# Working with neural networks at CERN

Roope Niemi

# Who I am & how I got here

- **Sotungin lukio 2008-2011**

- **2012-2017: non-IT work**

- **2017-2022: Studies at University of Helsinki, Kumpula (Computer Science, Data Science)**

- **2021-2024: Work at Nokia, applying to CERN a few times**

- **September 2024 → Software engineer / data scientist (QUEST) at CERN**

# What I do

- **Optimize neural networks for hardware, so they are fast and accurate → design methods to train compressed neural networks**

- **Programming with Python, read scientific papers, implement algorithms from them and compare them to other algorithms. Possibly improve them**

$$\mathcal{S}_g(w, s) := \text{sign}(w) \cdot \text{ReLU}(|w| - g(s))$$

```python
class STR(PruningLayer):
    def __init__(self, config, layer, out_size):
        super(STR, self).__init__()
        self.config = config
        threshold_size = get_threshold_size(config, out_size, layer.weight.shape)
        self.s = nn.Parameter(torch.ones(threshold_size) * -self.config.threshold_init)
        self.g = torch.sigmoid

    def forward(self, weight):
        """
        sign(W) * ReLu(|W| - g(s))
        """
        mask = self.get_mask(weight)
        return torch.sign(weight) * mask.view(weight.shape)

    def get_mask(self, weight):
        return torch.relu(torch.abs(weight).view(weight.shape[0], -1) - self.g(self.s))
```
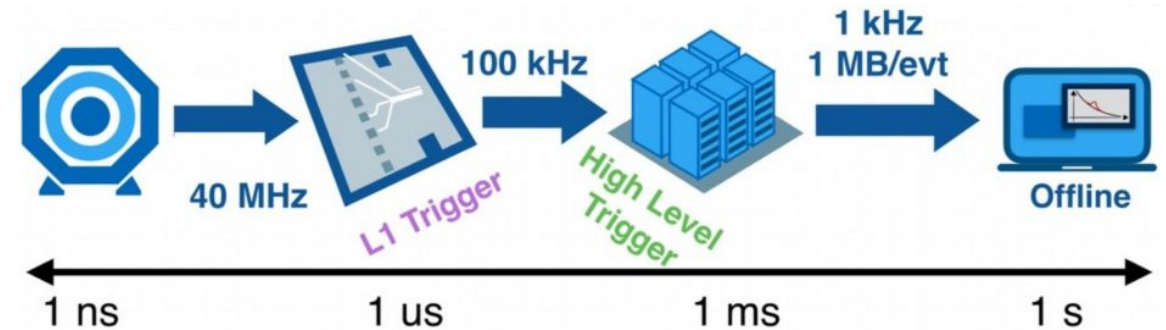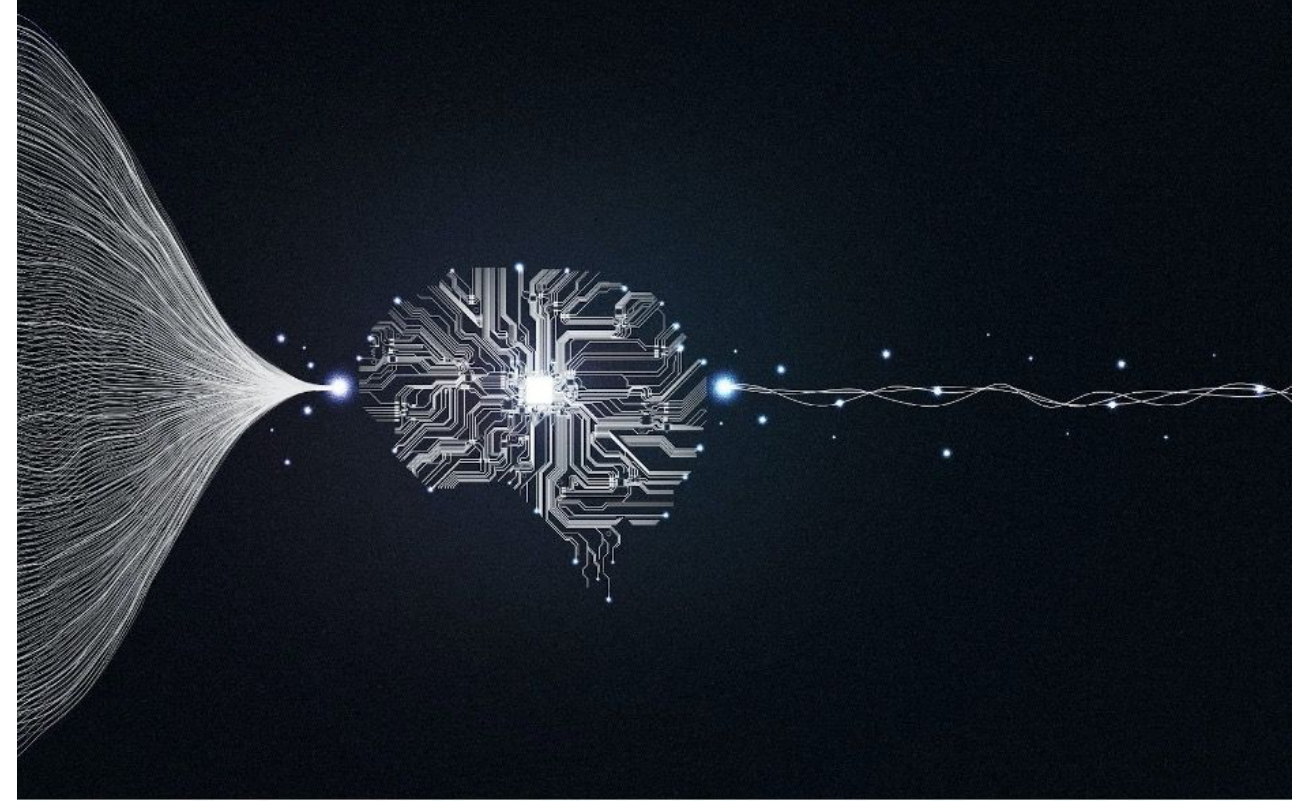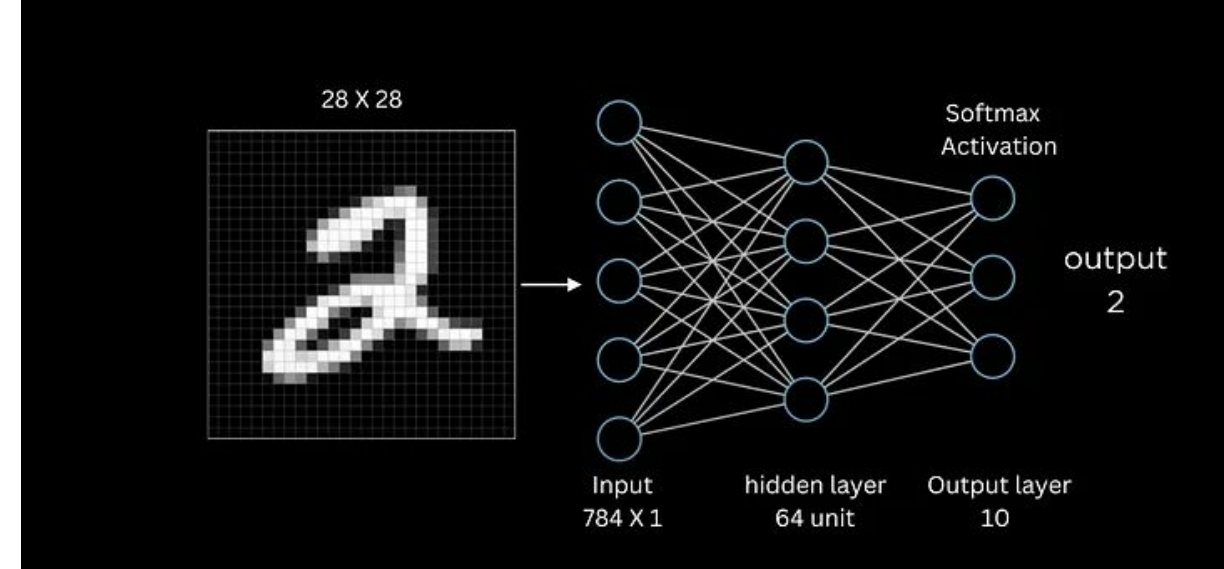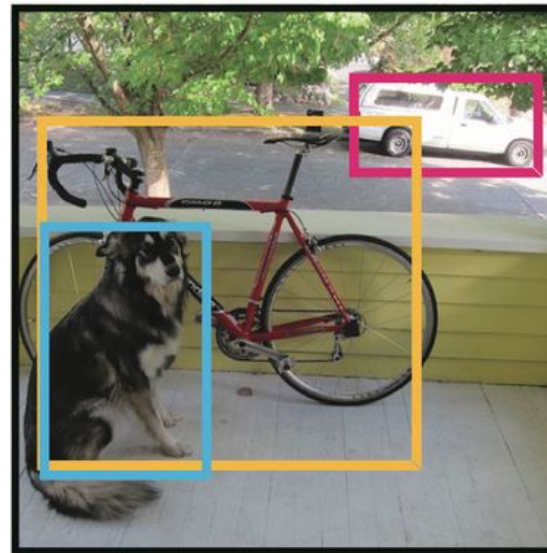
# Triggers

- **A lot of data from particle collisions. ATLAS has data volume of over 60TB/s**

- **Use triggers to save only relevant data. Has to be quick, but sensitive enough to signs of rare processes**

- **Selection based on heuristics such as energy, charge, direction, momentum**

- **Use neural networks to get better results than traditional rule-based methods?**

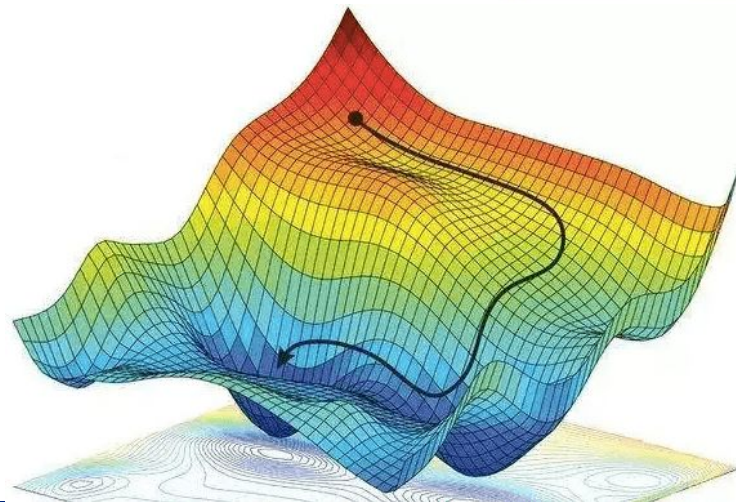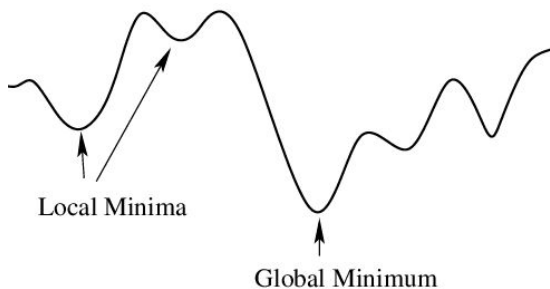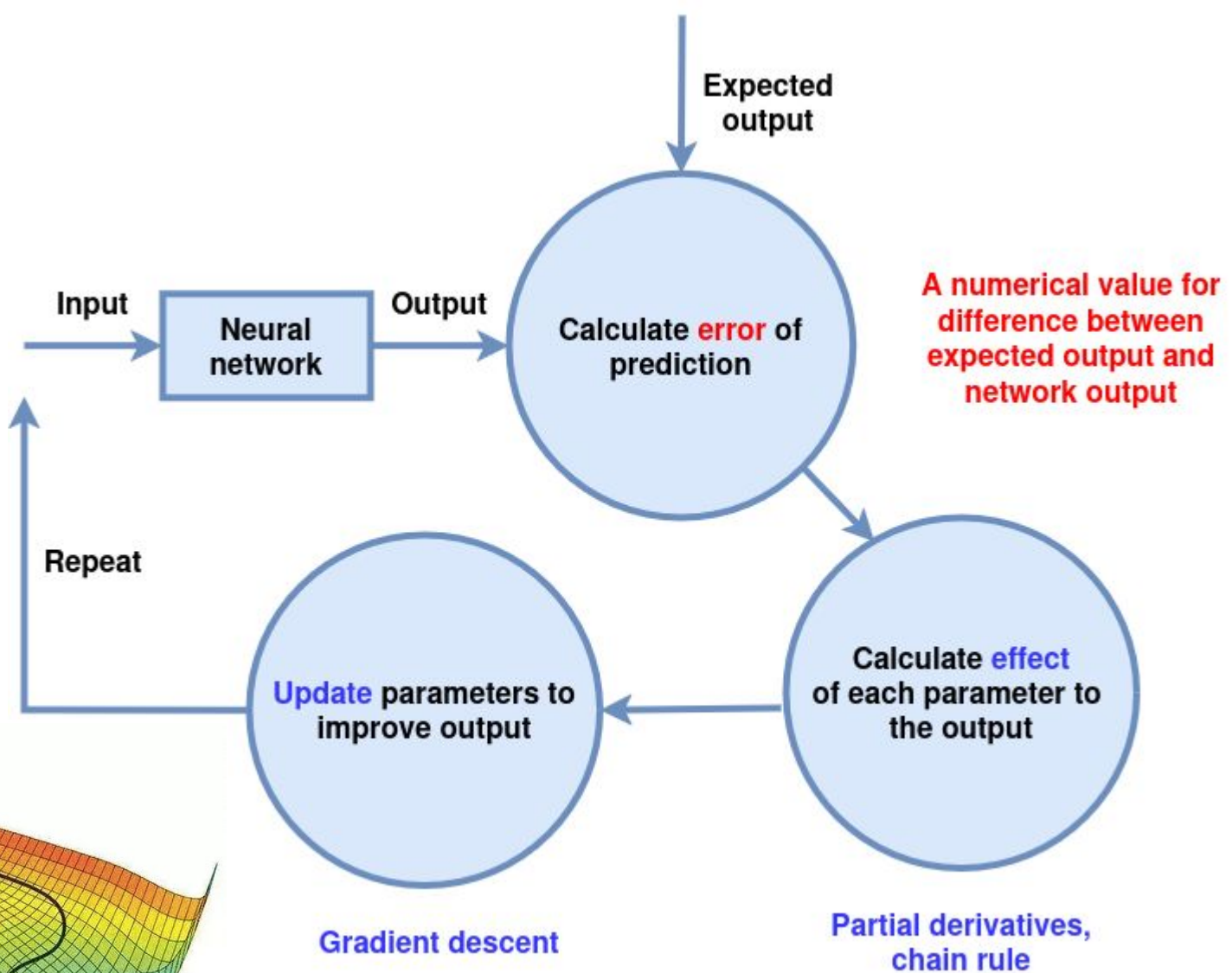- **A good neural network means nothing if it cannot be run efficiently in hardware**

# Neural networks

- **Neural networks learn from data**

- **Can be trained to do tasks such as classification, text, image or video generation, regression, pattern recognition**
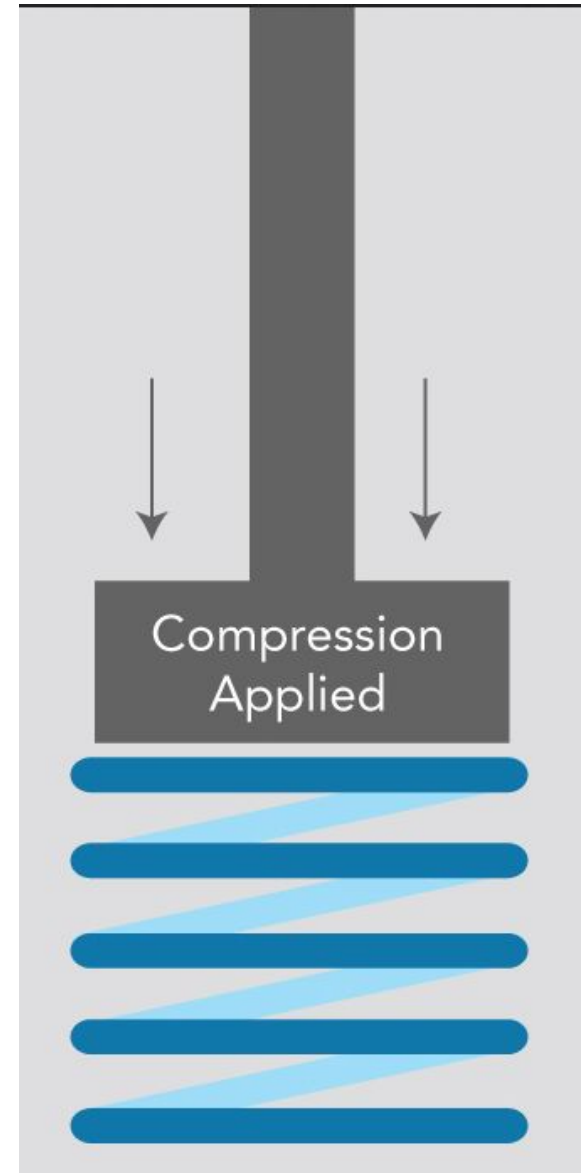
# Neural networks

- **Everything is numbers**

- **Neural networks produce an output by doing mathematical operations such as multiplication and addition**

# Neural networks

- **In the beginning, the neural network produces outputs that make no sense**

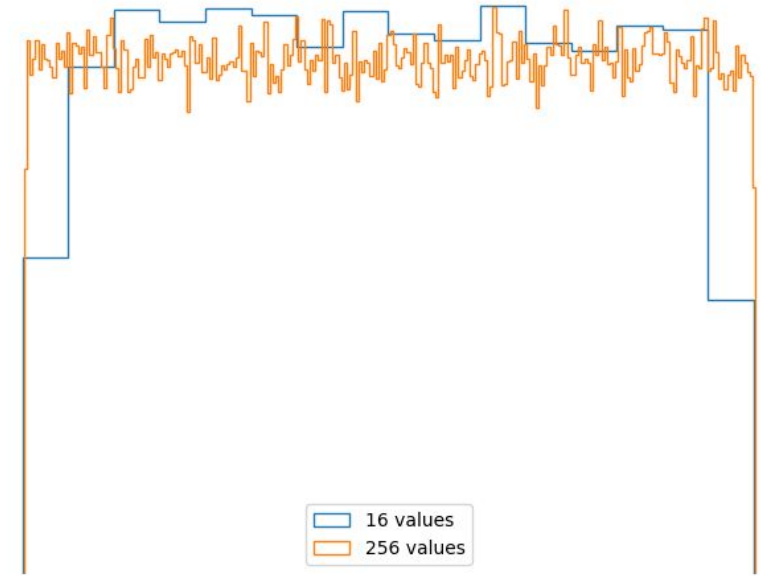- **During training the neural network learns to produce better outputs**



Input → Neural network → Output → Calculate **error** of prediction ← Expected output

A numerical value for difference between expected output and network output

Calculate **effect** of each parameter to the output

Partial derivatives, chain rule

Update parameters to improve output

Gradient descent

Repeat

Local Minima

Global Minimum

# Neural networks

- **Bigger neural networks can be more powerful, but slow**

- **In limited-resource or low-latency environments, this won't do.**

- **Use compression to make them faster and use fewer resources**



Compression Applied

# Compressing neural networks

- **Restrict parameters to be one of 2,4,16,256 etc. values**

- **Teach neural network to work with fewer parameters**

$$99 \times 85 + 93 \times 79 + 1 \times 55 = 15817$$

$$99 \times 85 + 93 \times 79 + 0 \times 55 = 15762$$

16 values
256 values

$$64\,112\,991 = 111101001001001001011011111 = 26 \text{ bits}$$
$$25\,812 = 110010011010100 = 15 \text{ bits}$$
$$54 = 110110 = 6 \text{ bits}$$

| | |
|---|---|
| min | -1 |
| max | 0.9921875 |
| mean | 0.03131510416666667 |
| std | 0.2122100147954903 |
| sparsity | 91.7% |

# Optimizing (compressed) neural networks

- **Hyperparameters: a set of parameters that configure the neural network architecture and how it is trained**

- **Neural networks are a black box. Training can take a long time. Have to wait until training ends to see how well a set of hyperparameters work ➡ automate**

- **With compressed neural networks that run directly on hardware, have to also consider requirements by hardware**

# Track reconstruction

**Offline:**
- spacepoint formation
- track seeding
- track following
- track fitting

**Online:**
- pattern recognition
- latency O(10)us

**Neural networks**



Current environment inside ATLAS at LHC
$\langle\mu\rangle = 23$

Expected environment inside ATLAS at HL-LHC
$\langle\mu\rangle = 140$

Hits — Graph Construction (1: Metric Learning or Module Map) — Graph — Edge Labeling (2: Graph Neural Network $v_b^{k+1} = \phi(e_{0j}^k, v_j^k, v_0^k)$) — Edge Scores — Graph Segmentation (3: Connected Components or Connected Components + Walkthrough) — Track Candidates

# Jet tagging

Identify the type of the particle that initiates the jet.

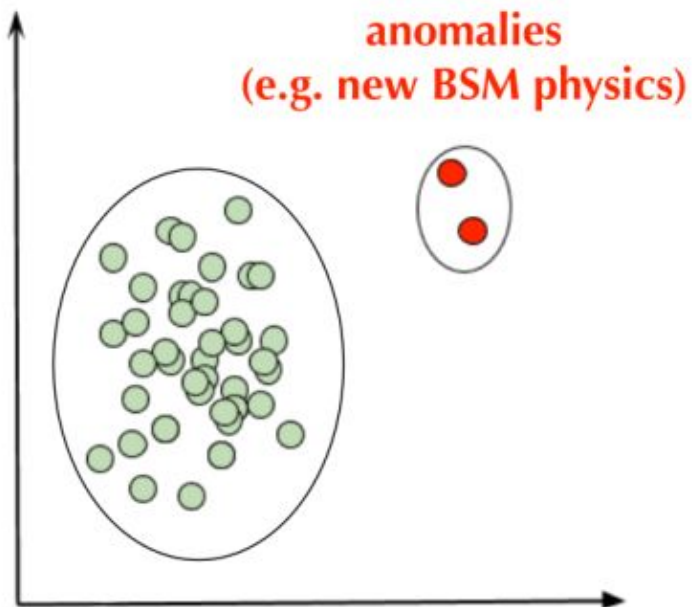Challenge: particles can radiate, radiated particles produce more particles

Neural network:
- use measured particle properties and particle pair interactions to identify particle types

# Anomaly detection

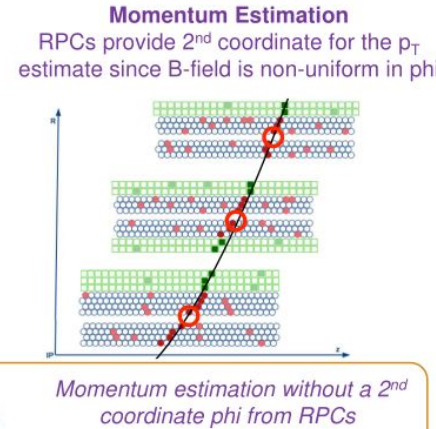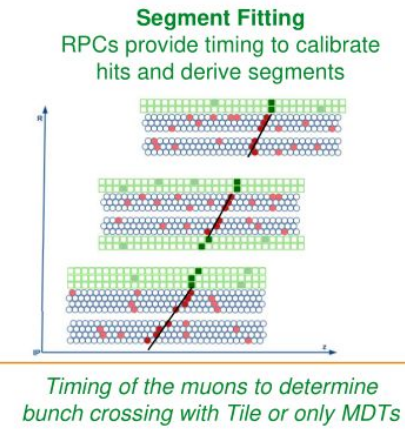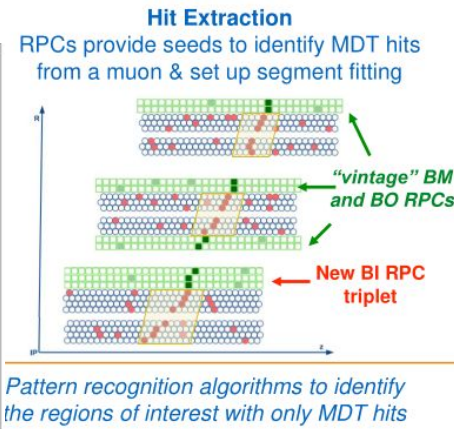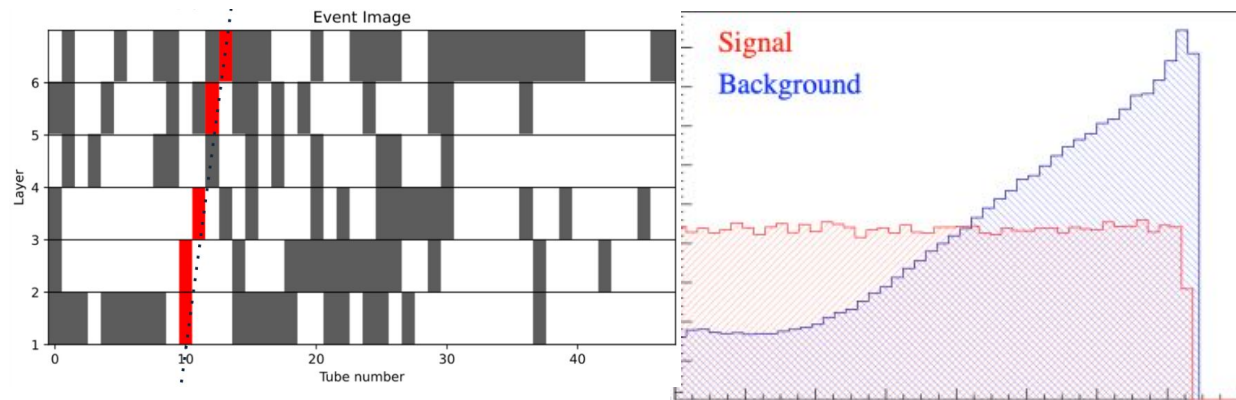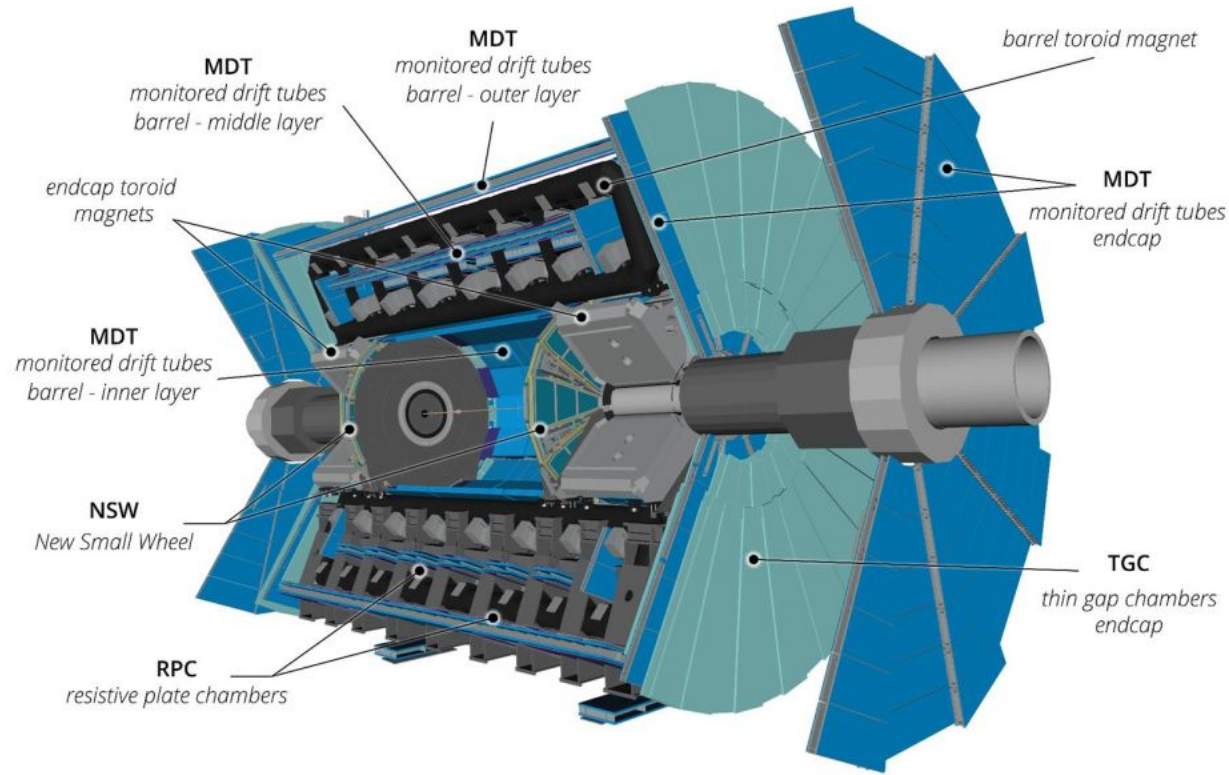- What if the trigger discards events that show new physics?
- Same input as Global Trigger, has to run in 50ns
- Find events that are very unusual



anomalies
(e.g. new BSM physics)



Anomalous Jet

$H \to bb$ Jet

CMS Experiment at the LHC, CERN
Data recorded: 2023-May-24 01:42:17.826112 GMT
Run / Event / LS: 367883 / 374187302 / 159

Most anomalous event!

# L0 Muon trigger



**For cases with reduced RPC performance:**
- Pattern recognition neural networks, distinguish muon hits from backgrounds
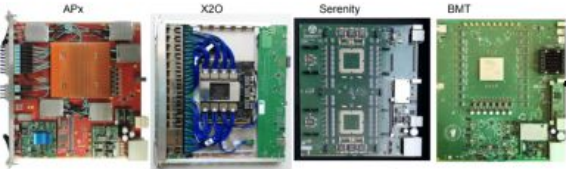- Momentum estimation

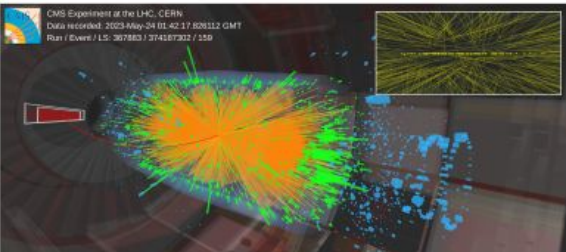# AI where?