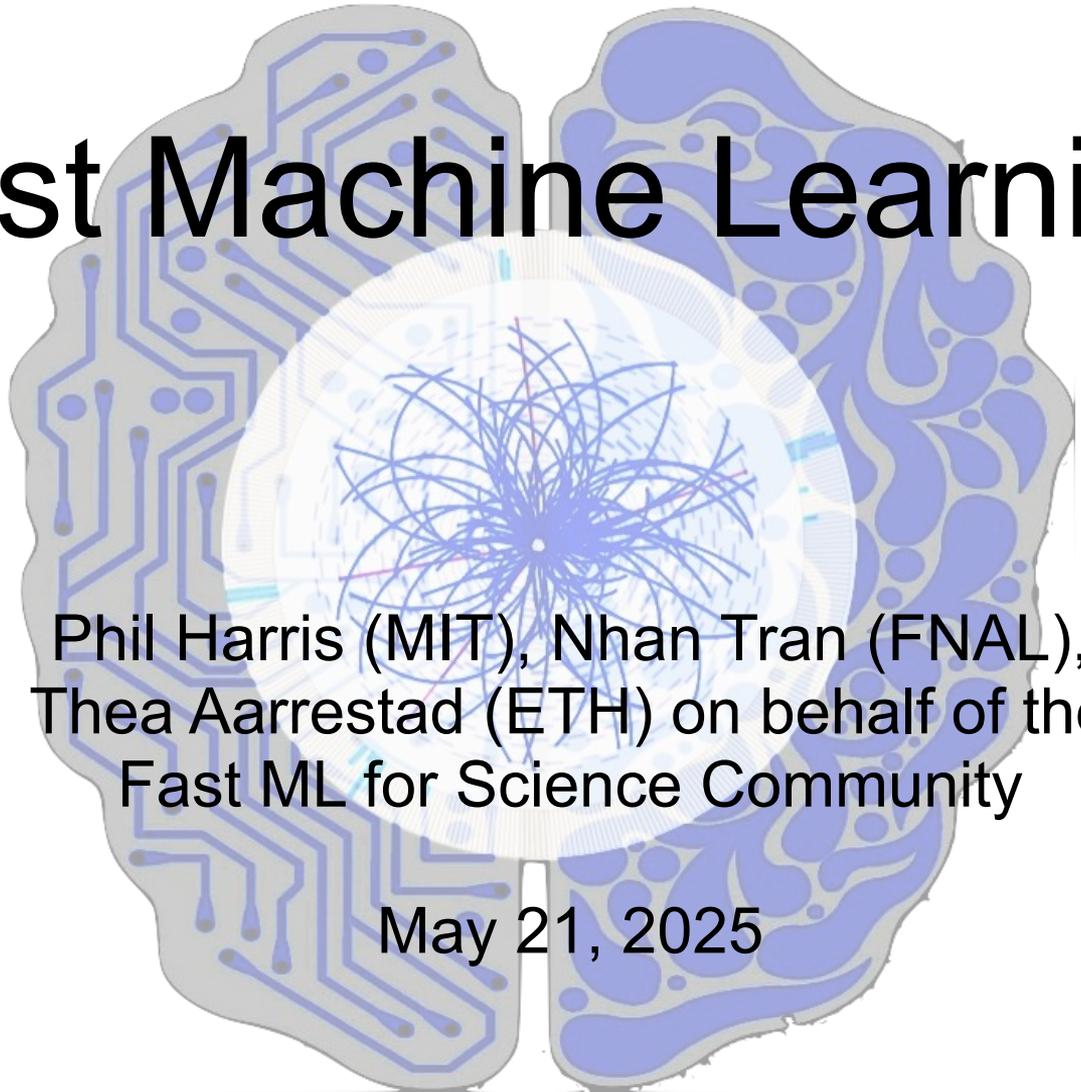


Fast Machine Learning



Phil Harris (MIT), Nhan Tran (FNAL),
Thea Aarrestad (ETH) on behalf of the
Fast ML for Science Community

May 21, 2025

ETH zürich

 **Fermilab**



What is Fast Machine Learning?

- Origin
 - Fast Machine Learning started not long after the HLS4ML paper
 - We were exploring a way to bring together the community around FastML
 - Clear interests were emerging from others at the LHC
- Fast ML was put together more concretely by end of 2019

fastmachinelearning.org was registered on:

Friday 16th of August 2019

4 years, 5 months and 28 days

general



Get Notifications for All Messages

About

Members 1,515

Agents

Tabs

Fast Machine Learning

10–13 Sept 2019
Fermi National Accelerator Laboratory
US/Central timezone

Enter your search term



<https://indico.cern.ch/event/822126/>

Mission

Cultivate resources in support of the multi-disciplinary Fast ML for Science community of domain science, machine learning (ML), and engineering researchers for the advancement of accelerated and autonomous scientific experimentation

- Organization and administration of regular meetings, events, and conferences to encourage the cross-pollination of ideas within the community and foster an inclusive environment for promoting new, multi-disciplinary collaborations
- Development and support of open-source software packages, firmware tools, hardware platforms, and benchmarks that increase productivity and accessibility to novel research techniques
- Community engagement and collection of feedback to ensure that activities are aligned with the needs and goals of community members

FastML Conferences and Workshops

J. Duarte

Timeline

1st annual @ Fermilab, 2019

2nd annual, online, 2020



4th annual @ Imperial, 2023

Satellite @ ICCAD, 2023

5th annual @ Purdue, 2024

Fast Machine Learning
September 10-12, 2017 at Fermilab

Sept. 10-11
HIS-HEP Blueprint Meeting

Sept. 12-13
Developer Bootcamp

Accelerating ML in science:
Ultrafast on-detector inference and real-time systems
Acceleration as-a-service
Hardware platforms
Distributed learning

Local Organization:
Charles Goble (Fermilab),
Gerrit Gottlieb (Fermilab),
Christine Gray (Fermilab),
Sally Gunther (Fermilab),
Alexey Kalinichev (Fermilab),
David Kuck (Fermilab),
Michael Neuman (Fermilab)

Scientific Organization:
Paul Adam (CERN),
Srinivas Aravamudan (Google),
Srinivas Aravamudan (Google)

Organizing Committee:
John DeNero (SMU),
Rohit Nayyar (SMU),
Thomas Coot (SMU),
Elizabeth Pfister (SMU)

Organizing Committee:
Javier Diaz (CERN),
Phil Hays (MIT),
Burt Holman (Fermilab),
Scott Hogg (Fermilab),
Srinivas Aravamudan (Google),
Mia Li (Fermilab),
Alison McCam (Fermilab),
Markus Neuman (CERN),
Mihai Trif (Fermilab)

REGISTER AND MORE INFORMATION
<http://indico.cern.ch/event/20200>

World Changes Shaped Here SMU

FAST MACHINE LEARNING FOR SCIENCE
A Virtual Event Hosted by Southern Methodist University at Dallas, Texas
November 30 to December 3

Organizing Committee:
John DeNero (SMU),
Rohit Nayyar (SMU),
Thomas Coot (SMU),
Elizabeth Pfister (SMU)

Organizing Committee:
Javier Diaz (CERN),
Phil Hays (MIT),
Burt Holman (Fermilab),
Scott Hogg (Fermilab),
Srinivas Aravamudan (Google),
Mia Li (Fermilab),
Alison McCam (Fermilab),
Markus Neuman (CERN),
Mihai Trif (Fermilab)

REGISTER AND MORE INFORMATION
<http://indico.cern.ch/event/20200>

World Changes Shaped Here SMU

FAST MACHINE LEARNING FOR SCIENCE
SOUTHERN METHODIST UNIVERSITY
OCTOBER 24-25, 2022

ORGANIZING COMMITTEE:
John DeNero (SMU),
Rohit Nayyar (SMU),
Thomas Coot (SMU),
Elizabeth Pfister (SMU)

SCIENTIFIC COMMITTEE:
Javier Diaz (CERN),
Phil Hays (MIT),
Burt Holman (Fermilab),
Scott Hogg (Fermilab),
Srinivas Aravamudan (Google),
Mia Li (Fermilab),
Alison McCam (Fermilab),
Markus Neuman (CERN),
Mihai Trif (Fermilab)

REGISTER AND MORE INFORMATION
<http://indico.cern.ch/event/20200>

Fast Machine Learning for Science
Real-time and accelerated ML for fundamental sciences
Imperial College London
25-28 September 2023

Organizing Committee:
John DeNero (SMU),
Rohit Nayyar (SMU),
Thomas Coot (SMU),
Elizabeth Pfister (SMU)

SCIENTIFIC COMMITTEE:
Javier Diaz (CERN),
Phil Hays (MIT),
Burt Holman (Fermilab),
Scott Hogg (Fermilab),
Srinivas Aravamudan (Google),
Mia Li (Fermilab),
Alison McCam (Fermilab),
Markus Neuman (CERN),
Mihai Trif (Fermilab)

REGISTER AND MORE INFORMATION
<http://indico.cern.ch/event/20200>

FAST MACHINE LEARNING FOR SCIENCE
Oct. 10 to 12, 2024, Purdue University

Fast ML meets fundamental sciences, quantum information science, semiconductors

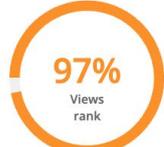
Organizing Committee:
John DeNero (SMU),
Rohit Nayyar (SMU),
Thomas Coot (SMU),
Elizabeth Pfister (SMU)

SCIENTIFIC COMMITTEE:
Javier Diaz (CERN),
Phil Hays (MIT),
Burt Holman (Fermilab),
Scott Hogg (Fermilab),
Srinivas Aravamudan (Google),
Mia Li (Fermilab),
Alison McCam (Fermilab),
Markus Neuman (CERN),
Mihai Trif (Fermilab)

REGISTER AND MORE INFORMATION
<http://indico.cern.ch/event/20200>

White paper

28,798
TOTAL VIEWS AND DOWNLOADS



<https://doi.org/10.3389/fdata.2022.787421>

frontiers | Frontiers in Big Data

REVIEW
SUBMITTED: 01 JULY 2022
doi: 10.3389/fdata.2022.787421

Applications and Techniques for Fast Machine Learning in Science

Allison McCarr Delana^{1*}, Nhan Tran^{1,2*}, Joshua Agar³, Michaela Blott⁴, Giuseppe Di Guglielmo⁵, Javier Duarte⁶, Philip Harris⁷, Scott Hauck⁸, Mia Liu⁹, Mark S. Neubauer¹⁰, Jennifer Ngalubwa¹¹, Seda Ogrenic-Memik¹², Maurizio Pierini¹³, Theo Ararates¹⁴, Steffen Bähr¹⁵, Jürgen Becker¹⁶, Anne-Sophie Berthold¹⁷, Richard J. Bonventre¹⁸, Tomás E. Müller Bravo¹⁹, Markus Diefenthaler²⁰, Zhen Dong²¹, Nick Fritzsche²², Amir Gholami²³, Ekaterina Goukova²⁴, Dongning Guo²⁵, Kyle J. Hazelwood²⁶, Christian Henning²⁷, Babar Khan²⁸, Sehoon Kim²⁹, Thomas Klitzmaier³⁰, Yaling Liu³¹, Kin Ho Lo³², Thi Nguyen³³, Gianantonio Pezzullo³⁴, Seyyidramin Rasoolimadadi³⁵, Ryan A. Rivera³⁶, Kate Scholberg³⁷, Justin Selig³⁸, Sougata Sen³⁹, Dmitri Strukov⁴⁰, William Tang⁴¹, Savannah Thies⁴², Kai Lukas Unger⁴³, Ricardo Vialta⁴⁴, Belina von Krosigk^{45,46}, Shen Wang⁴⁷ and Thomas K. Warburton⁴⁸

OPEN ACCESS

Edited by: Elena Cucco, European Gravitational Observatory, Italy

Reviewed by: Anshu Agrawal, IBM Research-Zurich, Switzerland; Ruty Rapoport, IBM, Israel; *Correspondence: Allison McCarr Delana, adelana@fdata.fri.uni-goettingen.de; Nhan Tran, ntran@fdata.fri.uni-goettingen.de

Specialty section: This article was submitted to Big Data and AI in High Energy Physics, a section of the journal Frontiers in Big Data

Received: 30 September 2021
Accepted: 21 January 2022
Published: 12 April 2022

Citation: Delana AM, Tran N, Agar J, Blott M, Di Guglielmo G, Duarte J, Harris P, Hauck S, Liu M, Neubauer MS, Ngalubwa J, Ogrenic-Memik S, Pierini M, Ararates T, Bähr S, Berthold A-S, Bonventre RJ, Bravo TEM, Becker J, Diefenthaler M, Dong Z, Fritzsche N, Gholami A, Goukova E, Guo D, Hazelwood K, Henning C, Kim S, Klitzmaier T, Lo KH, Nguyen T, Pezzullo G, Rasoolimadadi S, Rivera RA, Scholberg K, Selig J, Sen S, Strukov D, Tang W, Thies S, Unger N, Vialta R, von Krosigk B, Wang D and Warburton TK (2022) Applications and Techniques for Fast Machine Learning in Science. *Front. Big Data* 5:787421. doi: 10.3389/fdata.2022.787421

Frontiers in Big Data | www.frontiersin.org | April 2022 | Volume 5 | Article 787421

“

Scientific discoveries come from groundbreaking ideas and the capability to validate those ideas by testing nature at new scales—finer and more precise temporal and spatial resolution. This is leading to an explosion of data that must be interpreted, and ML is proving a powerful approach. The more efficiently we can test our hypotheses, the faster we can achieve discovery. To fully unleash the power of ML and accelerate discoveries, it is necessary to embed it into our scientific process, into our instruments and detectors.

”

Applications and Techniques for Fast Machine Learning in Science

Physics Community Needs, Tools, and Resources for Machine Learning

Snowmass <https://arxiv.org/pdf/2203.16255>

Snowmass <https://arxiv.org/pdf/2204.13223>

Transcendental Preprint
April 29, 2022

Smart sensors using artificial intelligence
for on-detector electronics and ASICs

Communication within FastML at different tiers

- Organize regular meetings, events, and conferences to encourage the cross-pollination of ideas within the community
 - <https://indico.cern.ch/category/11842/>

hls4ml Meeting (Type A)	One of these Every Friday @8am/10am/11am/17h PDT/CDT/EDT/CEDT	135 events	➡
Co-processor Meeting (Type B)		77 events	➡
General Meeting		53 events	➡
Workshops and Conferences	Once a year	6 events	➡
Special Meetings	Occasionally	25 events	➡
Training Events and Tutorials		1 event	➡

What are the goals of fast machine learning

Promote interdisciplinary collaborations

physicists, computer scientists, electrical and computer engineers, software engineers

Custom embedded systems

Off-the-shelf coprocessors

Build open-source, multi-technology
codesign workflows

Be nimble: abstraction, portability,
containerization

Novel ML research concepts: efficient, fault-tolerant, reliable

Open data, task-based, and data-based benchmarks

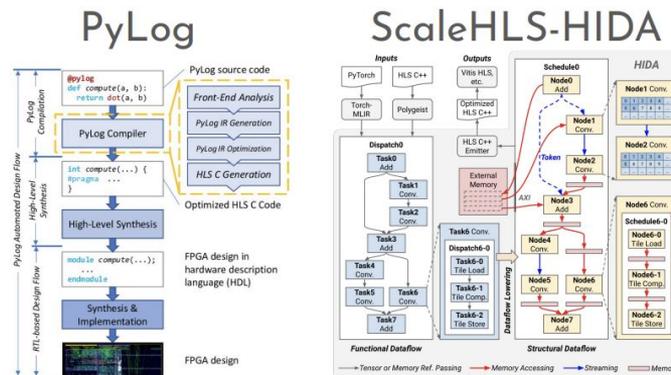
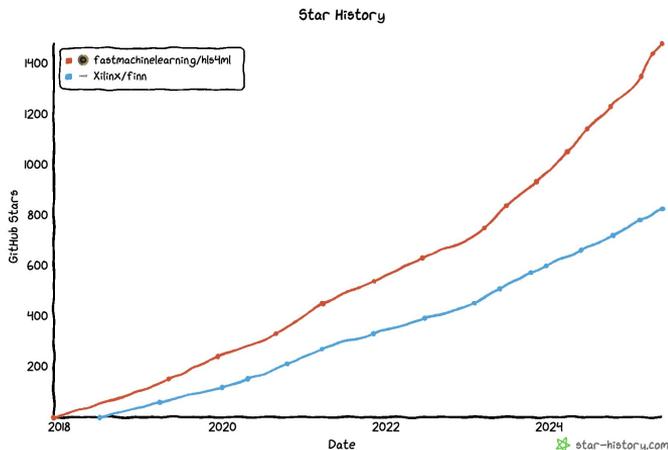
Support ecosystem integration and operation

SuperSONIC



HLS4ML(Embedded Systems) Subgroup

- **Conveners:** Benjamin Ramhorst (ETH-Z), Jan-Frederik Schulte (Purdue)
- Two types of meetings:
 - Type I : Survey the community for new ideas related to FPGA/ASIC programming
 - Here we often talk about community tools (not just HLS4ML)
 - Sometimes we hear about new scientific applications
 - Type II : Current status and planning of the HLS4ML software



Coprocessor Subgroup



- **Conveners:** Yongbin Feng (TTU), Yuantang Chou (UW), Ethan Marx (MIT)
- This meeting aims to cover our work with large scale compute clusters
 - How do we run Machine Learning Fast, but focusing on conventional tools and large scale

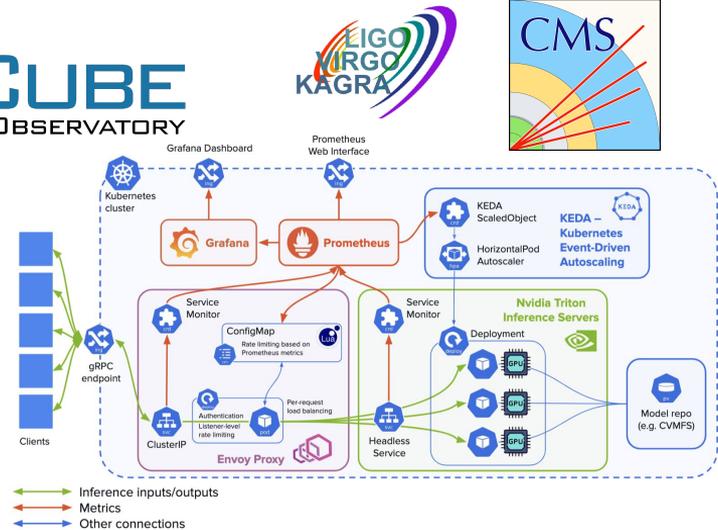
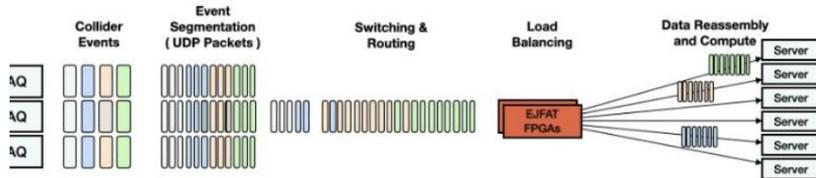


NVIDIA.

TRITON INFERENCE SERVER

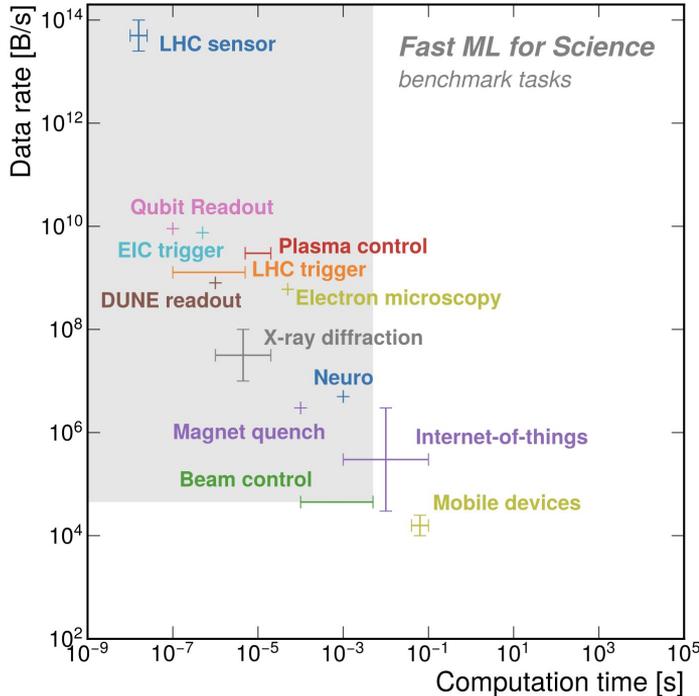


ICECUBE
NEUTRINO OBSERVATORY



FastML benchmarking

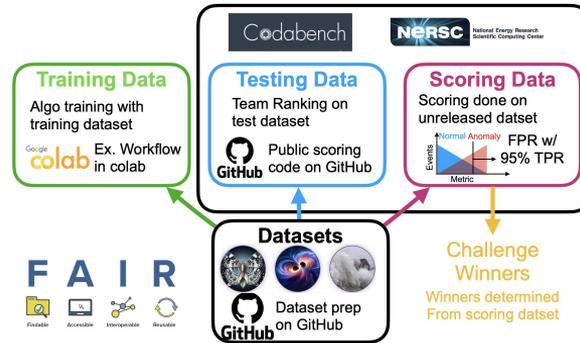
- Large variety of different domains and benchmarks
 - Actively working to connect this with greater ML benchmarking efforts



ML Commons



NSF HDR ML Challenge



mlcommons/
tiny_results_v1.1

This repository contains the results and code for the MLPerf™ Tiny Inference v1.1 benchmark.

5 Contributors 0 Issues 1 Star 2 Forks



Through A3D3 organized an challenge with 600 teams

FastML benchmarking

- Connecting FastML to the benchmarking community
 - Looking towards a website that can allow us to organize benchmarking
 - These elements are critical towards building robust Science drivers

MLPerf Tiny Benchmark

Colby Banbury^{*} Vijay Janapa Reddi^{*} Peter Torelli[†] Jeremy Holleman^{‡¶} Nat Jeffries[§]

Csaba Kiraly[¶] Pietro Montino^{*} David Kanter^{**} Sebastian Ahmed^{††} Danilo Pau^{†‡}

Urmish Thakker[‡] Antonio Torrin[¶] Peter Warden[§] Jay Cordaro[‡] Giuseppe Di Guglielmo^{¶¶}

Javier Duarte^{IV} Stephen Gibellini[†] Videet Parekh^V Honson Tran^V Nhan Tran^{VII}

Niu Wenxu^{VIII} Xu Xuesong^{VII}

Abstract

Advancements in ultra-low-power *tiny* machine learning (TinyML) systems promise to unlock an entirely new class of smart applications. However, continued progress is limited by the lack of a widely accepted and easily reproducible benchmark for these systems. To meet this need, we present MLPerf Tiny, the first industry-standard benchmark suite for ultra-low-power tiny machine learning systems. The benchmark suite is the collaborative effort of more than 50 organizations from industry and academia and reflects the needs of the community. MLPerf Tiny measures the accuracy, latency, and energy of machine learning inference to properly evaluate the tradeoffs between systems. Additionally, MLPerf Tiny implements a modular design that enables benchmark submitters to show the benefits of their product, regardless of where it falls on the ML deployment stack, in a fair and reproducible manner. The suite features four benchmarks: keyword spotting, visual wake words, image classification, and anomaly detection.

FASTML SCIENCE BENCHMARKS: ACCELERATING REAL-TIME SCIENTIFIC EDGE MACHINE LEARNING

Javier Duarte^{†1} Nhan Tran^{†2} Ben Hawks² Christian Herwig²
Jules Muhizi³ Shvetank Prakash¹ Vijay Janapa Reddi¹

ABSTRACT

Applications of machine learning (ML) are growing by the day for many unique and challenging scientific applications. However, a crucial challenge facing these applications is their need for ultra-low-latency and on-detector ML capabilities. Given the slowdown in Moore's law and Dennard scaling, coupled with the rapid advances in scientific instrumentation that is resulting in growing data rates, there is a need for ultra-fast ML at the extreme edge. Fast ML at the edge is essential for reducing and filtering scientific data in real-time to accelerate science experimentation and enable more profound insights. To accelerate real-time scientific edge ML hardware and software solutions, we need well-constrained benchmark tasks with enough specifications to be generically applicable and accessible. These benchmarks can guide the design of future edge ML hardware for scientific applications capable of meeting the nanosecond and microsecond level latency requirements. To this end, we present an initial set of scientific ML benchmarks, covering a variety of ML and embedded system techniques.

ML Commons Science

Evaluate, organize, curate, and integrate applications, models/algorithms, irreproducible data, and datasets. These artifacts are organized through the MLCommons GitHub, which includes independently funded activities at industry, government, and research.

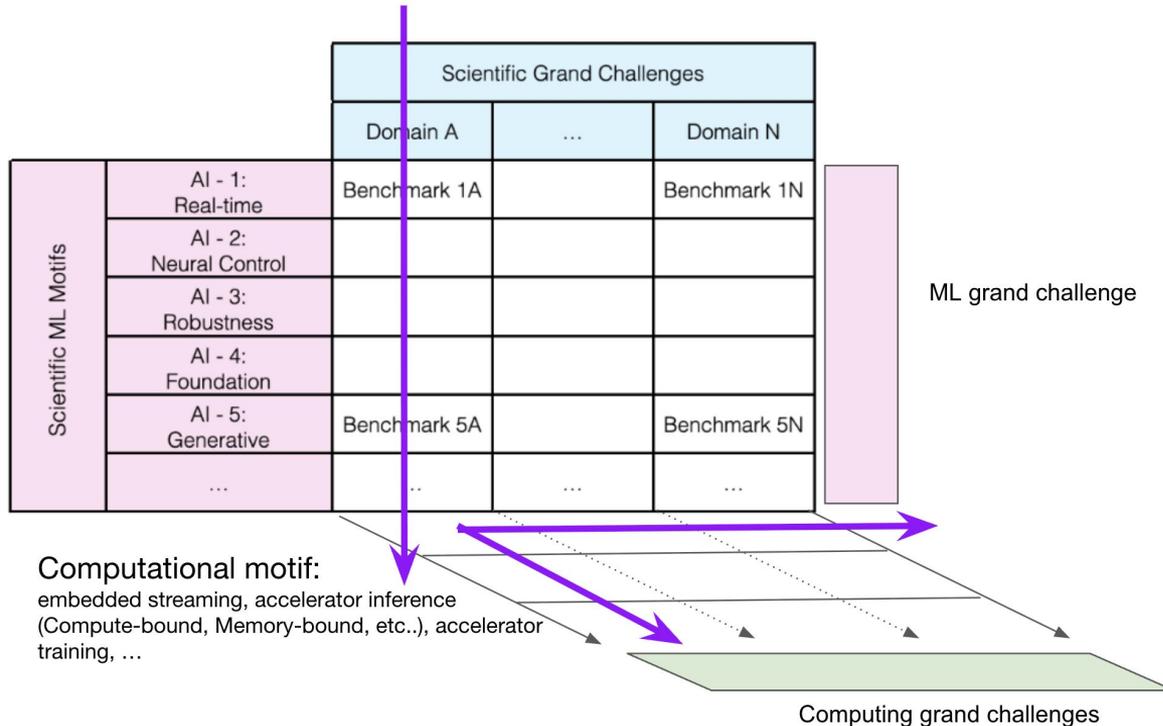
Join the Working Group

CONNECT WITH US:  

Building Machine Learning Challenges for Anomaly Detection in Science

Editors: Elizabeth G. Campolongo³, Yuan-Tang Chou², Ekaterina Govorkova¹, Wahid Bhimji¹⁰, Wei-Lun Chao³, Chris Harris¹⁰, Shih-Chieh Hsu², Hilmar Lapp⁴, Mark S. Neubauer⁶, Josephine Namayanja⁹, Aneesh Subramanian⁵, Philip Harris¹,

w/ML Commons Science



We are aiming to extend work on ML Challenges, benchmarking and performance to make a Computing axis for future scientific benchmarking

FastML Tutorials and Workshops

- We are regularly organizing tutorials and adjacent workshops
 - Recently hosted a workshop at ICCAD in 2023
 - Looking to organize a NeurIPS Wrokshop this year (fingers crossed)



Fast ML for Science @ ICCAD 2023

Home Submission Program Registration Committee

Fast Machine Learning for Science Workshop

Co-located with 2023 [International Conference on Computer-Aided Design \(ICCAD\)](#)

Date: November 2, 2023



Q Search

hls4ml-tutorial: Tutorial notebooks for hls4ml

Part 1: Getting started

Part 1: Getting started

```
from tensorflow.keras.utils import to_categorical
from sklearn.datasets import fetch_openml
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder, StandardScaler
import numpy as np
```



hls4ml AI Hardware Development Training

Arizona State University
April 30-May 1, 2025

Developed and taught by scientists leading research efforts in AI/ML ∩ Science!

FastML Engaging the public

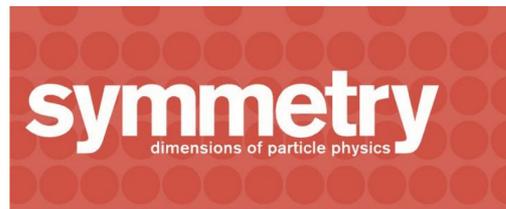
**Algorithms and Flow:
Lupe Fiasco's Creative
Use of LLMs**

September 26, 2025,
7:00 pm

Lecture Theatre 1,
Blackett Laboratory



Sign up here!
Fast Machine Learning Conference



FastML Engaging Industry

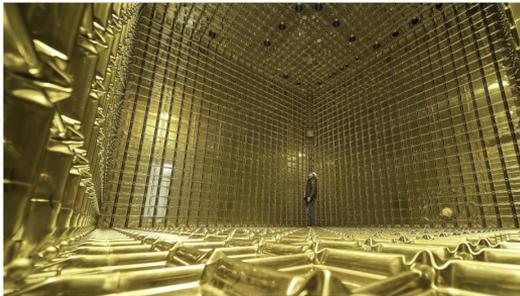
- Active collaborations with many different members of industry
 - Working with members of AMD, Nvidia, Siemens, Microsoft (previously), Intel, Graphcore,...
 - Much of our work has direct benefits with industry
- We see industry as a key player in helping scientific goals of FastML

Scaling Inference in High Energy Particle Physics at Fermilab Using NVIDIA Triton Inference Server

Apr 30, 2021

0 Like 0 Discuss (0)

By Shankar Chandrasekaran, Lindsey Gray, Farah Hariri, Kevin Pedro, Vartika Singh, Nhan Tran, Mike Wang and Tingjun Yang



RELEASE

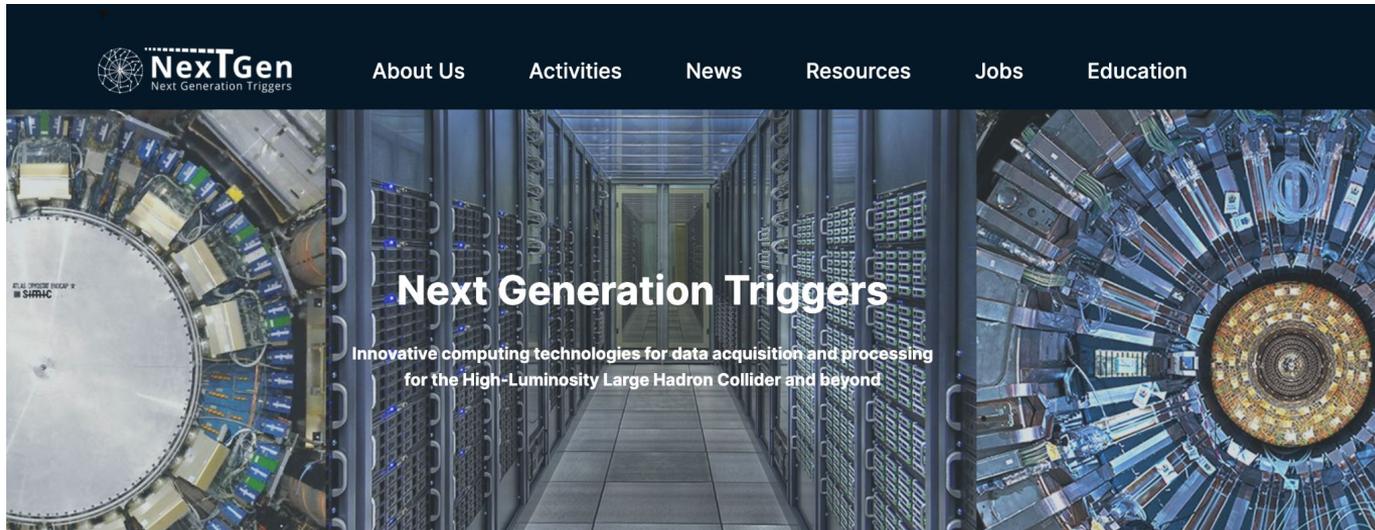
Siemens simplifies development of AI accelerators for advanced system-on-chip designs with Catapult AI NN

, 2024
Texas



FastML connections with the LHC

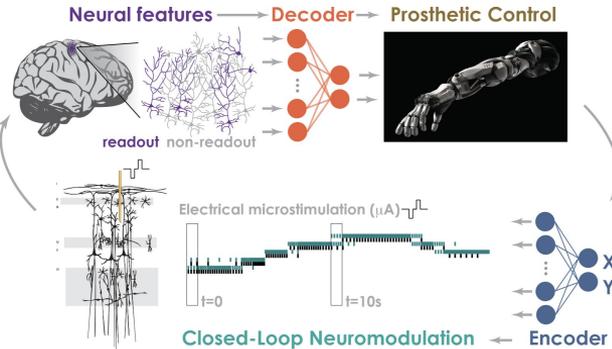
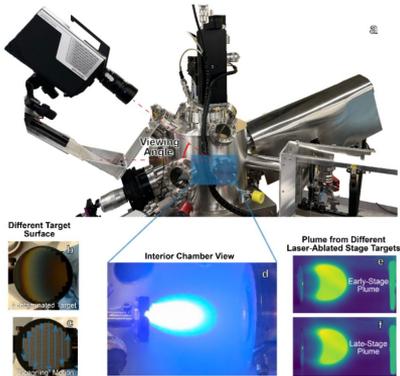
- Next Generation Trigger Project has elevated FastML at the LHC
 - <https://nextgentriggers.web.cern.ch/>
- Goal: establish working, cross-collaboration environment around triggering
 - Aligned workpackages and integration throughout all tiers of the LHC



FastML connections with Other Domains

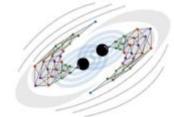
- Many domains are emerging that benefit from FastML
 - Astro, Plasma Physics, Neuroscience, Quantum Computing, Nuclear Physics , Health...
- We are constantly seeing interest to expand to many different domains

Materials Science



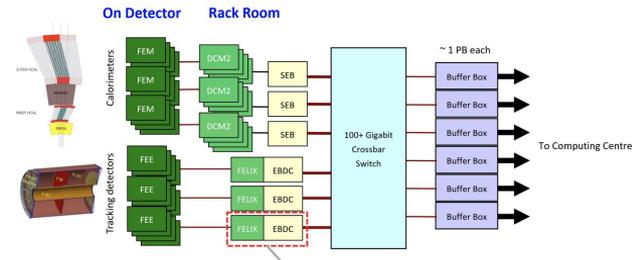
Neuroscience

Astrophysics



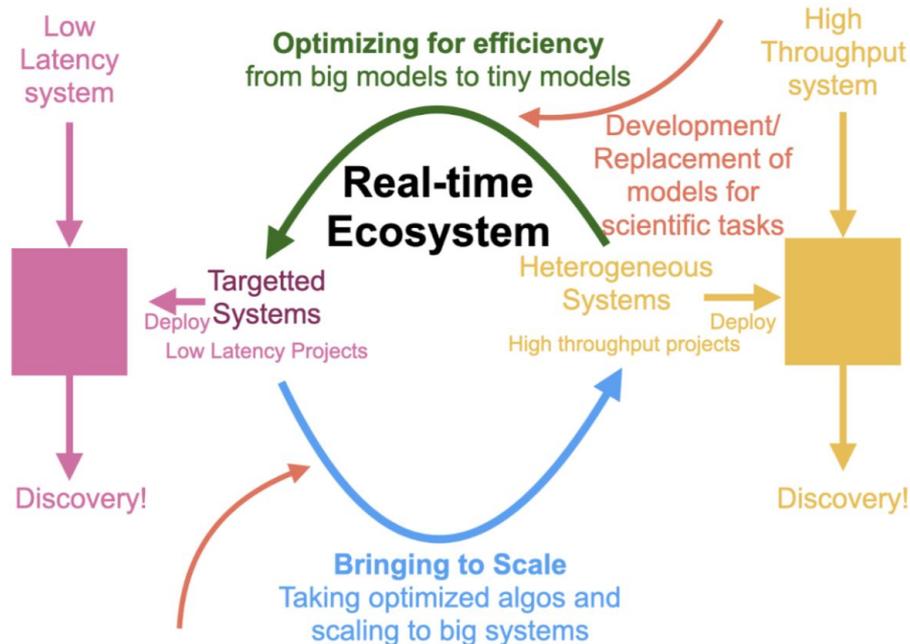
<https://github.com/ML4GW>

Nuclear Physics



FastML Sustainability Goals

- We are working to highlight the importance of FastML to the US



Cultivating an ecosystem is important

Workshops like this one are essential!

We are happy to see more interested people!

Conclusion and looking forward

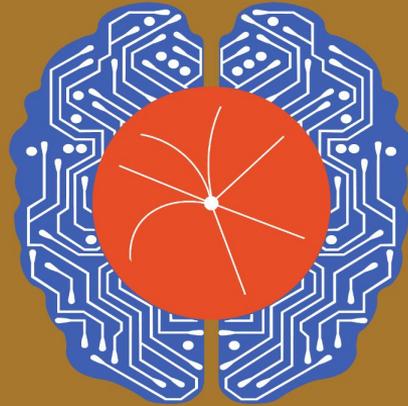
- Towards a more sustainable and transparent community
 - We would like to continue growing the community
 - Your input is essential to keeping the community vibrant and dynamic
- Discussions in the last 1.5 years about how to evolve the community
 - At its heart, we are open-source and open science driven
 - We want to continue to support the open source environment
 - Looking for away to sustain FastML for a long period of time
 - Possibility of support of projects and resources
 - A complicated process but we have learned a lot about other open-source models through this exploration –
 - <https://fastmachinelearning.org/pose/> – work in progress
- More Generally we are excited to hear from all of you
 - Keep in touch with the community throughout the year
 - Subscribe to the e-group hls-fml@cern.ch, join Slack, and attending meetings

Hope to see you there!

1-5 September 2025

fast machine learning for science

Real-time
and
accelerated
ML
for
fundamental
sciences



indi.to/fastml25

ETH zürich

