# **Development of real-time online** tracking system

Qi-Dong ZHOU Institute of Frontier and Interdisciplinary Science, Shandong Univ. (Qingdao)

21-23 Jul. 2027, Huizhou 2nd Workshop on Tracking in Particle Physics Experiments



# Readout system (Belle II vs. LHCb)

- Belle II: L1 trigger + HLT
  - Trigger efficiency:
    - •Had. B physics  $\sim$  100%  $\tau$  physics
      - 70~95%



- •LHCb: "triggerless" readout & DAQ
  - CPU+GPU based software trigger
  - Rate of physical process: ~MHz
    - •No hardware trigger available

- ALICE: continus readout
  - TPC w/ triggerless readout + others w/ hardware trigger •TPC signal: ~100 µs, physical event rate 50 kHz, TPC signal overlap
  - Very basic hardware+ more effective software trigger





# Readout and DAQ system(ALICE)





# Exp.Run timeData (PB)TotalBESIII2008-20280.510STCF-300-500

Exp.	Run time	Data (PB)	Total							
BESIII	2008-2028	0.5	10							
STCF	-	300-500	-							
CEPC	_	- 1.5-3(H) 500-50000 (Z)								
Data	High data rate	ASICs and systems		7.1						
density	New link tech	nologies (fibre, wireless,	wireline)	7.1						
	Power and rea	adout efficiency		7.1						
Intelligence	Front-end prog	Front-end programmability, modularity and configurability								
on the	Intelligent pov	Intelligent power management								
detector	Advanced dat	a reduction techniques	(ML/AI)	7.2						
40	High-performa	High-performance sampling (TDCs, ADCs)								
4U- tochniquos	High precisior	High precision timing distribution								
techniques	Novel on-chip	Novel on-chip architectures								
Extromo	Radiation hard	Radiation hardness								
environmer	Cryogenic ten	Cryogenic temperatures								
and longev	ity Reliability, fau	Reliability, fault tolerance, detector control								
	Cooling	Cooling								
	Novel microele	Novel microelectronic technologies, devices, materials								
Emeraina	Silicon photor	Silicon photonics								
technologie	s 3D-integration	and high-density interc	connects	7.5						
•	Keeping pace	Keeping pace with, adapting and interfacing to COTS								

Must happen or main physics goals cannot be met



\* LHCb Velo







Desirable to enhance physics reach

R&D needs being met

ECFA detector R&D

## Gain power of apparatus with data acceleration • Continues readout (less-hardware filtering) • Powered by hardware acceleration • Heterogeneous computing

## Typical TDAQ system



Trigger-less data readout system

Digital





# Roadmap of machine learning on "FPGA"

## Credit: Y. -T. Lai @ KEK



# Belle II trigger system

- Max. trigger rate: 30 kHz @ 6 x 10<sup>35</sup> cm<sup>-2</sup> s<sup>-1</sup> Challenges:  $\bullet$ 
  - Physics trigger ~15 kHz
- Latency limit: ~5 usec (SVD APV25 buffer structure)
  - A fixed latency of about 4.4 usec
- Event timing resolution: 10 nsec
- CDC, ECL: main triggers for tracks and clusters, KLM: trigger muon, TOP: event timing



- low multiplicity trigger vs. background
- High track trigger vs. crosstalk
- Drawback of track trigger at endcap
- Latency budget vs. transmission and logics

. . .

 $\approx 5\mu s$  after beam crossing





# Motivation of Neural Network for L1 Track trigger

- DAQ system is designed to handle 30 kHz
  - Physical trigger ~15 kHz, require S/N = 1
- L1 trigger rate depends significant on background condition
- Advanced CDC algorithm to further suppress background
- A fixed latency of about 4.4 usec







## Tracks $z_0$ distribution after trigger





Axial wire

Stereo wire



# **Basics of L1 CDC trigger**

# **Deep Neural Network for Z trigger**



- Inputs: Drift time  $t_{drift}$ , wires relative location  $\phi_{rel}$ , Crossing angle  $\alpha$  for priority wires + Drift time for all other wires
- Introduce the self-attention architecture to "focus" on certain inputs
- Output track vertex  $z_0$ , track  $\theta$  and signal/background classifier output (Q)

Parameter	#Attention value	#hidden nodes	#hidden layer	activate	precision	Total multiplier
Values	27	27	2	Leaky Relu	Float 16	4,185

# **Development flow of DNN on FPGA**





## Belle II UT4



XCVU080, XCVU160 25 Gbps with 64B/66B

![](_page_10_Picture_6.jpeg)

# Performance of DNN algorithm

![](_page_11_Figure_1.jpeg)

- Latency : 76 clock = 592.8 ns ;require: < 600ns</li>
- FPGA resource (UT4: Virtex UltraScale XCVU160) usage:
  - DSP: ~70%, LUT: ~50%, others <30%
- AUC do not get large drop comparing RTL and software simulation
- At signal efficiency ~95%
  - Background rejection rate ~85%
- DNN trigger with **HARDWARE** under commissioning, close to operate

## **Classifier output** Signal (RTL) Background (RTL) Signal (software) Background (software)

 $10^{-1}$ 

 $10^{-2}$ 

 $10^{-3}$ 

rate

Background rejection

![](_page_11_Figure_10.jpeg)

Q (%)

![](_page_11_Figure_15.jpeg)

![](_page_11_Picture_16.jpeg)

# **AMD Versal projects**

## AMD Versal<sup>™</sup> Adaptive SoCs

For any application from cloud to edge

New 2nd Gen Versal Portfolio

![](_page_12_Picture_4.jpeg)

Overview Portfolio Developers Resources

Overview

curved gradient divider

![](_page_12_Picture_10.jpeg)

#### **Heterogeneous Acceleration**

Highly integrated, multicore compute platform that can adapt with evolving and diverse algorithms.

![](_page_12_Picture_13.jpeg)

#### **Any Application**

Dynamically customizable at hardware and software levels to fit a wide range of applications and workloads.

![](_page_12_Picture_16.jpeg)

#### Any Developer

Architected around a programmable network on chip (NoC), Versal adaptive SoCs are programmable by software developers and hardware programmers alike.

![](_page_12_Picture_19.jpeg)

![](_page_12_Figure_20.jpeg)

![](_page_12_Picture_21.jpeg)

![](_page_12_Picture_23.jpeg)

# **DNN implementation on Versal ACAP**

- R&D of a new general FPGA device using the Versal ACAP
  - Heterogenous acceleration (VCK190, VCK5000 evaluation kit)
    - Al engine

![](_page_13_Figure_4.jpeg)

UG1079

Figure 2: AI Engine Array

![](_page_13_Figure_7.jpeg)

## sing the Versal ACAP D, VCK5000 evaluation kit)

Figure 4: AI Engine

![](_page_13_Figure_10.jpeg)

![](_page_14_Figure_0.jpeg)

- DNN implementation:
  - Model on a "graph"
  - Dense layer on a "kernel"
- Al engine: C++ based coding on Vitis
  - Al engine libraries
  - Al engine specific functions
  - Scaler, Vector engines, pipelining, etc.

	Layer 1	Layer 2	Layer 3	Layer 4	L
Input nodes	71	27	27	27	
Output nodes	27	27	27	27	
Active Func.	LeakyReLU	Softmax		LeakyReRU	

Al Engine Resource Utilization	
Tiles used for AI Engine Kernels:	5 of 400 (1.25 %)
Tiles used for Buffers:	7 of 400 (1.75 %)
Tiles used for Stream Interconnect:	8 of 450 (1.78 %)
DMA FIFO Buffers:	0
Interface Channels used for ADF Input/Output:	4 ( PLIO: 4 )
Interface Channels used for Trace data:	0

![](_page_14_Figure_11.jpeg)

![](_page_14_Picture_12.jpeg)

![](_page_14_Picture_13.jpeg)

# Latency optimization on Versal ACAP

						5.588 us	
NAME	VALUE	0.000000 us 1.00	0000 us 2.000000 us	s 3.000000 us	4.000000 us	5.000000 us 6.0	
> Tile(23,0)	_main	• • •		_main			
> Tile(23,1)	_main		_main	SO	ftm	_mai	
> Interface Tile(23)							
> Tile(24,0)	_main		_mai	n	hid4_27to2	7	
> Tile(24,1)	_main		main	hid3_27to27_no_a	ct(adf		
> Tile(24,2)	_main		_main	hid2_27to2		_main	
> Tile(25,0)	_main	hid1_71t	o27_leakyrelu(adf::io_b			_main	
					_	<b>T</b>	
	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Iotal	
Input nodes	71	27	27	27	27	—	
Output nodes	27	27	27	27	3	—	
Active Func.	LeakyReLU	Softmax		LeakyReRU	Tanh	—	
Ver.0 latency	~12us	~66us	~1.5us	~5.5us	~9.9us	~86us	
Ver.1 latency	~2.1us	~1.3us	~1.5us	0.9us	~0.2us	~5us	
Ver.2 latency	~0.48 us	~0.93us	~0.33us	~0.40us	~0.10us	~2.1us	

Total **305 clk cycles** one instance. Clock period **10ns**. Latency running on Versal ACAP is **3.05us** 

![](_page_15_Picture_4.jpeg)

![](_page_15_Picture_5.jpeg)

![](_page_15_Picture_6.jpeg)

![](_page_15_Picture_7.jpeg)

# **Real-time and AI integrated High Level Trigger**

![](_page_16_Figure_1.jpeg)

# **GNN for CDC track background filtering**

- for Belle II CDC hits clean up

![](_page_17_Figure_3.jpeg)

![](_page_17_Figure_5.jpeg)

# Acceleration on Versal ACAP platform

![](_page_18_Figure_1.jpeg)

# **Demo of GNN implementation on Versal ACAP**

![](_page_19_Figure_1.jpeg)

Versal is in order of ~ms

![](_page_19_Figure_4.jpeg)

![](_page_19_Picture_5.jpeg)

![](_page_19_Picture_7.jpeg)

## **CNN algorithm for PID** • DTOF as a PID subdetector of STCF CNN algorithm developed for Kaon/pion identification Kaon/Pion MC simple, 800w

![](_page_20_Figure_1.jpeg)

## EfficientNetV1

![](_page_20_Figure_3.jpeg)

## EfficientNetV2

![](_page_20_Figure_5.jpeg)

## Z. Yao et al.@SDU

## pion Signal Efficiency with 2% Misidentification Rate

36 -	100	100	100	100	100	100	100	100	100	100	100	100	99.3	98.6	95.5	91.9	87	75.8	68.4
35 -	100	100	100	100	100	100	100	100	100	100	100	100	100	99.5	98.1	95.8	94	84.8	79.2
34 -	100	100	100	100	100	100	100	100	100	100	100	100	99.7	99.8	99.2	96.5	97.2	86.9	87.5
[6] <sup>33</sup> -	100	100	100	100	100	100	100	99.7	100	100	100	100	100	99.8	99.5	98.5	96.8	94	83.7
р 32 -	100	100	100	100	100	100	100	100	100	100	100	100	100	99.7	99.9	99.3	98.7	97.3	91.9
- <sub>1</sub> د ص	100	100	100	100	100	100	100	100	100	100	100	100	99.4	100	100	99.8	99.1	96.6	94
<u>l</u> g	100	100	100	100	100	100	100	99.7	100	100	100	100	100	99.9	99.8	99.9	99	98.2	97.7
A 29 -	100	100	100	100	99.7	100	100	100	100	100	100	100	100	100	100	99.8	99.7	98.3	97.7
	99.7	100	100	100	100	100	100	100	100	100	99.7	100	100	99.9	99.7	99.8	99.3	98.3	97.2
G 27	99.3	100	100	100	100	100	100	100	99.7	100	100	100	99.7	100	100	99.7	99.4	98.2	98
20 -	100	100	99.7	100	100	100	100	100	100	100	100	100	99.7	100	100	99.7	99.3	98.9	98.5
23	99.4	100	100	100	100	100	100	100	100	100	100	100	100	99.9	99.8	99.8	99.3	98.3	97.7
24	99.8	99.7	99.7	100	100	100	100	100	99.7	100	100	99.7	100	99.6	99.5	99.7	99.1	98.6	98.2
23													5 2. <sup>0</sup>						
	Momentum [Gev/c]																		

![](_page_20_Picture_9.jpeg)

![](_page_20_Picture_10.jpeg)

![](_page_20_Picture_11.jpeg)

# Implementation on DPU with Vitis Al

#### AI Inference Development

If you are an AI developer, bring your TensorFlow and PyTorch trained models to directly infer on Versal using Mipsology Zebra and build, configure, and deploy computer vision applications on FPGA platforms with Aupera Video Machine Learning Streaming Server solution.

### Key Features

Explore partner solutions and articles, and learn about the key features for AI Inference Development with the VCK5000

![](_page_21_Figure_5.jpeg)

2x TCO Reduction vs Mainstream nVidia GPUs

- 2x perf/w and perf/\$ compared to Nvidia Ampere with standard MLPerf Models
- Achieves 90% compute efficiency
- Consume less than 100W at card level

![](_page_21_Figure_10.jpeg)

#### 2x End-to-End Video Analytics Throughput vs nVidia GPUs

- Full pipeline from H.264 decode to computer vision to up to 10 AI models
- Video decode and CV run on x86 CPU or discrete U30 Alveo card
- Plug-in based pipeline composition from FFmpeg / Gstreamer

![](_page_21_Figure_15.jpeg)

- board

## AMD Vitis<sup>™</sup> AI Integrated Development Environment

![](_page_21_Figure_21.jpeg)

![](_page_21_Figure_22.jpeg)

VCK5000

Easy to Use with Familiar Frameworks

• Easy-to-use software flow for any CPU & GPU users, no hardware programming required

Run inference from Tensorflow framework directly on

 State-of-the-art model supported with mainstream frameworks PyTorch, TensorFlow, TensoFlow 2 and Caffe

![](_page_21_Picture_29.jpeg)

Setting up explorer 's environment in the Docker container... Running as vitis-ai-user with ID 0 and group 0

![](_page_21_Picture_31.jpeg)

Docker Image Version: ubuntu2004-3.5.0.306 (CPU) Vitis AI Git Hash: 6a9757a Build Date: 2023-06-26 WorkFlow: pytorch

vitis-ai-user@dell-Precision-7960-Tower:/workspace\$ ls

![](_page_21_Picture_34.jpeg)

![](_page_21_Picture_35.jpeg)

# Implementation of CNN on DPU

![](_page_22_Figure_1.jpeg)

- Training performed same as typical software model development
- Quantization and compilation using the Vitis-AI tool (AMD)
- Vitis-AI Runtim and Xilinx Runtime tools used for deployment on Versal ACAP

![](_page_22_Picture_5.jpeg)

## Performance comparison Inference result based on 10000 samples CPU GPU

![](_page_23_Figure_1.jpeg)

## DPU

DPU based on Versal ACAP shows 8-15x(CPU)/2-4x(GPU) inference time

![](_page_23_Picture_5.jpeg)

# Summary and prospects

- collider experiments
- System design, algorithm and heterogeneous platform developments
- DNN algorithm was designed for L1 trigger, implemented to FPGA
- DNN/GNN algorithm for L1 trigger/HLT was approved able to be implemented to Versal ACAP
- A possible solution for latency-limit algorithm of HLT (e.g. track seed finding) DPU based on Versal ACAP shows 8-15x(CPU)/2-4x(GPU) inference time
- performance for the STCF PID CNN model
  - A possible solution for HLT (CNN algorithm)

Advance data reduction technique (real-time, ML/AI) is essential for new/upgraded

![](_page_24_Picture_9.jpeg)

![](_page_24_Picture_10.jpeg)