# Fast ML in the NSF HDR Institute A3D3

Shih-Chieh Hsu

University of Washington
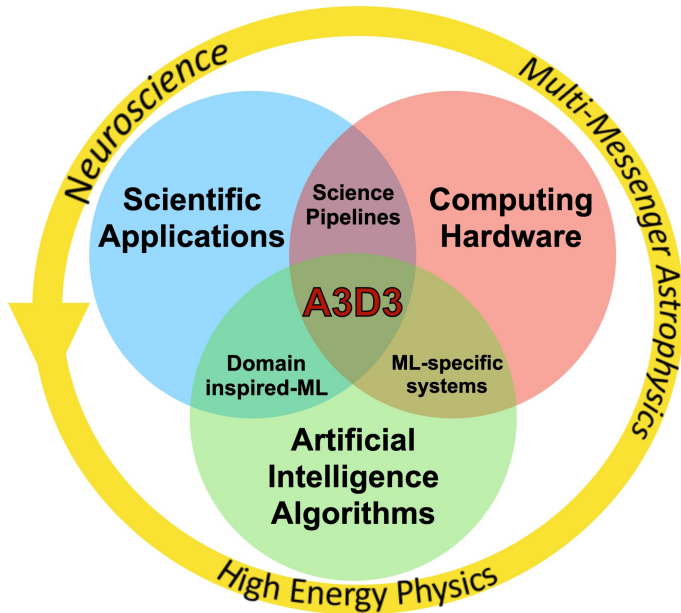
FastML Workshop ICCAD

Nov 2 2023

https://fastmachinelearning.org/iccad2023/program.html
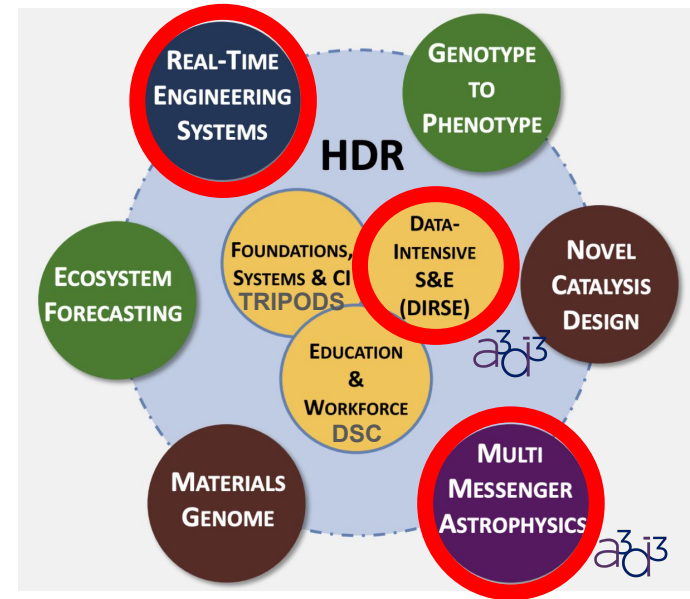
https://a3d3.ai/

# NSF HDR Institute: *A*ccelerated *A*rtificial Intelligence *A*lgorithms for *D*ata-*D*riven *D*iscovery (since 2021)



- **Our mission:**
  To enable real-time AI techniques for scientific and engineering discovery by uniting three core components: Scientific Applications, Artificial Intelligence Algorithms, and Computing Hardware

- **Our vision:**
  To make real-time AI accessible to the scientific and engineering community in order to accelerate discovery.

# Harnessing the Data Revolution

- A national-scale initiative to enable new modes of **data-driven discovery** addressing fundamental questions in science & engineering
- Three parallel tracks:
  - Institutes (**5** awards, $75M)
    - **A3D3**
    - I-GUIDE
    - iHARP
    - Imageonics
    - ID4
  - Ideas Labs + Frameworks (28, $53M)
  - TRIPODS (28, $42M) & DSC (19, $25M)

# Multi-disciplinary multi-institution

Spread across **16** institutions globally and **106** members (**70%** students + postdocs).

ICCAD FastML organizers associated to A3D3
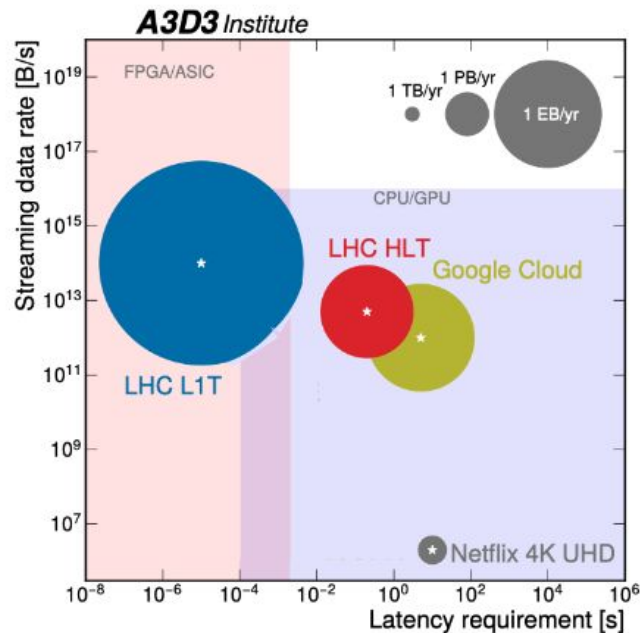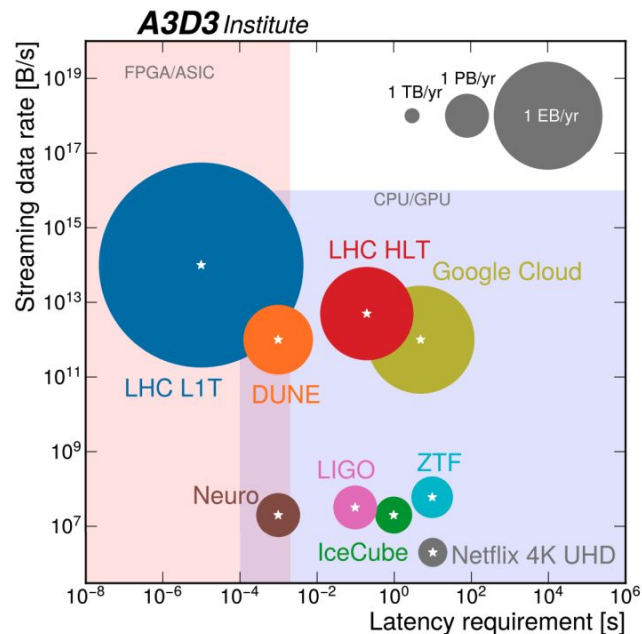- Nhan Tran (EAB)
- Mia Liu
- Javier Duarte

# Next generation of big data challenge

- The broader use of **AI/ML** in industry and academia is fueling rapid innovation in hardware accelerators.
- **High Energy Physics** at the LHC driving technology frontier
  - Both data size and streaming rates exceed those handled by industry leaders.
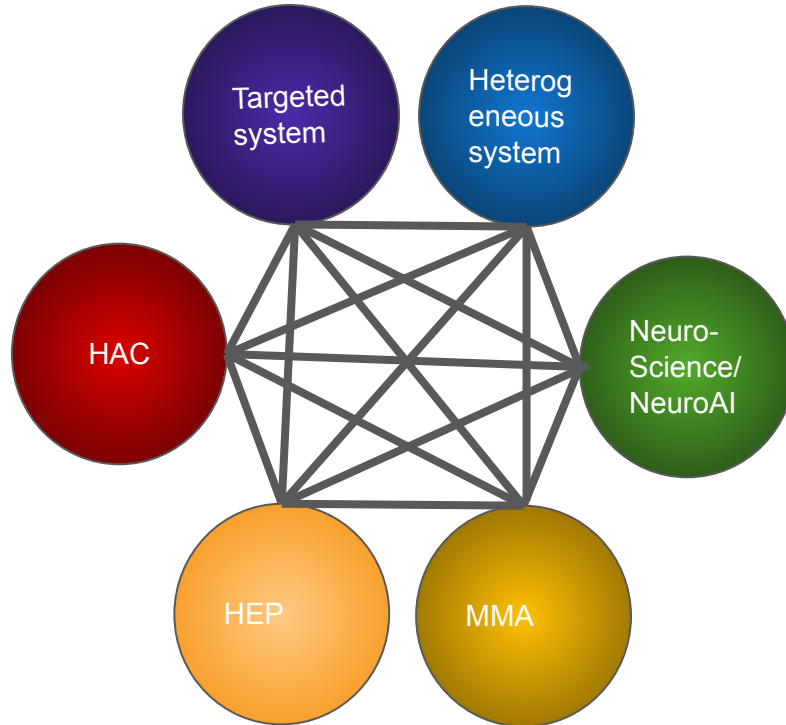
# Common challenge cross disciplinary

- **Multi-messenger Astrophysics** facilities rapidly increasing detection rates due to transformative network growth
- **Neuroscience** entering massive data analysis and interpretation thanks to neural recordings at scale
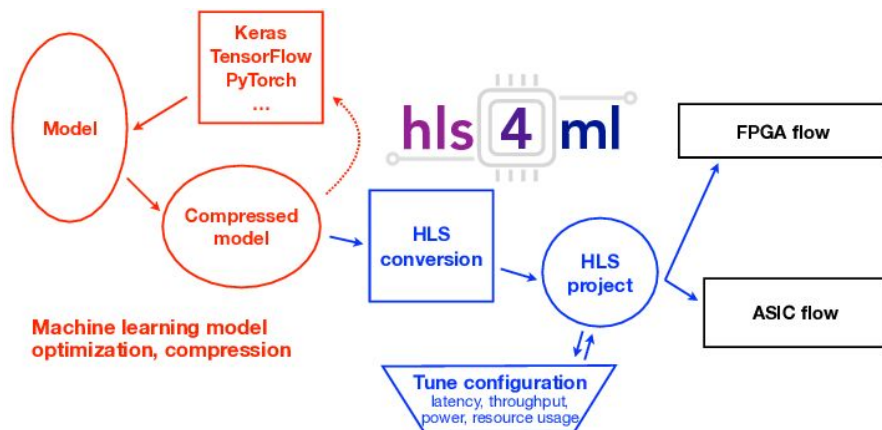
**Four** **focus areas** supported by core expertise for sustainability.

**Two** **Integrated** **systems** to facilitate integration and deployment.

Targeted system

Heterogeneous system

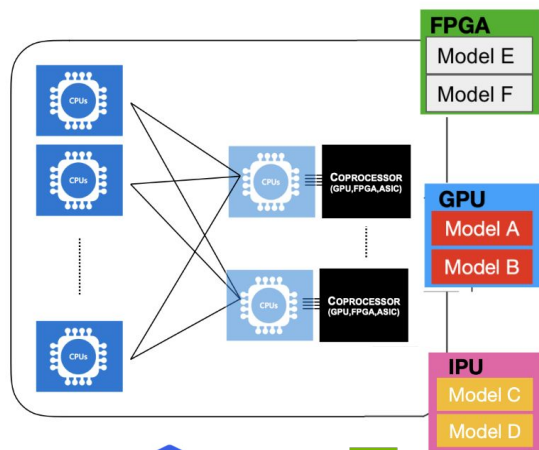Neuro-Science/ NeuroAI

HAC

HEP

MMA

# Targeted system for low latency/power

- [hls4ml](#): an open-source package enabling FPGAs & ASICs deployment of ML/AI algorithms
- A3D3 members are **core contributors and maintainers of package**, as well as **building a community of users**
    - AMD (FINN), TinyML, Imperial College London, University of Toronto, University of Zurich, CERN, FNAL, …, etc.
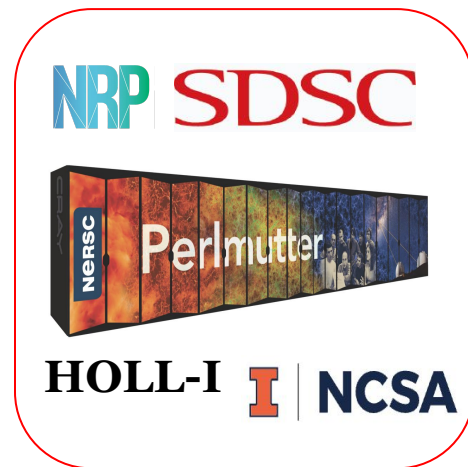
# Heterogeneous system for high throughput

- **ML as-a-Service** enabling users in sync with the most up-to-date AI model, and the inference server handling job execution in heterogeneous computing system.
  - A3D3 develops workflow platforms (SONIC, hermes) using standard industry tools and collaborates with IT Cloud providers & HPCs to evaluate performance
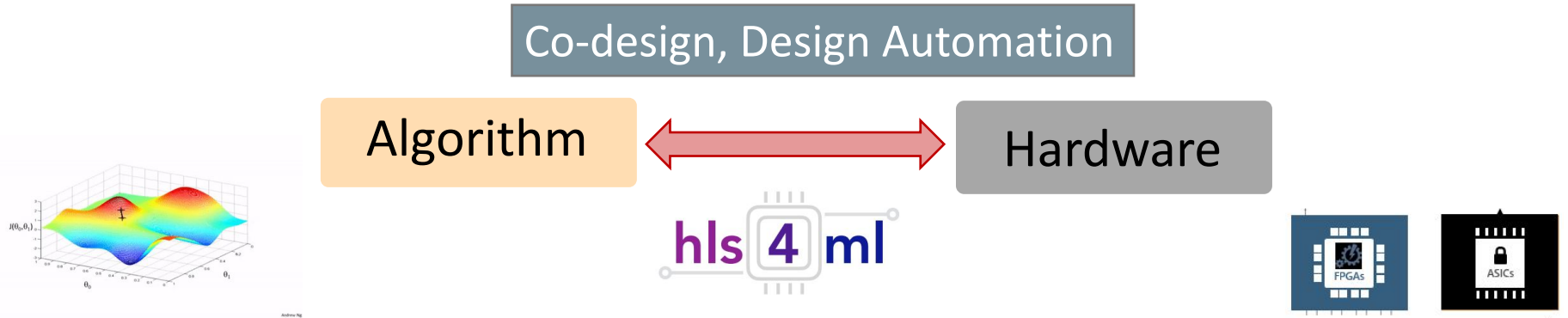


IT Cloud Providers

High Performance Computing

# Hardware-Algorithm Co-design (HAC)



Co-design, Design Automation

Algorithm ⟷ Hardware

hls 4 ml

**Challenges in Algorithm Design:**

- Irregular data (graphs, point clouds)
- Label scarsity
- AI models are hard to be interpreted
- …

**Challenges in Deployment in Hardware:**

- Computation efficiency issues (e.g. see Caroline Johnson's talk)
- Power/memory constraints
- Hard to be implemented on FPGA/ASIC

…

--> hardware design automation tools

# HAC: Innovative application

- New algorithms and hardware being prototyped with computational benchmark dataset and applied to domain science.

  - A3D3 researchers proactively seeks synergy cross different data



Torchsparse

Credit: Z. Liu

Torchsparse/ Torchsparse++ (Haotian Tang, et al. @ MLSys'22)

SPVCNN++ (Zhijian Liu et al . + HEP team)
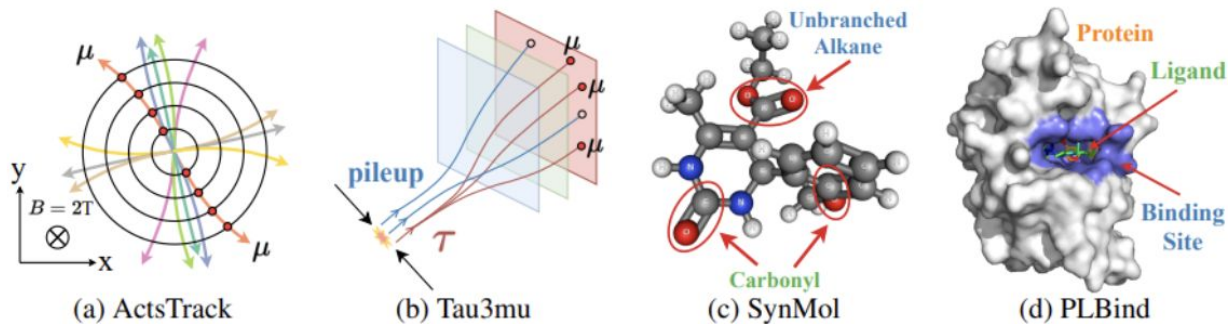
# HAC: ML Algorithms development

- **GSAT & LRI** (Siqi Miao, et al., @ ICML'22, ICLR'23)

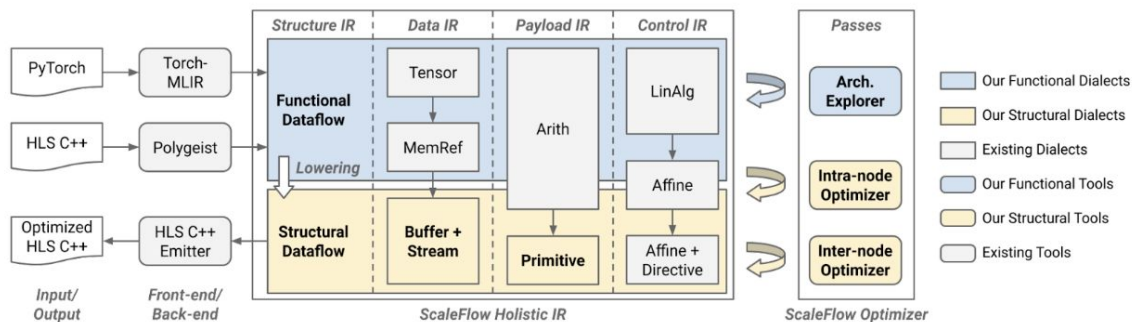   How to build interpretable and generalizable graph/geometric learning models?



A good model should capture the truly effective data patterns

(a) ActsTrack    (b) Tau3mu    (c) SynMol    (d) PLBind

- Theoretically grounded by the principle of information bottleneck
- Outperform baselines with a 10% improvement in detection accuracy of effective patterns and a 3% improvement in out-of-distribution generalization prediction accuracy
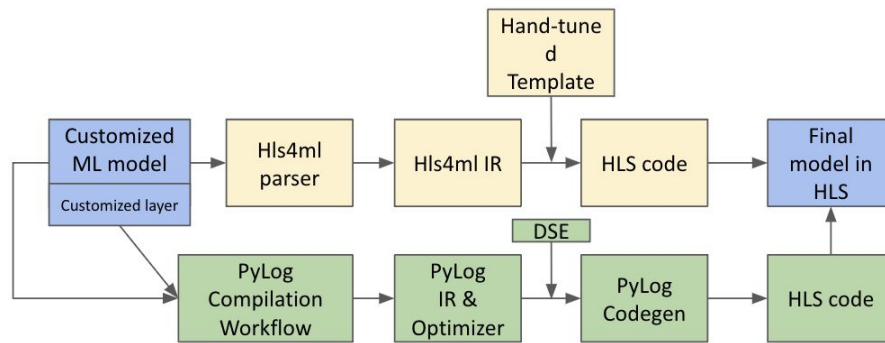
Credit: Pan Li

# Design Automation

- ### ScaleHLS / ScaleHLS 2.0 (Hanchen Ye, et al.)

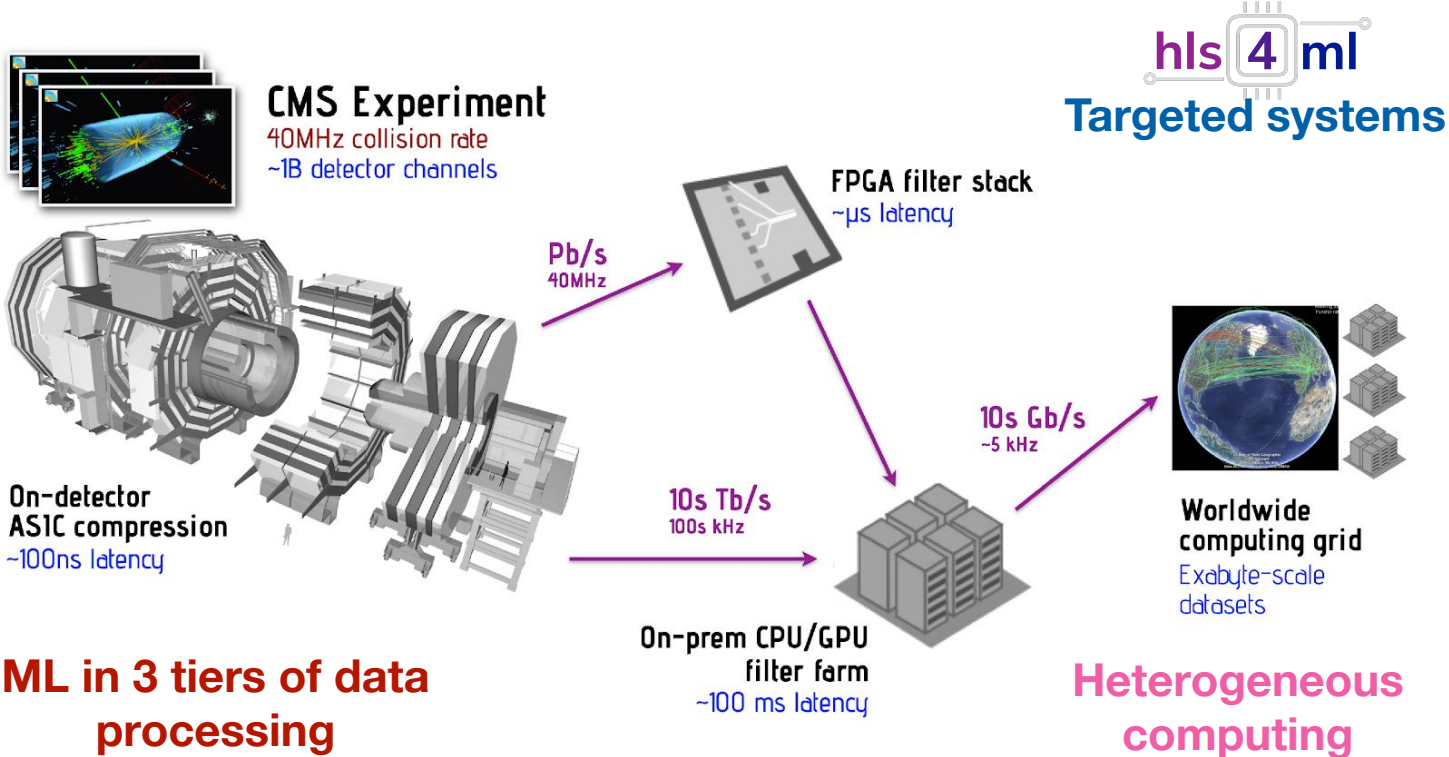  - generate highly-efficient hardware accelerators for scientific algorithms without much design effort



- ### PyLog + HLS4ML (Tim Zhang, et al.)

  - Integration of PyLog and HLS4ML enables significant code reduction in FPGA-oriented ML model development
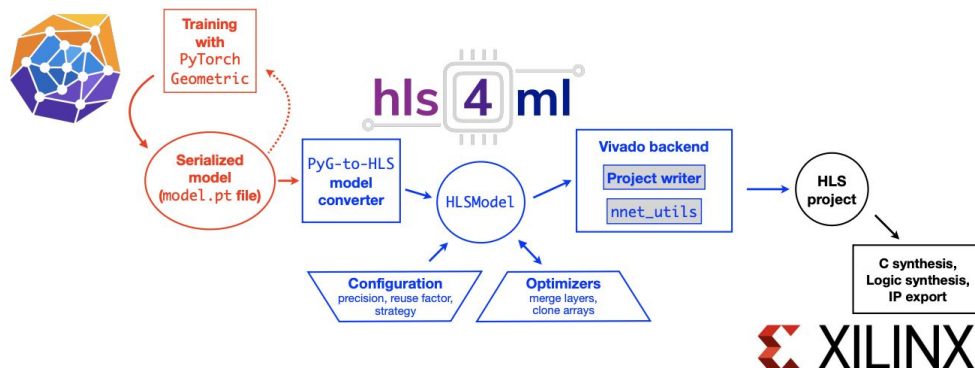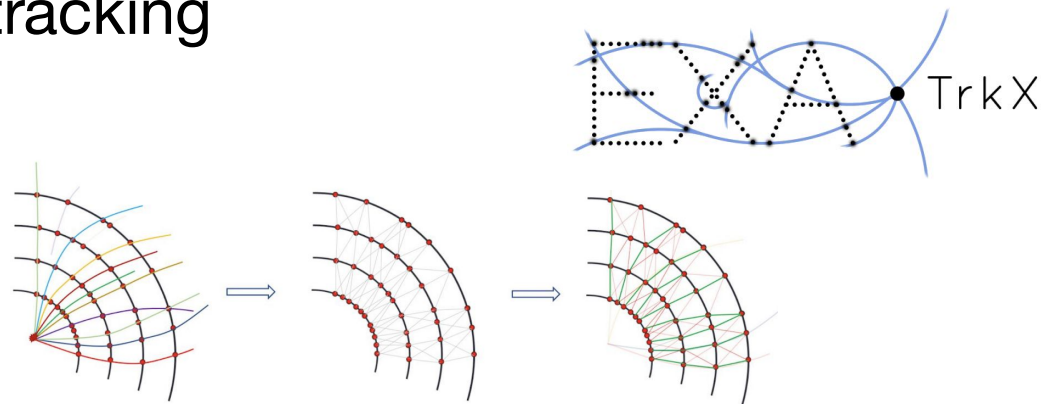


13

# High Energy Physics (HEP)



**hls 4 ml**
**Targeted systems**

**CMS Experiment**
40MHz collision rate
~1B detector channels

**FPGA filter stack**
~μs latency

Pb/s
40MHz

**On-detector ASIC compression**
~100ns latency

10s Tb/s
100s kHz

10s Gb/s
~5 kHz

**Worldwide computing grid**
Exabyte-scale datasets

**On-prem CPU/GPU filter farm**
~100 ms latency

**ML in 3 tiers of data processing**
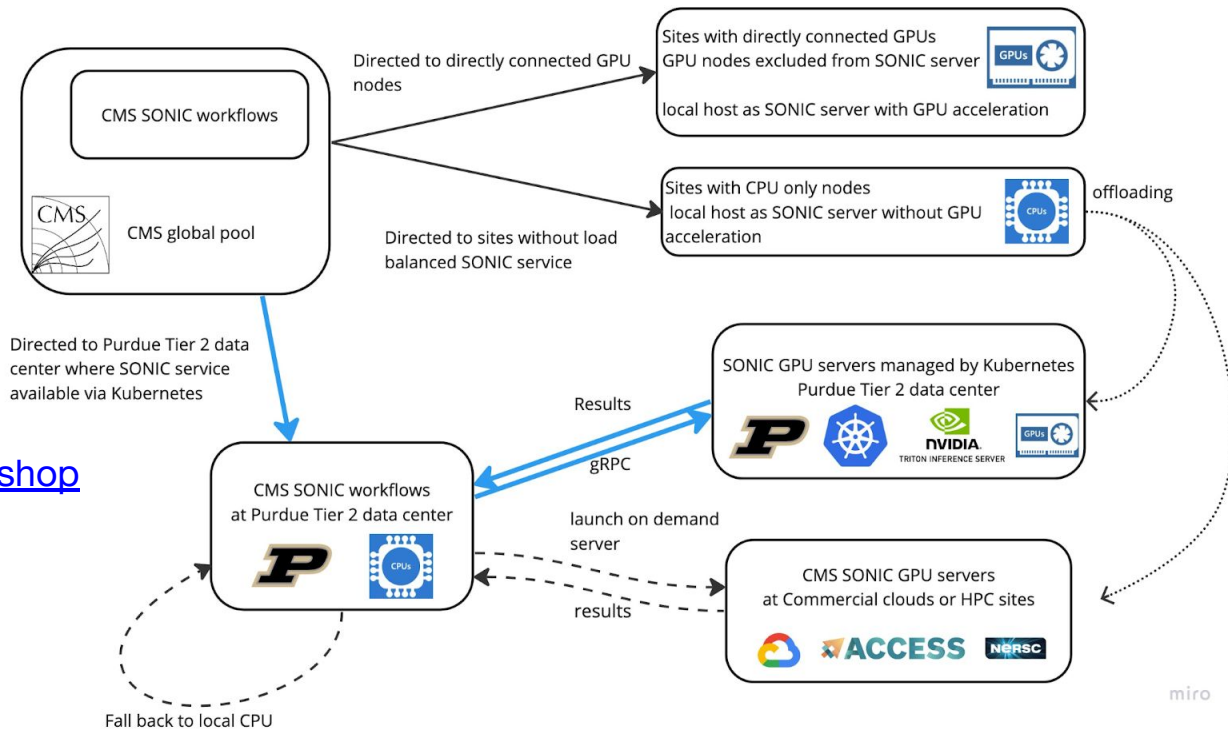
**Heterogeneous computing**

Credit M. Liu

# Graph Neural Network for tracking



- Algorithms making tracking highly parallelizable both low latency FPGA version and GPU version
    - Front. Big Data 5 (2022) 828666
    - 2306.11330
- Can be used at various tiers of track reconstruction
    - ExaTrkX as a service CTD2023

# Heterogeneous computing as-a-service (SONIC)

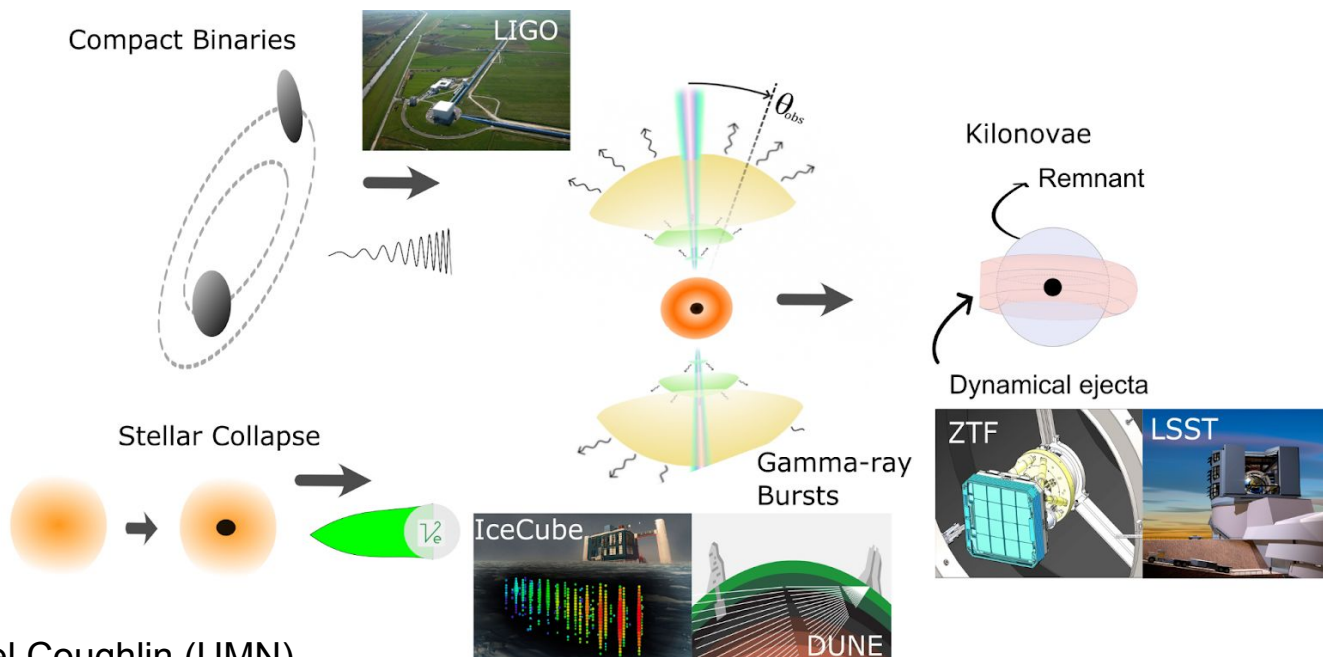Significant progress in integration of SONIC in CMS for minAOD production



[Talk at fast ml workshop](#)

[CHEP 2023](#)

# Multi-messenger Astrophysics

- Develop and deploy software within astronomical facilities to enable discovery
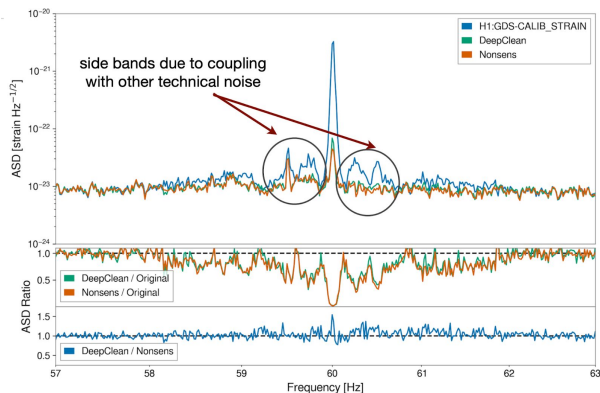


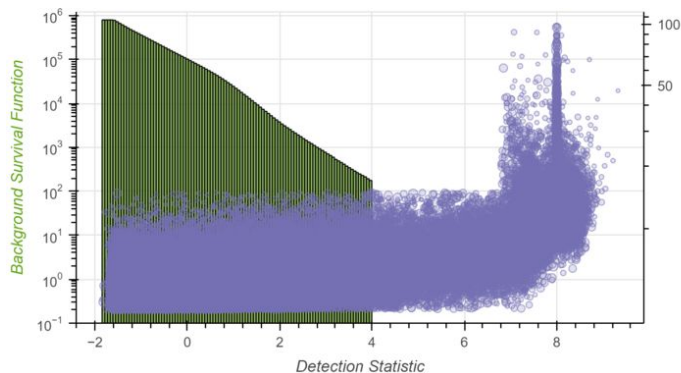Credit: Michael Coughlin (UMN)

# Gravitational Waves (LVK)

Github: <u>ML4GW</u>

All algorithms use our <u>inference-as-a-service</u> (IaaS) prototype to implement a real-time noise subtraction pipeline (DeepClean), detection (aframe/GWAK), and parameter estimation for use during the fourth observing run (O4) of LIGO-Virgo-KAGRA on dedicated hardware at the detector sites.
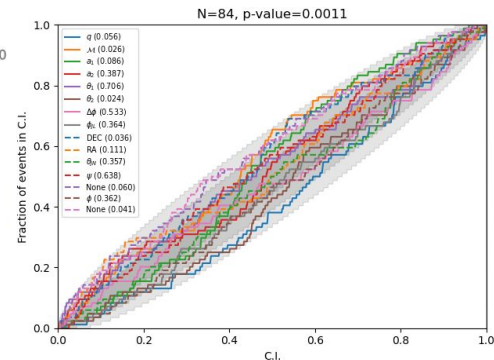
**Clean the Data: DeepClean (CNN)**
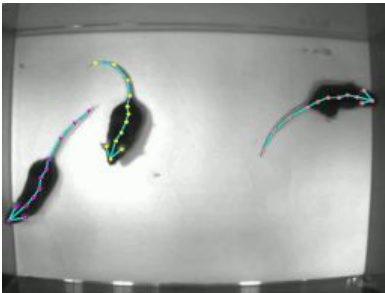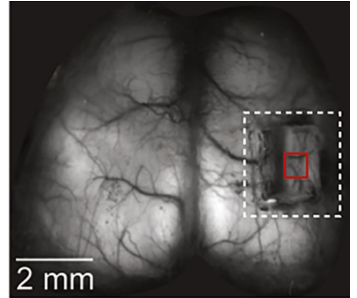
**Detect the GWs: aframe (CNN)/GWAK (autoencoders)**

**Characterize the GWs: (MAF*)**



18

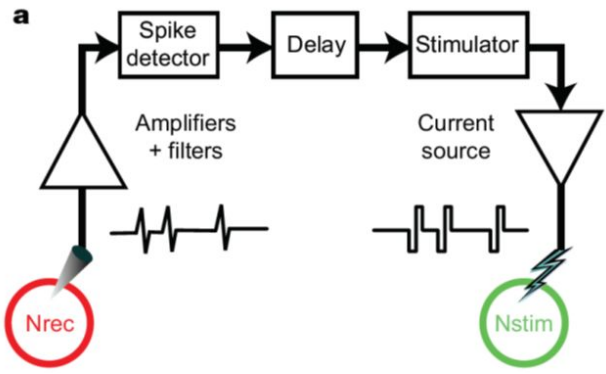# Neuroscience needs high-throughput & real-time AI

**Rapid increase in number, type of measurements**



**Brain**

**Behavior**

**Must *perturb* the system to disentangle causality, treat disorders.**



a

Spike detector → Delay → Stimulator

Amplifiers + filters

Current source

Nrec

Nstim

**Need: data-driven discovery of relevant features, structure in data**

**Need: low-latency algorithms (<1ms)**

# Improved time-series reconstruction methods

- Developed new Multi-block Recurrent Auto-Encoder (MRAE) to increase bandwidth more efficiently

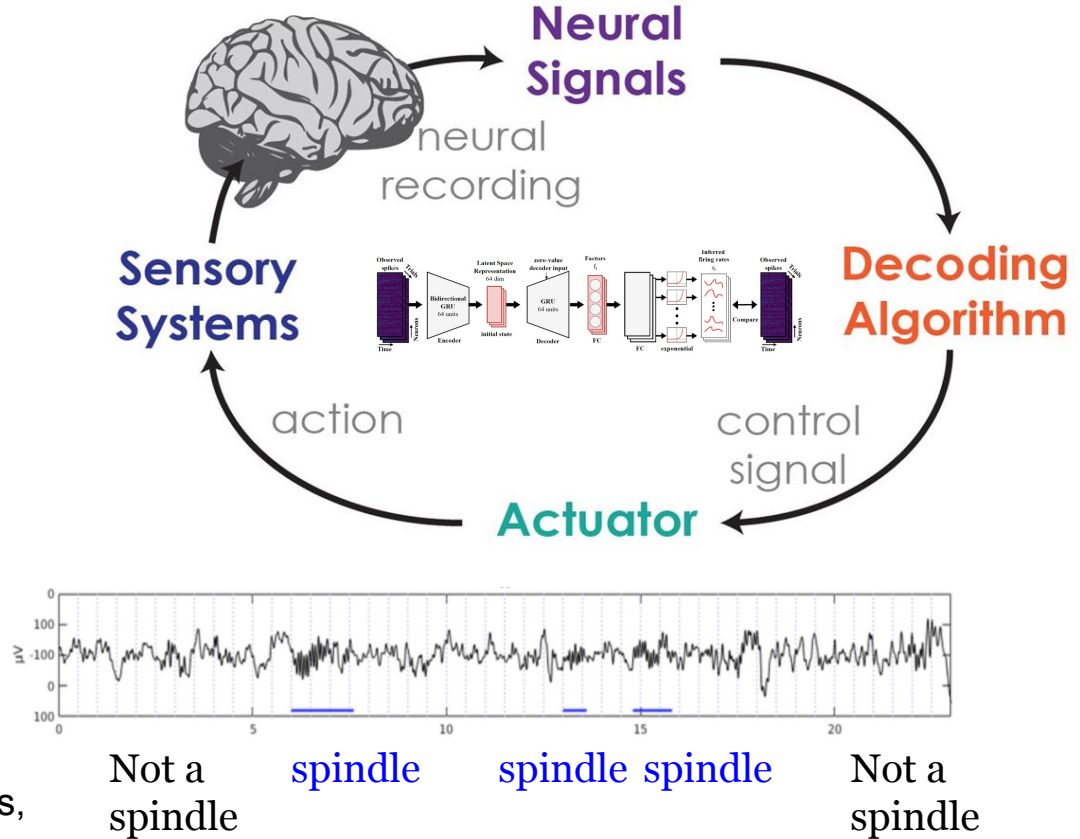- Developed Spatio-Temporal Transformer for Spiking Neural Data

Nolan, Pesaran, Shlizerman & Orsborn, *bioarxiv 2022*
*Le & Shlizerman, NeurIPS 2022*
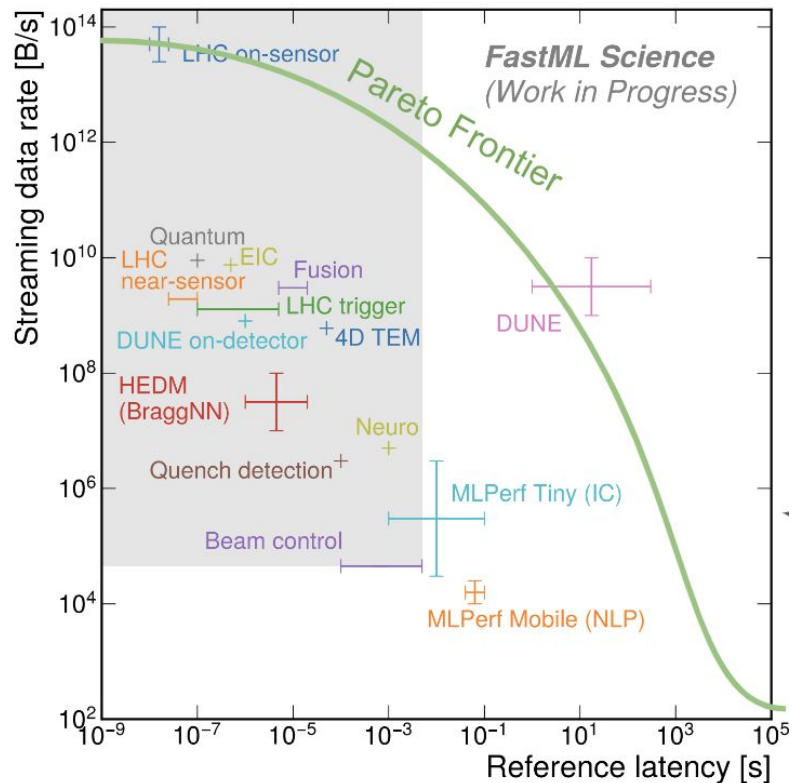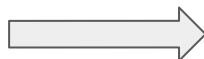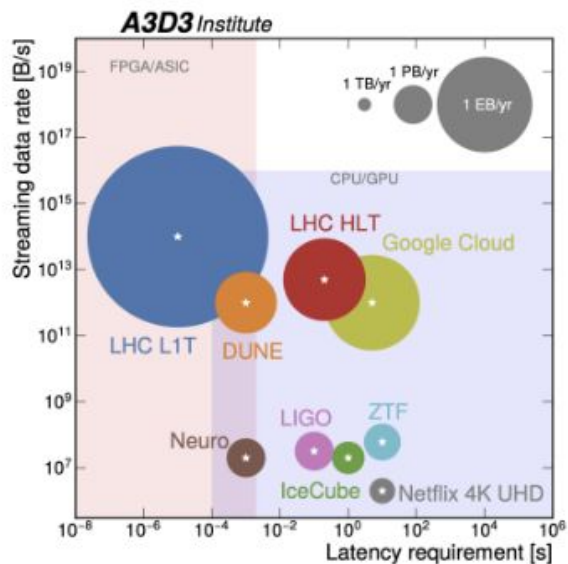
# NeuroAI Integration

- A popular autoencoder model used on neural data (LFADS) in FPGA, Elham Khoda's talk

- Neuro A3D3 develops methods for reconstruction, forecasting and clustering of time-series
- Potential applications/uses:
  - Detect noise and artifacts
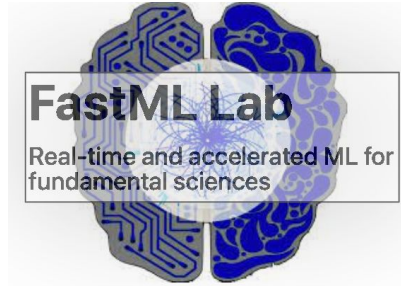  - Detect rare neural events of interest (e.g., seizures, spindles, etc)



Not a spindle    spindle    spindle spindle    Not a spindle

# Fast Machine Learning Community

Our aim is to build a large-scale public scheme to advertise this work

# Partnership and FastML Ecosystem

*Growing strong industry connections with support through the Fast ML community*



**Partner Projects**

FAIR4HEP

EXA·TrkX

sPHENIX

AstroAI

ats

OzGrav

+ Many more...

**Experiments**

ATLAS EXPERIMENT

ICECUBE NEUTRINO OBSERVATORY

CMS

DUNE DEEP UNDERGROUND NEUTRINO EXPERIMENT

ZTF

LIGO VIRGO KAGRA

VERA C. RUBIN OBSERVATORY

KITT PEAK NATIONAL OBSERVATORY

**National & Int'l Laboratories**

Fermilab

Brookhaven National Laboratory

CERN

Argonne NATIONAL LABORATORY

BERKELEY LAB

Los Alamos NATIONAL LABORATORY

**Coprocessors**

AMD, XILINX, intel

SambaNova SYSTEMS

NVIDIA

GRAPHCORE

habana An Intel Company

Cerebras

UNTETHER AI

**IT Cloud Providers**

Microsoft Azure

aws

Google Cloud Platform

CrusoeCloud

**High Performance Computing**

SDSC

NRP

NERSC Perlmutter

NCSA

DELTA

# A3D3 Ecosystem & Engagement

- *High-Throughput AI Methods and Infrastructure Workshop*

- *Postbaccalaureate Workshop*





**Fast Machine Learning for Science**

Real-time and accelerated ML for fundamental sciences

**Imperial College London**

**25-28 September 2023**

**Scientific Committee**
Thea Århestad (ETH Zurich)
Javier Duarte (UCSD)
Phil Harris (MIT)
Burt Holzman (Fermilab)
Scott Hauck (U. Washington)
Shih-Chieh Hsu (U. Washington)
Sergo Jindariani (Fermilab)
Mia Liu (Purdue University)
Allison McCarn Deiana (Southern Methodist University)
Mark Neubauer (U. Illinois Urbana-Champaign)
Jennifer Ngadiuba (Fermilab)
Maurizio Pierini (CERN)
Sioni Summers (CERN)
Alex Tapper (Imperial College)
Nhan Tran (Fermilab)

**Organising Committee**
Sunita Aubeeluck
Robert Bainbridge
David Colling
Patrick Dunne
Wayne Luk
Andrew Rose
Sioni Summers (co-chair)
Alex Tapper (co-chair)
Yoshi Uchida
Ioannis Xiotidis

indi.to/fastml23
fastmachinelearning.org

# Summary

- A3D3 focusing on accelerating real-time AI to solve common challenges through interdisciplinary collaboration
  - 4 focus areas: HAC, HEP, MMA, Neuros
  - 2 integrated systems: Targeted system, Hetereogenous computing
- A3D3 is closely connected with the FastML Community
  - Leverage our leadership in FastML to connect to main different domains
  - Touches on many fields in industry/science not part of A3D3 scope
    - Plasma Physics/Materials Science/…/ASIC design
- Welcome to participate in A3D3 activities
  - HDR Ecosystem Workshops
  - Postbac Program Enhancements
  - Machine Learning Challenges
    - Nov 17 planning meeting https://indico.cern.ch/event/1342015/

Shih-Chieh Hsu

http://faculty.washington.edu/schsu/

schsu@uw.edu

# Cross-discipline



**HEP**

Hsu
PI

Harris
co-PI

Neubauer
co-PI

Liu

Duarte

**MMA**

Coughlin
co-PI

Scholberg
co-PI

Graham

Hanson

Katsavounidis

**Neuros**

Orsborn

Shlizerman

Dadarlat Makin

**CS/EE**

Hauck

Li

Chen

Han

**17** Senior Personal

# A3D3 fully staffed

**106** Members (including 5 affiliate)



HS
0.9%

Undergraduate
15.1%

Master
7.5%

PhD
35.8%

Professor
5.7%

Assoc. Prof
2.8%

Assist. Prof
9.4%

Scientist
6.6%

Engineer/Specialist
1.9%

Postbacc
3.8%

Postdoc
10.4%

**74% trainee**

**Prog. Ope. Spec.**

Zhang (UW)

**Affiliate faculty/staff**

Rankin
(Upenn)
A3D3 Alumni

Sravan
(Drexel)
A3D3 Alumni

Ju
(LBNL)

Lai
(NYCU)

Carlson
(Westmont)

**PostBacc fellow**

Gray
(UMN)

Peterson
(UW)

Lian
(Duke)

Skivington
(UCSD)

# ML Challenge: Unifying across domains

- challenge across HDR domains

  - Try to find anomalies over many different datasets with one metric

Would be a FAIR workflow challenge?
Could extend this to semi/self-supervised learning (foundation models)
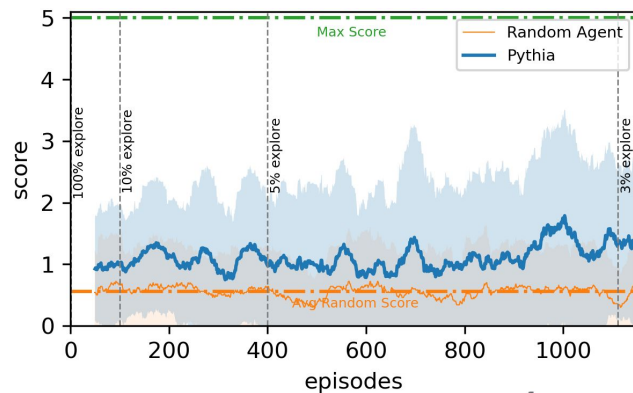


**Many Datasets covering whole HDR**

**Anomaly Algorithm**

**Anomaly Metric**

Dot product the input and output
Large Value : Good
Small Value : Anomaly

# Optical Astronomy - Overview

**Simulate Observations: NMMA (emulator)**



Github work areas:

NMMA SCOPE Pythia

**Optimize Observations: Pythia (RL)**



~4 faculty, 3 postdocs, 5 grad students, 3 postbac/undergraduat

**Classify the sources: Scope (CNN)**



**Main focus**: Deploy ML algorithms throughout the observation preparation and follow-up for source identification and characterization
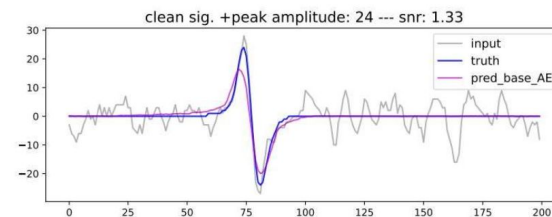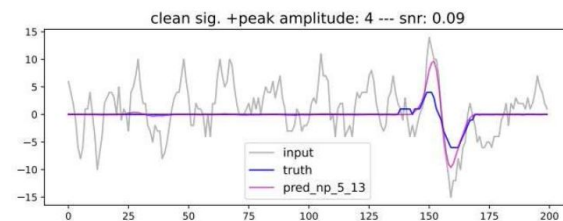
# Neutrinos - Overview

**PMT Voltage Picking (CNN)**



~2 faculty, 2 postdocs, 2 grad students, 2 postbac/undergraduates

**Supernova Reconstruction (1DCNN autoencoder + pointing)**



**Main focus**: Porting existing algorithms to GPUs and FPGAs for the purpose of detection and localization reconstruction.
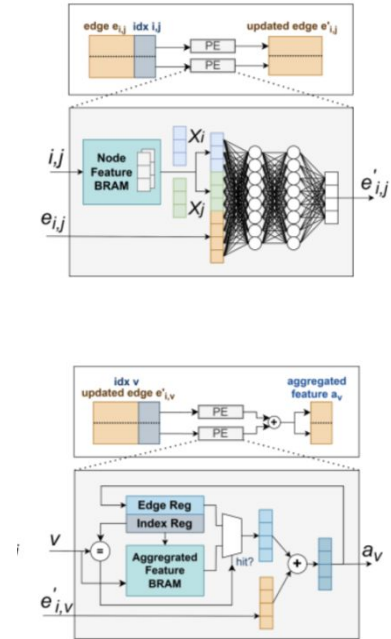
*See: See Pan's Talk in Hardware-Algorithm Co-Development*

# LOW LATENCY EDGE CLASSIFICATION GNN

Shi-Yu Huang, Yun-Chen Yang, Yu-Ru Si, et. al. FPL 2023

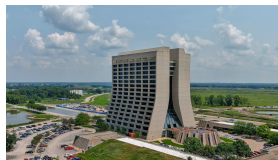Modularized parallel architecture for each computational pipelines



**Achieving 2.07 us Latency with 3.225 Throughput (MGPS)**
- Xilinx Virtex UltraScale+ VU9P  HLS 2019.2

# National Lab: HLS4ML for Analog AI

- Project: "*Democratizing AI Hardware with an Open Source, Automated AI-Chip Design Toolkit* "

- Joint initiative with Discovery Partners Institute and Fermilab
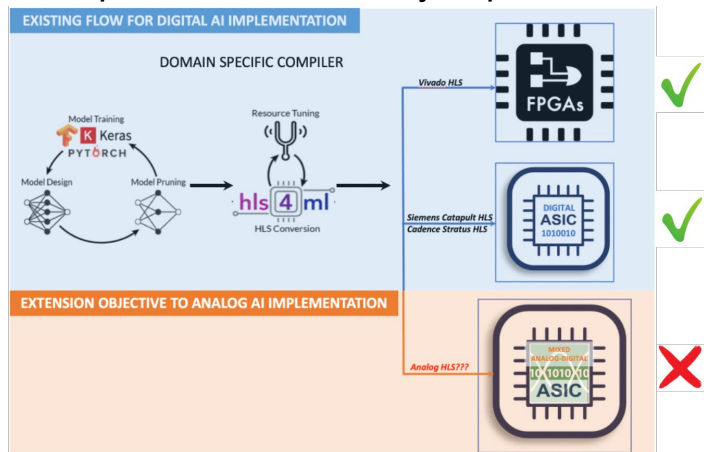


**Why Analog AI?**
More efficient, Better Latency, Less Area

**Why Automate Analog AI?**
Cheaper, faster, less risky implementation



Farah Fahim
Fermi Lab,
ASIC Research &
Development Head

Ben Parpillon
Fermi Lab,
Senior ASIC
Engineer

**AI-Chip Prototyping and Analog Primitive Automation**

Amit R. Trivedi
UIC,
Electrical and
Computer Engineering

Nhan Tran
Fermi Lab,
Accelerator-based
Experiments

**High-Level Synthesis and Digital Automation Flow**

Ahmet Cetin
UIC,
Electrical and
Computer Engineering

Mark Neubauer
UIUC,
High Energy
Physics

**Application Studies: Low Barrier Custom-AI for Small Businesses**
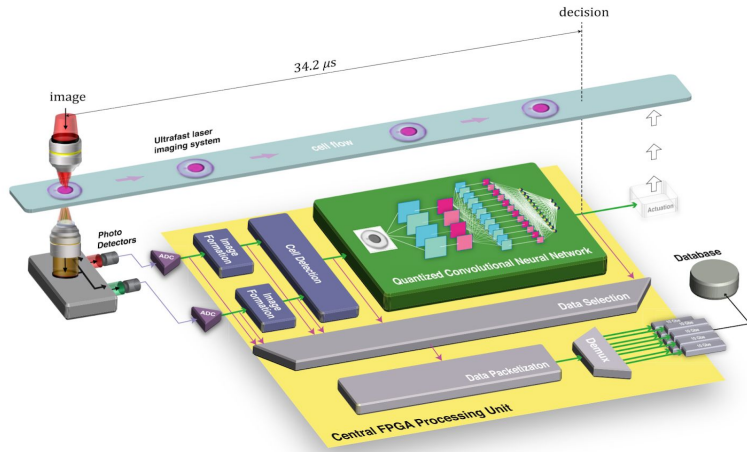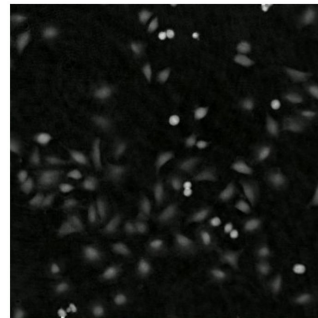
# Industry: Real-time  Blood Cell Id



Diagram from: ieee paper

- **Collaboration between MIT, CERN and Phiab**
  - Led/initiated by Vladimir Loncar
- Working to bring HLS4ML to cell identification
  - Working directly with industry to deploy
  - Builds on A3D3 AI initiatives

- Collaboration with https://phiab.com/

- Key Ideas
  - Real time tagging of blood cells
  - Can be used for cell therapy
    - Cancers/….
  - Non-invasive
    - No chemicals
    - All electronics based

Original holography info        Segmented cell instances