# Development of DNN trigger for Belle II experiment

Shuangshuang Zhang, Qi-Dong Zhou

Institute of frontier and interdisciplinary science, Shandong University
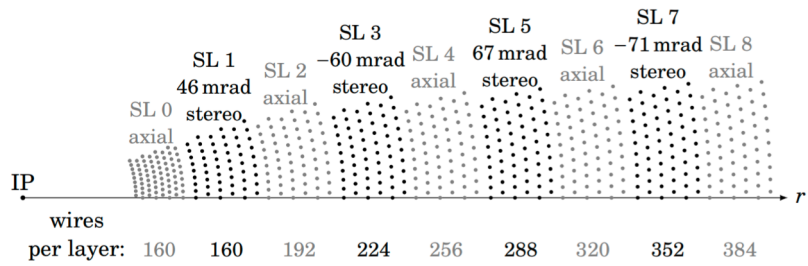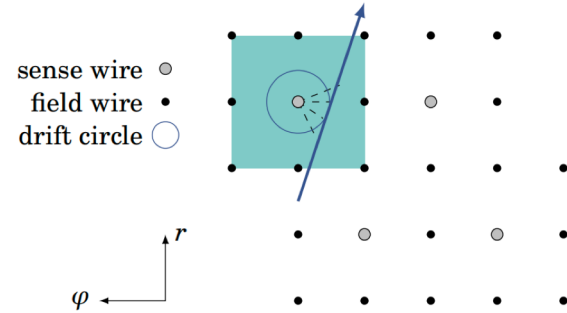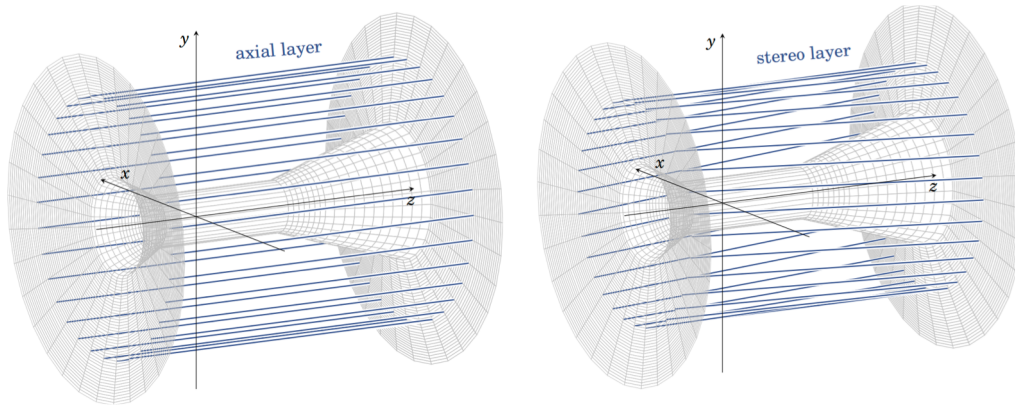
12 Dec. 2024

CEPC TDAQ meeting

# Outline

- Motivation

- Deep Neural Network (DNN) model

- Development of DNN model with python

- High Level Syntheis (HLS) with Vitis
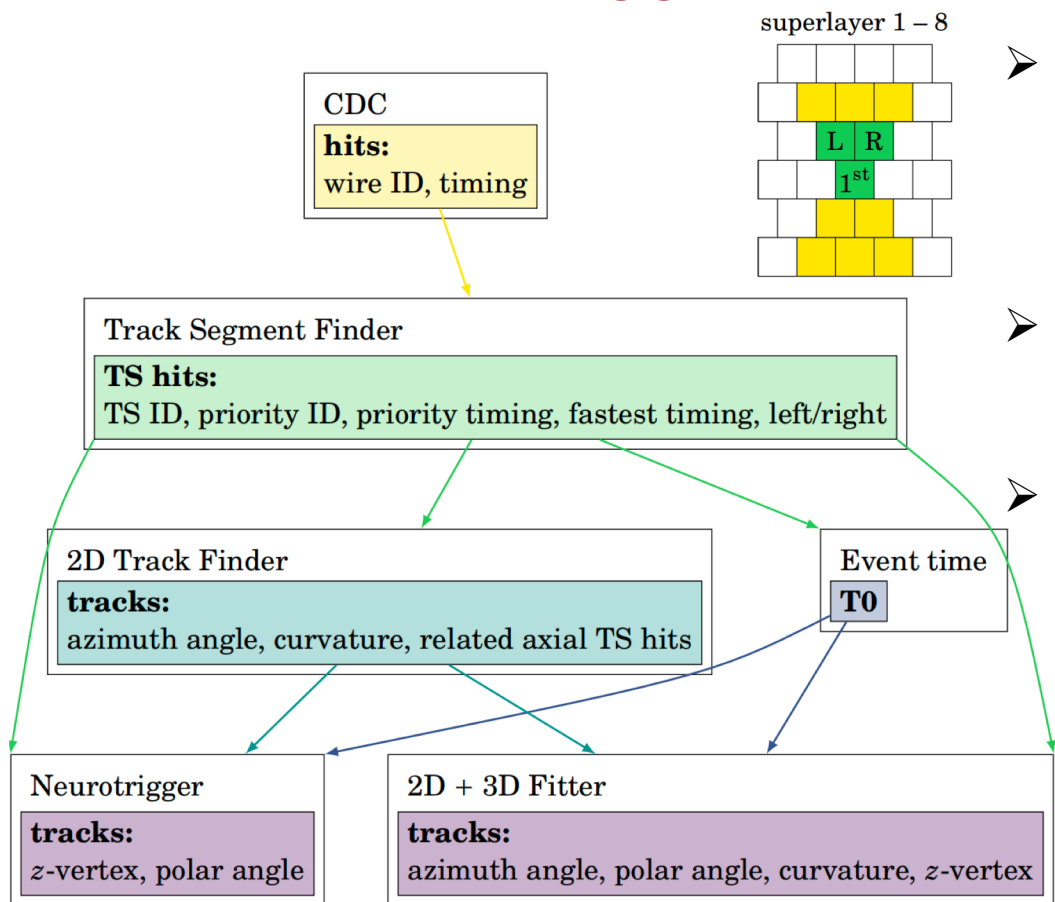
- Testing IP cores in Vivado

# Motivation

➢ Study the CDC DNN trigger and master the implementation process of

   DNN in FPGA

➢ Optimize the model and observe changes in performance

➢ PLAN: Investigating the possiblity of using more input variables in DNN

   models, such as ADC signals and momentum

➢ PLAN: Development of DNN model to Versal ACAP

# Central drift chamber





- ➤ A drift cell of the CDC is formed of one sense wire and eight field wires
- ➤ A hit on an axial wire provides coordinates in the transverse plane(2D)($\Phi$,r)
- ➤ By combining axial and stereo hits, a three dimensional track can be reconstructed(3D)

# The CDC track trigger



- ➢ The Belle II trigger consists of two levels: the first level trigger(hardware) and the high level trigger(software)

- ➢ CDC track trigger provides charged track information

- ➢ The track finding for the trigger is based on track segments("TS")

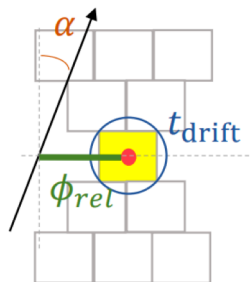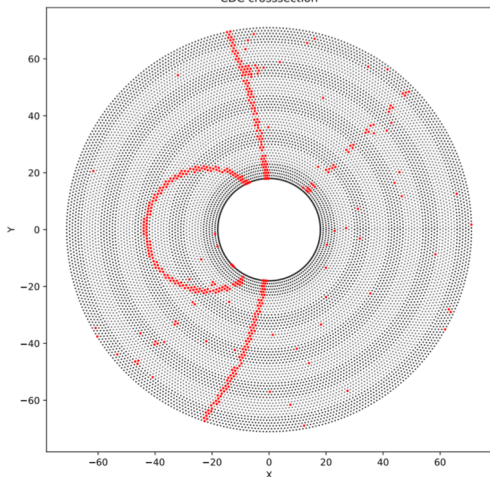# Deep Neural Network (DNN) model



CDC crosssection



$\alpha$, $t_{drift}$, $\phi_{rel}$



Single head self attention

Soft-max

weights

value

Track Fitting

$z_0$
$\theta$
$Q$

- Extract a track segment from each super layer

- Every track segment (TS) contains a set of $\alpha$, $t_{drift}$ and $\Phi$

- DNN trigger : 2D track candidates+ Drift time for all other wires in the stereo wire

- Input : $\alpha$, $t_{drift}$ , phi of priority wire, $t_{drift}$ of all other 10 wires for every TS($3*9+11*4=71$)

- Output : $z_0$, $\theta_0$, classifier output (Q)

2024-12.12

4

# Development of DNN model with python



- Using pytorch lib for model building and training
- Add two hidden layers
- Modify learning rates

Initial

Development

# Workflow of DNN model development on FPGA

```
Train DNN with    →    Extract weights    →    Import file in    →    C simulation
pytorch                and bias                vitis hls
                                                                          ↓
Using IP Core in  ←    Generate IP       ←    C/RTL           ←    C synthesis
vivado                                        cosimulation
```



```
▾ C SIMULATION
  ▶ Run C Simulation
  ▸ Reports & Viewers
▾ C SYNTHESIS
  ▶ Run C Synthesis
  ▸ Reports & Viewers
▾ C/RTL COSIMULATION
  ▶ Run Cosimulation
  ▸ Reports & Viewers
▾ IMPLEMENTATION
  ▶ Export RTL
  ▶ Run Implementation
  ▸ Reports & Viewers
```

```
|| + ap_ctrl
—  in_data_ap_vld
—  ap_clk                    layer15_out_ap_vld  —
—  ap_rst                    layer15_out[59:0]   —
—  in_data[925:0]

            Nntrg_hls
```

# Development flow with Vitis HLS

A typical flow using HLS has the following steps:

- Convert python based model to C/C++ model with hls4ml (Optional)

- Write the algorithm using C/C++ with a target architecture in mind

- Verify the functionality at the behavioral level

- Generate the RTL based model

- Verify the functionality of the generated RTL model

Optimizing performance: "#pragma HLS"

```
#pragma HLS array_partition variable = w2 complete dim = 0
```

- Generate multiple small memories

- Increases the amount of read and write ports for the storage

# Performance



Delta track z

baseline(NN)/software(DNN)/RTL(HLS)

Virtex Ultra
XCVU160

Versal
AI Core
XCVC1902

# Latency(testbench)



initial:
74clock

optimize:
81clock

2024-12.12

9

# Next plan

➢ Further optimize the model and perform pruning

➢ Investigating the possiblity of using more input variables in DNN models, such

  as ADC signals and momentum(Hope to receive more suggestions)

  • Check the size of the model scale(Utilization, latency)

  • Predicting the use of the next generation UT5 board

➢ Deploy the neutral networks to Versal ACAP

  • First step: deploy the DNN model into Versal AI engine

  • Further plan: deploy the NN model on Versal DPU

# Thanks for your listening

Back up

# Virtex UltraScale XCVU160



initial

```
+----------------------------------------------------------------------+
| Design Summary                                                       |
| impl_1                                                               |
| xcvu160_CIV-flgc2104-2-e                                             |
+----------------------------------------------------------------------+
| Criteria                                    | Guideline | Actual | Status |
+----------------------------------------------------------------------+
| LUT                                         | 70%    | 50.79% | OK     |
| FD                                          | 50%    | 10.22% | OK     |
| LUTRAM+SRL                                  | 25%    | 7.98%  | OK     |
| CARRY8                                      | 25%    | 29.50% | REVIEW |
| MUXF7                                       | 15%    | 0.00%  | OK     |
| LUT Combining                               | 20%    | 11.29% | OK     |
| DSP                                         | 80%    | 68.33% | OK     |
| RAMB/FIFO                                   | 80%    | 0.26%  | OK     |
| DSP+RAMB+URAM (Avg)                          | 70%    | 34.30% | OK     |
| BUFGCE* + BUFGCTRL                           | 24     | 2      | OK     |
| DONT_TOUCH (cells/nets)                     | 0      | 0      | OK     |
| MARK_DEBUG (nets)                           | 0      | 0      | OK     |
| Control Sets                                | 17370  | 266    | OK     |
| Average Fanout for modules > 100k cells     | 4      | 2.44   | OK     |
| Max Average Fanout for modules > 100k cells | 4      | 3.67   | OK     |
| Non-FD high fanout nets > 10k loads         | 0      | 4      | REVIEW |
+----------------------------------------------------------------------+
| TIMING-6 (No common primary clock between related clocks) | 0 | 0 | OK |
| TIMING-7 (No common node between related clocks)          | 0 | 0 | OK |
| TIMING-8 (No common period between related clocks)        | 0 | 0 | OK |
| TIMING-14 (LUT on the clock tree)                         | 0 | 0 | OK |
| TIMING-35 (No common node in paths with the same clock)   | 0 | 0 | OK |
+----------------------------------------------------------------------+
| Number of paths above max LUT budgeting (0.425ns) | 0 | 0 | OK     |
| Number of paths above max Net budgeting (0.298ns) | 0 | 1 | REVIEW |
+----------------------------------------------------------------------+
```

optimize

```
+----------------------------------------------------------------------+
| Design Summary                                                       |
| impl_1                                                               |
| xcvu160-flgb2104-2-e                                                 |
+----------------------------------------------------------------------+
| Criteria                                    | Guideline | Actual | Status |
+----------------------------------------------------------------------+
| LUT                                         | 70%    | 64.08% | OK     |
| FD                                          | 50%    | 12.59% | OK     |
| LUTRAM+SRL                                  | 25%    | 8.13%  | OK     |
| CARRY8                                      | 25%    | 37.98% | REVIEW |
| MUXF7                                       | 15%    | 0.00%  | OK     |
| LUT Combining                               | 20%    | 12.84% | OK     |
| DSP                                         | 80%    | 96.03% | REVIEW |
| RAMB/FIFO                                   | 80%    | 0.26%  | OK     |
| DSP+RAMB+URAM (Avg)                          | 70%    | 48.15% | OK     |
| BUFGCE* + BUFGCTRL                           | 24     | 2      | OK     |
| DONT_TOUCH (cells/nets)                     | 0      | 0      | OK     |
| MARK_DEBUG (nets)                           | 0      | 0      | OK     |
| Control Sets                                | 17370  | 255    | OK     |
| Average Fanout for modules > 100k cells     | 4      | 2.57   | OK     |
| Max Average Fanout for modules > 100k cells | 4      | 3.59   | OK     |
| Non-FD high fanout nets > 10k loads         | 0      | 6      | REVIEW |
+----------------------------------------------------------------------+
| TIMING-6 (No common primary clock between related clocks) | 0 | 0 | OK |
| TIMING-7 (No common node between related clocks)          | 0 | 0 | OK |
| TIMING-8 (No common period between related clocks)        | 0 | 0 | OK |
| TIMING-14 (LUT on the clock tree)                         | 0 | 0 | OK |
| TIMING-35 (No common node in paths with the same clock)   | 0 | 0 | OK |
+----------------------------------------------------------------------+
| Number of paths above max LUT budgeting (0.425ns) | 0 | 0 | OK     |
| Number of paths above max Net budgeting (0.298ns) | 0 | 0 | OK     |
+----------------------------------------------------------------------+
```

```
Implementation tool: Xilinx Vivado v.2023.2
Project:          DNN_1
Solution:         solution1
Device target:    xcvu160_CIV-flgc2104-2-e
Report date:      Sat Dec 07 01:58:02 CST 2024

#=== Post-Implementation Resource usage ===
SLICE:          0
LUT:          470542
FF:           189296
DSP:            1066
BRAM:             17
URAM:              0
LATCH:             0
SRL:           17550
CLB:           73455

#=== Final timing ===
CP required:                     4.000
CP achieved post-synthesis:      4.410
CP achieved post-implementation: 7.624
Timing not met
```

```
Implementation tool: Xilinx Vivado v.2023.2
Project:          new_hidelayer
Solution:         solution1
Device target:    xcvu160-flgb2104-2-e
Report date:      Sat Dec 07 02:53:23 CST 2024

#=== Post-Implementation Resource usage ===
SLICE:            0
LUT:          593640
FF:           233306
DSP:            1498
BRAM:             17
URAM:              0
LATCH:             0
SRL:           17871
CLB:           89668

#=== Final timing ===
CP required:                     4.000
CP achieved post-synthesis:      4.324
CP achieved post-implementation: 9.196
Timing not met
```

# Some questions