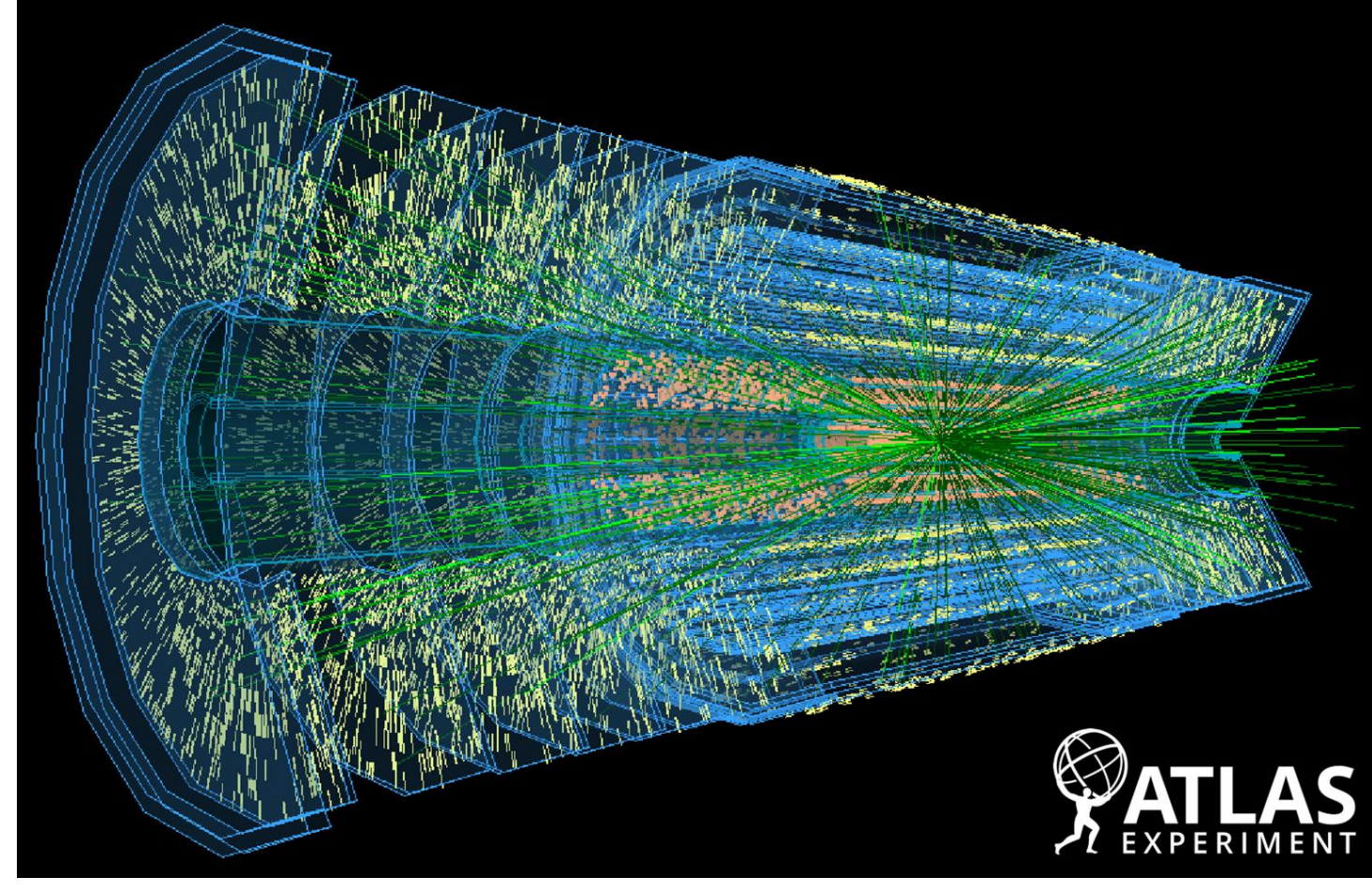# An Integrated FPGA Implementation of Complete GNN-Based Trajectory Reconstruction

Yun-Chen Yang*, Hao-Chun Liang*, Hsuan-Wei Yu*,

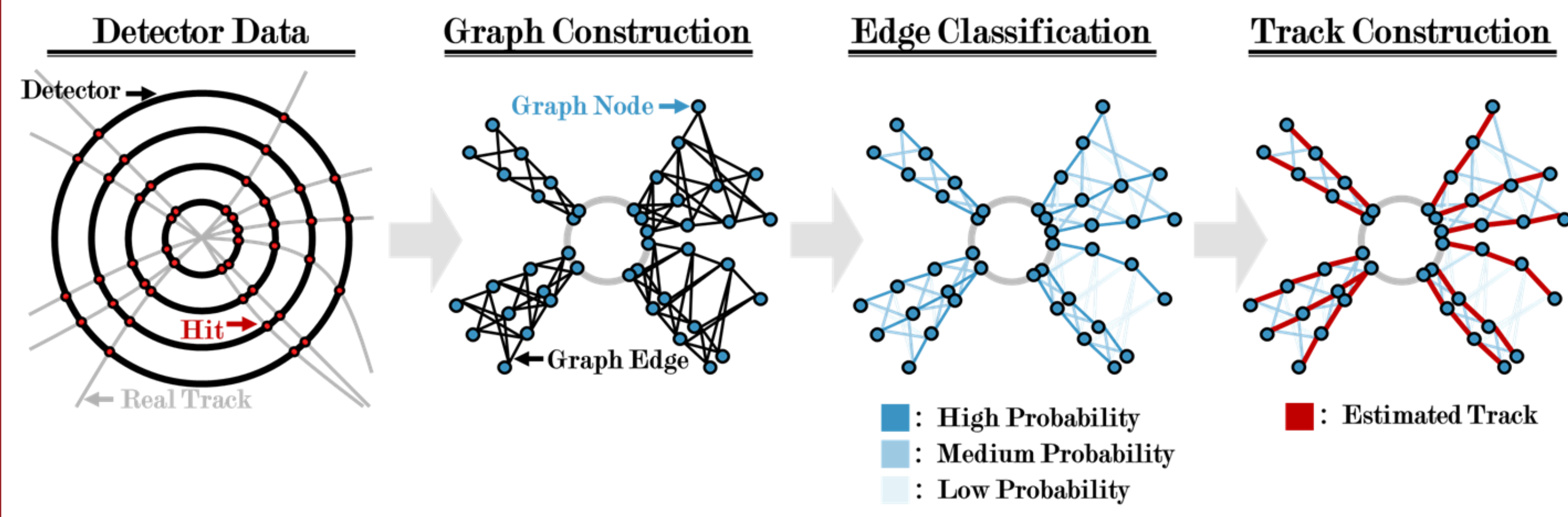Bo-Cheng Lai*, Shih-Chieh Hsu†, Mark Neubauer‡, Santosh Parajuli‡

National Yang Ming Chiao Tung University, Hsinchu, Taiwan*, University of Washington, USA† University of Illinois Urbana-Champaign, USA‡
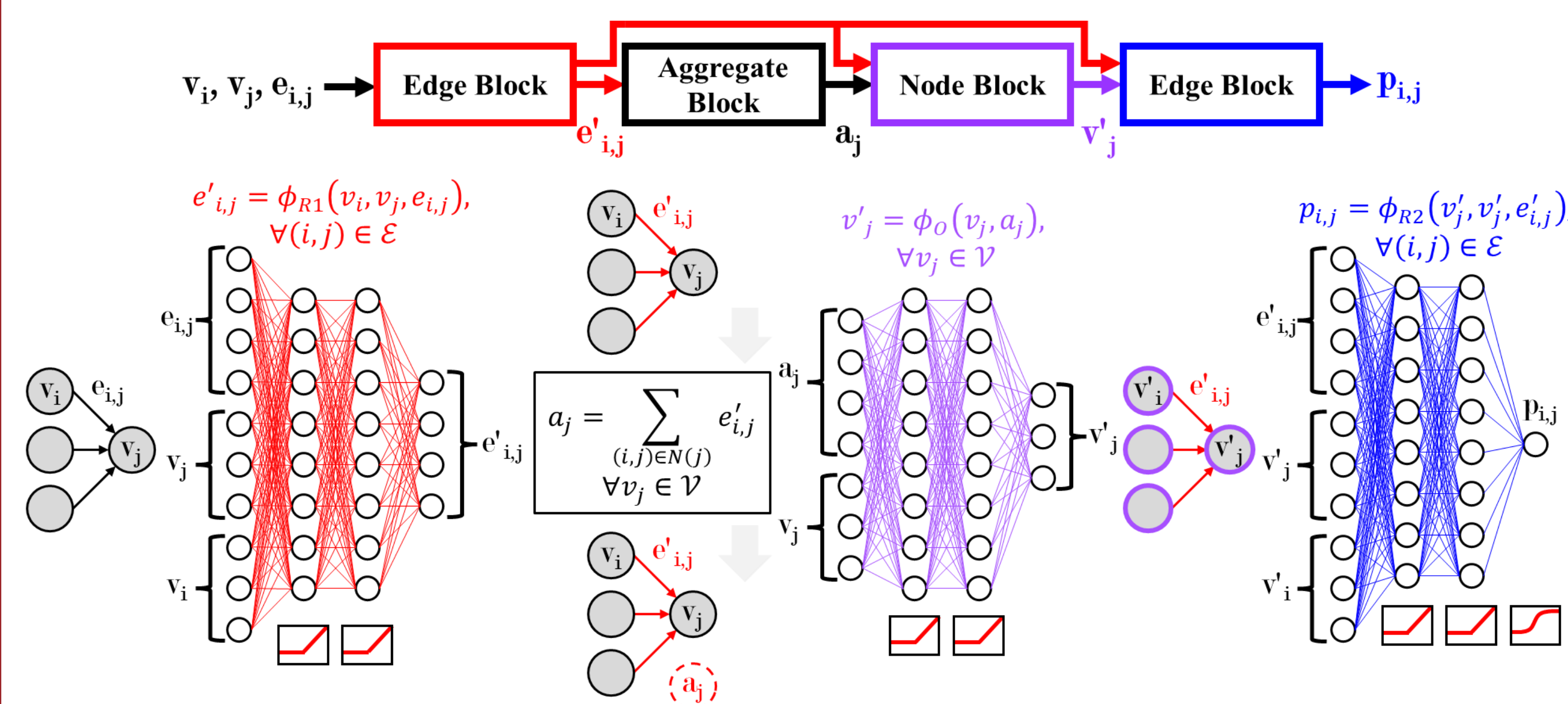
## Motivation & Challenge



- **Collaborate with CERN**
  - **Proton-Beam Collisions**
  - **Detector Produce Hits**
- **High-Volume Collision Hits**
  - **L1T System Filtering**
  - **via Track Reconstruction**
- **HL-LHC Requirement**
  - **Latency: 4 µs**
  - **Throughput: 2.22MHz**

- **Accuracy-Focused, Execution-Speed-Agnostic Approach** [7, 8]
  - **CPU-Based Limits Hindering Task-Specific Tuning**
  - **GPU-Based Inefficiency in Latency-Critical Scenarios**
    - **Missed Microsecond Targets By Milliseconds**
- **FPGA-Accelerated Framework**
  - **Scope Constrained to GNN Edge-Classification Stage** [9,10]
    - **Host-to-Device Data Transfers Latency**
    - **FPGA Resources Underutilization**
  - **Throughput-Driven Processing In Minor Graph Subregions** [11]
    - **Small Subgraph Lowers Accuracy**
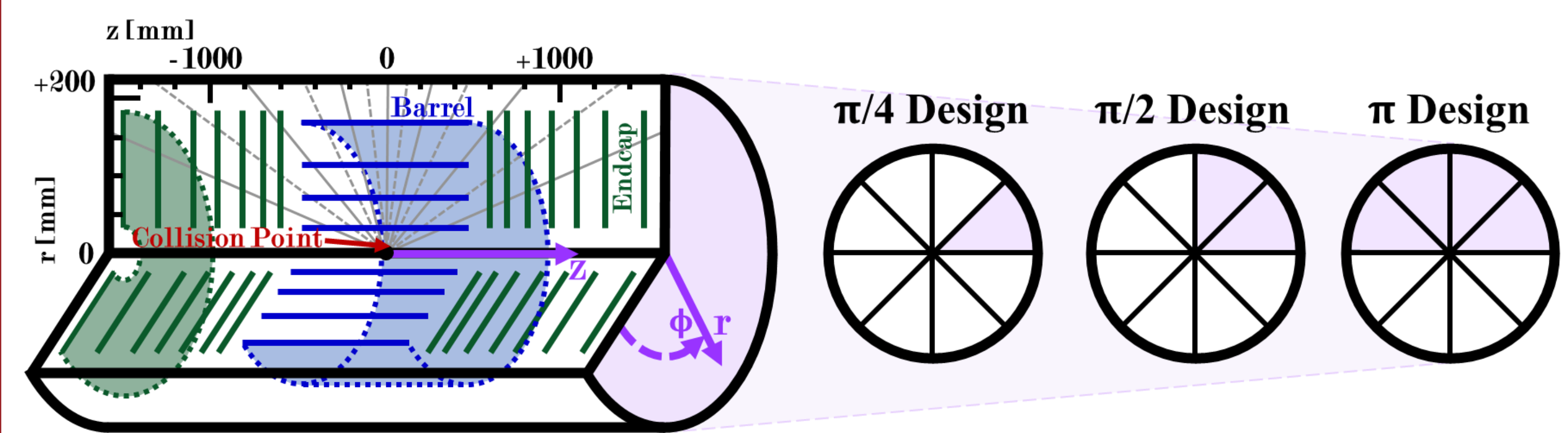
## Trajectory Reconstruction



- **Graph Construction – Map Hits to Nodes and Filter Directed Edges**
- **Edge Classification – Assign Probabilities to Edges**
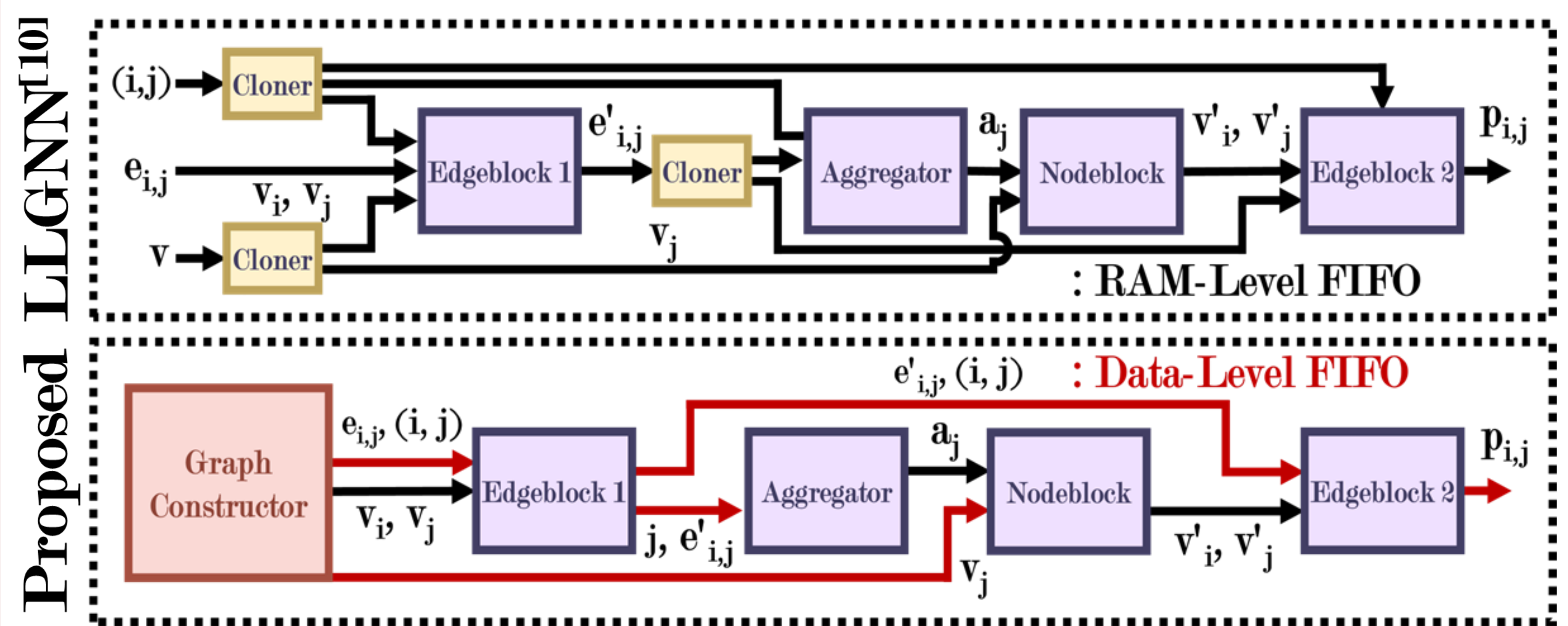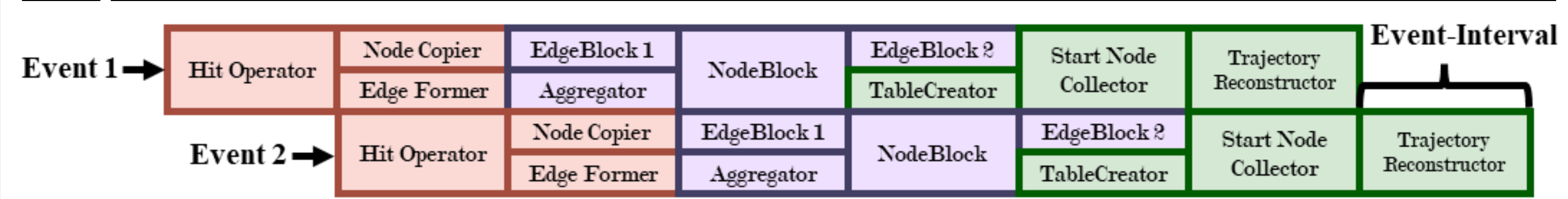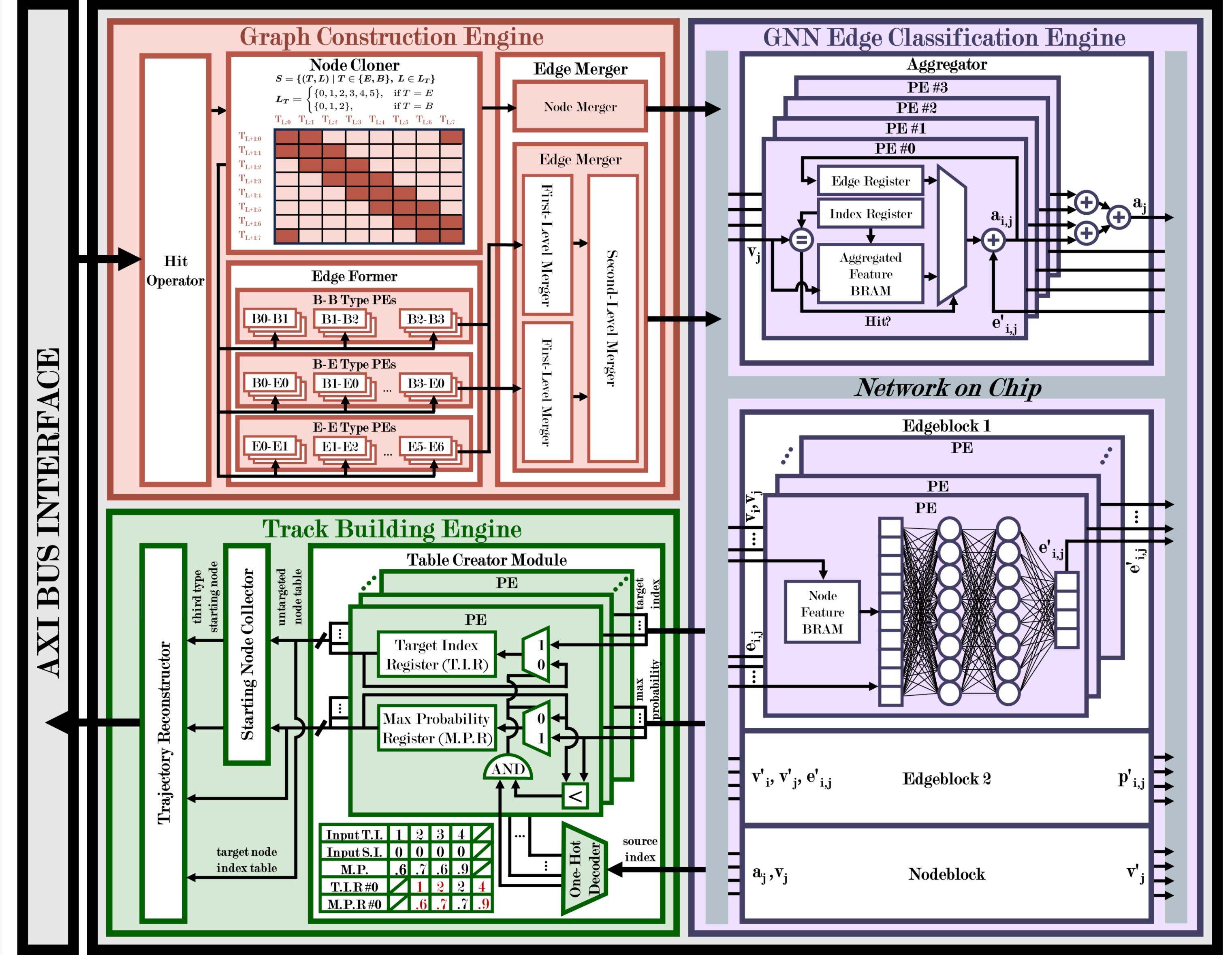- **Track Construction – Integrate Edges with High Probabilities**



$$e'_{i,j} = \phi_{R1}(v_i, v_j, e_{i,j}), \forall (i,j) \in \mathcal{E}$$

$$a_j = \sum_{(i,j) \in N(j)} e'_{i,j} \quad \forall v_j \in \mathcal{V}$$

$$v'_j = \phi_O(v_j, a_j), \forall v_j \in \mathcal{V}$$

$$p_{i,j} = \phi_{R2}(v'_i, v'_j, e'_{i,j}), \forall (i,j) \in \mathcal{E}$$

- **Interaction Network (IN) Framework** [6]
  - **Graph Neural Network for Object-Object Interaction Modeling**

## Configuration



- **Cross-Sectional View of Cylindrical Collider Detector Architecture**
  - **4× Cylindrical Barrels and 14× Planar Endcaps**
- **Hit Distribution across Segment**
  - **Longitudinal Segmentation along the Z Axis**
  - **Azimuthal Segmentation along the Φ-Axis into 2/4/8 Sectors**
    - **Adapted to Support Three Distinct Design Variants**
  - **Spatial-Multiplexed FPGA for Parallel Processing**

## Architecture



- **Edge Classification from Batch Processing** [10] **to Data Streaming**
  - **52.3% Latency Reduction from 2.86 µs to 1.365 µs**

## Algorithm & Conclusion

- **Geometry-Aware Edge Pruning in Graph Construction**
  - **Edge Count Reduction Achieving 37.5%–71.0% Pruning** [7]



$$L(T) \in \begin{cases} \{0,1,2,3,4,5\}, & \text{if } T = E \\ \{0,1,2\}, & \text{if } T = B \end{cases}$$

- **Probability-Based Sequential Building in Track Construction**
  - **LUT-Based Mapping from $B_0$ (Blue), $E_0$ (Brown), Other (Green)**
  - **Latency Reduction Ranging from 107× to 641×** [7]



- **First End-to-End FPGA of GNN-Based Trajectory Reconstruction**
  - **65024× Acceleration with Enhanced Accuracy Metrics** [7,9]
  - **Throughput 2.35 MHz with Latency 2.36 µs Meets L1T Criteria**