



# Real-time data processing with ML

周启东 (ZHOU Qi-Dong)

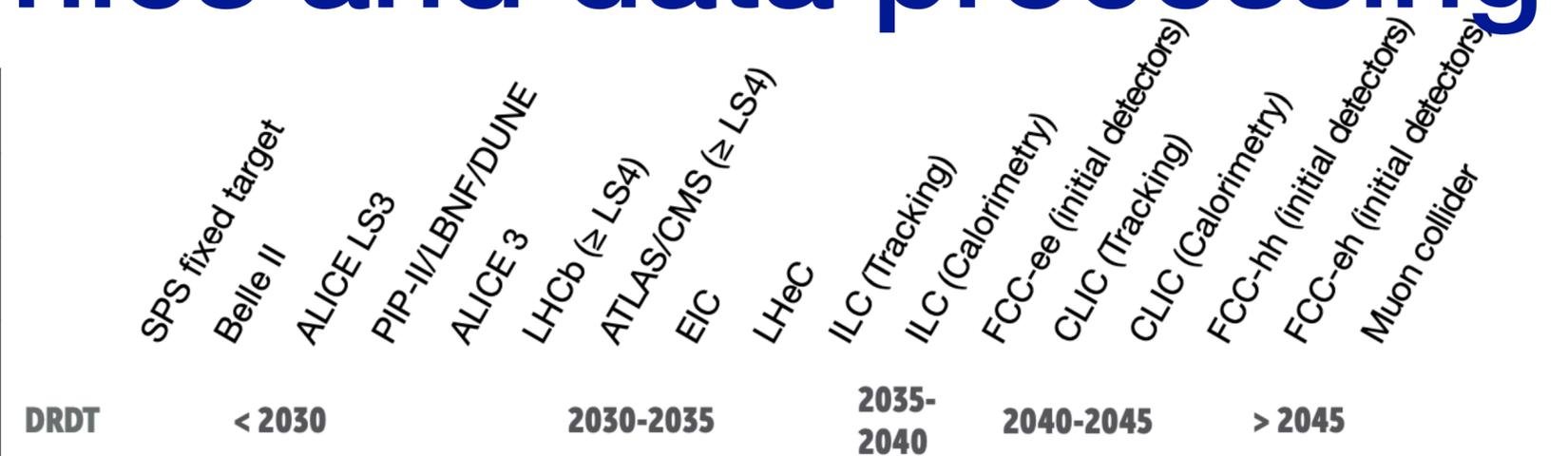
Institute of Frontier and Interdisciplinary Science,  
Shandong Univ. (Qingdao)

11-12 Jan. 2025, Hefei USTC

量子计算和人工智能与高能物理交叉研讨会

# Roadmap of electronics and data processing

Exp.	Run time	Data (PB)	Total
BESIII	2008-2028	0.5	10
STCF	-	300-500	-
CEPC	-	1.5-3(H) 500-50000 (Z)	-



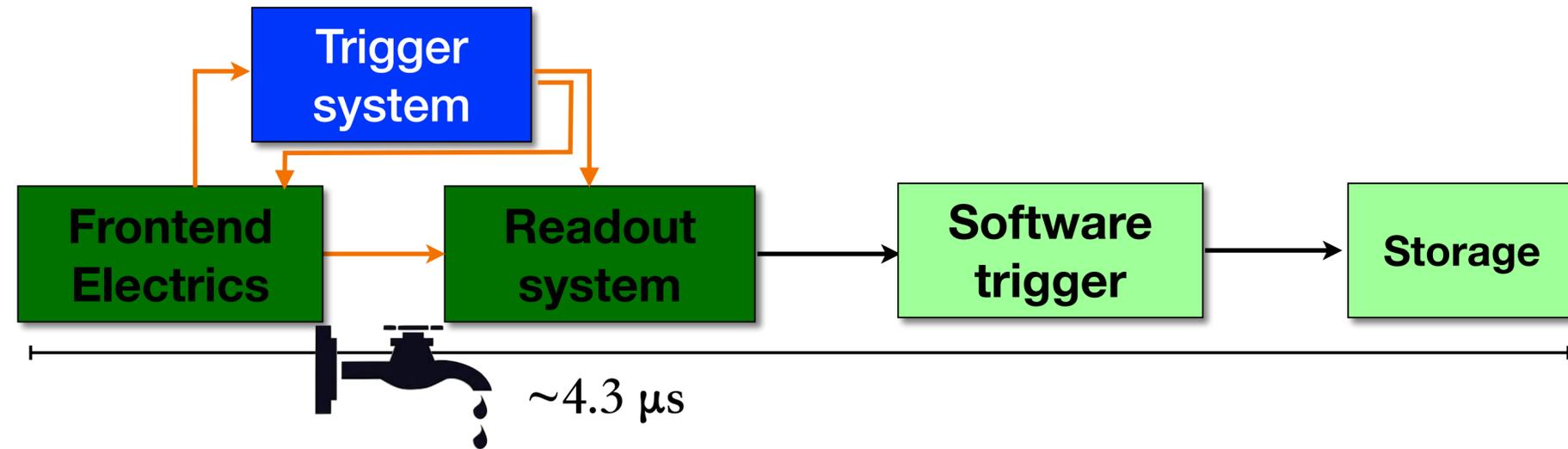
● Must happen or main physics goals cannot be met    
 ● Important to meet several physics goals    
 ● Desirable to enhance physics reach    
 ● R&D needs being met

\* LHCb Velo

# 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:  $\sim$ MHz
    - No hardware trigger available



Latency

Trigger rate 127 MHz

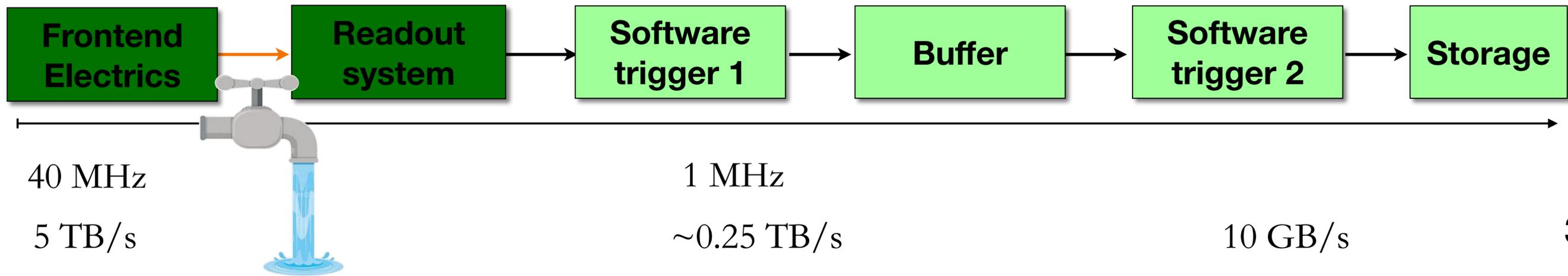
30 kHz

30 kHz

Throughput 3(33) GB/s

2 (32) GB/s

3 GB/s



Trigger rate 40 MHz

1 MHz

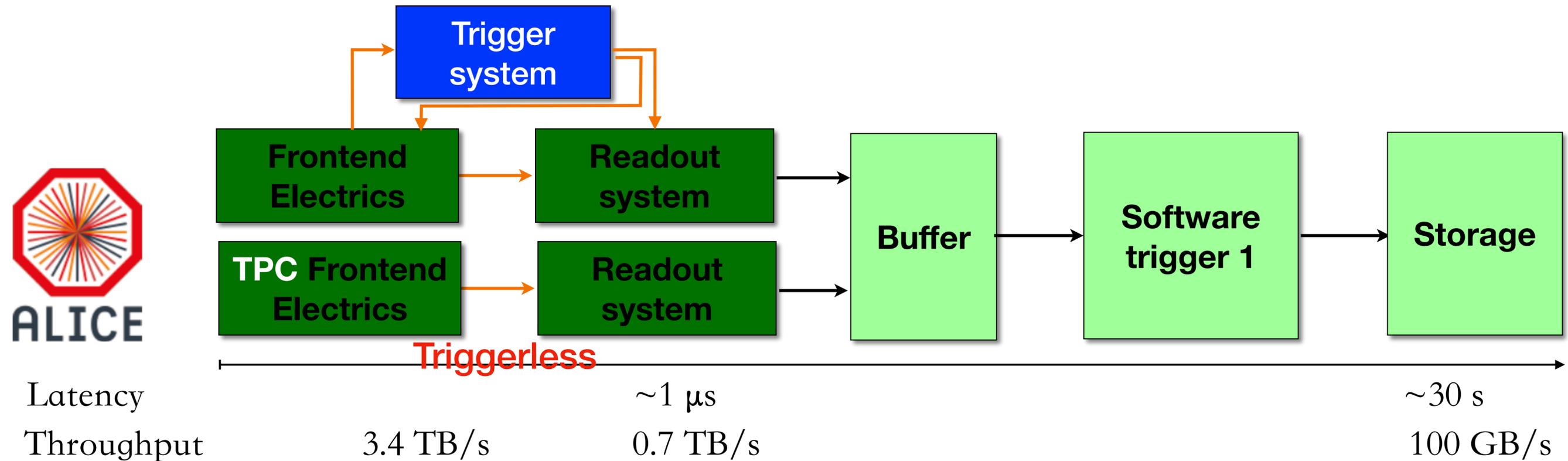
Throughput 5 TB/s

$\sim 0.25$  TB/s

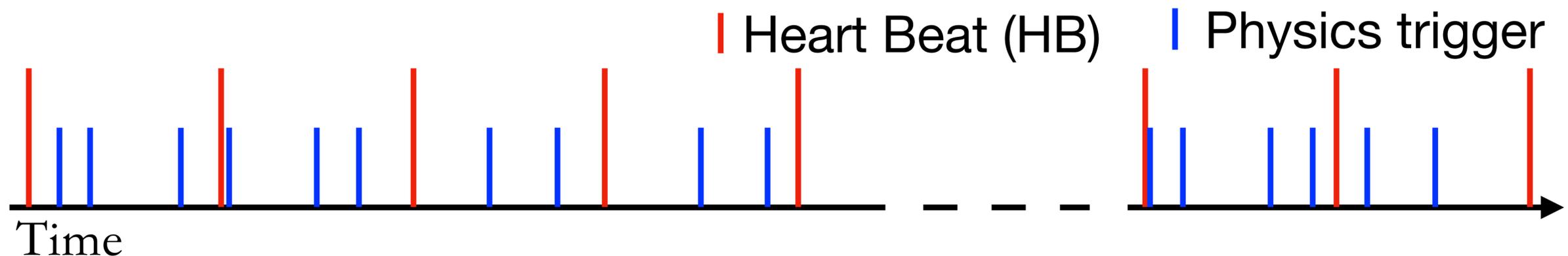
10 GB/s

# Readout and DAQ system(ALICE)

- ALICE: continuous readout
  - TPC w/ triggerless readout + others w/ hardware trigger
    - TPC signal:  $\sim 100 \mu\text{s}$ , physical event rate 50 kHz, **TPC signal overlap**
  - Very basic hardware+ more effective software trigger

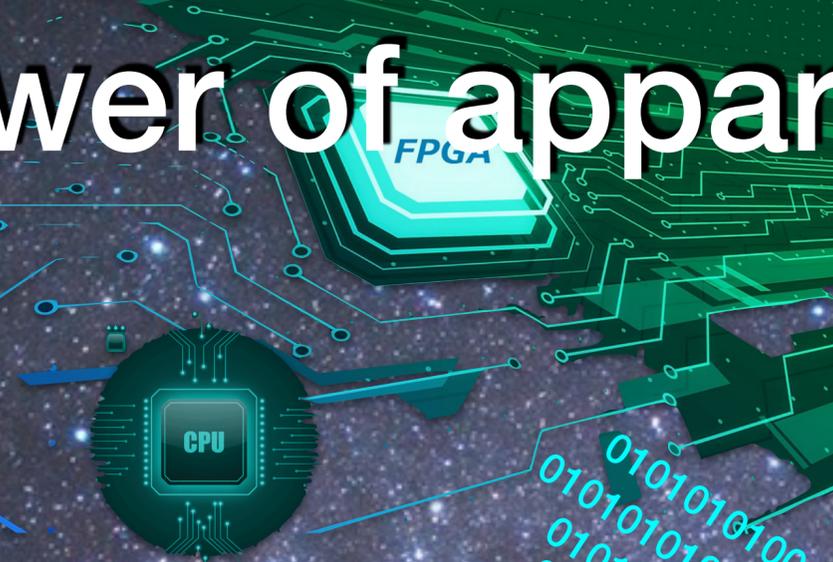
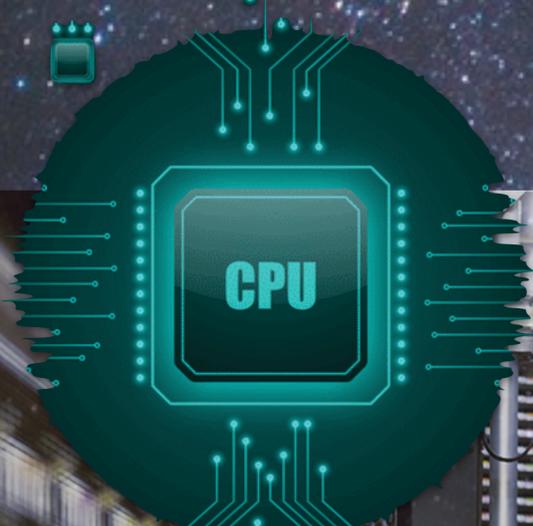


- HB:  $\sim 10 \text{ kHz}$
- Time Frame:  $\sim 50 \text{ Hz}$



# 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

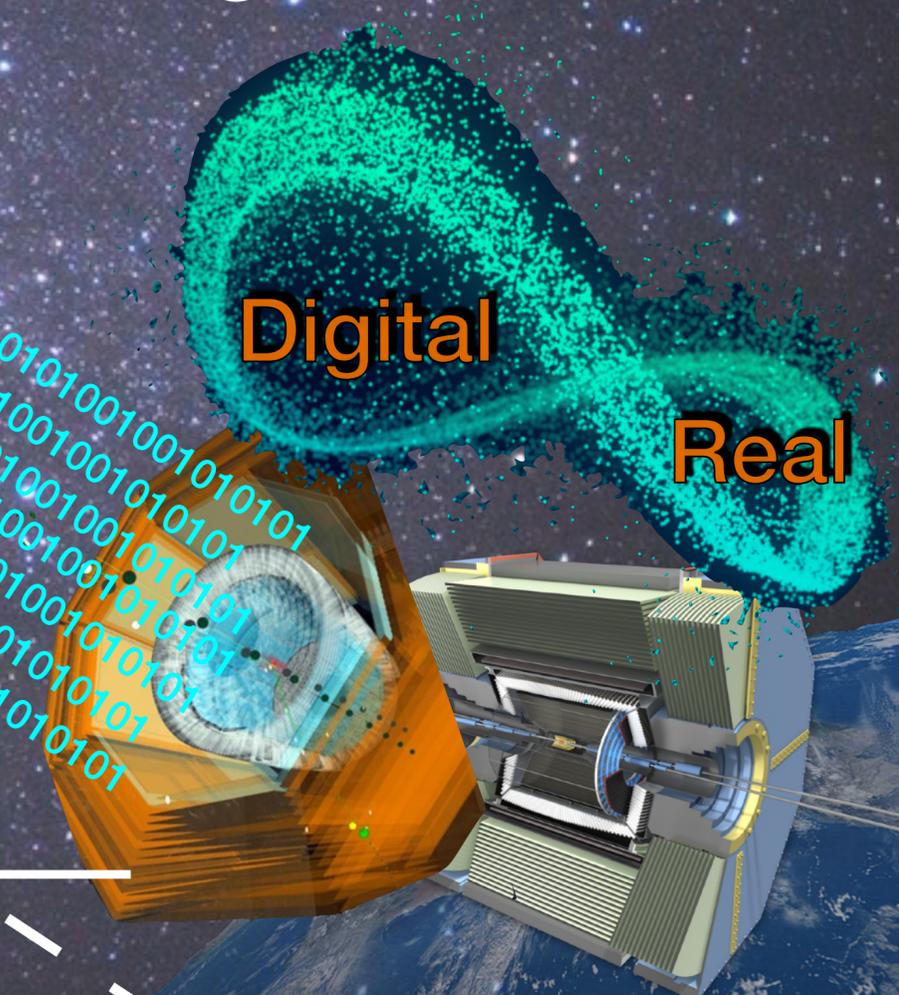
Data readout  
system

Decisions ↑

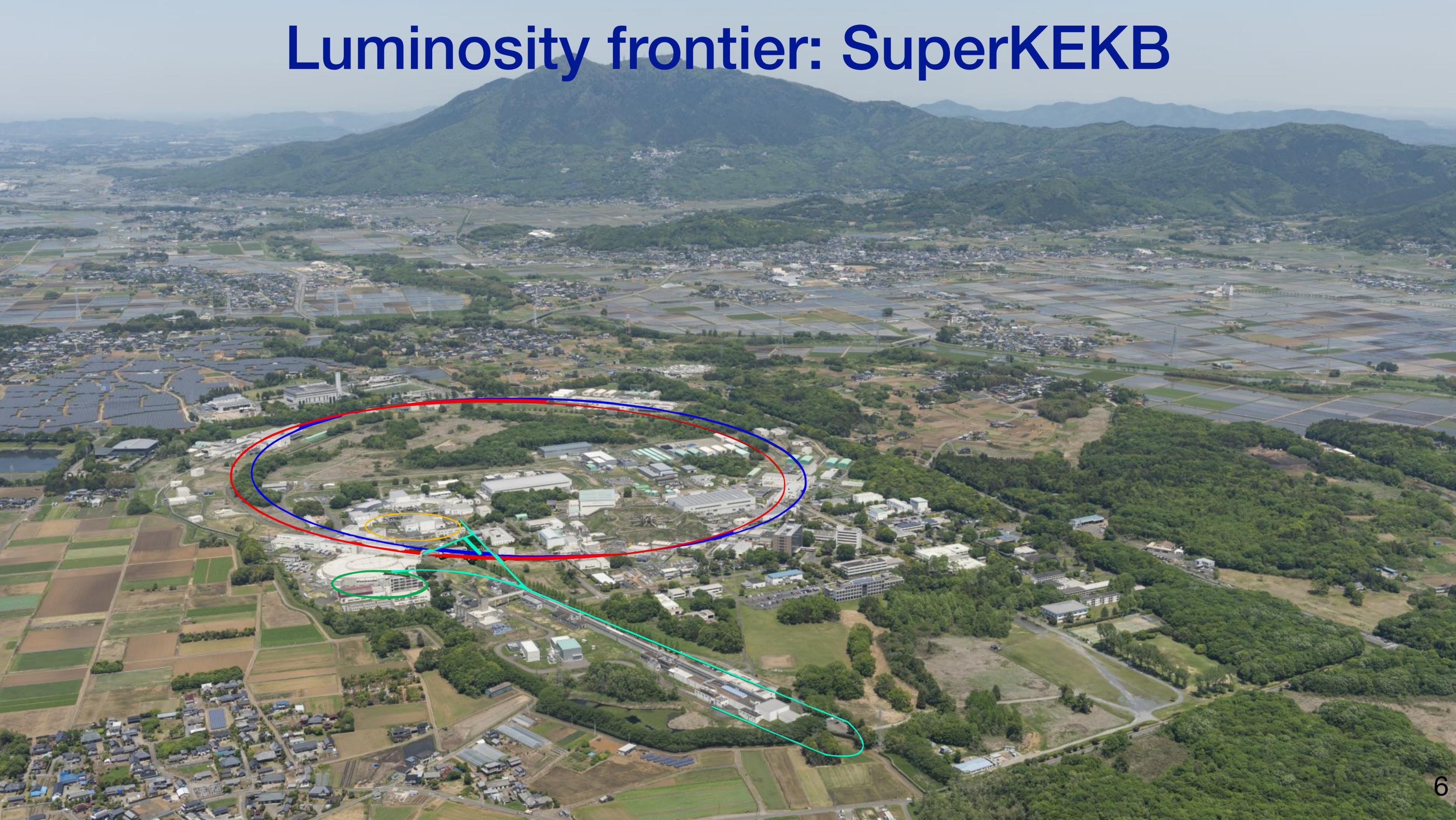
Trigger system (L1)  
(hardware filtering)

Digital

Real



# Luminosity frontier: SuperKEKB



# Luminosity frontier: SuperKEKB

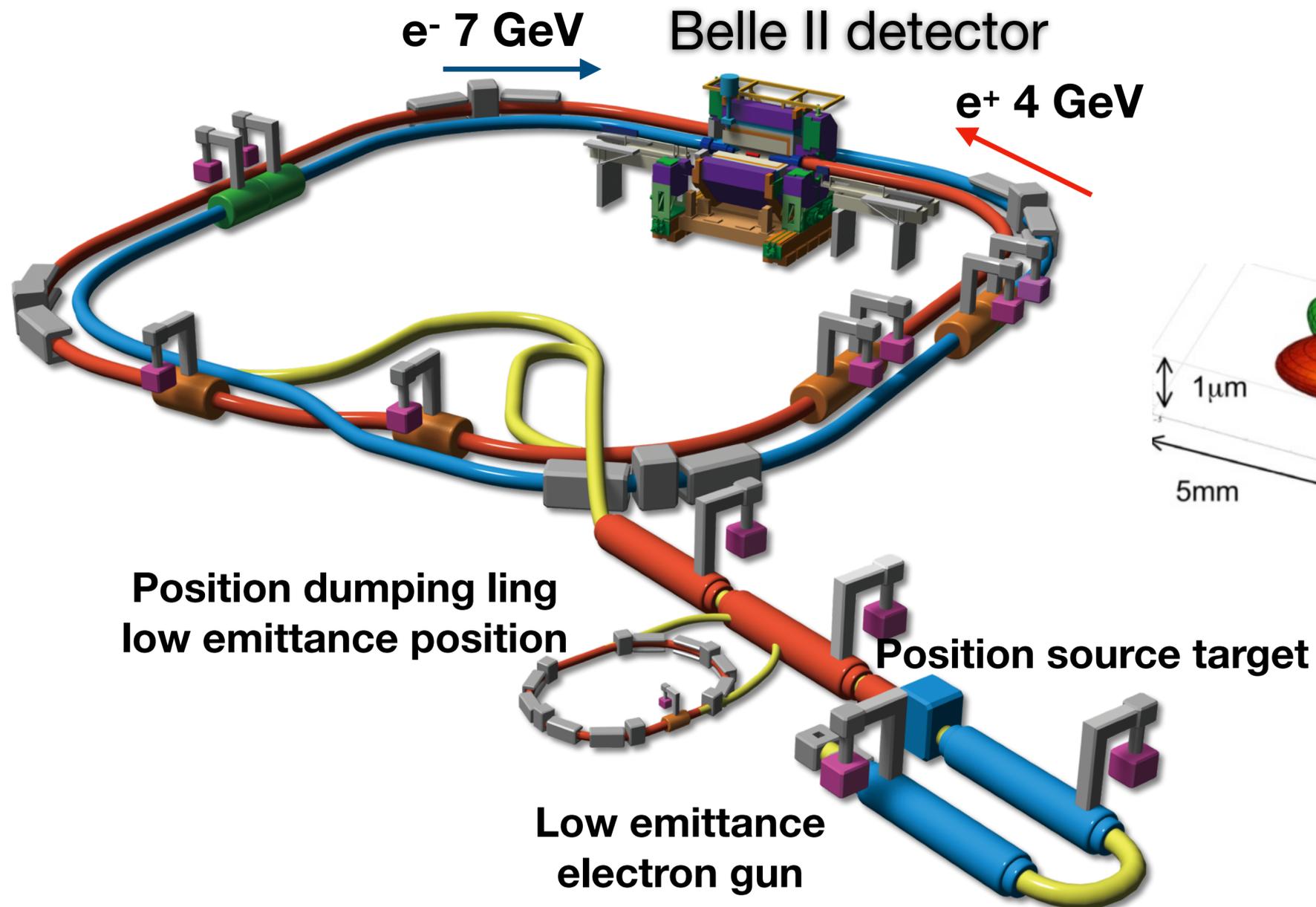
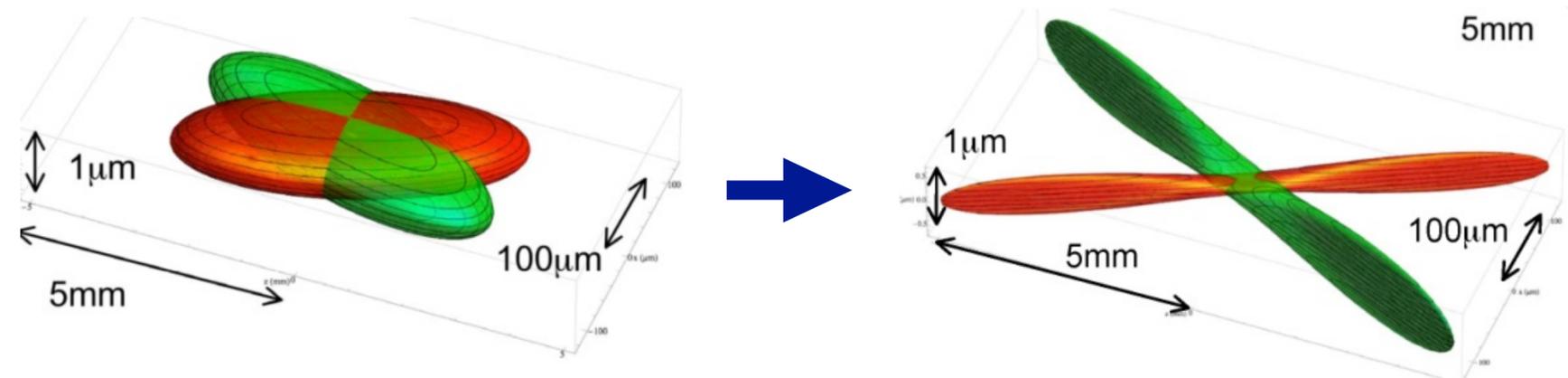
- Asymmetric  $e^+e^-$  collider
  - $e^+e^- \rightarrow \gamma(4S) \rightarrow B\bar{B}$
  - ▶ very clean and well-known initial state

Beam current: KEKB x ~1.5

$$L = \frac{\gamma_{\pm}}{2e r_e} \left(1 + \frac{\sigma_y^*}{\sigma_x^*}\right) \frac{I_{\pm} \xi_{\pm y}}{\beta_y^*} \left(\frac{R_L}{R_y}\right)$$

Beam squeeze: KEKB / ~20

## Nano beam scheme



**Target:  $L = 60 \times 10^{34} \text{ cm}^{-2} \text{ s}^{-1}$**   
**Achieved :  $5.1 \times 10^{34} \text{ cm}^{-2} \text{ s}^{-1}$  (Record)**

• Data:

- $575 \text{ fb}^{-1}$  (Belle II)  $\leftrightarrow$   $980 \text{ fb}^{-1}$  (Belle)

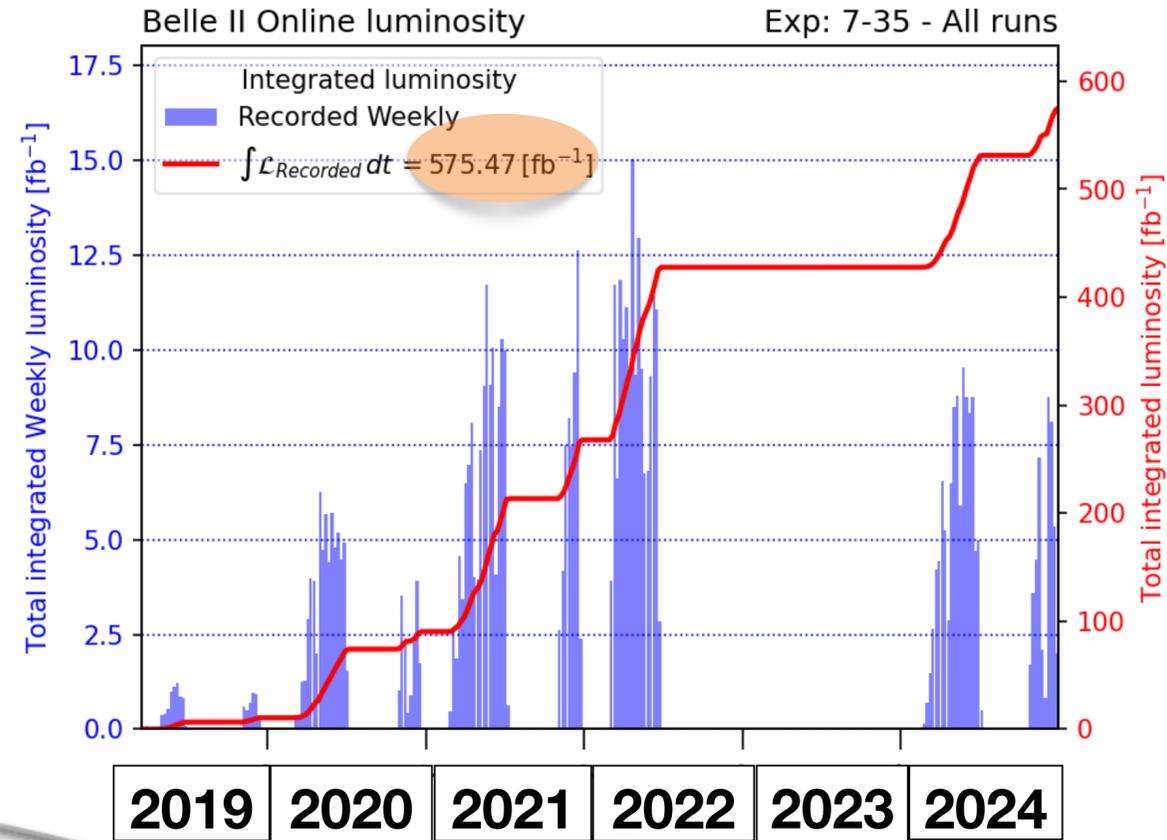
# Belle II detector and dataset

**Vertex detector (VXD)**  
 Inner 2 layers: pixel detector (PXD)  
 Outer 4 layers: strip sensor (SVD)

**Central Drift Chamber (CDC)**  
 He (50%), C<sub>2</sub>H<sub>6</sub> (50%), small cells, long lever arm

**Particle Identification**  
 Barrel: Time-Of-Propagation counters (TOP)  
 Forward: Aerogel RICH (ARICH)

**ElectroMagnetic Calorimeter (ECL)**  
 CsI(Tl) + waveform sampling



**e<sup>-</sup> (7GeV)**

**e<sup>+</sup> (4GeV)**

**K<sub>L</sub>/μ detector (KLM)**  
 Outer barrel: Resistive Plate Counter (RPC)  
 Endcap/inner barrel: Scintillator

- Features:
  - Near-hermetic detector
  - Vertexing and tracking:  $\sigma$  vertex  $\sim 15\mu\text{m}$ , CDC spatial res.  $100\mu\text{m}$   $\sigma(P_T)/P_T \sim 0.4\%$
  - Good at measuring neutrals,  $\pi^0$ ,  $\gamma$ ,  $K_L\dots$   $\sigma(E)/E \sim 2\text{-}4\%$

# Belle II trigger strategy

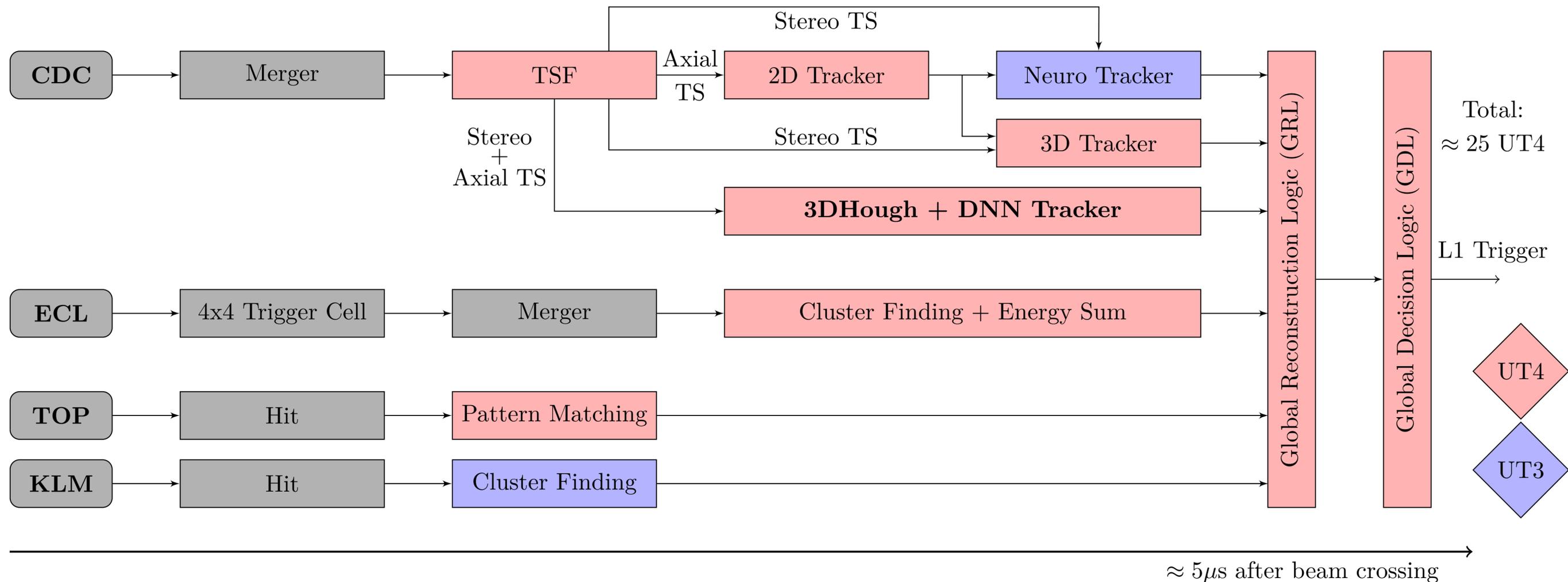
- Design requirements:  $\sim 100\%$  for  $\Upsilon(4S) \rightarrow BB$  (hadronic decay), Tau/Charm, Exotics
  - No dead-time  $\rightarrow$  pipeline
  - Single photon trigger
  - Single track trigger
- Max. trigger rate: 30 kHz @  $6 \times 10^{35} \text{ cm}^{-2} \text{ s}^{-1}$ 
  - Physics trigger  $\sim 15$  kHz
- Latency limit:  $\sim 5$  usec (SVD APV25 buffer structure)
  - A fixed latency of about 4.4 usec
- Event timing resolution: 10 nsec

Process	$\sigma(\text{nb})$	Rate@L= $6 \times 10^{35}$ (kHz)
Bunch. cross.	-	$2 \times 10^5$
Beam bkg	-	300-600
Bhabha	44	50
Total $\rightarrow$ L1	-	200350 $\rightarrow$ $\sim 15$

Process	$\sigma(\text{nb})$	L1@L= $6 \times 10^{35}$ (kHz)
Bhabha	44	0.35*
Two photon	13	10
Upsilon(4S)	1.2	0.96
Continuum	2.8	2.2
$\mu\mu$	0.8	0.64
$\tau\tau$	0.8	0.64
$\gamma\text{-}\gamma$	2.4	0.019*
Total	67	$\sim 15$

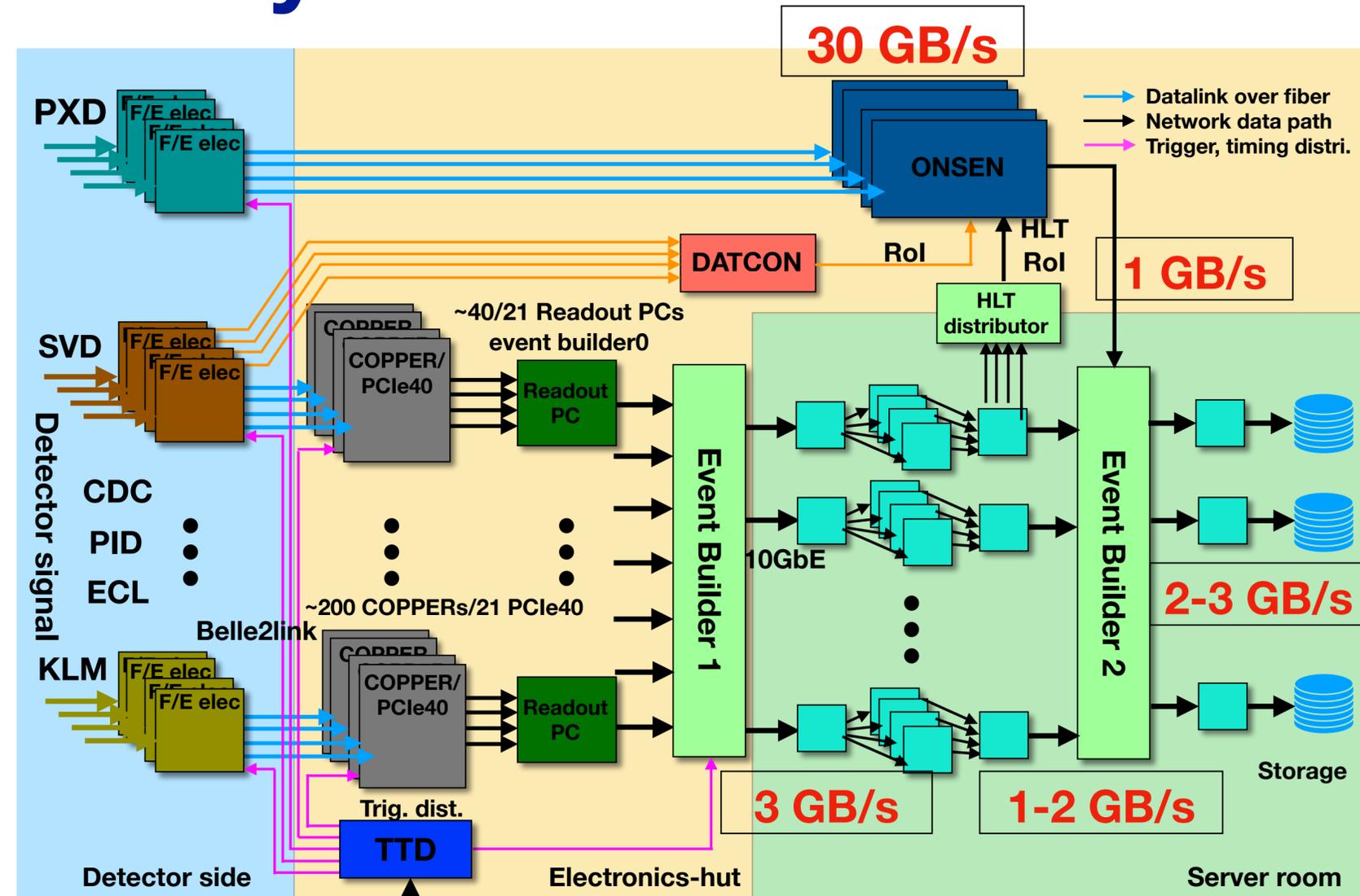
# Belle II trigger system

- CDC, ECL: main triggers for tracks and clusters
- KLM: trigger muon
- TOP: event timing
- GRL: matching of sub-triggers
- GDL: final trigger decision
- Challenges:
  - low multiplicity trigger vs. background
  - High track trigger vs. crosstalk
  - Drawback of track trigger at endcap
  - Latency budget vs. transmission and logics
  - ...

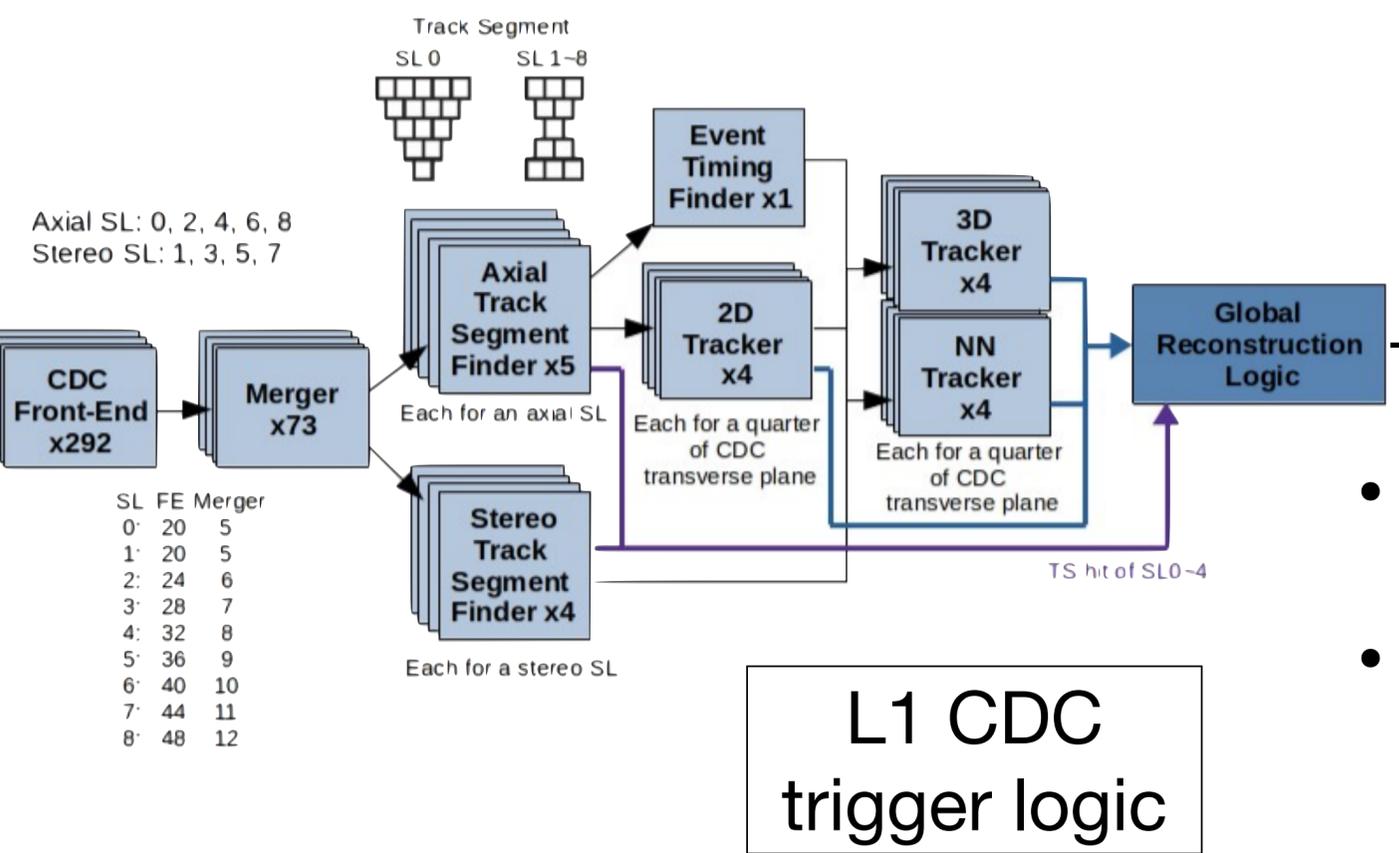


# Belle II TDAQ system

- Unified common readout system (except for PXD)
- Unified timing and trigger distribution (TTD) system
- A pipeline readout
- To handle 30 kHz level 1(L1) trigger with 0.1% dead time under raw event size of 1 MB



## Example: CDC

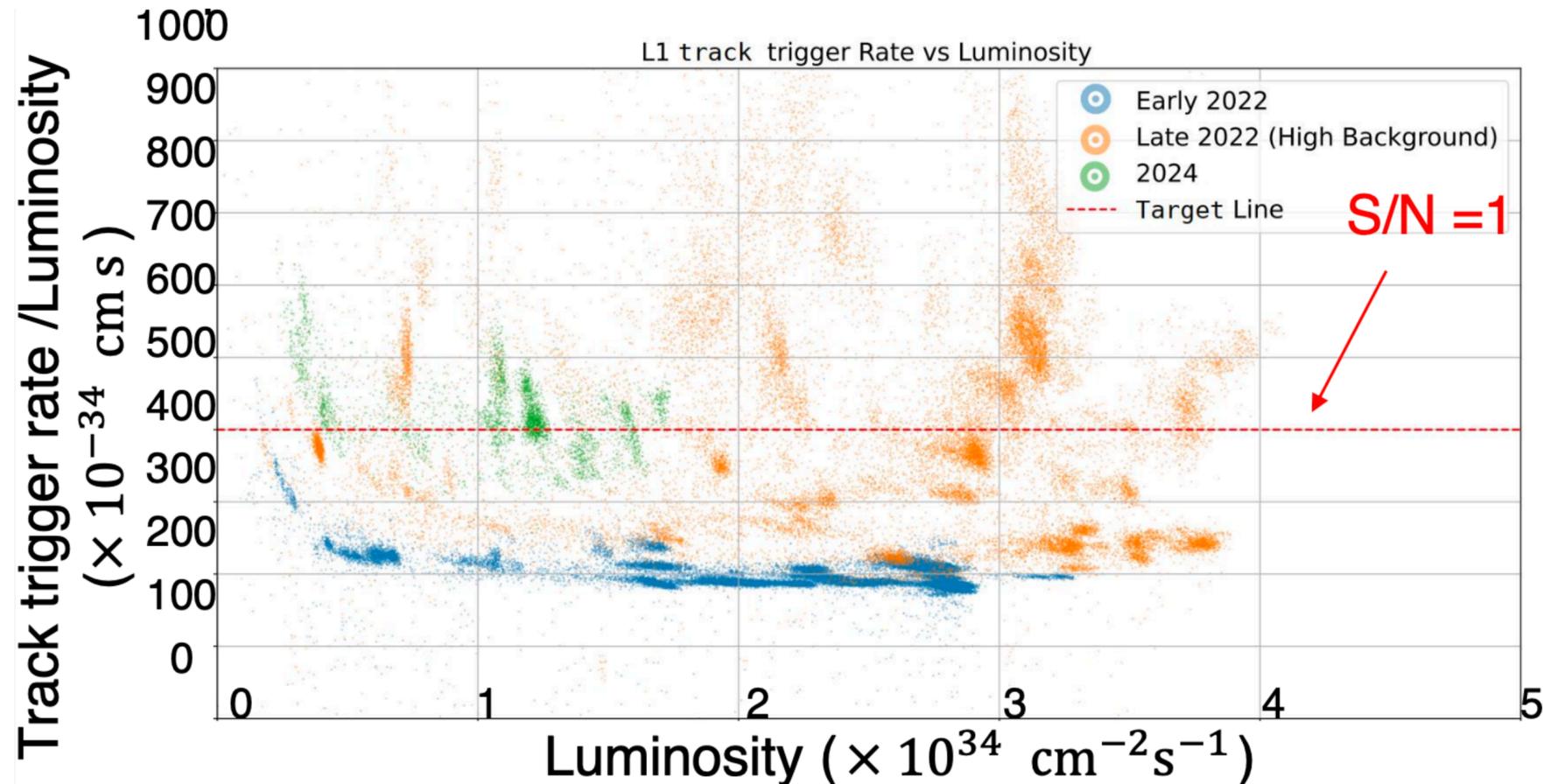
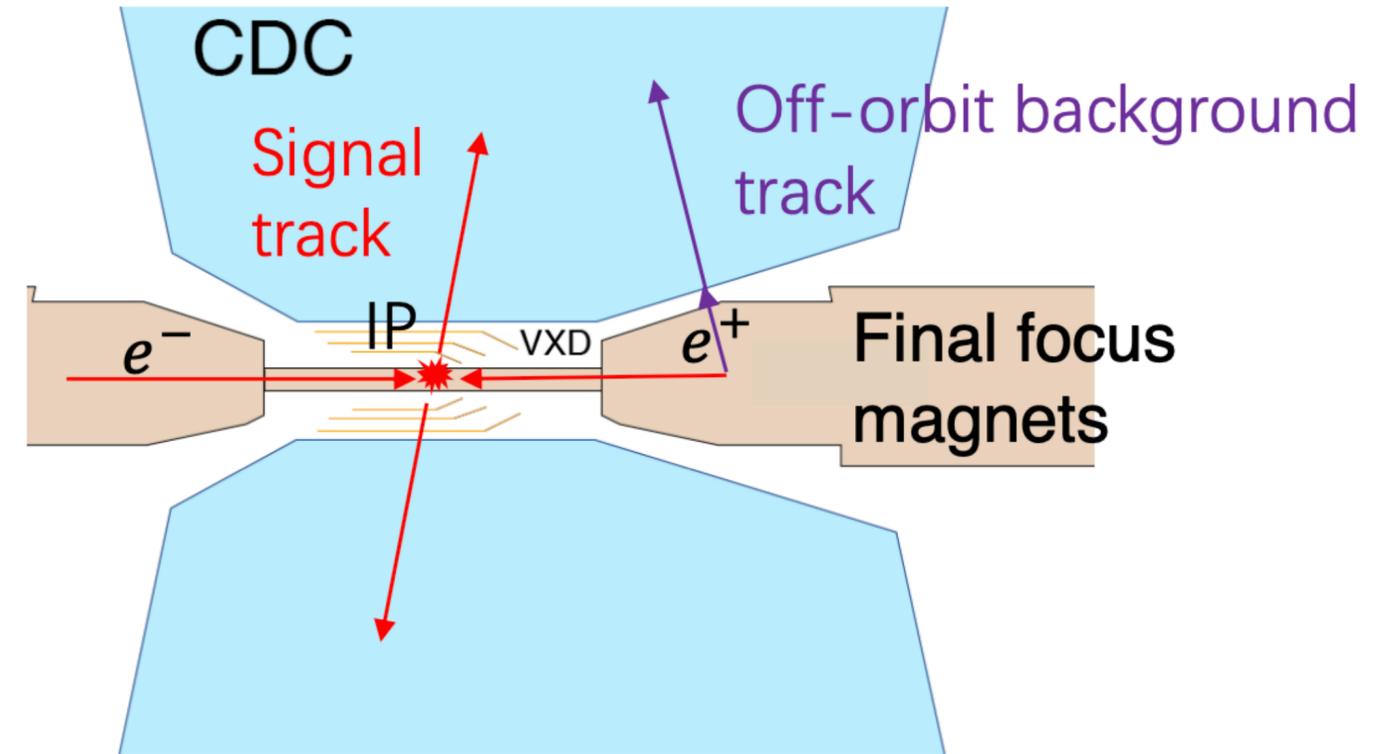


## L1 CDC trigger logic

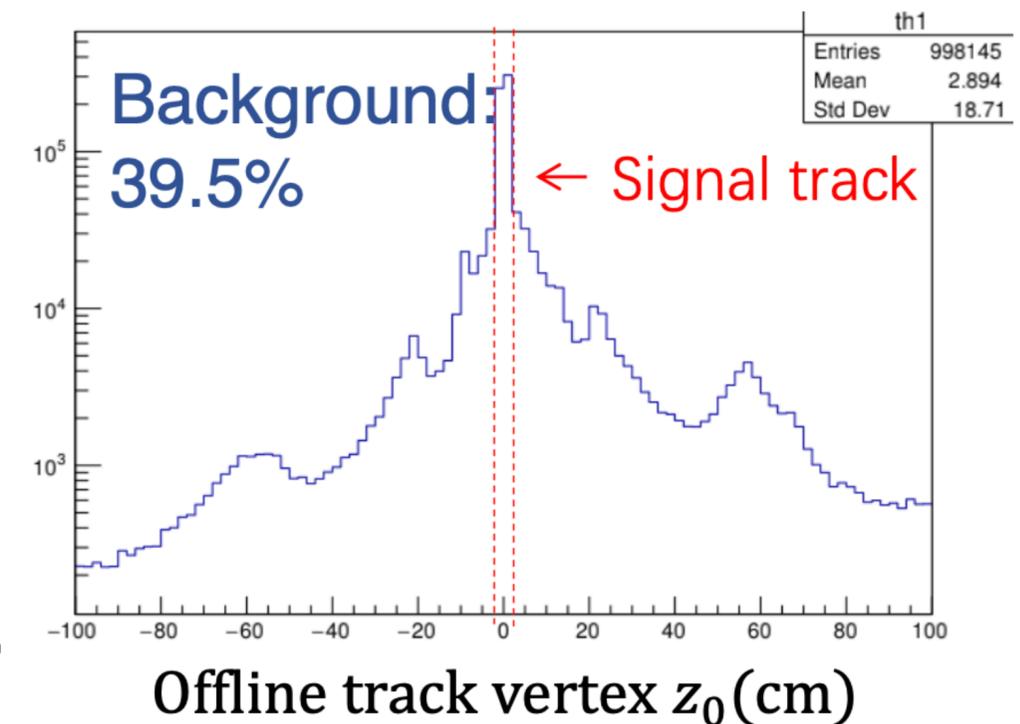
- Provide L1 trigger signal to DAQ using FPGA chips for real-time processing on detector raw data.
- HLT provide Region of Interest (RoI) to PXD for significantly reducing the data size.
  - Latency 0 sec.

# Motivation of Neural Network for L1 Track trigger

- DAQ system is designed to handle 30 kHz
  - Physical trigger  $\sim 15$  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

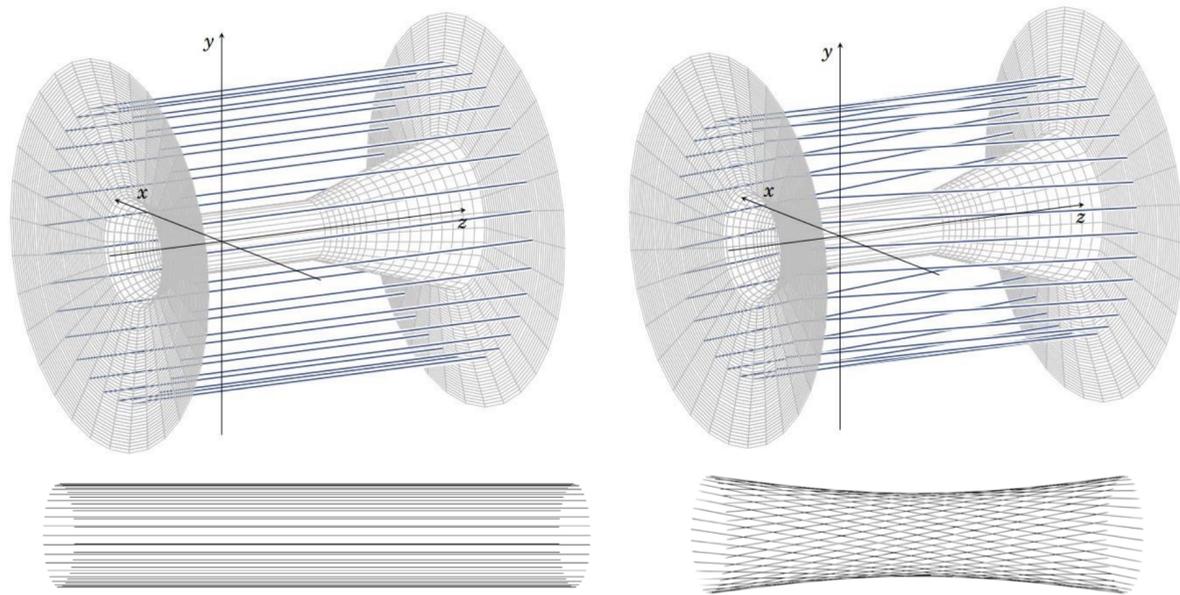


### Tracks $z_0$ distribution after trigger



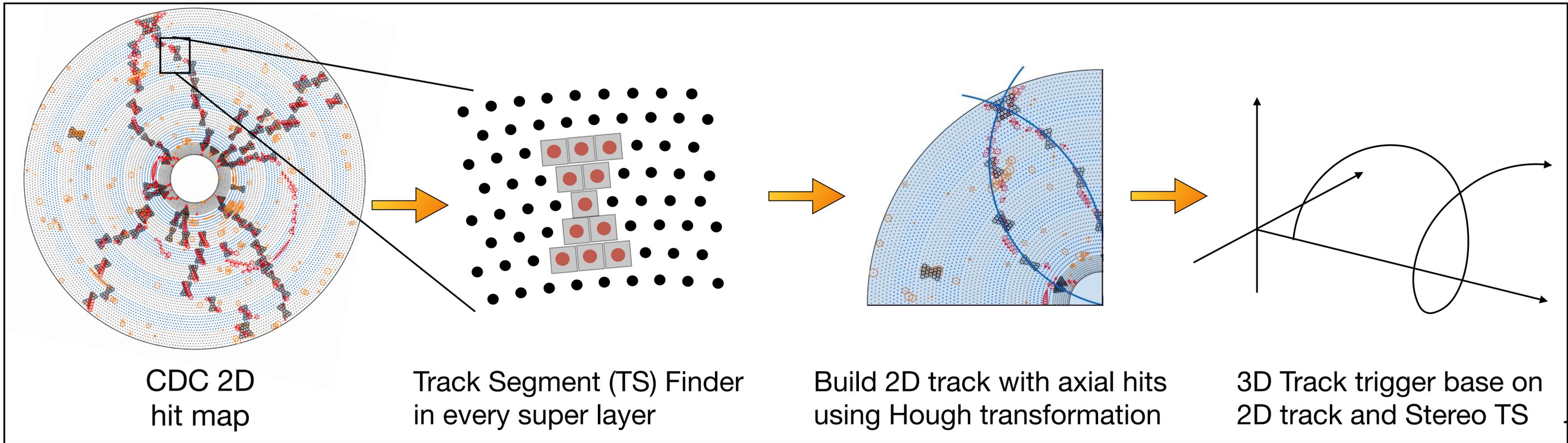
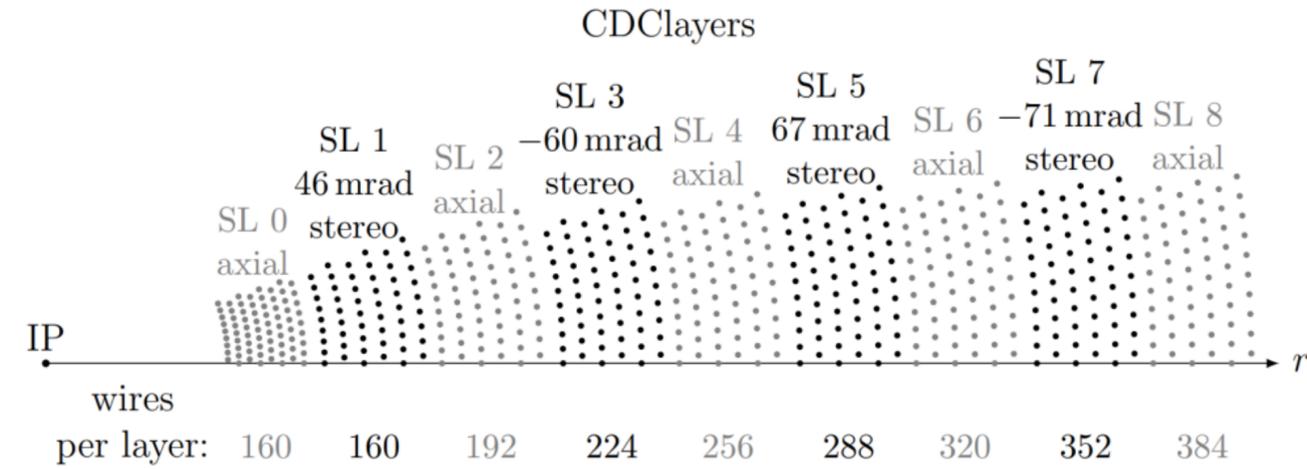
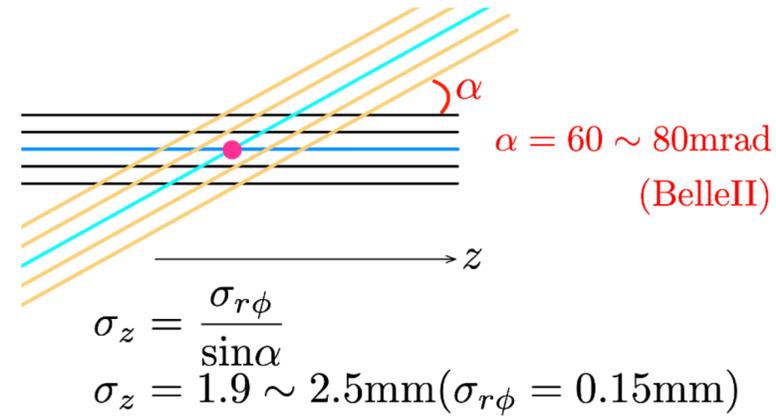
# Machine Learning for L1 Track trigger (HARDWARE)

# Basics of L1 CDC trigger



Axial wire

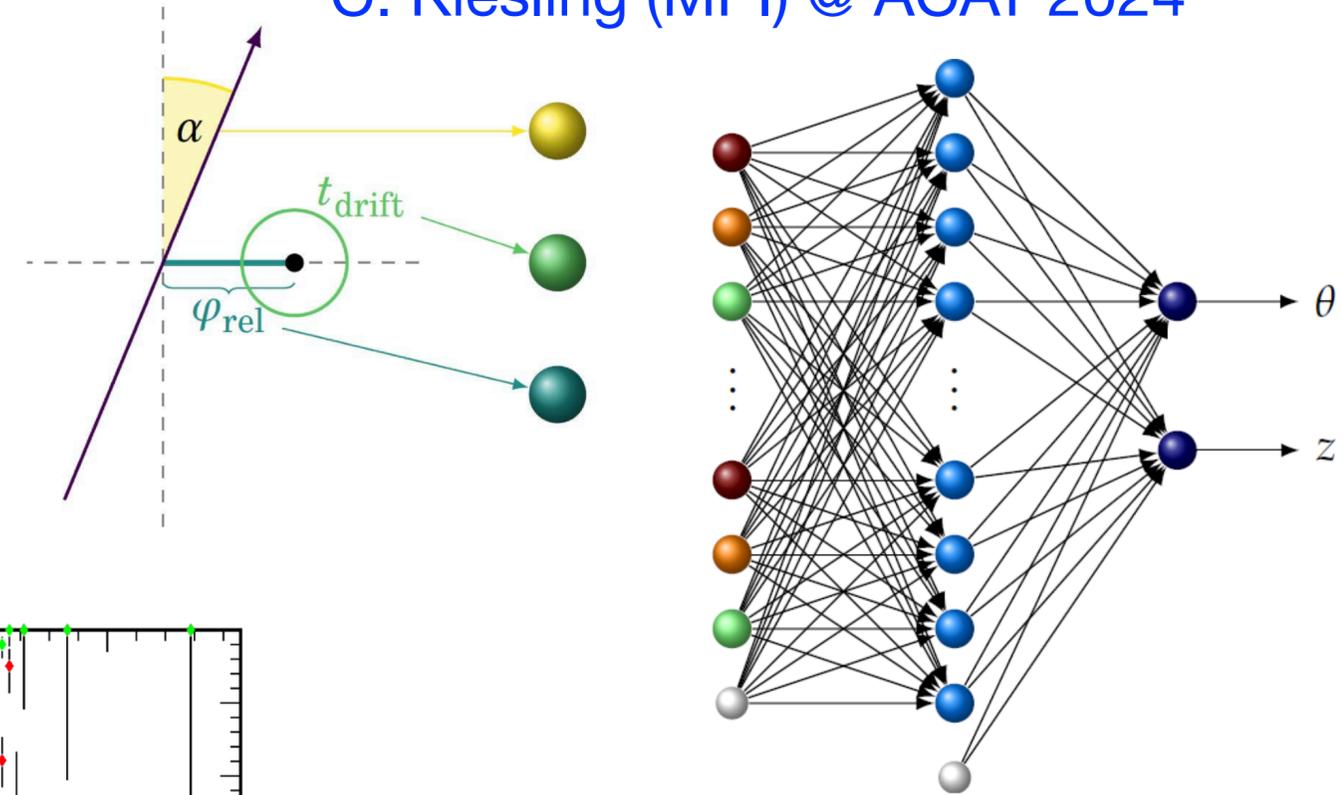
Stereo wire



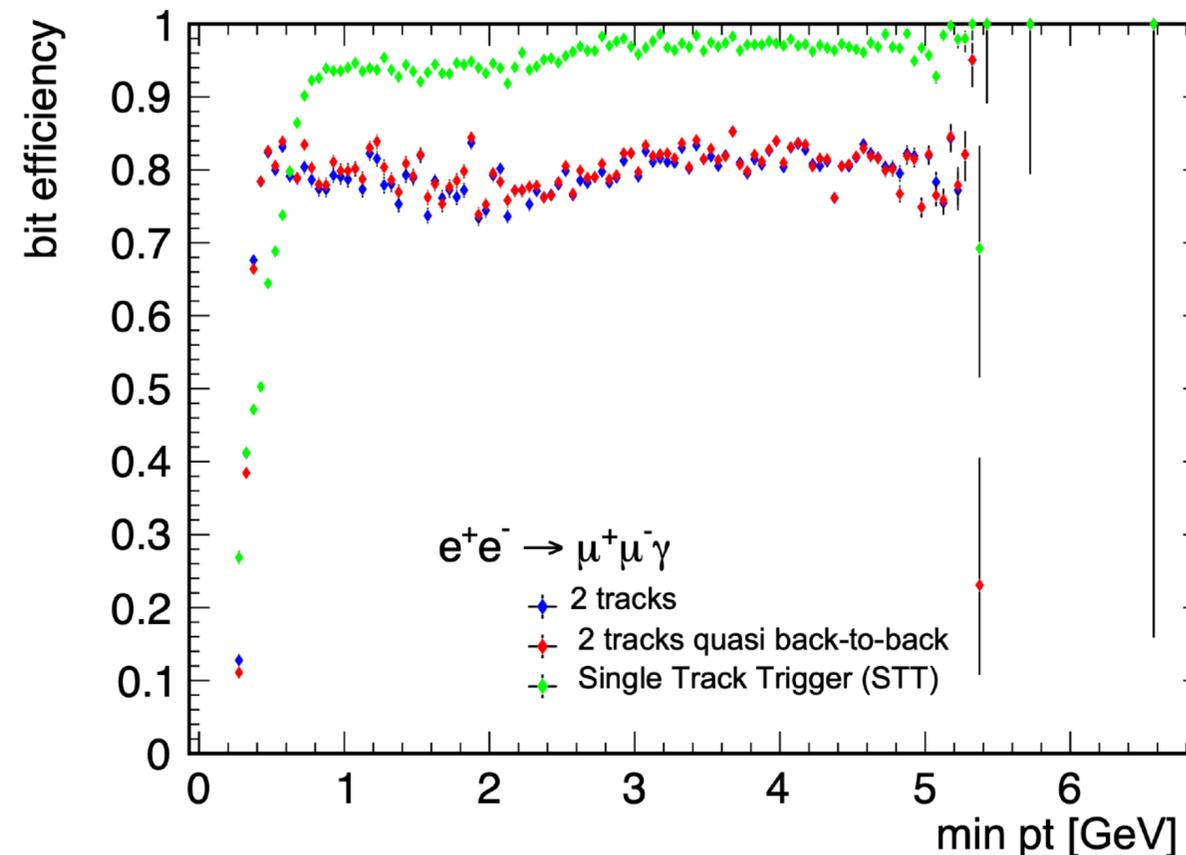
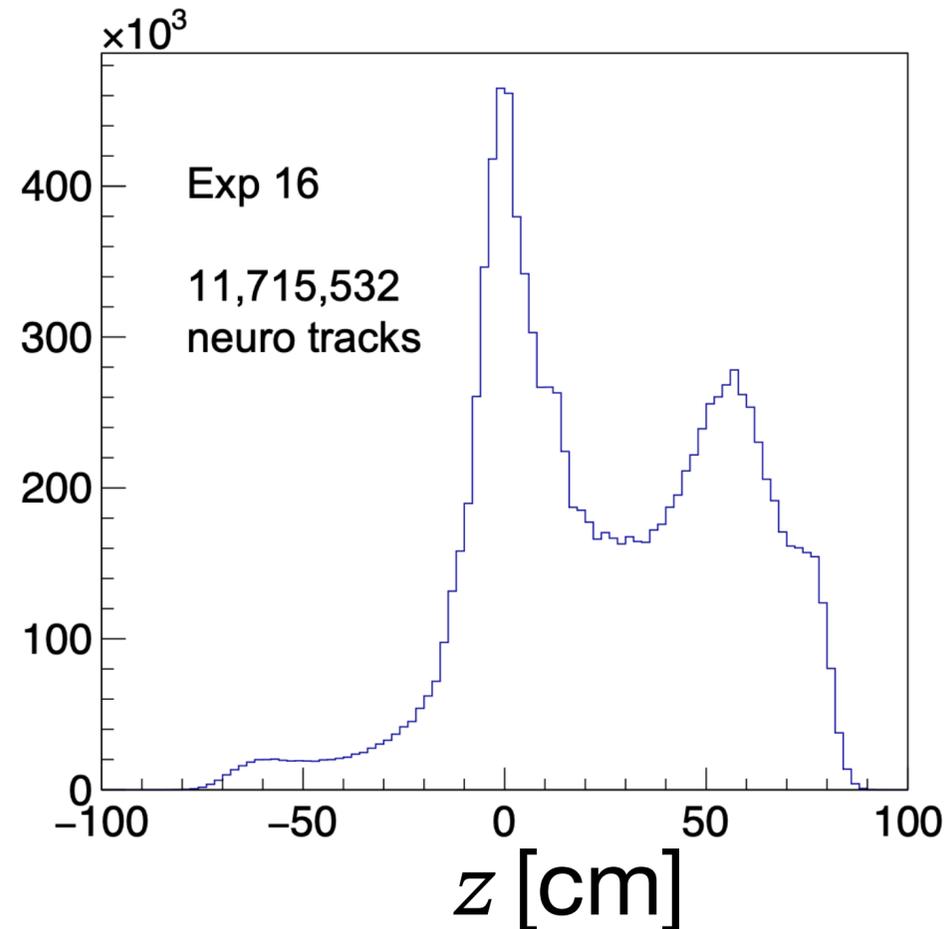
# Neural Network z-trigger

- Pioneer of NN on **HARDWARE** for Belle II trigger
- Inputs: **Drift time**  $t_{\text{drift}}$ , **wires relative location**  $\phi_{\text{rel}}$ , **Crossing angle**  $\alpha$  for priority wires
- Outputs: track vertex  $z_0$ , track  $\theta$
- Selected 1 Track Segment per one Super Layer
- Networks trained with real data from May-June 2020

C. Kiesling (MPI) @ ACAT 2024



[arXiv:2402.14962](https://arxiv.org/abs/2402.14962) (submitted to NIMA)

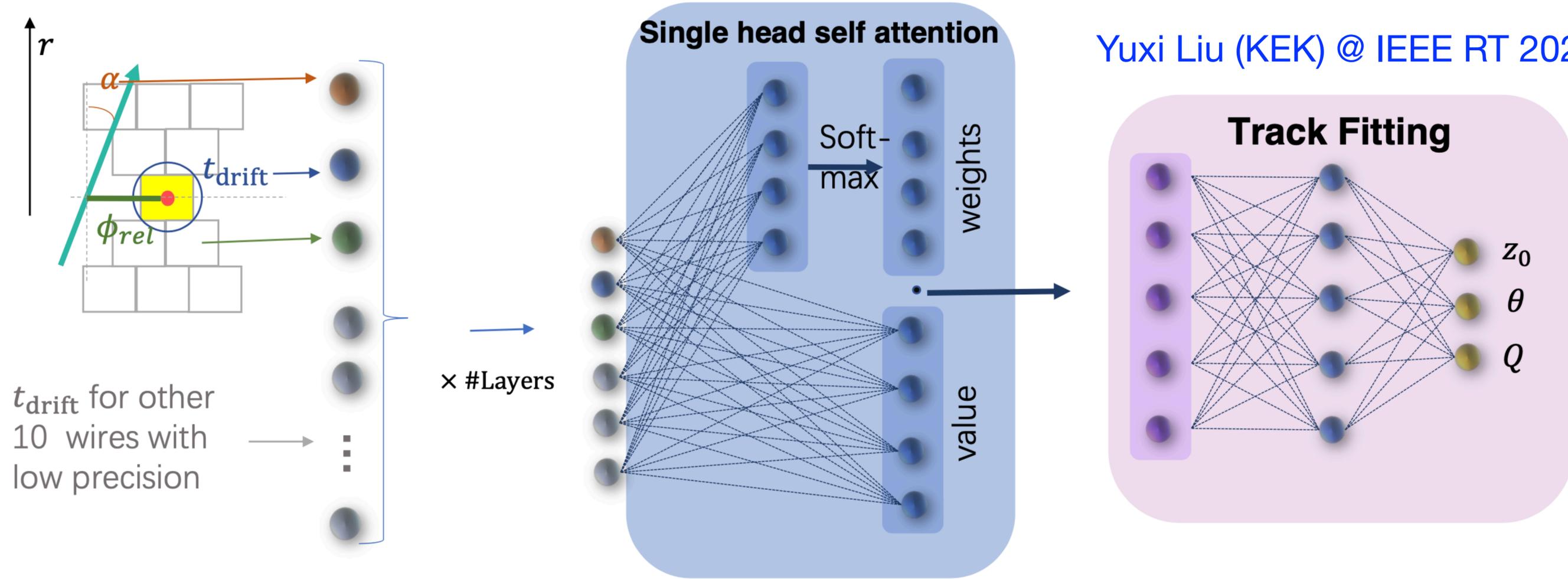


- Backgrounds still high
- Fake background

- Upgrade plan in trigger group:
- Track finding in **3D Hough space**
- network architecture: “deep-learning” + additional inputs

# Deep Neural Network

Yuxi Liu (KEK) @ IEEE RT 2024



- 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



- Machine Learning model
- Parameter



- C/C++ transition



- Translate into Verilog/VHDL FPGA language

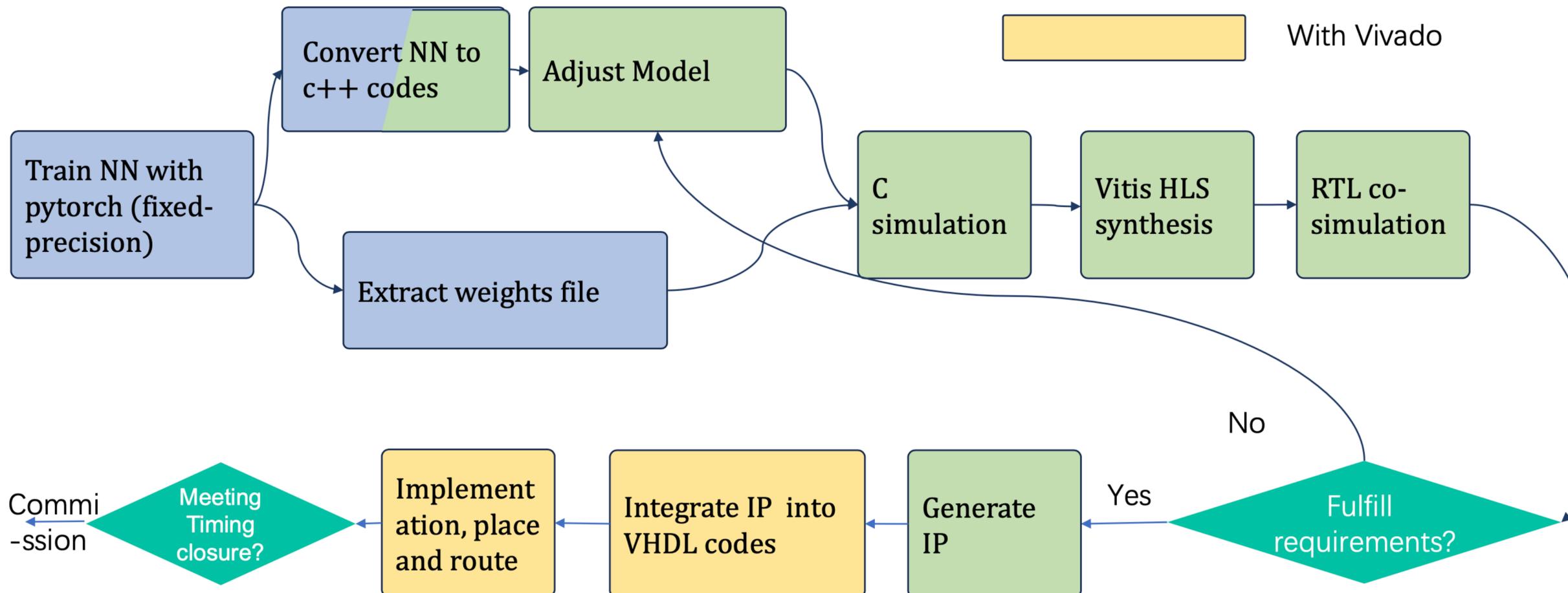


- Start fitter



- Evaluation

\*include some function from hls4ml lib



Belle II UT4

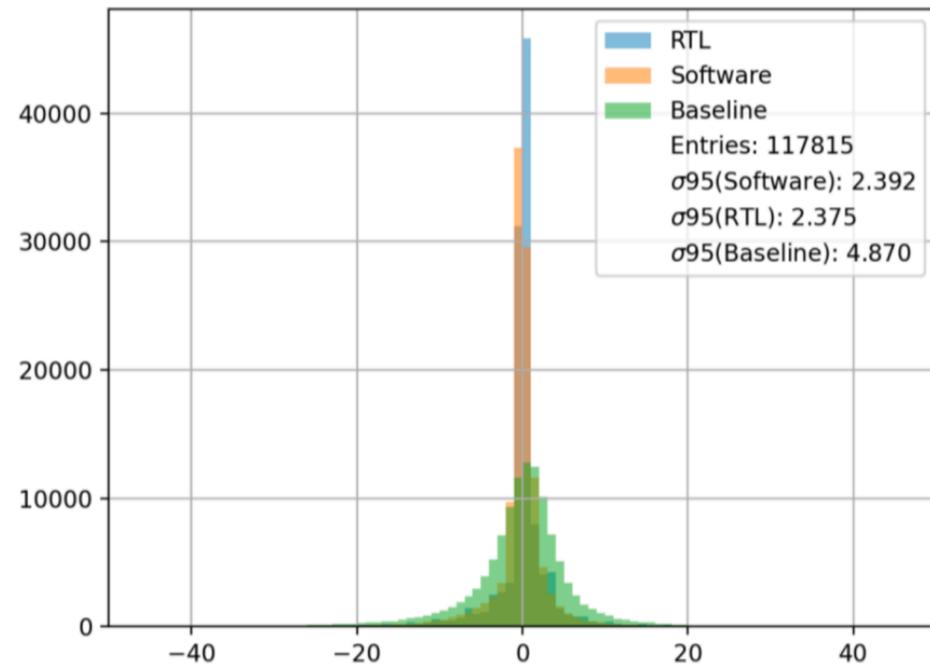


Xilinx UltraScale XCVU080, XCVU160 25 Gbps with 64B/66B

# Simulation performance of DNN

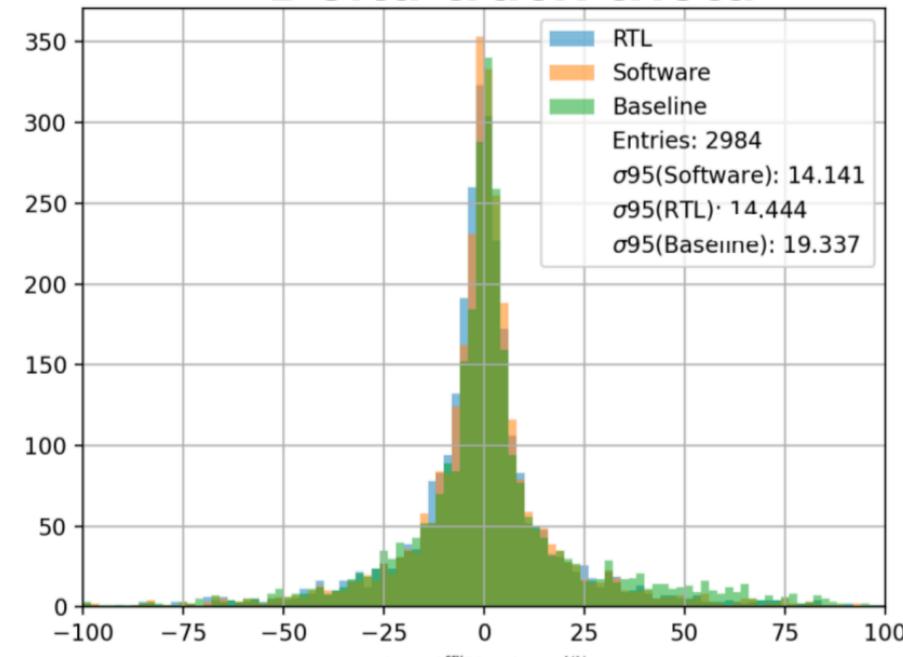
Yuxi Liu (KEK) @ IEEE RT 2024

### Delta track z



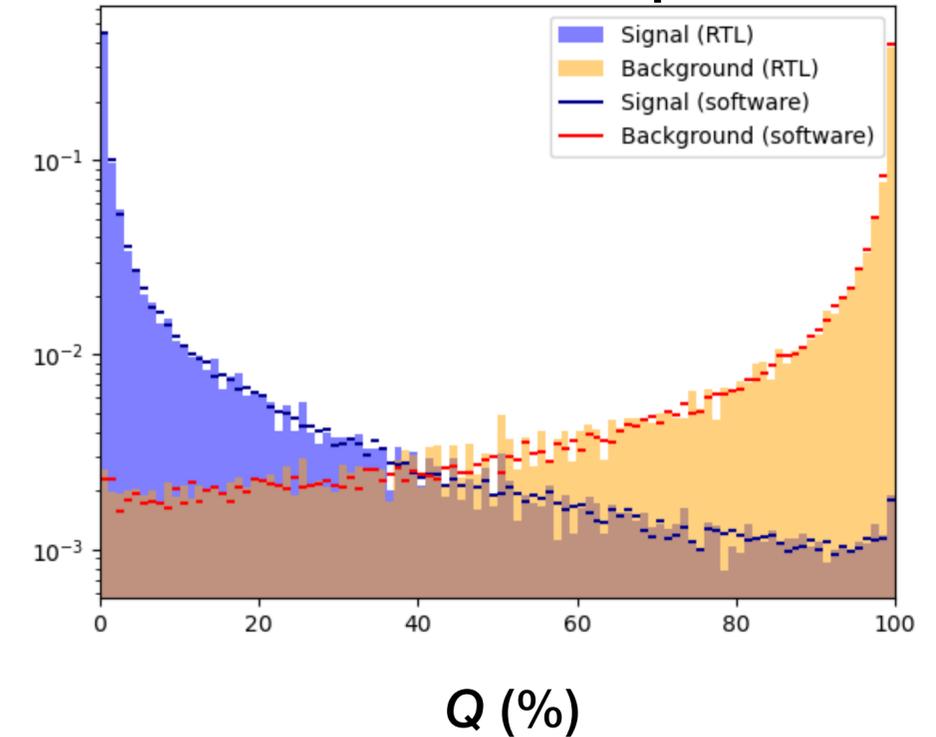
$$z_0^{NN} - z_0^{offline} (cm)$$

### Delta track theta



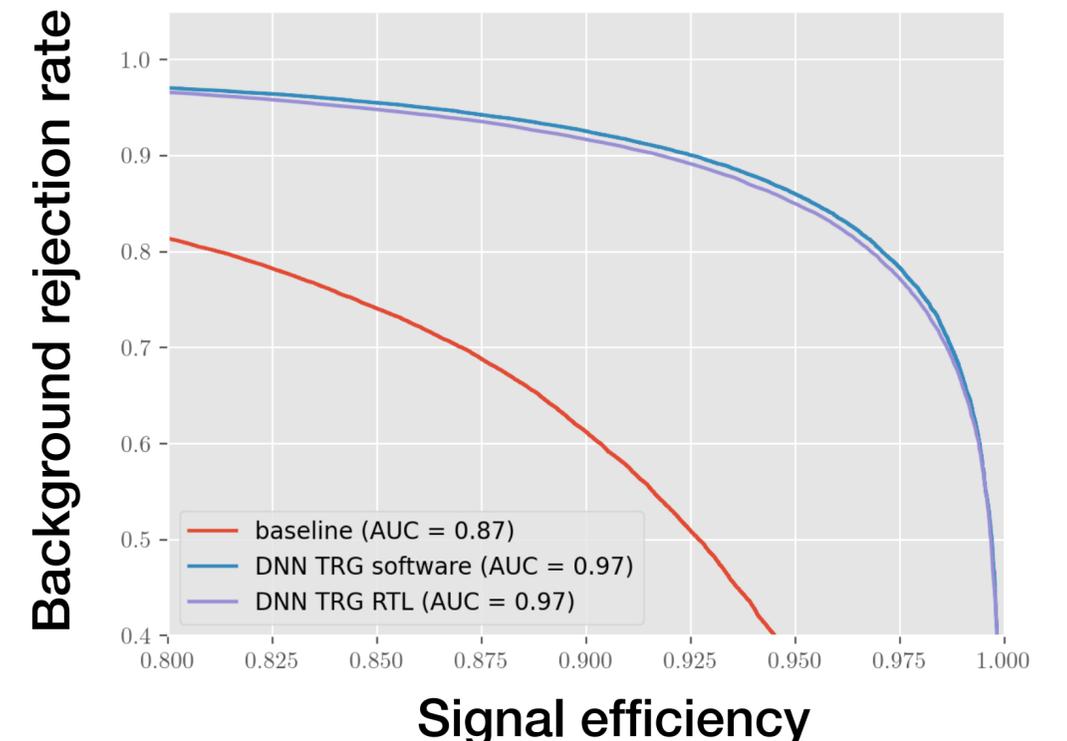
$$\theta_0^{NN} - \theta_0^{offline} (^\circ)$$

### Classifier output



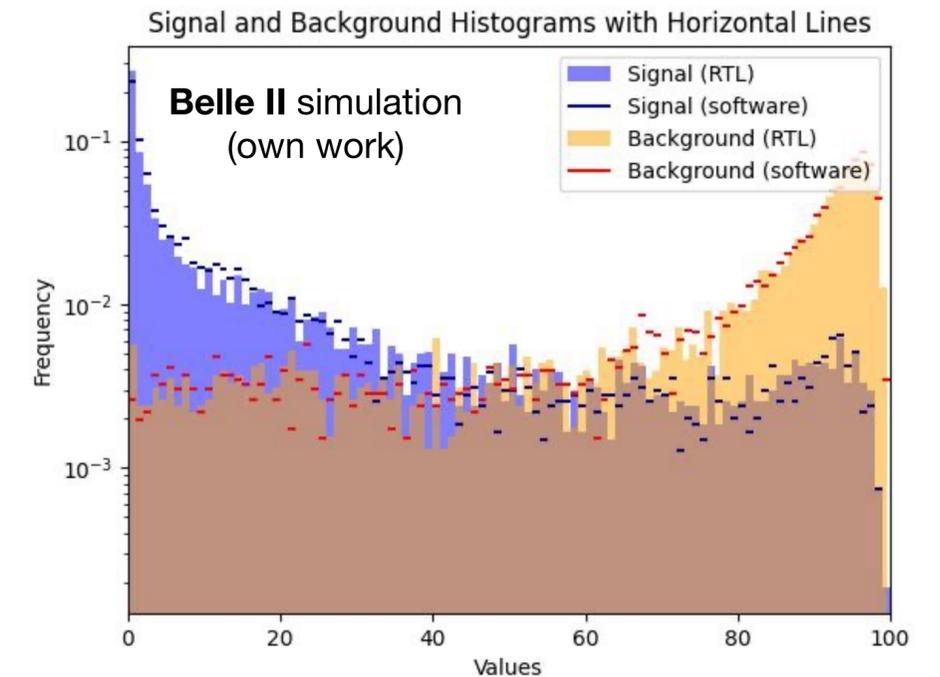
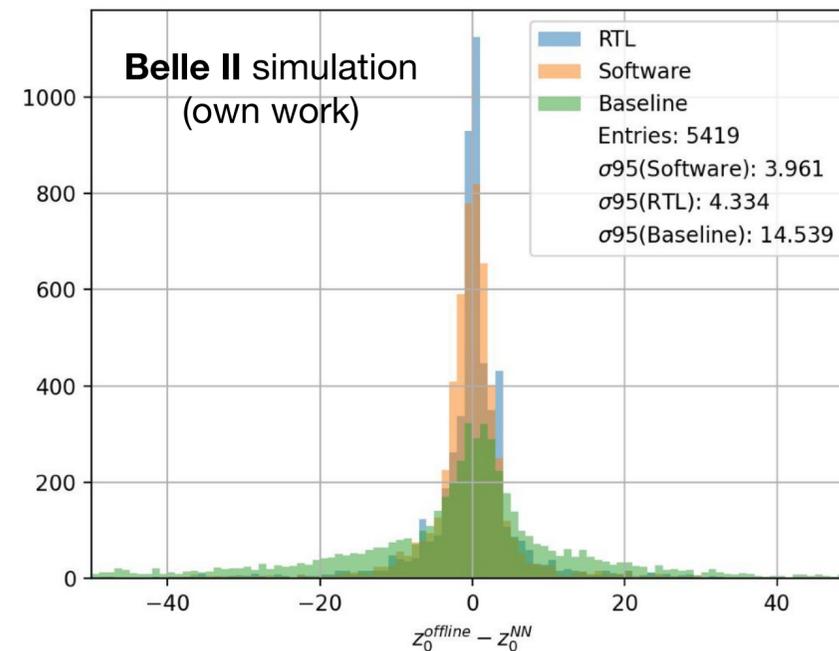
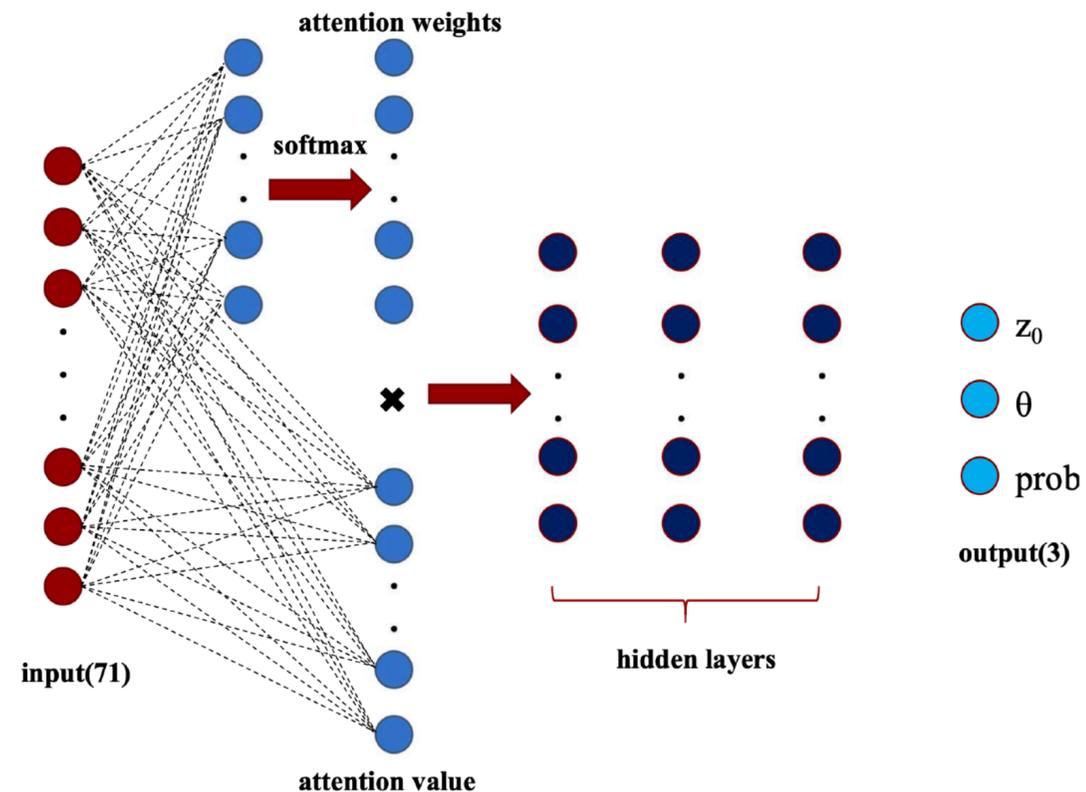
- Latency : 76 clock = 592.8 ns ;require: < 600ns
- FPGA resource (UT4: Virtex UltraScale XCVU160) usage:
  - DSP: ~70%, LUT: ~50%, others <30%
- AUC do not get large drop comparing RTL and software simulation
- At signal efficiency ~95%
  - Background rejection rate ~85%

- DNN trigger with **HARDWARE** under commissioning, close to operate



# Improvement try for CDC track trigger

- Develop a algorithm improve the performance for the upgrade (10 usec latency)
  - Start from optimization of DNN model
- Modify the number of hidden layers and learning rate
  - Hidden layer: 2 -> 4, learning rate:  $1e^{-2}$  ->  $1e^{-3}$
  - Others keep the same
  - Latency: 76 clock (592.8 ns) -> 82 clocks (640 ns)
- Next step, change the inputs (CDC hits info.), instead of 2D track parameters



# DNN implementation on Versal ACAP

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

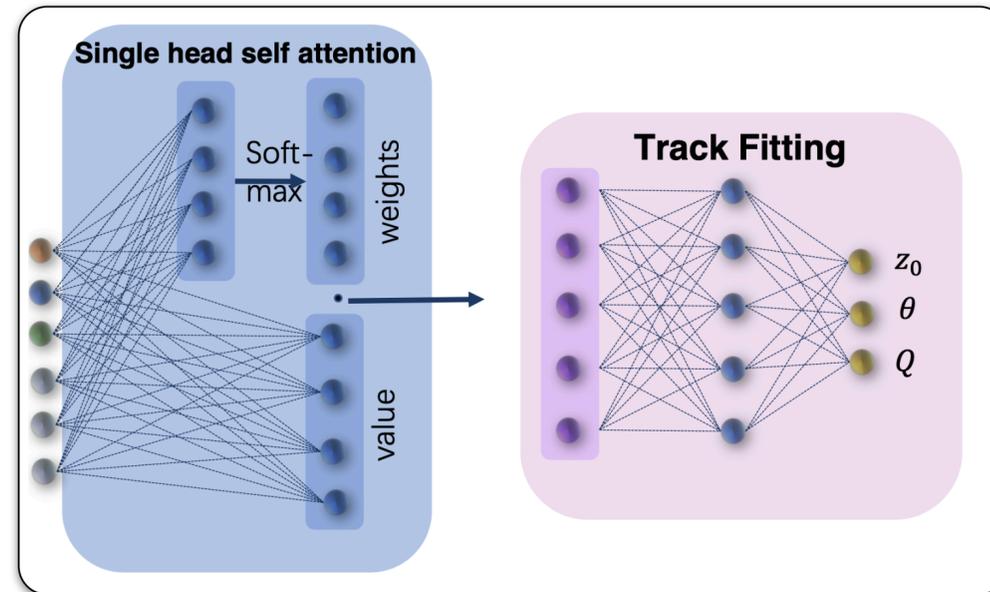


Figure 2: AI Engine Array

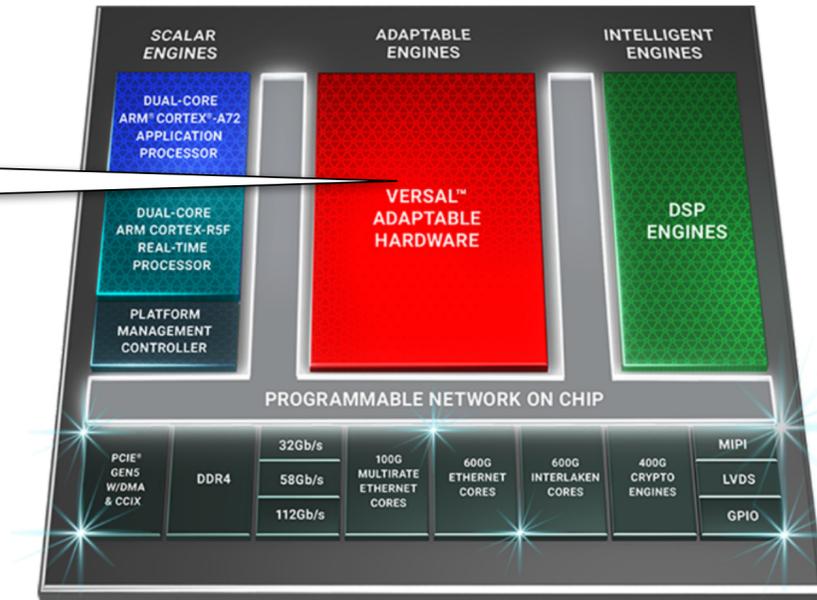
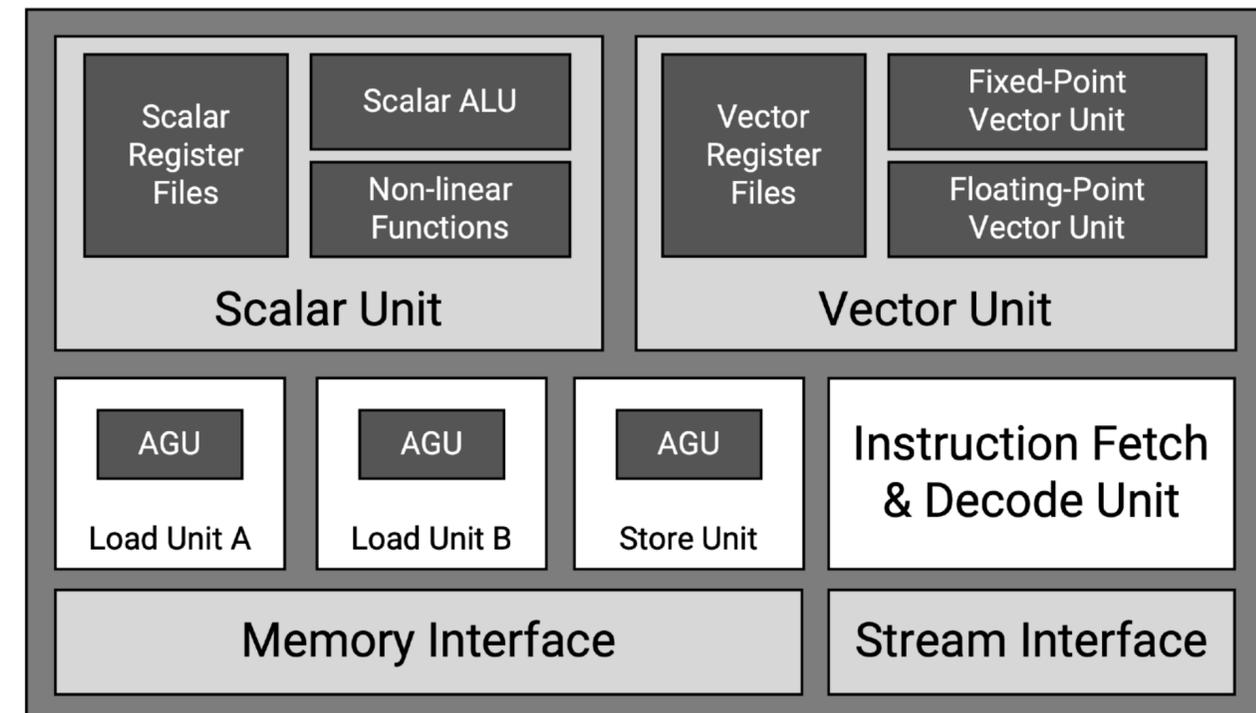
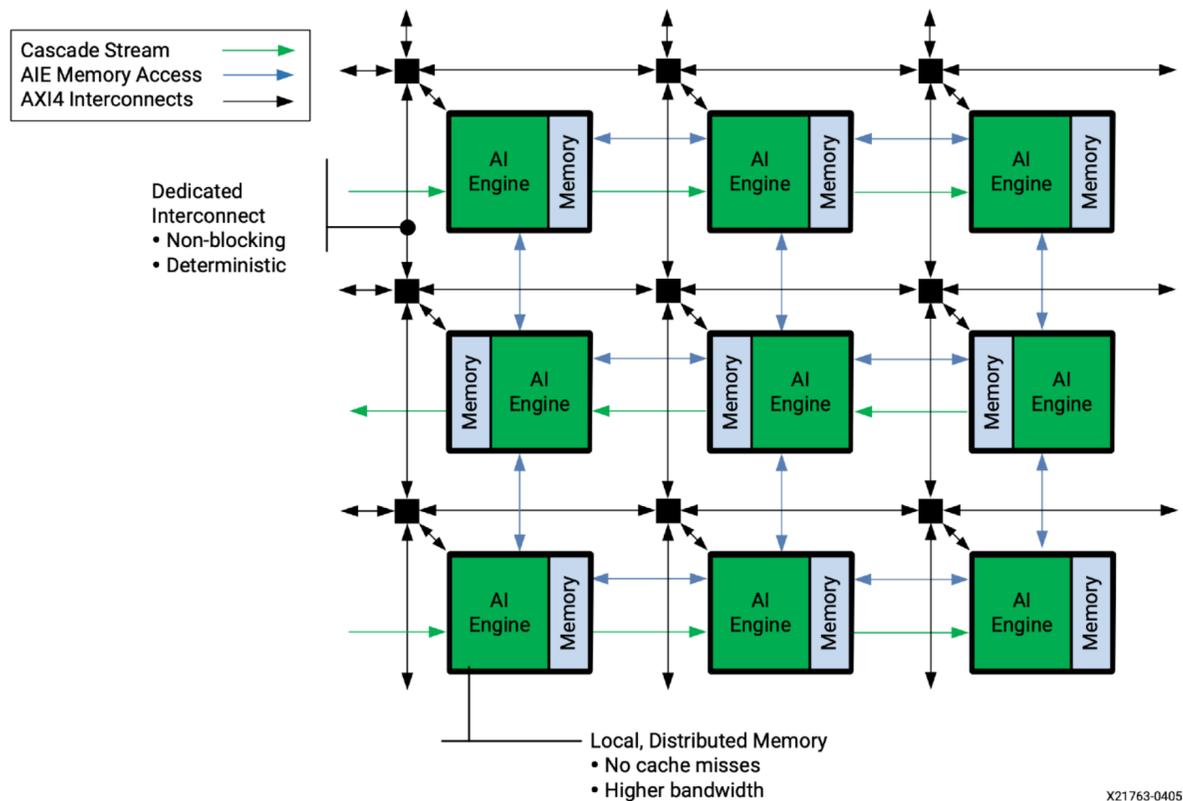
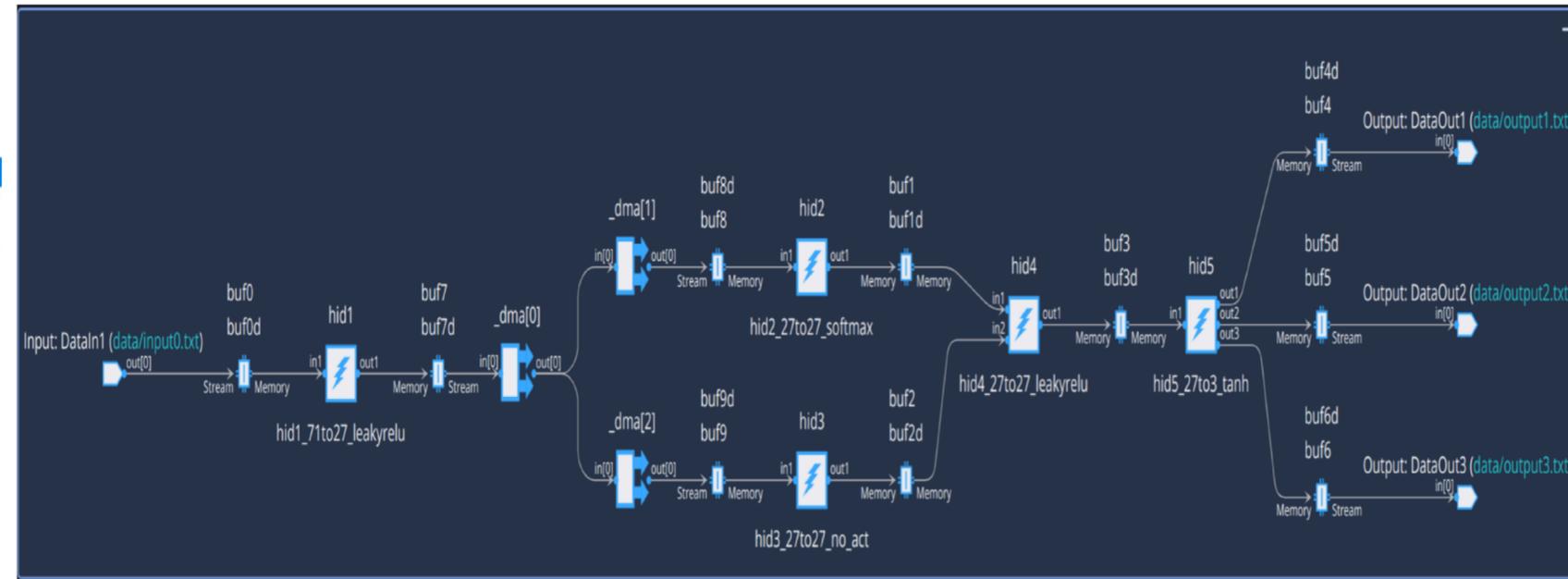
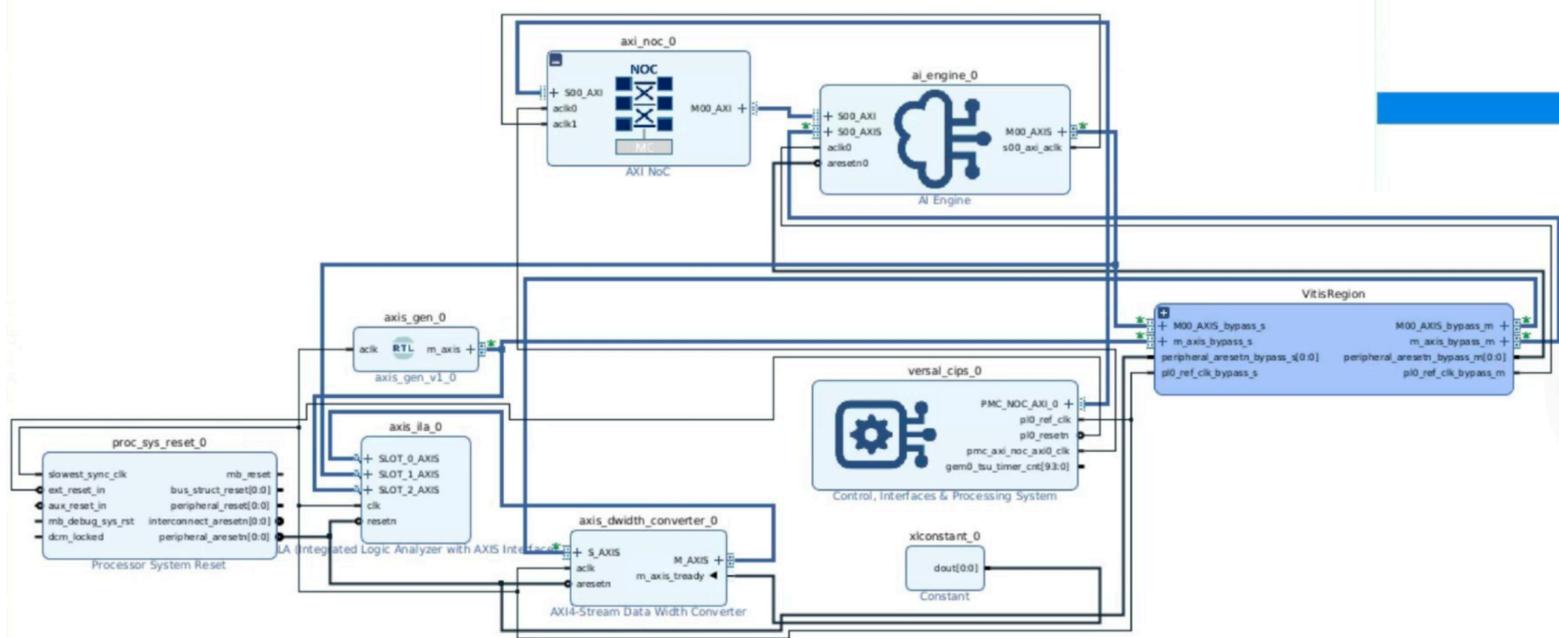


Figure 4: AI Engine

UG1079



# DNN acceleration on Versal ACAP



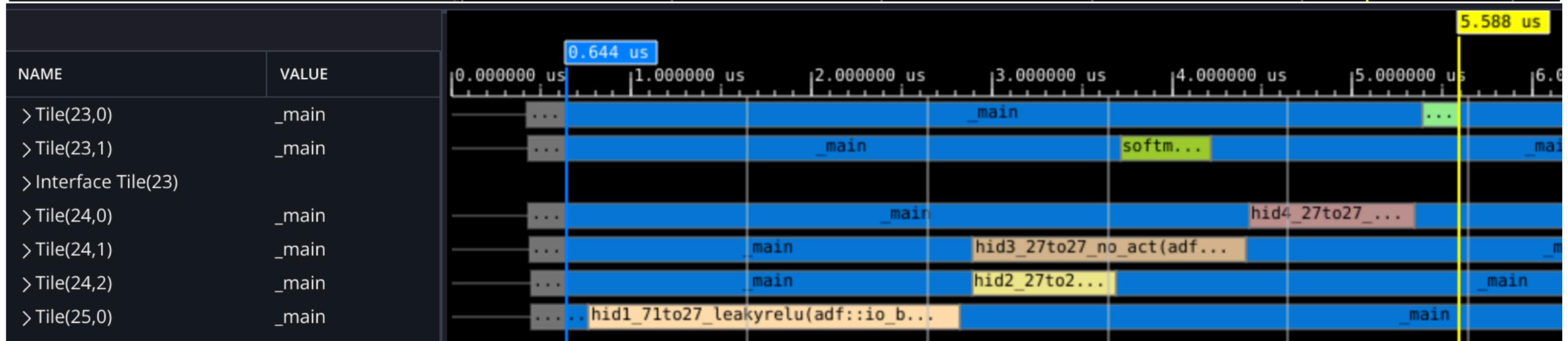
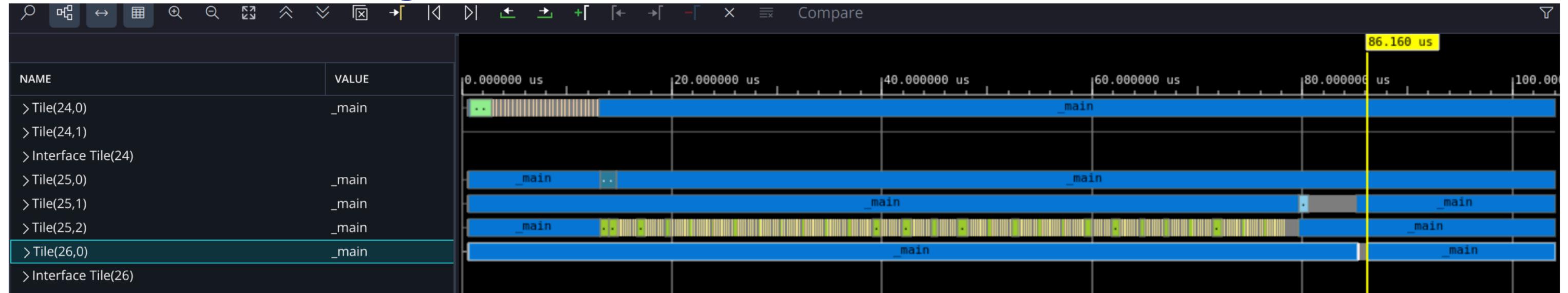
- DNN implementation:
  - Model on a “graph”
  - Dense layer on a “kernel”
- AI engine: C++ based coding on Vitis
  - AI engine libraries
  - AI engine specific functions
  - Scaler, Vector engines, pipelining, etc.

	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5
<b>Input nodes</b>	71	27	27	27	27
<b>Output nodes</b>	27	27	27	27	3
<b>Active Func.</b>	LeakyReLU	Softmax	—	LeakyReLU	Tanh

✓ **AI 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

# Latency optimization on Versal ACAP



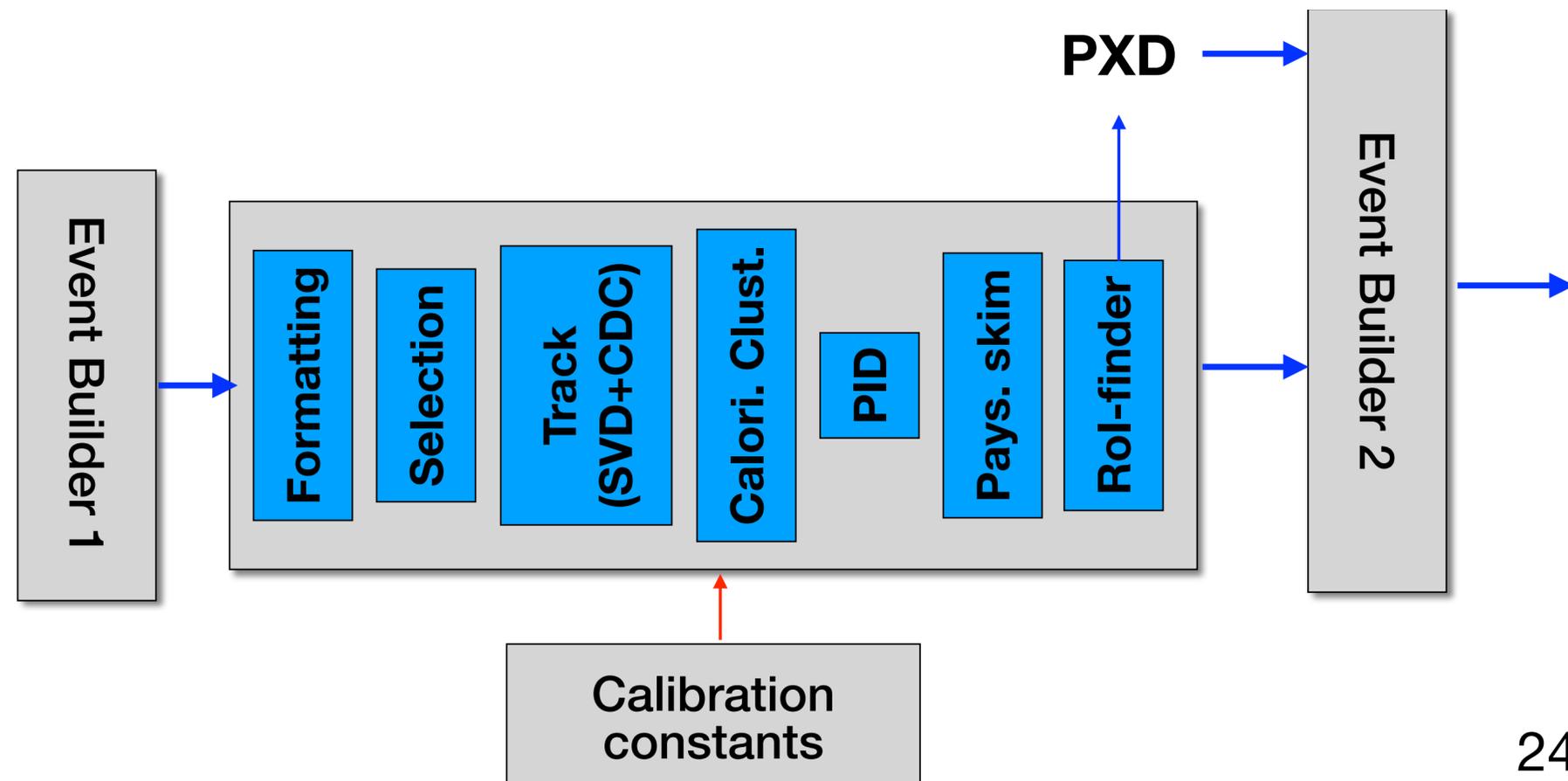
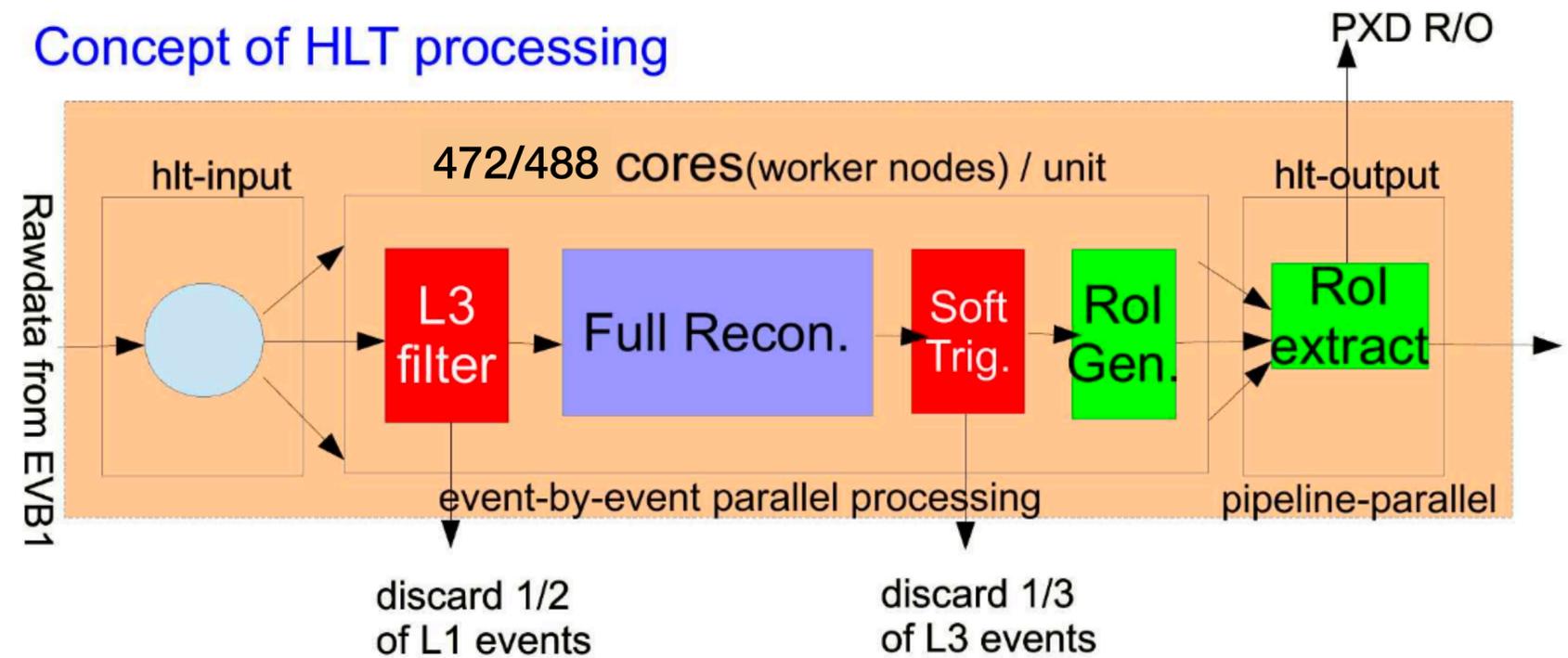
	Layer 1	Layer 2	Layer 3	Layer 4	Layer 5	Total
Input nodes	71	27	27	27	27	—
Output nodes	27	27	27	27	3	—
Active Func.	LeakyReLU	Softmax	—	LeakyReLU	Tanh	—
Ver.0 latency	~12us	~66us	~1.5us	~5.5us	~9.9us	~86us
Ver.1 latency	~2.1us	~1.3us	~1.5us	0.9us	~0.2us	~5us

# Machine Learning for software track trigger (SFOTWARE)

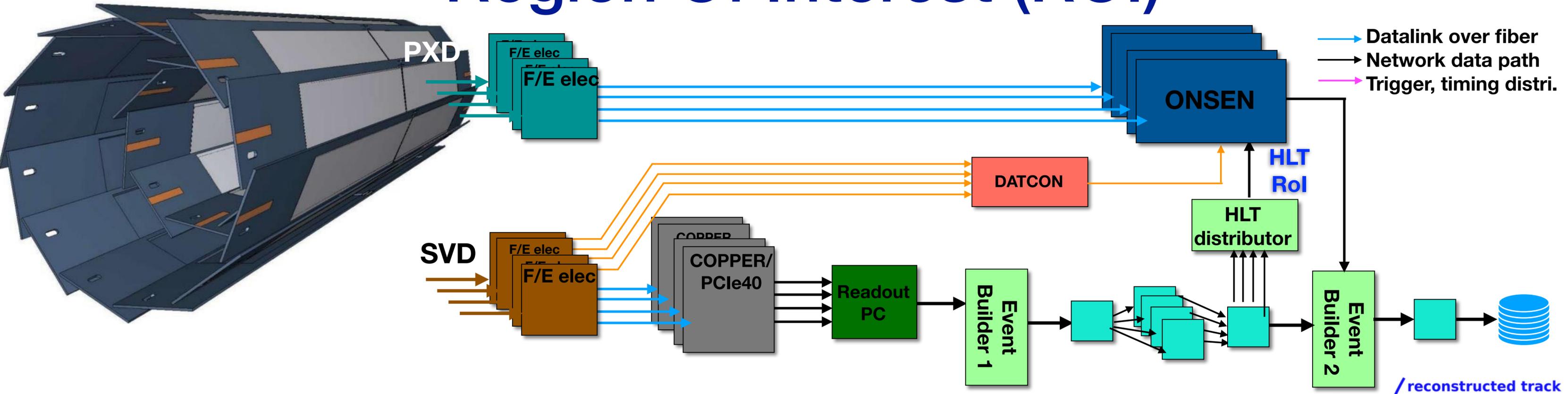
# Overview of high level trigger system at Belle II

- Full event reconstruction (same as offline processing)
- Crude calibration constant
- 13 HLT units, in total ~6200 CPU cores (design: 7000 cores)
- Data processing: ~ 2.1kHz/ HLT unit w/ hyper-threading
- Event size at HLT in the last run period: ~150 kB/event
- PXD event size = 1MB/event, 10 times larger than the rest of detectors
- Region of interest (RoI) method is effective to reduce the data size
- ROI
  - Tracking software running on HLT nodes

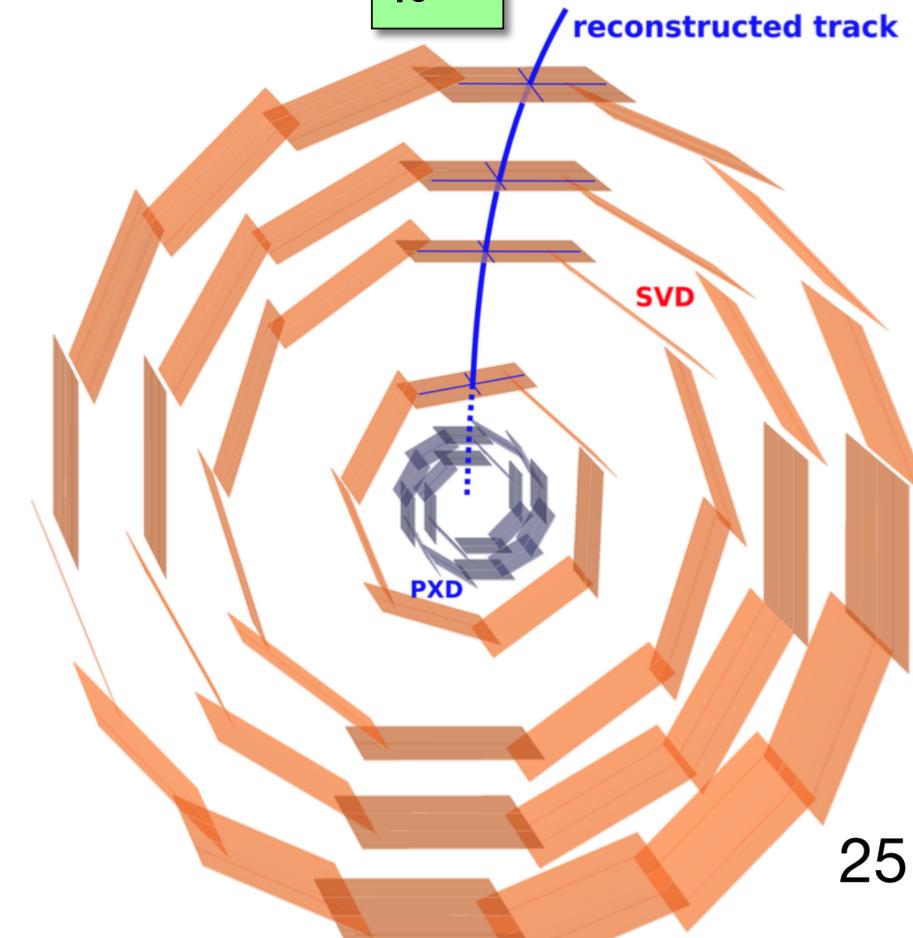
Concept of HLT processing



# Region Of Interest (ROI)



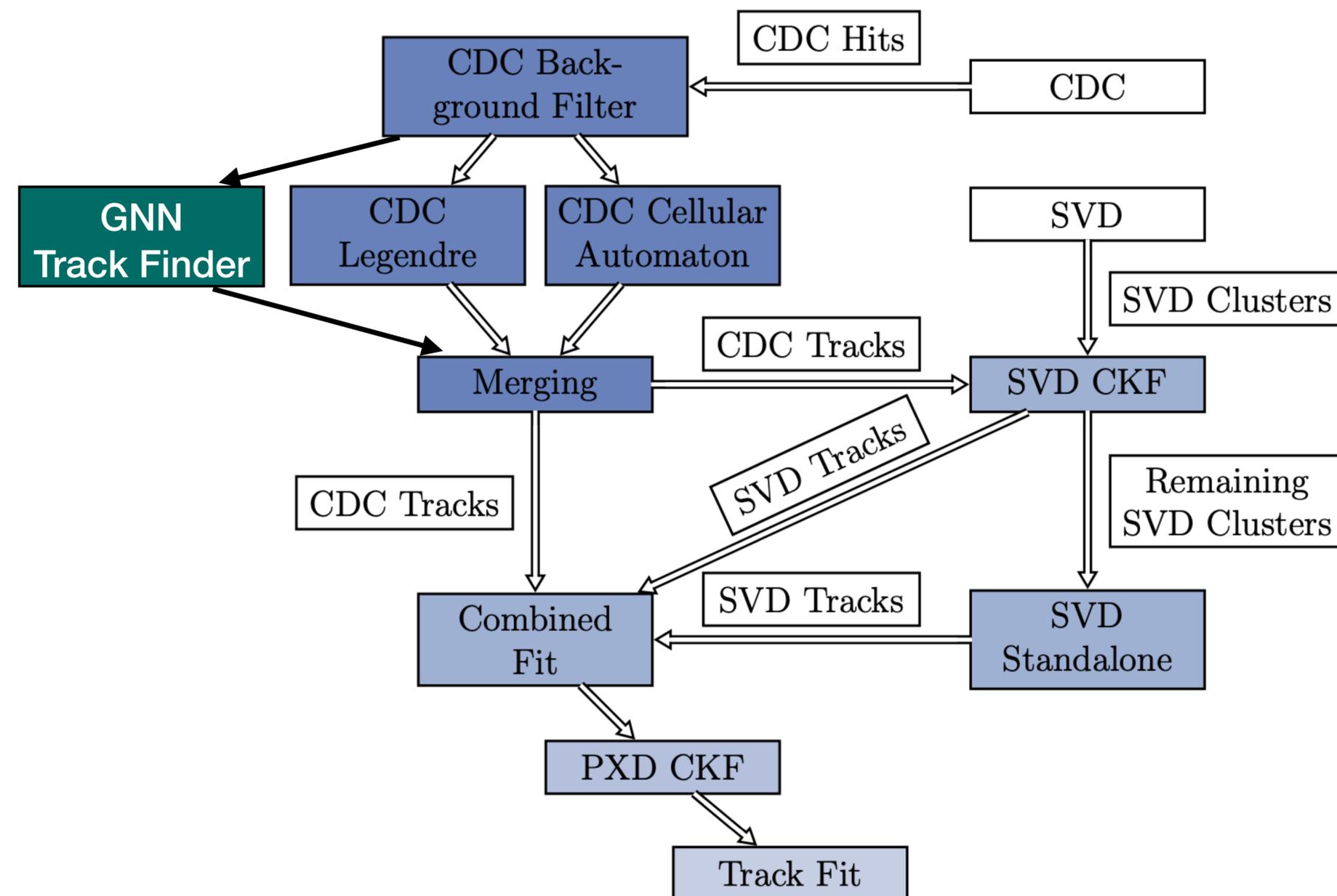
- PXD event size = 1MB/event, 10 times larger than the rest of detectors
- Region of interest method is effective to reduce the data size
- ROI
  - Tracking software running on HLT nodes
- PXD event data size reduced by 1/10 with ROI
  - In addition, trigger rate reduced by 1/3 with HLT ROI



# GNN based CDC track finder

- Motivations of introducing a GNN track finder (SOFTWARE)
- Low efficiency for displaced vertices
  - Efficiency decrease as displacement increase
  - Important signature for new physics search
- Higher background
- CDC wire inefficiencies
  - Bad wires or electric
  - Decreased efficiency

Comput.Phys.Commun. 259 (2021) 107610



- Modular structure for track finding, with flexible of reconstruction sequence

# GNN for offline track finding

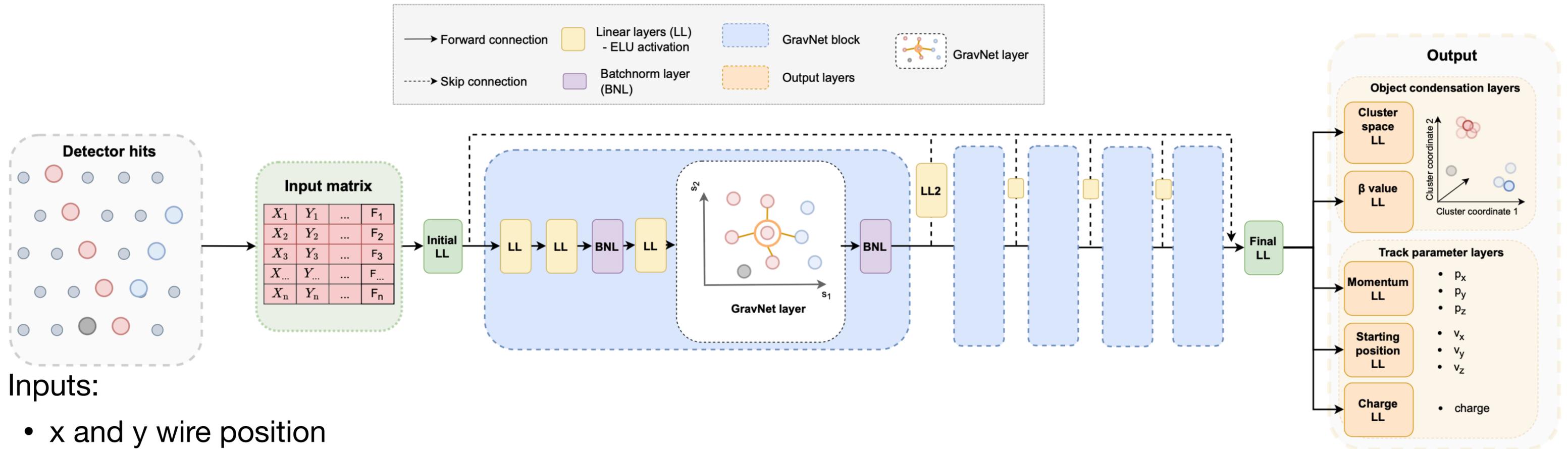
- Find track parameters: momentum, starting position and charge
- Find unknown number of tracks → Object Condensation ([arXiv:2002.03605](https://arxiv.org/abs/2002.03605))
- Computing resource and time constraint may be reducible

Noise filtering

Clustering

Fitting

[L. Reuter et. al. at \(KIT\) arXiv: 2411.13596](https://arxiv.org/abs/2411.13596)

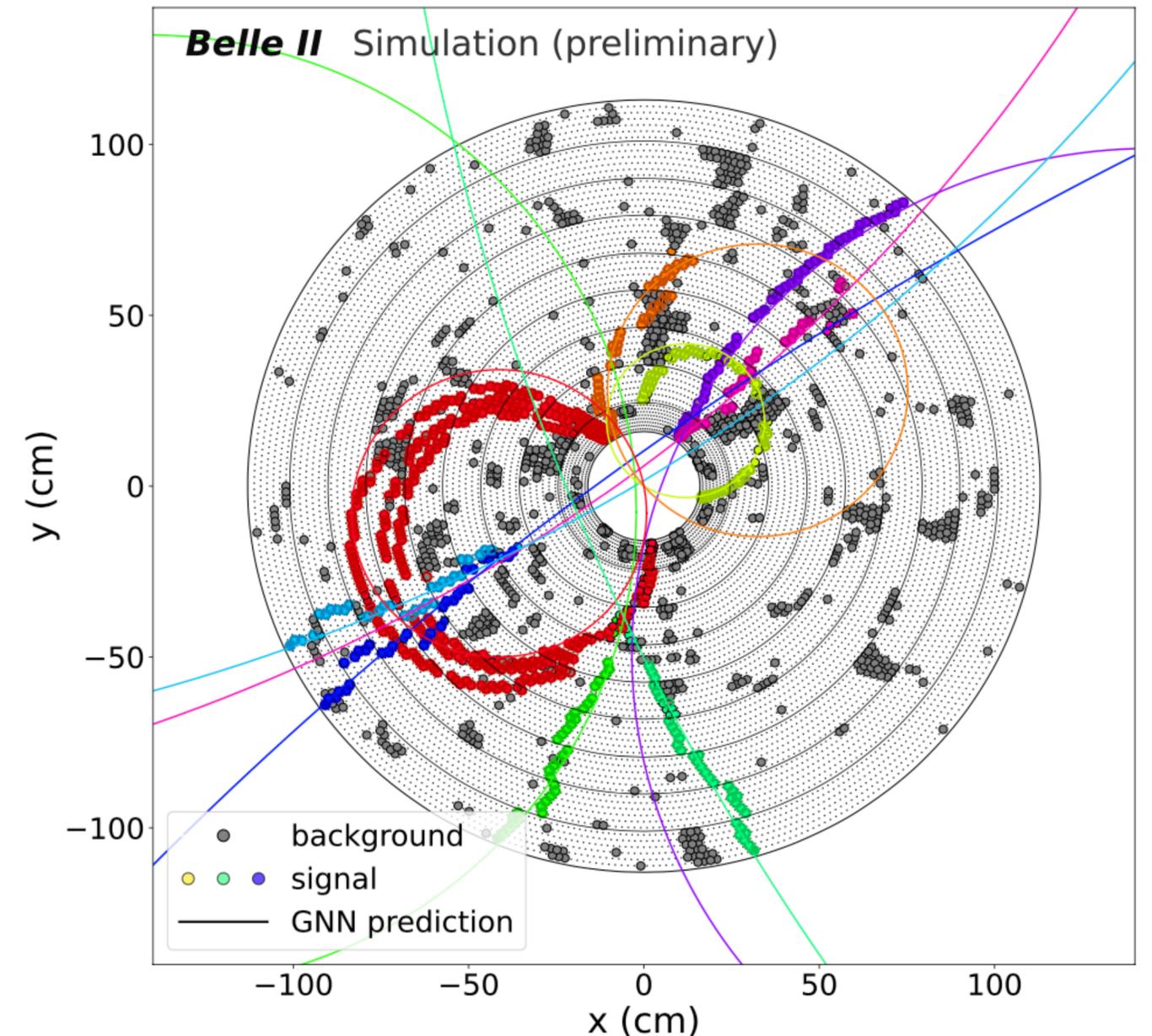
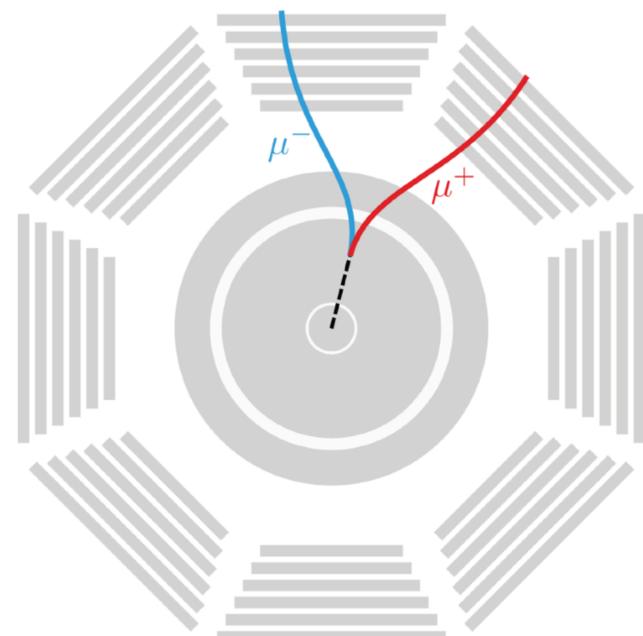


- Inputs:
  - x and y wire position
  - TDC and ADC of signal information
  - layer, superlayer, and layer info. with superlayer
- Adjustable Parameters
  - 797,812 trainable parameters (3MB weight files)

# Performance of GNN

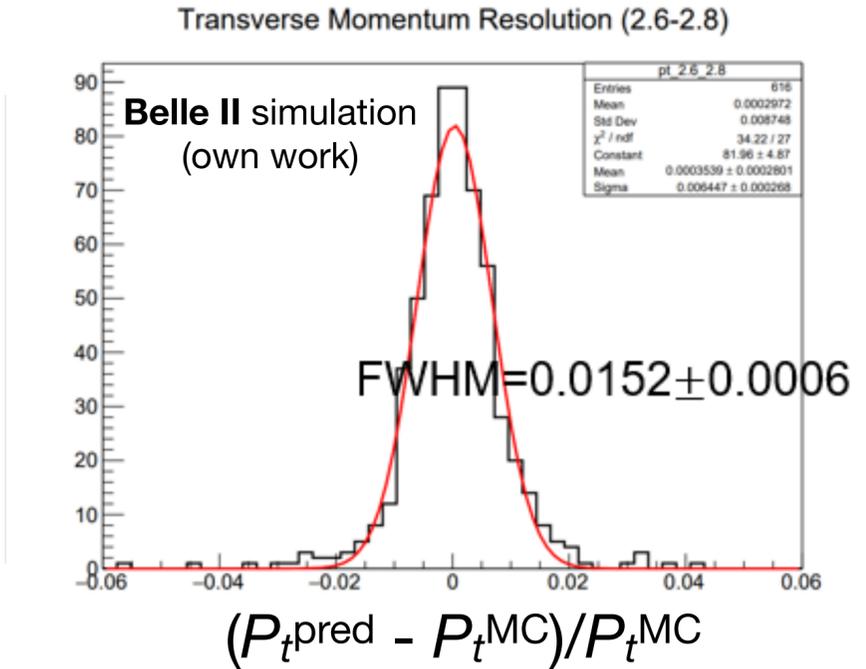
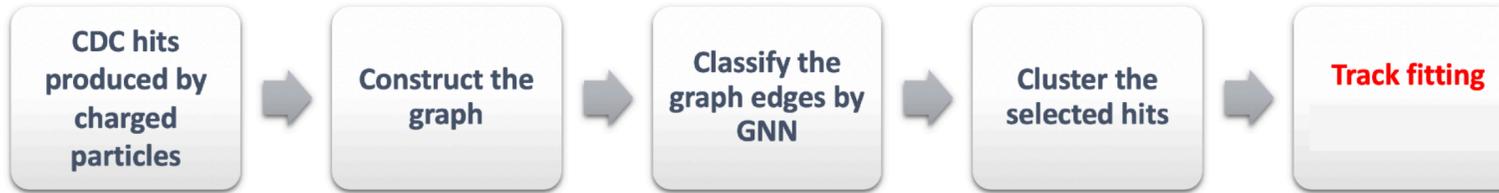
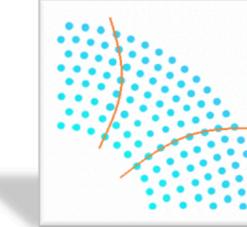
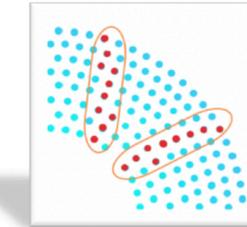
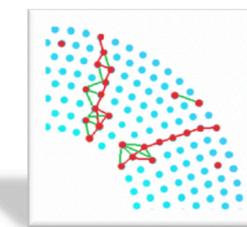
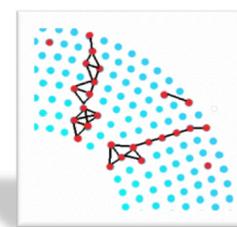
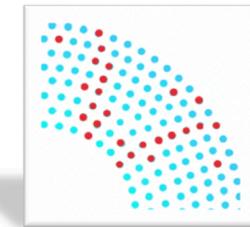
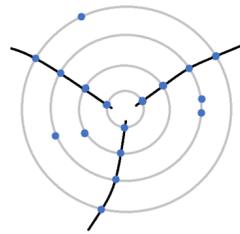
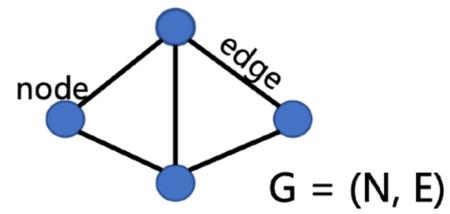
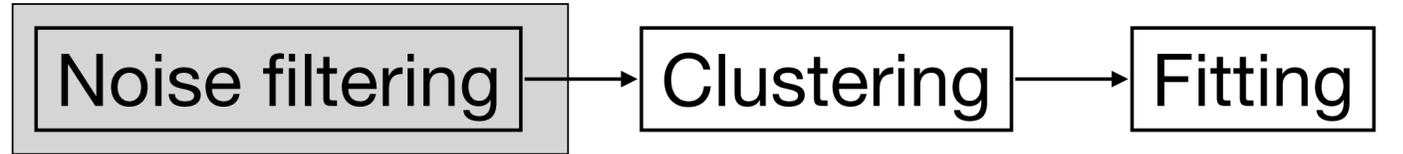
- Efficiency of displaced vertex tracks improved from 85.4% with a fake rate of 2.5%, compared to 52.2% and 4.1%
  - The other performance similar as original algorithm
- Momentum  $p_x$  ,  $p_y$  ,  $p_z$  starting position  $v_x$  ,  $v_y$  ,  $v_z$  , charge
  - Provide initial inputs for GENFIT
- GNN prediction is drawn according to the track parameters predicted by the GNN
- Plan to added as additional track finder for Belle II

[L. Reuter et. at \(KIT\) arXiv: 2411.13596](#)

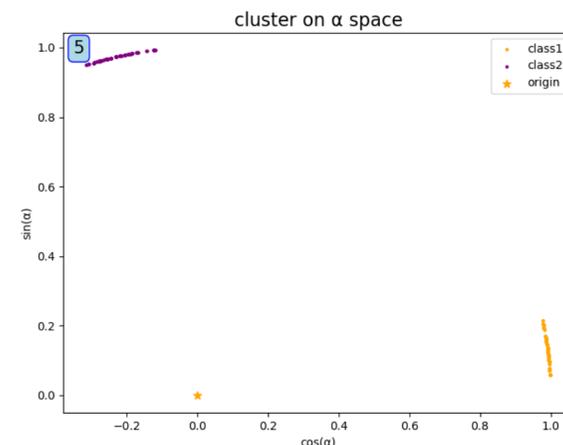
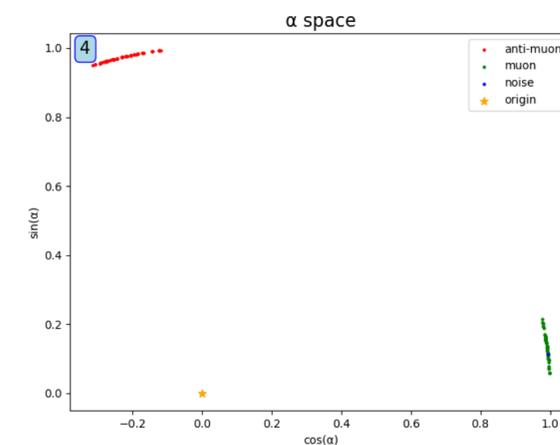
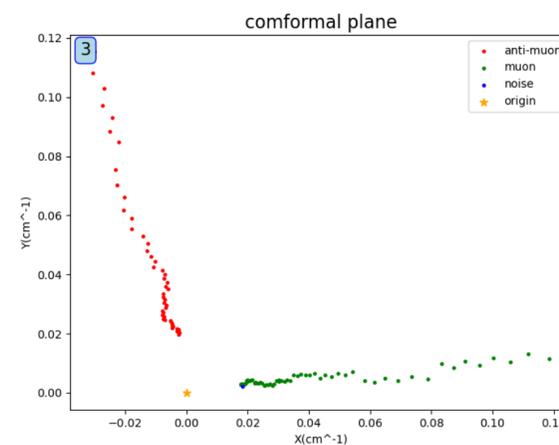
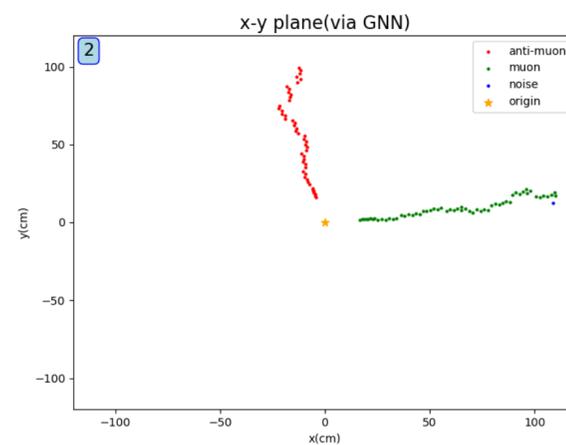
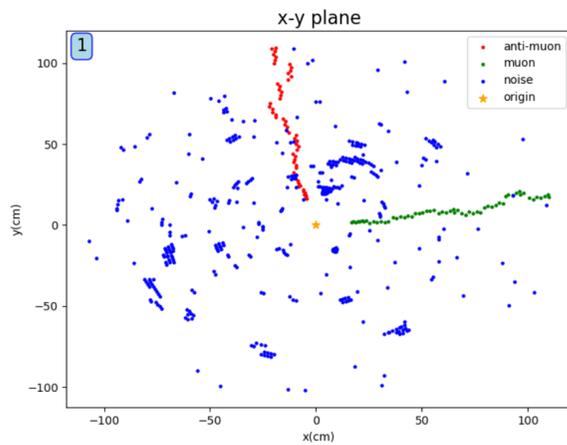


# GNN for CDC track background filtering

- Developed a GNN algorithm (based on [BESIII's algorithm](#)) for Belle II CDC hits clean up
  - Inputs: TDC, position coordinates  $r, \phi$



## Belle II simulation (own work)



$\mu^+ \mu^-$  (particle gun)

GNN noise filtering

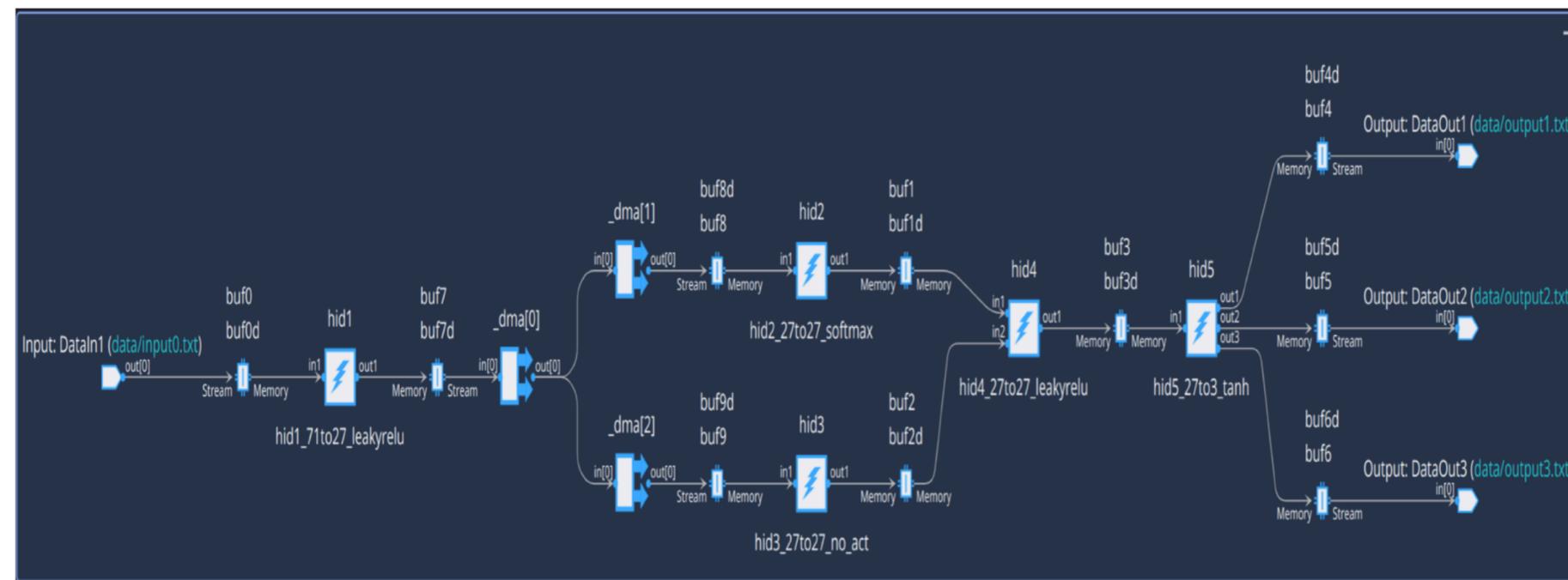
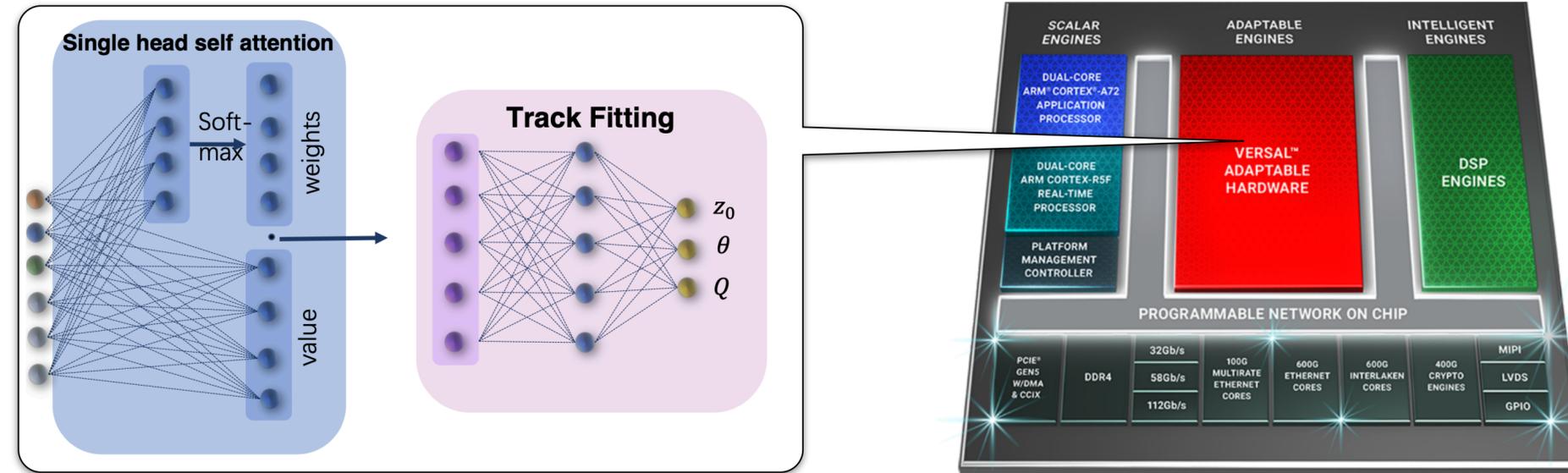
Transform space

Transform  $\alpha$  space

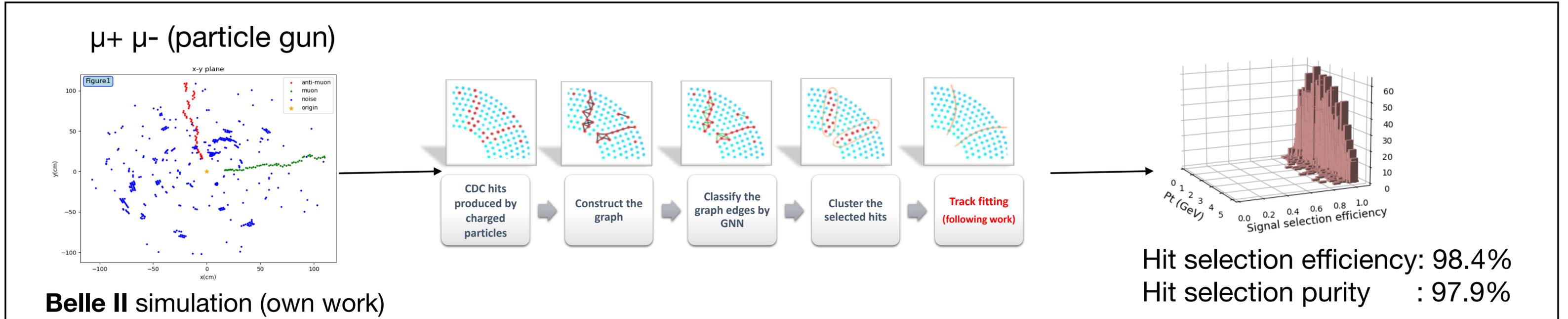
DBSCAN clustering

# NN acceleration on Versal ACAP

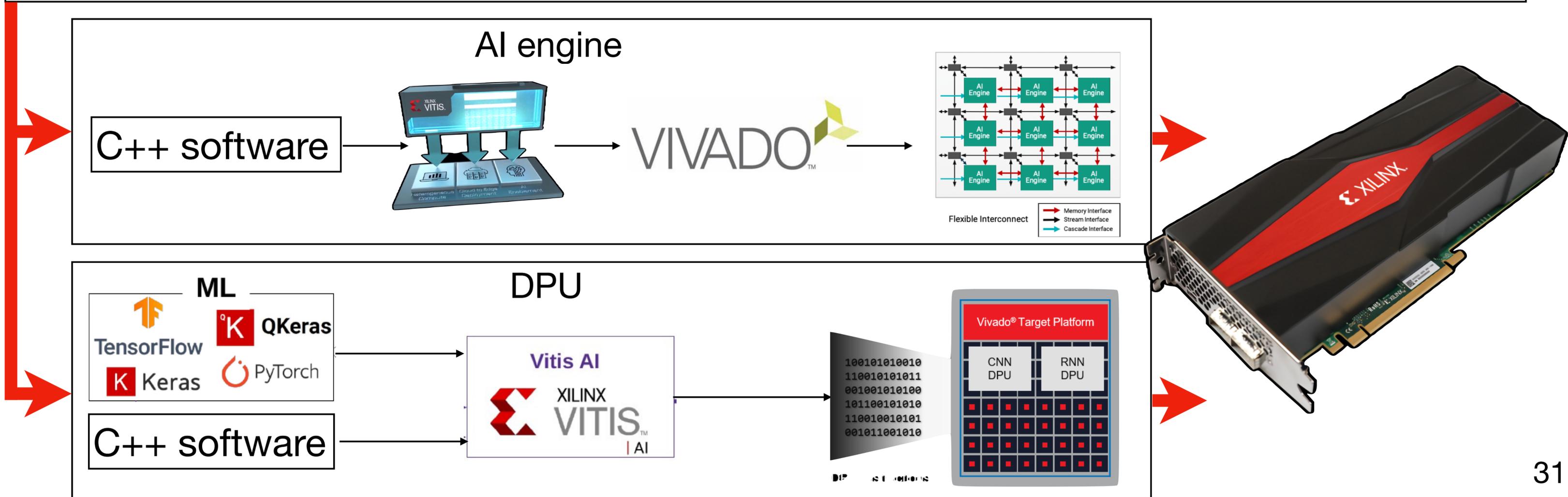
- Real-time graph building algorithm enables GNN implementation on FPGA for Belle II  
M. Neu et al. Comp. Soft. BigSci. 8, 8(2024)
- R&D of a new general FPGA device using the Versal ACAP
  - Heterogenous acceleration (VCK190, VCK5000 evaluation kit)
    - AI engine, DPU



# Acceleration on Versal ACAP platform



Belle II simulation (own work)



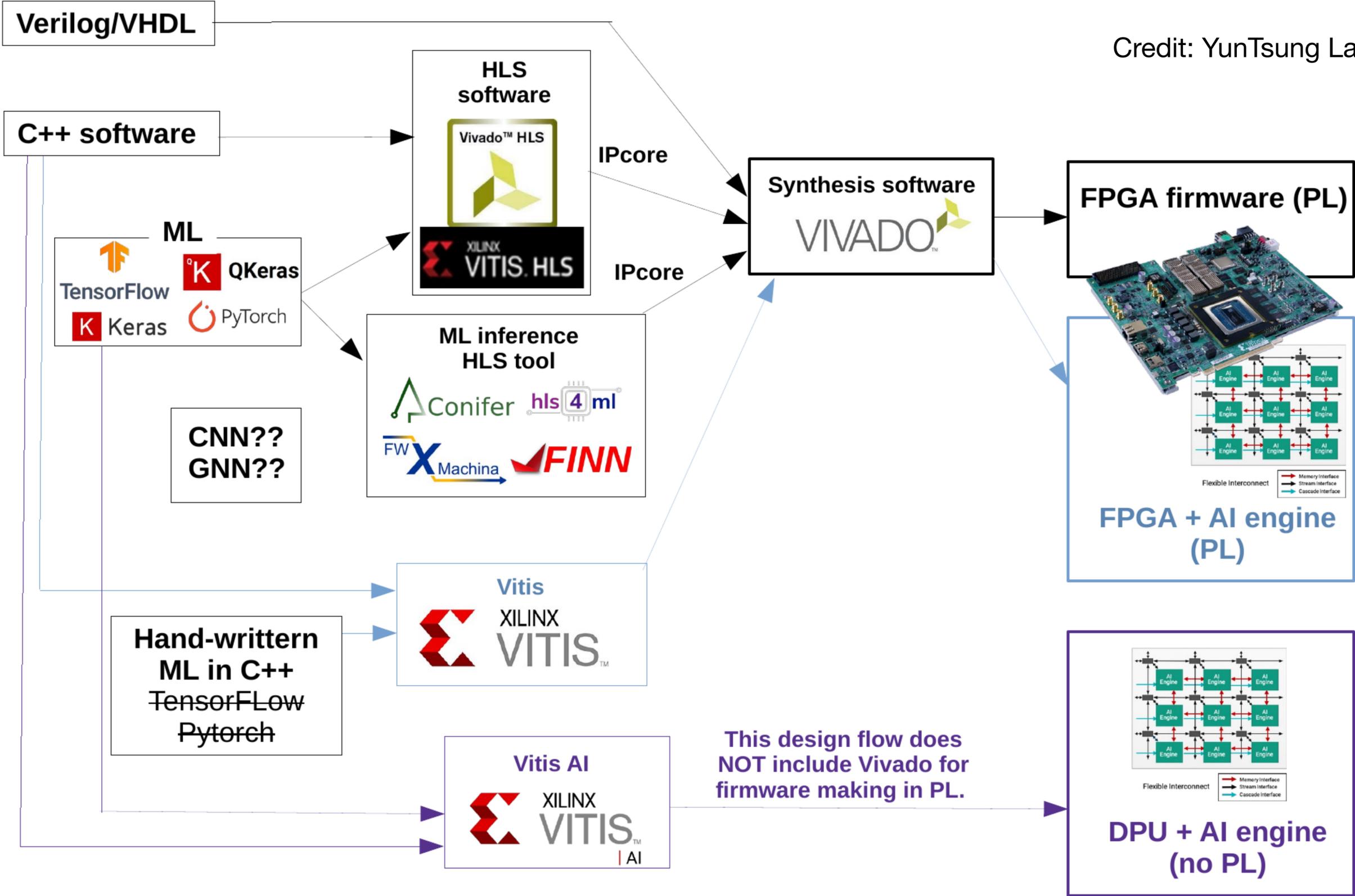
# Summary and prospects

- Belle II TDAQ system was designed to handle 30 kHz level 1 trigger
- NN and DNN with hardware based CDC L1 track trigger to improve background rejection
- GNN with software based offline CDC track finder to improve the efficiency of displaced vertex tracks
- Not covered in the talk: GNN with hardware based clustering trigger for Belle II is under commissioning
- Target the upgrade of ongoing and future collider projects
  - ML implementation on heterogenous computing system for acceleration

# Backup

# FPGA implementation path of ML algorithm

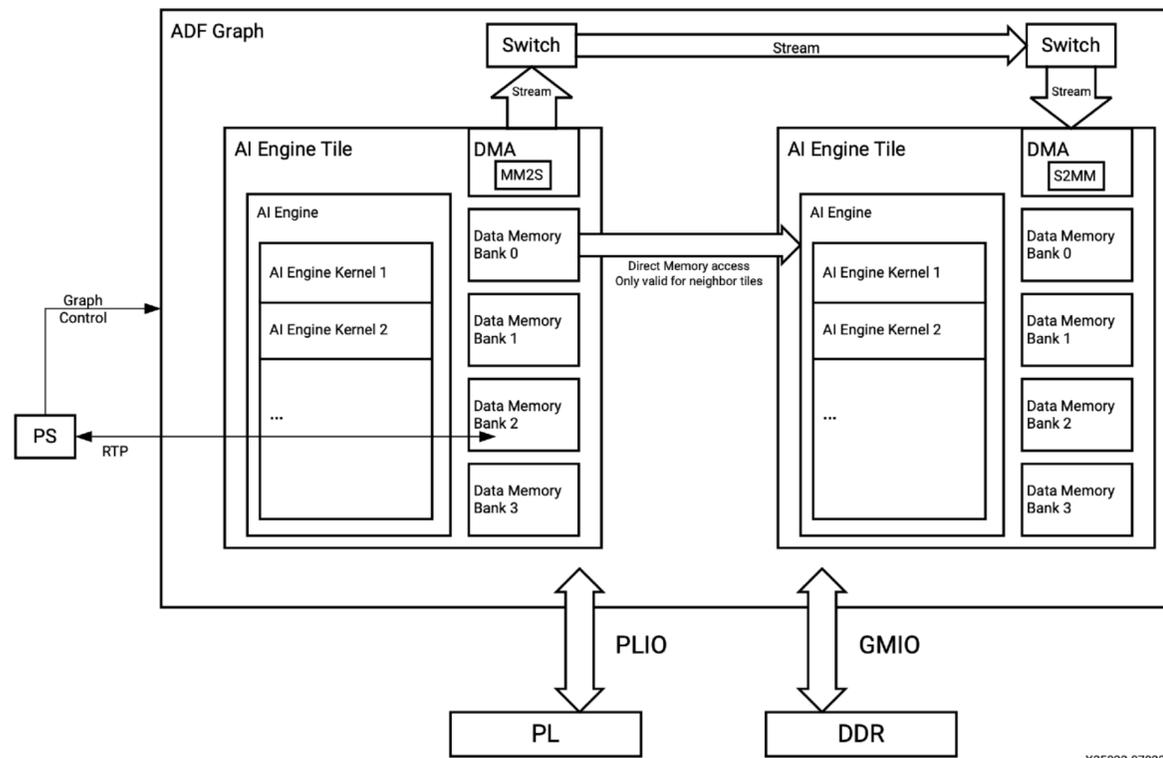
Credit: YunTsung Lai



# AI engine structure

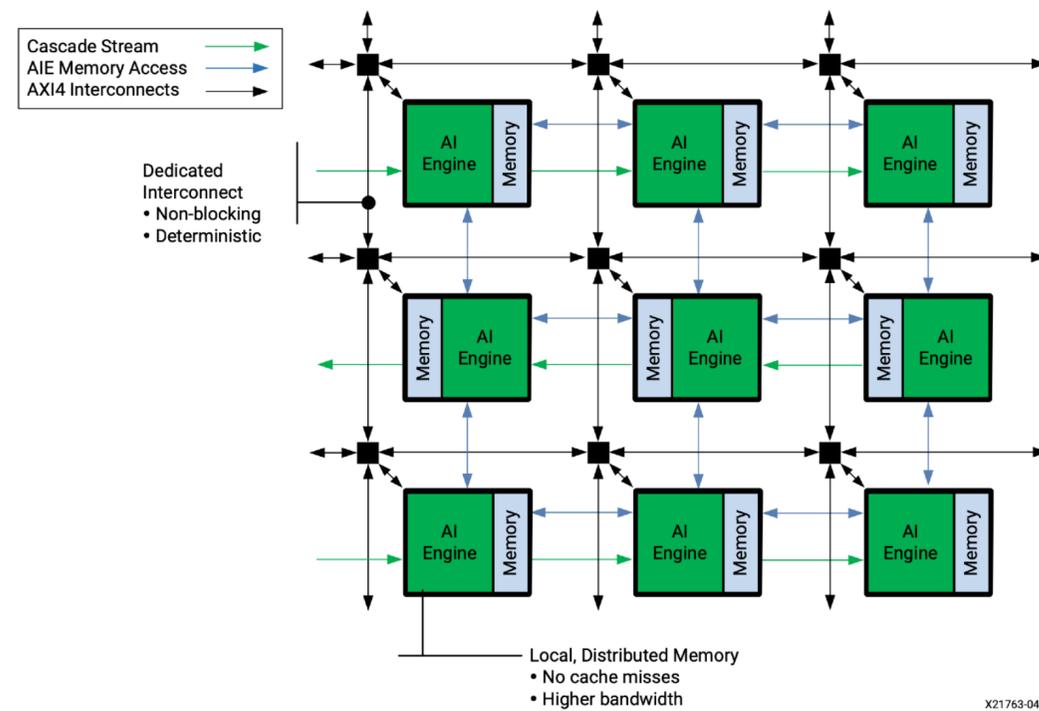
UG1079

Figure 1: Conceptual Overview of the ADF Graph



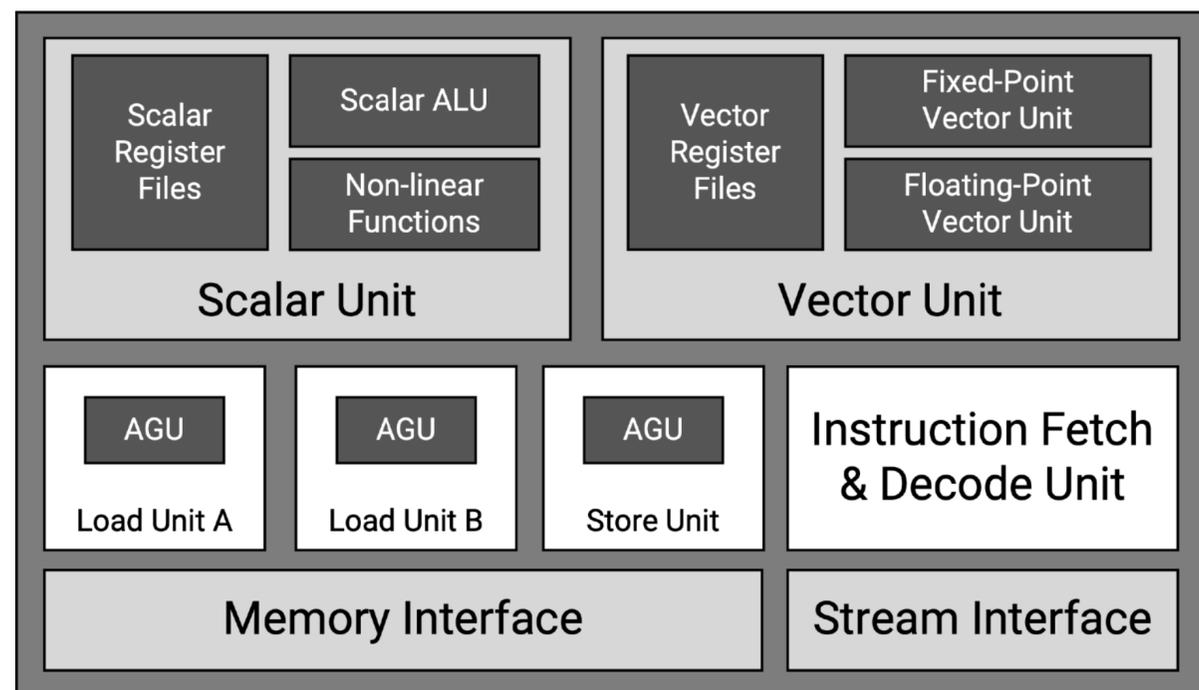
X25022-070221

Figure 2: AI Engine Array



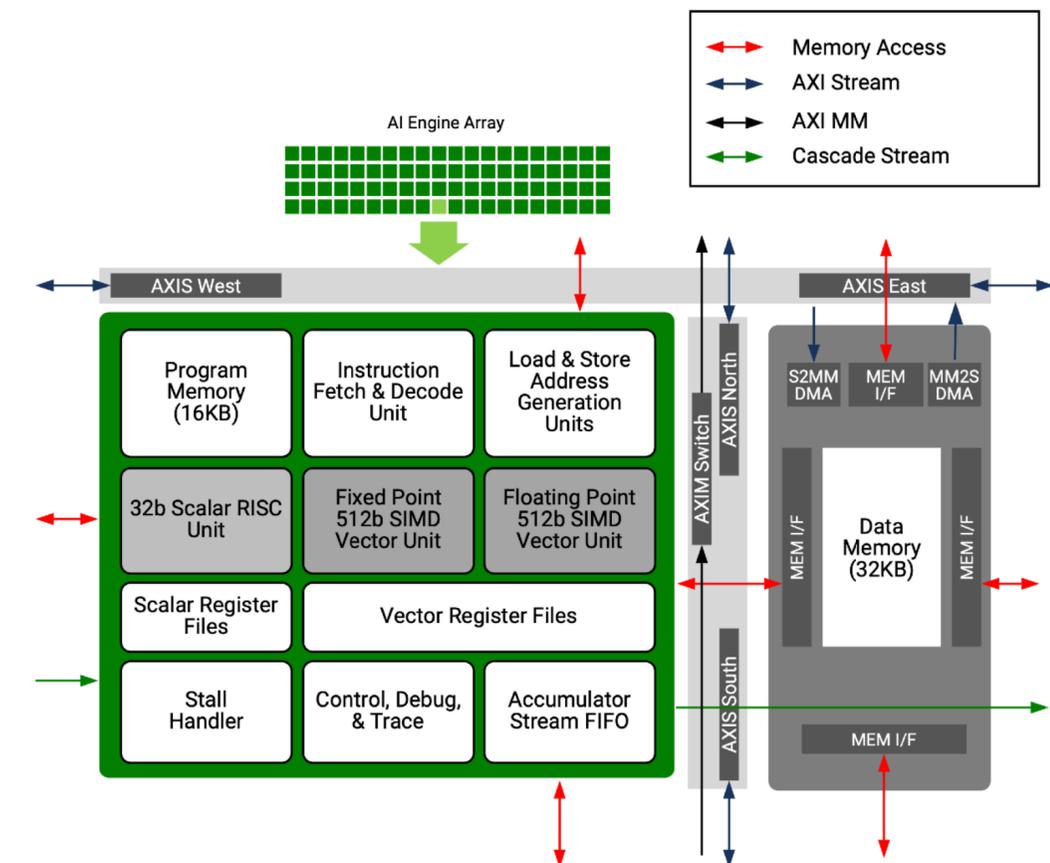
X21763-040519

Figure 4: AI Engine



X25020-011321

Figure 3: AI Engine Tile Details



# Kernel optimization for latency



before



after

Optimization	Before	After
Dense layer	Vector algorithm	Vector algorithm
Act. function	Scalar algorithm	Vector algorithm
Latency	~12us	~1.6us

# Motivations of trigger-DAQ upgrade

## Physics

- Tau trigger efficiency now is  $>95\%$  (to be pre-scaled if luminosity is high)
- Low multiplicity trigger efficiency (to be pre-scaled pre-scaled if luminosity is high)
- Low-momentum track trigger efficiency
- “Anomaly” trigger
  - Design a special trigger line for some specific physics channel
- Trigger efficiency of displaced vertex

## Current hardware limitation:

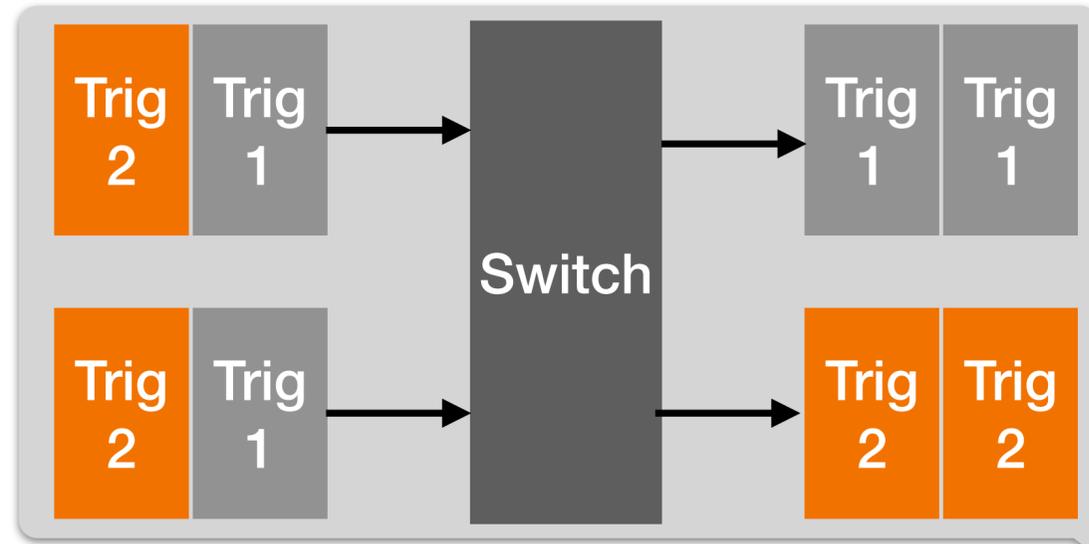
- DAQ system is designed to handle 30 kHz
  - L1 latency 4.4  $\mu\text{s}$  (SVD APV25 buffer)
    - CDC DNN trigger latency  $\sim 500$  ns, latency already limited more large model
- L1 trigger rate will reach to  $\sim 20$  kHz at  $0.9 \times 10^{-35} \text{ cm}^{-2} \text{ s}^{-1}$  (13 HLT units, w/o hyperthreading), planed full HLT: 15 units (7000 CPU cores)
- TTD system: VME bus limit, no more than 3 triggers within 80 clock (624ns)

## Vertex detector is planed to be upgraded during long shutdown 2 (after 2028)

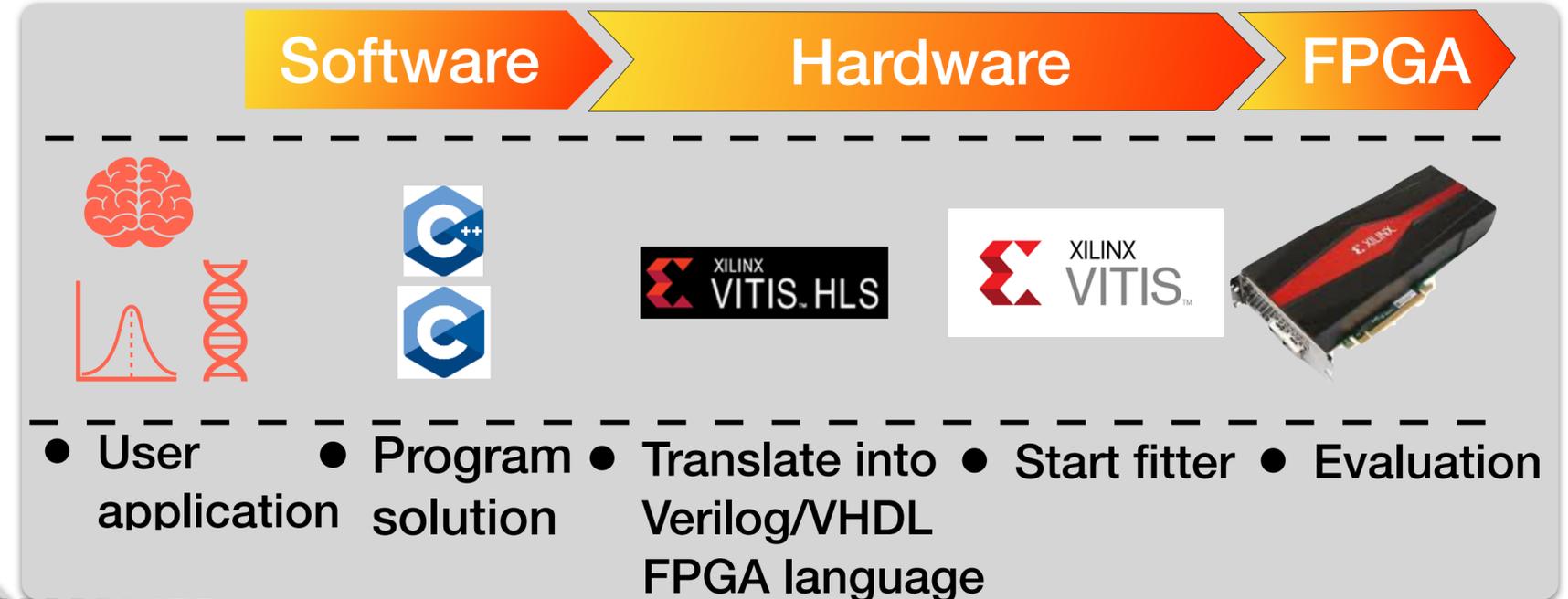
- Latency limit target: 5  $\mu\text{s}$   $\rightarrow$  10  $\mu\text{s}$  (5.2  $\mu\text{s}$  KLM, 9  $\mu\text{s}$  TOP, considering upgrade)
- New TTD hardware: VME bus  $\rightarrow$  Ethernet
- New trigger board (UT5): Versal ACAP

# Idea of upgrade trigger and DAQ system

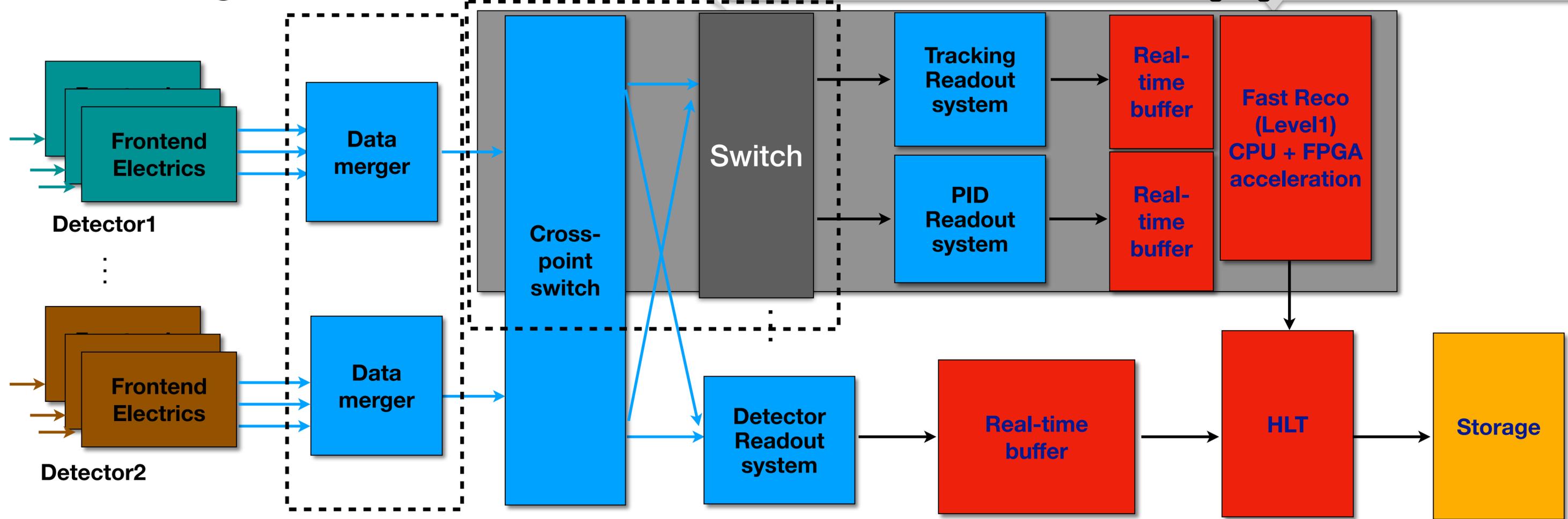
Proposed on Belle II trigger DAQ workshop 2022



Sorting

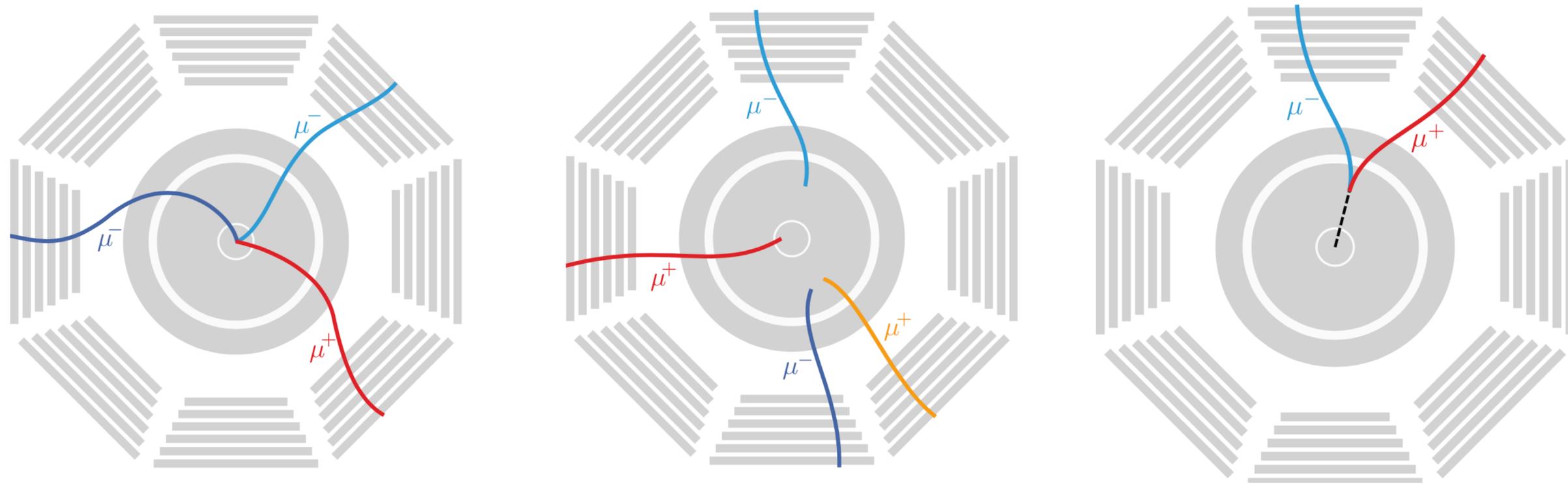


Optional



# Training of GNN

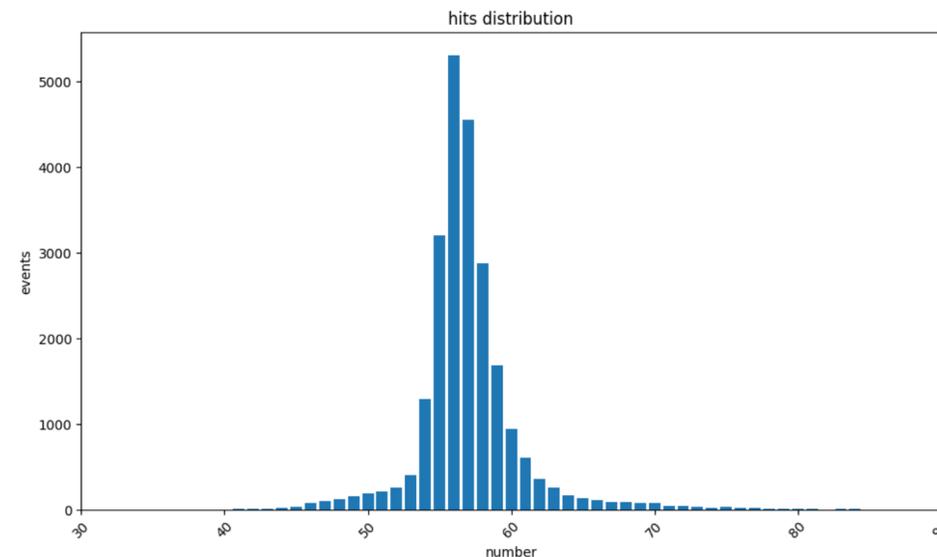
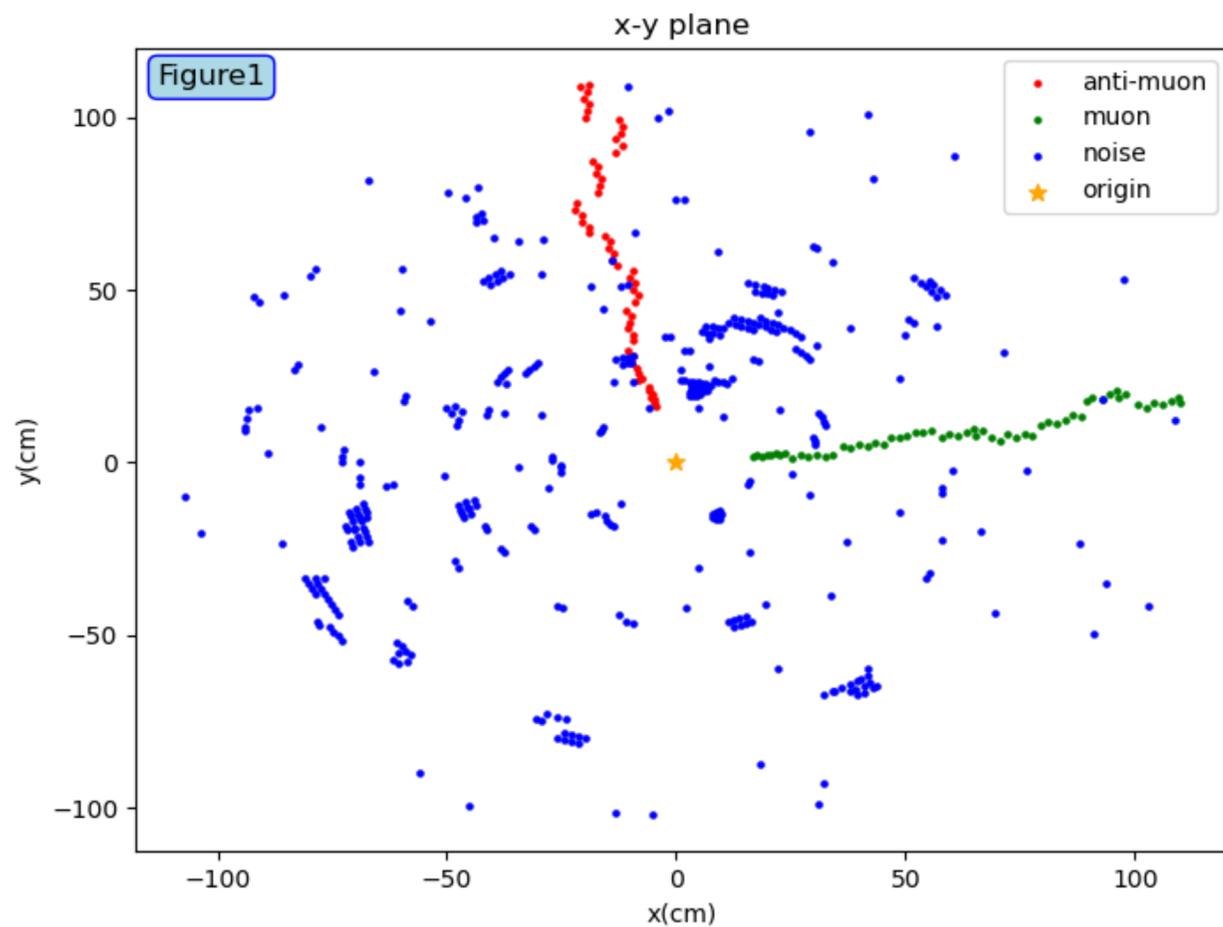
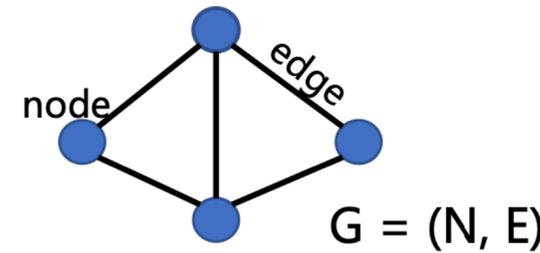
- Simulate 1 million events with over 4 million tracks
  - Train: Validation = 4 :1
- Training samples contain different topologies that cover all interested event features, to not bias the model, **no conservation laws involved here!**
  - crucial step to be agnostic about the physics processes
- Sample features
  - Low momentum tracks forming circles in the CDC ( $P_t < 0.4$  GeV)  $\leftrightarrow$  High momentum tracks
  - Short tracks  $\leftrightarrow$  tracks penetrate all CDC layers
  - Small opening angle  $\leftrightarrow$  well isolated two tracks
  - ...



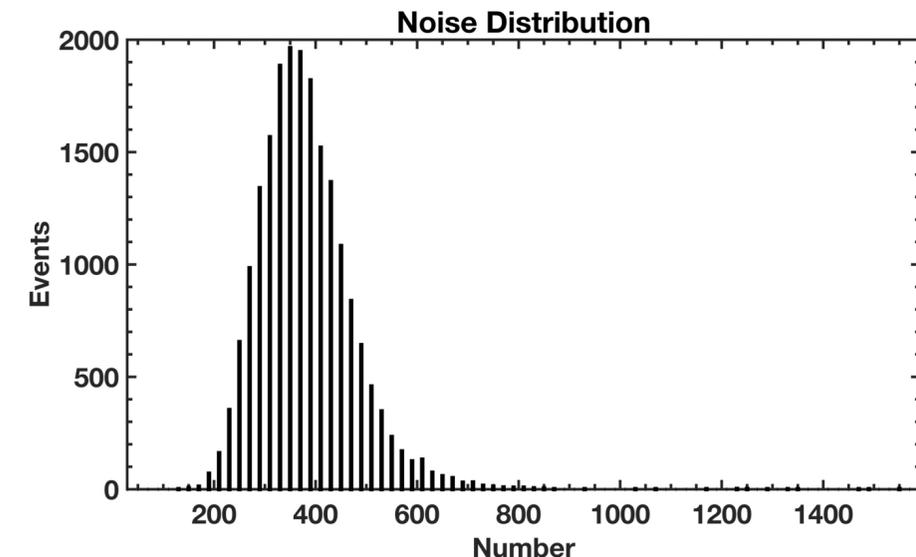
# Development of GNN tracking algorithm

- Belle II MC simulation data-set (Own simulation)

- $\mu^+ \mu^-$  (particle gan)
- $0.3 \text{ GeV}/c < P < 5.0 \text{ GeV}/c$
- Theta:  $30^\circ - 120^\circ$ , within on barrel CDC
- Phi:  $0 - 2\pi$
- Train: Validation: Test = 3: 1: 1
- noise : `/group/belle2/dataproduct/BGOverlay/early_phase3/release-06-00-05/overlay/BGx1/set0/`



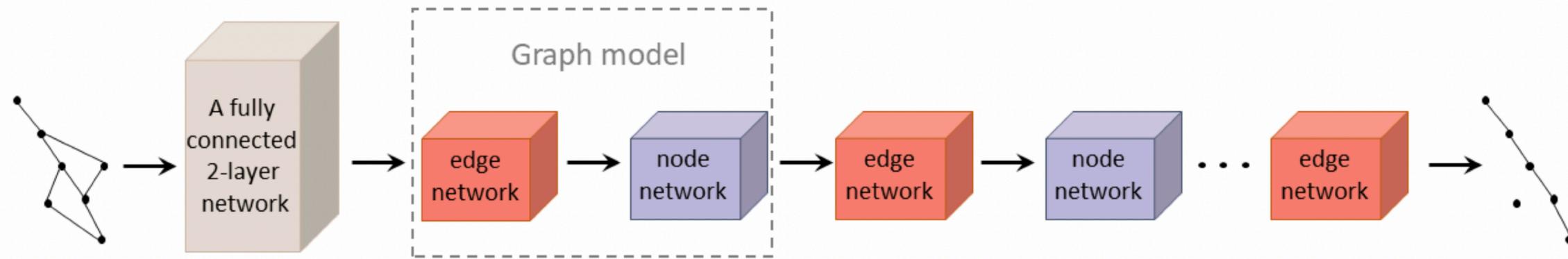
Signal hit No. distribution



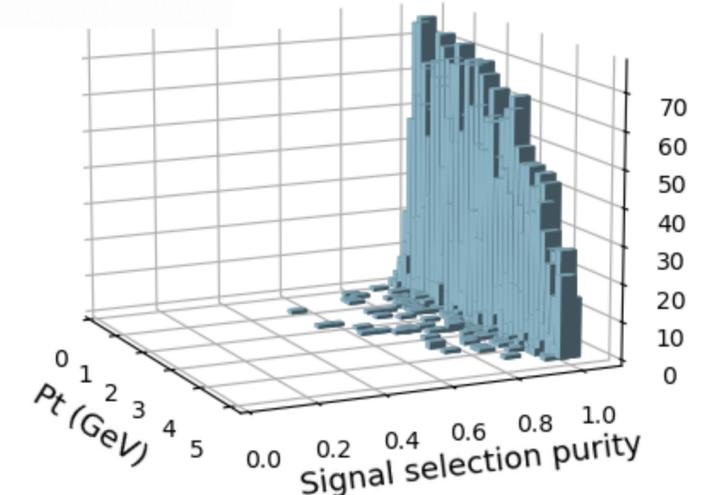
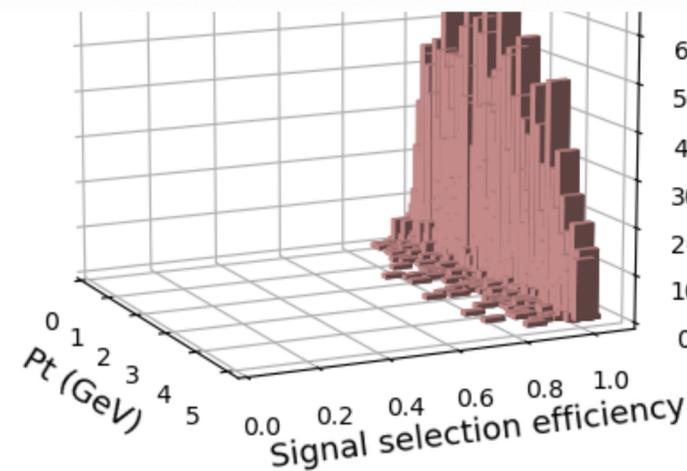
Noise hit No. distribution

# Development of GNN tracking algorithm

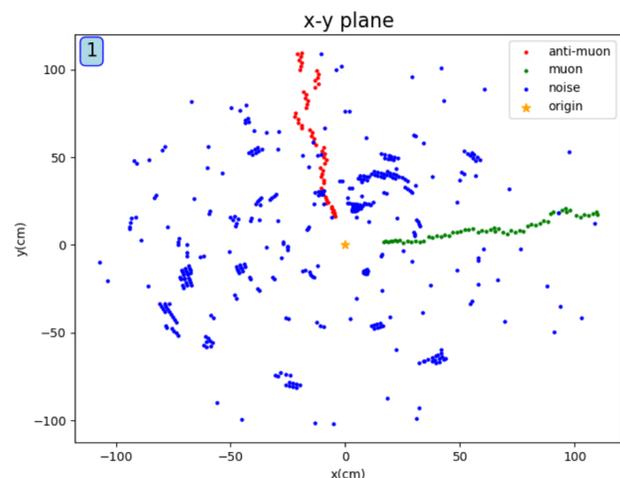
- Graph Neural Network edge classifier
- Input network
  - Node features embedded in latent space
- Graph model
  - Edge network computes weights for edges using the features of the start and end nodes
  - Node network computes new node features using the edge weight aggregated features of the connected nodes and the nodes' current features
  - MLPs
  - 8 graph iterations
- Strengthen important connections and weaken useless or spurious ones



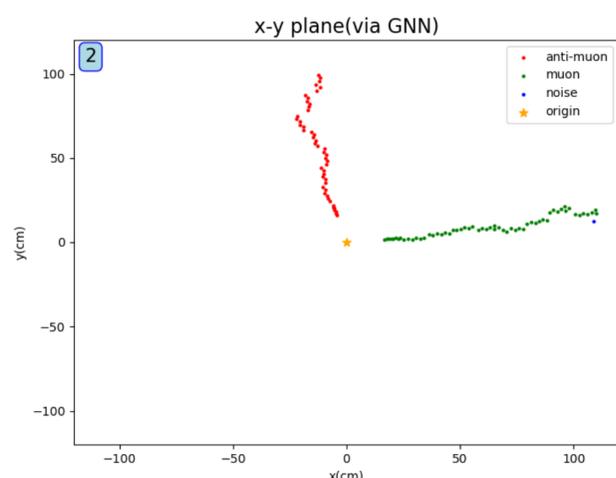
Hit selection efficiency: 98.4%  
Hit selection purity : 97.9%



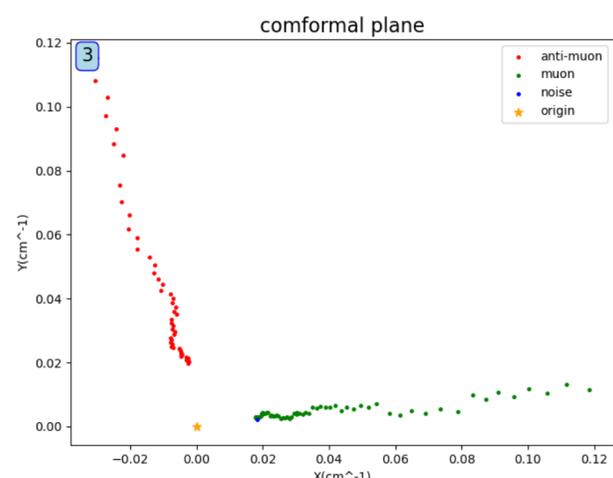
# Performance step-by-step



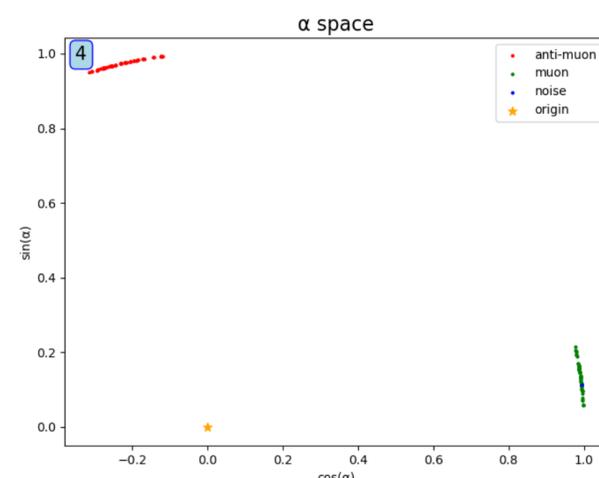
$\mu^+ \mu^-$  (particle gan)



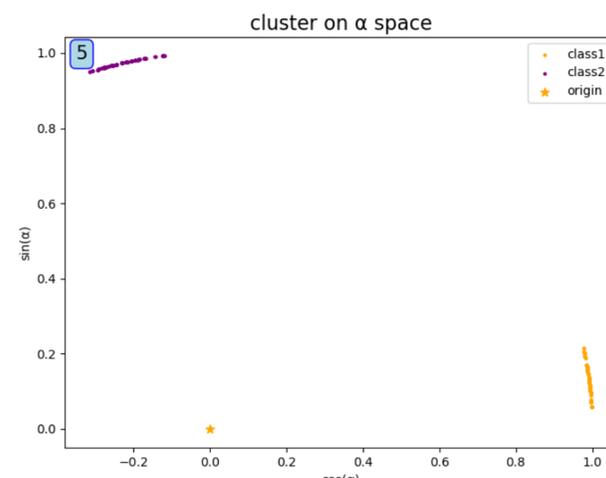
GNN noise filtering



Transform space



Transform  $\alpha$  space



DBSCAN clustering

## 1. Original MC data sample

- $\mu^+ \mu^-$  (use particle gan)
- $P$  (0.3GeV - 5.0GeV)

## 2. Remove noise via GNN

## 3. Transform to Conformal plane

- $X=2x/(x^2+y^2)$   $Y=2y/(X^2+y^2)$
- Circle passing the origin transform into a straight line

## 4. Transform to ' $\alpha$ ' parameter plane

- Hits connected in the X-Y plane in a straight line
- $\alpha$  as the angle between the straight line and X axis
- The parameter space as  $\cos\alpha$  and  $\sin\alpha$

## 5. DBSCAN clustering in ' $\alpha$ ' parameter plane

- Density-Based Spatial Clustering of Application with Noise
- Hits in a cluster are considered to be in the same track

Cluster efficiency: 97.7%  
Cluster purity : 96.9%