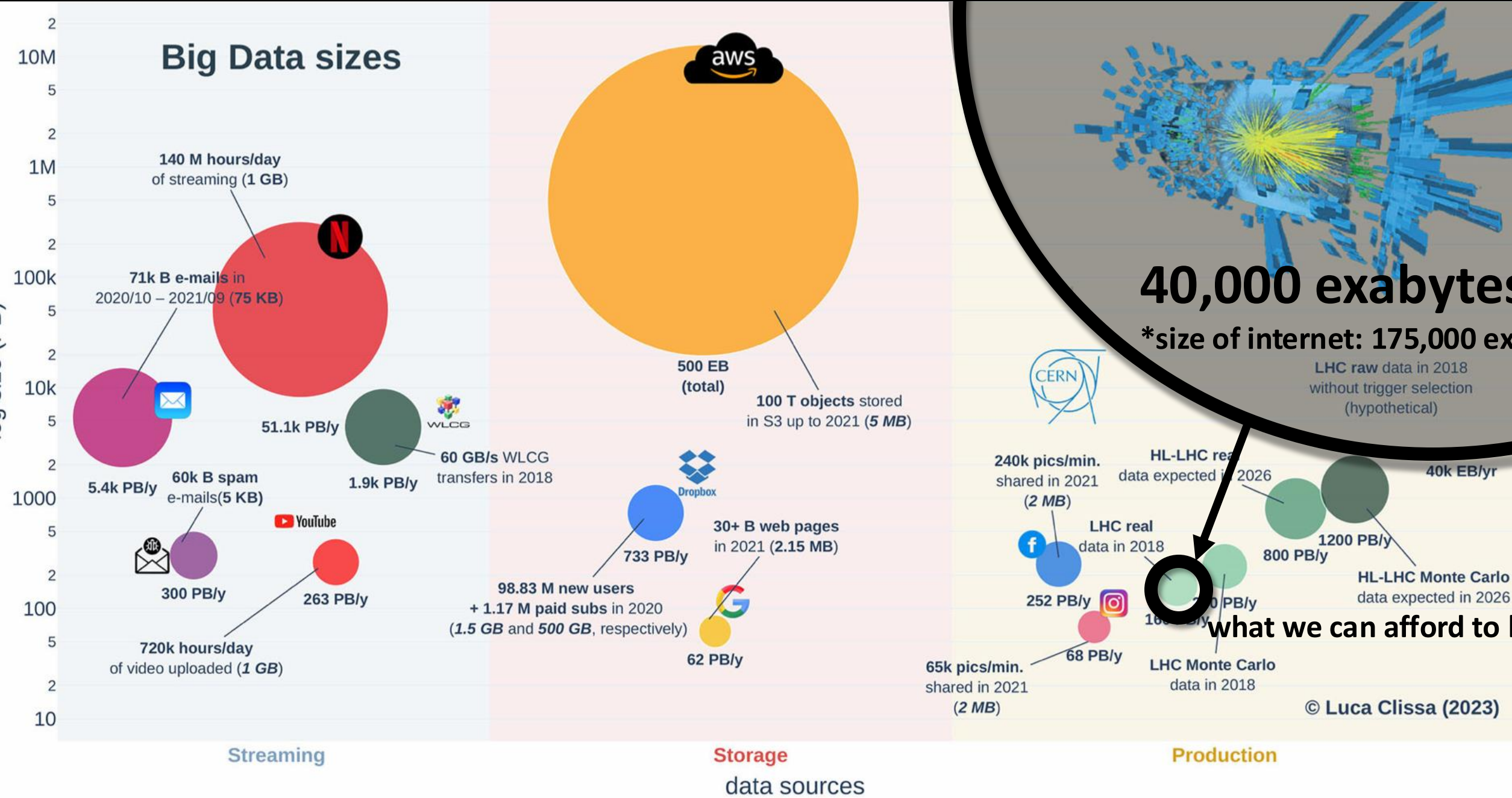
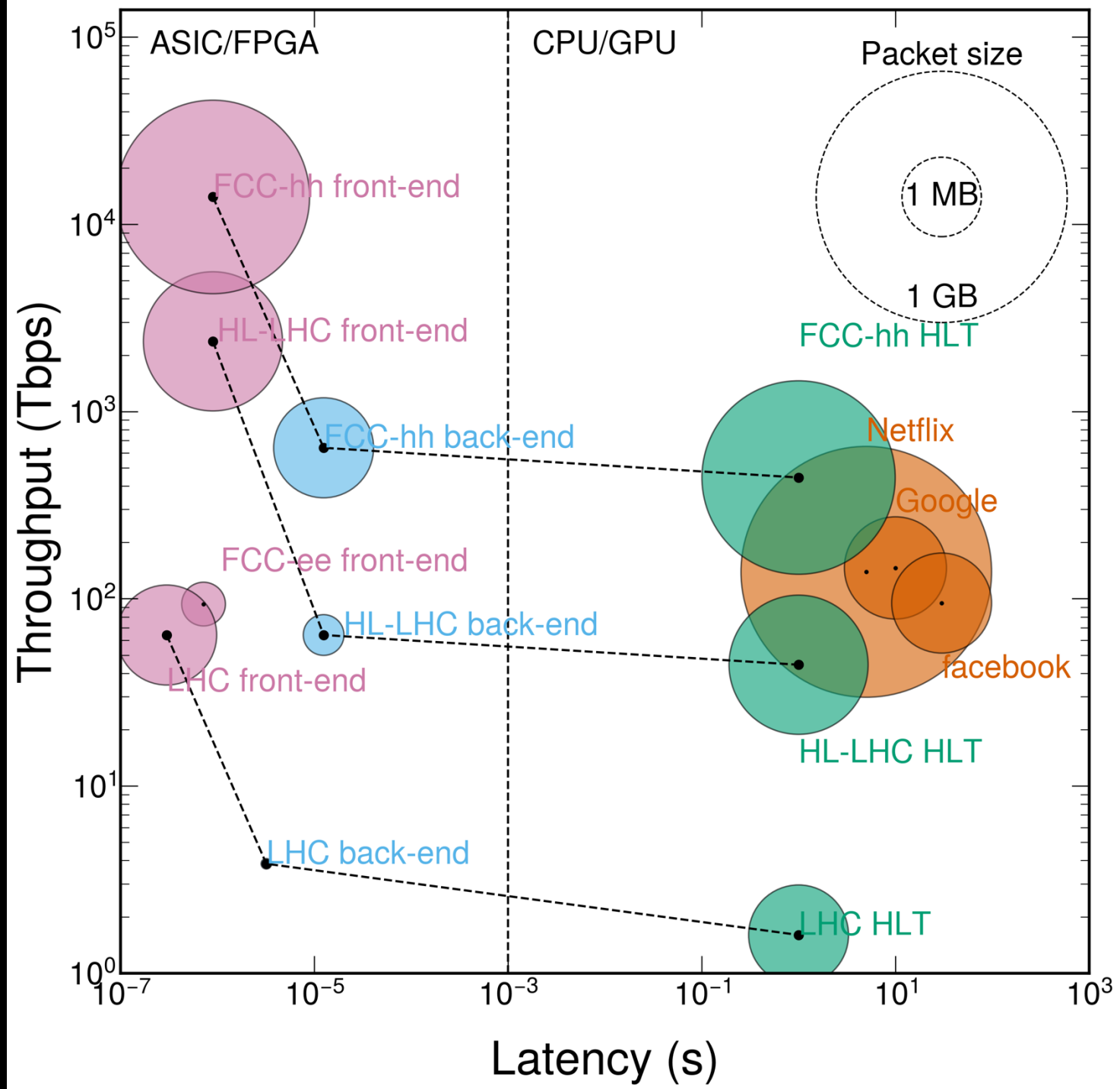




**ULTRA LOW-LATENCY INFERENCE  
ON FPGA AND ASIC FOR  
TRIGGERING AND DAO AT LHC**

Thea Klæboe Aarrestad  
(ETH Zürich)





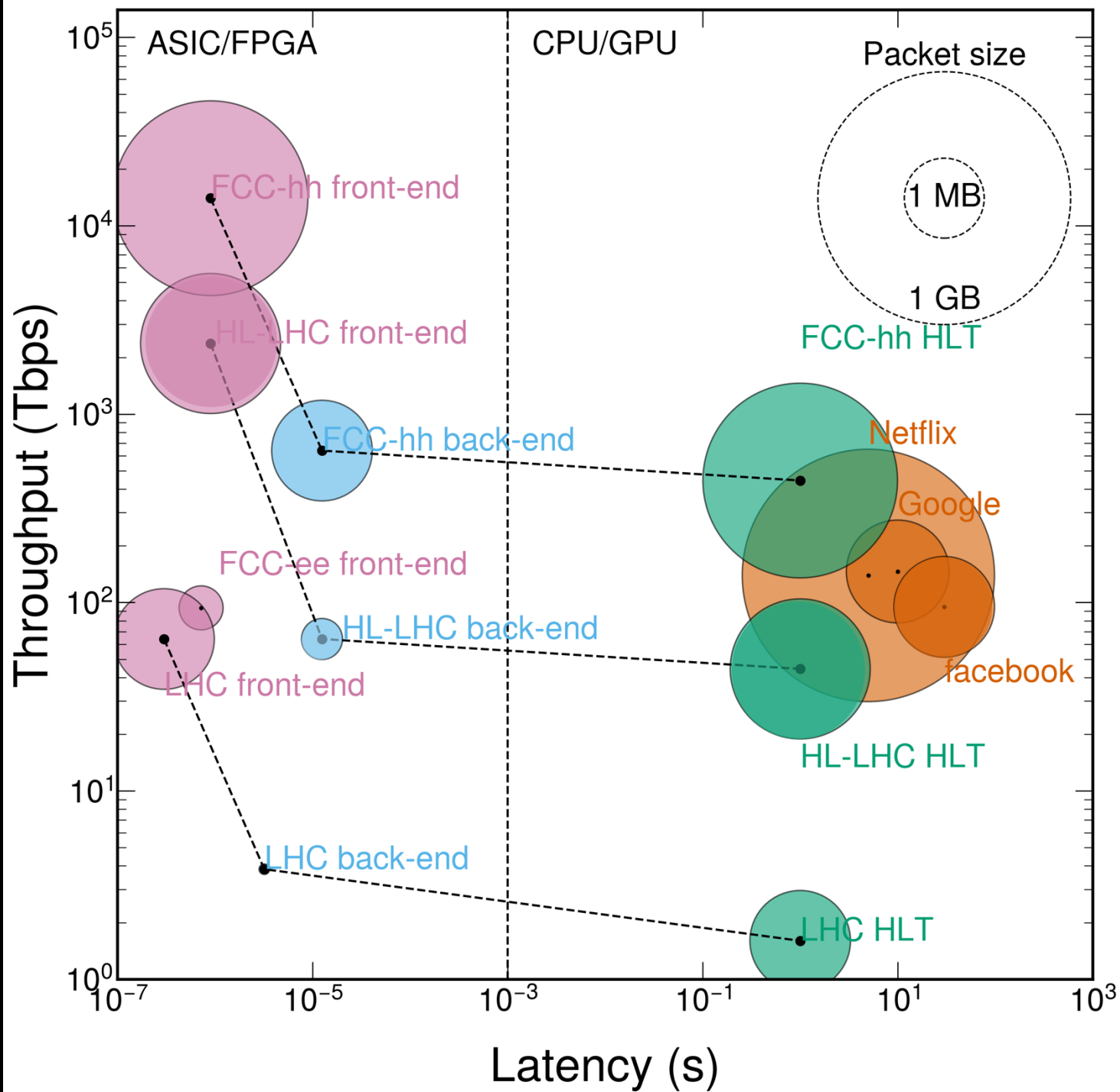
↑  
How much data goes through our system

○  
unit size of transmitted data (e.g., event size, webpage size)

→  
Time to process

ON-DETECTOR

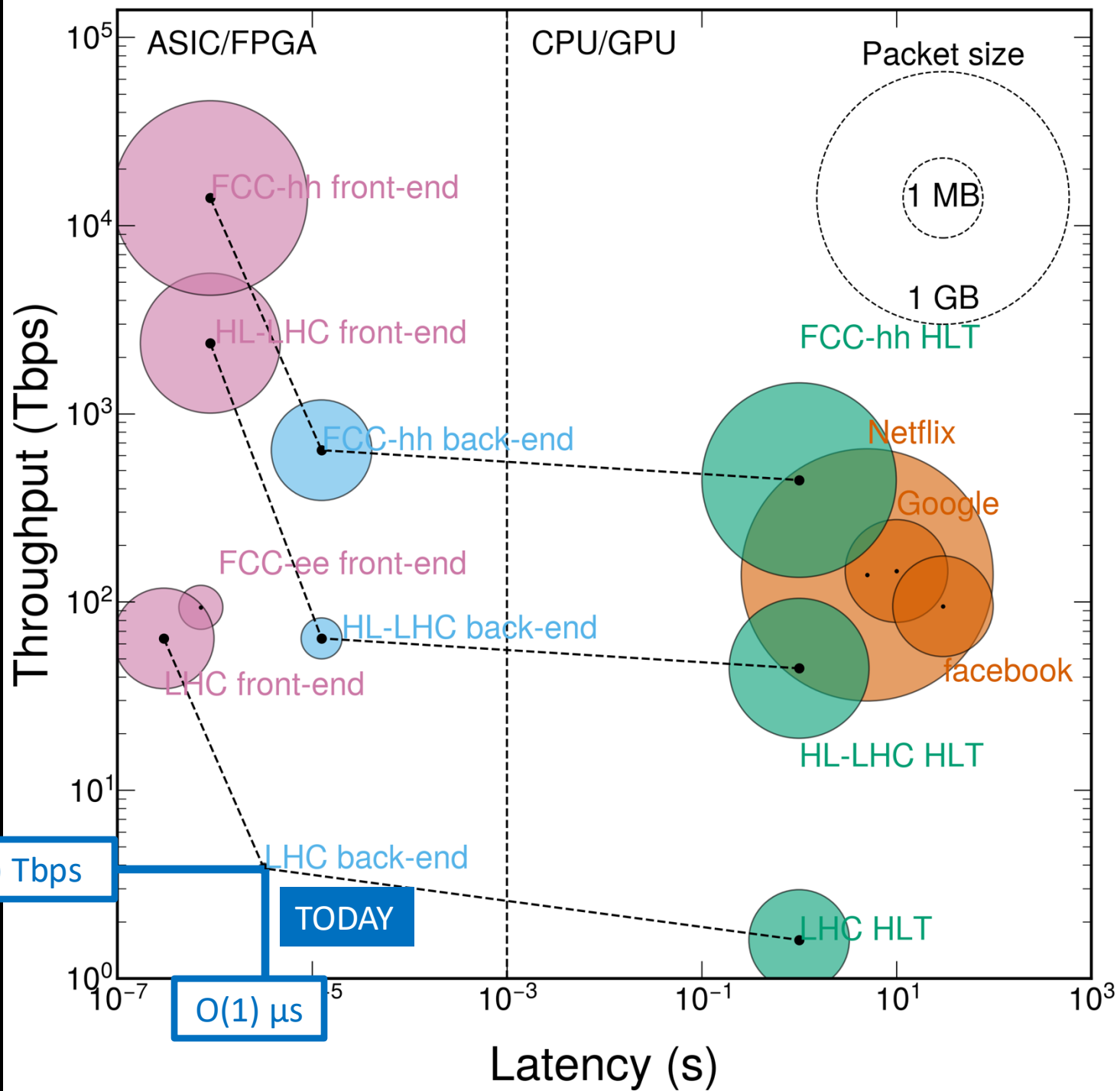
HARDWARE TRIGGERS



INDUSTRY

SOFTWARE TRIGGERS

How did we  
get to  
operating 50  
nanosecond  
ML triggers?

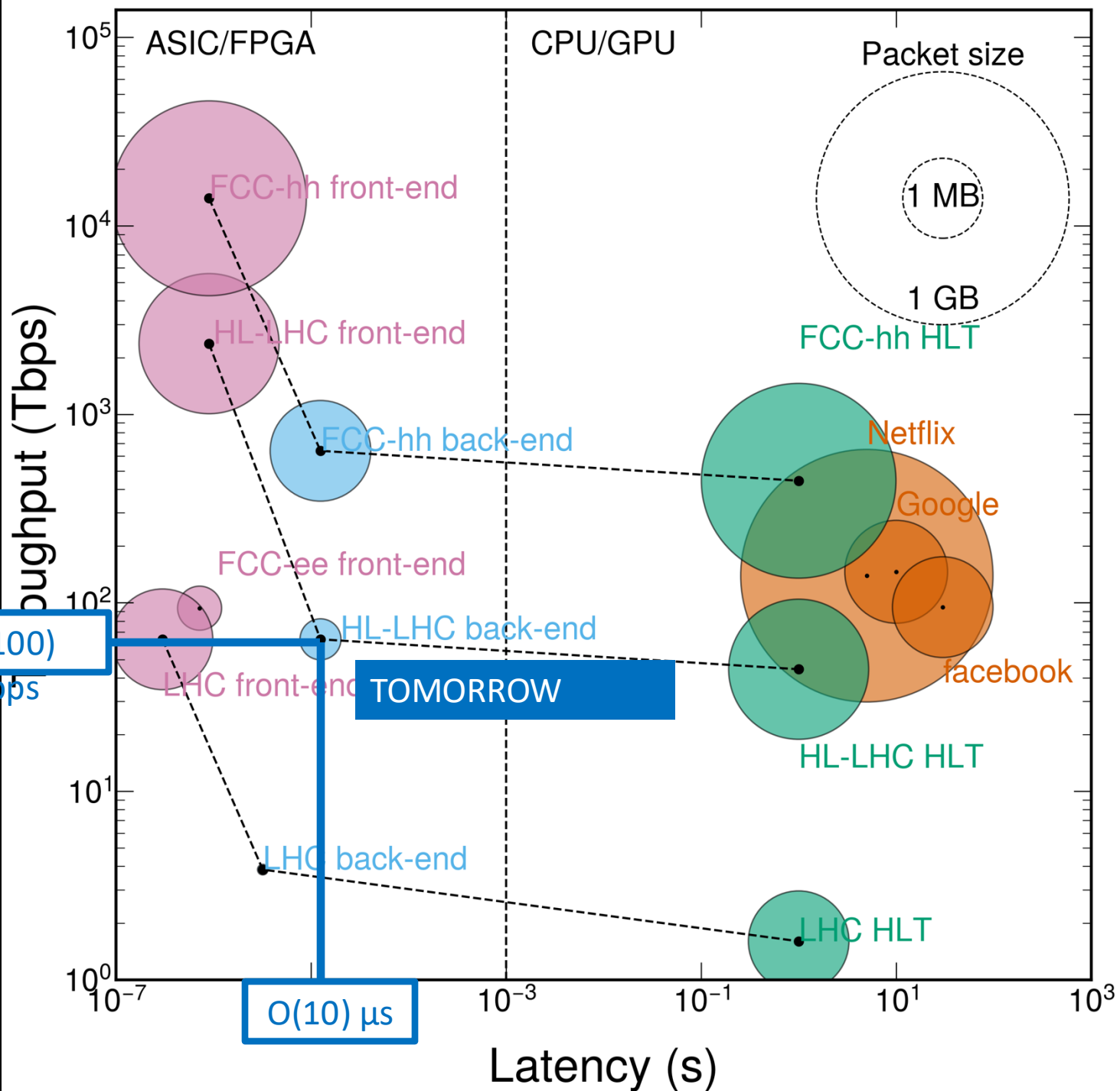


How will we maintain physics performance with x10 higher throughput?

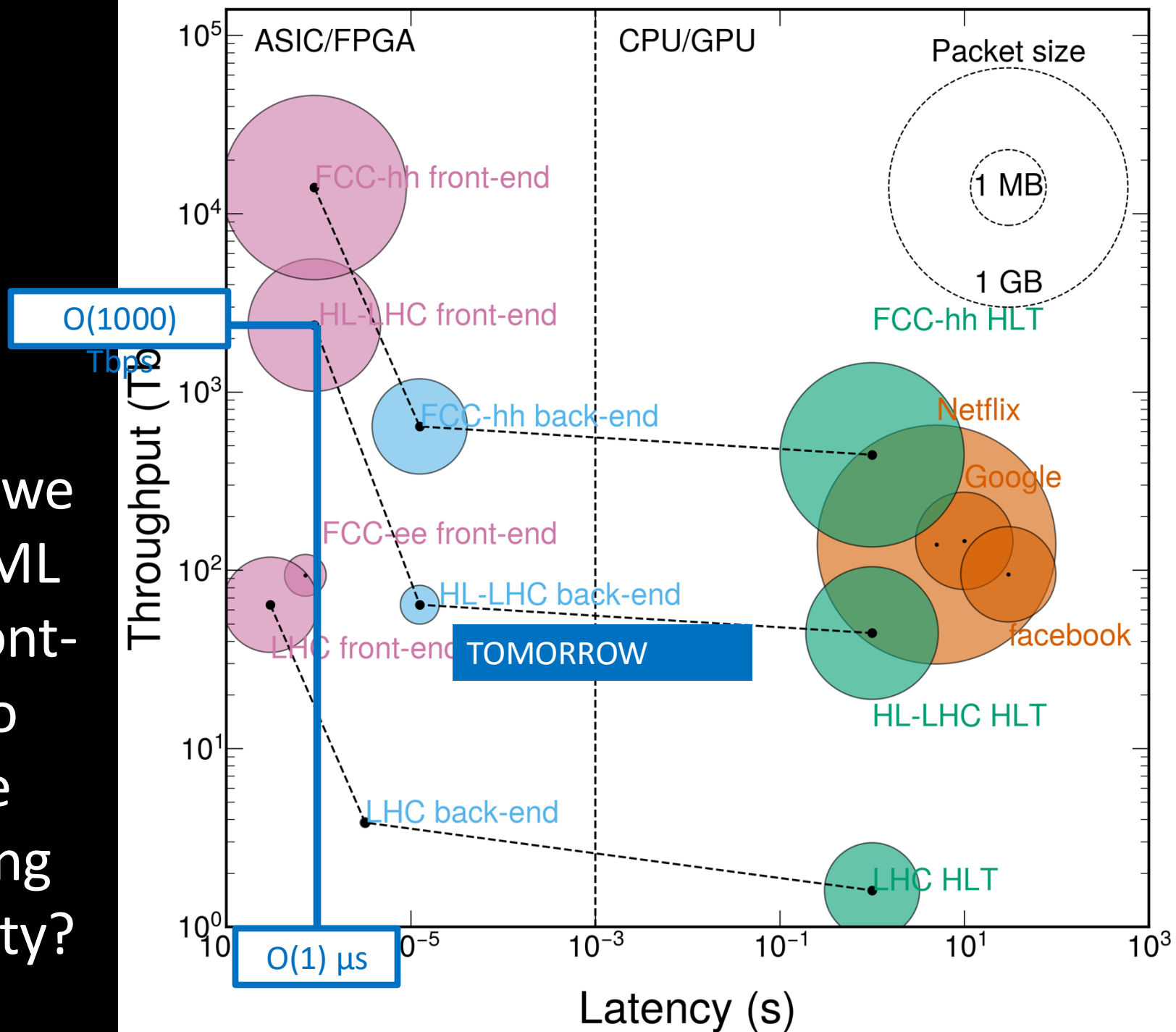
O(100) Tbps

O(10)  $\mu$ s

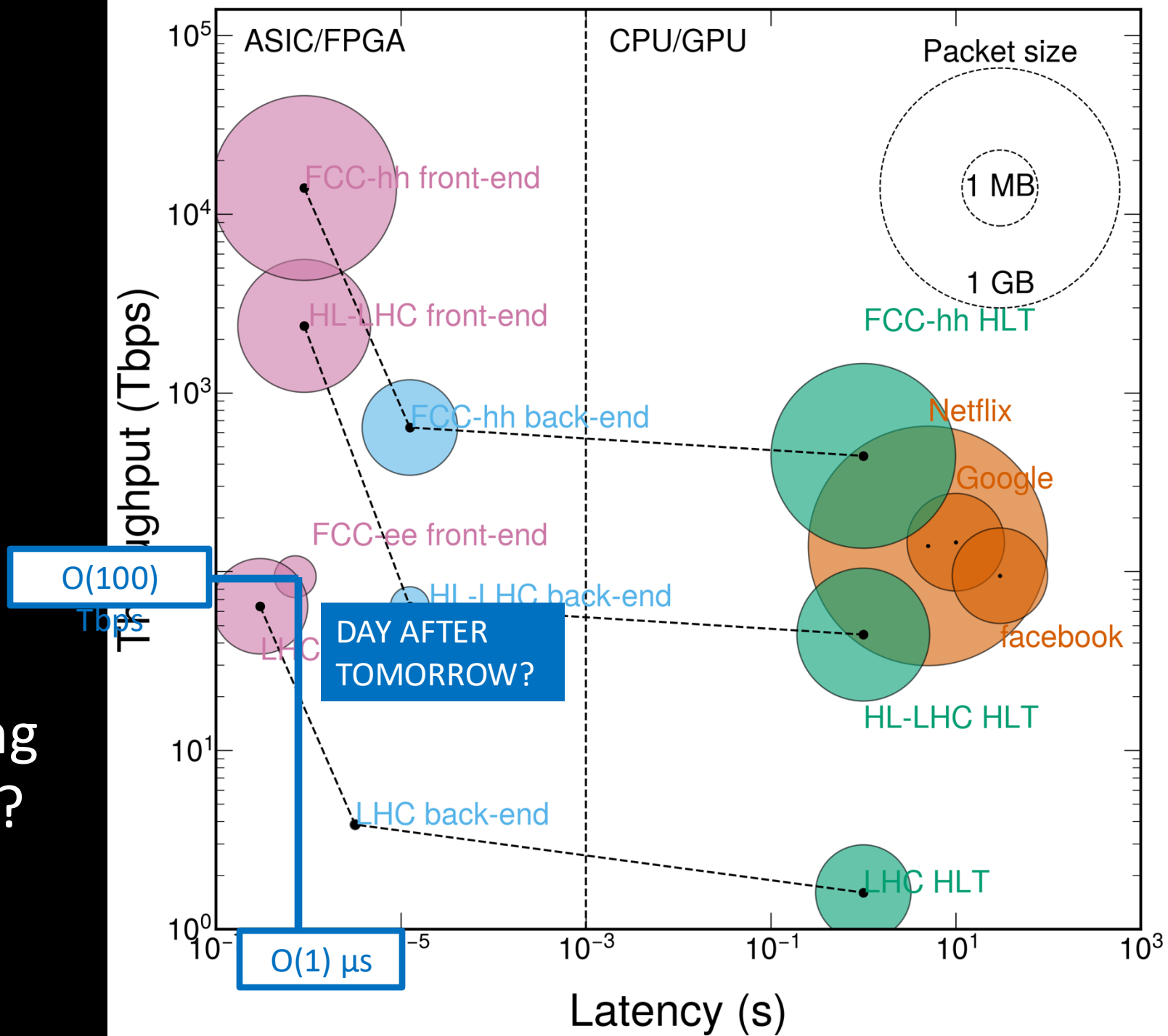
TOMORROW

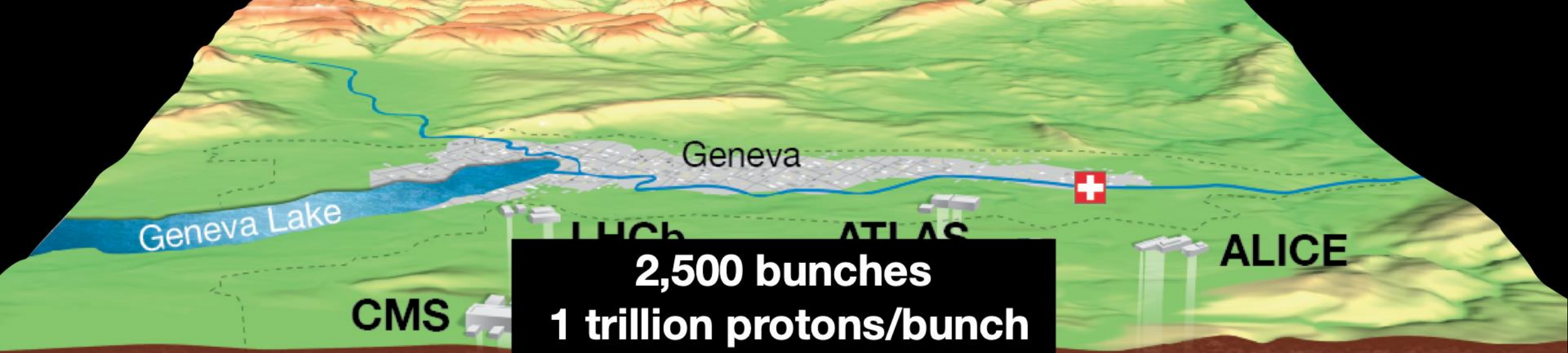


How are we getting ML in our front-ends to handle increasing granularity?

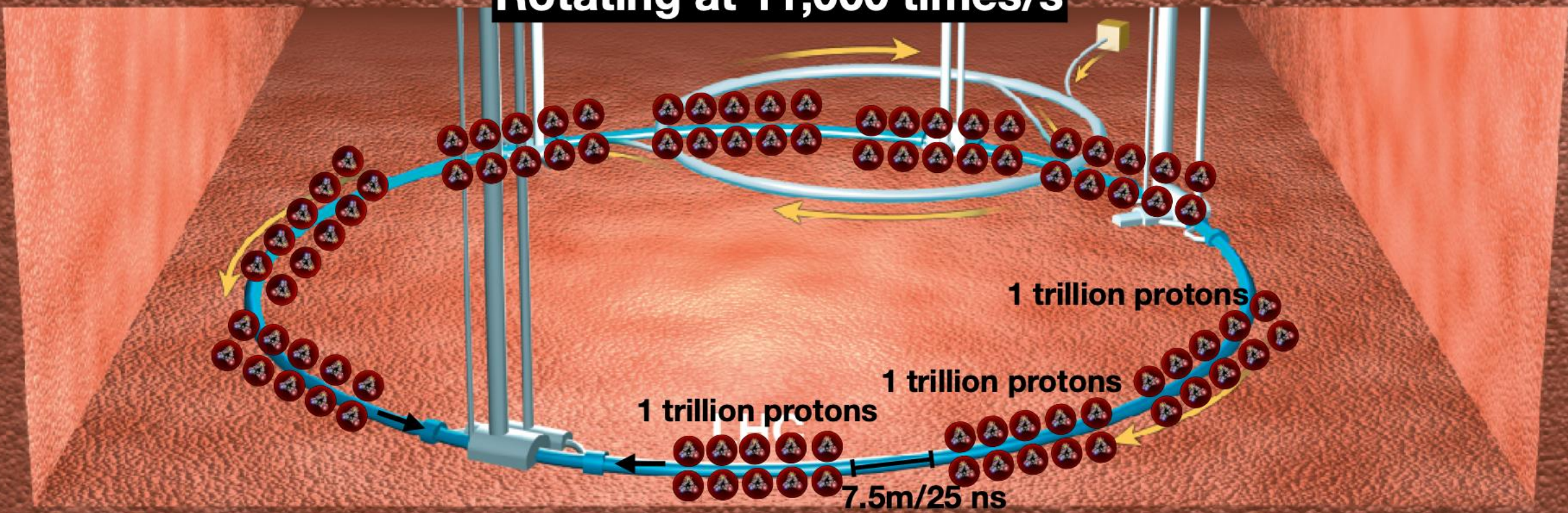


Streaming readout?





**2,500 bunches**  
**1 trillion protons/bunch**  
**Rotating at 11,000 times/s**





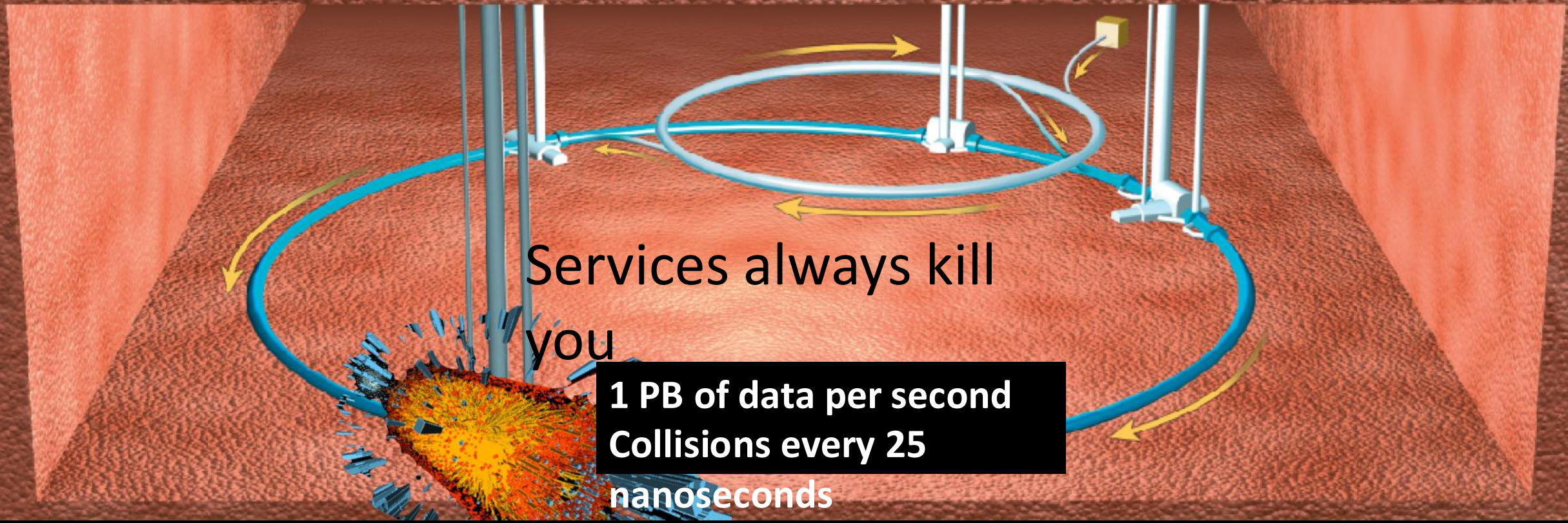
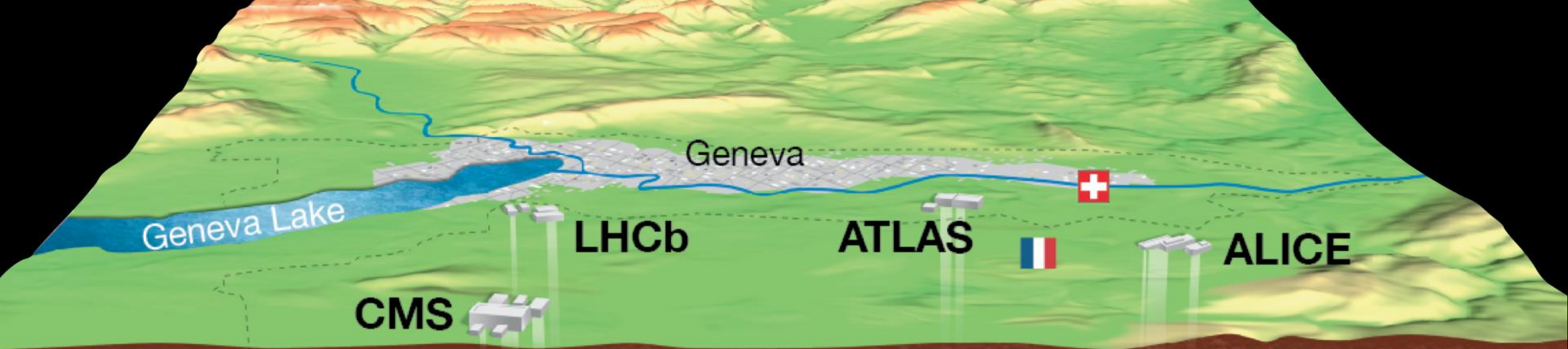
# CMS Experiment at the LHC, CERN

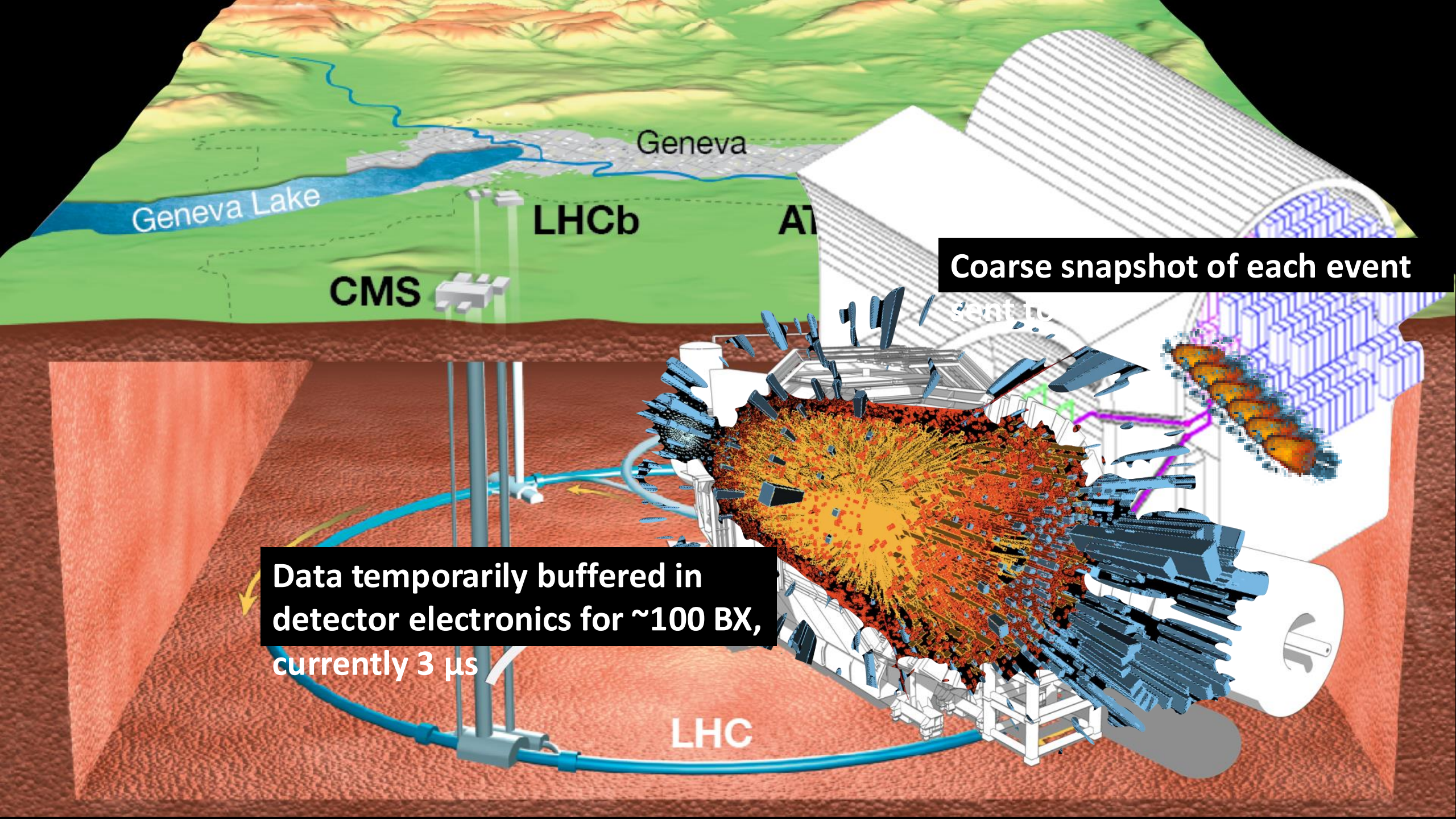
Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST)

Run / Event: 151076 / 1405388

**1 MB of data/collision**  
**1 billion**  
**collisions/second**  
**→ 1 PB of data/second**







Geneva Lake

Geneva

LHCb

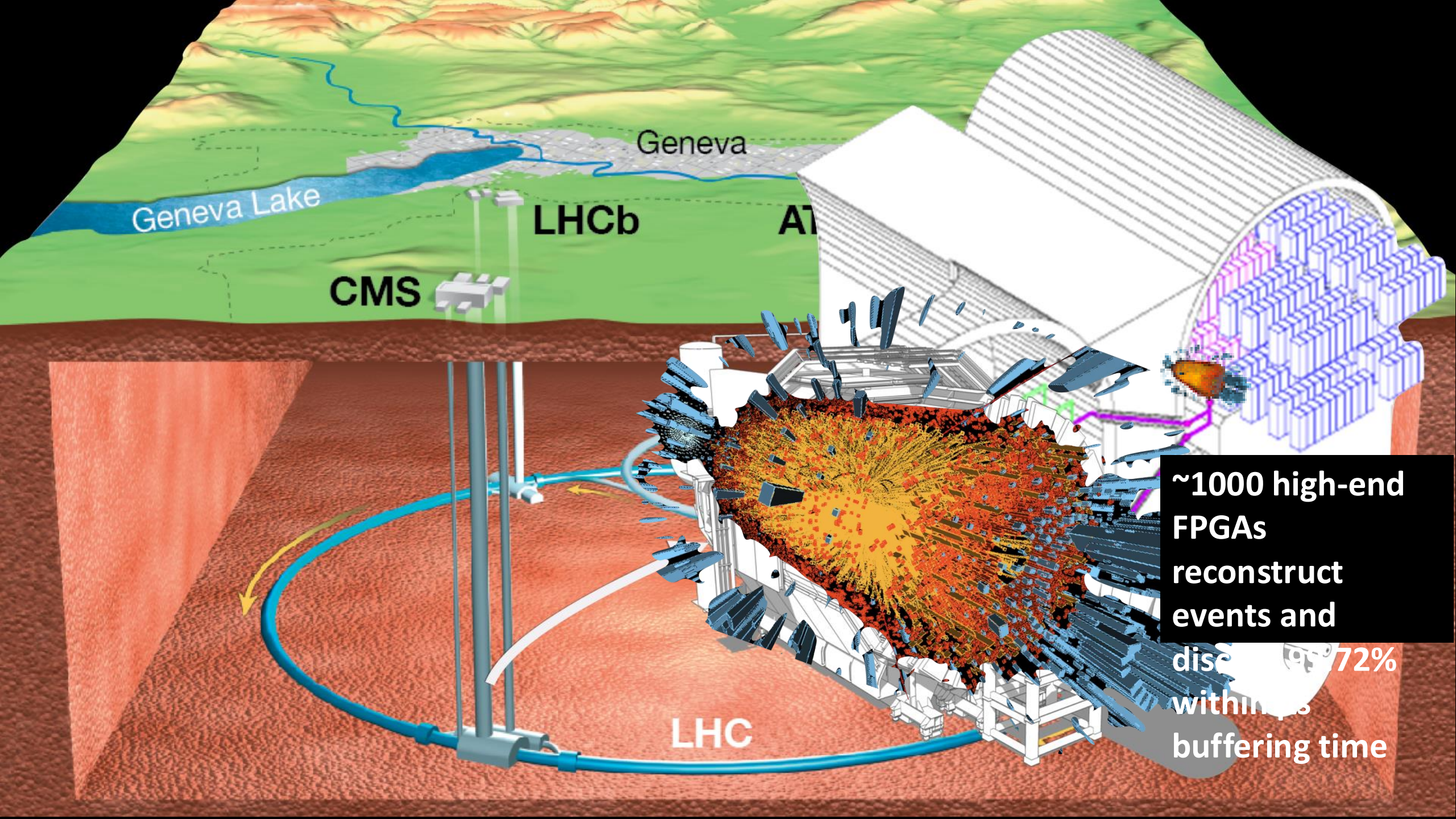
ATLAS

CMS

Coarse snapshot of each event

Data temporarily buffered in detector electronics for ~100 BX, currently 3  $\mu$ s

LHC



~1000 high-end  
FPGAs  
reconstruct  
events and  
discuss 99.72%  
within 100 ns  
buffering time

Geneva Lake

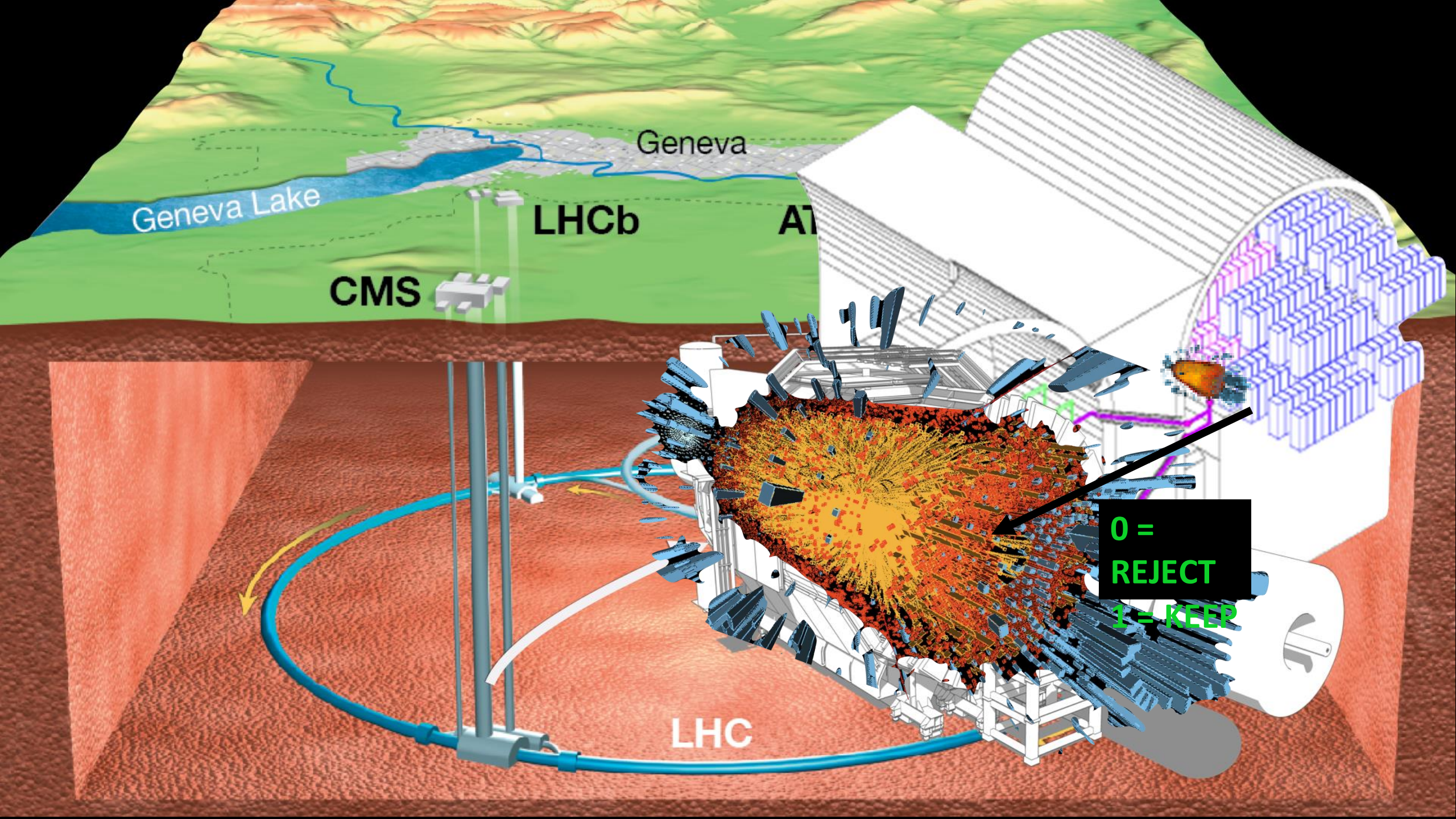
Geneva

CMS

LHCb

ATLAS

LHC



Geneva

Geneva Lake

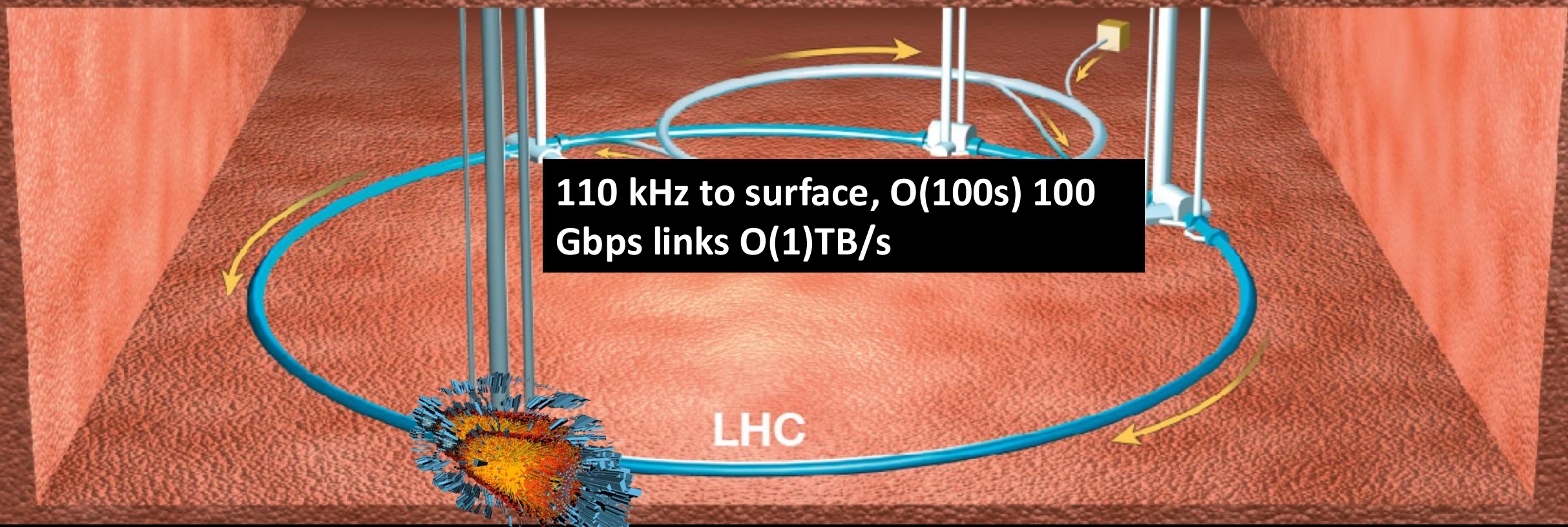
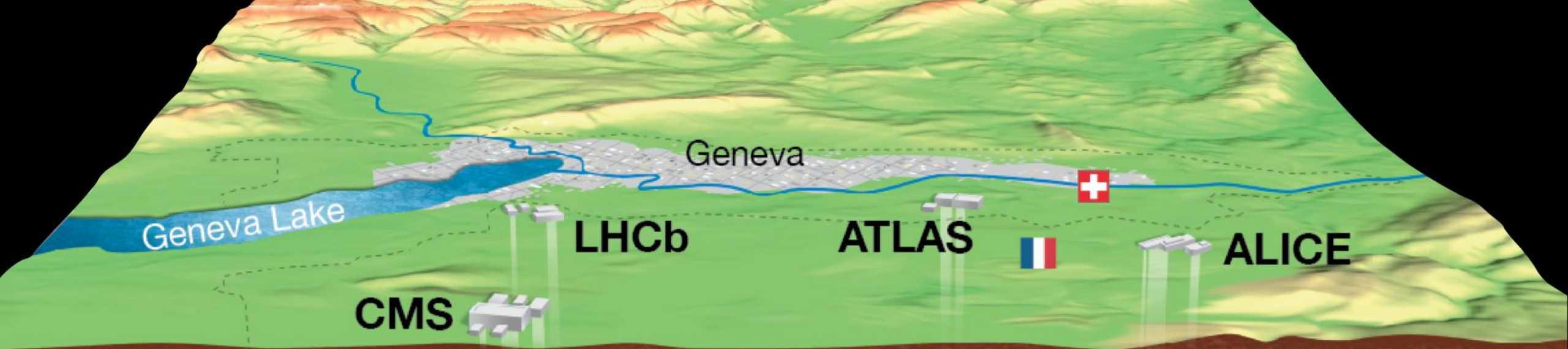
LHCb

ATLAS

CMS

LHC

0 =  
REJECT  
1 = KEEP



**HLT Software trigger:  
25,600 CPUs and 400  
GPUs**

**Offline  
reconstruction and  
permanent storage**

Reduce rate to  
100 MB/s

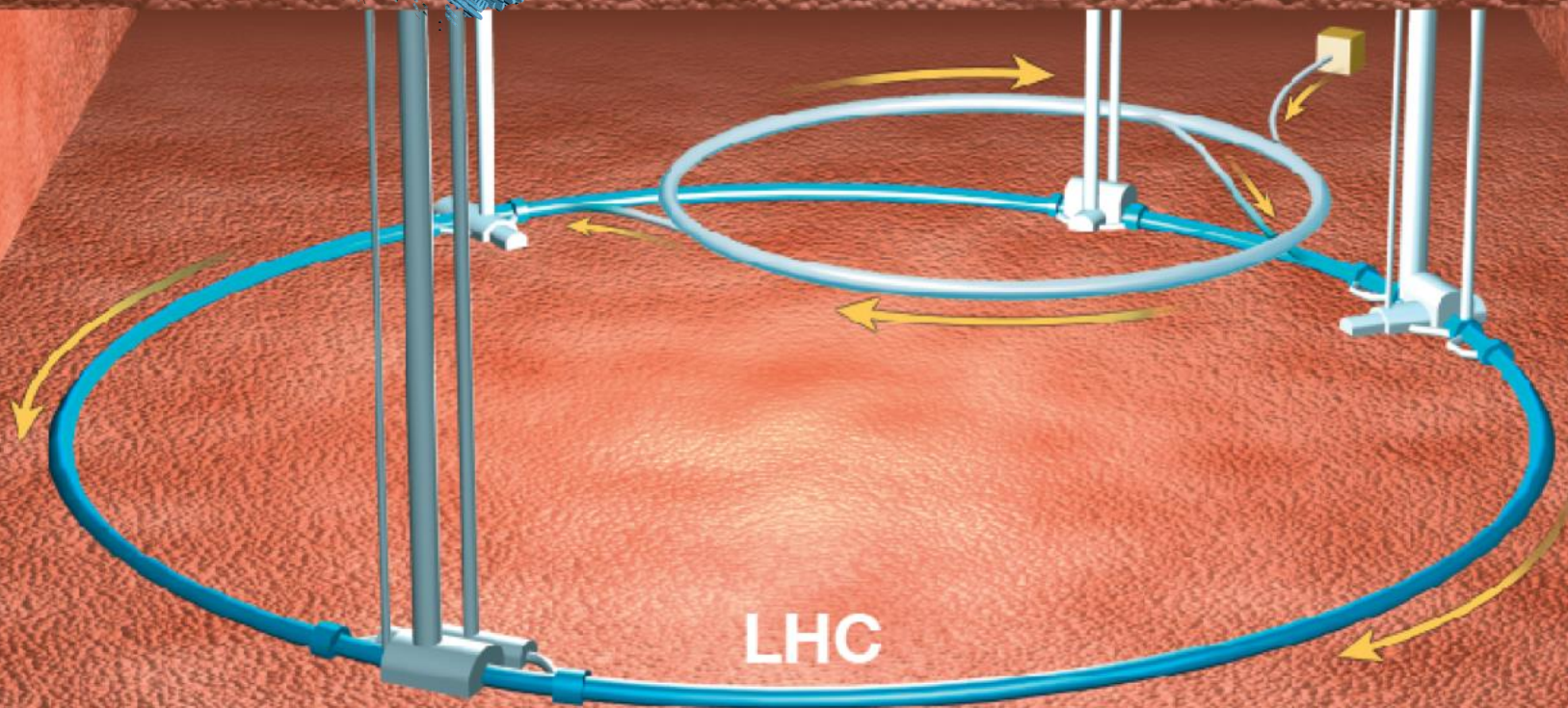
Geneva

ATLAS



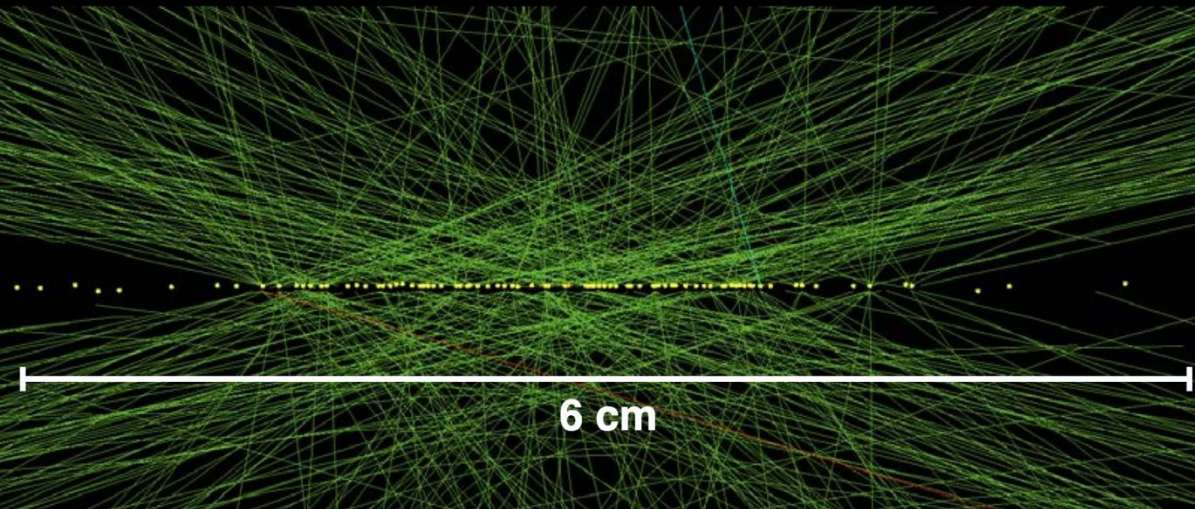
ALICE

Geneva La



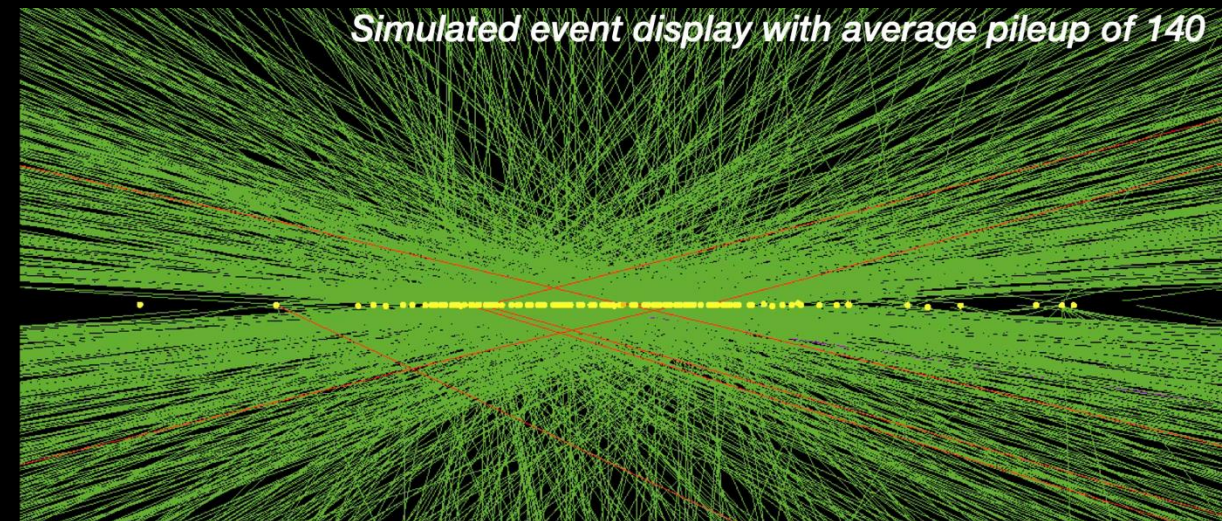
LHC

78 vertices  
(average 60)

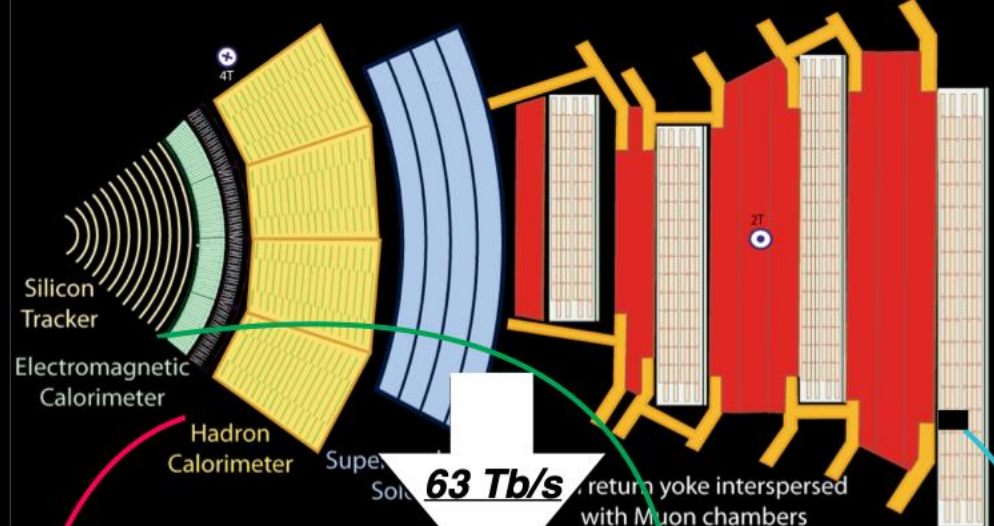


HL-LHC

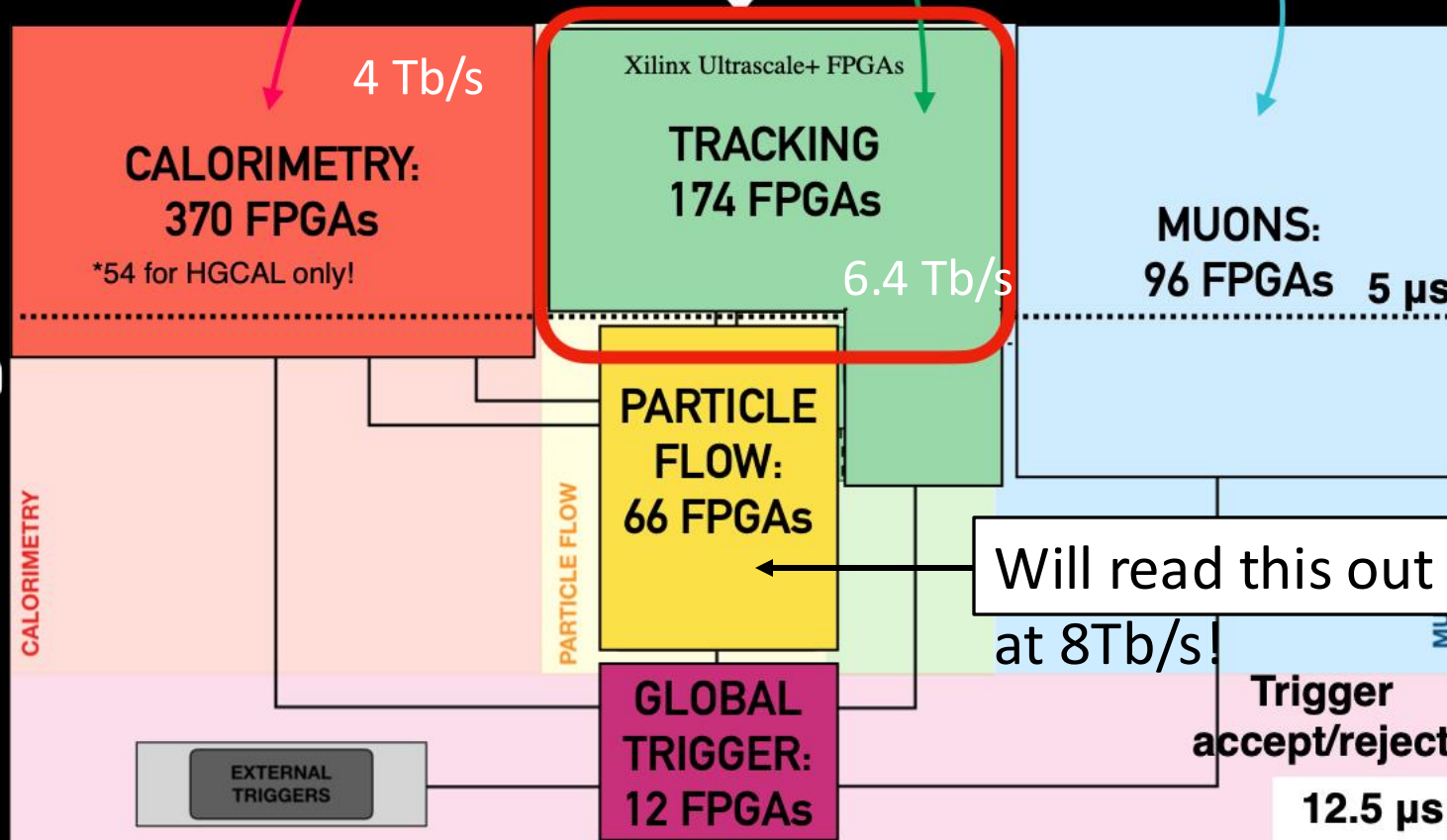
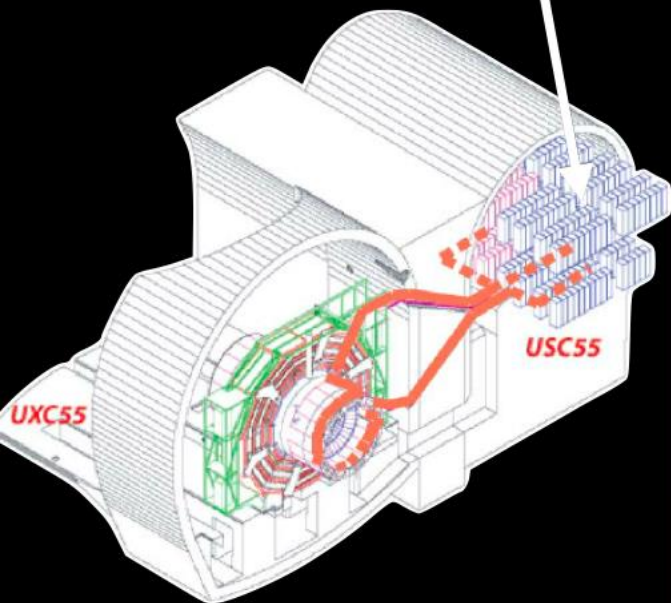
200 vertices (average  
140)

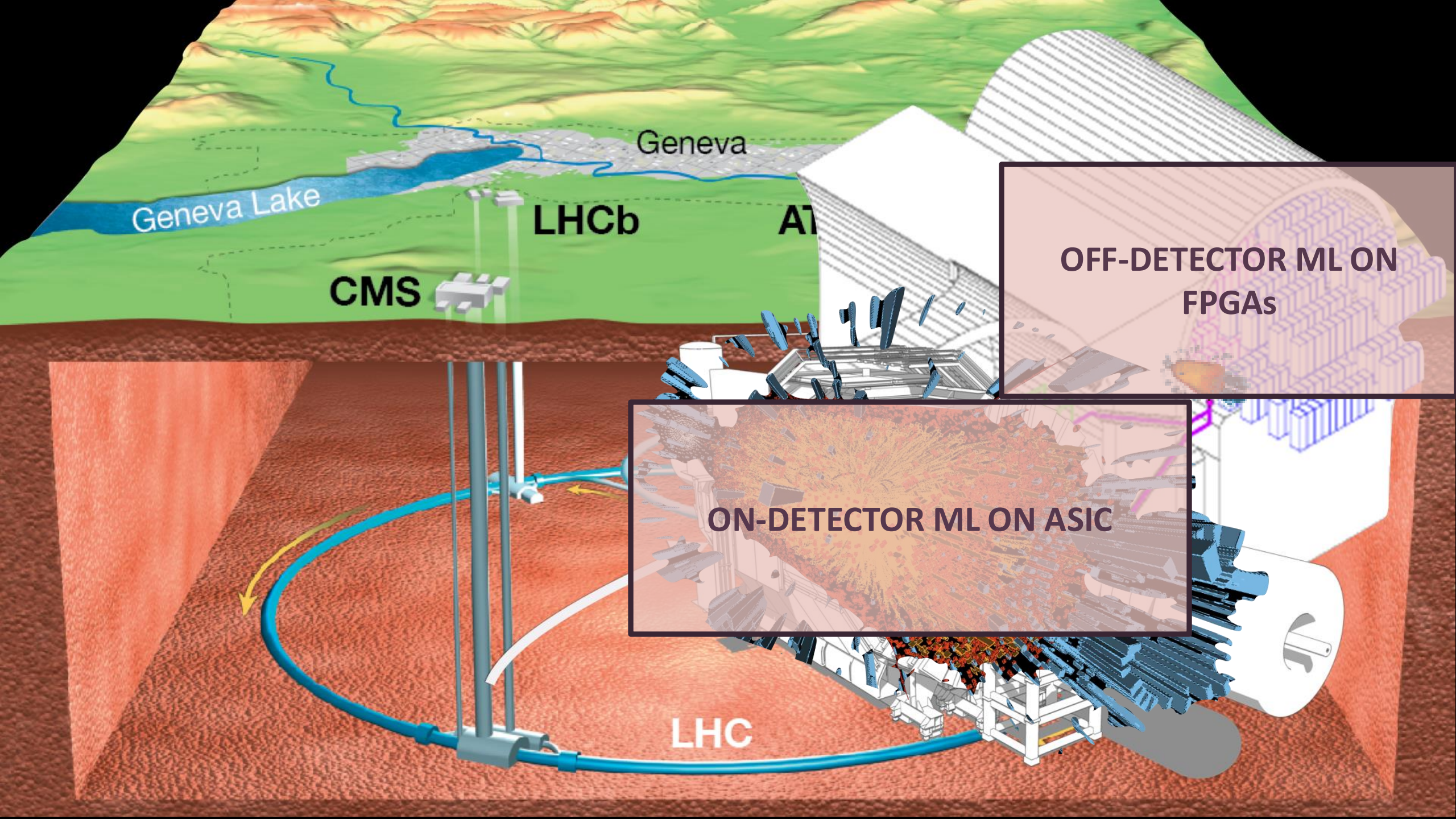


Need to upgrade detectors and algorithms to  
maintain physics performance in challenging  
environment



Input: 2 Tb/s  $\rightarrow$  63 Tb/s  
 Latency: 4  $\mu$ s  $\rightarrow$  12  $\mu$ s  
 Output: 100kHz  $\rightarrow$  750 kHz





Geneva Lake

Geneva

CMS

LHCb

ATLAS

OFF-DETECTOR ML ON  
FPGAs

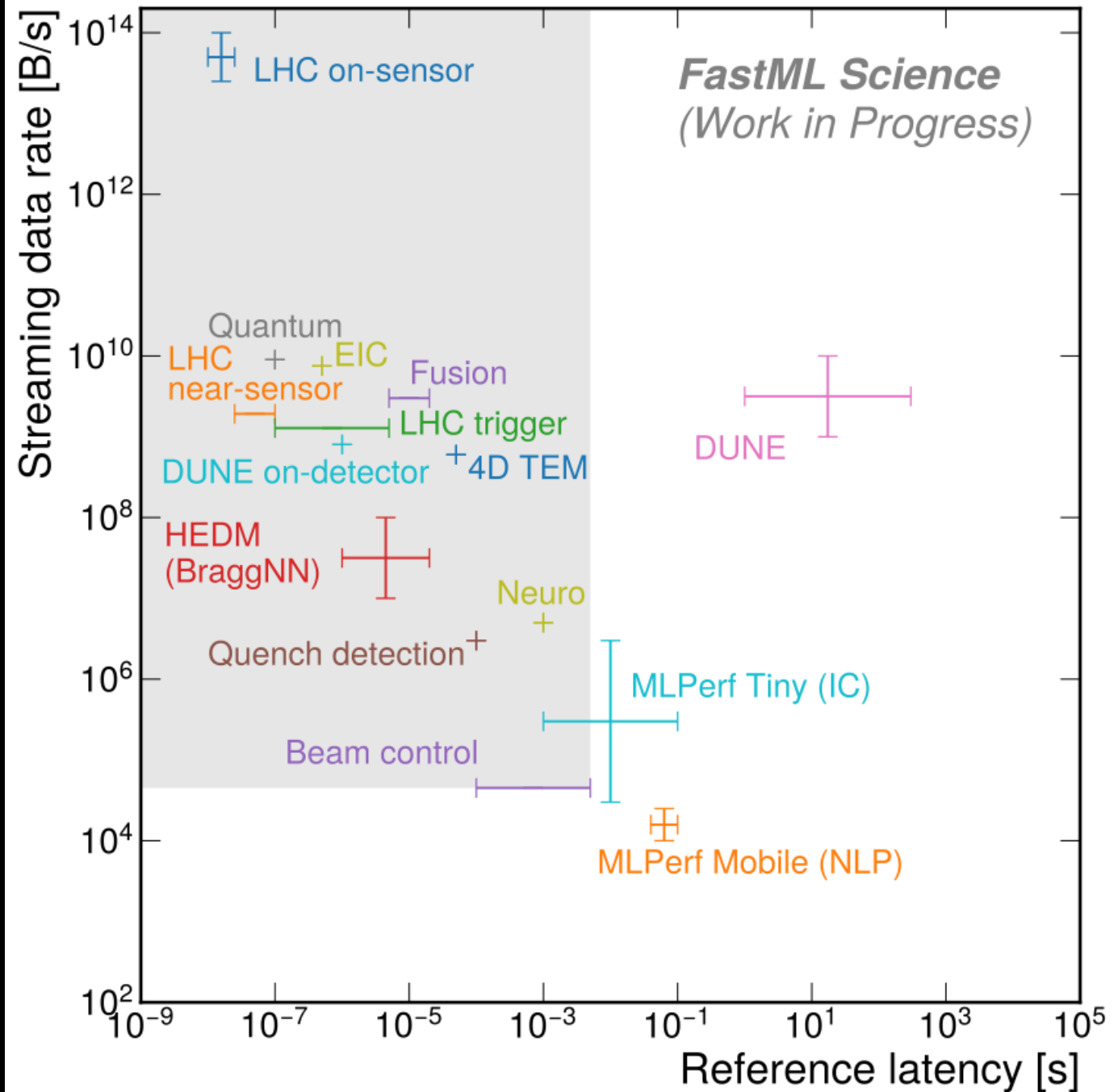
ON-DETECTOR ML ON ASIC

LHC

# Fast Machine Learning for HEP

Our tasks not well-represented by industry-driven benchmarks

Must develop our own tools, datasets and data challenges



# ML inference at low latency

## HEP Tools and Communities



Co-processing kernel  
(Xilinx accelerators/SoCs)

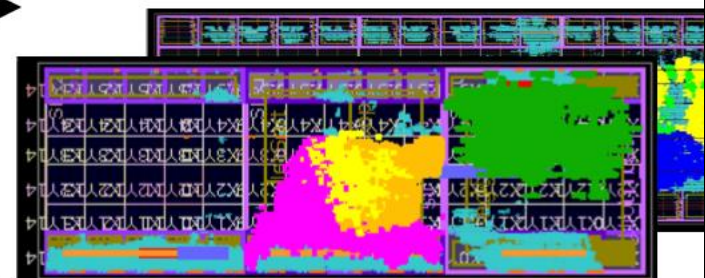
Model  
(quantized/pruned)



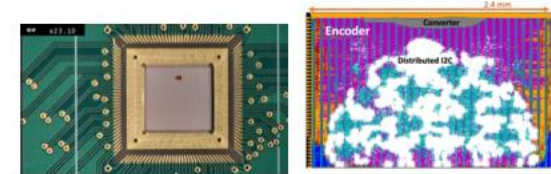
**HEP hardware ML libraries:**

hls4ml  
Conifer

FPGA custom designs  
(eg trigger algorithms)



ASICs

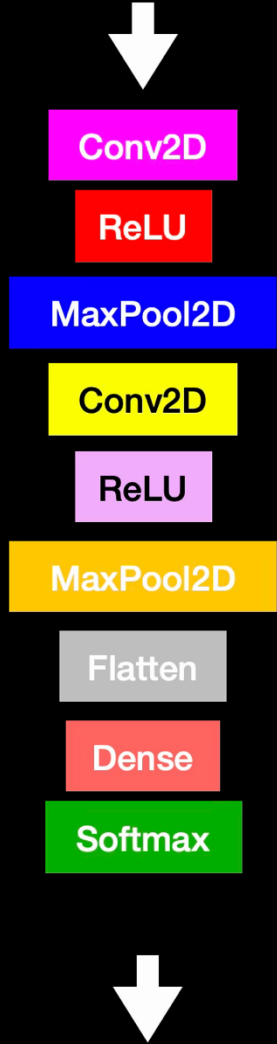
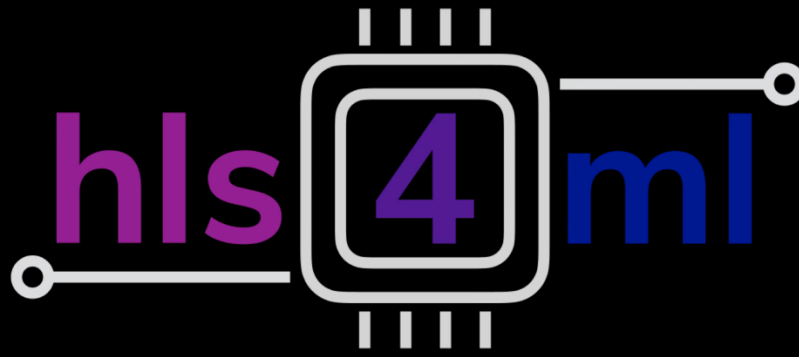
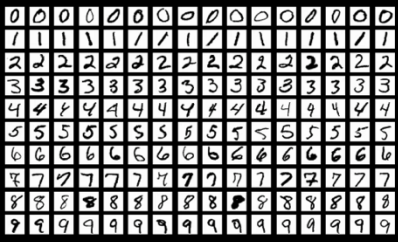


HEP quantization libraries:



HGQ

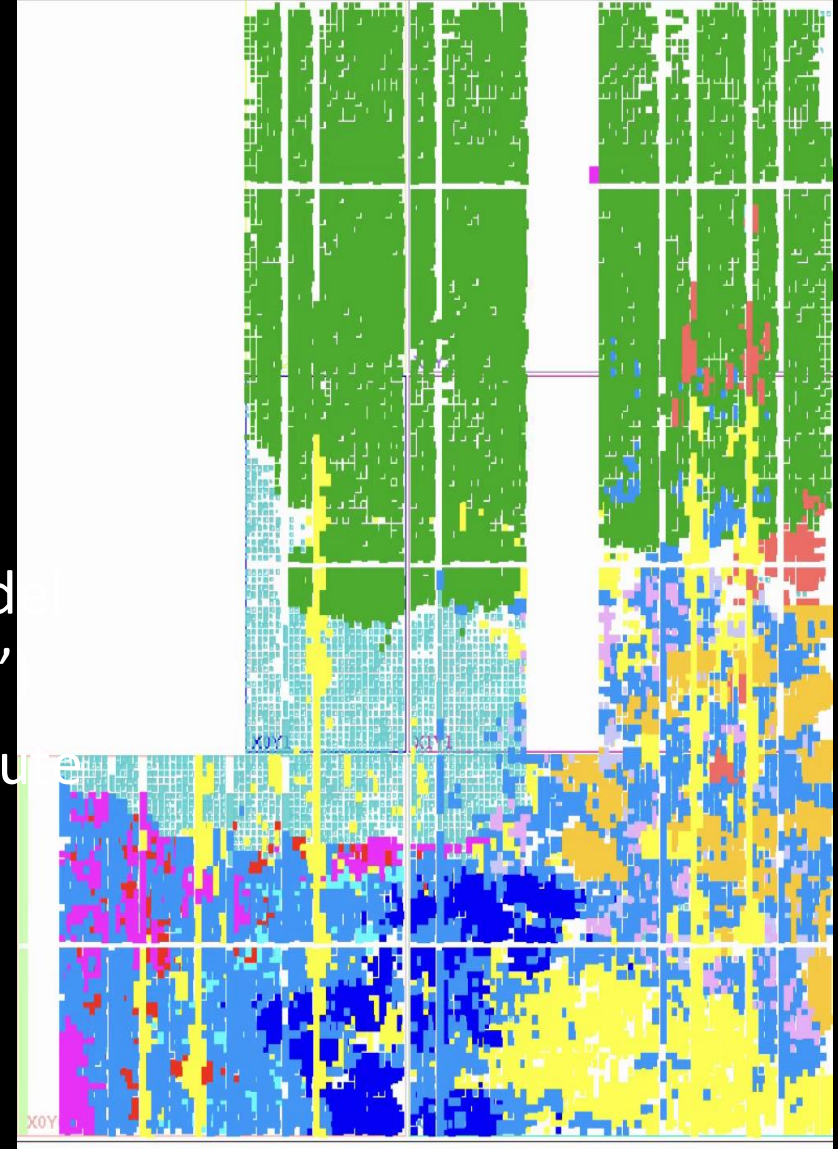




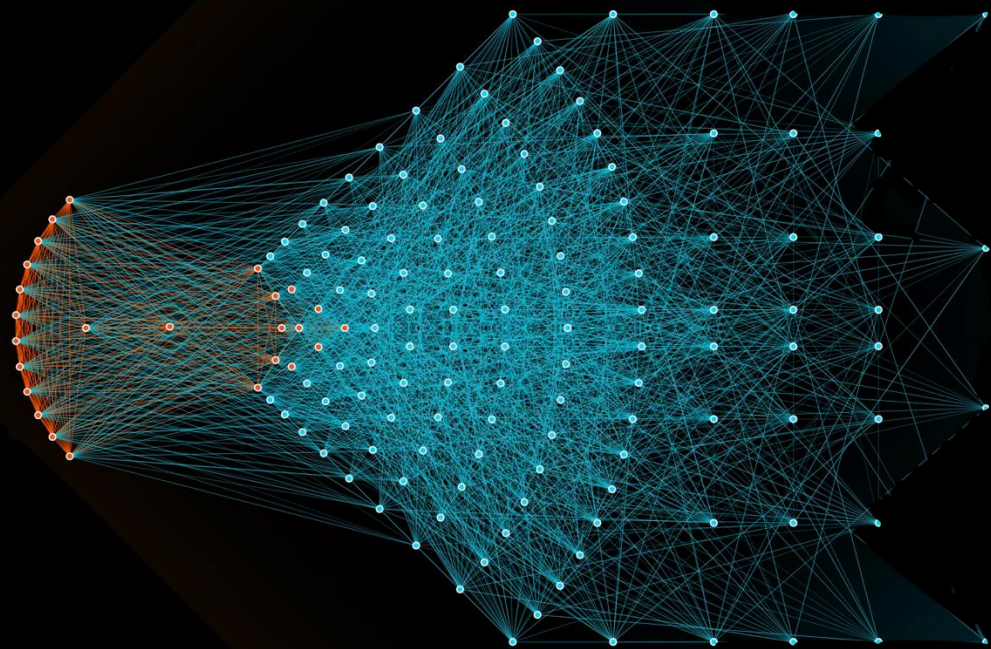
Prediction

### Data flow architecture

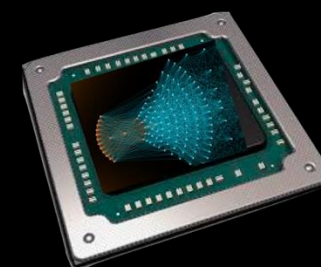
- Tailored hardware for a model (‘‘Decisions are design time’’)
- Each layer is separate computation unit
- Stay on-chip



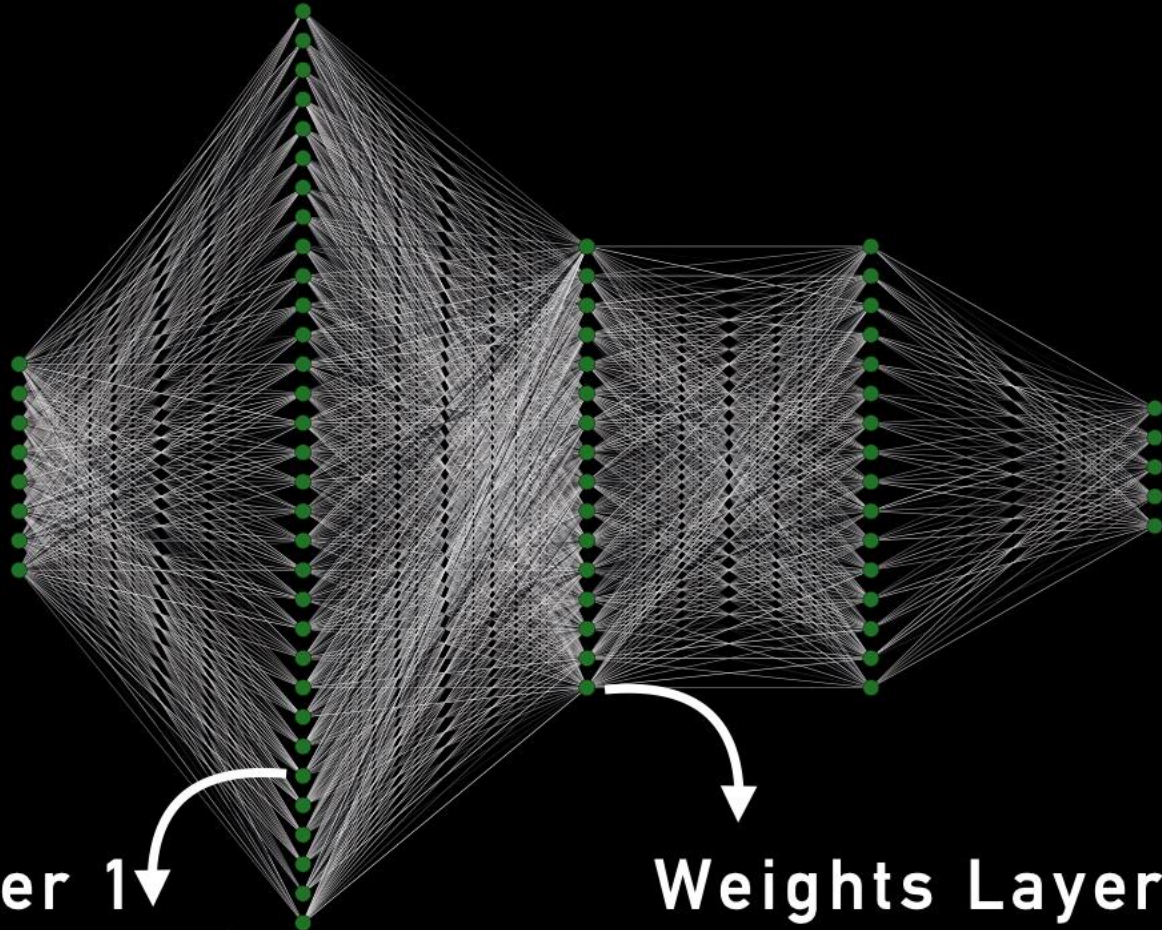
IDEALLY



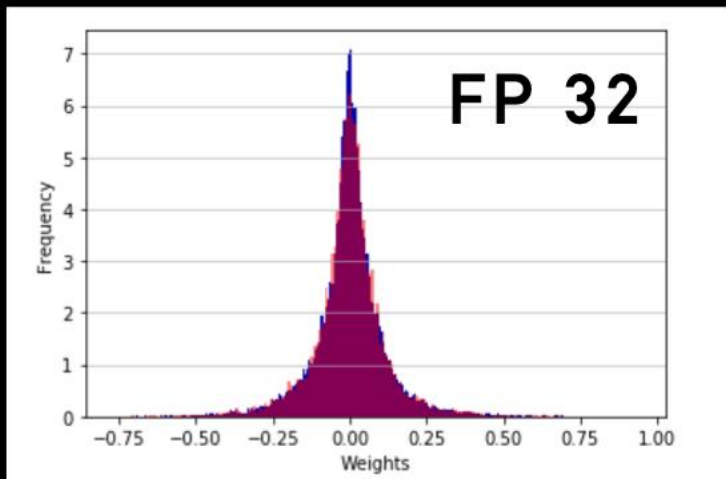
REALITY



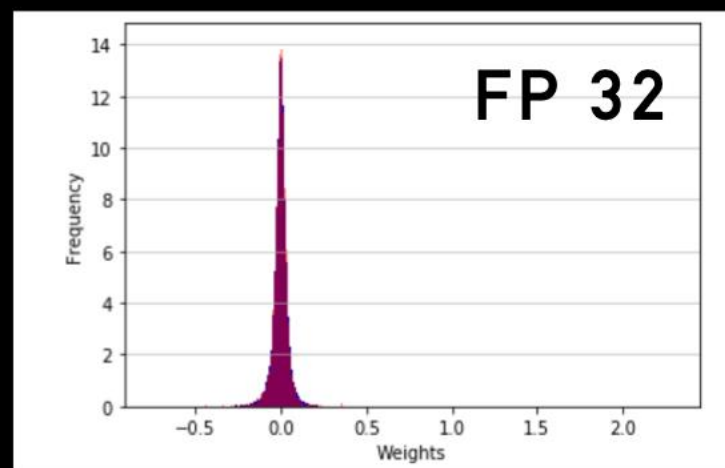
Quantization  
Pruning  
Parallelisation  
Knowledge  
Distillation



**Weights Layer 1**

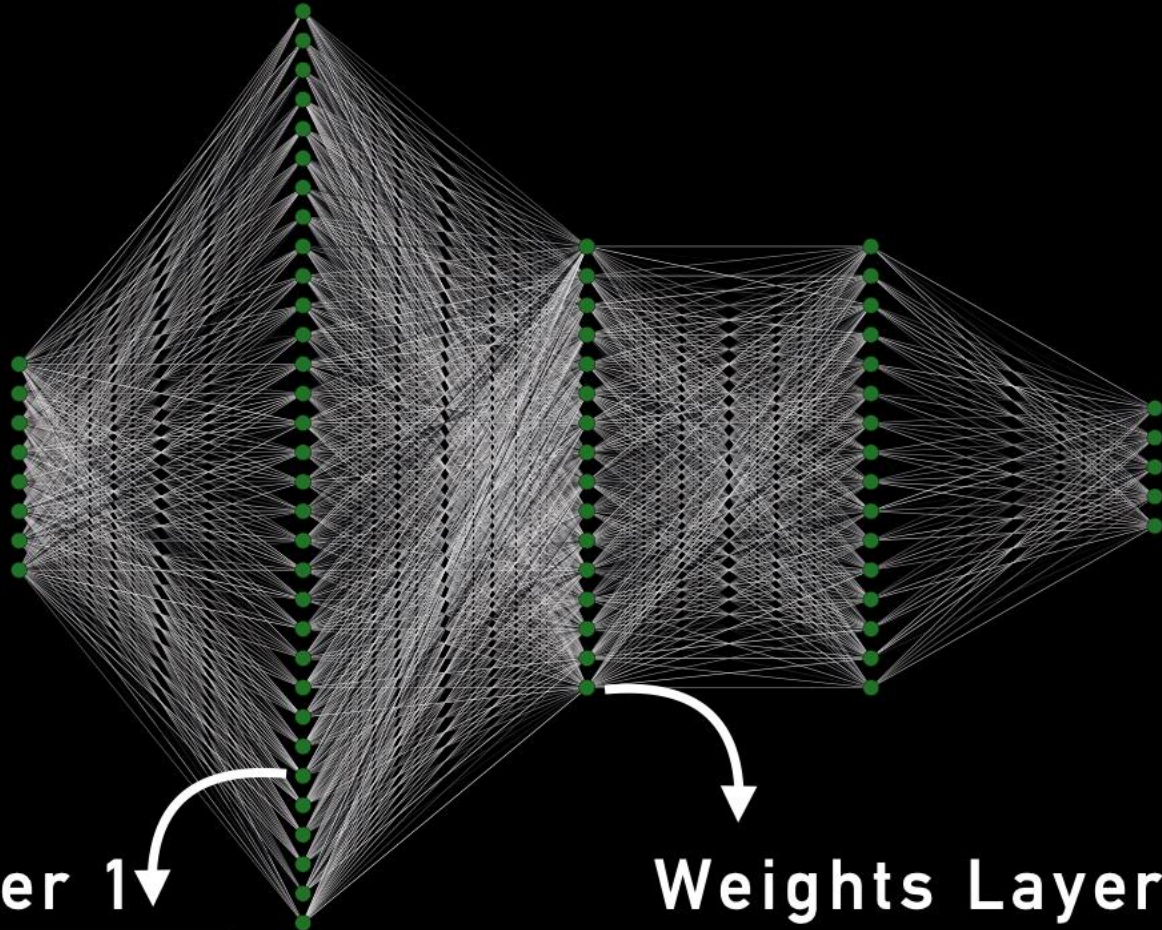


**Weights Layer 2**

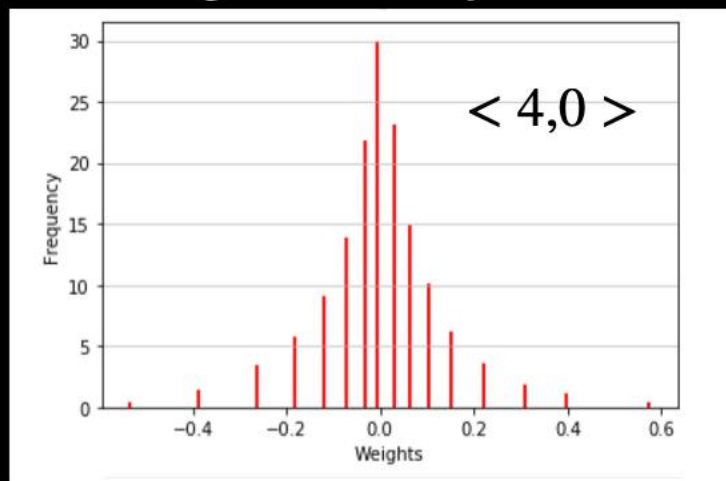


Fixed point

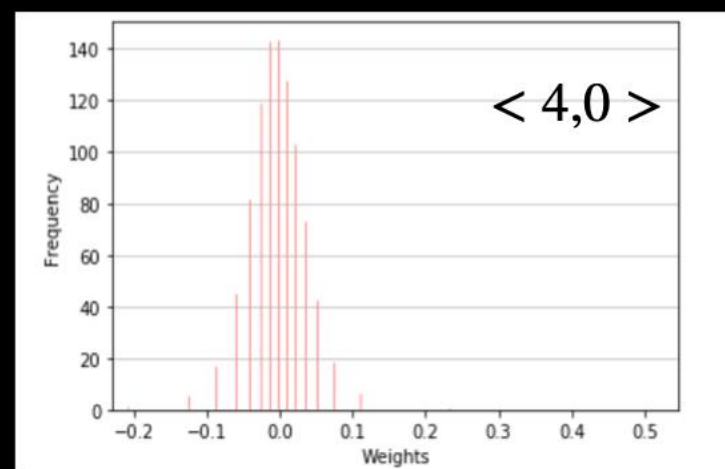
0101.1011101010



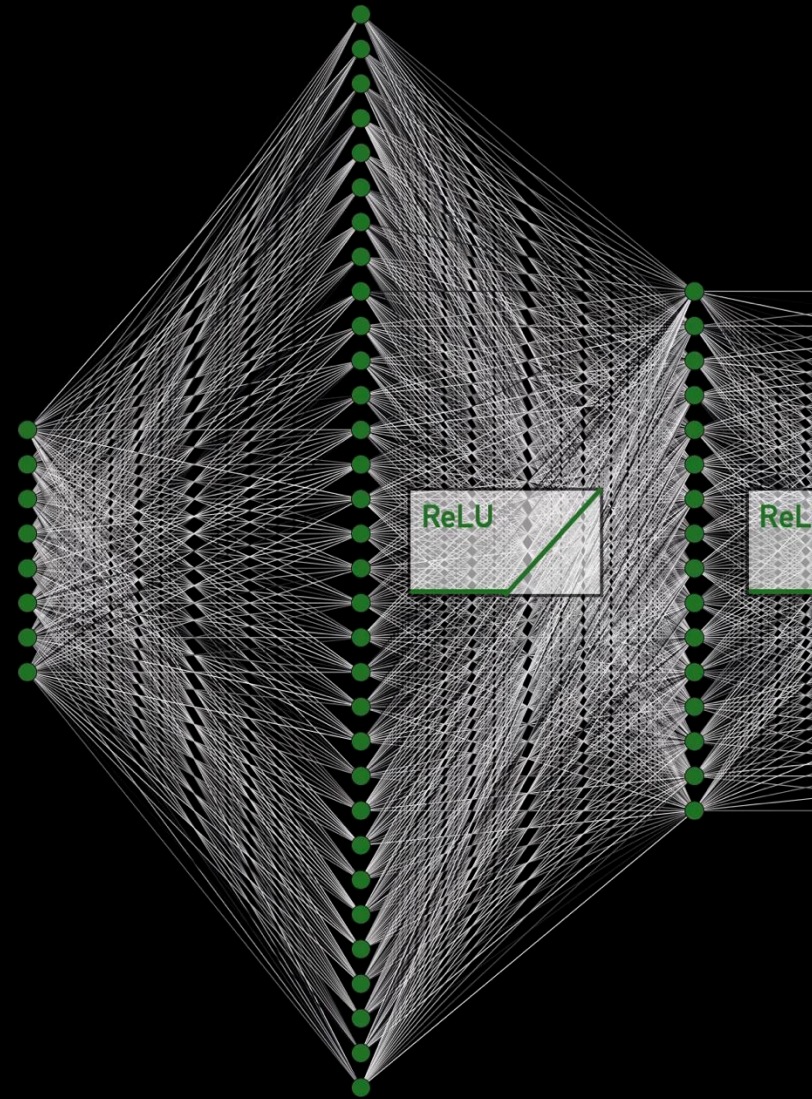
Weights Layer 1



Weights Layer 2



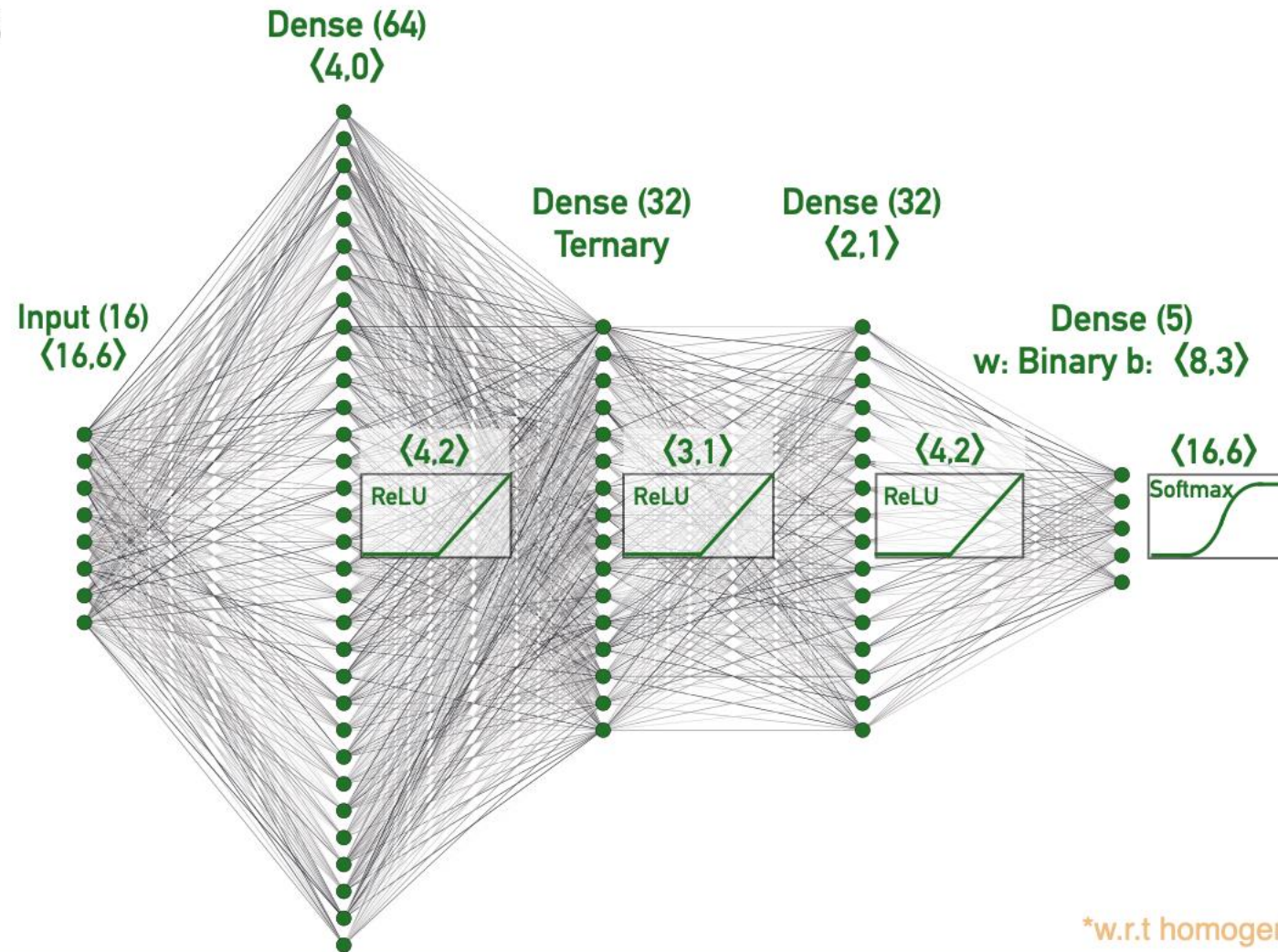
# QKERAS



```
from tensorflow.keras.layers import Input, Activation
from qkeras import quantized_bits
from qkeras import QDense, QActivation
from qkeras import QBatchNormalization
```

```
x = Input((16))
x = QDense(64,
          kernel_quantizer = quantized_bits(6,0, alpha=1),
          bias_quantizer   = quantized_bits(6,0, alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(32,
          kernel_quantizer = quantized_bits(6,0, alpha=1),
          bias_quantizer   = quantized_bits(6,0, alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(32,
          kernel_quantizer = quantized_bits(6,0, alpha=1),
          bias_quantizer   = quantized_bits(6,0, alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(5,
          kernel_quantizer = quantized_bits(6,0, alpha=1),
          bias_quantizer   = quantized_bits(6,0, alpha=1))(x)
x = Activation('softmax')(x)
```

# AutoQKeras



\*w.r.t homogeneously quantized 6 bit model

Model	Accuracy (%)	Precision								$\frac{E}{E_{Q6}}$	$\frac{\text{Bits}}{\text{Bits}_{Q6}}$
		Dense	ReLU	Dense	ReLU	Dense	ReLU	Dense	Softmax		
QE	72.3	$\langle 4, 0 \rangle$	$\langle 4, 2 \rangle$	Ternary	$\langle 3, 1 \rangle$	$\langle 2, 1 \rangle$	$\langle 4, 2 \rangle$	w: Stoc. bin. b: $\langle 8, 3 \rangle$	$\langle 16, 6 \rangle$	0.27	0.18*

# High Granularity Quantization

```
from tensorflow.keras.layers import Input
from HGQ import HQuantize, HDense

inp = Input((16,))
out = HQuantize(name='inp_q', beta=beta)(out)
out = HDense(64, activation='relu', beta=beta)(out)
out = HDense(32, activation='relu', beta=beta)(out)
out = HDense(32, activation='relu', beta=beta)(out)
out = HDense(5, activation='linear', beta=beta)(out)

hgq_model = Model(inp, out)
```

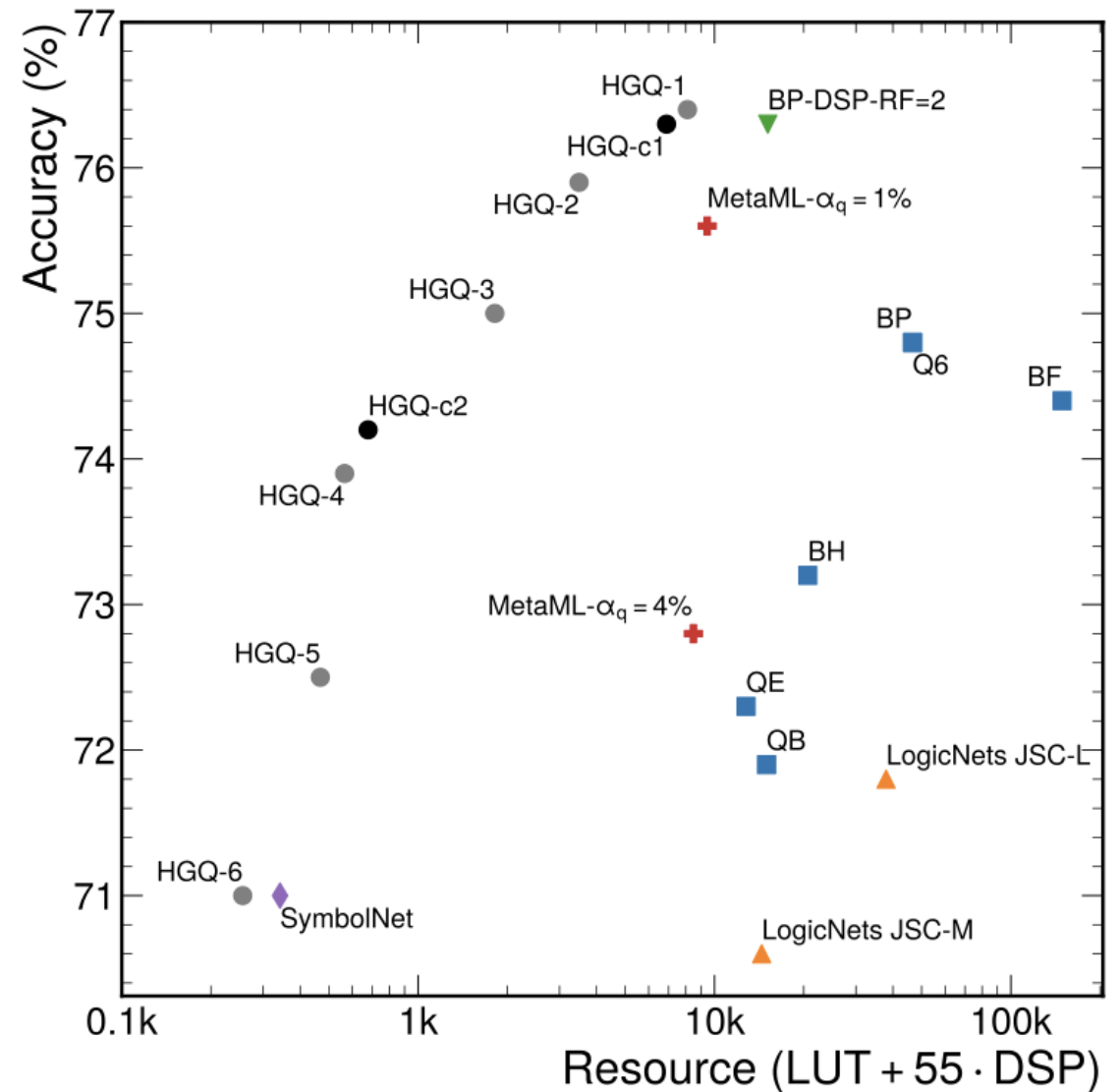
# High Granularity Quantization

```
from tensorflow.keras.layers import Input
from HGQ import HQuantize, HDense

inp = Input((16,))
out = HQuantize(name='inp_q', beta=beta)(out)
out = HDense(64, activation='relu', beta=beta)(o
out = HDense(32, activation='relu', beta=beta)(o
out = HDense(32, activation='relu', beta=beta)(o
out = HDense(5, activation='linear', beta=beta)(

hgq_model = Model(inp, out)
```

$$\mathcal{L} = \mathcal{L}_{\text{base}} + \beta \cdot \bar{E}$$



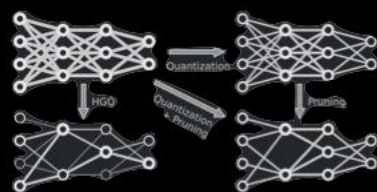
# hls4ml + AMD Research Labs: QONNX

Frontend (QAT)

Backend (compiler)

QK Keras

PYTORCH Brevitas



HQG

QONNX

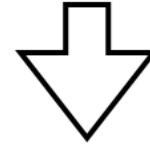
FINN AMD

hls4ml

Marius Köppel  
A. Pappalardo et al. (2022)

# Standard industry tools couldn't reach the specs that we require

**Better:**



Platform	Framerate (FPS)	Power (W)	pixels / s / W ( $\times 10^6$ )
<b>hls4ml</b>	<b>1500</b>	<b>3.7</b>	26.5
Vitis AI	40	9.4	1.1
<b>Ratio hls4ml / Vitis AI</b>	<b>38</b>	<b>0.4</b>	<b>24.0</b>

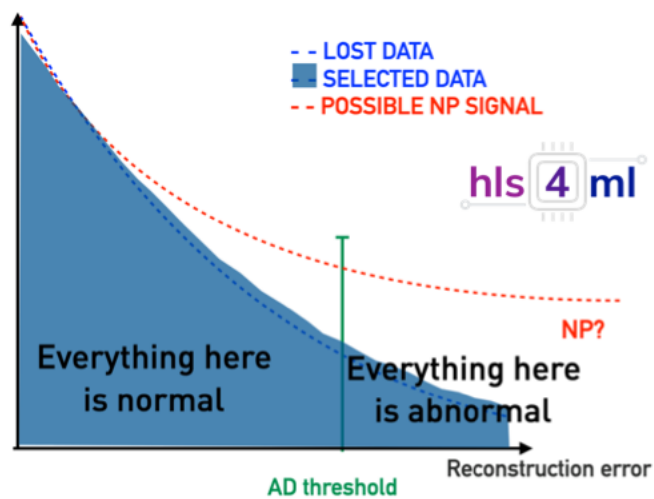


ZCU102  
(Zynq Ultrascale+ 9)

# 2024: Neural hardware triggers

## CMS:

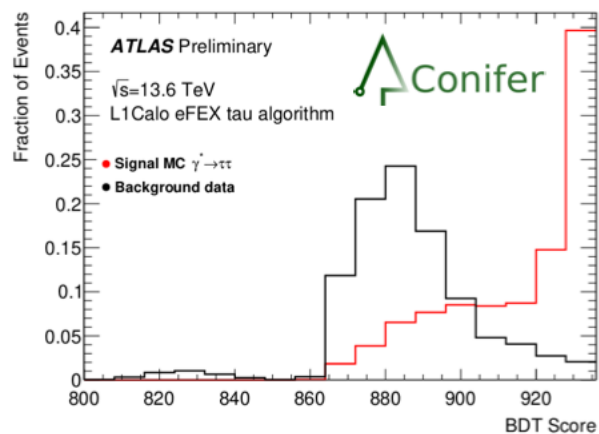
- Anomaly detection in 50 ns



CMS DP2023\_079

## ATLAS:

- BDT selecting candidate  $\tau$  events in <100ns

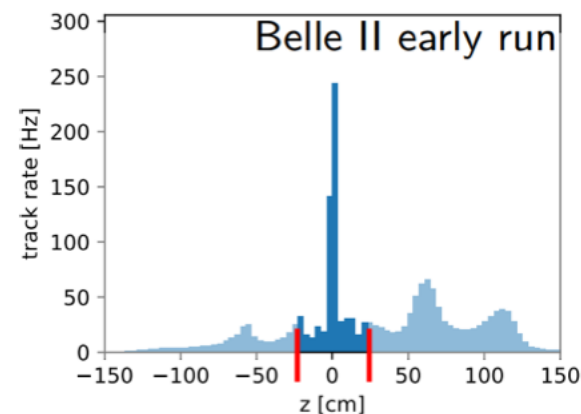


L1CaloTriggerPublicResults

## Belle-2:

- Neural track trigger since April 2021

Reject  $|z| \neq 0$  cm tracks



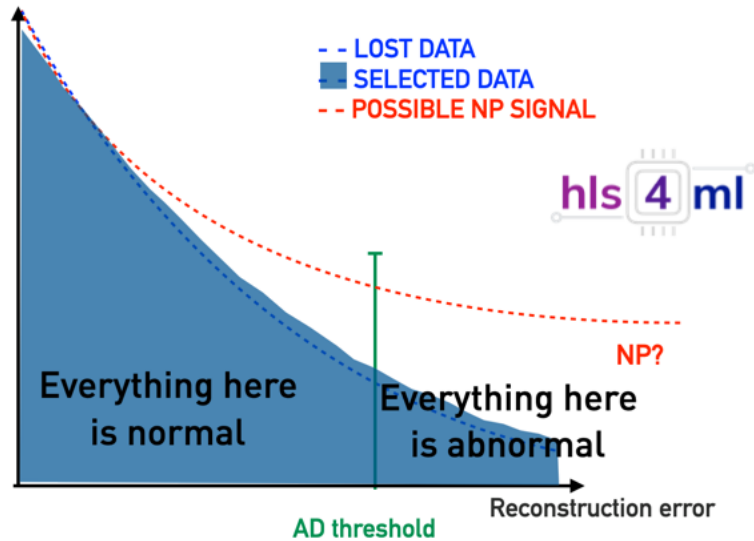
arxiv:2402.1496

2

# 2024: Neural hardware triggers making decisions in LHC experiments!

CMS Experiment:

- Anomaly detection in 50 ns
- 300 events/second

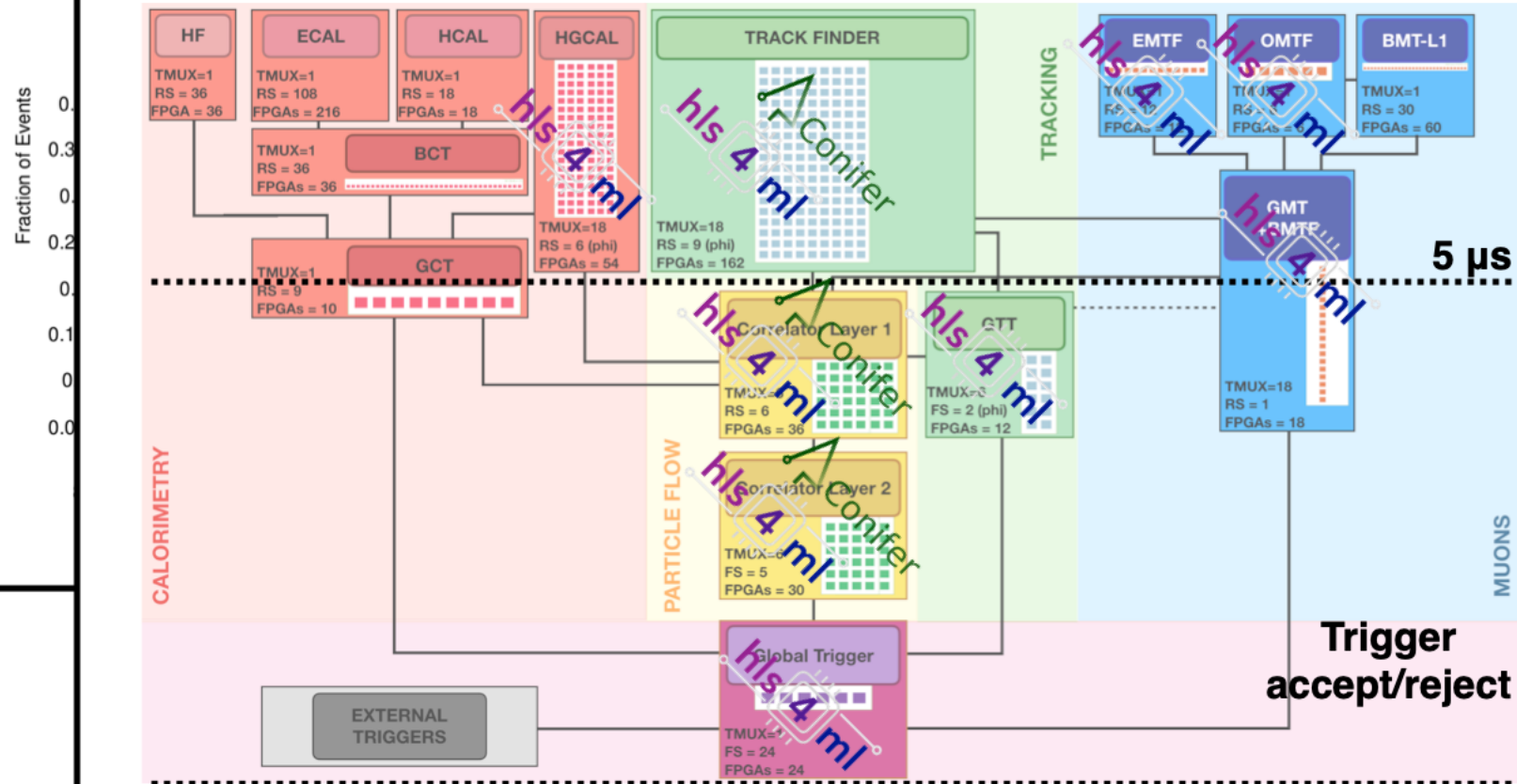


CMS DP2023\_079

ATLAS Experiment:

- BDT selecting candidate  $\tau$  lepton events in <100 ns

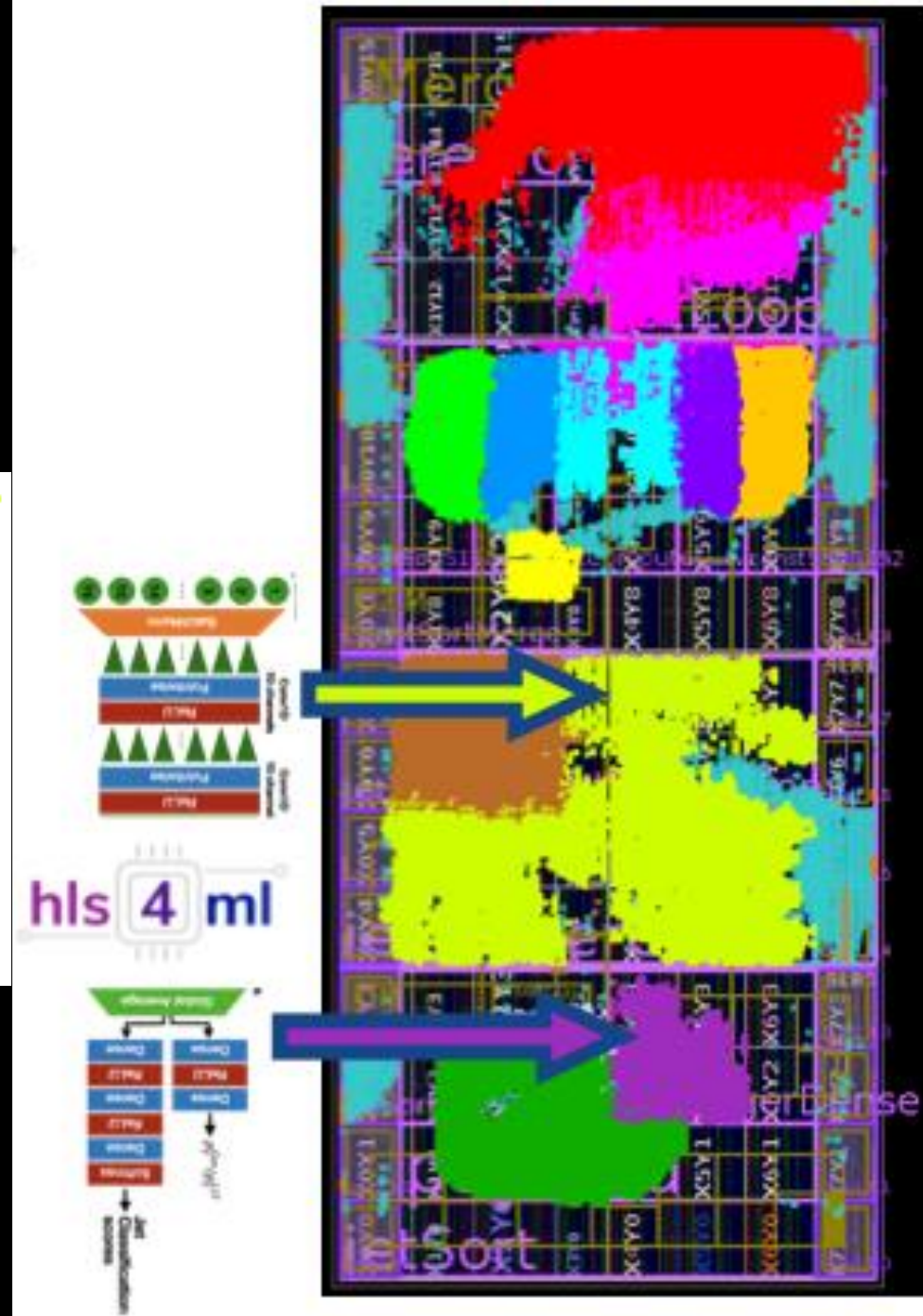
## 2031: ~20 billion inferences/s during HL-LHC

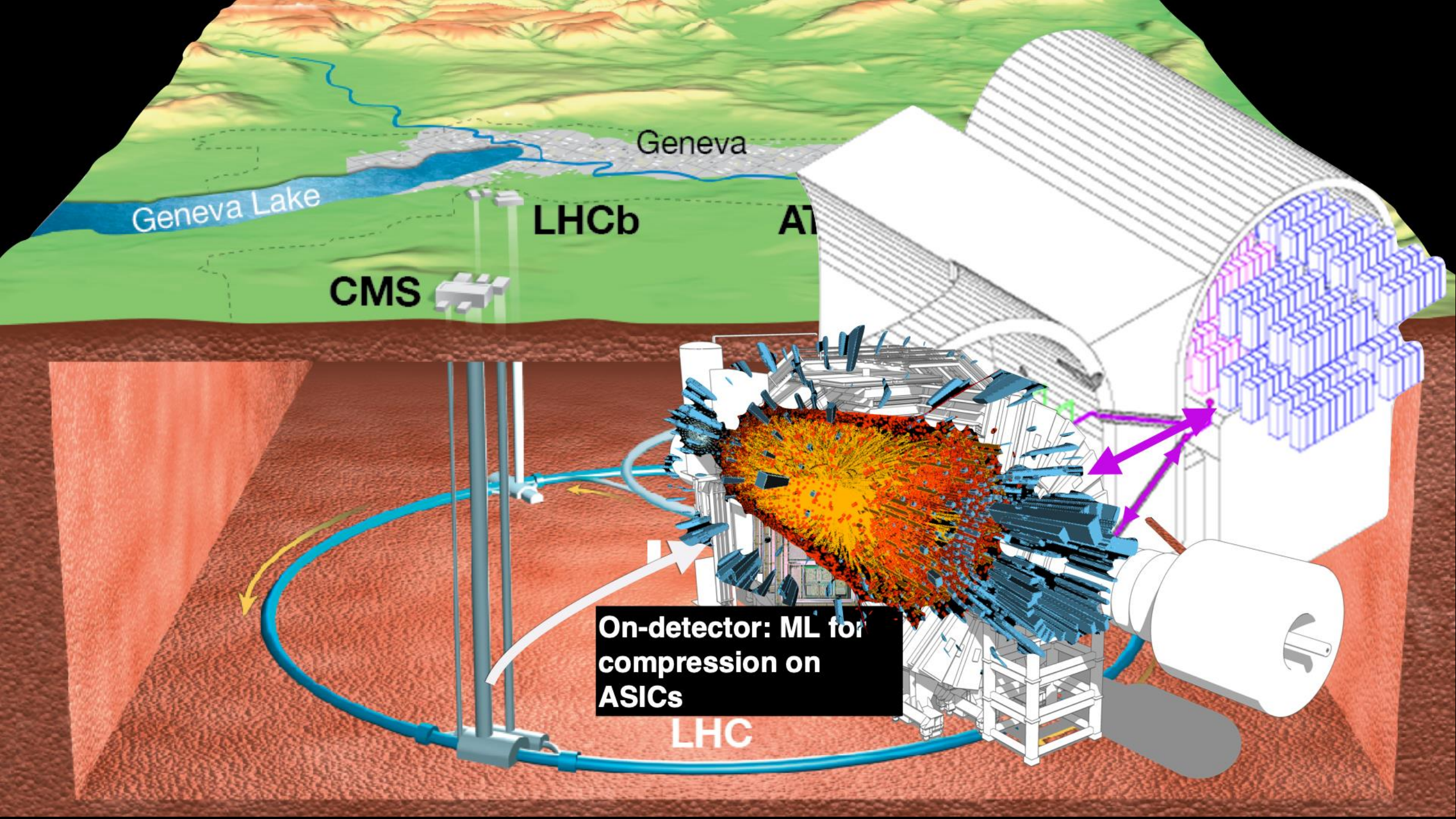


# HL-LHC CMS will run with FPGA track finding, particle flow, jet clustering and flavour tagging in < 12 $\mu$ s

- **Reco. jets**; sort the constituents by  $p_T$ ; compute input variables; per-constituent NN part; per-jet NN part; sort jets by  $p_T$

Component	Latency ( $\mu$ s)	LUT	FF	DSP	BRAM
Preprocessing (incl. sort)	0.10	6%	4%	1%	0.6%
Jet Tag NN	0.19	7%	5%	14%	1%
Jet Reco	0.74	9%	7%	5%	0%
Total System	1.01	25%	15%	19%	1%





Geneva

Geneva Lake

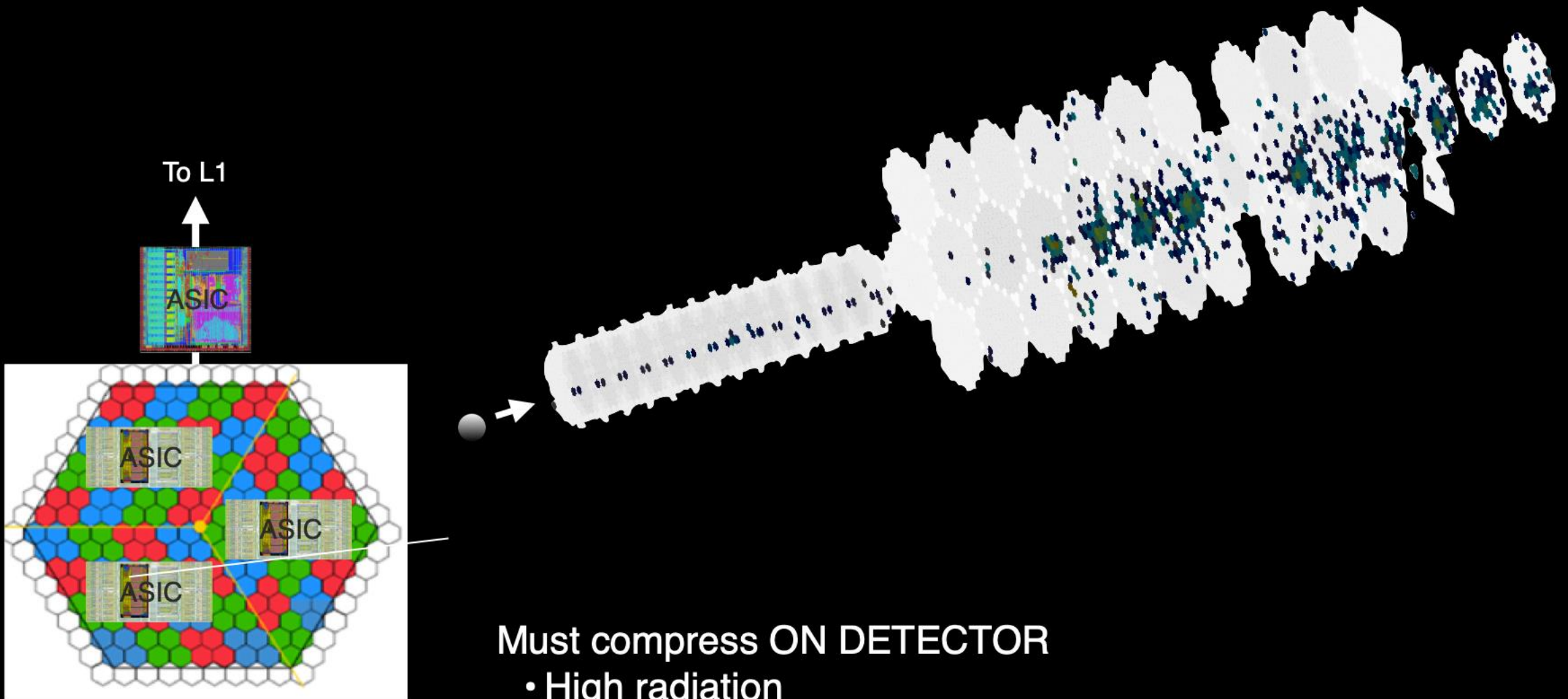
LHCb

ATLAS

CMS

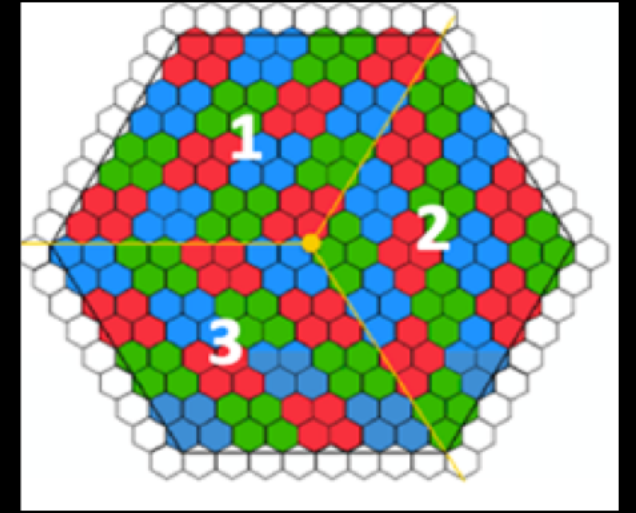
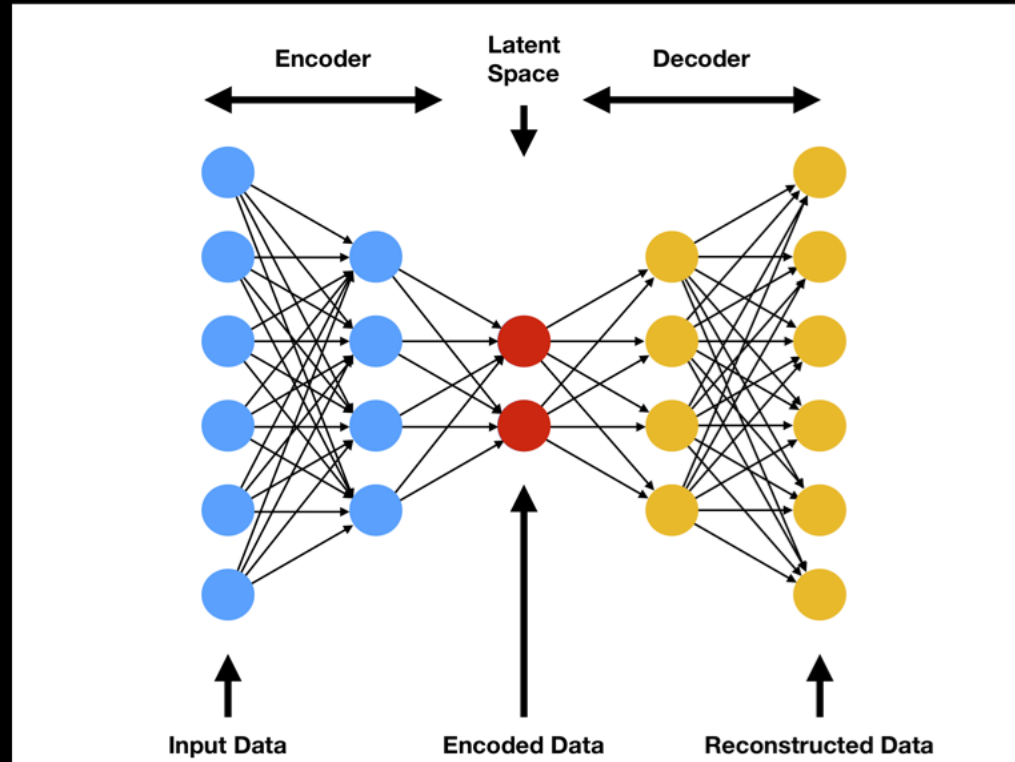
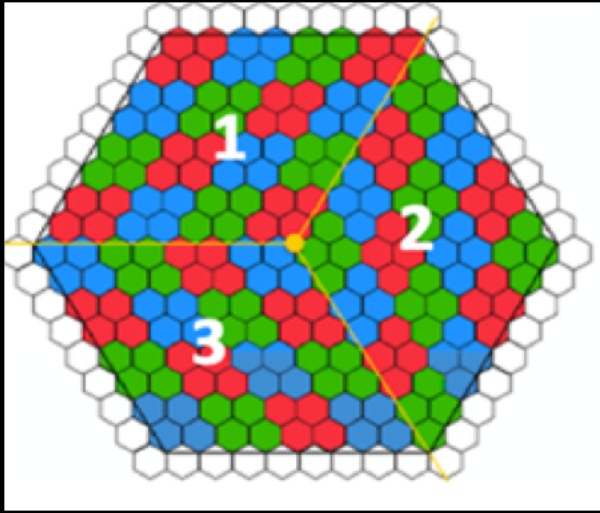
On-detector: ML for compression on ASICs

LHC



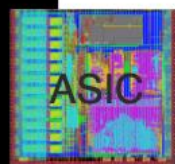
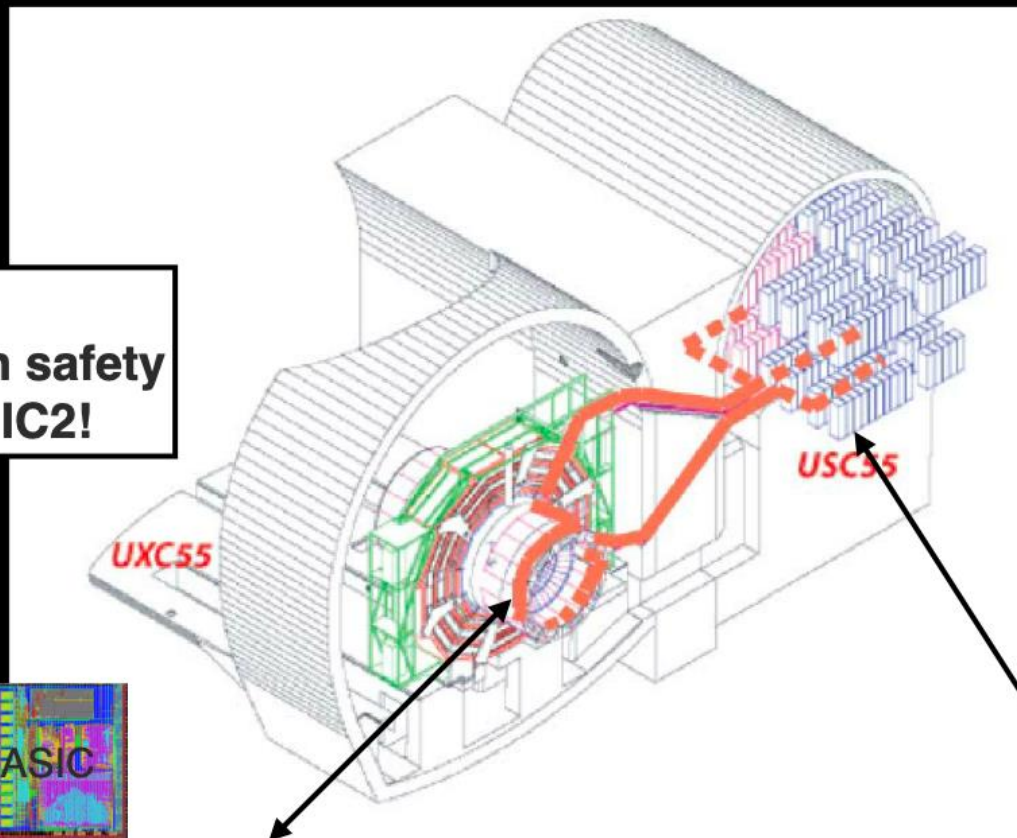
### Must compress ON DETECTOR

- High radiation
- Cooled to -30 → low power
- 1.5  $\mu$ s latency

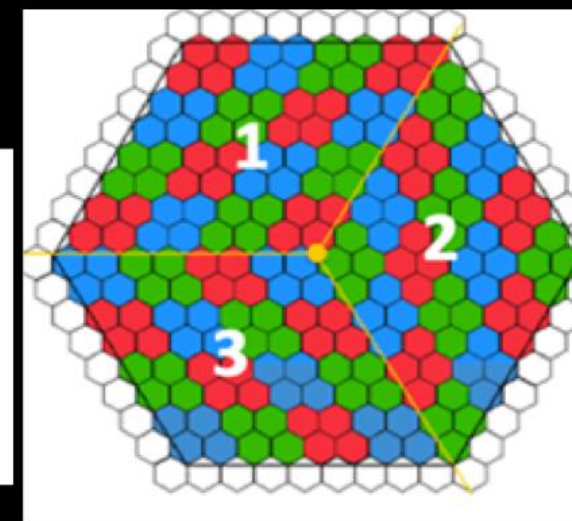
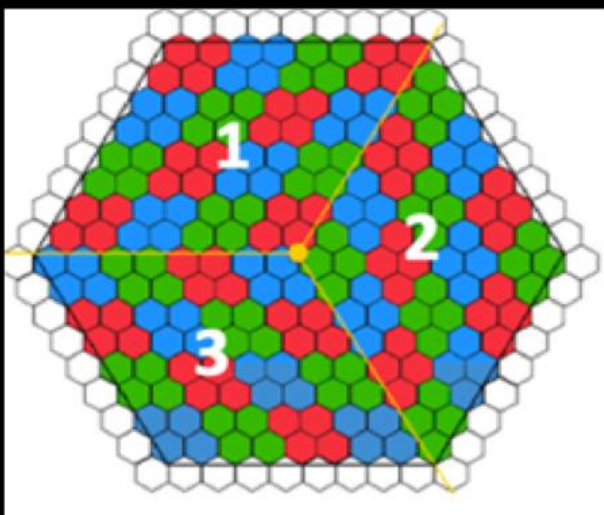
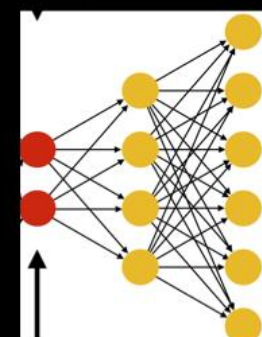
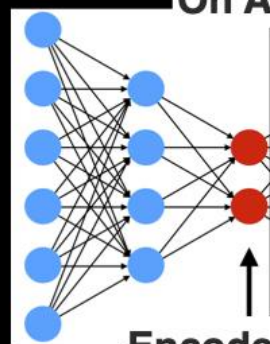


**Variational Autoencoder**

- 75-100 mW
- Triplicated w/b for radiation safety
- Reprogrammable w/b over IC2!



On ASIC



PRESS RELEASE

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

May 21, 2024

Plano, Texas



Catapult AI NN brings together hls4ml, an open-source package for machine learning hardware acceleration, and Siemens' Catapult™ HLS software for High-Level Synthesis. Developed in close collaboration with Fermilab, a U.S. Department of Energy Laboratory, and other leading contributors to hls4ml, Catapult AI NN addresses the unique requirements of machine learning accelerator design for power, performance, and area on custom silicon.



# Data readout bottleneck



Example vertex detector inner layer from before (#241 FCC-ee):

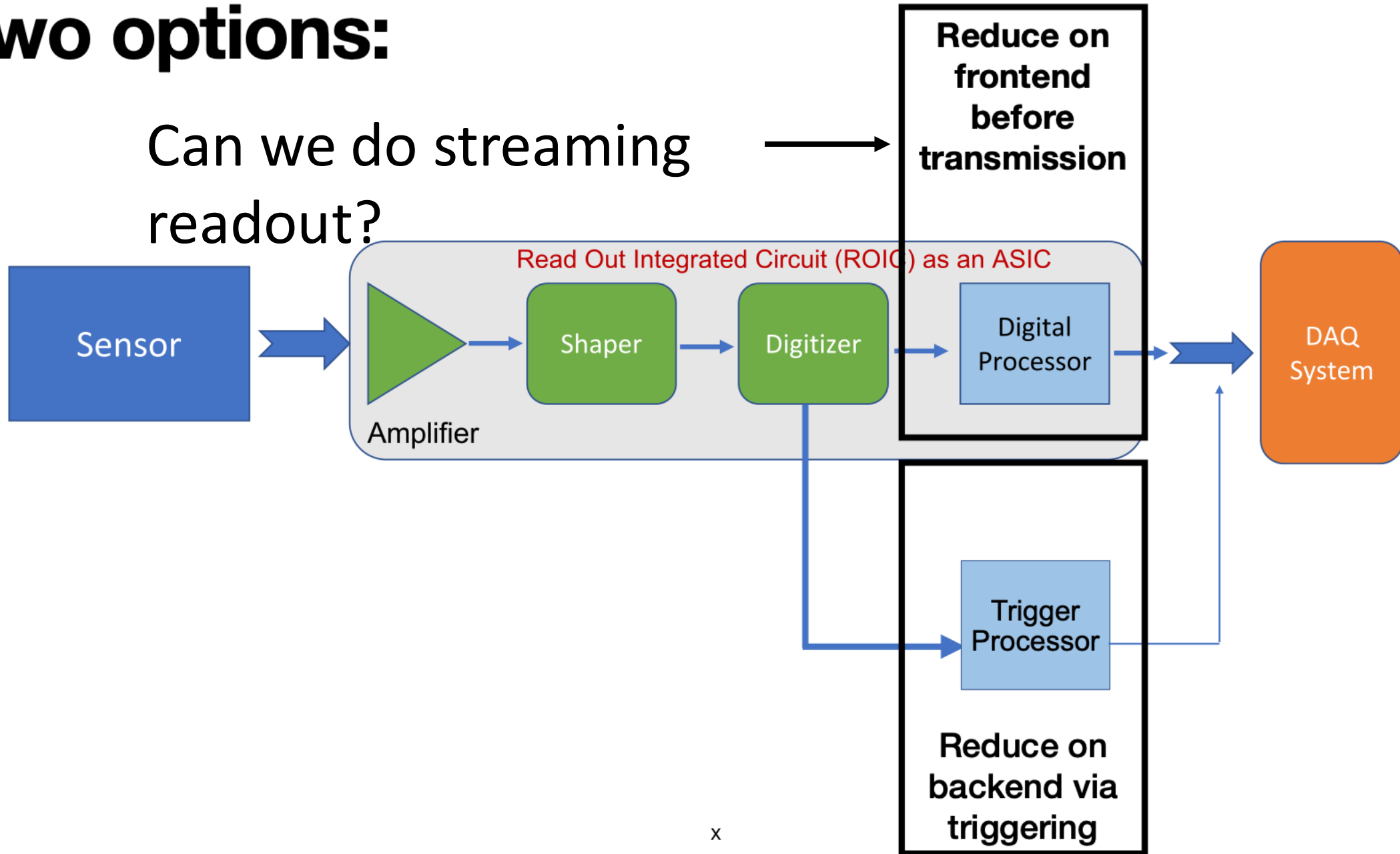
32 bits pixel data in inner layer @ 200 MHz/cm<sup>2</sup>

➔ **24.4 Gbit/s**

**2.2 Tb/s for Layer 1** (x10 less in Layer 2)

# Two options:

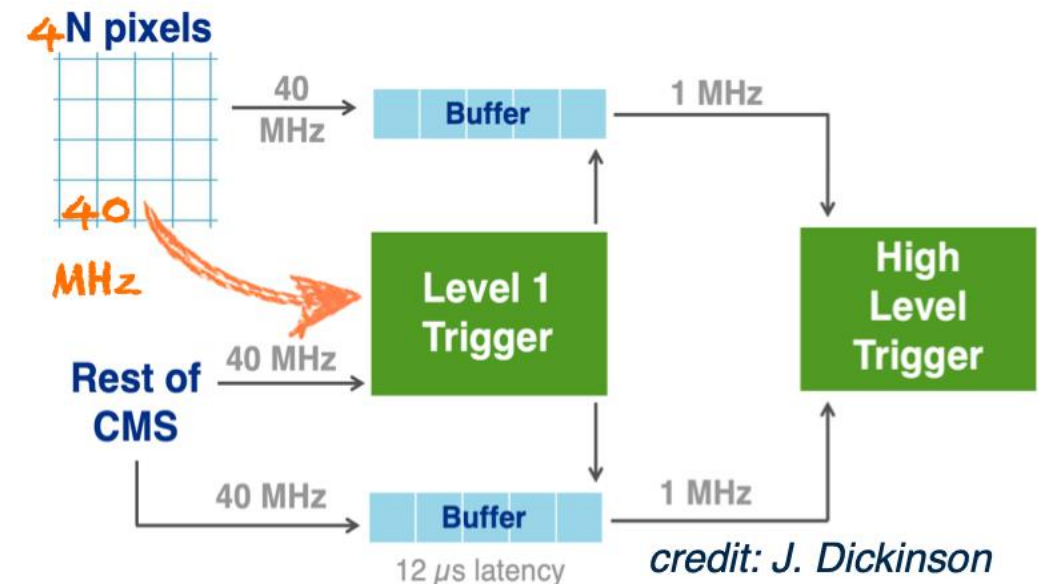
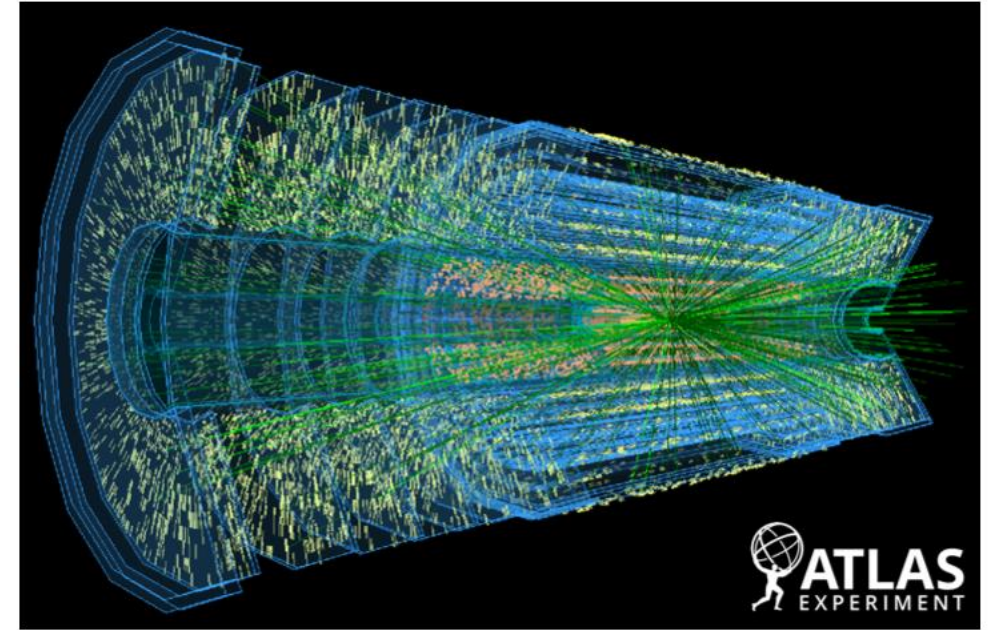
Can we do streaming readout?



# Smart Pixels

<https://fastmachinelearning.org/smart-pixels/>

- Reduce silicon data via **in-pixel intelligence**
  - **frontend filtering**: discard low- $p_T$  track data ( $< 2$  GeV)
  - **feature extraction**: Extract particle position and angle in pixel front-end ASICs from charge in **single pixel layer**
- Bandwidth savings of 57-75%!

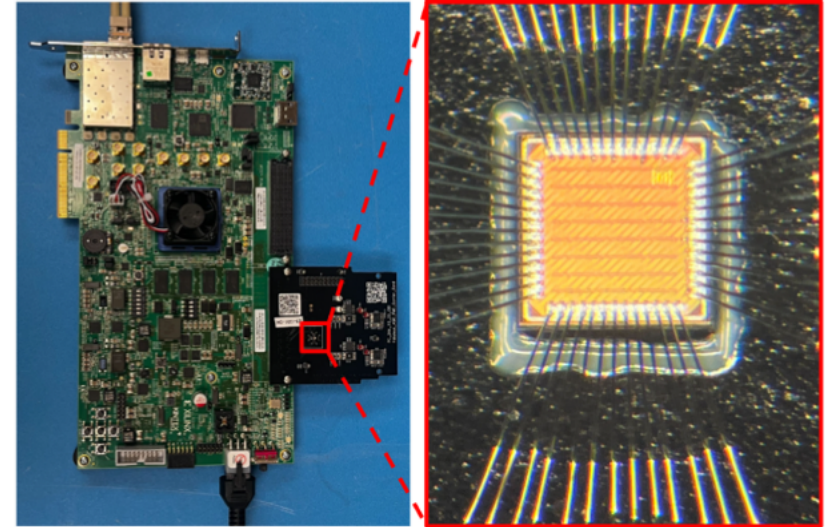


# Reconfigurable logic in ASIC design

## The Embedded FPGA framework

- Pathway to implementing ML "at source"
  - Fully reconfigurable logic on detector frontend
- Open source (FABulous, OpenFPGA)
  - potential to apply to variety of subsystems/ fields (SuperKEKB, FCC-ee, DUNE, free-electron lasers)

28nm CMOS ASIC (1x1mm)



# aie4ml

AI engines can offer higher compute density,  
provide necessary high speed I/O interfaces.  
Make them easier to use for trigger  
experiments!

hls4ml for AI engines

Status:

linear layers and ReLu implemented

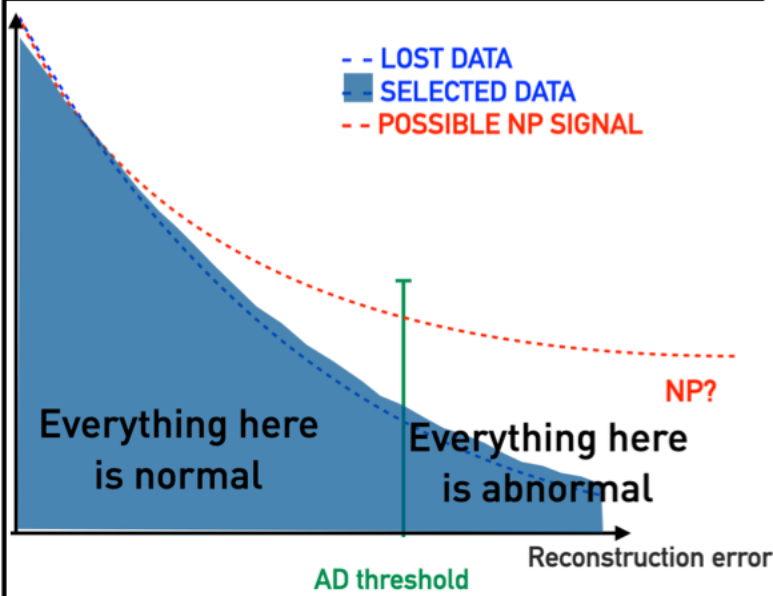
scales workloads across AIE array

support for Gen2/Gen3 devices

NGT WP 1.2

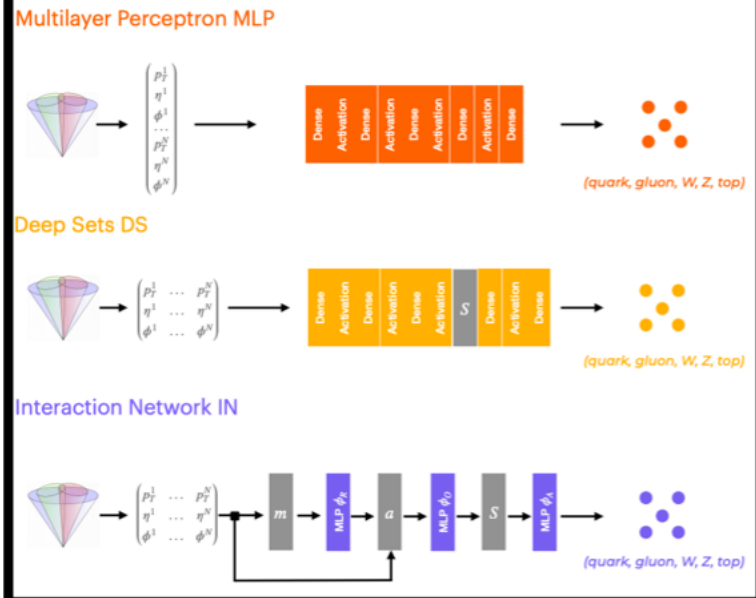


# Anomaly detection with VAEs in 50 ns



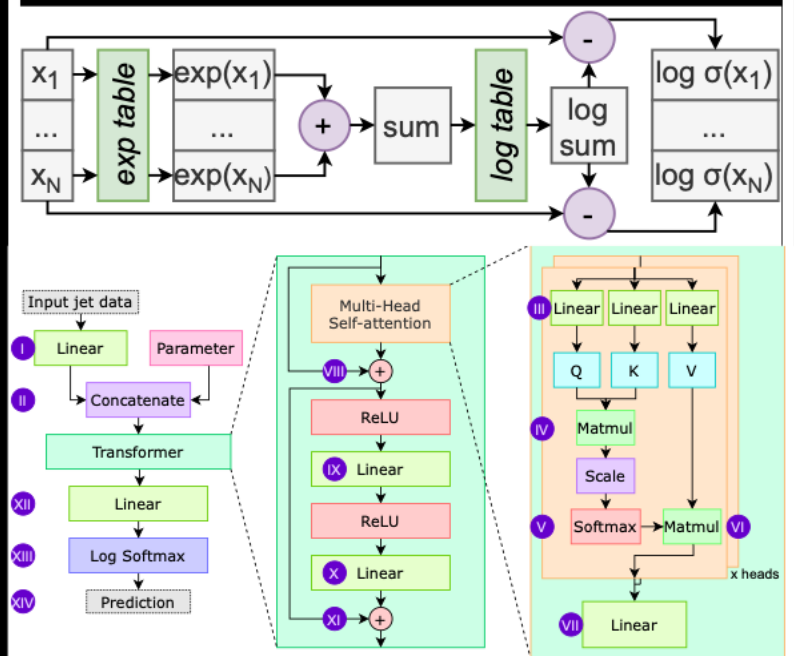
**CMS DP2023\_079**  
**E. Govorkova et al (2022)**

# Quantised Interaction Networks and Deep Sets in <160 ns



**P. Odagiu et al. 2024**

# Fully on-chip transformers in 90 ns (500k flops, 9k param)



**IEEE ICFTP 2022**  
**arxiv:409.05207**

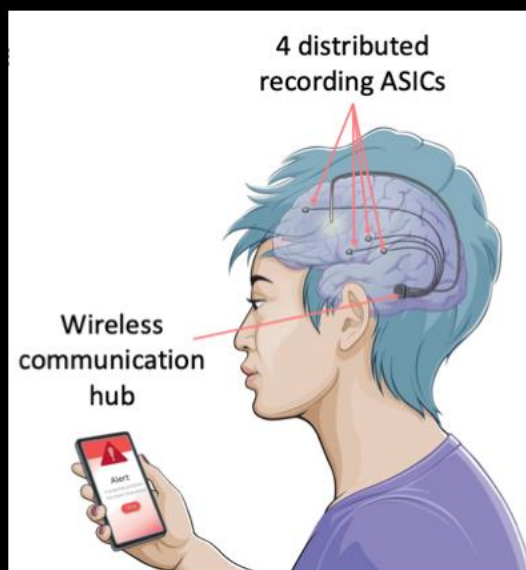
...and outside

## Semantic segmentation for autonomous vehicles



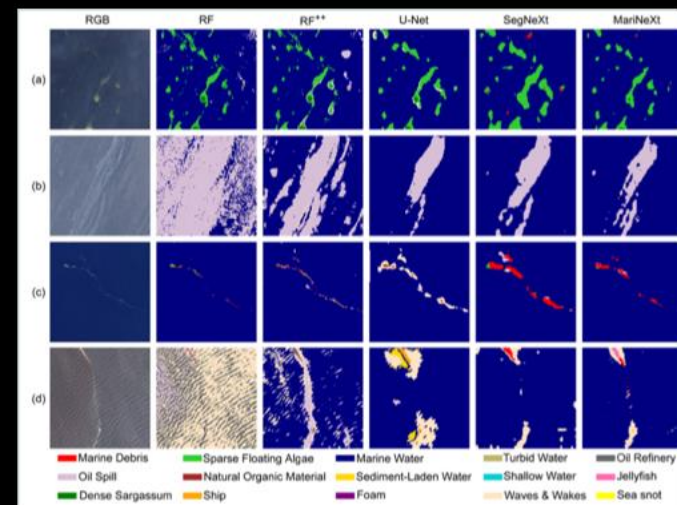
N. Ghilmetti et al. 2022

## Seizure Predicting Brain Implant



W. Lemaire et al. 2022

## Earth monitoring in satellites



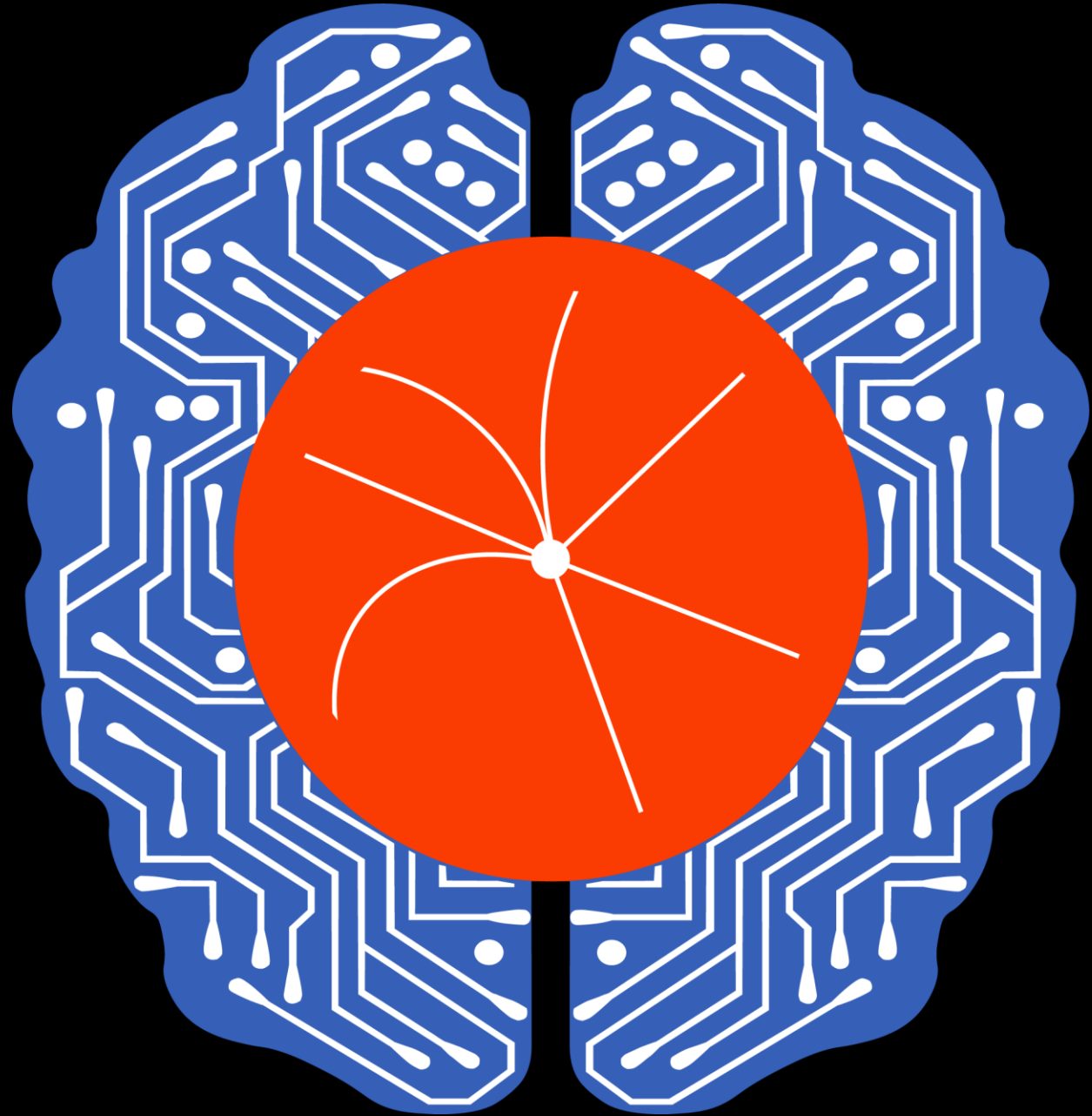
Edge SpAlce, S. Summers

- MLPerf tinyML benchmarking
- For fusion science phase/mode monitoring
- Crystal structure detection
- Triggering in DUNE

- Accelerator control
- Magnet Quench Detection
- Food contamination detection
- Quantum control etc....

# FastML Lab

Keep in touch with community by signing up to the e-group [hls-fml@cern.ch](mailto:hls-fml@cern.ch) ([sign up here!](#)) and attending [meetings \(listed here\)](#)



BACKUP

Currently  
being  
built

# Smart Pixels

Technology	65nm CMOS	28nm CMOS
Pixel ROIC size	50x50 $\mu\text{m}^2$	25x25 $\mu\text{m}^2$
Pixel Sensor size	100x25	50x12.5 $\mu\text{m}^2$
Pixels	394x400 = 157.6k	788x800 = 0.63M
Detection threshold	~1000e-	~500e-
Hit rate	< 3GHz/cm <sup>2</sup>	< 3GHz/cm <sup>2</sup>
Readout rate	1MHz	40MHz (?)
Digital buffer	12.5 $\mu\text{s}$	(?)
Readout Bandwidth	1-4 links @ 1.28Gbps	Photonic link @ 30-100 Gbps
Radiation tolerance	500Mrad at -15°C	1Grad at -15°C
Power	1 W /cm <sup>2</sup>	1 W /cm <sup>2</sup>

J. Dickinson