# Accelerating Artificial Intelligence for High Energy Physics

Shih-Chieh Hsu (徐士傑）

University of Washington

PHY-2117997

IAS HEP 2024 (https://indico.cern.ch/event/1335278/)

HKUST Jockey Club Institute for Advanced Study

Jan 24 2024

https://a3d3.ai/

## P5 Report (Draft Dec 2023)

https://www.usparticlephysics.org/2023-p5-report/
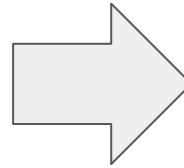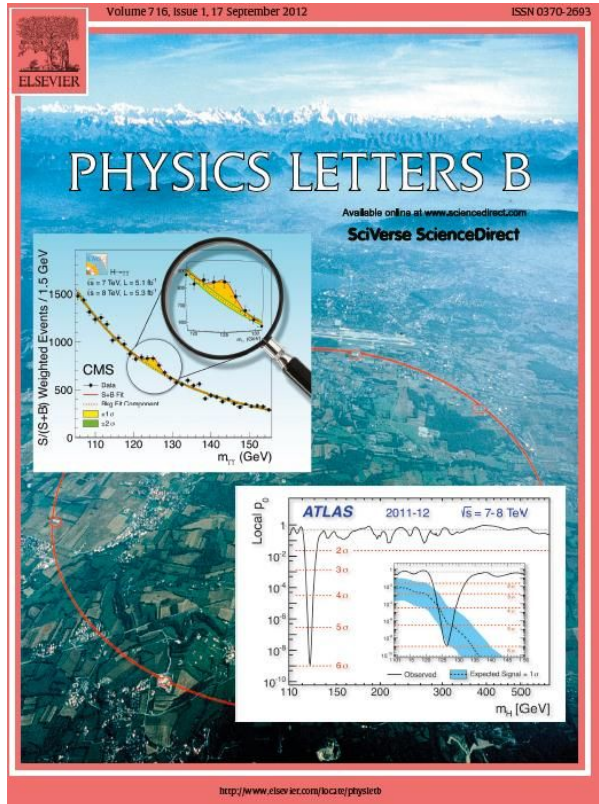
**Investing in the scientific workforce and enhancing computational and technological infrastructure is crucial.** To achieve this goal, funding agencies should support programs that foster a supportive, collaborative work environment; help recruit and retain diverse talent; and reinforce professional standards. Targeted increases in support for theory, general accelerator R&D (GARD), instrumentation, and computing will bolster areas where US leadership has begun to erode. These areas align with national initiatives in **artificial intelligence and machine learning (AI/ML)**, quantum information science (QIS), and microelectronics, creating valuable synergies. Such increased support maximizes the return on scientific investments, fosters innovation, and benefits society in domains from medicine to national security.

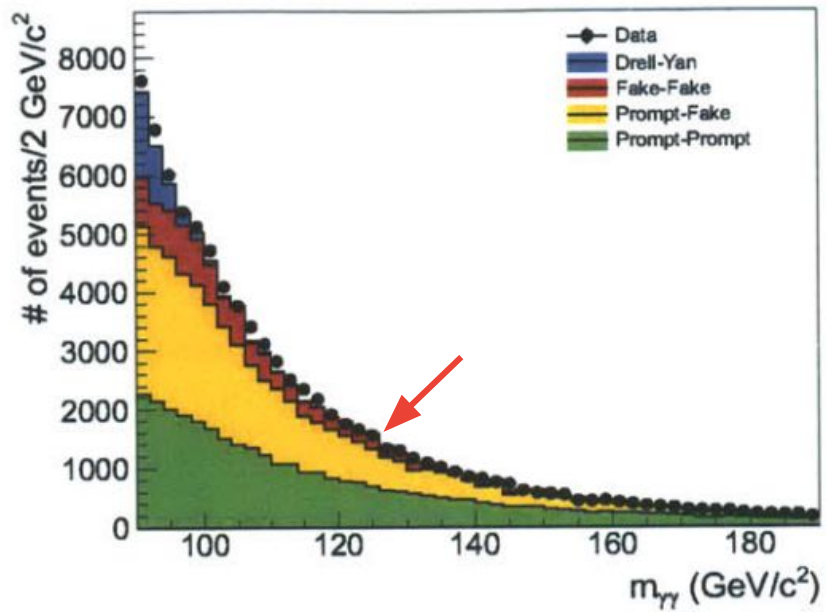# AI/ML has made critical contributions to the Higgs Discovery!



© Nobel Media AB. Photo: A. Mahmoud
François Englert

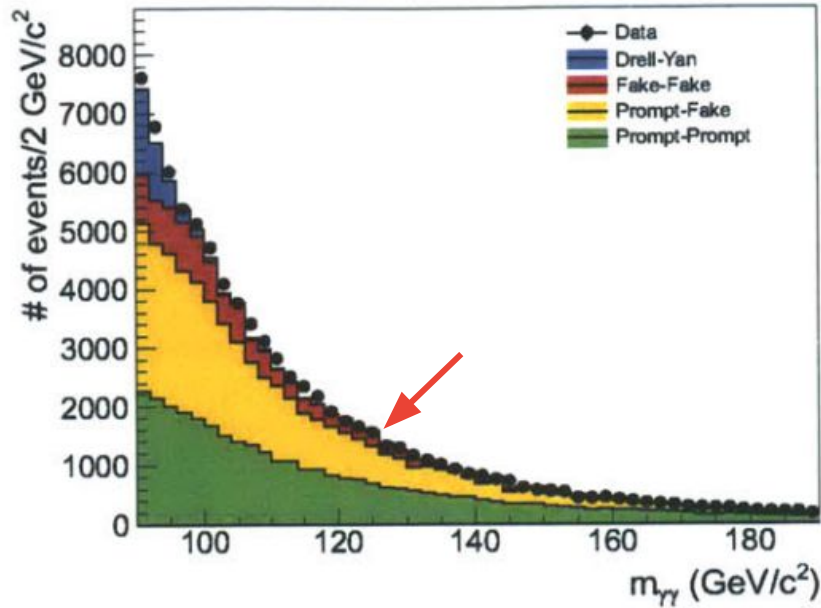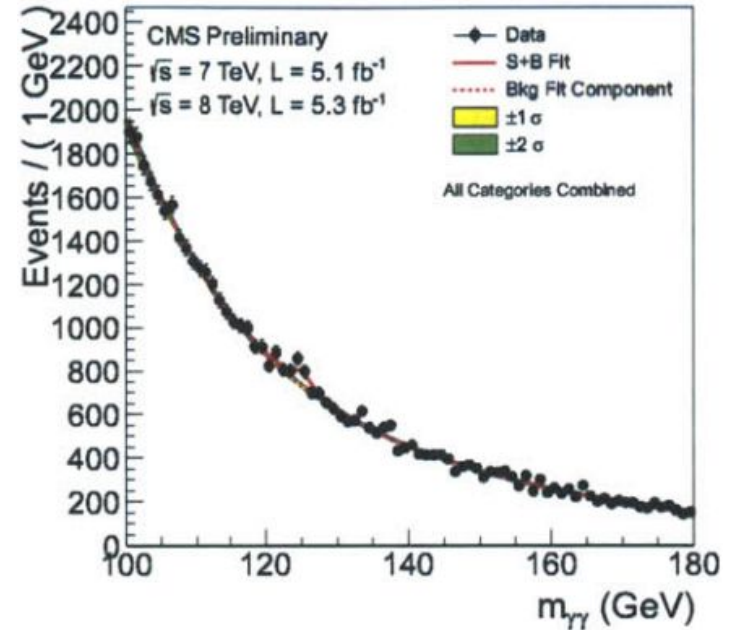© Nobel Media AB. Photo: A. Mahmoud
Peter W. Higgs

2013

J.L. Bendavid, THESIS-2013-079

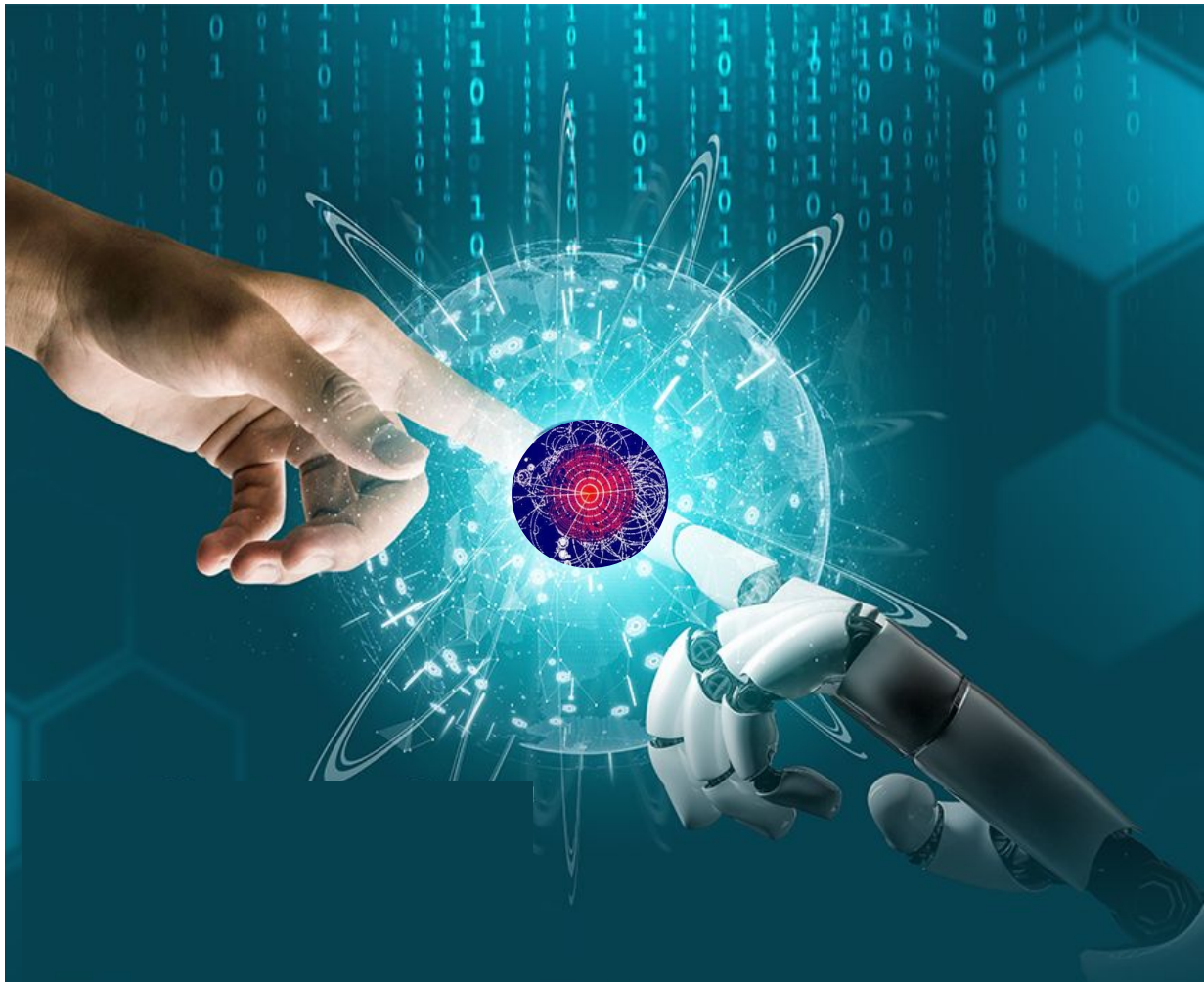**Key for discovery**

- Optimizing **signal-to-background** ratio



BDT

# 2012: A Breakthrough Year for Deep Learning



AlexNet Comm. ACM. 60 (6): 84–90



ACM 2018 Turing Award

# Exponential trend of computation need for AI



Training FLOPs Scaling for SOTA Models

7

8

Credit: Onpassive

Peak luminosity • Integrated luminosity — DATA VOLUME
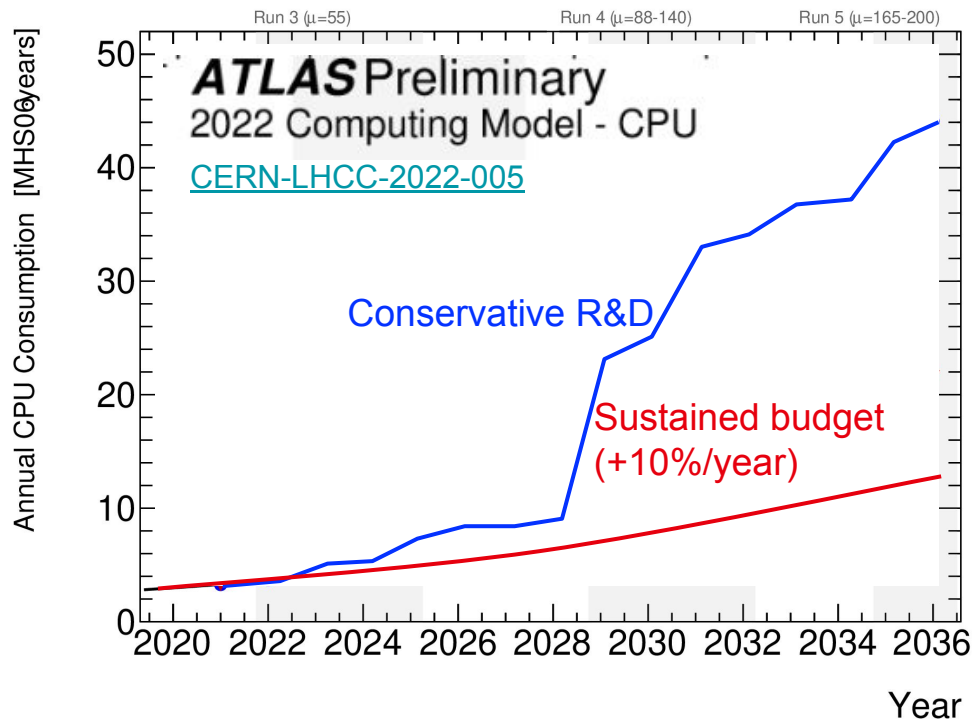
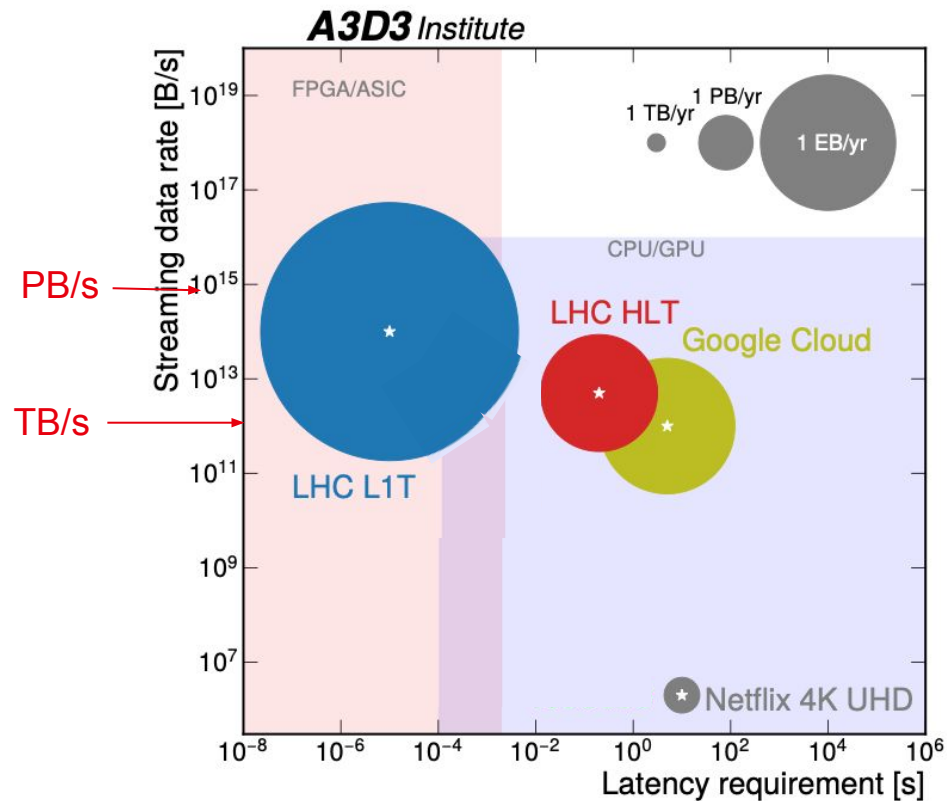# Critical computing challenge



- To preserve current physics we are upgrading the system
  - We will have to take data at 4 times the current rate
  - Our event size will have to be 10x larger

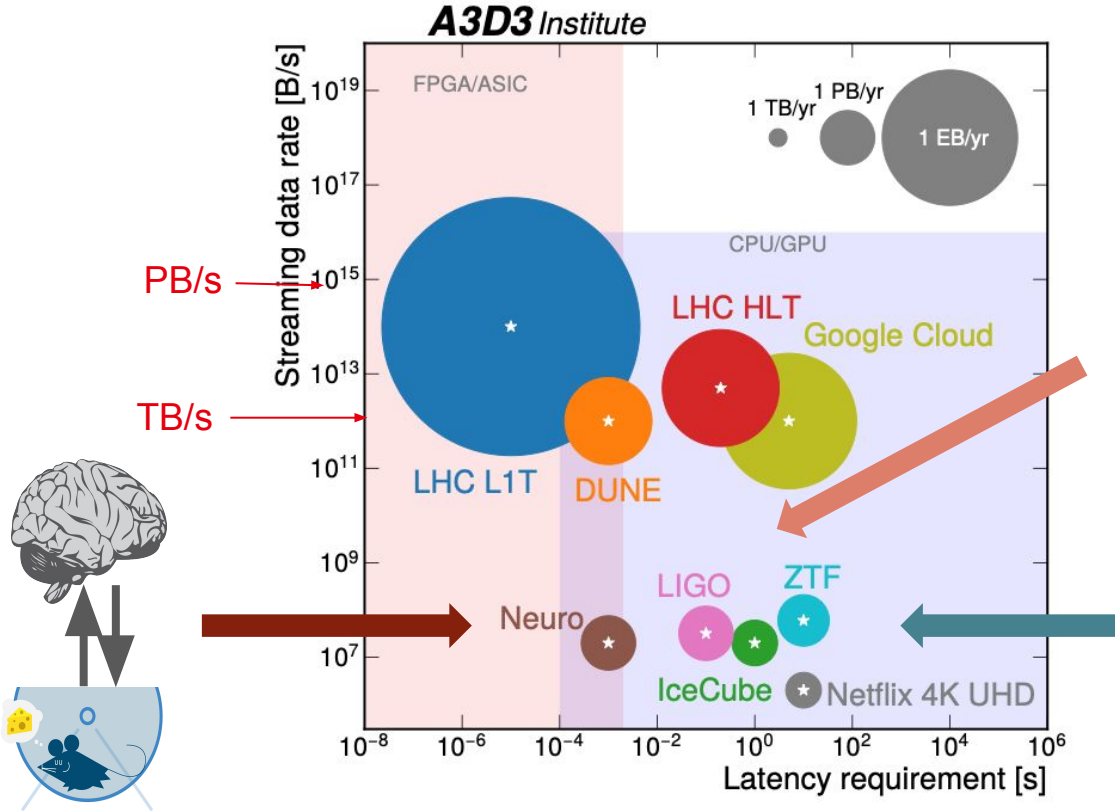# Critical computing challenge



- To preserve current physics we are upgrading the system
  - Our event size will have to be 10x larger
  - We will have to take data at 4 times the current rate
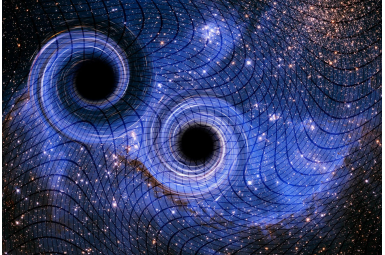- However, we are lacking of sufficient budget to sustain required computing
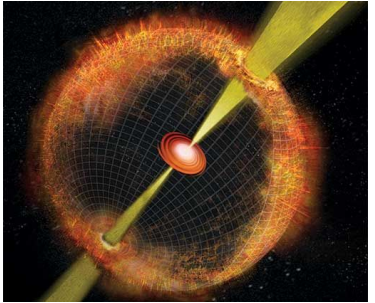
11

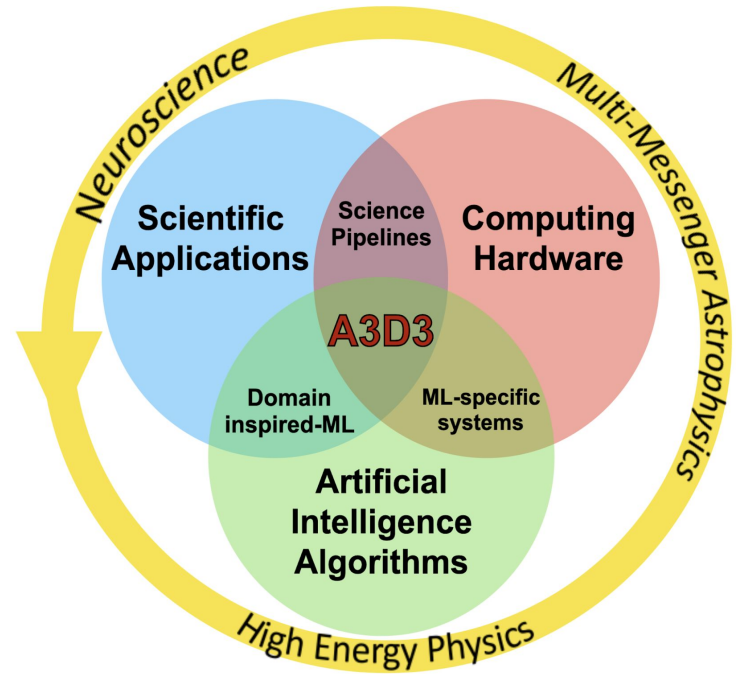# Critical computing challenges

# Common big data challenges

# NSF HDR Institute **A3D3**

***A**ccelerated **A**rtificial Intelligence **A**lgorithms for **D**ata-**D**riven **D**iscovery*

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

# Cross-institution

**16** institutions
**104** members

# Cross-discipline

**HEP**

Hsu
PI/Director

Harris
co-PI

Neubauer
co-PI

Liu

Duarte

Rankin

Aarastad

Gonski
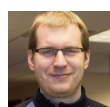
Carlsen

**MMA**

Coughlin
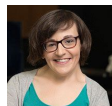co-PI

Scholberg
co-PI

Graham

Riedel

Katsavounidis

Li

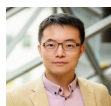Sravan

**Neuros**

Orsborn

Shlizerman

Dadarlat Makin

Sun
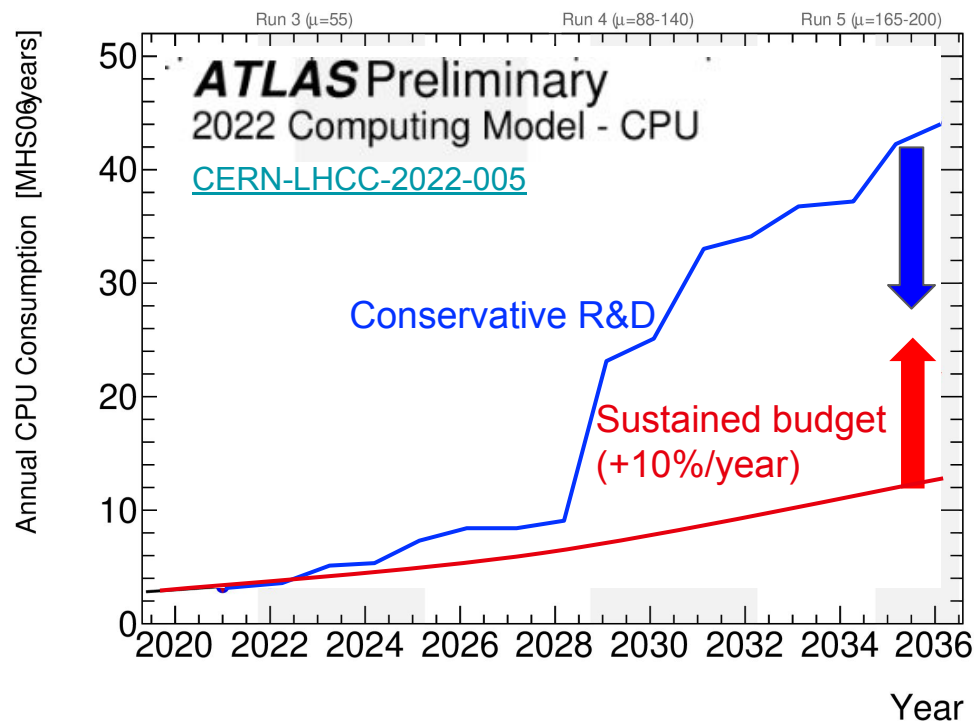
**CS/EE**

Hauck

Li

Chen

Han

Ju

Lai

**17** Senior Personal

**9** Affiliates
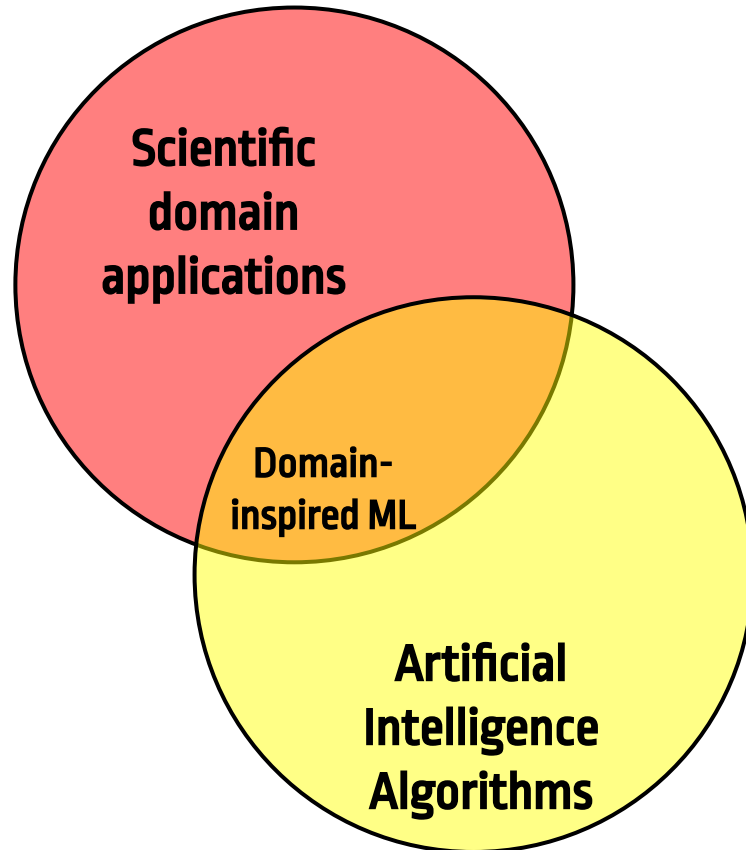
**10** Postdocs

**58** Graduates
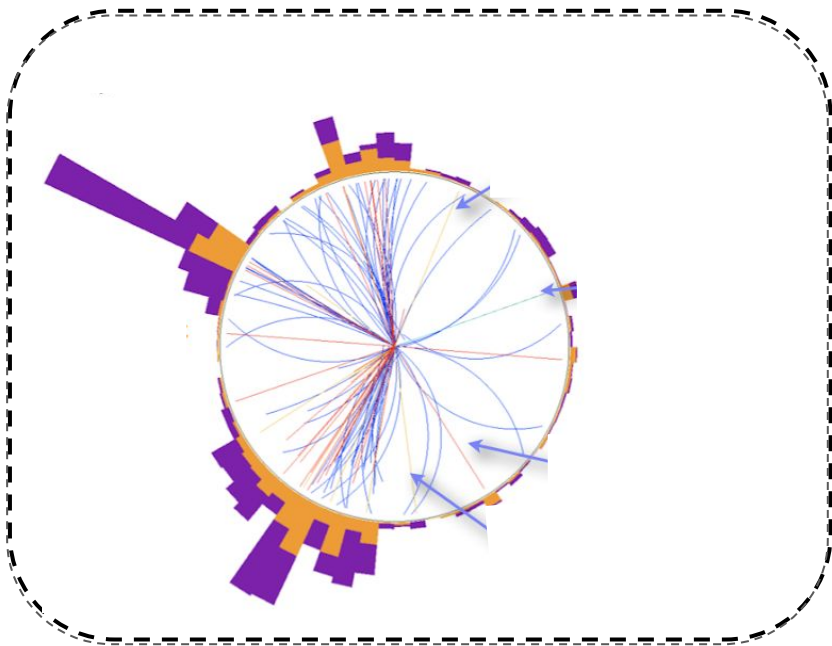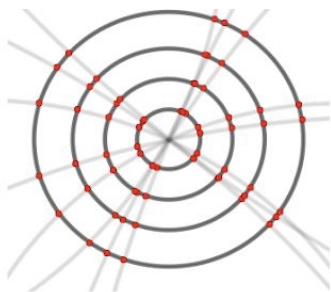
**10** Undergraduates

# Closing the gap
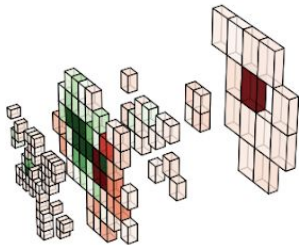


Smarter Algorithms - AI

Faster Hardware - Co-processor

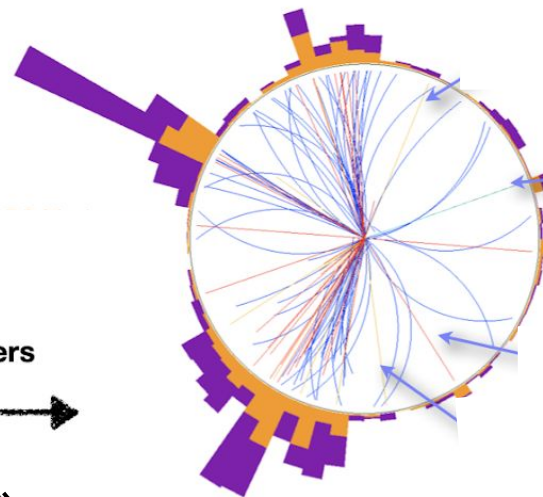# Smarter Algorithms
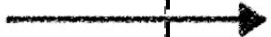
**Connecting the dots**

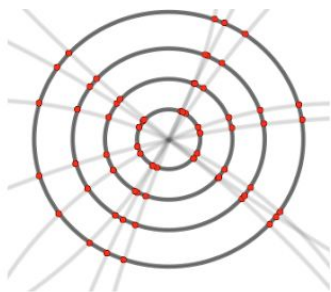**Energy Clusters**

**Charged particle tracks**

HCAL
deposit

**Energy clusters**

**Connecting the dots**

Charged particle tracks

HCAL deposit

Energy clusters

**Energy Clusters**
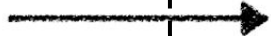
**Particle Flow Reconstruction**

Higgs?
Top Quark Pa
W boson?
Z boson?
Multiple boso
.
.
.
New Physics

**End-to-end**

# Track Reconstruction as Graph

IASHEP: GNN tracking

X. Ju, et. al. EPJC 81, 876 (2021)

# Track Reconstruction as Graph

TrkX

Graph Neural Network to identify correct edge connecting adjacent nodes

# Clustering with Sparse Point Voxel Convolutional Neural Network

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

  **2.9X** faster than MinkowskiEngine (NVIDIA)
  **1.8X** faster than SpConv (TuSimple).

**Energy Clusters**

# Clustering with Sparse Point Voxel Convolutional Neural Network

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

  **2.9X** faster than MinkowskiEngine (NVIDIA)
  **1.8X** faster than SpConv (TuSimple).

Particles are a set of 3D points and can be processed by our efficient 3D algorithms.

**4% higher** mIoU and **10+% higher** PQ



SPVCNN++ (Ours)    Groundtruth

25

Industry Focus

CPUs     GPUs     FPGAs     ASICs

FLEXIBILITY ← → EFFICIENCY

# The Need for the FastML



40 MHz

1 ns

latency c

# Heterogeneous Computing

**40 MHz**

ASIC

**1 ns**

latency

# The Need for the FastML

# The Need for the FastML



100 KHz

CPU Cluster

1 KHz

FPGA

1 MB/ev

40 MHz

High-Level trigger

L1 trigger

ASIC

1 ns    1 μs    100 ms

latency constraint

latency AND throughput constraint

# The Need for the FastML

# The Need for the FastML



100 KHz

40 MHz

1 KHz

1 MB/evt

FPGA

ASIC

L1 trigger

CPU Cluster

High-Level trigger

CPU GRID

Offline reconstruction

Computing time

Fast Machine Learning Collaboration

Core
(HLS4ML)

Co-processor
(SONIC)

Innovative Algorithms

# HAC Research Focus

Co-design, Design Automation

Algorithm ⟷ Hardware

hls 4 ml

FPGAs    ASICs

**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
- Power/memory constraints
- Hard to be implemented on FPGA/ASIC

…

--> hardware design automation tools

34

# HLS4ML translating ML into FPGA firmware

# HLS4ML translating ML into FPGA firmware

# HLS4ML translating ML into FPGA firmware

# Quantization

▶ Scan the bit width
until you reach
optimal performance

ap_fixed<width,integer>

0101.1011101010

integer | fractional

width

**Full performance
with 16 bits**

hls4ml



AUC / Expected AUC vs Fixed-point precision

Legend:
— Full
- - - Pruned

- g tagger
- q tagger
- w tagger
- z tagger
- t tagger

# Compression

- Remove **smallest** weigh
- Iterate

# Compression

- Remove **smallest** weights
- Iterate



70% REDUCTION OF WEIGHTS WITH NO LOSS IN PERF.

# CMS Level-1 trigger

**B-tagging**

**Autoencoder for Anomaly Detection**

**End-to-End Vertexing NN**



41

# High-Level Trigger (100 KHz, 100 ms latency)

**High-Level trigger**

Current 10K+

CPUs

# High-Level Trigger (100 KHz, 100 ms latency)

# ML-as-a-Service

- Simple support for mixed hardware
- Scaleable
- Throughput optimization for multiple-core
- Simple client-side



Clients     Server     Accelerators

FPGA
- Model E
- Model F

GPU
- Model A
- Model B

IPU
- Model C
- Model D

COPROCESSOR (GPU,FPGA,ASIC)

# Heterogeneous system for high throughput

- 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

# SONIC

- Within CMS software (CMSSW), the IaaS deployment scheme is called "Services for Optimized Network Inference on Coprocessors" (SONIC)

# Optimizing performance: CPU-to-GPU ratio

CMS mini-AOD
production



- Having explored server parameters, we can test the number of client jobs that a single GPU can handle
- We perform these tests in the cloud, as we need to synchronize jobs running on O(1000) CPU cores

# Summary

- **Artificial Intelligence heavily applied to Physics Discovery**
    - For examples, Higgs discovery!
- **HL-LHC confronted Big Data challenge**
    - Smart Machine Learning could offer partial solutions
- **A3D3 focusing on accelerating AI to solve common challenges through interdisciplinary collaboration**

**Fast Machine Learning for Science**

Real-time and accelerated ML for fundamental sciences

**Imperial College London**

25-28 September 2023

2024 TBA

Scientific Committee
Thea Årrestad (ETH Zurich)
Javier Duarte (UCSD)
Phil Harris (MIT)
Burt Holzman (Fermilab)
Scott Hauck
S[...]
S[...]
Mi[...]
Allis[...] (...odist University)
Mar[...] (...ana-Champaign)
Jennif[...] (...rmilab)
Mauri[...]ini (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

Shih-Chieh Hsu

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

schsu@uw.edu

Backup

# 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

51

# Studying SONIC at scale

- Inferences for three classes of algorithms were run through SONIC:
  - ONNX-based jet tagger
  - TensorFlow based missing energy calculation
  - TensorFlow based CNN for tau lepton ID
- These algorithms consume about 10% of total workflow latency

| Algorithm | Time [ms] | Fraction [%] | Input [MB] |
|---|---|---|---|
| PN-AK4 | 42.4 | 4.3 | 0.04 |
| PN-AK8 | 11.4 | 1.1 | 0.003 |
| DeepMET | 13.2 | 1.3 | 0.33 |
| DeepTau | 21.1 | 2.1 | 1.18 |
| ParticleNet+DeepMET+DeepTau | 88.1 | 8.8 | 1.55 |
| Total | 993.3 | 100.0 | — |

# Multi-messenger Astrophysics

- Develop and deploy software within astronomical facilities to enable discovery



Credit: Michael Coughlin (UMN)

# Gravitational Waves (LVK)

All algorithms use our underline inference-as-a-service (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*)**

# 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.**



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

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

# Improved time-series reconstruction methods
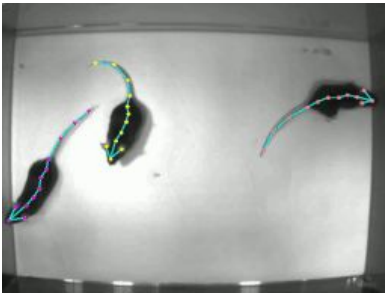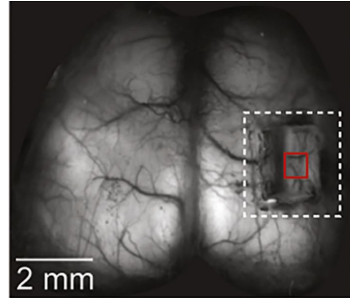
- Developed new Multi-block Recurrent Auto-Encoder (MRAE) to increase bandwidth more efficiently

- 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

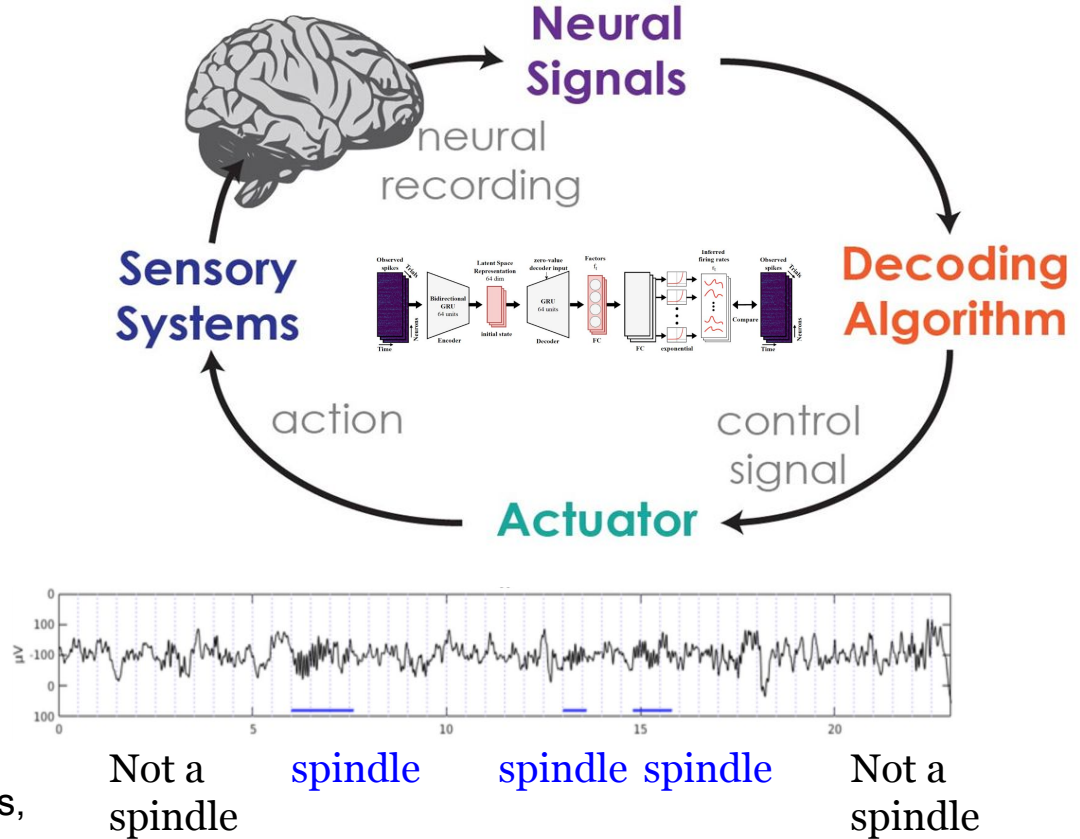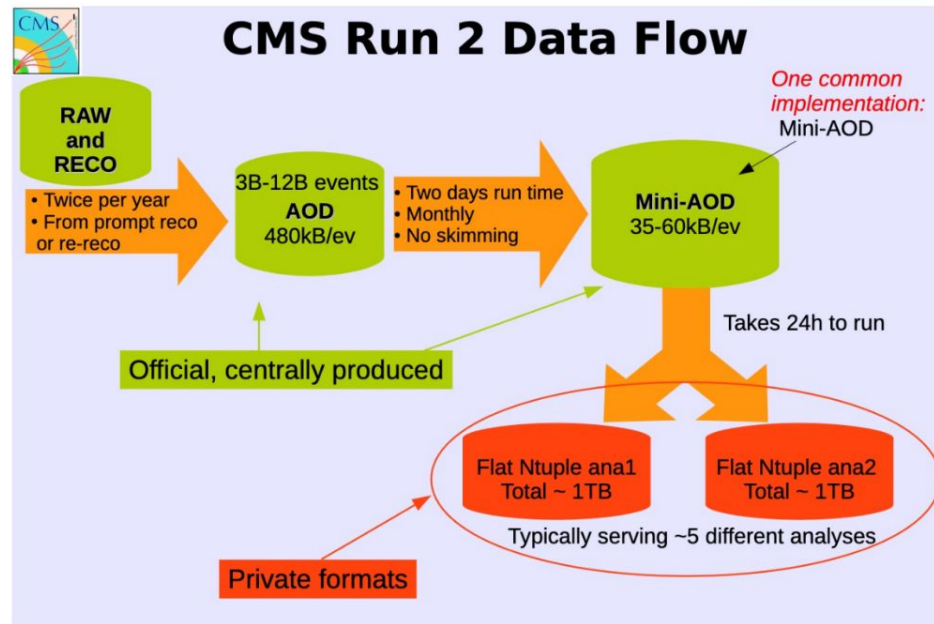- Neuro A3D3 develops methods for reconstruction, forecasting and clustering of time-series
- Potential applications/uses:
  - Detect noise and artifacts
  - Detect rare neural events of interest (e.g., seizures, spindles, etc)

# Studying SONIC at scale

- As a testbed for SONIC-enabled deployment, we created a MiniAOD demonstrator workflow
  - Runs a refinement and slimming step of CMS data processing
  - Full MiniAOD processing workflow typically run ~monthly



**CMS Run 2 Data Flow**

One common implementation: Mini-AOD

RAW and RECO
- Twice per year
- From prompt reco or re-reco

3B-12B events AOD 480kB/ev
- Two days run time
- Monthly
- No skimming

Mini-AOD 35-60kB/ev

Official, centrally produced

Takes 24h to run

Flat Ntuple ana1 Total ~ 1TB

Flat Ntuple ana2 Total ~ 1TB

Typically serving ~5 different analyses

Private formats

Mini-AOD production typically takes about 0.5 seconds per event on production grid nodes