

ML on FPGA for Event Selection

Sergey Furletov Jefferson Lab

Workshop VIII on Streaming readout

29 Apr 2021

04/29/21

Sergey Furletov

Outline

- Motivation
 - ✓ Examples from HEP experiments at LHC
 - ✓ EIC experiment
- > ANN in FPGA
- > ML in FPGA for physicists
 - ✓ Experimental setup
 - ✓ Offline PID analysis with ML (root / TMVA)
 - ✓ Moving on to FPGA
 - ✓ Run ML on FPGA
- Optimization ML in FPGA with hls4ml package
- Global PID with multiple PID detectors
- Hardware solutions
- Outlook

<u> Team :</u>

F. Barbosa, L. Belfore (ODU), C. Dickover, C. Fanelli (MIT),Y. Furletova, L. Jokhovets (Jülich Research Centre, Germany),D. Lawrence, D. Romanov

Jefferson Lab

Accelerator Facility



- With increase of luminosity for accelerator colliders as well as a granularity of detectors for particle physics, more challenges fall on the readout system and data transfer from detector front-end to computer farm and long term storage.
- Concepts of trigger-less readout and data streaming will produce large data volumes being read from the detectors.
- From a resource standpoint, it makes much more sense to perform data pre-processing and reduction at early stages of data streaming. Would allow to use information-rich data sets for event selection.



Example from CMS



- CMS bandwidth: Phase-2 ~50 Tb/s (1.8TB/s in Phase-1)
- The task of the real-time processing is to filter events to reduce data rates to manageable levels for offline processing.
- Level-1 typically uses custom hardware with ASICs or FPGAs (decision ~4 μs)
- The second stage of triggering, High Level Trigger (HLT), uses commercial CPUs to process the filtered data in software. (decision ~100 000 μs)

Machine Learning Inference with FPGAs

Reconstruction, Trigger, and Machine Learning for the HL-LHC @ MIT April 26th, 2018



Motivation



- The growing computational power of modern FPGA boards allows us to add more sophisticated algorithms for real time data processing.
- Many tasks could be solved using modern Machine Learning (ML) algorithms which are naturally suited for FPGA architectures.

Level 1 works with Regional and sub-detector Trigger primitives

Using ML on FPGA many tasks from Level 2 and/or Level 3 can be performed at Level 1



ML in physics

- ✦ Machine learning methods are widely used and have proven to be very powerful in particle physics.
- Although the methods of machine learning and artificial intelligence are developed by many groups and have a lot in common, nevertheless, the hardware used and performance is different:
 - CPU only
 CPU and GPU accelerator
 CPU and FPGA accelerator
 pure FPGA
- While the large numerical processing capability of GPUs is attractive, these technologies are optimized for high throughput, not low latency.
- FPGA-based trigger and data acquisition systems have extremely low, sub-microsecond latency requirements that are unique to particle physics.
- Definitely FPGA can work on a computer farm as an ML accelerator, but the internal FPGA performance will be degraded due to slow I/O through the computer and the PCIe bus. Not to mention the latency, which will increase by 2-3 orders of magnitude.
- Therefore, the most effective would be the use of ML-FPGA directly between the front-end stream and a computer farm, on which it is already more efficient to use the CPU and GPU for ML/AI.

efferson Lab

EIC readout





- The correct location for the ML on the FPGA filter is called "FEP" in this figure.
- This gives us a chance to reduce traffic earlier.
- Allows us to touch physics: ML brings intelligence to L1.
- However, it is now unclear how far we can go with physics at the FPGA.
- Initially, we can start in pass-through mode.
- Then we can add background rejection.
- Later we can add filtering processes with the largest cross section.
- In case of problems with output traffic, we can add a selector for low cross section processes.
- The ML-on-FPGA solution complements the purely computer-based solution and mitigates DAQ performance risks.

Filter design proposal

Jefferson Lab



Images from the Internet are for illustration of scale only.

Sergey Furletov



Beam test with GEMTRD and eCAL



9

Beam setup at JLab Hall-D

04/29/21



• Tests were carried out using electrons with an energy of 3-6 GeV, produced in the converter of a pair spectrometer at the upstream of GlueX detector.



GEM-TRD prototype



- A test module was built at the University of Virginia
- The prototype of GEMTRD/T module has a size of 10 cm × 10 cm with a corresponding to a total of 512 channels for X/Y coordinates.
- The readout is based on flash ADC system developed at JLAB (fADC125) @125 MHz sampling.
- GEM-TRD provides e/hadron separation and tracking





GEMTRD clusters on the track



GEM-TRD can work as mini TPC, providing 3D track segments



GEMTRD offline analysis





• For data analysis we used a neural network library provided by root /TMVA package : MultiLayerPerceptron (MLP)

- All data was divided into 2 samples: training and test samples
- Top right plot shows neural network output for single module:

Red - electrons with radiator

Blue - electrons without radiator

Moving forward



- Offline analysis using ML looks promising.
- Can it be done in real time ?
- Here are some of the possible solutions :
 - > Computer farm.
 - CPU + GPU
 - CPU + FPGA
 - FPGA only
- Steps to implement an FPGA solution:
 - Select FPGA for application in ML
 - > Export an offline trained neural network (NN) from root to C++ file.
 - Convert logical topology of NN coded in C++ to RTL structure of FPGA in VHDL or Verilog.
 - > Optimize the NN for application in FPGA.
 - > Create an I/O interface and configure FPGA.
 - > Perform the test with hardware.

Artificial Neural Network







It can perform logical operation in parallel

Inference on an FPGA





Image from: https://www.embeddedrelated.com/showarticle/195.php

Modern FPGA

Modern FPGAs have DSP slices - specialized hardware blocks placed between gateways and routers that perform mathematical calculations.

- The number of DSP slices can be up to 6000-12000 per chip.
- In addition, they often have ARM cores implemented using non-programmable gates.



Image from: https://www.embeddedrelated.com/showarticle/195.php

•

Sergey Furletov

Jefferson Lab

Accelerator Facility

Xilinx HLS: C++ to Verilog





The C/C++ code of the trained network is used as input for Vivado_HLS.

The Xilinx Vivado HLS (High-Level Synthesis) tool provides a higher level of abstraction for the user by synthesizing functions written in C,C++ into IP blocks, by generating the appropriate ,low-level, VHDL and Verilog code. Then those blocks can be integrated into a real hardware system.

18//	1// ===================================
10//	2// RTL generated by Vivado(TM) HLS - High-Level Synthesis from C, C++ and SystemC
2 // float_regex.sn:: converted to (tx_t)	3// Version: 2019.1
3 //	4// Copyright (C) 1986-2019 Xilinx. Inc. All Rights Reserved.
4 //cxx Tile	5//
5 #include "trd_ann.h"	6//
6 #include <cmath></cmath>	7
7⊜ /*	* timescale l ps / l ps
8 fx_t ann(int index,fx_t in0,fx_t in1,fx_t in2,fx_t in3,fx_t in4,fx_t in5,fx_t in6,fx_t in7,	These are ins / i ps
9 input0 = (in0 - (fx_t)1.96805)/(fx_t)7.63362;	(* CODE CENERATION INFO-M+sdopp blo in 2010 1 (MIC INDUIT TYDE-SWY MIC INDUIT FLOAT-1
10 input1 = (in1 - (fx_t)4.75766)/(fx_t)11.9138;	10 (* CORE_GENERATION_INFO="(I'damn, nts_1p_2019_1, thts_INPO1_ITPE=CXX, hts_INPO1_FLOAT=1
<pre>11 input2 = (in2 - (fx_t)4.40589)/(fx_t)11.4831;</pre>	
12 input3 = (in3 - (fx t)4.24519)/(fx t)11.2533;	12 module trdann (
13 input4 = $(in4 - (fx t)4.30175)/(fx t)11.2252$;	13 ap_clk,
14 input5 = $(in5 - (fx t)3.87414)/(fx t)10.1781$:	14 ap_rst_n,
15 input6 = (in6 - (fx t)3.75959)/(fx t)9.69367:	15 s_axi_AXILiteS_AWVALID,
16 input = (inf - (fx + 1)3, 84352)/(fx + 1)9, 66213	16 s_axi_AXILiteS_AWREADY,
17 inputs = (ing - (fx +)3 65047)/(fx +)9 00565.	17 s_axi_AXILiteS_AWADDR,
18 input9 - (ing - (fx 1)5 96775)/(fx 1)1 3203	18 s_axi_AXILiteS_WVALID,
switch (index) {	19 s_axi_AXILiteS_WREADY,
	20 s_axi_AXILiteS_WDATA,
zo case o.	21 s axi AXILiteS WSTRB,
21 Teturin neuronox3204C90();	22 s axi AXILiteS ARVALID,
22 default:	23 saxi AXILites ARREADY.
(++)	24 s axi AXILites ARADDR. VEITIOS
24 }	25 s axi AXILiteS RVALID.
25 }	26 saxi AXTLiteS REFADY.
26 */	27 s axi AXII iteS RDATA
27@fout_t trdann(int index, finp_t input[10]) {	28 saxi AXII ites RRESP
<pre>28 input0 = (fx_t(input[0]) - (fx_t)1.96805)/(fx_t)7.63362;</pre>	20 s avi ATLites RVALTD
<pre>29 input1 = (fx_t(input[1]) - (fx_t)4.75766)/(fx_t)11.9138;</pre>	20 S_BAT_AATLICE_DECADY
<pre>30 input2 = (fx_t(input[2]) - (fx_t)4.40589)/(fx_t)11.4831;</pre>	30 S_AAT_AATLIES_DREADI,
<pre>31 input3 = (fx_t(input[3]) - (fx_t)4.24519)/(fx_t)11.2533;</pre>	si s_axi_AxiLites_bresp,
<pre>32 input4 = (fx_t(input[4]) - (fx_t)4.30175)/(fx_t)11.2252;</pre>	32 Interrupt
<pre>33 input5 = (fx_t(input[5]) - (fx_t)3.87414)/(fx_t)10.1781;</pre>	33);
<pre>34 input6 = (fx_t(input[6]) - (fx_t)3.75959)/(fx_t)9.69367;</pre>	
<pre>35 input7 = (fx t(input[7]) - (fx t)3.84352)/(fx t)9.66213;</pre>	35 parameter ap_SI_TSm_state1 = 23°d1;
<pre>36 input8 = (fx t(input[8]) - (fx t)3.65047)/(fx t)9.09565;</pre>	36 parameter ap_SI_tsm_state2 = 23'd2;
<pre>37 input9 = (fx t(input[9]) - (fx t)5.96775)/(fx t)11.3203;</pre>	37 parameter ap_ST_tsm_state3 = 23'd4;
38 switch(index) {	38 parameter ap_ST_tsm_state4 = 23'd8;
39 case 0:	39 parameter ap_ST_tsm_state5 = 23'd16;
40 return neuron0x32b4c90():	40 parameter ap_ST_fsm_state6 = 23'd32;
41 default:	41parameter ap_ST_fsm_state7 = 23'd64;
42 return (fx t)0.:	42parameter ap_ST_fsm_state8 = 23'd128;
	43 parameter ap_ST_fsm_state9 = 23'd256;
Note: fixed point calculation	44parameter ap_ST_fsm_state10 = 23'd512;
Note: fixed point calculation	45 parameter ap_ST_fsm_statel1 = 23'd1024;
45 fx + nouronax32hf85a() {	46parameter ap_ST_fsm_state12 = 23'd2048;
	47 parameter ap_ST_fsm_state13 = 23'd4096;
40 L	48 parameter ap ST fsm state14 = 23'd8192;
40	49 parameter ap_ST_fsm_state15 = 23'd16384;
455	50 parameter ap_ST_fsm_state16 = 23'd32768;
	51parameter ap ST fsm state17 = 23'd65536;
or return input;	52 parameter ap ST fsm state18 = 23'd131072;
Thanks to Dan Davida far halm	53 parameter ap ST fsm state19 = 23'd262144;
Indriks to ben kayuu ior help.	54 parameter ap ST fsm state20 = 23'd524288:
54 TX_T neuronux32DT4dU() {	55 parameter ap ST fsm state21 = 23'd1048576:
collection inputz;	
56 3	

Sergey Furletov

ML FPGA Core for TRD

Jefferson Lab

• Using HLS significantly decreases development time. (at the cost of lower efficiency of use of FPGA resources)



04/29/21

Test ML FPGA



Test tools:

- 1. Vivado SDK
- 2. Petalinux

ev=0	out=0.192	out0=0.197
ev=1	out=0.192	out0=0.197
ev=2	out=0.233	out0=0.236
ev=3	out=0.192	out0=0.197
ev=4	out=0.165	out0=0.169
ev=5	out=0.192	out0=0.196
ev=6	out=0.462	out0=0.470
ev=7	out=0.187	out0=0.191



C++ code for test : XTrdann ann; // create an instance of ML core.





• A package hls4ml is developed based on High-Level Synthesis (HLS) to build machine learning models in FPGAs.



Sergey Furletov

GEMTRD neural network optimization



Full size neural network, accuracy-optimized.

+	Timing (r * Summa	ns): ary:						
	Clo	ck Ta	arget	Estima	ated	Uncert	ainty	
	ap_cl +	× (5.00	3.	968		0.62	
+	Latency * Summa	(clock ary:	cycle	5):	L	aten	cy =	75r
	+ Late min	ency max	Inte min	erval max	Pip T	eline ype		
	15		1	+ 1 +	fun	ction		

±	+		L	L	L4
Name	BRAM_18K	DSP48E	FF	LUT	URAM
+ DSP	–	2			 –
Expression		- 1	0	24	i –i
FIFO	-	-	-	-	-
Instance	19	692	3737	16446	-
Memory	2	-	0	0	-
Multiplexer	-	-	–	36	-
Register		-	1532	-	-
Total	21	694	5269	16506	0
 Available SLR	1440	2280	788160	394080	320
Utilization SLR (%)	1	30	~0	4	0
Available	4320	6840	2364480	1182240	960
Utilization (%)	~0	10	~0	1	0

DSP utilization 10%

Size-optimized neural network

+ Tir ,	ming (ns) ∗ Summary	:							
	Clock	Targ	et	Estima	ated	Uncerta	inty		
	ap_clk	5.	00	3	.883	+ +	0.62		
+ La ⁻	' tency (cl ∗ Summary	ock cy	cle	es):	L	Late	ency =	85ns	
	Latenc min m	y 3 ax m	Int in	terval max	Pi 	peline Type			
(17		3	3 3 -+	fu +	nction +			
+ !	Name		+ 	BRAM_	18K	DSP48E	FF	LUT	+ URAM
DSP Expres	ssion		 		- - -	2	- 0 -	_ 24	– –
Instar Memory	nce / plexer				- 2	177	3132 0 -	10696 0 81	- -
Regist	ter				-	_	1423	-	-
Total			l		2	179	4555	10801	0
Availa	able SLR			1	440	2280	788160	394080	320
Utiliz	zation SL	R (%)		~0		7	~0	2	0
Availa	able			4	320	6840	2364480	1182240	960
Utiliz	zation (%	5)	+ 	~0		2	~0	~0	0

DSP utilization 2%

ML for Calorimeter e/pi separation





Classification		Last-layer activation	Loss function
single-label		softmax	categorical_crossentropy
multi-label (scores for cand	lidates)	sigmoid	binary_crossentropy





Examples of events with e and π^- showers and μ^- passing through.

by D. Romanov

04/29/21

Sergey Furletov

Calorimeter NN implementation report





+ Latency (clock cycles):

* Summary:

Late	ncy	Inte	erval	Pipeline
min	max	min	max	Type
12	12	1	1	function

Latency = 60ns

Name	BRAM_18K	DSP48E	FF	LUT	URAM
DSP	++ 	3	+ –	+ –	
Expression	i –i	- i	0	42	-
FIF0	–	-	-1	-	-
Instance	–	208	931	4426	-
Memory	3	-	0	0	-
Multiplexer	–	-	-1	36	-
Register	-	-	946	-	-
Total	3	211	1877	4504	0
Available SLR	1344	3072	864000	432000	320
Utilization SLR (%)	~0	6	~0	1	0
Available	5376	12288	3456000	1728000	1280
Utilization (%)	~0		~0	~0	0

DSP utilization 1%

Possible hardware solutions for eic @ IP6



- ★ There a multiple way to implement ML-FPGA in DAQ.
- The main idea: to collect information from all detectors in one unit to process the full event.
- For data acquisition in physics, ATCA and OpenVPX standards have become widespread.

Compute Node (PXD, Belle II)



- The pixel detector of Belle II with its ~ 8 million channels will deliver data at rate of 22 Gbytes/s for a trigger rate of 30 kHz
- A hardware platform capable of processing this amount of data is the ATCA based Compute Node. (Advanced Telecommunications Computing Architecture).
- A single ATCA crate can host up to 14 boards interconnected via a full mesh backplane.
- Each AMC board is equipped with 4 Xilinx Virtex-5 FX70T FPGA.







ADC based DAQ for PANDA STT



Level 0 Open VPX Crate

ADC based DAQ for PANDA STT (one of approaches):

- 160 channels (shaping, sampling and processing) per payload slot, 14 payload slots+2 controllers;
- totally 2200 channels per crate;
- time sorted output data stream (arrival time, energy,...)
- noise rejection, pile up resolution, base line correction, ...







 All information from the straw tube tracker is processed in one unit.

 Allows to build a complete STT event.

 This unit can also be used for calorimeters readout and processing.



٠



Sergey Furletov

Unified hardware solution (ATCA or OpenVPX)





Outlook

- An FPGA-based Neural Network application would offer online event preprocessing and allow for data reduction based on physics at the early stage of data processing.
- The ML-on-FPGA solution complements the purely computer-based solution and mitigates DAQ performance risks.
- FPGA provides extremely low-latency neural-network inference on the order of 100 nanoseconds.
- Open-source hls4ml software tool with Xilinx[®] Vivado[®] High Level Synthesis (HLS) accelerates machine learning neural network algorithm development.
- The ultimate goal is to build a real-time event filter based on physics signatures.



Figure 2.1: Feynman diagrams of the Quark Parton Model, QCD-Compton and Boson Gluon Fusion processes in NC DIS.

Published in 2007

Measurement of multijet events at low \$x_{Bj}\$ and low \$Q^2\$ with the ZEUS detector at HERA

T. Gosau



Jefferson Lab

ccelerator Facility

04/29/21





Backup



29

Example from LHCb





LHCb 'Tubo'

Turbo



⁴arXiv:1604.05596, NEW arXiv:1903.01360





11/20

Jefferson Lab

04/29/21

analysis

15kB

 $15 \rightarrow 70 \text{kB}$

the raw event: 70kB



Readout electronics, shaping time optimization



32

Jefferson Lab

Readout electronics, sampling rate optimization



a preamplifier (GAS2 ASIC chip) with shaping times of~10-12ns. The flash ADC has a sampling rate of 125 MHz and 12 bit resolution but provides only pipe-lined triggered readout (price ~50\$/channel)









Sergey Furletov