



Machine Learning at the Edge of Scale and Speed

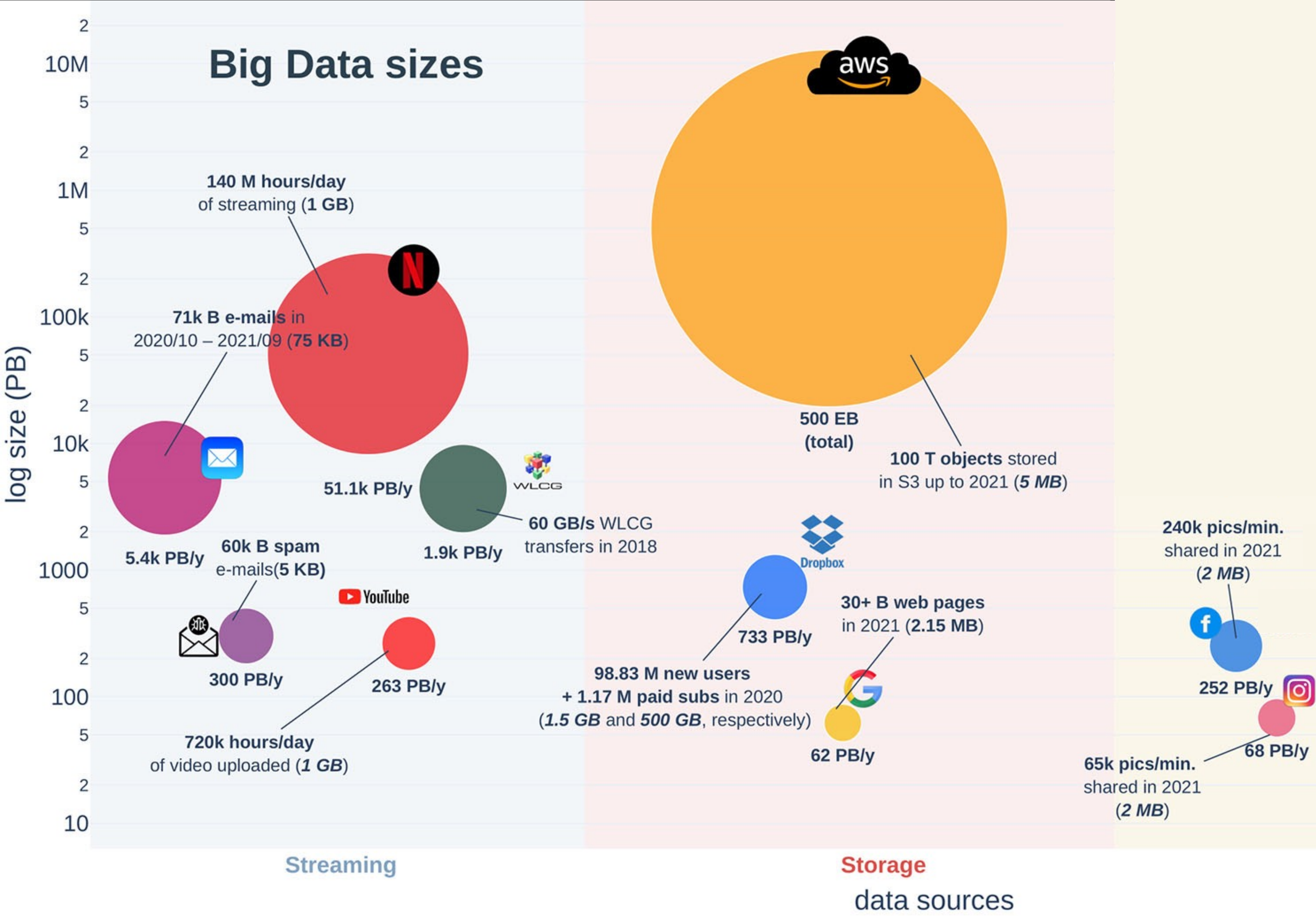
Theo Klæboe Årrestad (ETH Zürich)

The Zettabyte Era

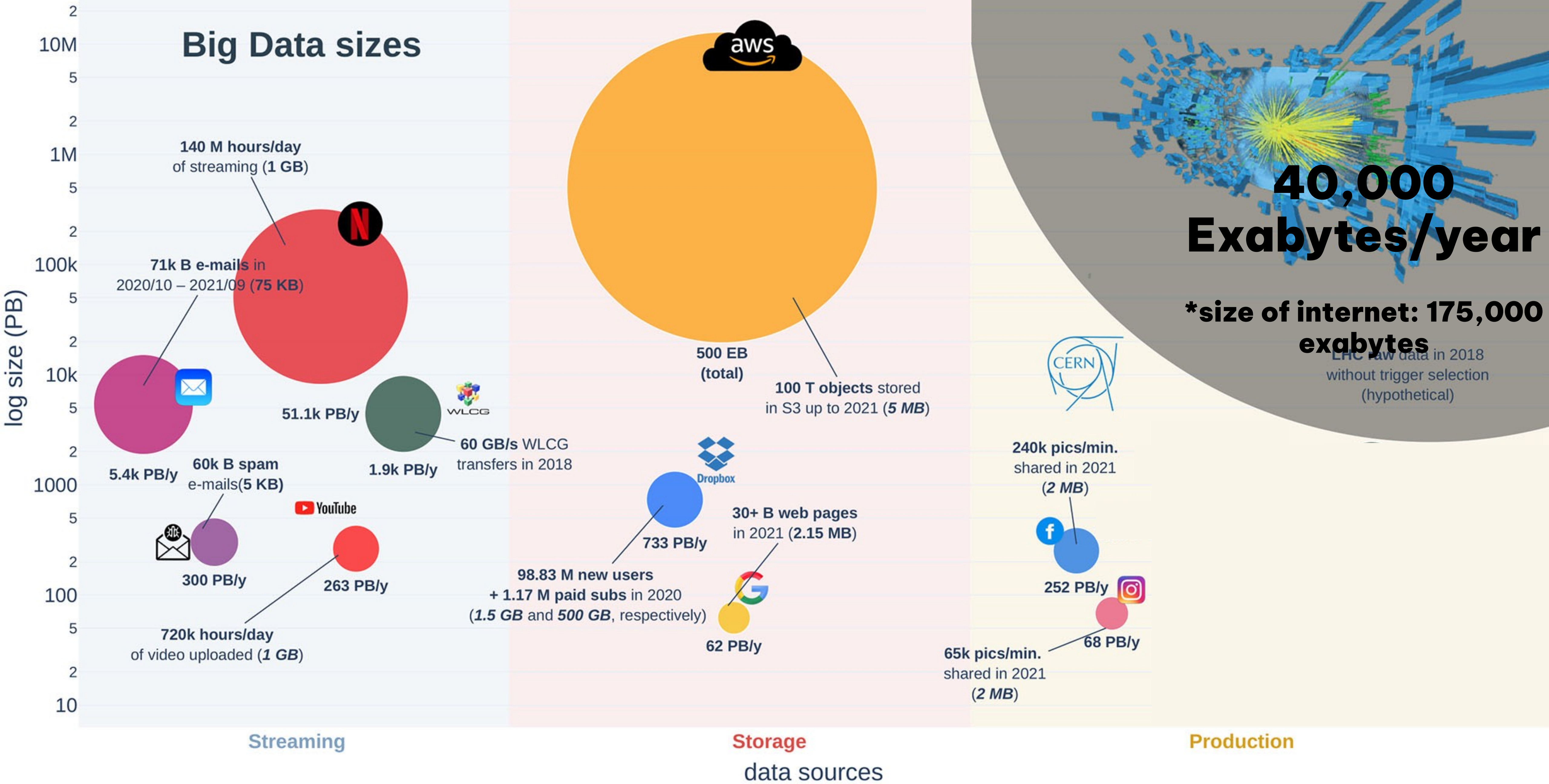
When will global Internet traffic reach an annual run rate of one Zettabyte? Well, that day has finally come. According to our estimate, the world's collective Internet use will reach the Zettabyte threshold for this calendar year on September 9, 2016. Finally... Consider the fact that the Internet essentially began to scale for global consumer

... (~size of internet ~175 Zettabytes)

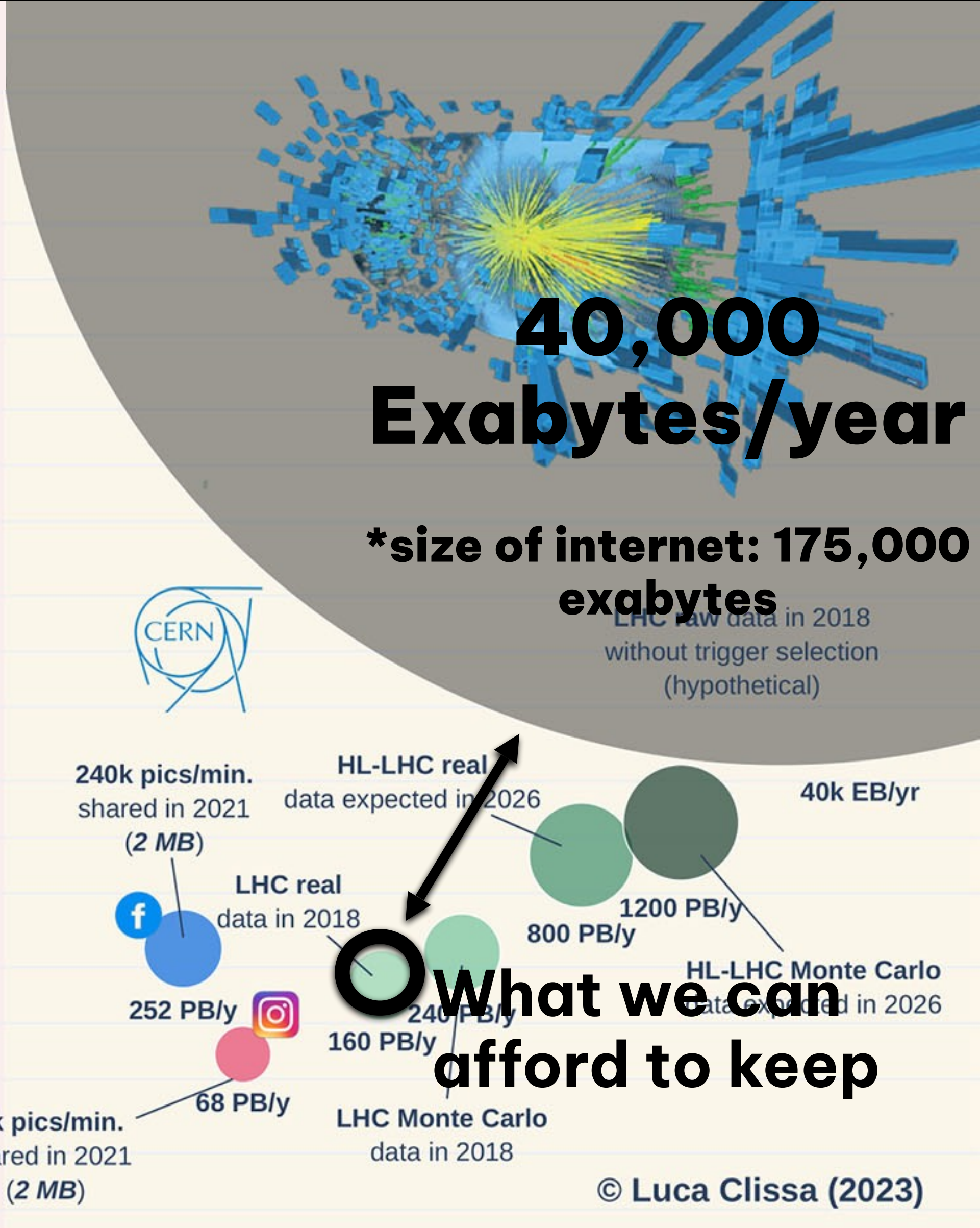
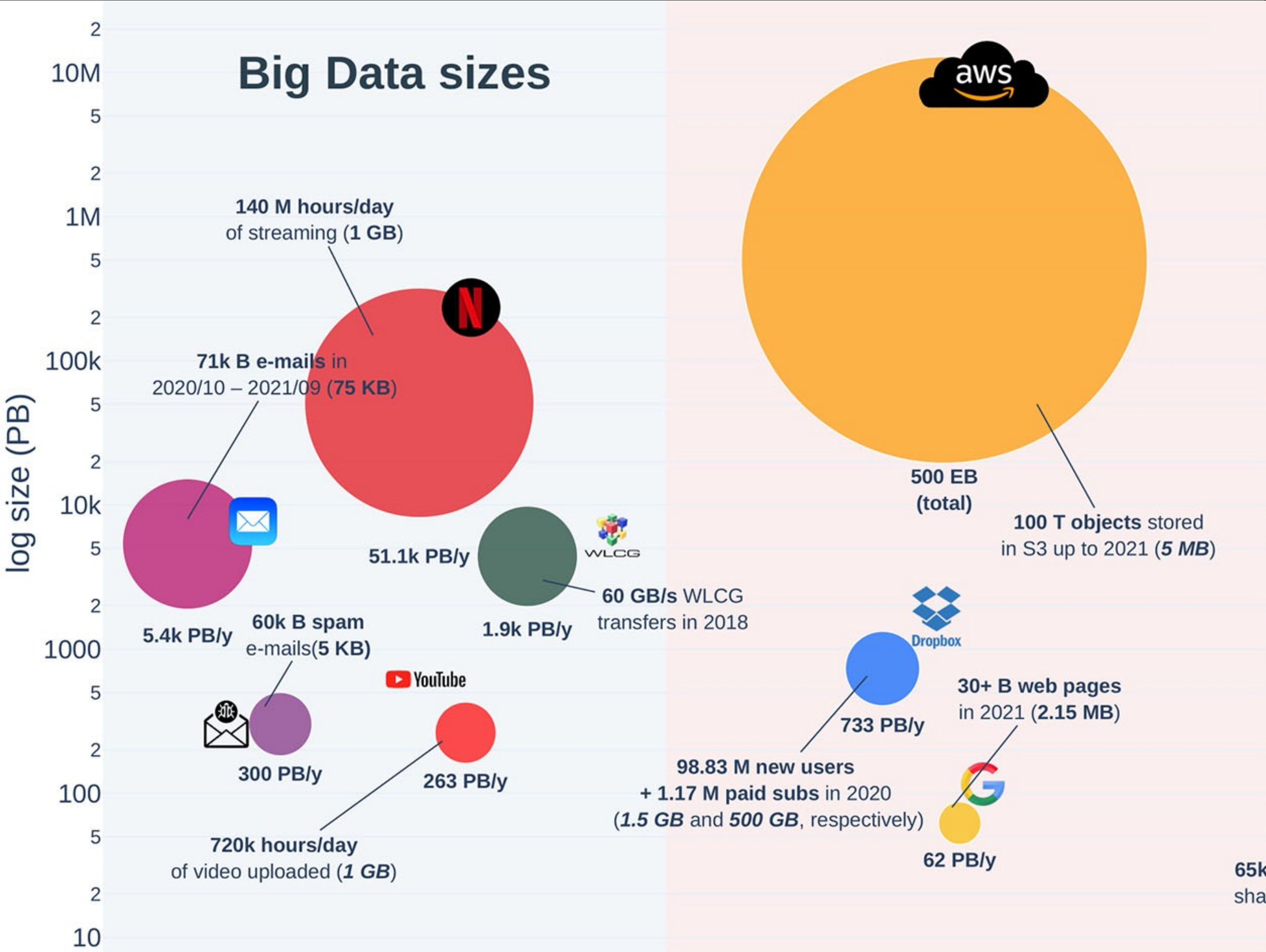
Big Data sizes



Big Data sizes



Big Data sizes



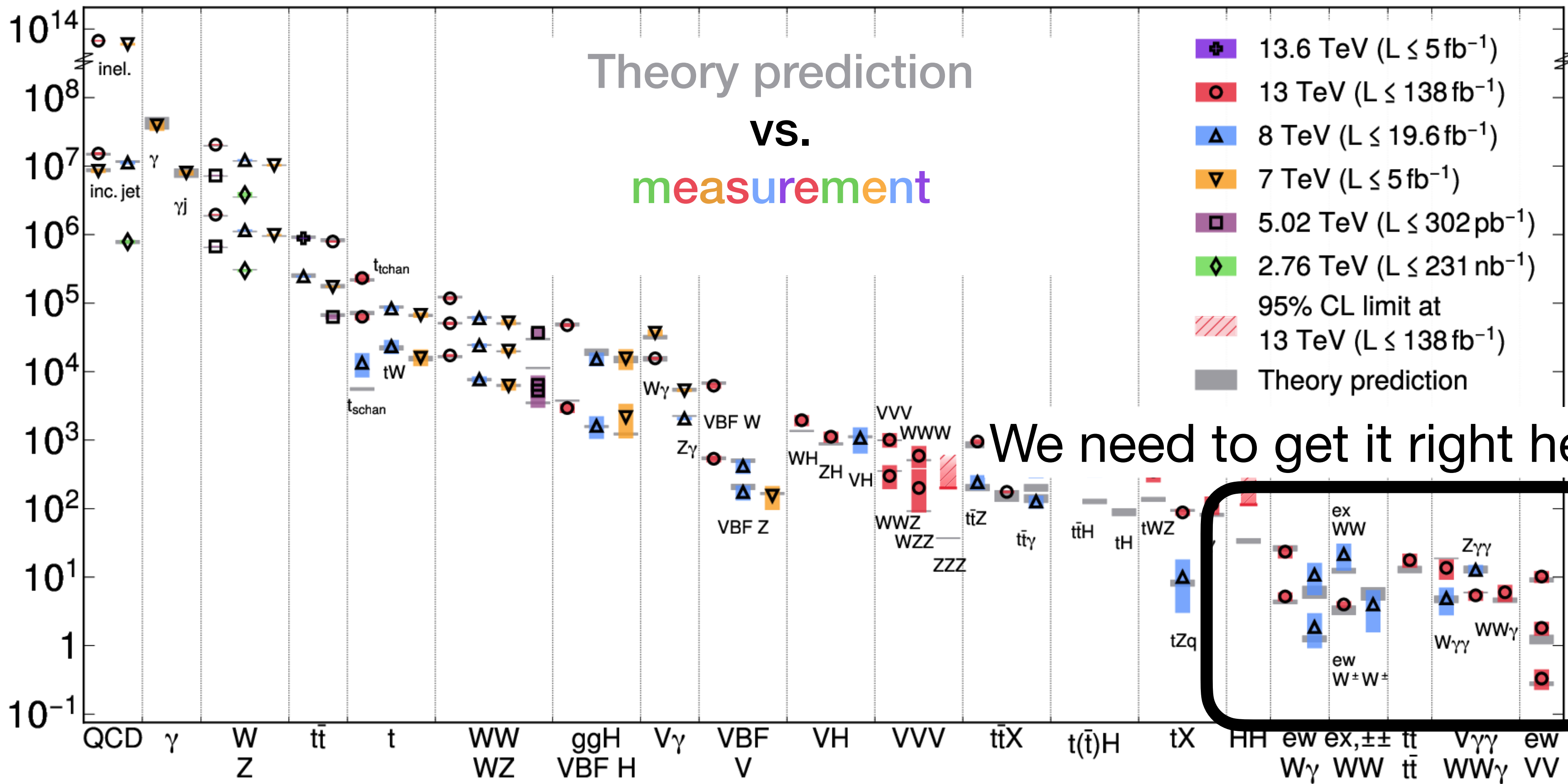


“While the humans guessed correctly 70 percent of the time, RankBrain had an 80 percent success rate [...]”

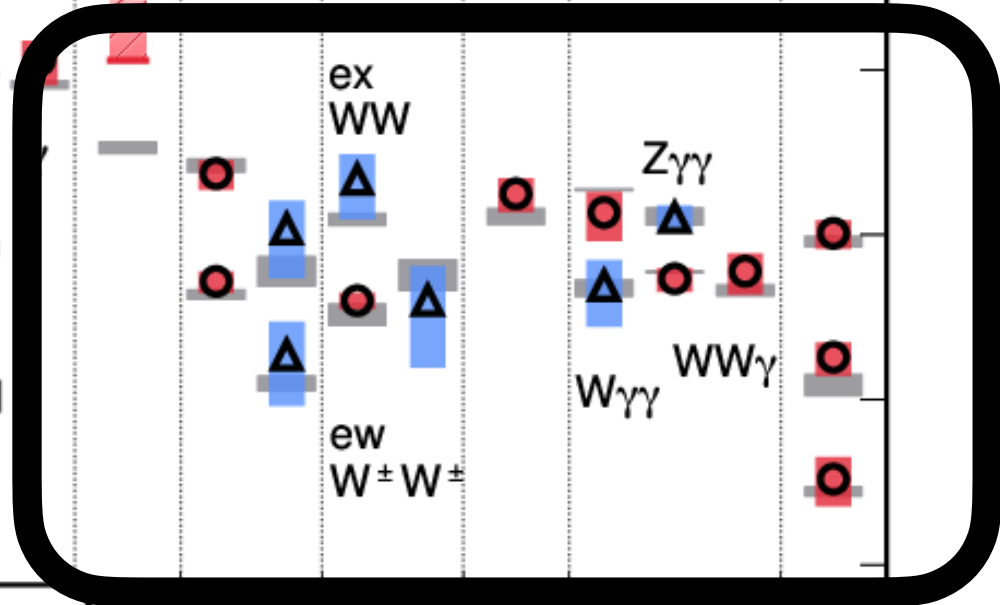
CMS


Theory prediction
vs.
measurement

Production frequency



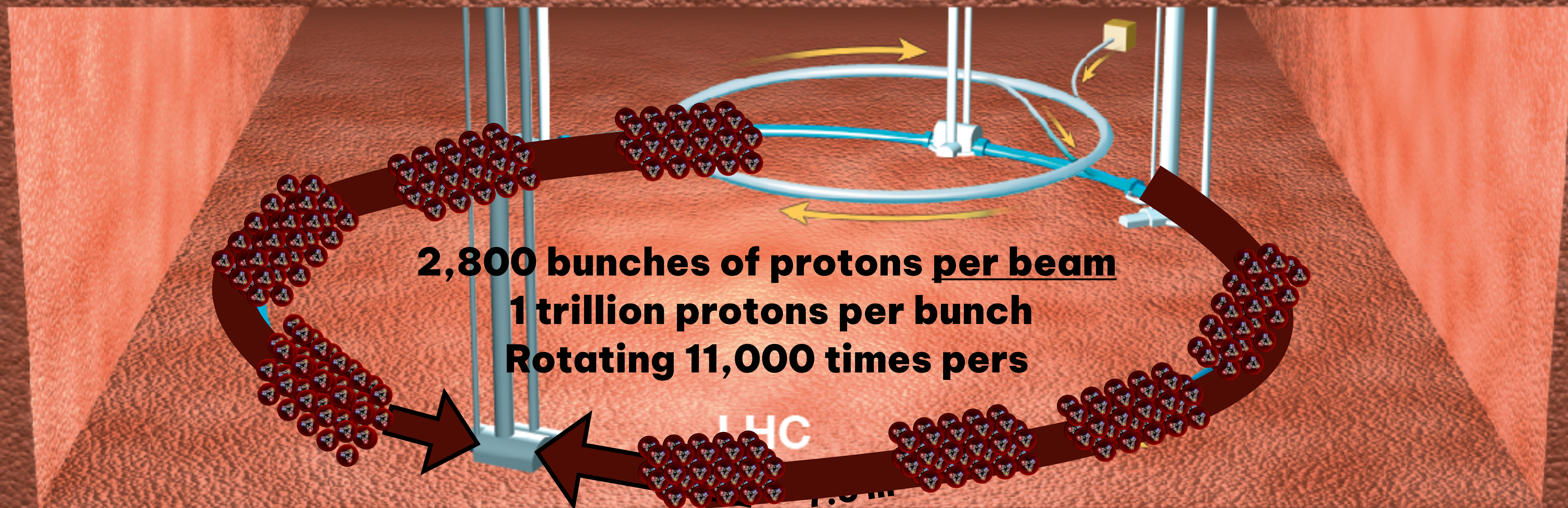
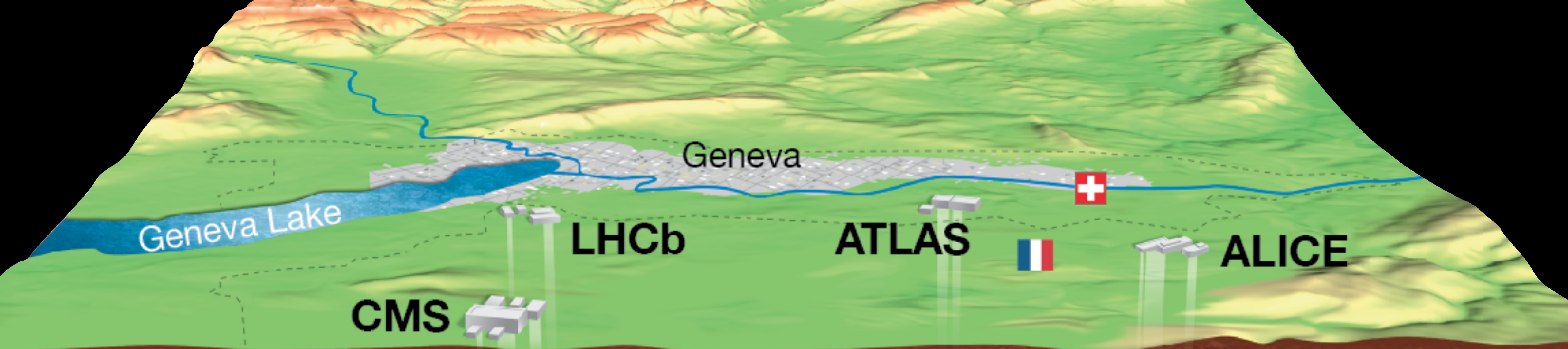
We need to get it right here



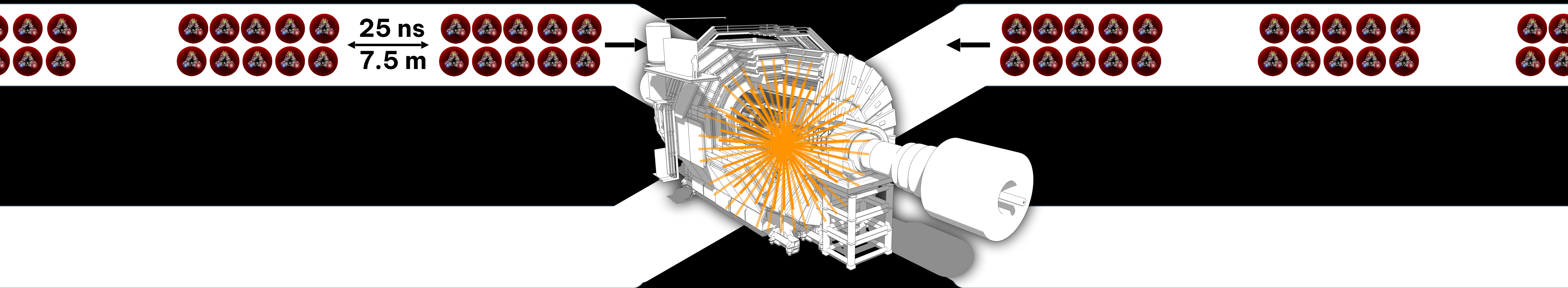
A close-up photograph of a large stack of needles, viewed from above. The needles are arranged in a dense, circular pattern, creating a radial effect. One needle, located in the lower right quadrant, is highlighted in a bright blue color, standing out from the otherwise uniform silver needles. The background is dark and slightly blurred, showing some structural elements of the container.

Get it right one in a one hundred trillionth

Machine Learning and the **needle in the needle stack**



The bunches “collide” every 25 ns (40 MHz)



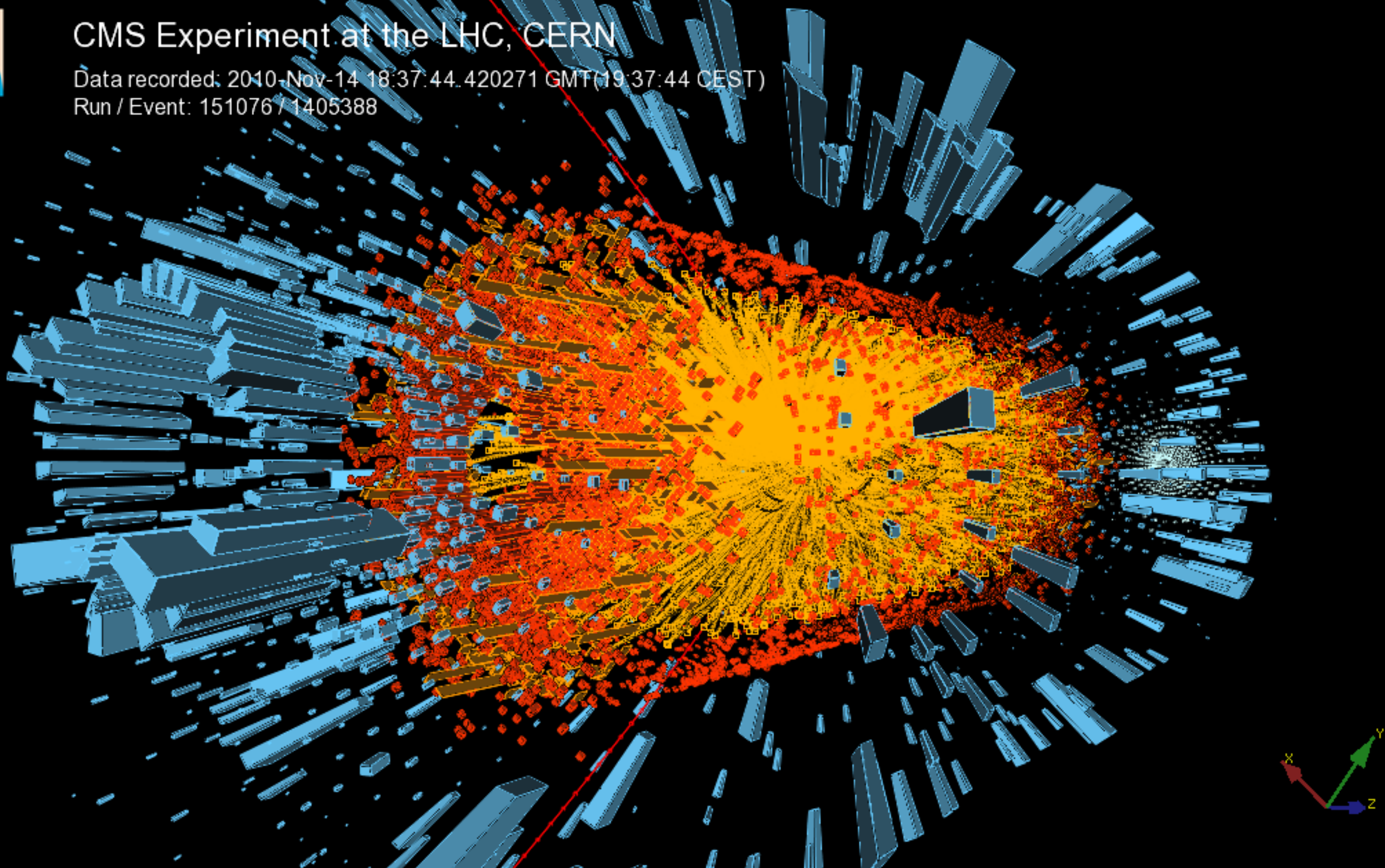
~Out of 10^{10} protons, only 60 collisions per crossing
Interaction is extremely rare!



CMS Experiment at the LHC, CERN

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

Run / Event: 151076 / 1405388



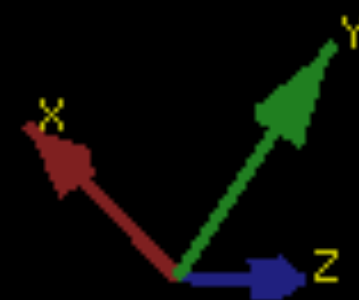


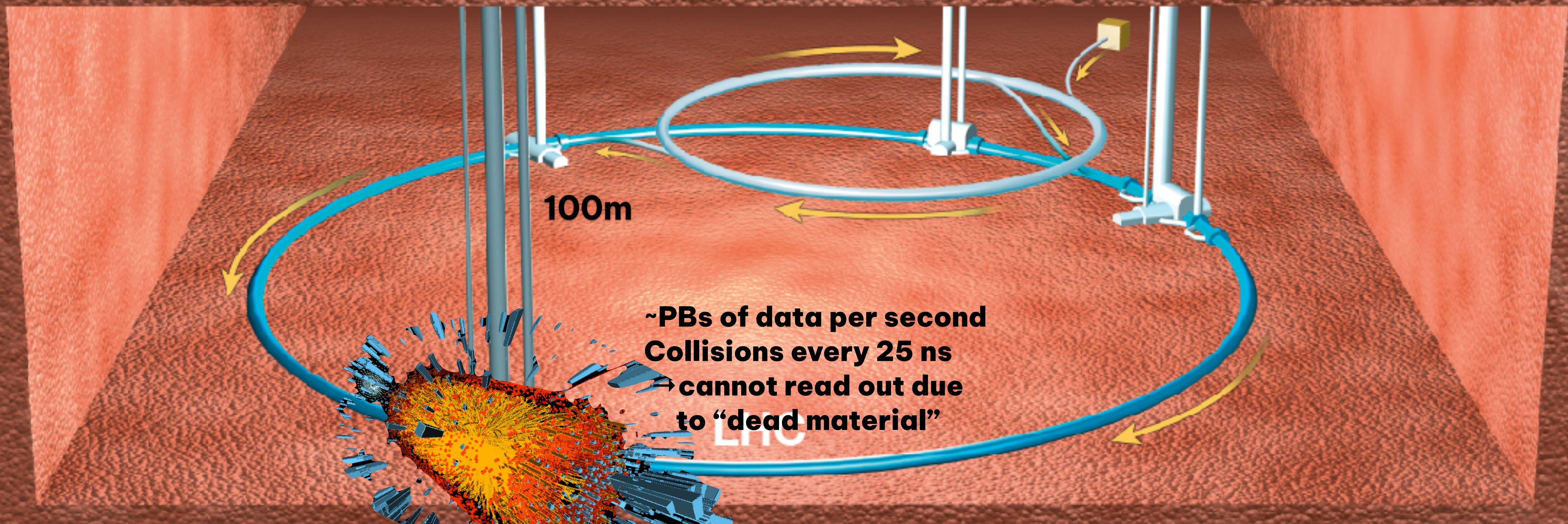
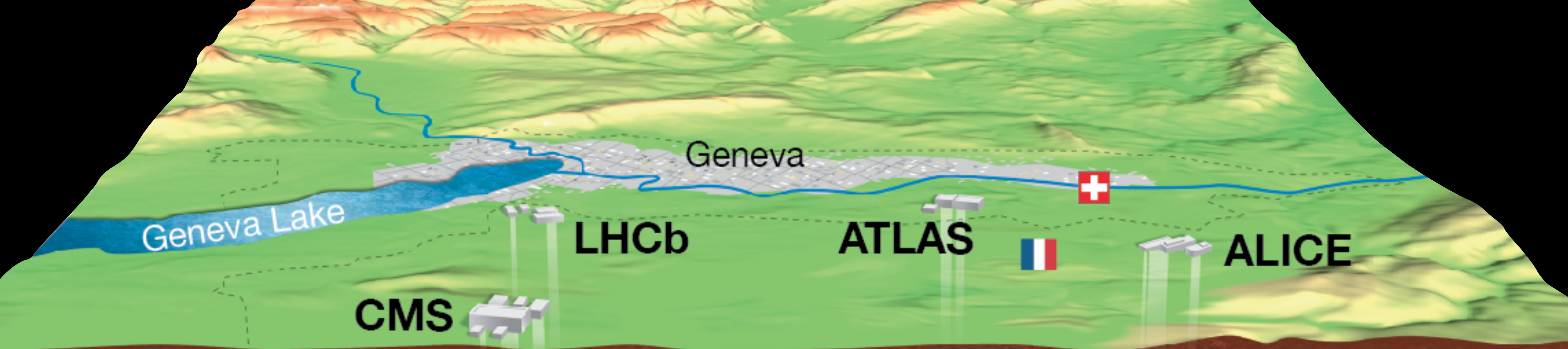
CMS Experiment at the LHC, CERN

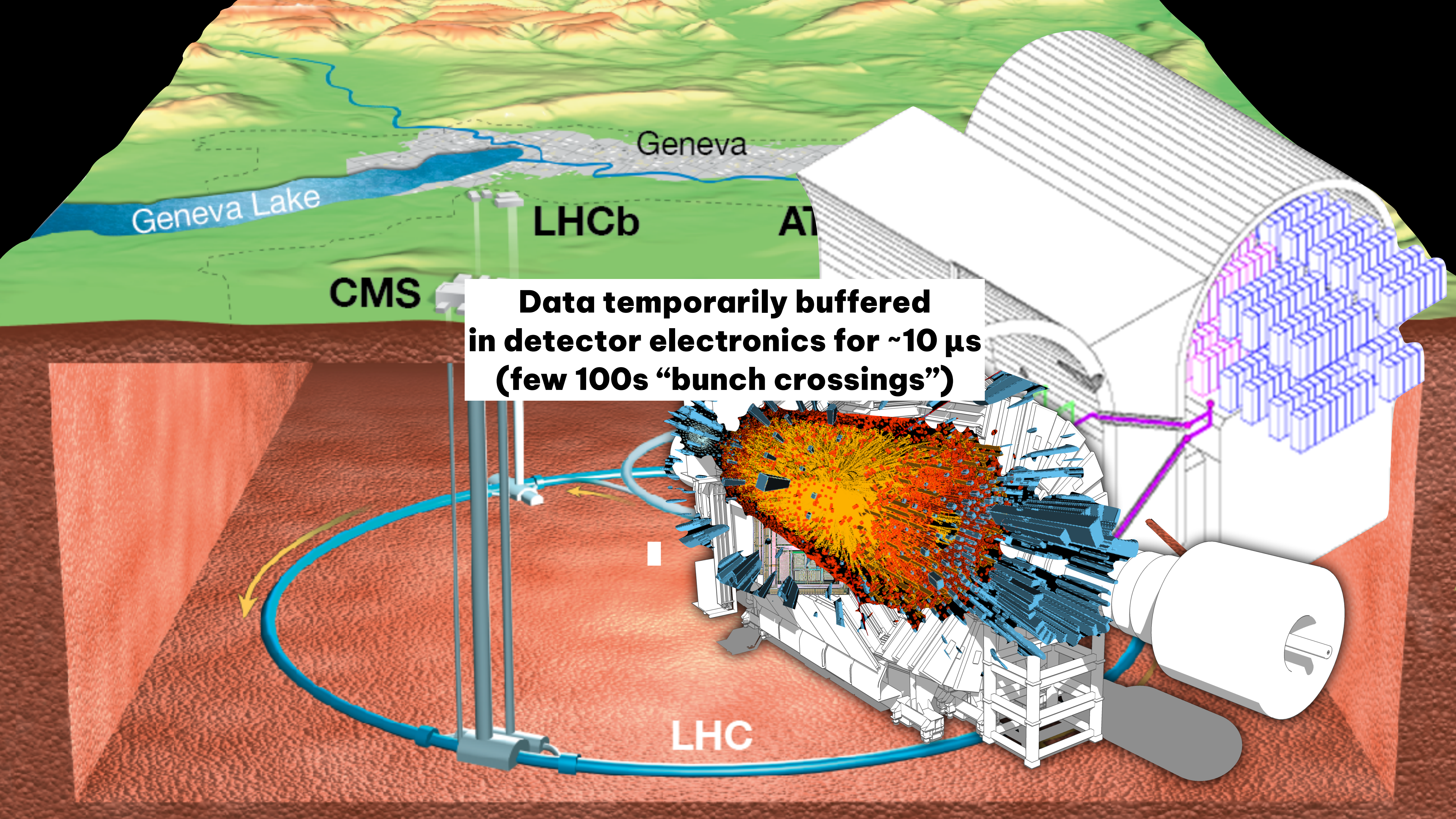
Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST)

Run / Event: 151076 / 1405388

1 collision \cong $O(1)$ MB
 $\sim O(1)$ billion collisions per second
 $\rightarrow O(1)$ PB / second







Geneva Lake

Geneva

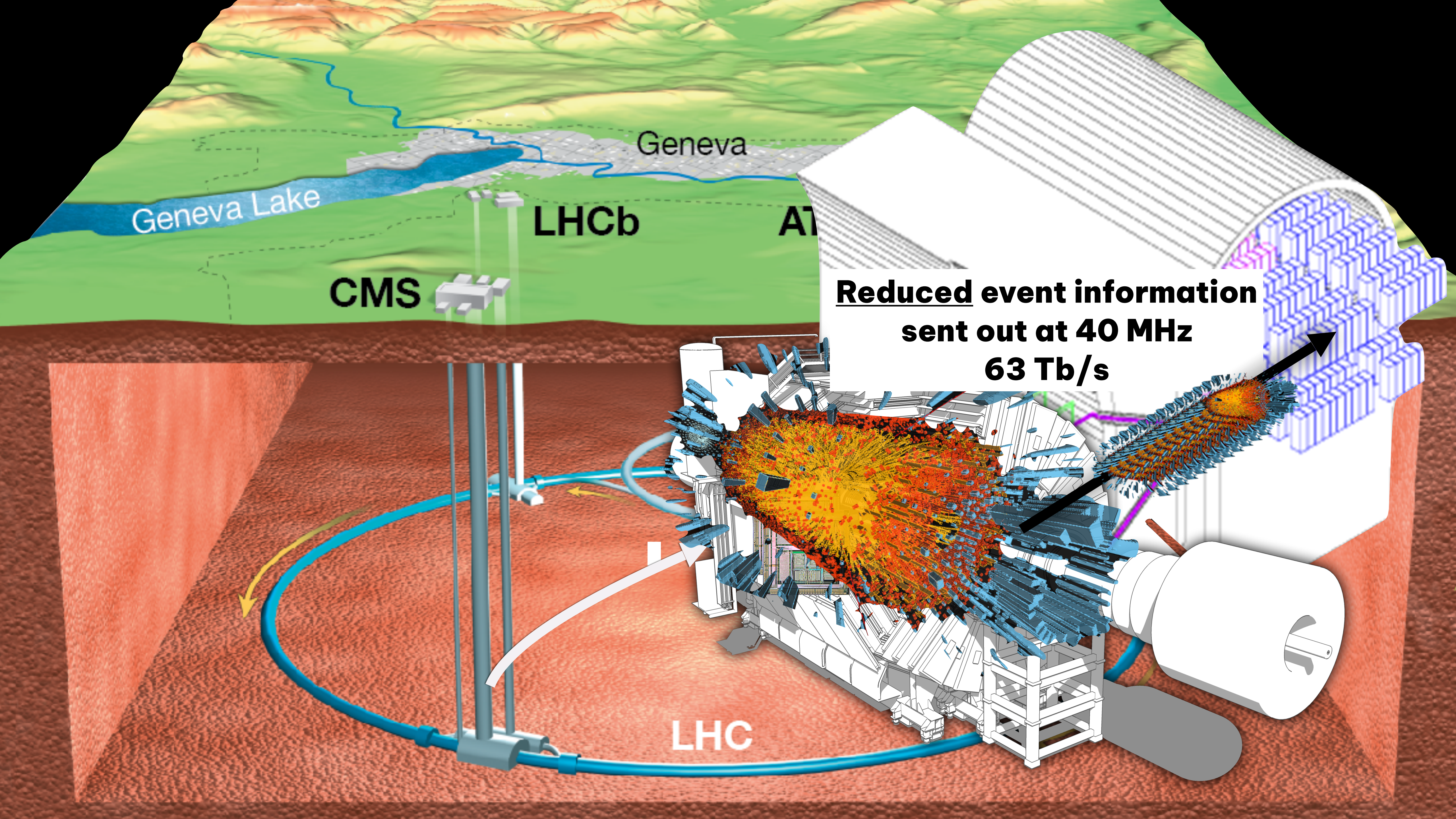
CMS

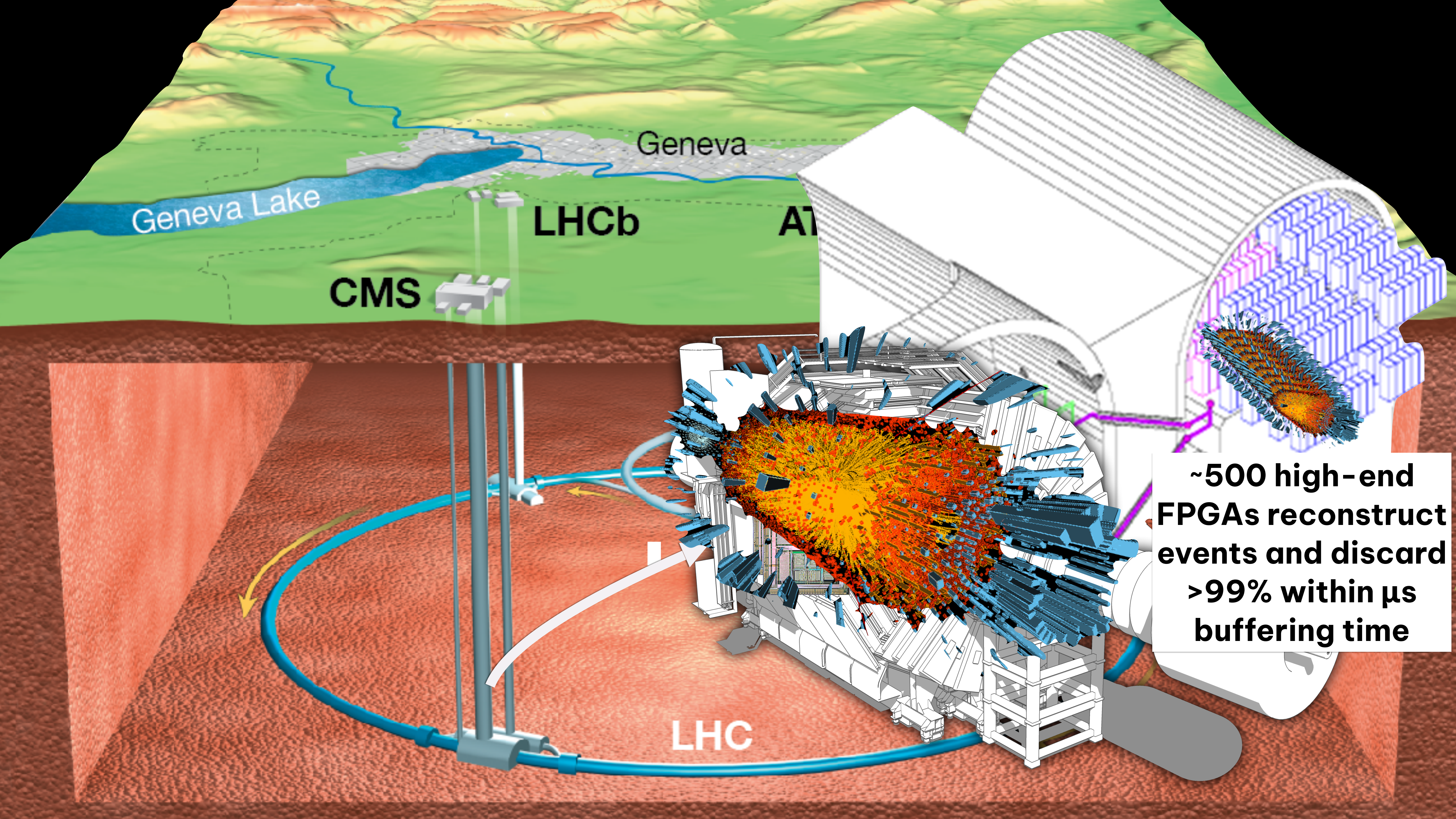
LHCb

ATLAS

**Data temporarily buffered
in detector electronics for $\sim 10 \mu\text{s}$
(few 100s “bunch crossings”)**

LHC





Geneva

Geneva Lake

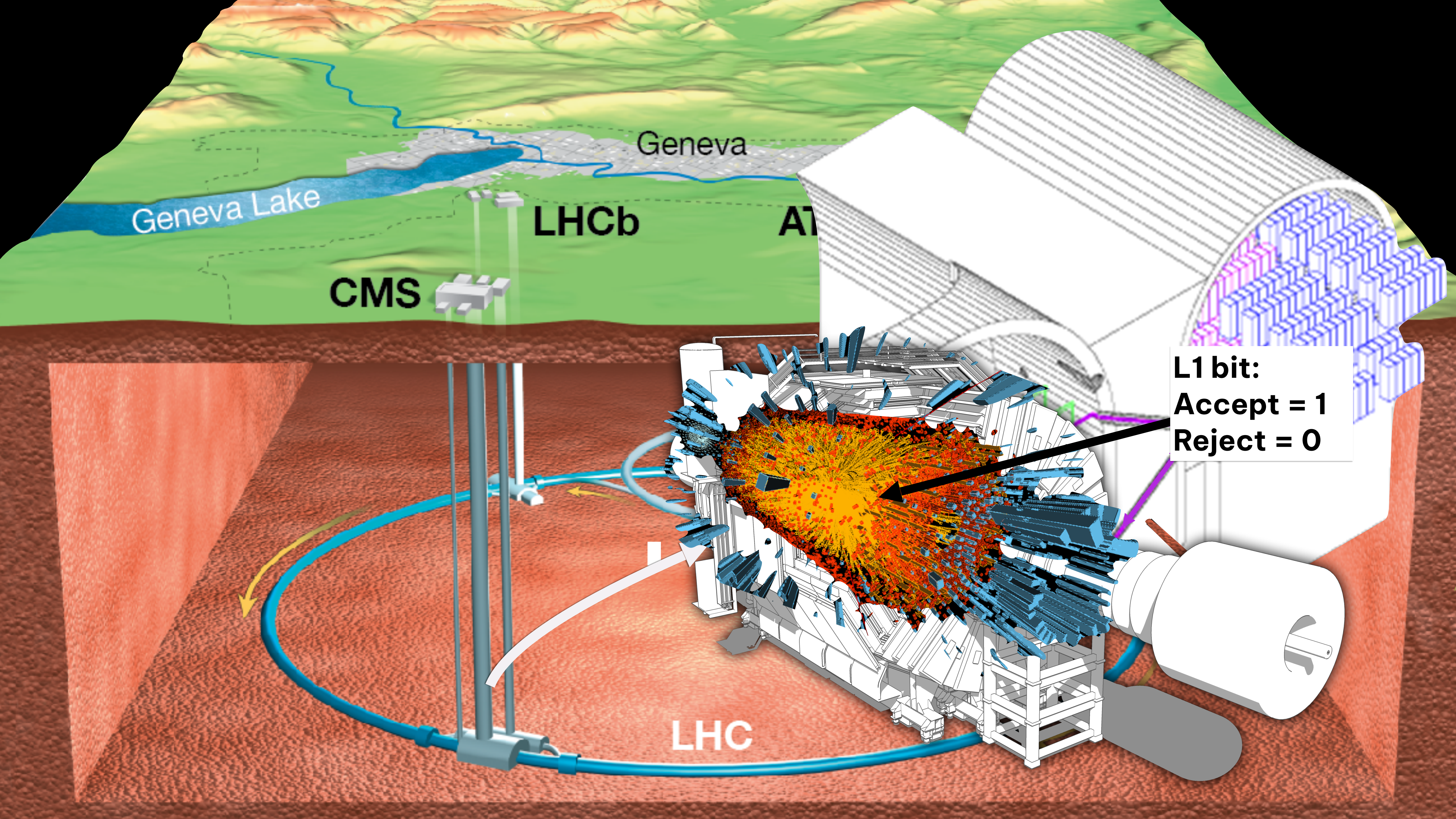
CMS

LHCb

ATLAS

LHC

**~500 high-end
FPGAs reconstruct
events and discard
>99% within μ s
buffering time**



Geneva

Geneva Lake

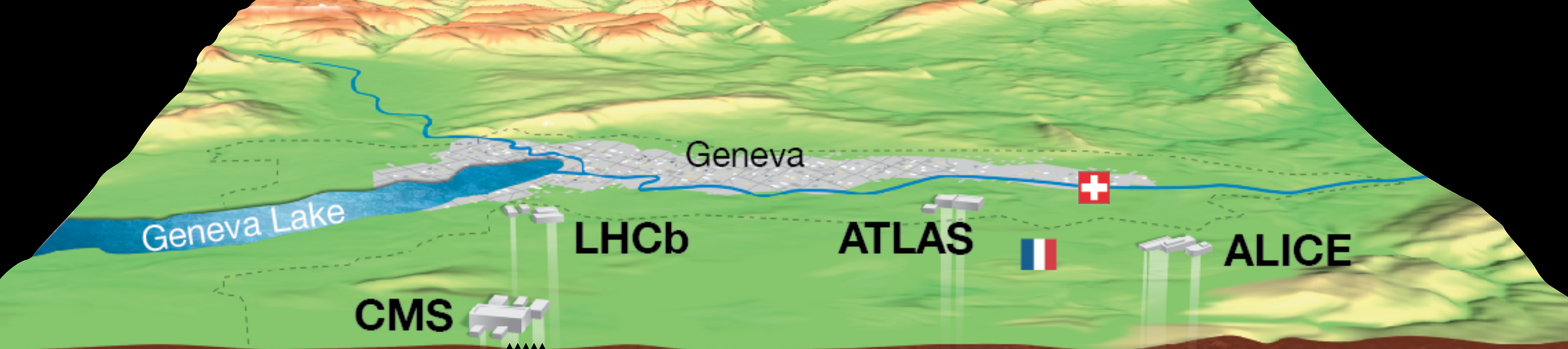
LHCb

ATLAS

CMS

LHC

L1 bit:
Accept = 1
Reject = 0

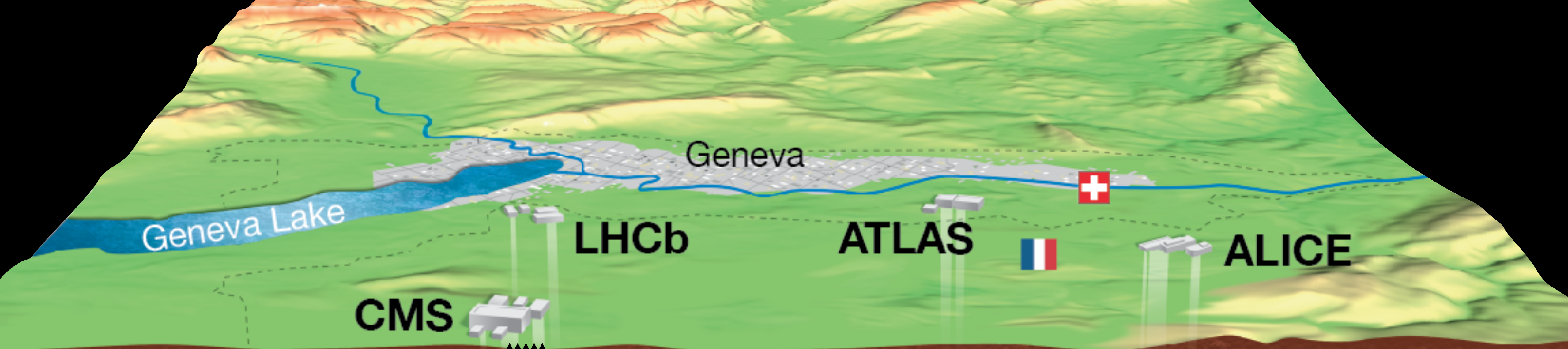


x560 100 Gbps links
200m long fibers

0.2% of events left
750 kHz, O(1)TB/s

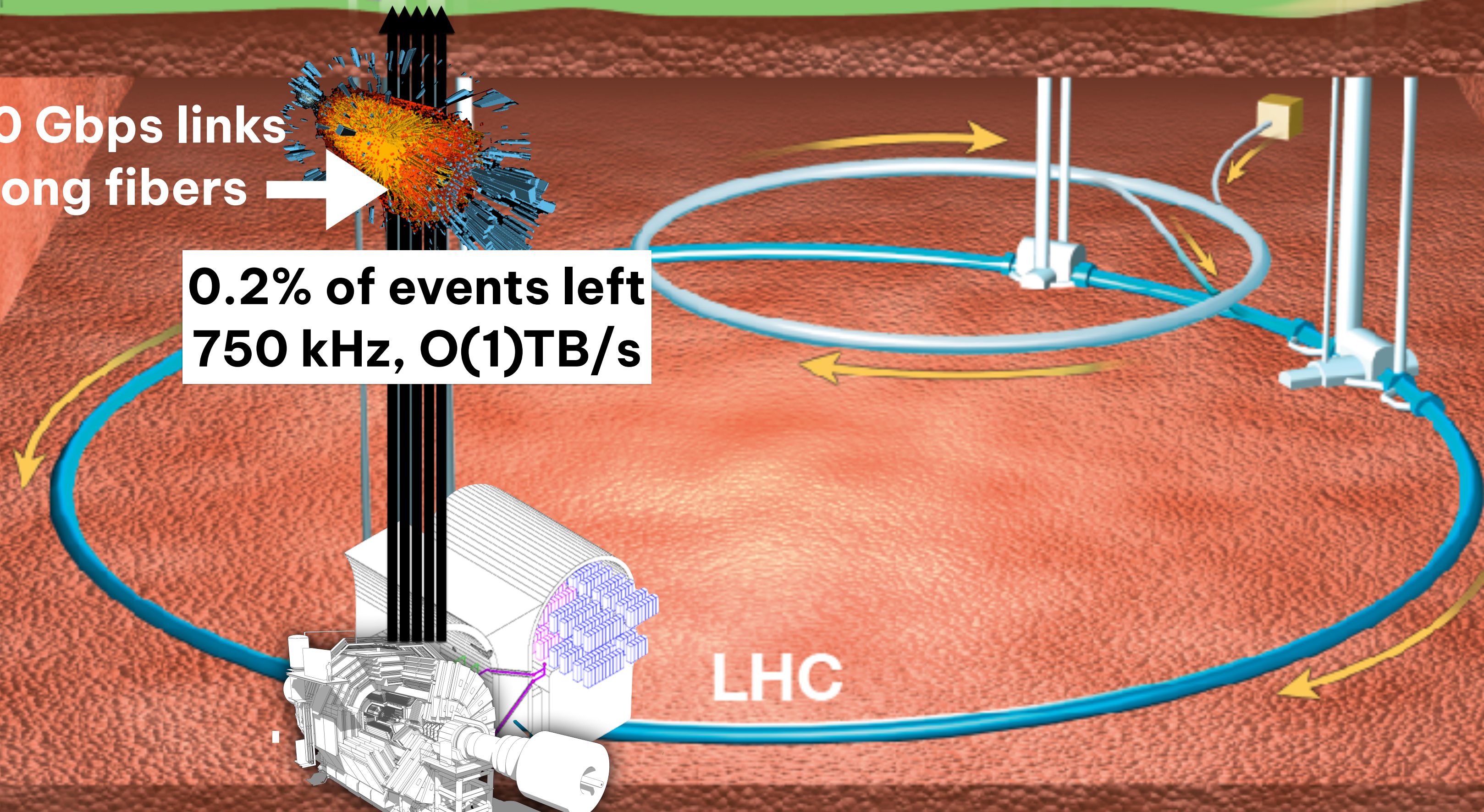
LHC

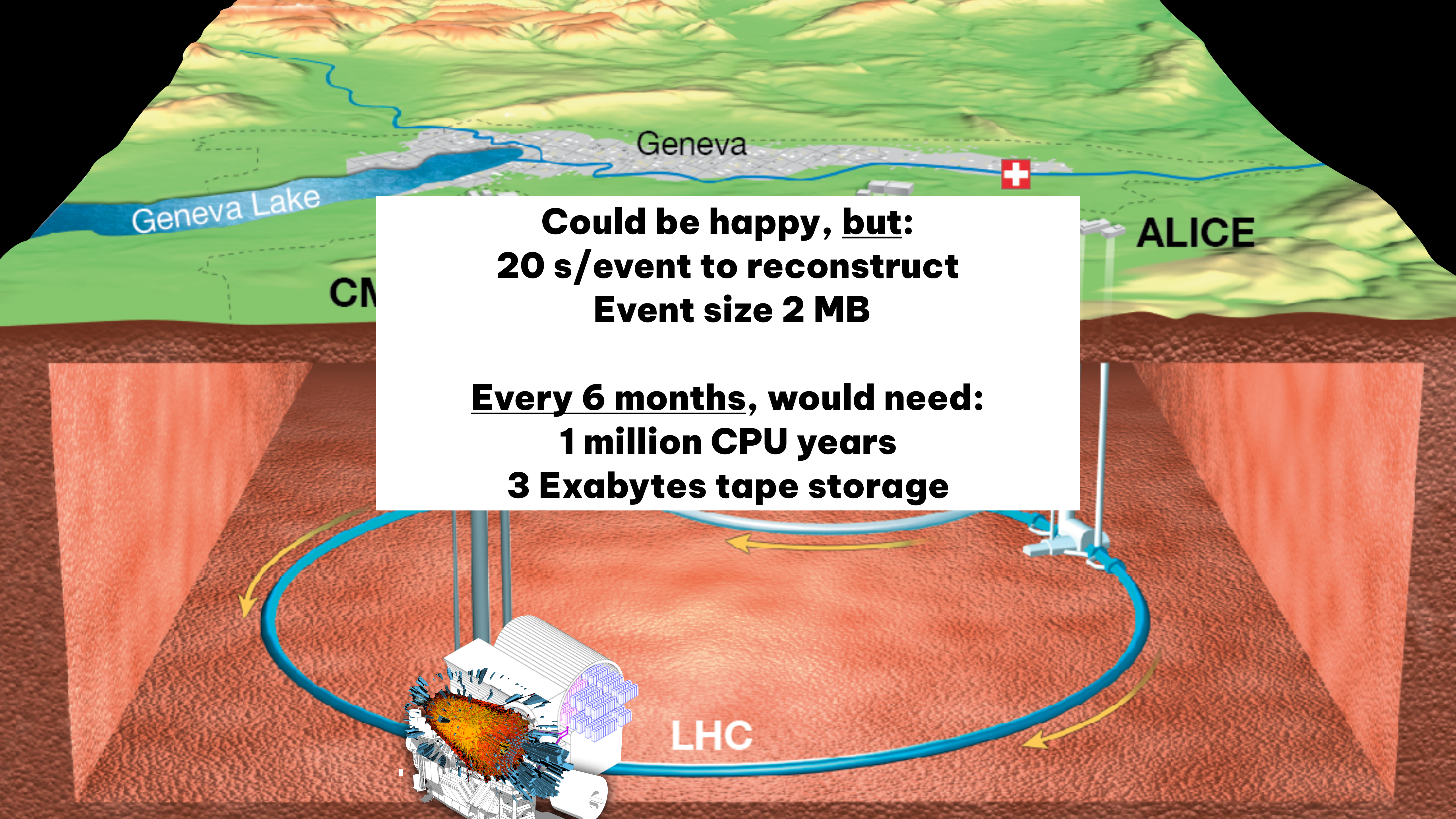
This diagram illustrates the LHC tunnel and the fiber optic links connecting it to the experiments. The LHC is shown at the bottom, with a large number of fiber optic links extending from it. The links are shown as blue lines forming a circular path around the LHC. Yellow arrows indicate the direction of data flow. The text 'x560 100 Gbps links' and '200m long fibers' is shown on the left, and '0.2% of events left' and '750 kHz, O(1)TB/s' is shown in a white box in the center. The LHC is labeled at the bottom.



x560 100 Gbps links
200m long fibers

0.2% of events left
750 kHz, O(1)TB/s





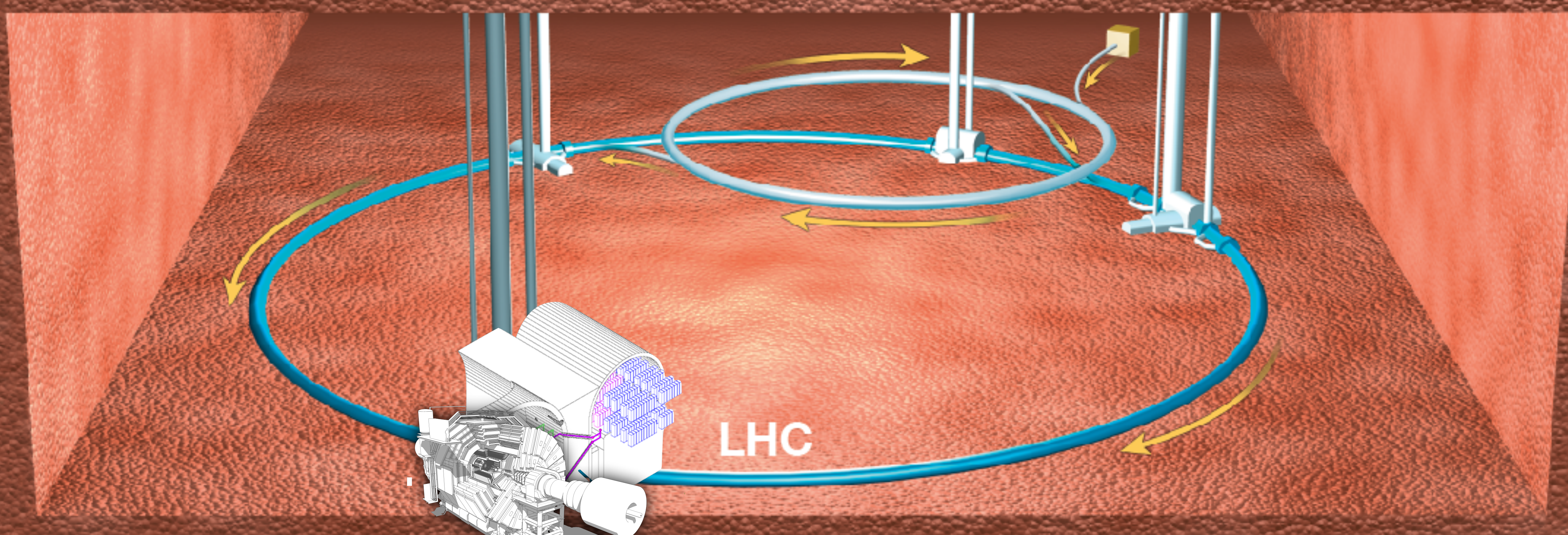
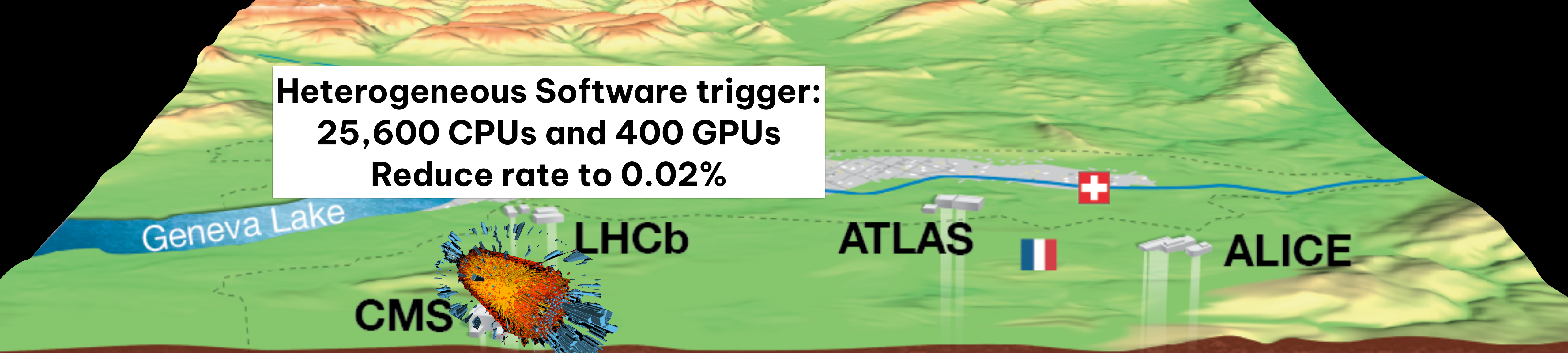
**Could be happy, but:
20 s/event to reconstruct
Event size 2 MB**

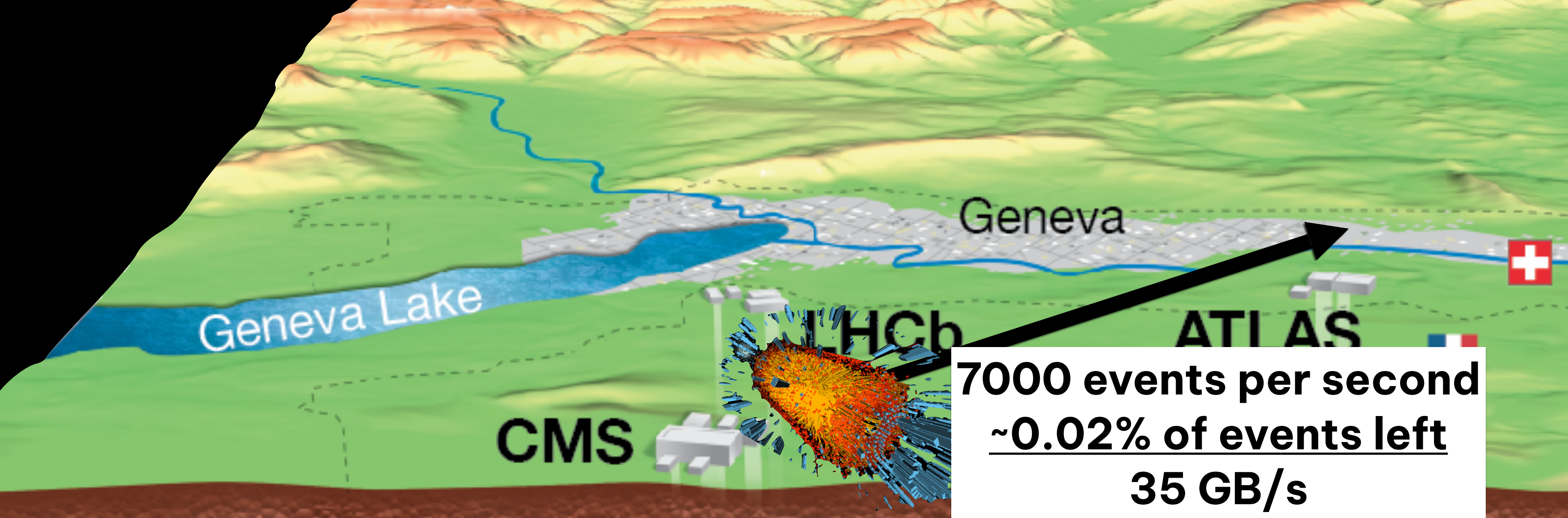
**Every 6 months, would need:
1 million CPU years
3 Exabytes tape storage**

ALICE

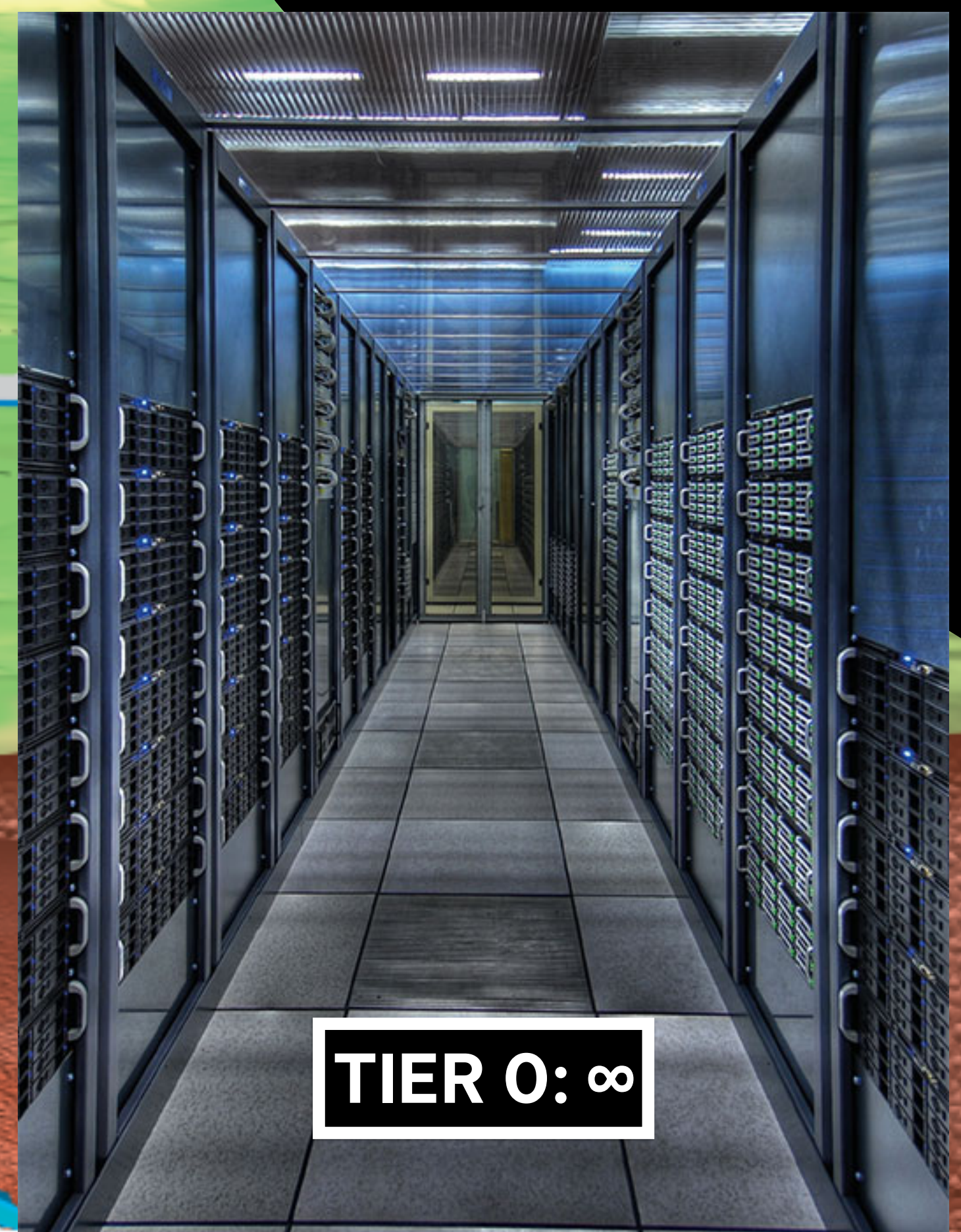
LHC

**Heterogeneous Software trigger:
25,600 CPUs and 400 GPUs
Reduce rate to 0.02%**



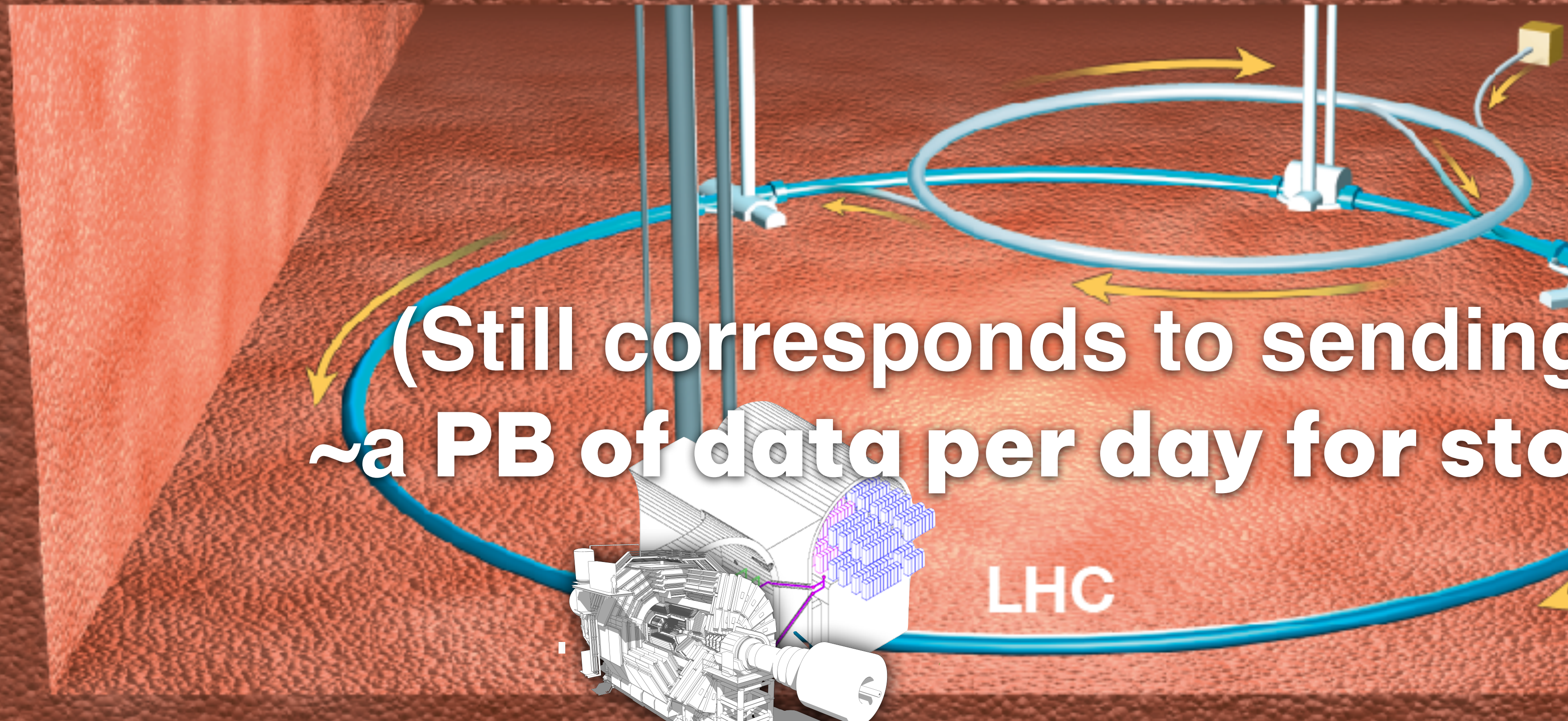


7000 events per second
~0.02% of events left
35 GB/s

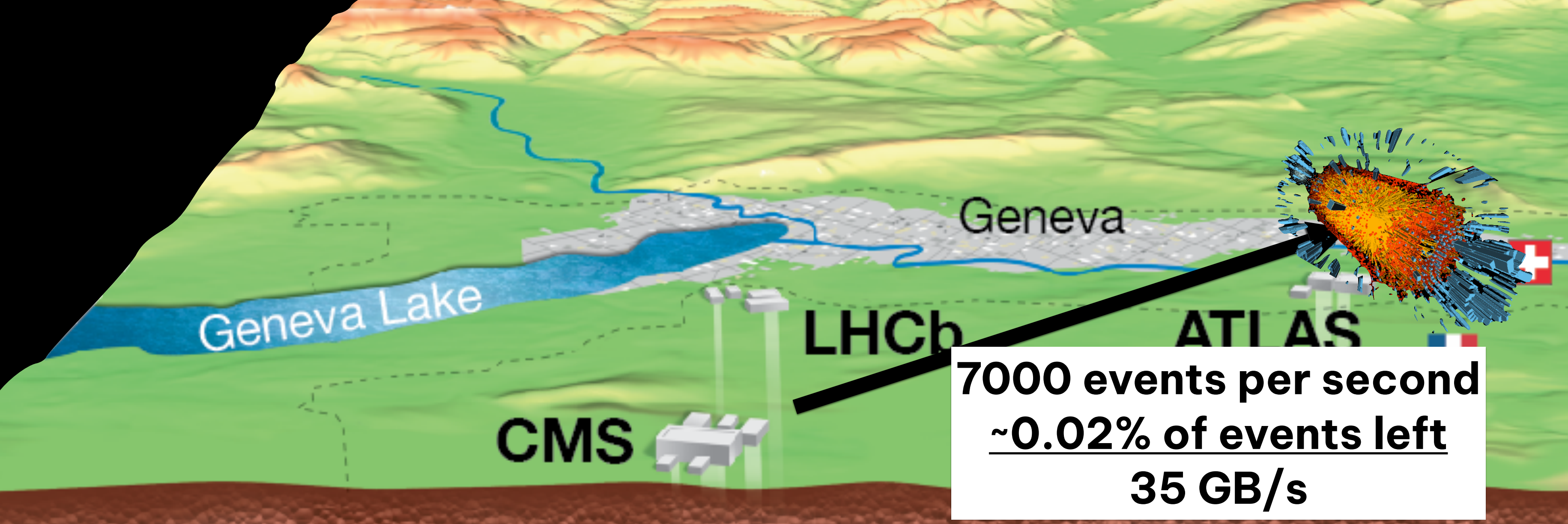


TIER 0: ∞

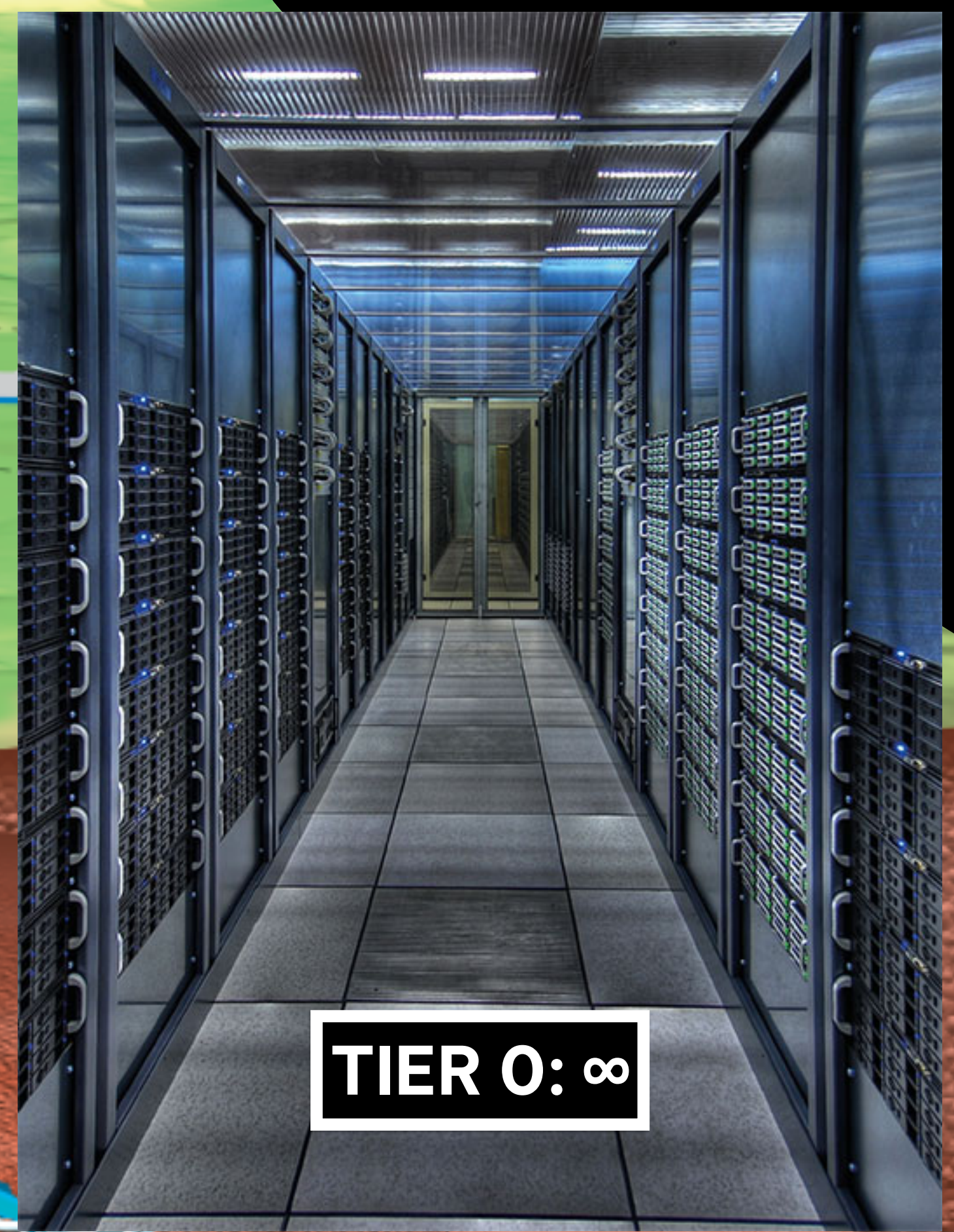
**(Still corresponds to sending out
~a PB of data per day for storage)**



LHC

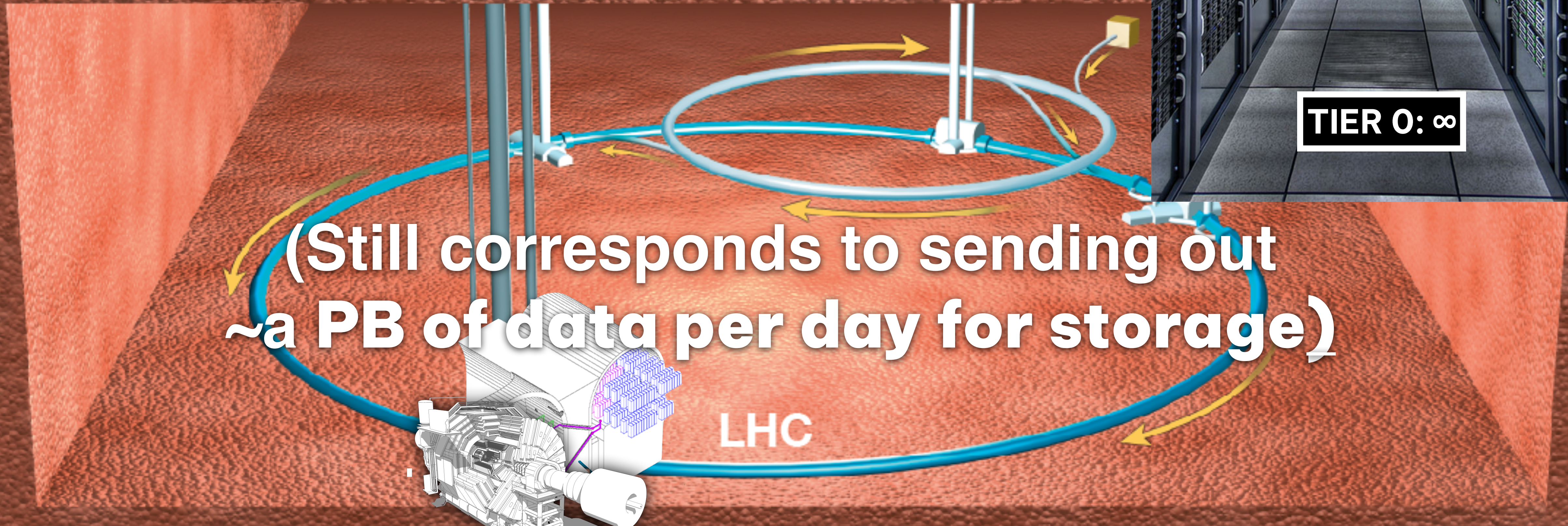


7000 events per second
~0.02% of events left
35 GB/s

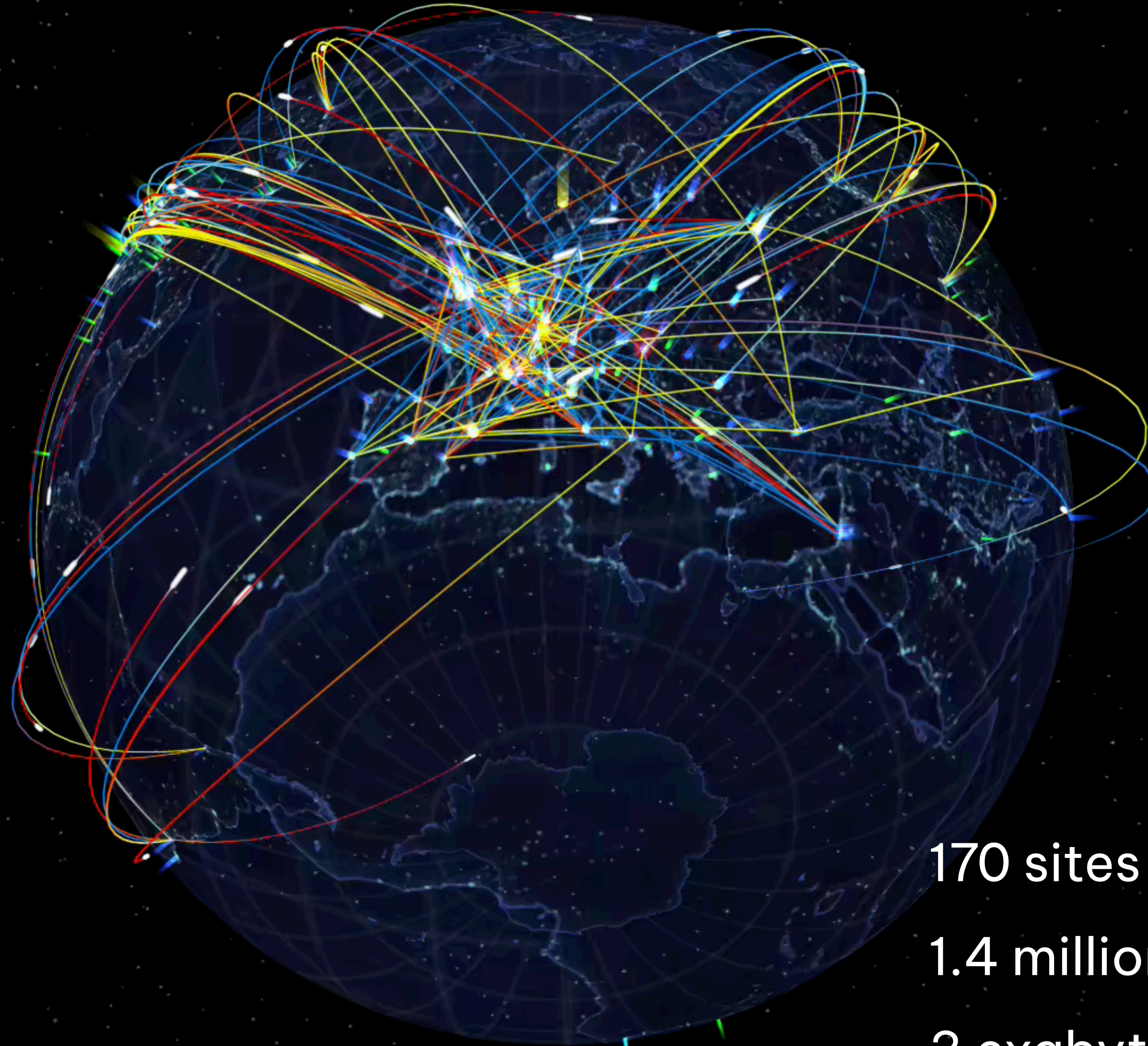


TIER 0: ∞

(Still corresponds to sending out
~a PB of data per day for storage)



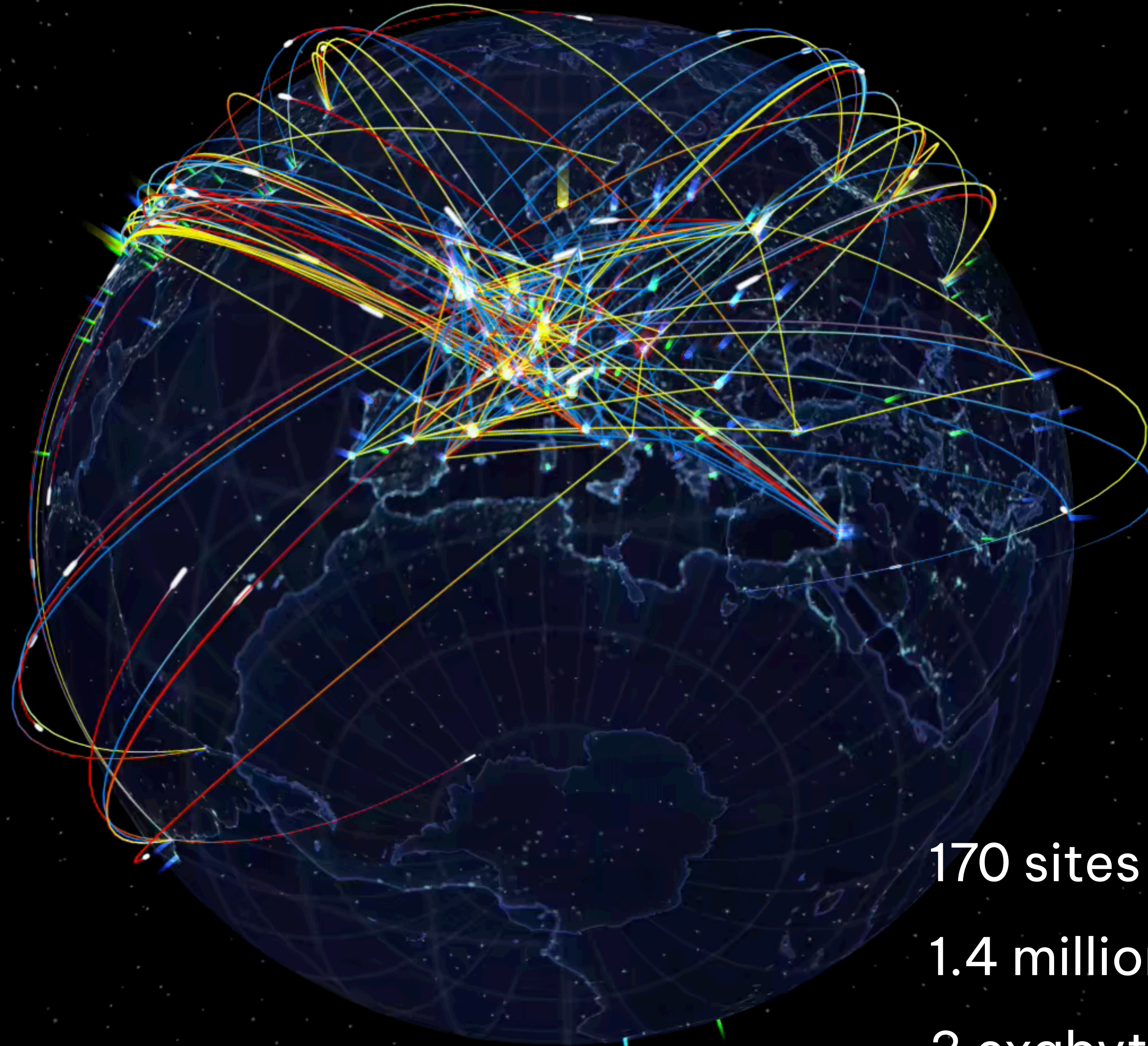
The Worldwide LHC Computing Grid



170 sites in 42 countries
1.4 million computer cores
3 exabytes of storage

Running jobs: 365644
Active CPU cores: 807139
Transfer rate: 21.54 GiB/sec

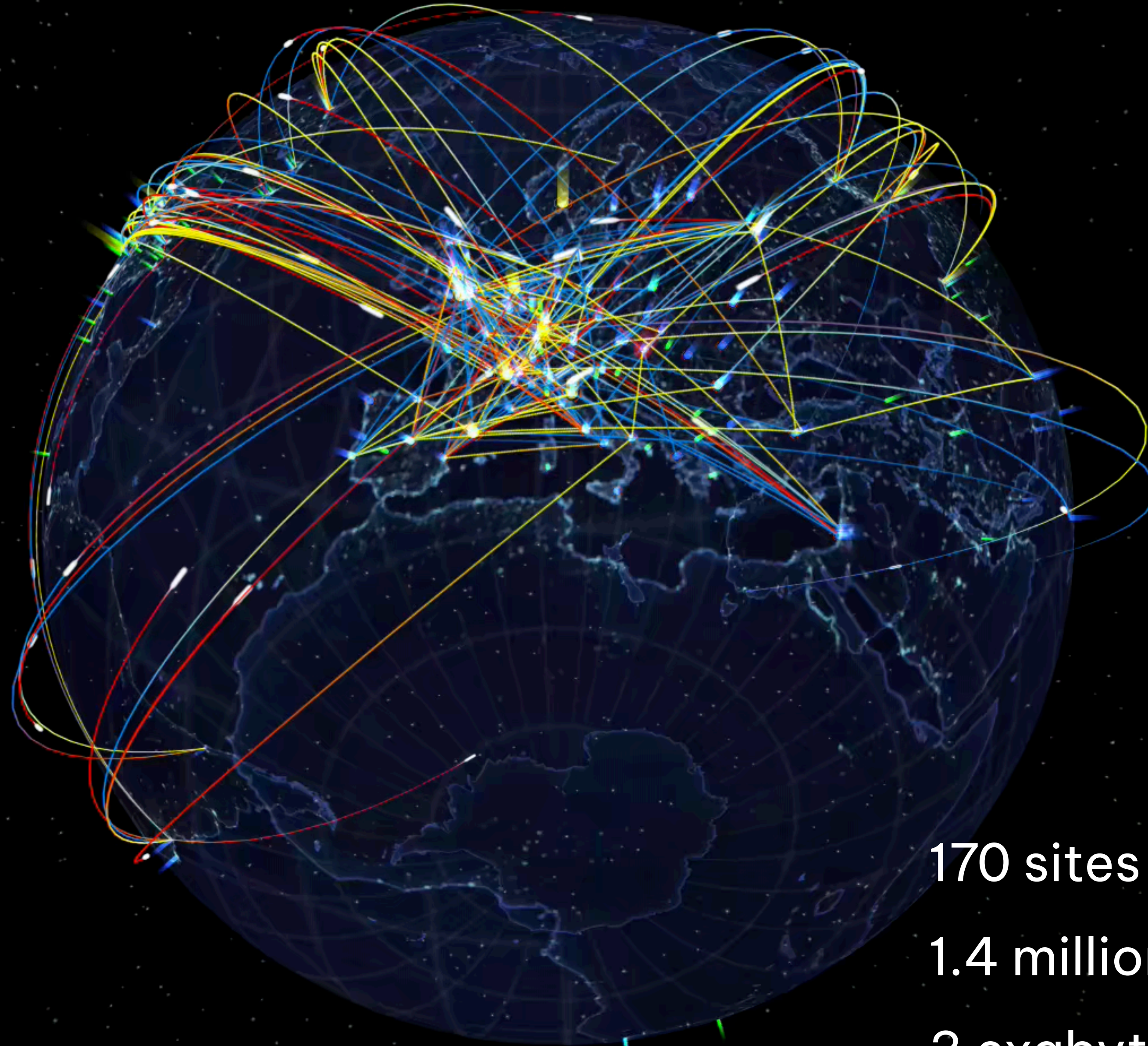
The Worldwide LHC Computing Grid



170 sites in 42 countries
1.4 million computer cores
3 exabytes of storage

Running jobs: 365644
Active CPU cores: 807139
Transfer rate: 21.54 GiB/sec

The Worldwide LHC Computing Grid



170 sites in 42 countries
1.4 million computer cores
3 exabytes of storage

Running jobs: 365644
Active CPU cores: 807139
Transfer rate: 21.54 GiB/sec

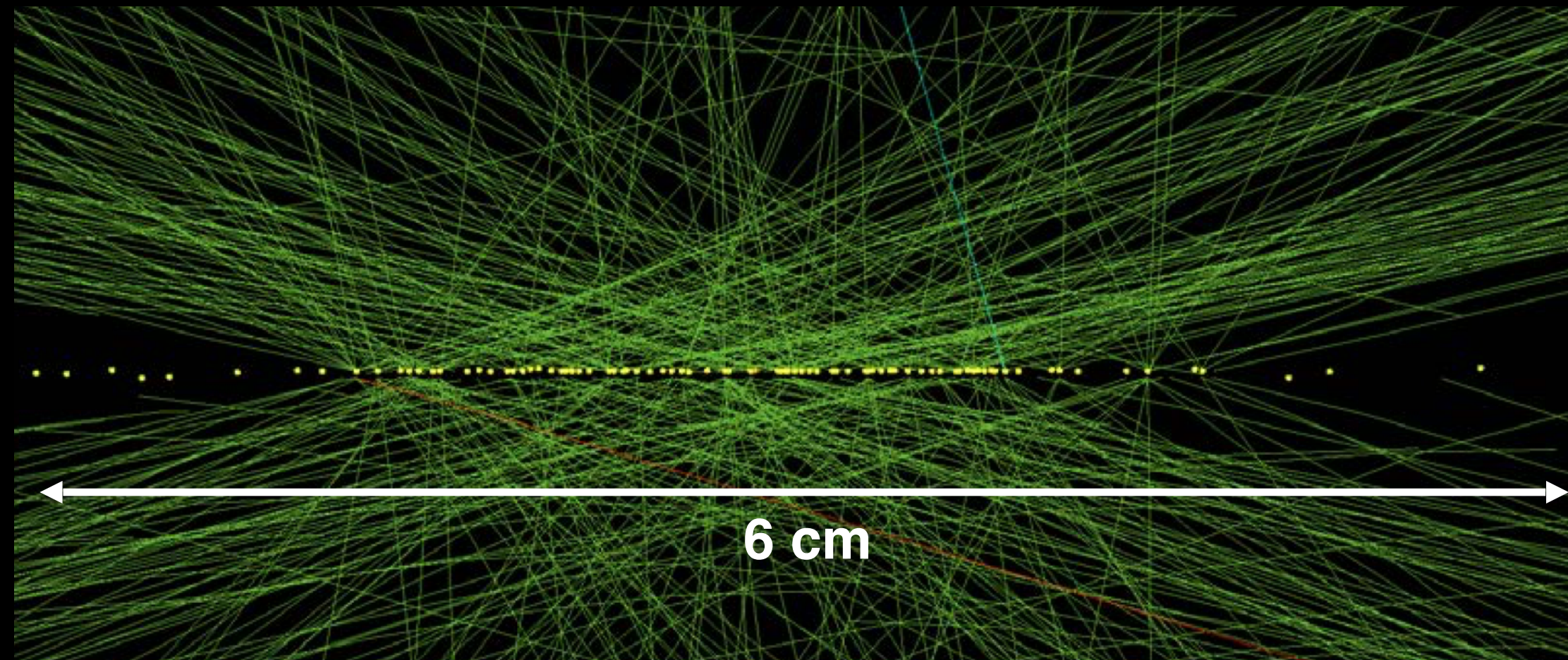
The Standard Model

$$\begin{aligned}
& -\frac{1}{2}\partial_\nu g_\mu^a \partial_\nu g_\mu^a - g_s f^{abc} \partial_\mu g_\nu^a g_\mu^b g_\nu^c - \frac{1}{4}g^2 f^{abc} f^{ade} g_\nu^b g_\nu^c g_\mu^d g_\nu^e + \\
& \frac{1}{2}ig_s^2(\bar{q}_i^\sigma \gamma^\mu q_j^\sigma)g_\mu^a + \bar{G}^a \partial^2 G^a + g_s f^{abc} \partial_\mu \bar{G}^a G^b g_\mu^c - \partial_\nu W_\mu^+ \partial_\nu W_\mu^- - \\
& M^2 W_\mu^+ W_\mu^- - \frac{1}{2}\partial_\nu Z_\mu^0 \partial_\nu Z_\mu^0 - \frac{1}{2c_w^2}M^2 Z_\mu^0 Z_\mu^0 - \frac{1}{2}\partial_\mu A_\nu \partial_\mu A_\nu - \frac{1}{2}\partial_\mu H \partial_\mu H - \\
& \frac{1}{2}m_h^2 H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - M^2 \phi^+ \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \frac{1}{2c_w^2}M\phi^0 \phi^0 - \beta_h[\frac{2M^2}{g^2} + \\
& \frac{2M}{g}H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-)] + \frac{2M^4}{g^2}\alpha_h - igc_w[\partial_\nu Z_\mu^0(W_\mu^+ W_\nu^- - \\
& W_\nu^+ W_\mu^-) - Z_\nu^0(W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) + Z_\mu^0(W_\nu^+ \partial_\nu W_\mu^- - \\
& W_\nu^- \partial_\nu W_\mu^+)] - ig s_w[\partial_\nu A_\mu(W_\mu^+ W_\nu^- - W_\nu^+ W_\mu^-) - A_\nu(W_\mu^+ \partial_\nu W_\mu^- - \\
& W_\mu^- \partial_\nu W_\mu^+) + A_\mu(W_\nu^+ \partial_\nu W_\mu^- - W_\nu^- \partial_\nu W_\mu^+)] - \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\nu^+ W_\nu^- + \\
& \frac{1}{2}g^2 W_\mu^+ W_\nu^- W_\mu^+ W_\nu^- + g^2 c_w^2(Z_\mu^0 W_\mu^+ Z_\nu^0 W_\nu^- - Z_\mu^0 Z_\nu^0 W_\mu^+ W_\nu^-) + \\
& g^2 s_w^2(A_\mu W_\mu^+ A_\nu W_\nu^- - A_\mu A_\nu W_\mu^+ W_\nu^-) + g^2 s_w c_w[A_\mu Z_\nu^0(W_\mu^+ W_\nu^- - \\
& W_\nu^+ W_\mu^-) - 2A_\mu Z_\mu^0 W_\nu^+ W_\nu^-] - g\alpha[H^3 + H\phi^0 \phi^0 + 2H\phi^+ \phi^-] - \\
& \frac{1}{8}g^2 \alpha_h[H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2] - \\
& gMW_\mu^+ W_\mu^- H - \frac{1}{2}g\frac{M}{c_w^2}Z_\mu^0 Z_\mu^0 H - \frac{1}{2}ig[W_\mu^+(\phi^0 \partial_\mu \phi^- - \phi^- \partial_\mu \phi^0) - \\
& W_\mu^-(\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)] + \frac{1}{2}g[W_\mu^+(H\partial_\mu \phi^- - \phi^- \partial_\mu H) - W_\mu^-(H\partial_\mu \phi^+ - \\
& \phi^+ \partial_\mu H)] + \frac{1}{2}g\frac{1}{c_w}(Z_\mu^0(H\partial_\mu \phi^0 - \phi^0 \partial_\mu H) - ig\frac{s_w^2}{c_w}MZ_\mu^0(W_\mu^+ \phi^- - W_\mu^- \phi^+) + \\
& ig s_w MA_\mu(W_\mu^+ \phi^- - W_\mu^- \phi^+) - ig\frac{1-2c_w^2}{2c_w}Z_\mu^0(\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) + \\
& ig s_w A_\mu(\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \frac{1}{4}g^2 W_\mu^+ W_\mu^- [H^2 + (\phi^0)^2 + 2\phi^+ \phi^-] - \\
& \frac{1}{4}g^2 \frac{1}{c_w^2}Z_\mu^0 Z_\mu^0 [H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)^2 \phi^+ \phi^-] - \frac{1}{2}g^2 \frac{s_w^2}{c_w}Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + \\
& W_\mu^- \phi^+) - \frac{1}{2}ig^2 \frac{s_w^2}{c_w}Z_\mu^0 H(W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\
& W_\mu^- \phi^+) + \frac{1}{2}ig^2 s_w A_\mu H(W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{s_w}{c_w}(2c_w^2 - 1)Z_\mu^0 A_\mu \phi^+ \phi^- - \\
& g^1 s_w^2 A_\mu A_\mu \phi^+ \phi^- - \bar{e}^\lambda (\gamma \partial + m_e^\lambda) e^\lambda - \bar{\nu}^\lambda \gamma \partial \nu^\lambda - \bar{u}_j^\lambda (\gamma \partial + m_u^\lambda) u_j^\lambda - \\
& \bar{d}_j^\lambda (\gamma \partial + m_d^\lambda) d_j^\lambda + ig s_w A_\mu [-(\bar{e}^\lambda \gamma^\mu e^\lambda) + \frac{2}{3}(\bar{u}_j^\lambda \gamma^\mu u_j^\lambda) - \frac{1}{3}(\bar{d}_j^\lambda \gamma^\mu d_j^\lambda)] + \\
& \frac{ig}{4c_w}Z_\mu^0[(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{e}^\lambda \gamma^\mu (4s_w^2 - 1 - \gamma^5) e^\lambda) + (\bar{u}_j^\lambda \gamma^\mu (\frac{4}{3}s_w^2 - \\
& 1 - \gamma^5) u_j^\lambda) + (\bar{d}_j^\lambda \gamma^\mu (1 - \frac{8}{3}s_w^2 - \gamma^5) d_j^\lambda)] + \frac{ig}{2\sqrt{2}}W_\mu^+[(\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + \\
& (\bar{u}_j^\lambda \gamma^\mu (1 + \gamma^5) C_{\lambda\kappa} d_j^\kappa)] + \frac{ig}{2\sqrt{2}}W_\mu^-[(\bar{e}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{d}_j^\kappa C_{\lambda\kappa}^\dagger \gamma^\mu (1 + \\
& \gamma^5) u_j^\lambda)] + \frac{ig}{2\sqrt{2}}\frac{m_d^\lambda}{M}[-\phi^+ (\bar{\nu}^\lambda (1 - \gamma^5) e^\lambda) + \phi^- (\bar{e}^\lambda (1 + \gamma^5) \nu^\lambda)] - \\
& \frac{g}{2}\frac{m_u^\lambda}{M}[H(\bar{e}^\lambda e^\lambda) + i\phi^0 (\bar{e}^\lambda \gamma^5 e^\lambda)] + \frac{ig}{2M\sqrt{2}}\phi^+ [-m_d^\kappa (\bar{u}_j^\lambda C_{\lambda\kappa} (1 - \gamma^5) d_j^\kappa) + \\
& m_u^\lambda (\bar{u}_j^\lambda C_{\lambda\kappa} (1 + \gamma^5) d_j^\kappa) + \frac{ig}{2M\sqrt{2}}\phi^- [m_d^\lambda (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 + \gamma^5) u_j^\kappa) - m_u^\kappa (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 - \\
& \gamma^5) u_j^\kappa) - \frac{g}{2}\frac{m_d^\lambda}{M}H(\bar{u}_j^\lambda u_j^\lambda) - \frac{g}{2}\frac{m_d^\lambda}{M}H(\bar{d}_j^\lambda d_j^\lambda) + \frac{ig}{2}\frac{m_d^\lambda}{M}\phi^0 (\bar{u}_j^\lambda \gamma^5 u_j^\lambda) - \\
& \frac{ig}{2}\frac{m_d^\lambda}{M}\phi^0 (\bar{d}_j^\lambda \gamma^5 d_j^\lambda) + \bar{X}^+ (\partial^2 - M^2) X^+ + \bar{X}^- (\partial^2 - M^2) X^- + \bar{X}^0 (\partial^2 - \\
& \frac{M^2}{c_w^2}) X^0 + \bar{Y} \partial^2 Y + igc_w W_\mu^+ (\partial_\mu \bar{X}^0 X^- - \partial_\mu \bar{X}^+ X^0) + ig s_w W_\mu^+ (\partial_\mu \bar{Y} X^- - \\
& \partial_\mu \bar{X}^+ Y) + igc_w W_\mu^- (\partial_\mu \bar{X}^- X^0 - \partial_\mu \bar{X}^0 X^+) + ig s_w W_\mu^- (\partial_\mu \bar{X}^- Y - \\
& \partial_\mu \bar{Y} X^+) + igc_w Z_\mu^0 (\partial_\mu \bar{X}^+ X^+ - \partial_\mu \bar{X}^- X^-) + ig s_w A_\mu (\partial_\mu \bar{X}^+ X^+ - \\
& \partial_\mu \bar{X}^- X^-) - \frac{1}{2}gM[\bar{X}^+ X^+ H + \bar{X}^- X^- H + \frac{1}{c_w^2}\bar{X}^0 X^0 H] + \\
& \frac{1-2c_w^2}{2c_w}igM[\bar{X}^+ X^0 \phi^+ - \bar{X}^- X^0 \phi^-] + \frac{1}{2c_w}igM[\bar{X}^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-] + \\
& igMs_w[\bar{X}^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-] + \frac{1}{2}igM[\bar{X}^+ X^+ \phi^0 - \bar{X}^- X^- \phi^0]
\end{aligned}$$

...and measuring its 19 free parameters to 5σ (= 0.00003% chance of statistical fluctuation)...

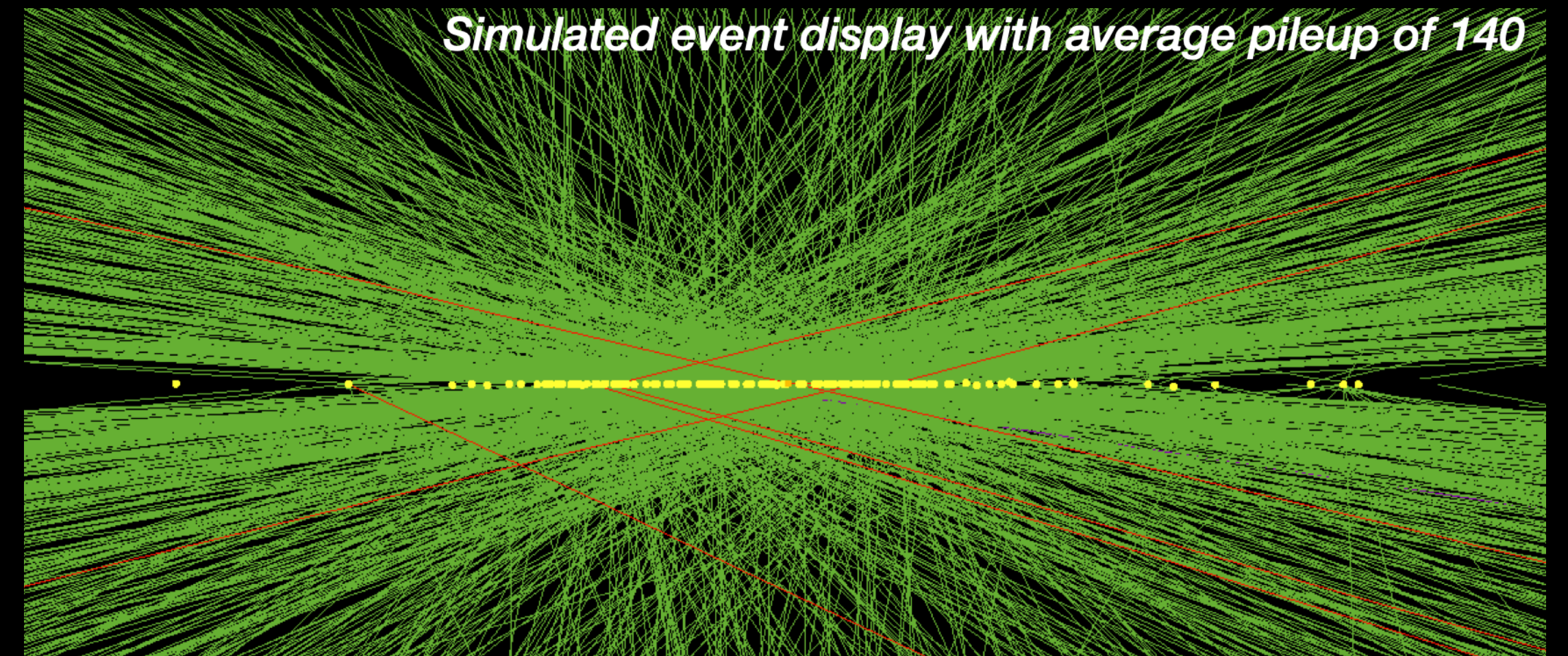
LHC (currently)

78 vertices
(average 60)

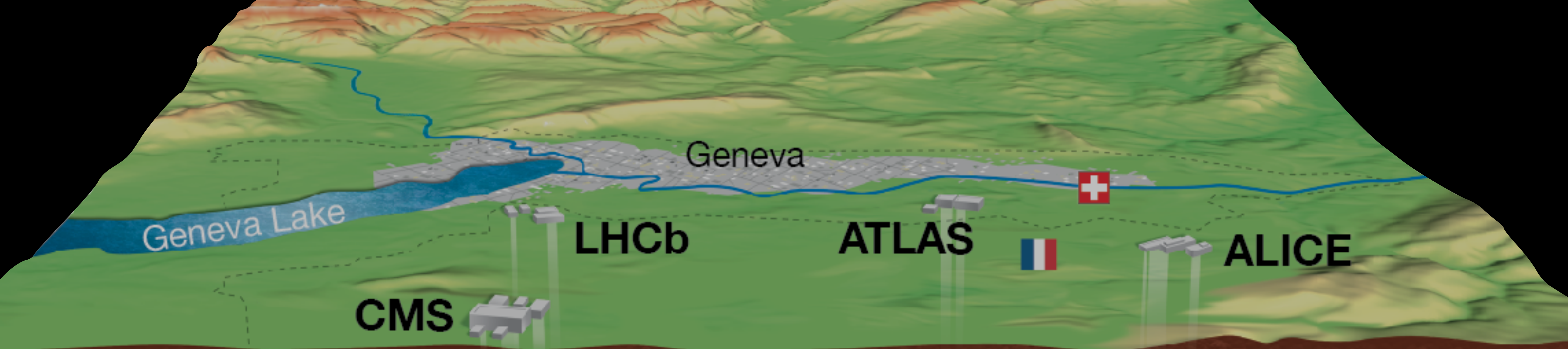


High Luminosity LHC (2030-2041)

200 vertices
(average 140)

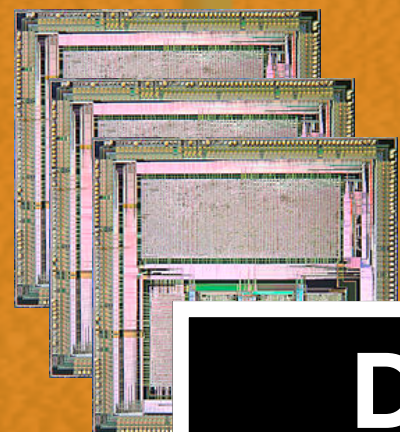


→ **event size: 2 → 8 MB,**
→ **throughput: 4 → 63 Tb/s**



ML on specialised hardware

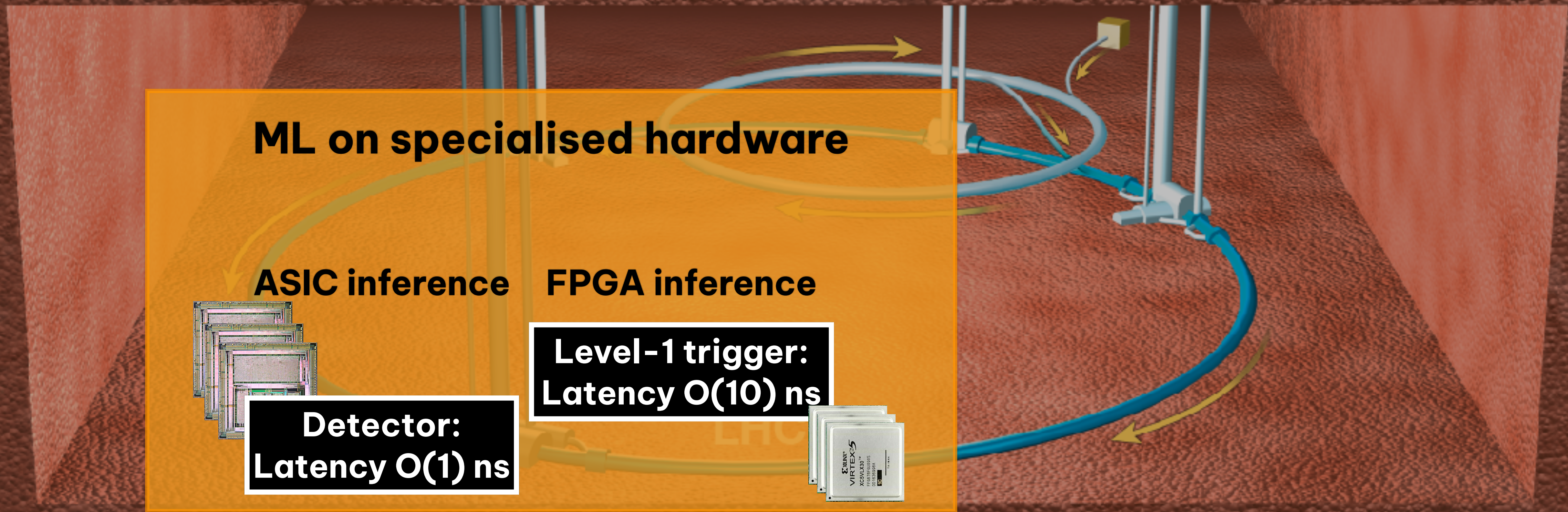
ASIC inference

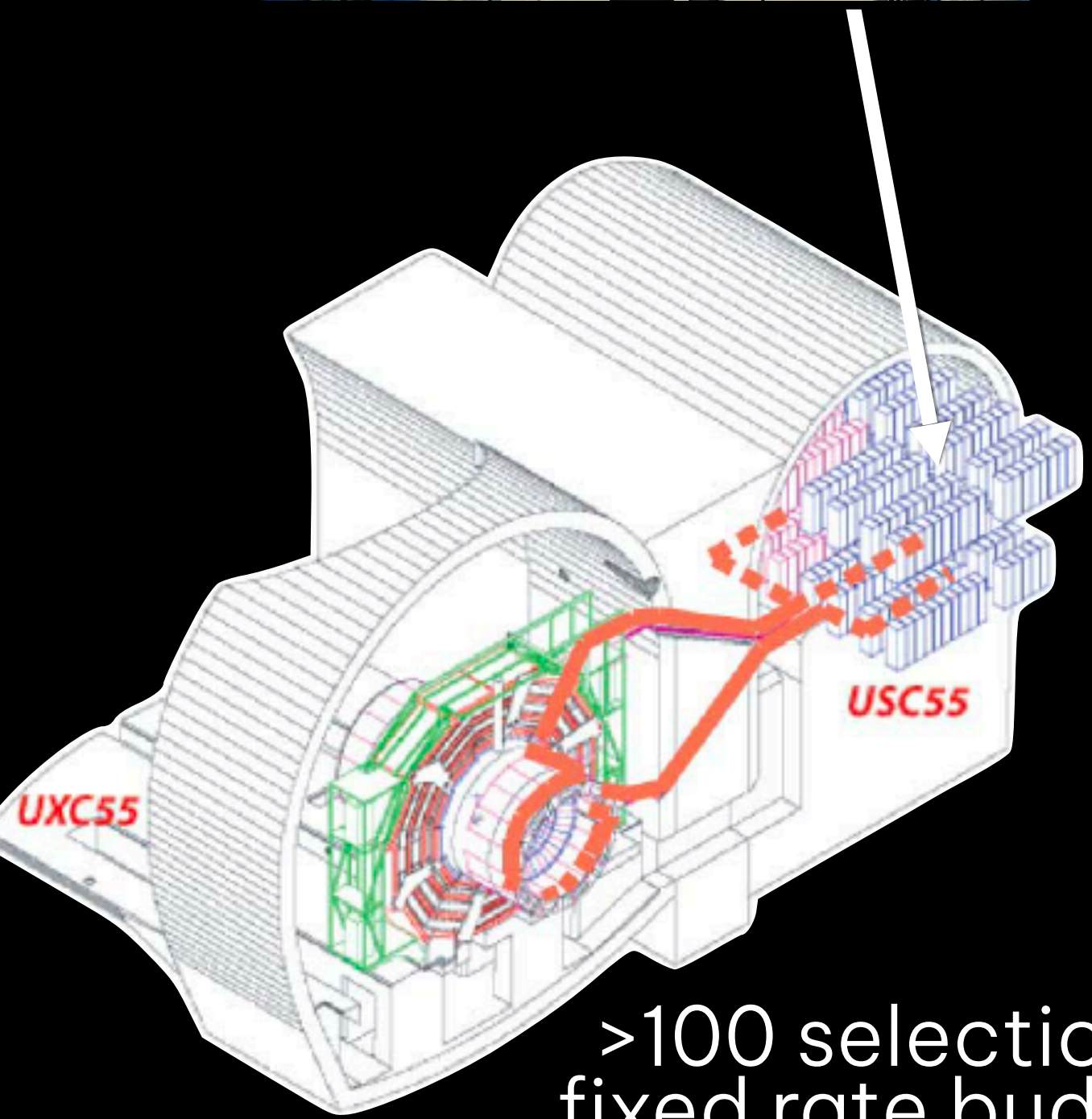
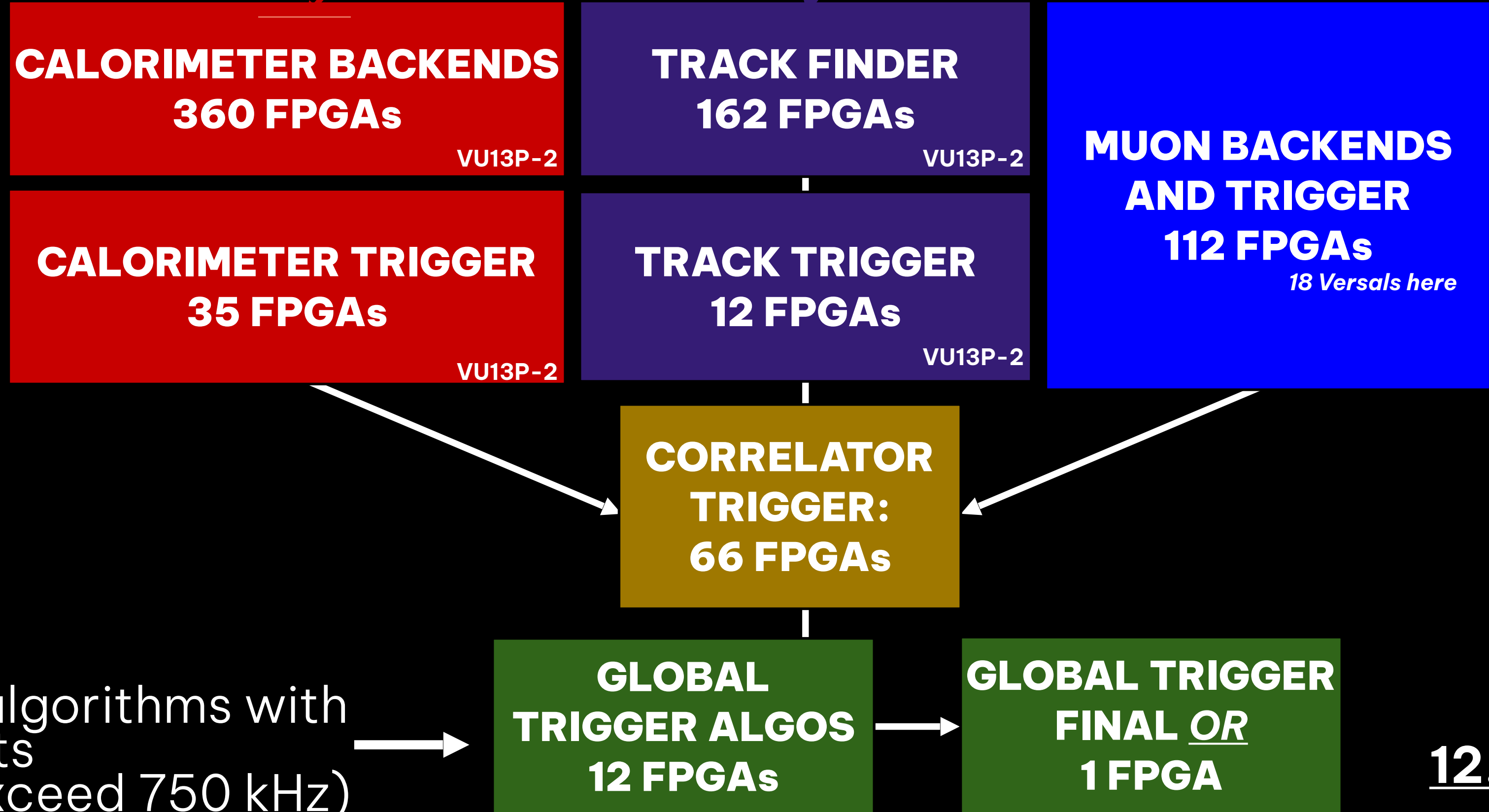
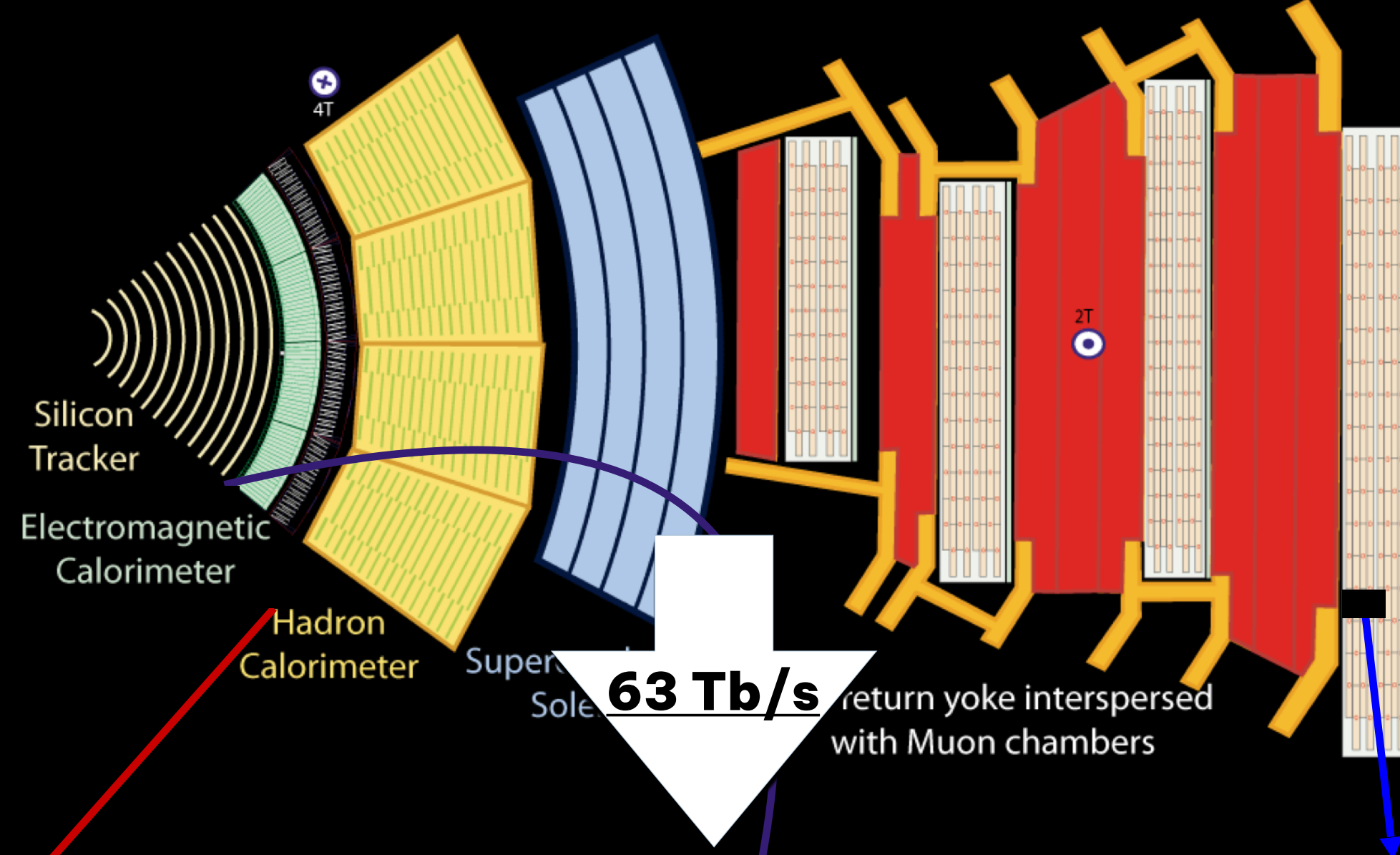


Detector:
Latency $O(1)$ ns

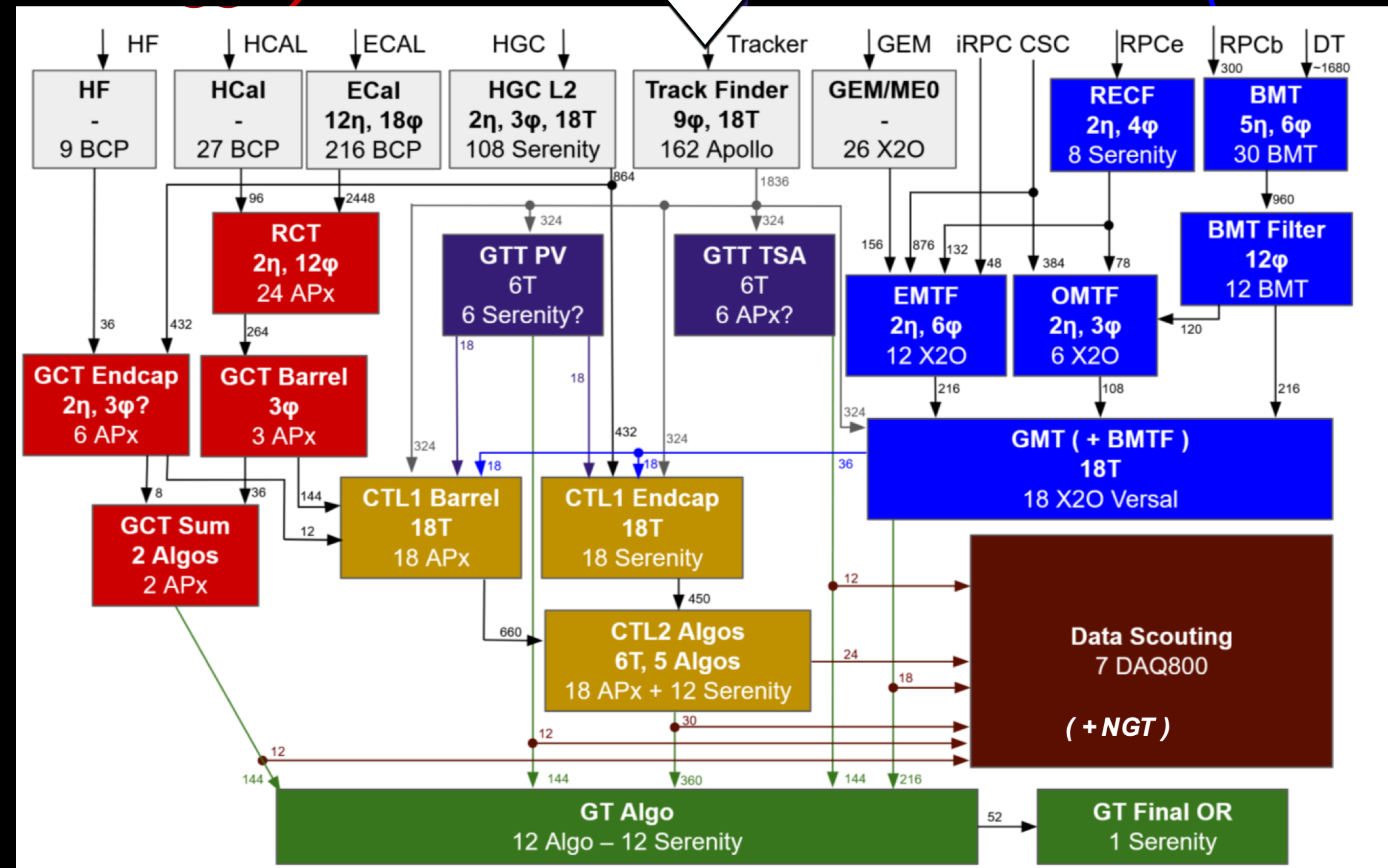
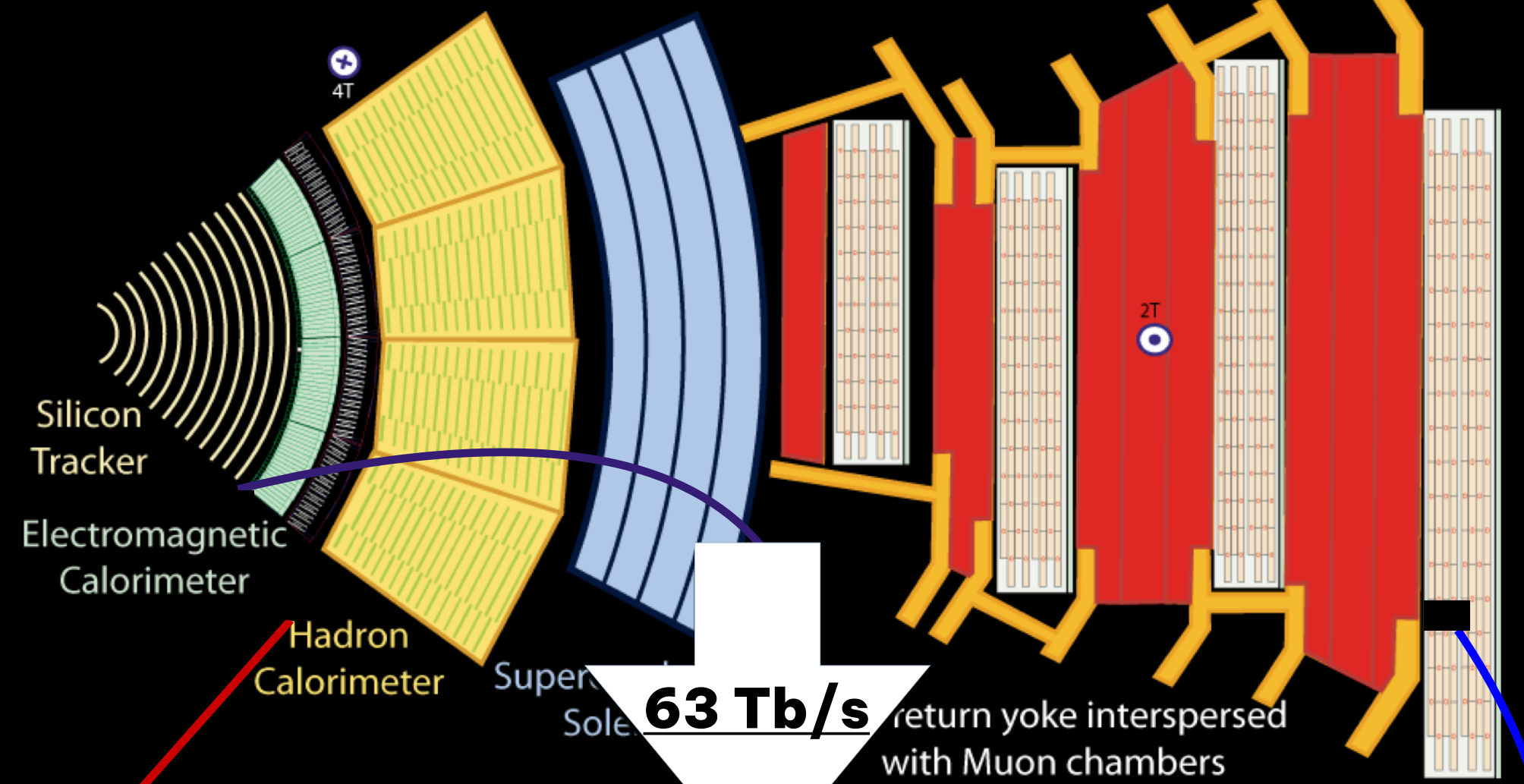
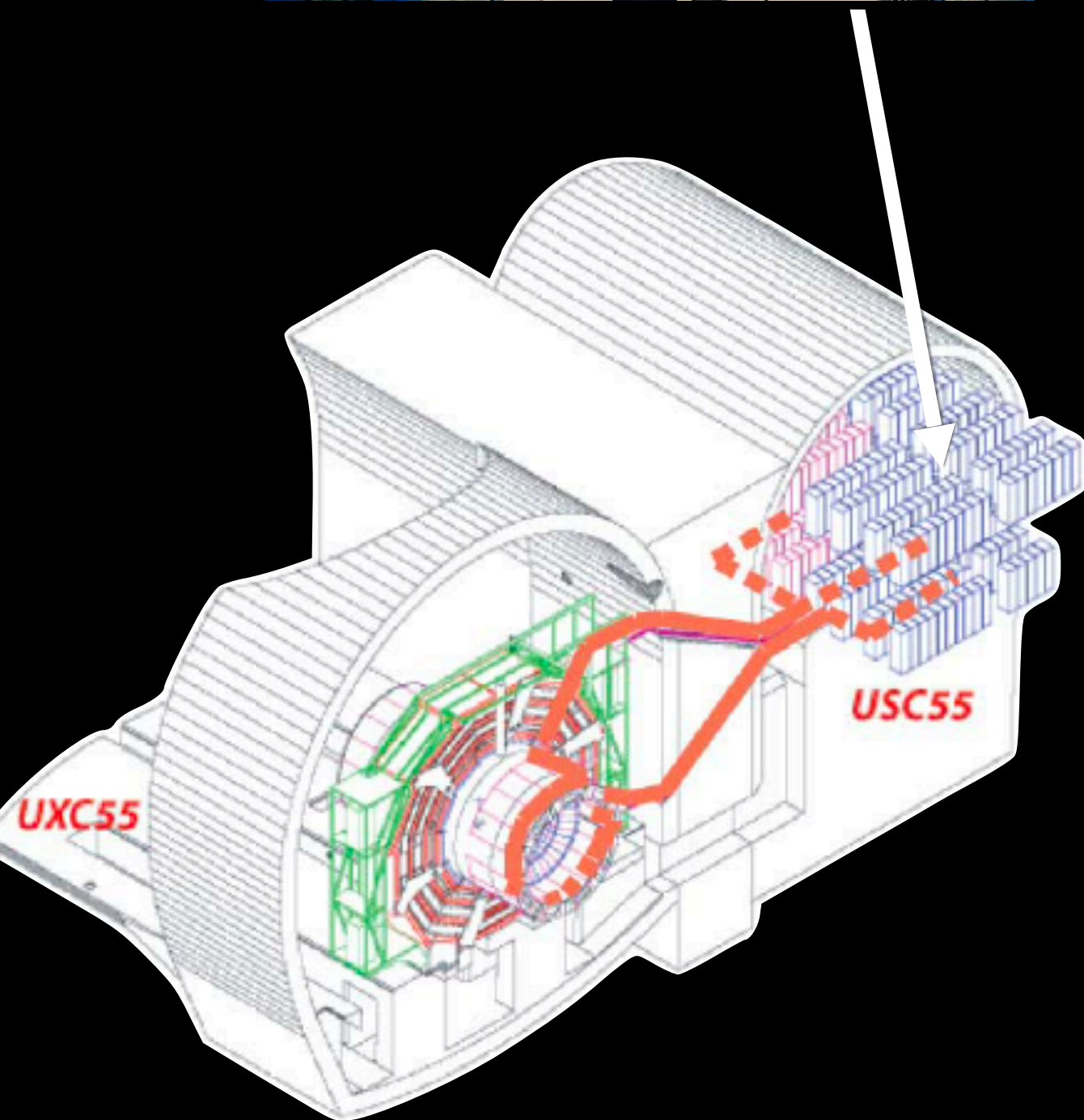
FPGA inference

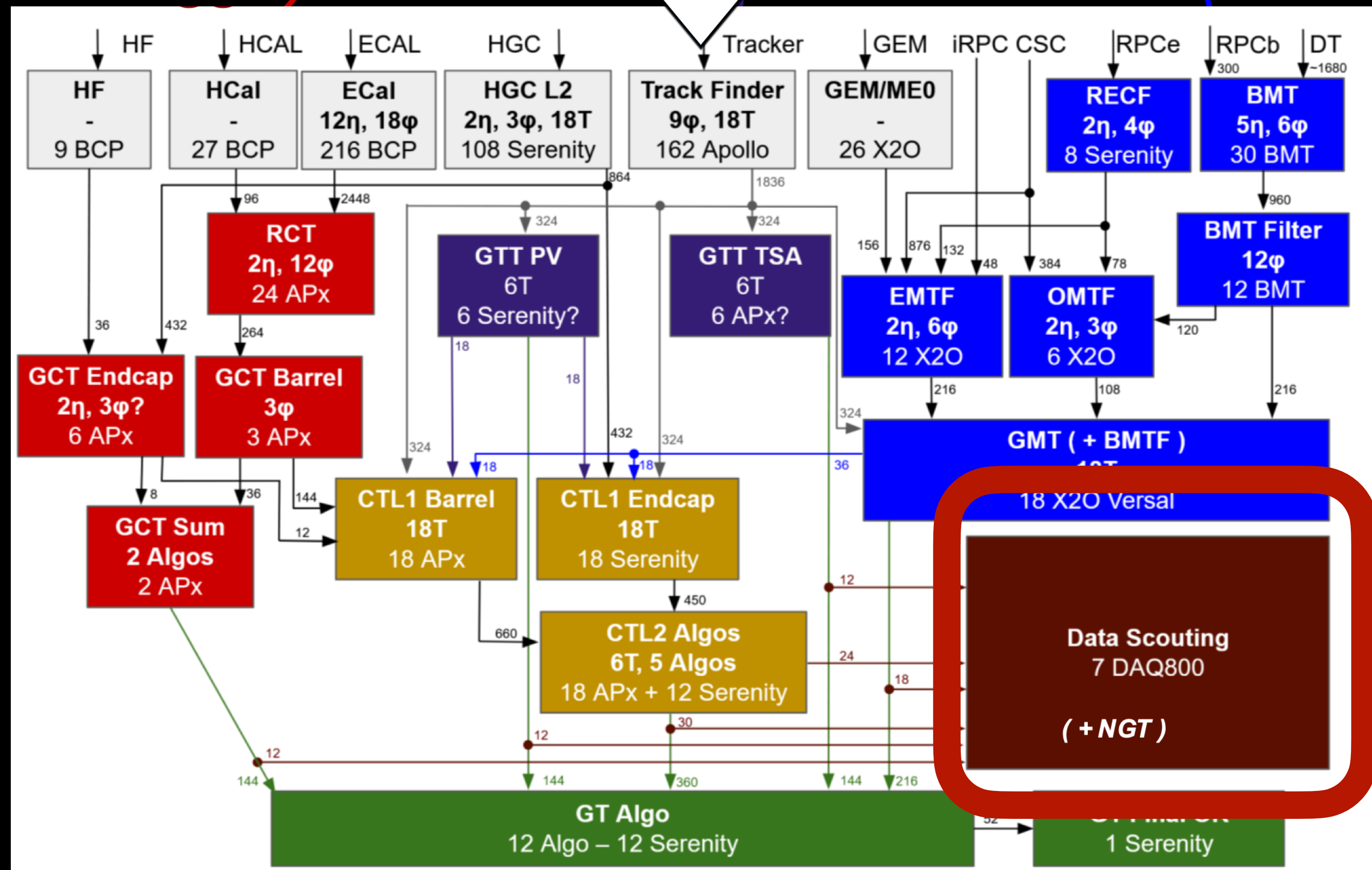
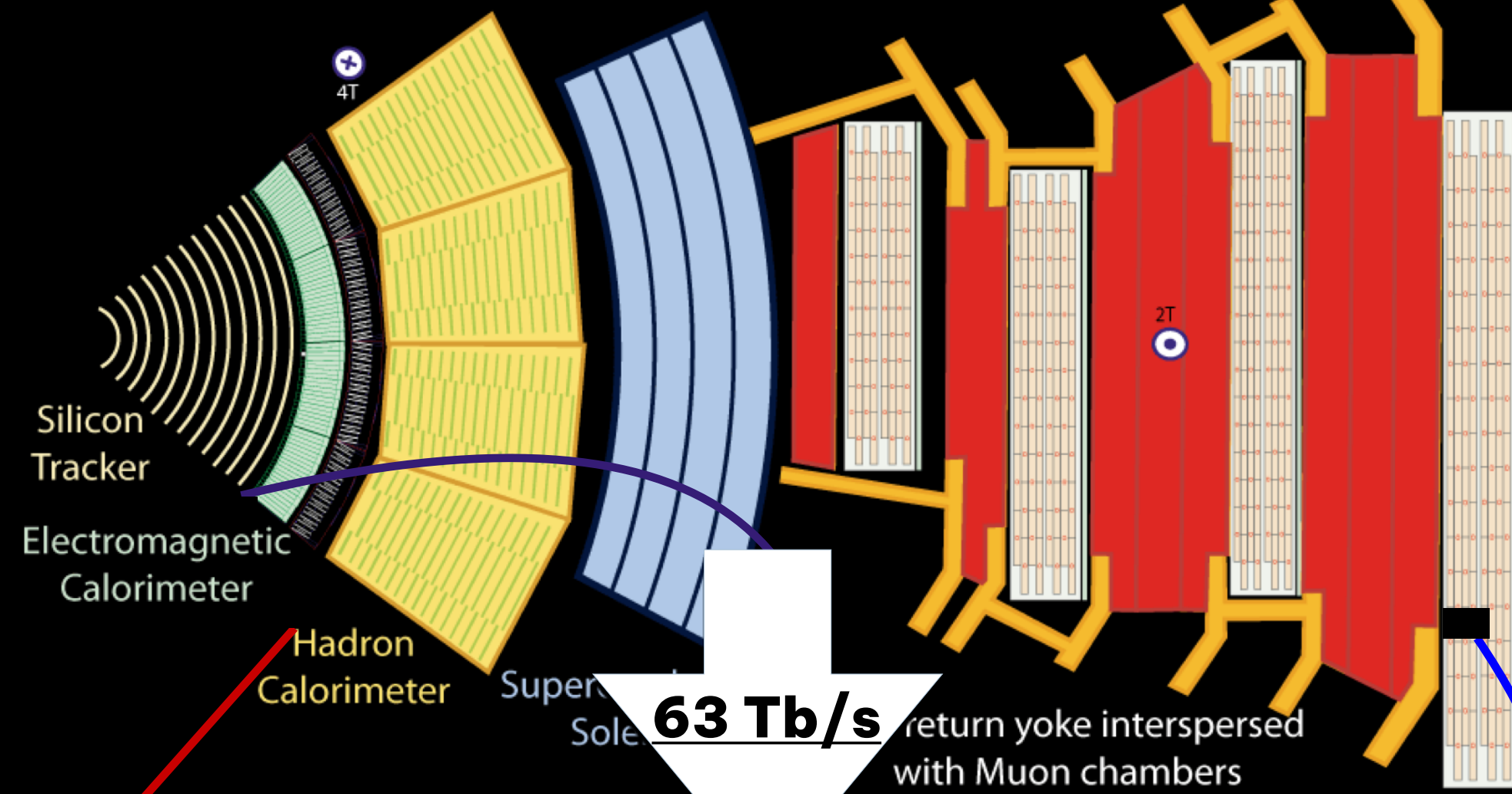
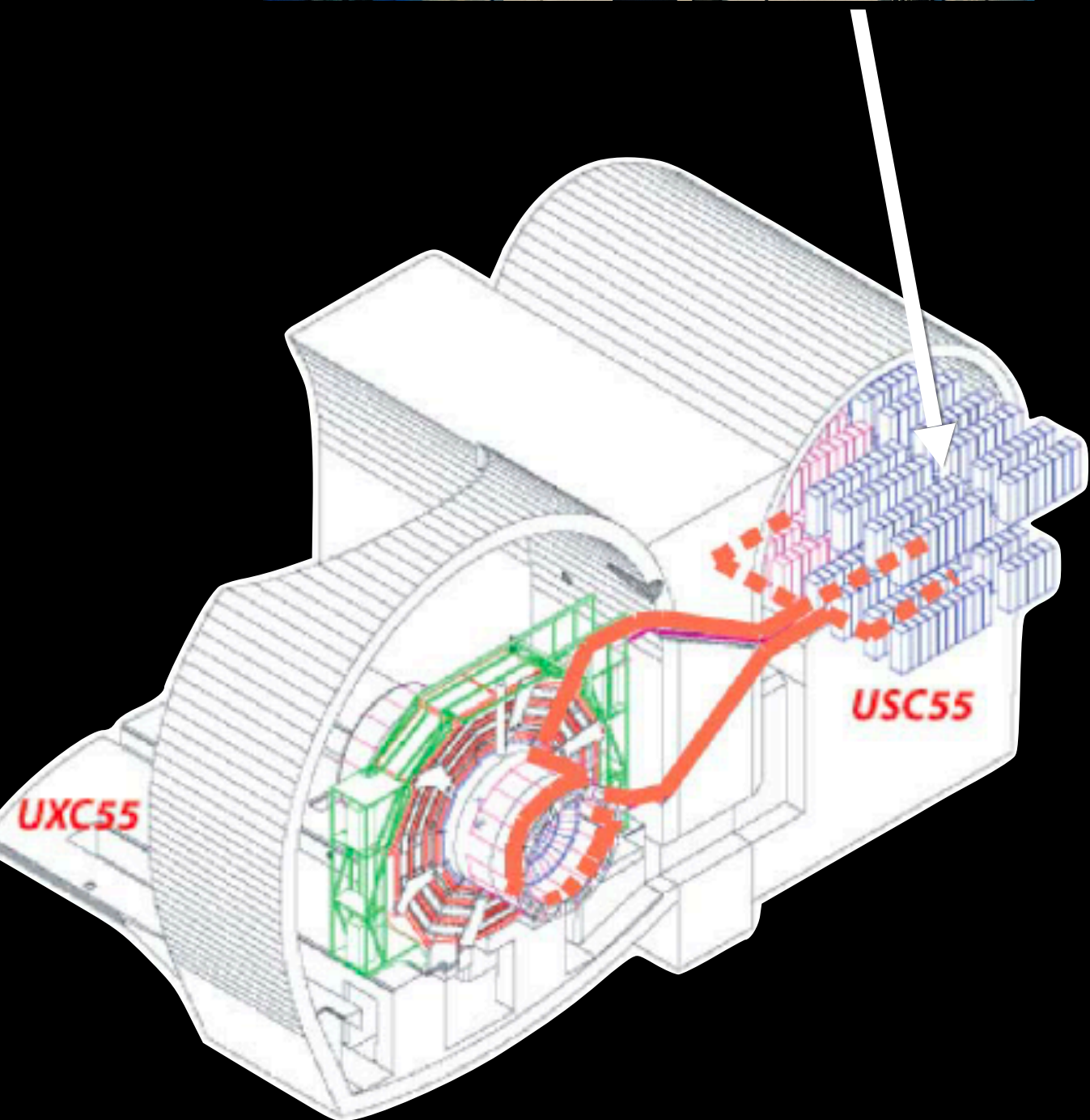
Level-1 trigger:
Latency $O(10)$ ns





>100 selection algorithms with fixed rate budgets (total cannot exceed 750 kHz)

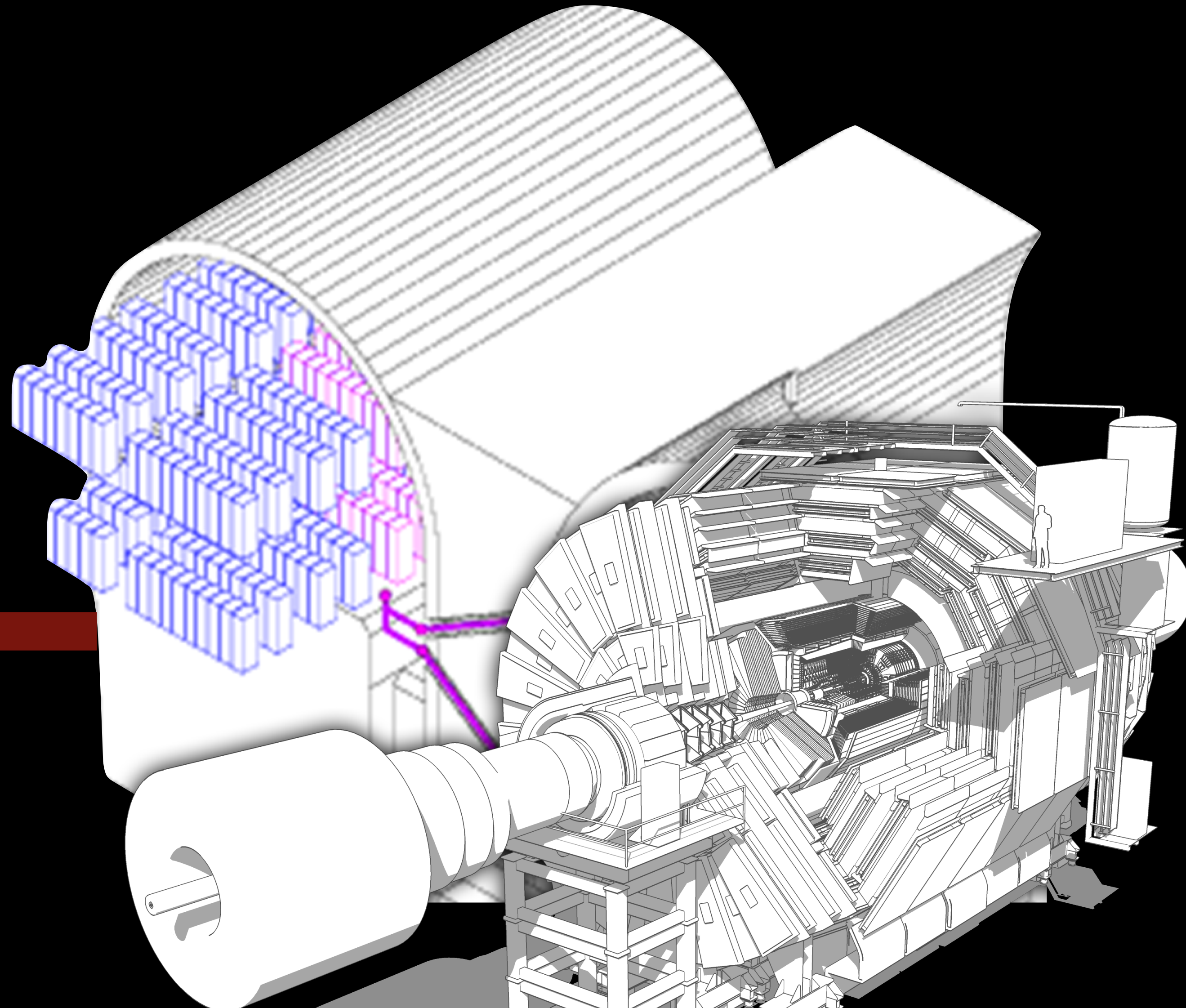


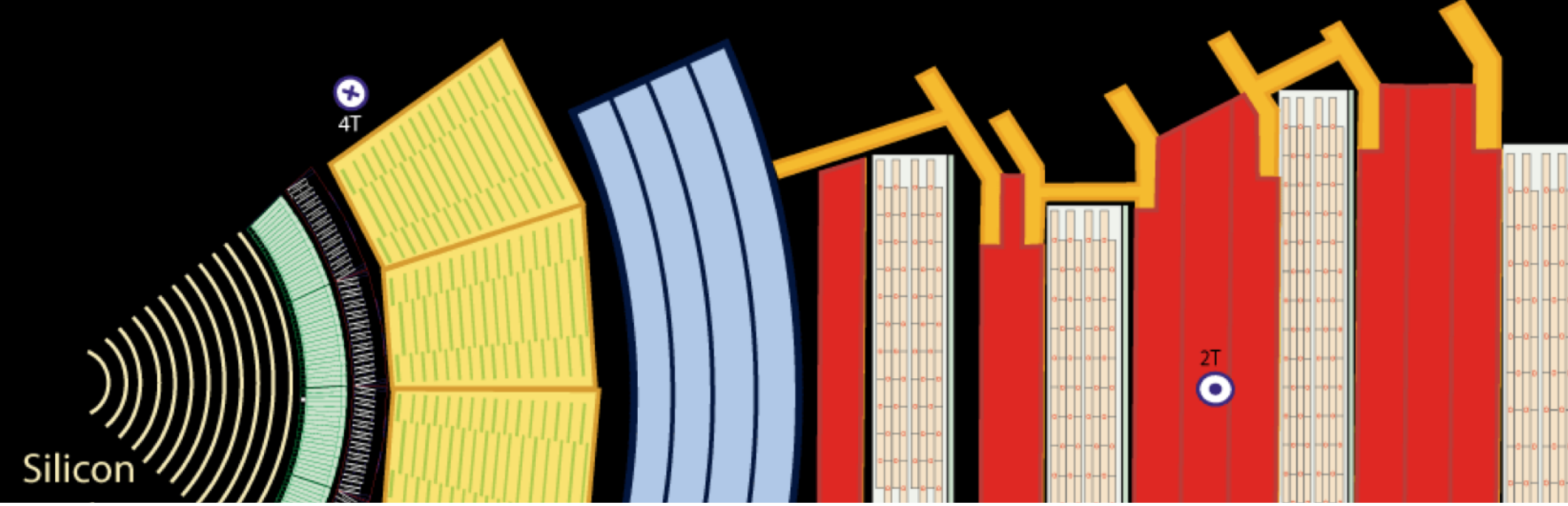


New system to read out all hardware trigger data at 40 MHz
Enables analysis of every collision event.

On-GPU real-time analysis in $O(1)$
 μs , reduce rate to 5 GB/s!

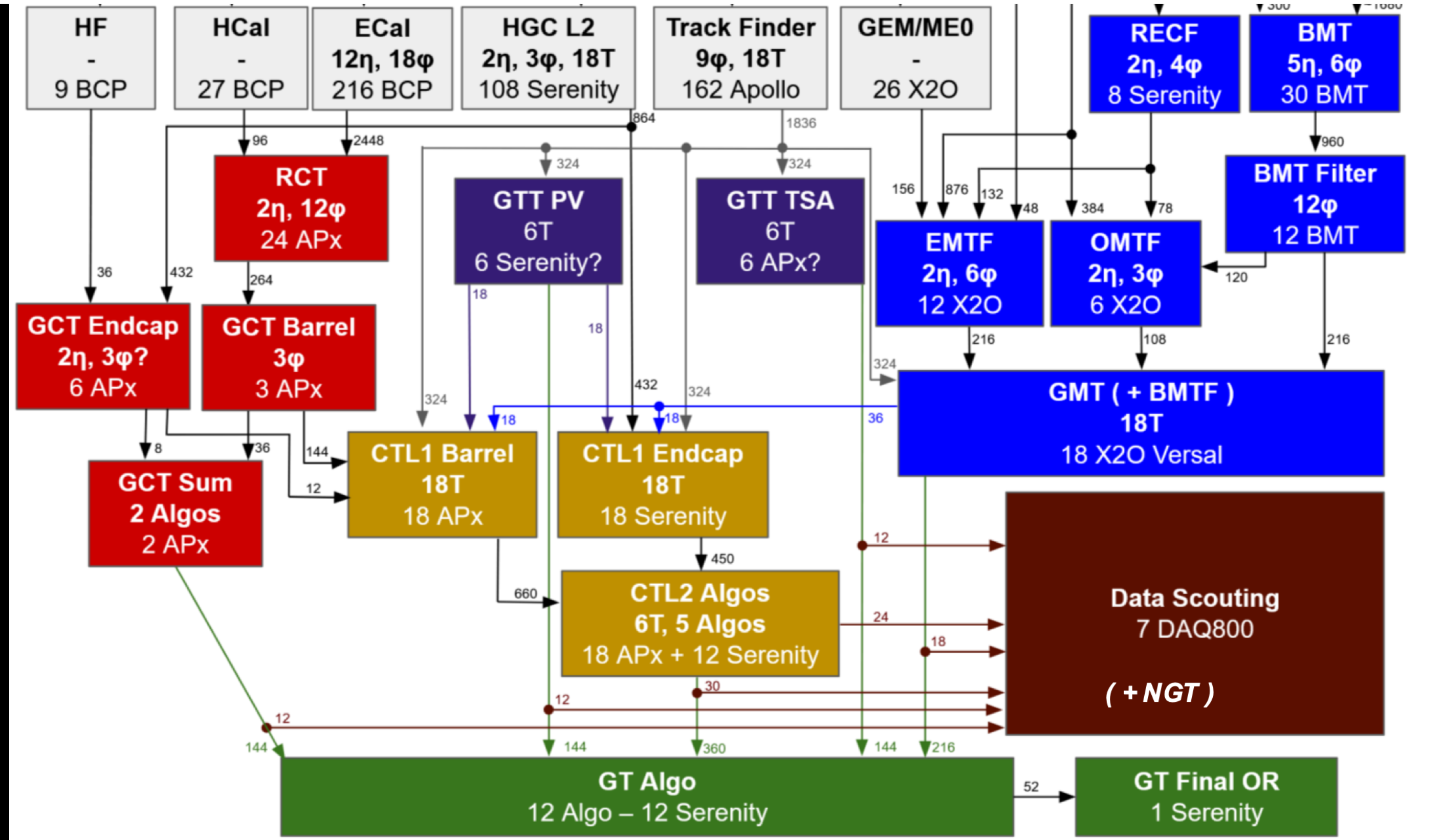
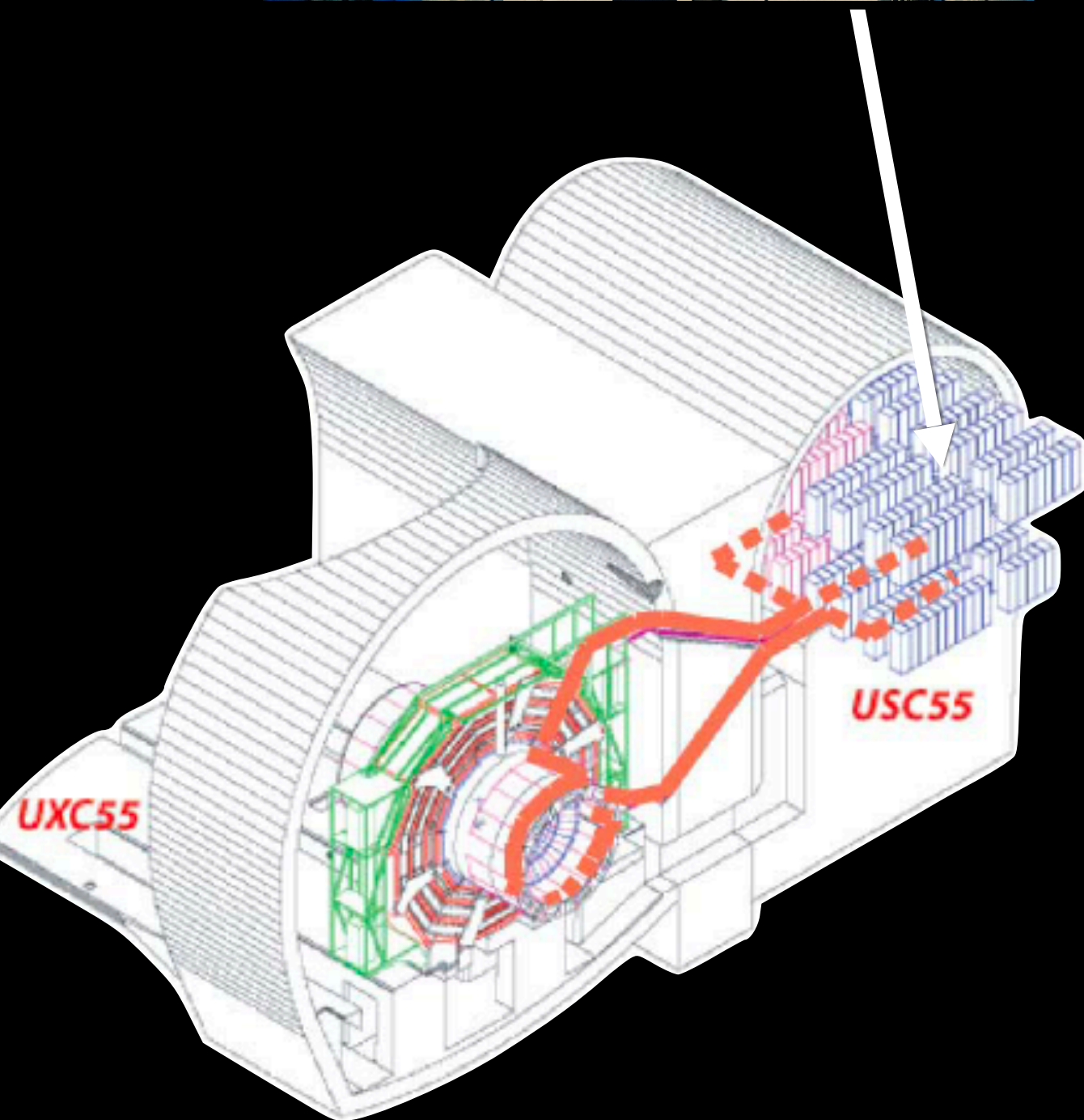
On-FPGA processing in <25 ns
7.65 Tb/second



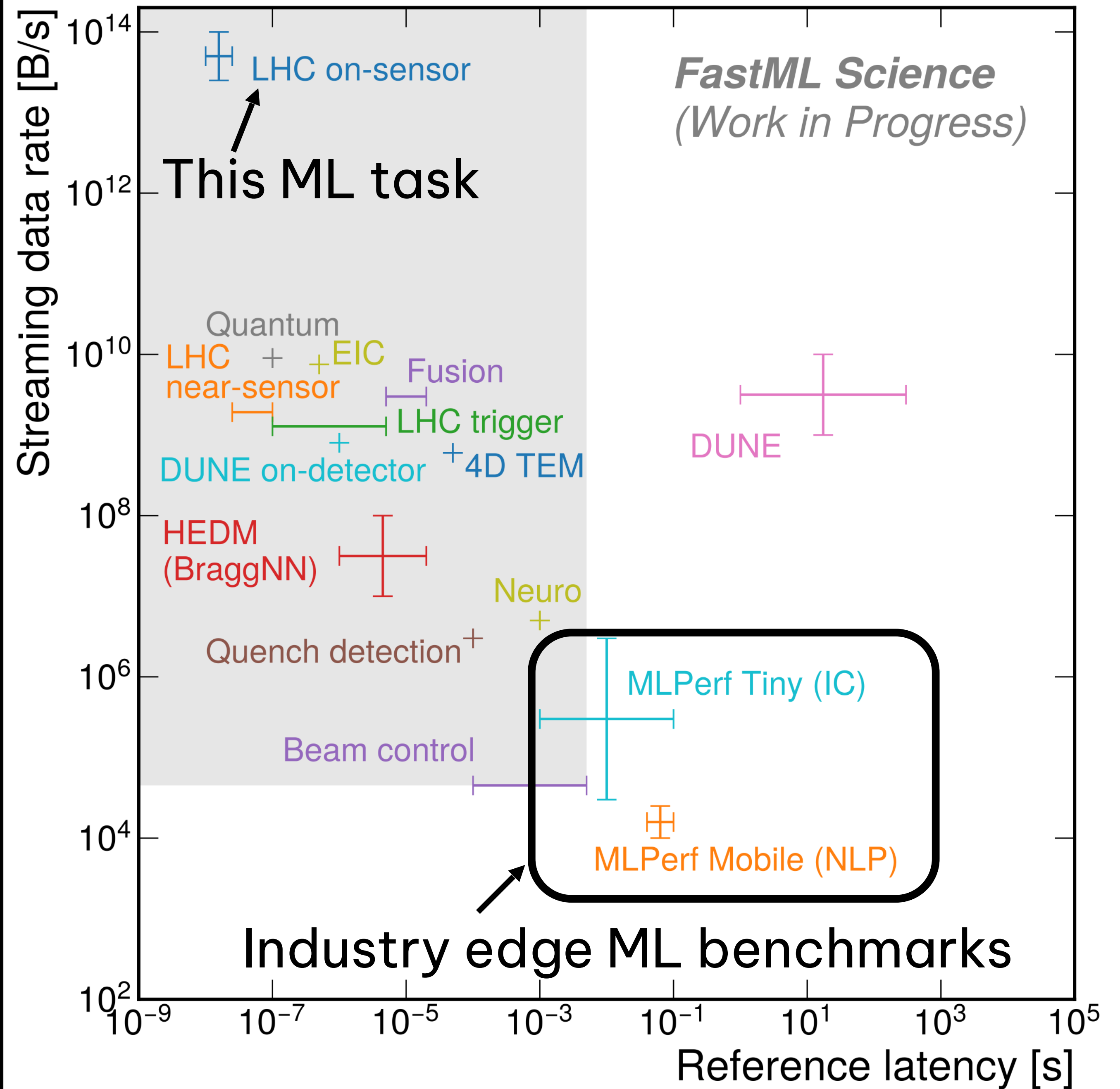


12 microseconds latency

Processing 5% of internet traffic

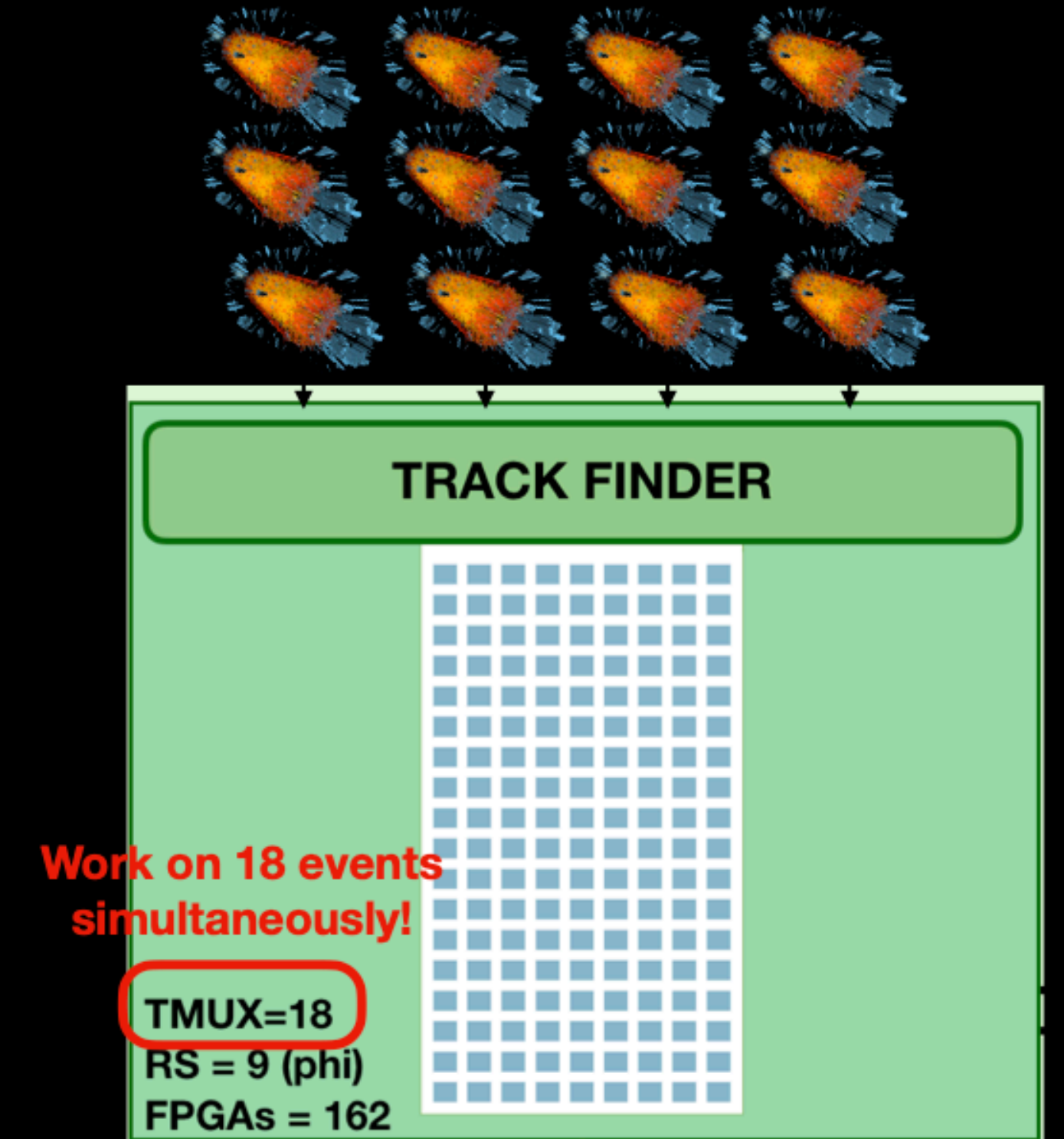


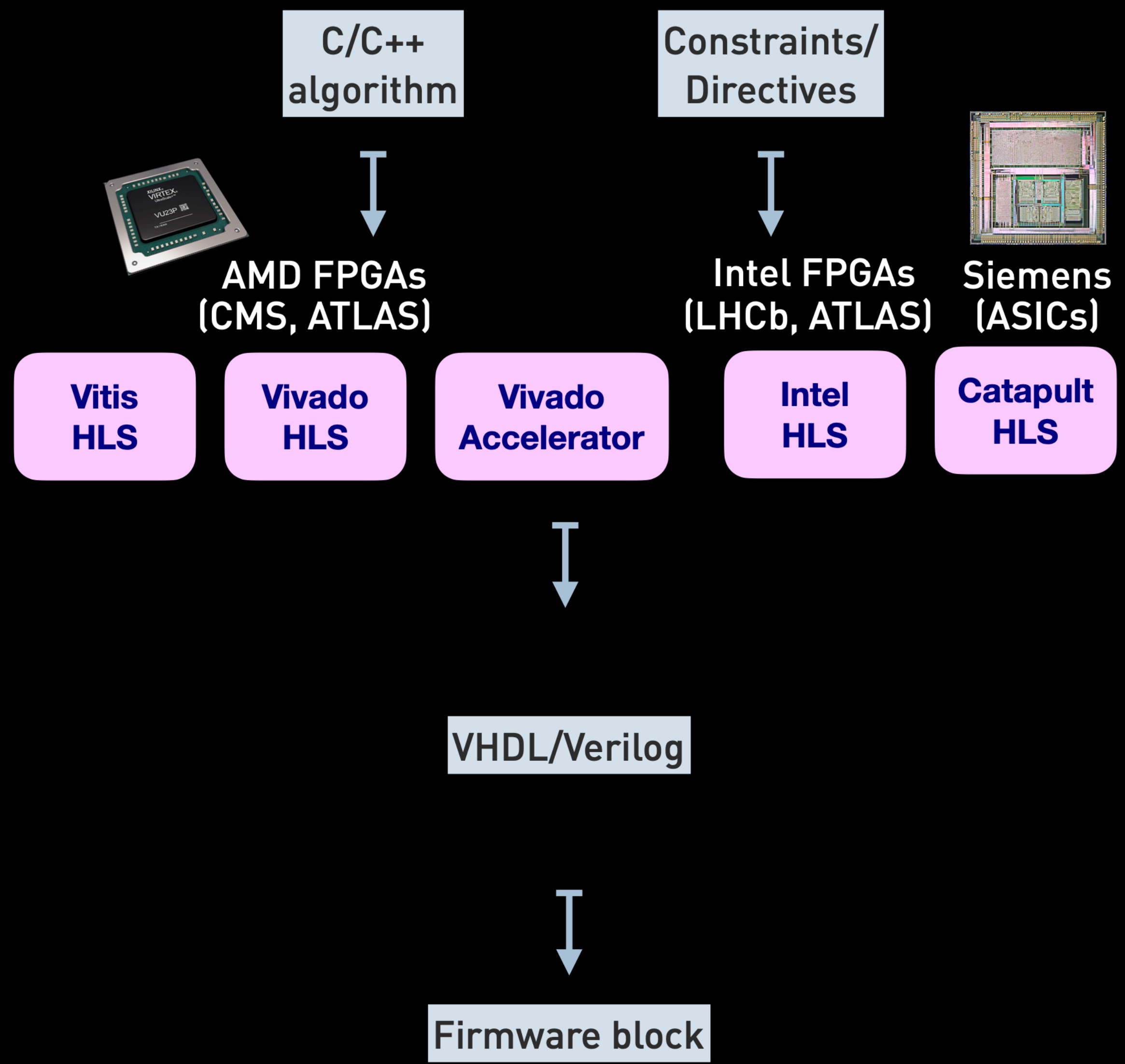
***Had to design our own tools
(to do inference in <100 ns)***

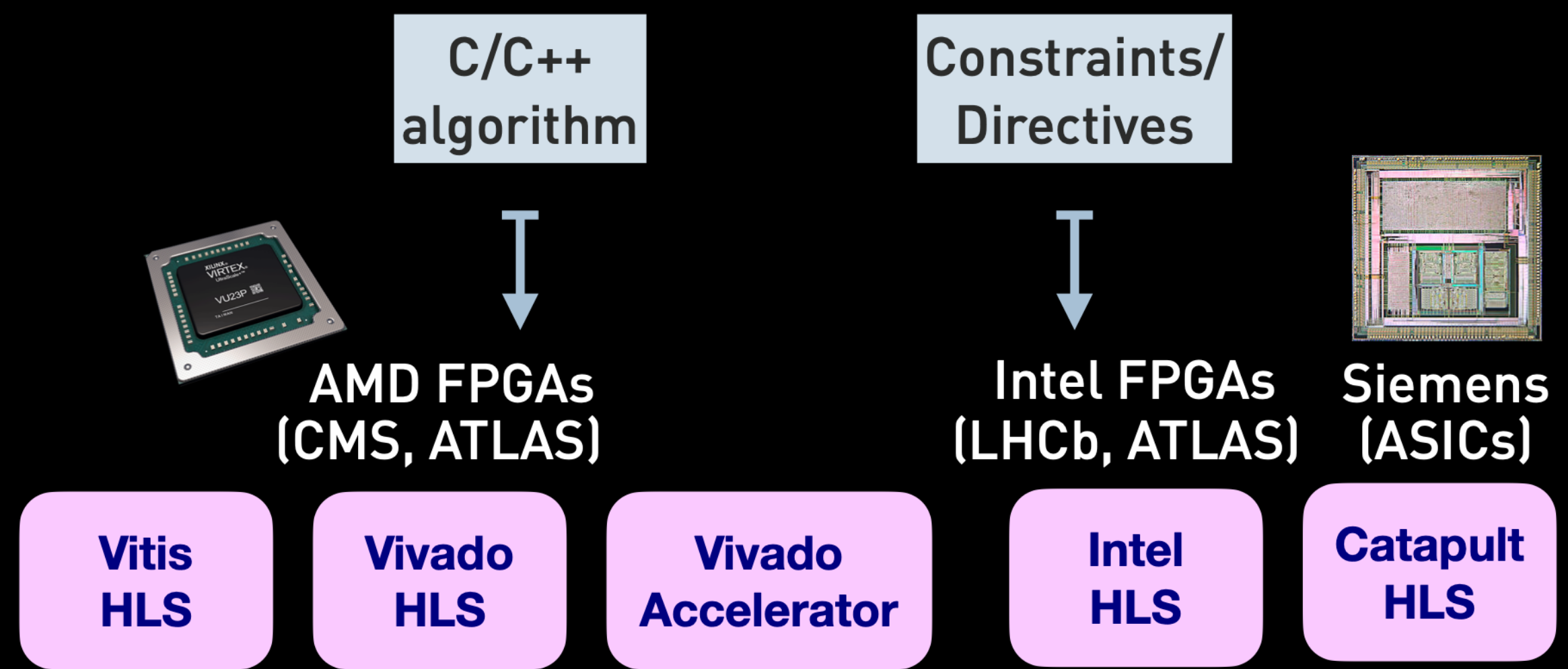


Why FPGAs

- Resource parallelism ↔ latency
- Pipeline parallelism ↔ throughput
- Latency deterministic (40 MHz collisions)
- Specialized blocks for I/O
- High-speed transceivers with Tb/s total bandwidth (PCIe, 100Gbps ethernet, InfiniBand)
- Relatively low power per op







Hardware Description Languages:
programming languages which describe
electronic circuits

VHDL/Verilog

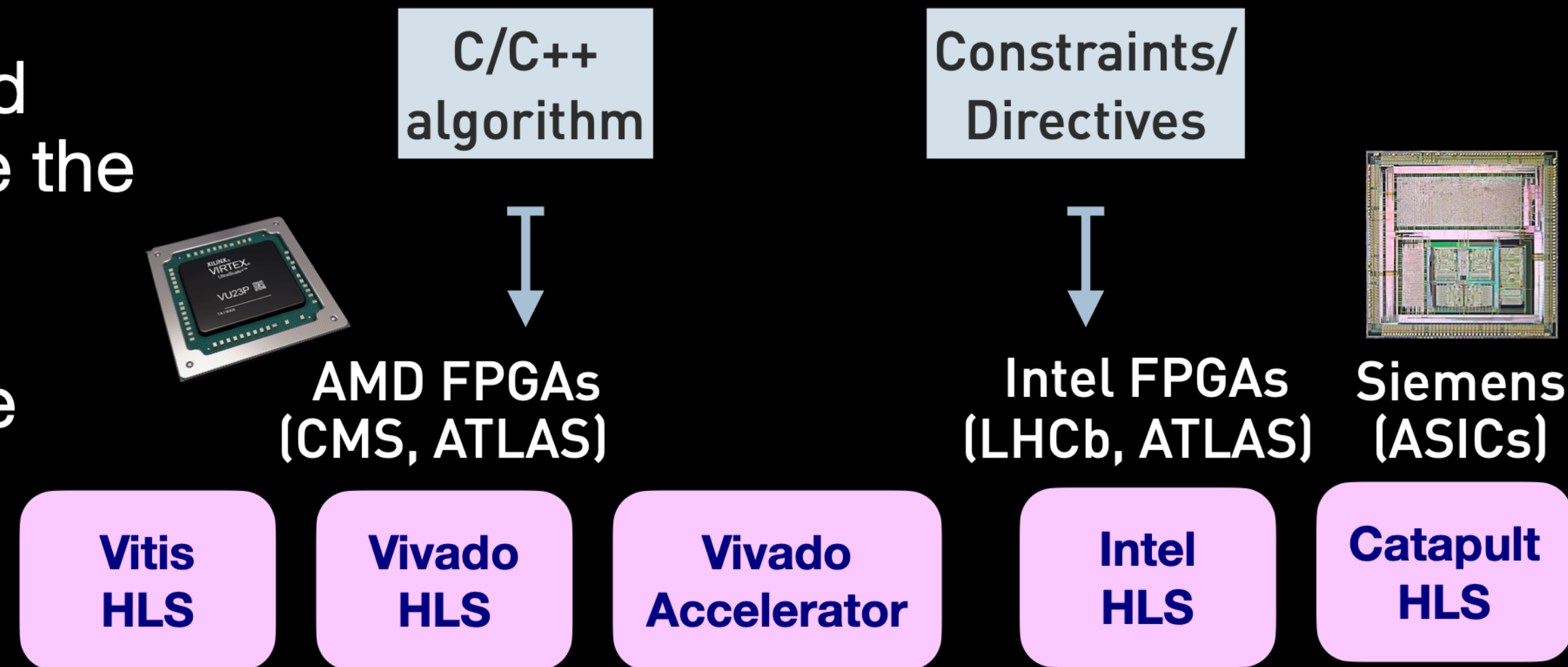
Firmware block

High Level Synthesis:

Compile from C/C++ to VHDL

Pre-processor directives and constraints used to optimize the design

Drastic decrease in firmware development time!



Hardware Description Languages:

programming languages which describe electronic circuits

VHDL/Verilog

Firmware block

hls 4 ml

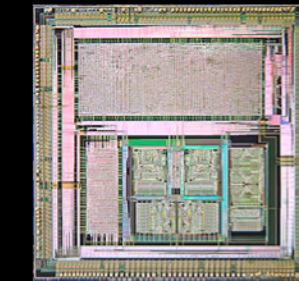
C/C++
algorithm

Constraints/
Directives

Let's tune this for
nanosecond ML!



AMD FPGAs
(CMS, ATLAS)



Intel FPGAs
(LHCb, ATLAS)

Siemens
(ASICs)

Vitis
HLS

Vivado
HLS

Vivado
Accelerator

Intel
HLS

Catapult
HLS

VHDL/Verilog

Firmware block

Floating point:

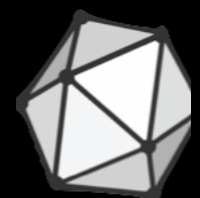


Model

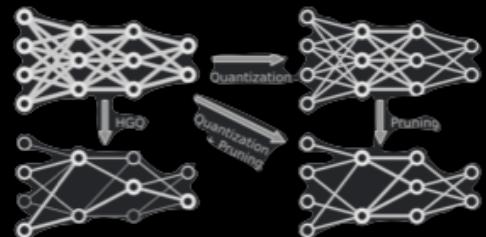
Quantized:



Brevitas



QONNX

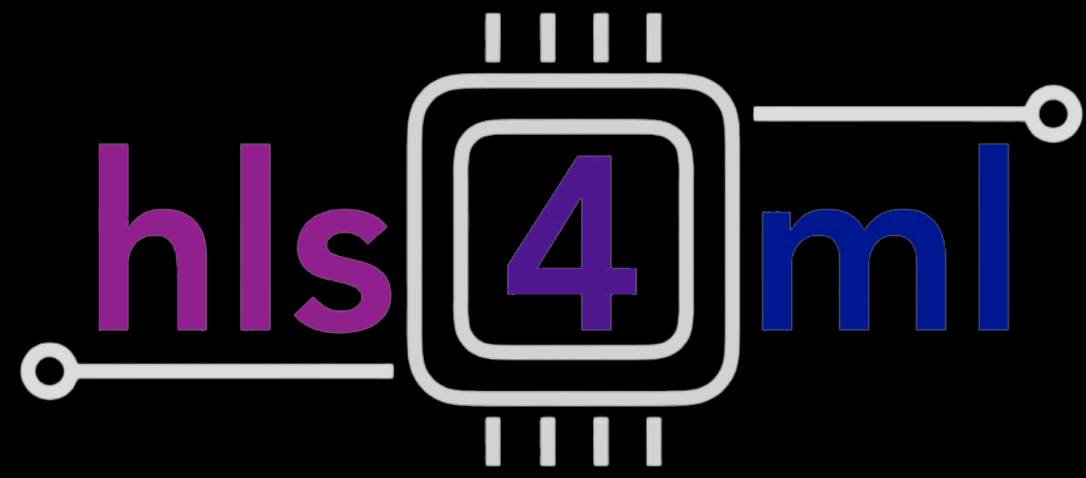
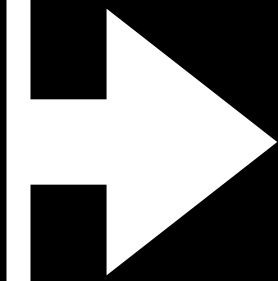


HQG

Floating point:



Model



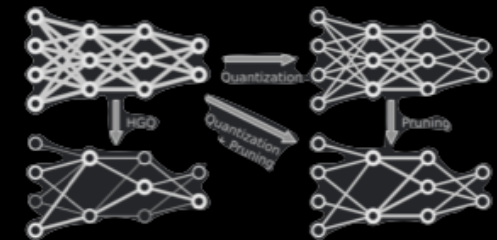
HLS Model

Write HLS project
Emulate
Run synthesis

Quantized:



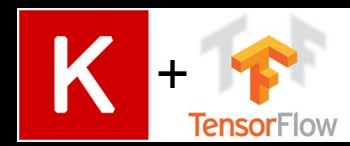
Brevitas



HQQ



Floating point:

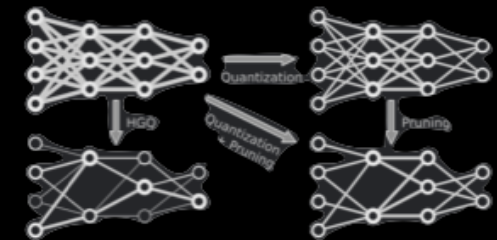


Model

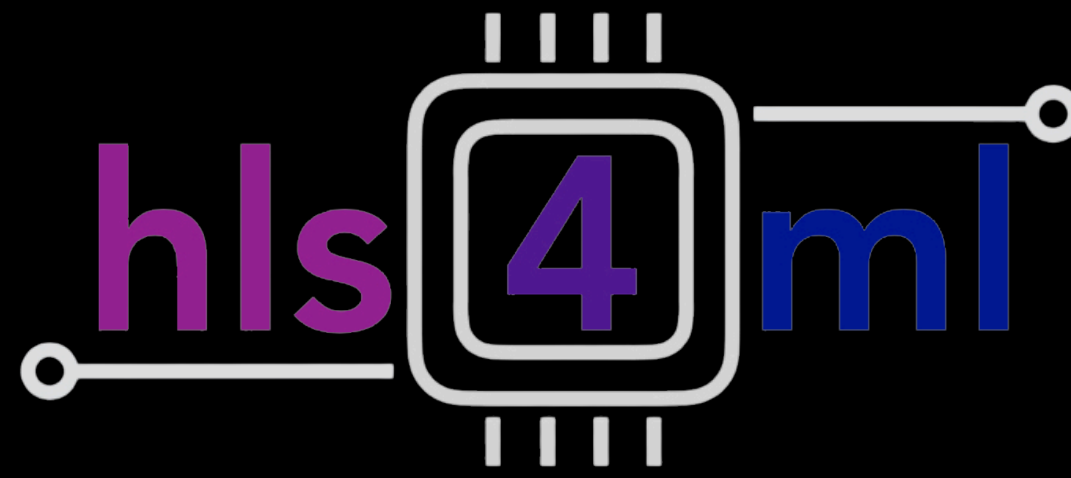
Quantized:



Brevitas

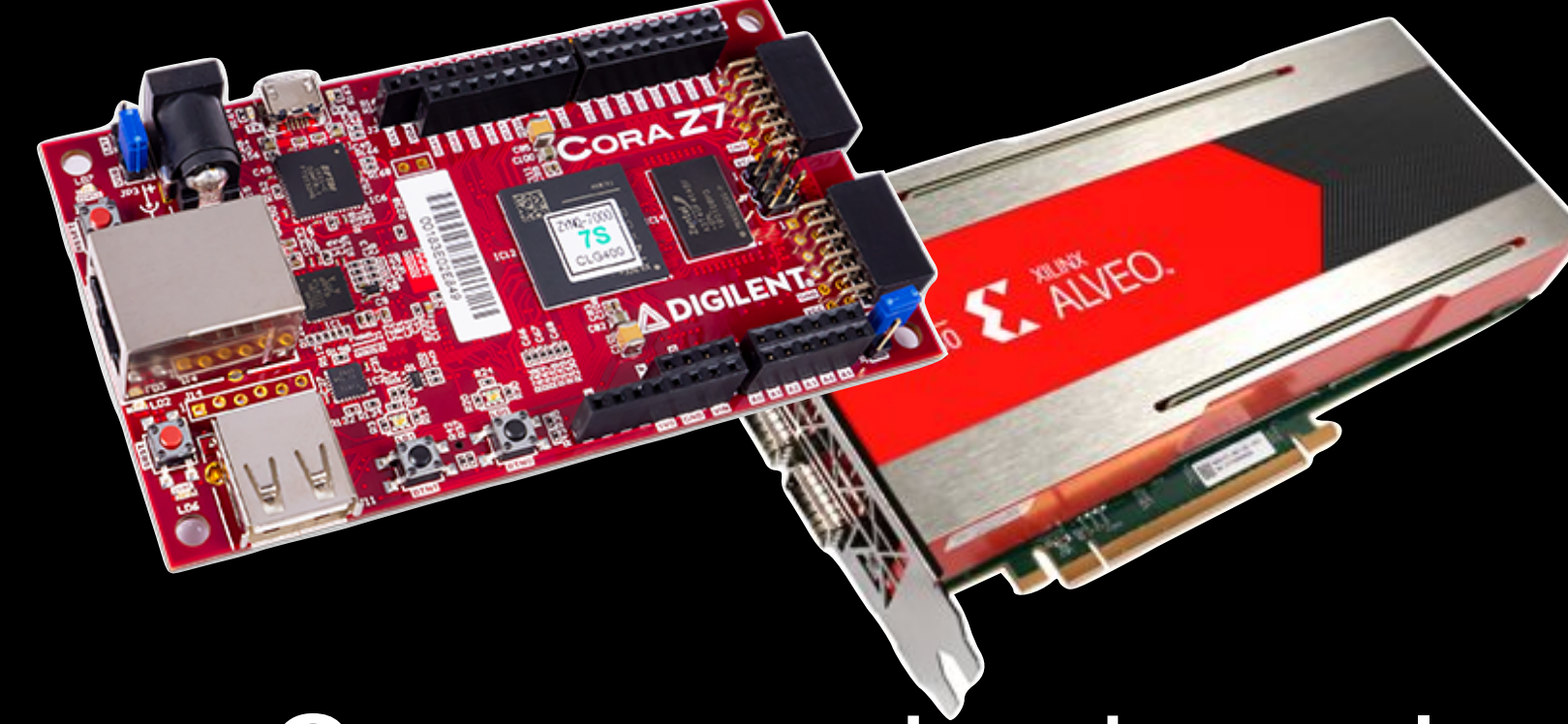


HQG



HLS Model

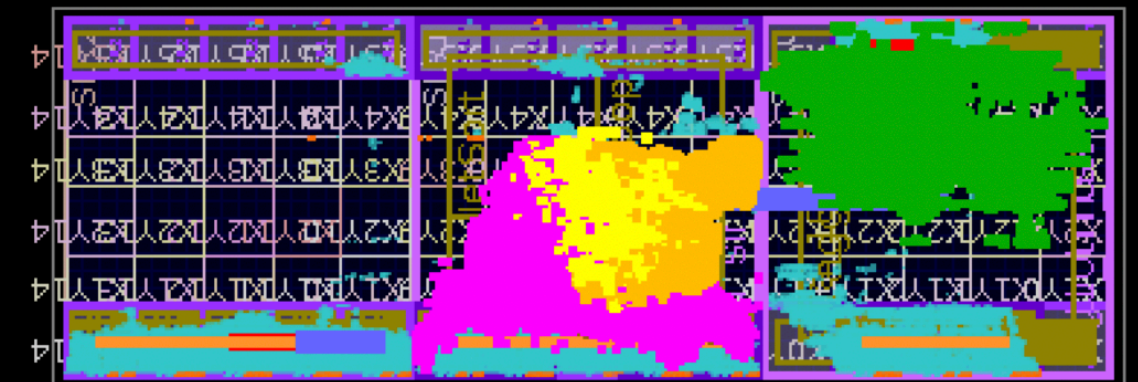
Write HLS project
Emulate
Run synthesis



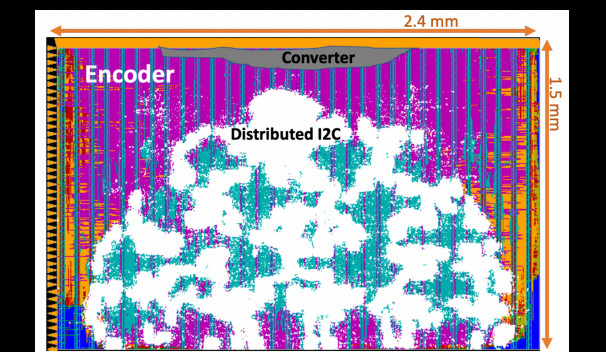
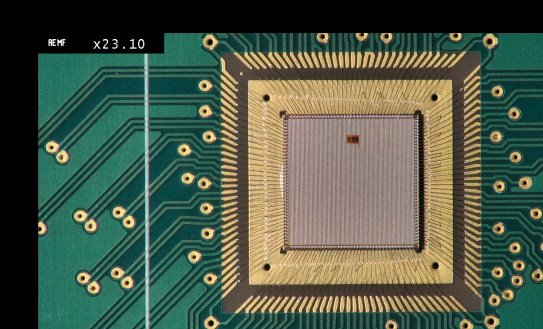
Co-processing kernel
(Accelerators/SoCs)

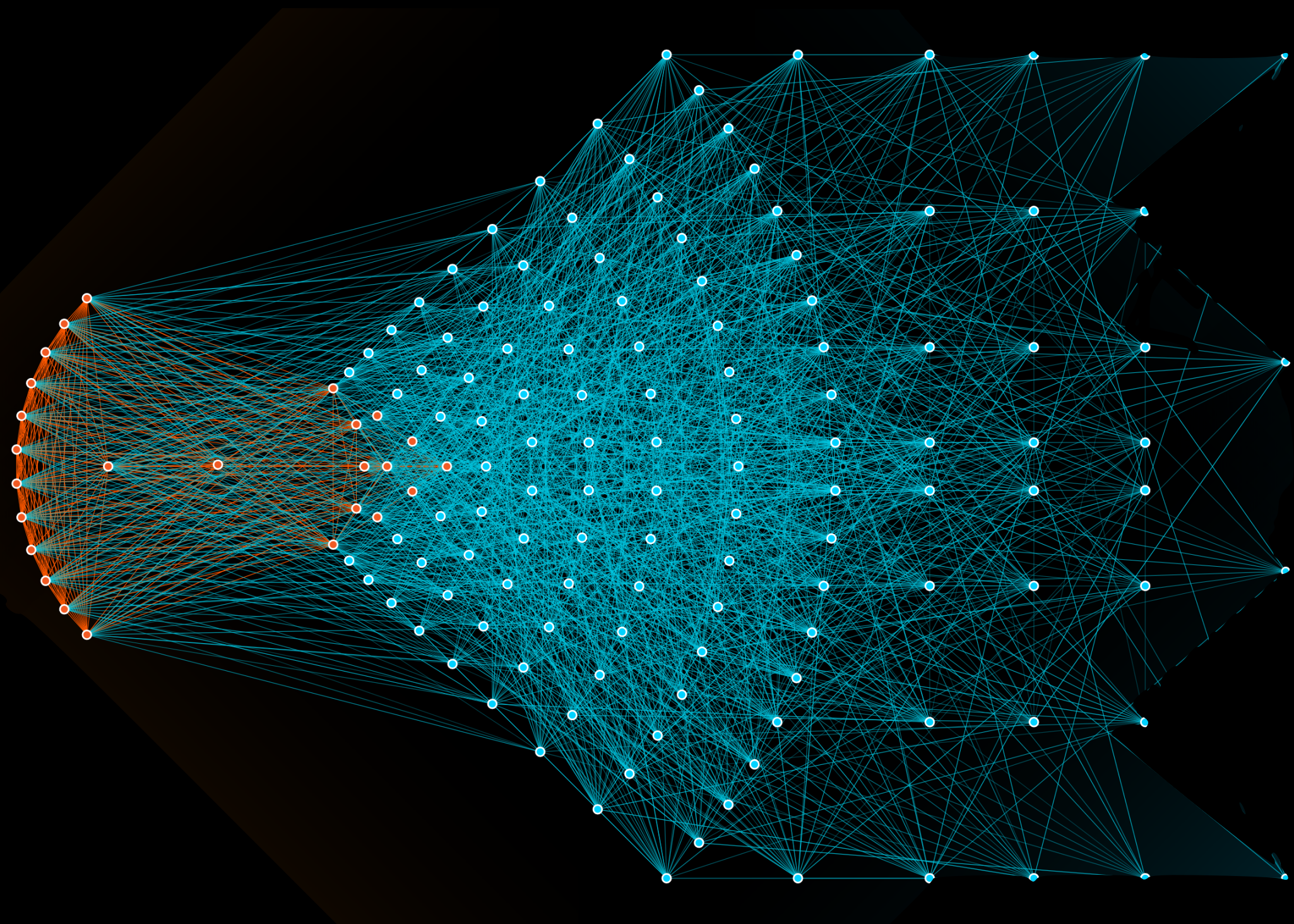
Hardware Model

FPGA custom designs

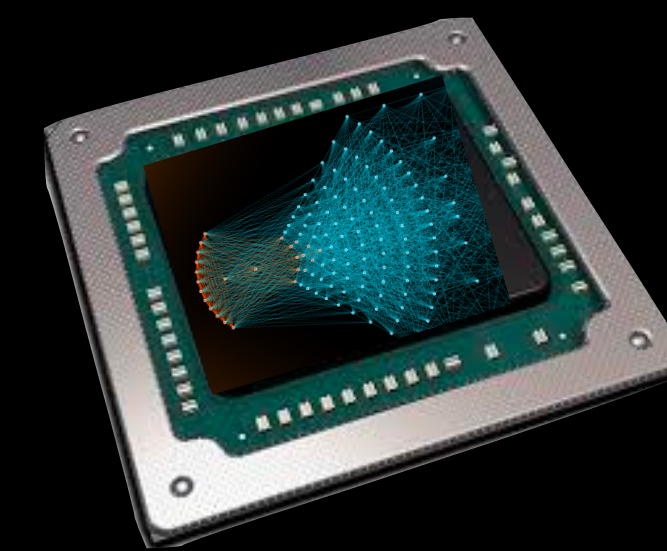


ASICs

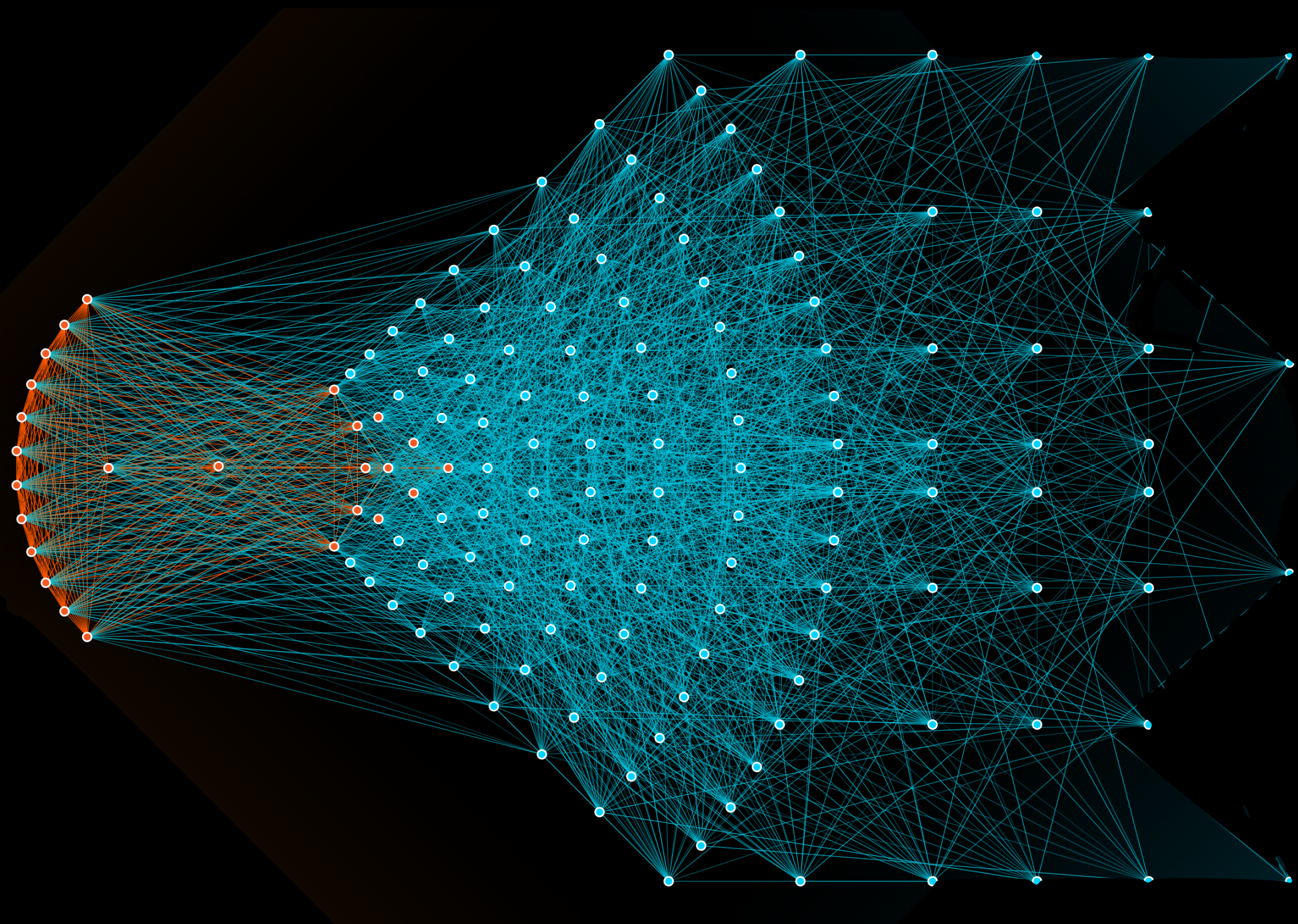




Ideally



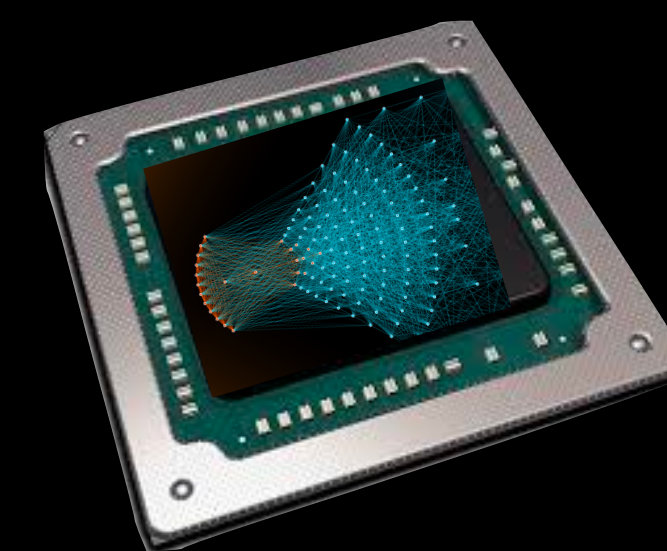
Reality



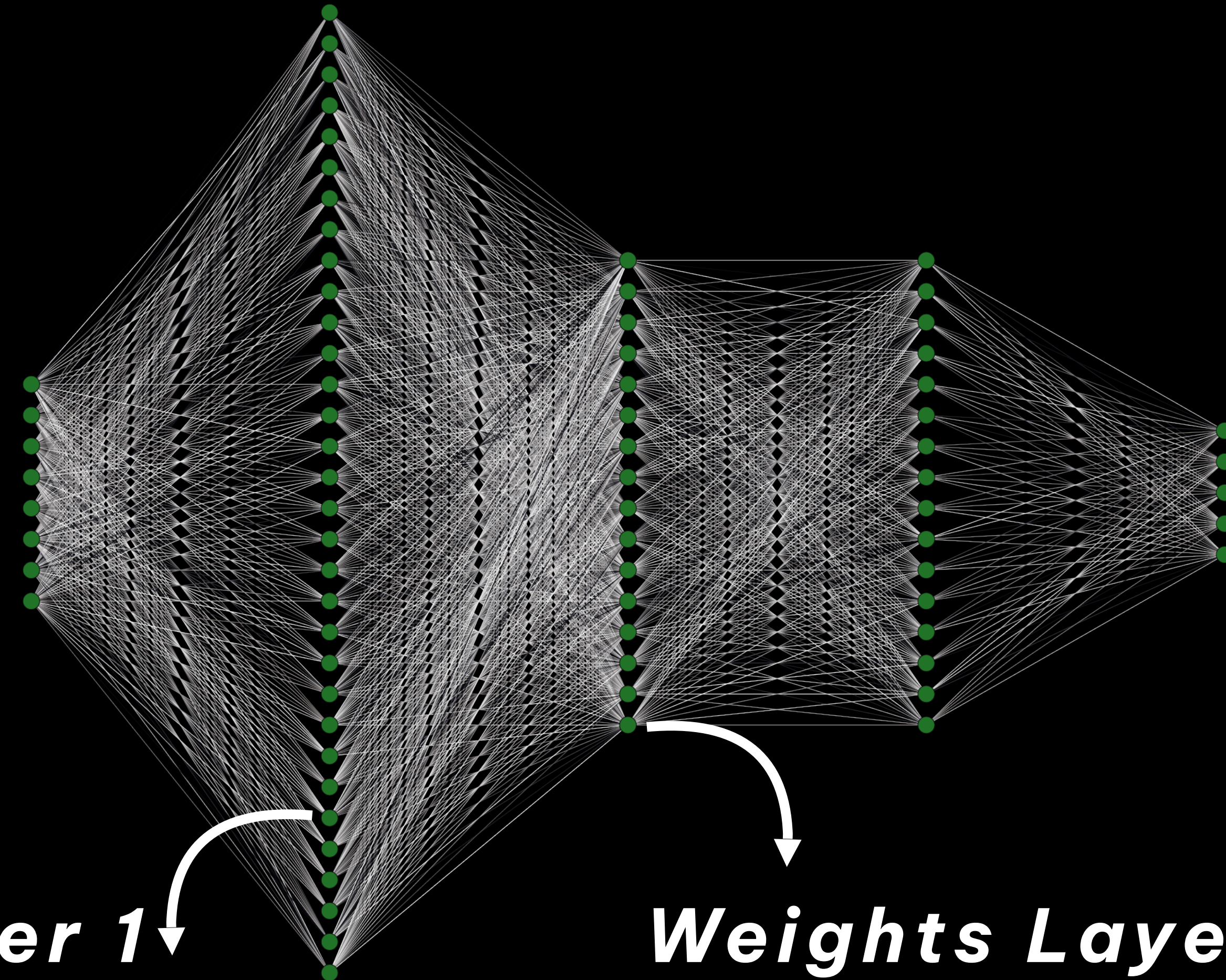
Ideally



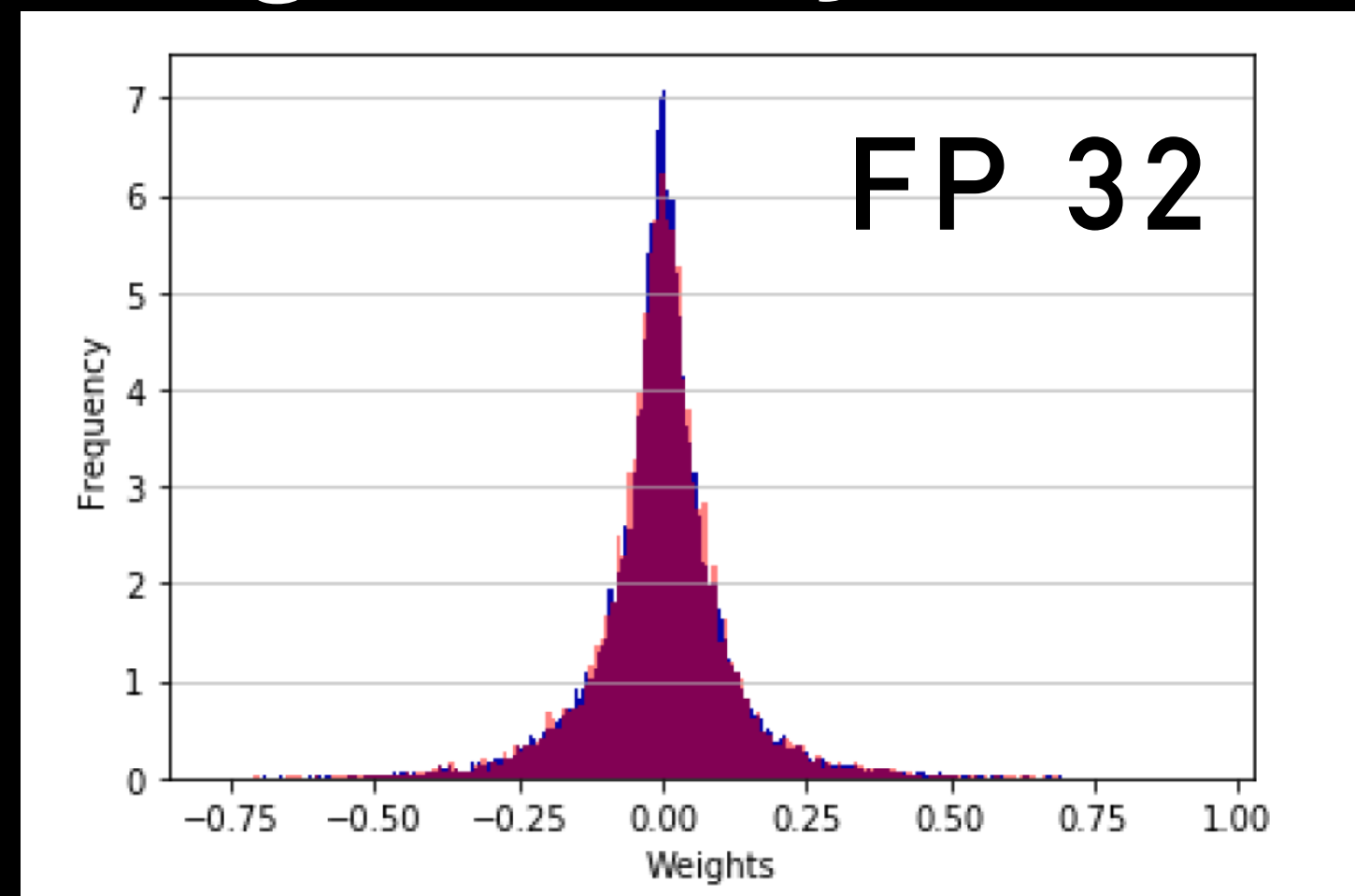
- Quantization
- Pruning
- Parallelisation
- Knowledge distillation



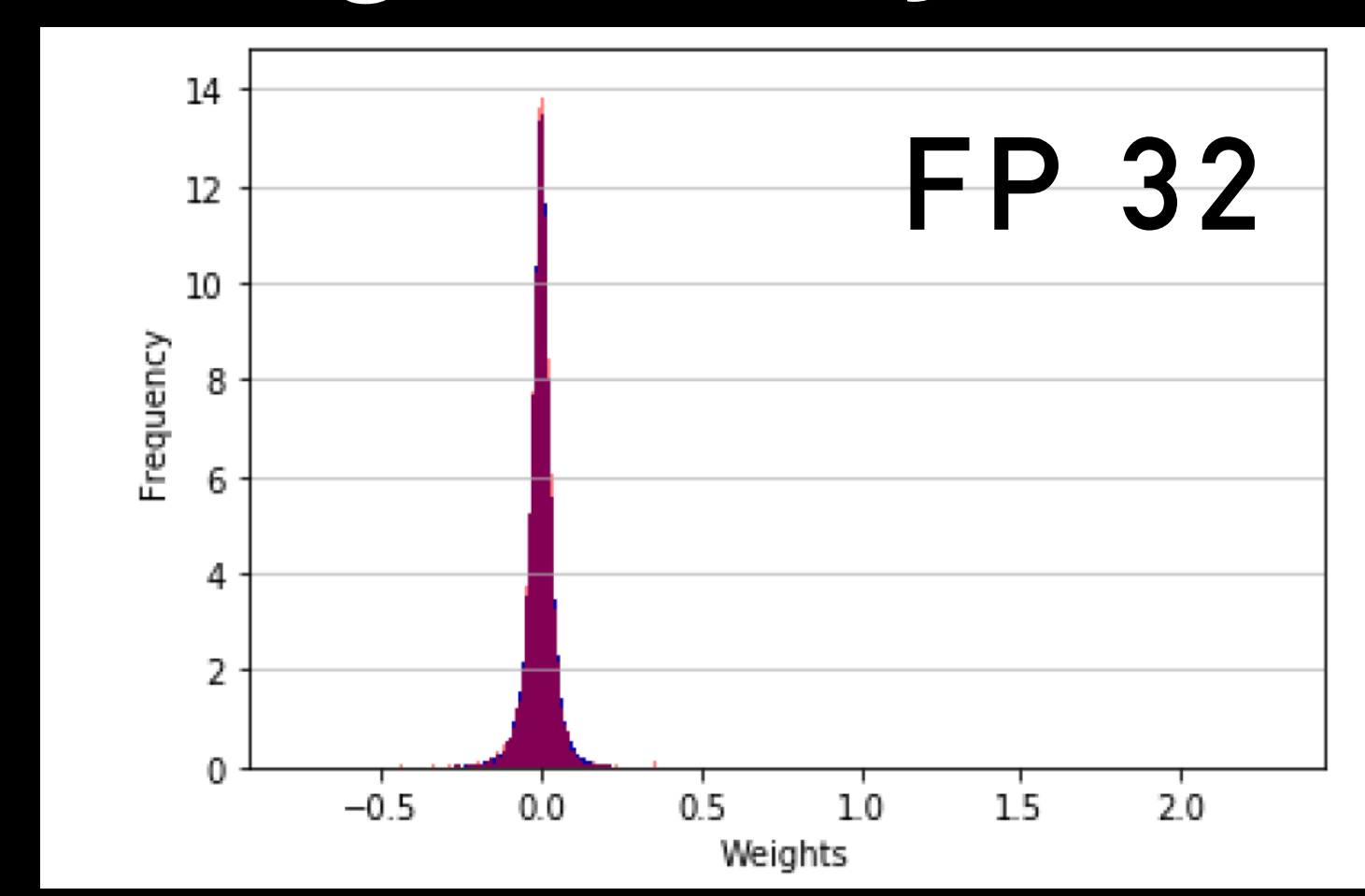
Reality

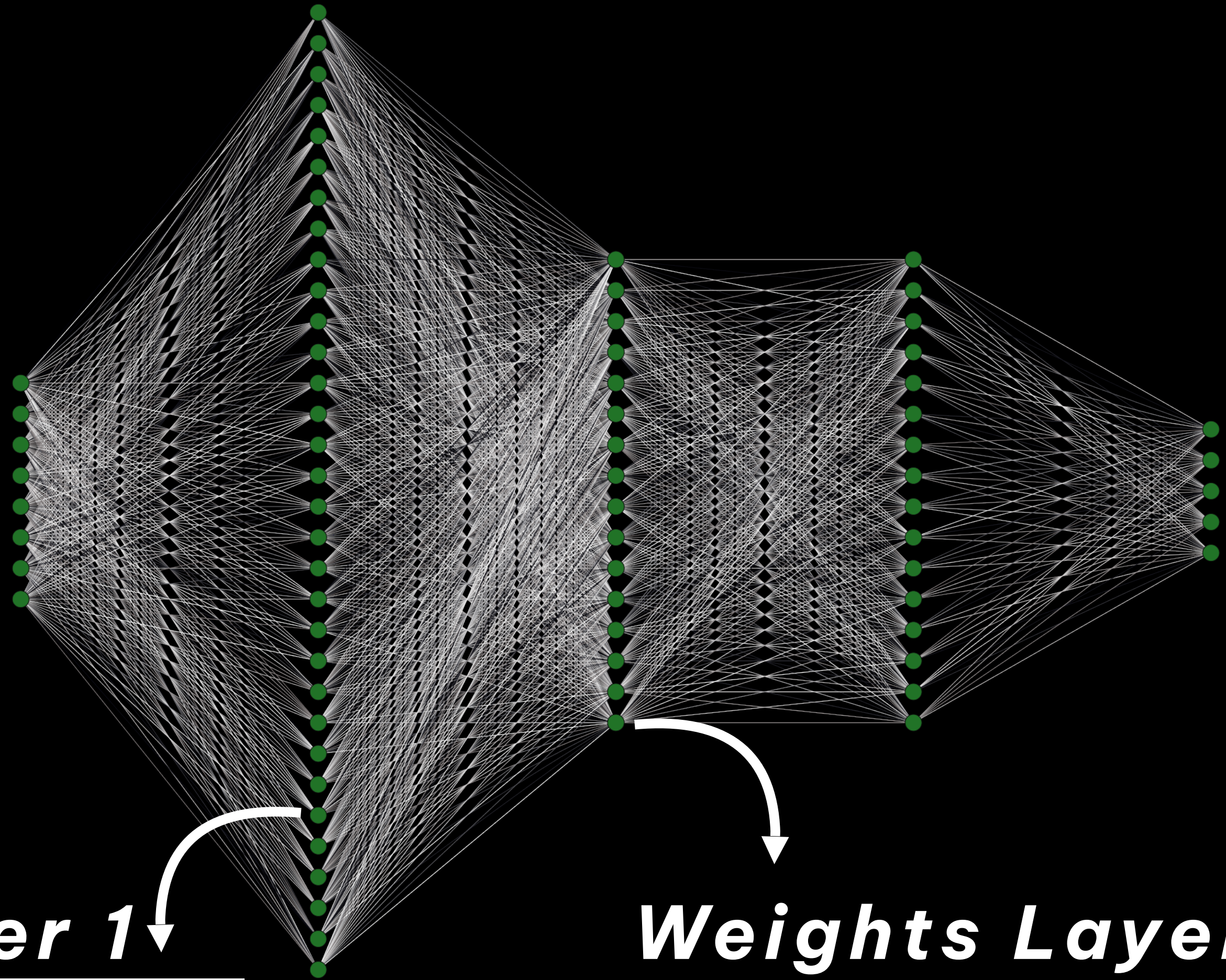


Weights Layer 1



Weights Layer 2



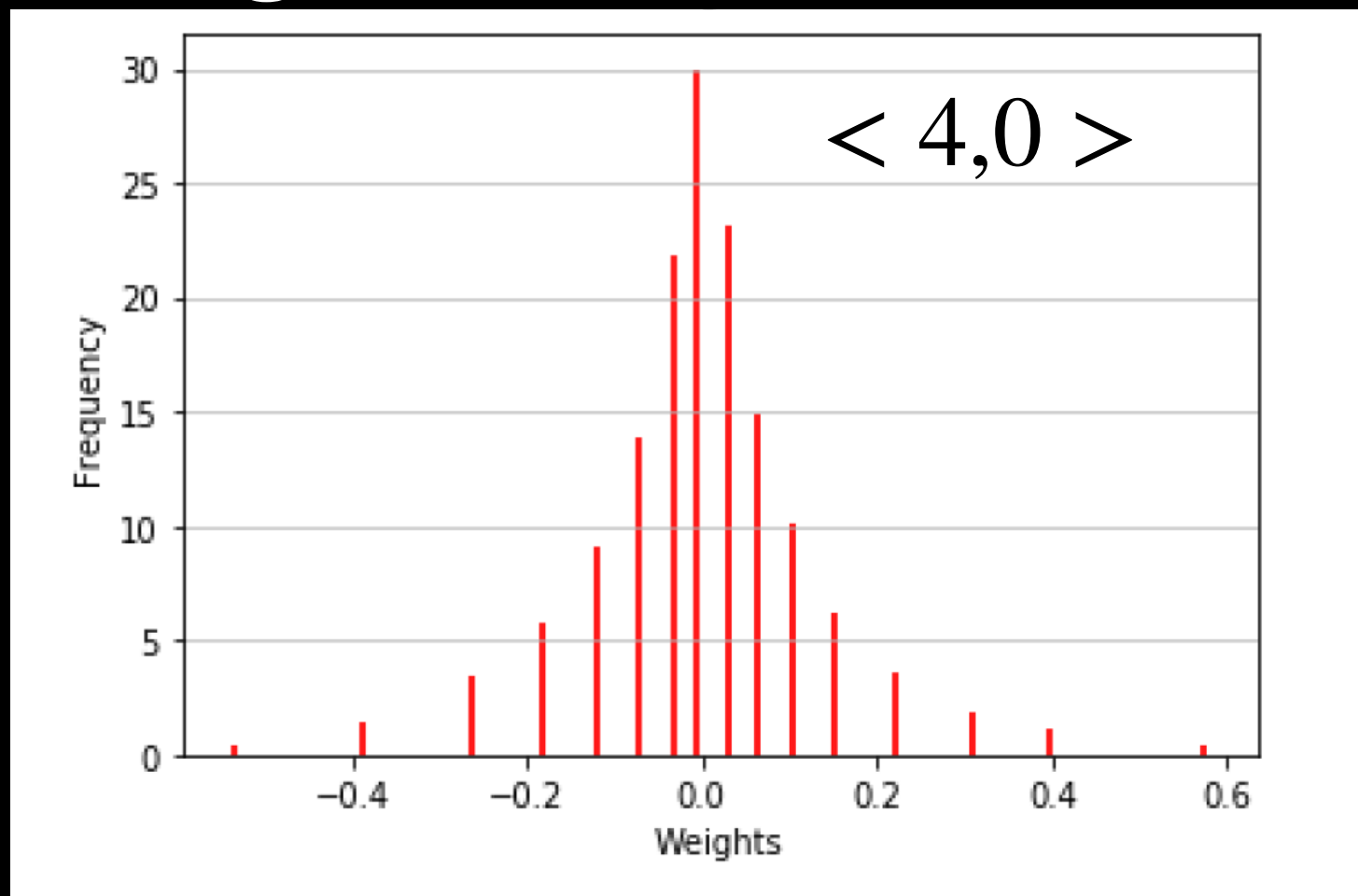


Fixed point

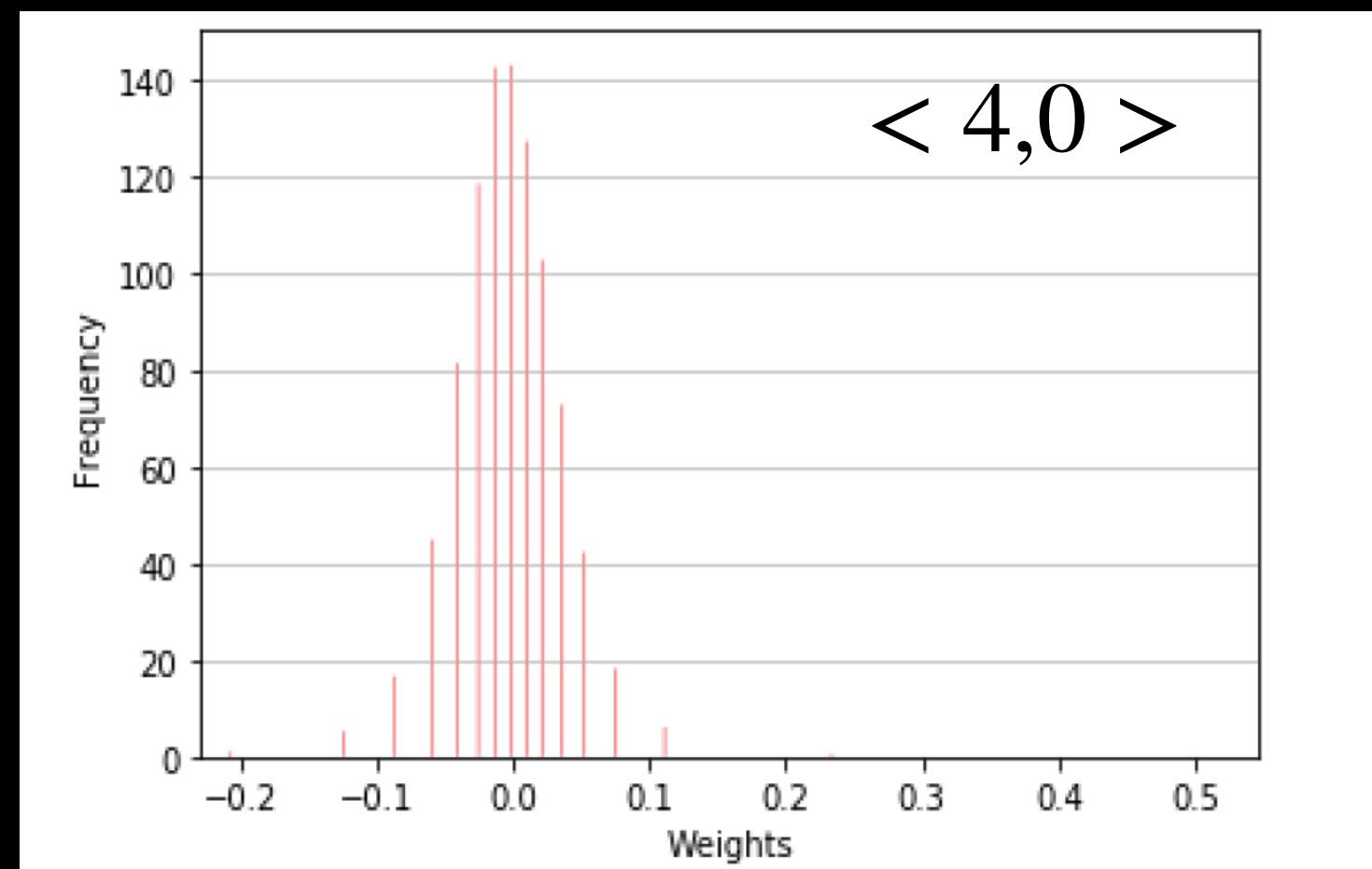
0101.1011101010



Weights Layer 1

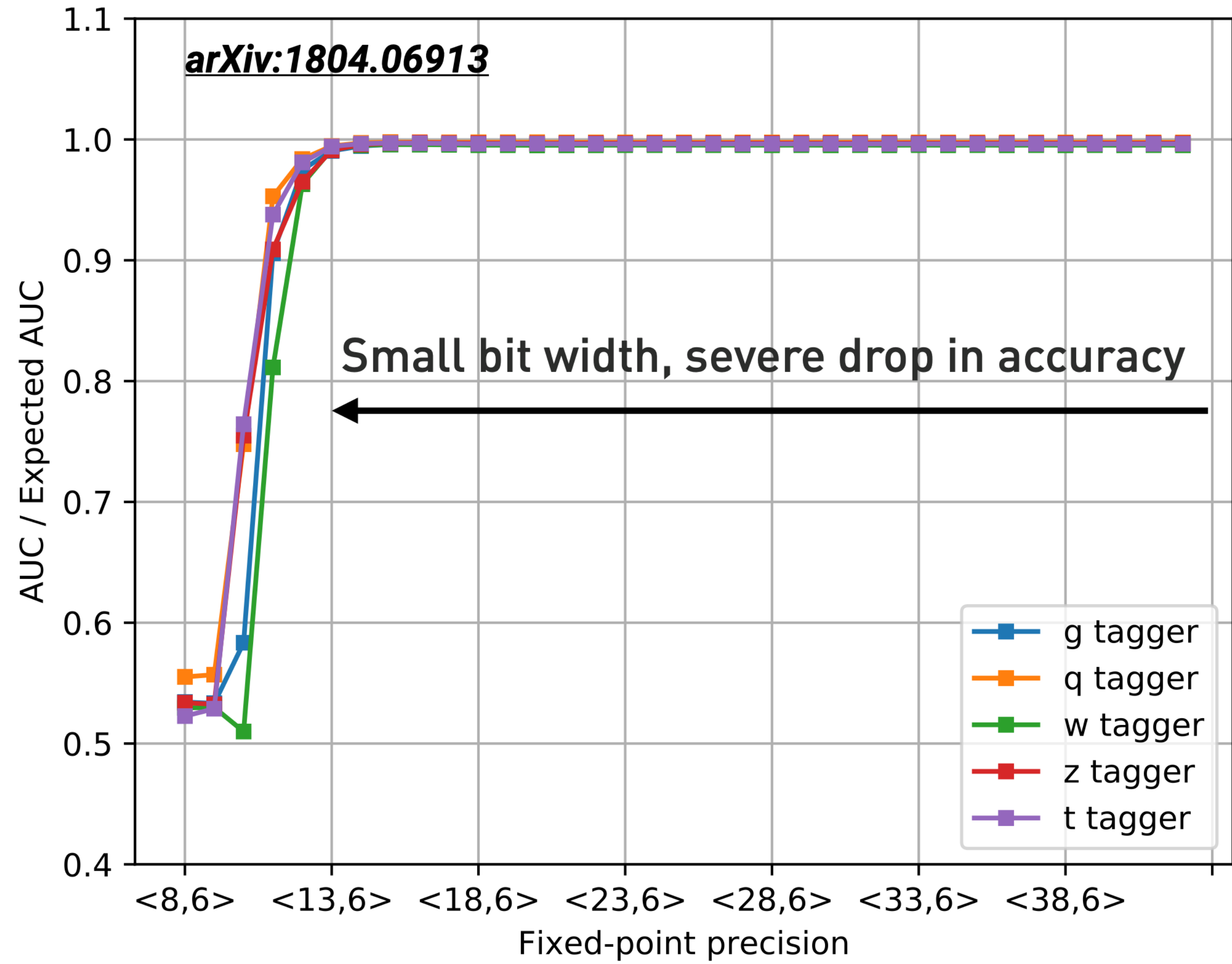


Weights Layer 2



hls4ml

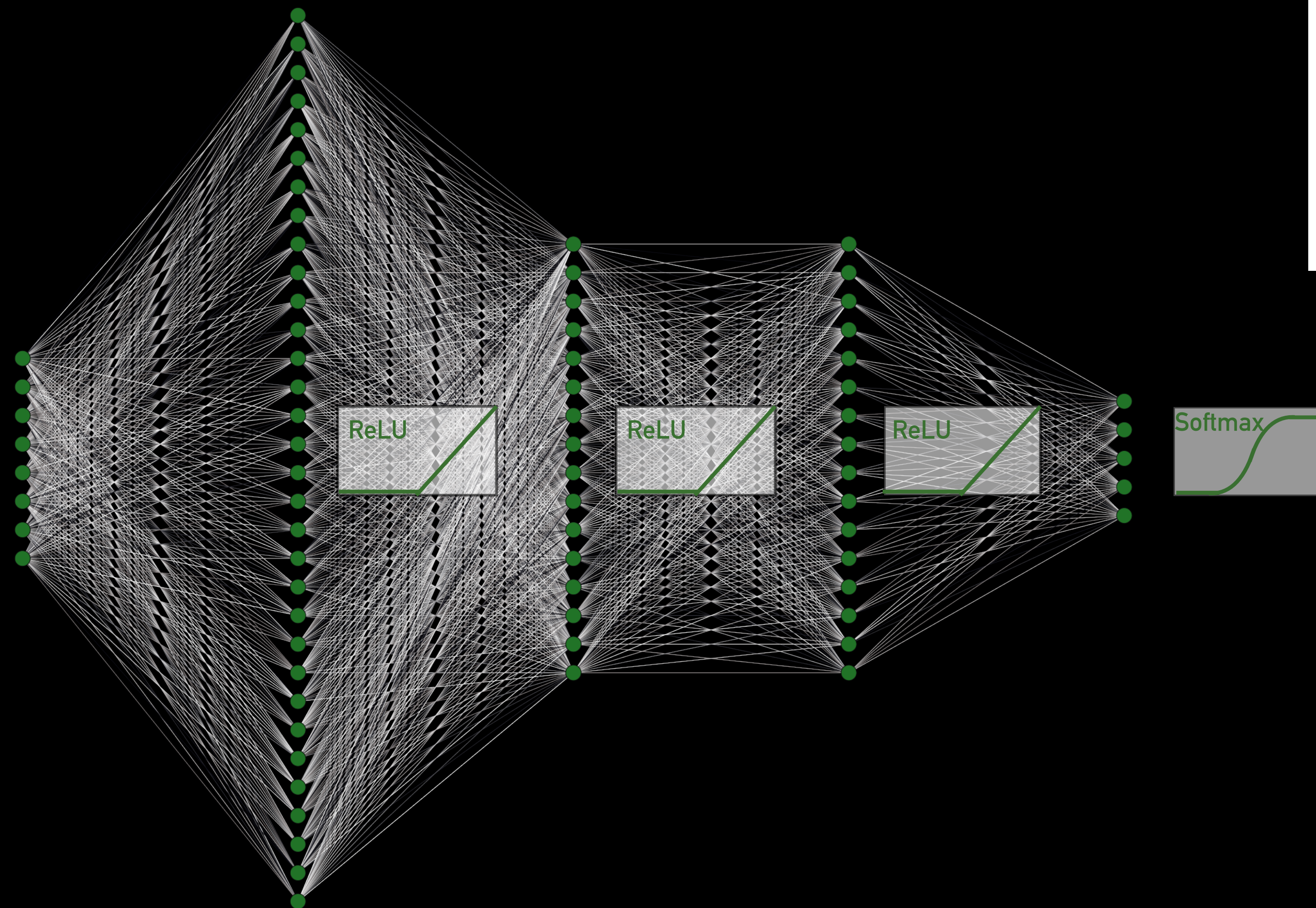
[arXiv:1804.06913](#)



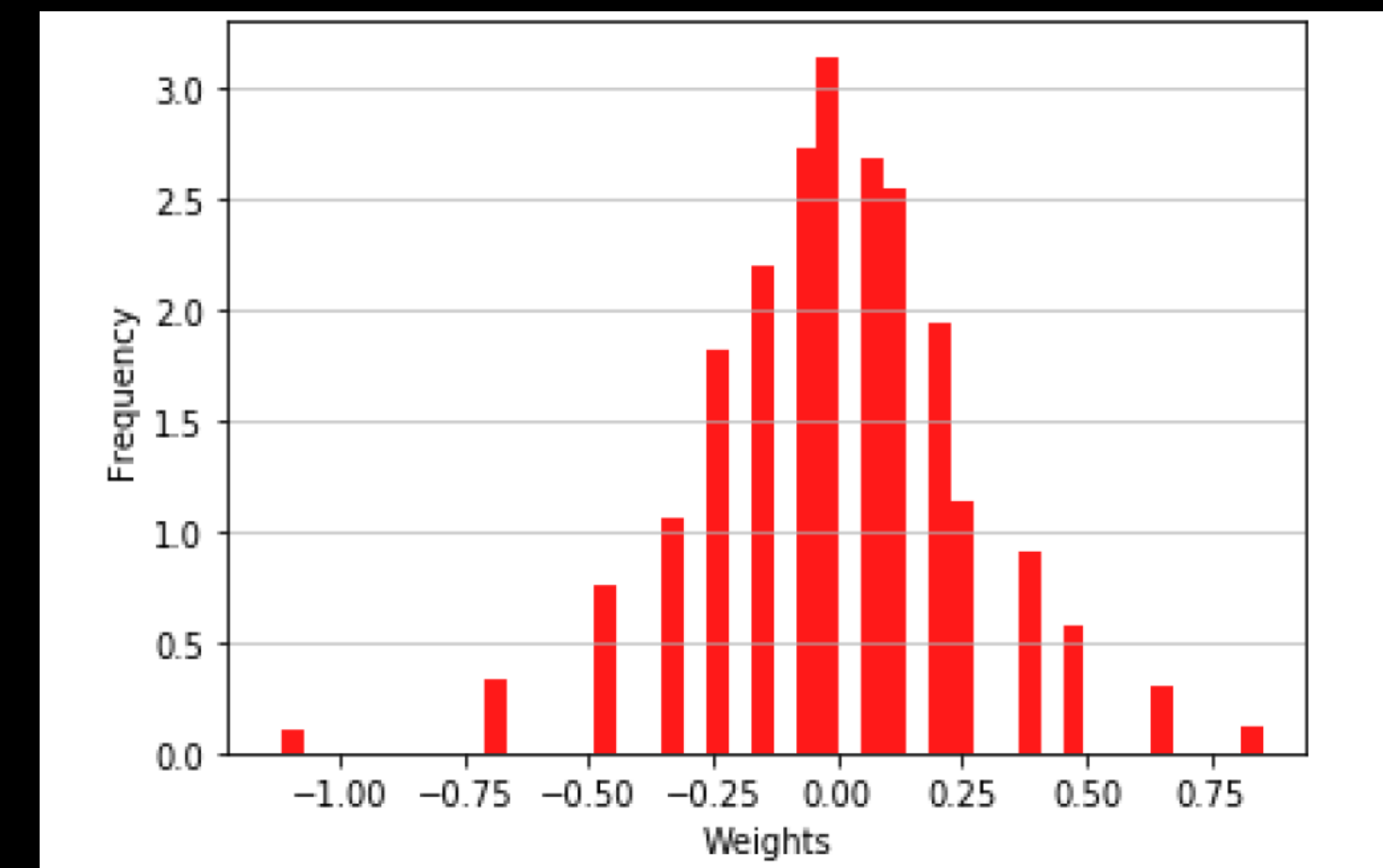
Small bit width, severe drop in accuracy

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

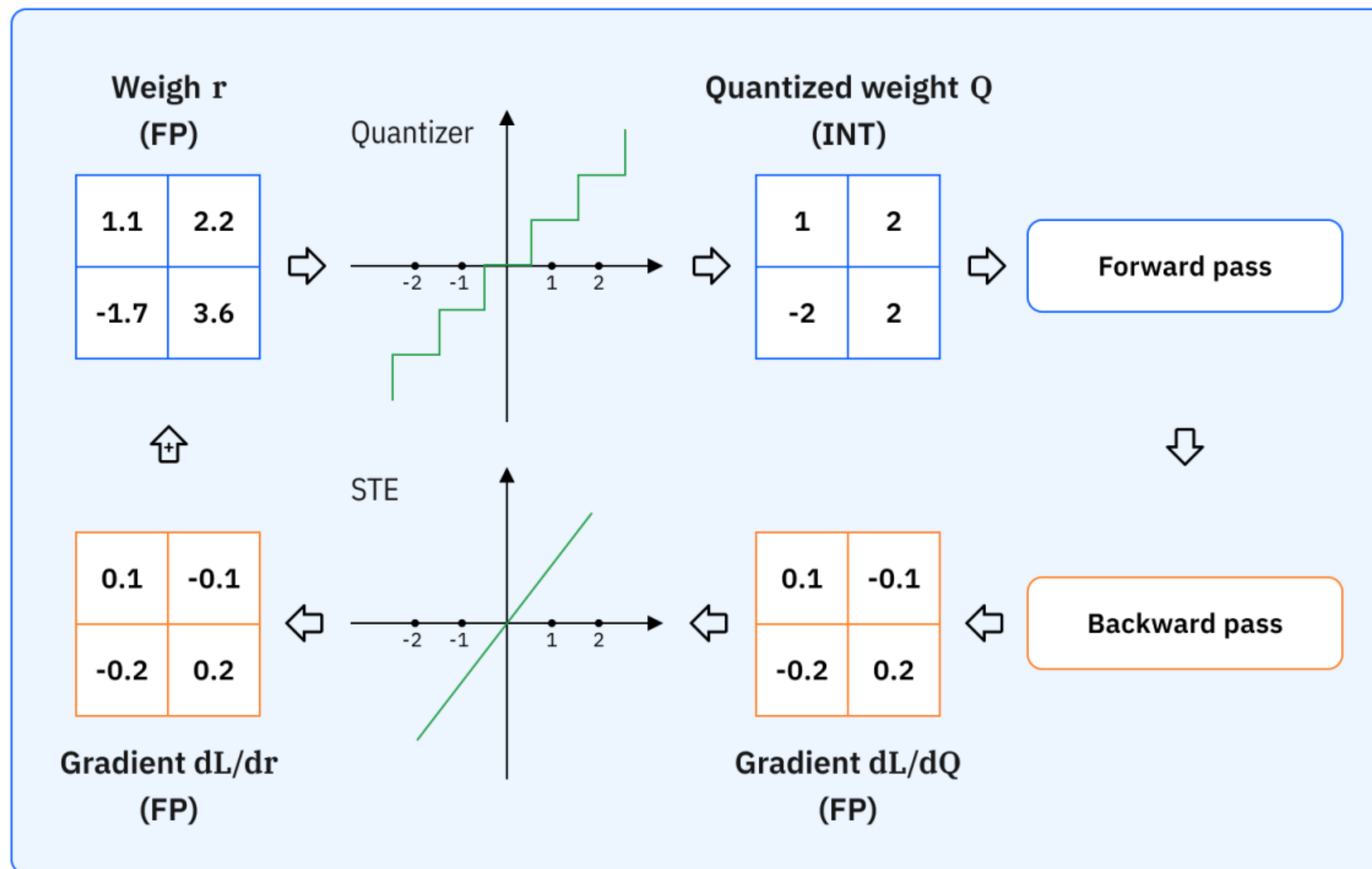
hls4ml + Google
Quantization-aware training



Forward pass →



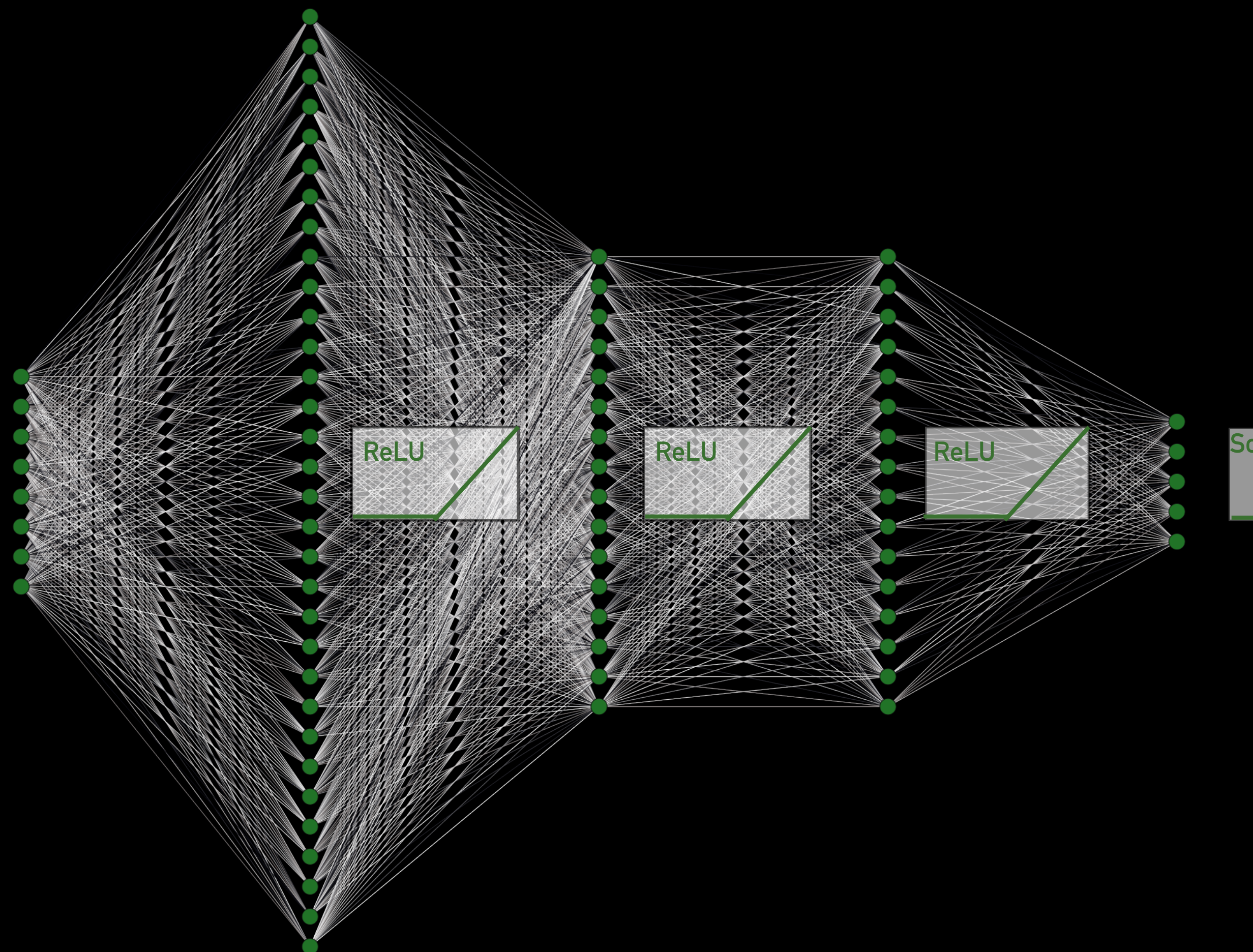
Straight-through estimator



hls4ml + Google Quantization-aware training

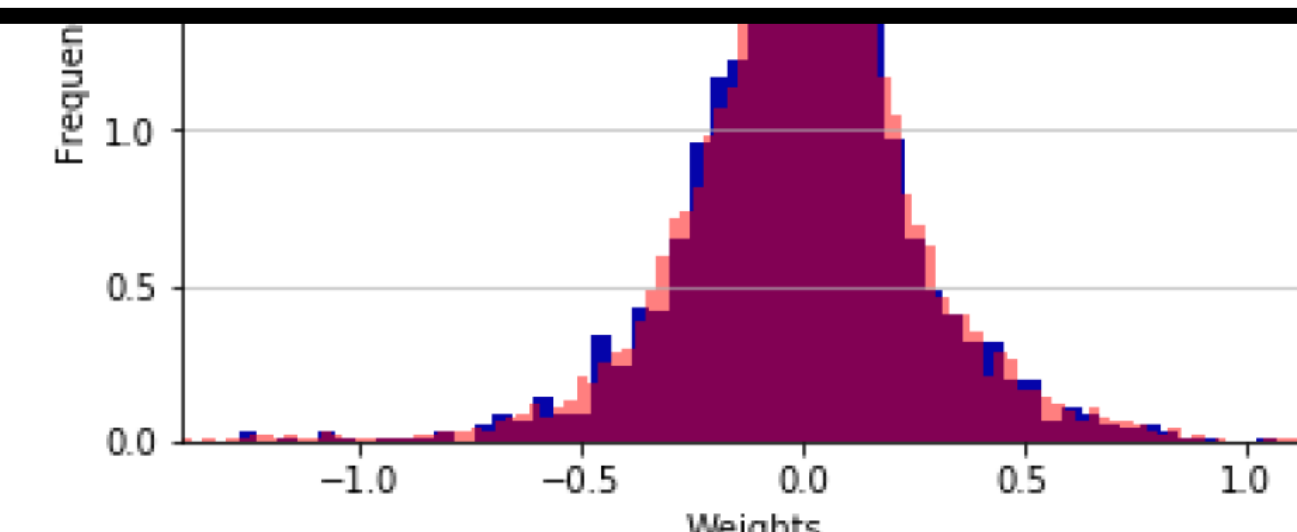
Sundays IEEE RT pre-conference program!

Forward pass →

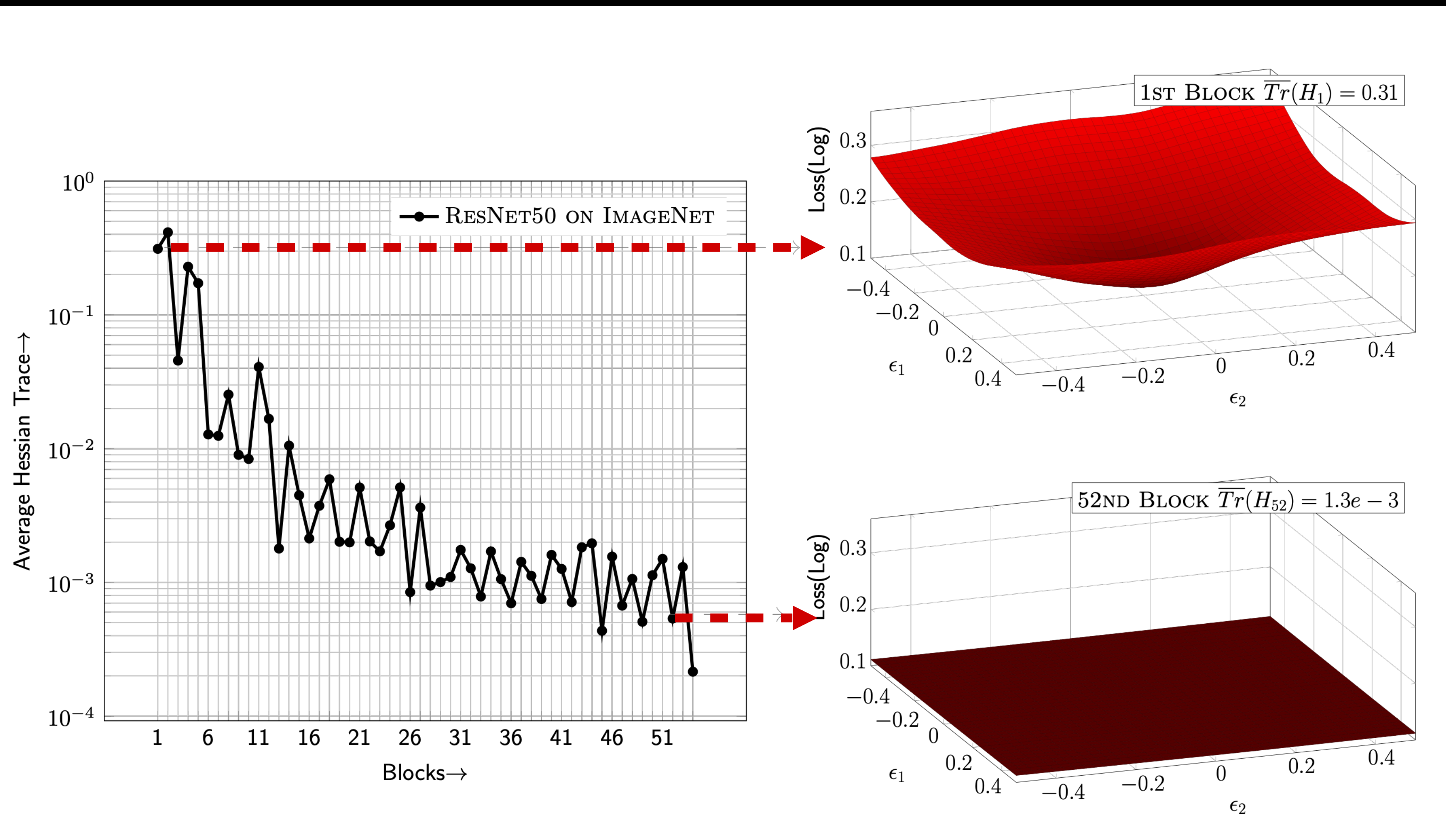


```
from tensorflow.keras.layers import Input, Activation
from qkeras import quantized_bits
from qkeras import QDense, QActivation
from qkeras import QBatchNormalization

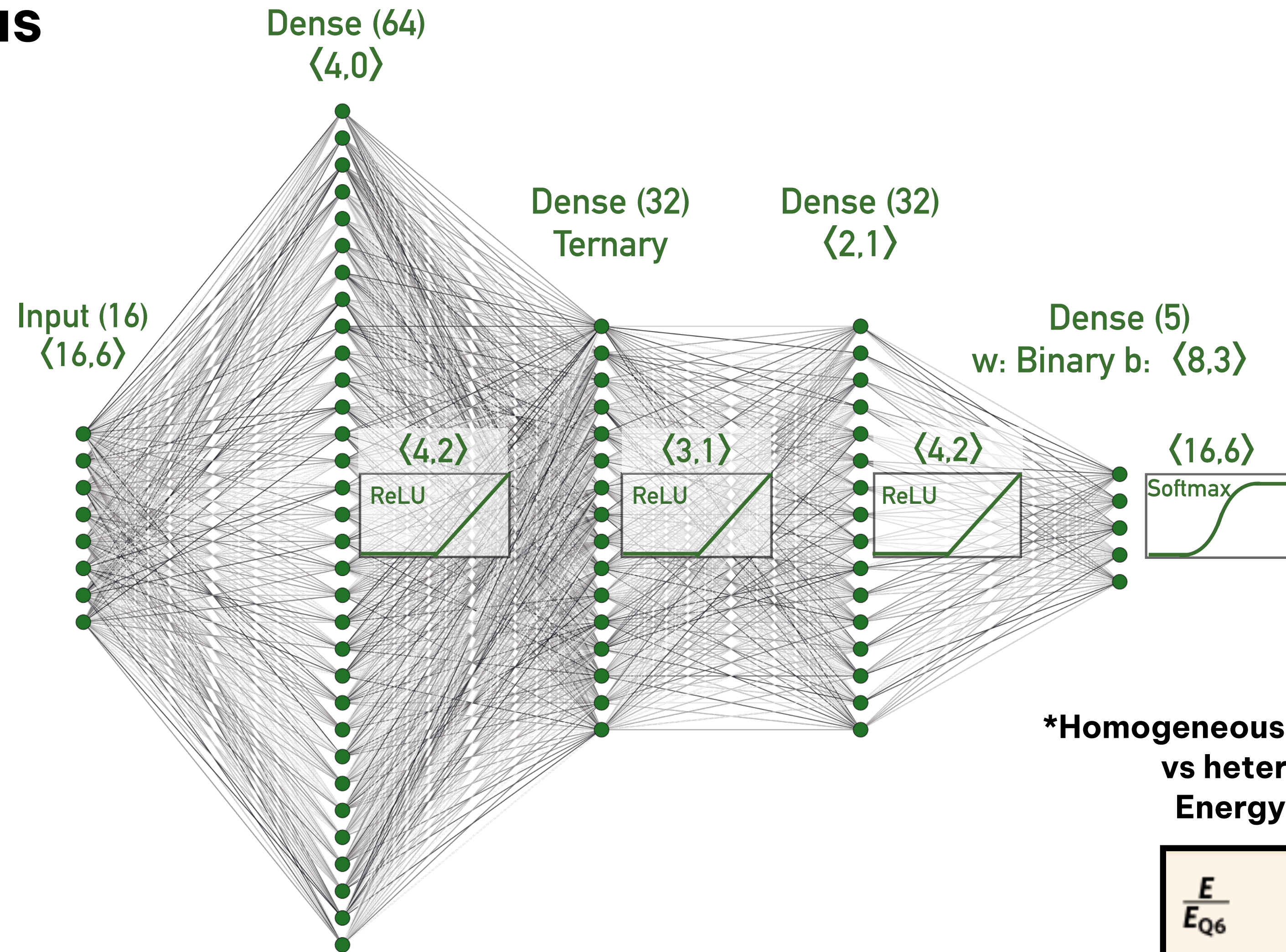
x = Input((16))
x = QDense(64,
          kernel_quantizer = quantized_bits(6,0,alpha=1),
          bias_quantizer   = quantized_bits(6,0,alpha=1))(x)
x = QBatchNormalization()(x)
x = QActivation('quantized_relu(6,0)')(x)
x = QDense(32,
          kernel_quantizer = quantized_bits(6,0,alpha=1),
      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



Some layer more accommodating to aggressive quantization than others!



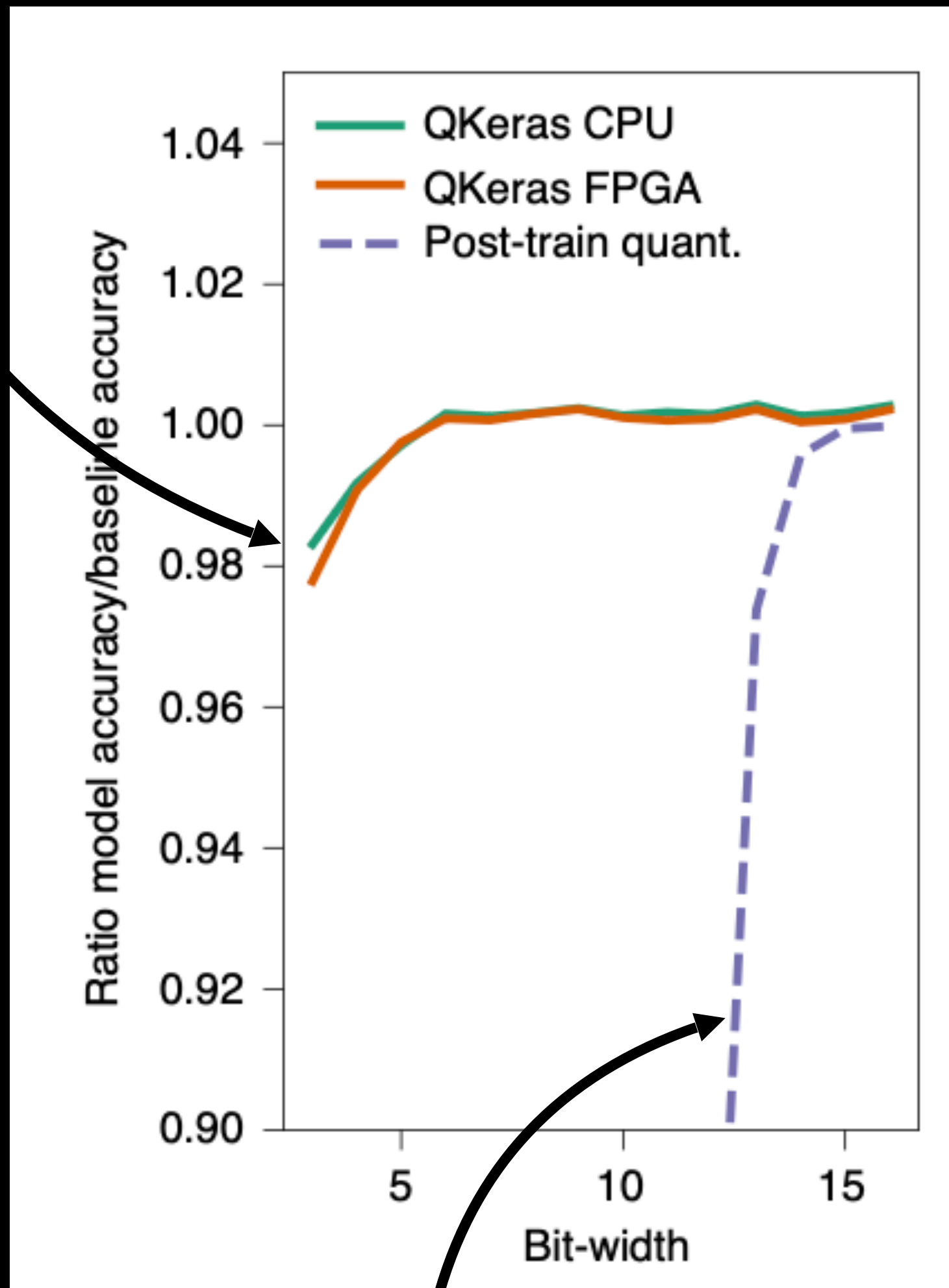
AutoQKeras



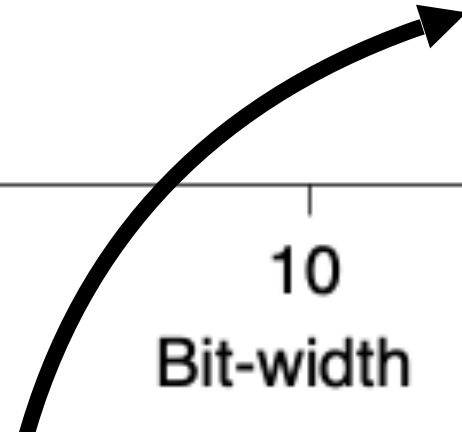
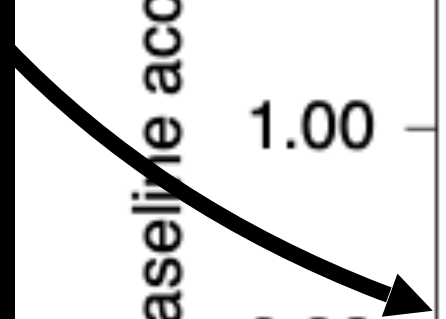
***Homogeneously quantized 6 bit model vs heterogeneous model
Energy cost, chip area**

$\frac{E}{E_{Q6}}$	$\frac{\text{Bits}}{\text{Bits}_{Q6}}$
0.27	0.18

Accuracy

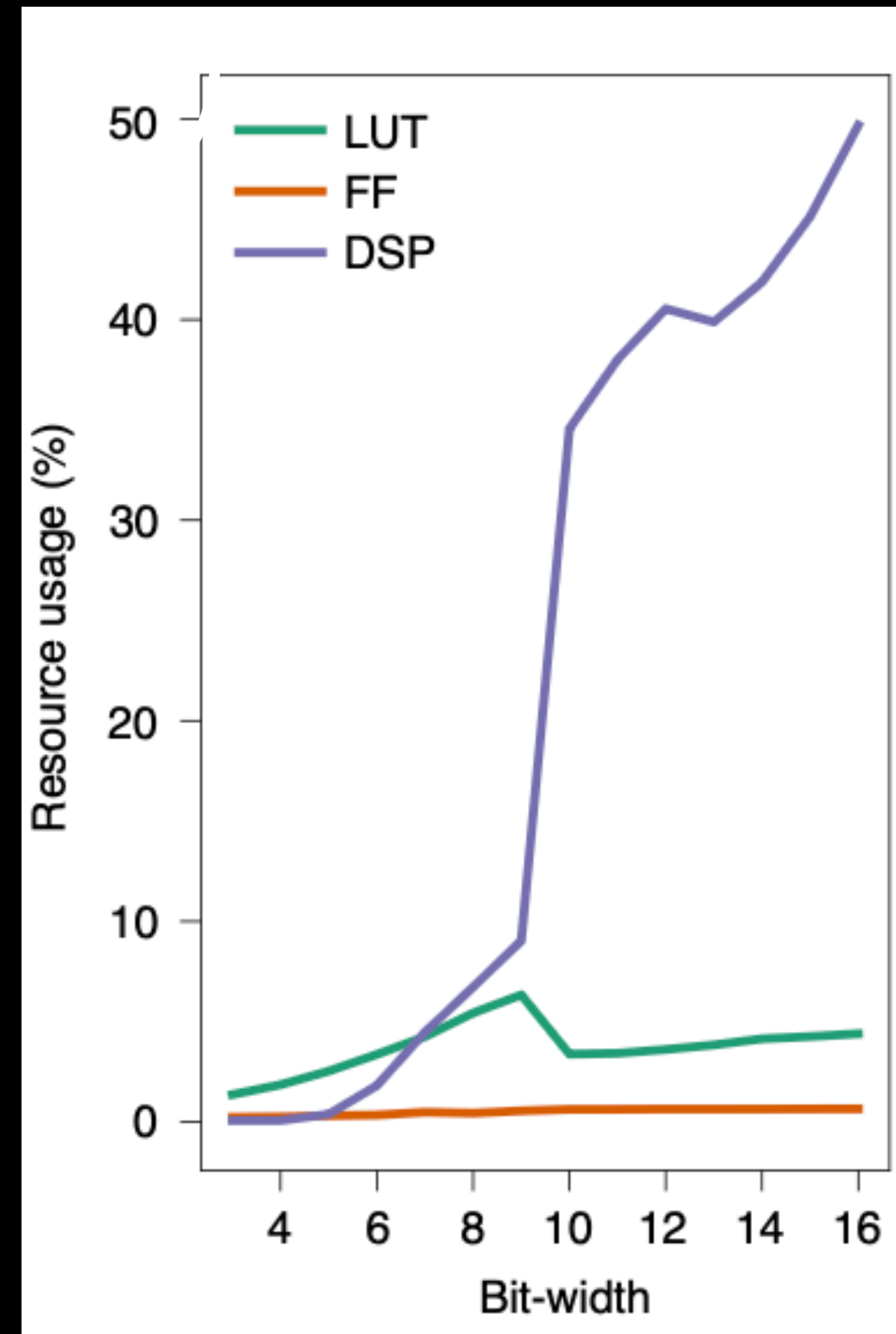


QAT



Post-training quantization

Resource cost



But why stop there?

```
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)
```

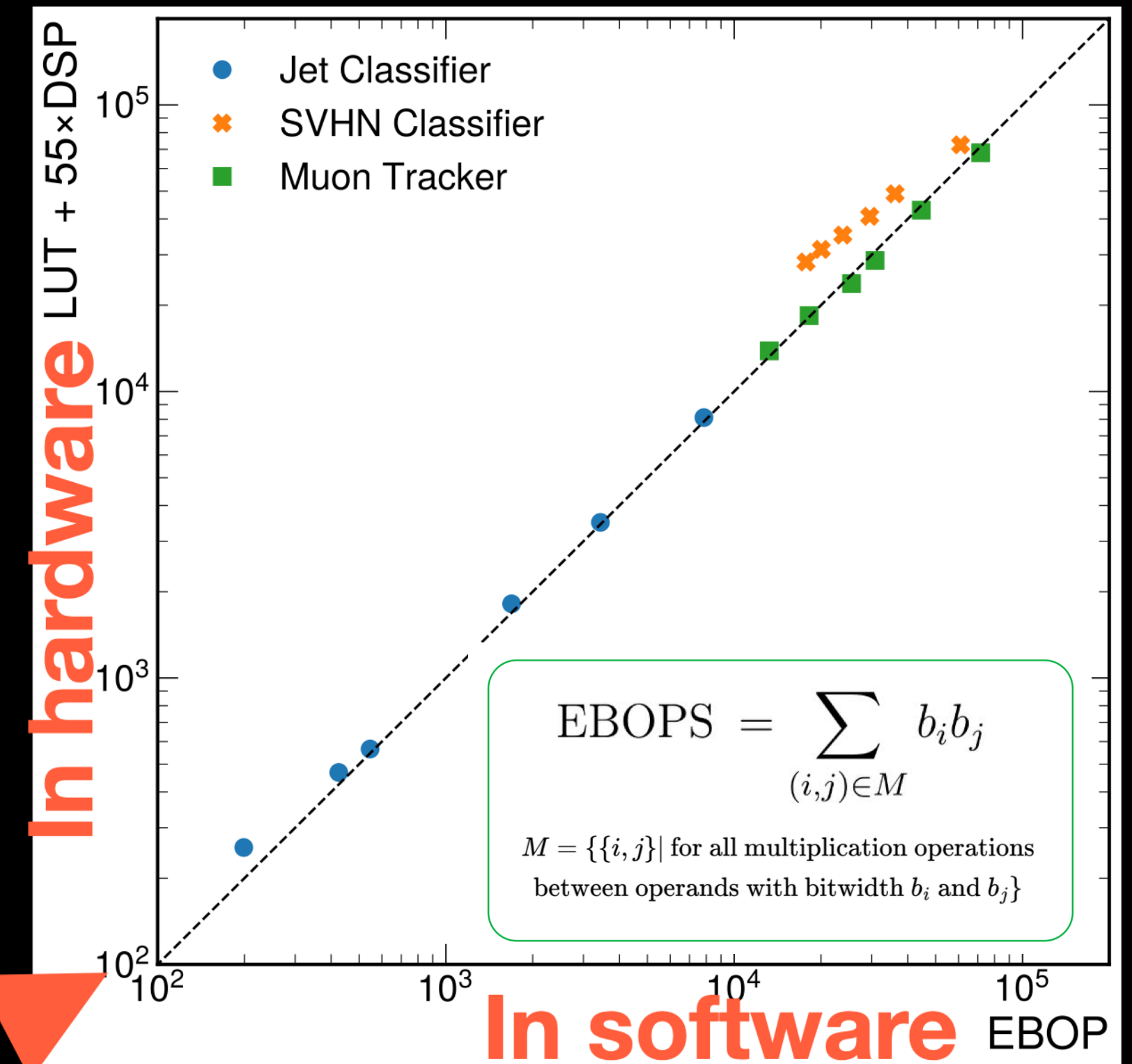
HGQ: High Granularity Quantization

C.Sun et al. 2024

```
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)

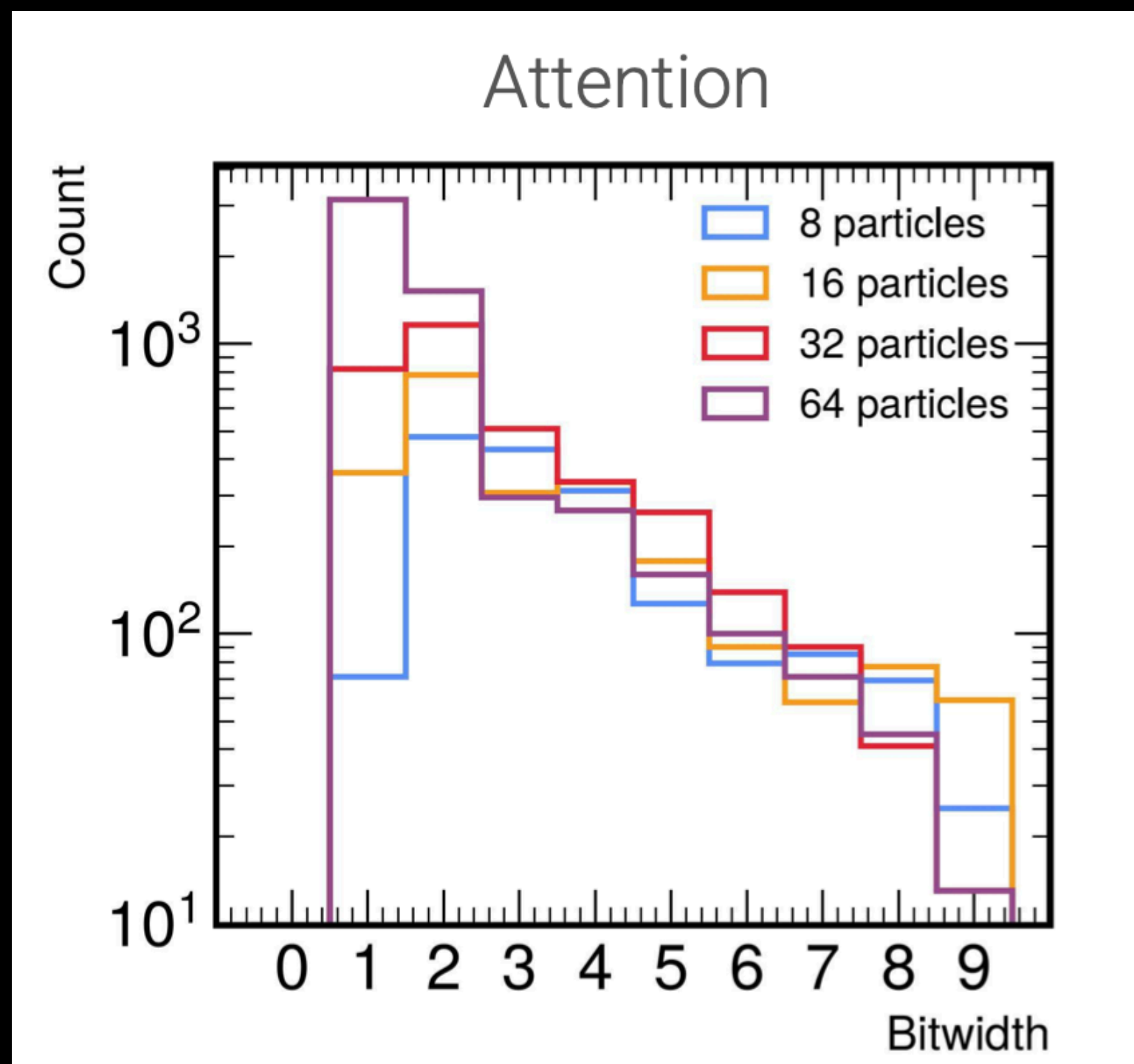
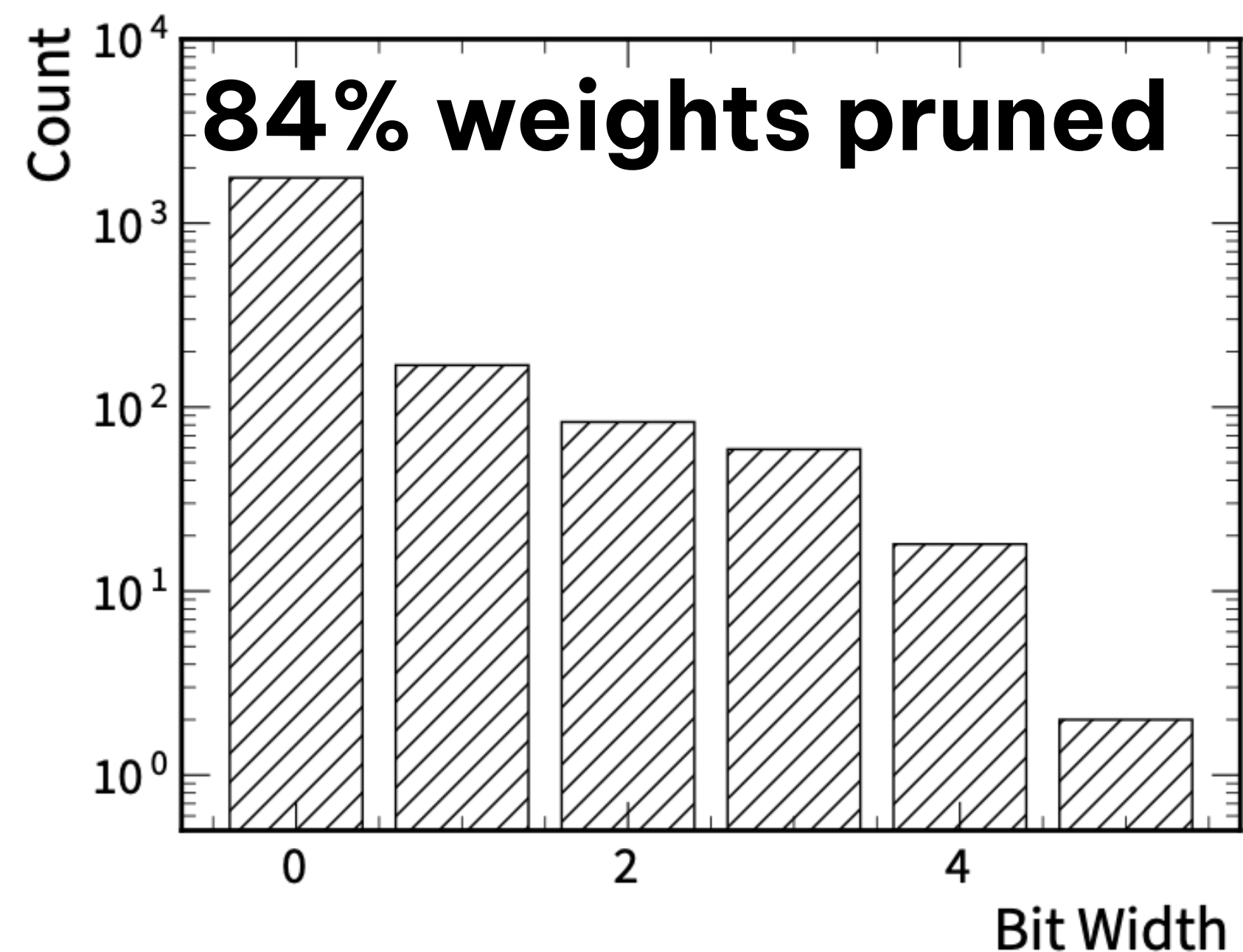
hgq_model = Model(inp, out)
```



Proxy for resource cost

$$\mathcal{L} = \mathcal{L}_{\text{base}} + \beta \cdot \overline{\text{EBOPs}} + \gamma \cdot \text{L1}_{\text{norm}}$$

Making bitwidths differentiable! Fully heterogeneous quantization!



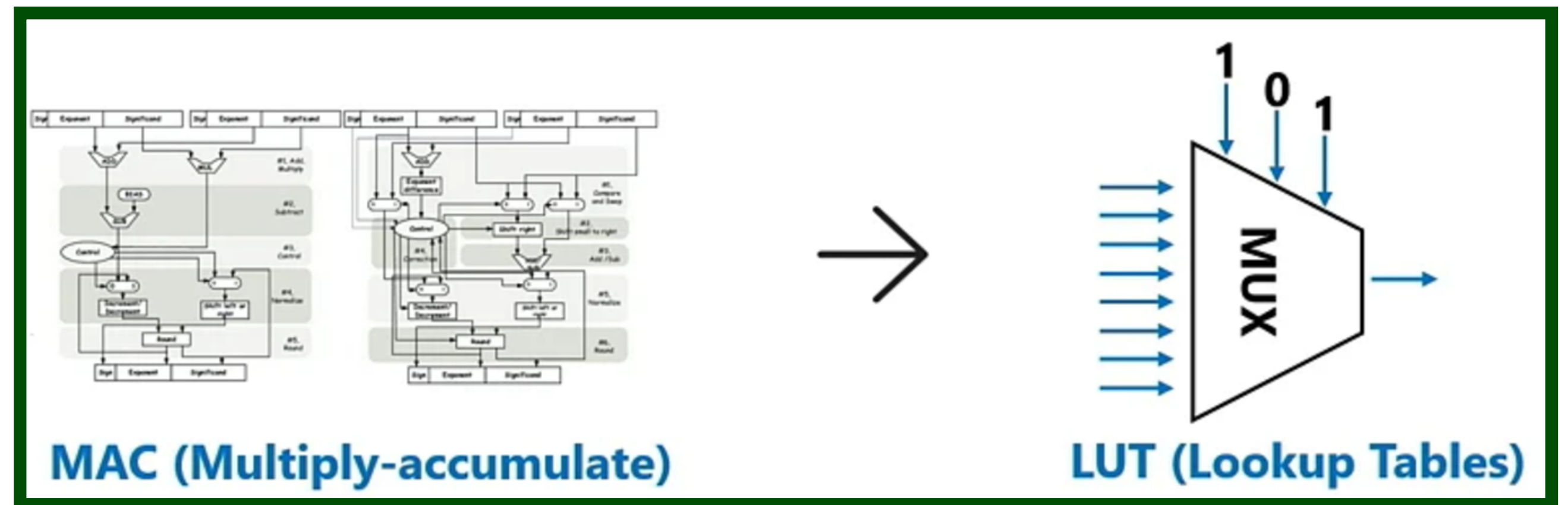
**It's not a FAD:
first results in using Flows for unsupervised
Anomaly Detection at 40 MHz at the Large Hadron
Collider**

**Normalizing Flows on
chip with 35 ns latency**

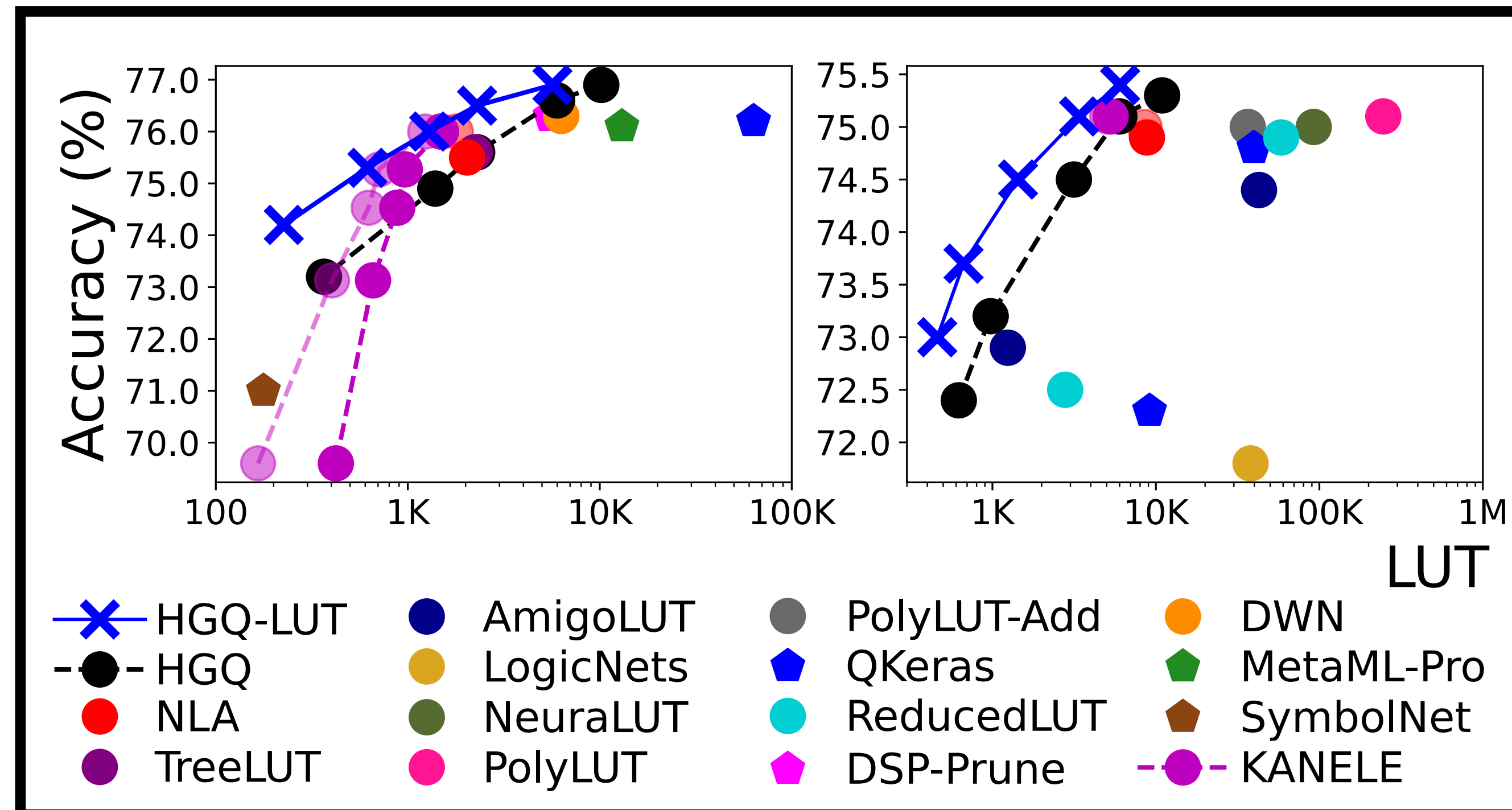
**Heterogeneously
quantized
transformers with 100
ns latency!**

HGQ-LUT

C. Sun et al (2026)




Blue is LUT dense layers
Black is without!



Why do tree-based models still outperform deep learning on typical tabular data?

Leo Grinsztajn, Edouard Oyallon, Gael Varoquaux

06 Jun 2022 (modified: 16 Jan 2023) NeurIPS 2022 Datasets and Benchmarks Readers:  Everyone [Show Bibtex](#) [Show Revisions](#)

Abstract: While deep learning has enabled tremendous progress on text and image datasets, its superiority on tabular data is not clear. We contribute extensive benchmarks of standard and novel deep learning methods as well as tree-based models such as XGBoost and Random Forests, across a large number of datasets and hyperparameter combinations. We define a standard set of 45 datasets from varied domains with clear characteristics of tabular data and a benchmarking methodology accounting for both fitting models and finding good hyperparameters. Results show that tree-based models remain state-of-the-art on medium-sized data ($\sim 10K$ samples) even without accounting

Computer Science > Machine Learning

[Submitted on 11 Oct 2022 (v1), last revised 25 Oct 2022 (this version, v3)]

Neural Networks are Decision Trees

[Caglar Aytakin](#)

In this manuscript, we show that any neural network with any activation function can be represented as a decision tree. The representation is equivalence and not an approximation, thus keeping the accuracy of the neural network exactly as is. We believe that this work provides better understanding of neural networks and paves the way to tackle their black-box nature. We share equivalent trees of some neural networks and show that besides providing interpretability, tree representation can also achieve some computational advantages for small networks. The analysis holds both for fully connected and convolutional networks, which may or may not also include skip connections and/or normalizations.

Subjects: **Machine Learning (cs.LG)**

Cite as: [arXiv:2210.05189 \[cs.LG\]](#)

(or [arXiv:2210.05189v3 \[cs.LG\]](#) for this version)

<https://doi.org/10.48550/arXiv.2210.05189> 

Submission history

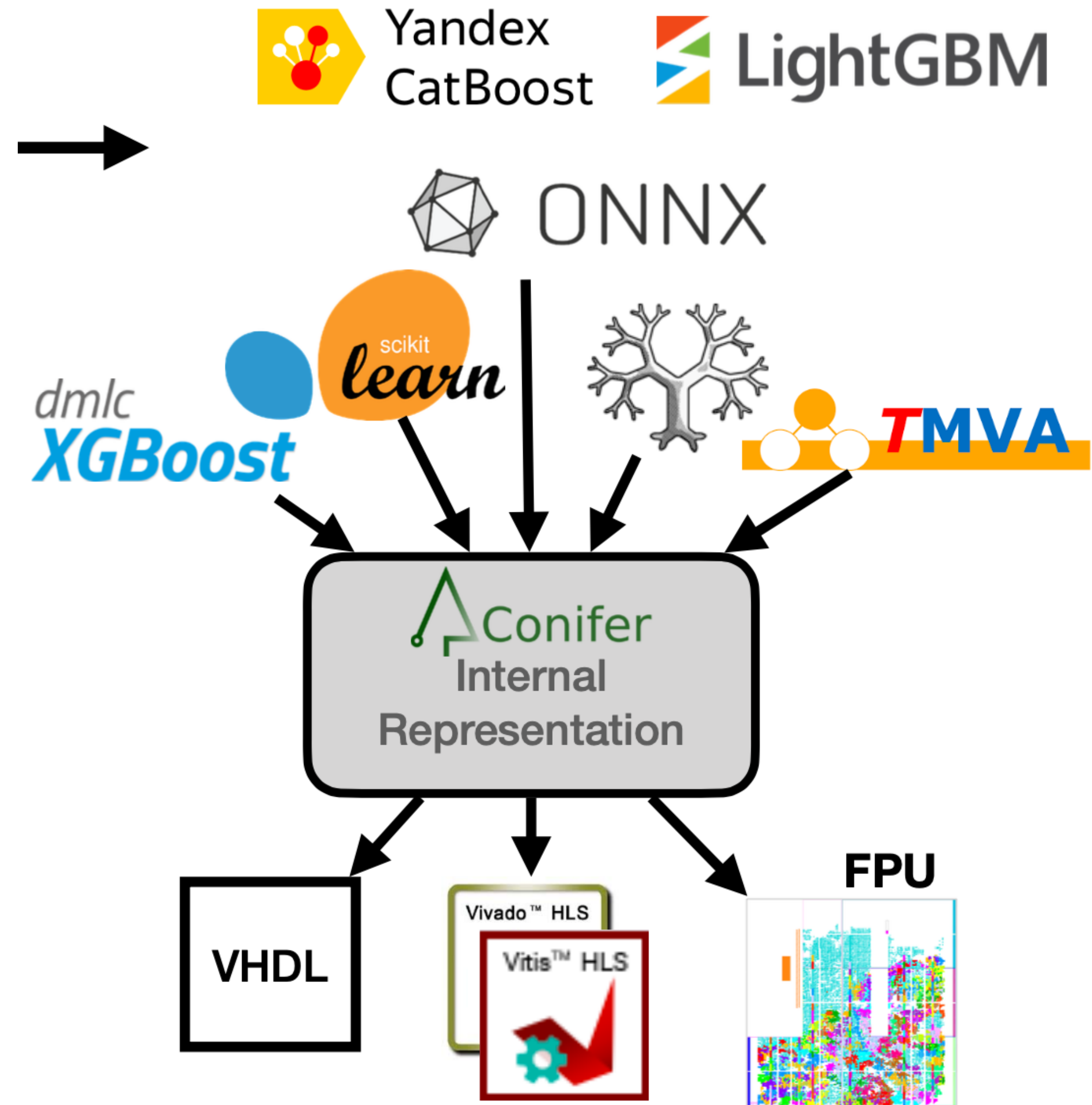
From: Çağlar Aytakin [[view email](#)]

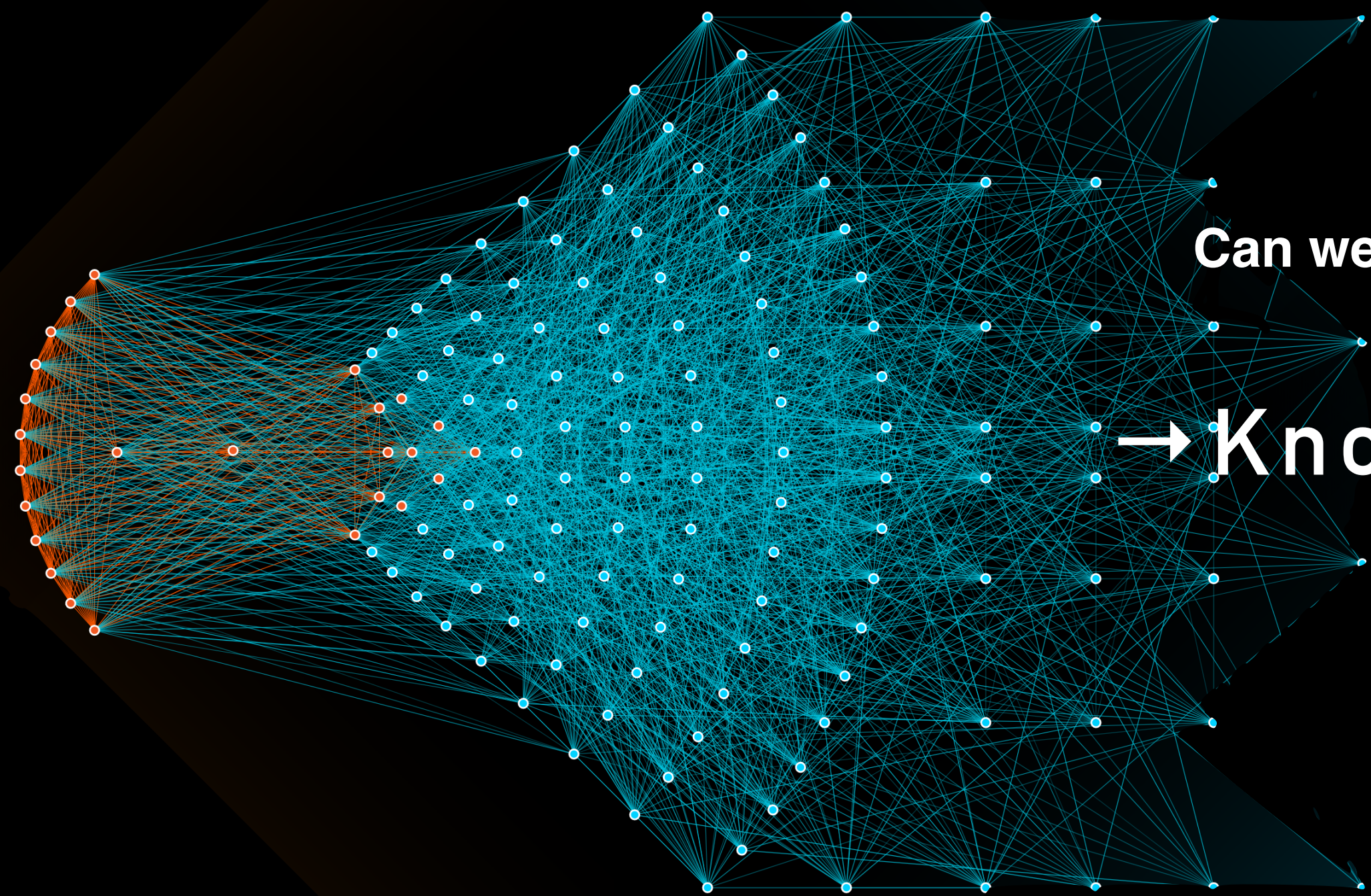
[v1] Tue, 11 Oct 2022 06:49:51 UTC (216 KB)

[v2] Mon, 17 Oct 2022 15:18:14 UTC (224 KB)

[v3] Tue, 25 Oct 2022 17:32:33 UTC (240 KB)

%VU9P	Accuracy	Latency	DSP	LUT
qDNN	75.6%	40 ns	22 (~0%)	1%
BDT	74.9%	5 ns	-	0.5%

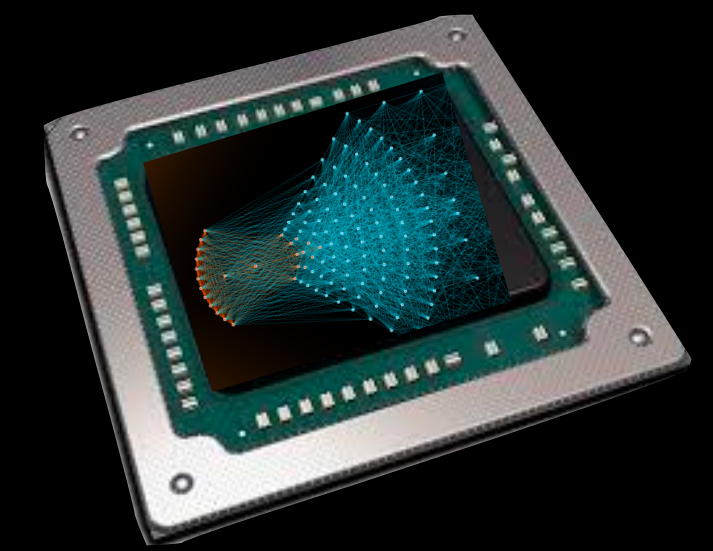




Train

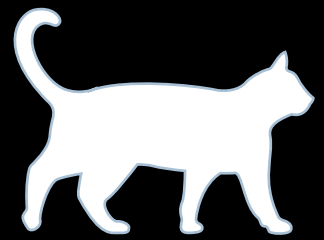
Can we have the best of both worlds?

→ Knowledge Distillation

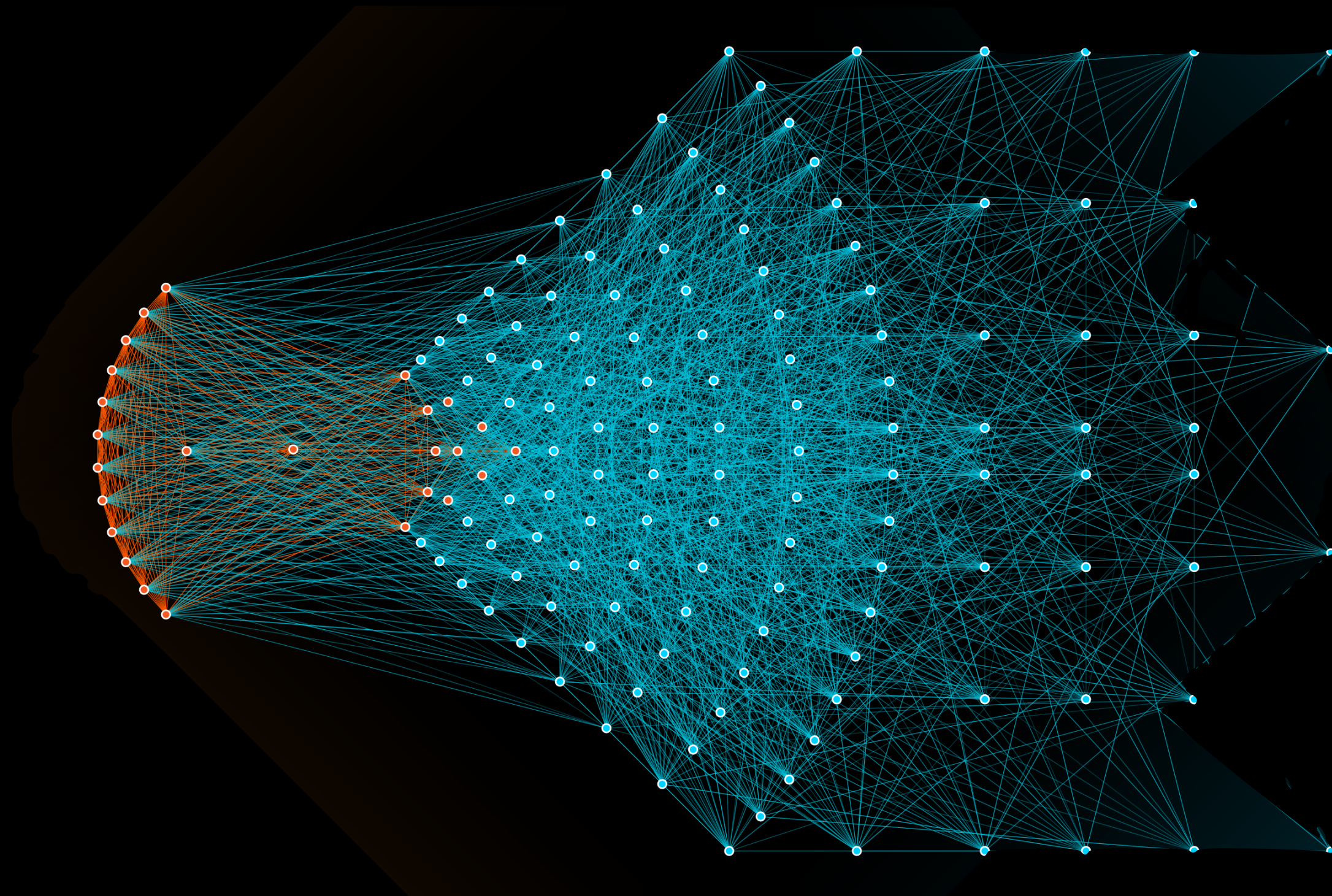
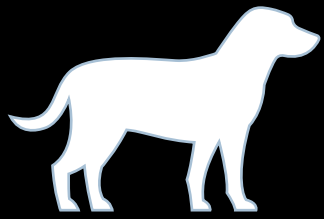


Inference

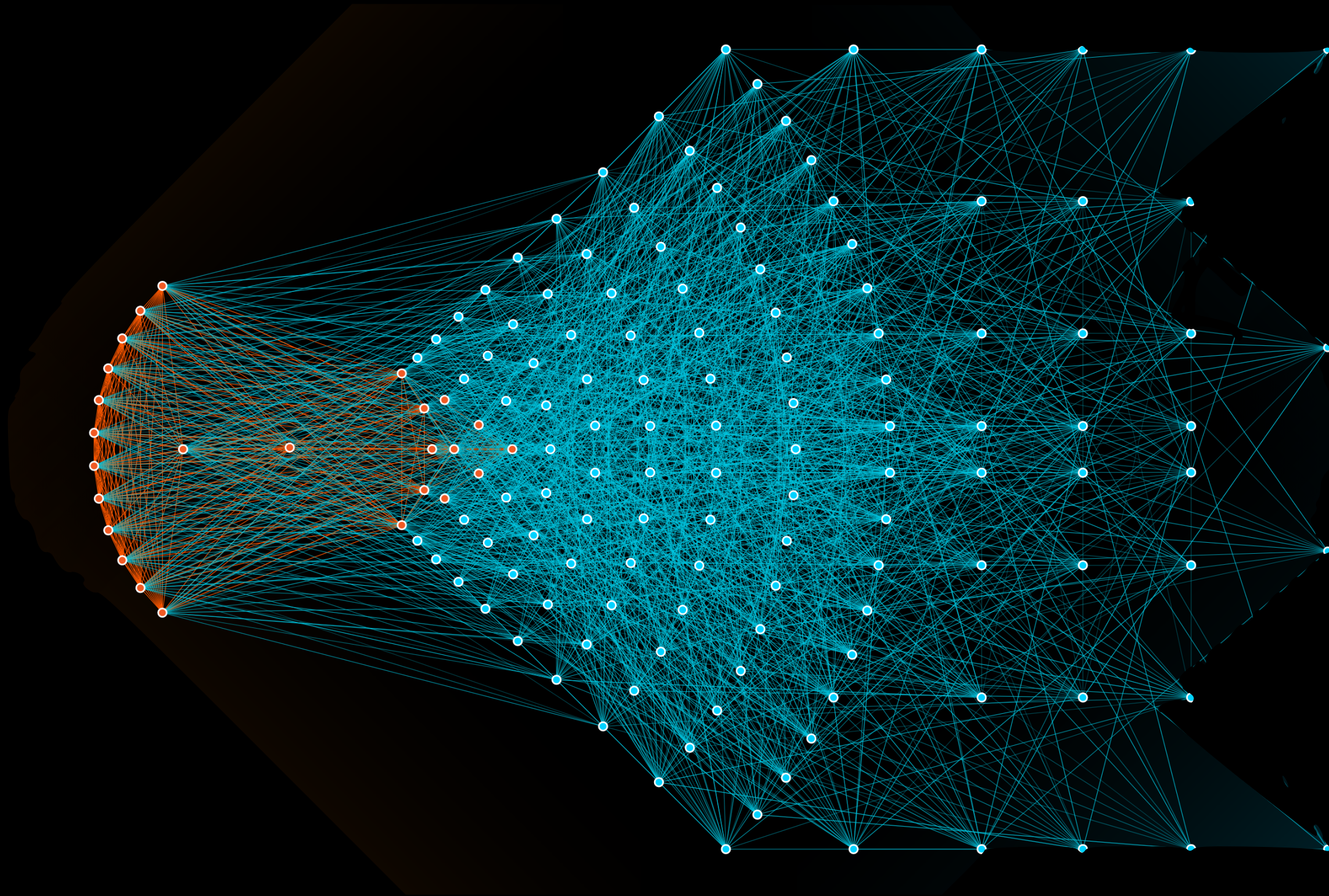
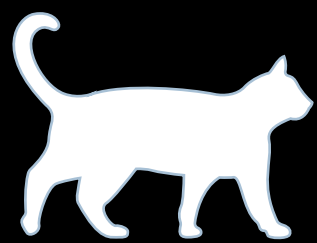
Cat



Dog

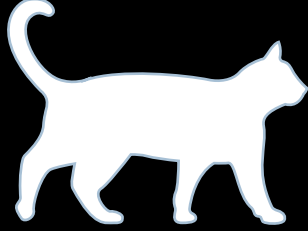


Cat



is cat

is dog

Cat 

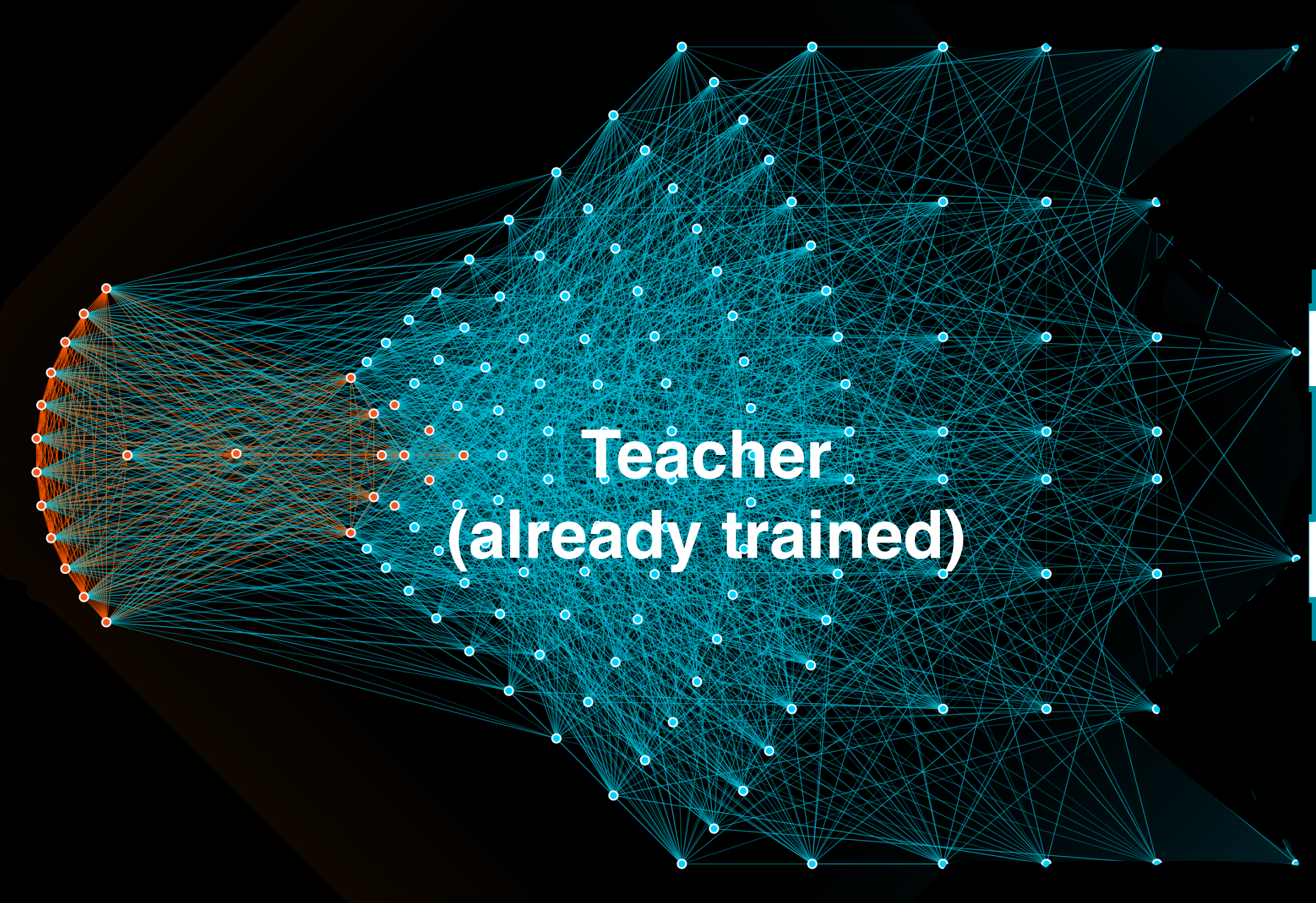


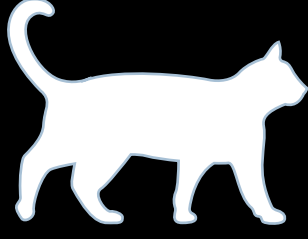
Teacher
(already trained)

Predicted labels

is cat = 0.89

is dog = 0.11



Cat 

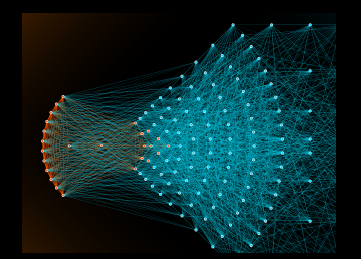
True labels

is cat = 1
is dog = 0

**Teacher
(already trained)**

Predicted labels

is cat = 0.89
is dog = 0.11





Predicted labels

is cat = 0.46

is dog = 0.54

True labels

is cat = 0

is dog = 1

Soft labels contain information!!



Predicted labels

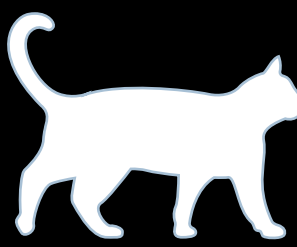
is cat = 0.03

is dog = 0.97

True labels

is cat = 0

is dog = 1

Cat 

True labels

is cat = 1

is dog = 0

Predicted labels

is cat = 0.89

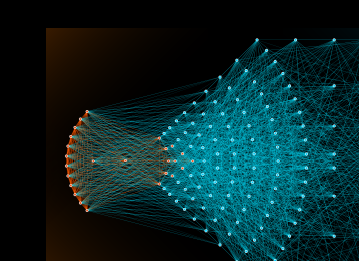
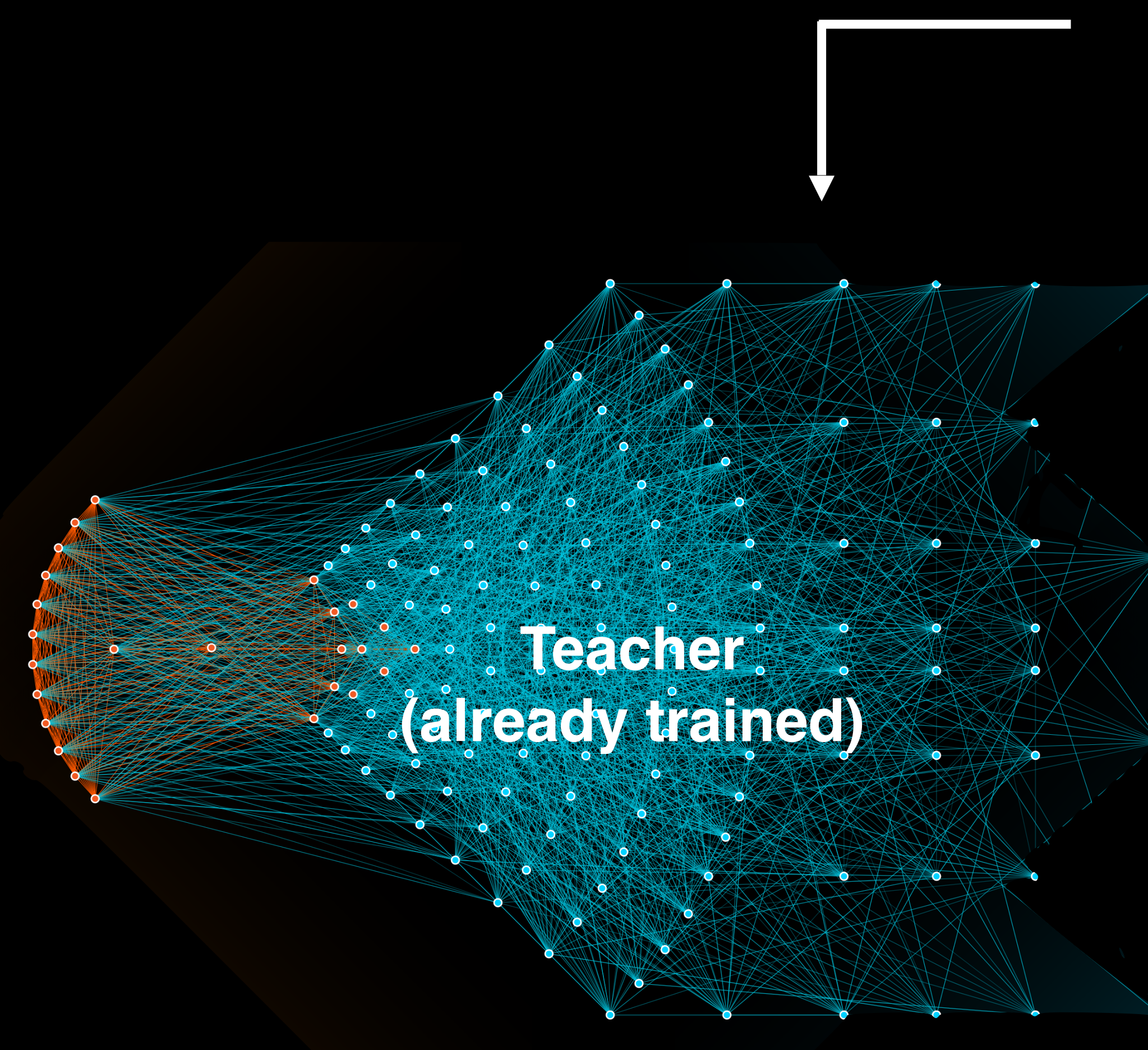
is dog = 0.11

Distilled knowledge

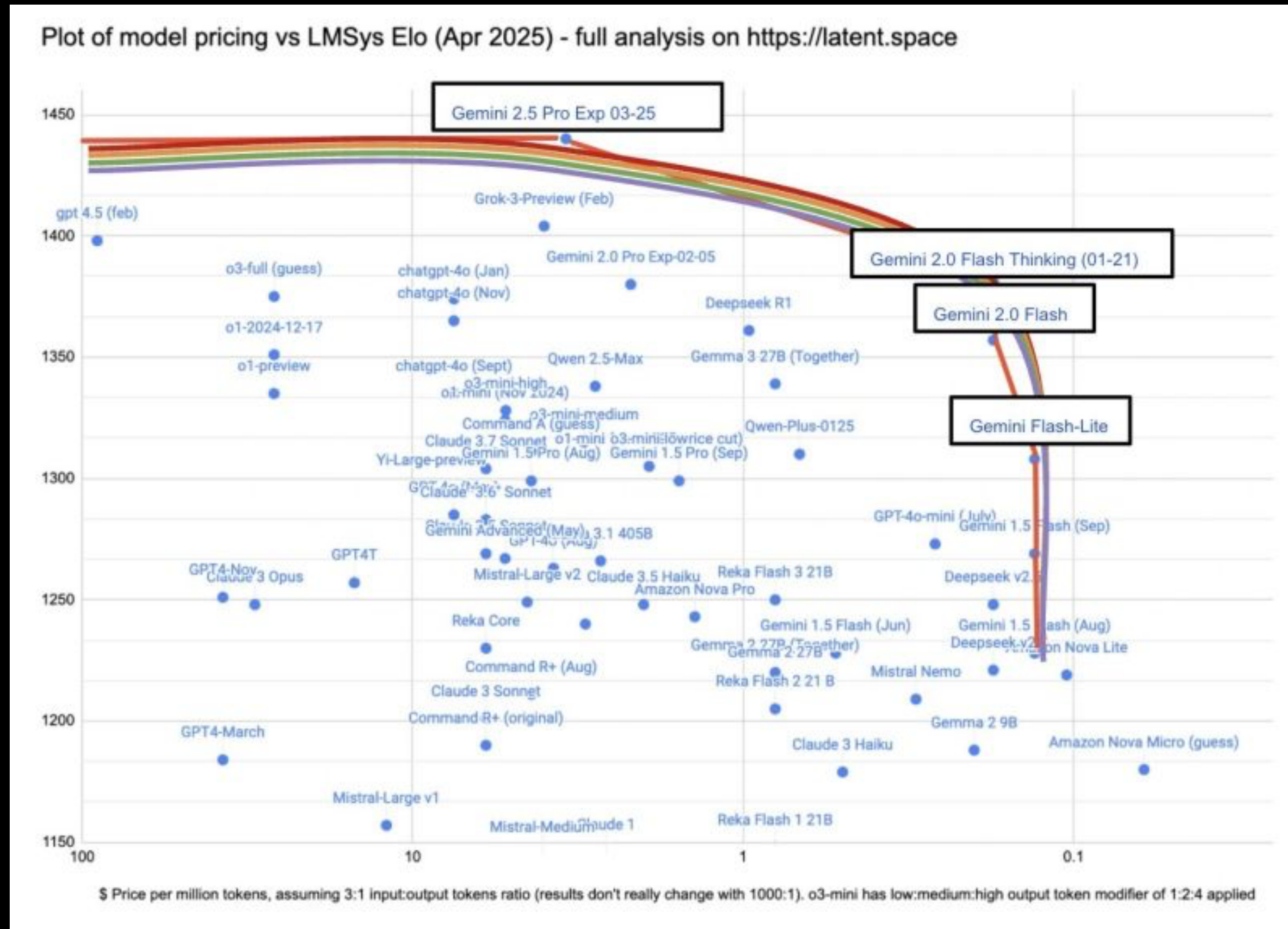
Train student to learn both true and predicted (teacher) labels!

$$L_{total} = \beta \times L_{Distillation} + \alpha \times L_{student}$$

Student learns subtle learned features from teacher!



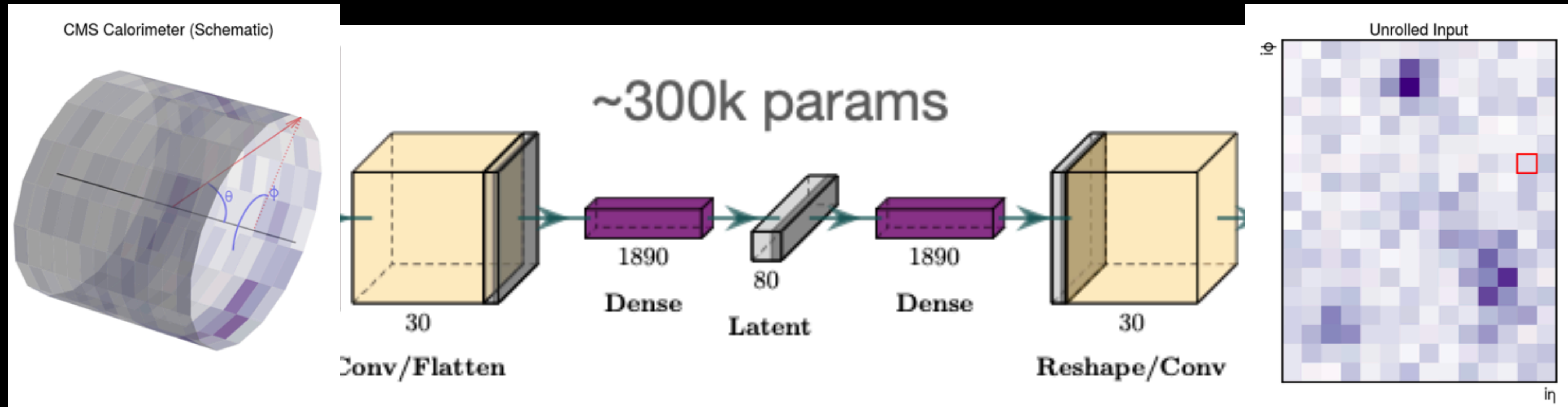
Used for Gemini Flash



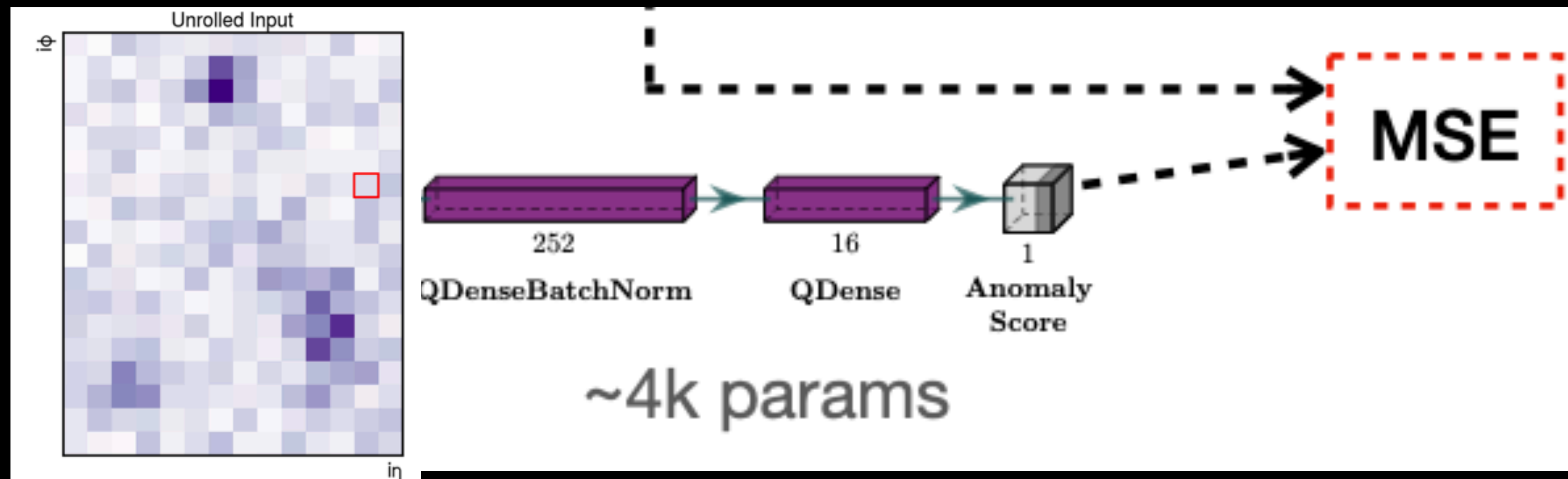
Gemini 3 Flash

www.latent.space/,
Jeff Dean HLF

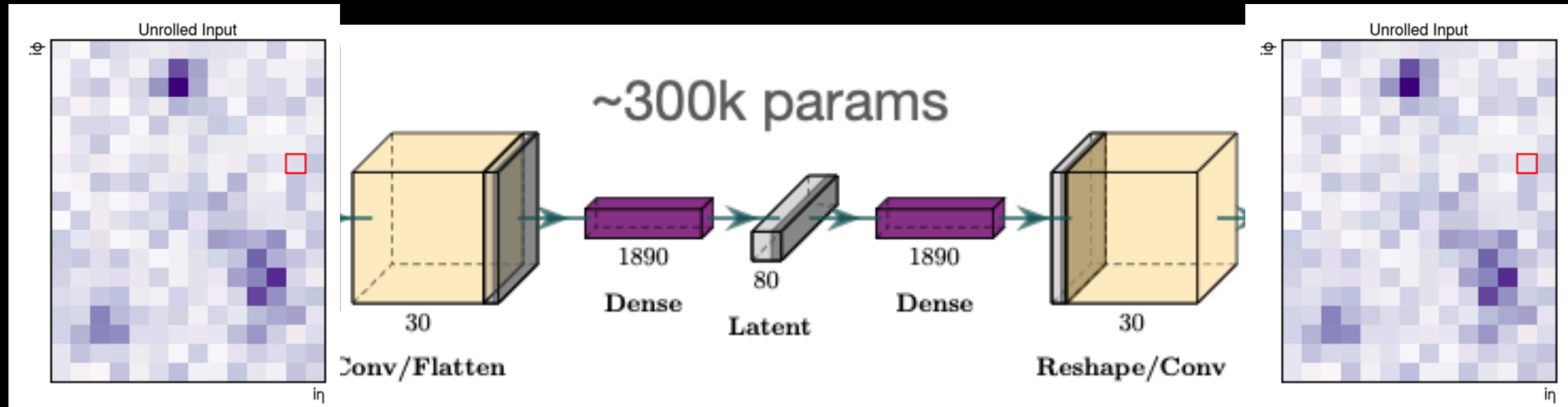
Used in CMS trigger!



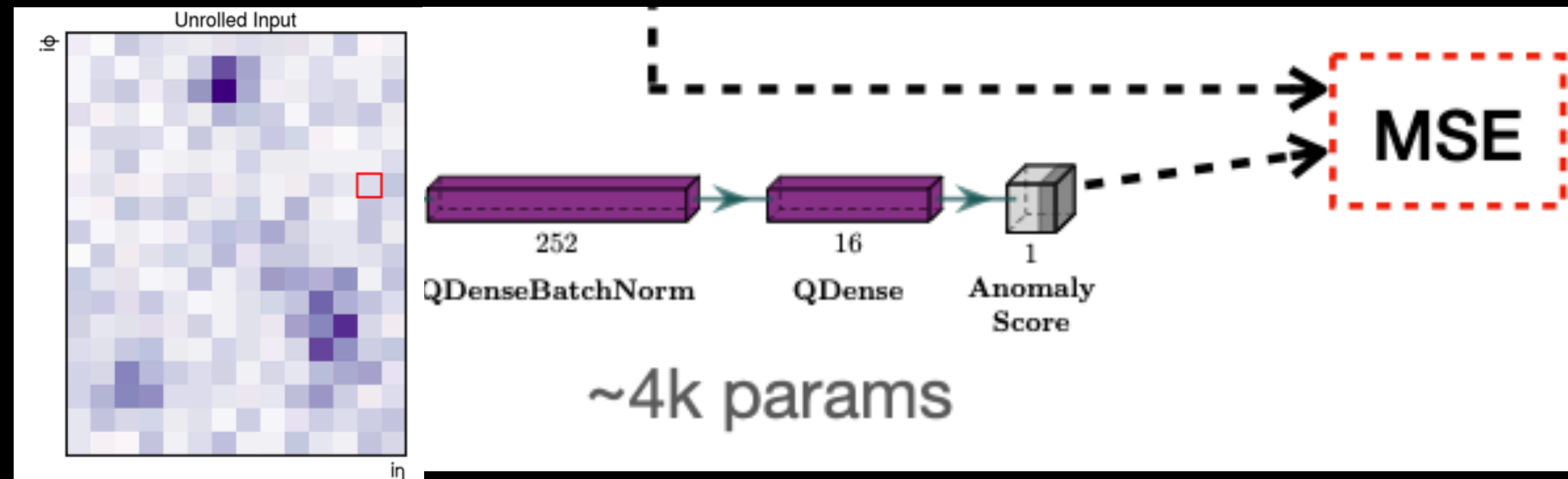
Knowledge distillation: input to anomaly score

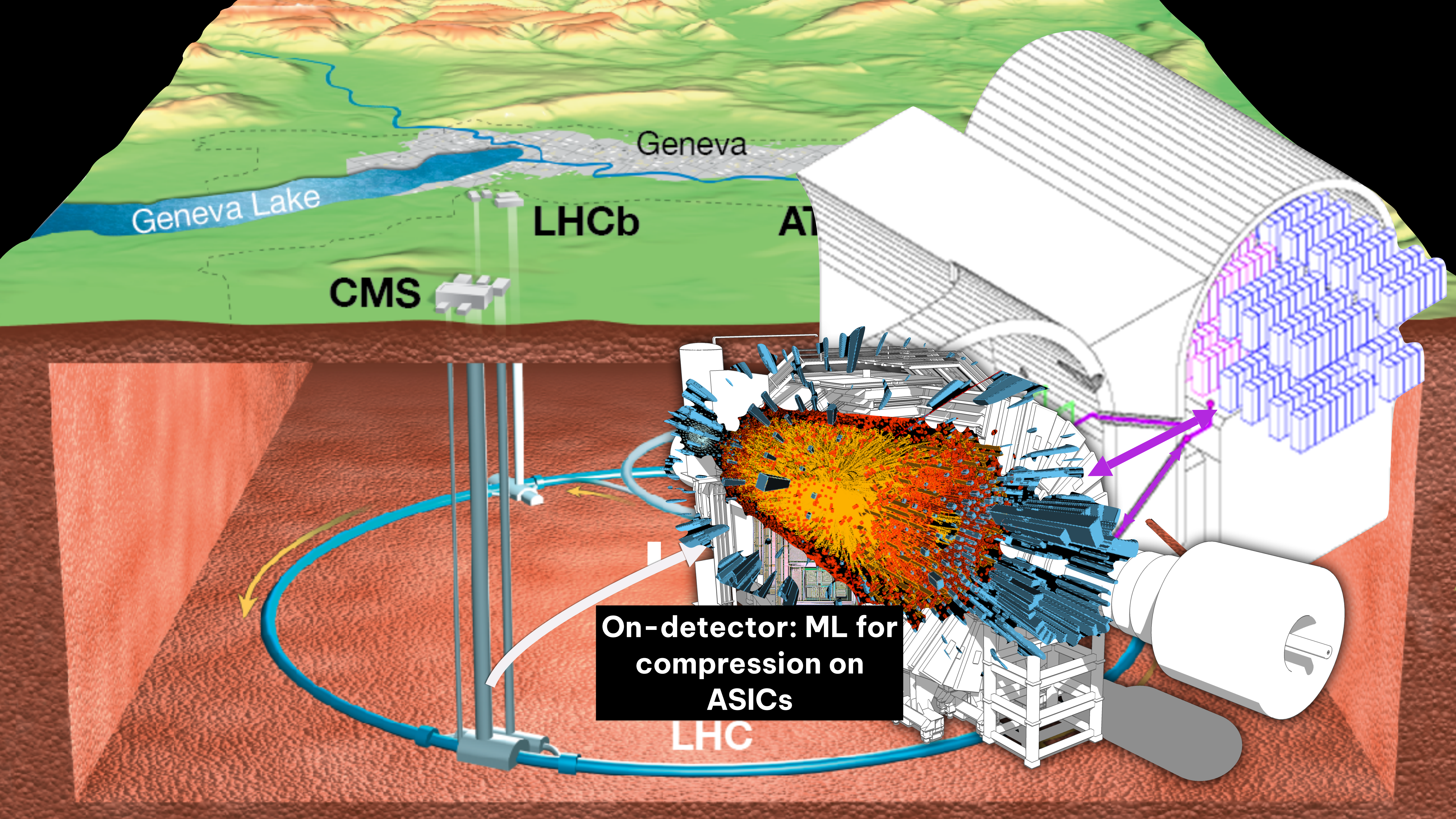


Used in CMS trigger!



Knowledge distillation: input to anomaly score





Geneva

Geneva Lake

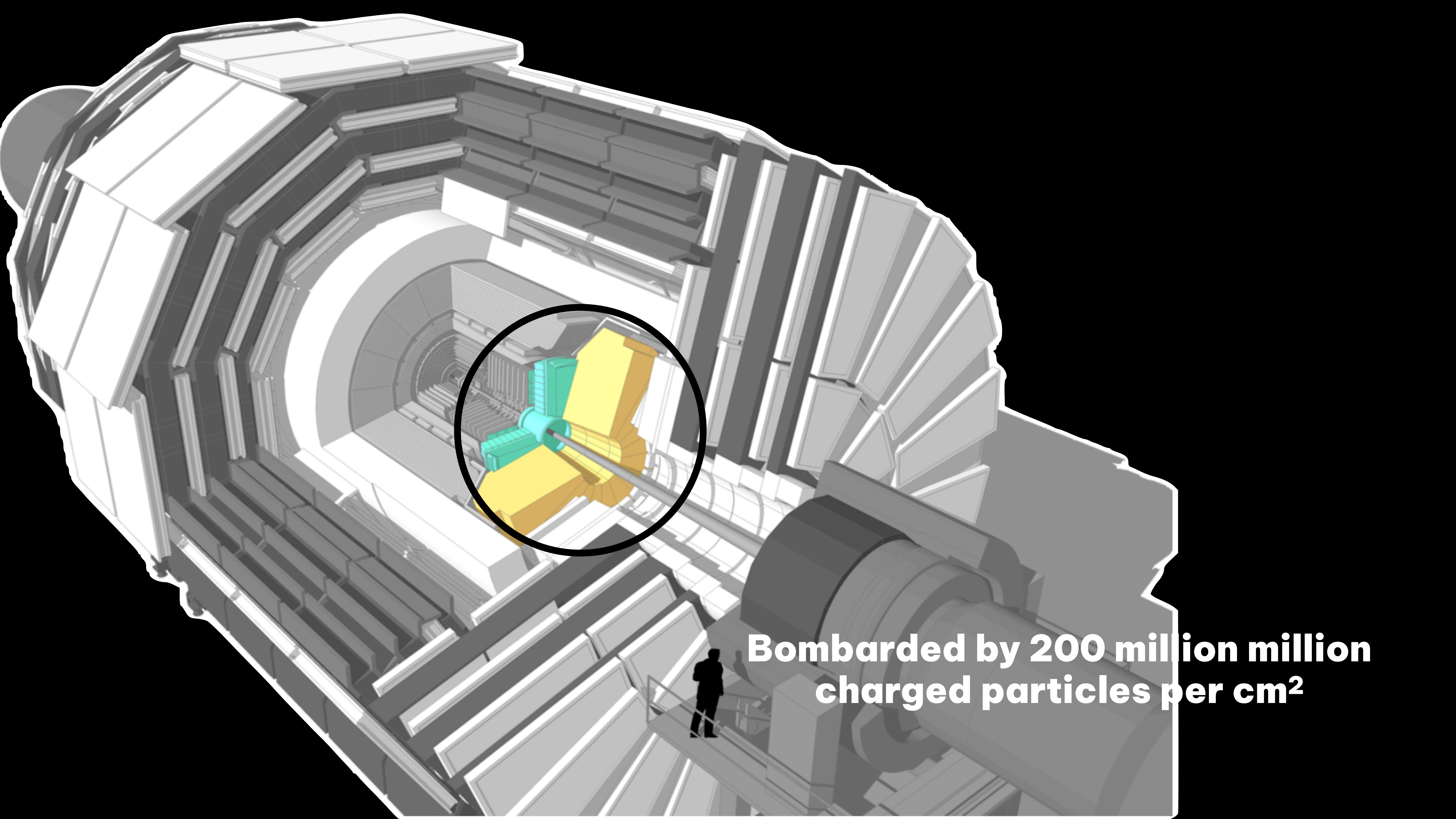
LHCb

ATLAS

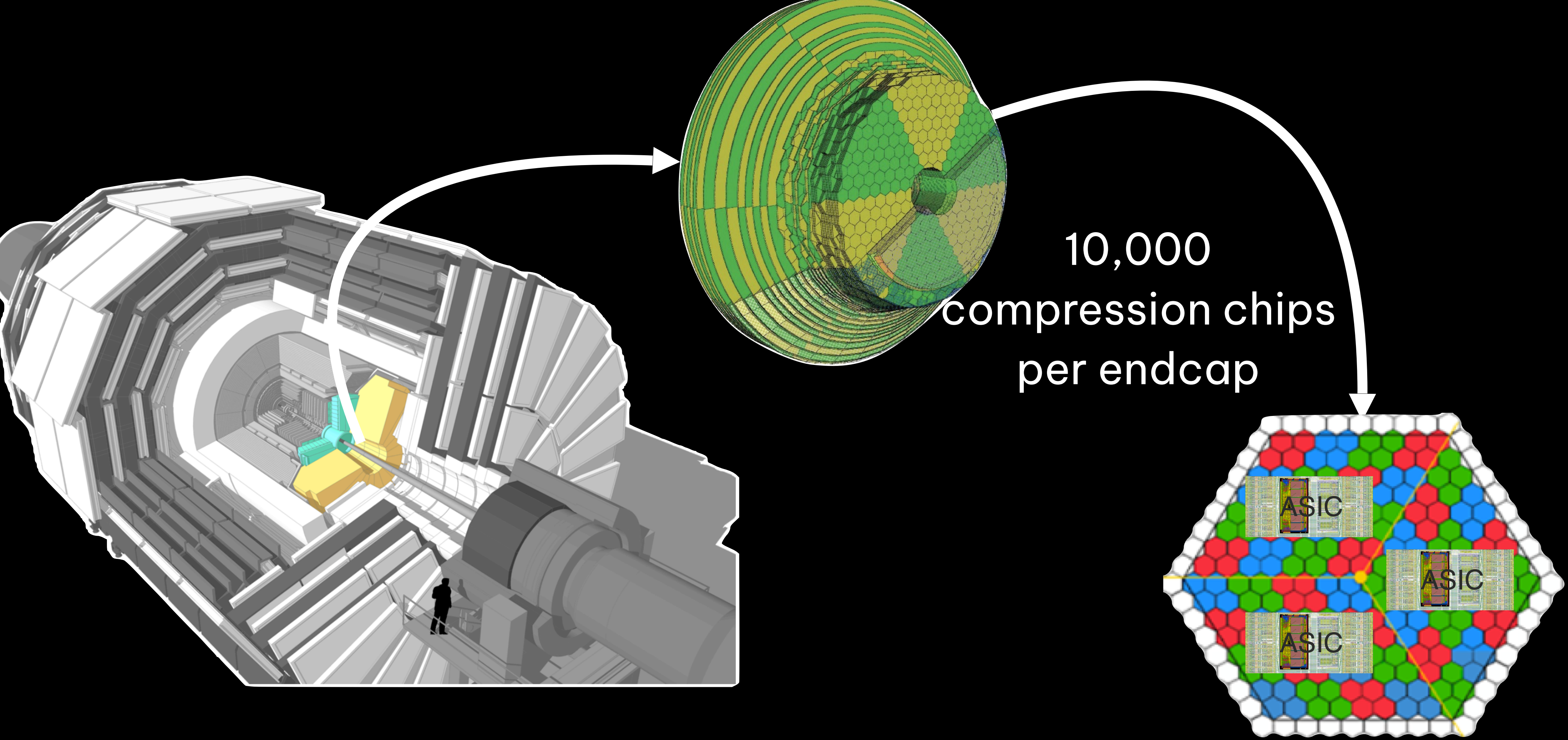
CMS

On-detector: ML for compression on ASICs

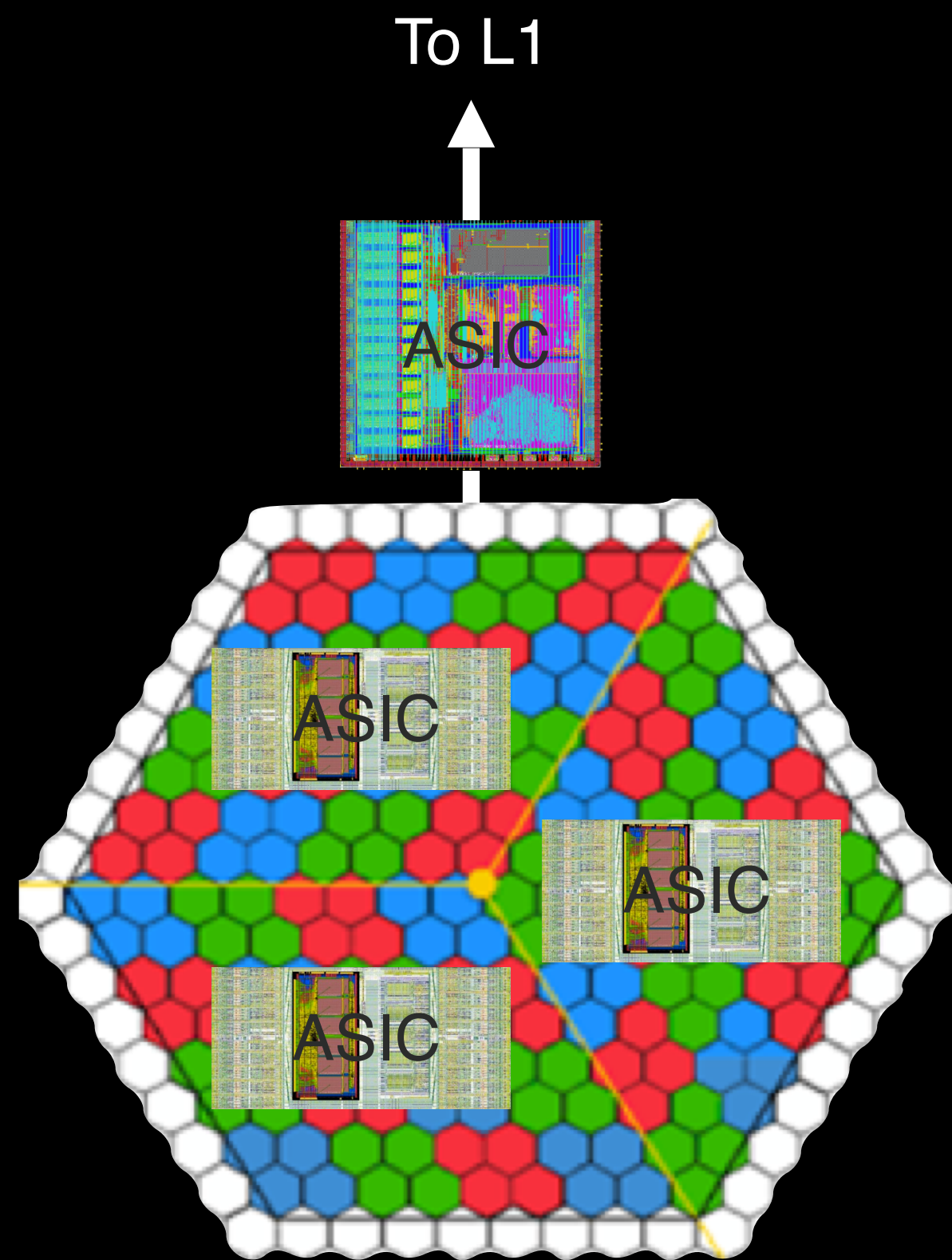
LHC



**Bombarded by 200 million million
charged particles per cm²**

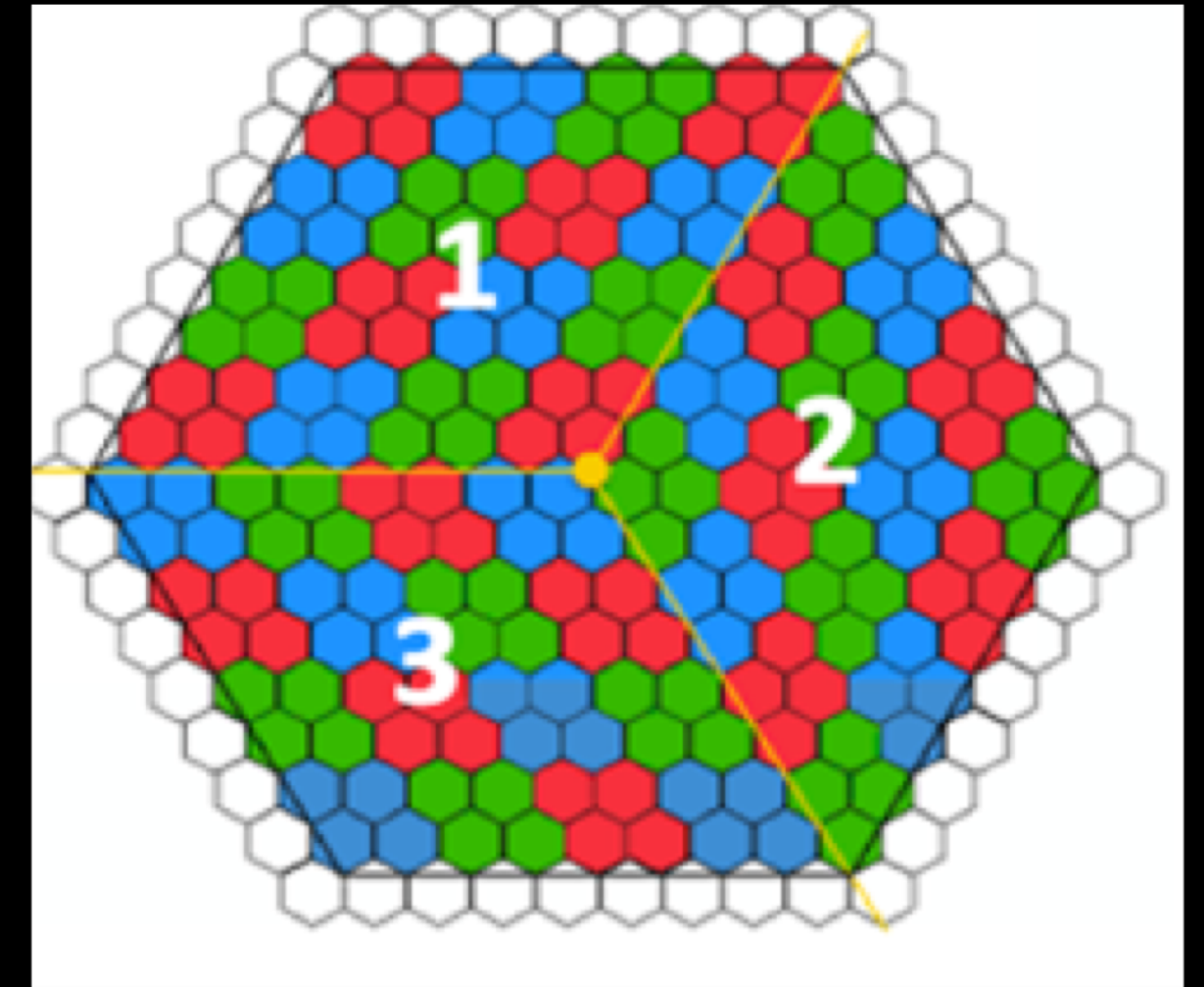
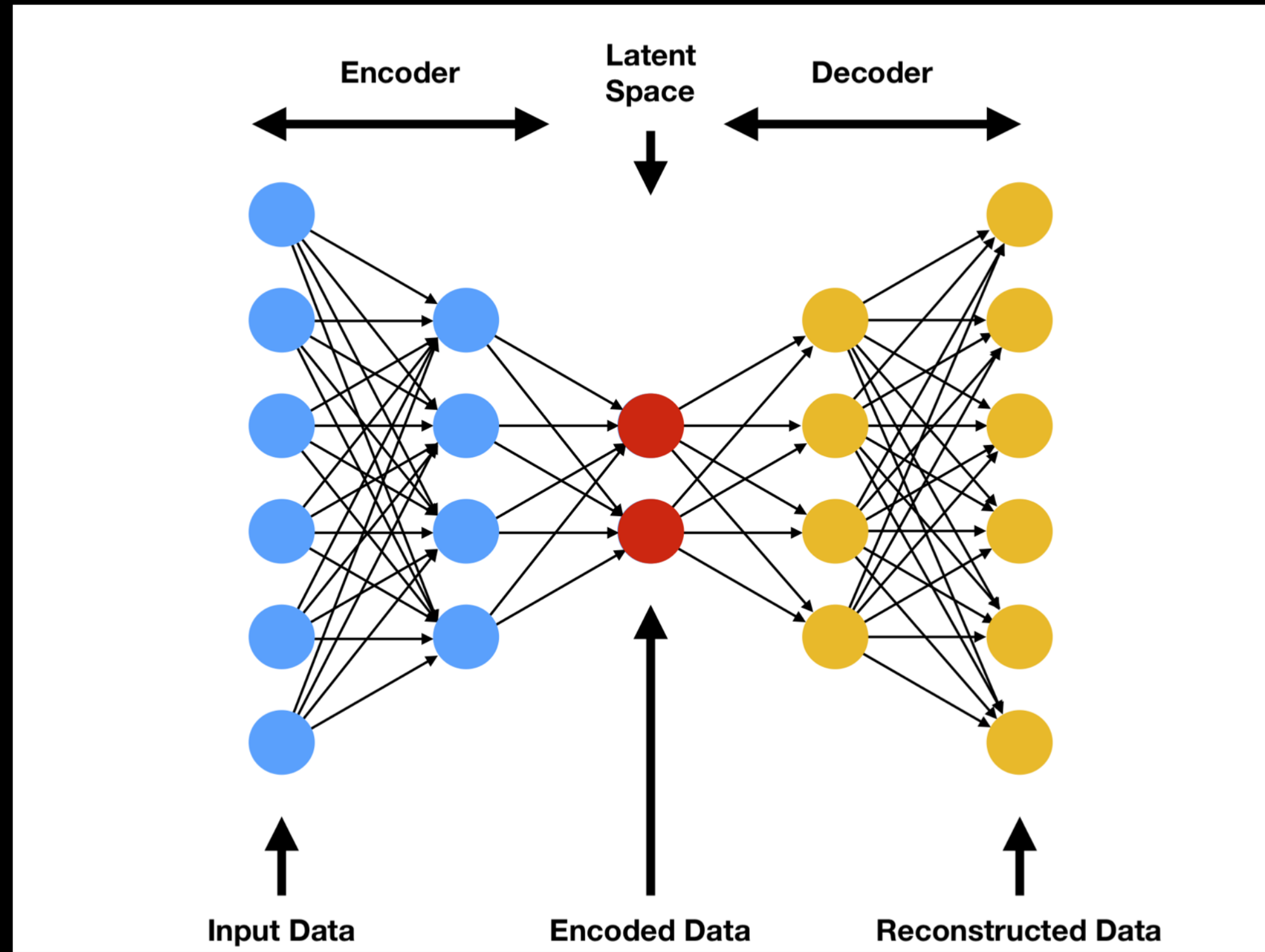
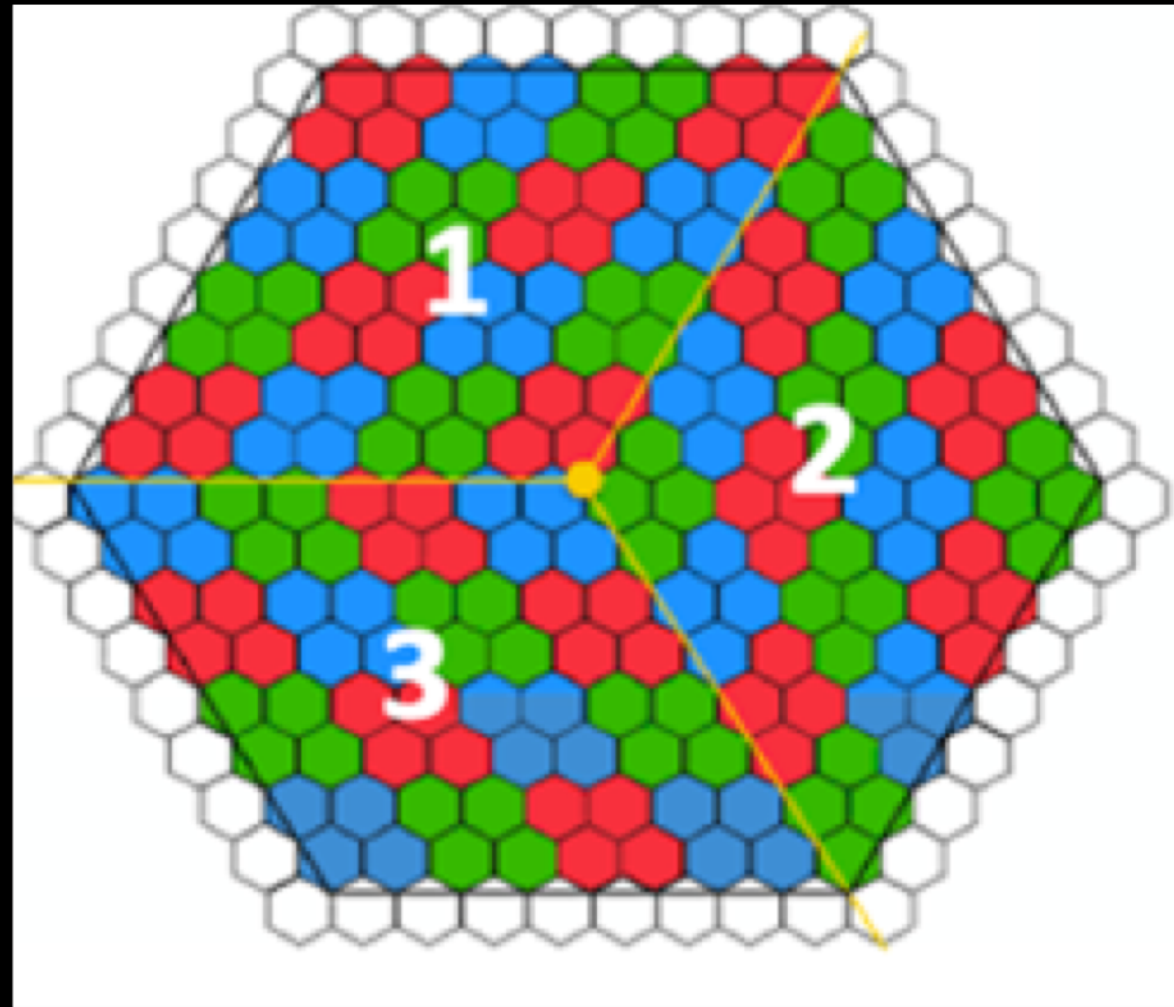


Total 6.5 million readout-channels.
Can not read out all at 40 MHz, need to compress

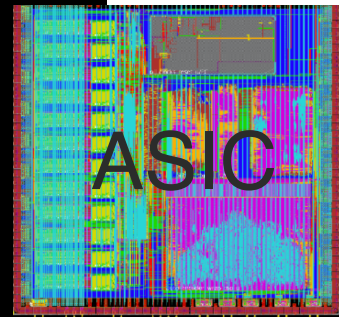
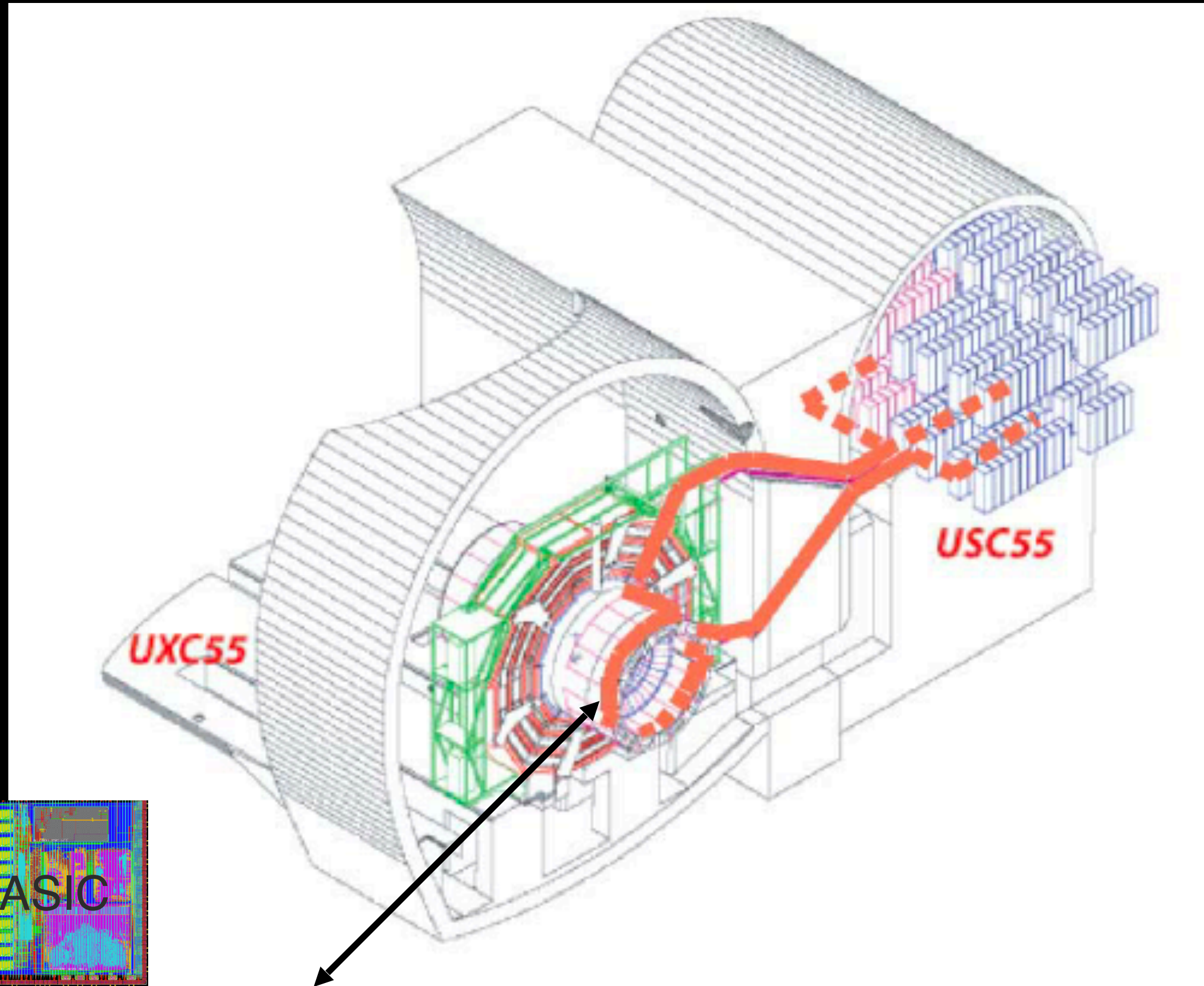


Must compress ON DETECTOR

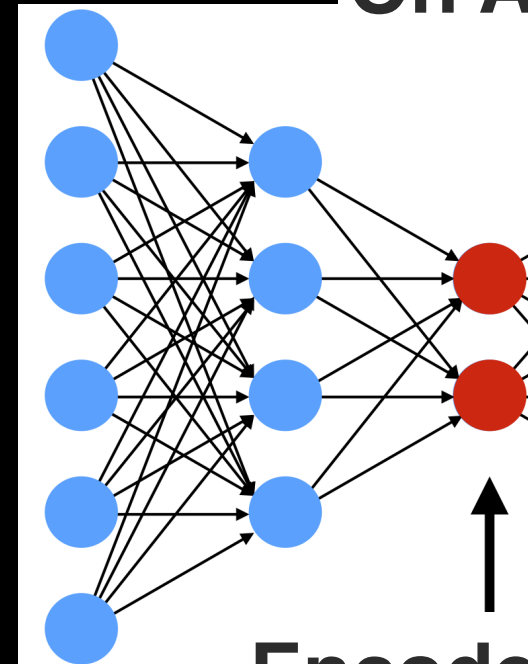
- High radiation
- Cooled to -30 degrees
- 400 ns latency budget



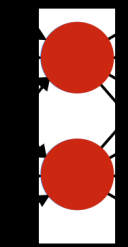
Variational Autoencoder



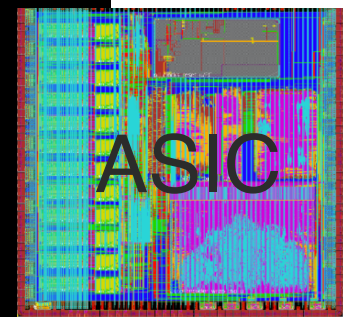
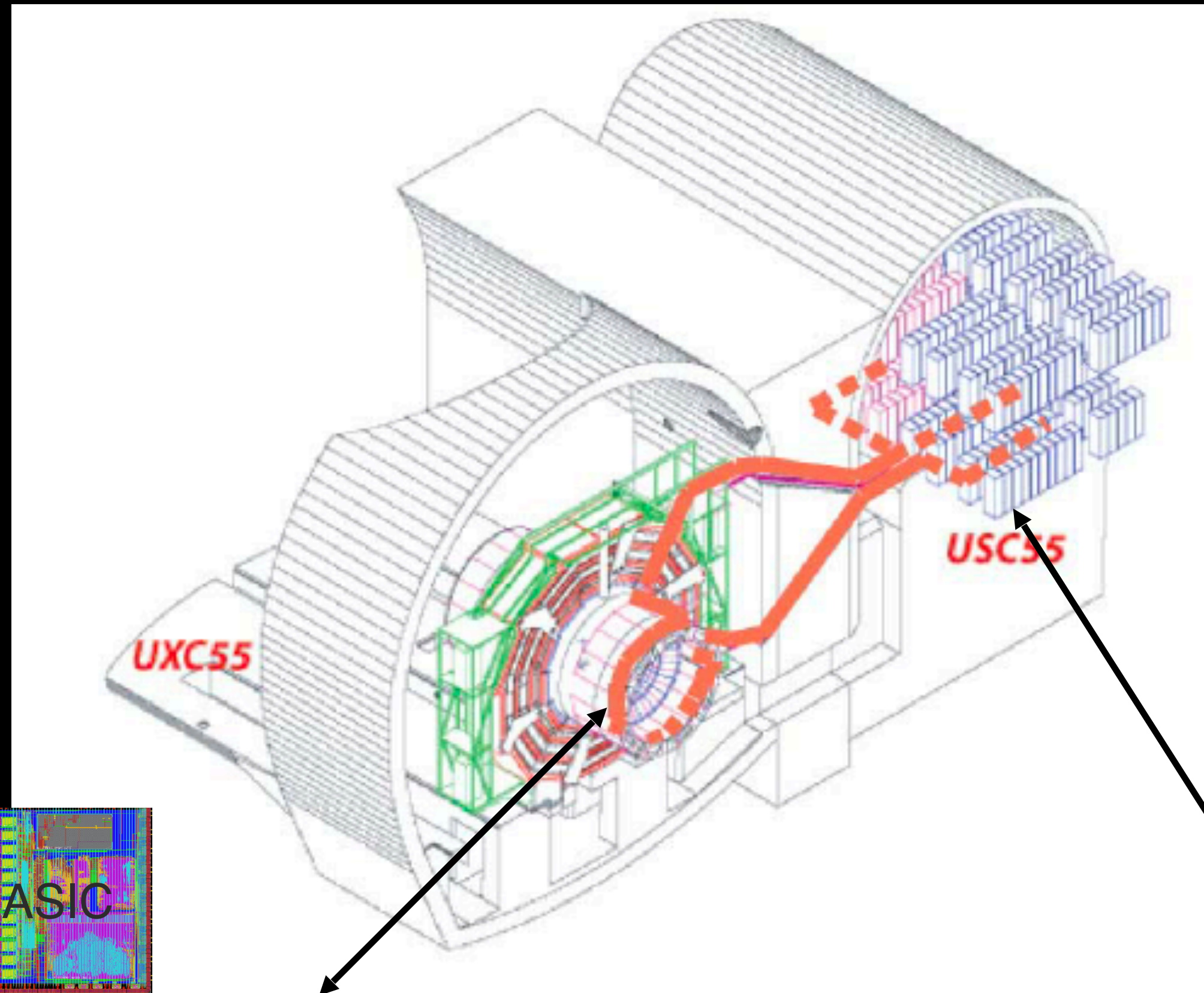
On ASIC



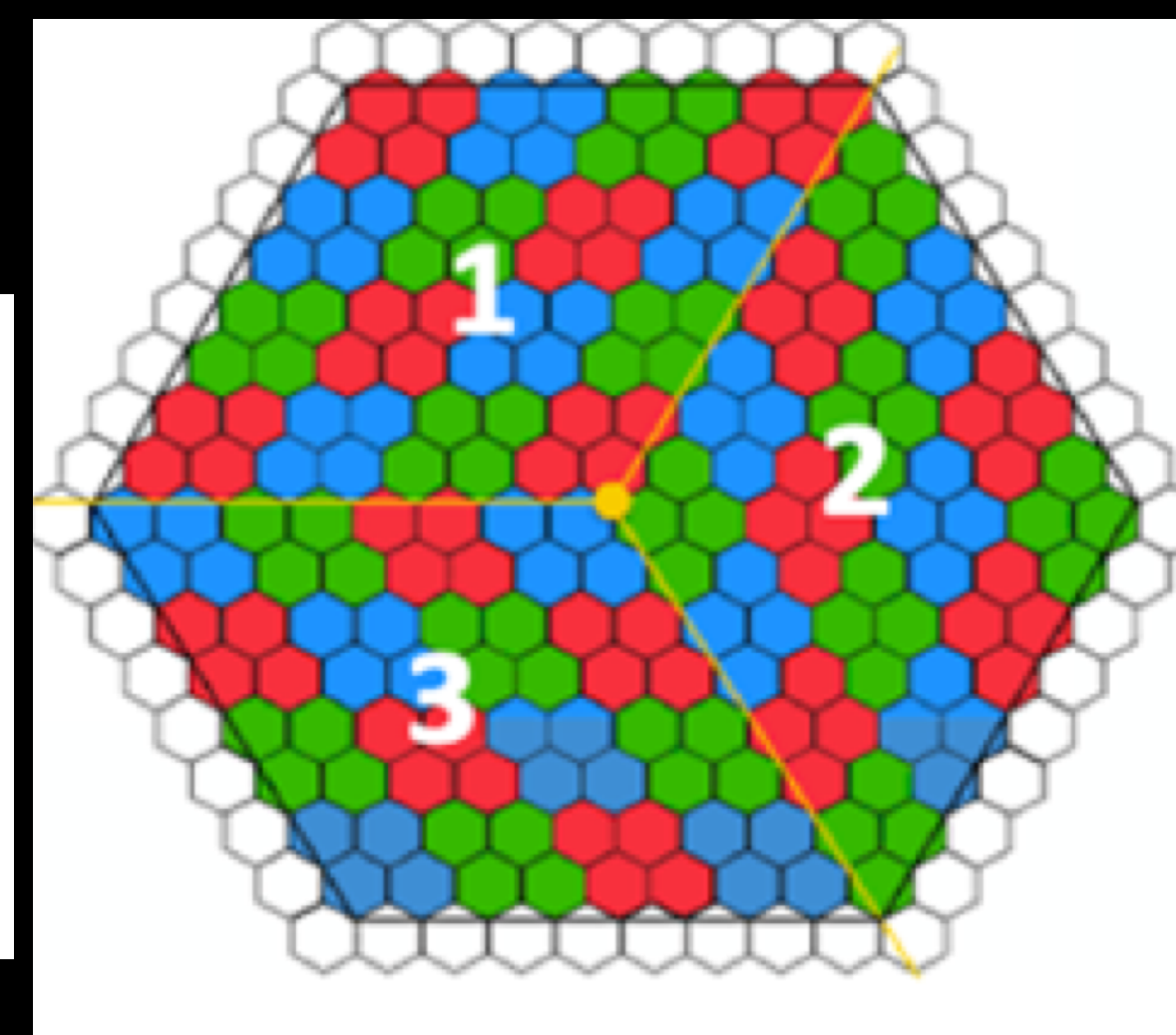
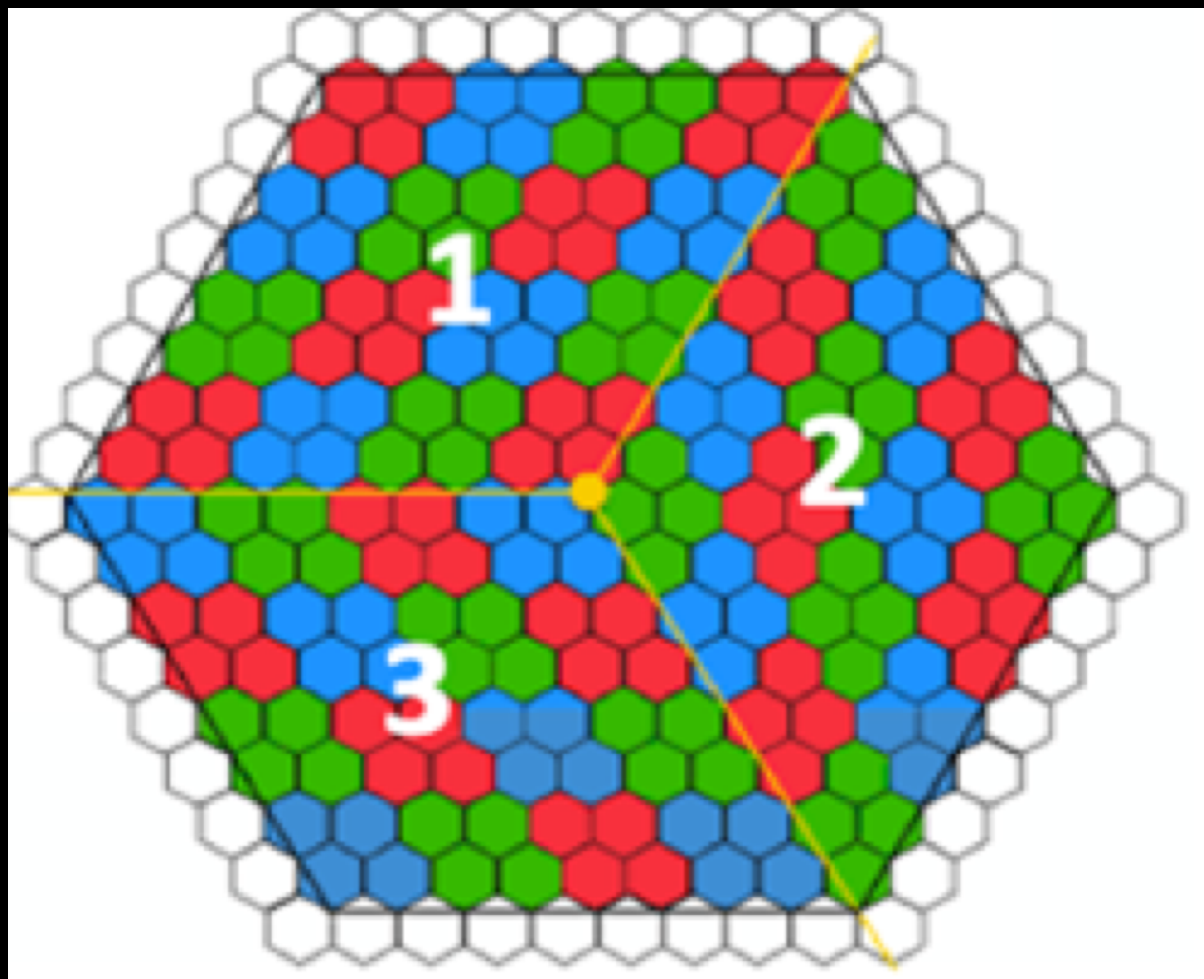
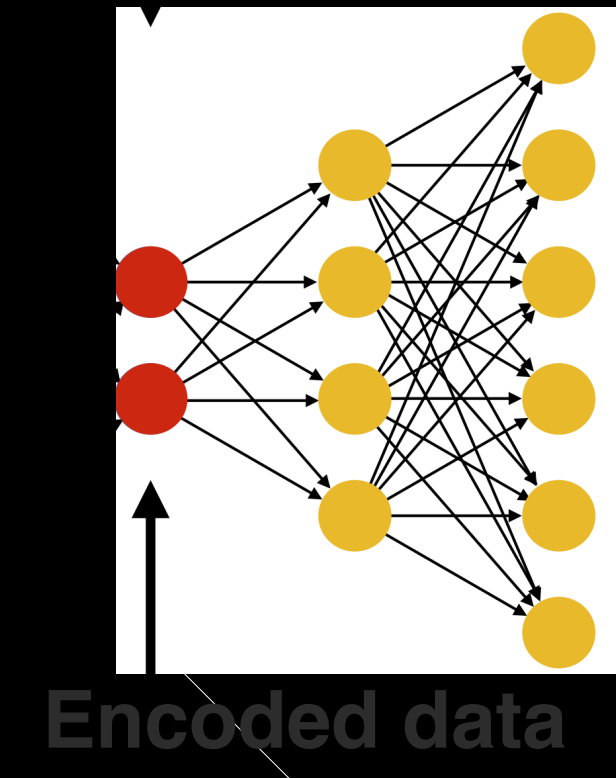
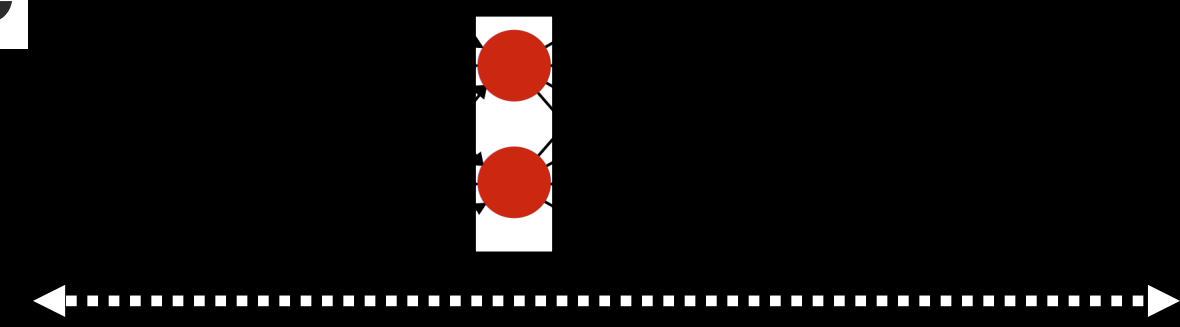
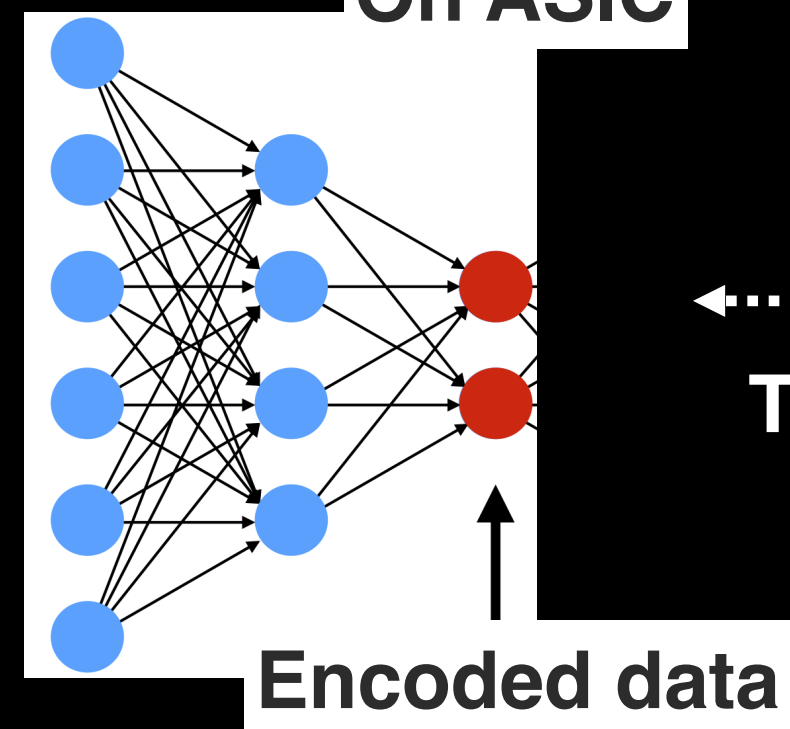
Encoded data



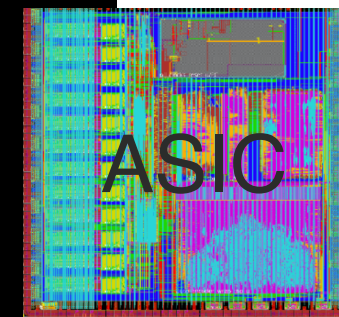
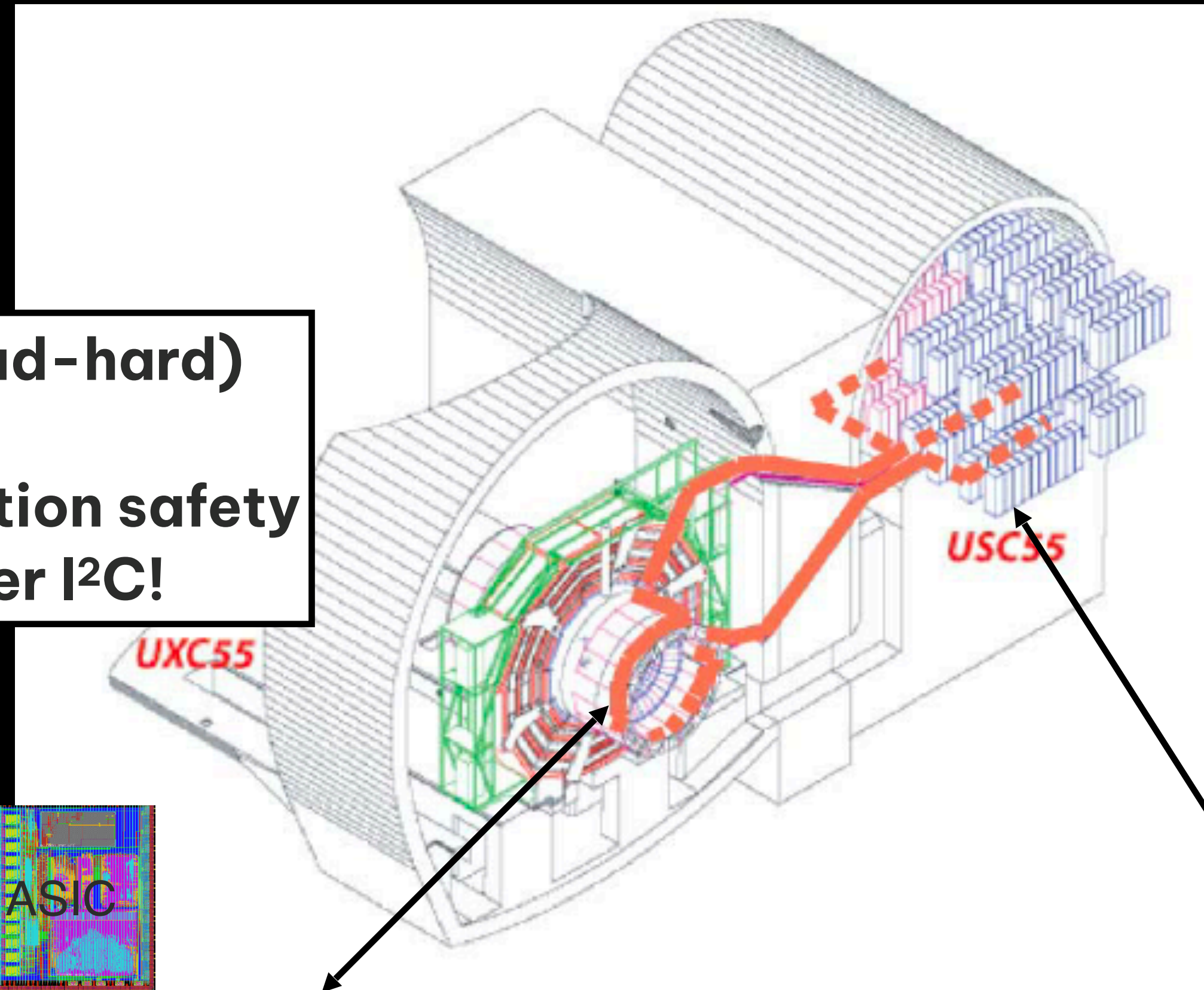
Transmit encoded data!



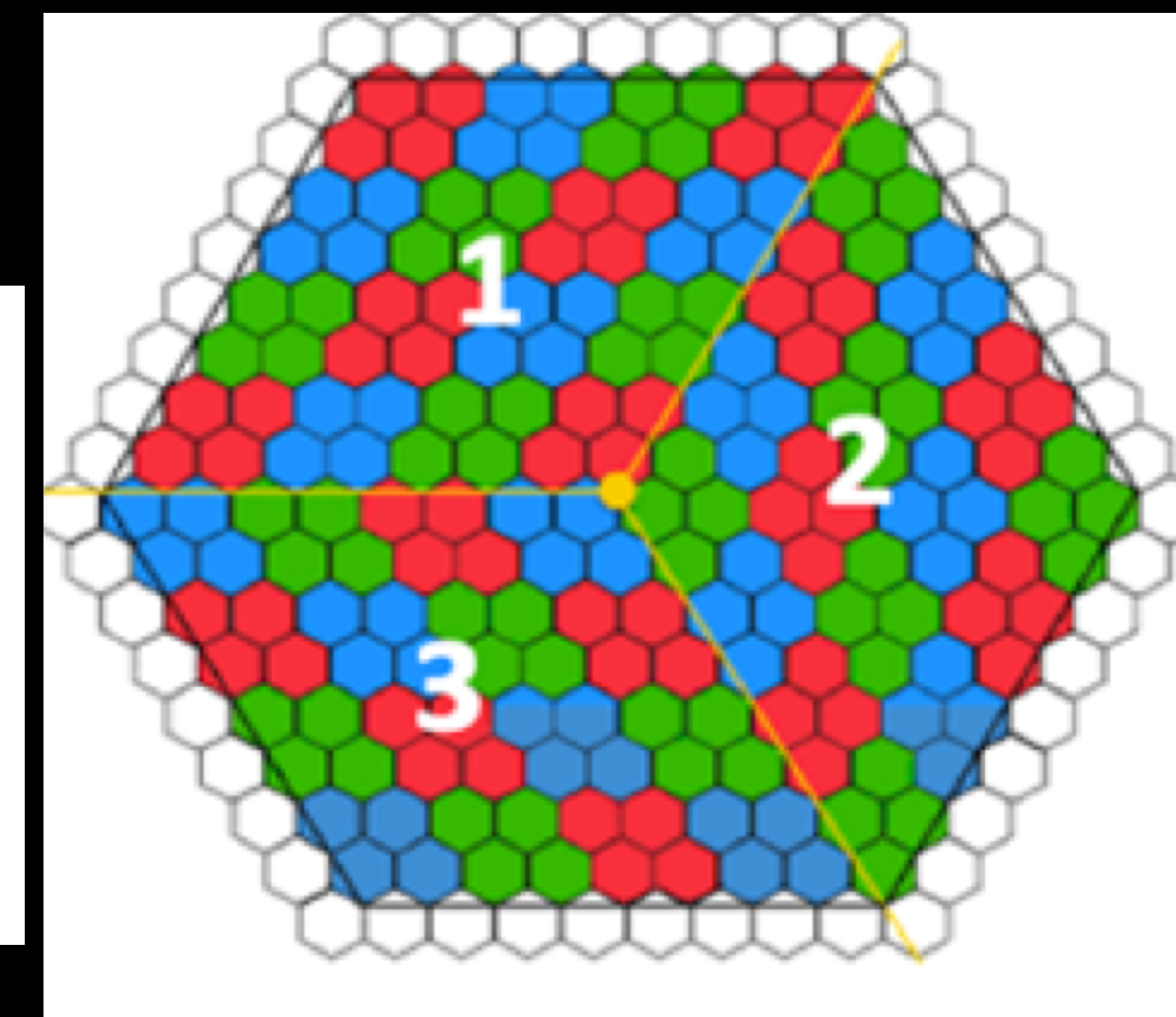
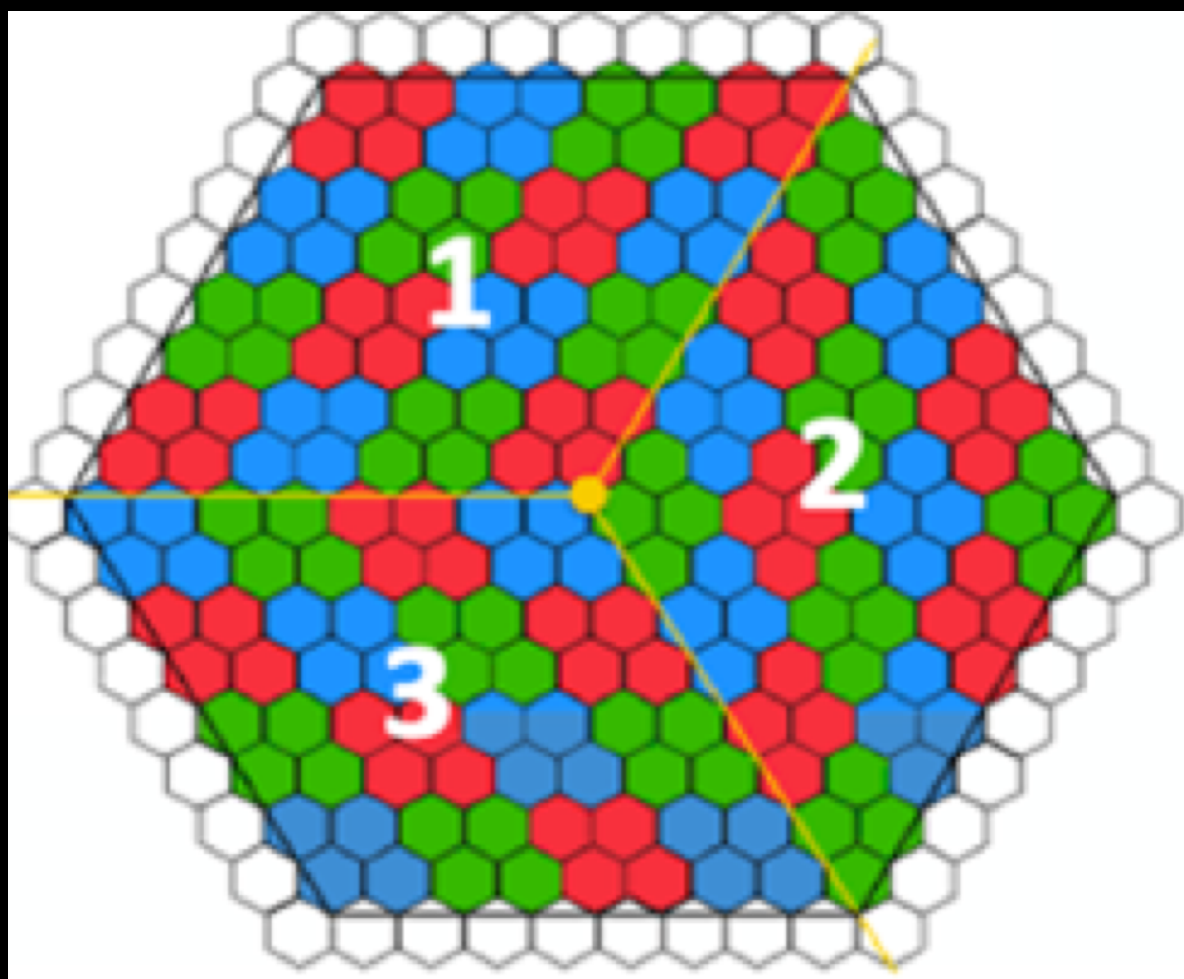
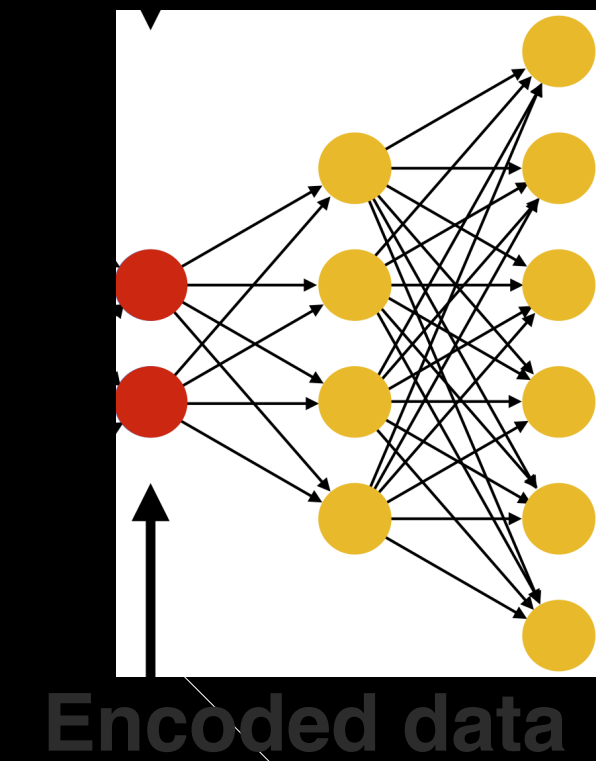
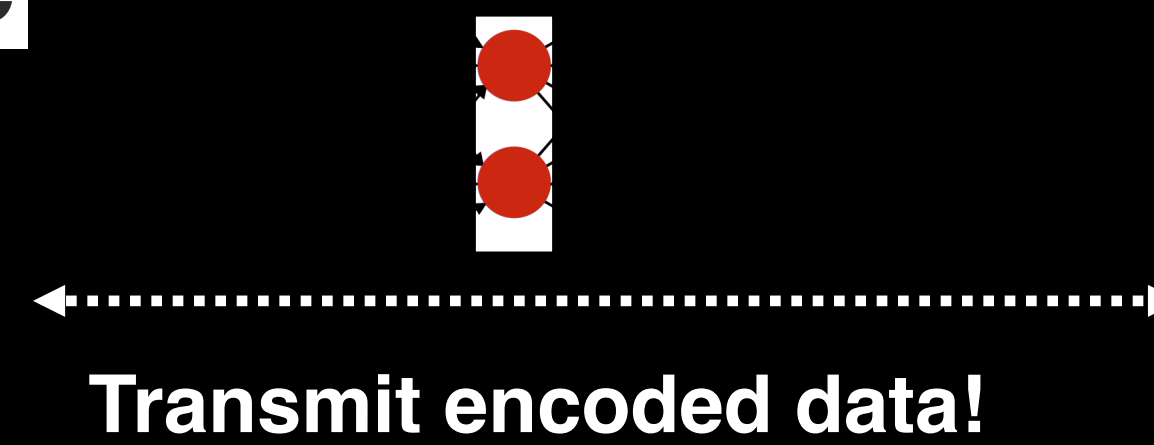
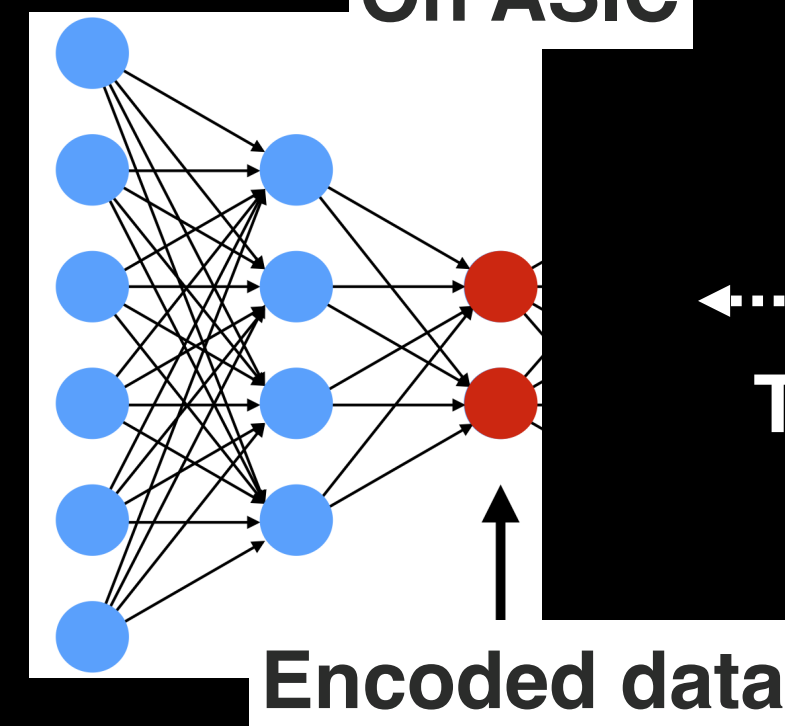
On ASIC

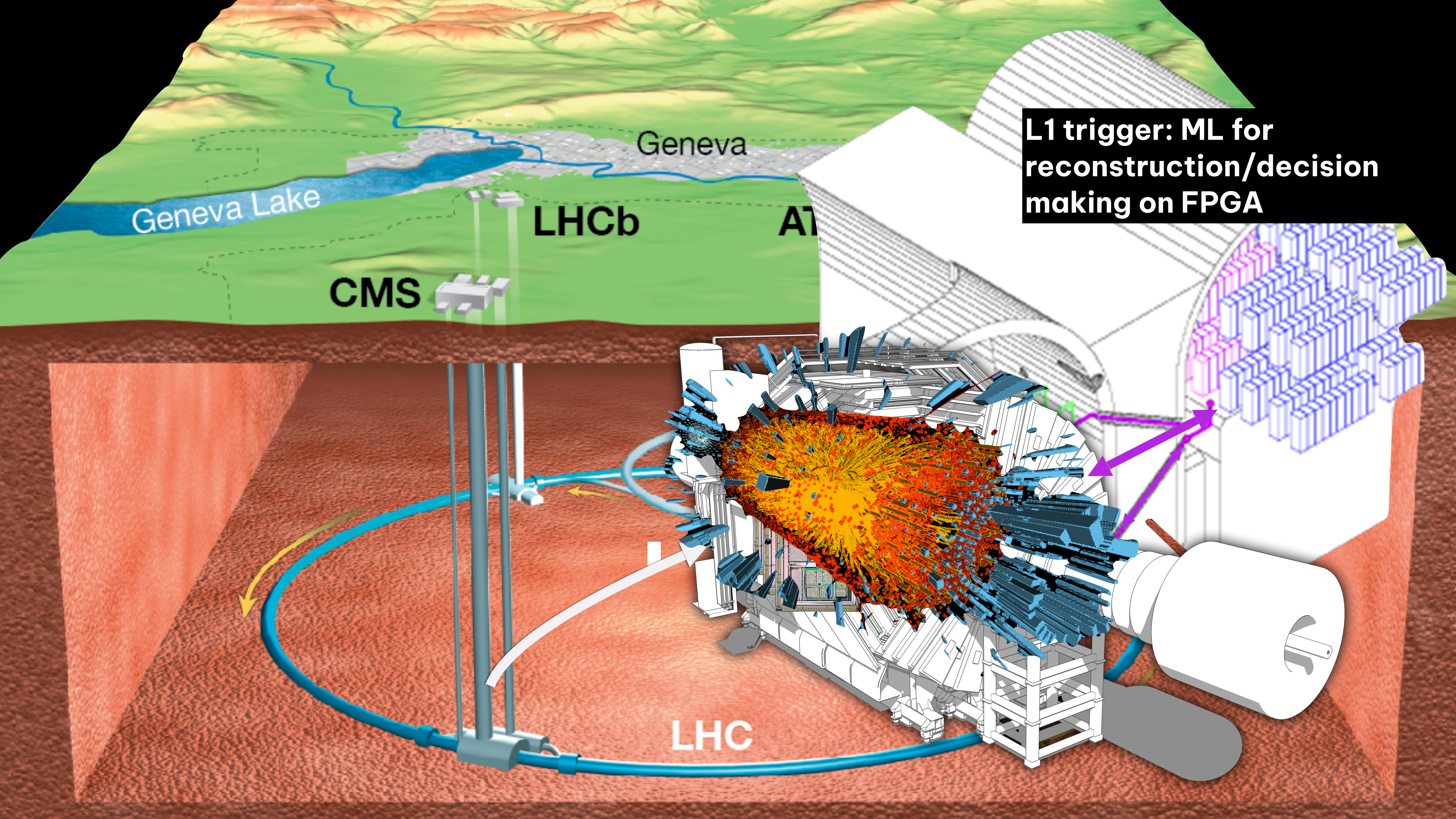


- ASIC inference (65 nm, rad-hard)
- < 4.5 mW/channel
- Triplicated w/b for radiation safety
Reprogrammable w/b over I²C!



On ASIC





Geneva

Geneva Lake

LHCb

ATLAS

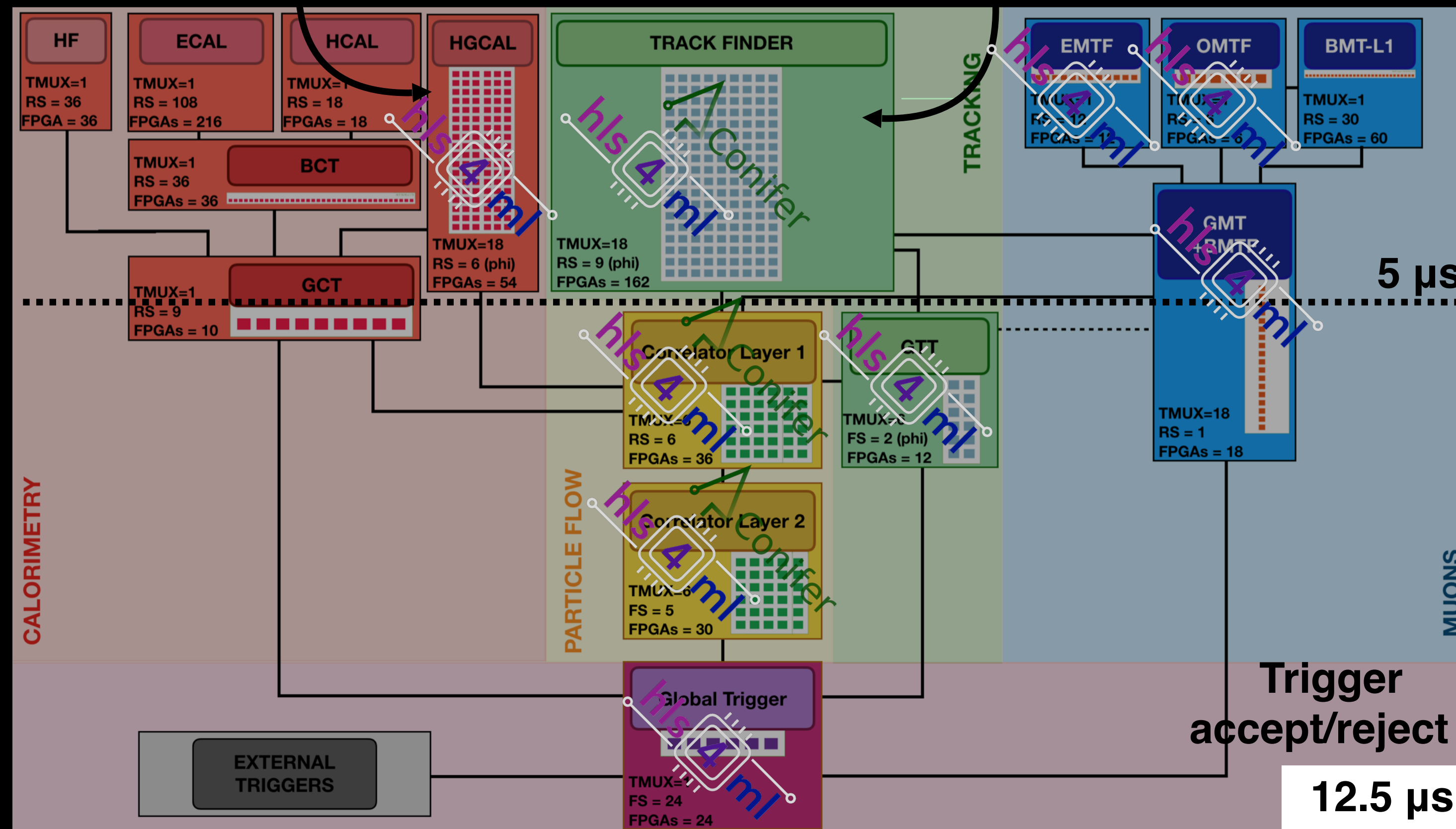
CMS

L1 trigger: ML for reconstruction/decision making on FPGA

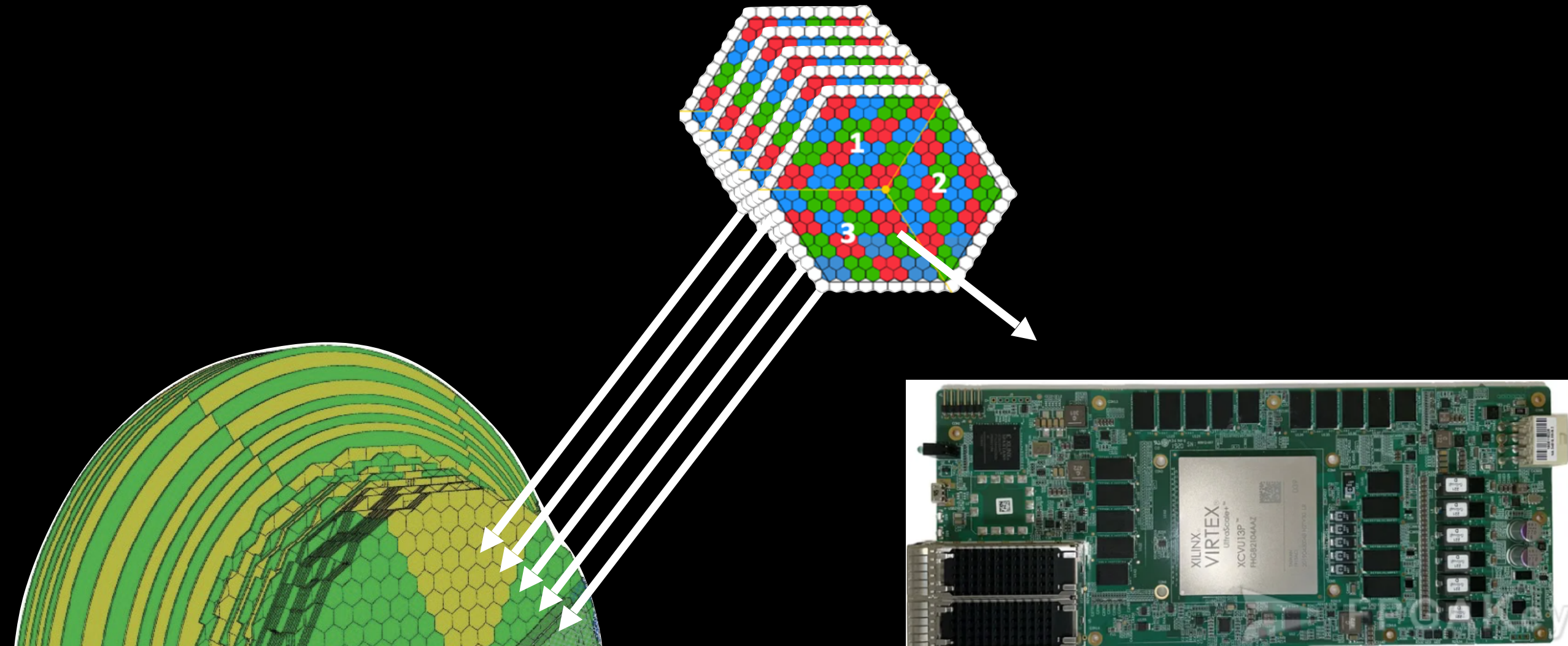
LHC

Nanosecond ML inference on FPGAs!

~40 billion inferences/s during HL-LHC

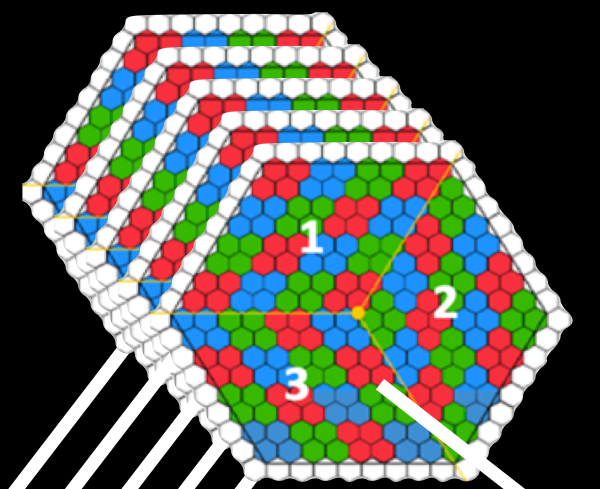
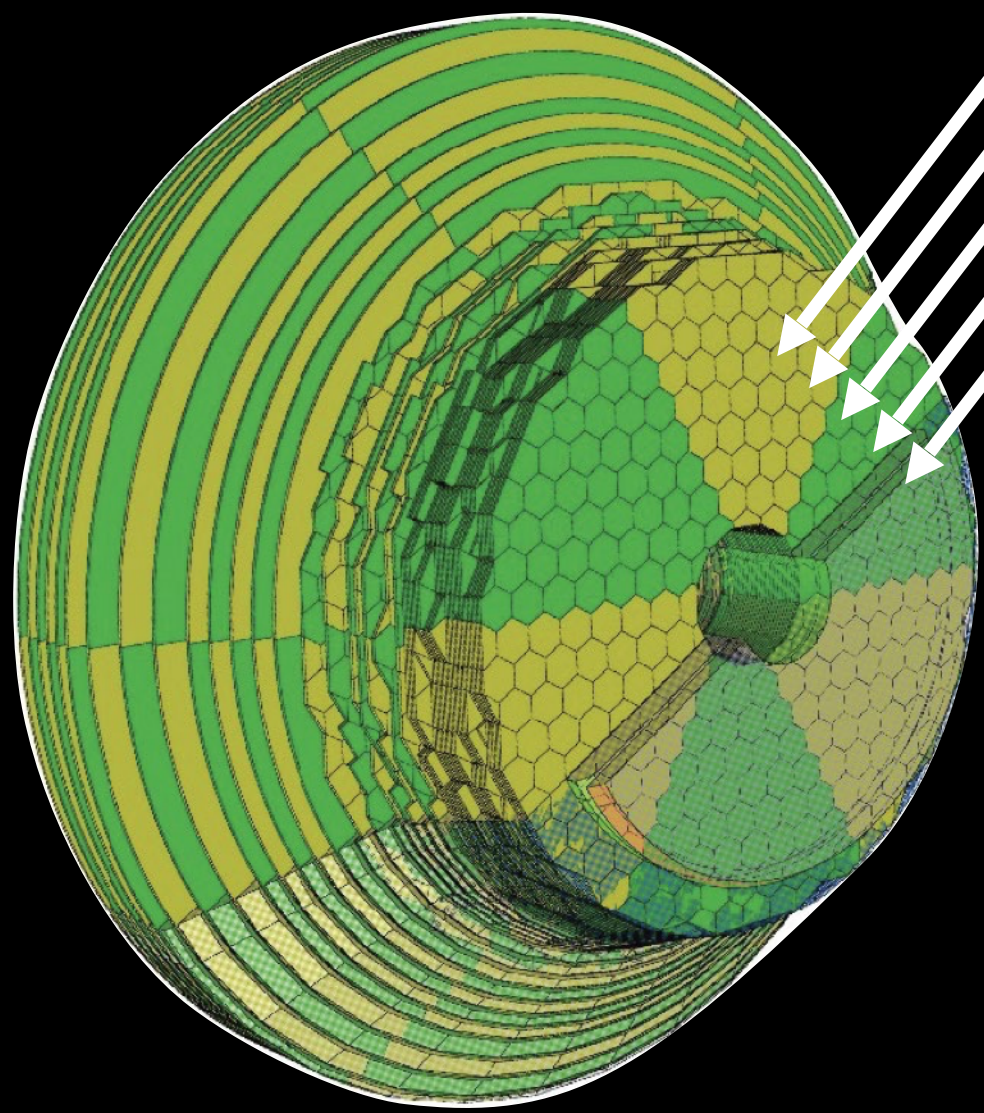


HEP developed libraries for fast ML on FPGAs

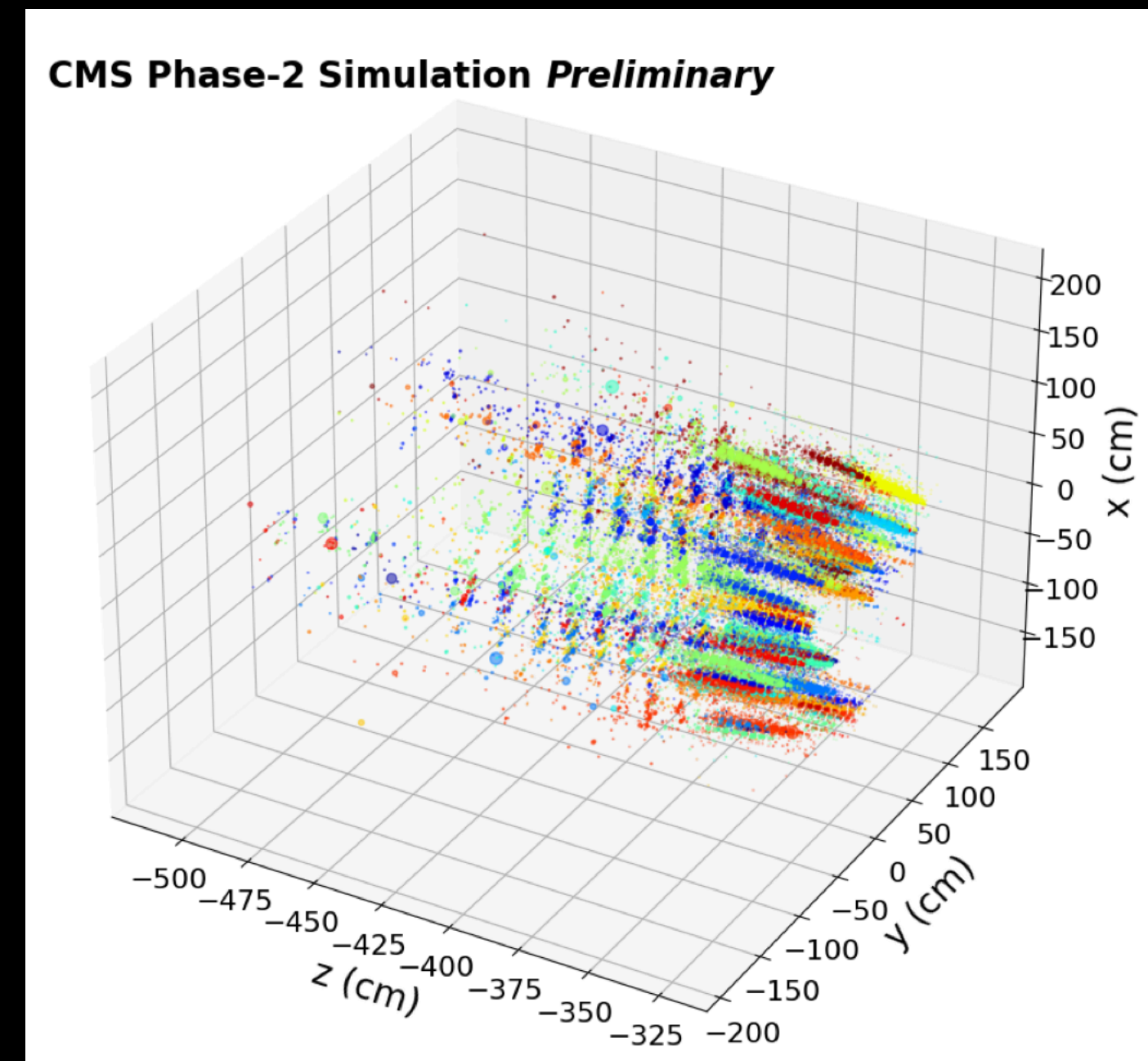


**5,000 “sensor
input” per FPGA
(and 18 of them to
process 18 events
in parallel)**

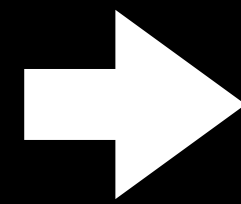
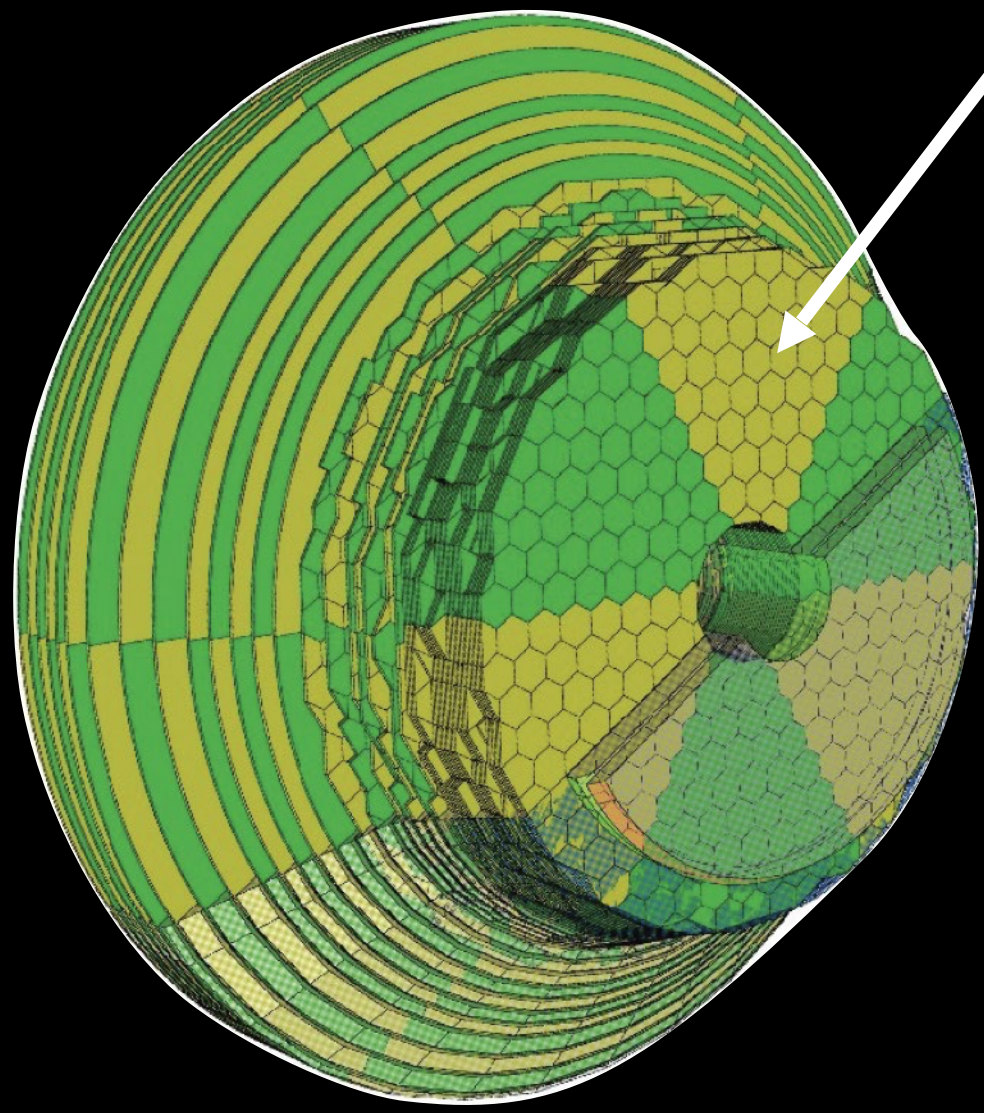
5,000 “sensor input” per FPGA



2 μ s to cluster into distinct particle energy clusters!

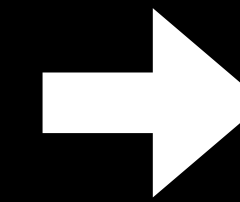
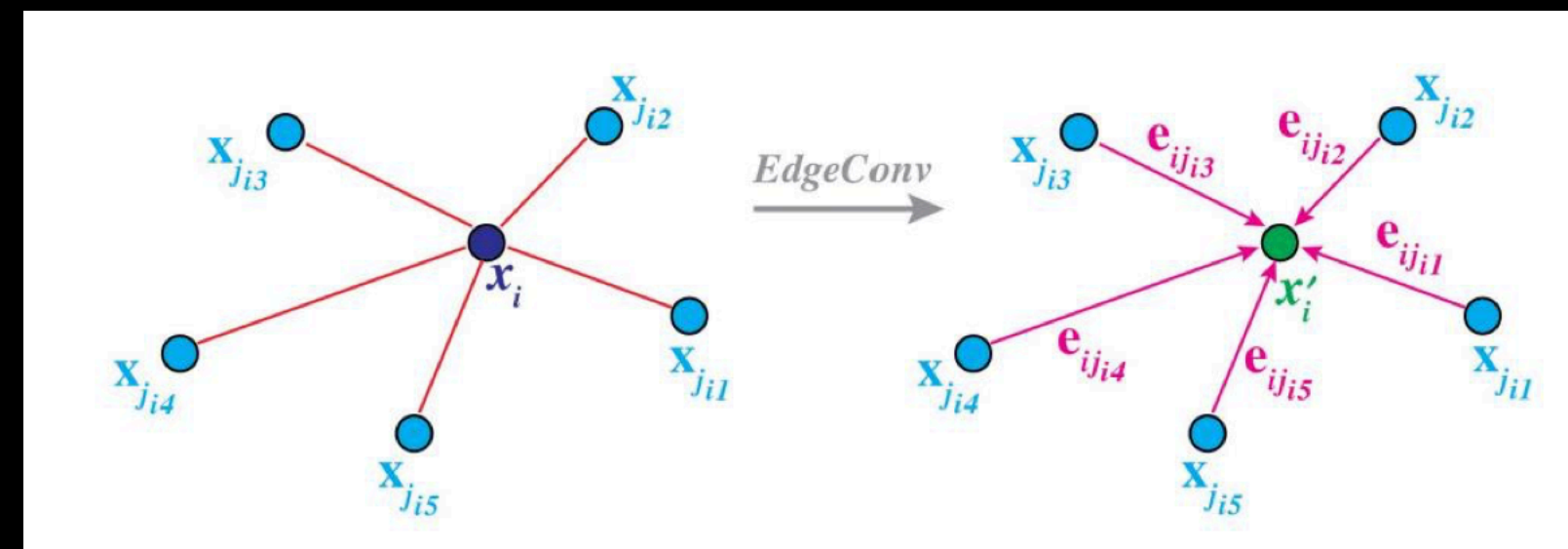


5,000 "vertices"
per FPGA

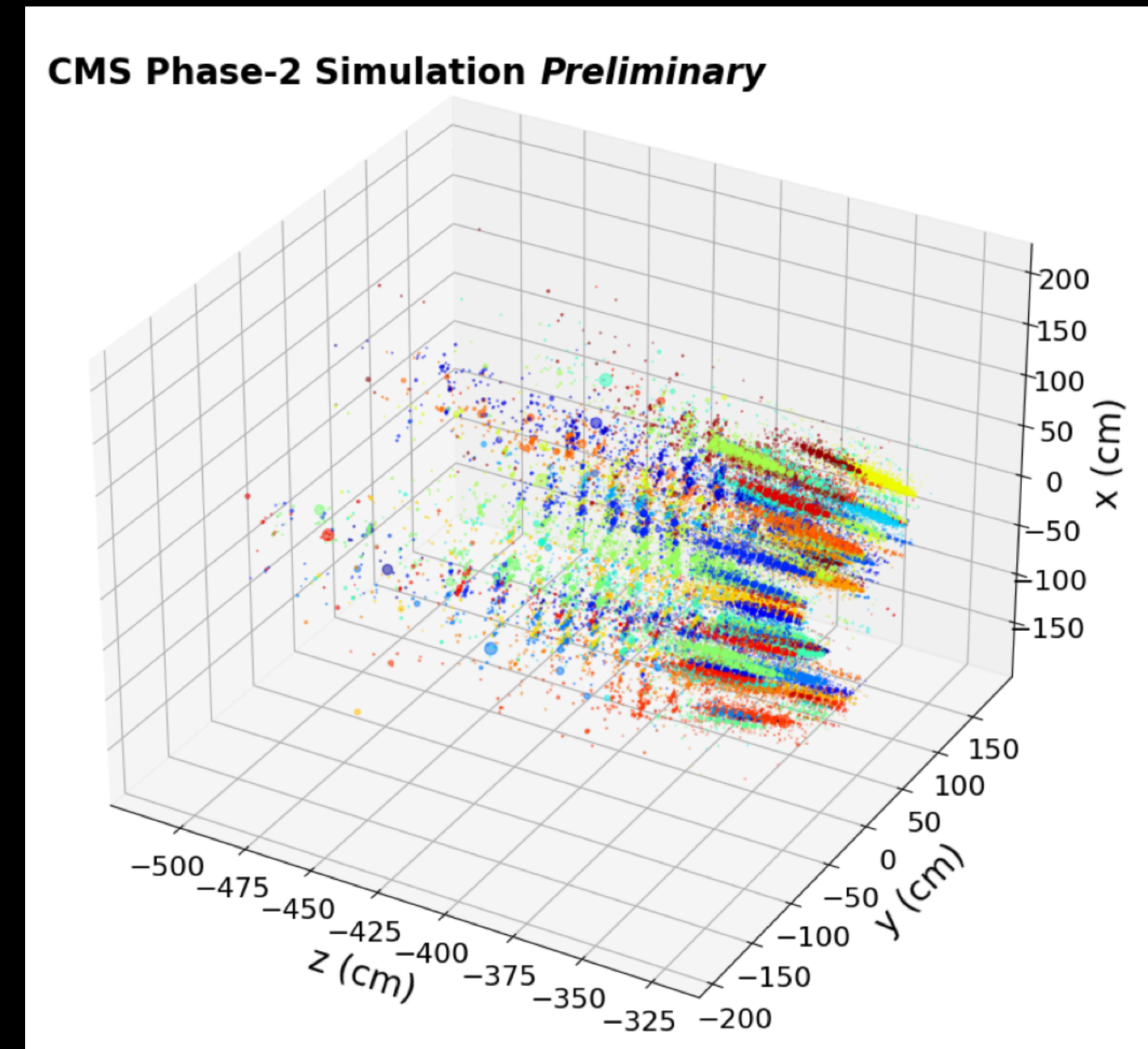


GNN with dynamic graph building is
clustering surrogate:

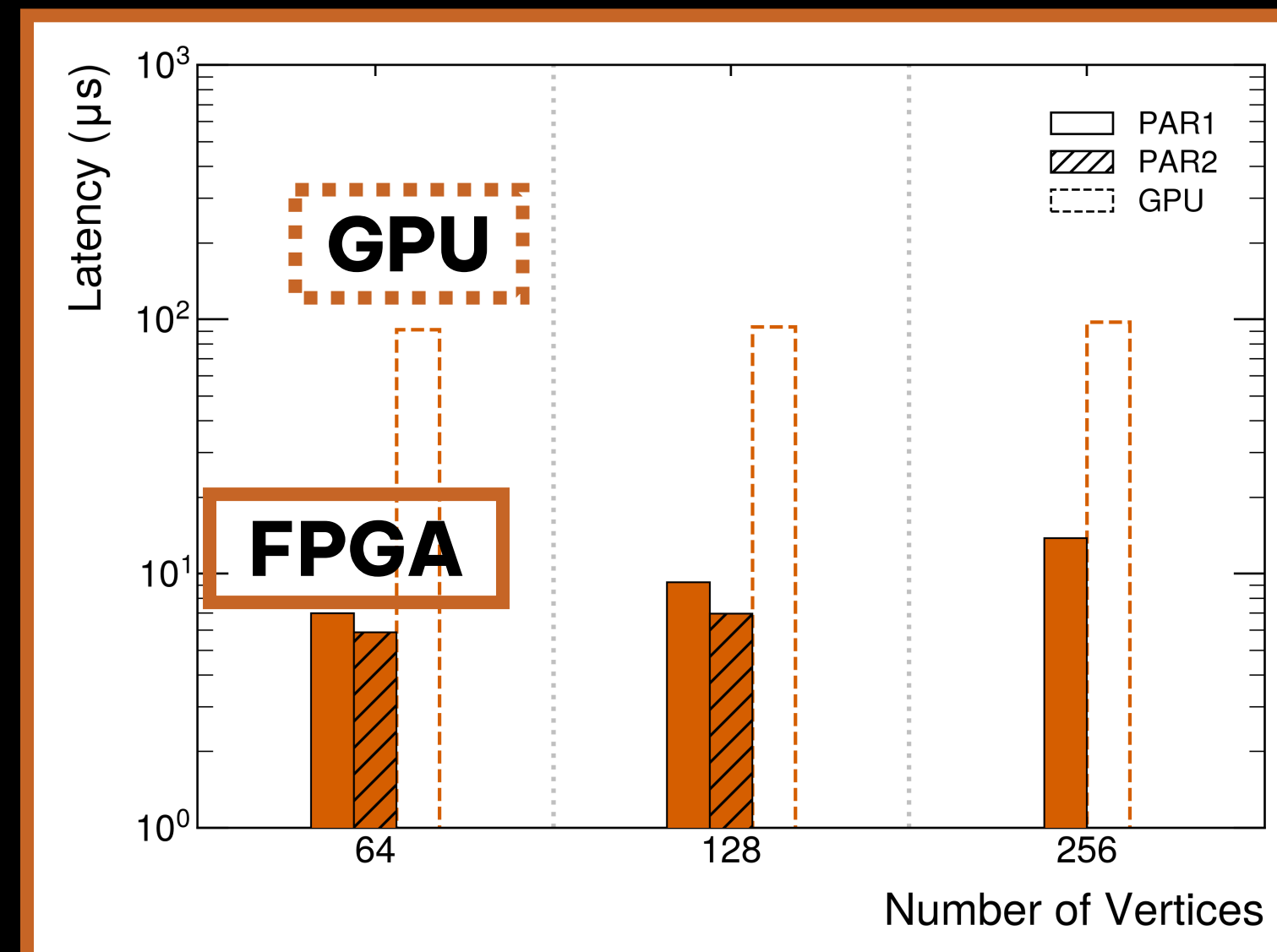
Hits belonging to same object clustered



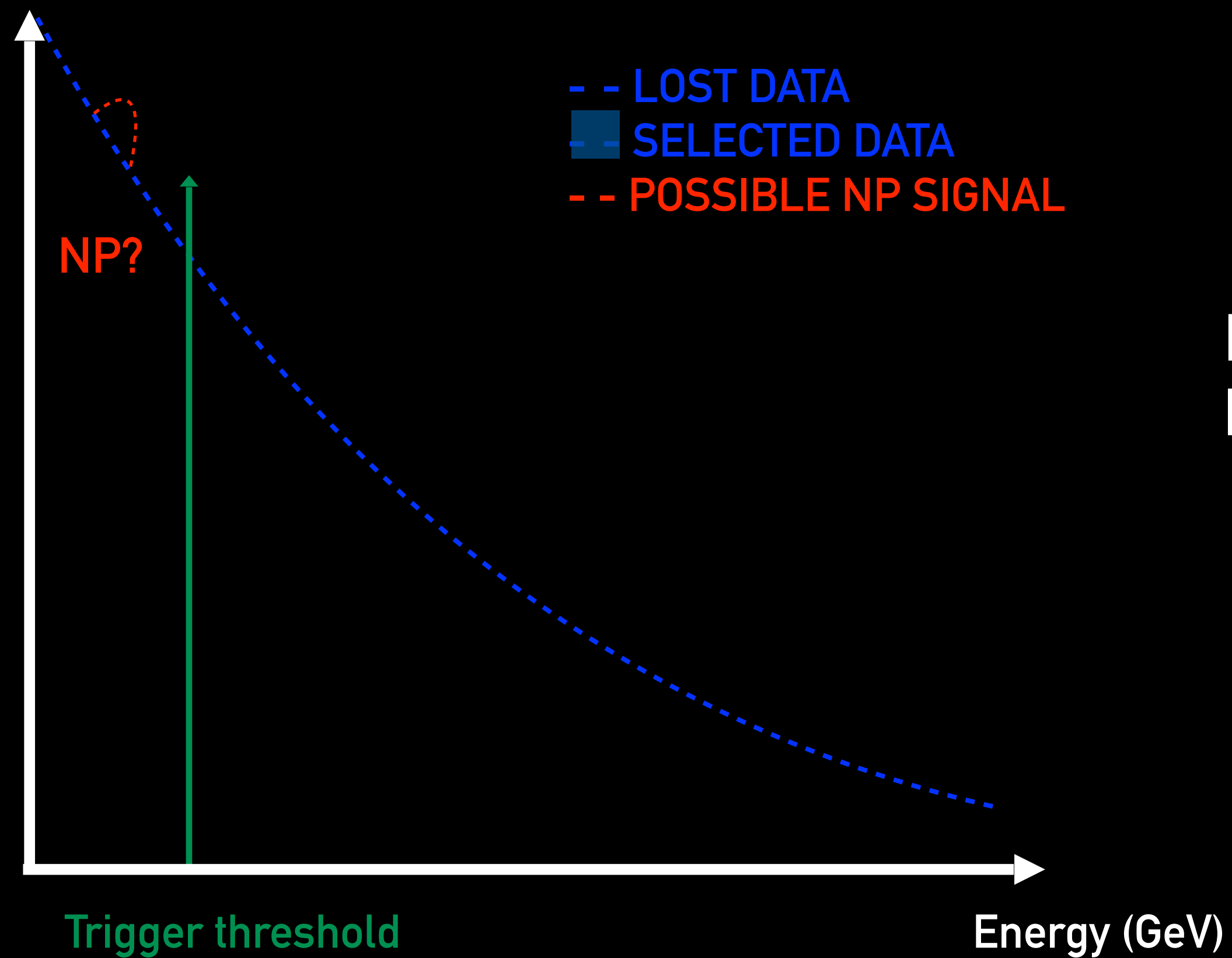
Cluster into distinct particle
energy clusters



$\sim 2 \mu\text{s}$ latency constraints!

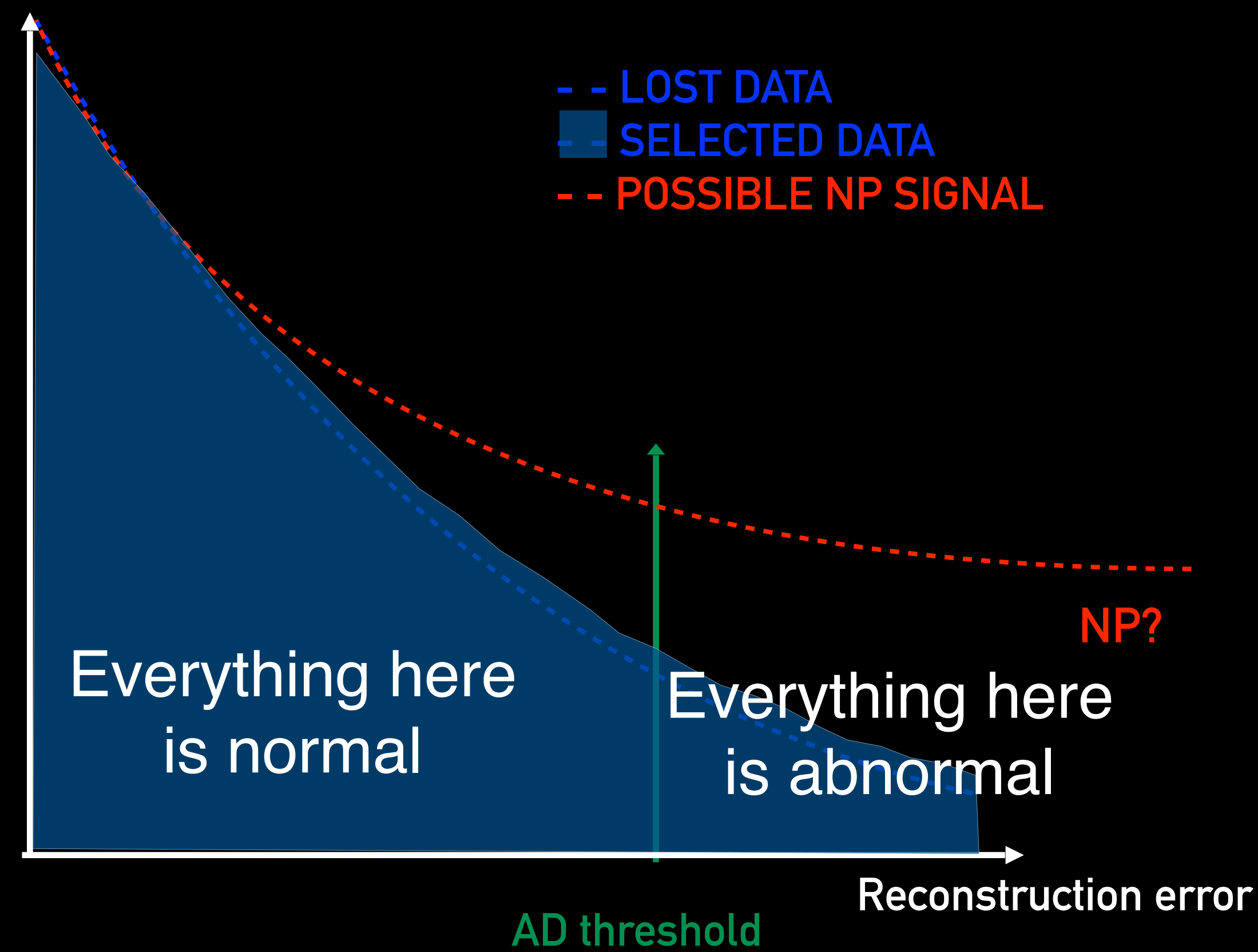
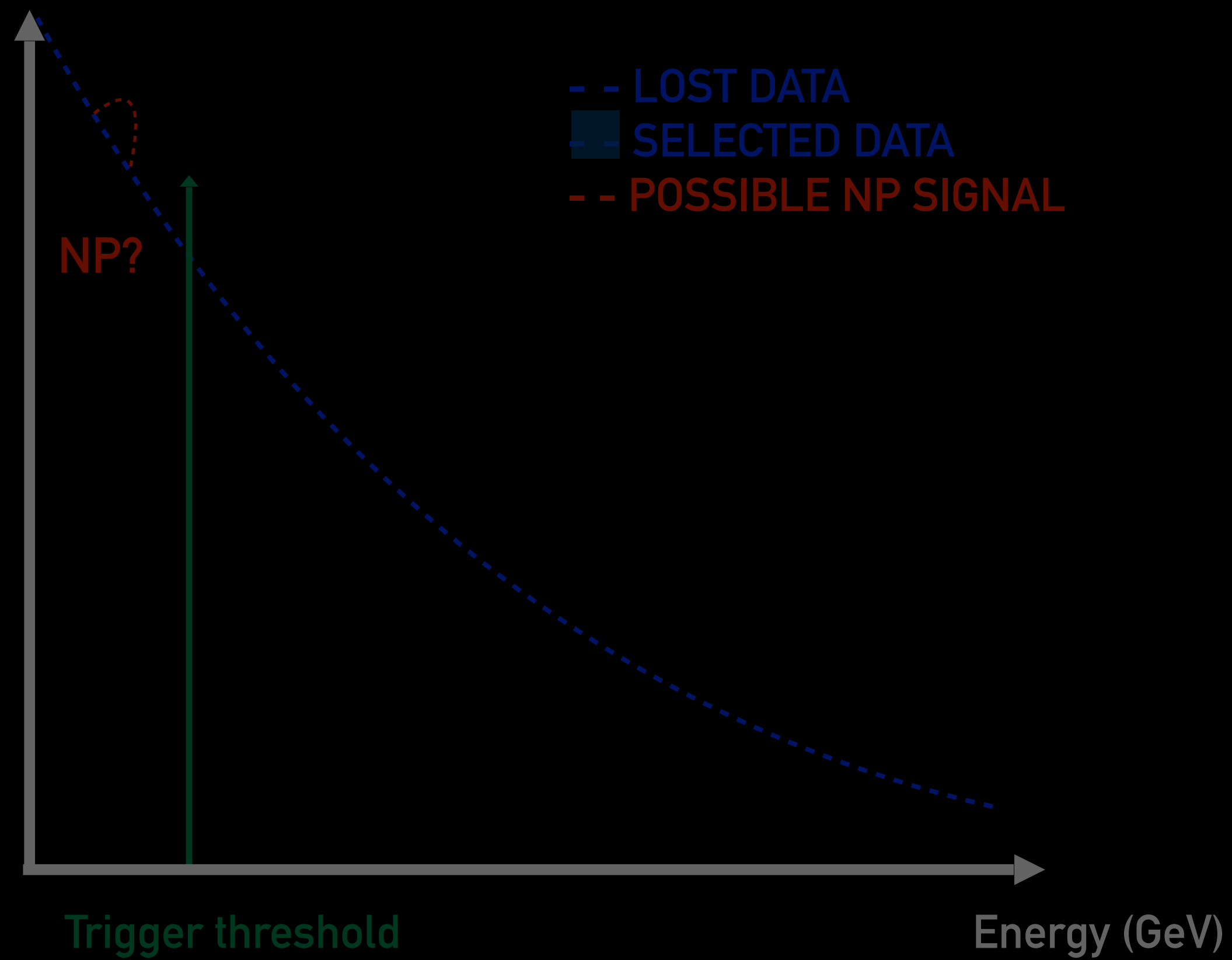


Anomaly Detection on FPGA

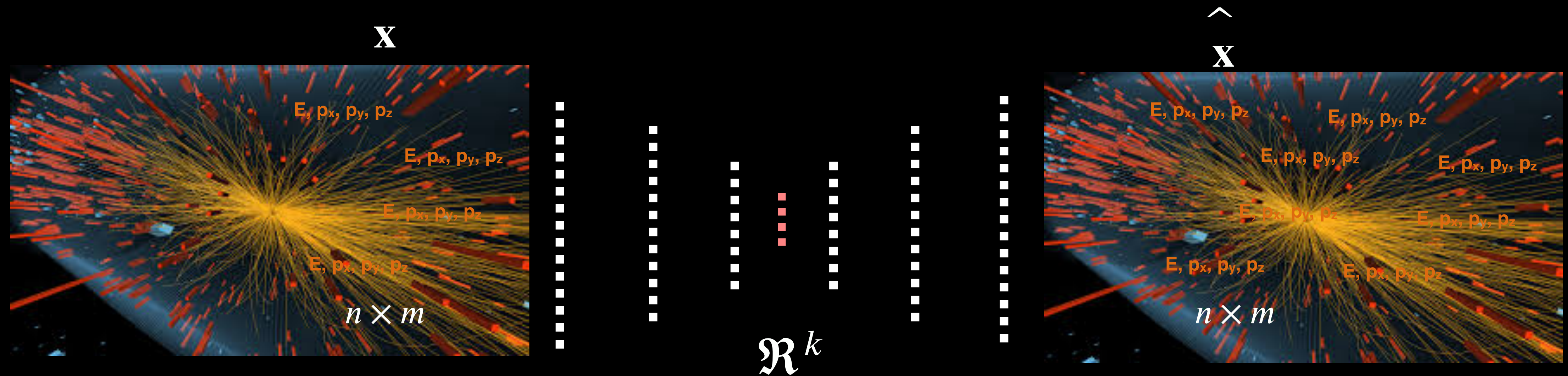


Level-1 rejects >99% of events!
Is there a smarter way to select?

Anomaly Detection triggers

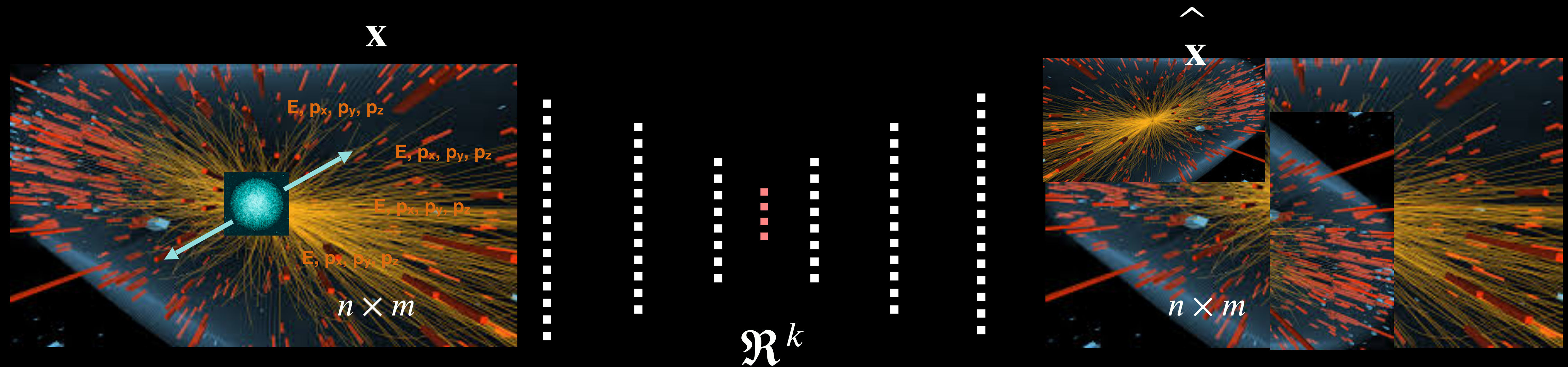


Outlier detection

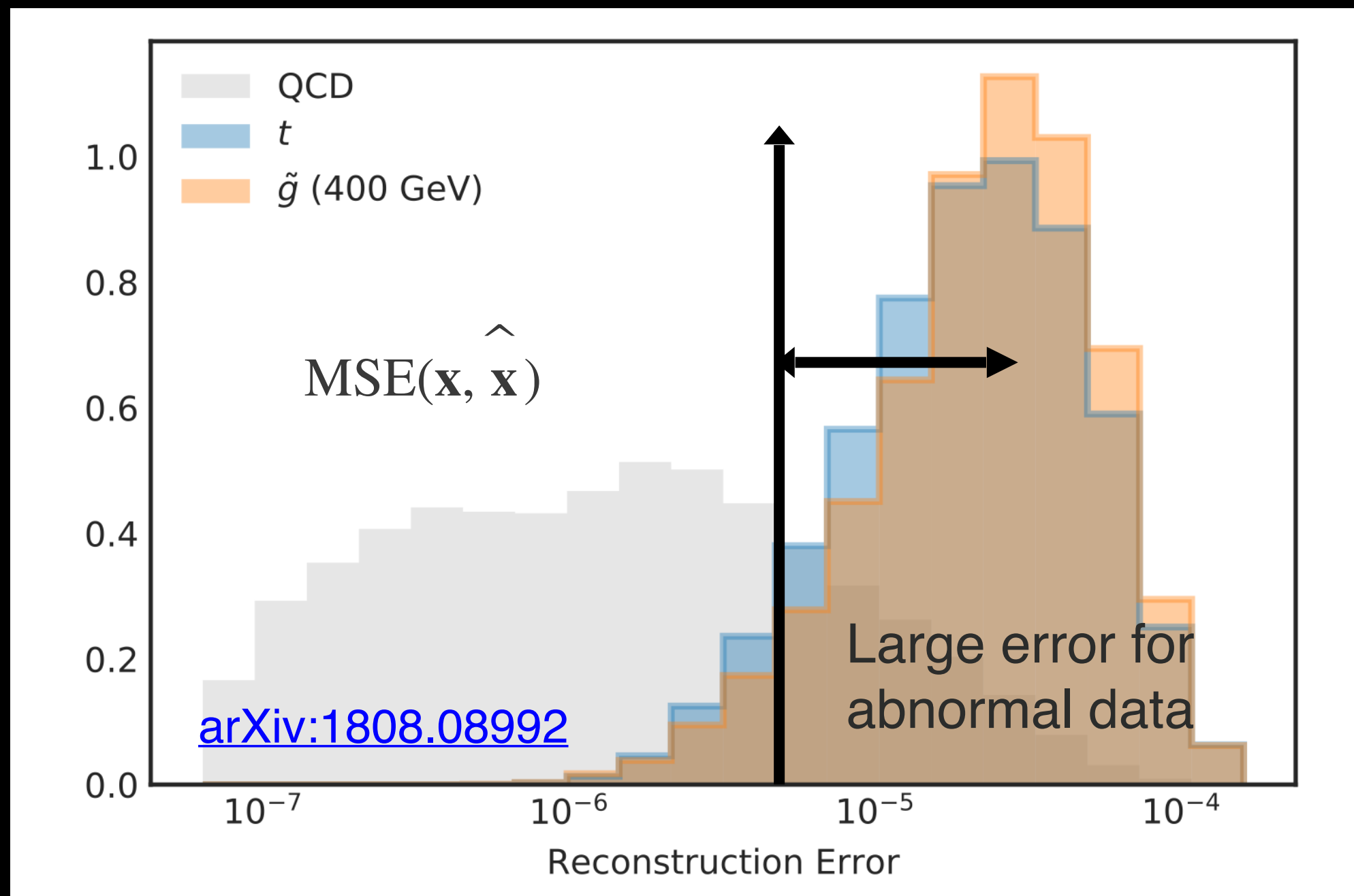
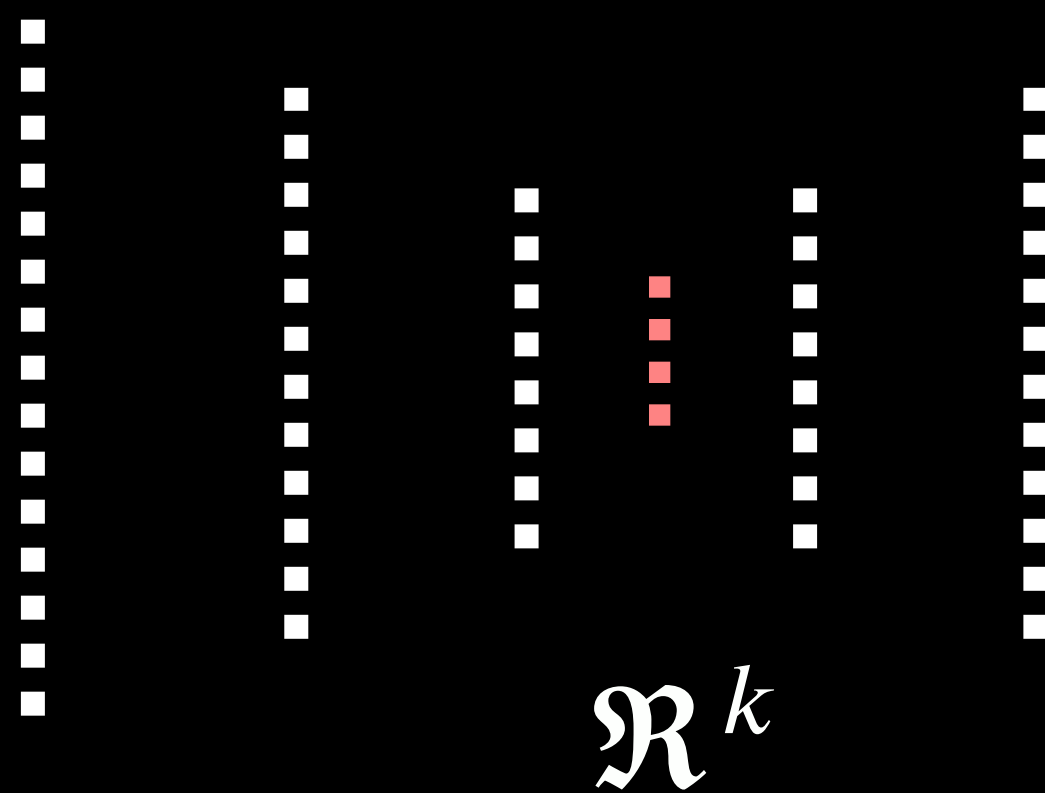
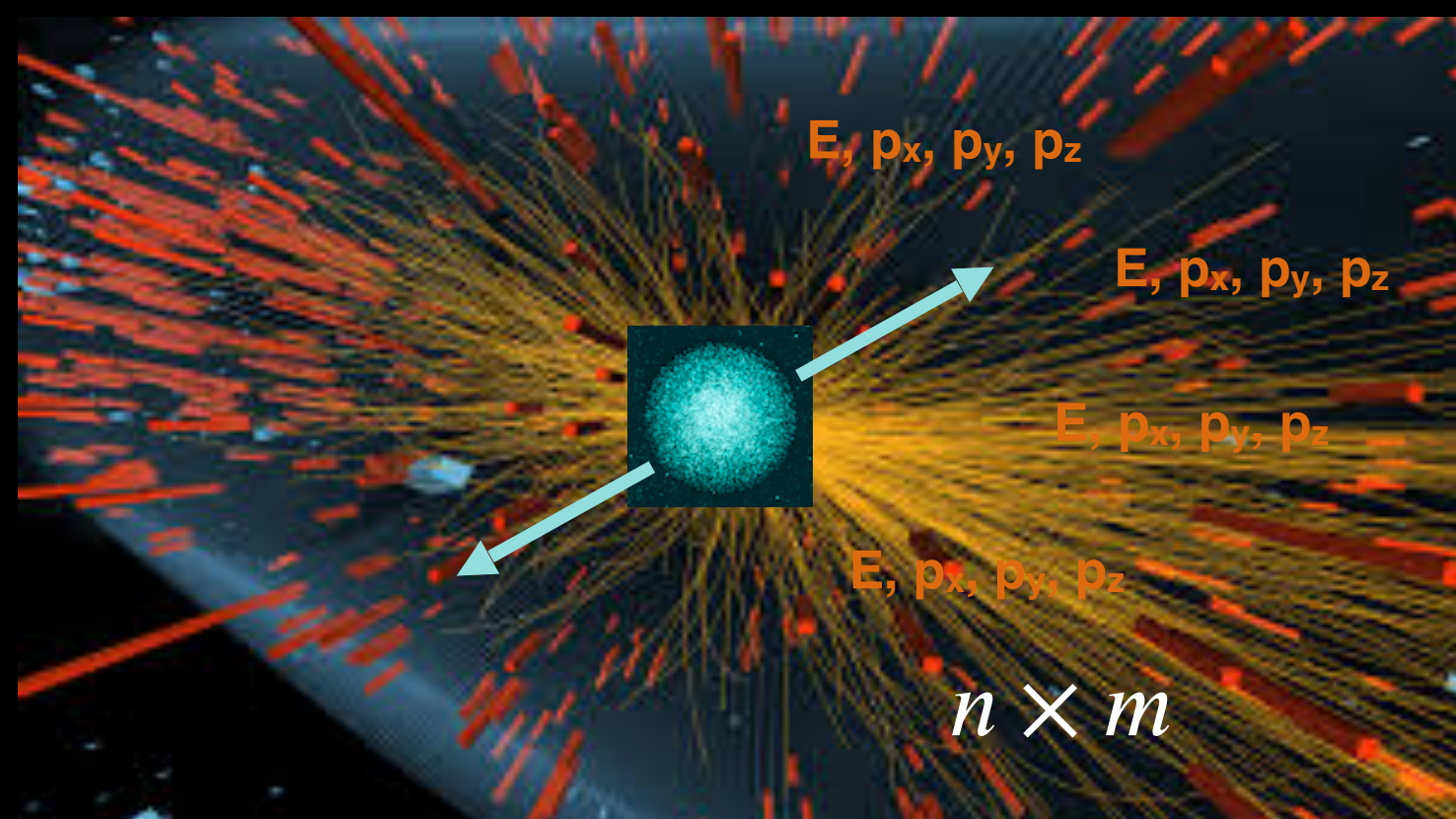


Compressed representation of \mathbf{x} .
Latent space \mathcal{R}^k , $k < m \times n$
prevents memorisation of input, must learn

Outlier detection

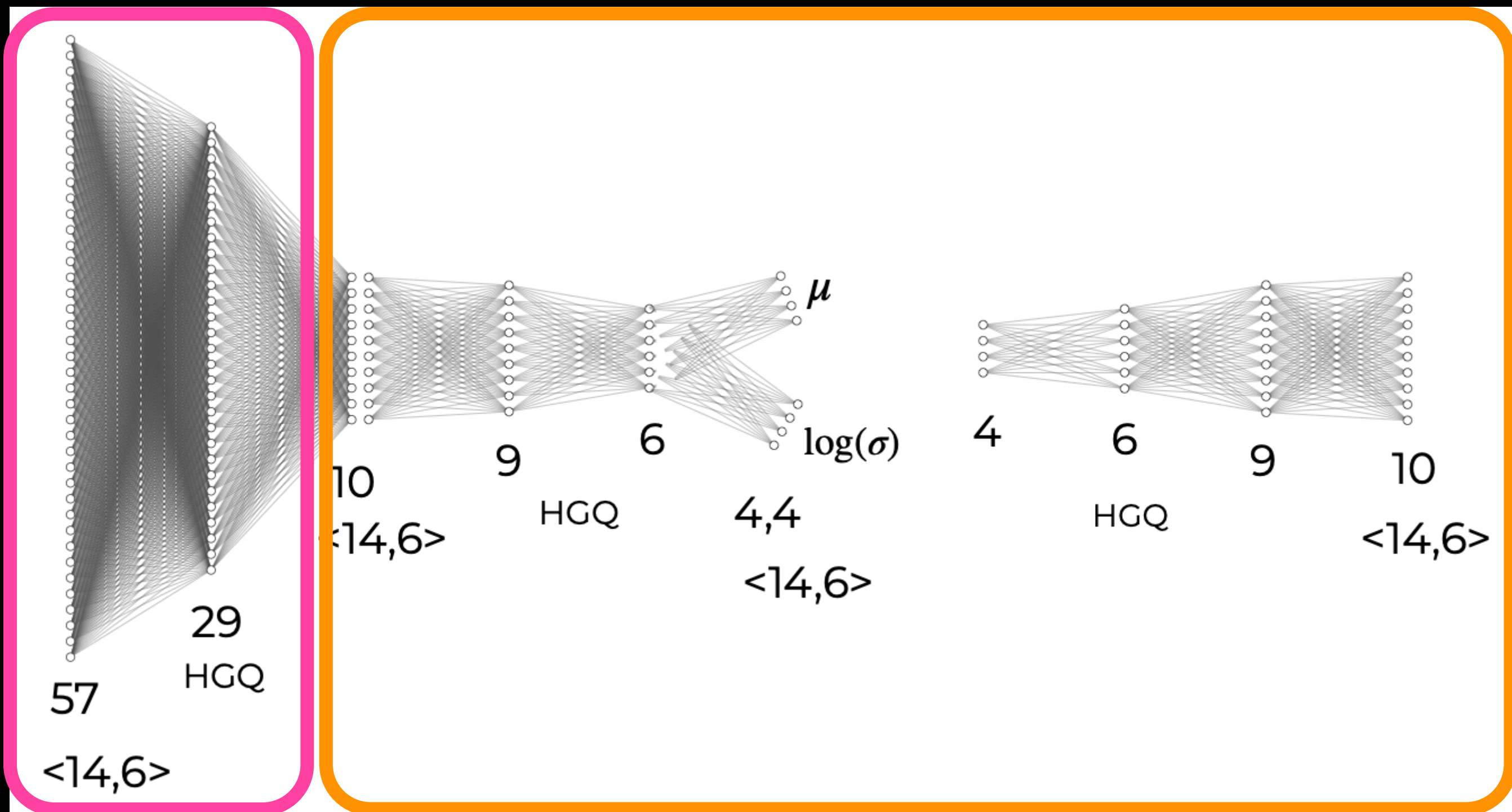


$\mathcal{L}(\mathbf{x}, \hat{\mathbf{x}})$ is Mean Squared Error($\mathbf{x}, \hat{\mathbf{x}}$), “high error events” proxy for “degree of abnormality”



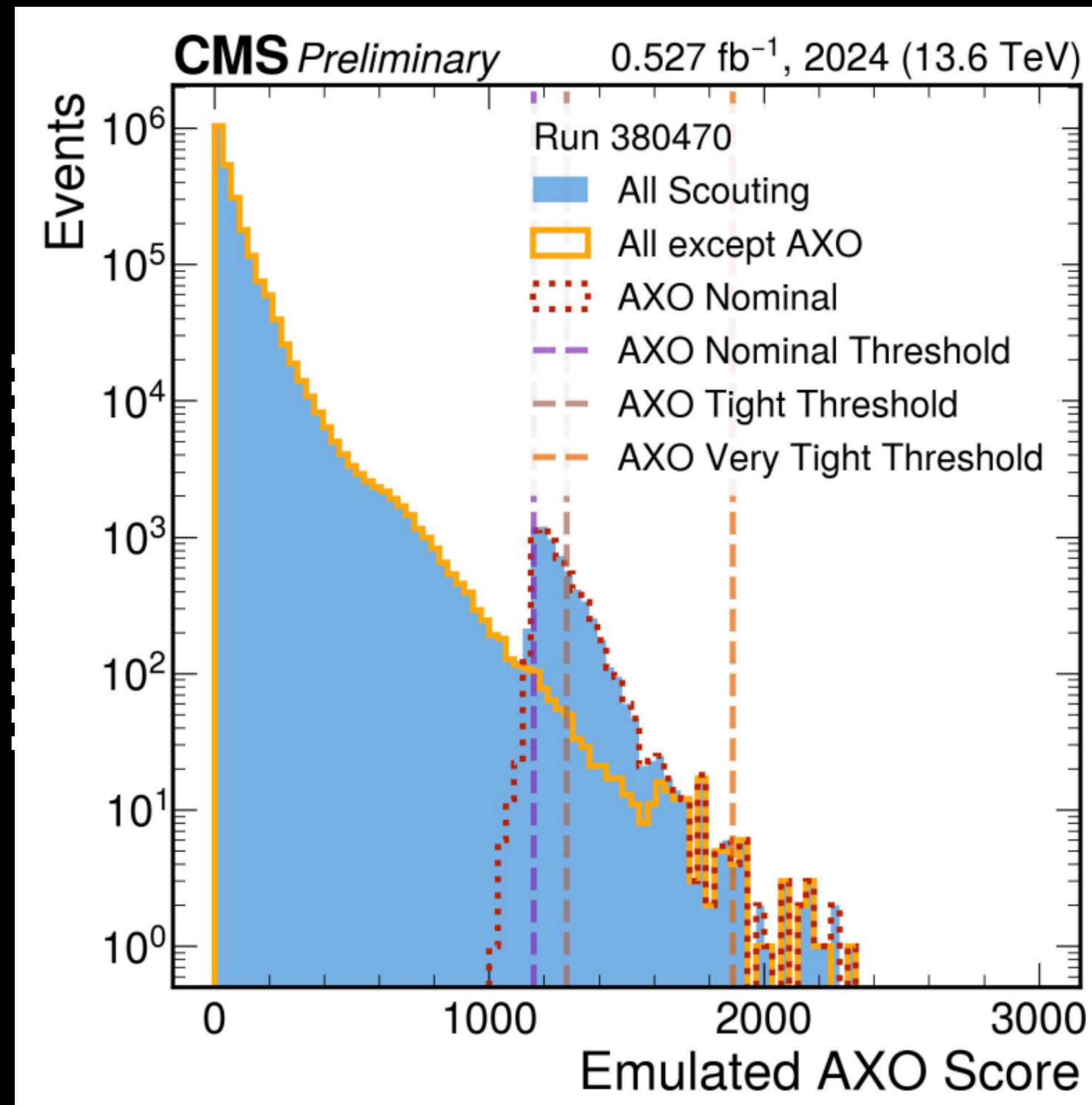


Neural embedding with VicReg contrastive learning

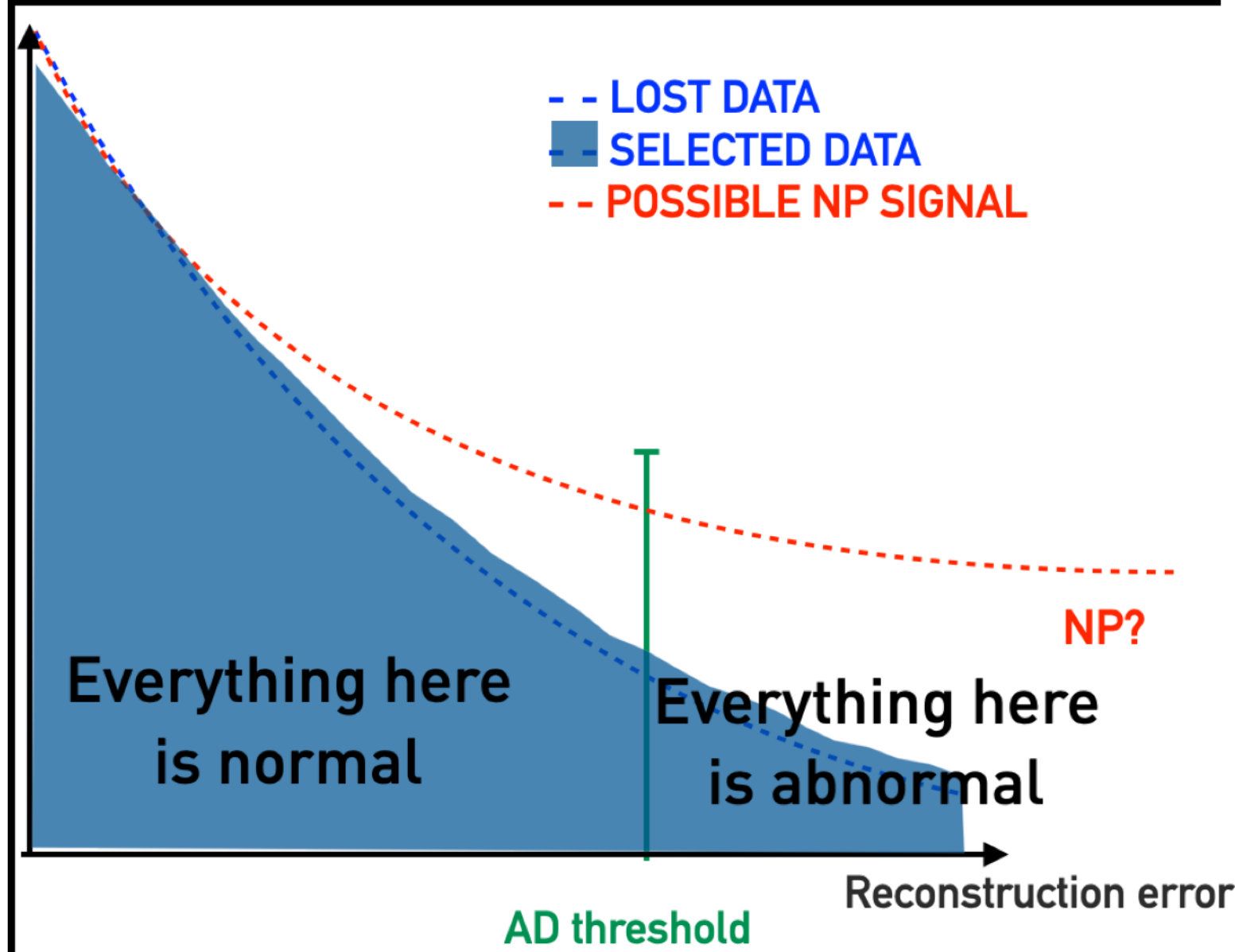


Variational Autoencoder

8 layers in 50 ns on FPGA,
~1% of total rate selected by AD

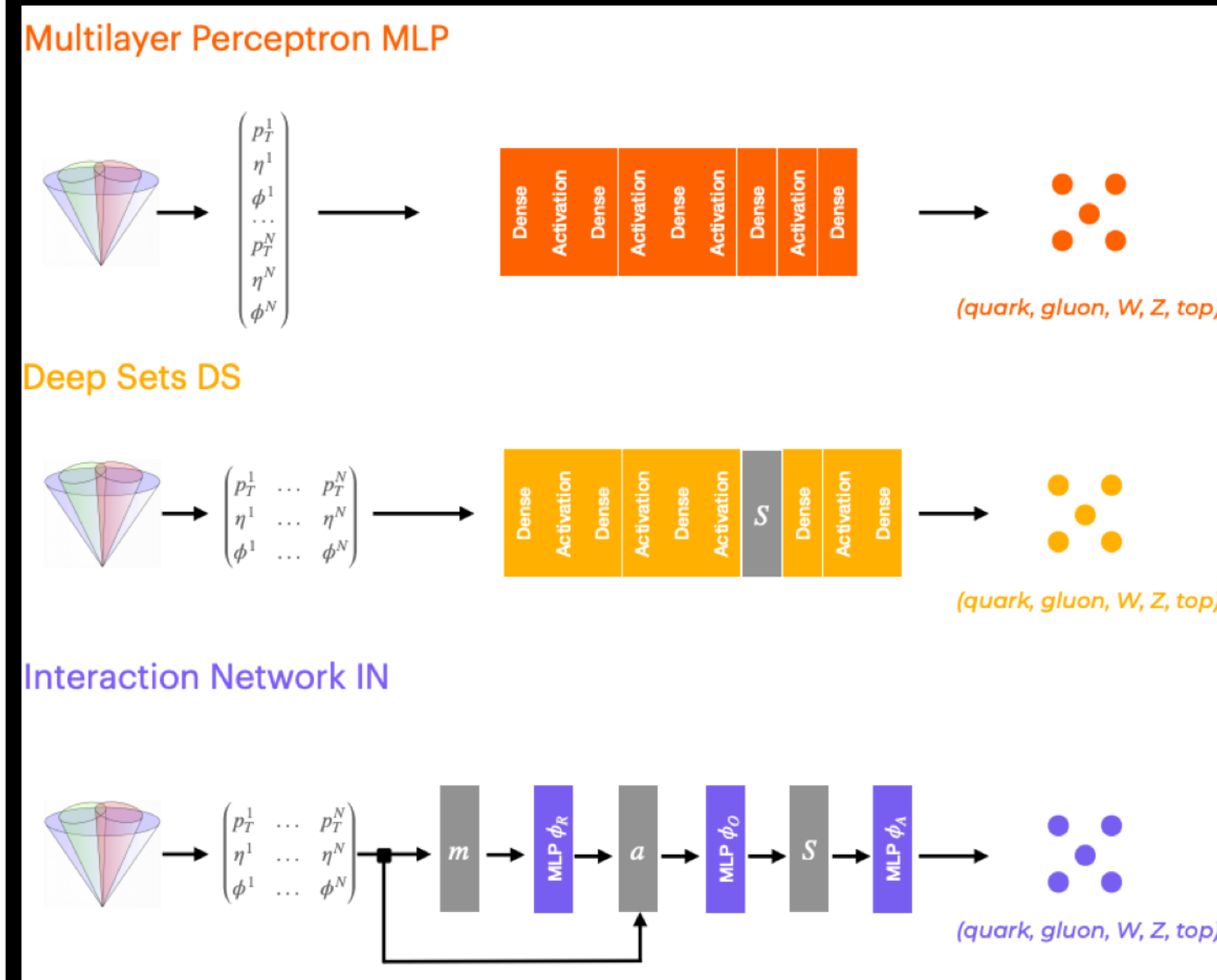


Anomaly detection with VAEs in 50 ns



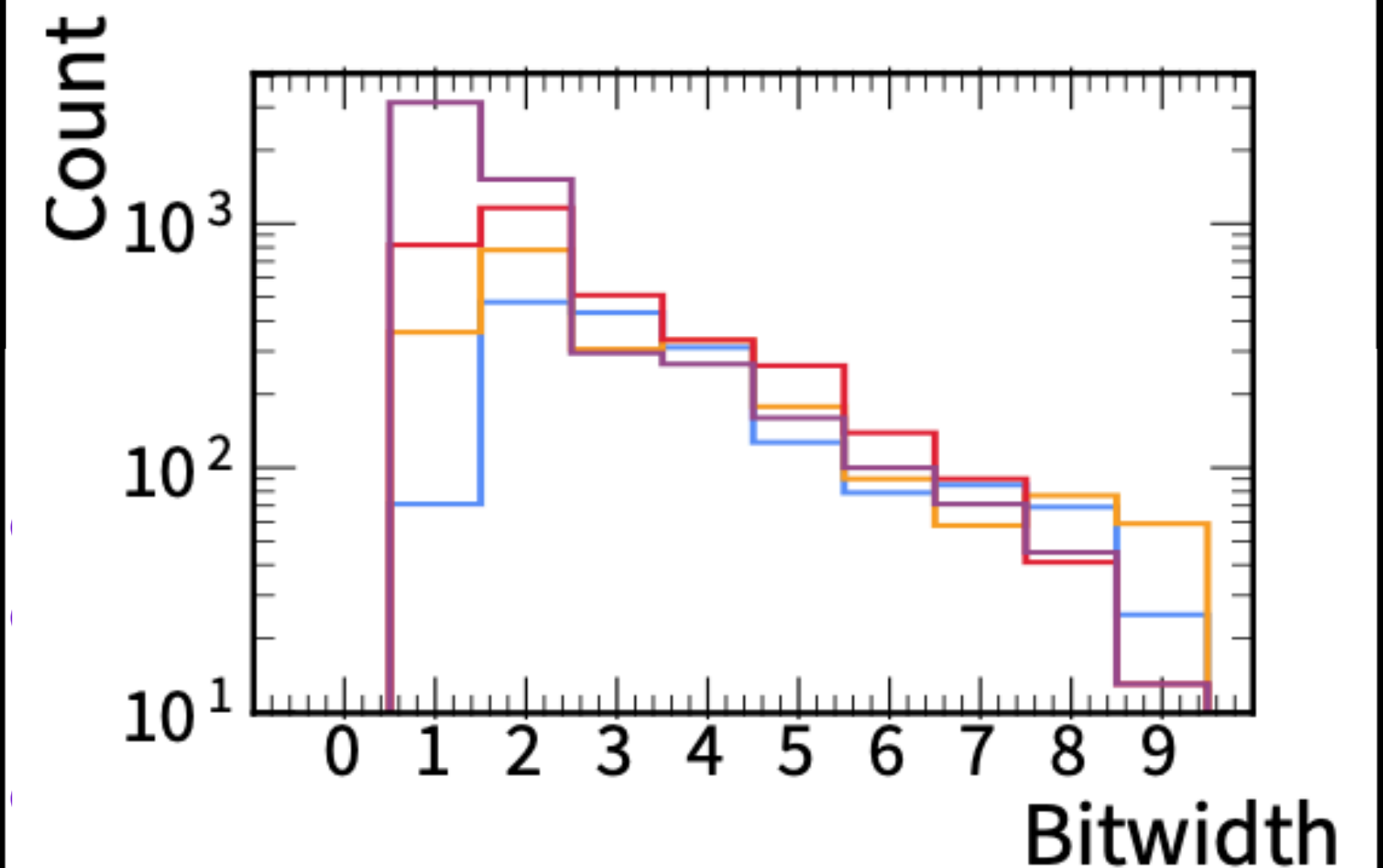
CMS DP2023_079
 E. Govorkova et al (2022)

Quantised Interaction Networks and Deep Sets in <160 ns



P. Odagiu et al. 2024

Fully heterogeneously quantized transformers in <100 ns



(a) Attention Layers

[arxiv:2510.24784](https://arxiv.org/abs/2510.24784)

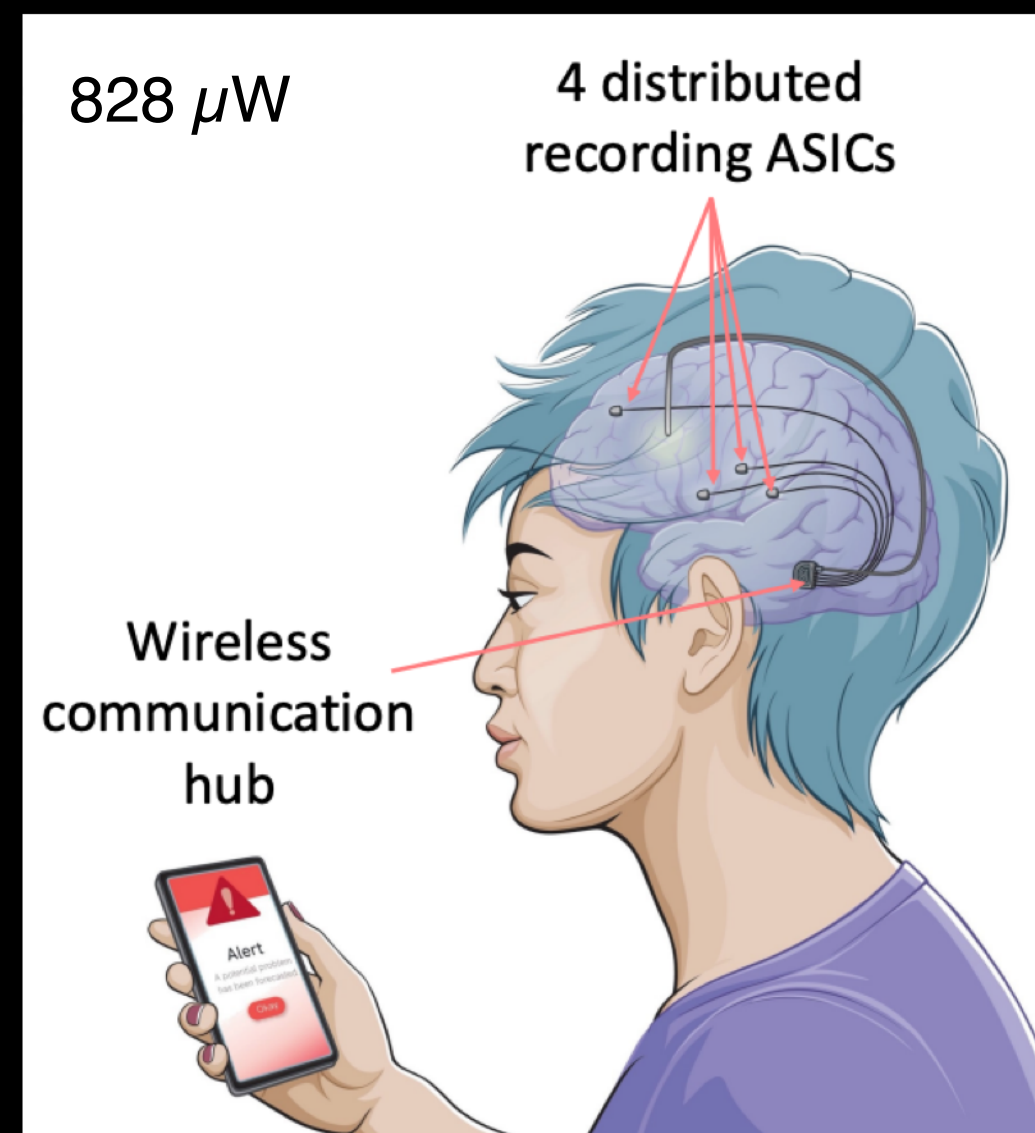
...and outside

Semantic segmentation for autonomous vehicles



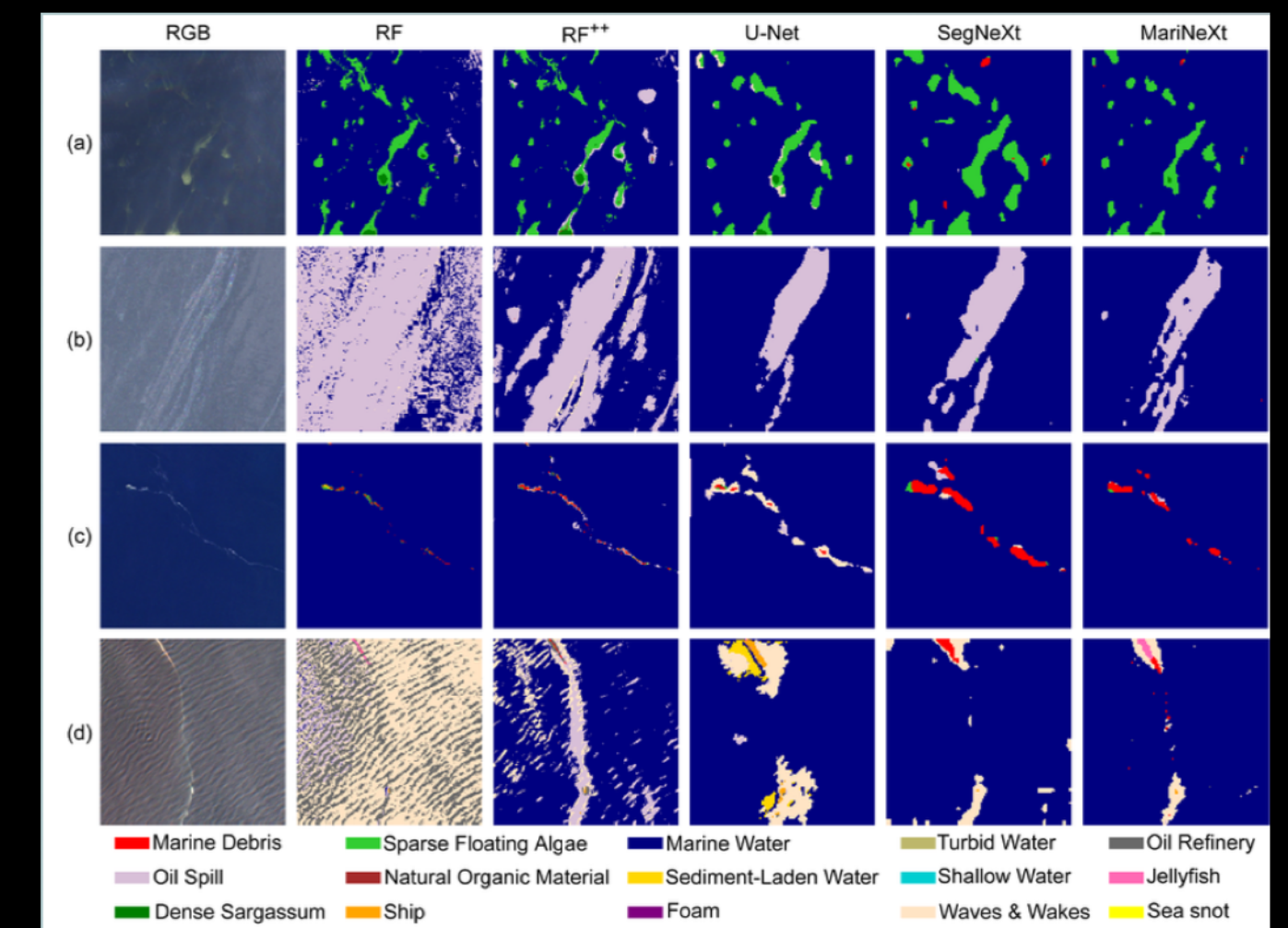
N. Ghielmetti et al. 2022

Seizure Predicting Brain Implant



W. Lemaire et al. 2022

Earth monitoring in satellites



Edge SpAlce, S. Summers

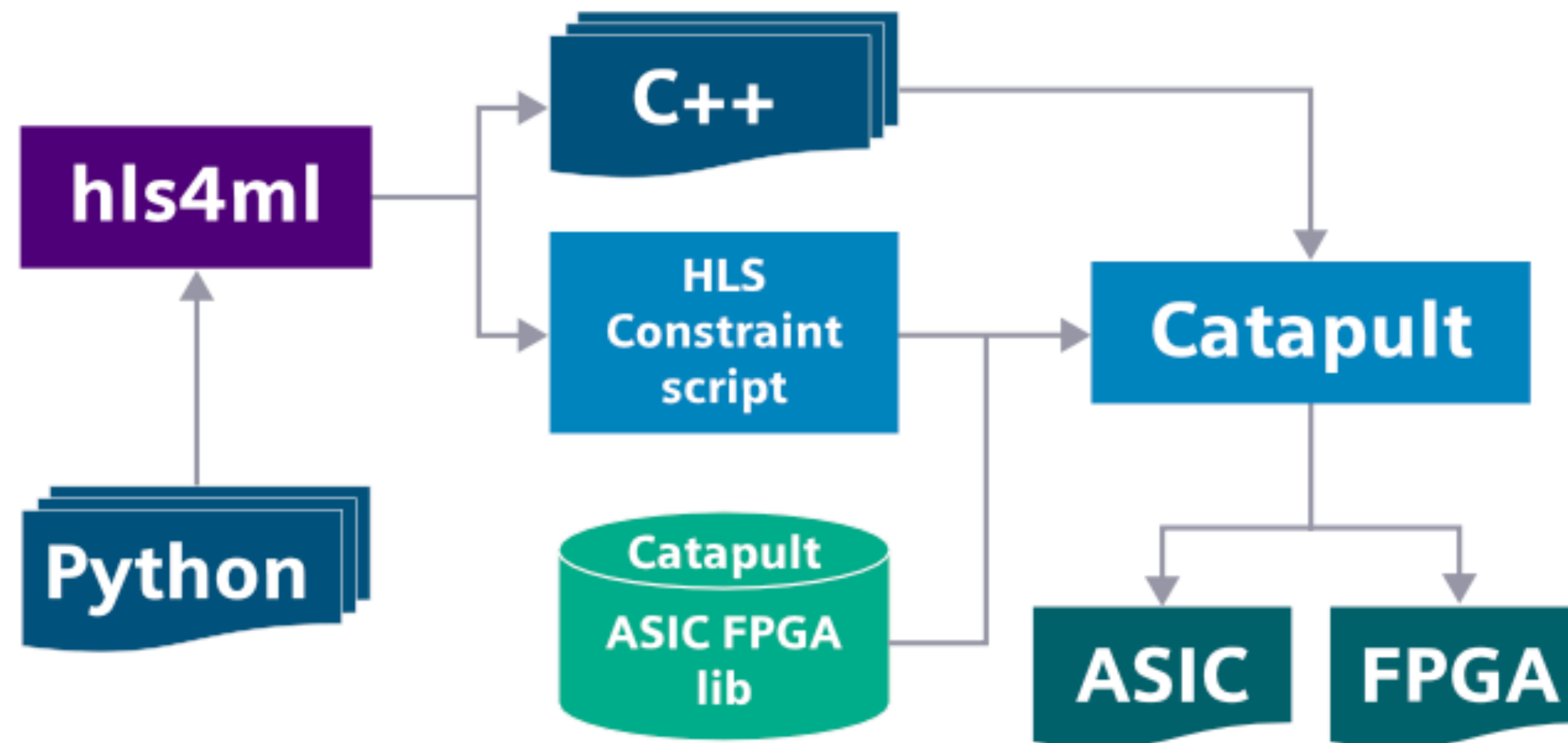
- MLPerf tinyML benchmarking
- For fusion science phase/mode monitoring
- Crystal structure detection
- Triggering in DUNE

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

PRESS RELEASE

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

May 21, 2024
Plano, Texas



“...brings together hls4ml and Siemens' Catapult™ HLS software...

... addresses the unique requirements of

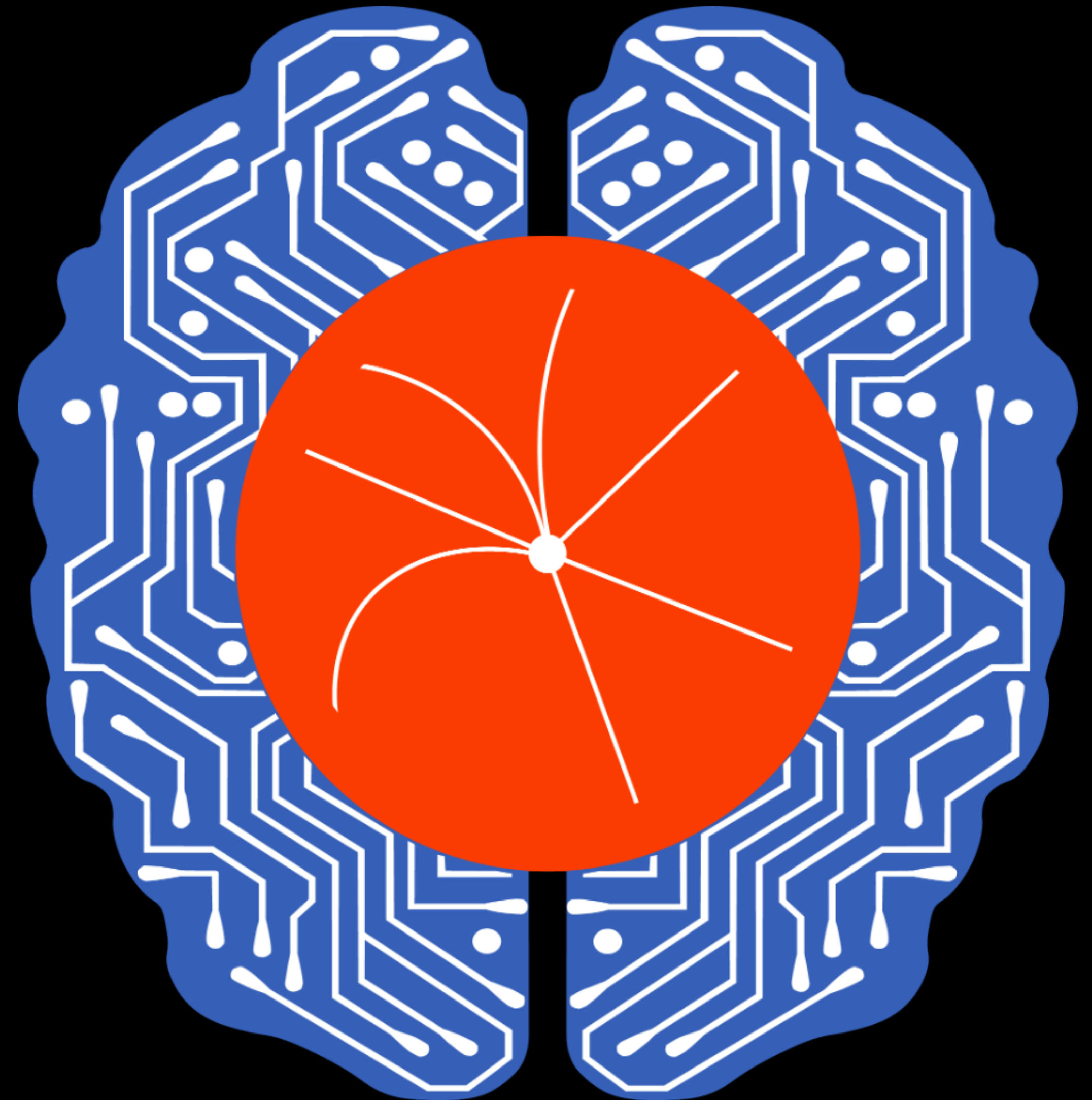
ML accelerator design for power, performance, and area on custom silicon”

fast machine learning foundation

fastmachinelearning.org

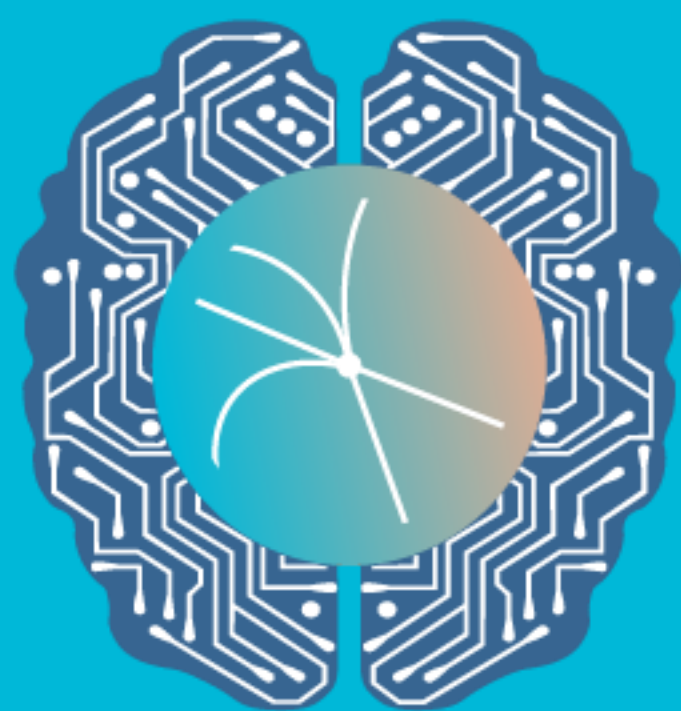
Collaborating Institutes

Argonne National Laboratory	Pacific Northwest National Laboratory
CERN	Princeton University
California Institute of Technology	Purdue University
Carnegie Mellon University	SLAC
ETH Zurich	Southern Methodist University
Fermi National Accelerator Laboratory	Technical University of Denmark
Georgia Institute of Technology	Tohoku University
HEPIA - Haute école du paysage, d'ingénierie et d'architecture	UC San Diego
Imperial College London	UCL
Jefferson Lab	University of Bristol
Jožef Stefan Institute	University of Chicago
KIT	University of Colorado Boulder
LPNHE, CNRS/IN2P3, France	University of Geneva
Los Alamos National Laboratory	University of Illinois at Urbana-Champaign
Massachusetts Institute of Technology	University of Michigan
National Yang Ming Chiao Tung University	University of Pennsylvania
Northwestern University	University of Washington
Norwegian University of Science and Technology	University of Wisconsin-Madison
	Université de Sherbrooke



Fast Machine Learning for Science

August 31 - September 4, 2026



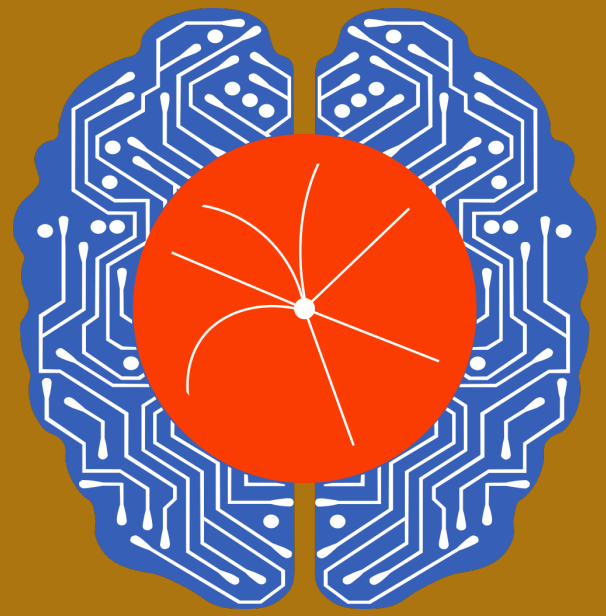
indi.to/fastml26

UC San Diego



The Fast ML community

<https://fastmachinelearning.org/>



Join our Slack to get updates!