwX

• It is a software package to design nanosecond implementation of machine learning / artificial intelligence algorithms on FPGA for use in high energy physics

Some figures

Nature Communications page



Caption Illustrative example of *****coder as two visual representations of the same decision tree. Deep decision tree (left) rendered as the decision tree grid (center) and implemented by the paralle decision paths (right). Two-depth deep decision tree (DDT) is the encoder (step 1) shown as a conventional binary split diagram; the latent space is the bin number (step 2); the latent space data is decoded using the decision tree grid (DTG) (step 3); and the simultaneous encoding and decoding with **±**coder (star-coder) architecture (right) represented by parallel decisio paths (PDP) of Ref. [79]. The DTG is the visualization as a grid of partitions in V-dimensional space. In this example, the input x = (55, 70) yields the output x = (27, 25) without needing to explicitly produce the latent laver Demonstration of decision tree-based autoencoder and a demonstration of data transmission / anomaly detection using the MNIST dataset, which is a set of images of handwritten numbers converted to 28 × 28 pixels, or 784-length input vector

ative to the

<u>⊻</u> ப ≡



	3	Anomaly detection with end-to- end decision tree-based autoencoder in HLS Application in ATLAS Upgrade ks / Posters Date Type: Title 2021-05-24 Talk: Comparison: results 2021-06-06 Poster: Nanoseco		used in v1 of the paper draft [arXiv:2304.03836v1] bencoder in HLS bencoder in HLS		 scholarship.pitt.edu/id/eprint/4443 testbench is used in v1 of the paper [arXiv:2304.03836v1] IP testbench: Xilinx inputs for nanos detection with decision trees for two jets, http://d-scholarship.pitt.edu/id (2024-02-01). This testbench is used of the paper. 	pprint/44431 (2023-04-23). This of the paper draft its for nanosecond anomaly trees for two photons and two p.pitt.edu/id/eprint/45784 ench is used in the final version		
	4	Application in	n ATLAS Upgrade	0 -		0 -			
1	Talks	/ Posters							
	#	Date	Type: Title		Ven	ue / Link	Speaker		
	1	2021-05-24	Talk: Comparisons results	to hls4ml's boosted decision tree	Phe indi	nomenology Symposium, Pheno 2021, co	T.M. Hong		
	2	2021-06-06	Poster: Nanosecon high energy physic:	d machine learning with BDT for s	Virte Offs	ual HEP conference on Run4@LHC, hell 2021, indico	B.T. Carlson		
	3	2021-07-13	Talk: Nanosecond r high energy physic:	nd machine learning with BDT for ysics		sion of Particles and Fields (DPF) in the erican Physical Society (APS), indico	B.T Carlson		
	4	2021-09-28	Seminar: Invisible H at the LHC	liggs decays & trigger challenges	University of Geneva, Switzerland		T.M. Hong		
	5	2021-10-18	Talk: Presentation	of fwX BDT	18th Exp ICAI	n Int'l Conf. on Accelerator and Large erimental Physics Control Systems, LEPCS 2021, indico	S.T. Roche		
	6	2021-10-22	Seminar: Machine l the LHC: A discussi decision trees, Rea	earning in real-time triggers at on on Machine learning, Boosted I-time trigger, and ML on FPGA	Dep Ten	artment of Physics, University of nessee, Knoxville	T.M. Hong		
	7	2021-10-20	Poster: Presentatio	n of fwX BDT	IEEE Nuclear Science Symposium and Medical Imaging Conference, 2021 IEEE NSS MIC, link		S.T. Racz		
-									

PIKIMO 11 indice

2023, indic

Phenomenology Symposium, Pheno 2023

Fast Machine Learning for Science Workshop

fwXmachina example: Anomaly detection, Mendeley Data, doi:

(2023-04-11). This sample is

10 176224

Hong

Pittsburgh

Python: Available upon request

detection with decision trees. http://

IP testbench: Xilinx inputs for nanosecond anomaly

T.M. Hong

S.T. Roche

T.M. Hong

9

Tutorial

SMARTHEP Edge ML School 9/24/24

Slides

indico.cern.ch/event/1405026/contributions/6103378/

Videos on synthesizing & test bench

indico.cern.ch/event/1405026/contributions/6103386/

fwX_Tutorial - [C/Users/pas218/fwxHDL Elle Edit Flow Tools Repor	/fwX_Tutorial ts <u>W</u> indov	/fwX_Tutoria v Lagout	Lxpr] - Vivado 2019.2 View Bun Help Q. Qui	ck Access							-	□ × Ready
BAR BEX .		Σ %	1 X H	10 us	✓ 王	C					III Default La	yout 🗸
Flow Navigator 🗄 🔍 –	SIMULATIO	N - Behavio	ral Simulation - Functional - sim_1 - a	e_testbench								? :
PROJECT MANAGER					_							
© Settings		- 0 8	fwX_ae_behavioral_tb.vhd ×	ntitled 1 ×								? 🗆 🖸
Add Sources	Q."	۹."	Q 🖬 Q Q 💥 📲	I PI	12 27 4	Fe of H						٥
	Nam^	Narr ^						423.965 ns				^
Language Templates	~ 8	и	Name	Value	0.000 ns	200.00	ns40	0.000 ns 6	00.000 ns	000.	.000 ns	1,000.0
P Catalog		1	> Vevent0Temp[7:0]	62	84	6 X	62	X 41	X	90	X 146	
		14	> Vevent1Temp[7:0]	176	(193)	254	176	137) I	81	X 201	
IP INTEGRATOR		> 10	> Vevent2Temp[7:0]	219	210	120	219	235	2	117	206	
Create Block Design		2.0	> Vevent3Temp[7:0]	255				255				
Open Block Design	-	2.8	> VexpectedDistTemp[7:0]	251		255	251	179		2	27	
Generate Block Design	_		> V prevEvent0Temp[7:0]	62	0 84	6	X 62	41		90	146	
		2.1	# prevEvent1Temp(7:0)	176	0 193	254	176	137	X	181	X 201	
SIMULATION		> 11	> V prevEvent2Temp[7:0]	219	0 210	120	X 219	235	X	217	206	
Run Simulation		> 10	> ♥ prevEvent3Temp[7:0]	255	(a)			255				
		> 1	> vevent0[7:0]	62	0, 84	6	62	41		90	146	
RTL ANALYSIS		> 10	> event1[7:0]	176	0 193 X	254	176	117		181	201	\rightarrow
> Open Elaborated Design		> 10	> vevent2[7:0]	219	0 210 j	120	<u>_}</u>	A 235		217	<u></u> × × × × × × × × × × × × × × × × × × ×	
		5.10) Was return[7:0]	255	<u> </u>	Vavav v	255	251	179	¥	2	\Rightarrow
 SYNTHESIS 		> 8	> weeter output[70]	251	0Y 7 Y	255	¥ 251	179	Y	2	27	\rightarrow
Run Synthesis		> 10	Ve addrTemp	3		1 Y	2	A 3	Y 4			$ \rightarrow $
A Constant Series		5.0	> ¥ addr(9:0)	3	(0)		2	3	1	÷	5	\equiv
 Open symmesized Design 		и			<u></u>	A			MI.	A		
	<>> v	> N < > Y		<	> <							
- INFLEMENTIATION	-			_								
 Run implementation 	Tcl Conso	le × Mes	ssages Log									
 Open Implemented Design 	Q ¥	€ II	B 38 B									
	t rur	1000ns									Salty	
PROGRAM AND DEBUG	INFO:	[USF-XSI:	n-96] XSim completed. Design	snapshot	ae_testbench	behav' loaded.					B	
Generate Bitstream	- INPO:	tosr-xsis	n-97] xSim simulation ran fo ion: Time (s): cpu = 00:00:0	7 ; elapse	i = 00:00:08	. Memory (MB): pee	a = 1064.094 ; gai	n = 6.12			-	
> Open Hardware Manager	1	-									20	Maria h
	<							1		1	The Same	No. of Concession, name
											a start of	
P Type here to search												4.6
											4	- 11

Talk: Comparisons of fwX's BDT to hls4ml's neural

Talk: Decision tree autoencoder anomaly detection

Talk: fwXmachina part 1: Classification with boosted

decision trees on FPGA for L1 trigge

network results

on FPGA at L1 triggers

8 2021-12-04

2023-05-12

2023-09-25 10