# Machine Learning for real-time data processing at Belle II

Qi-Dong ZHOU
Institute of Frontier and Interdisciplinary Science,
Shandong Univ. (Qingdao)
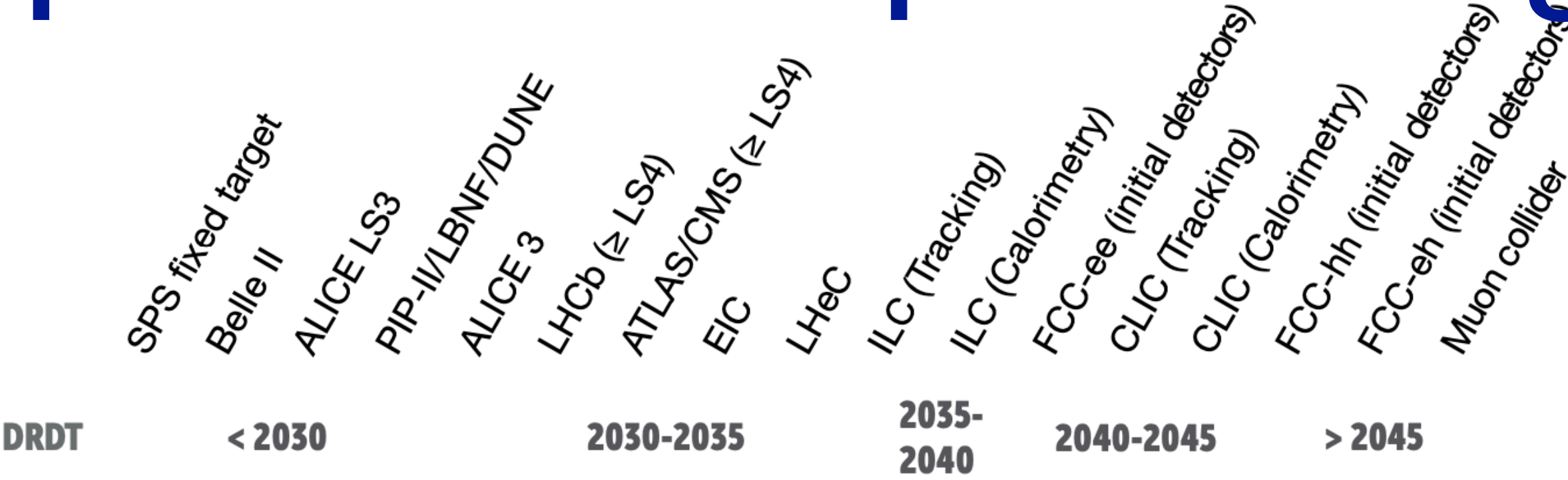
20-23 Aug. 2025, Qingdao
Quantum computing and machine learning workshop 2025

**Belle II**

# Roadmap of techniques for 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) | - |



**Data density**
- High data rate ASICs and systems — 7.1
- New link technologies (fibre, wireless, wireline) — 7.1
- Power and readout efficiency — 7.1

**Intelligence on the detector**
- Front-end programmability, modularity and configurability — 7.2
- Intelligent power management — 7.2
- Advanced data reduction techniques (ML/AI) — 7.2

**4D-techniques**
- High-performance sampling (TDCs, ADCs) — 7.3
- High precision timing distribution — 7.3
- Novel on-chip architectures — 7.3

**Extreme environments and longevity**
- Radiation hardness — 7.4
- Cryogenic temperatures — 7.4
- Reliability, fault tolerance, detector control — 7.4
- Cooling — 7.4

**Emerging technologies**
- Novel microelectronic technologies, devices, materials — 7.5
- Silicon photonics — 7.5
- 3D-integration and high-density interconnects — 7.5
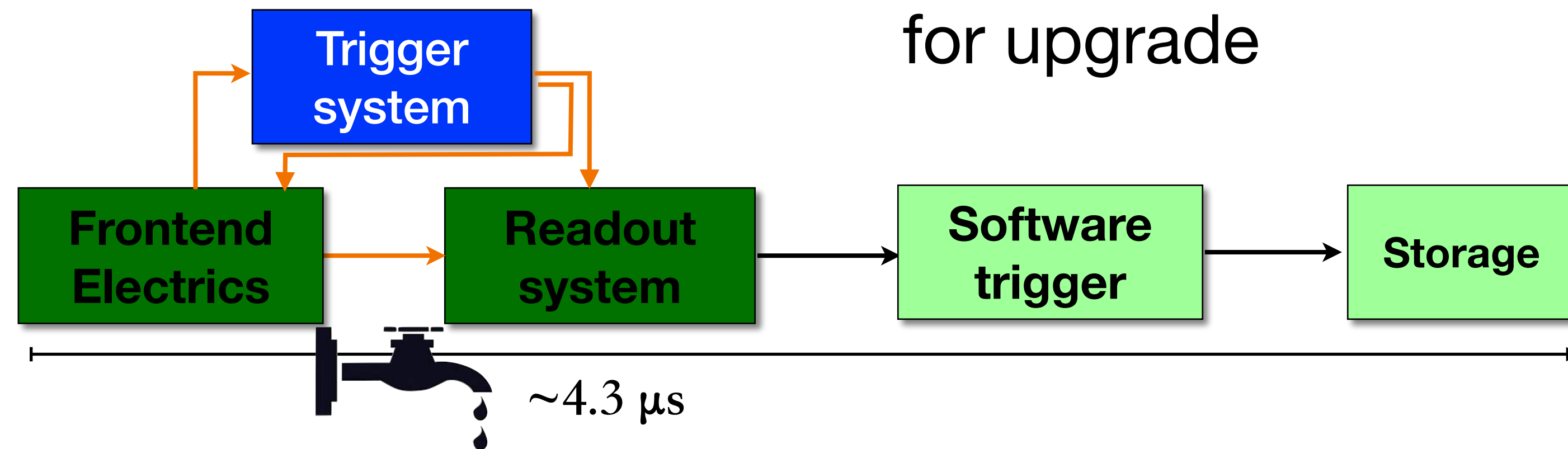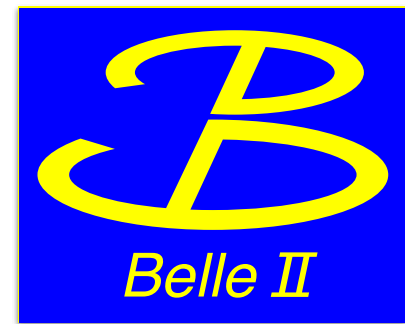- Keeping pace with, adapting and interfacing to COTS — 7.5

Column headers: SPS fixed target, Belle II, ALICE LS3, PIP-II/LBNF/DUNE, ALICE 3, LHCb (≥ LS4), ATLAS/CMS (≥ LS4), EIC, LHeC, ILC (Tracking), ILC (Calorimetry), FCC-ee (initial detectors), CLIC (Tracking), CLIC (Calorimetry), FCC-hh (initial detectors), FCC-eh (initial detectors), Muon collider

Timeline: < 2030 | 2030-2035 | 2035-2040 | 2040-2045 | > 2045

Legend:
- 🔴 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

ECFA detector R&D

2

# Data processing system (Belle II vs. LHCb)

- Belle II: L1 trigger + HLT
  - Trigger efficiency:
    - Had. B physics~100% τ physics
    70~95%

- LHCb: "triggerless" readout & DAQ
  - CPU+GPU based software trigger
  - ~350 GPU RTX A5000
  - Part of online data processing with FPGA for upgrade



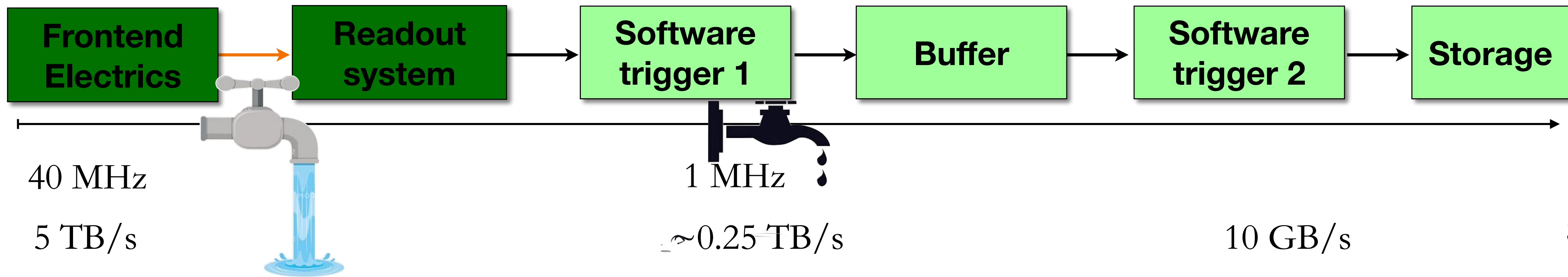| | | | |
|---|---|---|---|
| Latency | | ~4.3 µs | |
| Trigger rate | 127 MHz | 30 kHz | 30 kHz |
| Throughput | 3(33) GB/s | 2 (32) GB/s | 3 GB/s |

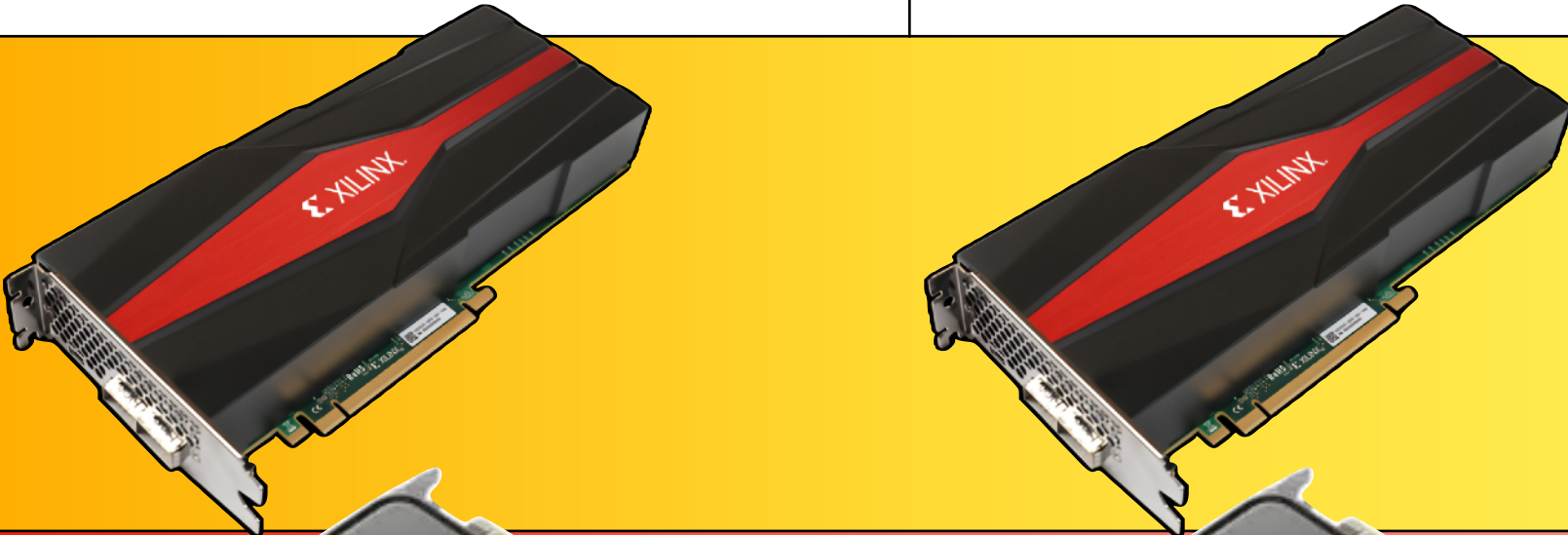| | | | |
|---|---|---|---|
| Trigger rate | 40 MHz | 1 MHz | |
| Throughput | 5 TB/s | ~0.25 TB/s | 10 GB/s |

# Heterogenous computing system

- System integrated with different devices
  - System level, chip level
- FPGA widely used in frontend and trigger electrics
- FPGA, AI engine, DPU, System On Chip, Network On Chip for computing acceleration

| Devices | Specifics | | | | | |
|---|---|---|---|---|---|---|
| | | Trigger system → Frontend Electrics → Readout system | | Software trigger | Buffer | Offline Software |
| **FPGA** | • Level 1 trigger: low latency w/ **Ous** <br> • Simple ML algorithm available | | | | | |
| **FPGA/ AI engine/ DPU** | • Software trigger w/ **Oms** <br> • CPU/FPGA/GPU system on chip <br> • CNN easily run on DPU | | | | | |
| **GPU** | • Software trigger & offline software <br> • User friendly and widely supported | | | | | |



4

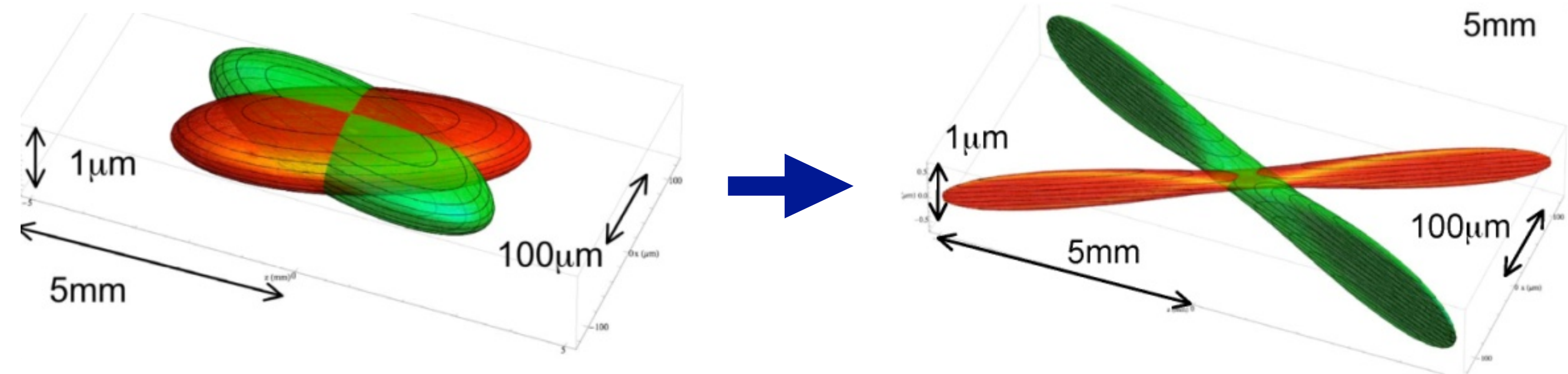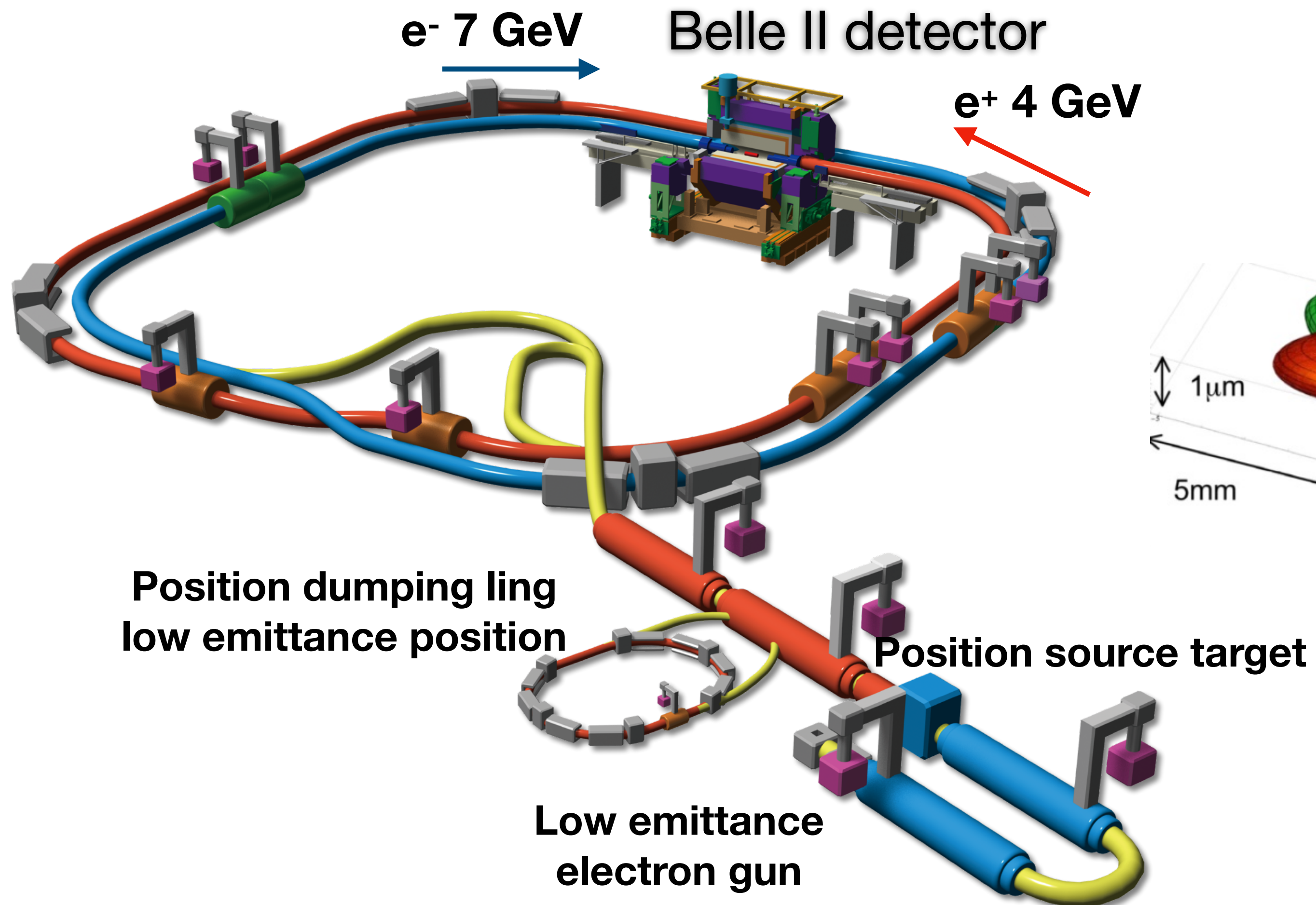# Luminosity frontier: SuperKEKB

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

Beam current: KEKB x ~1.5

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

Beam squeeze: KEKB / ~20

**e⁻ 7 GeV**    Belle II detector

**e⁺ 4 GeV**

**Nano beam scheme**



**Position dumping ling
low emittance position**

**Position source target**

**Low emittance
electron gun**

Target: L = 60 x 10³⁴ cm⁻² s⁻¹
Achieved : 5.1 x 10³⁴ cm⁻² s⁻¹ (Record)
- Data:
  - **575** fb⁻¹ (Belle II)  <->  980 fb⁻¹ (Belle)

# Belle II detector and dataset

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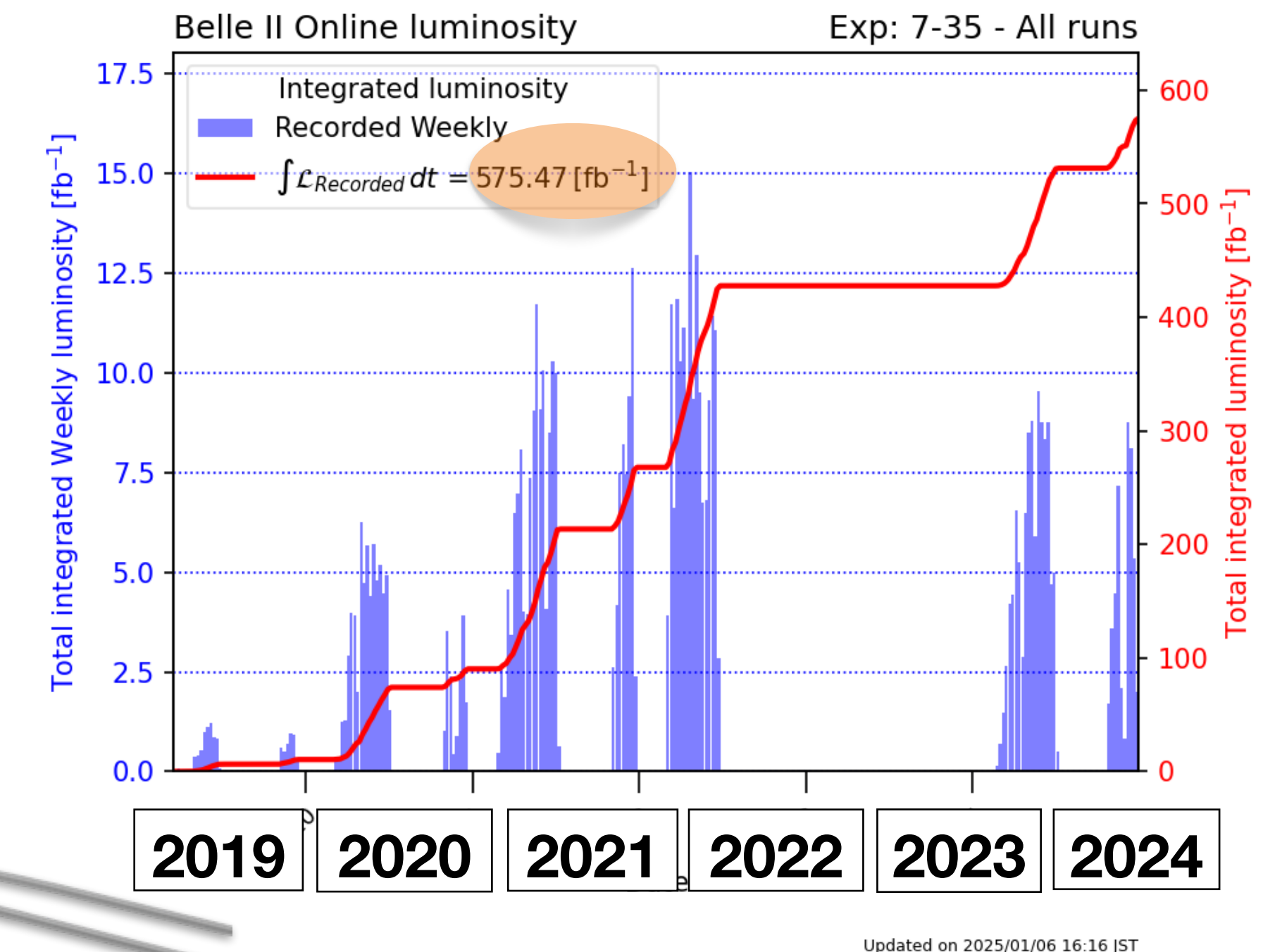Outer 4 layers: strip sensor (SVD)

**Central Drift Chamber (CDC)**
He (50%), $C_2H_6$ (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
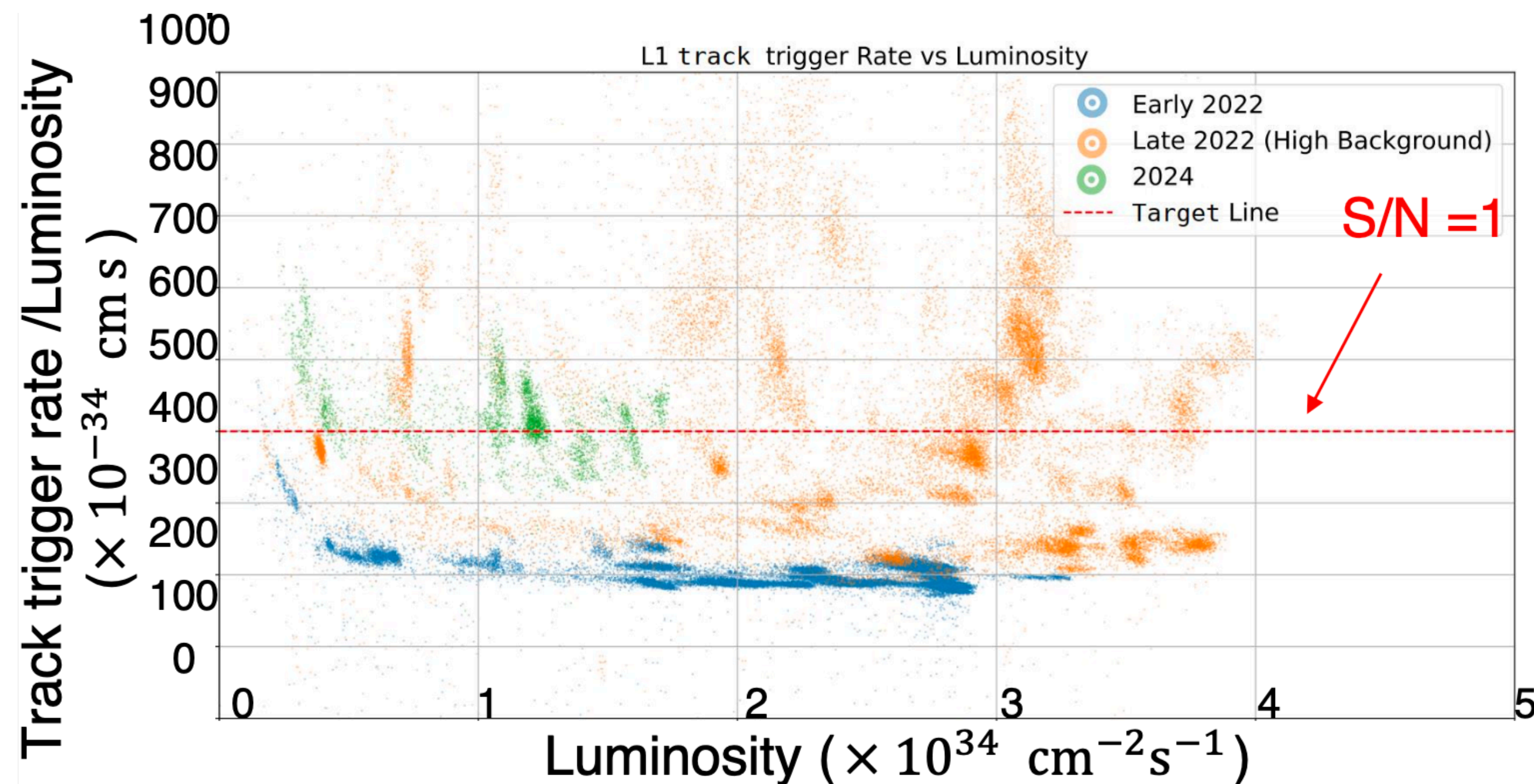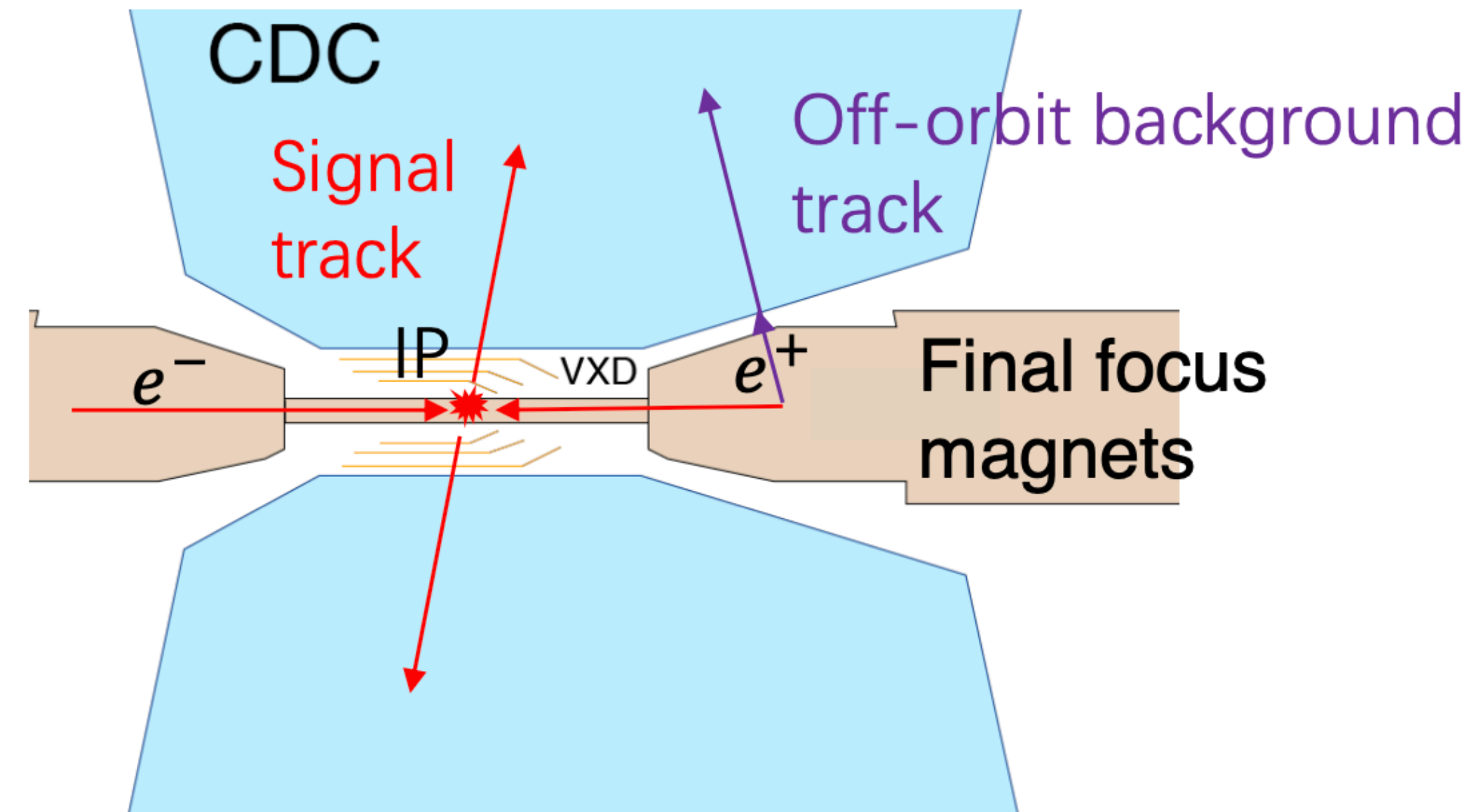
$e^-$ (7GeV)

$e^+$ (4GeV)

$K_L/\mu$ detector (KLM)
Outer barrel: Resistive Plate Counter
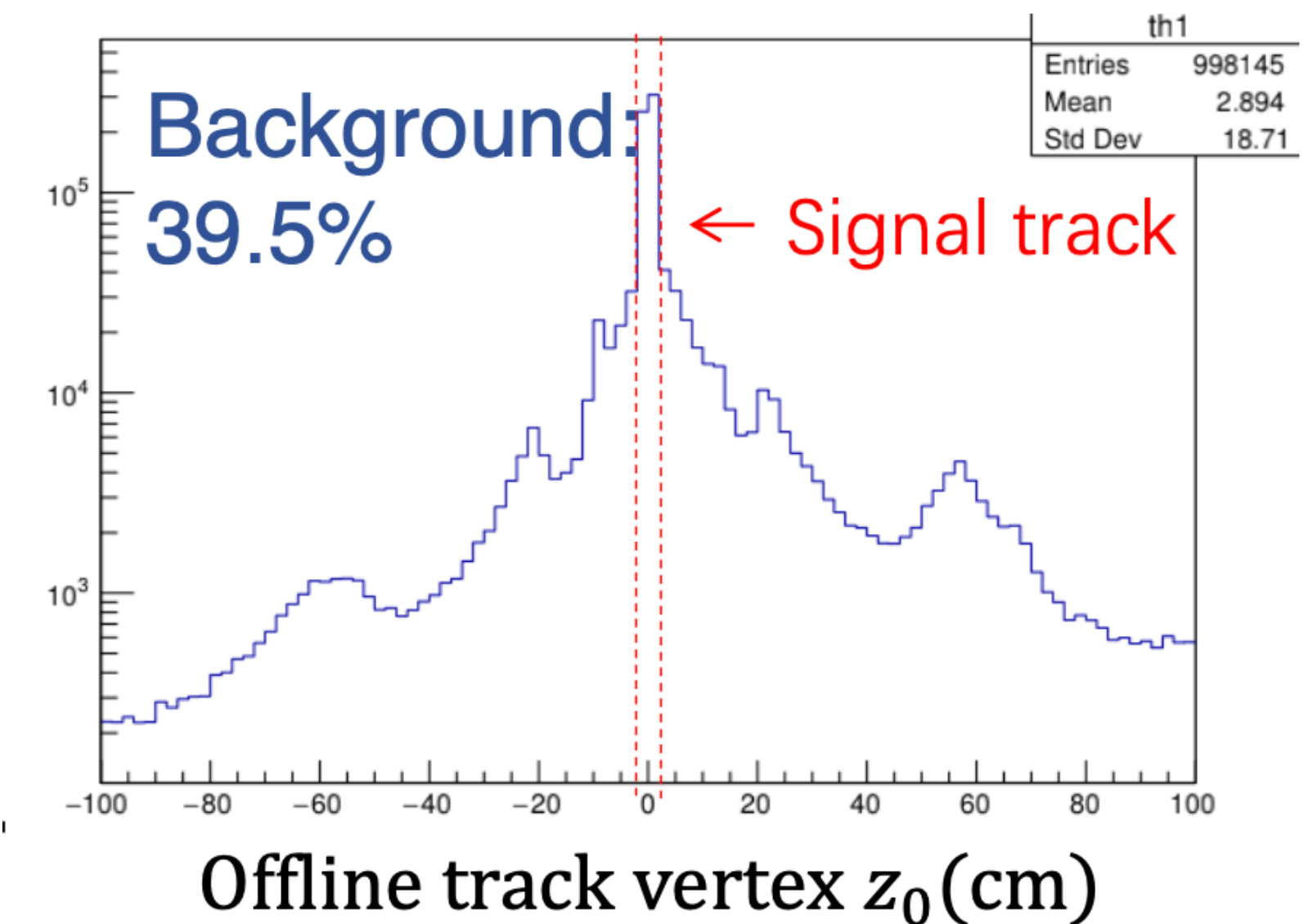(RPC)
Endcap/inner barrel: Scintillator



Belle II Online luminosity        Exp: 7-35 - All runs

Integrated luminosity
Recorded Weekly
$\int \mathcal{L}_{Recorded} \, dt = 575.47 \, [fb^{-1}]$

2019  2020  2021  2022  2023  2024

Updated on 2025/01/06 16:16 JST

- Features:
  - Near-hermetic detector
  - Vertexing and tracking: σ vertex ~ 15μm, CDC spatial res. 100μm $\sigma(P_T)/P_T$ ~ 0.4%
  - Good at measuring neutrals, $\pi^0$, $\gamma$, $K_L$... $\sigma(E)/E$ ~ 2-4%

6

# Motivation of Neural Network for L1 Track trigger
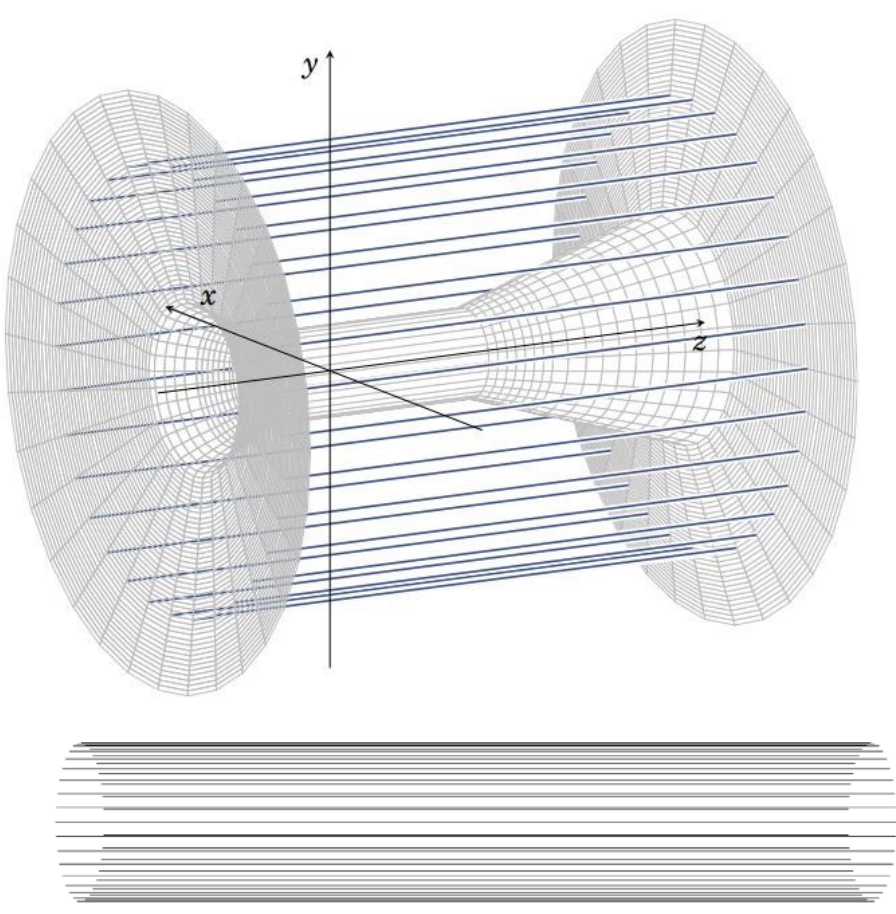
- DAQ system is designed to handle 30 kHz
  - Physical trigger ~15 kHz, require S/N = 1
  - 200350 kHz (total rate) -> ~15 kHz (physics rate)
- L1 trigger rate depends significant on background condition
- Advanced CDC algorithm to further suppress background
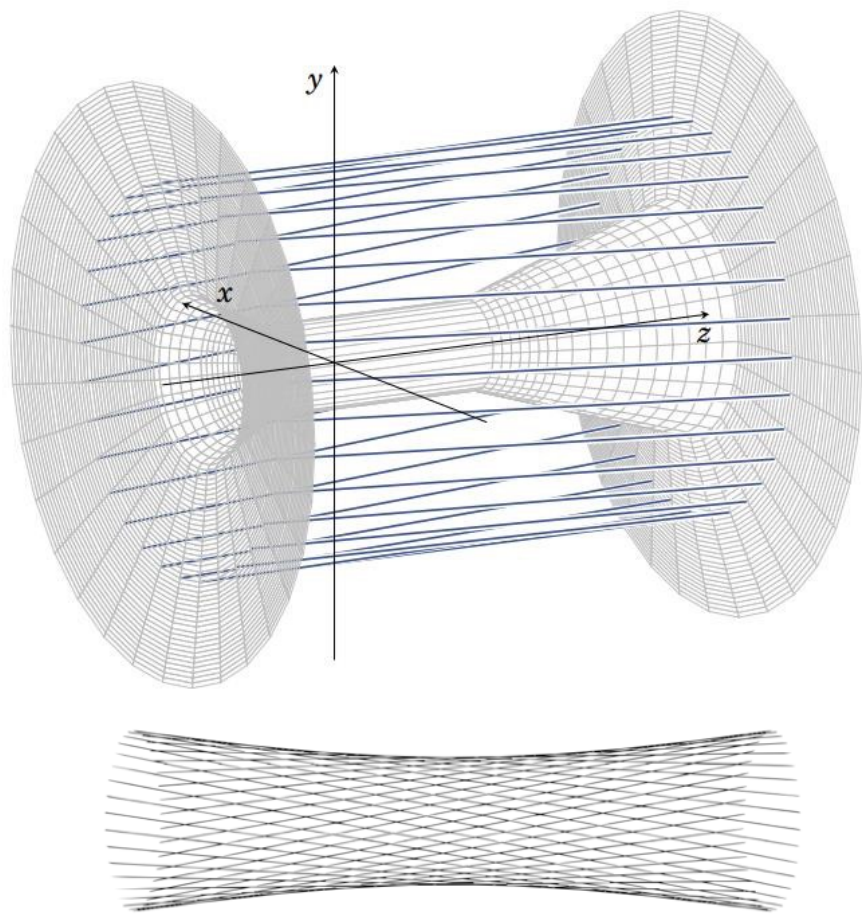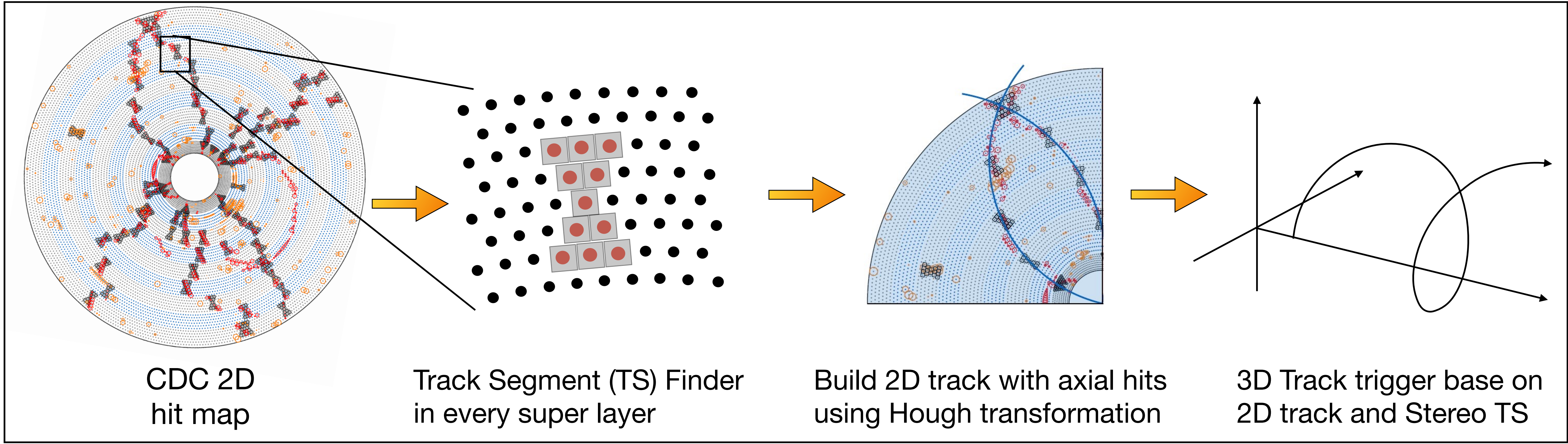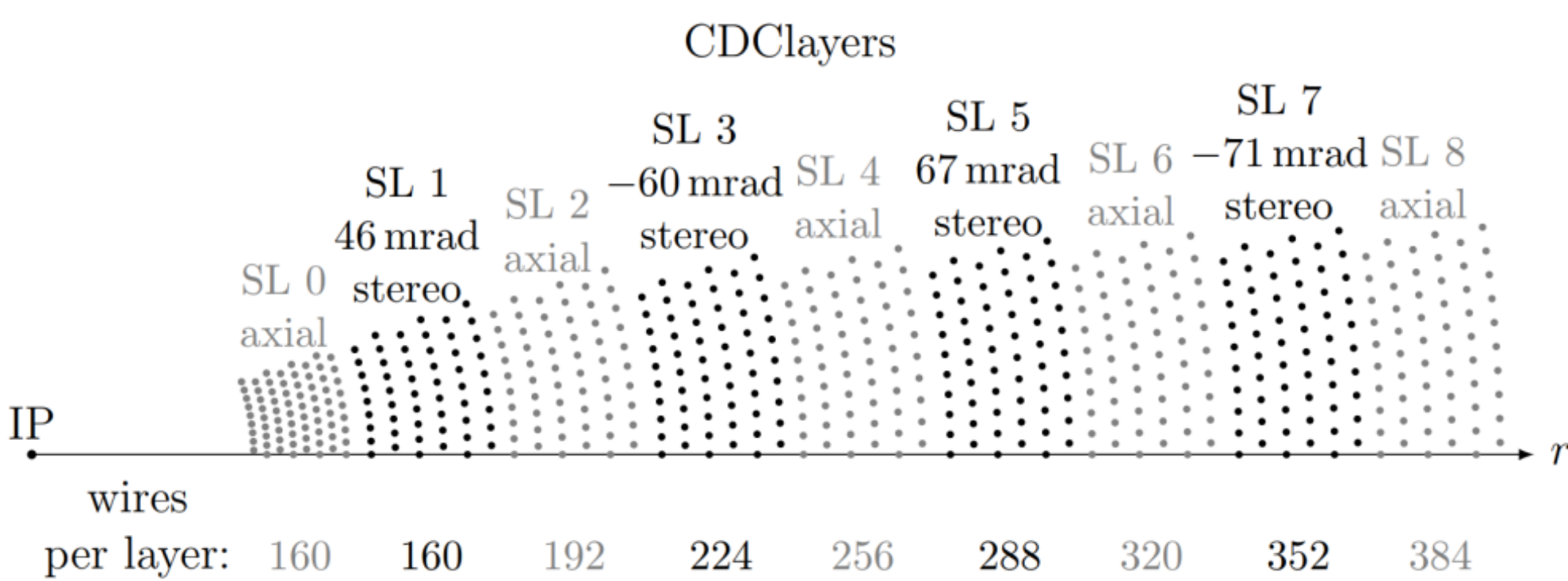- A fixed latency of about **4.4 usec**





L1 track trigger Rate vs Luminosity

- Early 2022
- Late 2022 (High Background)
- 2024
- ---- Target Line

S/N =1

Track trigger rate /Luminosity $(\times 10^{-34}$ cm s $)$

Luminosity $(\times 10^{34}$ cm$^{-2}$s$^{-1})$

**Tracks $z_0$ distribution after trigger**

Background: 39.5%

← Signal track

| th1 | |
|---|---|
| Entries | 998145 |
| Mean | 2.894 |
| Std Dev | 18.71 |

Offline track vertex $z_0$(cm)

# Basics of L1 CDC trigger

$\alpha = 60 \sim 80\text{mrad}$

(BelleII)

$$\sigma_z = \frac{\sigma_{r\phi}}{\sin\alpha}$$
$$\sigma_z = 1.9 \sim 2.5\text{mm}(\sigma_{r\phi} = 0.15\text{mm})$$

CDClayers

SL 7
−71 mrad SL 8
stereo axial

SL 5
67 mrad SL 6
stereo axial

SL 3
−60 mrad SL 4
stereo axial

SL 1
46 mrad SL 2
stereo axial

SL 0
axial

IP

wires
per layer: 160   160   192   224   256   288   320   352   384

Axial wire

Stereo wire

CDC 2D
hit map

Track Segment (TS) Finder
in every super layer

Build 2D track with axial hits
using Hough transformation

3D Track trigger base on
2D track and Stereo TS

8

# Deep Neural Network for Z trigger



- Inputs: **Drift time $t_{\text{drift}}$, 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 |

9

# Development flow of DNN on FPGA



Software → Hardware → FPGA

**Belle II UT4**

- Machine Learning model
- Parameter

- C/C++ transition

- Translate into Verilog/VHDL FPGA language

- Start fitter

- Evaluation

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

*include some function from hls4ml lib

With Python
With Vitis HLS
With Vivado

Convert NN to c++ codes

Adjust Model

Train NN with pytorch (fixed-precision)

Extract weights file

C simulation

Vitis HLS synthesis

RTL co-simulation

Fulfill requirements?

No

Yes

Generate IP

Integrate IP into VHDL codes

Implementation, place and route

Meeting Timing closure?

Commi-ssion

# Quantization-Aware Training (QAT) for FPGA implementation

- Quantization is essential technique to speed up inference, reduce resource usage
  - Embedded system, edge device
  - Fake quantization during training
  - For example, convert 32-bit floating to 8-bit integer
    - Reduction in the model size, memory bandwidth 4x
  - Optimization item: performance vs. latency vs. resource usage

|  | No QAT | QAT |
|---|---|---|
| LUT | ~46% | ~27% |
| DSP | ~64% | ~56% |
| Latency | 551 ns (70 clock) | 488 ns (62 clock) |

# Performance of DNN algorithm

**Delta track z**



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

**Delta track theta**



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

**Classifier output**



*Q* (%)

- Latency : 76 clock = **551.2 ns** ;require: < 600ns
- FPGA resource (UT4: Virtex UltraScale XCVU160) usage:
  - DSP: ~75%, LUT: ~45%, others <30%
- AUC do not get large drop comparing RTL and software simulation
- At signal efficiency ~95%
  - Background rejection rate ~85%

**Background rejection rate**



**Signal efficiency**

12

# 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 electrics
  - Decreased efficiency

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

# Model I: GNN for CDC track background filtering

- Developed a GNN algorithm (based on X. Q. Jia (SDU) et al. BESIII's algorithm)    Xiaoqian Hu (SDU)
  for Belle II CDC hits clean up
  - Inputs: TDC, position coordinates r, φ



node — edge

G = (N, E)

CDC hits produced by charged particles → Construct the graph → Classify the graph edges by GNN → Cluster the selected hits → **Track fitting**

A fully connected 2-layer network → Graph model [ edge network → node network ] edge network → node network ··· edge network →

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

**Belle II** simulation (own work)



μ+ μ- (particle gun)       GNN noise filtering       Transform space       Transform α space       DBSCAN clustering

14

# Model II: GNN for offline track finding

- Find track parameters: momentum, starting position and charge

- Find unknown number of tracks → **Object Condensation** (arXiv:2002.03605)

- Computing resource and time constraint may reducible

Noise filtering ⟶ Clustering ⟶ Fitting

L. Reuter et. al. (KIT) arXiv: 2411.13596



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

# Performance of GNN

- Model II (Object condensation) shows better clustering performance than model I
- Track finding efficiency do not increase that much
  - Failed in the track finding, even clustering of hits shows better performance
  - Further improvement is needed

Belle II simulation (own work)          Model I GNN          Model II GNN

# Training samples for 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!**
  $\rightarrow$ crucial step to be agnostic about the physics processes

- Sample features

  - Low momentum tracks forming circles in the CDC ($P_t$ < 0.4 GeV) <-> High momentum tracks

  - Short tracks <-> tracks penetrate all CDC layers

  - Small opening angle <-> well isolated two tracks

  - …

# 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

L. Reuter et. al. (KIT) arXiv: 2411.13596



18

# CNN algorithm for STCF PID

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

Z. Yao et al.@SDU



EfficientNetV2

# Heterogenous computing platform

- R&D of a new general FPGA device using the AMD Versal ACAP
  - Heterogenous acceleration (VCK190, VCK5000 evaluation kit)
    - AI engine (AIE)
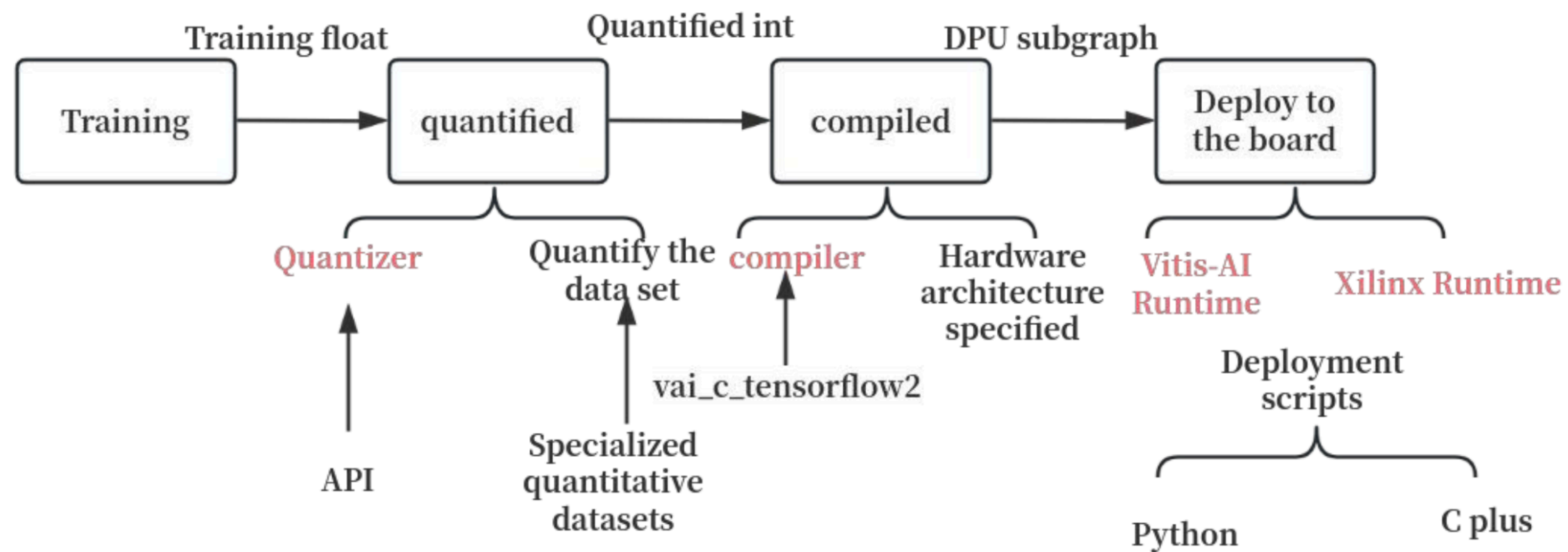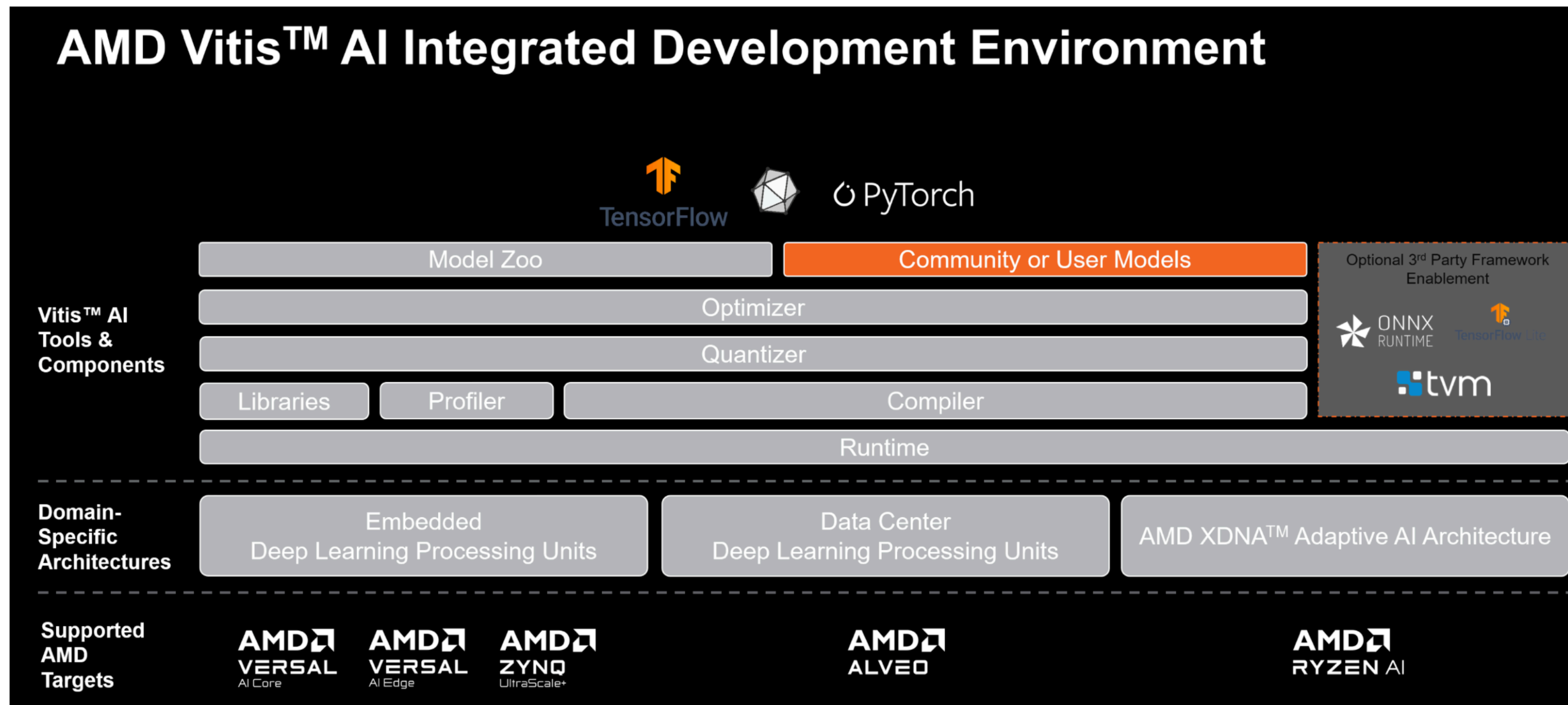
**UG1079**

*Figure 2:* **AI Engine Array**



*Figure 4:* **AI Engine**
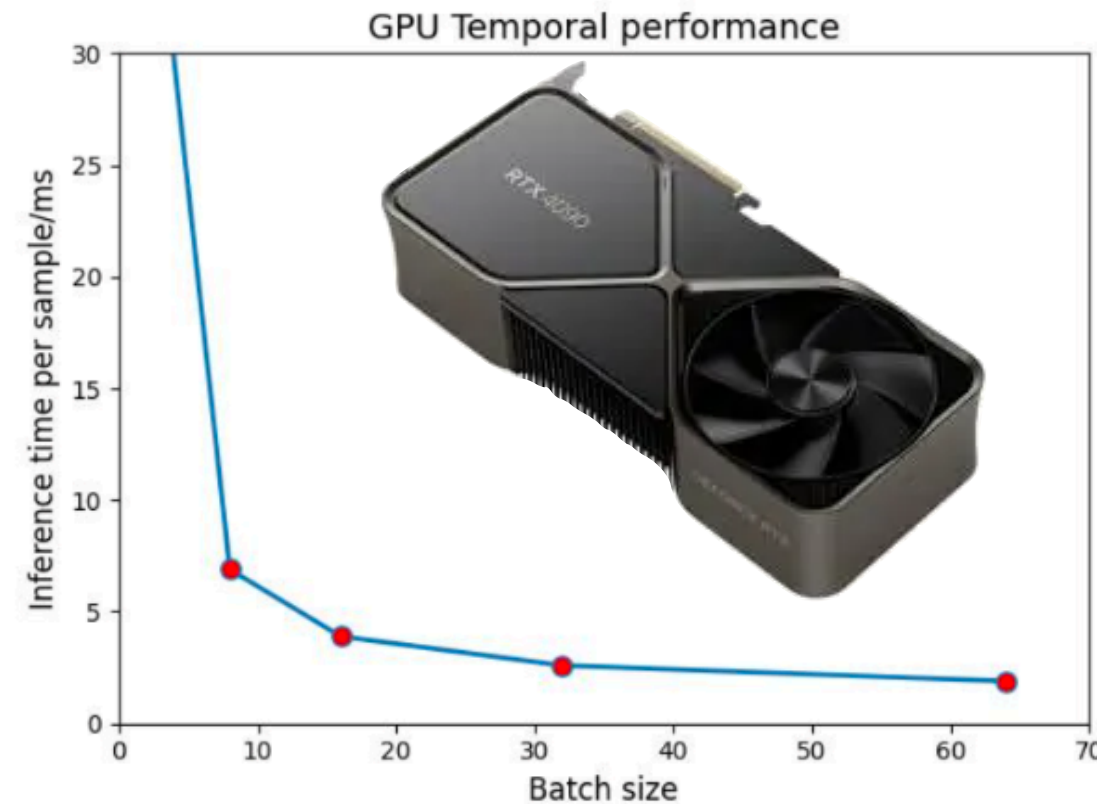
# CNN algorithm implementation

# CNN algorithm implementation

Inference result based on 10000 samples



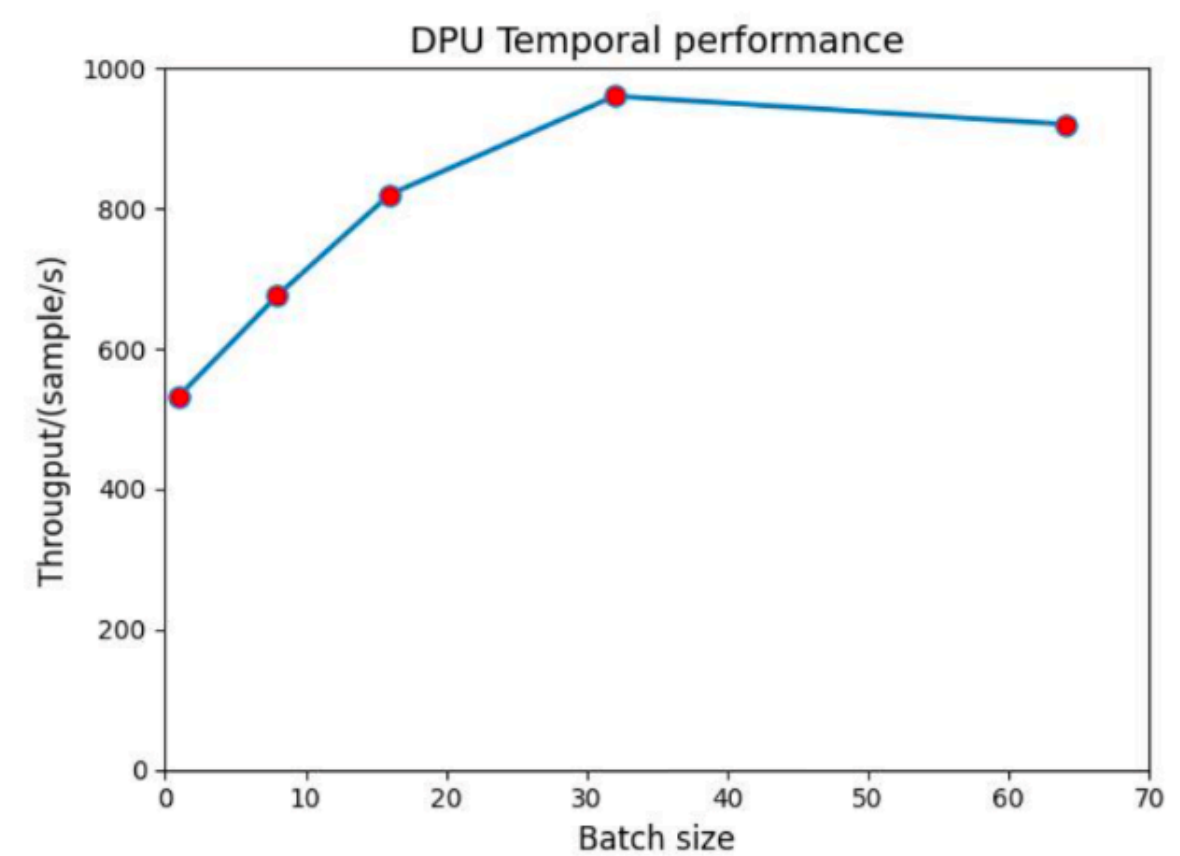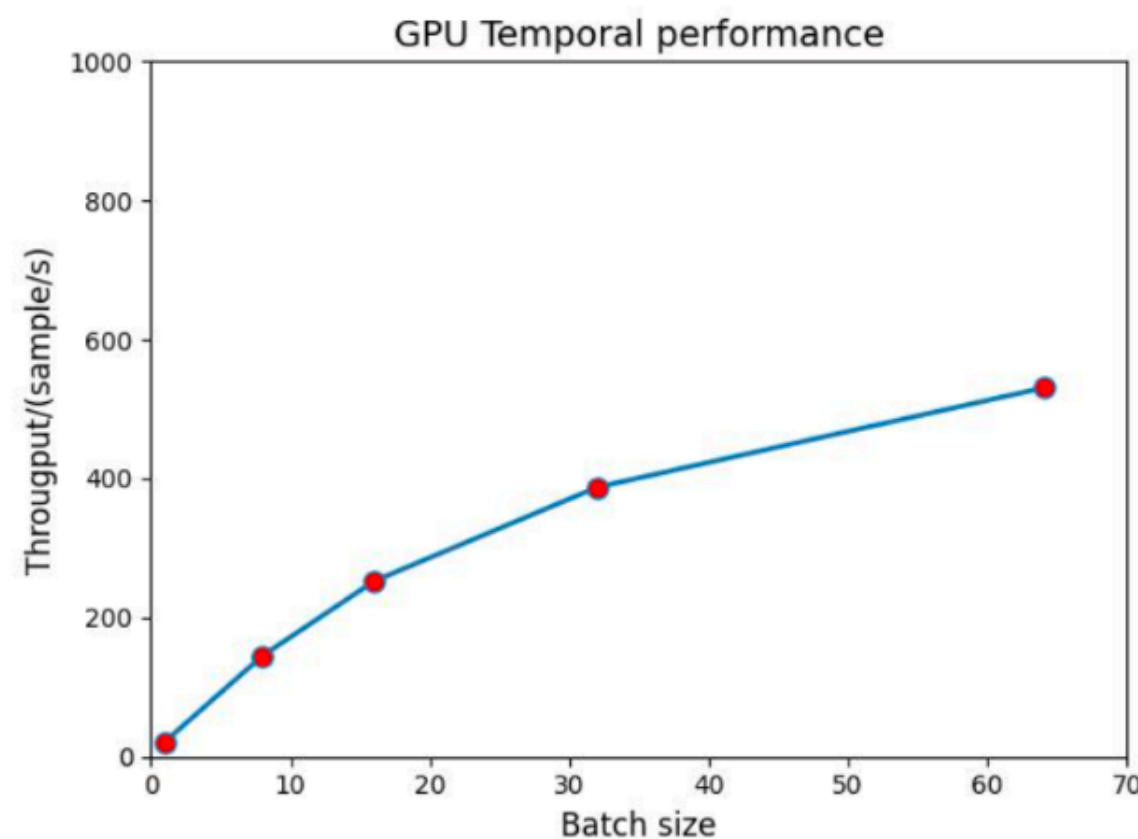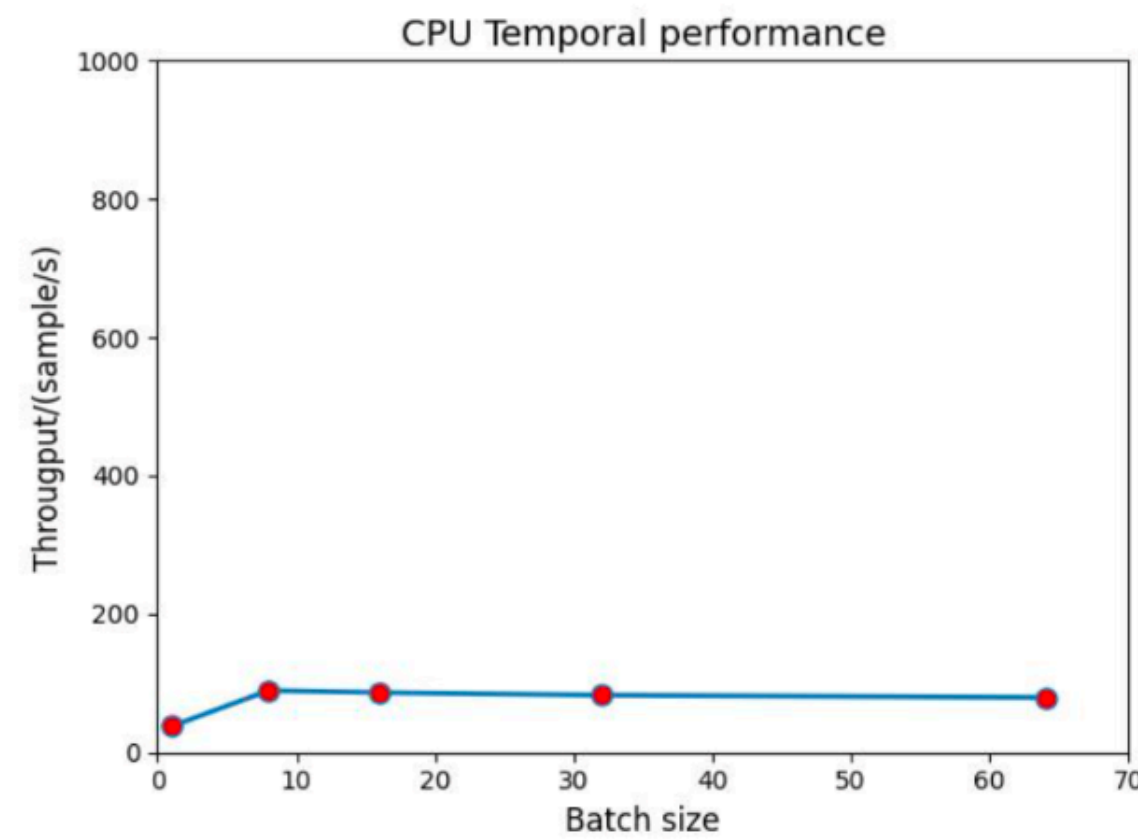DPU based on AMD Versal ACAP shows ~13 times(CPU)/~3(GPU) faster inference time

# Summary and prospects

- Advanced data reduction technique is essential for next-generation HEP experiment.
  - AI/ML integrated with heterogenous computing acceleration
- DNN with hardware based L1 track trigger for improving background rejection
- GNN based hit filter, and DNN were implemented on AMD Versal ACAP
- CNN based PID algorithm for STCF running on DPU has advanced performance than CPU/GPU
  - Both GPU and DPU can be a solution to accelerate data processing for online or offline data processing system