# Design Ideas for an Online Data Reduction System for the ePIC dRICH Detector

Alessandro Lonardo

INFN Roma, APE Lab

for the ePIC Roma1/2 team

ePIC General Meeting
January 21th 2025

# ePIC dRICH

Compact cost-effective solution for particle identification in the high-energy endcap at EIC

## dRICH
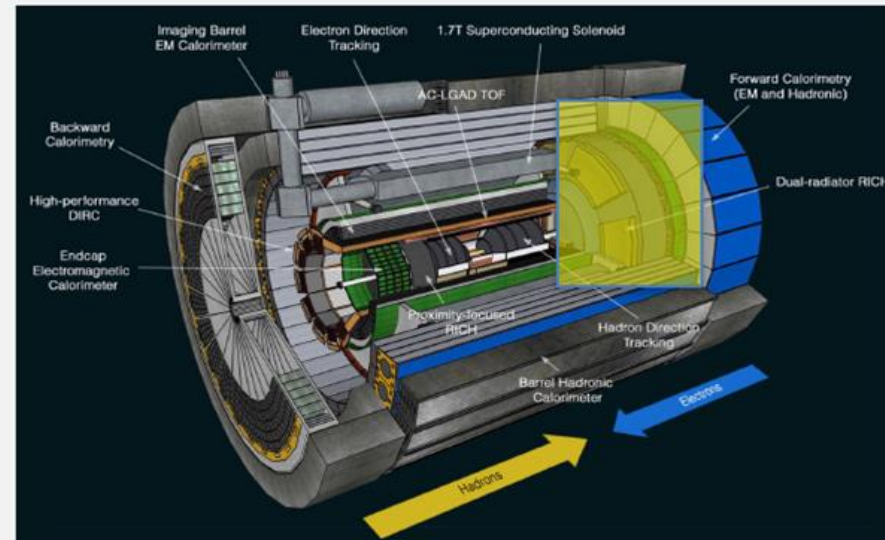


BA, BO, CS, CT, FE, GE, LNS, RM1, RM2, SA, TO, TS

## EPIC



## EIC RICH Consortium
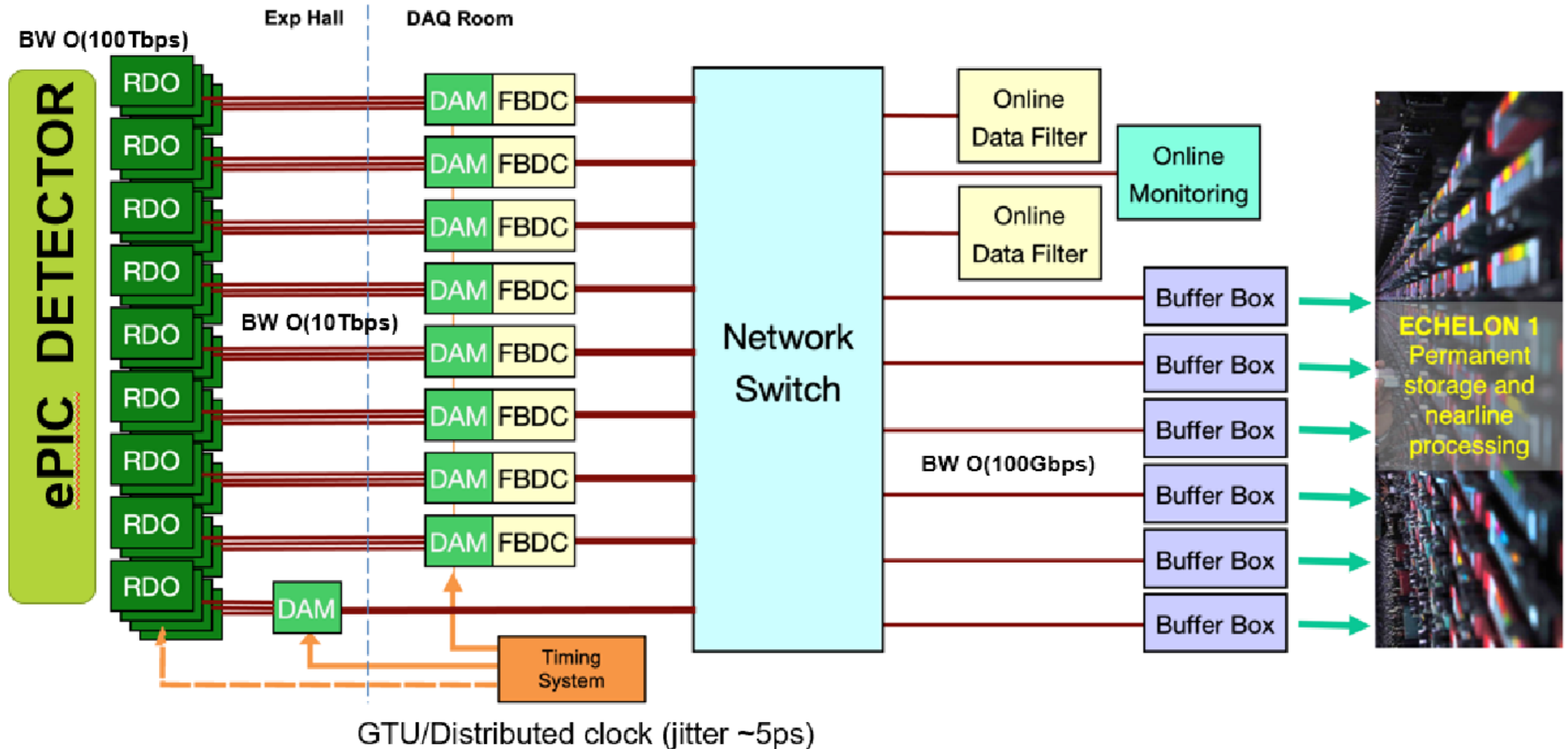


Forward particle detection

Hadron ID in the extended 3-50 GeV/c interval

Support electron ID up to 15 GeV/c

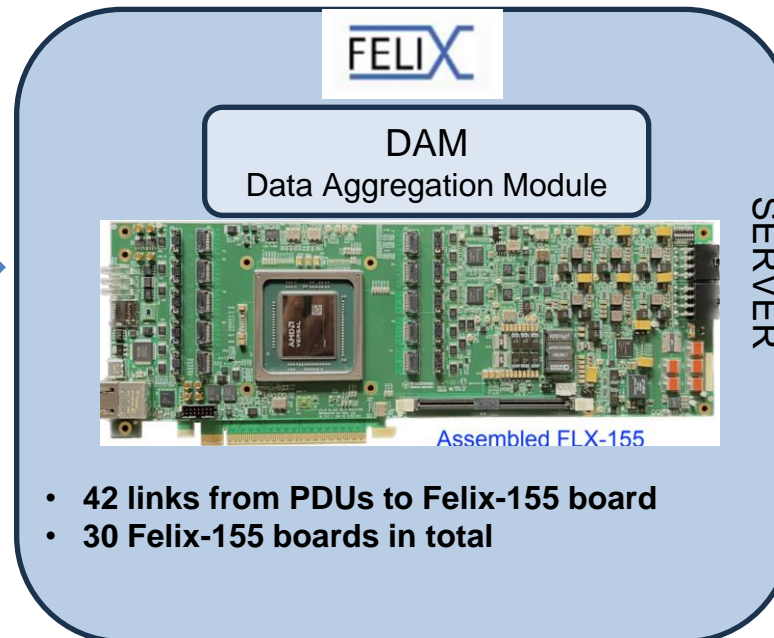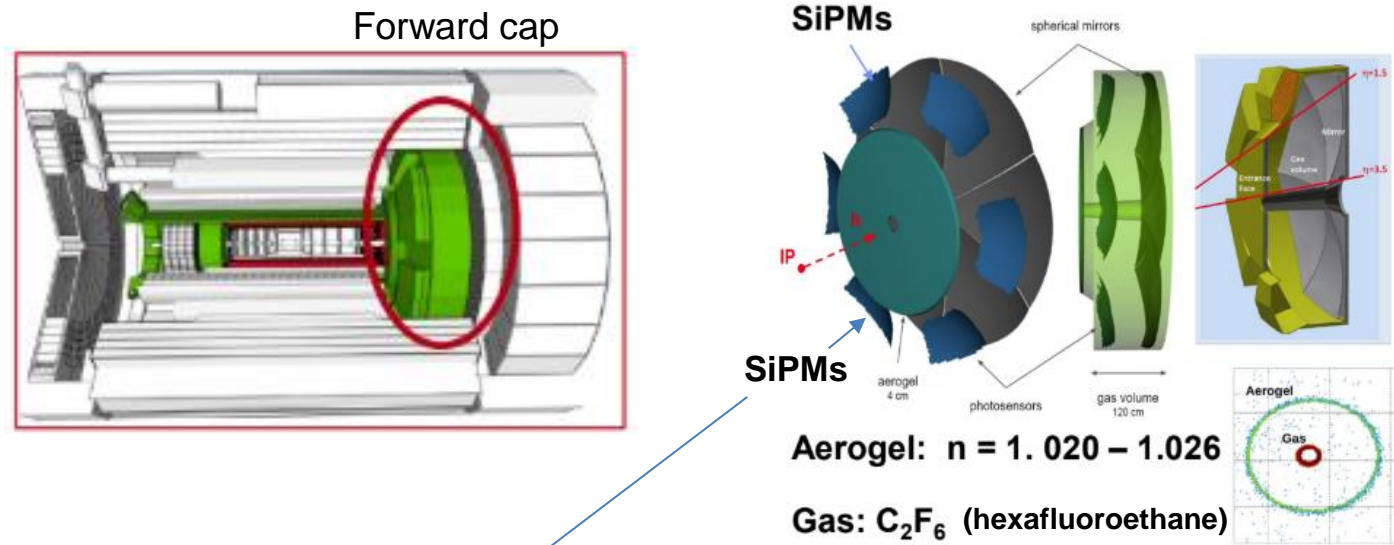**Main challenges:**

**Cover wide momentum range  3 - 50 GeV/c**    -> dual radiator
**Work in high (~ 1T) magnetic field**              -> SiPM
**Fit in a quite limited (for a gas RICH) space**   -> curved detector

# ePIC General DAQ Scheme

# Dual Radiator RICH (dRICH)

Forward cap

SiPMs

spherical mirrors

SiPMs

aerogel 4 cm   photosensors   gas volume 120 cm

Aerogel: $n = 1.020 - 1.026$

Gas: $C_2F_6$ (hexafluoroethane)

Aerogel

Gas

PDU

FELIX

DAM
Data Aggregation Module

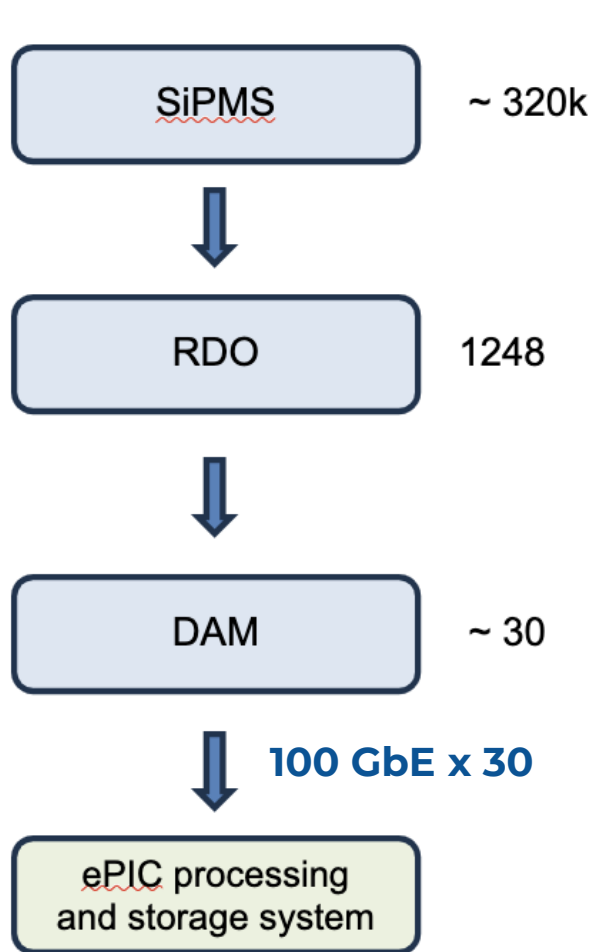SERVER

Assembled FLX-155

ePIC processing
and storage system

- 1 photodetector unit PDU: 4x64 SiPM array device (256 channels), 4 FEBs, 1 RDO
- 1248 PDUs for full dRICH readout
- 319488 readout channels divided in six sectors

- 42 links from PDUs to Felix-155 board
- 30 Felix-155 boards in total

4

# Analysis of dRICH Output Bandwidth

The dRICH DAQ chain in ePIC ➜ **the throughput issue**

SiPMS — ~ 320k

RDO — 1248

DAM — ~ 30

**100 GbE x 30**

ePIC processing and storage system

| dRICH DAQ parameters | |
|---|---|
| RDO boards | 1248 |
| ALCOR64 x RDO | 4 |
| dRICH channels (total) | 319488 |
| Number of DAM L1 | 27 |
| Input link in DAM L1 | 47 |
| Output links in DAM L1 | 1 |
| Number of DAM L2 | 1 |
| Input link to DAM L2 | 27 |
| Link bandwidth [ Gb/s] (assumes VTRX+) | 10 |
| Interaction tagger reduction factor | 1 |
| Interaction tagger latency [s] | 2,00E-03 |
| **EIC parameters** | |
| EIC Clock [MHz] | 98,522 |
| Orbit efficiency (takes into account gap) | 0,92 |

| Bandwidth analysis | | Limit |
|---|---|---|
| Sensor rate per channel [kHz] | 300,00 ▾ | 4.000,00 |
| Rate post-shutter [kHz] | 55,20 | 800,00 |
| Throughput to serializer [ Mb/s] | 34,50 | 788,16 |
| Throughput from ALCOR64 [Mb/s] | 276,00 | |
| Throughput from RDO [ Gb/s] | 1,08 | 10,00 |
| Input at each DAM I [Gbps] | 50,67 | 470,00 |
| Buffering capacity at DAM I [MB] | 12,97 | |
| Output from every DAM | 50,67 | 10,00 |
| **Total throughput** | 1.368,14 | 270,00 |

- Sensors DCR: 3-300 kHz (increasing with radiation damage ➜ with experiment lifetime).
- Full detector throughput (FE): 14-1400Gbps
- **A reduction is needed** to match 30 channels aggregated bandwidth (and safety margin)
- EIC beams bunch spacing: 10 ns ➜ bunch crossing rate of 100 MHz
- For the low interaction cross-section (DIS) ➜ one interaction every ~100 bunches ➜ interaction rate of ~1MHz.
- **A system tagging the (DIS) interacting bunches** could solve the issue reducing down to ~1/100 the data throughput

**Two complementary approaches are possible:**
1. **Develop a dedicated sub-detector tagging relevant interactions.**
2. **This proposal.**

# dRICH: Data Reduction

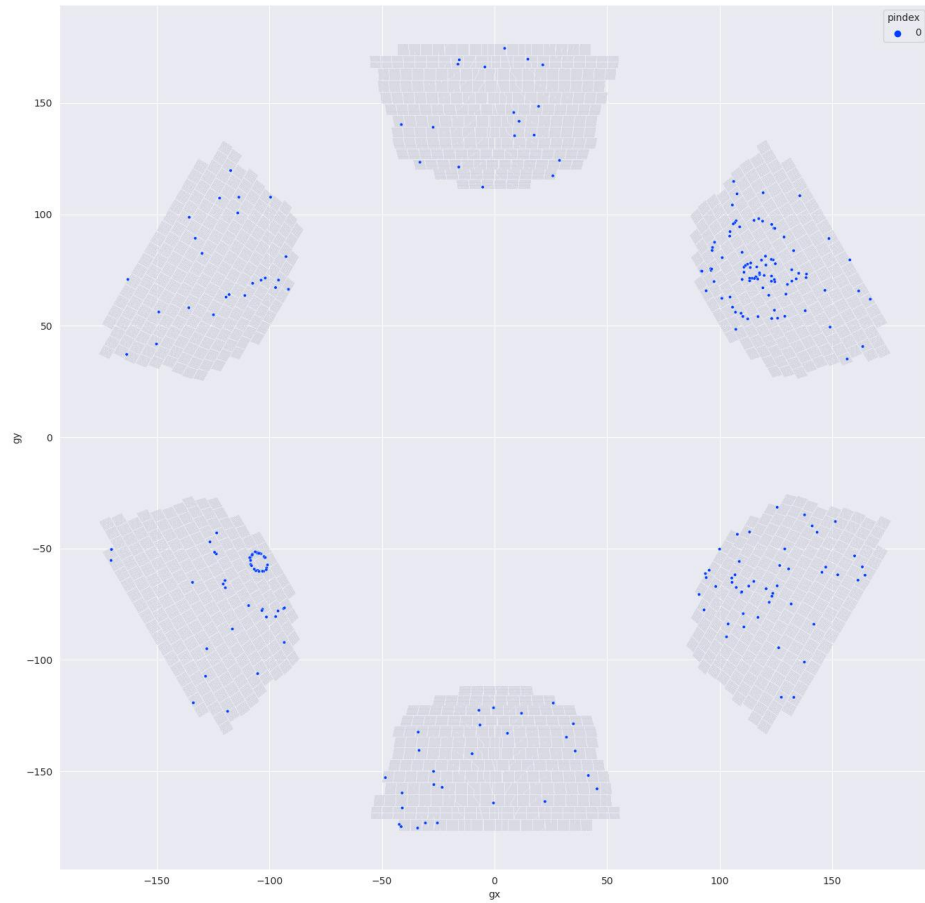**Online Signal/ Noise discrimination using ML**

- **Signal (i.e. Merged Phys Signal + Bkg)**:
  - **Physics Signal:**
    - e.g DIS
  - **Physics Background:**
    - e/p with beam pipe
    - Synchrotron radiation (currently not including it)

- **SiPM Noise:**
  - Dark current rate (DCR) modelled in the reconstruction stage (*recon.rb* eic-shell method)
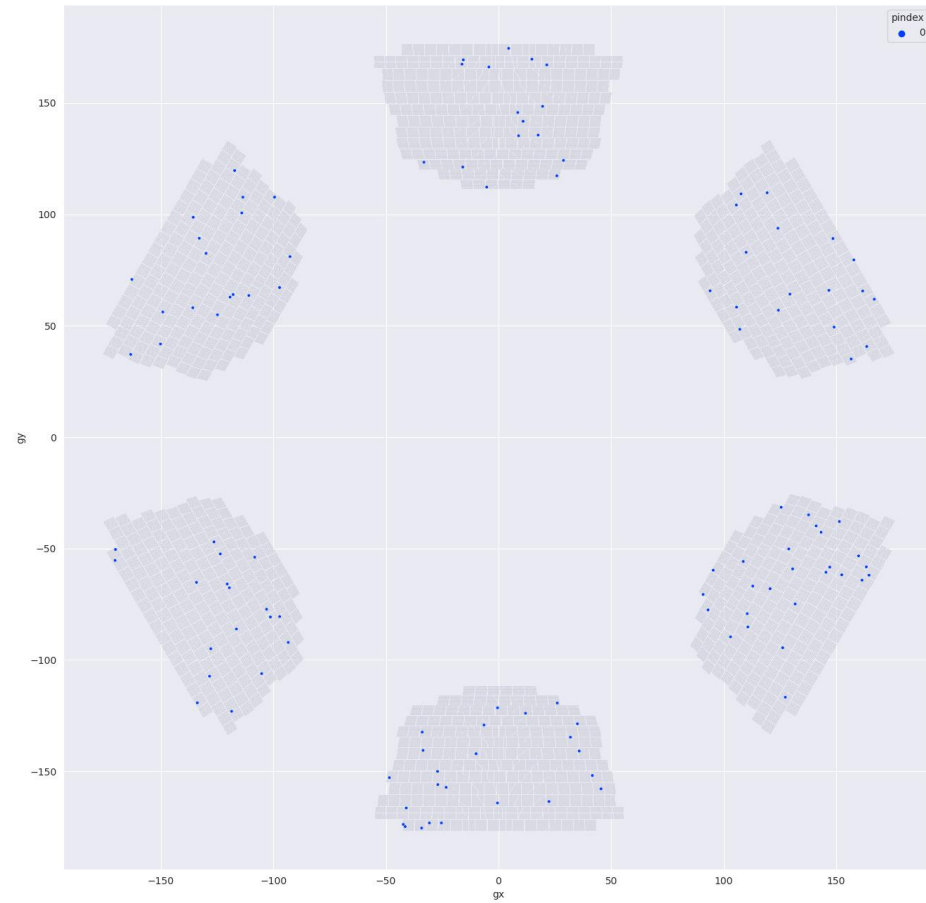
**ML task:**
Discriminate between **Noise Only** and **Signal + Noise** events

# dRICH: Dataset for training, classes
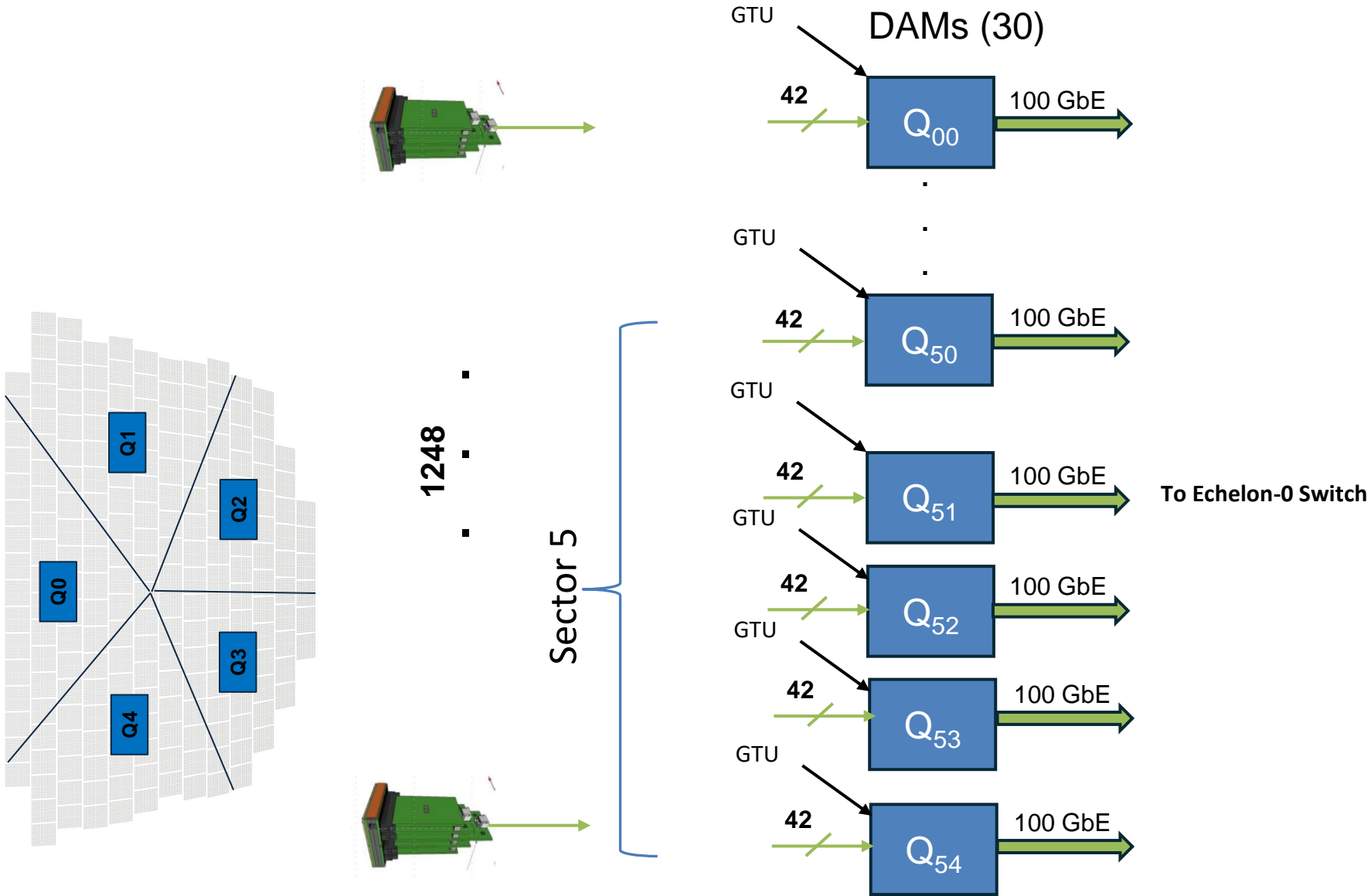
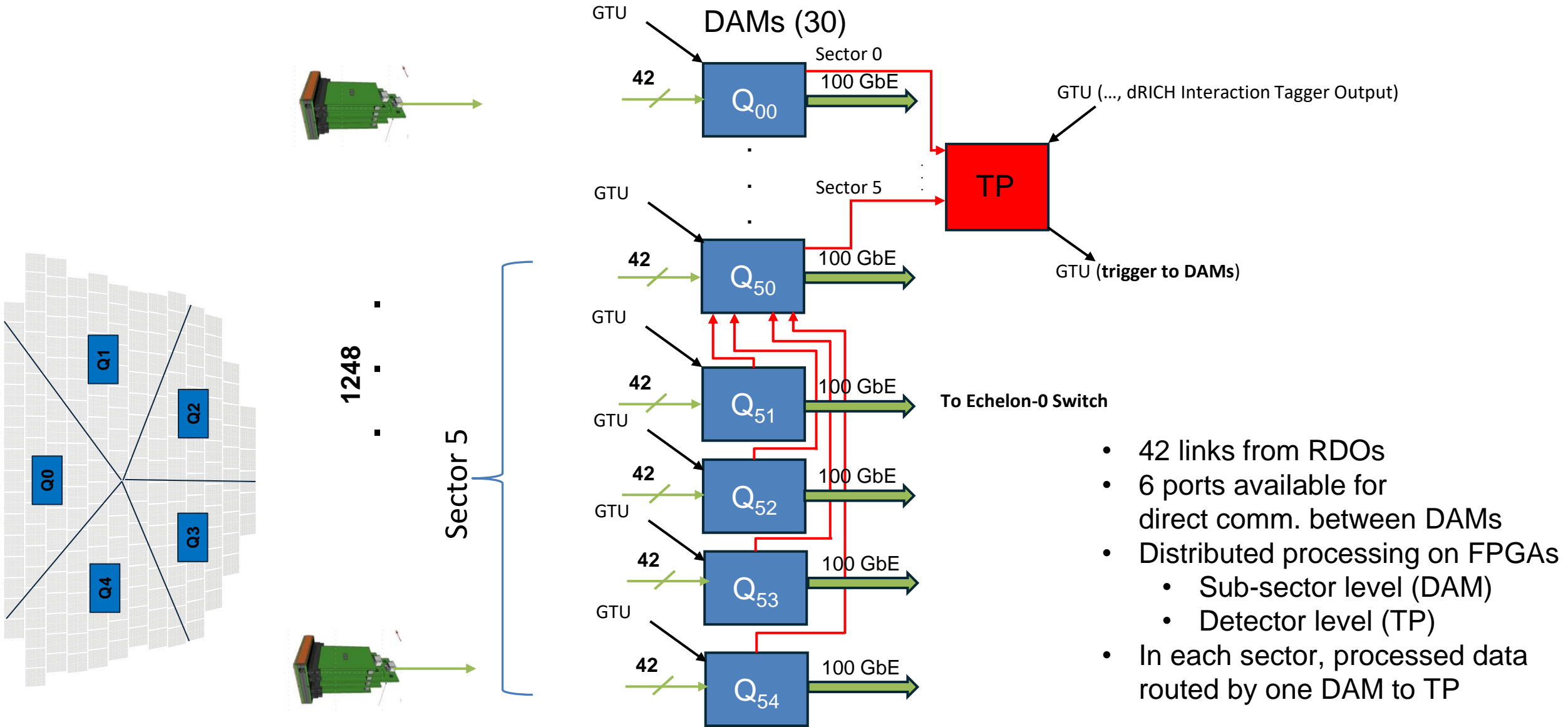**Phys Signal+Phys Background+Noise**



**Noise Only**

# dRICH Data Reduction Stage on FPGA

- Online «Noise only» classifier driving a dRICH local trigger.
  - Study of Inference Models
    - Restricting our study to inference models that can be deployed on FPGA with reasonable effort (using a High-Level Synthesis workflow)
      - MLP, CNN, GNN Models (HLS4ML)
    - Inference throughput (98.5 MHz) is the main challenge.
    - HDL optimized implementation is an option.
    - Not necessarily ML-based…
- Deployment on multiple Felix DAMs and on an additional FPGA (TP – Trigger Processor) directly interconnected with the APE communication IP.
- Possibly integrate with the **dRICH Interaction Tagger** to boost performance.

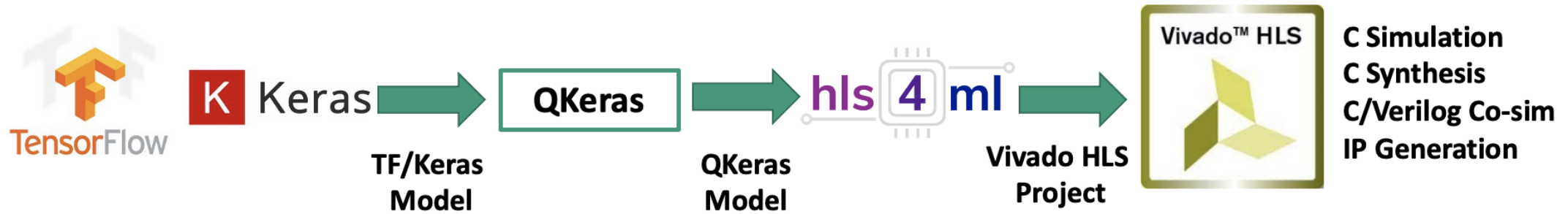# dRICH DAQ

# dRICH DAQ & Data Reduction



DAMs (30)

GTU

42 → Q_00 → Sector 0, 100 GbE

GTU

42 → Q_50 → Sector 5, 100 GbE

1248

Sector 5

GTU

42 → Q_51 → 100 GbE → To Echelon-0 Switch

GTU

42 → Q_52 → 100 GbE

GTU

42 → Q_53 → 100 GbE

GTU

42 → Q_54 → 100 GbE

TP

GTU (…, dRICH Interaction Tagger Output)

GTU (**trigger to DAMs**)

Q1 Q2 Q0 Q3 Q4

- 42 links from RDOs
- 6 ports available for direct comm. between DAMs
- Distributed processing on FPGAs
  - Sub-sector level (DAM)
  - Detector level (TP)
- In each sector, processed data routed by one DAM to TP

# Some Background Activities

- INFN APE Lab @ Roma1/2 **https://apegate.roma1.infn.it**
- Design and development of 4 generations of parallel computing architectures (mainly) dedicated to LQCD (1986-2010)
- Two recent research activities are relevant for this presentation:
  - APEIRON: a framework offering hardware and software support for the execution of real-time dataflow applications on a system composed by interconnected FPGAs. [https://doi.org/10.1051/epjconf/202429511002]
  - FPGA-RICH: online ring counting system based on FPGA for the RICH detector of the NA62 experiment at CERN. In publication.
- Other research activities of possible interest
 - APENet: a high-throughput network interface card based on FPGA used in hybrid, GPU-accelerated clusters with a 3D toroidal mesh topology. [http://doi.org/10.1088/1742-6596/898/8/082035]
 - NaNet: a family of FPGA-based PCIe Network Interface Cards (with GPUDirect/RDMA capability) for High Energy Physics to bridge the front-end electronics and the software trigger computing nodes. [https://doi.org/10.1088/1742-6596/1085/3/032022]

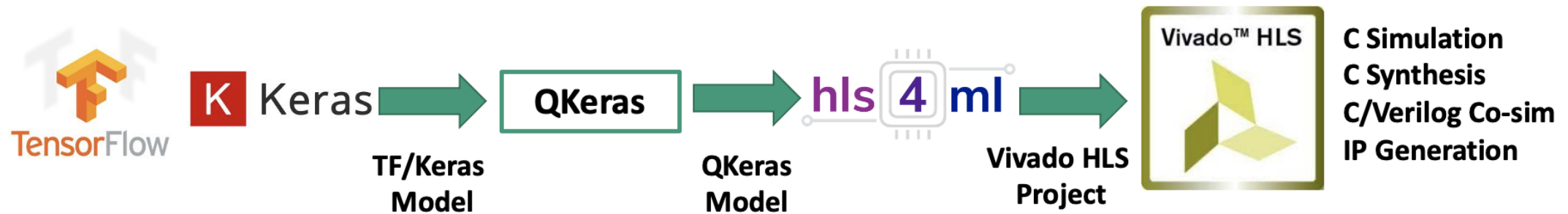# dRICH data reduction of FPGA: How?
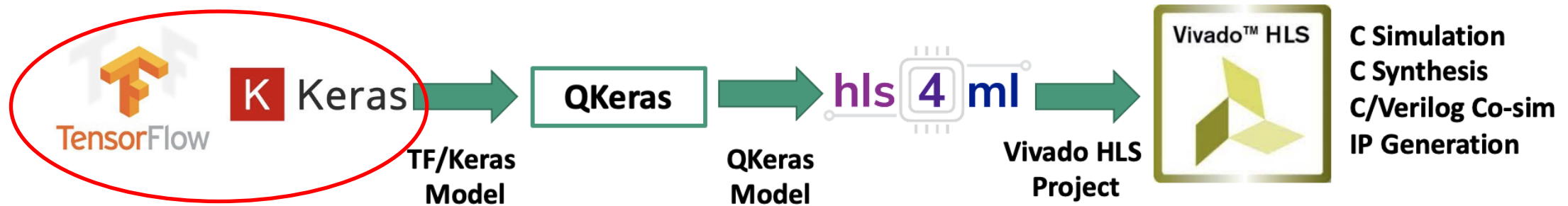# ➔ Design and Implementation Workflow



**Design targets** (accuracy, precision, recall, throughput, latency) and hardware constraints (mainly FPGA resource usage) must be taken into account and verified at any stage:

# dRICH Data reduction: How?
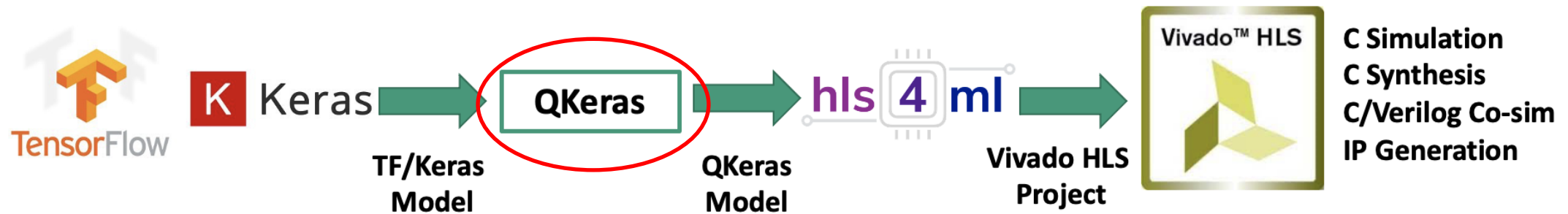# ➔ Design and Implementation Workflow



**Design targets** (accuracy, precision, recall, throughput, latency) and hardware constraints (mainly FPGA resource usage) must be taken into account and verified at any stage:

- **Generation strategy of training and validation data sets.**

# dRICH Data reduction: How?
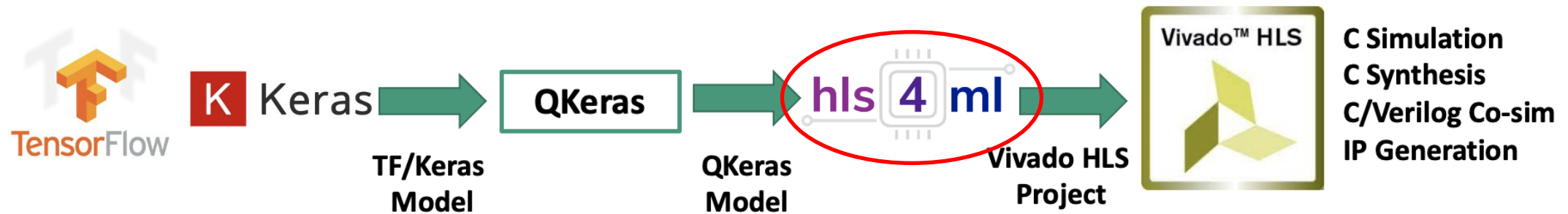# ➔ Design and Implementation Workflow



**Design targets** (accuracy, precision, recall, throughput, latency) and hardware constraints (mainly FPGA resource usage) must be taken into account and verified at any stage:

- **TensorFlow/Keras**
  ➔ NN architecture (number and kind of layers) and **representation of the input**
  ➔Training strategy (class balancing, batch sizes, optimizer choice, learning rate,…).

# dRICH Data reduction: How?
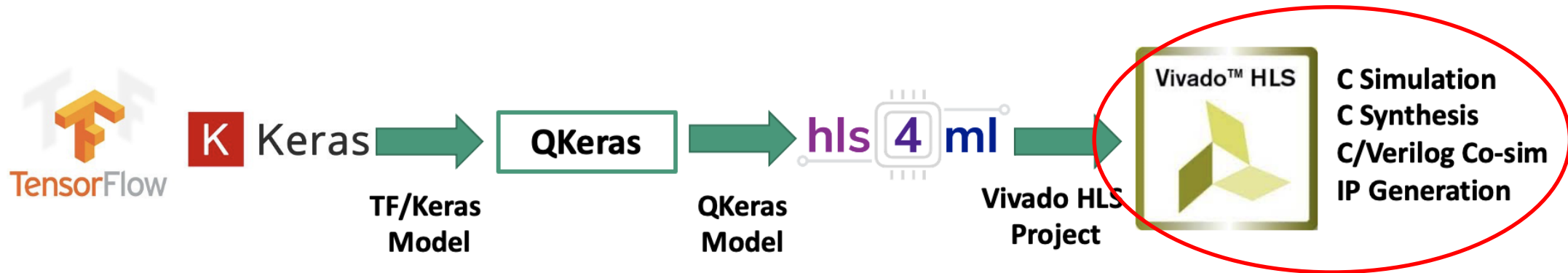# ➔ Design and Implementation Workflow



**Design targets** (accuracy, precision, recall, throughput, latency) and hardware constraints (mainly FPGA resource usage) must be taken into account and verified at any stage:

- **Qkeras** ➔ Search iteratively the minimal representation size in **bits** of input, weights, biases and activations.

# dRICH Data reduction: How?
# ➔ Design and Implementation Workflow



**Design targets** (accuracy, precision, recall, throughput, latency) and hardware constraints (mainly FPGA resource usage) must be taken into account and verified at any stage:

- **hls4ml** ➔ Tuning of REUSE FACTOR config param (low values ➔ low latency, high throughput, high resource usage), clock frequency.

# dRICH Data reduction: How?
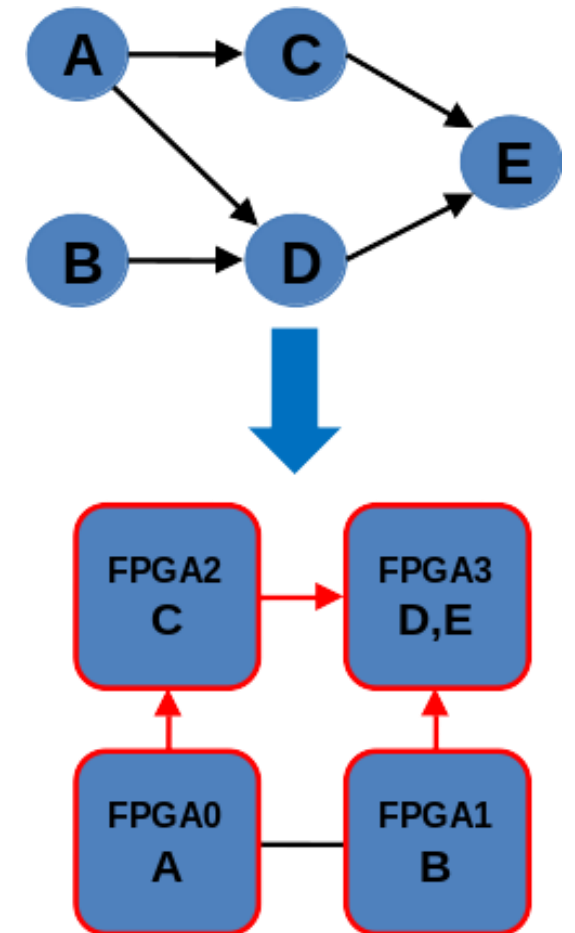# → Design and Implementation Workflow



**Design targets** (accuracy, precision, recall, throughput, latency) and hardware constraints (mainly FPGA resource usage) must be taken into account and verified at any stage:

- **Vivado HLS** → co-simulation for verification of performance (experimented very good agreement with QKeras Model)
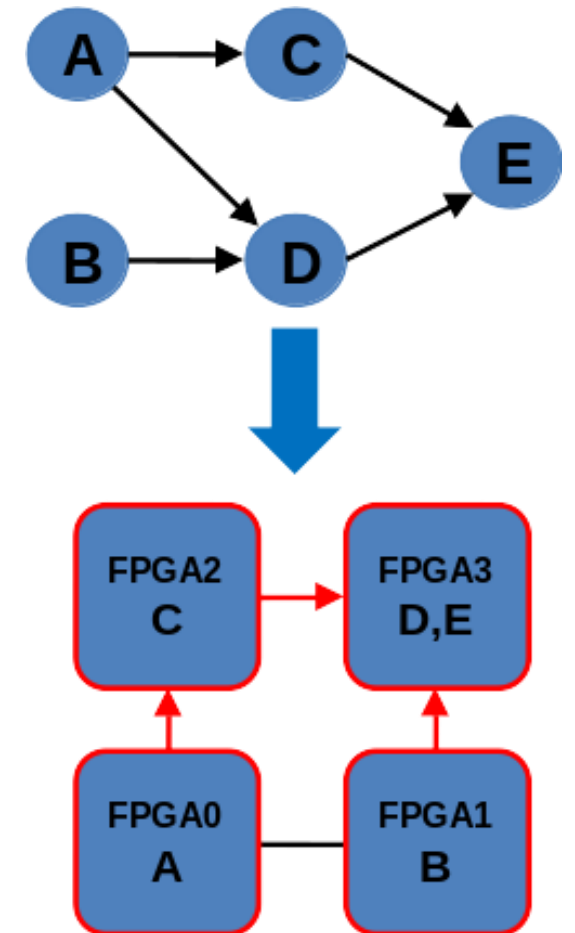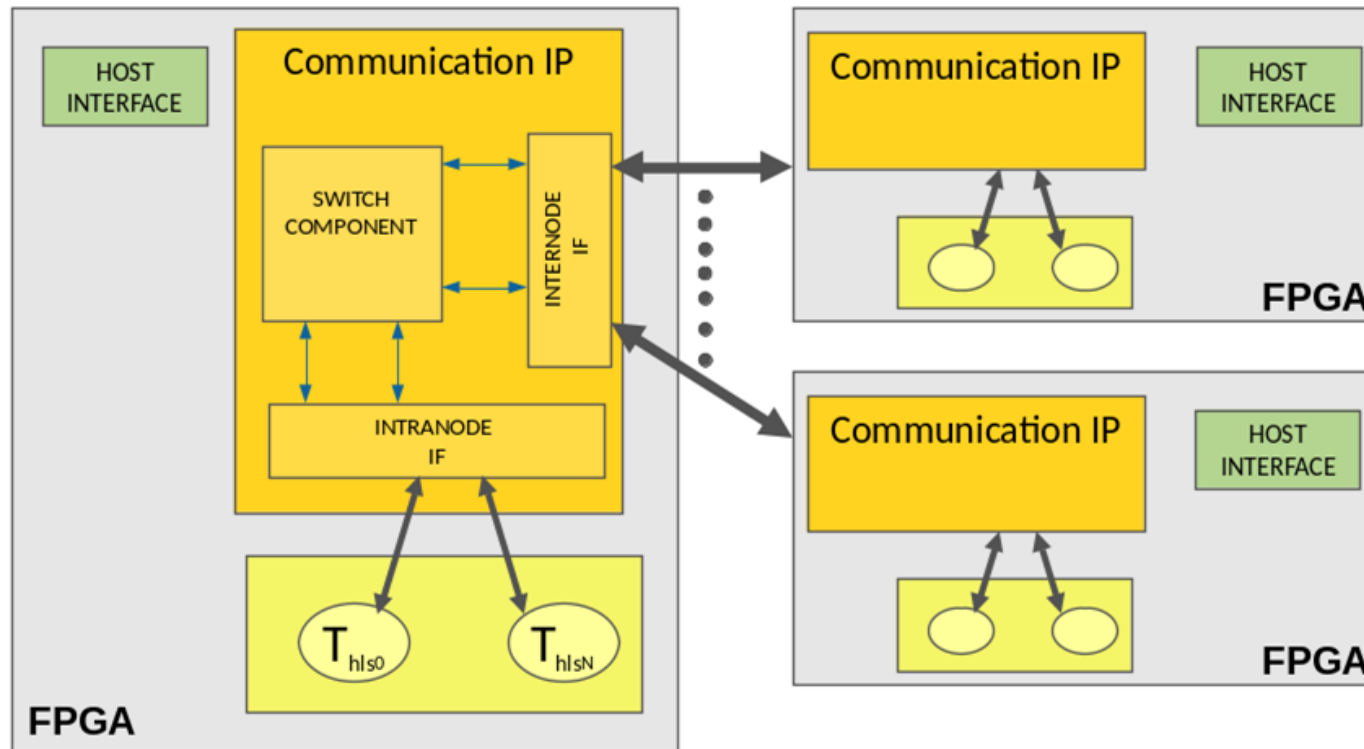
# APEIRON: an Overview

▪ Framework offering hardware and software support for the execution of real-time dataflow applications on a system composed by interconnected FPGAs .

- Map the dataflow graph of the application on the distributed FPGA system and offers runtime support for the execution.
- Allow users with no (or little) experience in hardware design tools, to develop their applications on such distributed FPGA-based platforms
  – Tasks are implemented in C++ using High Level Synthesis tools (Xilinx Vitis).
  – Lightweight C++ communication API
    • Non-blocking *send()*
    • Blocking *receive()*

▪ **APEIRON is based on Xilinx Vitis High Level Synthesis framework and on INFN Communication IP (APE Router)**
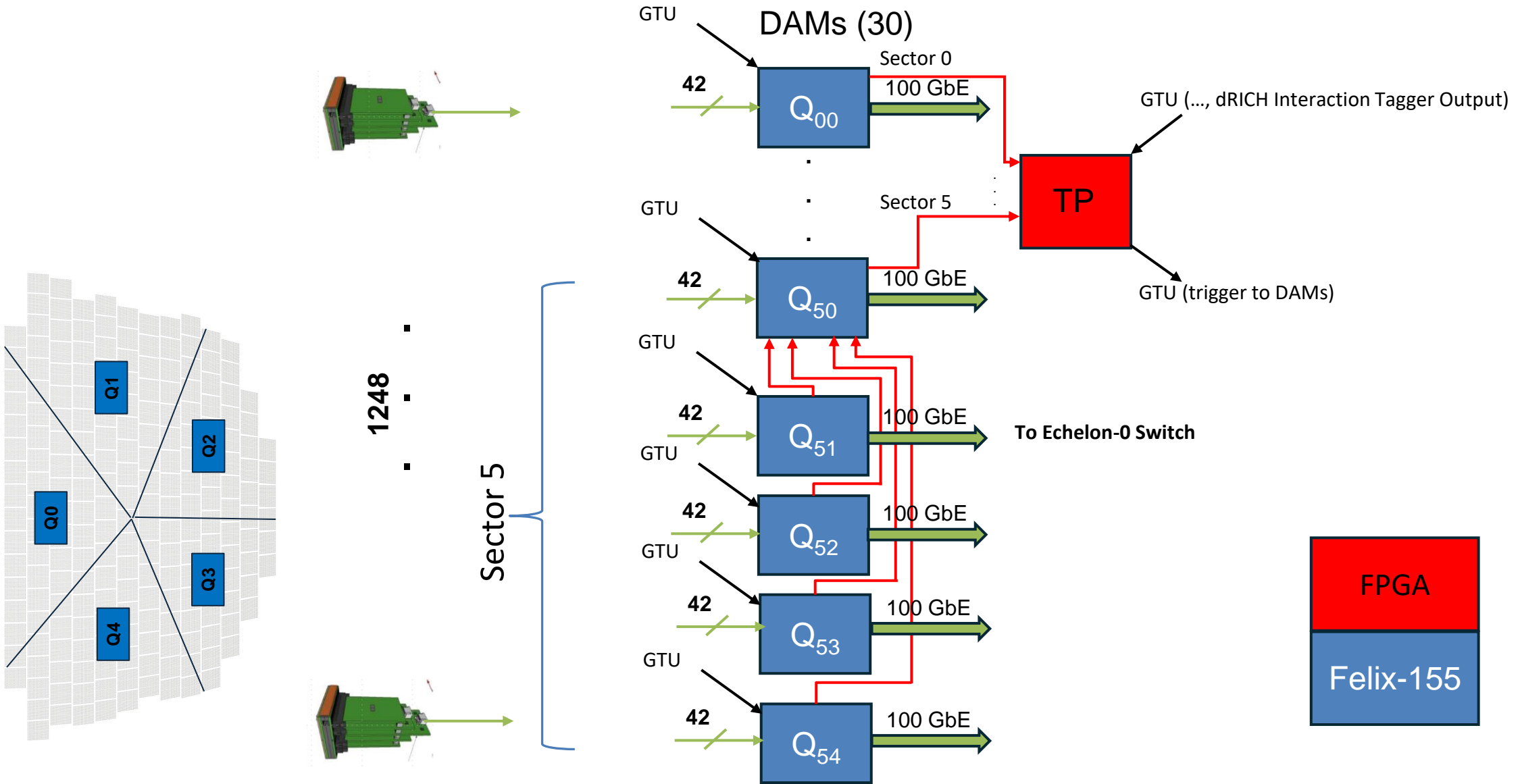
# APEIRON: INFN Communication IP

- INFN developed the IPs implementing a **direct network** that allows **low-latency** data transfer between processing **High Level Syntesis (C++) tasks** deployed on the same FPGA (intra-node communication) and on different FPGAs (inter-node communication).
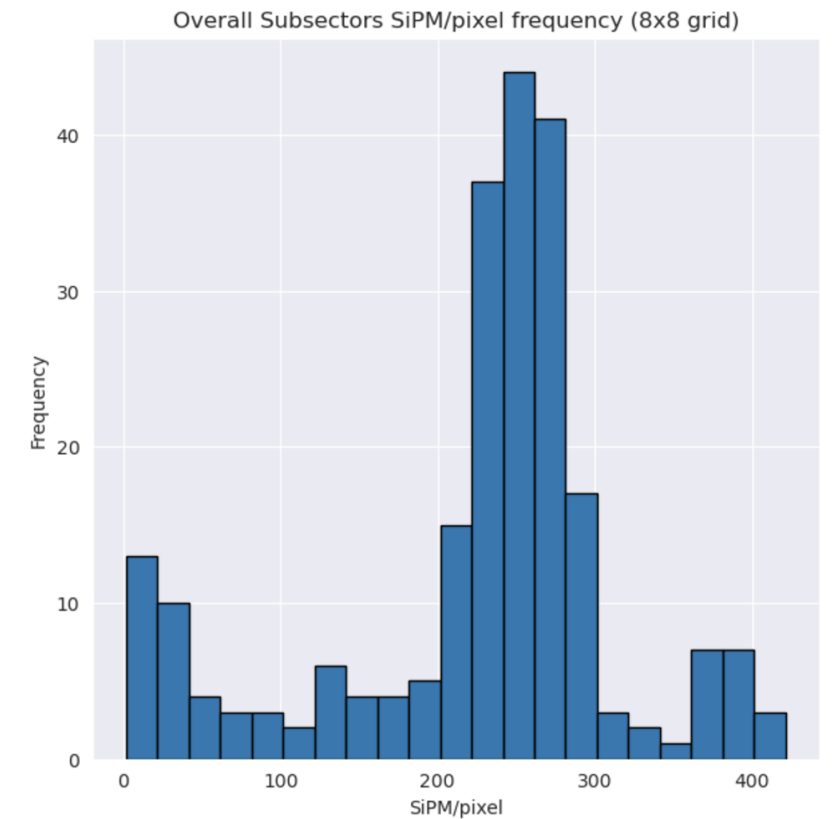
- Inter-node Latency < 1us for packet sizes up to 1kB (source and destination buffers in BRAM)
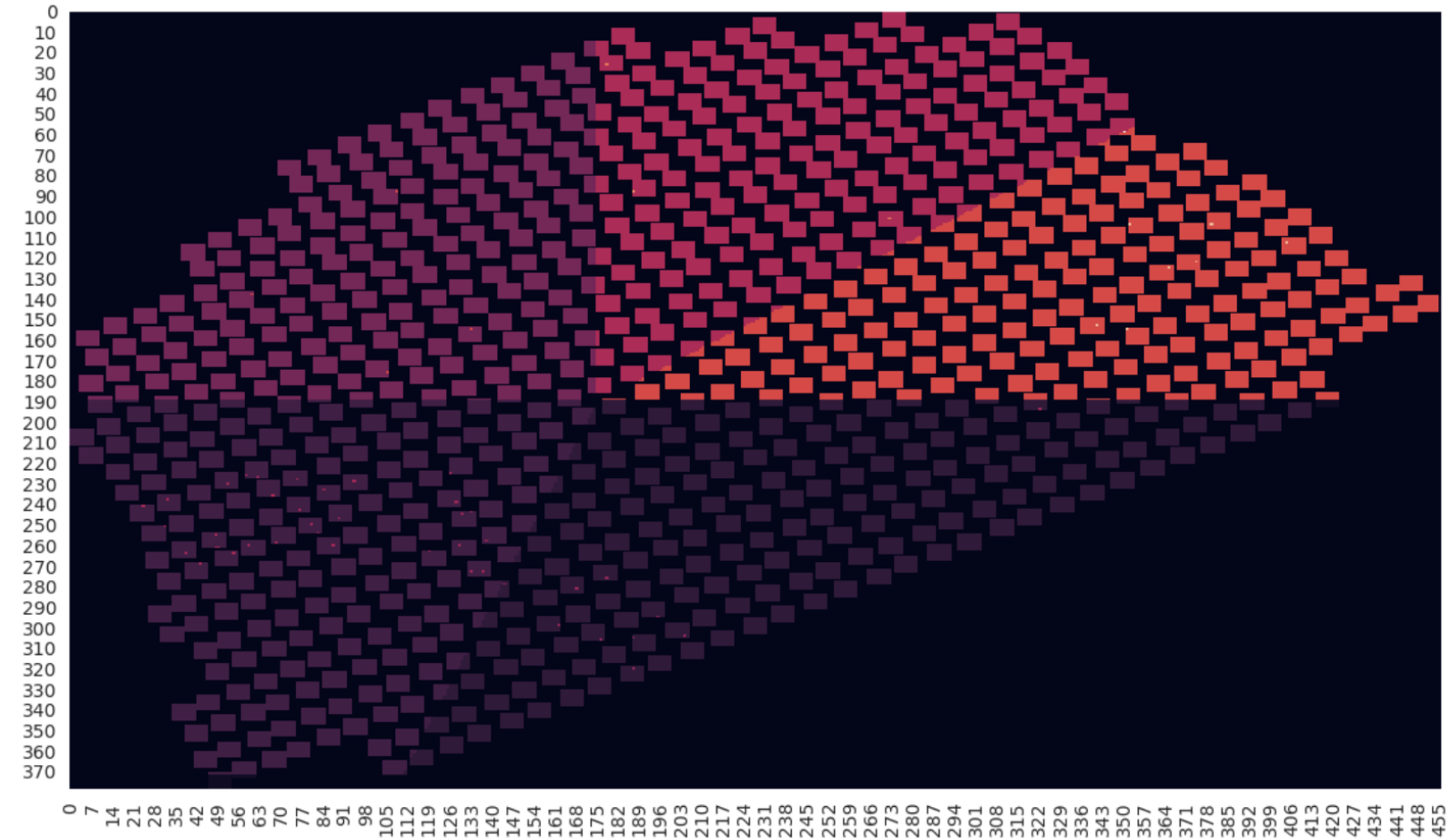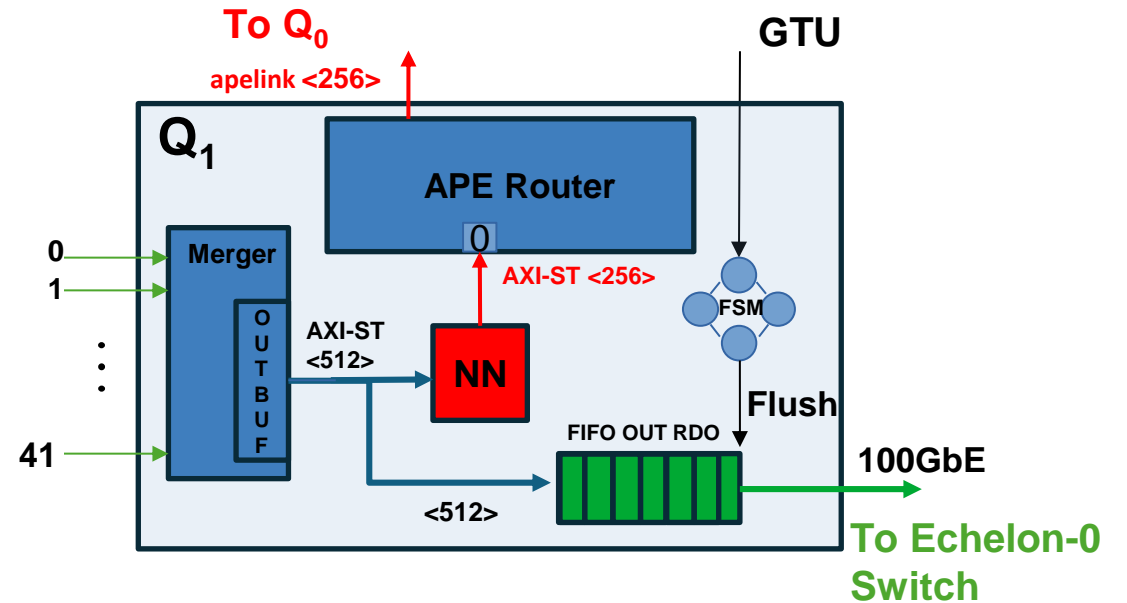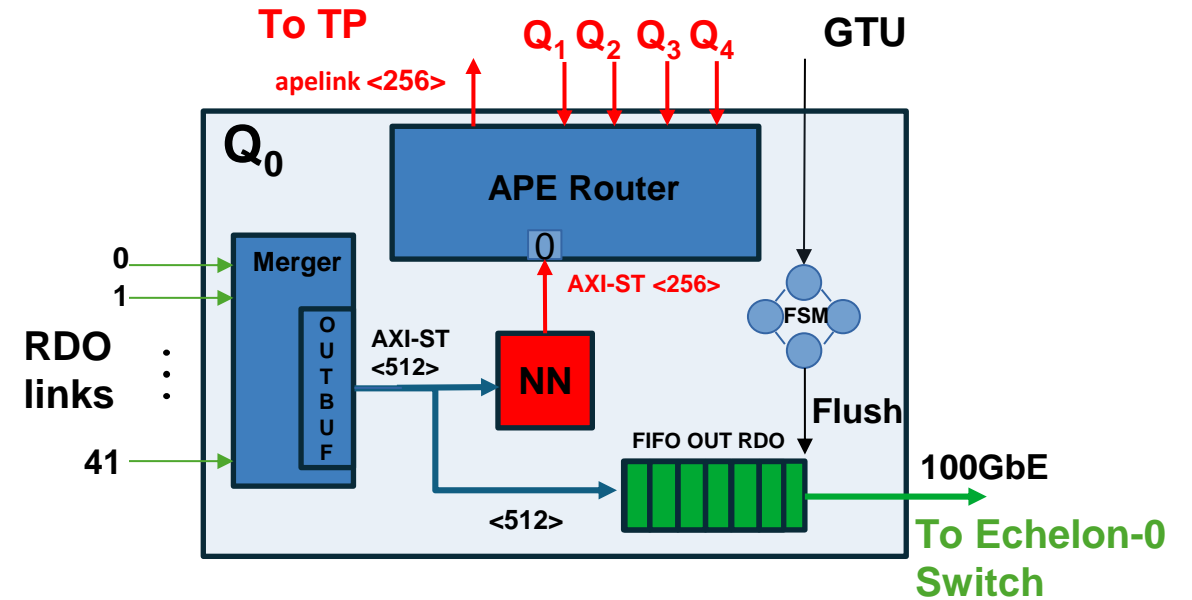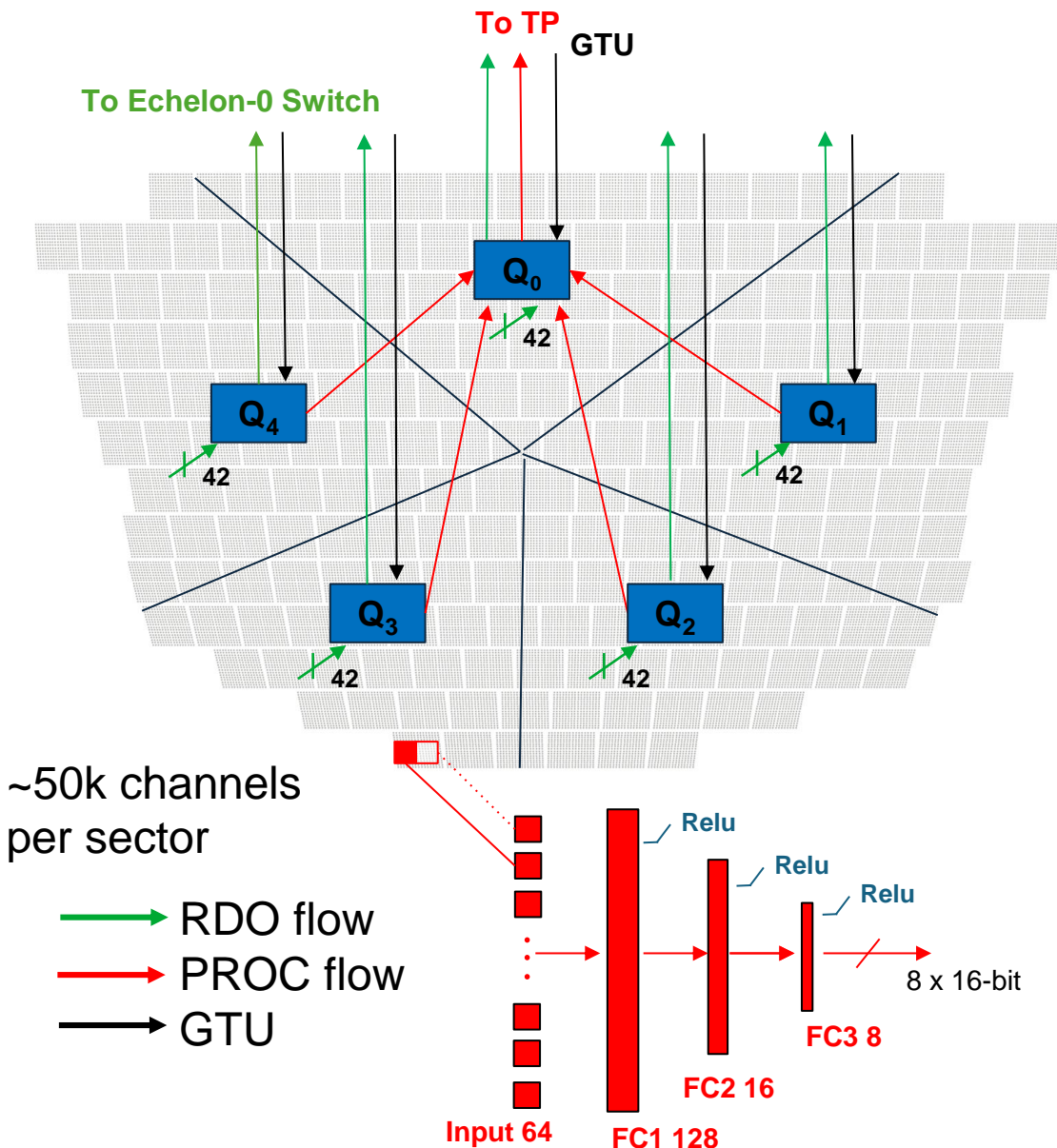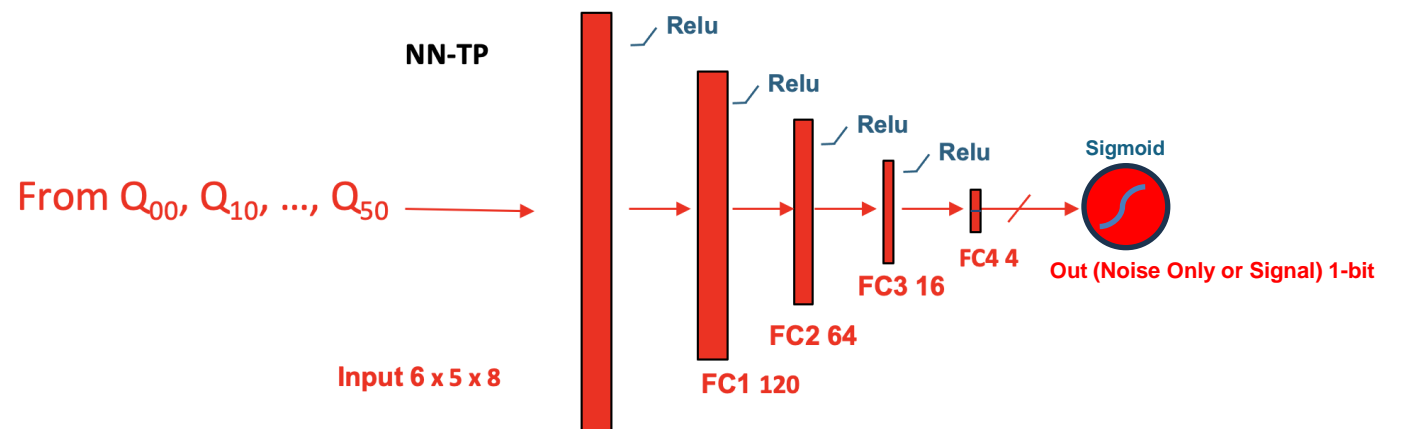
# dRICH DAQ & Data Reduction

# dRICH: Data reduction ➜ Subsectors

each subsector readout information discretized to a **8x8 grid ➜ 64 inputs to NN**



Overall Subsectors SiPM/pixel frequency (8x8 grid)

# dRICH Data Reduction on FPGA - Deployment

# dRICH Data Reduction on FPGA - Deployment

# dRICH: Data Reduction Dataset

**ePIC simulation pipeline:**

```
┌─────────────┐        ┌─────────────┐        ┌─────────────┐
│   HEPMC     │        │   FULL      │        │   RECO      │
│  MC events  │───────▶│   Geant4    │───────▶│  EICRecon   │
└─────────────┘        └─────────────┘        └─────────────┘
```

**PHYS signal**   e.g. DIS

**HEPMC_merger**

**PHYS bkg**   e - beamgas
             p - beamgas

Geant4 hits on detectors sensitive volume

In dRICH case, applies digitization, quantum efficiency, safety factor. **Possibly SiPM noise**

```
./SignalBackgroundMerger --signalFile dis.hepmc --signalFreq 0 \
    --bg1File hgas.hepmc --bg1Freq 31.9 \
    --bg2File egas.hepmc --bg2Freq 3177.25 \
```

root://dtn-eic.jlab.org//work/eic2/EPIC/FULL/24.05.0/epic_craterlake/BACKGROUNDS/MERGED/HEPMC_merger-1.0.2/10x100/RealisticSignalPerFrame/

**By default SiPM noise disabled** and not present on root://dtn-eic.jlab.org

Options:

➢ Start from Merged FULL root files available on server and enable noise at RECO stage using *drich-dev/recon.rb* with configs (but only ~ 7k events present on dtn-eic)

➢ **Run the entire simulation pipeline ourselves, starting from HEPMC files.**
  ○ Up to now we have produced 600k events to train and test our ML models.

# dRICH Data reduction:
# Input Data (Features Definition)

- **Gaussian** dark current SiPM **noise hits distribution,** obtained by modifying EICRecon source:
- avg = noiseRate*noiseTimeWindow
- sigma = 0.1*avg
- noiseTimeWindow = 2 ns



SiPMs Dark Current noise hits distribution

noiseRate = 200 KHz

# of hits

# dRICH Data reduction: Input Data (Features Definition)

➢ **Signal+Background+Noise**

➢ **Noise Only**

# dRICH Data Reduction: Input Data

➢ **Signal+Background+Noise**

➢ **Noise**

# dRICH Data Reduction: 8x8 Grid

➤ **Signal+Background+Noise**          ➤ **Noise**



**8x8 Grid ➔ 64 input NN neurons**

# dRICH Data Reduction: Tensorflow Model

- Coherently with the hardware design composition of the proposed system, we trained **30** (# of subsectors x #number of sectors) **concatenated MLP networks** into a single MLP model to be deployed on 30 DAM FPGAs + 1 TP FPGA



**DAM NN**

**Trigger Processor NN**

«Distributed MLP Model»

# Distributed MLP Tensorflow Model



Each MLP DAM output (**embedding**) is concatenated to the others to feed the final stage of the MLP (deployed on TP)

# dRICH Data reduction: model training & validation

o We trained the 30 MLP DAM models concatenated to the single MLP TP model by using 100k Signal+Background+Noise and 100k Noise Only event

o **200k balanced dataset** (90% training set, 10% validation set) for any of the considered values of noiseRate (100 KHz, 200 KHz, 400 KHz)

o We minimize a typical Binary CrossEntropy loss function in 1000 epochs, **backpropagating** the result to all the DAM input models ➜ in this way, trained 30 MLP DAM models result are **uncorrelated**

o Training and validation has been repeated after quantization

# Model training & validation: Loss

# Model training & validation: Accuracy

# Model performance @ noiseRate = 100 KHz



Confusion Matrix

- Accuracy = (TP+TN)/(TP+TN+FP+FN) = 0.997
- Precision = TP/(TP+FP) = 0.994
- Recall = TP/(TP+FN) = 1.000

# Model performance @ noiseRate = 200 KHz



Confusion Matrix

- Accuracy = (TP+TN)/(TP+TN+FP+FN) = 0.997
- Precision = TP/(TP+FP) = 0.994
- Recall = TP/(TP+FN) = 1.000

# Model performance @ noiseRate = 400 KHz



Confusion Matrix

- **Accuracy = (TP+TN)/(TP+TN+FP+FN) = 0.986**
- **Precision = TP/(TP+FP) = 0.974**
- **Recall = TP/(TP+FN) = 1.000**

# Quant. model performance @ noiseRate = 100 KHz



Confusion Matrix

- Accuracy =
  (TP+TN)/(TP+TN+FP+FN) = 0.990
- Precision = TP/(TP+FP) = 0.989
- Recall = TP/(TP+FN) = 0.991

**Model Quantization**
- **Inputs, Activations: fixed point<16,6>**
- **Weights, Biases: fixed point<8,1>**

# Quant. Model performance @ noiseRate = 200 KHz

Confusion Matrix



- Accuracy =
  (TP+TN)/(TP+TN+FP+FN)  = 0.985
- Precision = TP/(TP+FP) = 0.981
- Recall = TP/(TP+FN)     = 0.990

**Model Quantization**
- **Inputs, Activations: fixed point<16,6>**
- **Weights, Biases: fixed point<8,1>**

# Quant. model performance @ noiseRate = 400 KHz



Confusion Matrix

- Accuracy = $(TP+TN)/(TP+TN+FP+FN)$ = 0.903
- Precision = $TP/(TP+FP)$ = 0.909
- Recall = $TP/(TP+FN)$ = 0.895

**Model Quantization**
- **Inputs, Activations: fixed point<16,6>**
- **Weights, Biases: fixed point<8,1>**

# Summary of Distributed MLP Performance



Tensorflow Model Performance

Quantized Model Performance

# dRICH Data Reduction: HLS4ML ➜ HW Synthesis for 8x8 Grid DAM NN

➜ To correctly synthetize the model at 200 MHz of operational clock, we used a **REUSE FACTOR = 1**, obtaining an instantiation interval **II = 5 clock cycles**

➜ **Throughput = 40MHz (< 100 MHz)**

```
+ Timing:
    * Summary:
    +--------+-----------+-----------+------------+
    | Clock  |  Target   | Estimated | Uncertainty|
    +--------+-----------+-----------+------------+
    |ap_clk  |  5.00 ns  |  4.374 ns |   0.62 ns  |
    +--------+-----------+-----------+------------+

+ Latency:
    * Summary:
    +-----------------+---------------------+-------------+----------+
    |  Latency (cycles) |   Latency (absolute) |  Interval  | Pipeline |
    |  min  |   max   |    min    |    max    | min | max  |   Type   |
    +-----------------+---------------------+-------------+----------+
    |    14 |      14 | 70.000 ns | 70.000 ns |   5 |   5  | dataflow |
    +-----------------+---------------------+-------------+----------+
```

# dRICH Data Reduction: HLS4ML ➜ HW Synthesis for 8x8 Grid DAM NN

➜ The possible **overhead** in the full II pipepline **introduced by the communication between DAMs and TP** will be considered in further developments

➜ We synthetized the model at 200 MHz of operational clock, setting a **REUSE FACTOR = 1** and obtaining an instantiation interval **II = 5 clock cycles**

➜ **Throughput = 40MHz (< 100 MHz)**

**STILL LOW, BUT PROMISING! (can be improved via modifying part of HLS4ML code)**

```
+ Timing:
    * Summary:
    +---------+---------+-----------+-----------+
    |  Clock  |  Target | Estimated| Uncertainty|
    +---------+---------+-----------+-----------+
    |ap_clk   | 5.00 ns | 4.374 ns |   0.62 ns  |
    +---------+---------+-----------+-----------+

+ Latency:
    * Summary:
    +---------+---------+-----------+-----------+-----+-----+----------+
    | Latency (cycles) |  Latency (absolute) | Interval | Pipeline |
    |  min    |   max   |    min    |    max    | min | max |   Type   |
    +---------+---------+-----------+-----------+-----+-----+----------+
    |      14 |      14 | 70.000 ns | 70.000 ns |   5 |   5 | dataflow |
    +---------+---------+-----------+-----------+-----+-----+----------+
```

# Conclusions

o We sketched a data reduction system design to be deoployed on DAM's FPGAs as a risk-mitigation action to the possible problem of an excessive data bandwidth requirement from the dRICH to Echelon-0 due to SiPMs DCR.

o We showed results of the initial activities we made to proof the design concept.

o The design is based on a **distributed MLP NN** model, that can reach **near-optimal performance (using simulated data), and promising performance in terms of throughput at least in the first part of the distributed pipeline.**
**<u>These results need to be confirmed with a more realistic noise model.</u>**

o Next steps:
   o Deploy the distributed NN on two FPGAs already available in our lab (Xilinx Alveo U200) representing a DAM and the TP, integrating the communication in the pipeline and assessing its impact on pipeline throughput (and latency).
   o In addition different NN models (CNNs, GNN,...) and data reduction tasks/ideas (Cherenkov ring detection...) can be explored
   o Acquaint ourselves with the FELIX board HW and FW (we received a FLX-182 on loan from JLab) to start devising the integration of our design in its FW.
   o A initial «parasitic mode» deployment would allow the tuning and assessment of performace of the system, with periodic re-training of the NN with real data.

# Backup Slides

# APEIRON: the Node



APEIRON node in a 3D Torus network topology

- **Host Interface IP**: Interface the FPGA logic with the host through the system bus.
  – Xilinx XDMA PCIe Gen3
- **Routing IP**: Routing of intra-node and inter-node messages between processing tasks on FPGA.
- **Network IP**: Network channels and Application-dependent I/O
  – APElink 40 Gbps
  – UDP/IP over 10 GbE
- **Processing Tasks**: user defined processing tasks (Xilinx Vitis HLS Kernels)

# APEIRON: Communication Latency



**Test modes**
- Local-loop (red arrow)
- Local-trip (green arrows)
- Round-trip (blue arrows)

**Test Configuration**
- IP logic clock @ 200 MHz
- 4 intranode ports
- 2 internode ports
- 256-bit  datapath width
- 4 lanes inter-node channels

**Inter-node LATENCY (orange line) < 1us  for packet sizes up to 1kB (source and destination buffers in BRAM)**

# FELIX Hardware Development at BNL



FLX-711 @2017

FLX-712 @2018

FLX Demo @2019

FLX-181 @2021

FLX-182 @2022

FLX-182B @2023

FLX-155 @2024

PCIe Gen3 x16
48x 10 Gb/s
DDR3

PCIe Gen3 x16
48x 10 Gb/s

PCIe Gen3 x16
25 Gb/s with various optical modules

PCIe Gen4 x16
FMC+ 16x 25 Gb/s
12x 16 Gb/s
SoC / DDR4

PCIe Gen4 x16
28x 10/25 Gb/s
100 GbE
SoC / DDR4

PCIe Gen5 x16
48x 10/25 Gb/s
100/400 GbE
SoC / DDR4

# FLX-182B Hardware



Assembled FLX-182B

- FPGA: Xilinx Versal Prime XCVM1802
- PCIe Gen4 x16, 256 GT/s
- 24 FireFly links with 3 possible configurations
  - 24 links up to 25 Gb/s
  - 24 links up to 10 Gb/s (CERN-B FireFly)
  - 12 links up to 25 Gb/s + 12 links up to 10 Gb/s
- 4 FireFly links with 2 possible configurations with 14 or 25 Gb/s FireFly TRx
  - LTI interface
  - 100 GbE
- Built-in self test, online configuration and monitoring
- White Rabbit
- DDR4 Mini-UDIMM
- GbE/SD3.0/PetaLinux

Brookhaven
National Laboratory

# FLX-155 Hardware


Assembled FLX-155


3D model of FLX-155

- AMD/Xilinx Versal Premium FPGA: XCVP1552-2MSEVSVA3340
- 2 x PCIe Gen5 x8 512 GT/s
- 56 FireFly optical links
  - Compatible with various options
  - Default configuration for ATLAS
    - 48 data links up to 25 Gb/s
    - 4 links for LTI
  - Optional 4 links for 100 GbE
- Electrical IOs
- Built-in self test, online configuration and monitoring
- 1 16GB DDR4 Mini-UDIMM
- USB-JTAG/USB-UART
- GbE/SD3.0/PetaLinux
- Optional White Rabbit

| | VP1002 | VP1052 | VP1102 | VP1202 | VP1402 | VP1502 | VP2502 | VP1552 | VP1702 | VP1802 | VP2802 | VP1902 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| System Logic Cells | 833,000 | 1,185,800 | 1,574,720 | 1,969,240 | 2,233,280 | 3,763,480 | 3,737,720 | 3,836,840 | 5,557,720 | 7,351,960 | 7,326,200 | 18,506,880 |
| CLB Flip-Flops | 761,600 | 1,084,160 | 1,439,744 | 1,800,448 | 2,041,856 | 3,440,896 | 3,417,344 | 3,507,968 | 5,081,344 | 6,721,792 | 6,698,240 | 16,920,576 |
| LUTs | 380,800 | 542,080 | 719,872 | 900,224 | 1,020,928 | 1,720,448 | 1,708,672 | 1,753,984 | 2,540,672 | 3,360,896 | 3,349,120 | 8,460,288 |
| Distributed RAM (Mb) | 12 | 17 | 22 | 27 | 31 | 53 | 52 | 54 | 78 | 103 | 102 | 258 |
| Block RAM Blocks | 535 | 751 | 1,405 | 1,341 | 1,981 | 2,541 | 2,541 | 2,541 | 3,741 | 4,941 | 4,941 | 6,808 |
| Block RAM (Mb) | 19 | 26 | 49 | 47 | 70 | 89 | 89 | 89 | 132 | 174 | 174 | 239 |
| UltraRAM Blocks | 345 | 489 | 453 | 677 | 645 | 1,301 | 1,301 | 1,301 | 1,925 | 2,549 | 2,549 | 2,200 |
| UltraRAM (Mb) | 97 | 138 | 127 | 190 | 181 | 366 | 366 | 366 | 541 | 717 | 717 | 619 |
| Multiport RAM (Mb) | 80 | 80 | – | – | – | – | – | – | – | – | – | – |
| DSP Engines | 1,140 | 1,572 | 1,904 | 3,984 | 2,672 | 7,440 | 7,392 | 7,392 | 10,896 | 14,352 | 14,304 | 6,864 |
| AI Engines (AIE) | – | – | – | – | – | – | 472 | – | – | – | 472 | – |
| AIE Data Memory (Mb) | – | – | – | – | – | – | 118 | – | – | – | 118 | – |
| APU | Dual-core Arm Cortex-A72; 48 KB/32 KB L1 Cache w/ parity & ECC; 1 MB L2 Cache w/ ECC | | | | | | | | | | | |
| RPU | Dual-core Arm Cortex-R5F; 32 KB/32 KB L1 Cache; TCM w/ECC | | | | | | | | | | | |
| Memory | 256 KB On-Chip Memory w/ECC | | | | | | | | | | | |
| Connectivity | Ethernet (x2); UART (x2); CAN-FD (x2); USB 2.0 (x1); SPI (x2); I2C (x2) | | | | | | | | | | | |
| NoC to PL Master / Slave Ports | 22 | 22 | 30 | 28 | 42 | 52 | 52 | 52 | 76 | 100 | 100 | 192 |
| DDR Bus Width | 128 | 128 | 192 | 256 | 192 | 256 | 256 | 256 | 256 | 256 | 256 | 896 |
| DDR Memory Controllers (DDRMC) | 2 | 2 | 3 | 4 | 3 | 4 | 4 | 4 | 4 | 4 | 4 | 14 |
| PCIe w/DMA (CPM4) | 2 x Gen4x4 | 2 x Gen4x4 | – | – | – | – | – | – | – | – | – | – |
| PCIe w/DMA (CPM5) | – | – | – | 2 x Gen5x8 | – | 2 x Gen5x8 | 2 x Gen5x8 | 2 x Gen5x8 | 2 x Gen5x8 | 2 x Gen5x8 | 2 x Gen5x8 | – |
| PCIe (PL PCIE4) | 1 x Gen4x8 | 1 x Gen4x8 | – | – | – | – | – | – | – | – | – | – |
| PCIe (PL PCIE5) | – | – | 2 x Gen5x4 | 2 x Gen5x4 | 2 x Gen5x4 | 2 x Gen5x4 | 2 x Gen5x4 | 8 x Gen5x4 | 2 x Gen5x4 | 2 x Gen5x4 | 2 x Gen5x4 | 16 x Gen5x4 |
| 100G Multirate Ethernet MAC | 3 | 5 | 6 | 2 | 6 | 4 | 4 | 4 | 6 | 8 | 8 | 12 |
| 600G Ethernet MAC | 2 | 3 | 7 | 1 | 11 | 3 | 3 | 1 | 5 | 7 | 7 | 4 |
| 600G Interlaken | 1 | 2 | – | – | – | 1 | 1 | – | 2 | 3 | 3 | 0 |
| High-Speed Crypto Engines | 1 | 1 | 3 | 1 | 4 | 2 | 2 | 2 | 3 | 4 | 4 | 0 |
| GTY Transceivers[1] | 8 | 8 | – | – | – | – | – | – | – | – | – | – |
| GTYP Transceivers[1] | – | – | 8 | 28[3] | 8 | 28[3] | 28[3] | 68[3] | 28[3] | 28[3] | 28[3] | 128 |
| GTM Transceivers[1] 58Gb/s (112 Gb/s) | 24 (12) | 36 (18) | 64 (32) | 20 (10) | 96 (64)[2] | 60 (30) | 60 (30) | 20 (10) | 100 (50) | 140 (70) | 140 (70) | 32 (16) |