

2 **A demonstrator for a real-time AI-FPGA-based** 3 **triggering system for sPHENIX at RHIC**

4 **J. Kvapil,^{a,1} G. Borca-Tasciuc,^b H. Bossi,^c K. Chen,^d Y. Chen,^d Y. Corrales Morales,^c H.**
5 **Da Costa,^a C. Da Silva,^a C. Dean,^c J. Durham,^a S. Fu,^e C. Hao,^f P. Harris,^c O. Hen,^c H.**
6 **Jheng,^c Y. Lee,^c P. Li,^f X. Li,^a Y. Lin,^a M. X. Liu,^a A. Olvera,^e M. Purschke,^g M. Rigatti,^h**
7 **G. Roland,^c J. Schambach,ⁱ Z. Shi,^a N. Tran,^h N. Wuerfel,^j B. Xu,^f D. Yu,^k H. Zhang^f**

8 ^a*Los Alamos National Laboratory,*
9 *Bikini Atoll Rd, Los Alamos, NM 87545, United States*

10 ^b*Rensselaer Polytechnic Institute,*
11 *110 8th St, Troy, NY 12180, United States*

12 ^c*Massachusetts Institute of Technology,*
13 *77 Massachusetts Ave, Cambridge, MA 02139, United States*

14 ^d*Central China Normal University,*
15 *No.152, Luoyu Rd, Wuhan 430079, China*

16 ^e*University of North Texas,*
17 *1155 Union Cir, Denton, TX 76205, United States*

18 ^f*Georgia Institute of Technology,*
19 *225 North Ave, Atlanta, GA 30332, United States*

20 ^g*Brookhaven National Laboratory,*
21 *PO Box 5000 Upton, NY 11973, United States*

22 ^h*Fermilab,*
23 *PO Box 500. Batavia IL 60510, United States*

24 ⁱ*Oak Ridge National Laboratory,*
25 *P.O. Box 2008. Oak Ridge, TN 37831, United States*

26 ^j*University of Michigan,*
27 *500 S State St, Ann Arbor, MI 481091, United States*

28 ^k*New Jersey Institute of Technology,*
29 *323 Dr Martin Luther King Jr Blvd, Newark, NJ 07102, United States*

30 *E-mail: jakub.kvapil@lanl.gov; jakub.kvapil@cern.ch*

¹Corresponding author.

31 **ABSTRACT:** The RHIC interaction rate at sPHENIX will reach around 3 MHz in pp collisions and
32 requires the detector readout to reject events by a factor of over 200 to fit the DAQ bandwidth
33 of 15kHz. Some critical measurements, such as heavy flavor production in pp collisions, often
34 require the analysis of particles produced at low momentum. This prohibits adopting the traditional
35 approach, where data rates are reduced through triggering on rare high momentum probes. We
36 explore a new approach based on real-time AI technology, adopt an FPGA-based implementation
37 using a custom designed FELIX-712 board with the Xilinx Kintex Ultrascale FPGA, and deploy
38 the system in the detector readout electronics loop for real-time trigger decision.

39 **KEYWORDS:** Detector control systems (detector and experiment monitoring and slow-control sys-
40 tems, architecture, hardware, algorithms, databases); Trigger algorithms; Trigger concepts and
41 systems (hardware and software)

42 **ARXIV EPRINT:** [not yet](#)

43 1 Motivation

44 Realizing the science potential of modern nuclear physics (NP) experiments at colliders relies on
45 the collection and processing of very large datasets, with sustained data rates from future detectors
46 exceeding Terabits per second. Critical measurements in NP data often require the analysis of
47 complex final states of particles produced at low momentum. This prohibits adopting the traditional
48 approach used in high energy physics, where data rates are reduced through online event selection
49 (“triggering”) on rare high momentum probes that can be readily distinguished from the background.
50 On the other hand, archiving the full detector data stream exceeds current DAQ bandwidth limits
51 and would lead to cost-prohibitive offline storage and analysis computing requirements. We propose
52 to develop real-time AI technologies, implemented in the detector readout electronics loop, that
53 address these challenges for the next generation of NP experiments at RHIC and EIC. First, we
54 will deploy a demonstrator that is being developed under the current Fast-ML project for the pp
55 running in the sPHENIX experiment in 2024, and then generalize our approach for applications in
56 experiments at the EIC, using future generations of EIC detector technologies.

57 The triggered readout rate of sPHENIX is limited to 15 kHz due to the design of the calorimeter
58 readout system, which places limitations on the overall data volume and the expected collision rate.
59 In pp collisions, RHIC delivers collision rates of 3 MHz, limiting sPHENIX to collect less than 1%
60 of heavy-flavour (HF) events of the total pp (and p+Au) rate when using triggered readout. The
61 extended Streaming readout (SRO) of the tracking detectors can further improve the statistics up to
62 10% of the total luminosity. The goal of this project is to sample the remaining luminosity further
63 enhancing the collected data samples. The aim is to deploy a future system on the Electron-Ion
64 Collider (EIC) to identify the (non)interesting Deep-Inelastic-Scattering processed in the e+p/A
65 collisions.

66 2 sPHENIX detector

67 The sPHENIX detector [1] is located at the RHIC accelerator complex in BNL. The solenoid magnet
68 provides a magnetic field of 1.4 T and the detectors have a pseudorapidity coverage of $|\eta| < 1.1$.
69 The sPHENIX running period is 2023 - 2025, where 2023 was dedicated to commissioning,
70 2024 to pp collisions, and 2025 to Au+Au collisions. The main central barrel tracking detectors
71 are the Microvertex Detector (MVTX), Intermediate Tracker (INTT), Time Projection Chamber
72 (TPC), and TPC Outer tracker (TPOT), and the central barrel calorimetry system - Electromagnetic
73 Calorimeter (EMCAL) and Hadronic Calorimeter (HCAL). The sPHENIX detector schematic is
74 shown in figure 1 (left). The tracking detectors are capable of SRO. Even-though they are able to
75 record all data, the data volume of the TPC detector exceeds the capability of the computing center
76 and therefore a down-selection of what to save must be done. In the absence of online reconstruction
77 to identify interesting HF events, a random selection of events based on the minimum bias trigger
78 detector is done. The aim of this project is to reconstruct tracklets from the silicon detectors in
79 order to search for signatures of HF decays based on unique topology, and provide an additional
80 trigger to sPHENIX.

81 **Silicon detectors** The three innermost tracking layers (the MVTX detector) are based on Mono-
82 lithic Active Pixel Sensors (MAPS), the ALPIDE, originally developed for the ALICE experiment.

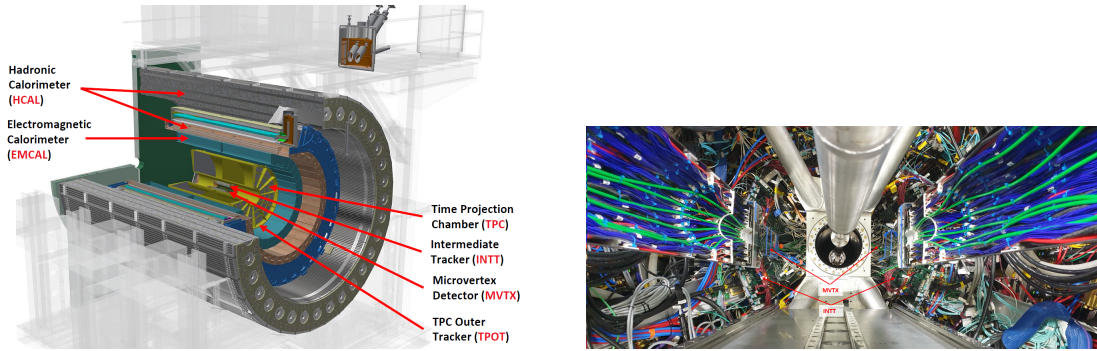


Figure 1. Left: The sPHENIX detector with highlighted central barrel detectors. Right: Installation of MVTX and INTT detectors.

83 The ALPIDE offers very fine pitch of $27 \mu\text{m} \times 29 \mu\text{m}$, collision event time resolution of $5 \mu\text{s}$; the
 84 MVTX contains a total of 270M channels. The next two layers (the INTT detector) contain silicon
 85 strip sensors (manufactured by Hamamatsu) with a pitch of $78 \mu\text{m} \times 16$ (or 20) mm with 360k
 86 channels in total. Both detectors are shown in figure 1 (right).

87 **sPHENIX readout, trigger and timing distribution** The schematic of the sPHENIX readout
 88 chain is shown in figure 2. The tracking detectors' Front-End Electronics (FEE) sends data to
 89 the Event Buffer and Data Compressor (EBDC) through the FELIX (FLX-712) [2] interface card.
 90 The calorimetry detectors' Front-End Modules (FEM) send data to the SubEvent Buffer (SEB)
 91 through the Data Collection Module (DCM2). The trigger and timing information is distributed
 92 via the Granule Timing Module (GTM). The Global Level 1 Trigger (GL1) and machine clock
 93 are transmitted to the FELIX cards for the tracking detectors and to the FEM for the calorimeters.
 94 There can be up to 4 LEMO and 4 Fiber connections to the GTM with 64 trigger inputs total.
 95 The most used trigger inputs are from the Hadronic Calorimeters for the cosmic and high energy
 96 jet triggers (normally not pre-scaled) and Minimum-Bias Detector for the beam collision trigger
 97 (heavily pre-scaled). The goal is to provide an additional high efficiency trigger input for HF events
 98 in pp collisions.

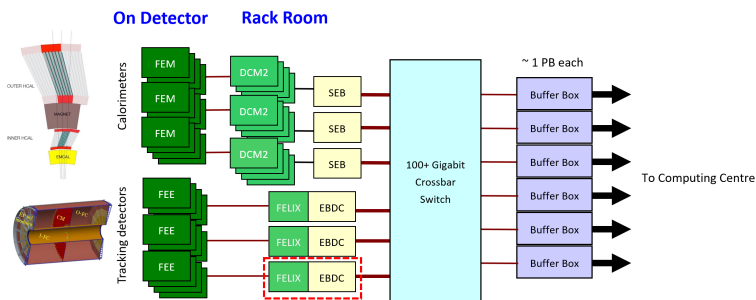


Figure 2. A schematic of sPHENIX readout chain. The calorimetry data are sent from Front-End Modules (FEM) to the SubEvent Buffer (SEB) through the Data Collection Module (DCM2), and the tracking data from Front-End Electronics (FEE) are sent to the Event Buffer and Data Compressor (EBDC) through a FELIX interface card.

99 **2.1 The DAQ-AI Data Flow**

100 Data from the MVTX and INTT are transmitted from the FELIX-DAQ board to the AI-Engine. It
 101 was decided that the AI-Engine will be hosted on a widely HEP community supported FELIX board
 102 for the following 3 reasons. First, FELIX offers 48 high-speed optical links to receive data; second,
 103 to reuse its Wupper module for the PCIe communication and associated software tools; and third,
 104 to use the sPHENIX infrastructure for tracking detectors that was built around the FELIX cards.
 105 The AI-Engine houses the raw data decoder, event builder, clusterizer, and GNN models to provide
 106 fast tracking and the trigger decision. The trigger decision is sent via a LEMO cable to the GTM.
 107 Since the decision is based on a event topology, a reference point (beam spot), which changes in
 108 time, must be precisely known and monitored. A GPU based feed-back system will be in place to
 109 process the data from the buffer boxes, reconstruct the beam spot position and update the position
 110 in the AI-Engine. The schematic of the data flow is shown in figure 3.

111 There are 144 optical links running at 3.2 Gbps per fiber for the MVTX alone, thus two engines
 112 will be used, one for each MVTX/INTT hemisphere. This will allow to have 24 links for MVTX and
 113 24 links for INTT. Since the DAQ, AI-Engine, and GTM sit in the Counting house, away from the
 114 radiation environment, it is not necessary to use radiation-hardened protocols. The FELIX optical
 115 links have been tested up to 14 Gbps with BER < 10⁻¹⁶ with an external loop-back measurement.
 116 INTT offers excellent time resolution to tag 100 ns RHIC bunch-crossing time to assign a unique
 timing to each collision event.

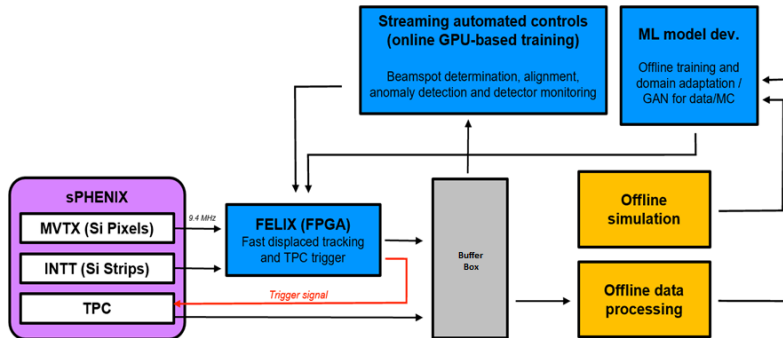


Figure 3. The schematic of the DAQ-AI data flow.

117

118 **2.2 Latency breakdown**

119 The TPC buffers can hold up to 30 μs of data. As we need to capture both sides around the triggered
 120 event, the aim is to deliver the trigger decision with 10 μs. The MVTX contributes up to 5 μs
 121 latency, the cables between the interaction region and counting house contribute around 0.3 μs
 122 (81m fibres). Forwarding data to the AI-Engine and decoding them takes up to 0.6 μs, which results
 123 from the maximum depth of FIFO holding the decoded hits (128). This leaves around 4 μs for
 124 the AI engine to perform tracking and trigger decision. The real latency of the decoding and data
 125 transmission depends on the occupancy, which is estimated to be around 50 physics hits per chip per
 126 event. The AI-Engine needs to ensure the latency is fixed, either by delaying the trigger decision,
 127 or vetoing the processing by decoding the bunch-crossing numbers from MVTX and INTT.

Table 1. Efficiency and Background Rejection with 1% and 0.1% signal/noise Ratios.

1% signal/background ratio			0.1% signal/background ratio		
Bg. rejection	Efficiency	Purity	Bg. rejection	Efficiency	Purity
90%	72.5%	7.25%	90%	78%	0.78%
99%	15.0%	15.0%	99%	17%	1.7%

128 3 Model description

129 The model is based on Graph Neural Network (GNN) using PyTorch and PyTorch geometric. The
 130 aim is to reconstruct the decay topology of the tracklets and search for secondary (displaced) vertices,
 131 which is one of the prominent features of HF decays. The displaced vertex will be $\sim O(100 \mu\text{m})$
 132 away from the primary vertex which has resolution $\sim O(10 \mu\text{m})$. We propose a novel method for
 133 treating events as graphs consisting of tracks as nodes and interconnection between tracks when they
 134 belong to the same particle decay rather than hit graphs [3]. Further improvement can be achieved
 135 by estimating the transverse momentum, p_T , based on the tracks. A 15% improvement is observed
 136 in the trigger decision performance, as expected. There are three stages in event processing: hit
 137 clustering, track reconstruction, and trigger decision. We use the GNN models in the second stage to
 138 reconstruct tracks, remove outliers (TrackGNN), and regress the track momentum onto the learned
 139 tracks. In the third stage, we design a bipartite GNN to reconstruct displaced vertices based on
 140 the corresponding tracks generated from decays and the estimated track momentum and estimate
 141 the probability of event being a trigger. The flowchart is shown in figure 4 and the efficiency and
 142 background rejection with 1% and 0.1% signal/noise ratios are summarized in Tab. 1.

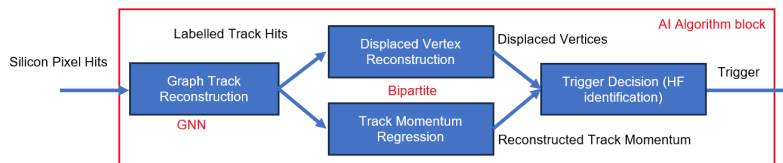


Figure 4. Flowchart of the GNN network.

143 4 Generation of the GNN IP core

144 We propose two parallel efforts to implement the software sPHENIX model onto FPGA, both use
 145 High-Level Synthesis (HLS) tools to generate synthesizable Register Transfer Level (RTL) code,
 146 i.e., Verilog. The first effort, manually translates the sPHENIX model into synthesizable C code
 147 and feeds it into the HLS tool, Vitis HLS [5], and then performs targeted optimization of the model
 148 following the FlowGNN architecture [6], which is the state-of-the-art GNN architecture on FPGA.
 149 The second effort is based on the hsl4ml framework [4], which is a generalized package to translate
 150 neural networks, such as deep neural networks (DNNs) and GNN, into an IP core. The target is to
 151 have a model implemented on FPGA that can process 100-200 nodes (hits) and 200-500 edges (hit
 152 connections) within $\sim O(10 \mu\text{s})$.

153 **FlowGNN Architecture with on-FPGA Implementation** The current TrackGNN model in
 154 sPHENIX we are using has one GNN layer, which includes 4 multi-layer perceptron (MLP) layers
 155 for both node and edge embedding with a dimension of 8. The proposed architecture follows the
 156 message-passing framework in FlowGNN [6]: the node embeddings are processed first, followed
 157 by an adapter to orchestrate the node information to the correct edge processing units for edge
 158 embedding computation and message aggregation. We also use quantization to reduce the data
 159 precision and to reduce the memory and computation requirements. We use `ap_fixed<18, 6>` for
 160 node embeddings, edge embeddings, and model weights and bias, and use `ap_fixed<21, 9>` for
 161 messages and input node features. The target FPGA board is the Alveo U280, approximately twice
 162 as big as FELIX FLX-712. The resource utilization is as follows: 194K (14.9%) LUT, 214K (8.2%)
 163 FF, 406 (20.2%) BRAM, and 488 (5.4%) DSP. The processing latency is measured on-board in an
 164 end-to-end fashion, including graph and weight loading, model computation, and results readback.
 165 The latency for an average-sized input (92 nodes and 142 edges) is $8.82 \mu\text{s}$ at a 285 MHz clock.
 166 The time-size scatter plot for nodes and edges are shown in figure 5, where the x-axis is the time
 167 spent processing one graph, and the y-axis is the graph size in terms of the number of nodes and
 168 edges. Compared to CPU calculations, the FPGA was 99.86% accurate.

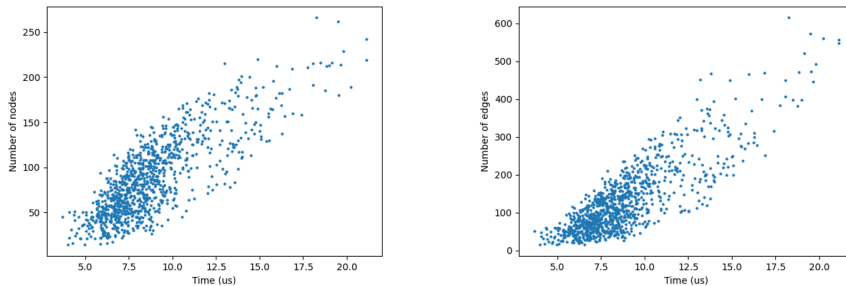


Figure 5. The time-size scatter plot of the TrackGNN model measured at Alveo U250. Left: for the nodes size. Right: for the edges size.

169 **hls4ml** A parallel effort to translate the model using hls4ml just started. The framework has been
 170 extensively used at CMS and we expect to have a first version of the inference at the end of October.

171 5 Summary and Outlook

172 Different modules of the AI-Engine has been tested independently, the final development is focused
 173 towards combining them into a single FPGA board. A further iterations with the model developers
 174 is done to account for the FPGA utilization. The beam test is expected in 2024 during RHIC pp run.

175 Acknowledgments

176 This project is funded by the United States Department of Energy, funding calls FOA-0002490,
 177 FOA-0002875 and Los Alamos National Laboratory LDRD program.

178 **References**

- 179 [1] sPHENIX Collaboration *sPHENIX Beam Use Proposal*, **sPH-TRG-2020-001** (2020).
- 180 [2] K. Chen, H. Chen, J. Huang, F. Lanni, S. Tang and W. Wu, *A Generic High Bandwidth Data*
181 *Acquisition Card for Physics Experiments*, **IEEE Trans. Instrum. Meas.** **69** (2020) no.7, 4569-4577,
- 182 [3] Xuan, T., Borca-Tasciuc, G., Zhu, Y., Sun, Y., Dean, C., Shi, Z. & Yu, D. *Trigger Detection for the*
183 *sPHENIX Experiment via Bipartite Graph Networks with Set Transformer.*, **Machine Learning And**
184 **Knowledge Discovery In Databases** (2023) pg. 51-67.
- 185 [4] F. Fahim, B. Hawks, C. Herwig, J. Hirschauer, S. Jindariani, N. Tran, L. P. Carloni, G. Di Guglielmo,
186 P. Harris and J. Krupa, *et al. hls4ml: An Open-Source Codesign Workflow to Empower Scientific*
187 *Low-Power Machine Learning Devices*, arXiv:2103.05579.
- 188 [5] Vitis High-Level Synthesis, Xilinx.
189 <https://www.xilinx.com/products/design-tools/vivado/high-level-design.html>
- 190 [6] R. Sarkar, S. Abi-Karam, Y. He, L. Sathidevi and C. Hao, *FlowGNN: A Dataflow Architecture for*
191 *Real-Time Workload-Agnostic Graph Neural Network Inference*, **2023 IEEE International**
192 **Symposium on High-Performance Computer Architecture (HPCA)**, Montreal, QC, Canada, 2023,
193 pp. 1099-1112, doi: 10.1109/HPCA56546.2023.10071015.