# A3D3 Workshop HEP Hackathon Project

# Real-Time Anomaly Detection with HEP Open Data

Daniel Diaz, Elham E Khoda, Melissa Quinnan
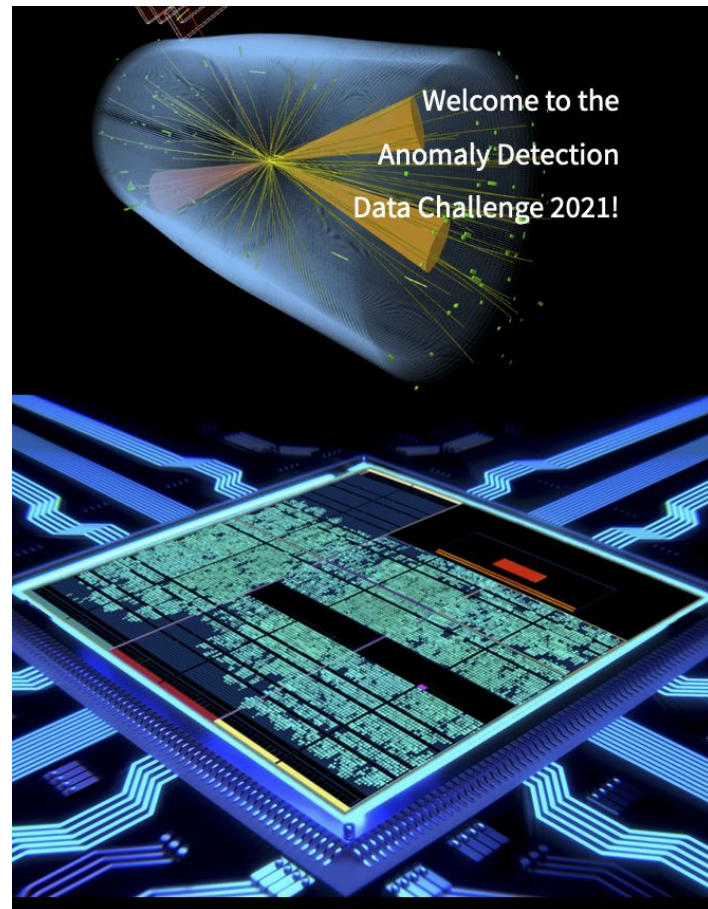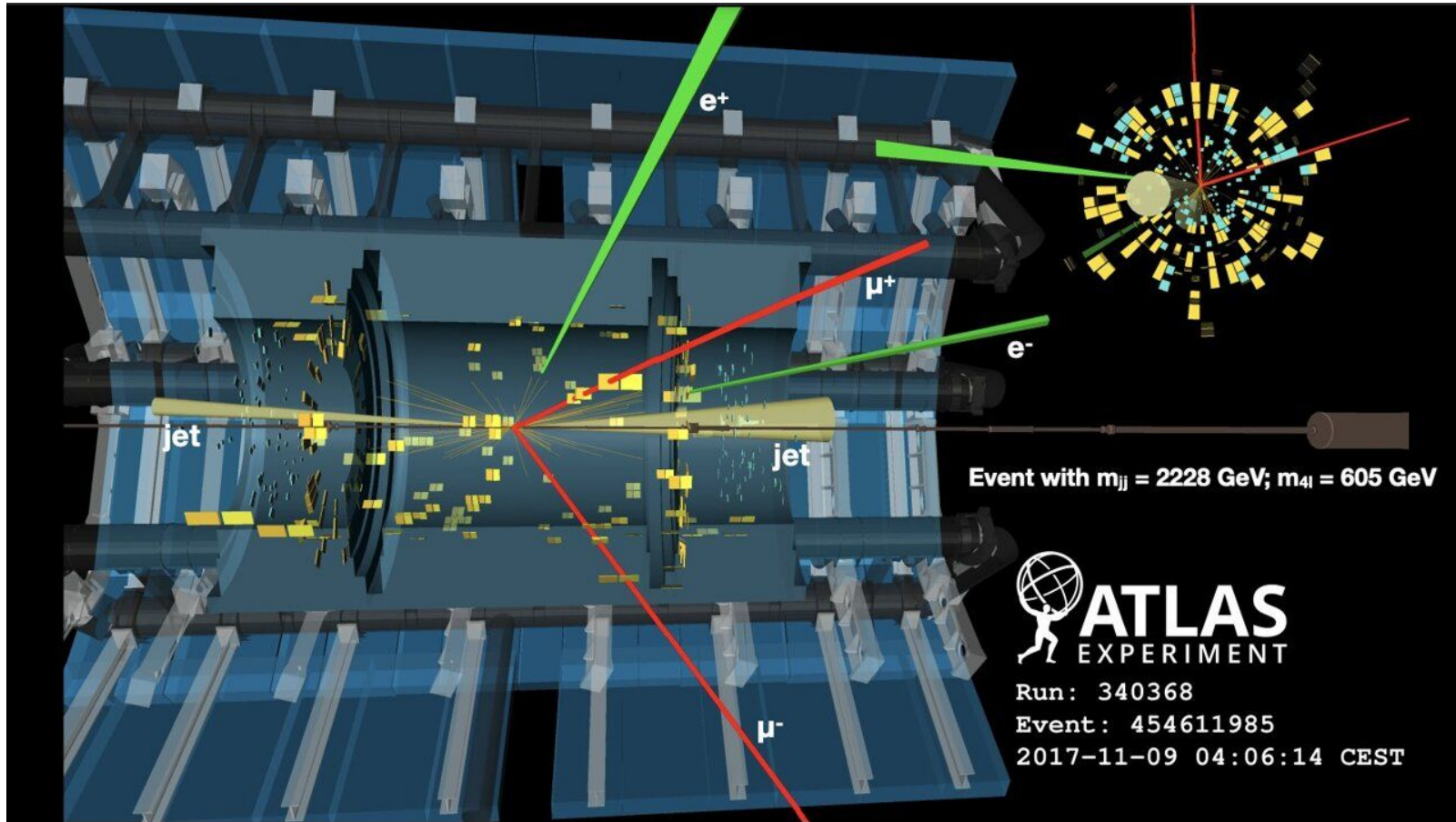
July 12-14 2023

# Introduction

- We are going to work on the Anomaly Detection Data Challenge 2021

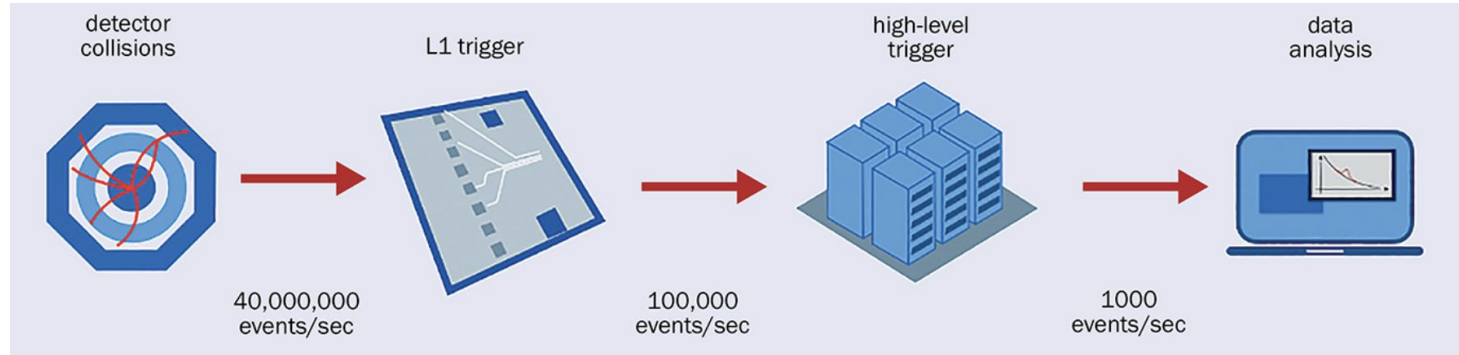  *Unsupervised new physics detection at 40 MHz*


- Resources:

  - [Challenge website](#)
  - [Challenge introduction from ML4Jets 2021](#)
  - [Challenge example code](#)

# Particle physics collisions



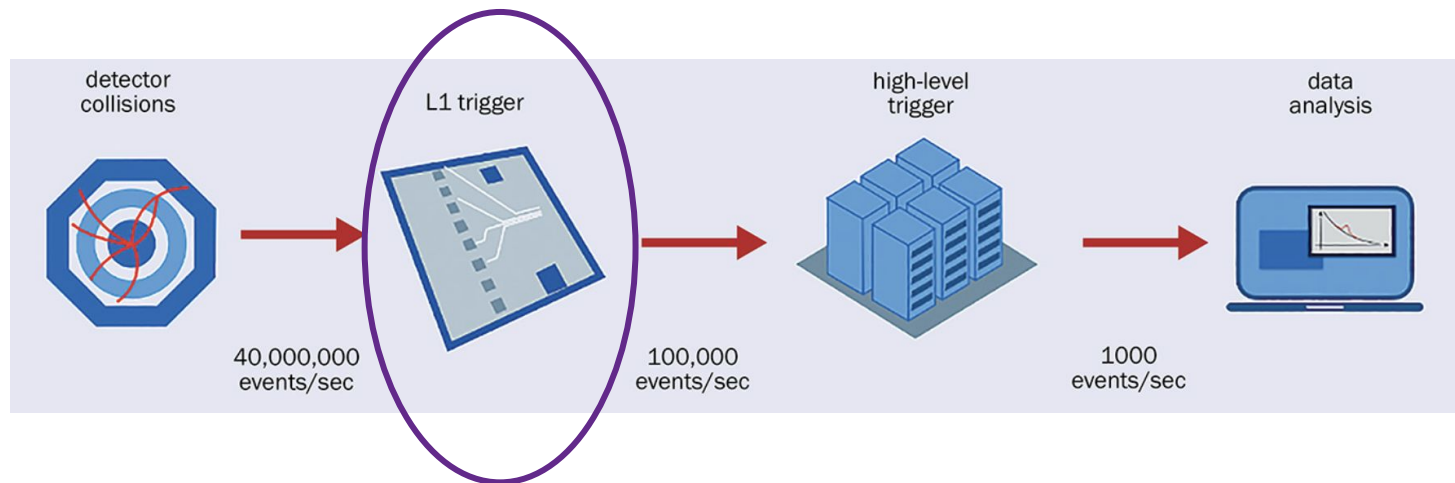Event with $m_{jj}$ = 2228 GeV; $m_{4l}$ = 605 GeV

ATLAS
EXPERIMENT

Run: 340368
Event: 454611985
2017-11-09 04:06:14 CEST

# HEP Data Processing



**L1 Trigger** (hardware: FPGAs) – *O(μs) hard latency*

- Typically coarse selections are applied

**High Level Trigger** (software: CPUs) – *O(100 ms) soft latency*

- More complex algorithms (full detector information available), some BDTs and DNNs used

# Focus of the challenge: L1 Trigger



**L1 Trigger** (hardware: FPGAs) – *O(μs) hard latency*

- Typically coarse selections are applied

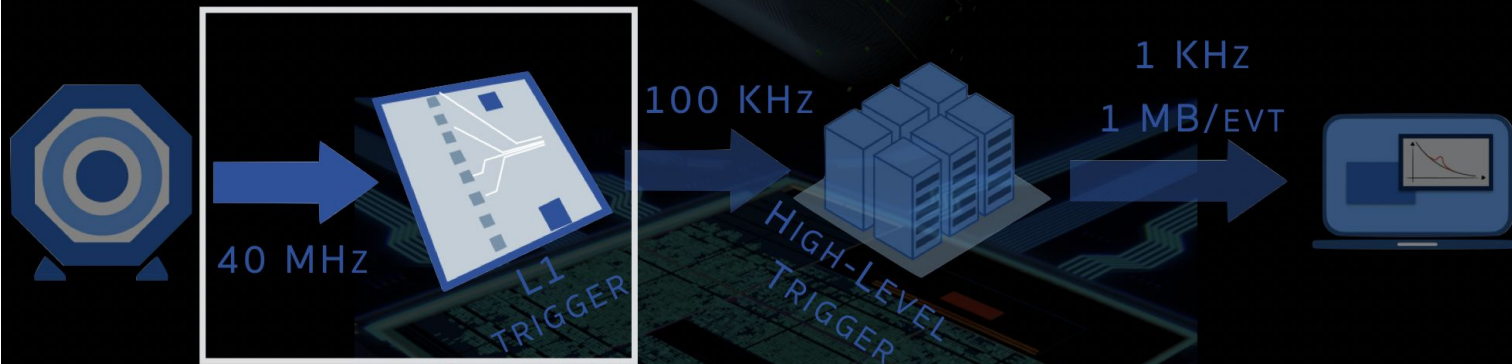**High Level Trigger** (software: CPUs) – *O(100 ms) soft latency*

- More complex algorithms (full detector information available), some BDTs and DNNs used

# Unsupervised new physics detection at 40 MHz

**Idea** is to **look for** something **very rare and unusual** directly in the **Level-1 Trigger** without any signal hypothesis in mind

**The challenge** is to find a-priori **unknown** and **rare New Physics** hidden in a data sample dominated by ordinary Standard Model processes

**40 MHz** → **L1 TRIGGER** — **100 KHz** → **HIGH-LEVEL TRIGGER** — **1 KHz 1 MB/EVT** →

**The deliverable** is a developed **algorithm** that can be deployed and run **in L1** with strict **latency** requirement of **< 1 microsecond**

**The task** is therefore to design an architecture that maximises the **sensitivity for New Physics** but at the **lowest** possible **resource** and **latency** budget

Taken from

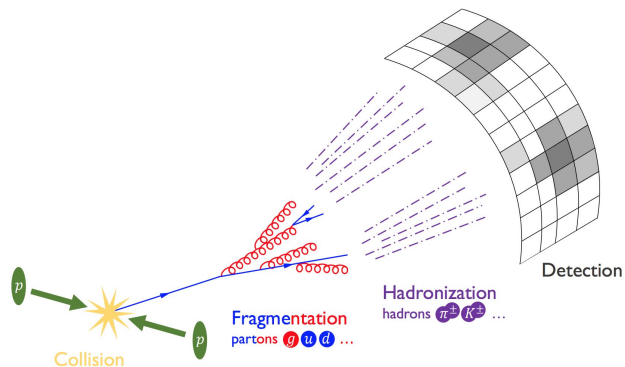Challenge introduction from ML4Jets 2021

# Hackathon Dataset

- Get the dataset from here: https://mpp-hep.github.io/ADC2021/
- There are 5 dataset files
  - Background dataset
  - 4 different signal datasets

- **Dataset dimensionality:** 57 = 19x3 "particles"
  - 10 jets
  - 4 electrons
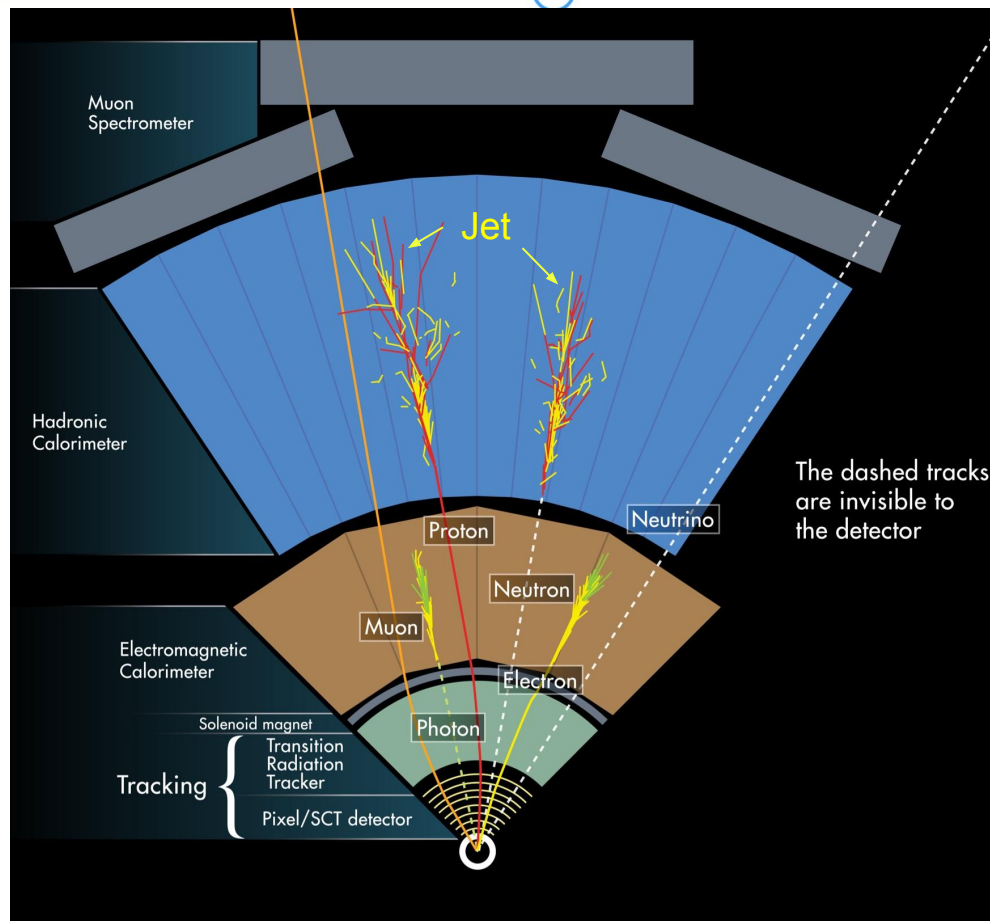  - 4 muons
  - MET

# Relevant Particle Responses

A quark always forms a spay of particle before getting detected

→ **Jet**



**MET = Missing Transverse Energy**

*corresponds to all the missing particles, invisible to the detector*

# Hackathon Dataset

- Get the dataset from here: https://mpp-hep.github.io/ADC2021/
- There are 5 dataset files
  - Background dataset
  - 4 different signal datasets

- **Dataset dimensionality:**

57 = 19x3 "particles"

  - 10 jets
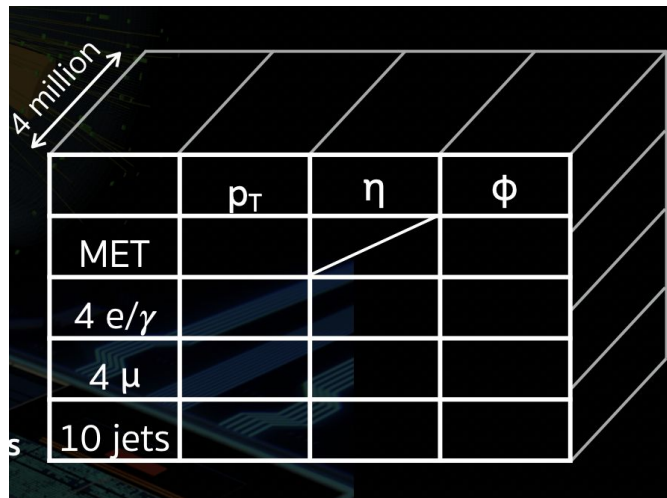  - 4 electrons
  - 4 muons
  - MET

# Eta and Phi

# Hackathon Training Dataset

Train with [4 million background](#)-like events 📄

The file contains:
- Inclusive W production, with $W \rightarrow l\nu$ (59.2%)
- Inclusive Z production, with $Z \rightarrow ll$ (6.7%)
- tt production (0.3%)
- QCD multijet production (33.8%)



Paper describing the dataset: [https://arxiv.org/abs/2107.02157](https://arxiv.org/abs/2107.02157)
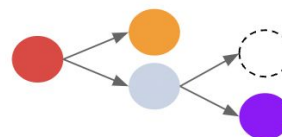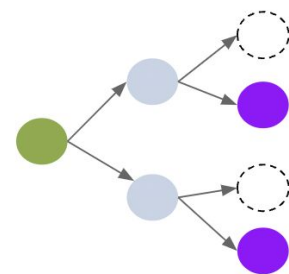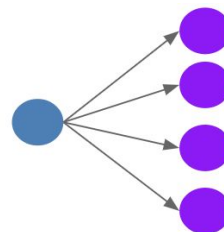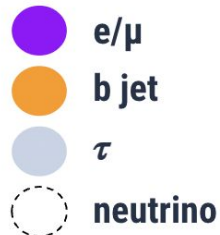
# Model Development and Evaluation

Evaluate performance on several different New Physics simulated samples

**New physics benchmarks**
- Neutral scalar boson (A), 50 GeV → 4 l
- Leptoquark (LQ), 80 GeV → b τ
- Scalar boson (h$^0$), 60 GeV → τ τ
- Charged scalar boson(h$^+$), 60 GeV → τ $\nu$

| | |
|---|---|
| 🔵 A | 🟣 e/μ |
| 🟢 h$^0$ | 🟠 b jet |
| 🔴 LQ | ⚪ $\tau$ |
| 🟡 h$^+$ | ⭕ neutrino |

# Autoencoder: One of the popular choice

- Train the model with **background-enriched data**

- Encode the inputs to a low dimensional representation and try to decode it back to the input set

- **Anomalous events** are often **poorly reconstructed** given low, if any, examples present during training

# Some Published Solutions:

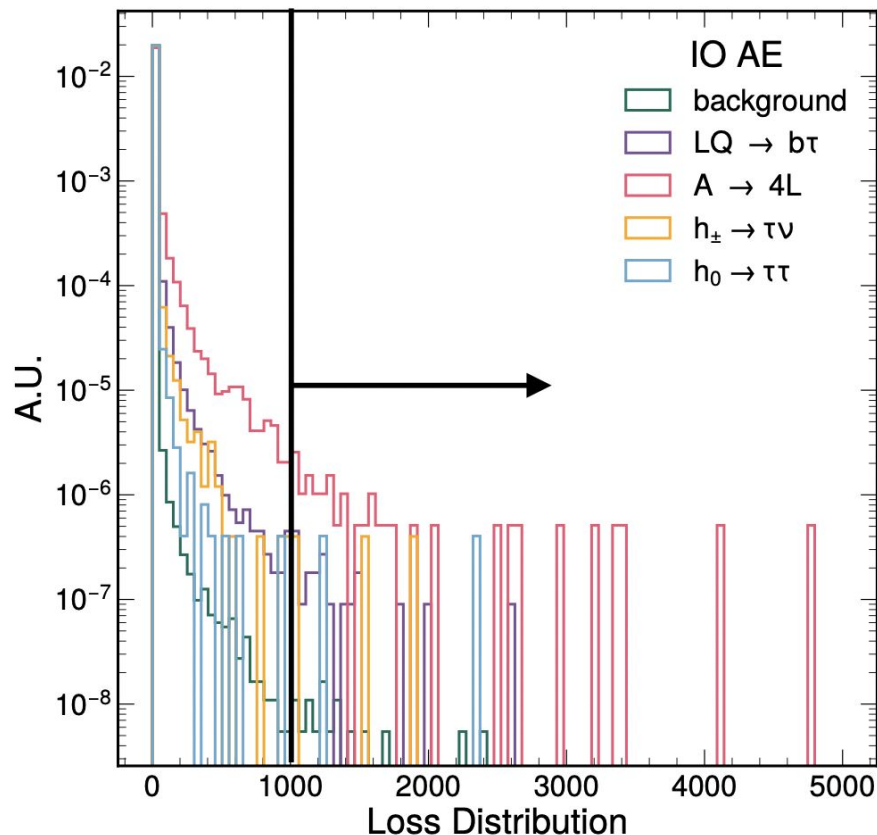**The first paper by the ADC2021 challenge organizers:**

https://arxiv.org/abs/2108.03986

⇒ Studies Autoencoders and Variational Autoencoders (VAE)

Another solution-based on contrastive learning (using Autoencoders):

https://arxiv.org/abs/2301.04660

# Challenges and Expectations

**Real-life application:** Low-Latency inference (~ 1 $\mu$s)

⇒ *Make sure your model is not too huge and can obey this latency constraints*

- An estimate of the algorithm efficiency can be obtained by calculating the floating- point operations per second (FLOPs)
- Example code to compute FLOPs: computeFLOPs.ipynb

**Some results are already published based on Autoencoders and VAEs**

⇒ Do not use vanilla Autoencoders and Variational Autoencoders

# Useful Resources

We will use [a3d3-hackathon/hep-ad-2023](#) repository to develop our code

**Other Resources:**

- [Challenge website](#)
- [Challenge introduction from ML4Jets 2021](#)
- [Challenge example code](#)
- Toolkit for implementing ML inference on FPGAs: [hls4ml](#)