

# **Intelligent Experiment Through Real-Time AI: Fast Data Processing and Autonomous Detector Control for sPHENIX and Future EIC Detectors**

**Ming Liu  
Los Alamos National Lab  
For the Fast-ML Team**

**sPHENIX Collaboration Meeting  
Jan. 5-7, 2022**

# The Problem

## – Rare events hard to trigger by conventional methods

### sPHENIX challenge p+p@200GeV:

- Very high p+p collision rate:  $\sim 10\text{MHz}$ 
  - **Charm production rate:  $\sim 100\text{kHz}$** 
    - $0.5\text{mb}/42\text{mb} \sim 1\%$
  - **Beauty production rate:  $\sim 500\text{Hz}$** 
    - $2\text{ub}/42\text{mb} \sim 0.005\%$
- *No effective trigger to select low  $p_T$  HF events*
  - Triggered MB rate  $\sim 1\text{kHz} \ll 10\text{MHz}$ 
    - Lost most of the HF events at low  $p_T$
  - High  $p_T$  jet trigger,  $p_T > 10\text{GeV}$
  - Streaming readout  $\rightarrow$  huge data volume, DAQ/tape cost

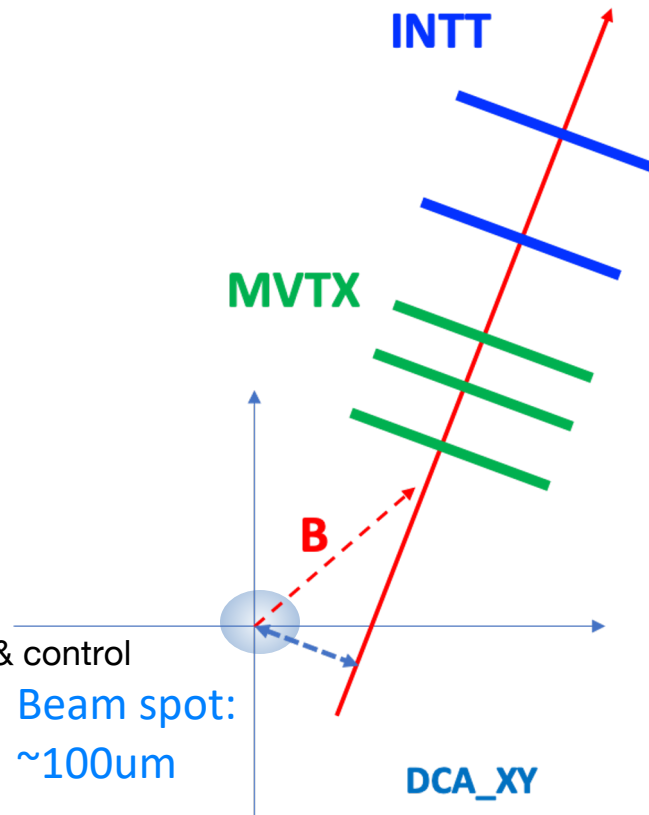
### EIC challenge:

- high rate e+p and e+A collisions

### Our approach:

Develop effective HF (or other rare events) triggers for p+p and e+p

- Streaming readout key detectors for high efficiency
- AI-based beam/detector monitoring and autonomous feedback & control
- ML-trained algorithm for HF tagging
  - **Rare physics don't require “slow detectors” for measurements**



# The DOE FOA Call in 2021

- Proposals called on 3/16, 2021
  - Short deadline, 4/30/2021**
  - Very intense work**



DEPARTMENT OF ENERGY  
OFFICE OF SCIENCE  
NUCLEAR PHYSICS



**DATA ANALYTICS FOR AUTONOMOUS OPTIMIZATION AND  
CONTROL OF ACCELERATORS AND DETECTORS**

- Initial team of NP, HEP and CS
  - LANL, MIT, FNAL and NJIT**
    - ORNL, CCNU and NTU joined later**

FUNDING OPPORTUNITY ANNOUNCEMENT (FOA) NUMBER:  
**DE-FOA-0002490**

ANNOUNCEMENT TYPE: INITIAL  
CFDA NUMBER: 81.049

## Intelligent experiments through real-time AI: Fast Data Processing and Autonomous Detector Control for sPHENIX and future EIC detectors

A proposal submitted to the DOE Office of Science  
April 30, 2021

- Embed AI/ML algorithms on fast FPGA-based trigger system
  - **Low trigger decision latency ~10us**
- Streaming readout key inner trackers to FPGAs to identify HF events through track topology
  - **High efficiency in HF tagging with AI/ML**
  - **HLS4ML package developed by HEP**
- Monitor and update beam-spot and detector alignment in real time
  - **Update geometry in real time**
- Send HF-trigger signal to the rest of other detectors
  - **Initiate readout if not already in the data stream**

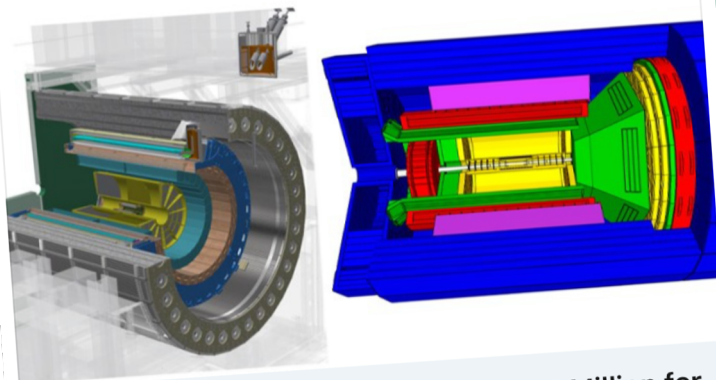


# DOE Awards Announced 12/2/2021

**\$1.5M for our proposal  
FY22-23**

 Brookhaven National Laboratory  
26,734 followers  
3w • 🌐

New funding from the **U.S. Department of Energy (DOE)** will support artificial intelligence advancements at **#RHIC** and the **#ElectronIonCollider**.



**Department of Energy Announces \$5.7 Million for  
Research on Artificial Intelligence and Machine...**

bnl.gov • 2 min read

Office of Science

## Department of Energy Announces \$5.7 Million for Research on Artificial Intelligence and Machine Learning (AI/ML) for Nuclear Physics Accelerators and Detectors

DECEMBER 2, 2021

Office of Science »

Department of Energy Announces \$5.7 Million for Research on Artificial Intelligence and Machine Learning (AI/ML) for Nuclear Physics Accelerators and Detectors

*Projects will advance understanding of atomic structure and the nature of matter and antimatter*

**WASHINGTON, D.C.** - Today, the **U.S. Department of Energy (DOE)** announced \$5.7 million for six projects that will implement artificial intelligence methods to accelerate scientific discovery in nuclear physics research. The projects aim to optimize the overall performance of complex accelerator and detector systems for nuclear physics using advanced computational methods.

"Artificial intelligence has the potential to shorten the timeline for experimental discovery in nuclear physics," said Timothy Hallman, DOE Associate Director of Science for Nuclear Physics.  
"Particle accelerator facilities and nuclear physics instrumentation face a variety of technical challenges in simulations, control, data acquisition, and analysis that artificial intelligence holds promise to address."

# The Team



- LANL (NP)
  - **Yasser Corrales, Cameron Dean, Zhaozhong Shi, Noah Wuerfel, Kun Liu, Cesar da Silva, Hugo Pereira da Costa, Ming Liu ... new PDs**
- MIT (NP, HEP)
  - **Gunther Roland, Philip Harris (HLS4ML), Yen-Jie Lee, Or Hen, Cristiano Fanelli et al**
- FNAL(HEP)
  - **Nhan Tran(HLS4ML), Engineer, Yu-Dai Tsai (Theorist, ML) et al**
- NJIT(CS)
  - **Dantong Yu, students + PDs**
- ORNL(NP)
  - **Jo Schambach**
- CCNU(EE, NP)
  - **Kai Chen(FELIX), Yaping Wang et al**
- NTU (CS)
  - **Fu Song, students + PDs**

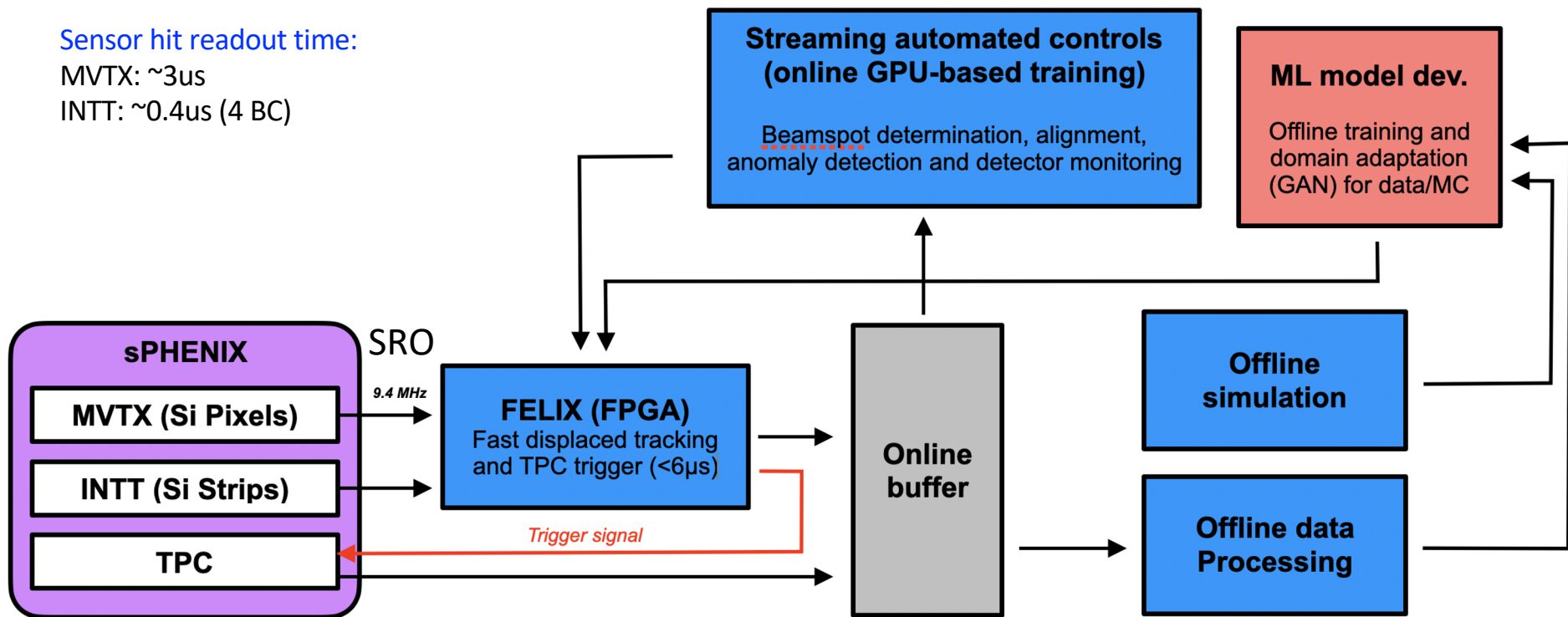
In collaboration with experts from BNL - Jin Huang, Martin Purschke, John Haggerty et al

# HF AI Trigger: sPHENIX as a Test Ground

Sensor hit readout time:

MVTX:  $\sim 3\mu\text{s}$

INTT:  $\sim 0.4\mu\text{s}$  (4 BC)



sPHENIX DAQ & Trigger integration challenge

# Timeline

2021

2022

2023

2024

2030+

- Project started
- Initial simulations constructed
- First data for algorithm training

- MVTX & INTT SRO
- Fast tracking algorithms in place
- GPU feedback machine R&D
- Initial FPGA bitstream

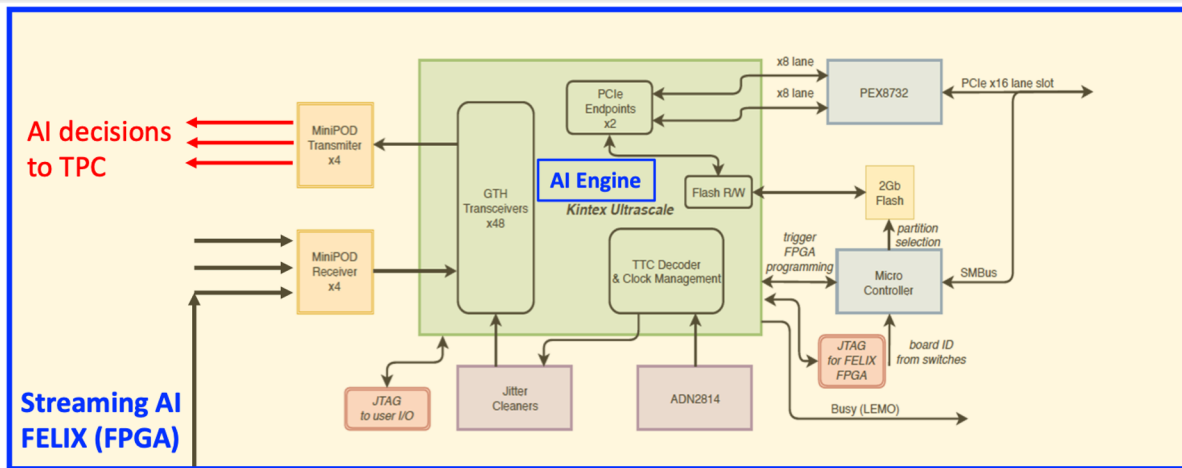
- Refine interface between system and detectors
- Improve algorithms with latest data stream
- Pre-commissioning

- Deploy device at sPHENIX
- pp/pA run

- Design updated system for EIC
- Take advantage of new technology if required

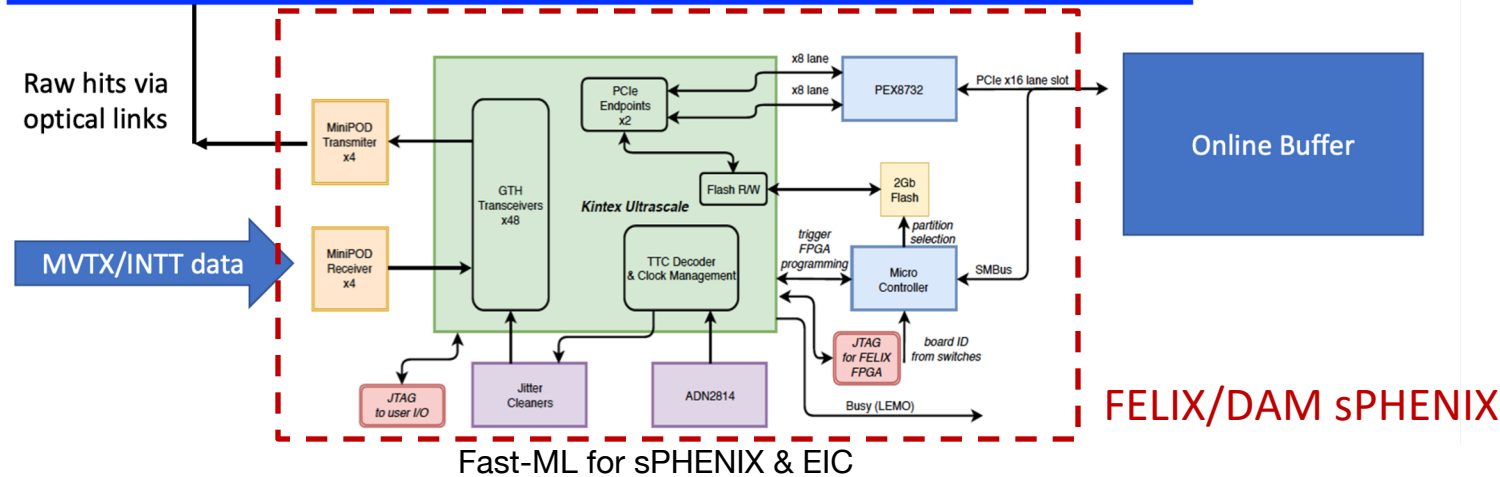
- Deploy device at EIC

# A Prototype Implementation Proposal



Trigger boards, TBD

- FELIX
- CMS trigger boards
- Other options



FELIX/DAM sPHENIX

Fast-ML for sPHENIX & EIC

(Pythia + GEANT) → MVTX/INTT hit maps → raw data (pixel hits) in JSON

- 1) trigger AI-ML training
- 2) generate raw real data like electronics bit stream for hardware simulations

Cameron, Zhaozhong, Jin,  
Yasser, Noah et al

# MVTX + INTT: 3 + 2 layers

## INTT

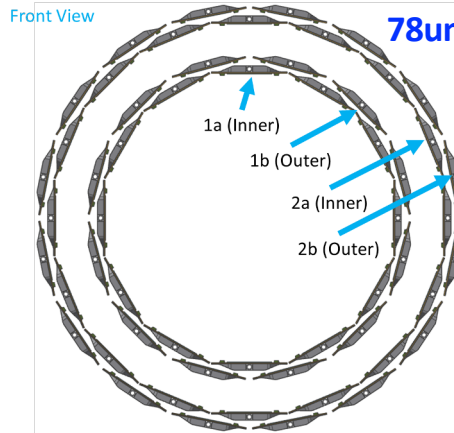
Barrel	Center of Sensor Tangent Radius (mm)
1	-
1a (Inner)	71.88
1b (Outer)	77.32
2	-
2a (Inner)	96.80
2b (Outer)	102.62

## MVTX

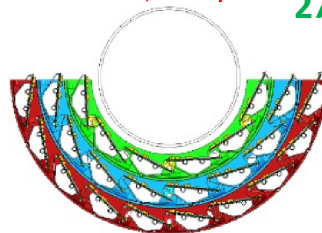
R (mm)	min	mid	max
Layer 0	24.61	25.23	27.93
Layer 1	31.98	33.35	36.25
Layer 2	39.93	41.48	44.26

## Silicon Strips:

78um x 16mm (A)/20mm (B) INTT

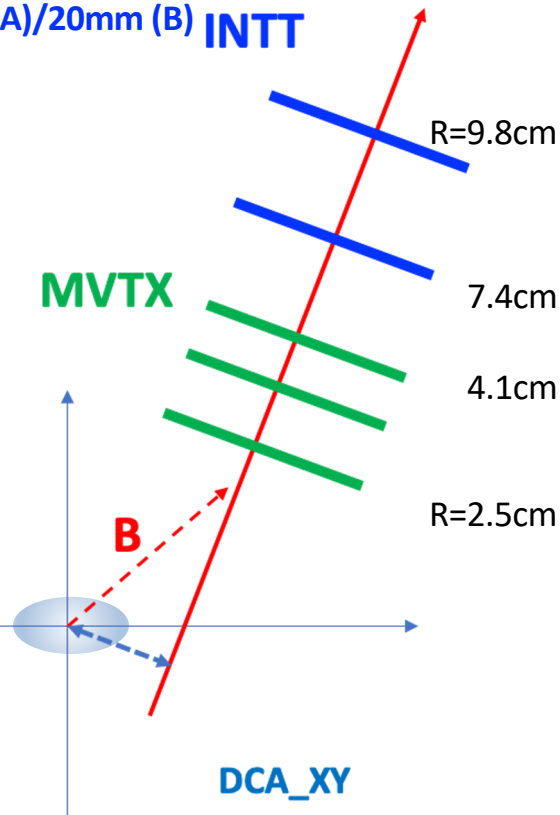


3-layer sensor barrel  
- 48 staves, 432 chips

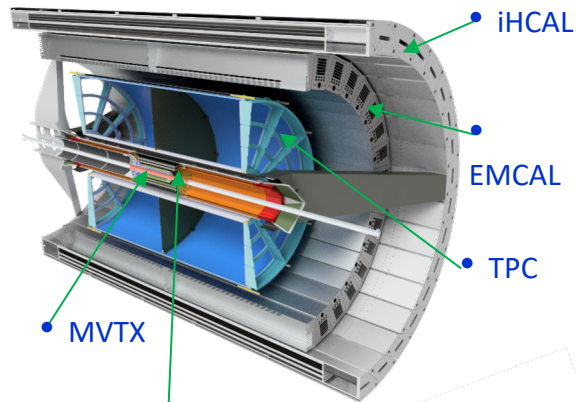


Silicon pixels:  
27um x 29um

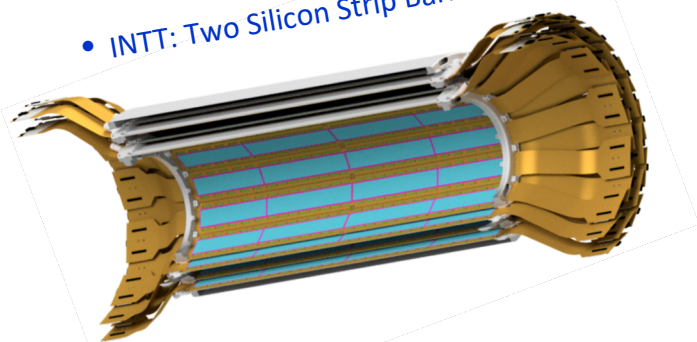
## MVTX



# Event Timing: Reject out of time MVTX hits with INTT<sub>sPHENIX</sub>



• INTT: Two Silicon Strip Barrels



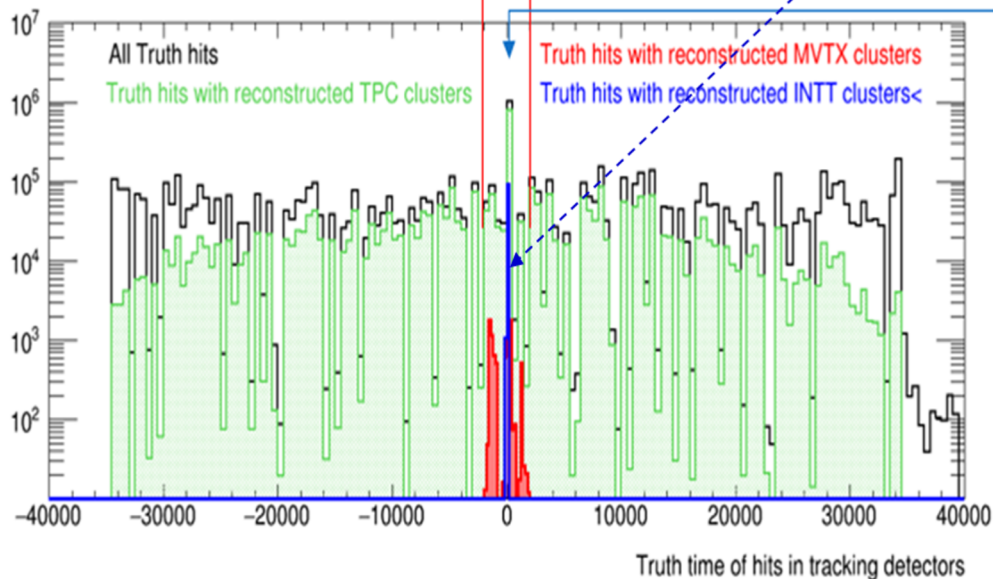
Collisions:  $\pm 35 \mu\text{s}$

TPC:  $\pm 35 \mu\text{s}$

MVTX:  $\pm 2 \mu\text{s}$

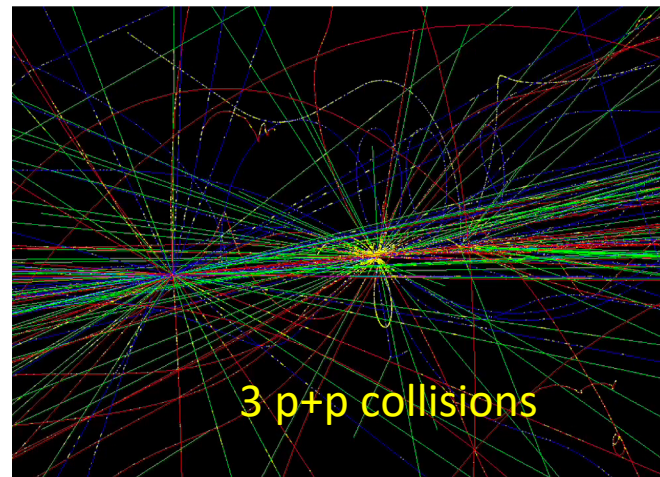
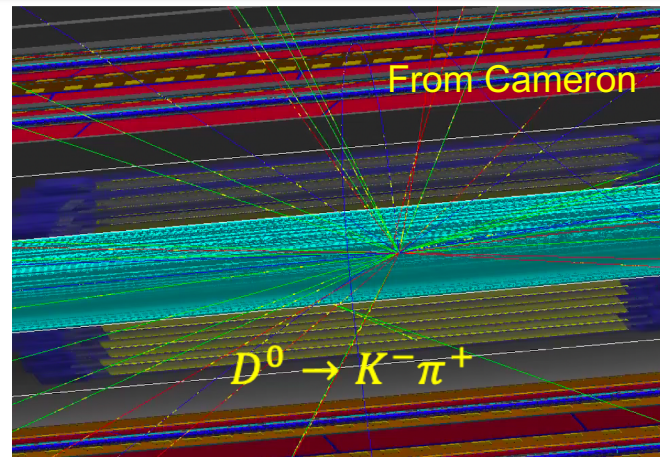
INTT:  $[-20 \text{ ns}, 80 \text{ ns}]$

A.U.



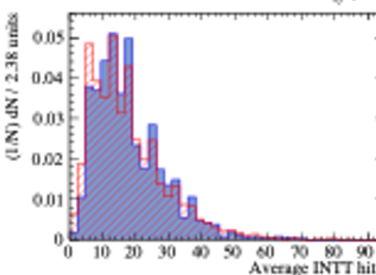
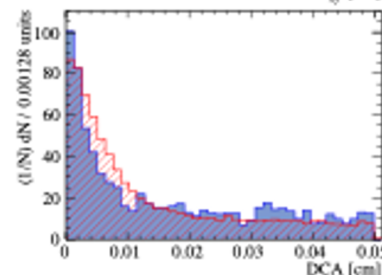
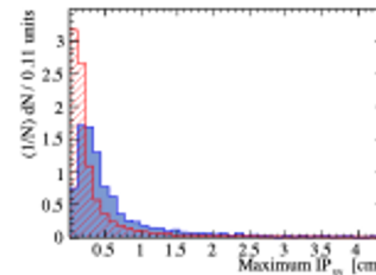
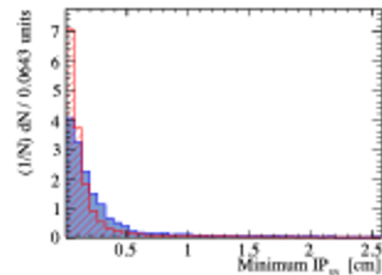
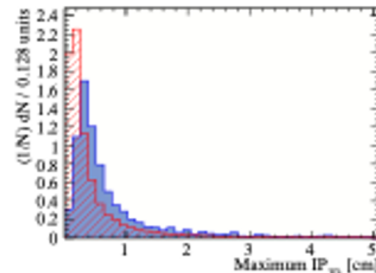
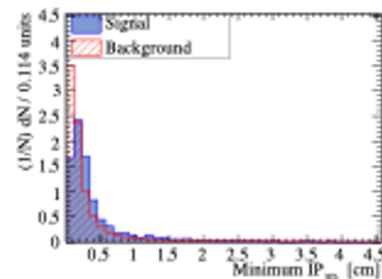
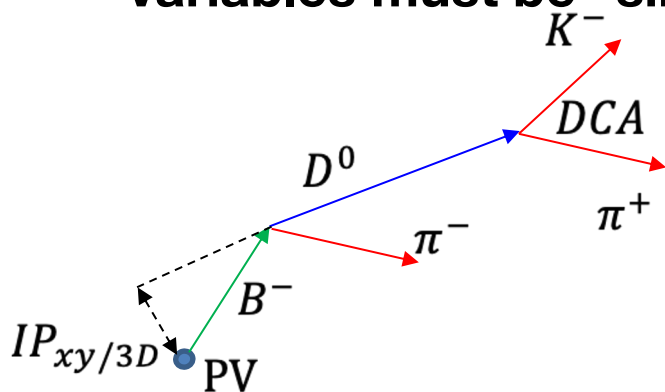


- Can simulate any number of signal and background events with full digitization
- Package developed to extract raw hit information, used for
  - **algorithm training (JSON output)**
  - **sim data to raw data bit pattern**



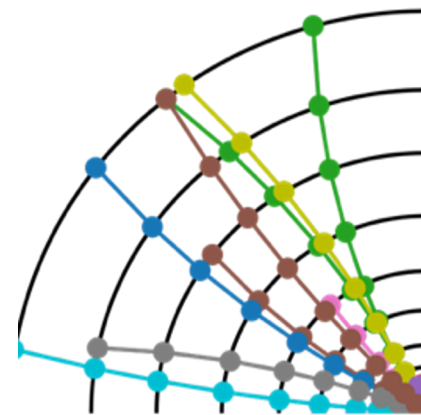
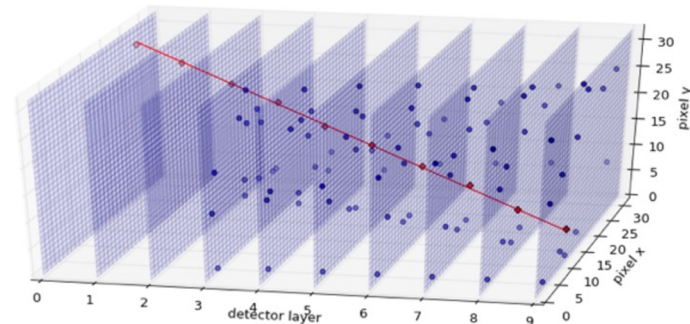
# Case study: AI HF selections

- Is ML better for selecting HF decays over conventional selections?
- Challenge:
  - Must run online, in FPGA. Hence variables must be “simple”



## Moving from images to points

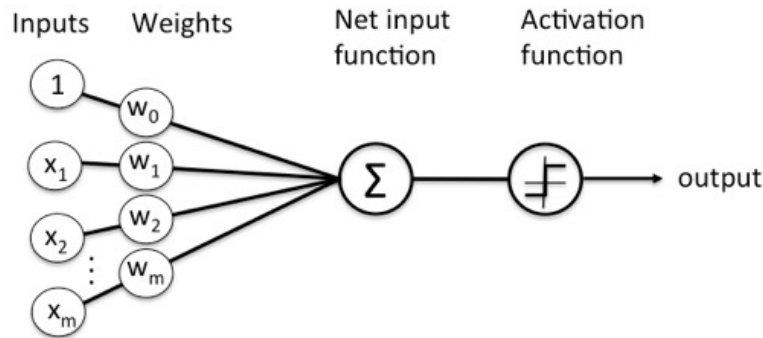
- **Image-based methods face challenges scaling up to realistic HL-LHC conditions**
  - High dimensionality ( $1k * 0.5K * 9$  per MVTX stave alone) and sparsity
  - Irregular detector geometry
- **Instead of forcing the data into an image, use the space point representation**
  - Harder to design models (variable-sized inputs/outputs, MVTX + INTT)
  - But now we can exploit the structure of the data with full precision
- **What ML models are appropriate for the event**
  - Recurrent Neural Networks and Graph Neural Networks



Implemented several models to solve the trigger detection problem:

- Directly applied GNN model to trigger detection problem (GNN)
- Added a global vector to the GNN model to represent some global feature (VPGNN)
- DiffPool model (DiffPool)
- VpGNN + DiffPool (GNNDiffPool)
- ParticleNet , Georgian

Another model we tried: Set2Graph (Affinity Matrix Prediction)



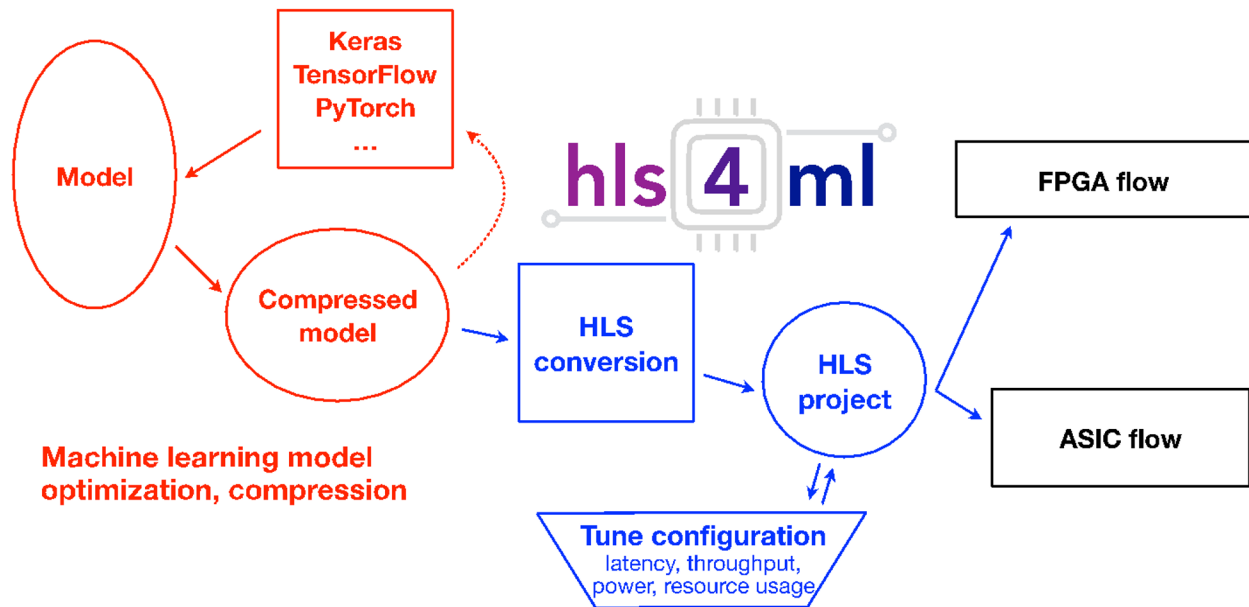
# Initial Study from NJIT: Efficiency and Purity

Proper mix of simulated data: 100(BG) + 1 HF(charm) events



	Current performance		Goal 1	Goal 2
Efficiency	50%	25%	20-50%	90%
Purity	5%	5%	5%	5%
Background Rejection	90%	95%	99%	95%

- Algorithms must have low latency and resource usage
- hls4ml** translates NN algorithms into high level synthesis
- Also generates IP cores for easy implementation



Red – typical ML algorithm development stages

Blue – HLS conversion to FPGA IP

Black – typical implementation onto chips

# AI-Trigger SW/FW Pipeline - Status

1. Fetch events from event buffer (Work in Progress, sim to raw data)



2. Data Pre-processing Clustering (Work in Progress on FPGA implements)



3. Tracking + Outlier hits Removal (Done in FPGA)



4. Triggering (Done in FPGA, need performance tuning)



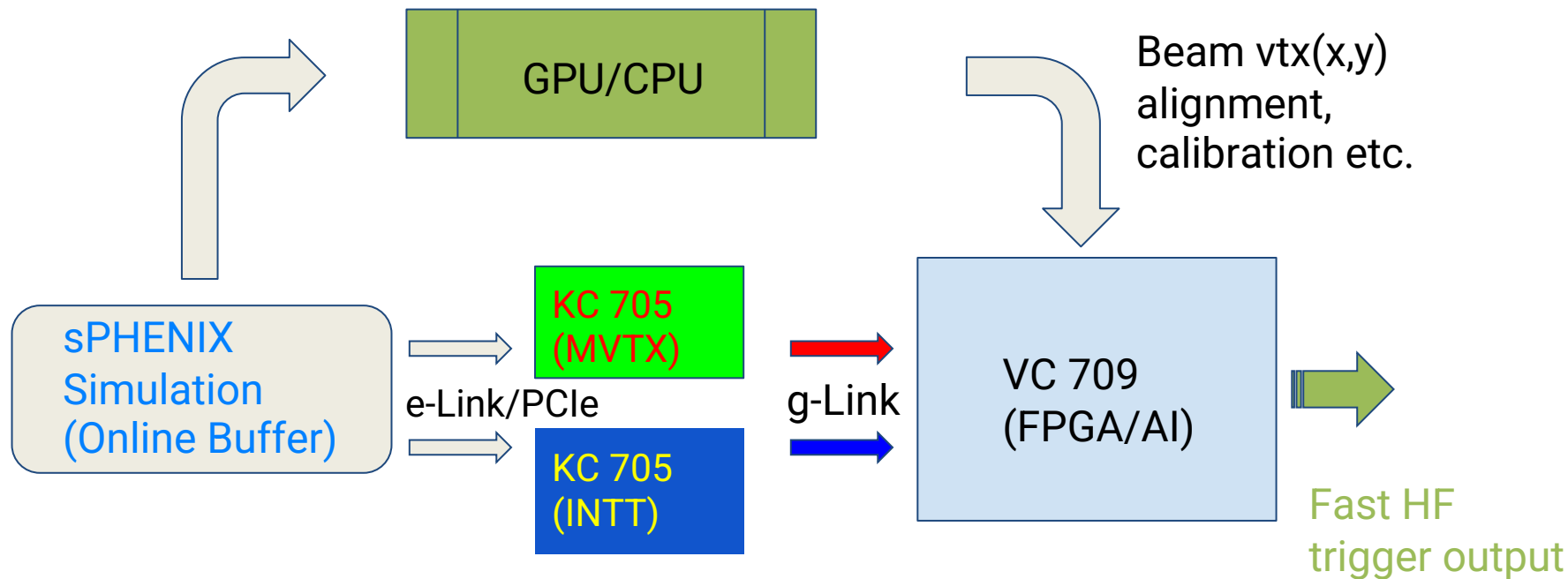
5. Triggers on TPC (Interface and integration with sPhenix Detector, last step)

## Hardware implementation

1. data stream processing
2. Test AI algorithms on FPGA



# A Toy Model – Hardware Implementation (sPHENIX)

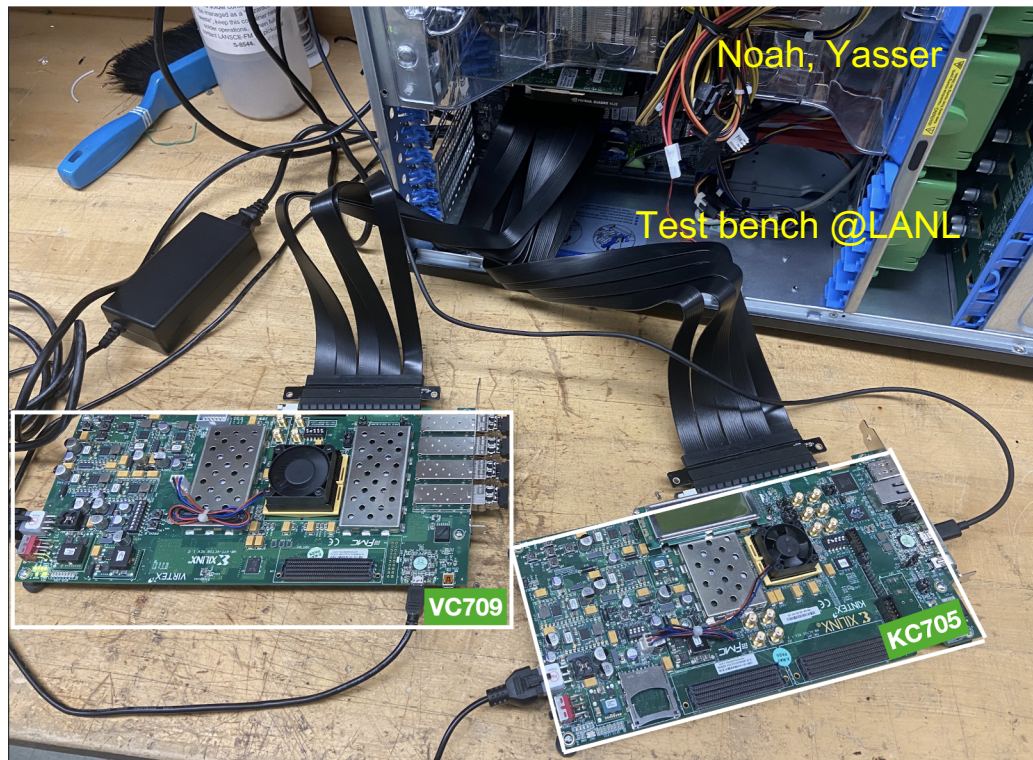


Streaming readout sim data:

8b/10b MVTX/INTT data (KC705) to FPGA/AI Engine (VC709)

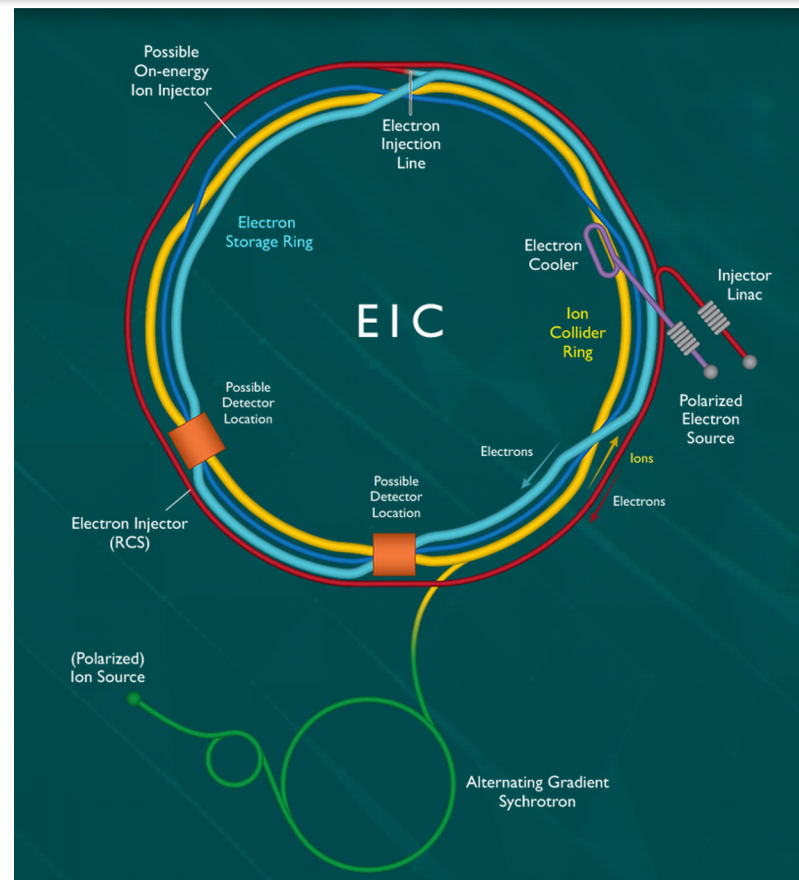
# Realizing in Firmware

- FELIX card shares same FPGA as Xilinx VC709, ideal testing ground
- KC705 represents our MVTX+INTT Data Aggregation Module
- Successfully transmit data from host PC to DMA
  - **Convert MVTX sim data to real-data-like bit-stream in progress**
- Next:
  - **Transmit MVTX/INTT sim data to VC709(AI-Engine) through G-Links**



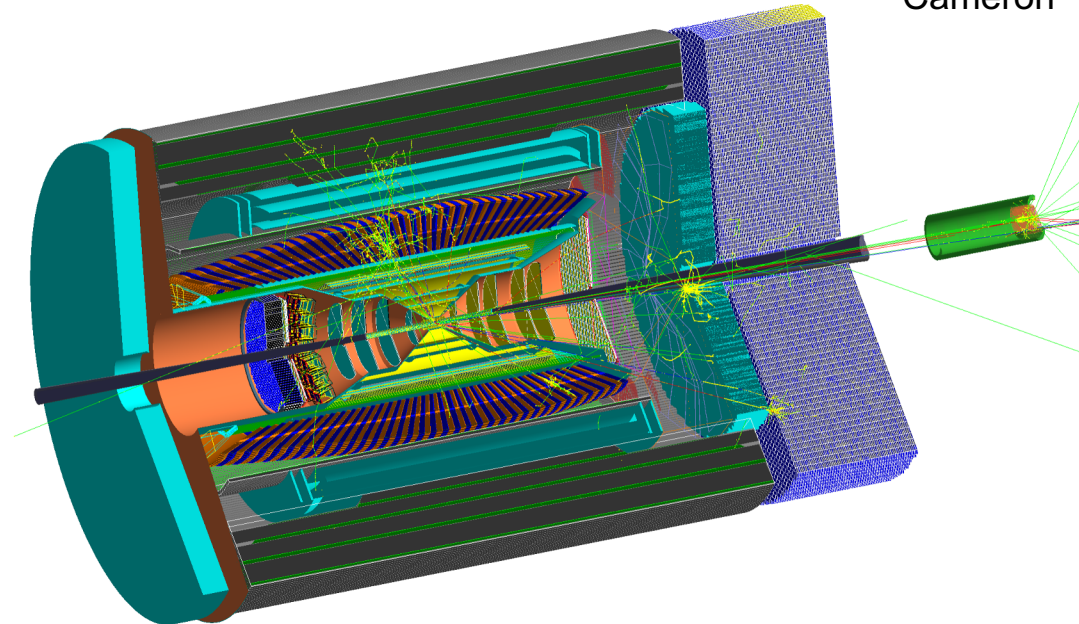
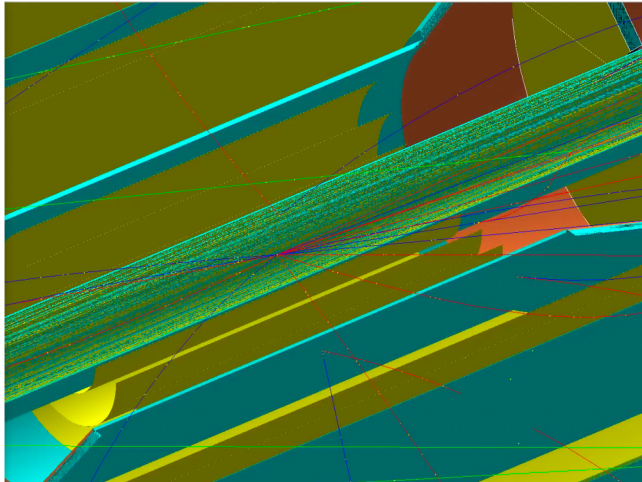
# The Electron-Ion Collider

- Next generation accelerator
  - To be operating at BNL from the early 2030s
  - the future of nucleon structure probes and many other studies
- Three collaborations have submitted detector proposals:
  1. ATHENA
  2. CORE
  3. ECCE



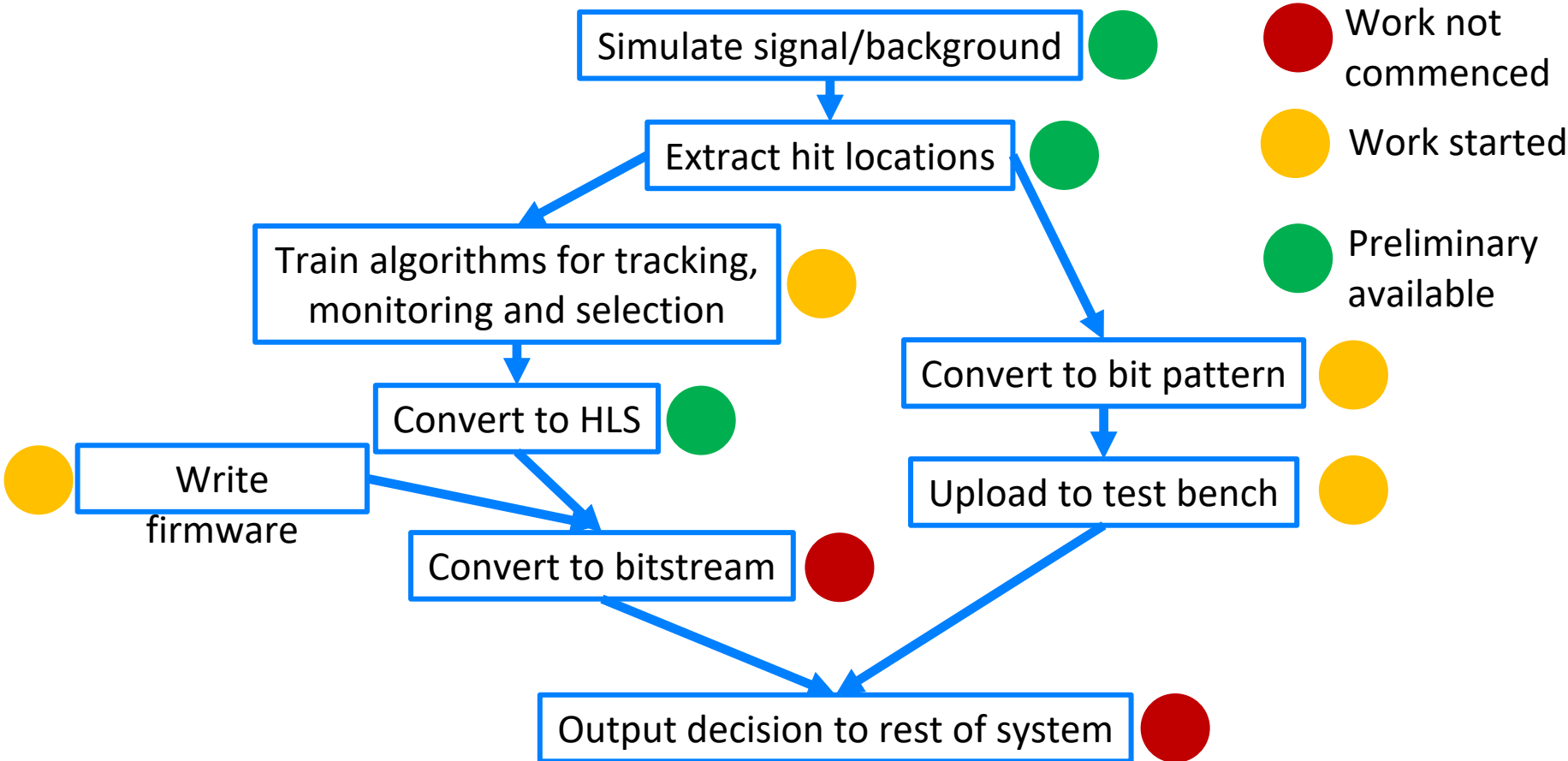
# Simulating events (ECCE)

- EIC physics simulations progressed rapidly in 2021
- Large volume of data already at hand (>800M events)
- No digitization yet but we can use smeared hits to understand potential





# Project Workflow Status



# Summary of Progress

---

Success stories since proposal approved

1. Full Geant4 simulations of MVTX and INTT plus Geant4 simulation of EIC detectors
2. Tracking GNN algorithms are being developed at NJIT
3. Prototype hardware set up at LANL with host-to-client transfers running
4. Second lab (FELIX) being set up at MIT
5. HLS4ML development at Fermilab, MIT and NJIT
6. FELIX FW development at ORNL and LANL

Next several months:

1. **Convert simulation output to equivalent bit pattern through G-Link**
2. **Develop initial tracking and selection algorithms**
3. **Convert algorithms to HLS code to run on FPGA**
4. **Pass simulated data to FPGA as if it were real data**

Goals:

1. **Build a full prototype and benchmark performance with simulations by 2023;**
  2. **Install device in sPHENIX before 2024 (RHIC *pp* run)**
- Project will significantly improve sPHENIX HF capabilities
  - Project relies on inner tracker MVTX and INTT SRO
  - After successful deployment at sPHENIX, focus shifts to future EIC detectors

# Backup

---



# from sPHENIX to EIC



- sPHENIX takes data from 2023
  - **Can be used as a proof-of-principle (as well as a real use case)**
- EIC has lower average multiplicity
  - **relatively easier to select**
  - **likely to use similar tracker technology to MVTX (ITS-2 vs ITS-3)**
- Large overlap of team between sPHENIX and EIC/ECCE
  - **knowledge preservation**
  - **share a simulation framework**

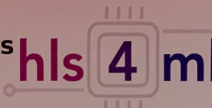
# Discussions at MIT, 12/14/2021



Phil, Yasser and Ming @MIT 12/2021



## Fast Inference of Deep Neural Networks for Real-Time Physics Applications

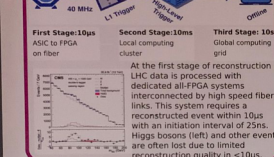


Giuseppe Di Guglielmo<sup>a</sup>, Javier Duarte<sup>a</sup>, Song Han<sup>a</sup>, Philip Harris<sup>b</sup>, Burt Holzman<sup>c</sup>, Sergio Jindariani<sup>d</sup>, Edward Kreinar<sup>d</sup>, Benjamin Kreis<sup>e</sup>, Vladimir Loncar<sup>f</sup>, Jennifer Nagdiub<sup>g</sup>, Maurizio Perini<sup>h</sup>, Dylan Rankin<sup>i</sup>, Ryan Rivera<sup>j</sup>, Sioni Summers<sup>k</sup>, Nhan Tran<sup>l</sup>, Zhenbin Wu<sup>m</sup>

<sup>a</sup>Columbia University, New York, NY 10027, USA, <sup>b</sup>Fermi National Accelerator Laboratory, Batavia, IL 60510, USA, <sup>c</sup>Massachusetts Institute of Technology, Cambridge, MA 02139 USA, <sup>d</sup>Travis Eye360, Herndon, VA 20176, USA, <sup>e</sup>IBM, CH-1211 Geneva 23 Switzerland, <sup>f</sup>Imperial College London, London, SW7 2BZ, UK, <sup>g</sup>University of Illinois at Chicago, Chicago, IL 60607, USA

### High Energy Physics data flow

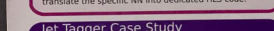
LHC data rates are >100 Tbps. Data is filtered down and saved at a rate of 1 GB/s. The filtering involves 3 tiers of reconstruction.



At the first stage of reconstruction LHC data is processed with dedicated all-FPGA systems interconnected by high speed fiber links. This system requires a reconstructed event within 10μs with an initiation interval of 25ns. High-p bosons (left) and other events are often lost due to limited reconstruction quality in <10μs.

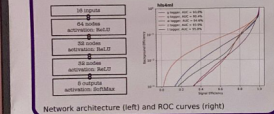
### HLS4ML Flow

To enhance algorithmic complexity we introduce HLS4ML, a tool to translate DNN to Vivado HLS then to an FPGA.



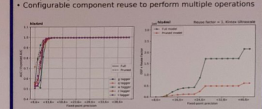
To ensure the lowest latency and initiation interval we translate the specific NN into dedicated HLS code.

### Jet Tagger Case Study



To achieve the smallest initiation interval (below 25ns) and latency (below 1μs) the whole algorithm is unrolled onto the FPGA. Also, arrays are fully partitioned and no block RAM is used. We aim to minimize FPGA resources through features:

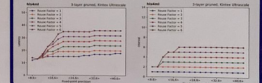
- Fixed-point number representation and reduced precision
- Compressing NN to reduce the size, and number of operations
- Configurable component reuse to perform multiple operations



With the jet tagger model, inference with 16-bit fixed-point numbers achieves equivalent performance to 32-bit floating-point (left). Using L1 Regularization, 70% of weights can be removed, without impacting on performance, reducing the resource utilization of the model.

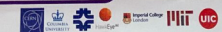


Reuse trade resources for latency and throughput by explicitly restricting the DSP to perform multiple operations (left). The impact of reuse on DSP utilization is shown in the right figure. It scales linearly with reuse. Inference latency increases with a larger reuse factor (left). With the targeted 200 MHz clock frequency, latencies as low as 75ns are achieved with an initiation interval of 1 clock (5ns).

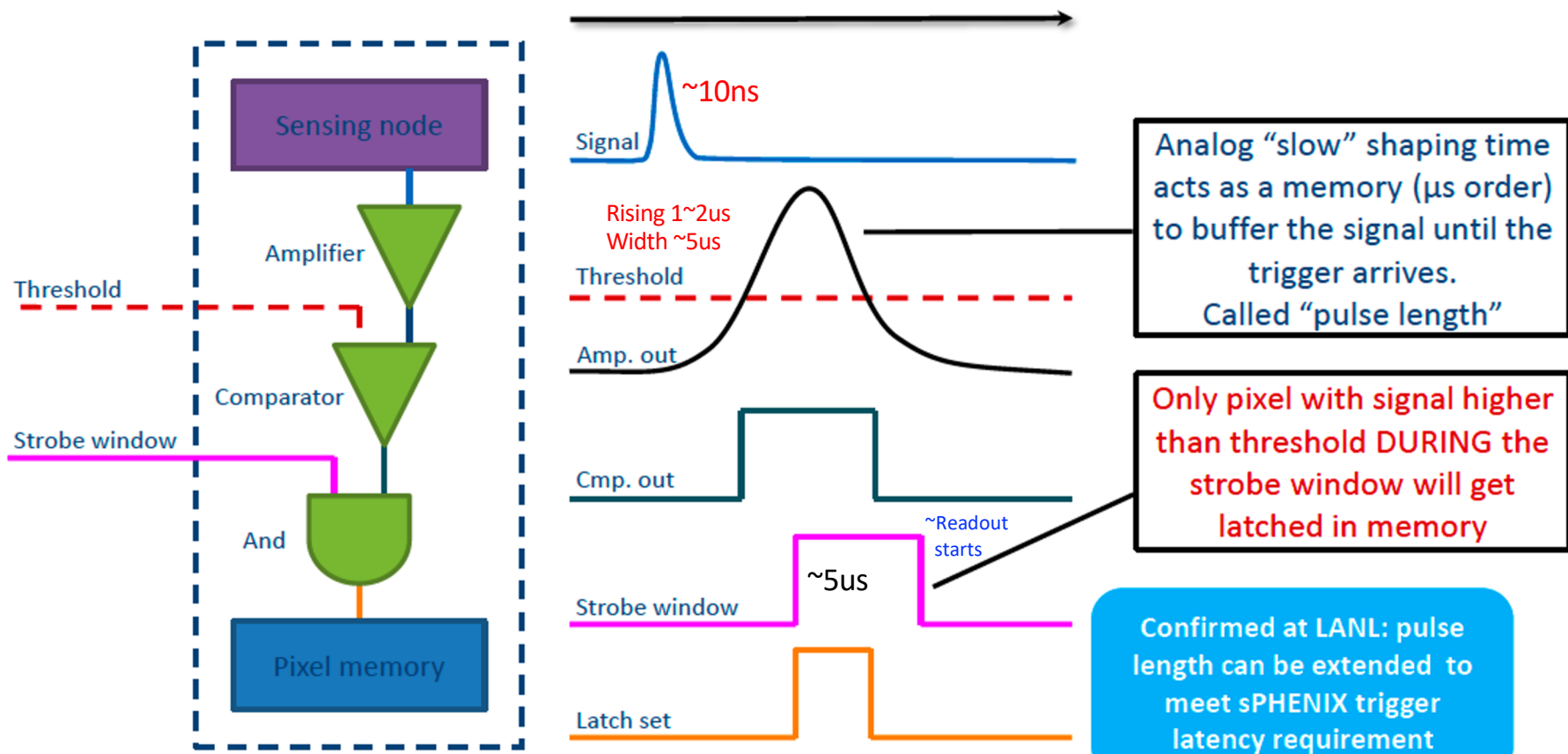


### Outlook

HLS4ML project has been tested on dedicated boards and will be used in the next running of the LHC. The project is expanding to cover most modern DNN architectures and support for network architectures up to several million weights. Full status is shown below. (\*denotes large NN support is being developed) <https://hls4ml-machine-learning.github.io/status/>

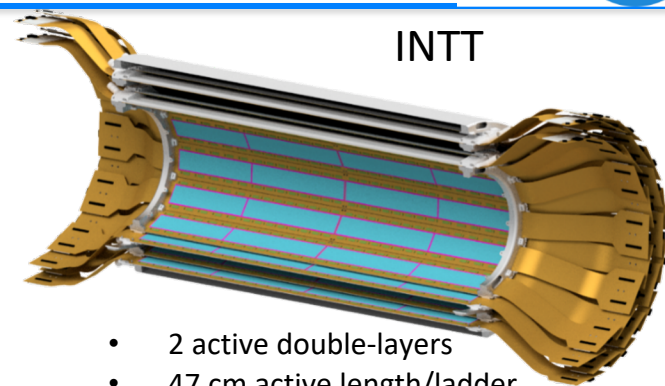


# ALPIDE Timing

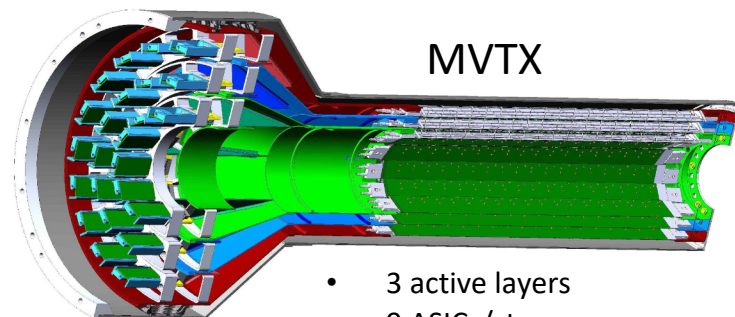


# Tracking at sPHENIX

- Tracking consists of 3 sub-detectors:
  - Pixel Vertex Detector (MVTX)
  - Intermediate Silicon Tracker (INTT)
  - Time Projection Chamber (TPC)
- MVTX and INTT are both capable of streaming readout
- Combined tracking to  $r = 10.3$  cm



- 2 active double-layers
- 47 cm active length/ladder
- Silicon strip detector



- 3 active layers
- 9 ASICs/stave
- 27 cm active length/stave
- Pixel detector

# sPHENIX HF constraints



- sPHENIX has great tracking and calorimetry
- However, limited by calorimetry backend readout rate (15kHz) in triggered mode
- RHIC pp rate is  $\sim 10$  MHz
- Plan: Use tracker SRO to recover some heavy flavor physics potential

Year	Species	$\sqrt{s_{NN}}$ [GeV]	Cryo Weeks	Physics Weeks	Rec. Lum. $ z  < 10$ cm	Samp. Lum. $ z  < 10$ cm
2023	Au+Au	200	24 (28)	9 (13)	3.7 (5.7) nb $^{-1}$	4.5 (6.9) nb $^{-1}$
2024	$p^{\uparrow}p^{\uparrow}$	200	24 (28)	12 (16)	0.3 (0.4) pb $^{-1}$ [5 kHz] 4.5 (6.2) pb $^{-1}$ [10%-str]	45 (62) pb $^{-1}$
2024	$p^{\uparrow}$ +Au	200	–	5	0.003 pb $^{-1}$ [5 kHz] 0.01 pb $^{-1}$ [10%-str]	0.11 pb $^{-1}$
2025	Au+Au	200	24 (28)	20.5 (24.5)	13 (15) nb $^{-1}$	21 (25) nb $^{-1}$

sPHENIX beam-use proposal. 5 kHz refers to final rate with triggered readout, 10%-str refers to 10% streaming readout

# How are FPGAs programmed?

## Hardware Description Languages

HDLs are programming languages which describe electronic circuits

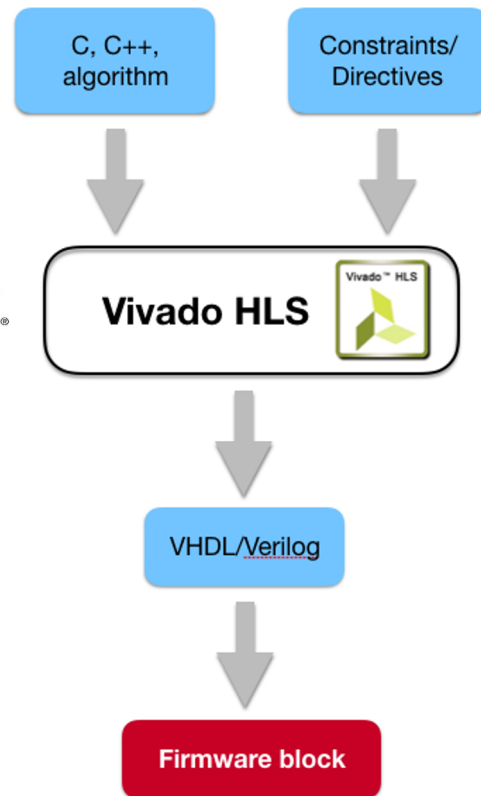
## High Level Synthesis

Compile from C/C++ to VHDL

Pre-processor directives and constraints used to optimize the design

**Drastic decrease in firmware development time!**

Today we'll use Xilinx Vivado HLS [\*]



[\*] [https://www.xilinx.com/support/documentation/sw\\_manuals/xilinx2020\\_1/ug902-vivado-high-level-synthesis.pdf](https://www.xilinx.com/support/documentation/sw_manuals/xilinx2020_1/ug902-vivado-high-level-synthesis.pdf)

- Aim to develop algorithms as Graph Neural Networks (GNN)
- Advantageous over Convolutional Neural Networks (CNN) by adding edge information
- Detector and physics knowledge will improve predictions
- Algorithms deployed at several points:
  1. **Fast tracking on FPGA**
  2. **Topological separation of HF signals on FPGA**
  3. **Beam-spot and anomaly detection on GPU**
    - Part of feedback system to improve 1 & 2 plus inform detector operators

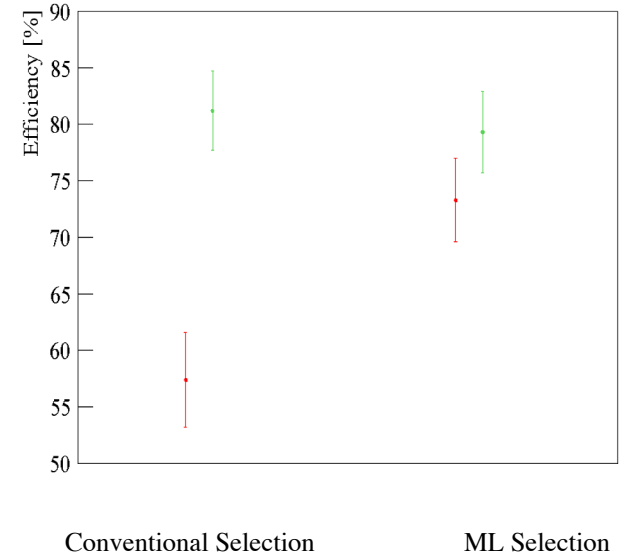
# Constructing ML algorithms

- Anomaly detection is important in all experiments
- RHIC experiments cannot be accessed during beam to fix issues
- Aim to use variational autoencoders
- Incoming data can be compared for:
  1. Noisy pixels
  2. Dead strips
  3. Change in beam spot or alignment
- Pass info. back to selection system to improve yield
- Pass back to control room



# Case study: AI HF selections

- Several algorithms trained using TMVA
  - Fast turnaround due to proposal time constraints
  - Algorithms used “out-of-the-box”, no optimizations
- Trained using samples with no HF signal and with  $D^0 \rightarrow K^- \pi^+$  signal
- Selection tuned for approx. equal signal efficiency



Green – The signal selection efficiency  
Red – The background rejection efficiency

# What we actually used in Json Data Files? SPHENIX

---

‘RawHit’ contains two part: [u'MVTXHits', u'Description']

- ‘MVTXHits’ contains all the hits information.

Each hit contains: [u'Coordinate', u'ID']

- ‘ID’ contains: [u'Layer', u'PixelZIndex', u'Chip', u'Stave',  
u'PixelPhiIndexInLayer', u'HitSequenceInEvent',  
u'PixelPhiIndexInHalfLayer', u'PixelHalfLayerIndex',  
u'Pixel',  
u'HalfLayer']

# Training Dashboard

