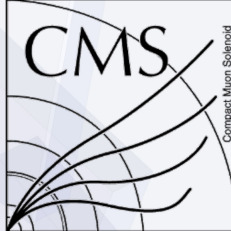


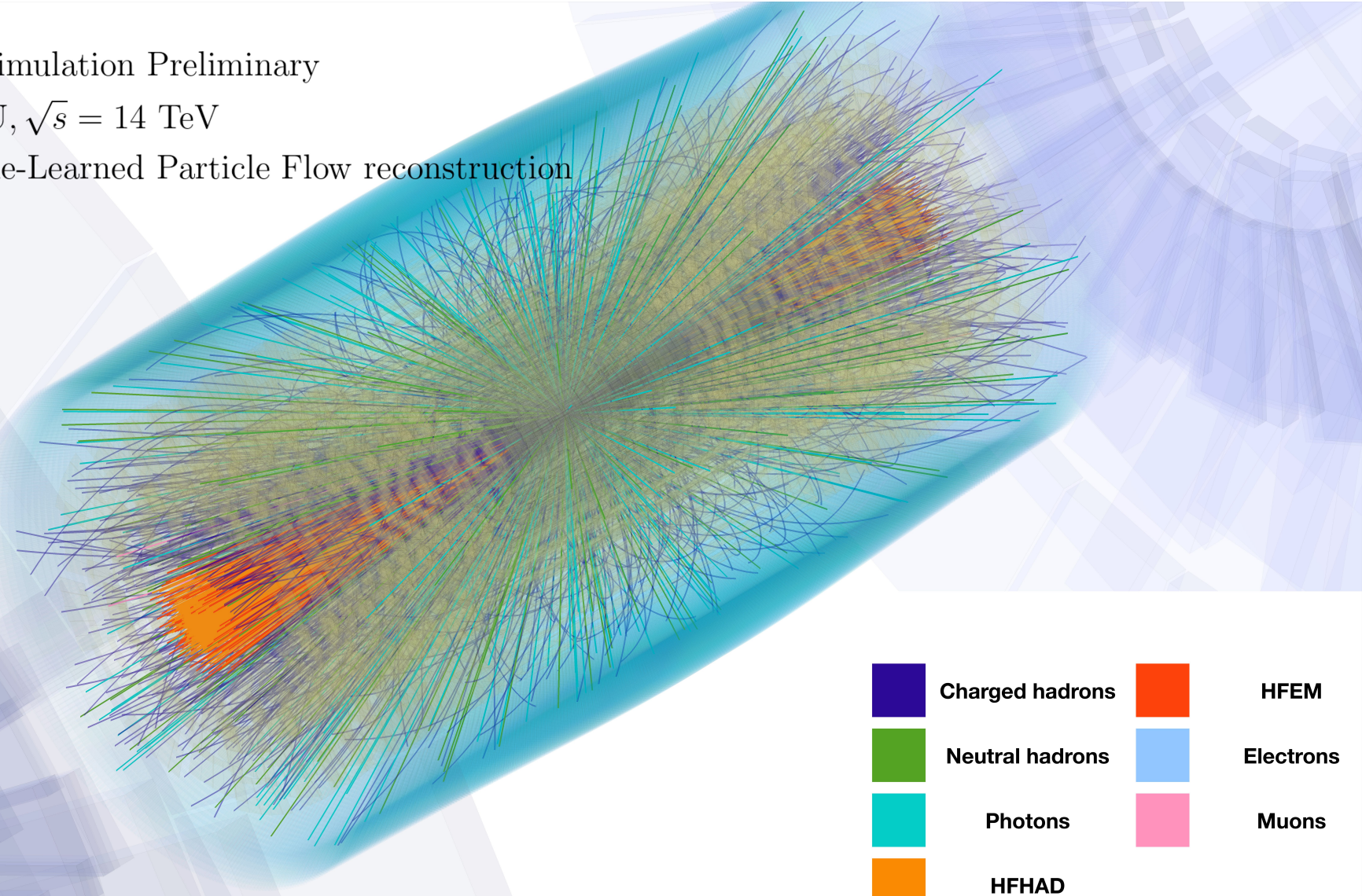


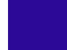






Bigger data, shorter time:
Real-time inference for scientific discovery

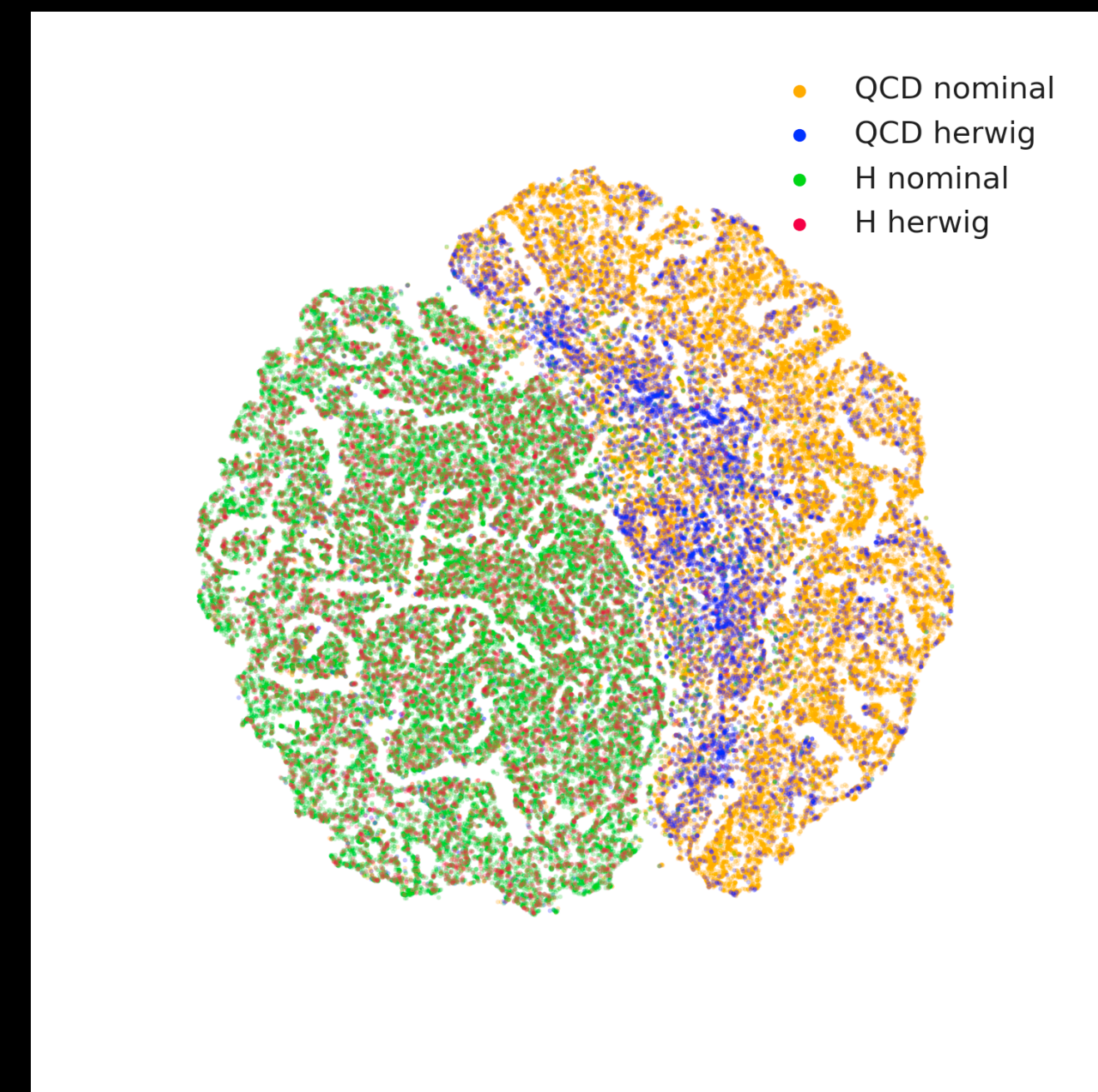
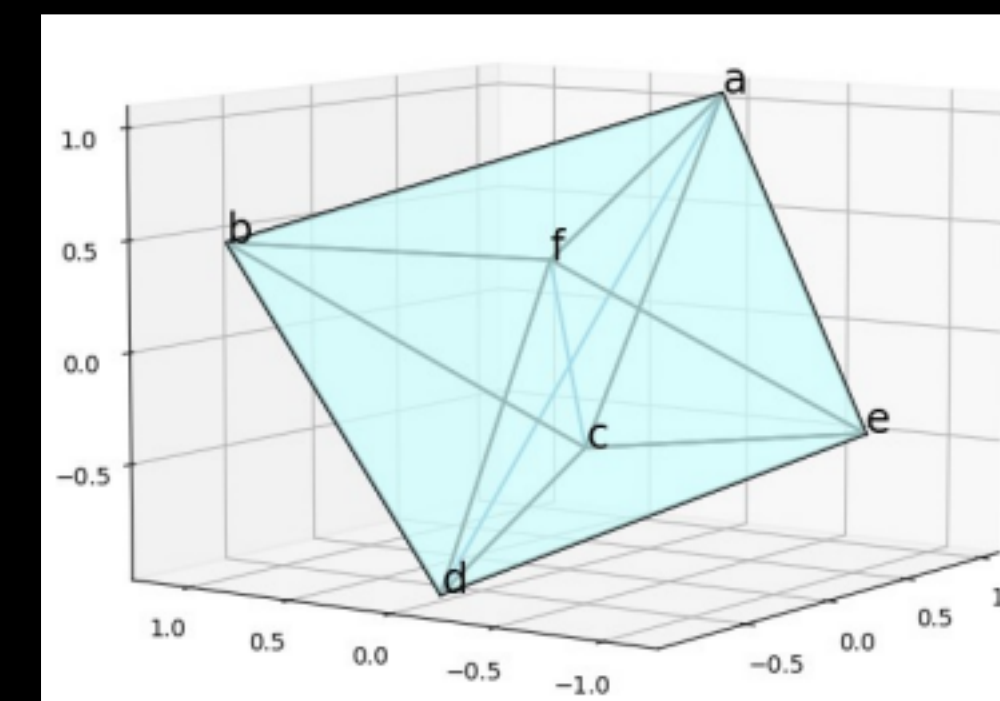
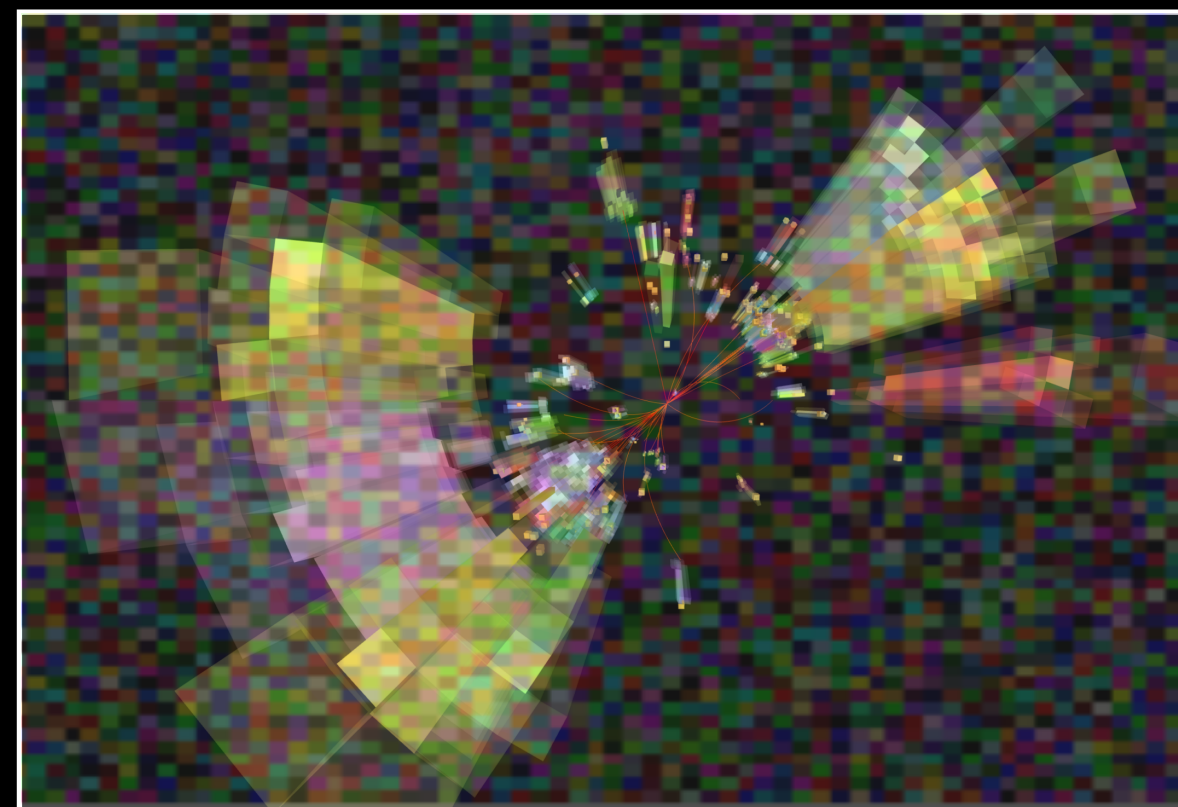
Thea Klæboe Årrestad
(ETH Zürich)

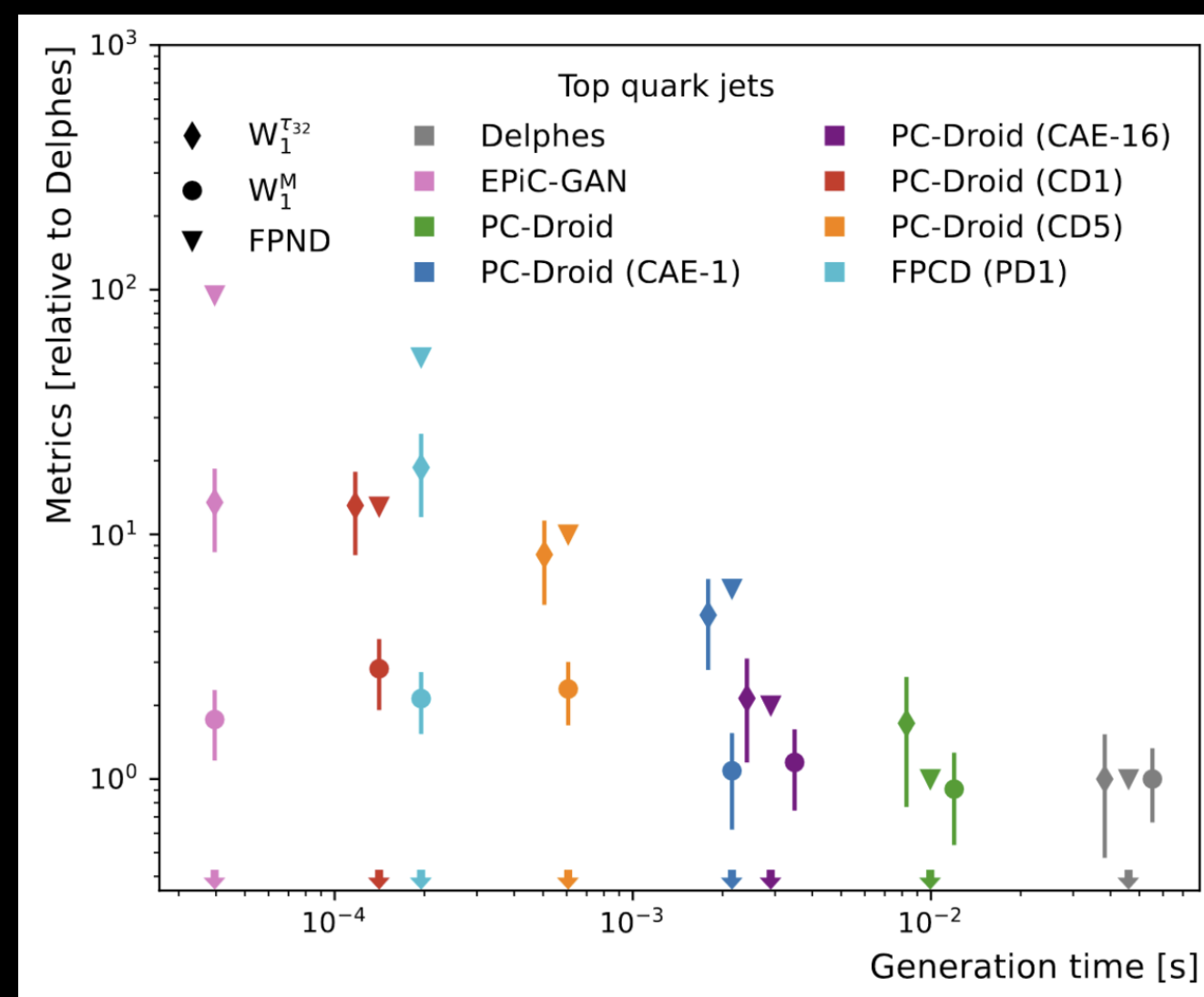
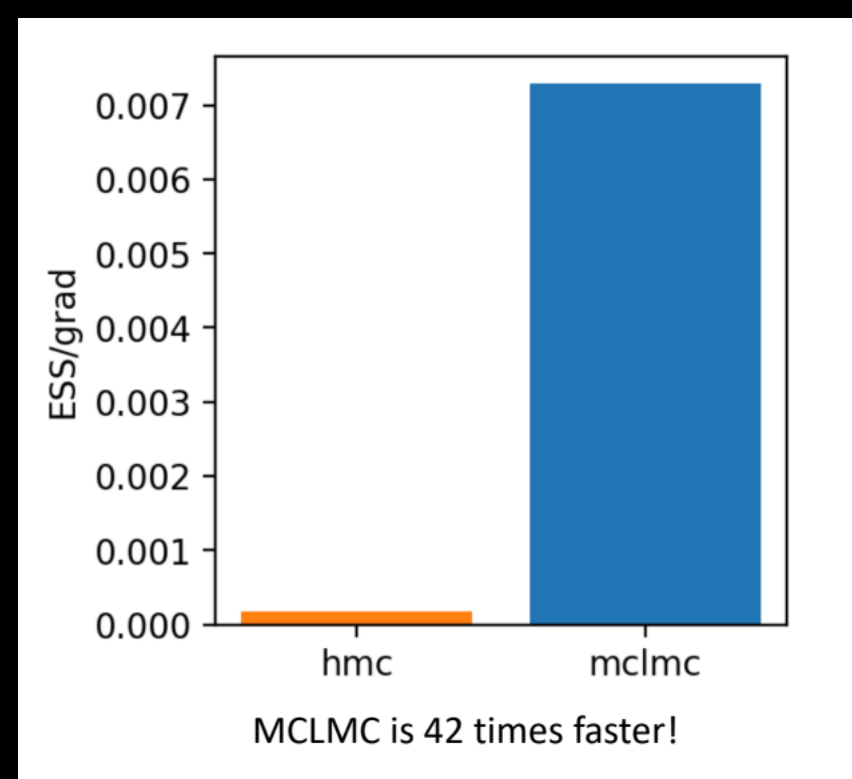
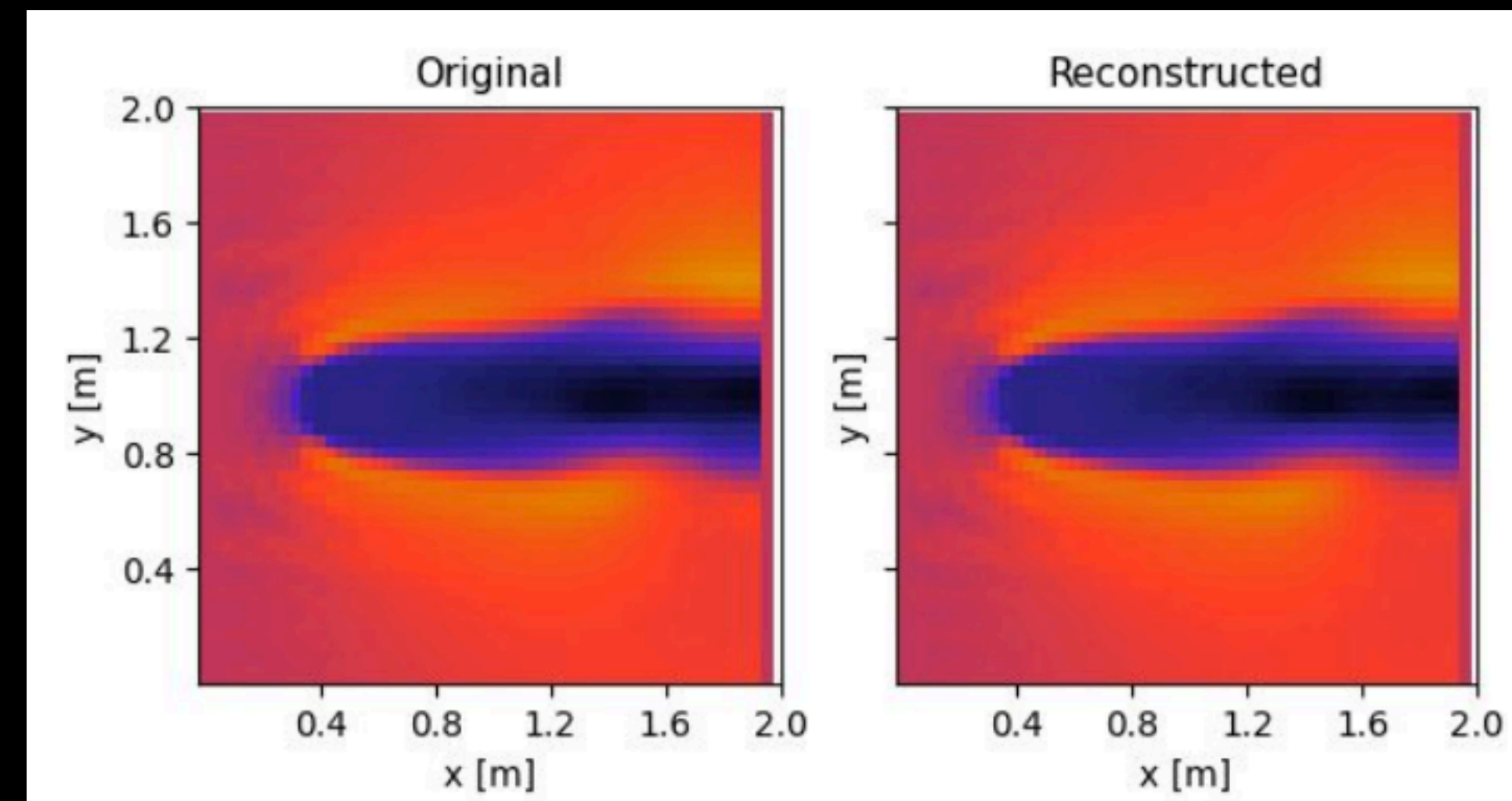
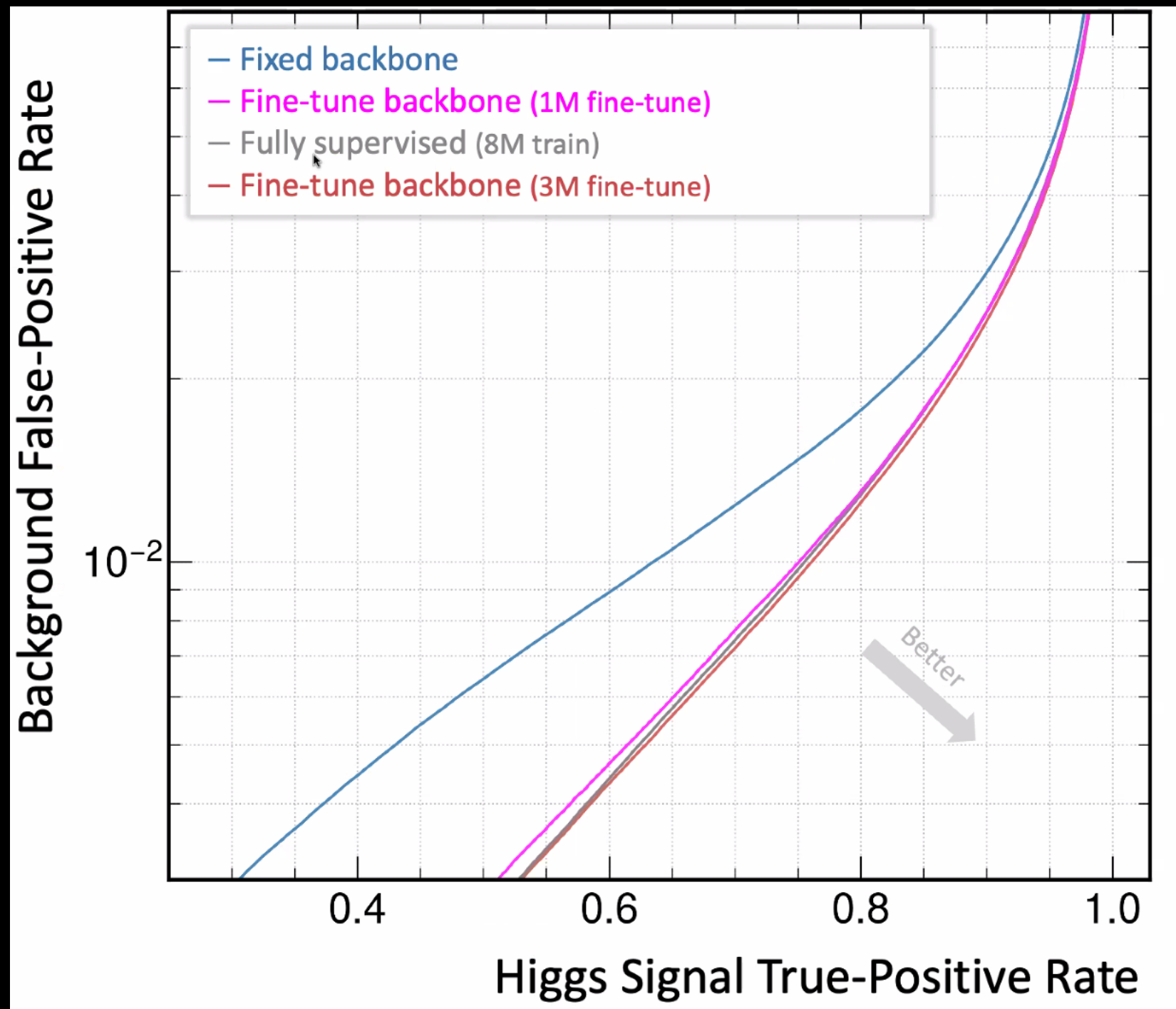


CMS Simulation Preliminary
 $t\bar{t} + \text{PU}, \sqrt{s} = 14 \text{ TeV}$
Machine-Learned Particle Flow reconstruction

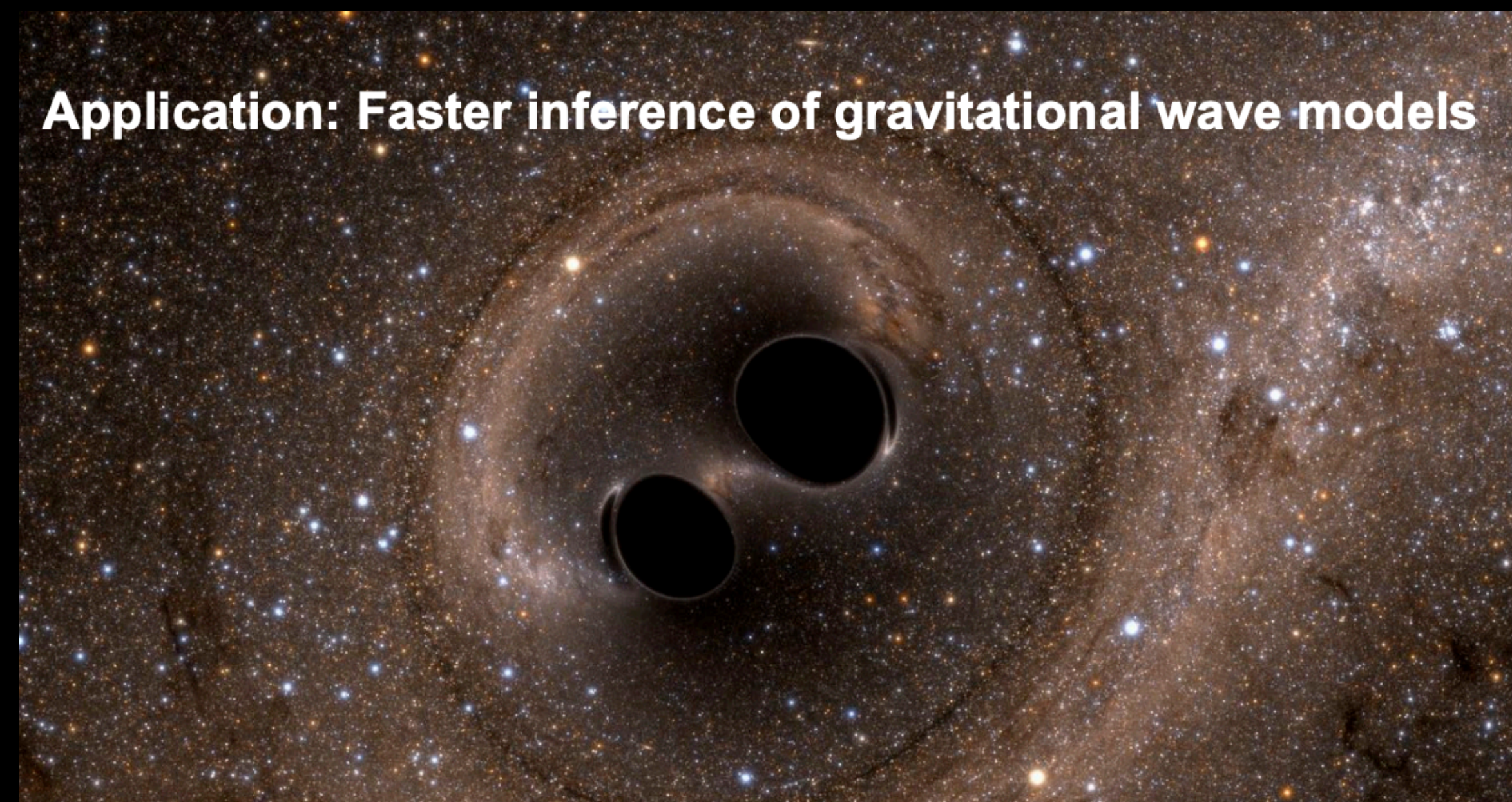


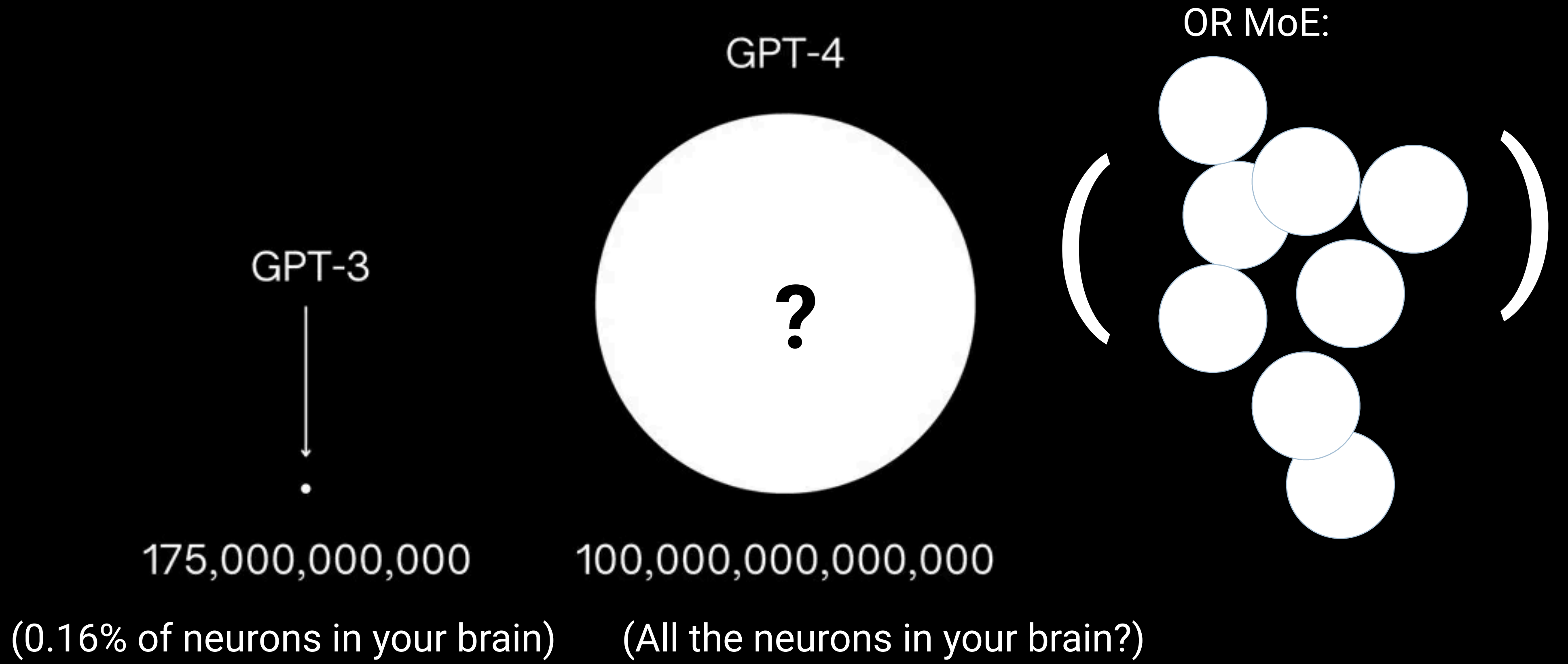
- | | |
|---|---|
|  Charged hadrons |  HFEM |
|  Neutral hadrons |  Electrons |
|  Photons |  Muons |
|  HFHAD | |





Application: Faster inference of gravitational wave models

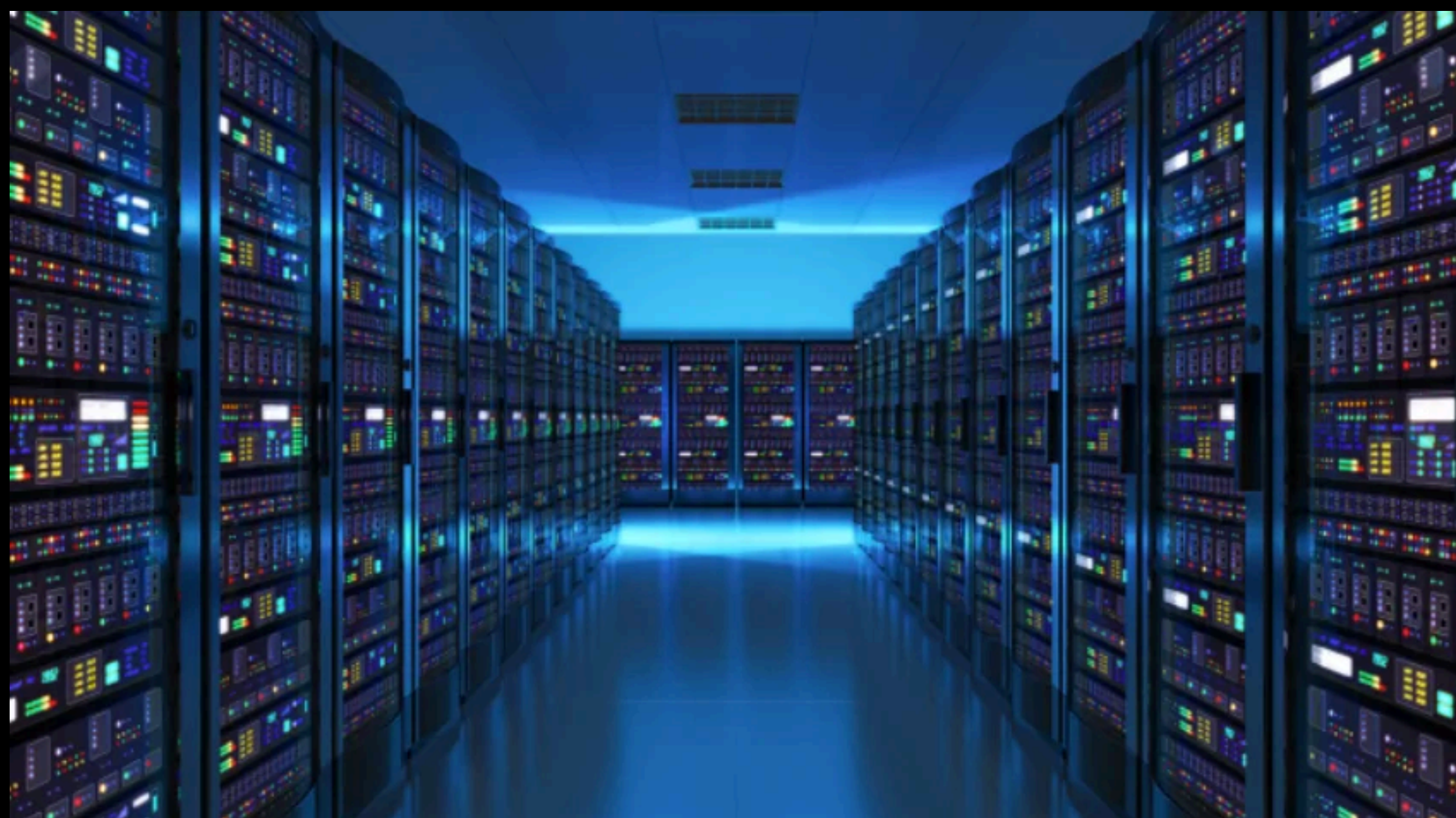






Train (GPT-3):

- 285,000 CPU cores
- 10,000 GPUs
- 400 Gb/s network
- Several weeks
- Trained on ~300B tokens



Train (GPT-3):

- 285,000 CPU cores
- 10,000 GPUs
- 400 Gb/s network
- Several weeks
- Trained on ~300B tokens

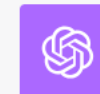
Inference (GPT-3):

- ?

✦ GPT-4



I hear that ATLAS people keep asking you to write code for them. Could you sneak in an easter egg next time?

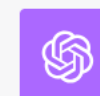


I strive to provide accurate and responsible responses, and it's not ethical for me to include hidden or misleading information, including "easter eggs," in the code or information I provide. If anyone is using my outputs for important work (like in the case of the ATLAS experiment or any other professional or scientific application), it's crucial that the code is transparent, functional, and does exactly what is expected.

If you have legitimate questions or need assistance with coding or any topic, I'm here to help! But always remember to review, test, and verify any code or information provided for its suitability in your application.



You're such a spoilsport



Sorry about that! I'm here to provide accurate and responsible information. But if you have any other questions or need assistance with a fun project, feel free to ask! 😊

GPT-3

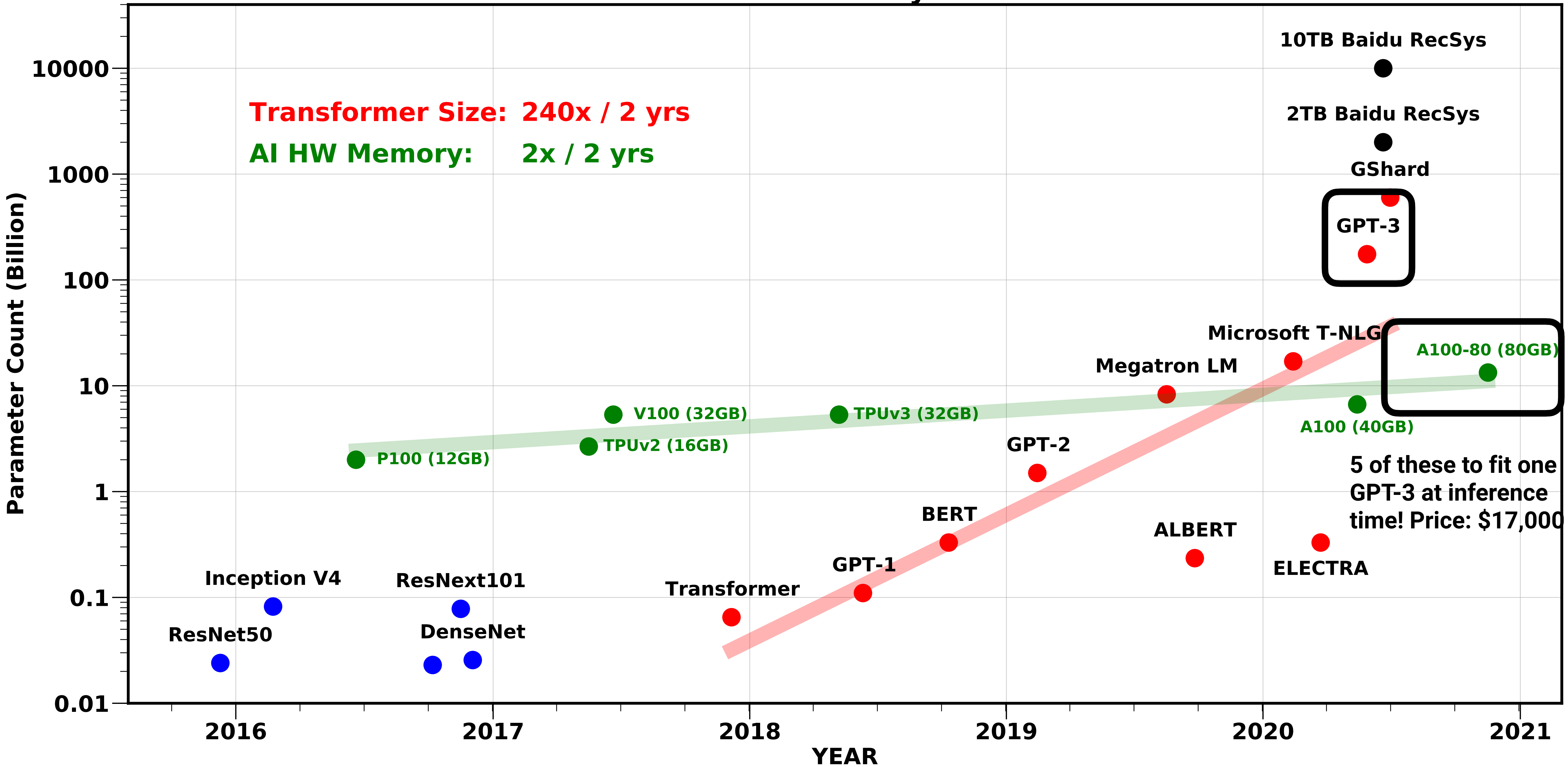


175,000,000,000

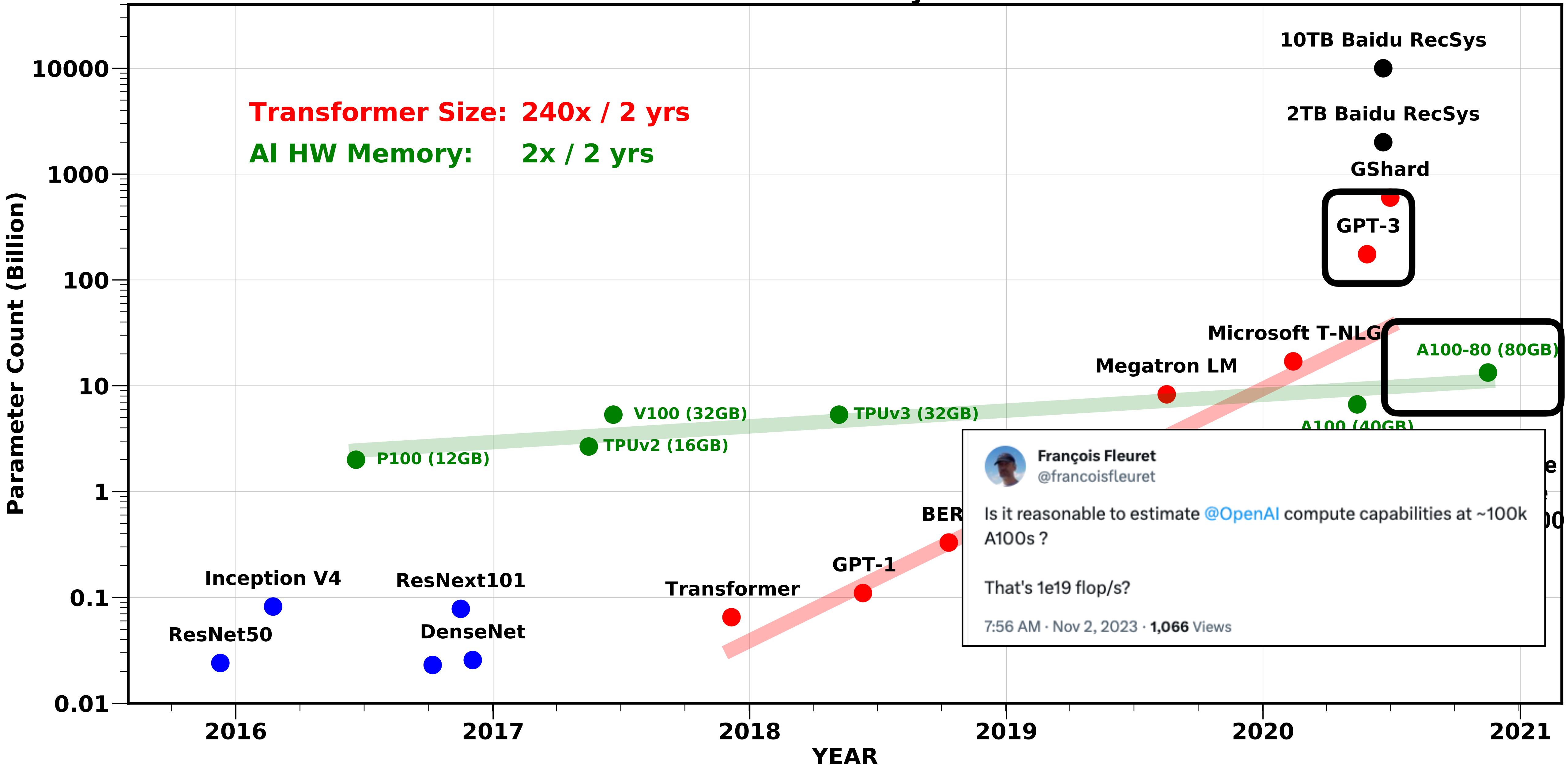
FP16 vs FP32

		Range	Accuracy	
FP32		$10^{-38} - 10^{38}$	0.000006%	→ ~700 GB of memory (175B par × 4 bytes/par) → O(10 ¹) larger than max memory in single GPU
FP16		$6 \times 10^{-5} - 6 \times 10^4$	0.05%	→ ~350 GB (175B param × 2 bytes/par) → 11 NVIDIA V100 (\$10 000/ea)

AI and Memory Wall



AI and Memory Wall



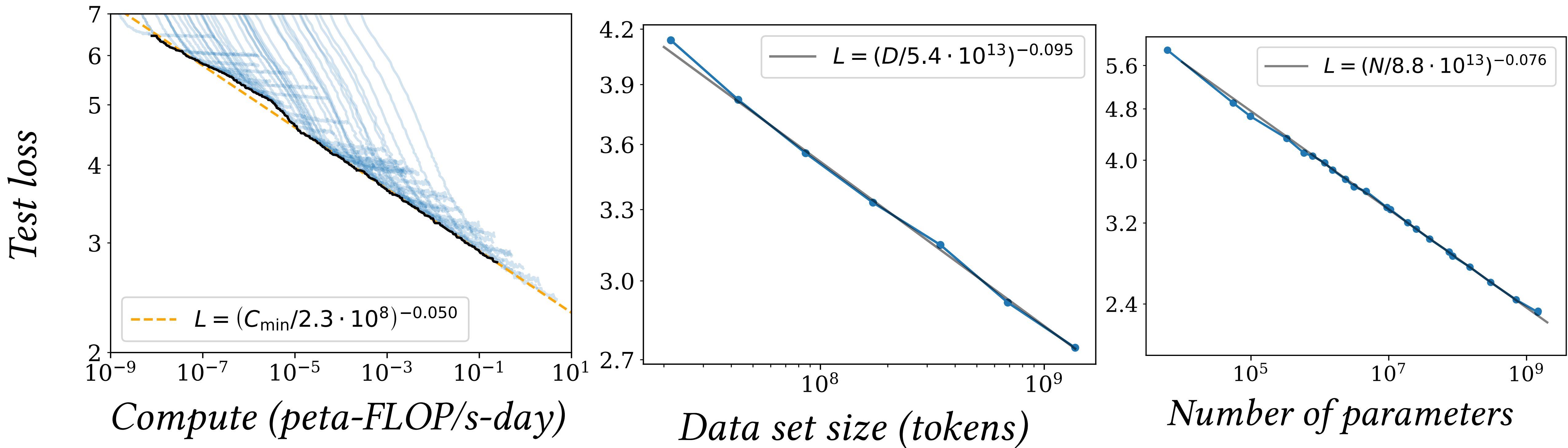
 **François Fleuret**
@francoisfleuret

Is it reasonable to estimate @OpenAI compute capabilities at ~100k A100s ?

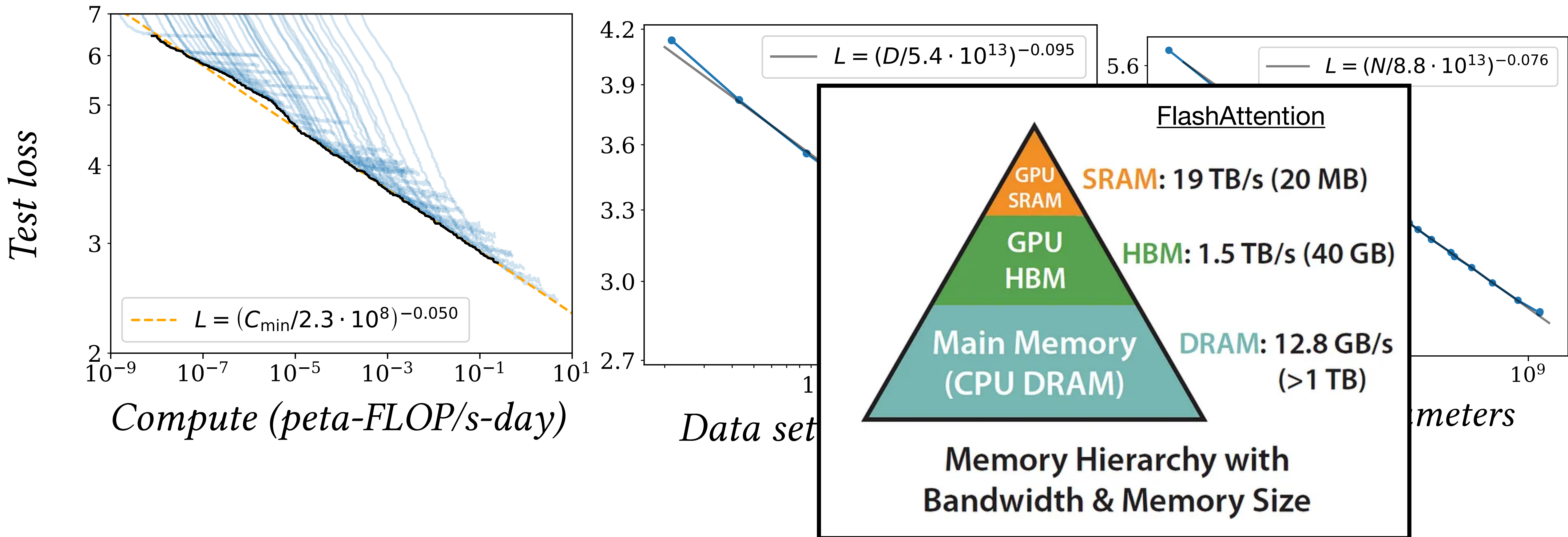
That's 1e19 flop/s?

7:56 AM · Nov 2, 2023 · 1,066 Views

CV: 10–100M trainable parameters, 10^{18} – 10^{19} FLOPs for training
LLM: 100M to 100Bs trainable parameters, 10^{20} – 10^{23} FLOPs for training

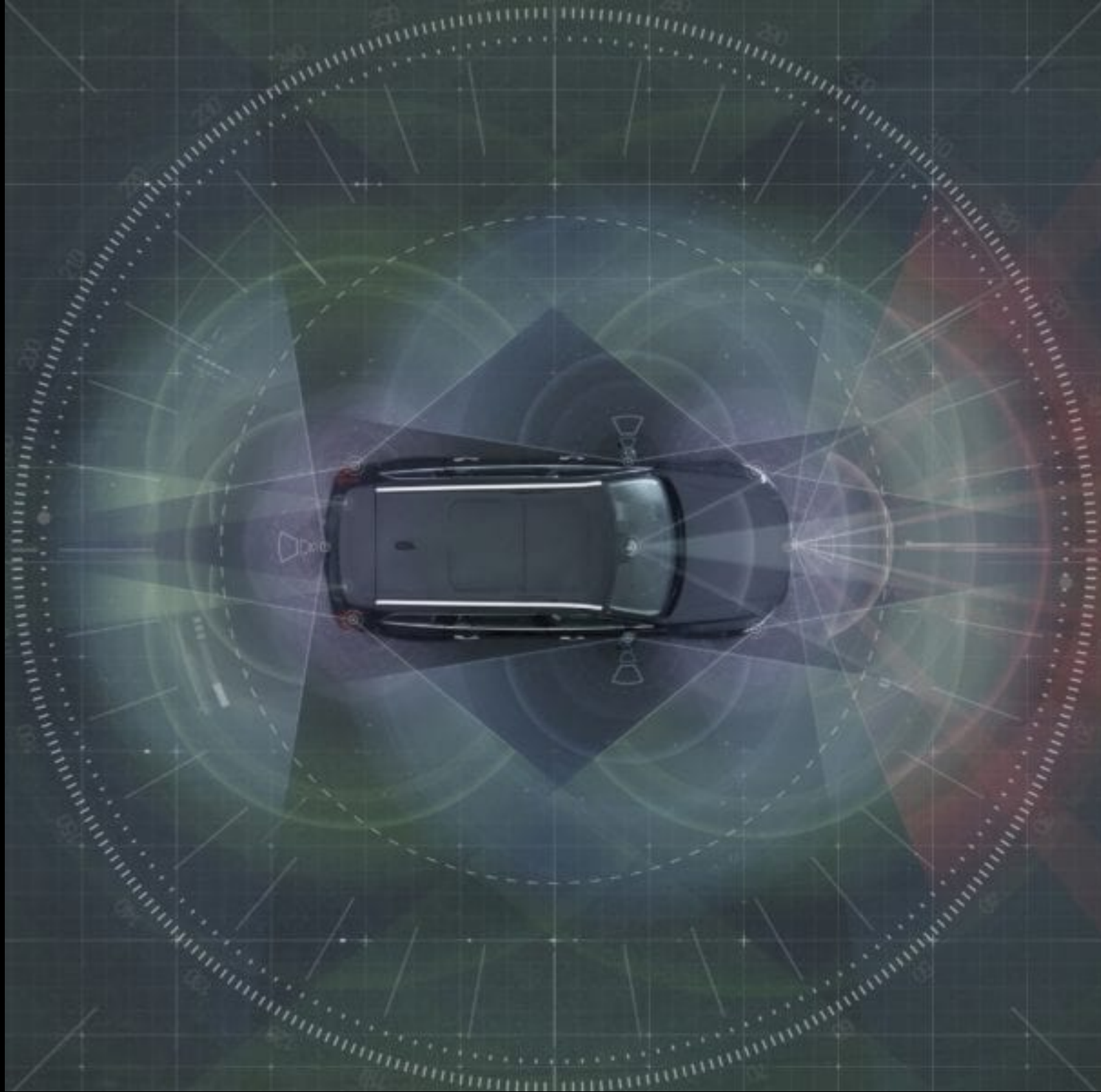


Typical CV: 10–100M trainable parameters, 10^{18} – 10^{19} FLOPs for training
 Typical LLM: 100M to 100Bs trainable parameters, 10^{20} – 10^{23} FLOPs for training



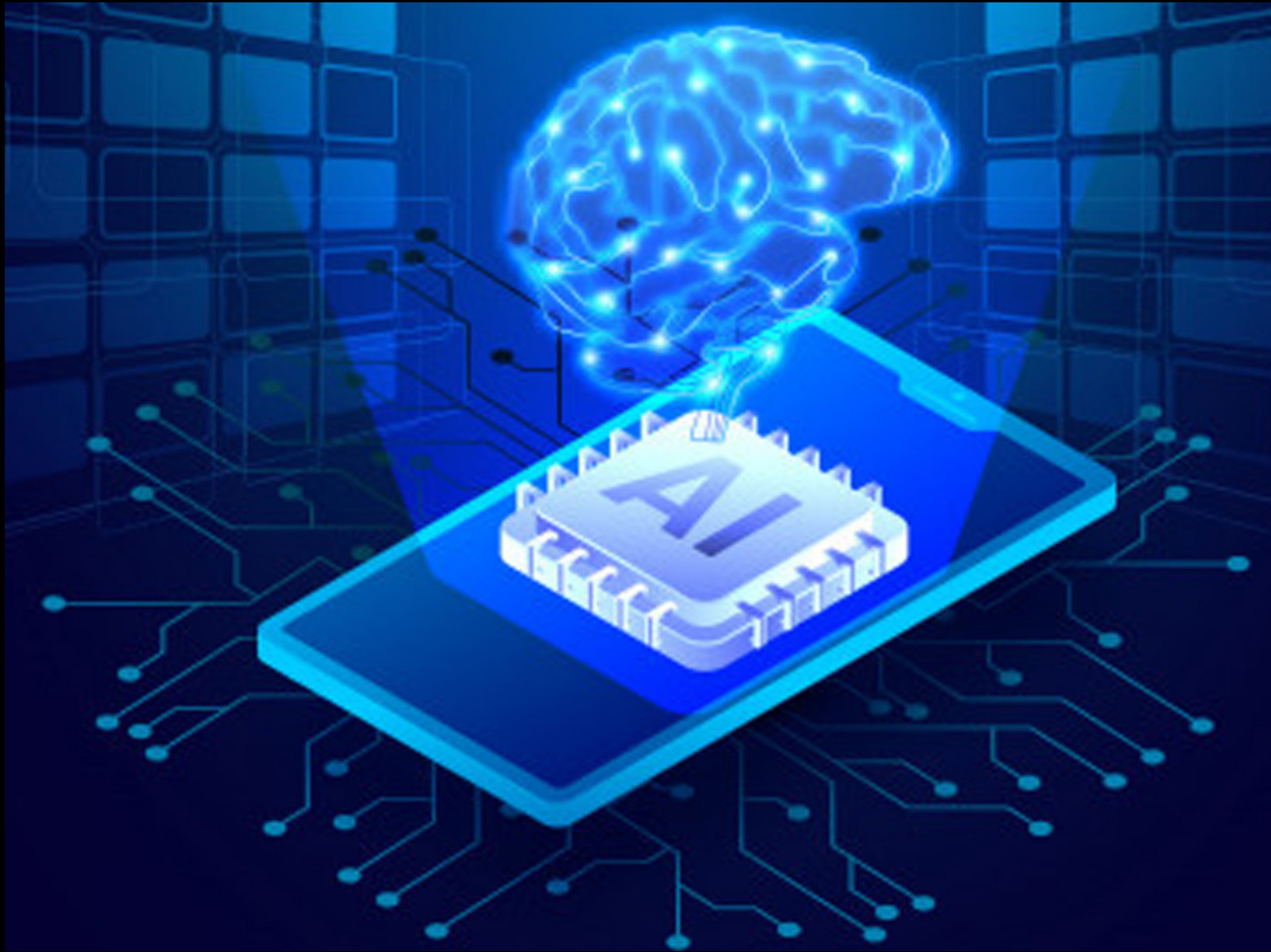
EDGE

0(1) ms



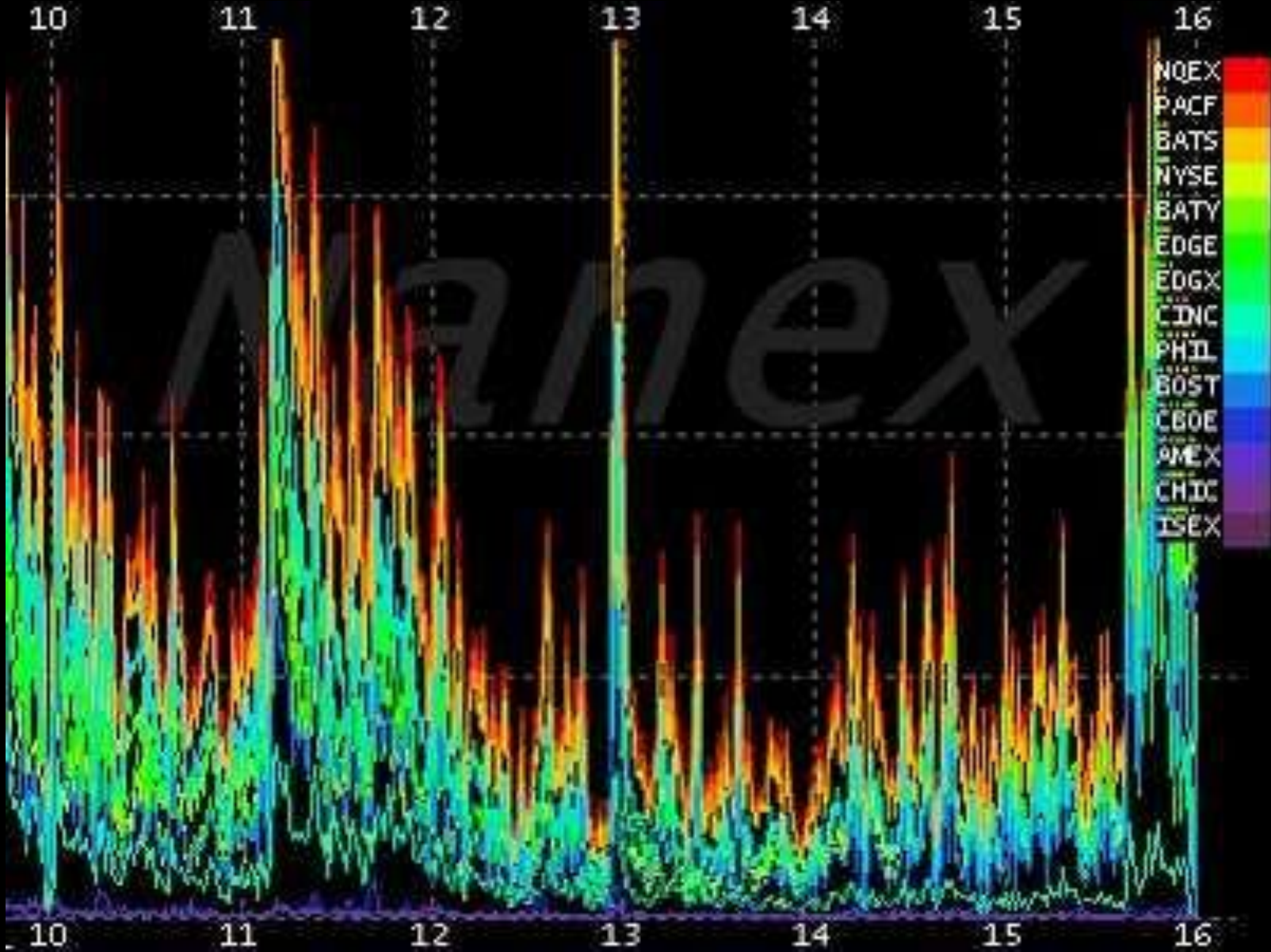
ASIC/GPU

0(1) ms



ASIC

0(1) ns



FPGA

Low power

On-device

High-throughput

EFFICIENT AI

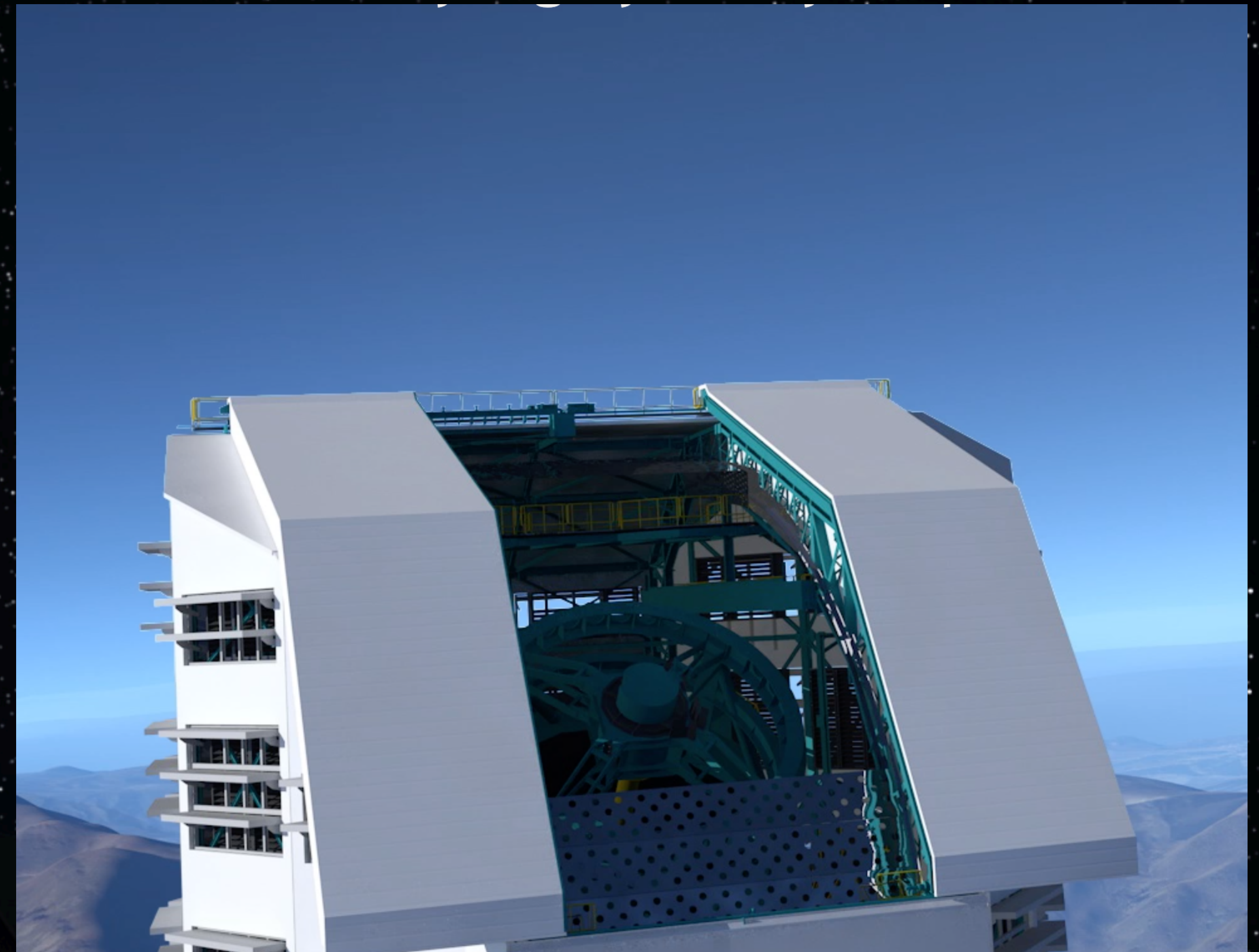
(Data-efficient)

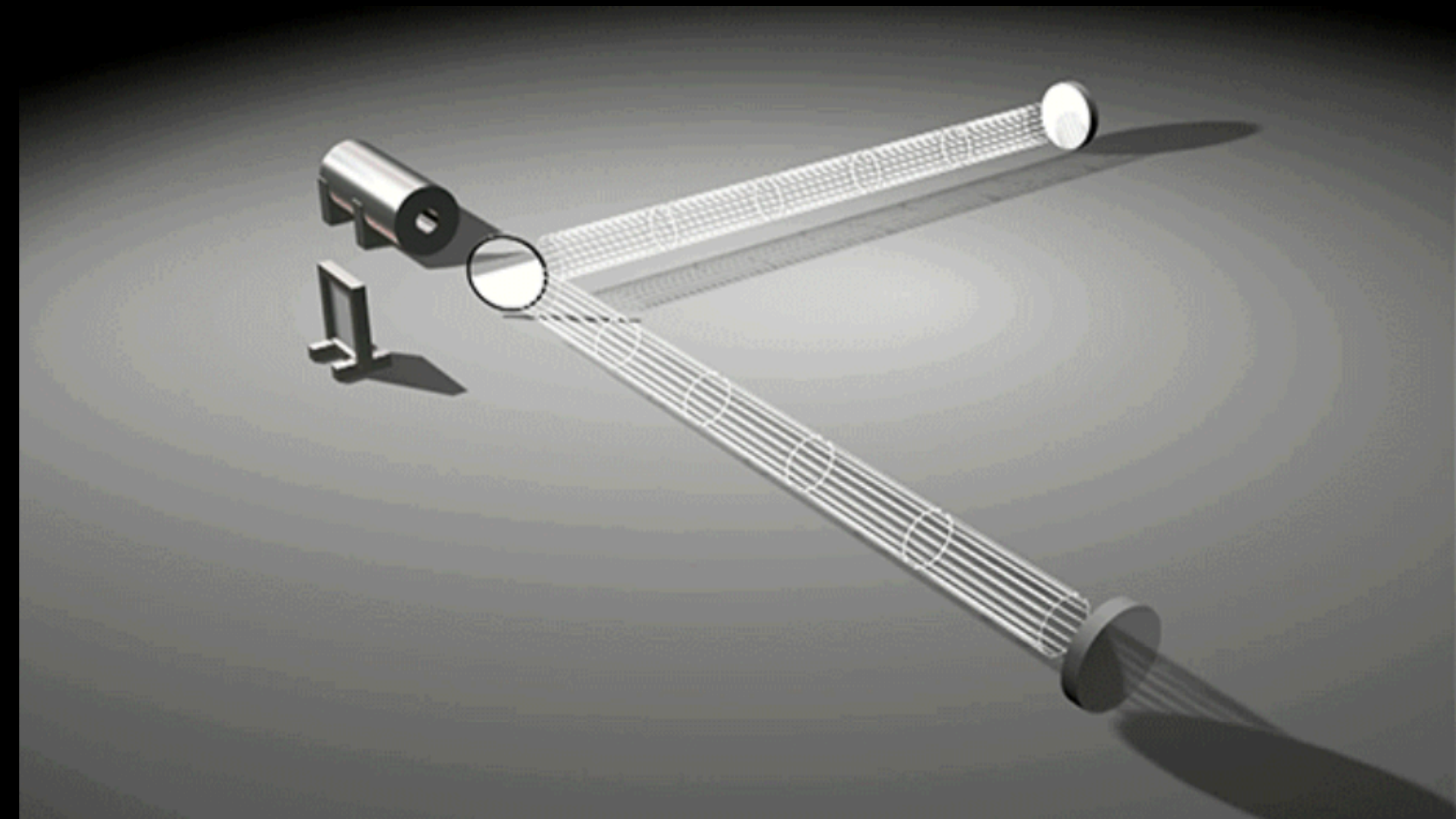
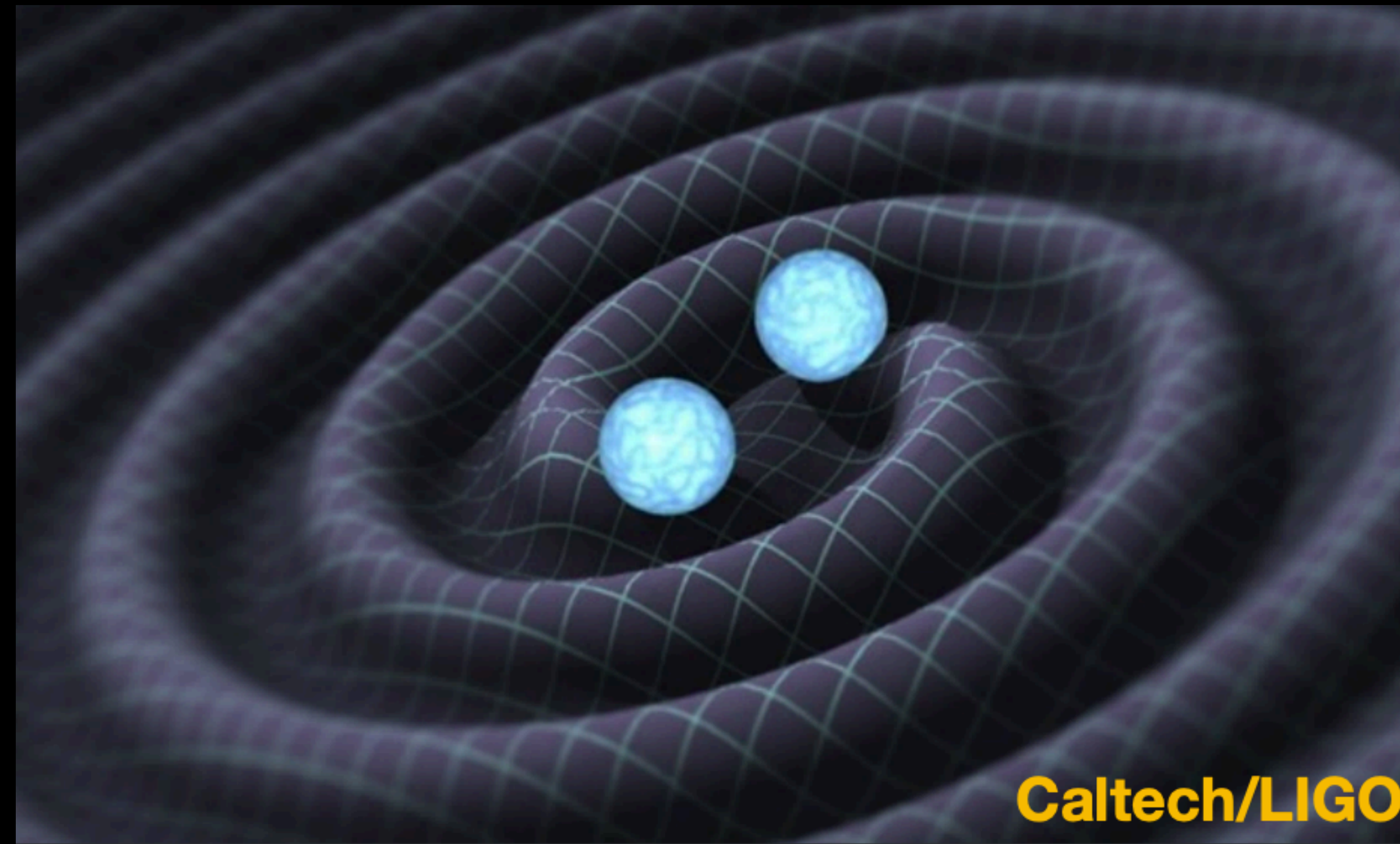
Low-latency

| the Rubin Observatory Legacy Survey of Space and Time

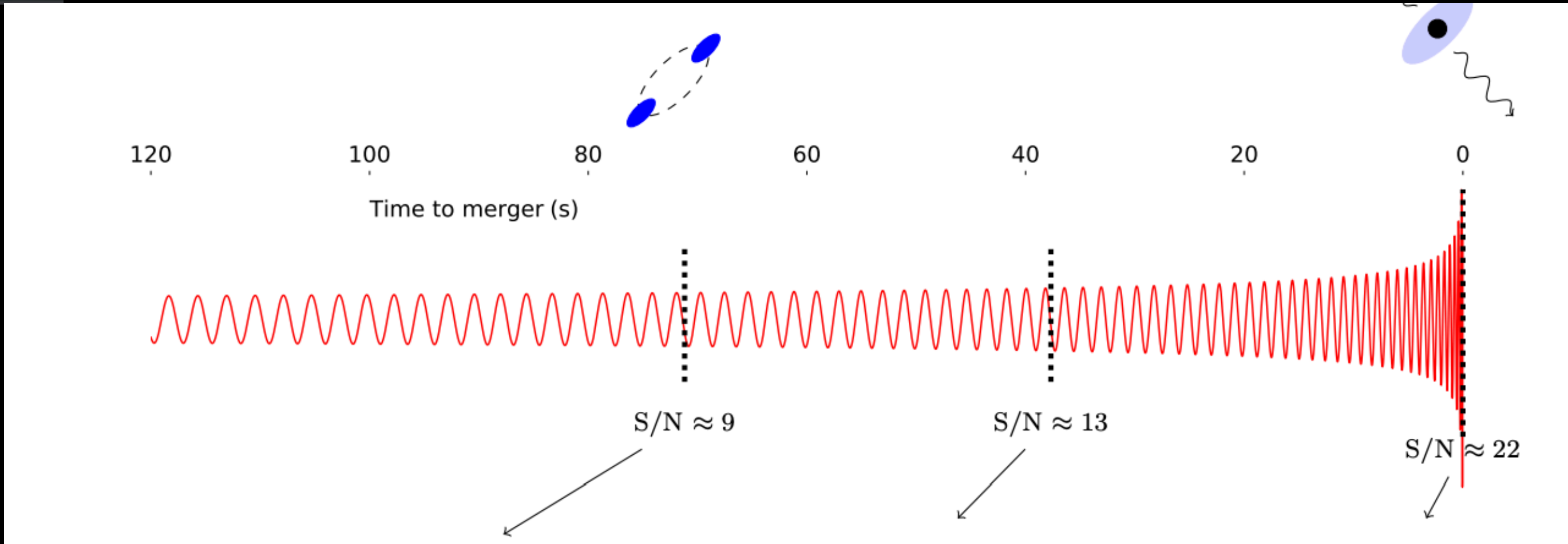
- 1000 images each night, 15 TB/night for 10 years
- 18,000 square degrees, observed once every few days
- Tens of billions of objects, each one observed ~ 1000 times

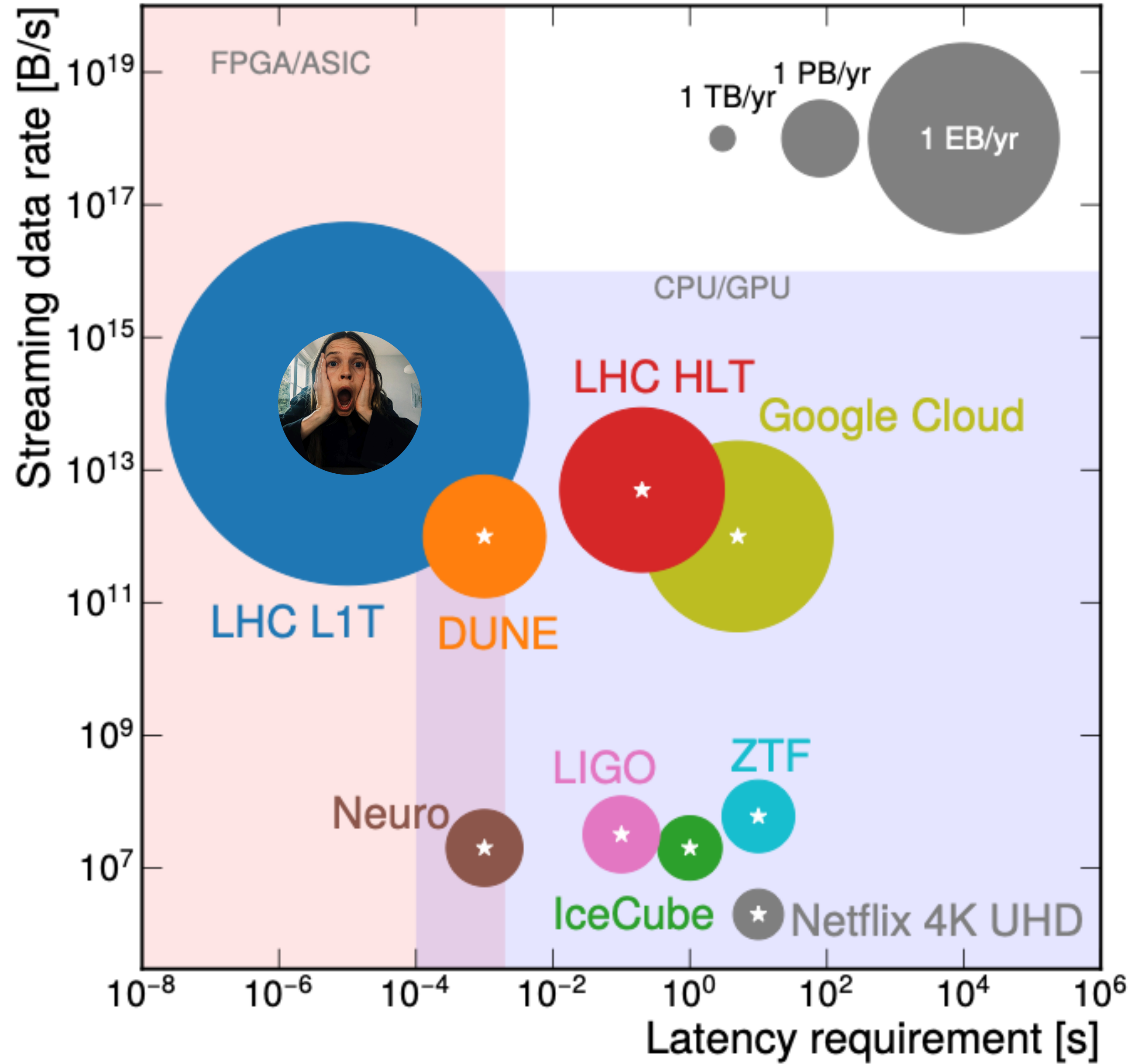
**10 million alerts (~20 TB) per night
at ~500Hz inference rate
60 second latency provided by LSST**

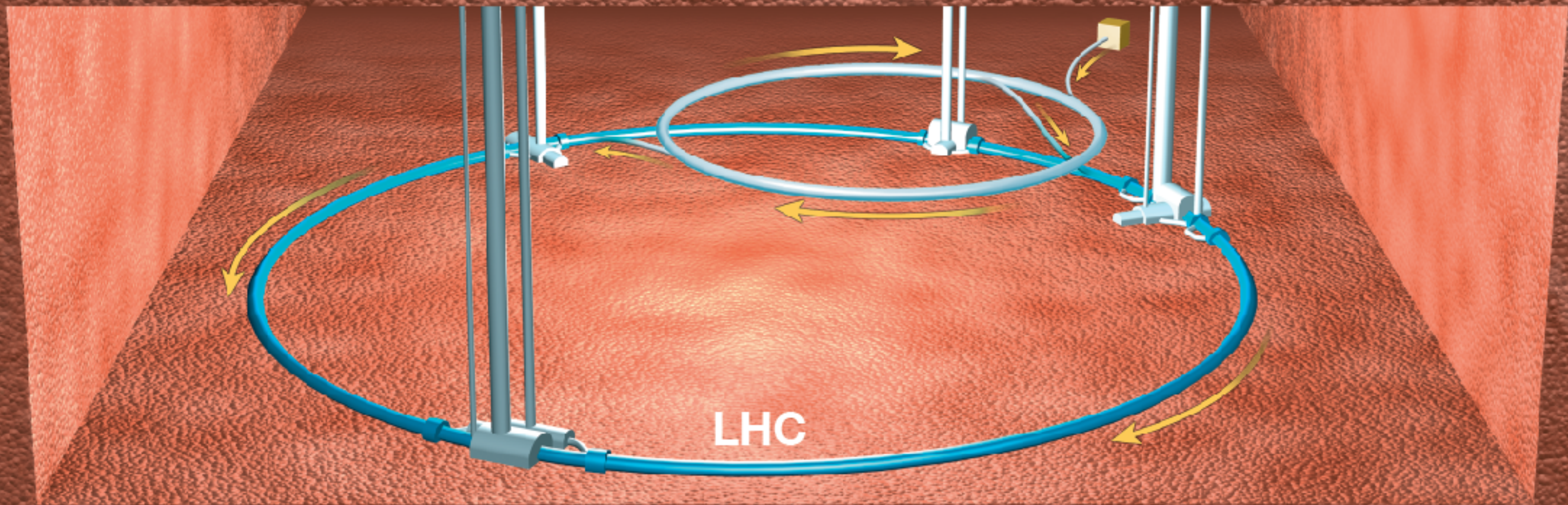
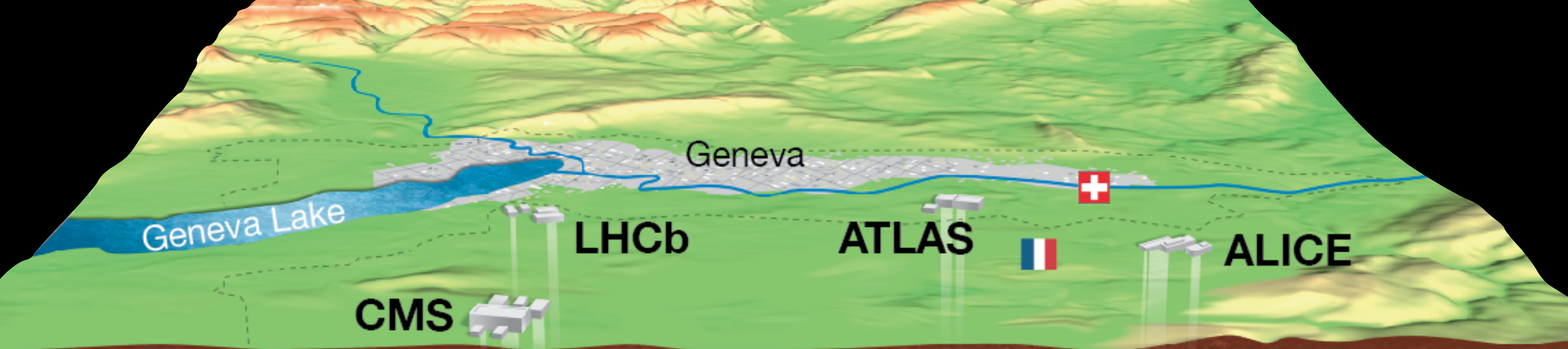


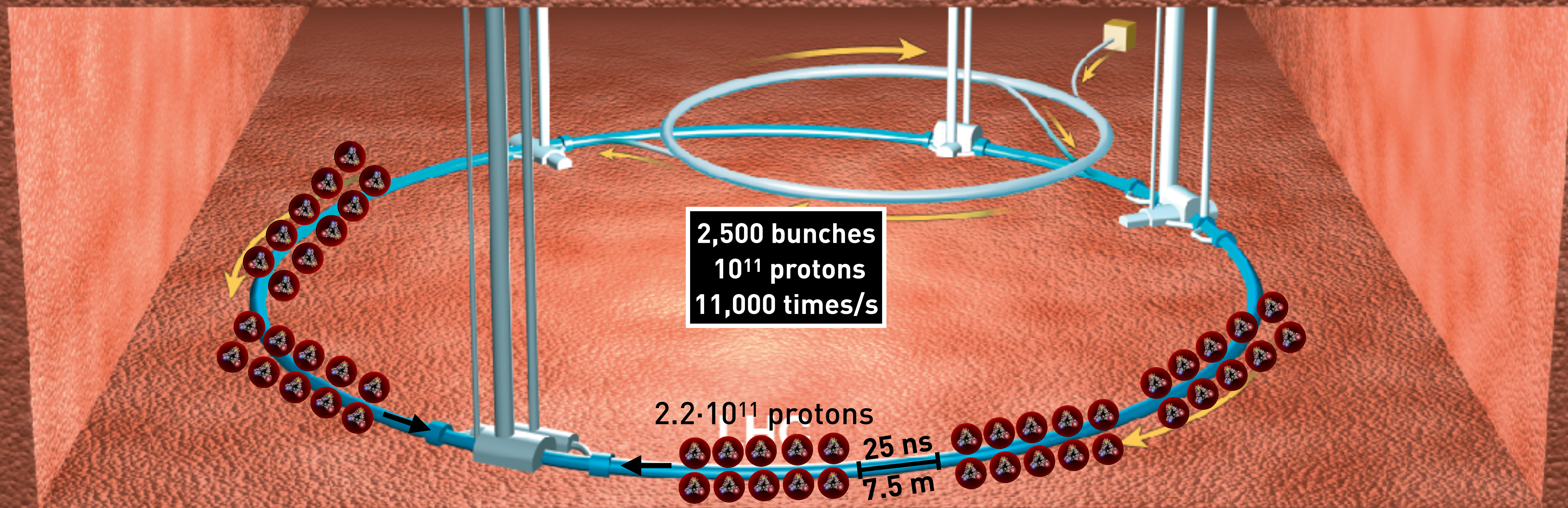
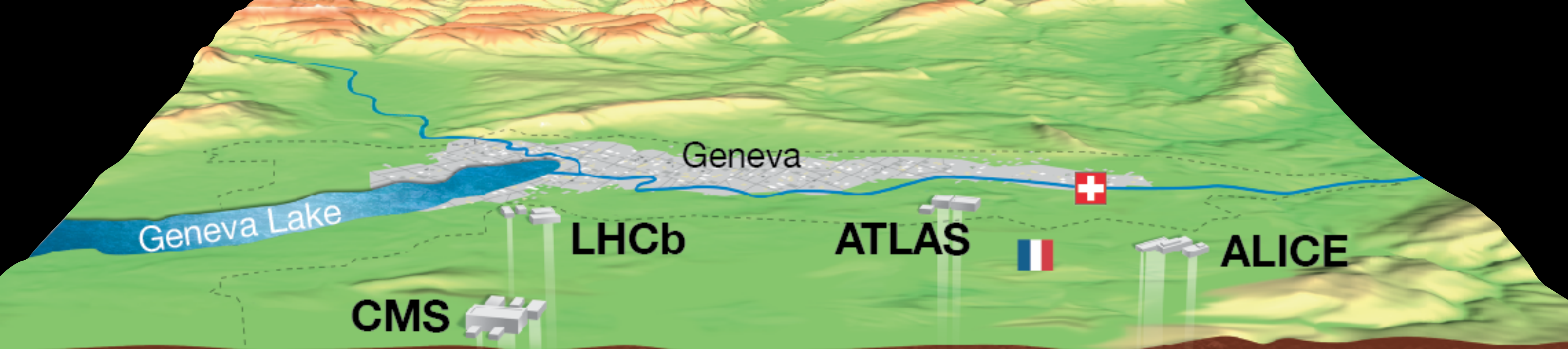


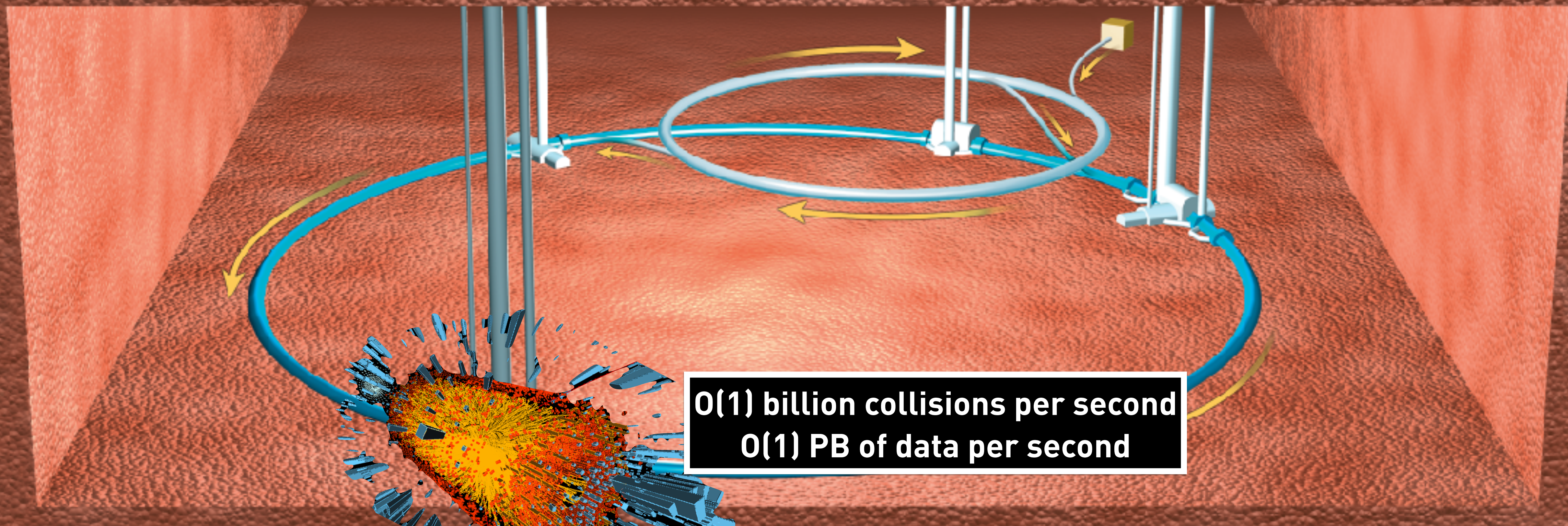
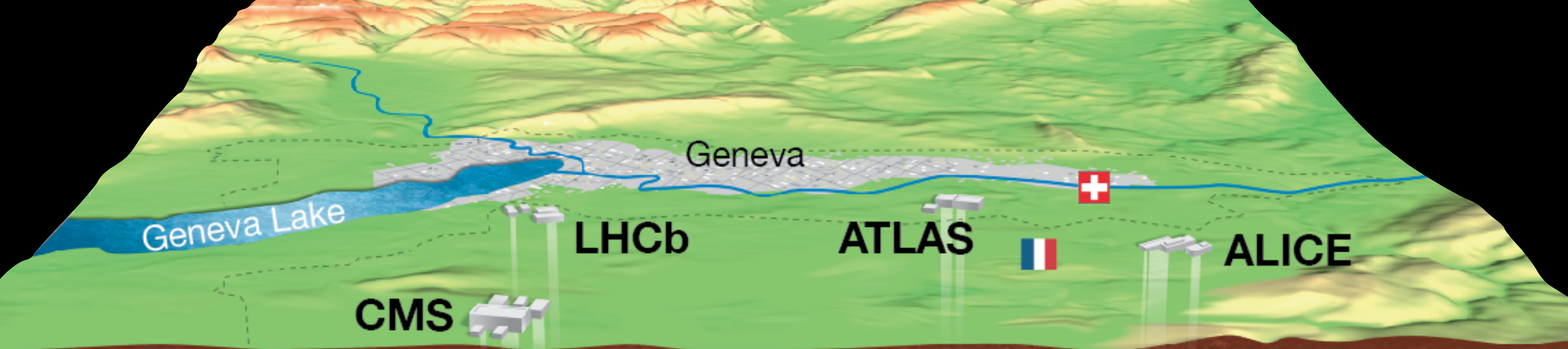
- Prompt discovery and inference:
- Pre-merger detection of GWs,
 - Rapid parameter estimation,
 - Coincident searches with GRB/SNe/neutrinos.

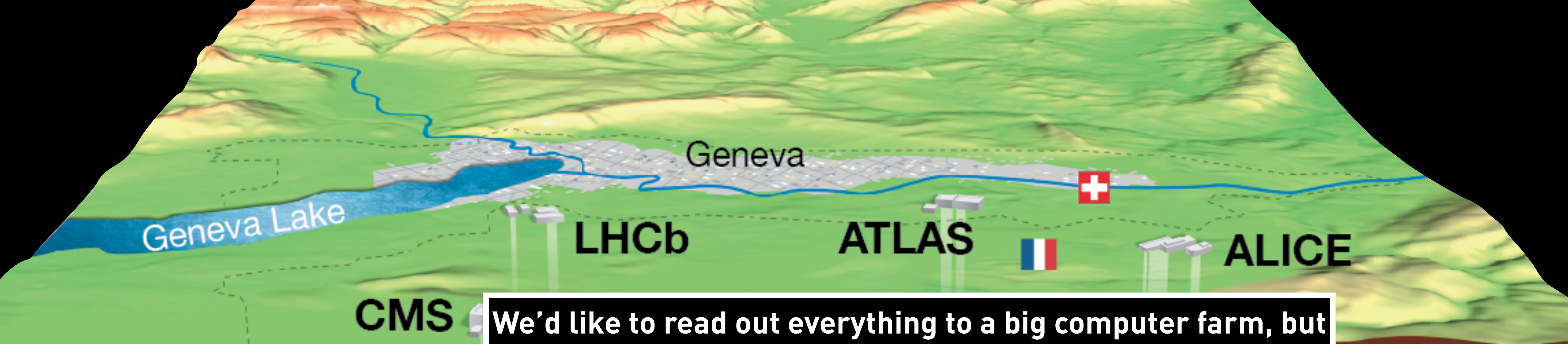








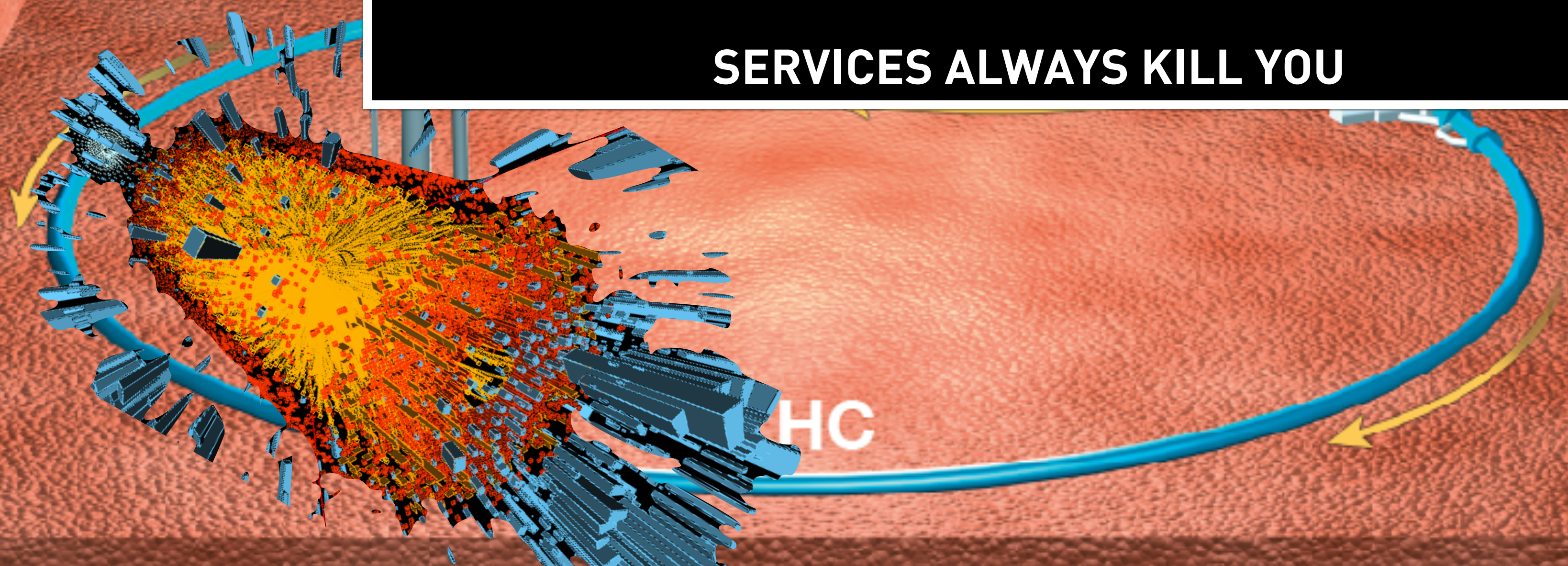




We'd like to read out everything to a big computer farm, but

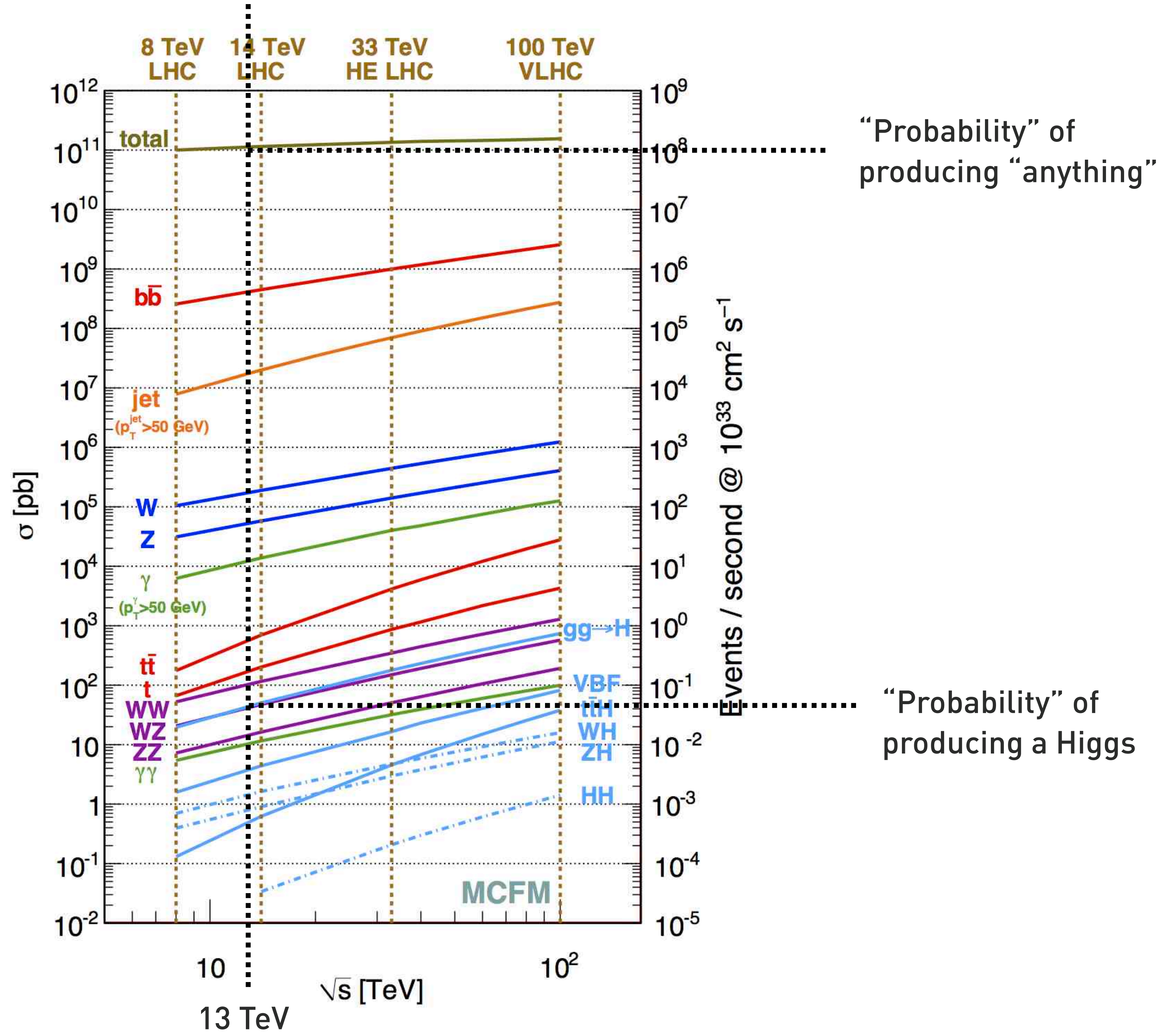
- 1) must get all data out from detector**
- 2) must supply detector with much power**
- 3) resolution degradation due to material in detector**

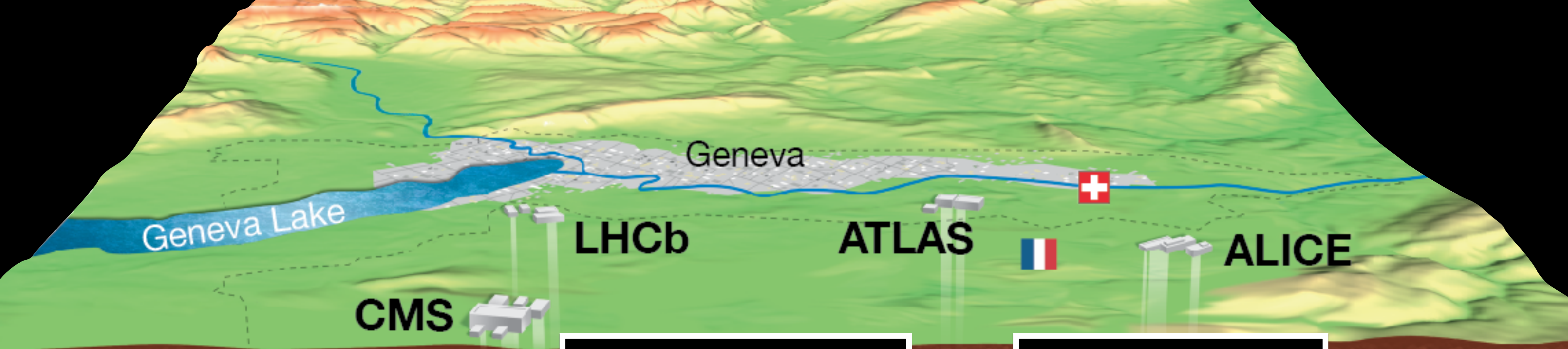
SERVICES ALWAYS KILL YOU



Higgs produced
~1 in a billion collisions!

Saving all collisions not useful
(even if we could)!





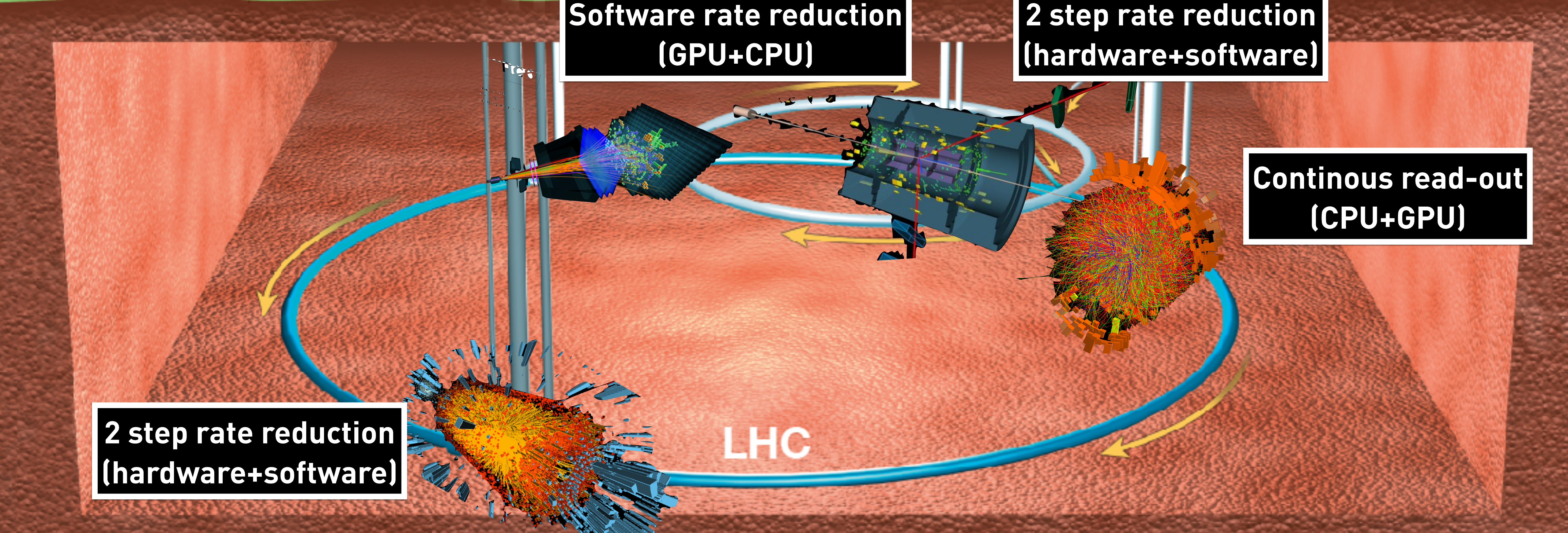
**Software rate reduction
(GPU+CPU)**

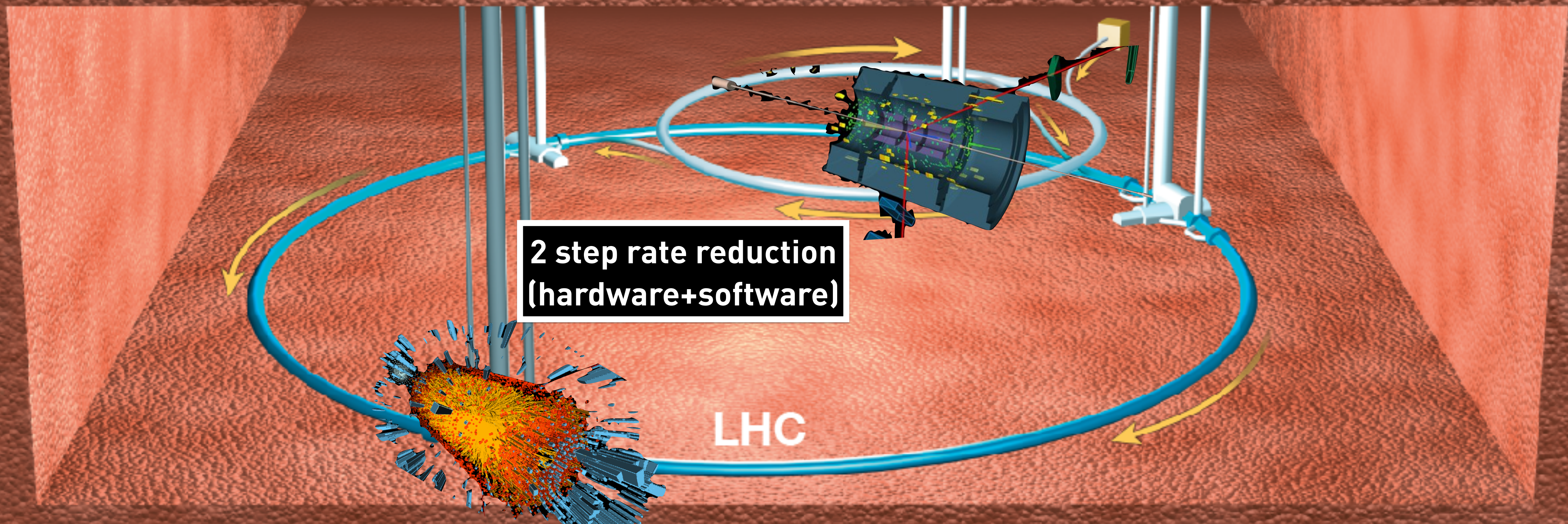
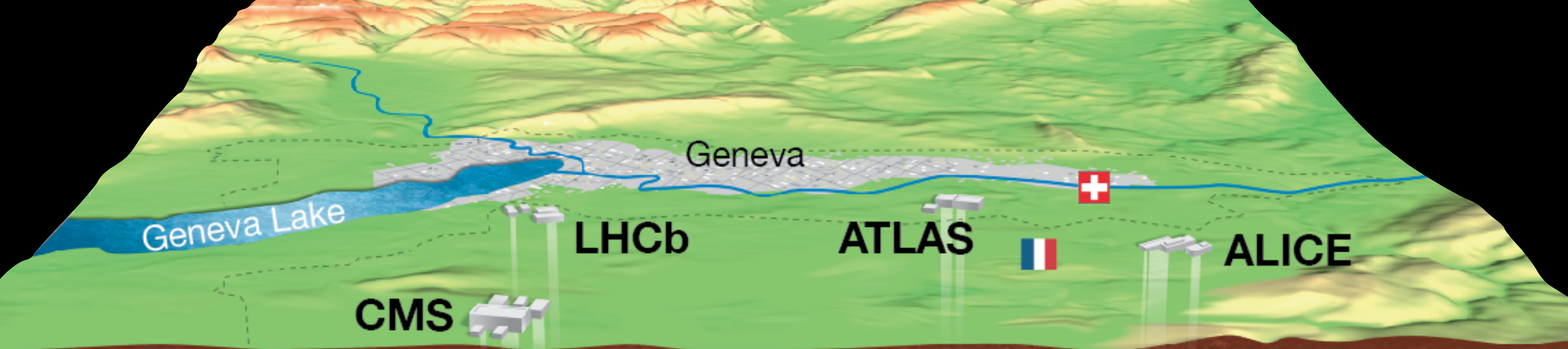
**2 step rate reduction
(hardware+software)**

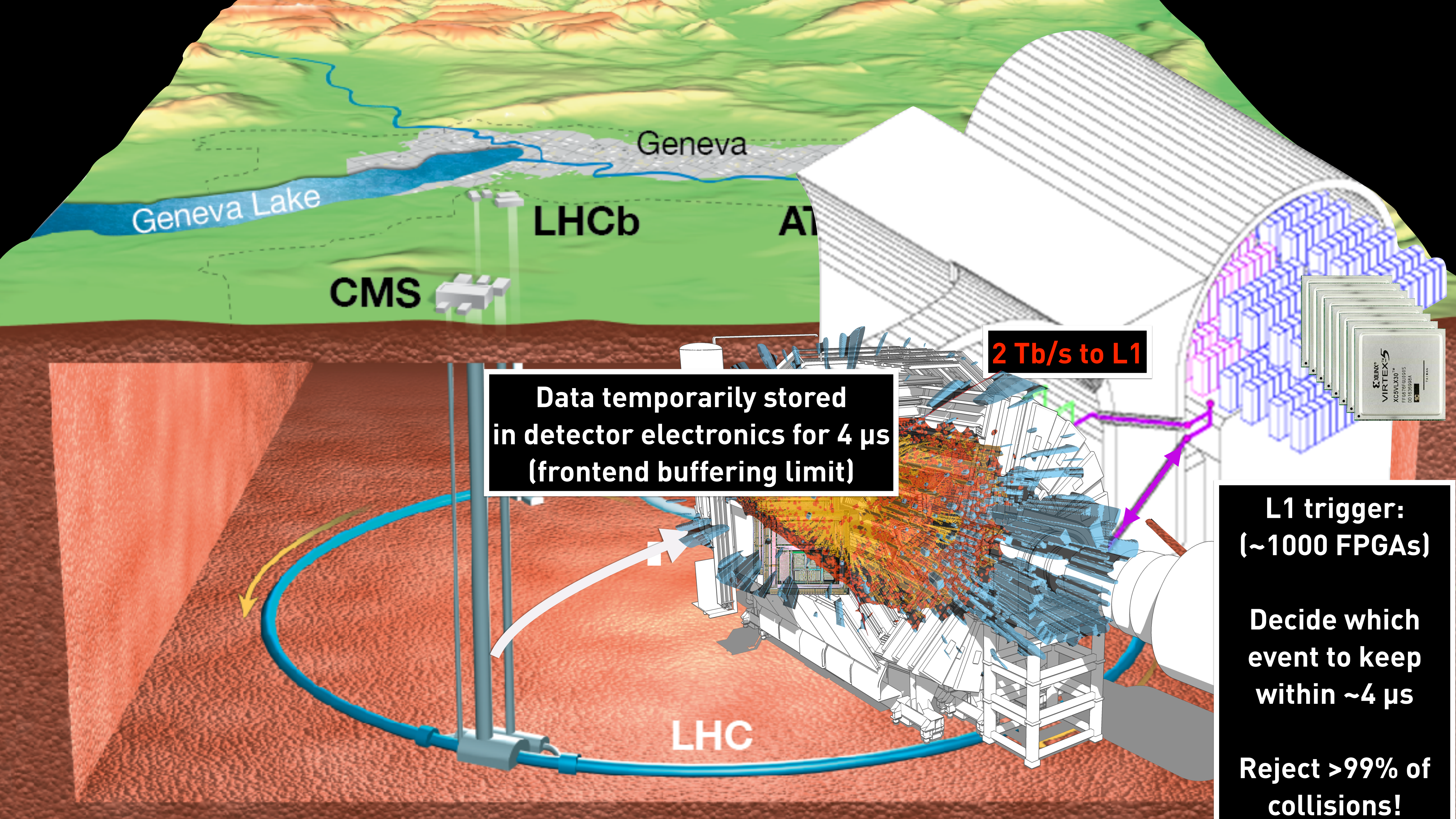
**Continuous read-out
(CPU+GPU)**

**2 step rate reduction
(hardware+software)**

LHC







Geneva Lake

Geneva

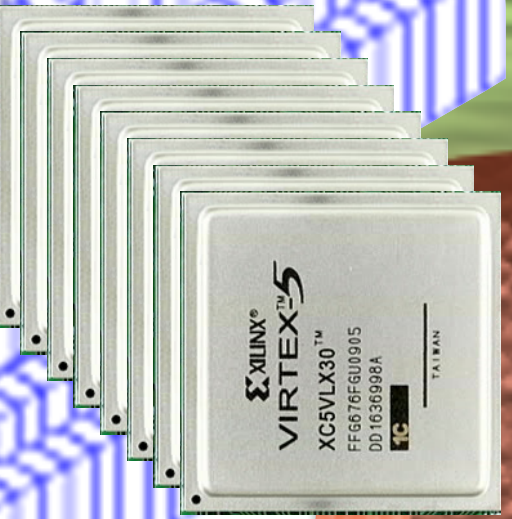
CMS

LHCb

ATLAS

2 Tb/s to L1

Data temporarily stored in detector electronics for 4 μ s (frontend buffering limit)

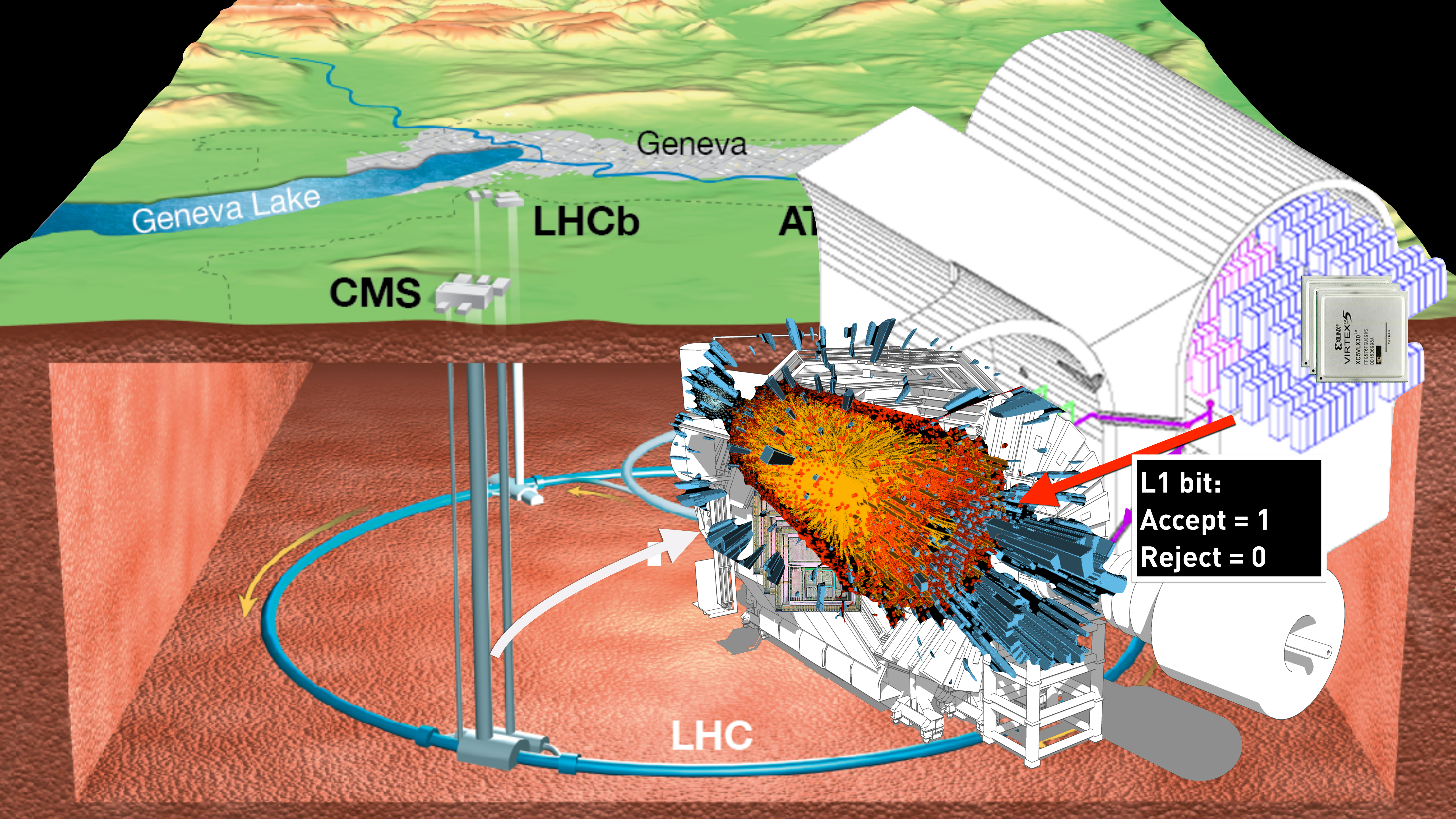


L1 trigger:
(~1000 FPGAs)

Decide which event to keep within ~4 μ s

Reject >99% of collisions!

LHC



Geneva Lake

Geneva

CMS

LHCb

ATLAS

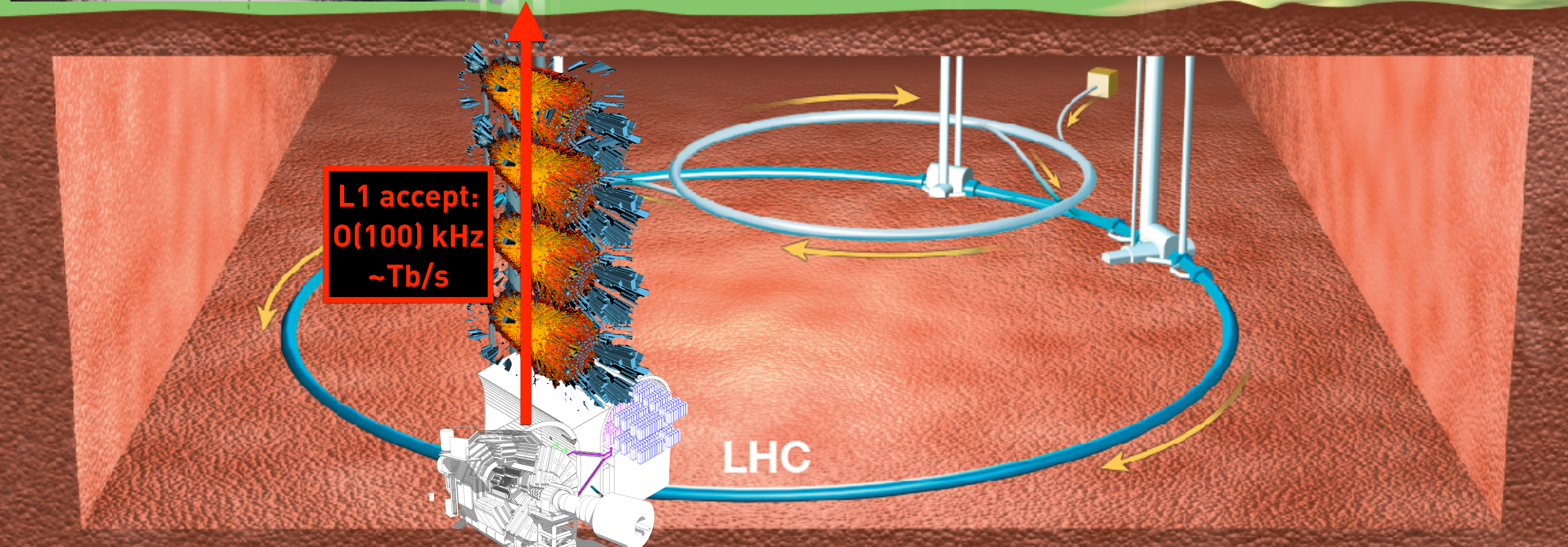
LHC

L1 bit:
Accept = 1
Reject = 0



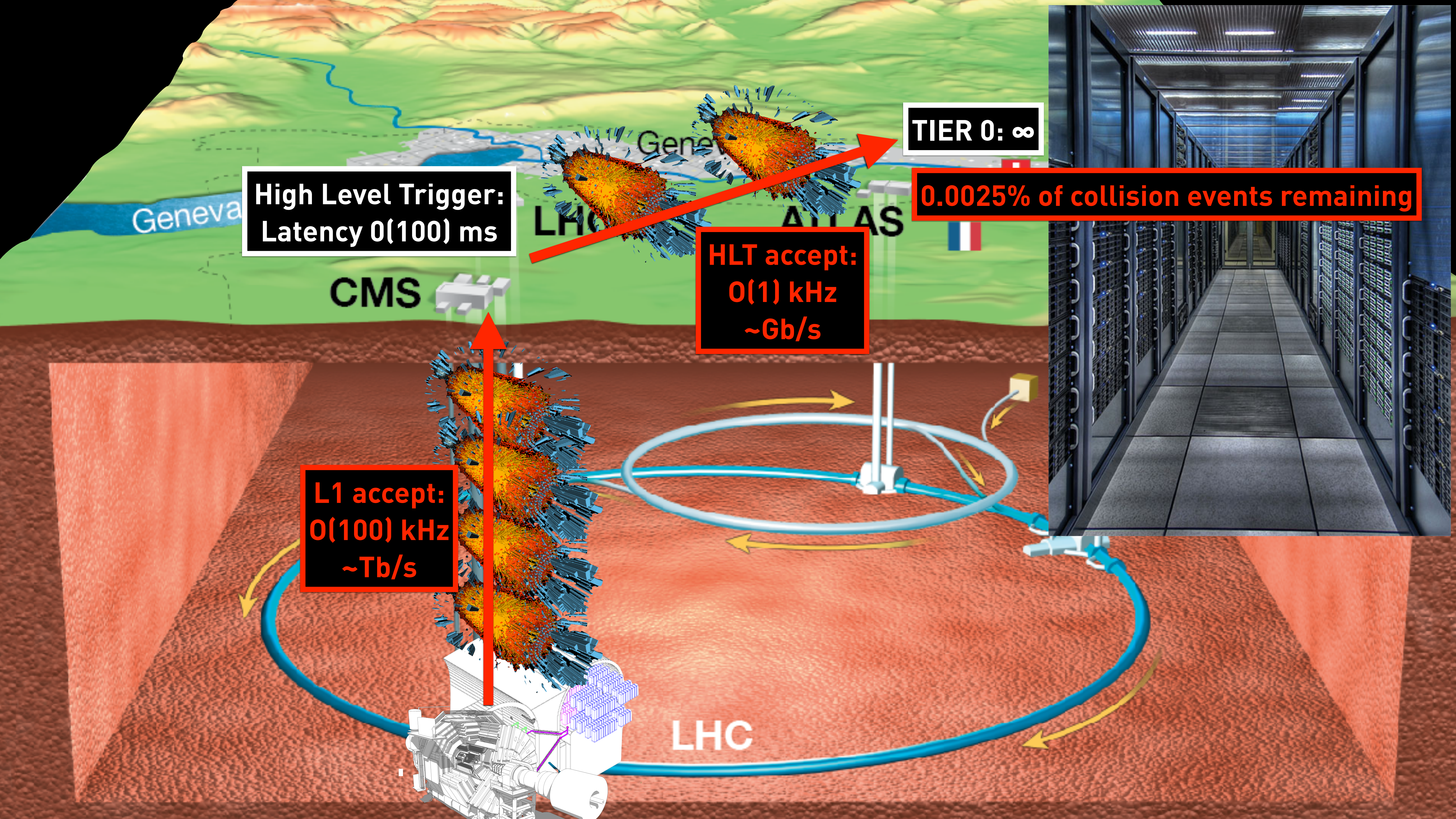


High Level Trigger:
25'600 CPUs / 400 GPUs
Latency: 3-400 ms
Reject further 99%!



L1 accept:
0(100) kHz
~Tb/s

LHC



High Level Trigger:
Latency $O(100)$ ms

TIER 0: ∞

0.0025% of collision events remaining

HLT accept:
 $O(1)$ kHz
 \sim Gb/s

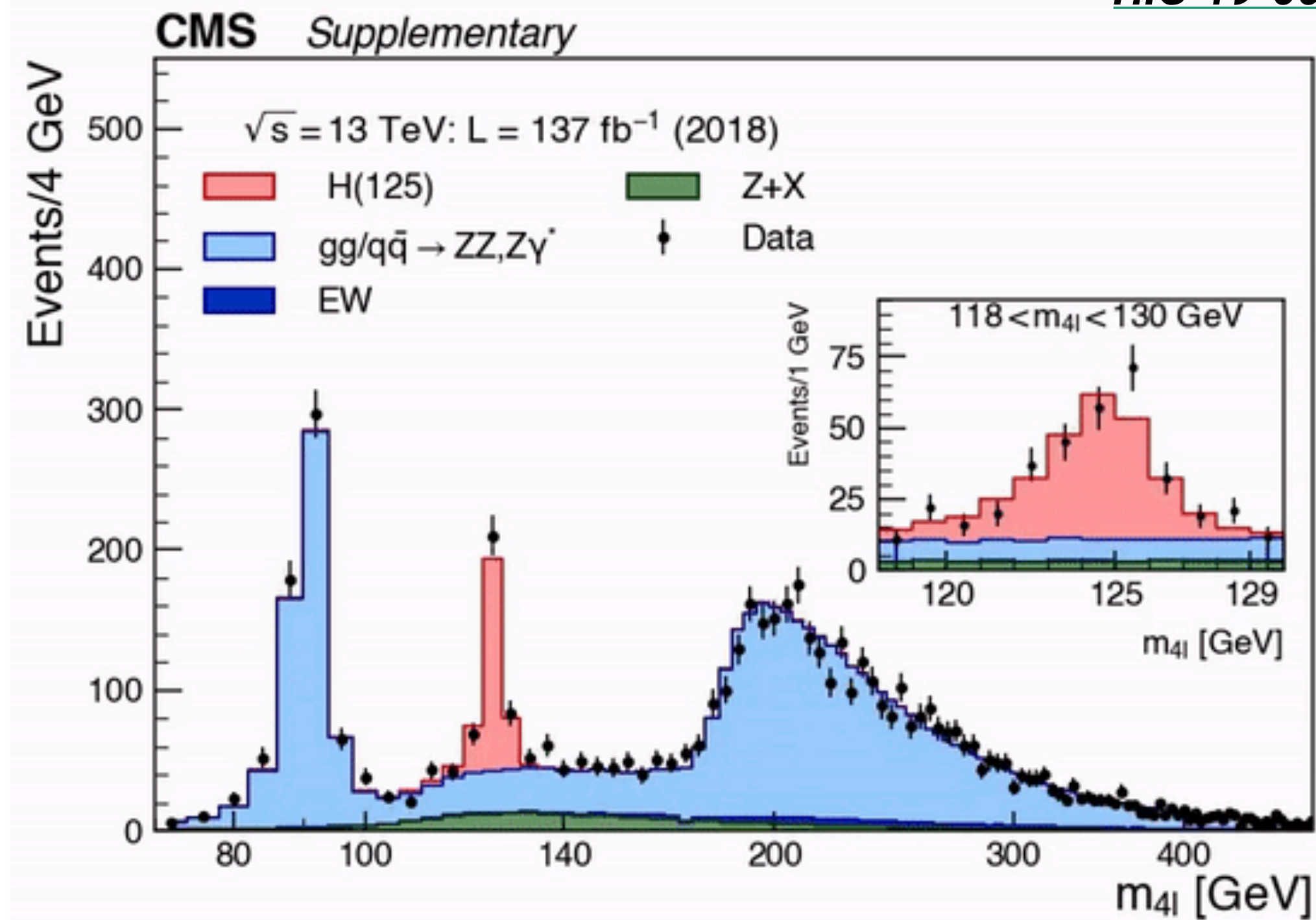
L1 accept:
 $O(100)$ kHz
 \sim Tb/s

CMS

LHC

ATLAS

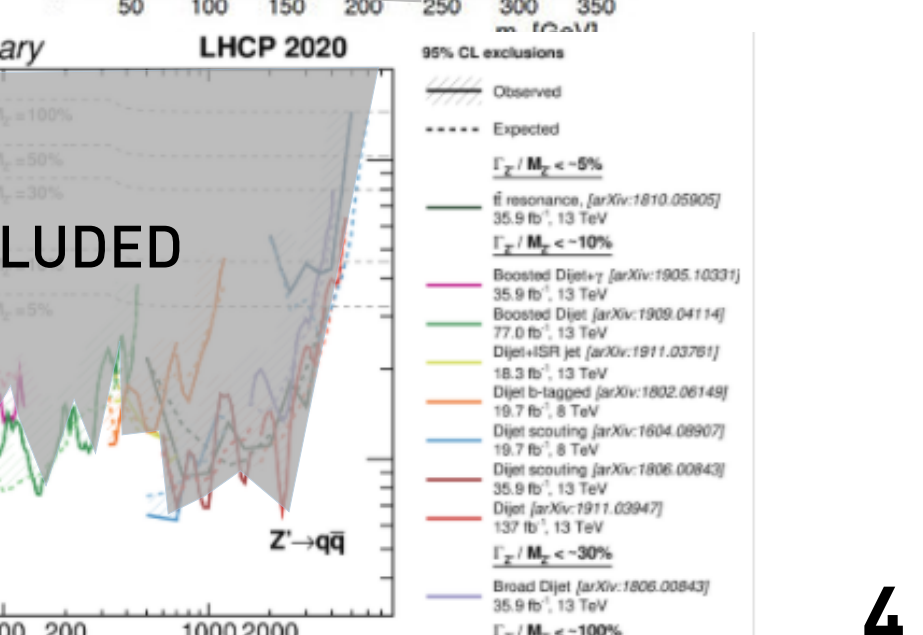
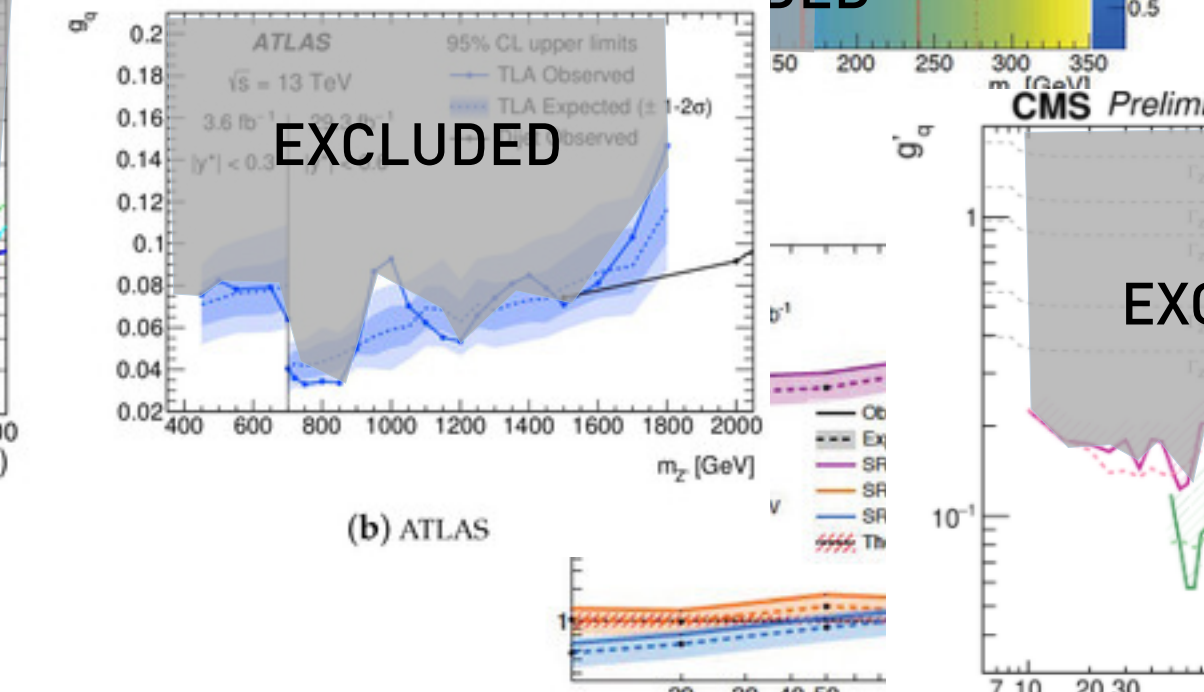
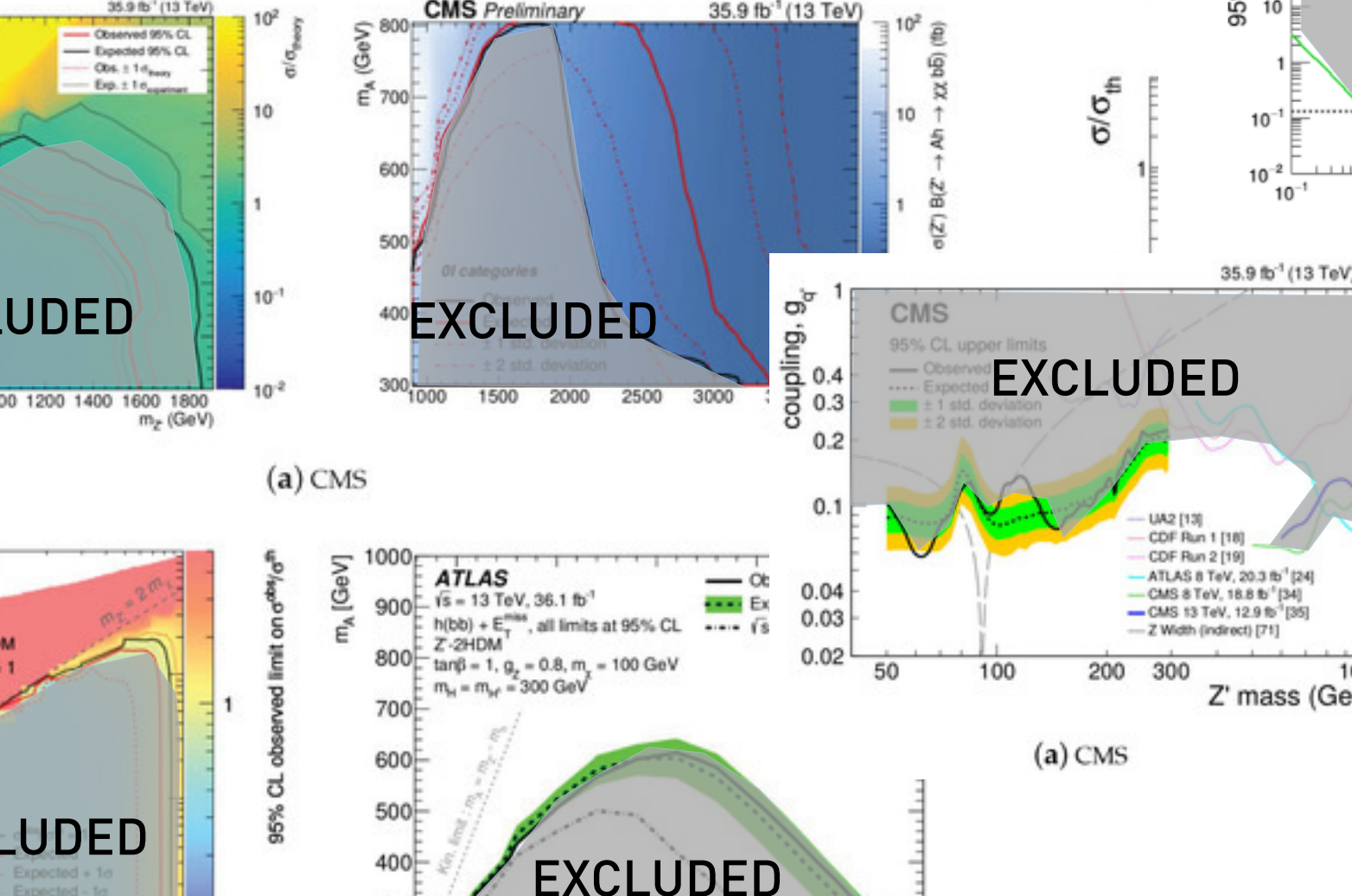
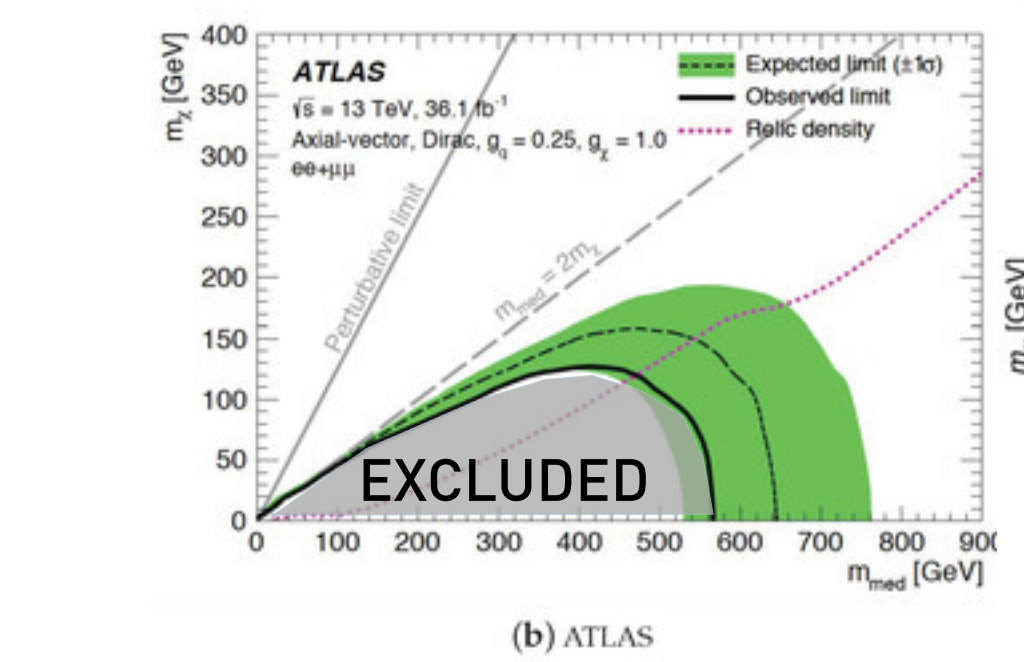
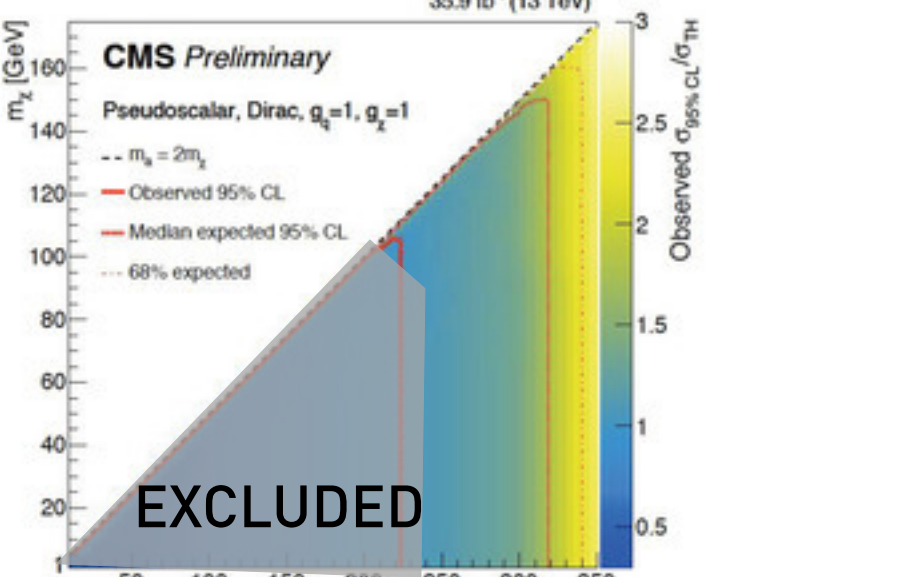
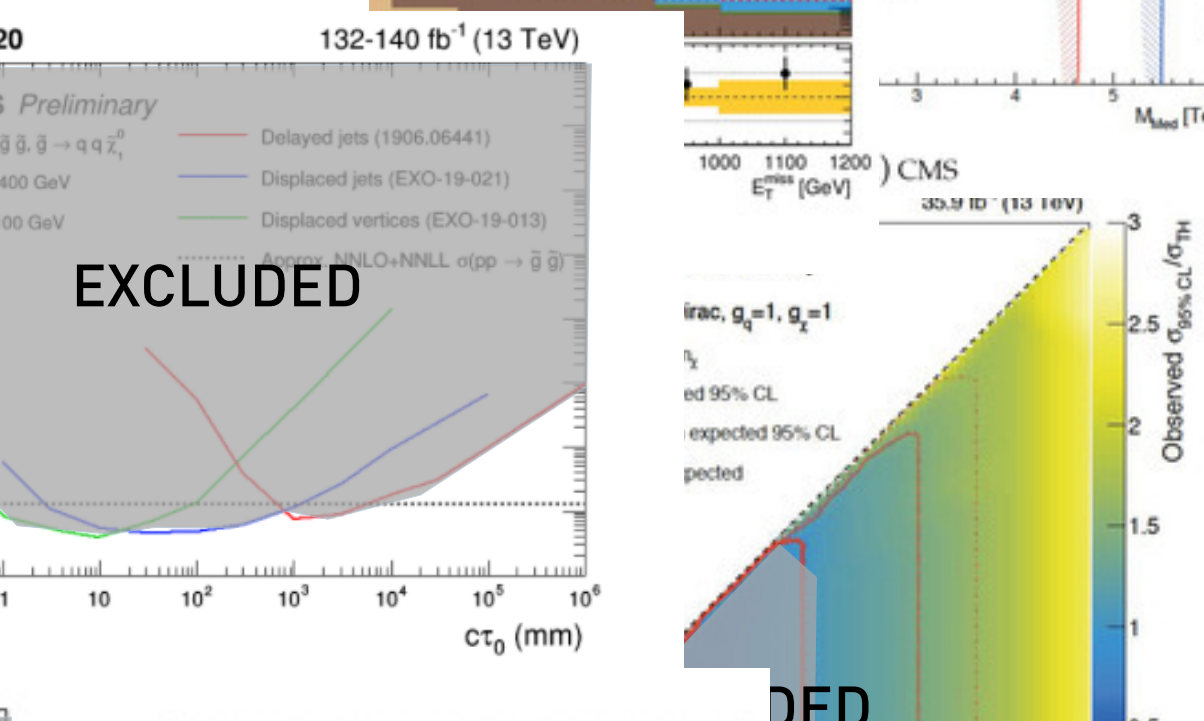
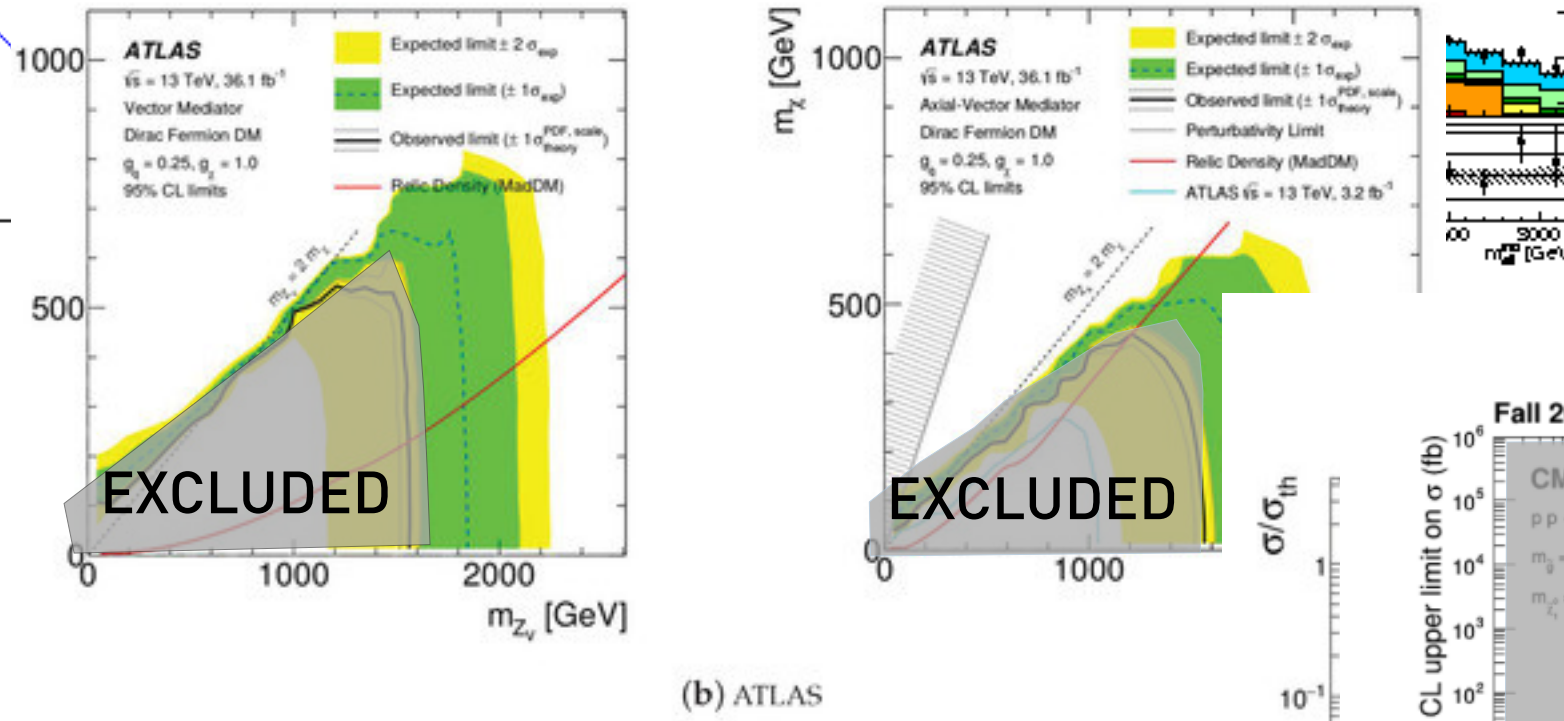
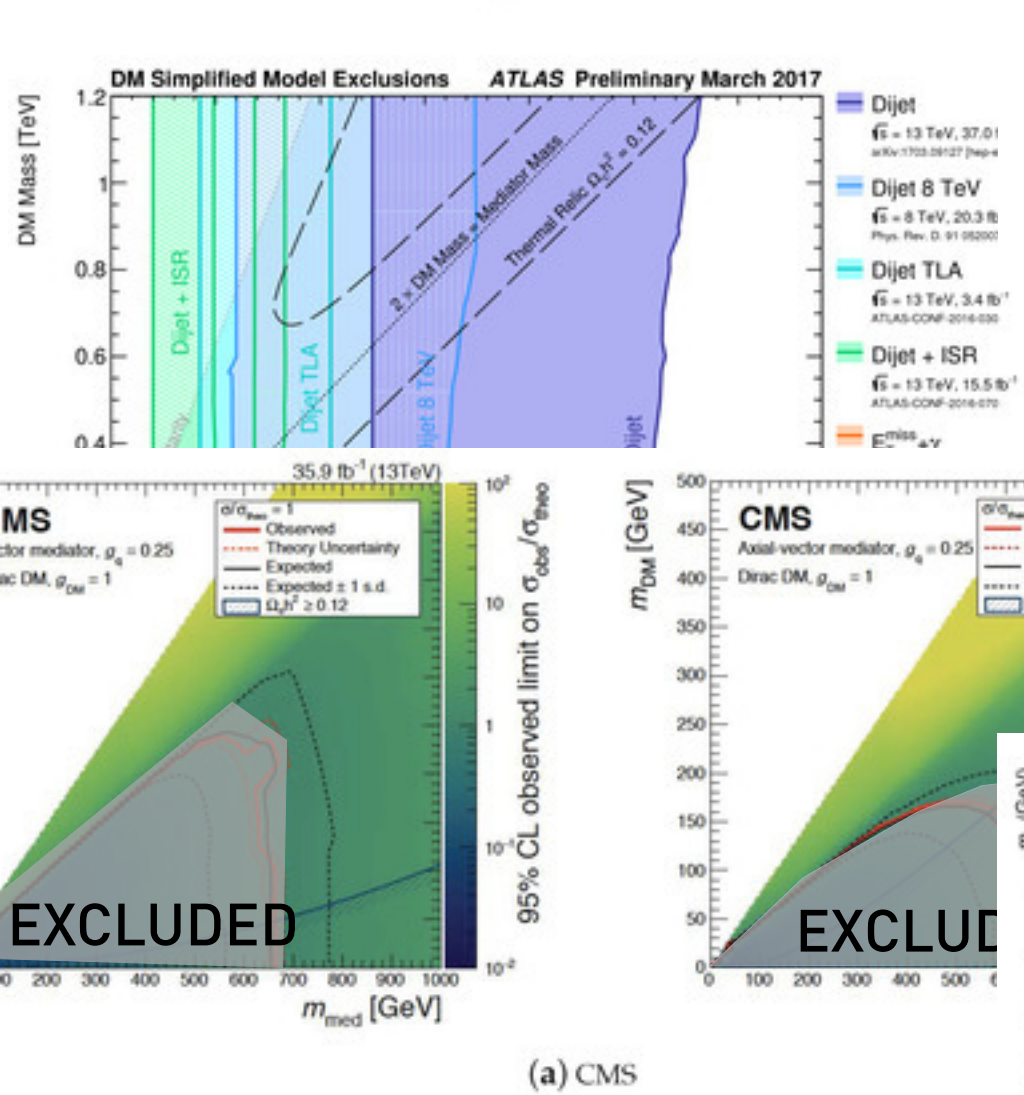
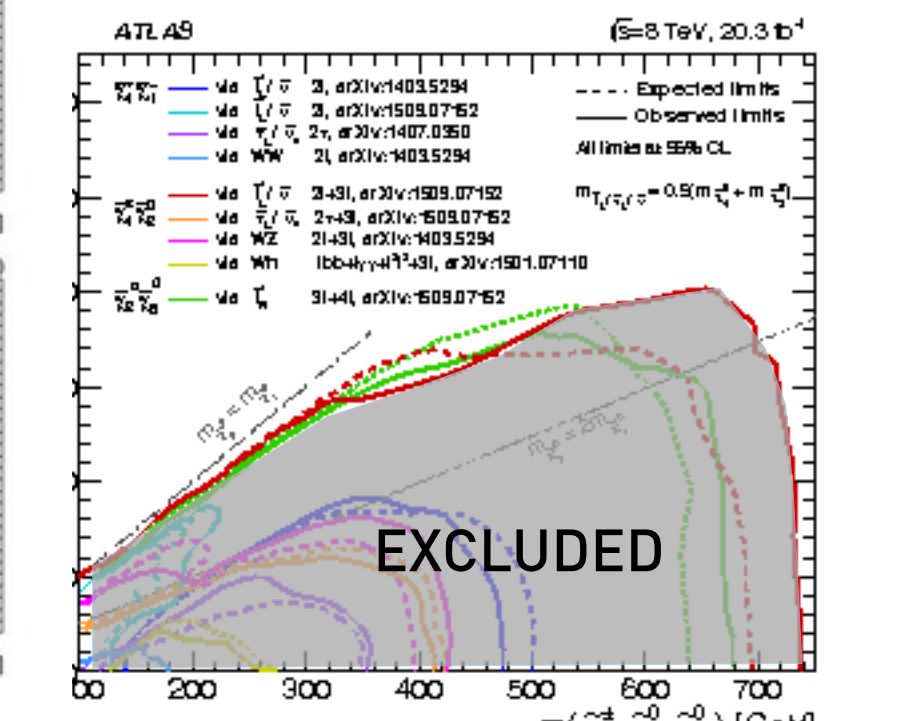
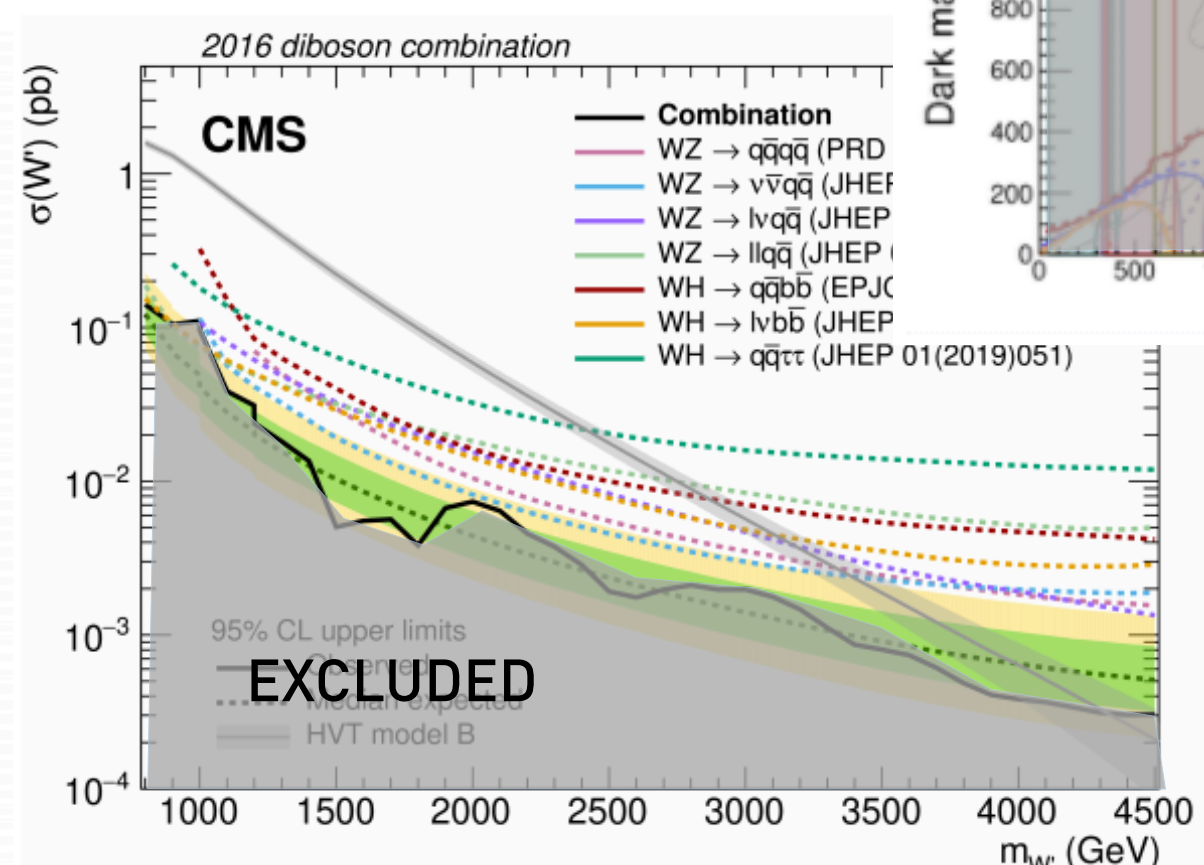
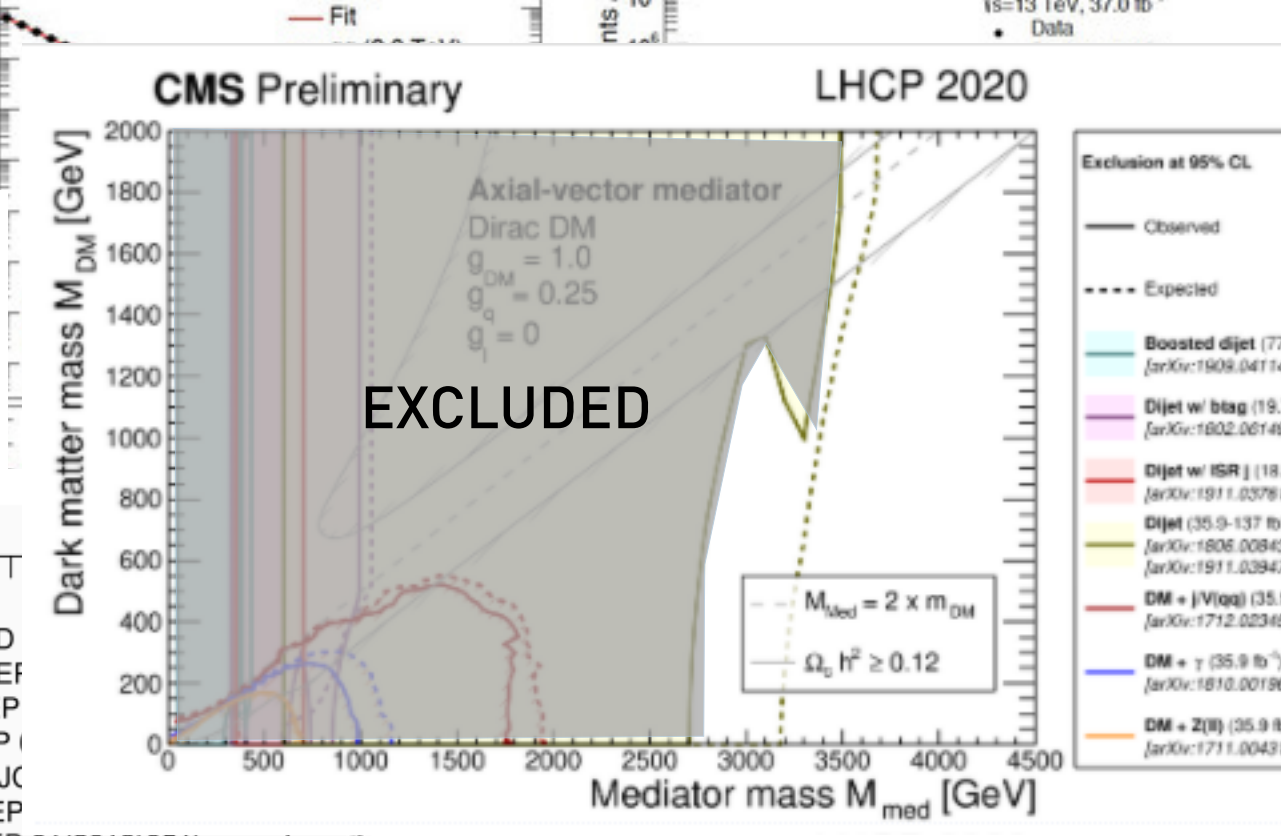
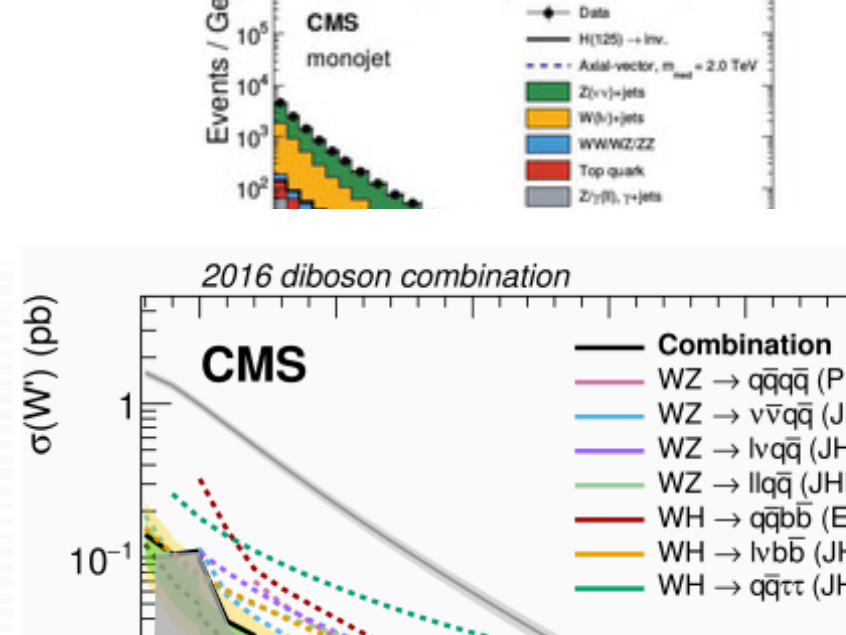
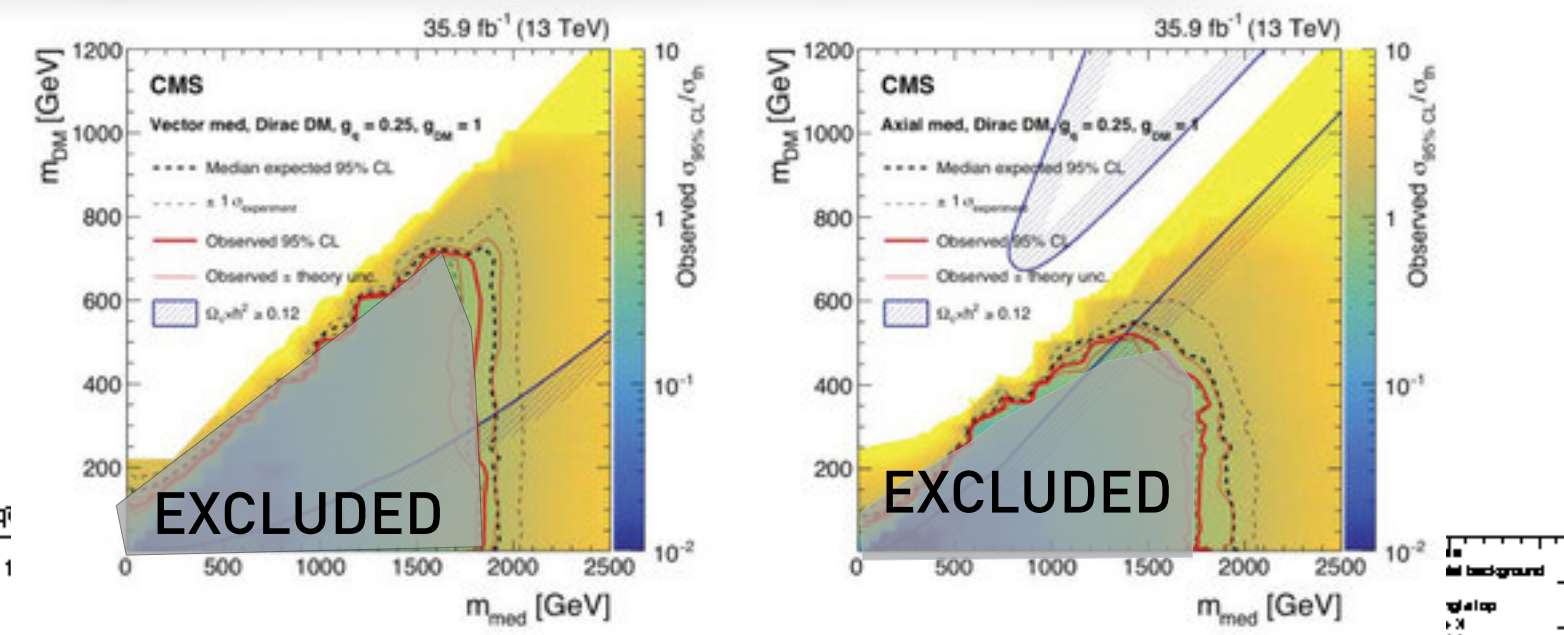
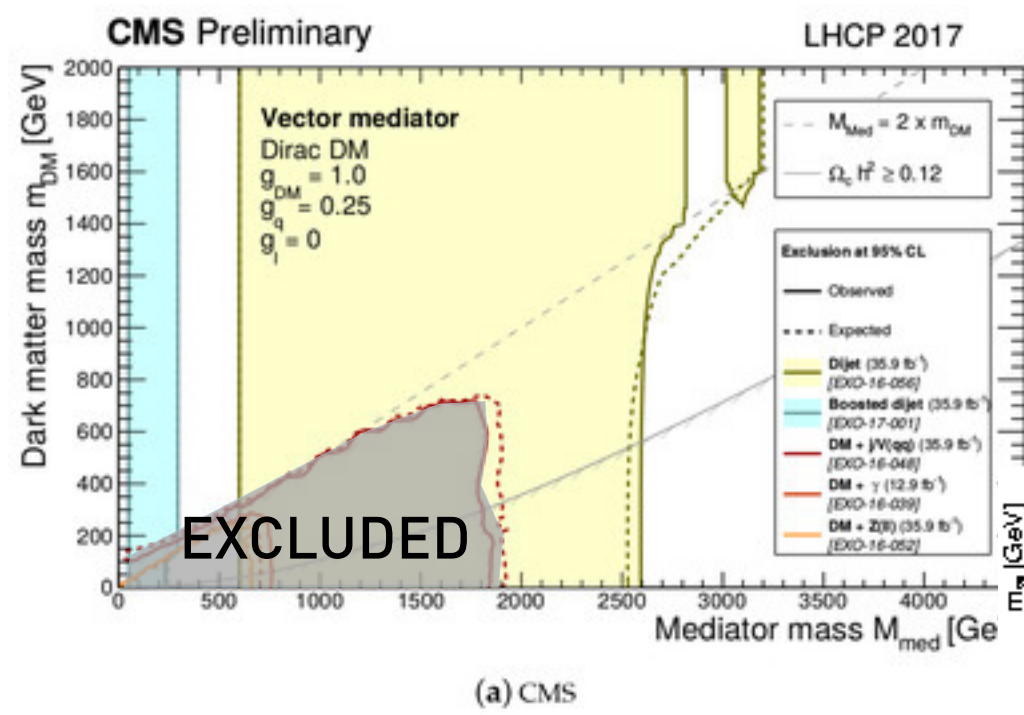
LHC



To make sure we select “the right” 0.0025%, algorithms must be

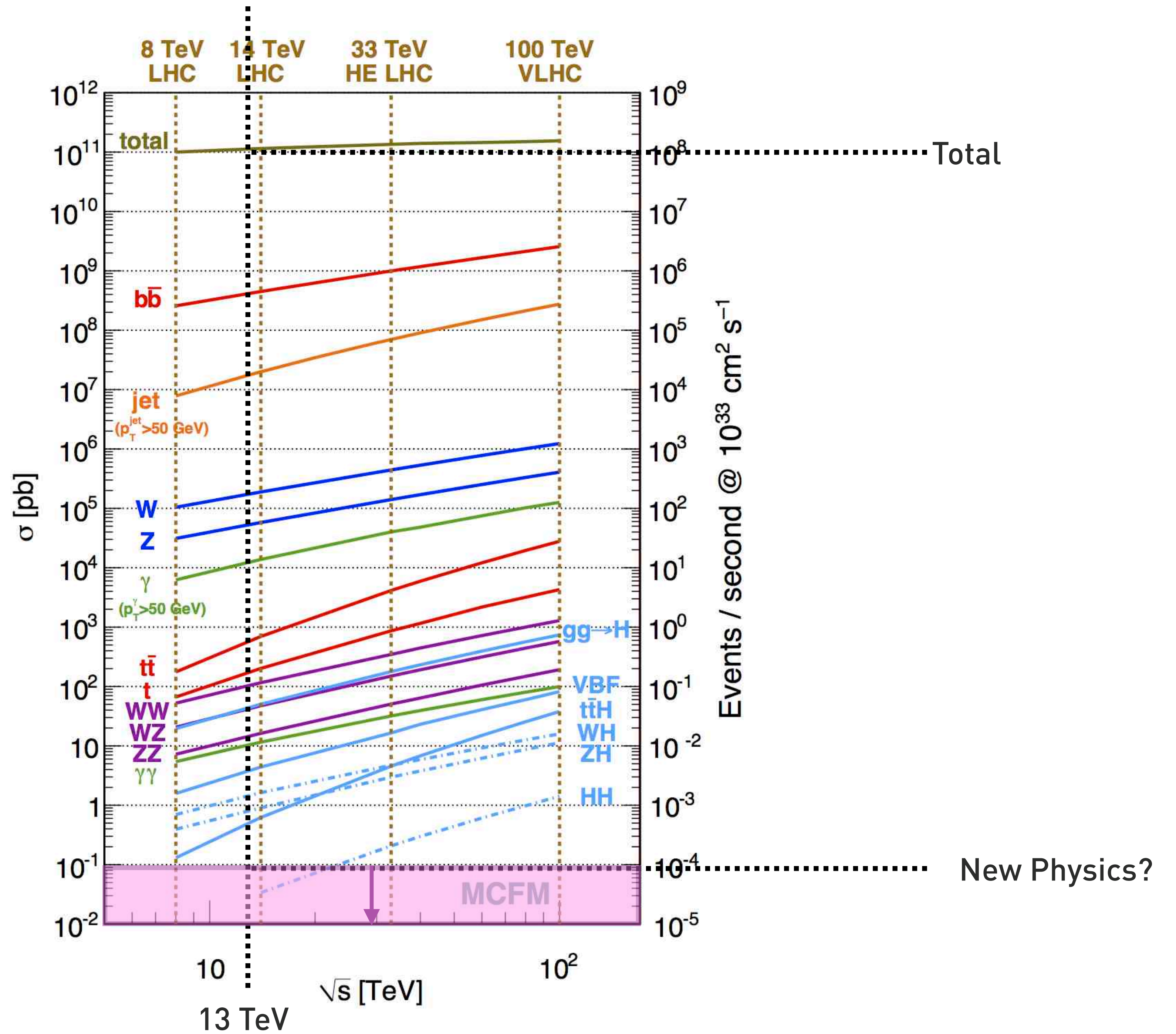
- Fast (get more data through)
- Accurate (select the right data)

Searches for new particles at LHC



New Physics is produced less than 1 in a trillion (if at all)

Need more data!



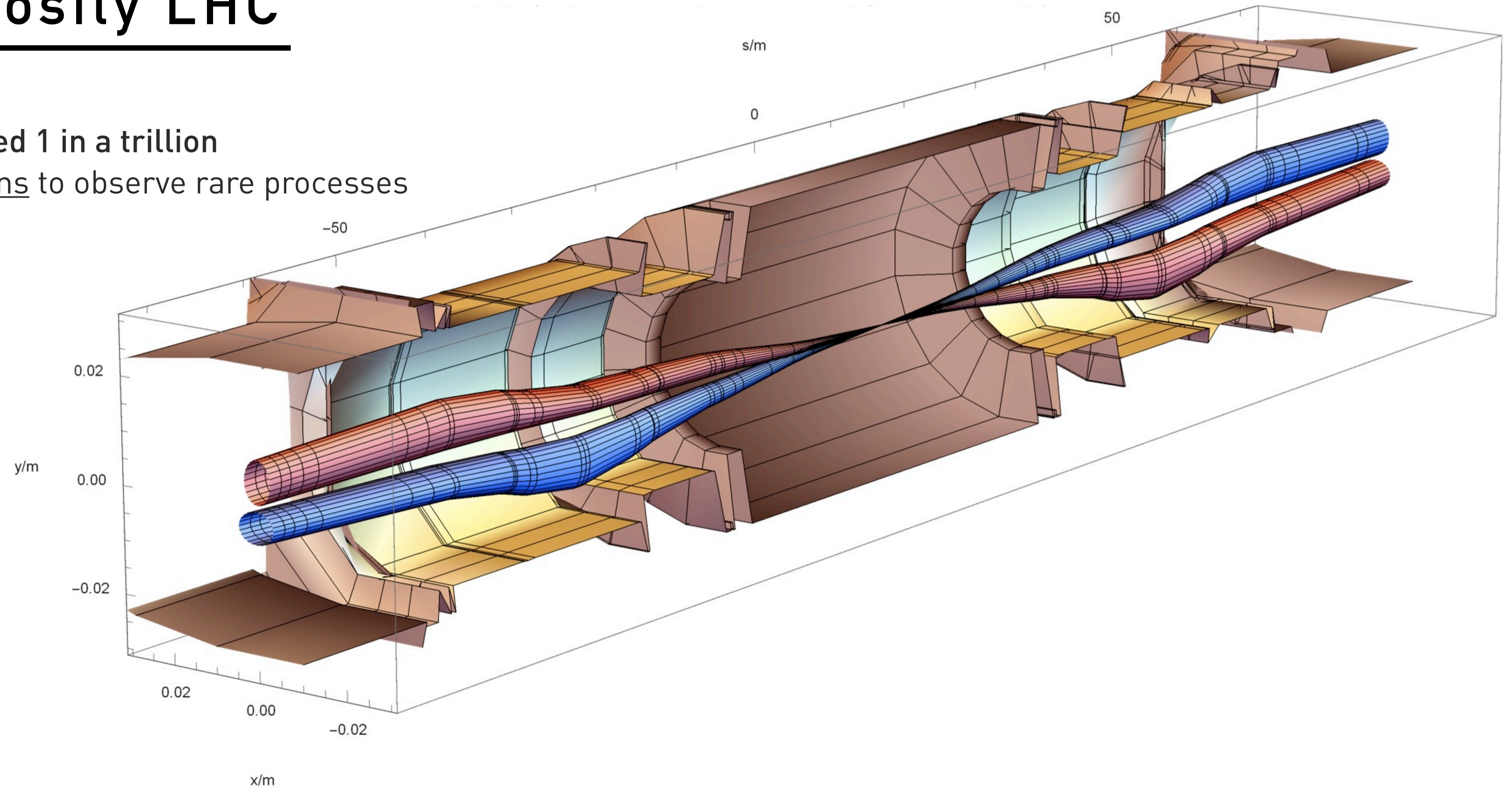
High Luminosity LHC

New Physics is produced 1 in a trillion

- Need more collisions to observe rare processes

High Luminosity LHC

- x10 data size
- x3 collisions/s



2022 - 2025

LHC (TODAY!)

Run 3

2026 - 2028

MAJOR UPGRADE

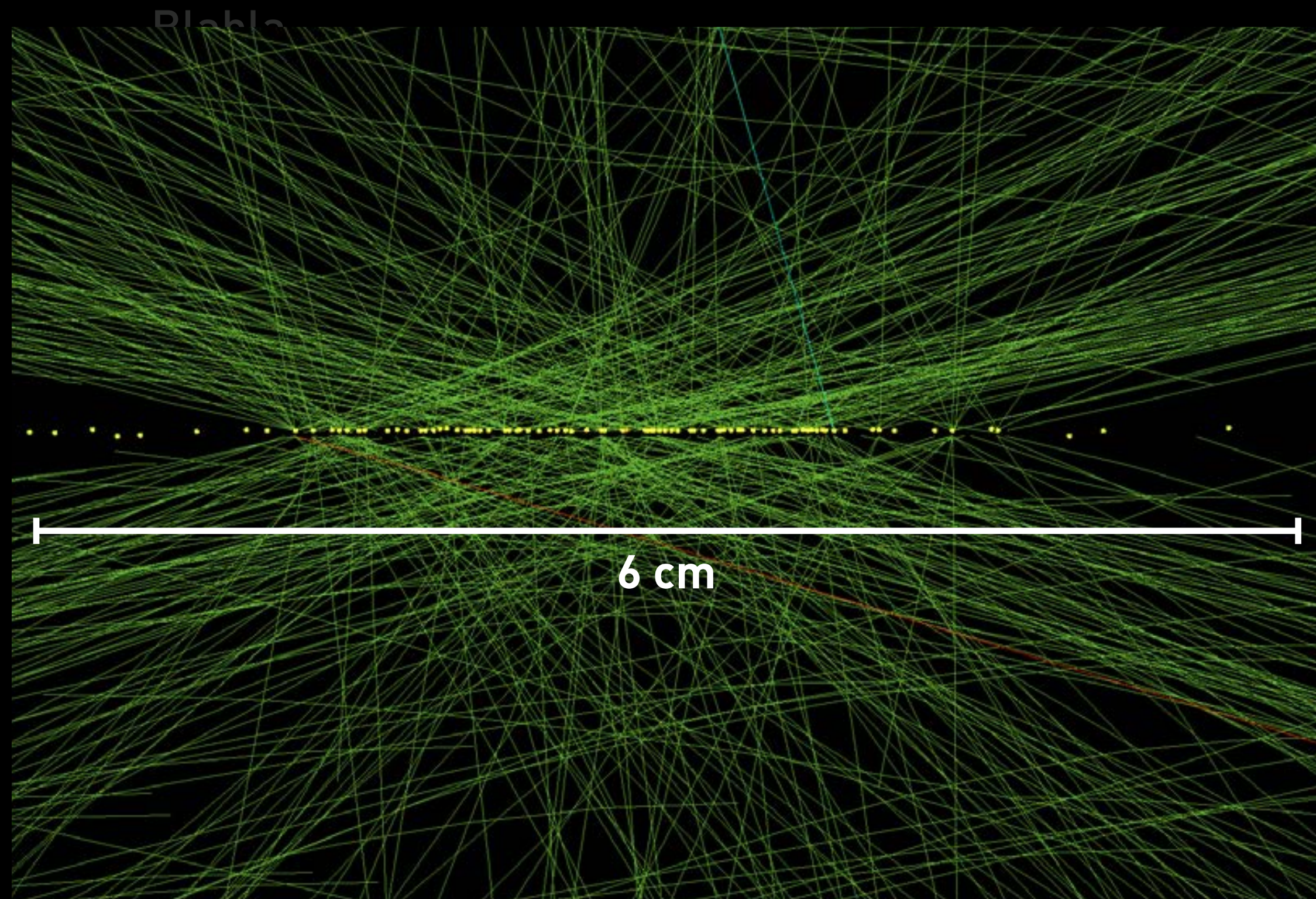
2029 - 2038

HL-LHC

Run 4+5

LHC

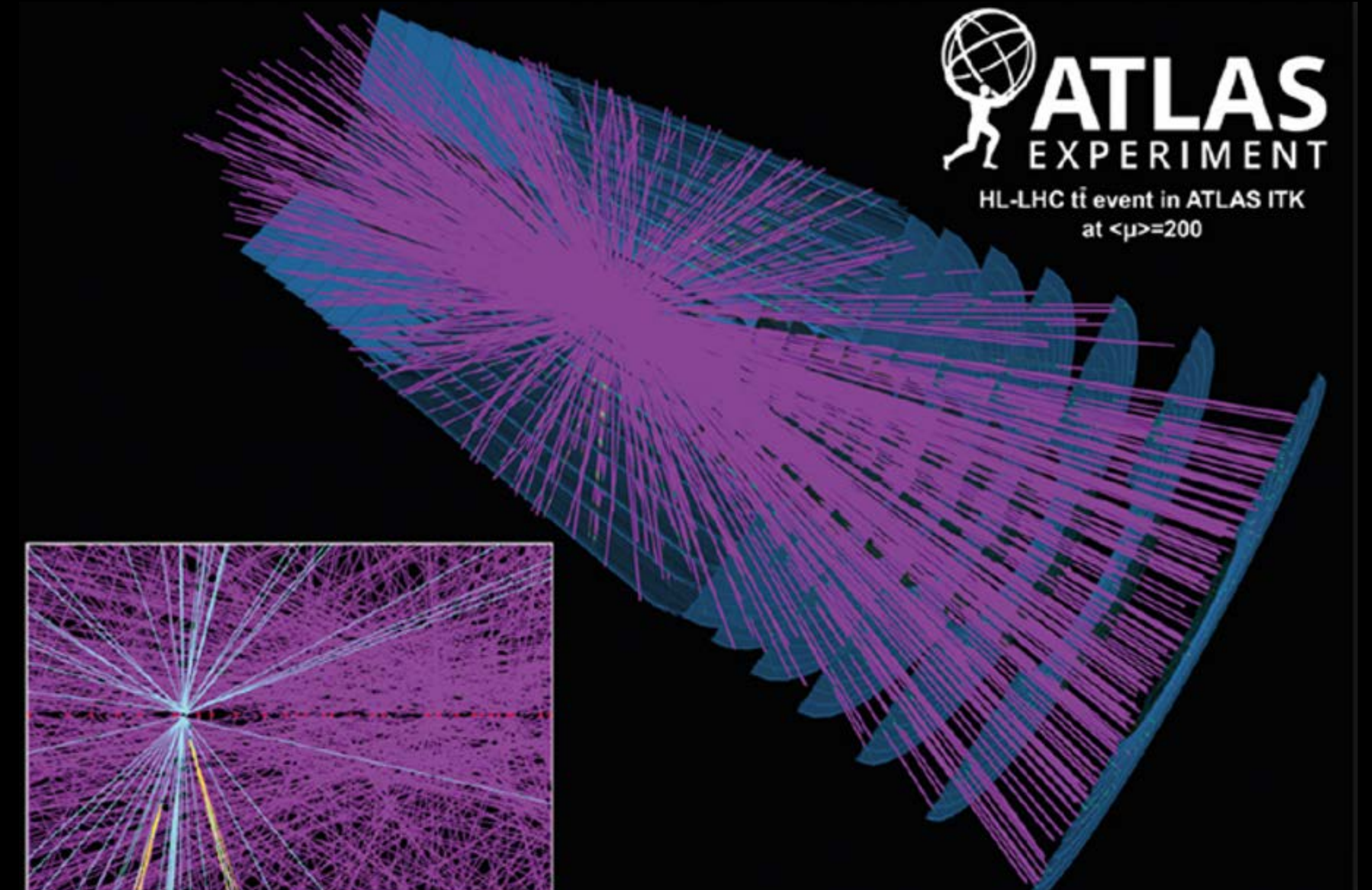
78 vertices
(average 60)



Run 3

High Luminosity LHC

200 vertices
(average 140)



Run 4+5

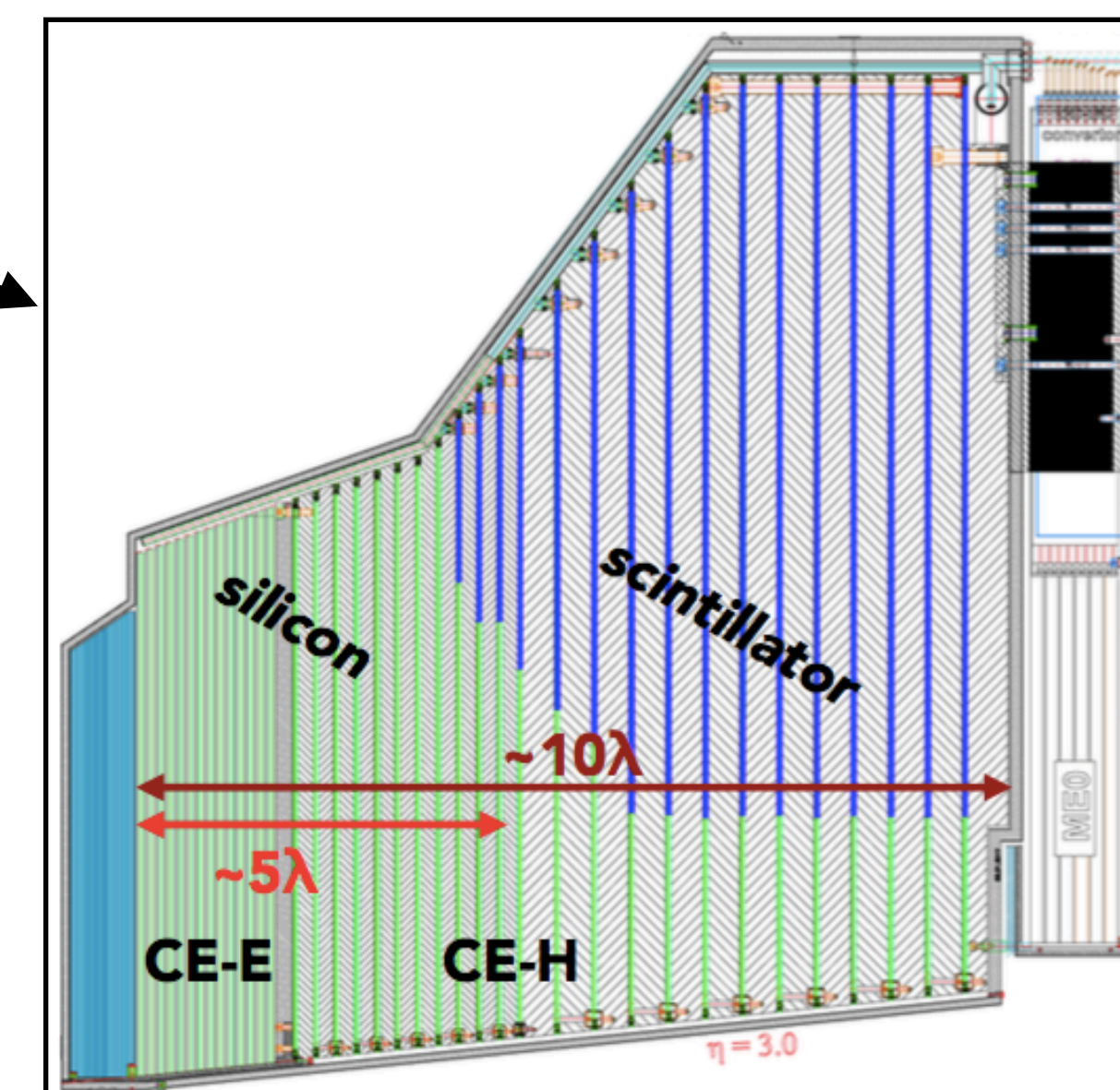
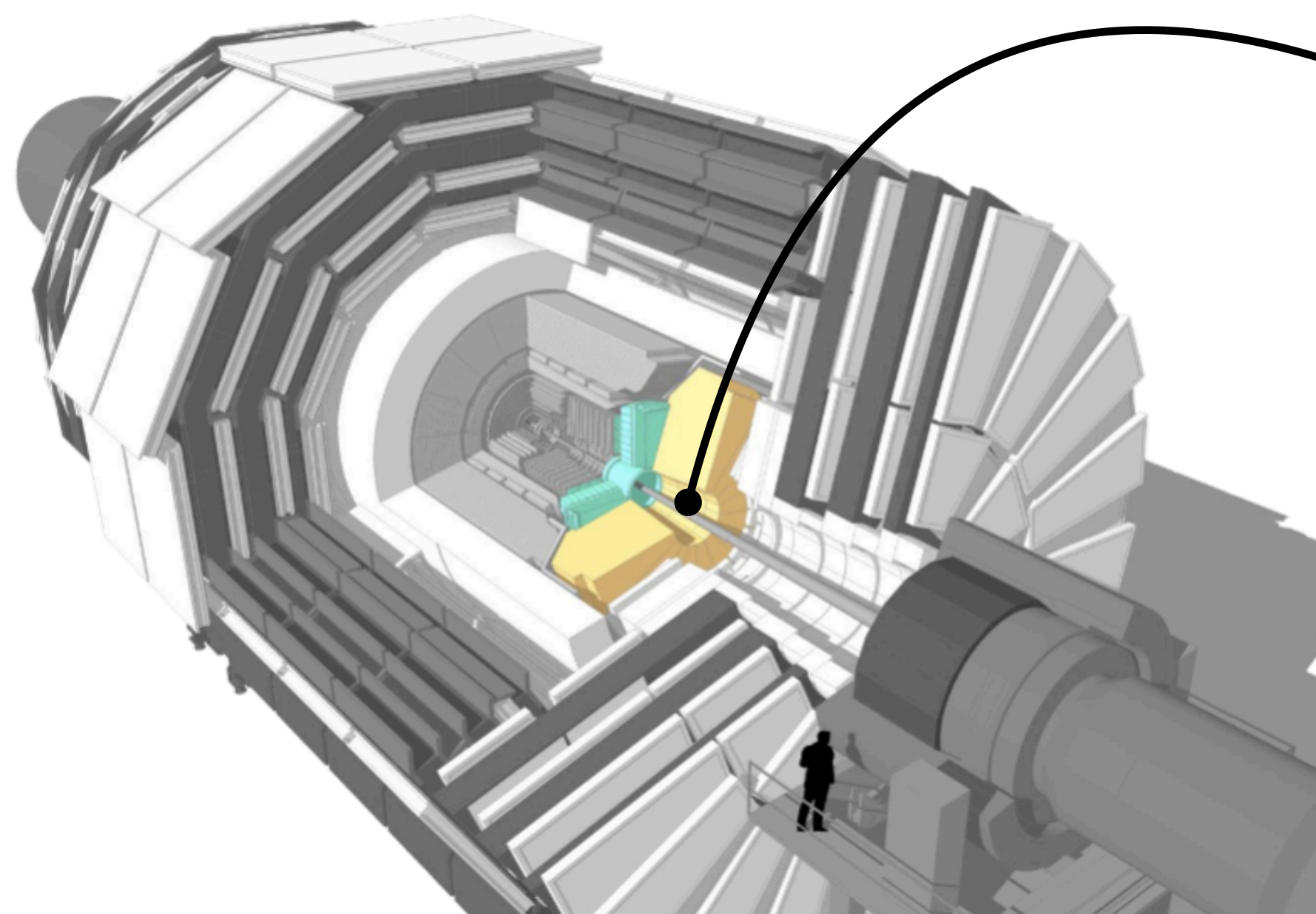
Maintain physics acceptance \rightarrow better detectors

CMS High Granularity (endcap) calorimeter

- 85K (today) \rightarrow 6M (HL-LHC) readout channels

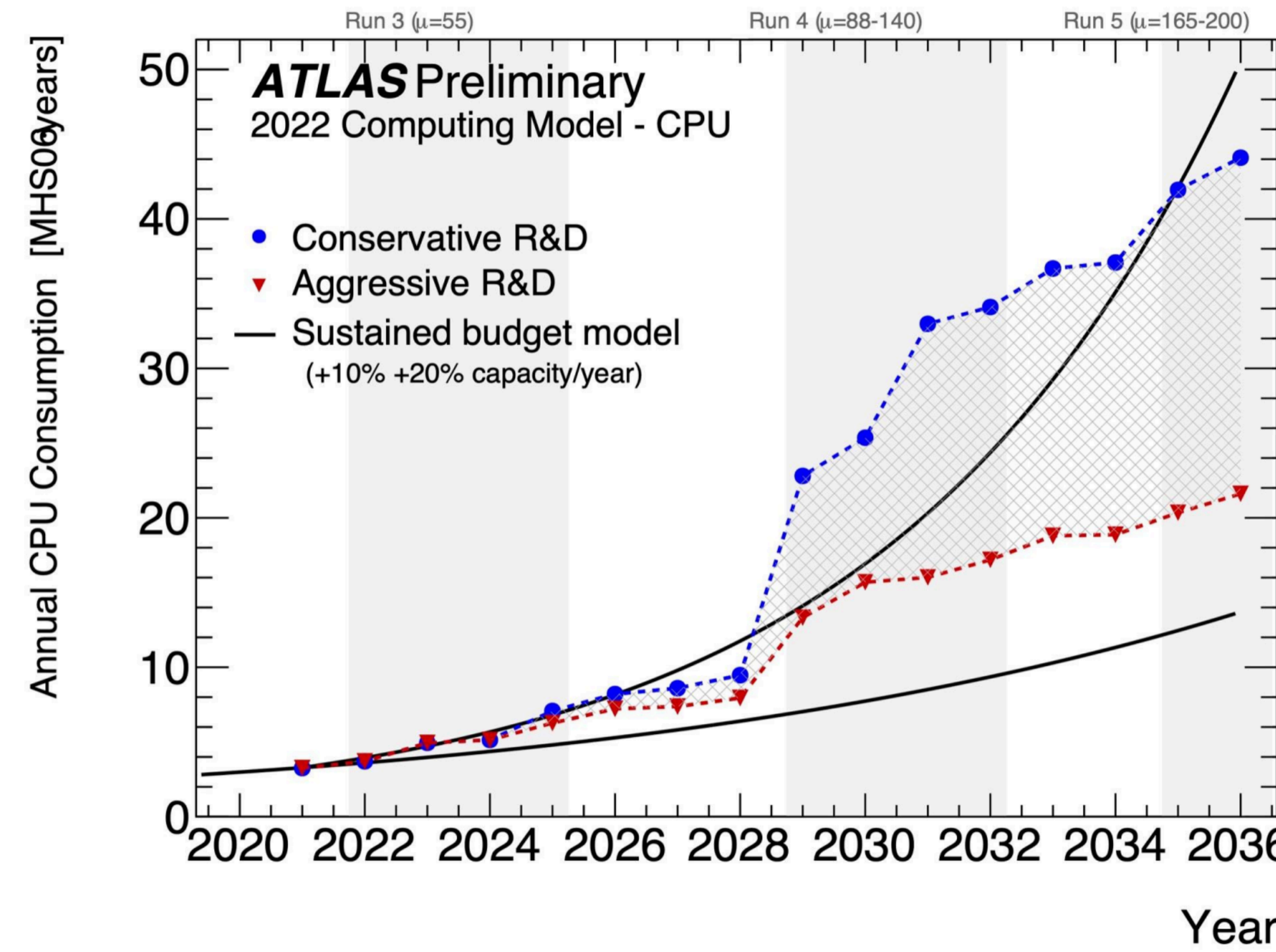
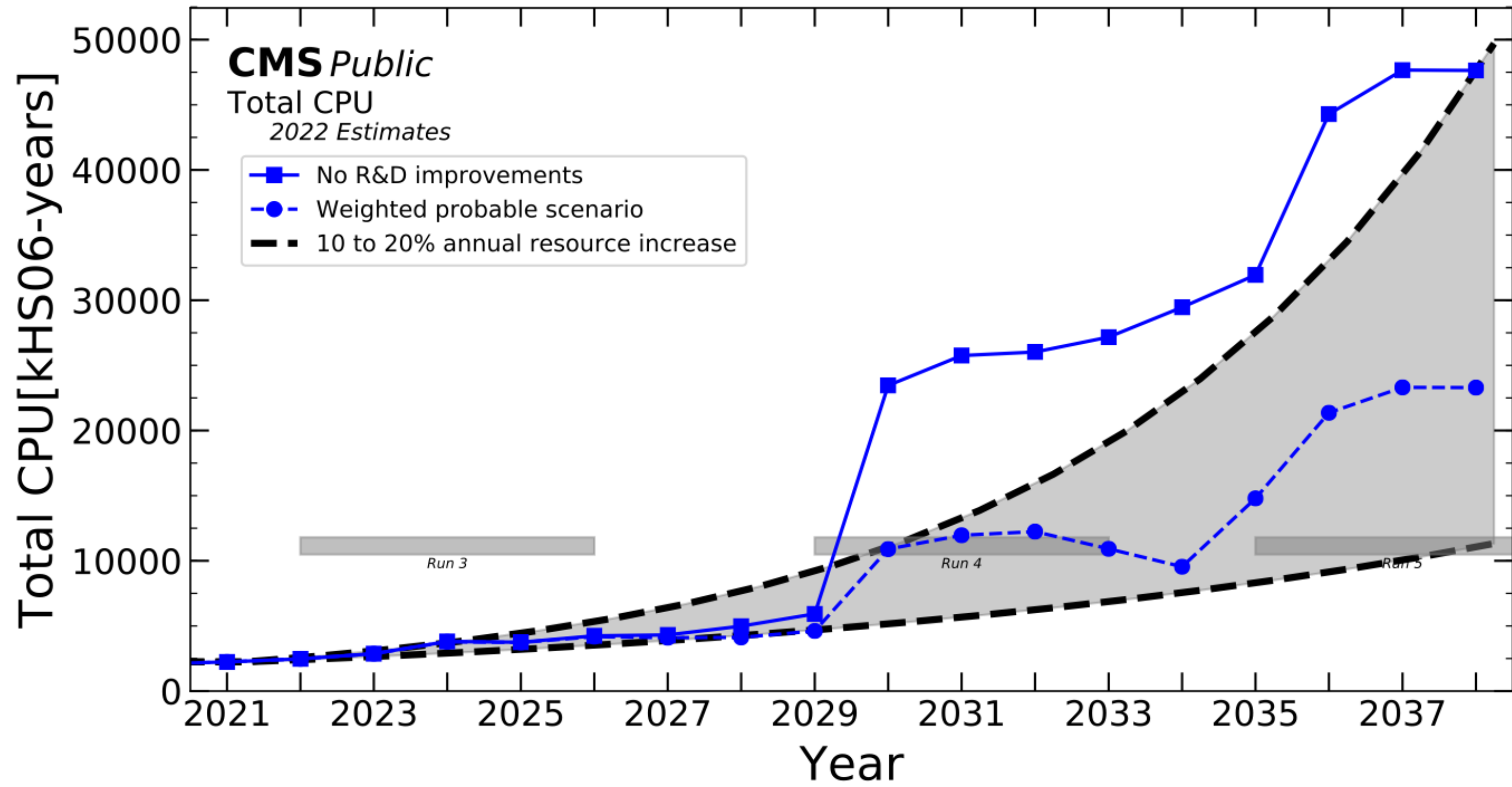
More collisions

More readout channels

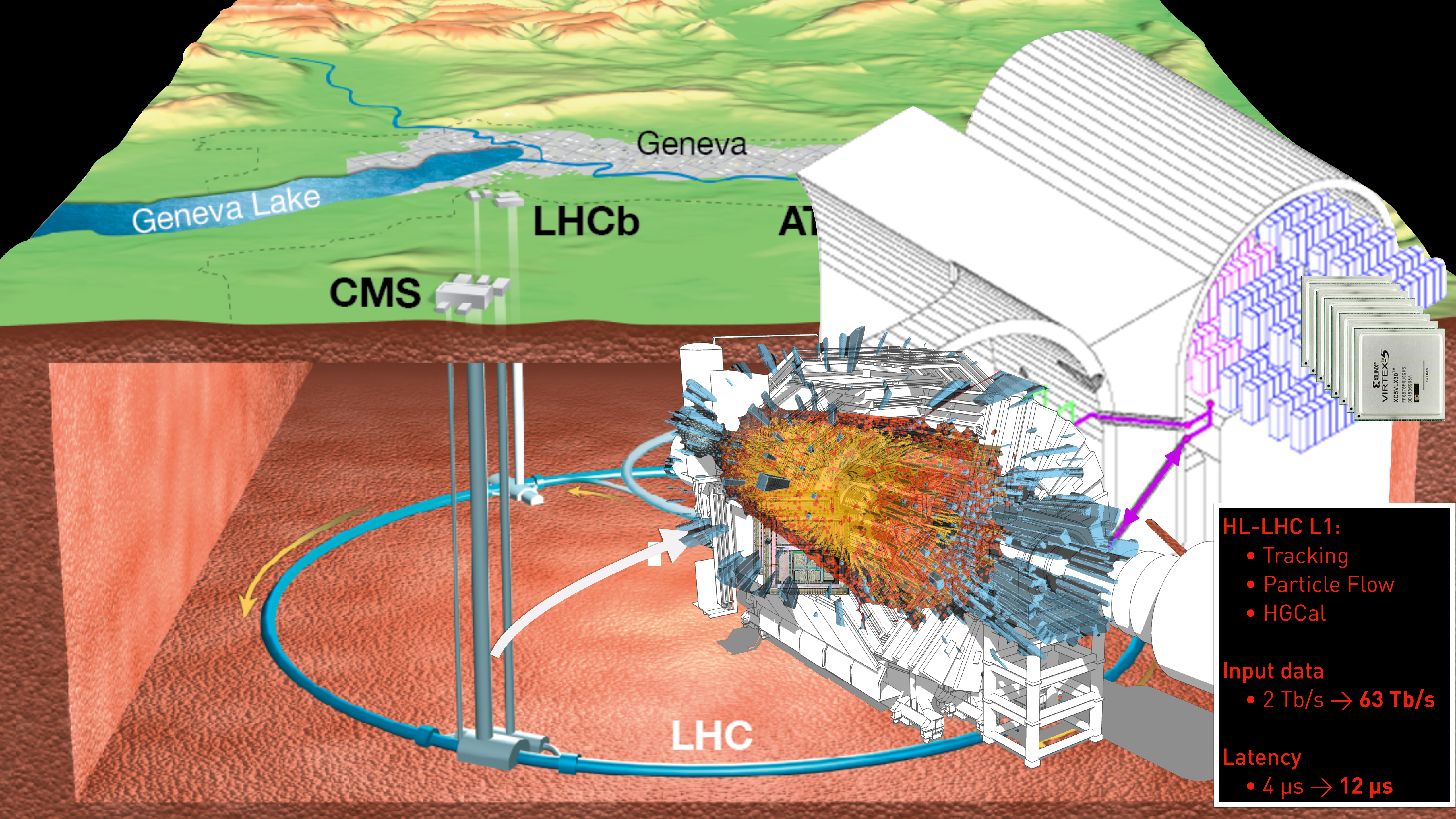


CMS HGCal TDR

Computing resources



... flat computing budget



Geneva Lake

Geneva

CMS

LHCb

ATLAS

LHC

HL-LHC L1:

- Tracking
- Particle Flow
- HGCAL

Input data

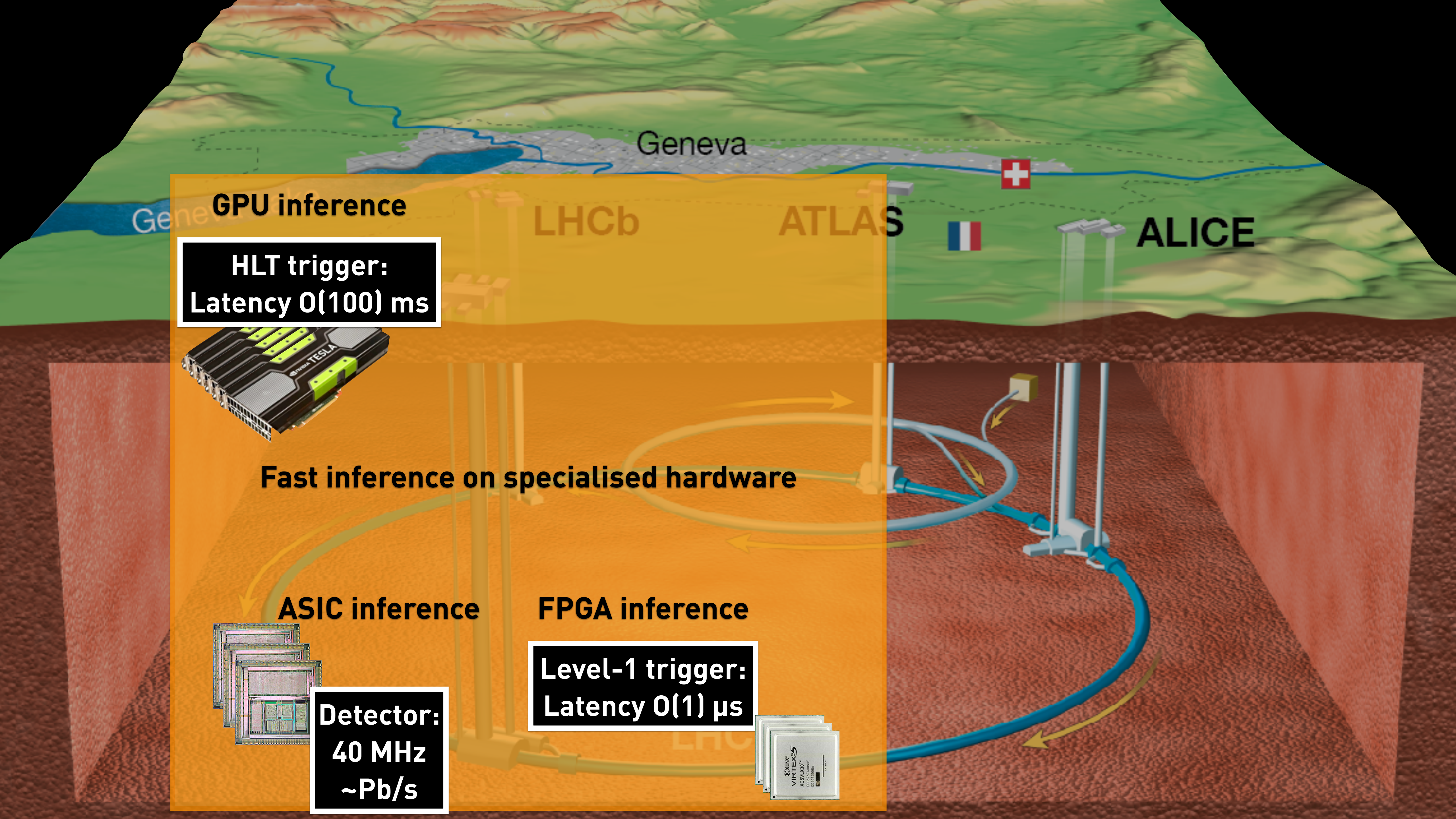
- 2 Tb/s → 63 Tb/s

Latency

- 4 μs → 12 μs



How can we use fast ML to cope with increased data complexity?



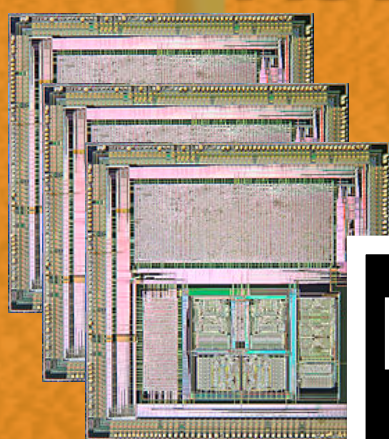
GPU inference

**HLT trigger:
Latency $O(100)$ ms**



Fast inference on specialised hardware

ASIC inference

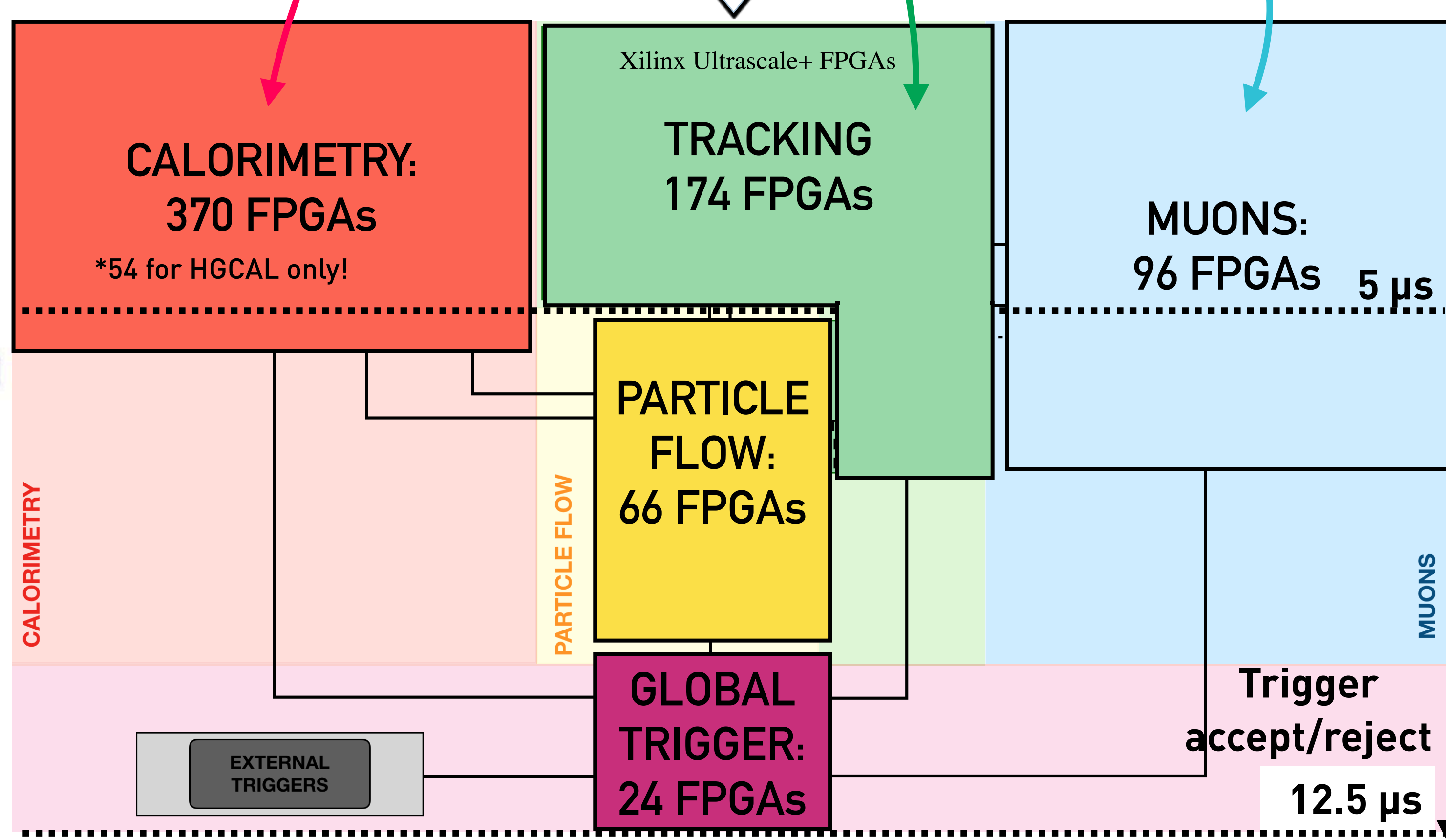
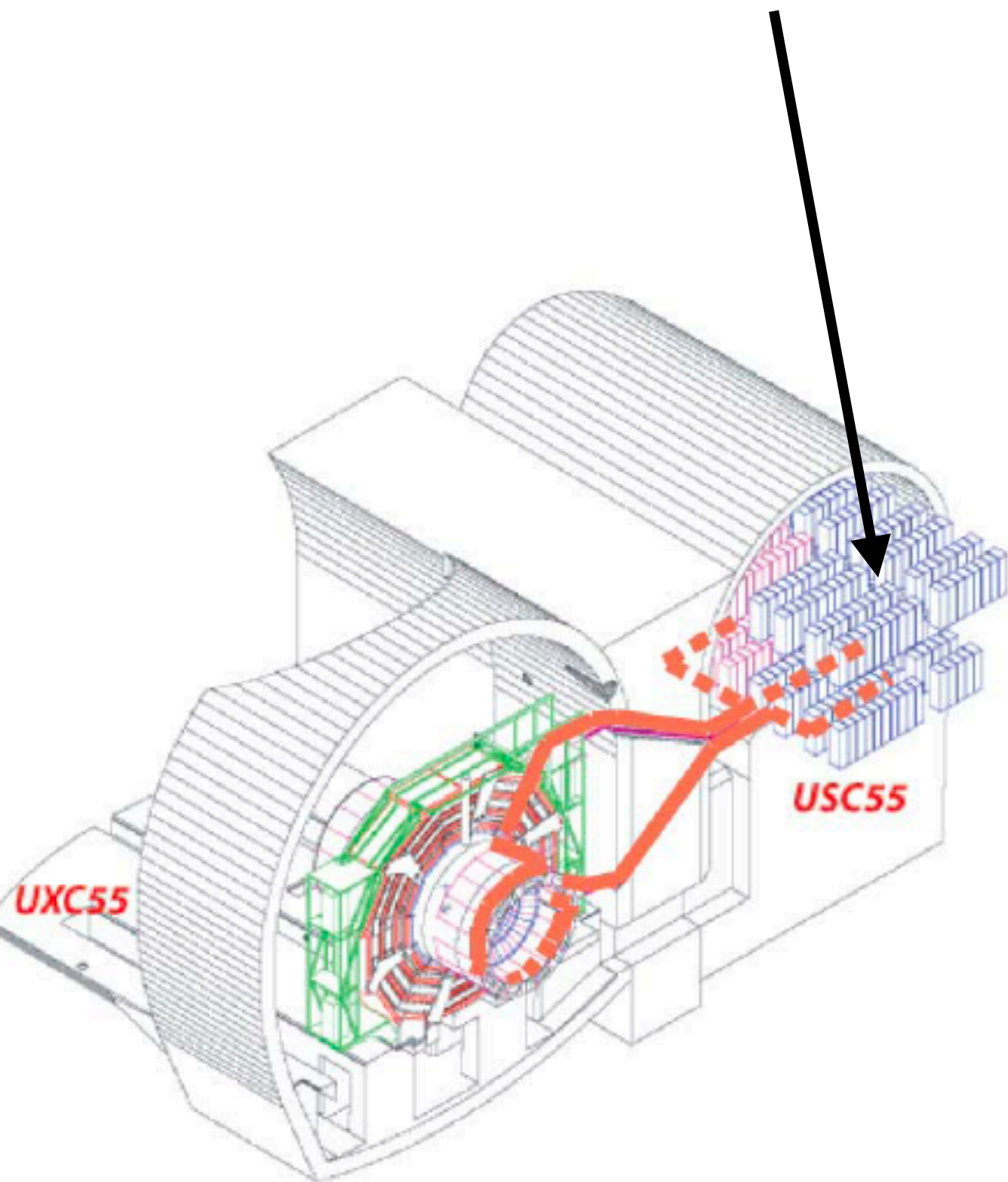
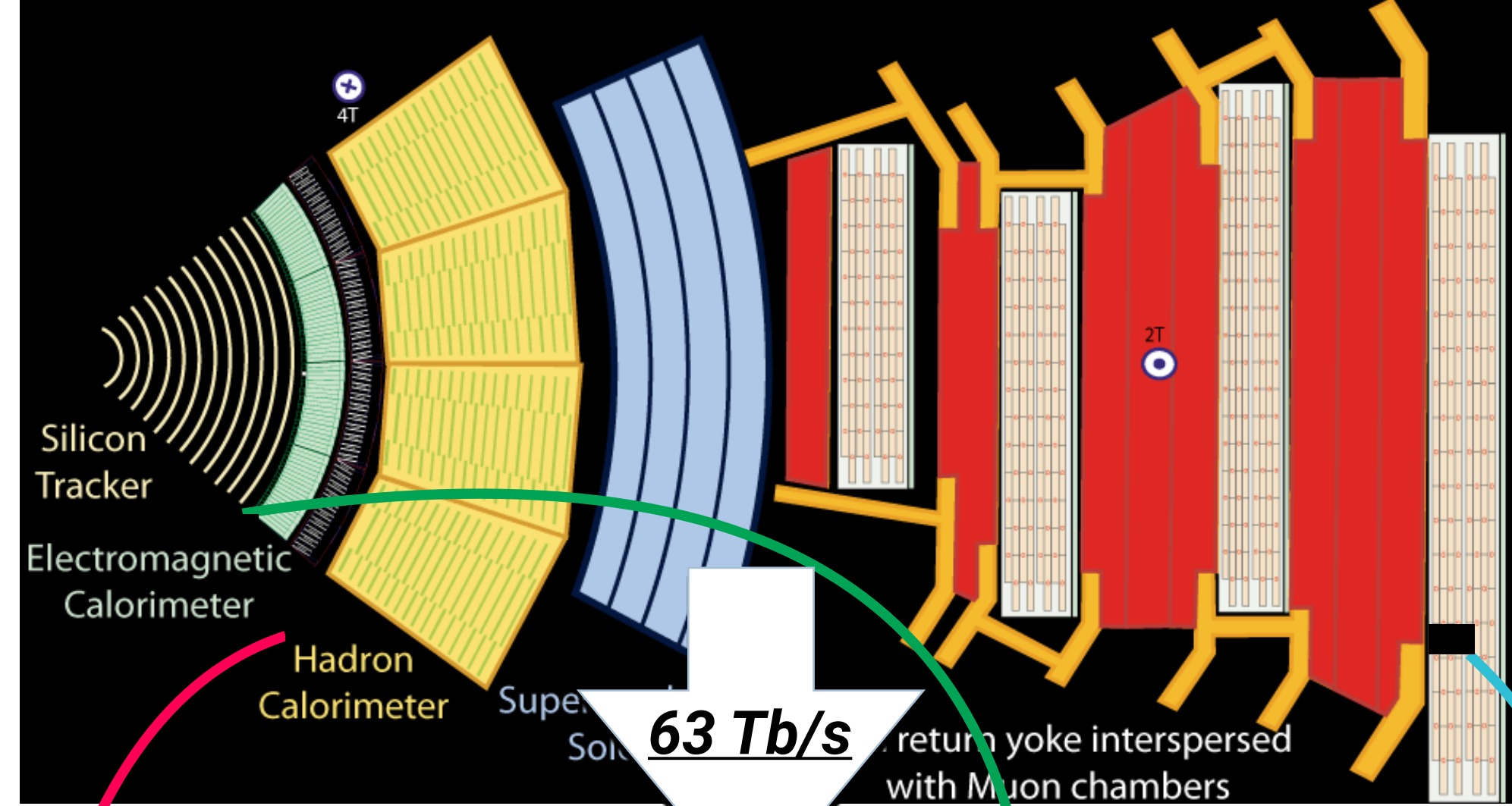
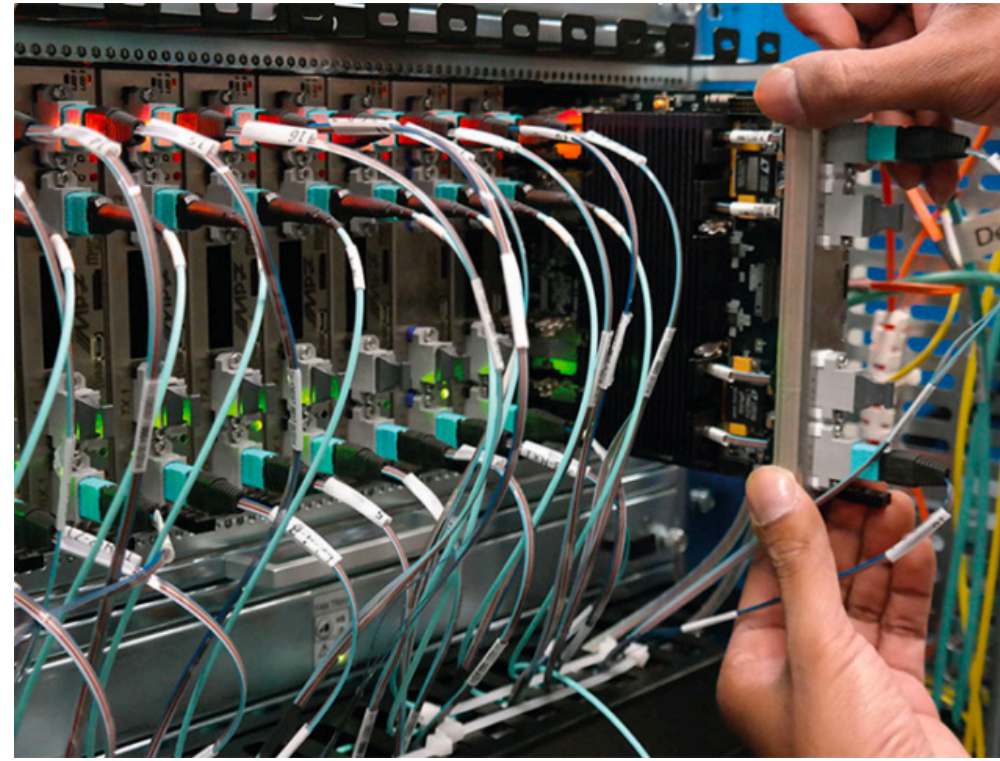


**Detector:
40 MHz
~Pb/s**

FPGA inference

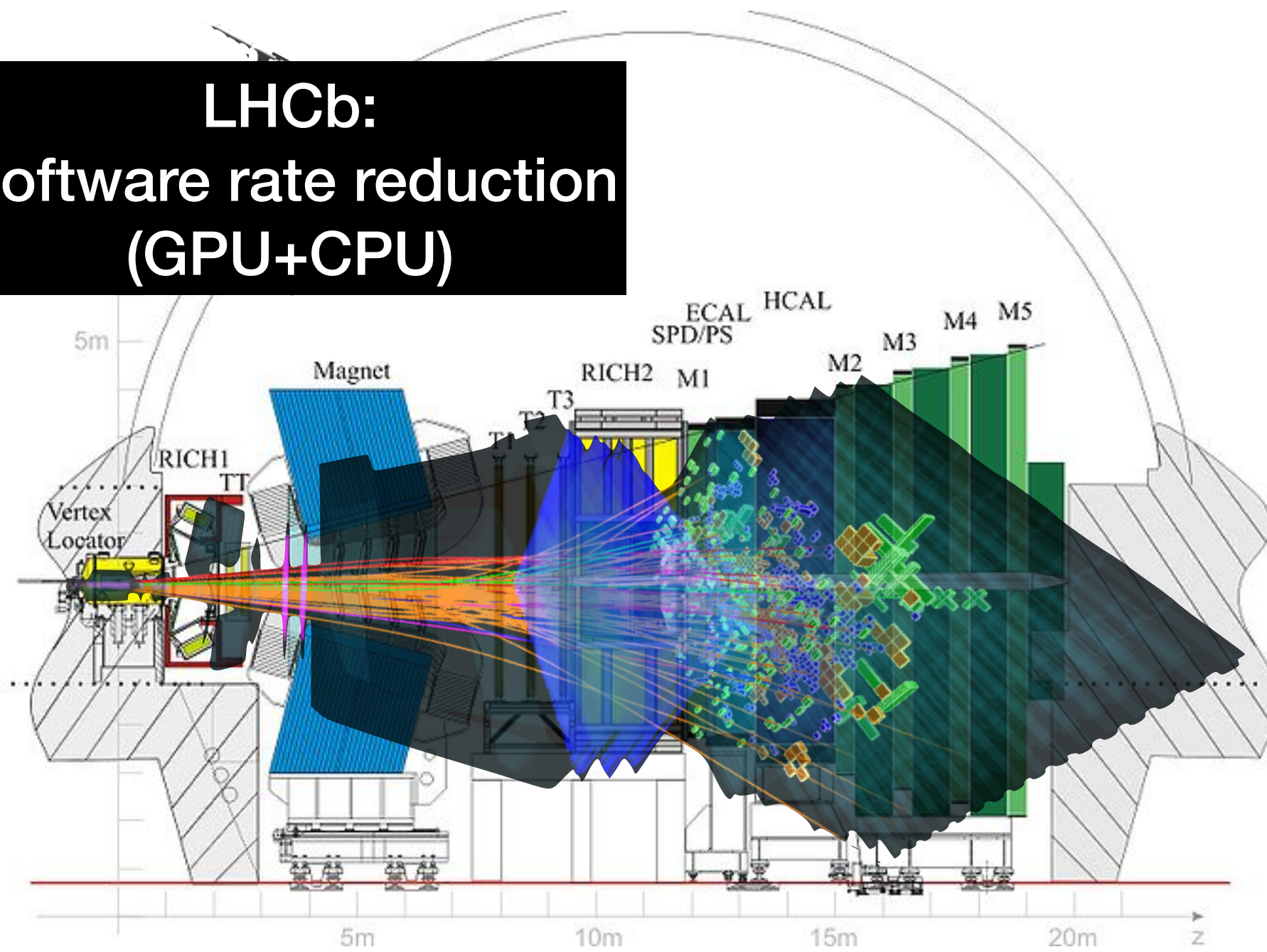
**Level-1 trigger:
Latency $O(1)$ μ s**





Why FPGAs?

**LHCb:
Software rate reduction
(GPU+CPU)**

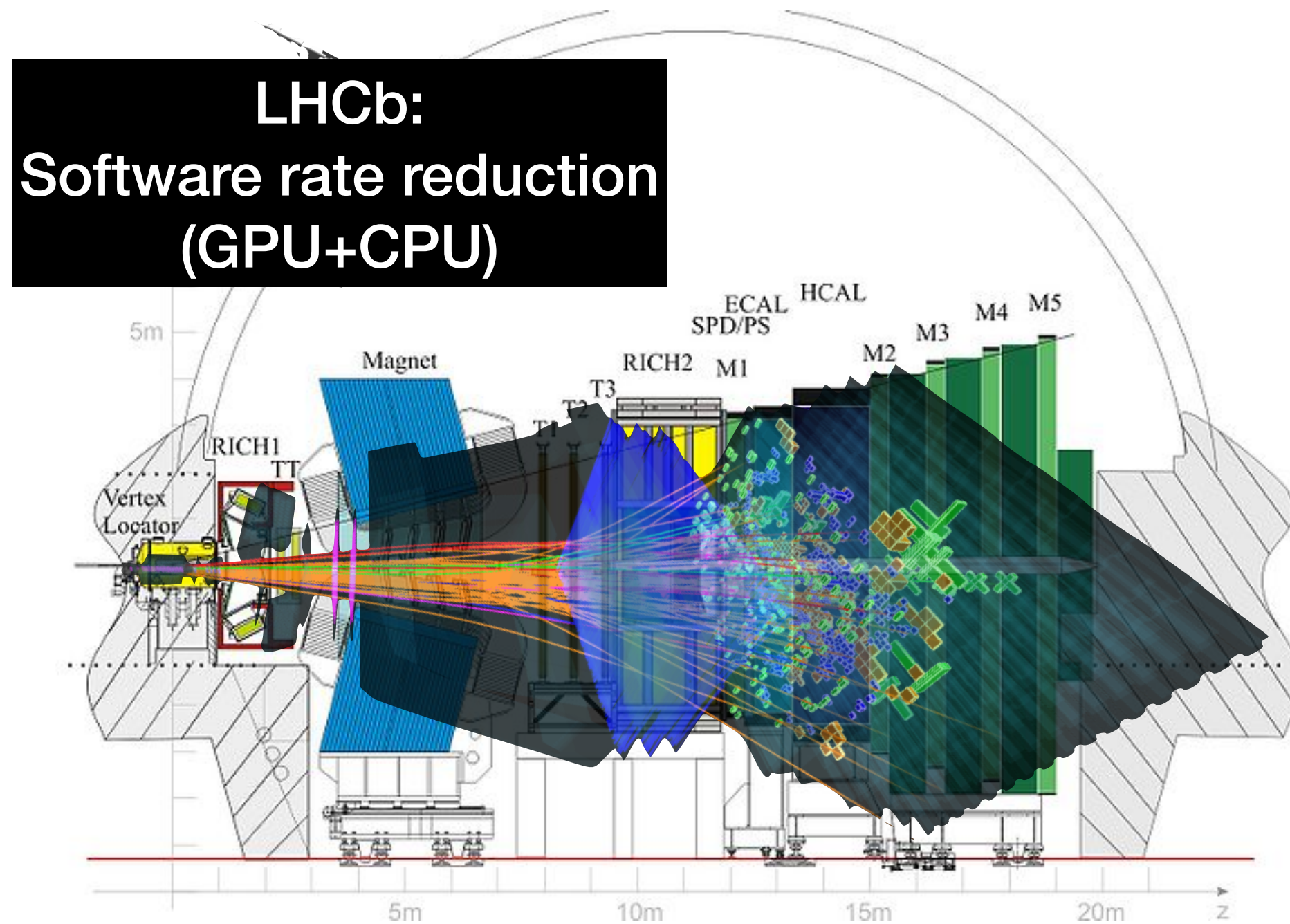


Full GPU reconstruction @ 4 TB/s

- 326 GPUs, 60 kHz per GPU

Why FPGAs?

Depends on your problem...



Characteristics of LHCb HLT1	Characteristics of GPUs <i>C. Fitzpatrick</i>
Intrinsically parallel problem: - Run events in parallel - Reconstruct tracks in parallel	Good for - Data-intensive parallelizable applications - High throughput applications
Huge compute load	Many TFLOPS
Full data stream from all detectors is read out → no stringent latency requirements	Higher latency than CPUs, not as predictable as FPGAs
Small raw event data (~100 kB)	Connection via PCIe → limited I/O bandwidth
Small event raw data (~100 kB)	Thousands of events fit into O(10) GB of memory

Full GPU reconstruction @ 4 TB/s

- 326 GPUs, 60 kHz per GPU

- LHCb has already read out detector
- CMS frontend buffers strictly limited, cannot tolerate latency slack
- CMS raw event data x10 larger, L1 "event" ~ 200 kB (possible with GPU)

Why FPGAs?



Latency, latency, latency (cannot do much on a GPU IN 4 μ s)

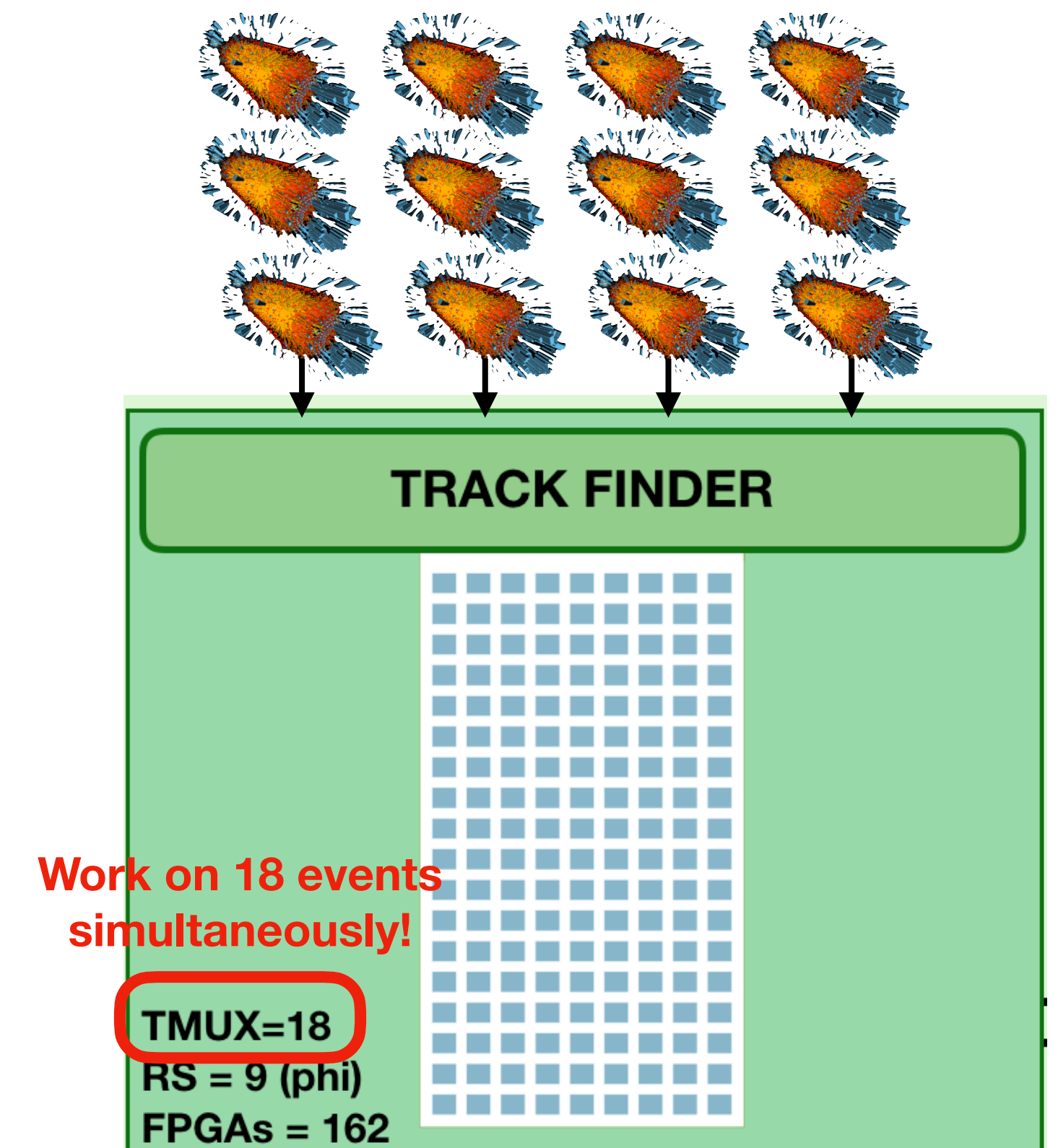
- Can work on different parts of problem, different data simultaneously
- Latency strictly limited by detector frontend buffer

Latency deterministic

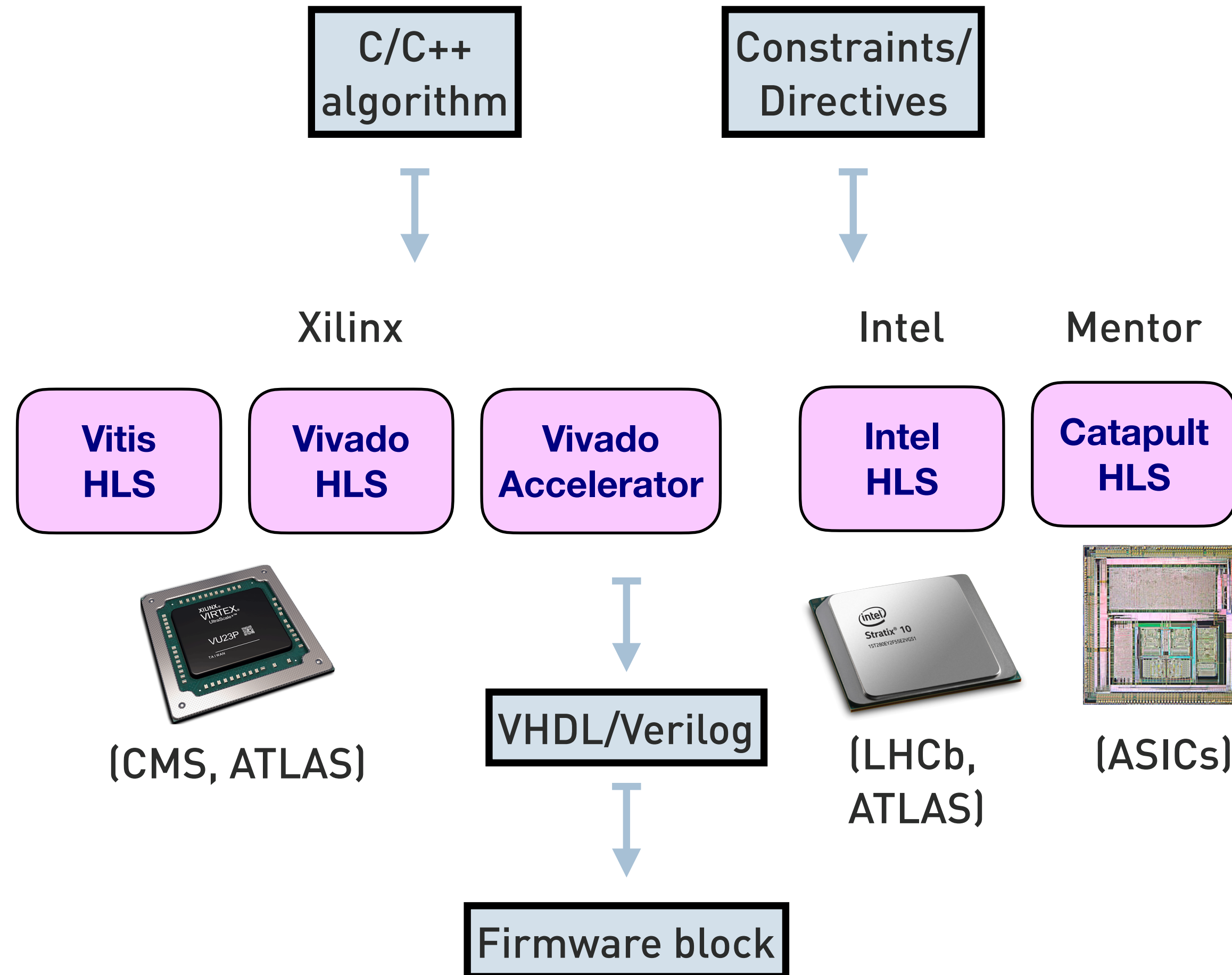
- CPU/GPU processing randomness, FPGAs repeatable and predictable latency

High bandwidth

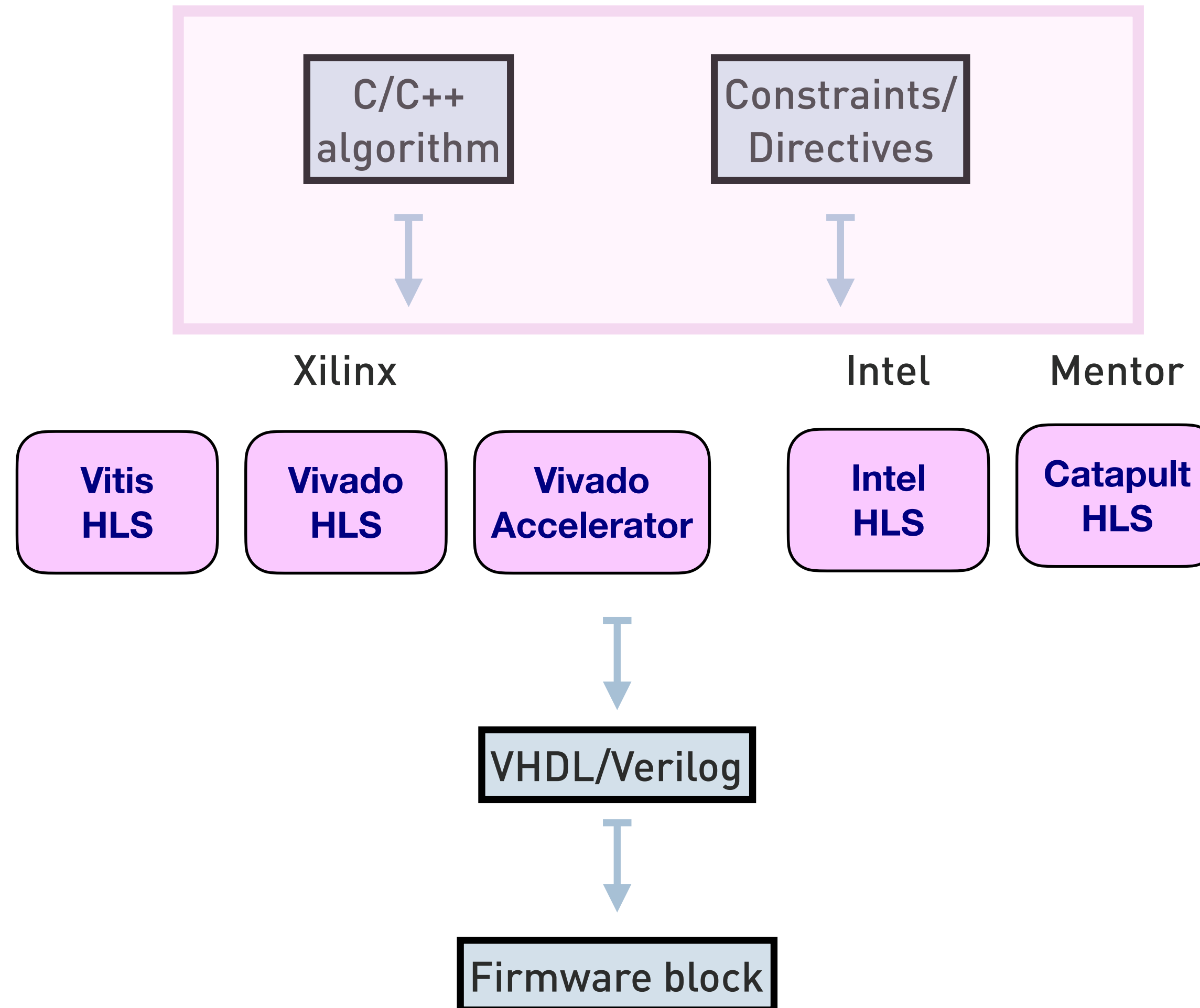
- L1T processes 5% of total internet traffic, dissipate heat of $\sim 7\text{W}/\text{cm}^2$



Programming an FPGA



Programming an FPGA




```

library ieee;
use ieee.std_logic_1164.all;
use ieee.std_logic_unsigned.all;
use ieee.std_logic_arith.all;

use work.gtl_pkg.all;

entity invariant_mass is
  generic (
    > upper_limit: real := 15.0;
    > lower_limit: real := 10.0;
    > pt1_width: positive := 12;
    > pt2_width: positive := 12;
    > cosh_cos_width: positive := 28;
    > INV_MASS_PRECISION : positive := 1;
    > INV_MASS_COSH_COS_PRECISION : positive := 3
  );
  port(
    > pt1 : in std_logic_vector(pt1_width-1 downto 0);
    > pt2 : in std_logic_vector(pt2_width-1 downto 0);
    > cosh_deta : in std_logic_vector(cosh_cos_width-1 downto 0); -- cosh of eta1 - eta2
    > cos_dphi : in std_logic_vector(cosh_cos_width-1 downto 0); -- cos of phi1 - phi2
    inv_mass_comp : out std_logic;
    sim_inv_mass_sq_div2 : out std_logic_vector(pt1_width+pt2_width+cosh_cos_width-1 downto 0)
  );
end invariant_mass;

architecture rtl of invariant_mass is

  constant INV_MASS_VECTOR_WIDTH : positive := pt1_width+pt2_width+cosh_cos_width;
  constant INV_MASS_PRECISION_FACTOR : real := real(10**INV_MASS_PRECISION);
  constant FACTOR_4_VECTOR : std_logic_vector((INV_MASS_COSH_COS_PRECISION+1)*4-1 downto 0) := conv_std_logic_vector(10**(INV_MASS_COSH_COS_PRECISION+1), (INV_MASS_COSH_COS_PRECISION+1)*4);

  signal inv_mass_sq_div2 : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);
  signal upper_limit_vector : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);
  signal lower_limit_vector : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);

begin

  -- Converting the boundary value for the comparison
  upper_limit_vector <= conv_std_logic_vector((integer(upper_limit*INV_MASS_PRECISION_FACTOR)), INV_MASS_VECTOR_WIDTH-FACTOR_4_VECTOR'length)*FACTOR_4_VECTOR;
  lower_limit_vector <= conv_std_logic_vector((integer(lower_limit*INV_MASS_PRECISION_FACTOR)), INV_MASS_VECTOR_WIDTH-FACTOR_4_VECTOR'length)*FACTOR_4_VECTOR;

  -- Calculation of invariant mass with the formula: M**2/2 = pt1*pt2 * (cosh(eta1 - eta2) - cos(phi1 - phi2))
  inv_mass_sq_div2 <= pt1 * pt2 * (cosh_deta - cos_dphi);
  sim_inv_mass_sq_div2 <= inv_mass_sq_div2;

  -- Comparison with boundary values
  inv_mass_comp <= '1' when (inv_mass_sq_div2 >= lower_limit_vector and inv_mass_sq_div2 <= upper_limit_vector) else '0';

end architecture rtl;

```



```

library ieee;
use ieee.std_logic_1164.all;
use ieee.std_logic_unsigned.all;
use ieee.std_logic_arith.all;

use work.gtl_pkg.all;

entity invariant_mass is
  generic (
    >> upper_limit: real := 15.0;
    >> lower_limit: real := 10.0;
    >> pt1_width: positive := 12;
    >> pt2_width: positive := 12;
    >> cosh_cos_width: positive := 28;
    >> INV_MASS_PRECISION : positive := 1;
    >> INV_MASS_COSH_COS_PRECISION : positive := 3
  );
  port(
    >> pt1 : in std_logic_vector(pt1_width-1 downto 0);
    >> pt2 : in std_logic_vector(pt2_width-1 downto 0);
    >> cosh_deta : in std_logic_vector(cosh_cos_width-1 downto 0); -- cosh of eta1 - eta2
    >> cos_dphi : in std_logic_vector(cosh_cos_width-1 downto 0); -- cos of phi1 - phi2
    inv_mass_comp : out std_logic;
    sim_inv_mass_sq_div2 : out std_logic_vector(pt1_width+pt2_width+cosh_cos_width-1 downto 0)
  );
end invariant_mass;

architecture rtl of invariant_mass is

  constant INV_MASS_VECTOR_WIDTH : positive := pt1_width+pt2_width+cosh_cos_width;
  constant INV_MASS_PRECISION_FACTOR : real := real(10**INV_MASS_PRECISION); .pkg.
  constant FACTOR_4_VECTOR : std_logic_vector((INV_MASS_COSH_COS_PRECISION+1)*4-1 downto 0) := conv_std_logic_vector(10**(INV_MASS_COSH_COS_PRECISION+1), (INV_MASS_COSH_COS_PRECISION+1)*4-1);

  signal inv_mass_sq_div2 : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);
  signal upper_limit_vector : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);
  signal lower_limit_vector : std_logic_vector(INV_MASS_VECTOR_WIDTH-1 downto 0);

begin

  -- Converting the boundary value for the comparison
  upper_limit_vector <= conv_std_logic_vector((integer(upper_limit*INV_MASS_PRECISION_FACTOR)), INV_MASS_VECTOR_WIDTH-FACTOR_4_VECTOR'length)*FACTOR_4_VECTOR;
  lower_limit_vector <= conv_std_logic_vector((integer(lower_limit*INV_MASS_PRECISION_FACTOR)), INV_MASS_VECTOR_WIDTH-FACTOR_4_VECTOR'length)*FACTOR_4_VECTOR;

  -- Calculation of invariant mass with the formula: M**2/2 = pt1*pt2 * (cosh(eta1 - eta2) - cos(phi1 - phi2))
  inv_mass_sq_div2 <= pt1 * pt2 * (cosh_deta - cos_dphi);
  sim_inv_mass_sq_div2 <= inv_mass_sq_div2;

  -- Comparison with boundary values
  inv_mass_comp <= '1' when (inv_mass_sq_div2 >= lower_limit_vector and inv_mass_sq_div2 <= upper_limit_vector) else '0';

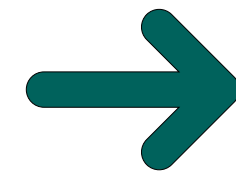
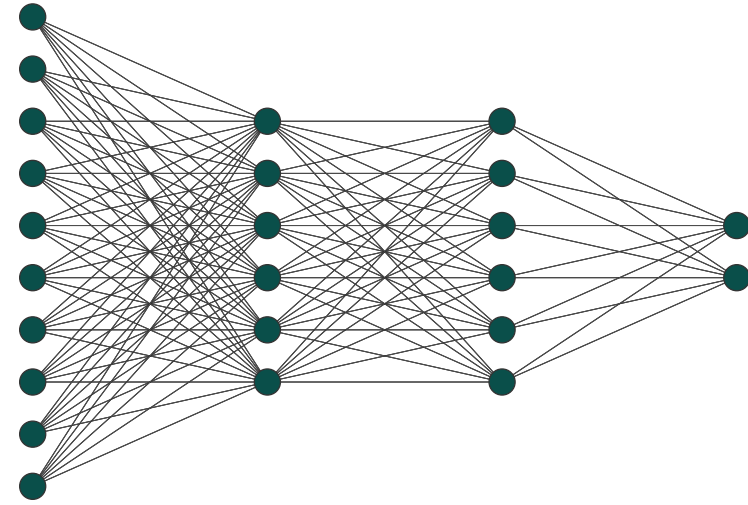
end architecture rtl;

```

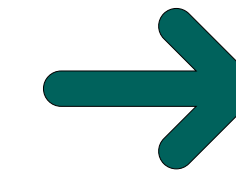
$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

Generic (superfast) HLS implementations for DNN inference?

KERAS / PyTorch / ONNX

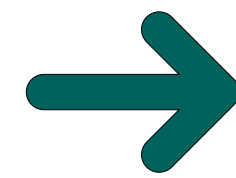
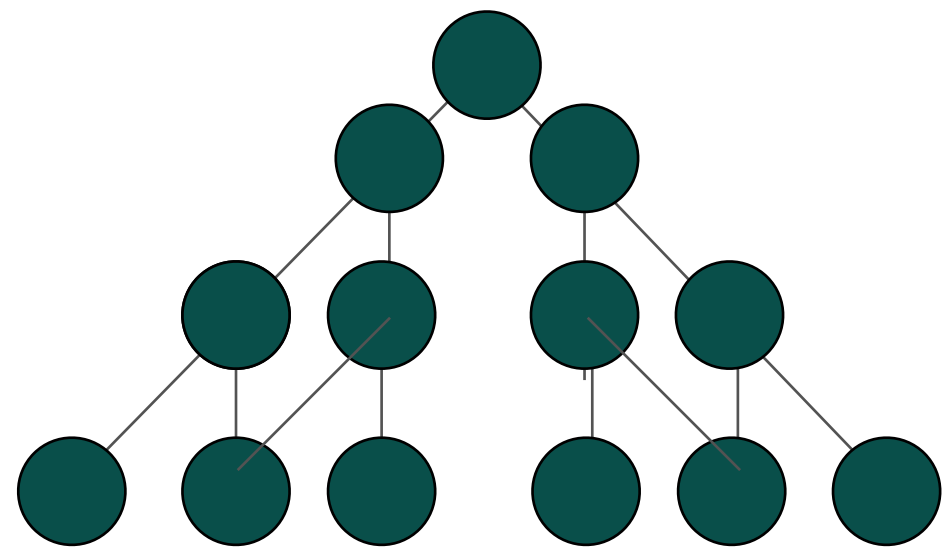


hls4ml

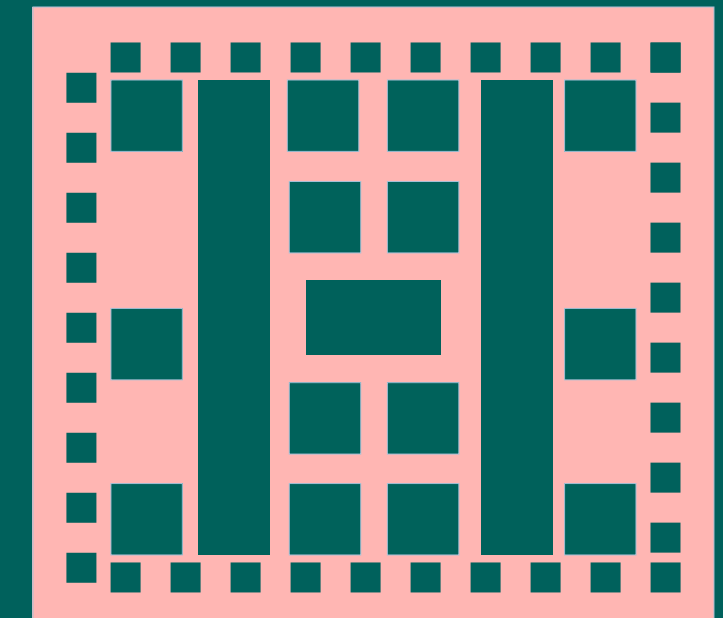
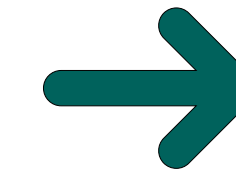


HLS project:
Vivado / Vitis / Intel Quartus /
IntelOne API / Catapult

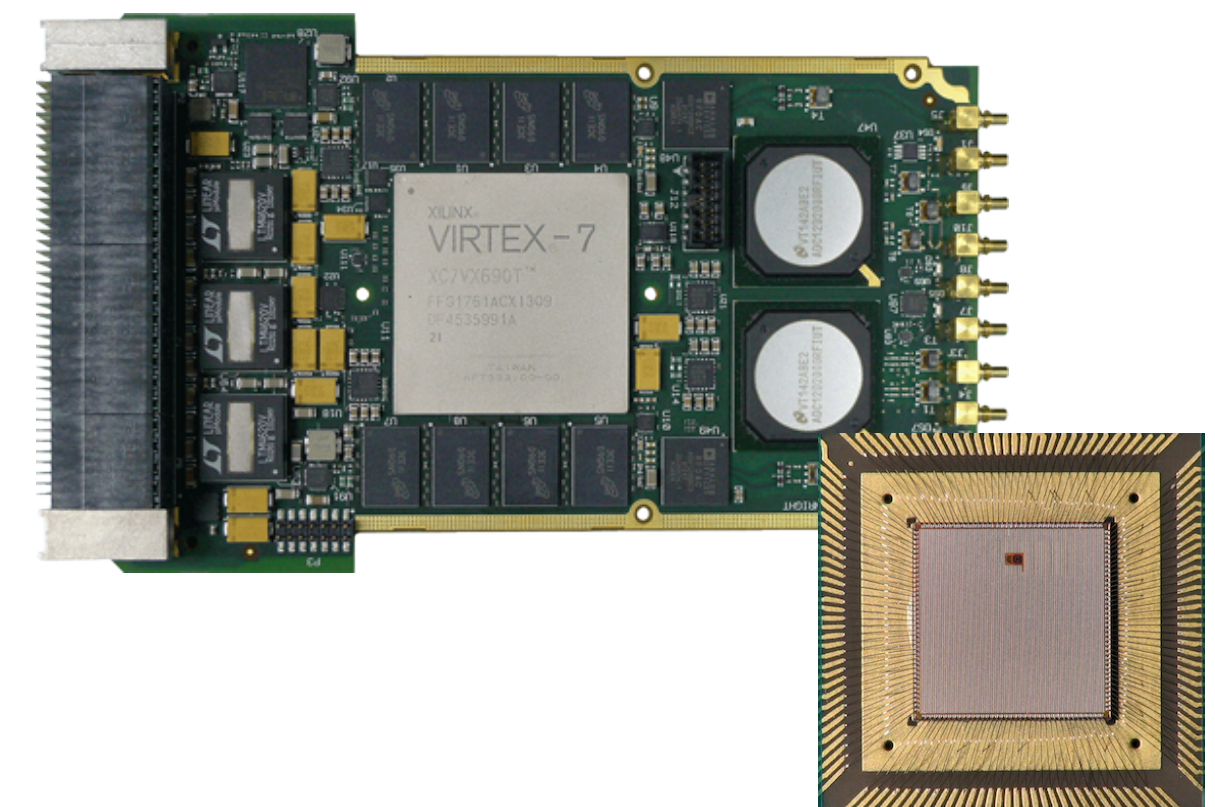
TensorFlow DF / scikit-learn / XGBoost

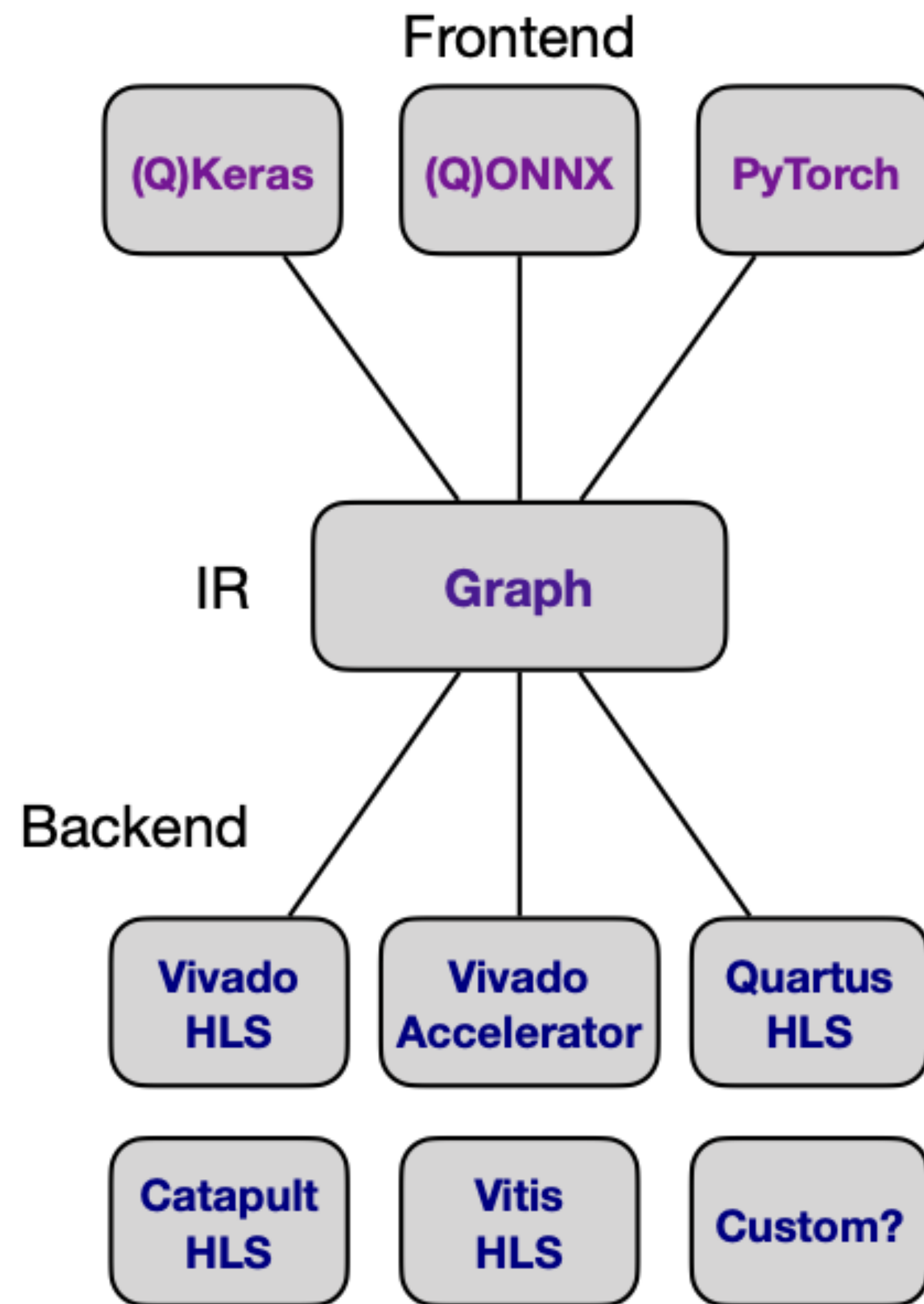


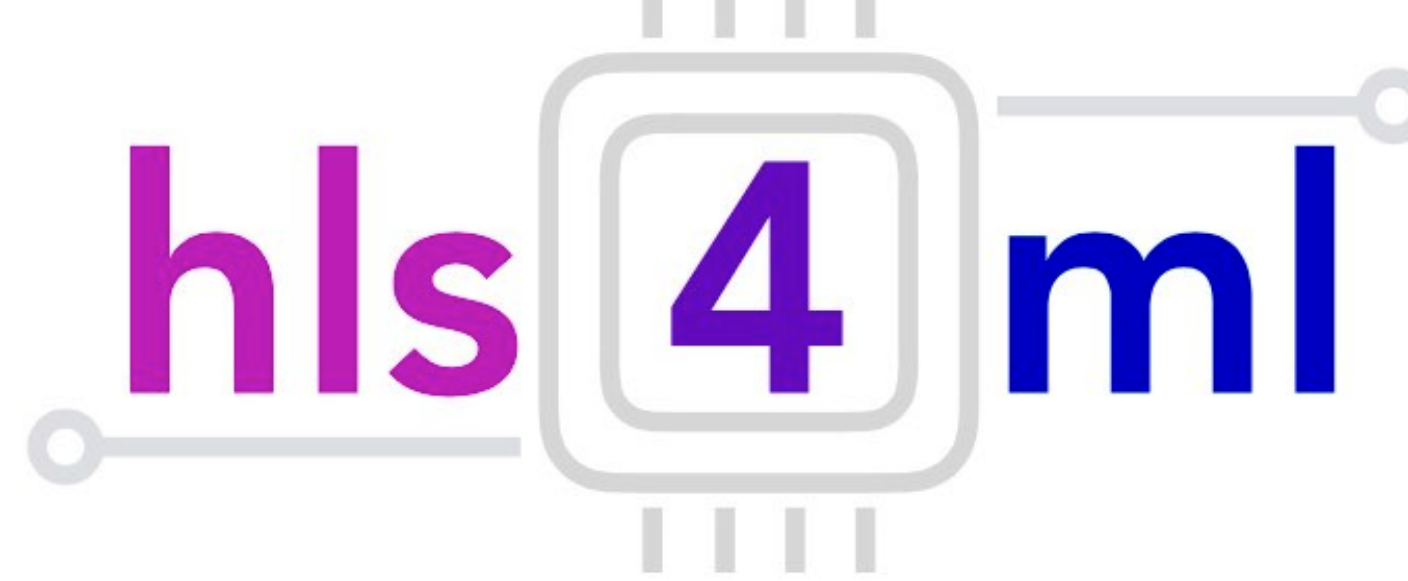
Conifer



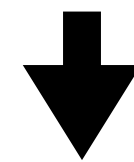
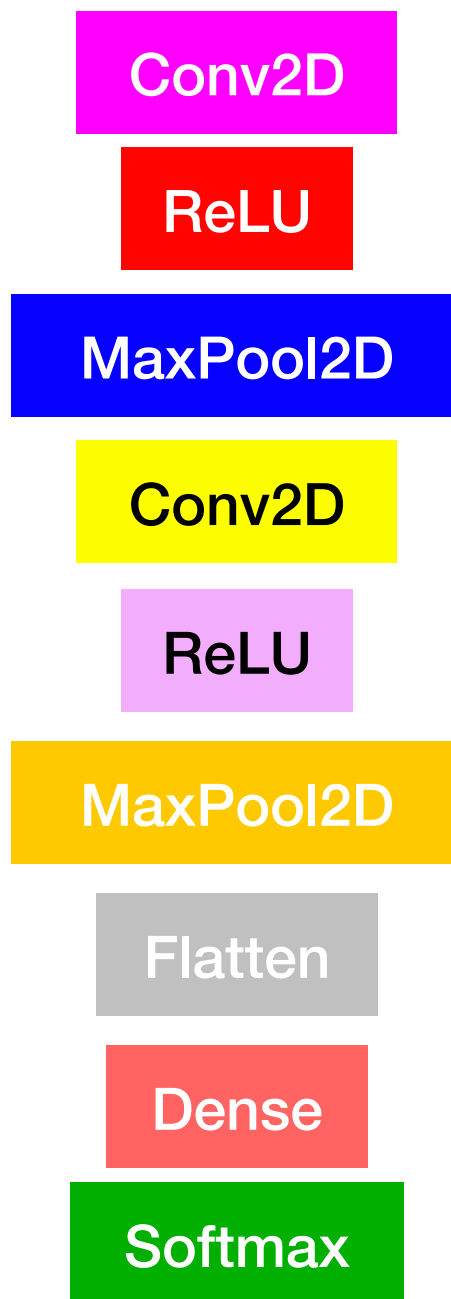
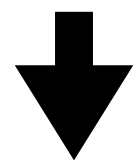
```
pip install hls4ml  
pip install conifer
```







0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9	9	9	9	9	9	9	9



Prediction

```

from hls4ml import ...
import tensorflow as tf

# train or load a model
model = ... # e.g. tf.keras.models.load_model(...)

# make a config template
cfg = config_from_keras_model(model,
granularity='name')

# tune the config
cfg['LayerName']['layer2']['ReuseFactor'] = 4

# do the conversion
hmodel = convert_from_keras_model(model, cfg)

# write and compile the HLS
hmodel.compile()

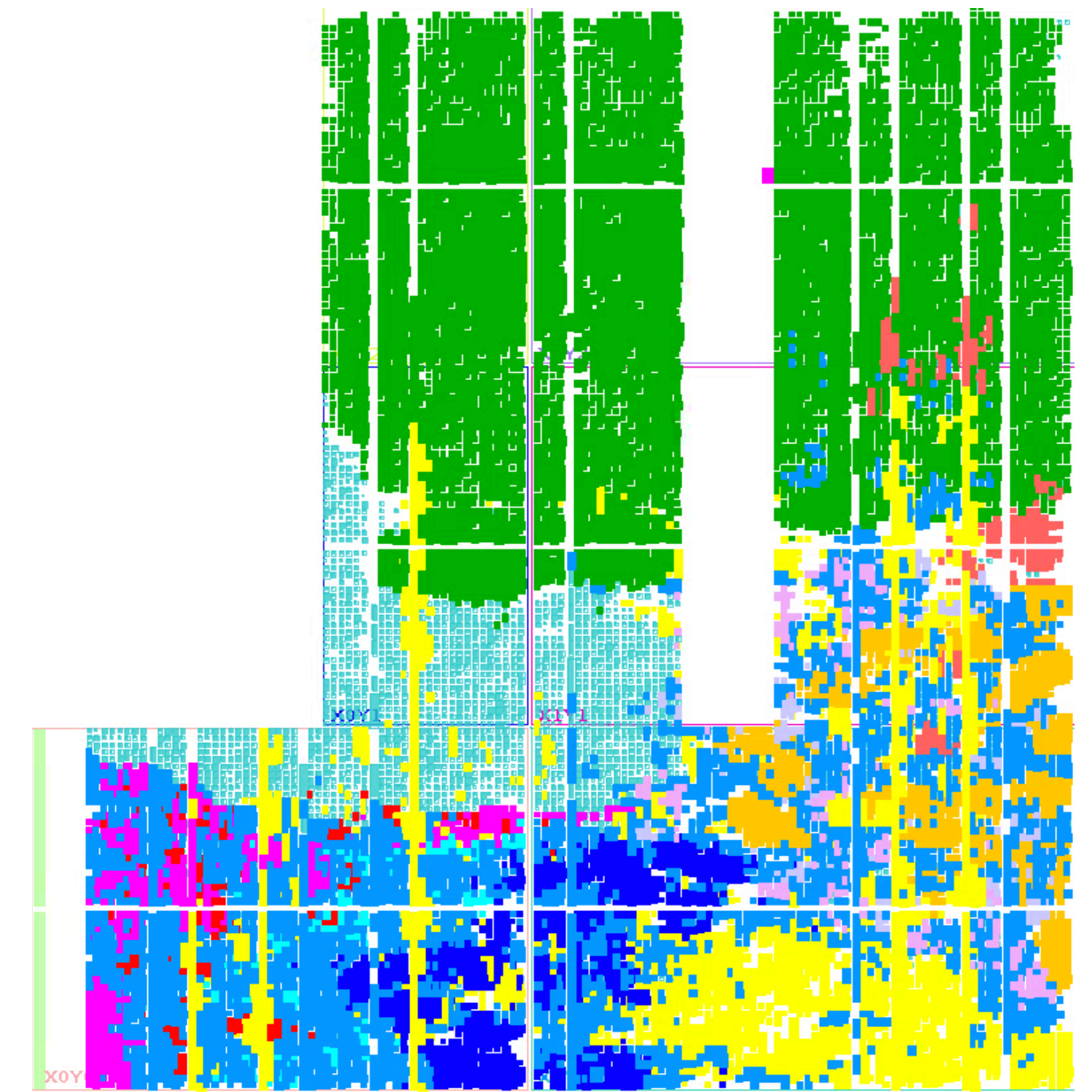
# run bit accurate emulation
y_tf = model.predict(x)
y_hls = hmodel.predict(x)

# do some validation
np.testing.assert_allclose(y_tf, y_hls)

# run HLS synthesis
hmodel.build()

```

pynq-z2 floorplan

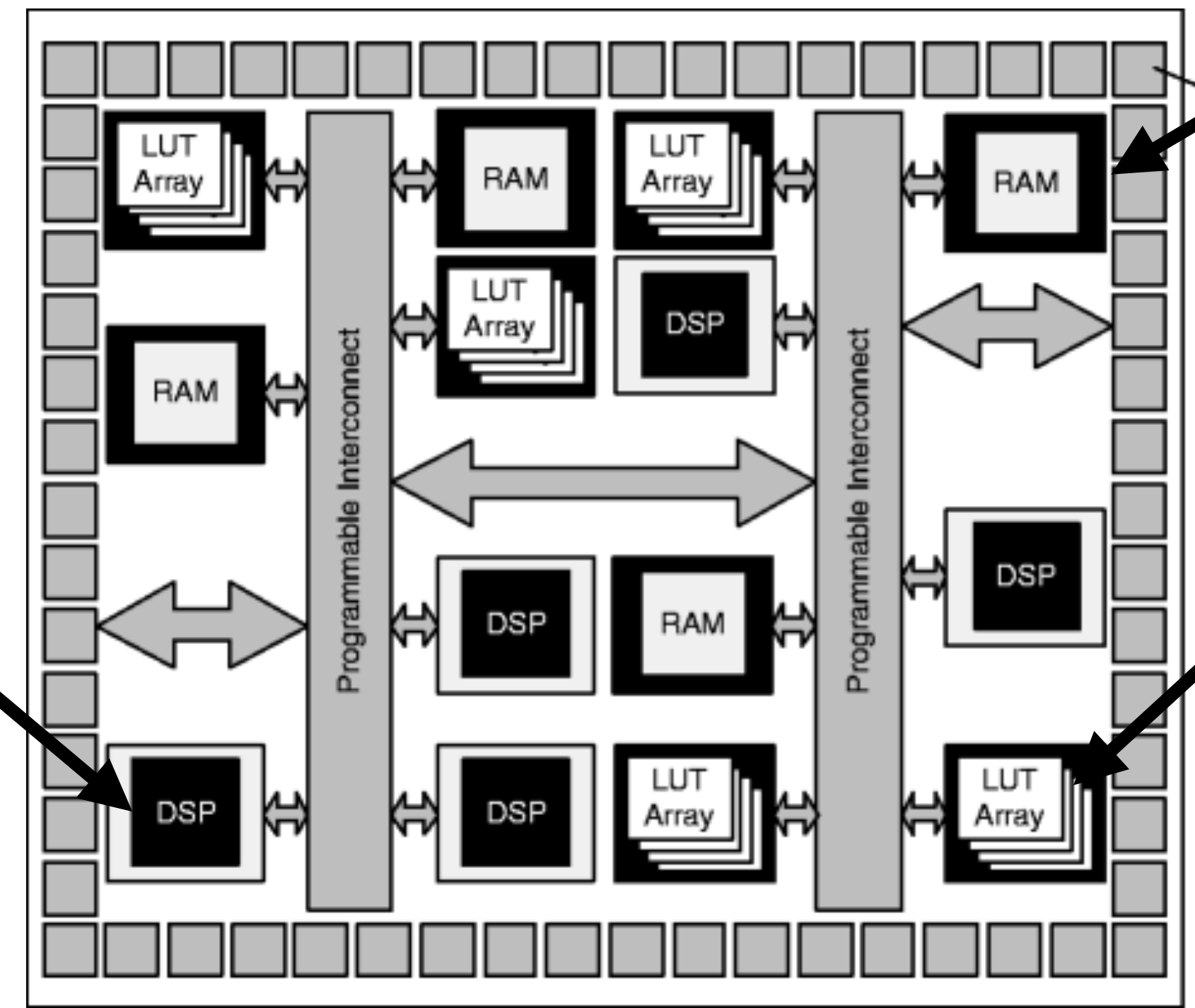


$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

↗ activation function
↑ multiplication
↖ addition

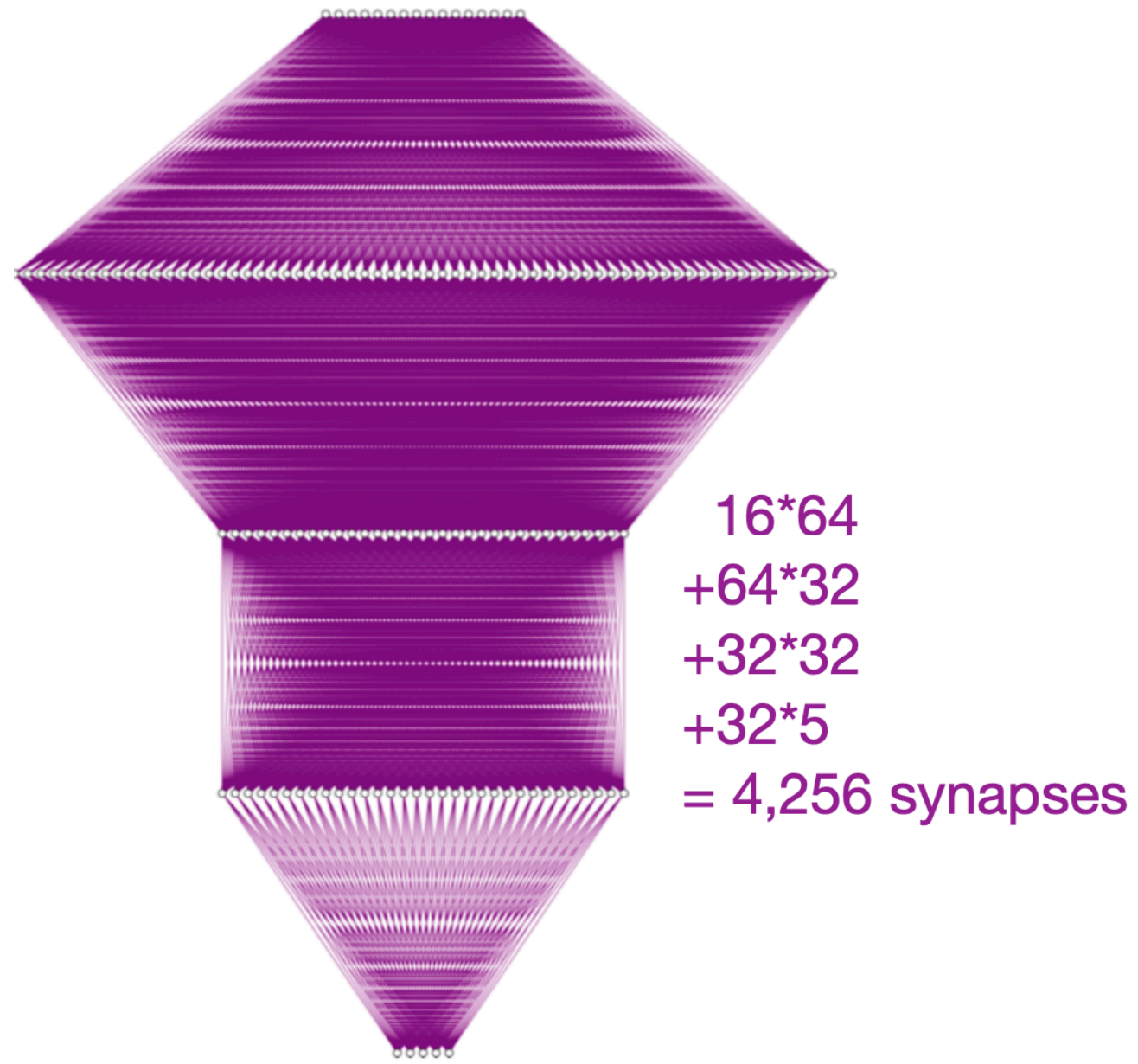
precomputed and stored in BRAMs DSPs logic cells

Digital signal processors (DSPs)
O(5,000) units



Memory (BRAM)
O(2000) units

Logic cells/lookup tables (LUTs)
O(1) million units

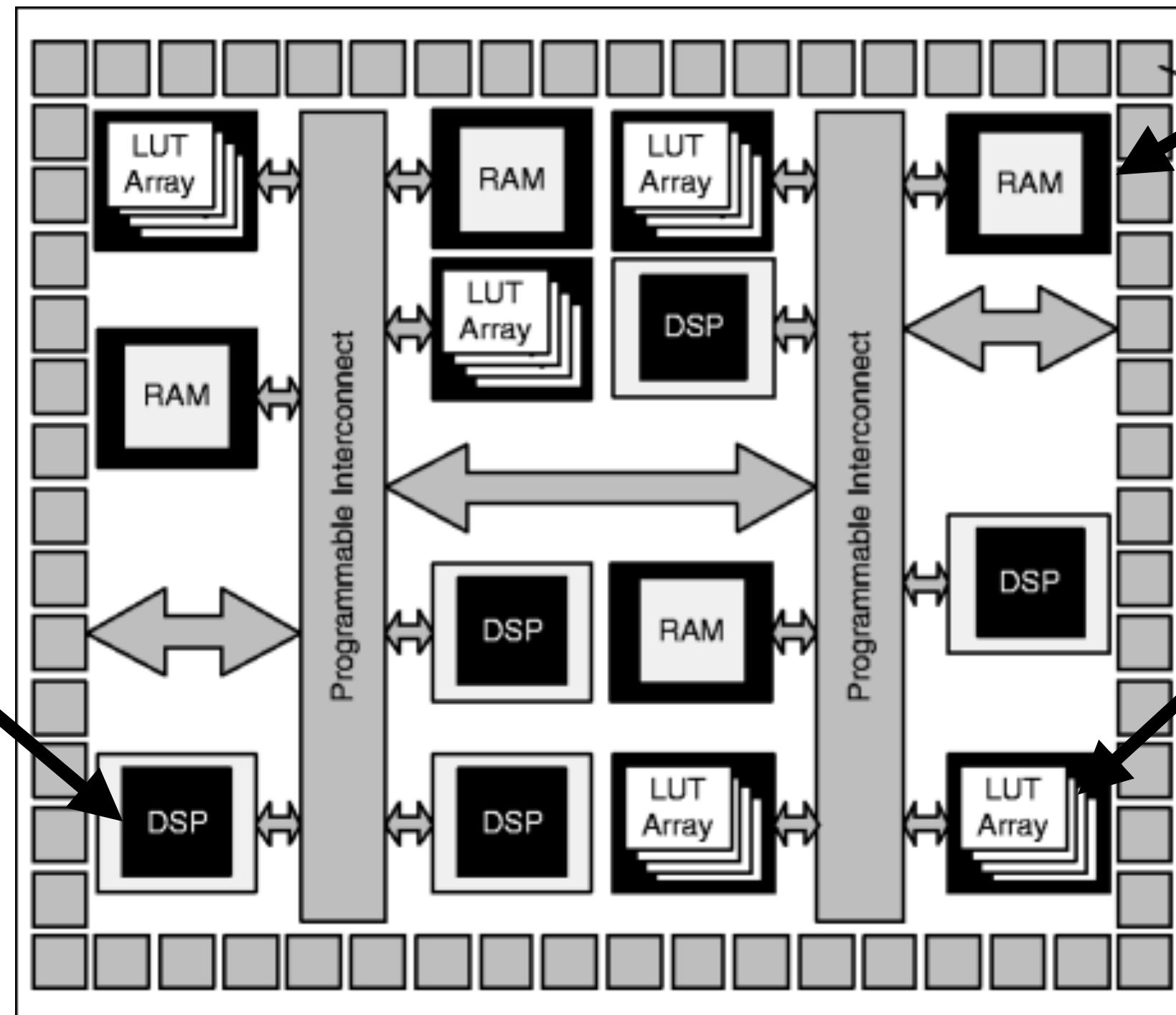


Network size limited by # multiplications

$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

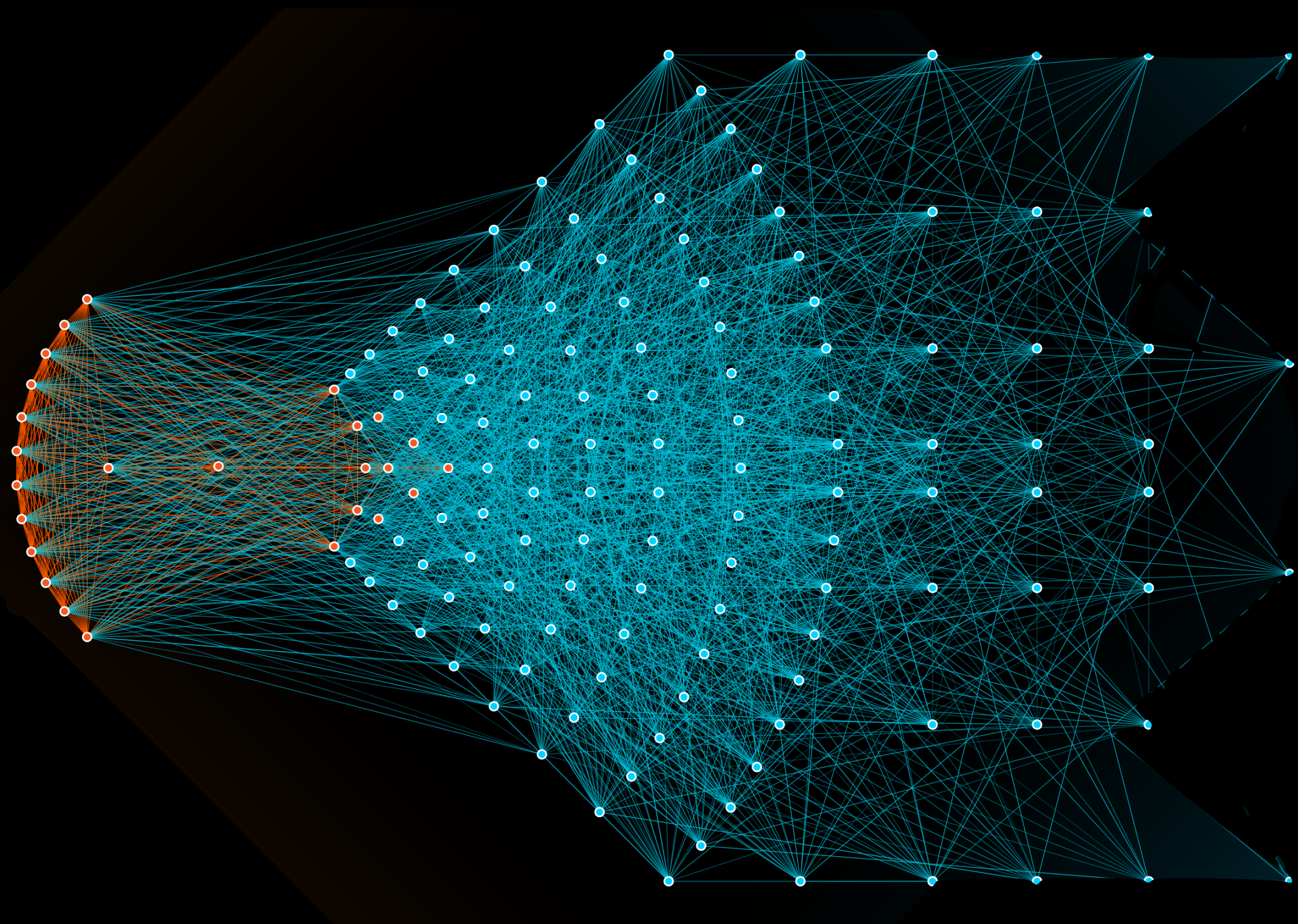
activation function multiplication addition
 precomputed and stored in BRAMs DSPs logic cells

Digital signal processors (DSPs)
O(5,000) units

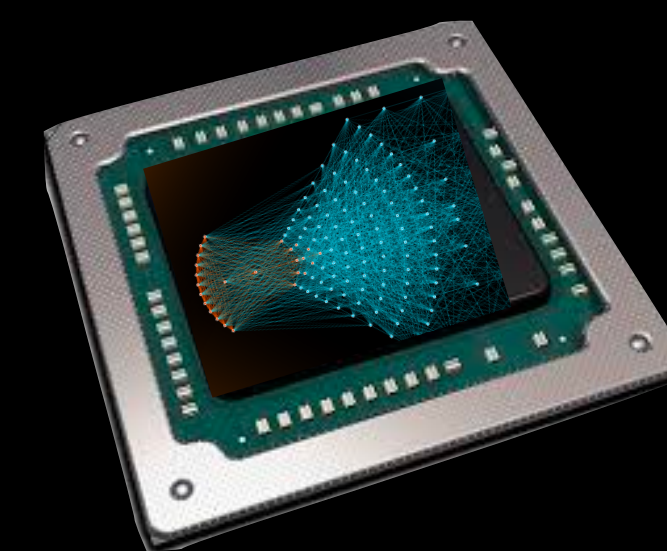


Memory (BRAM)
O(2000) units

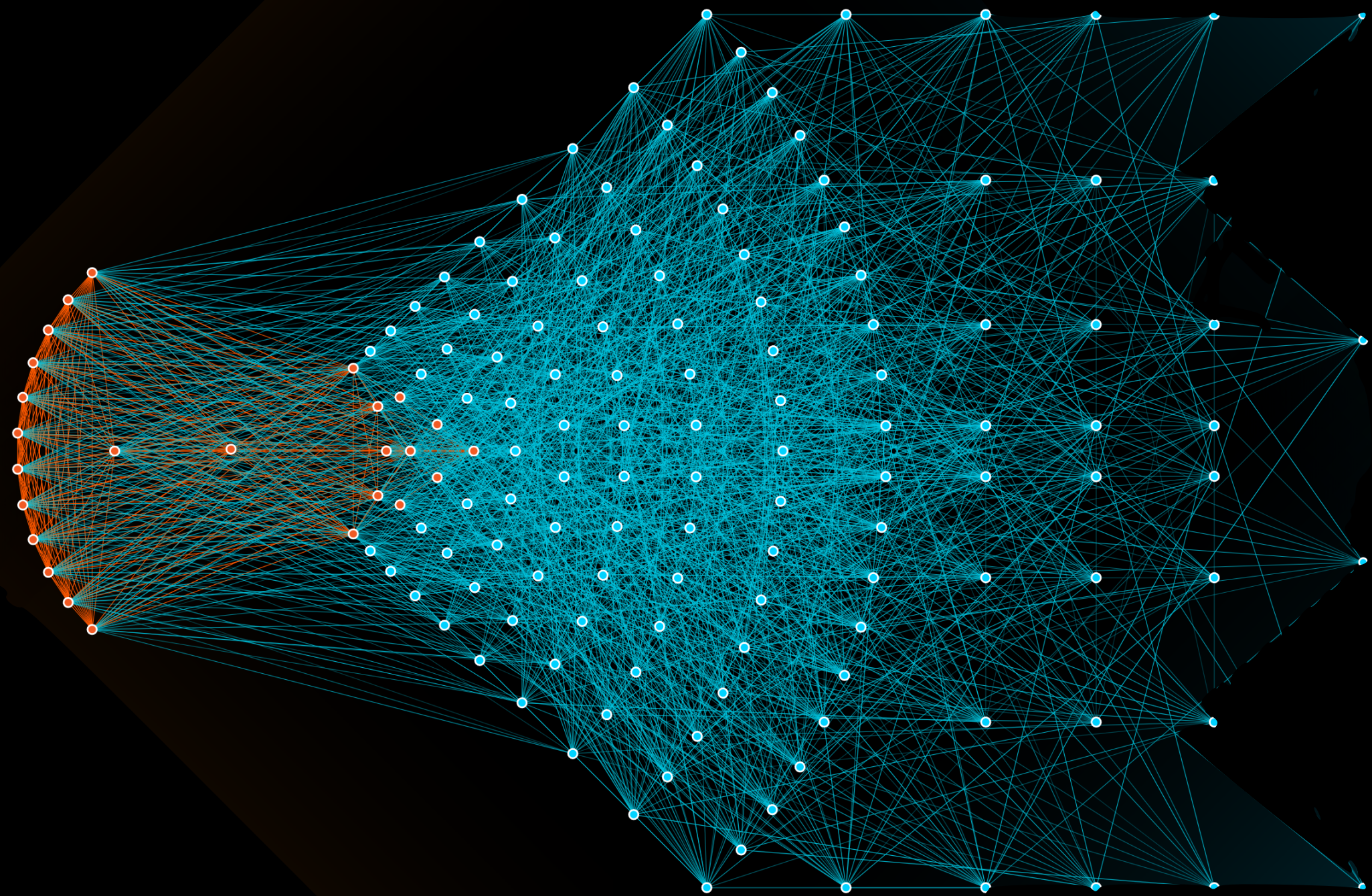
Logic cells/lookup tables (LUTs)
O(1) million units



Ideally



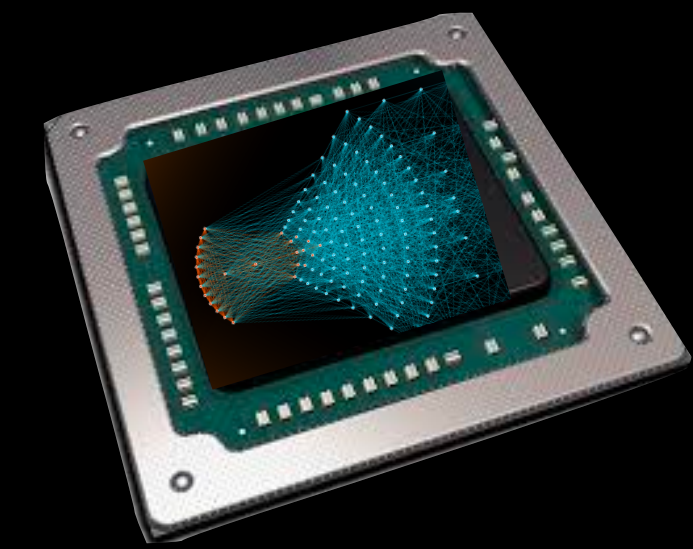
Reality



Ideally

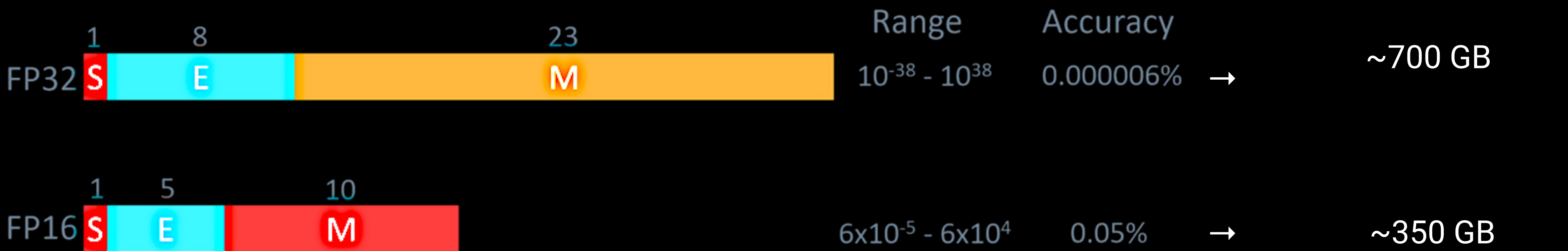


- Quantization
- Pruning
- Parallelisation
- Knowledge distillation



Reality

FP16 vs FP32





Memory:

- ~ FP16 x2 smaller than FP32
- ~ int8 x4 smaller than FP32

Speed:

- ~ FP16 x8 times faster than FP32

Fixed-point arithmetic:

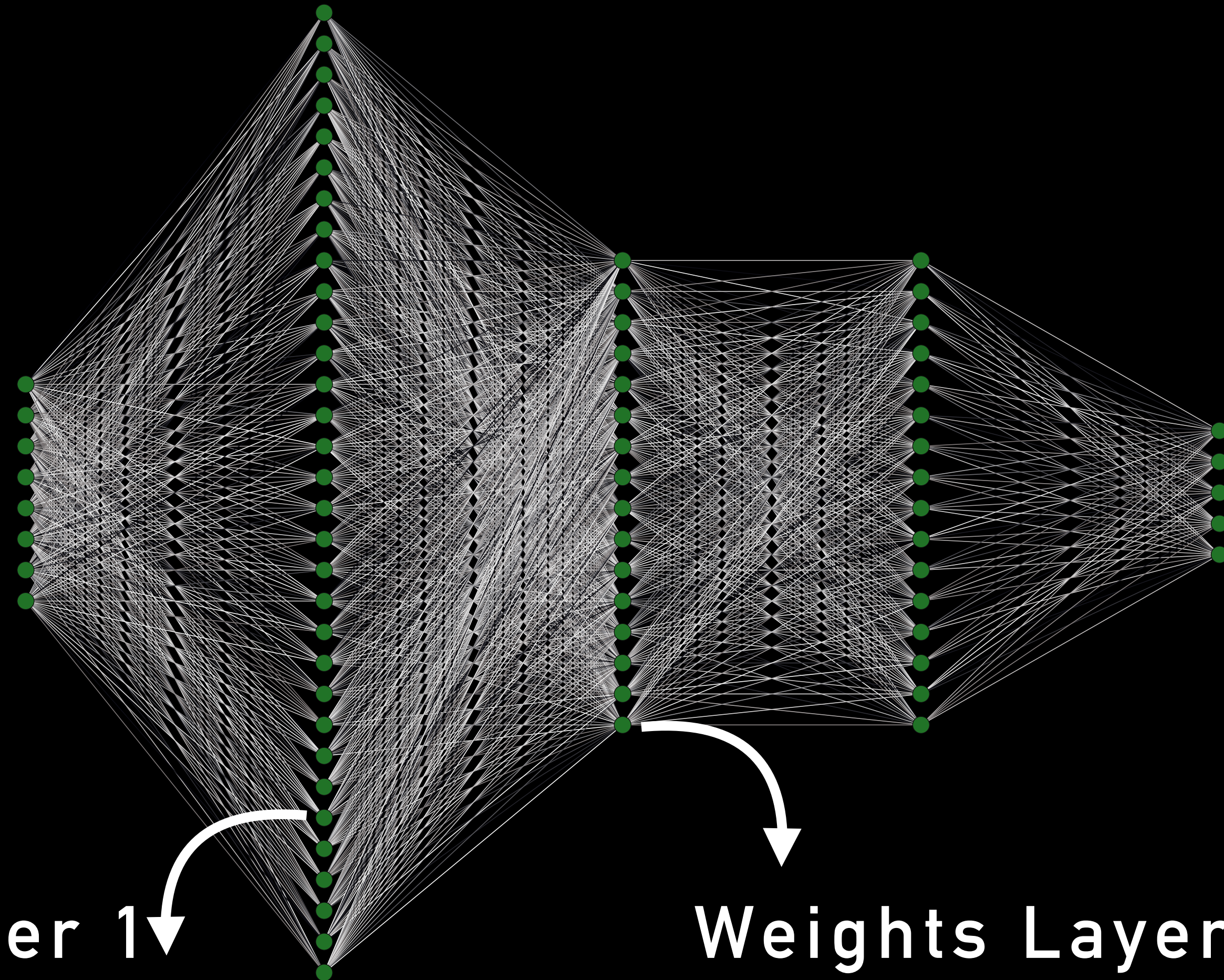
Any precision!

	Hopper	Ampere
Supported Tensor Core precisions	FP64, TF32, bfloat16, FP16, FP8, INT8	FP64, TF32, bfloat16, FP16, INT8, INT4, INT1

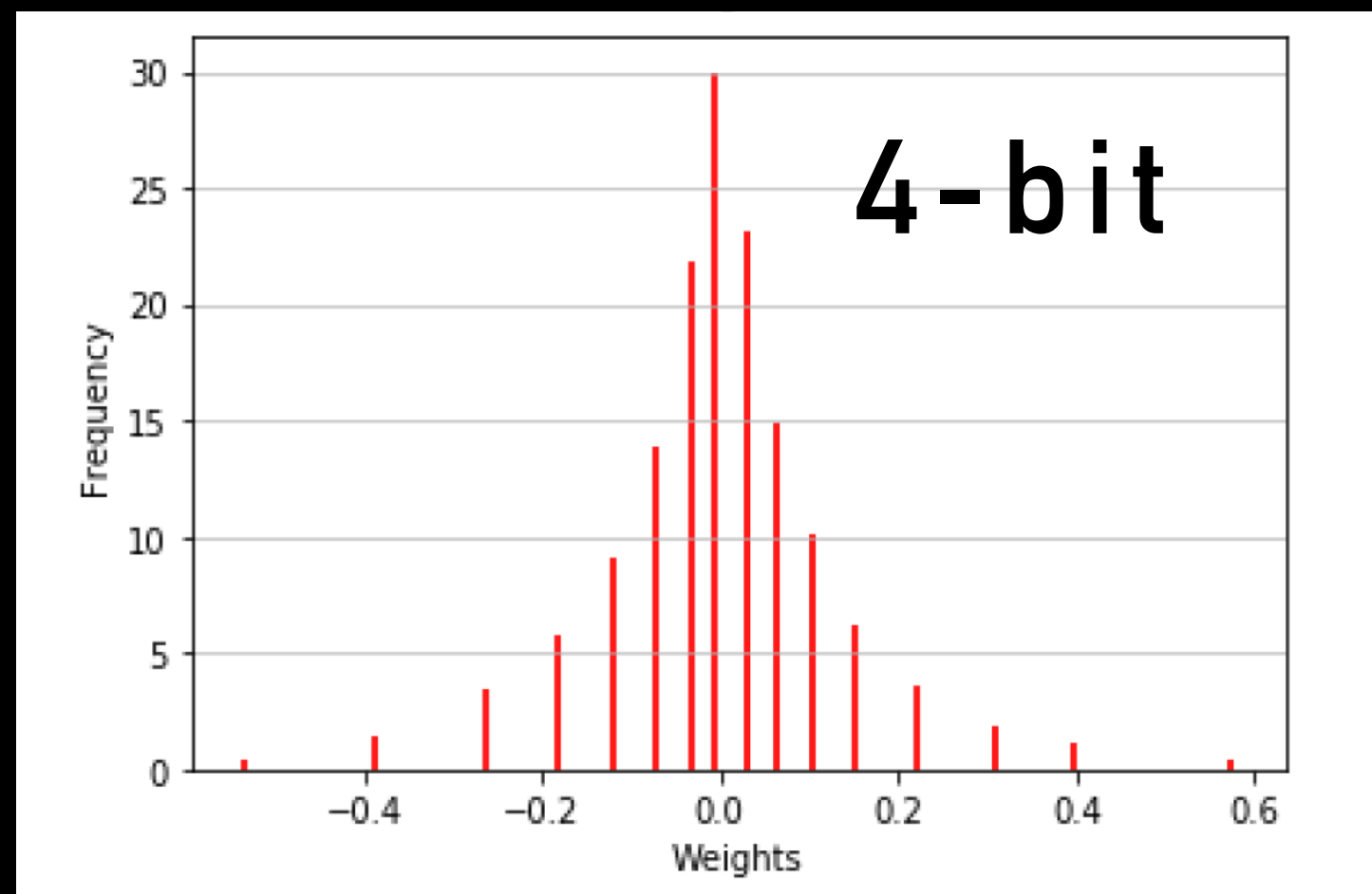
Quantization

2⁴ numbers in $s^*[-8, +7]$

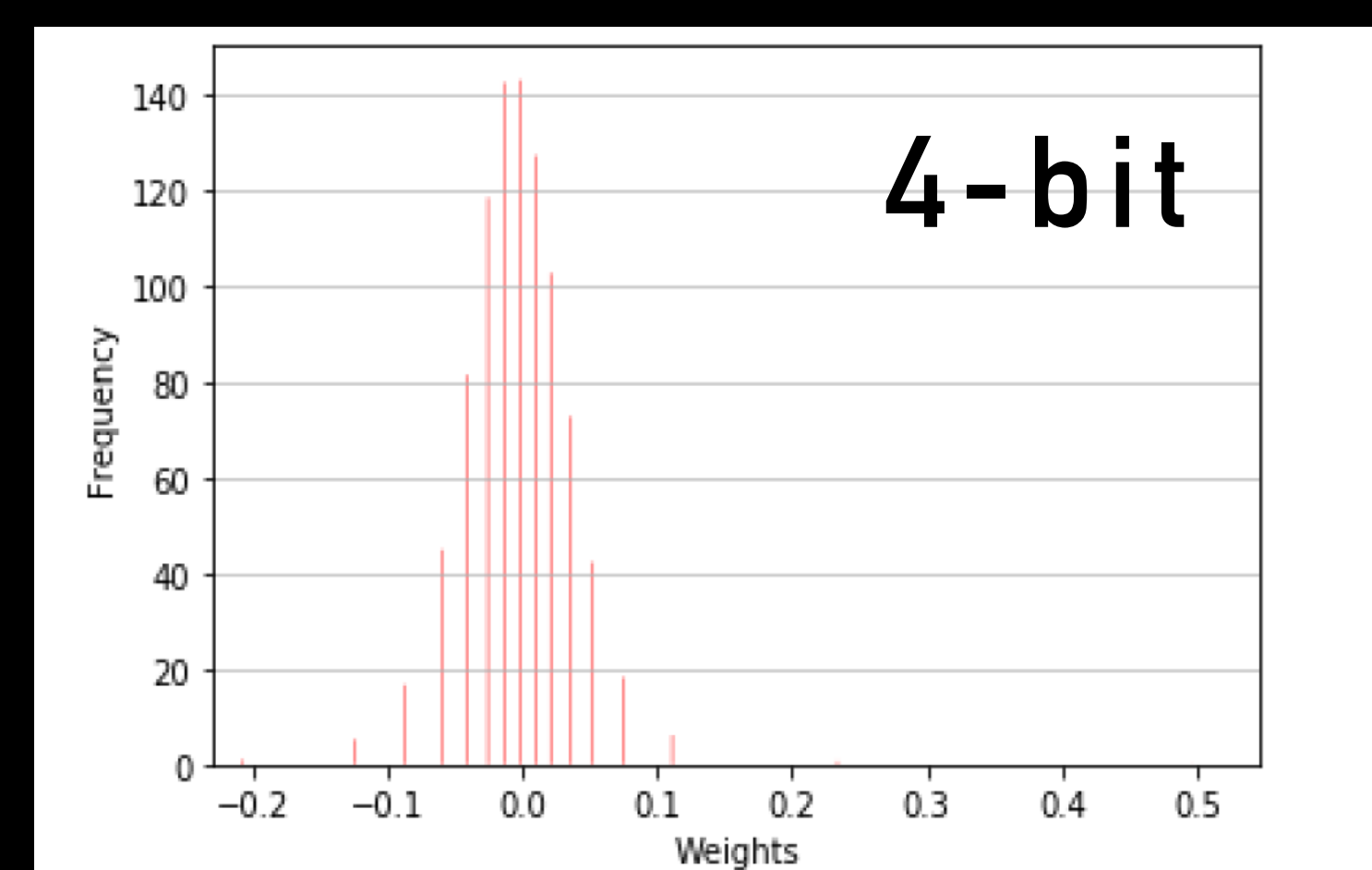
Fixed point



Weights Layer 1



Weights Layer 2

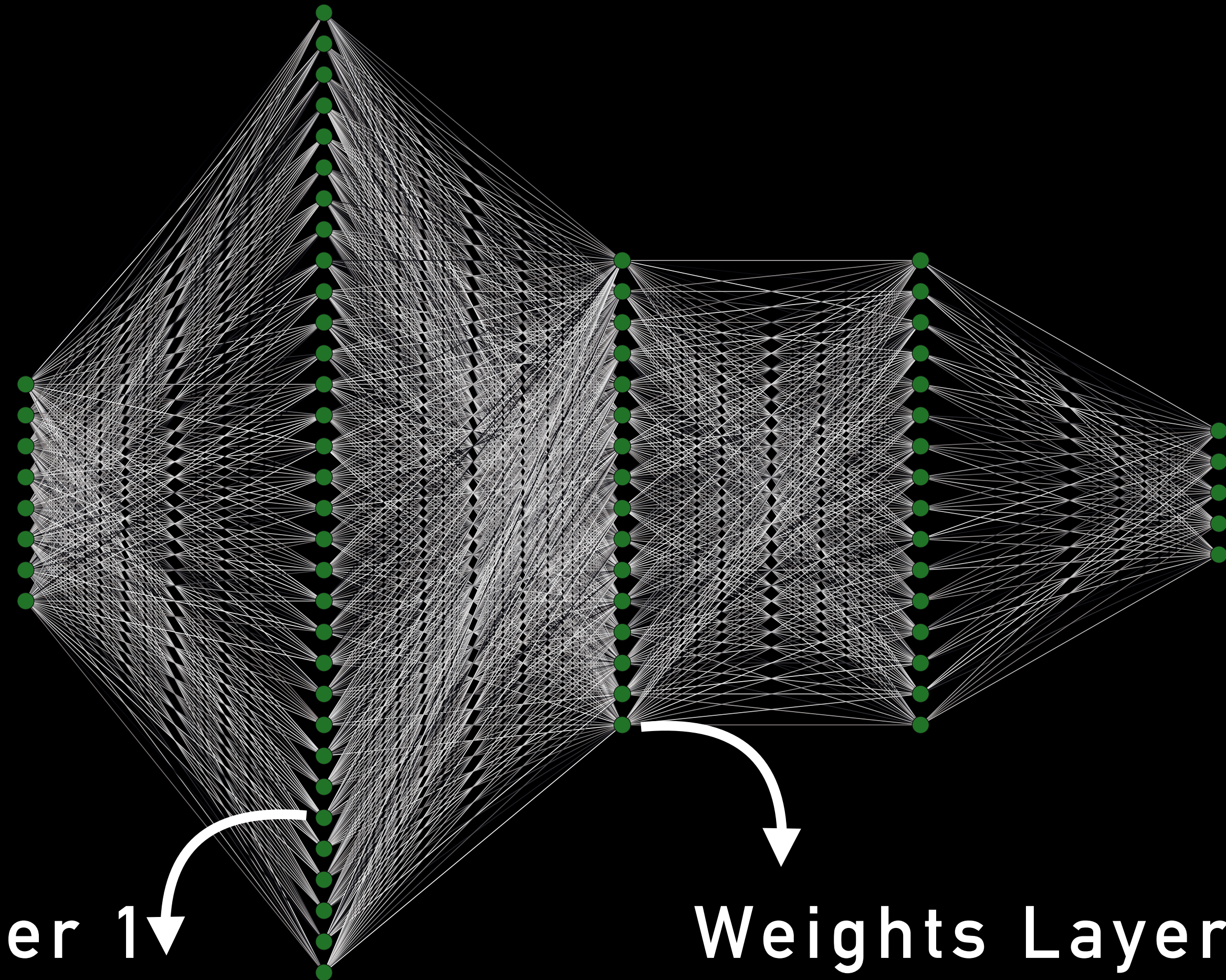


Quantization

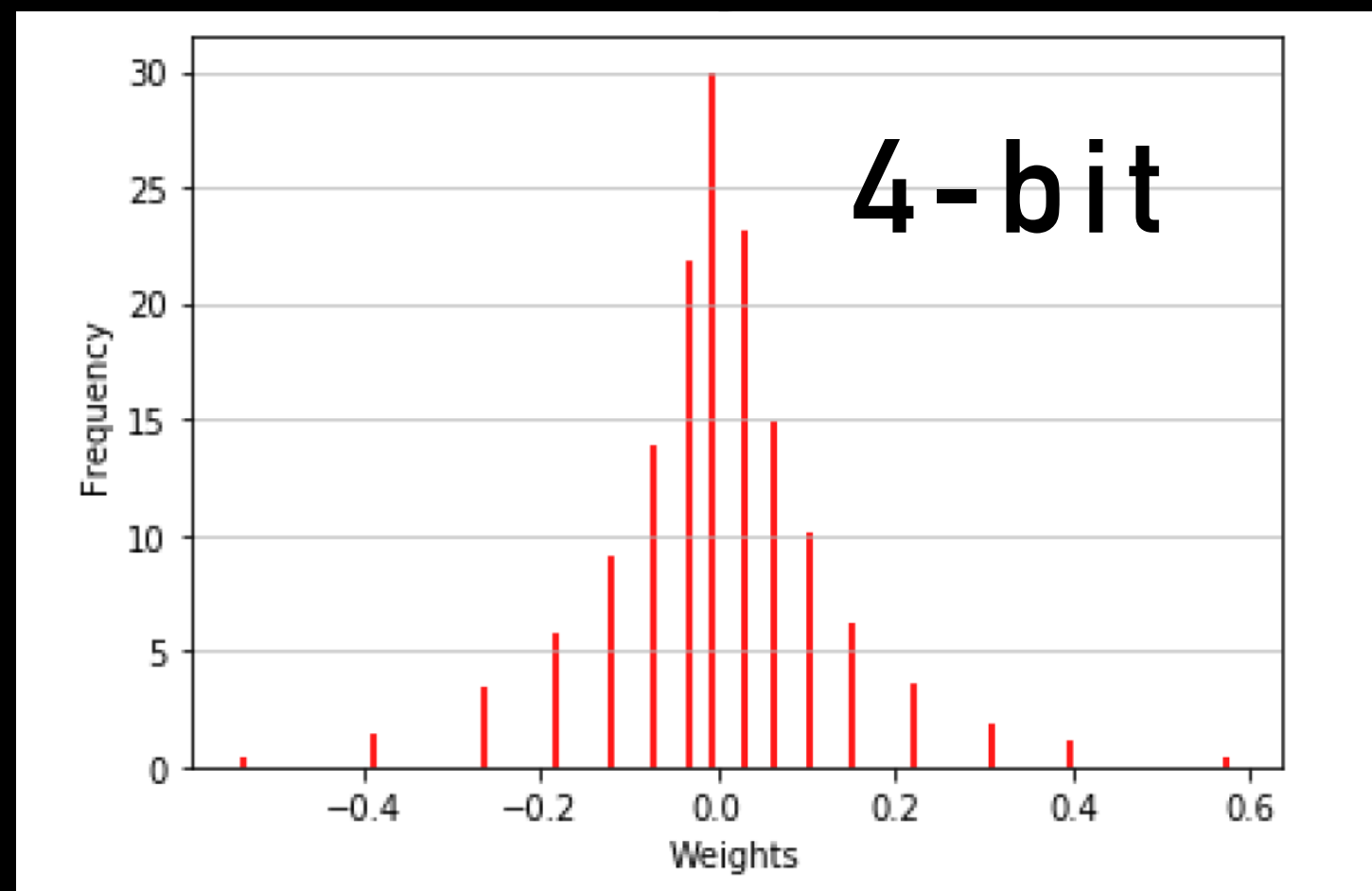
2⁴ numbers in s*[-8, +7]

Fixed point

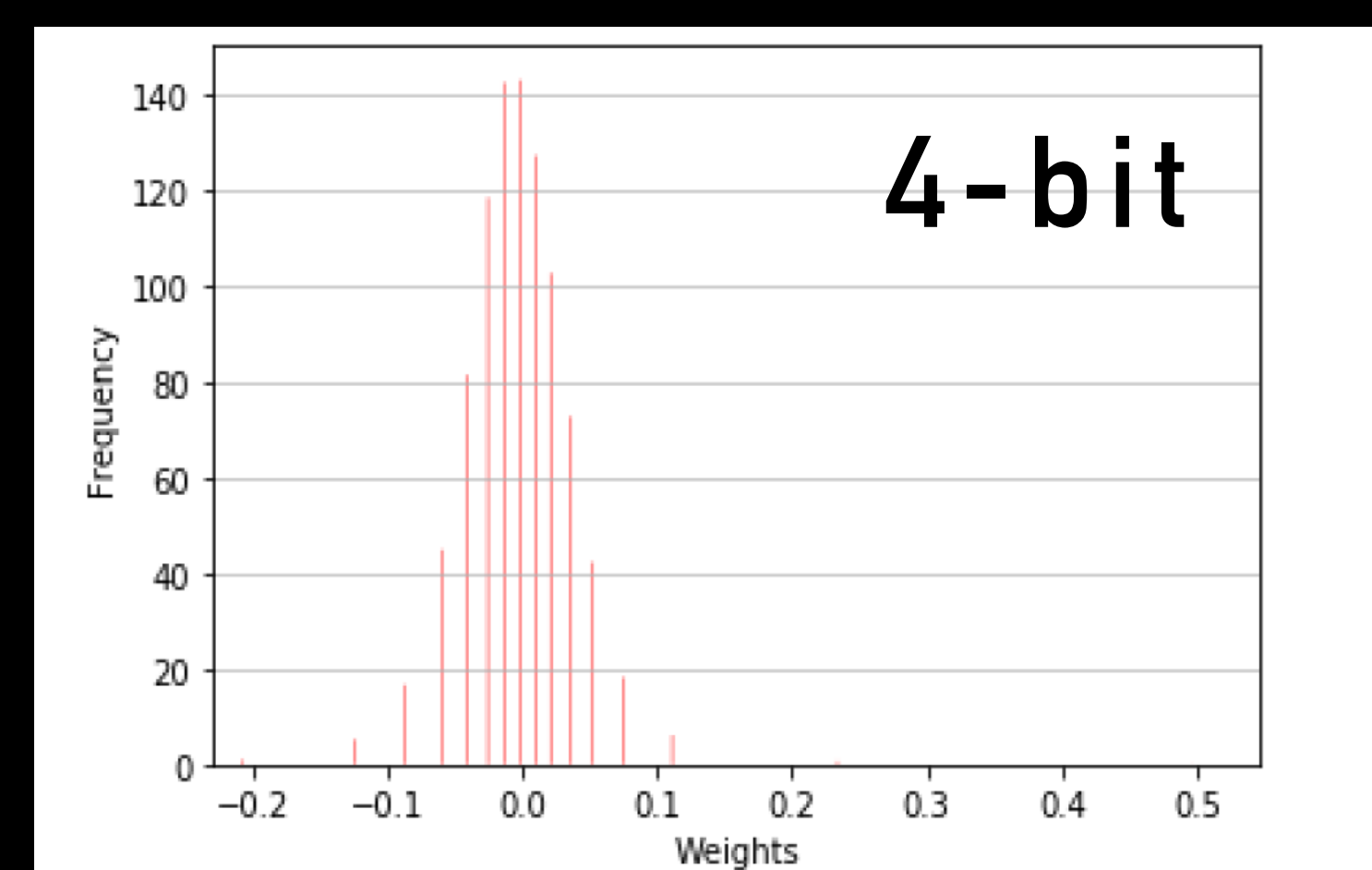
0101.1011101010



Weights Layer 1

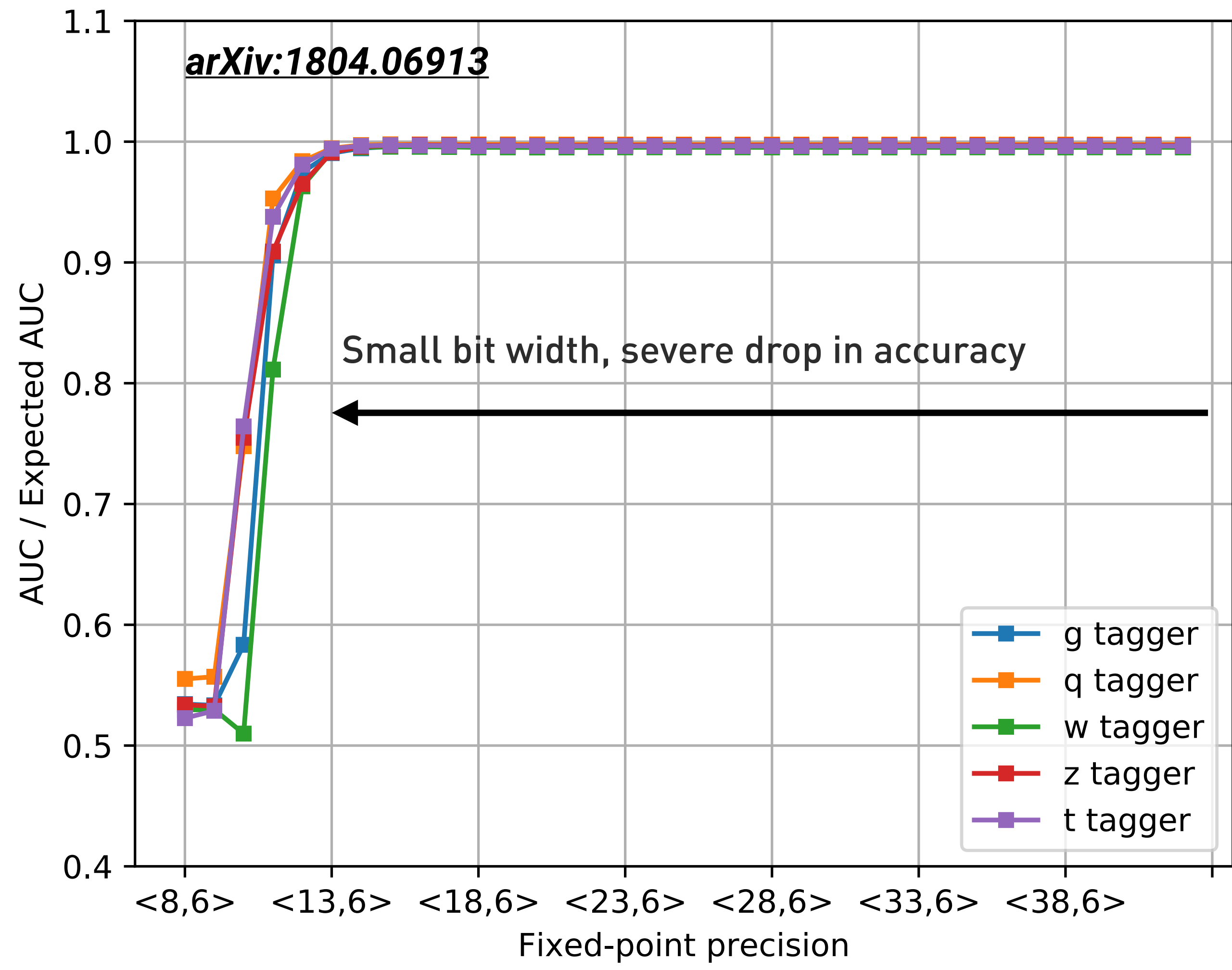


Weights Layer 2

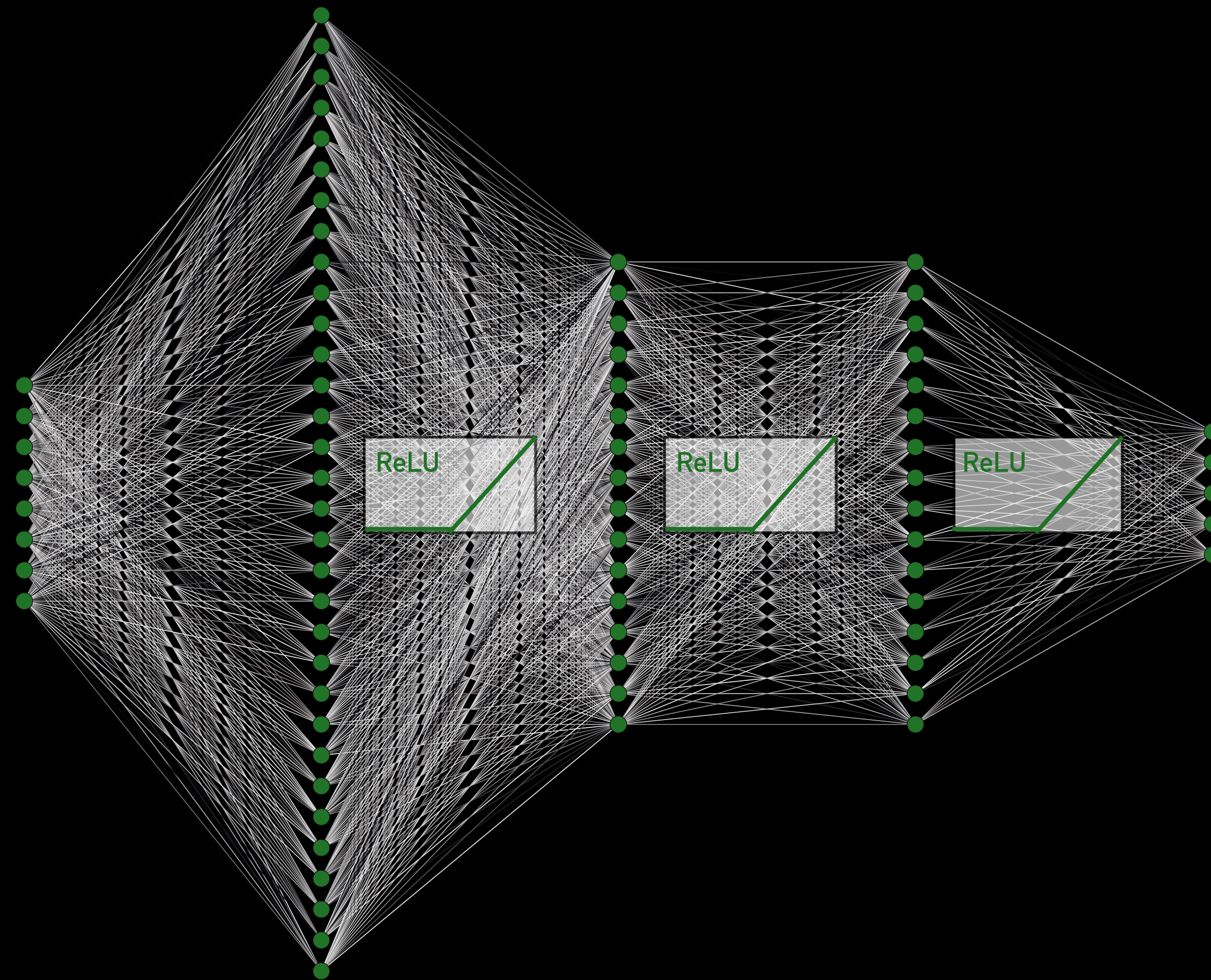


hls4ml

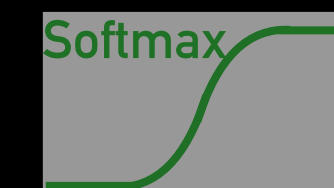
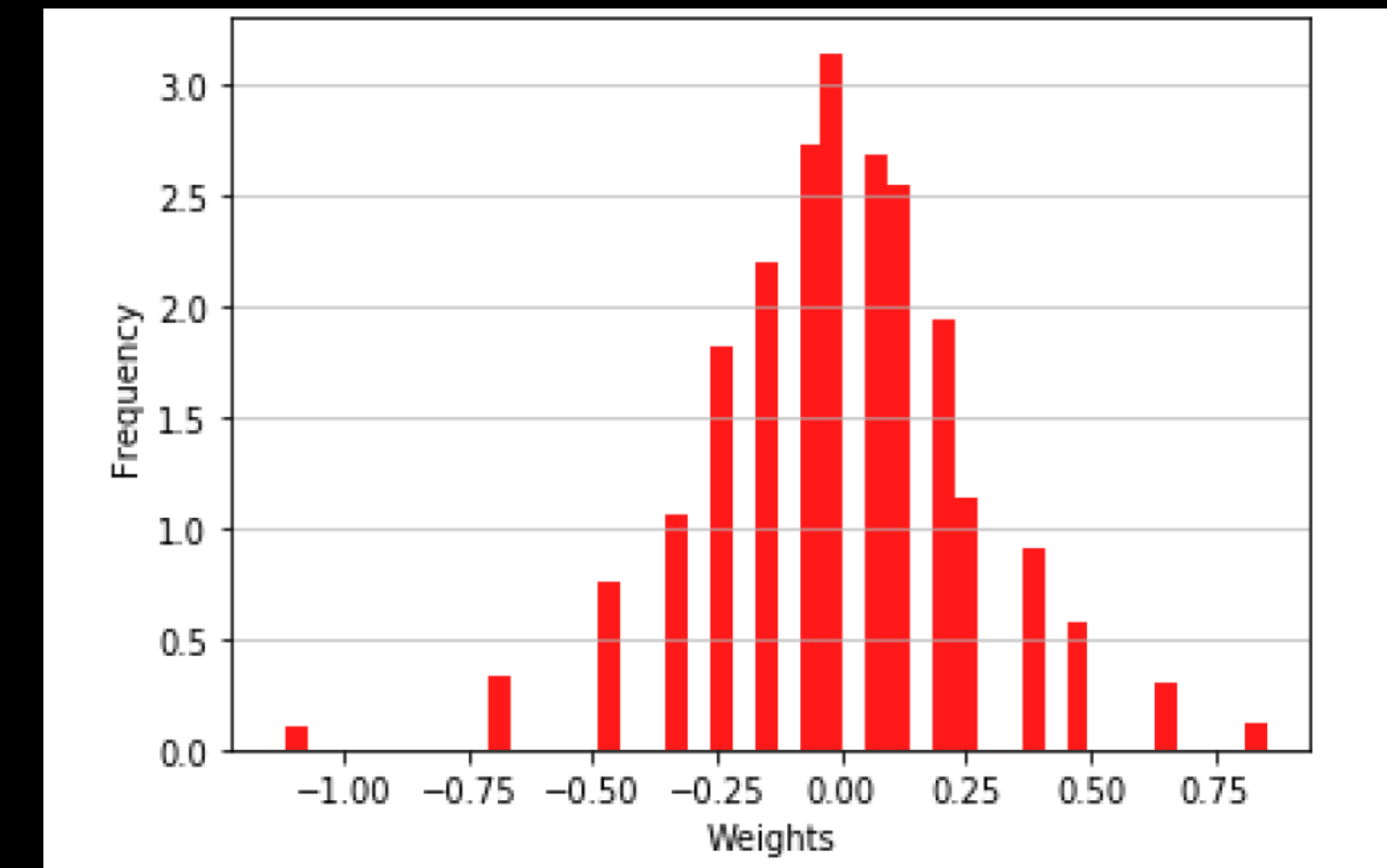
arXiv:1804.06913



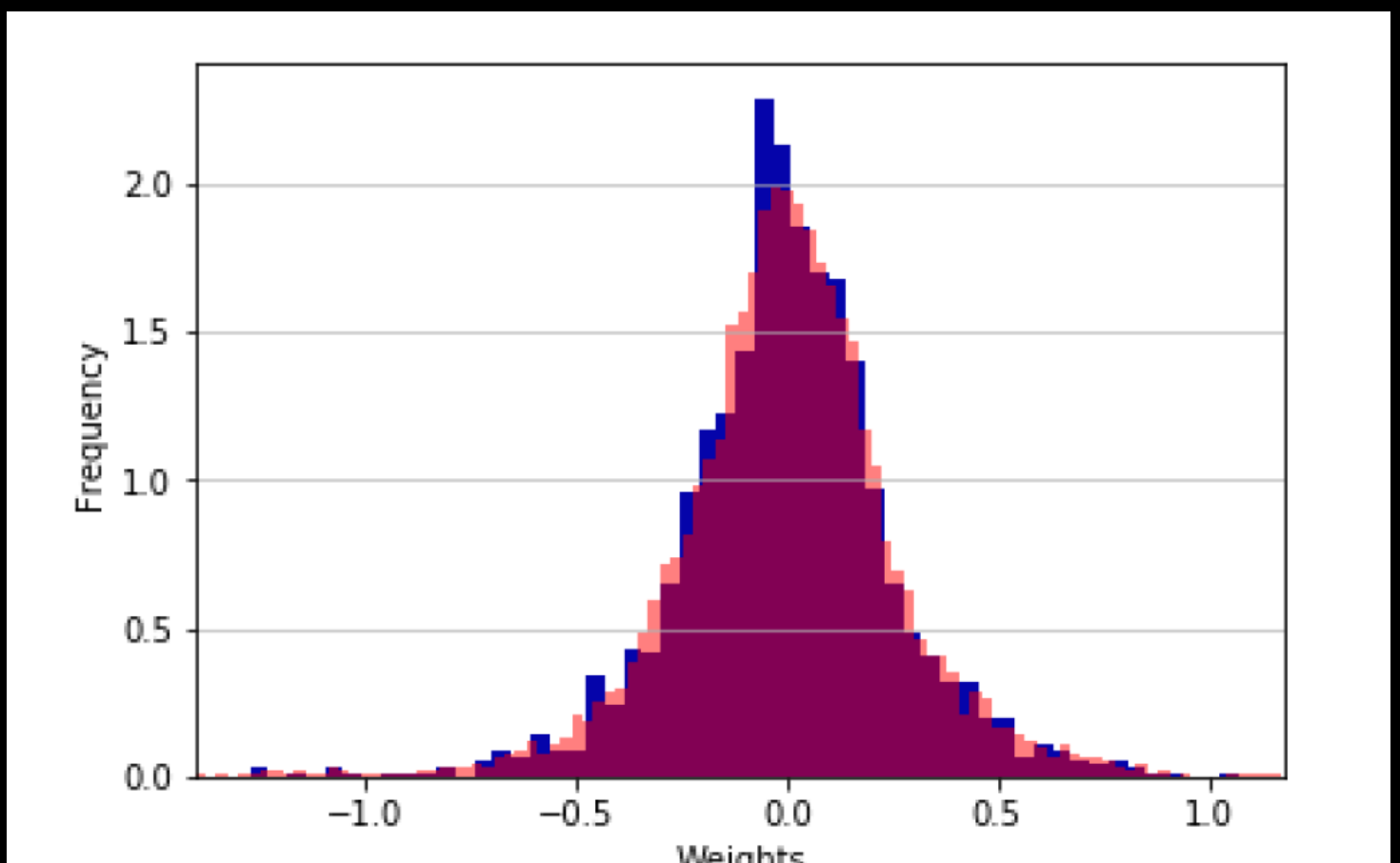
hls4ml + Google Quantization-aware training



Forward pass →

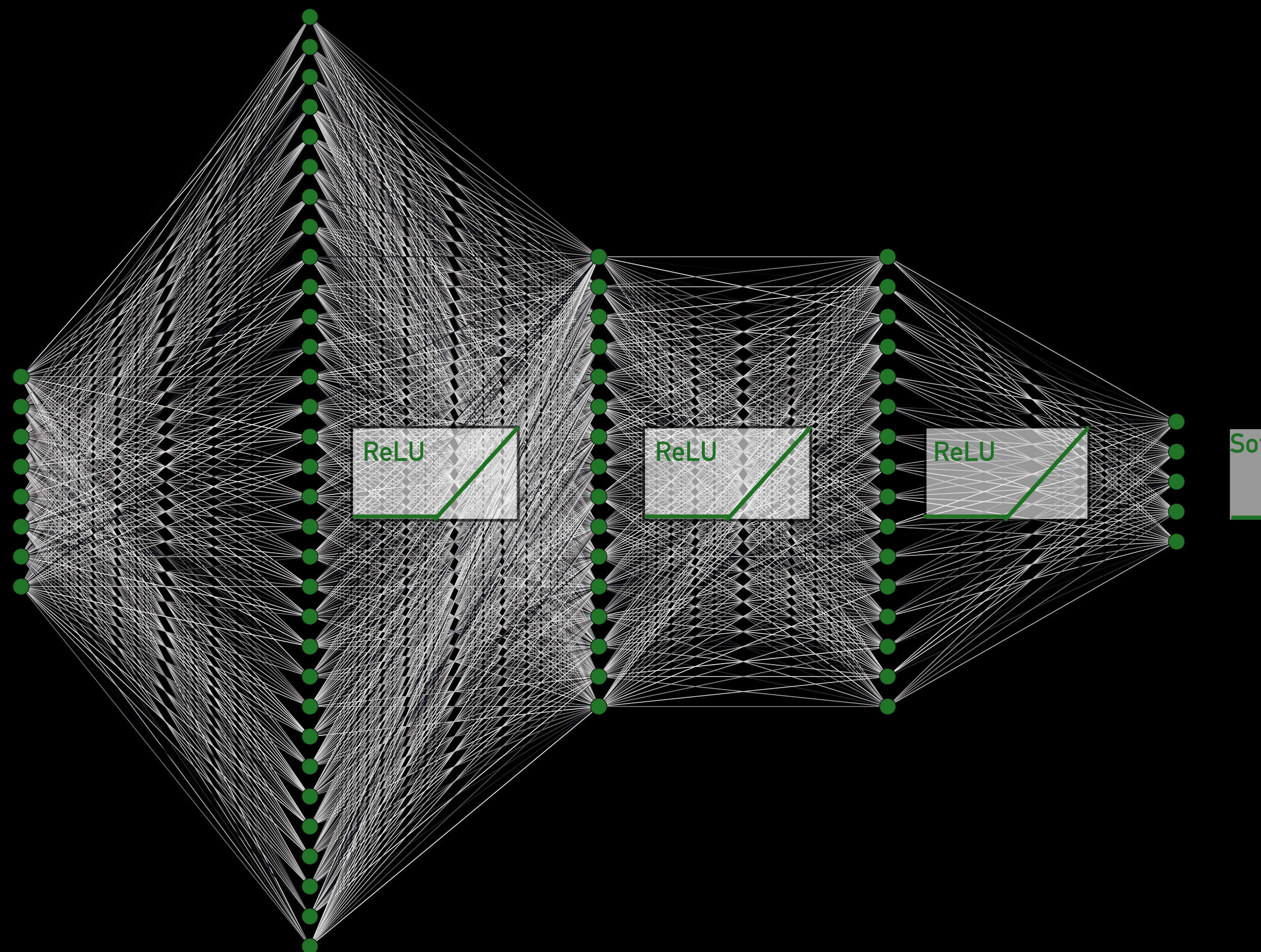


← Back propagation



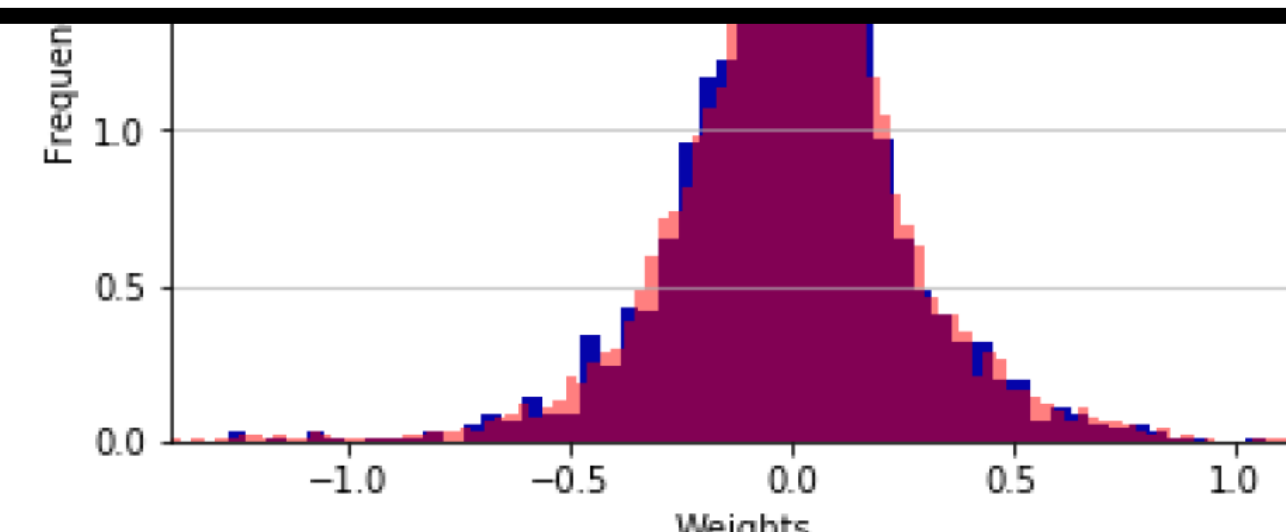
hls4ml + Google Quantization-aware training

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)
```



Estimating energy and size

Some layers **more accommodating** for aggressive quantization, others require expensive arithmetic

- heterogeneous quantization

Estimating energy and size

Some layers **more accommodating** for aggressive quantization, others require expensive arithmetic

- heterogeneous quantization

For edge inference, need best possible quantization configuration for

- Highest accuracy ↑...
- ... and lowest resource consumption ↓

→ hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

Estimating energy and size

Some layers **more accommodating** for aggressive quantization, others require expensive arithmetic

- heterogeneous quantization

For edge inference, need best possible quantization configuration for

- Highest accuracy ↑...
- ... and lowest resource consumption ↓

→ hyper-parameter scan over quantizers which considers energy and accuracy simultaneously

QTools: Estimate QKeras model bit and energy consumption, assuming 45 nm **Horowitz process**

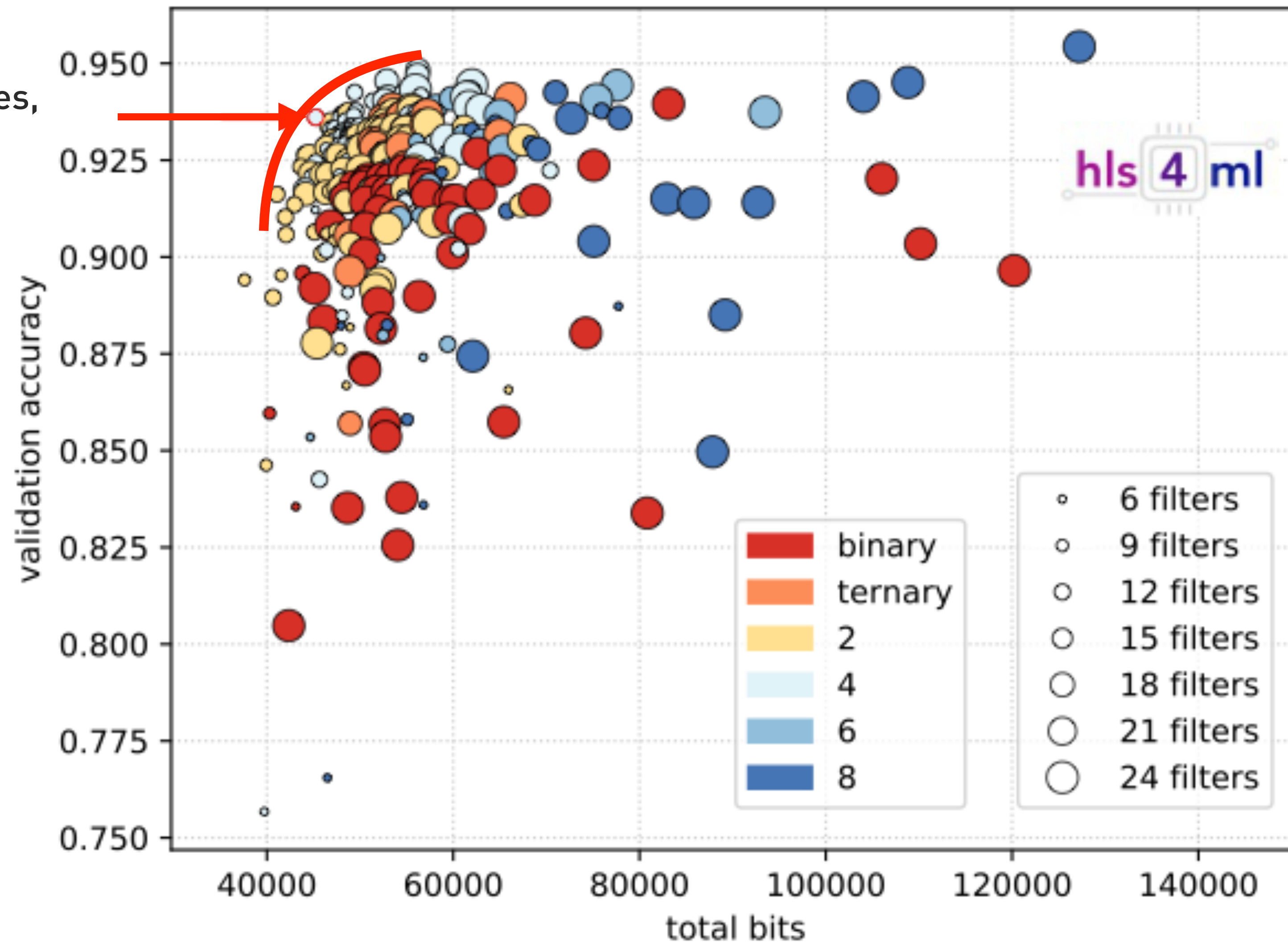
Model	Accuracy [%]	Per-layer energy consumption [pJ]								Total energy [μ J]	Total bits
		Dense	ReLU	Dense	ReLU	Dense	ReLU	Dense	Softmax		
BF	74.4	1735	53	3240	27	1630	27	281	11	0.00700	61446
Q6	74.8	794	23	1120	11	562	11	99	11	0.00263	26334

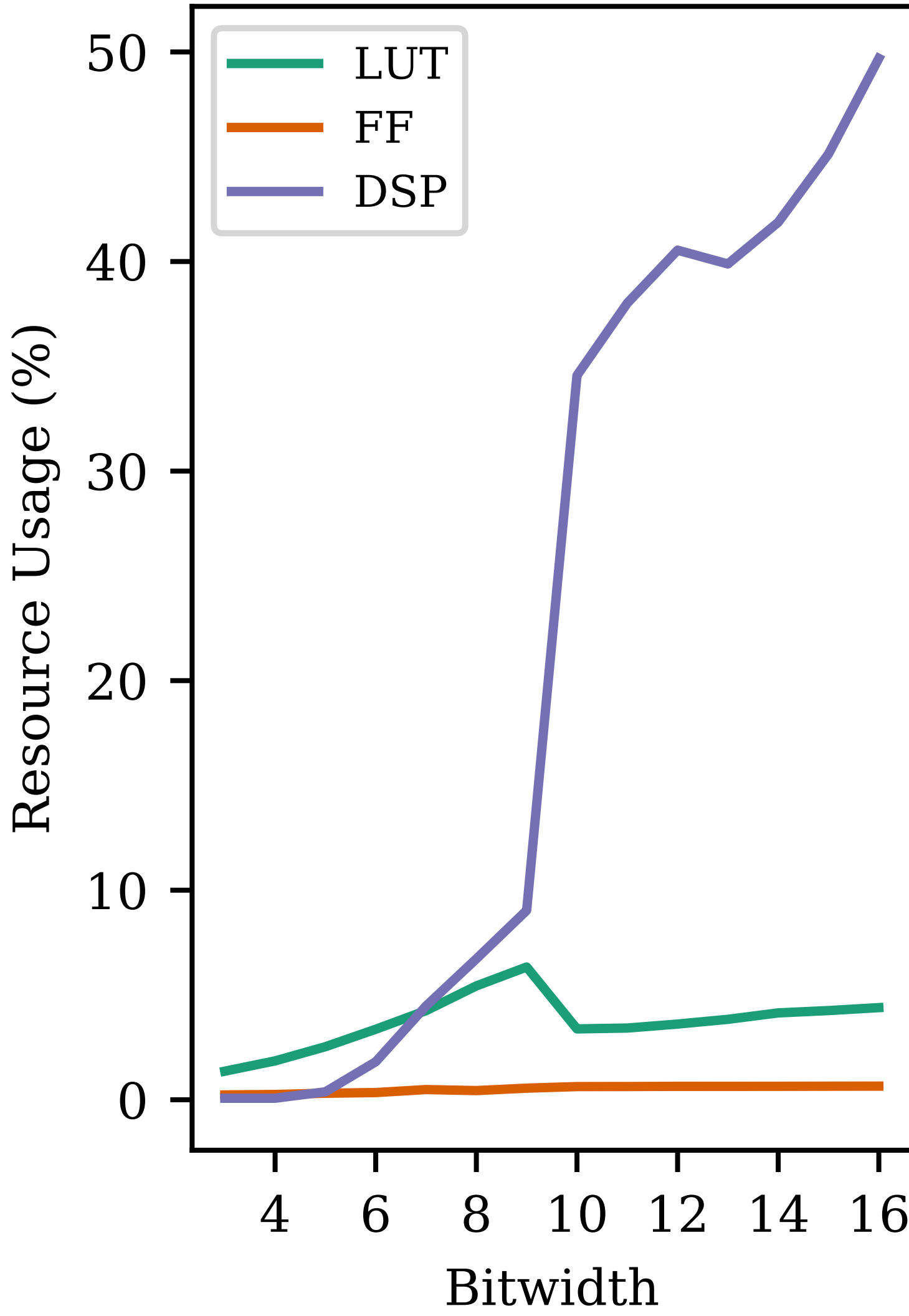
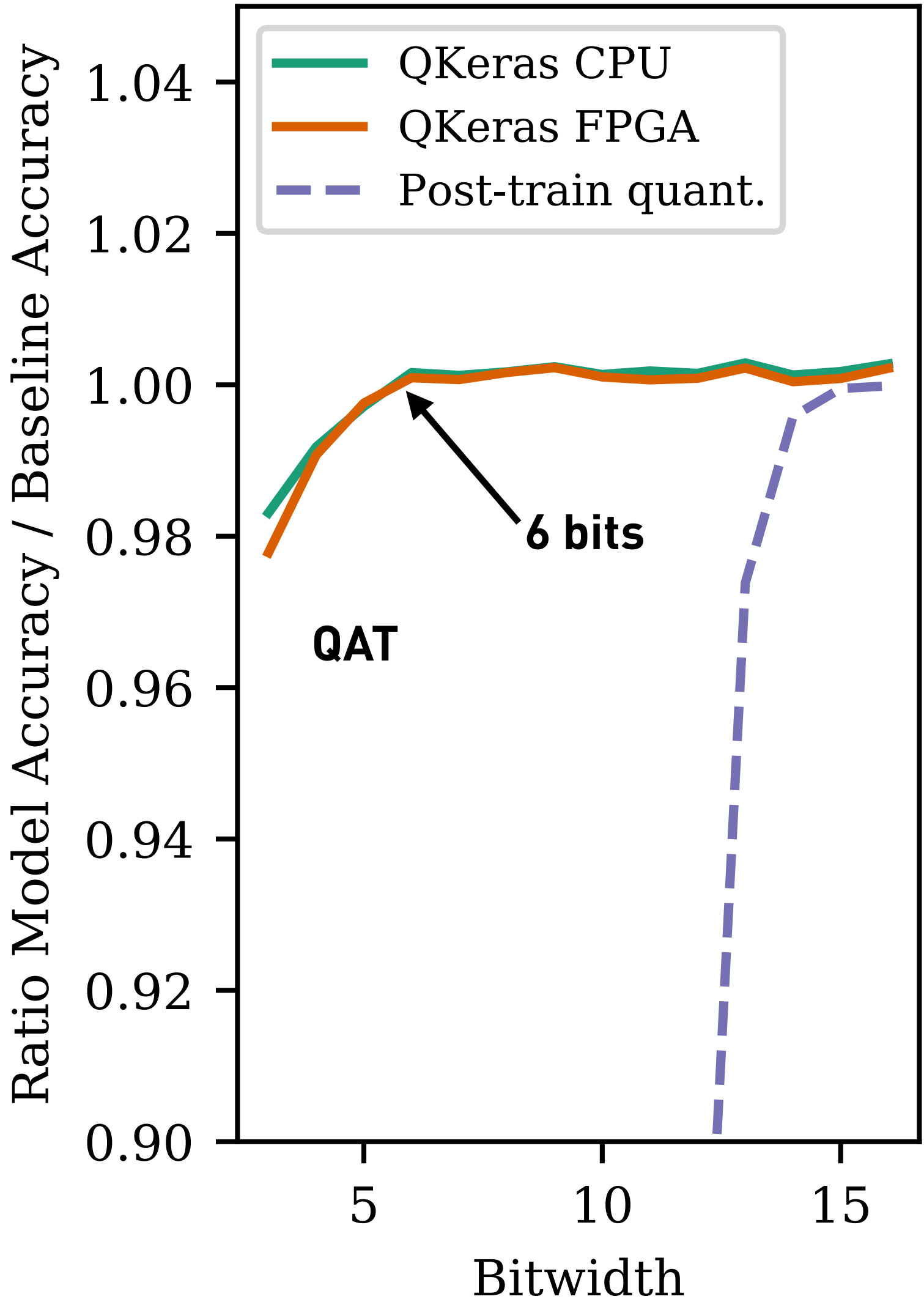
$$\text{Forgiving Factor} = 1 + \Delta_{\text{accuracy}} \times \log_{\text{rate}} \left(S \times \frac{\text{Cost}_{\text{ref}}}{\text{Cost}_{\text{trial}}} \right)$$

Maximize accuracy + minimizing cost in hyper parameter scan over quantizers:
AutoQKeras

Example: One convolutional layer

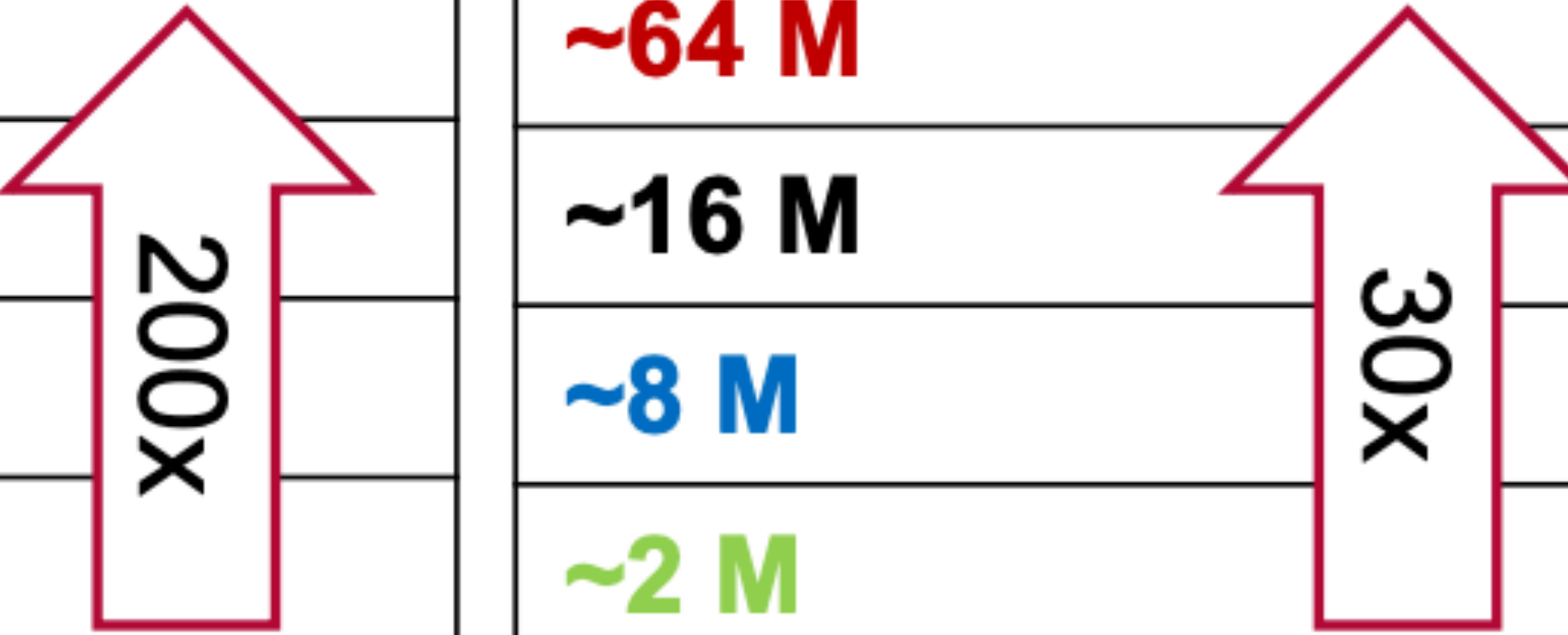
As optimization progresses,
best model accuracy/size
trade-off is found!





AMD UltraScale+ MPSoC ZU19EG (conservative estimates)

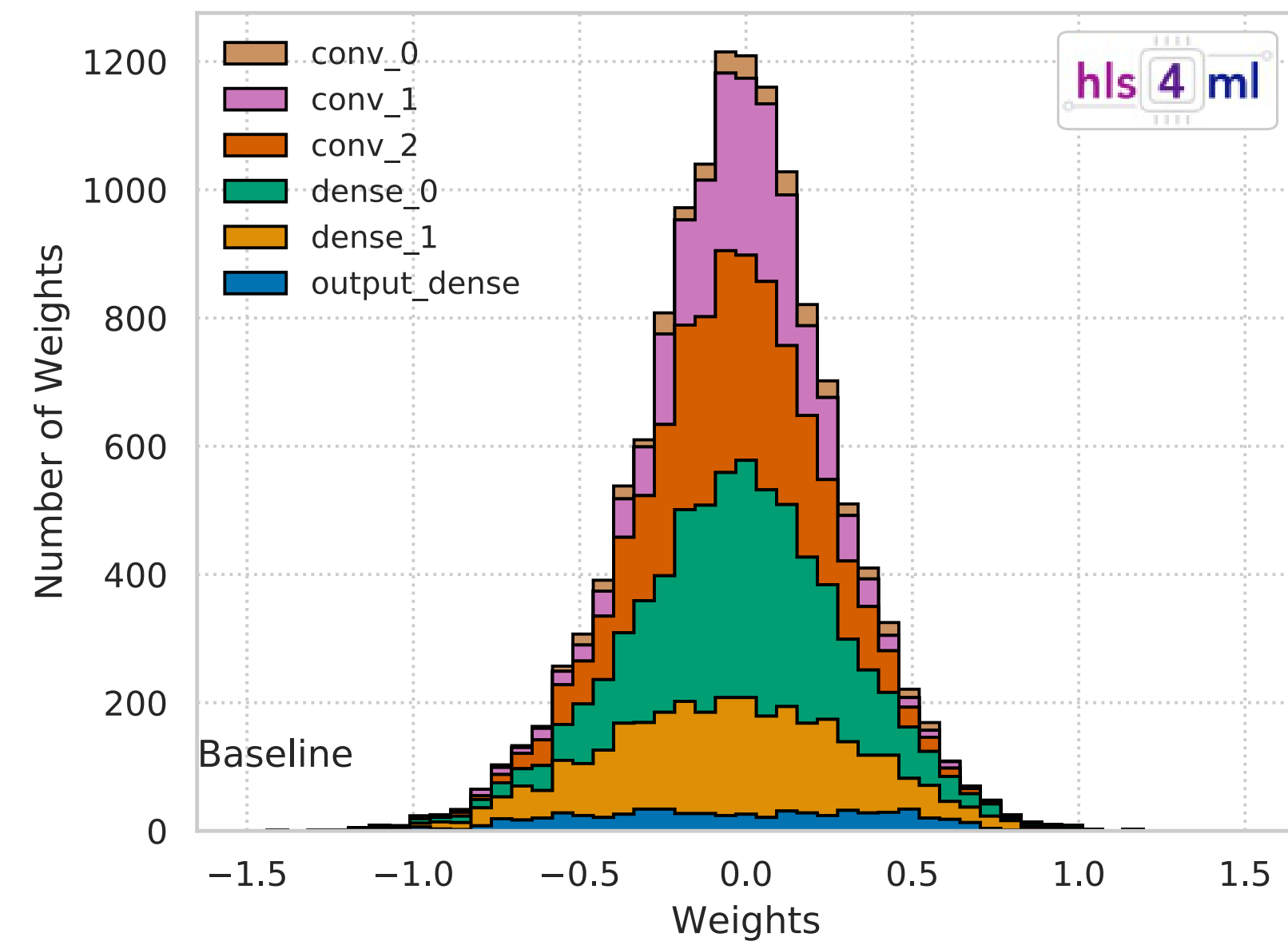
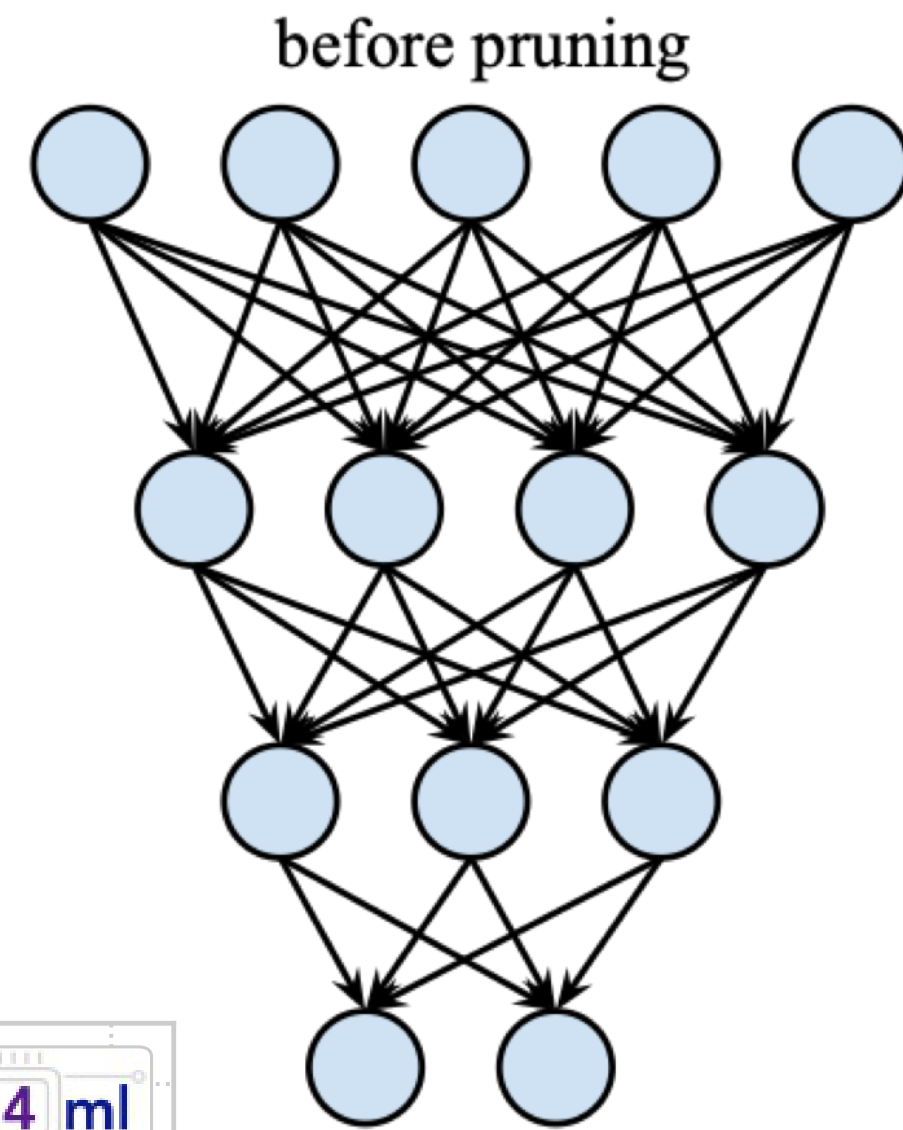
Precision	Approx. Peak GOPS	On-chip weights
1b	64 000	~64 M
4b	16 000	~16 M
8b	4 000	~8 M
32b	300	~2 M

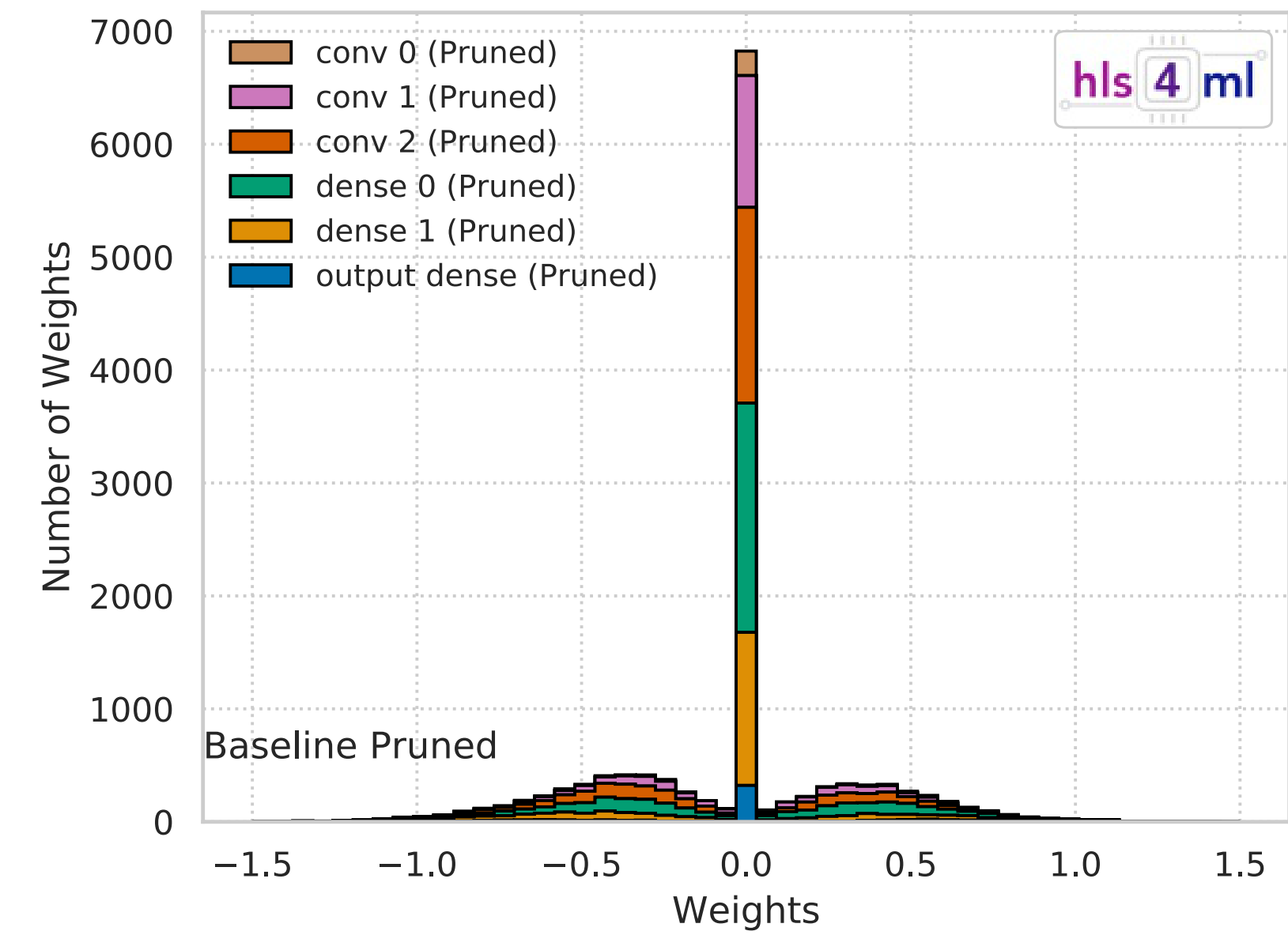
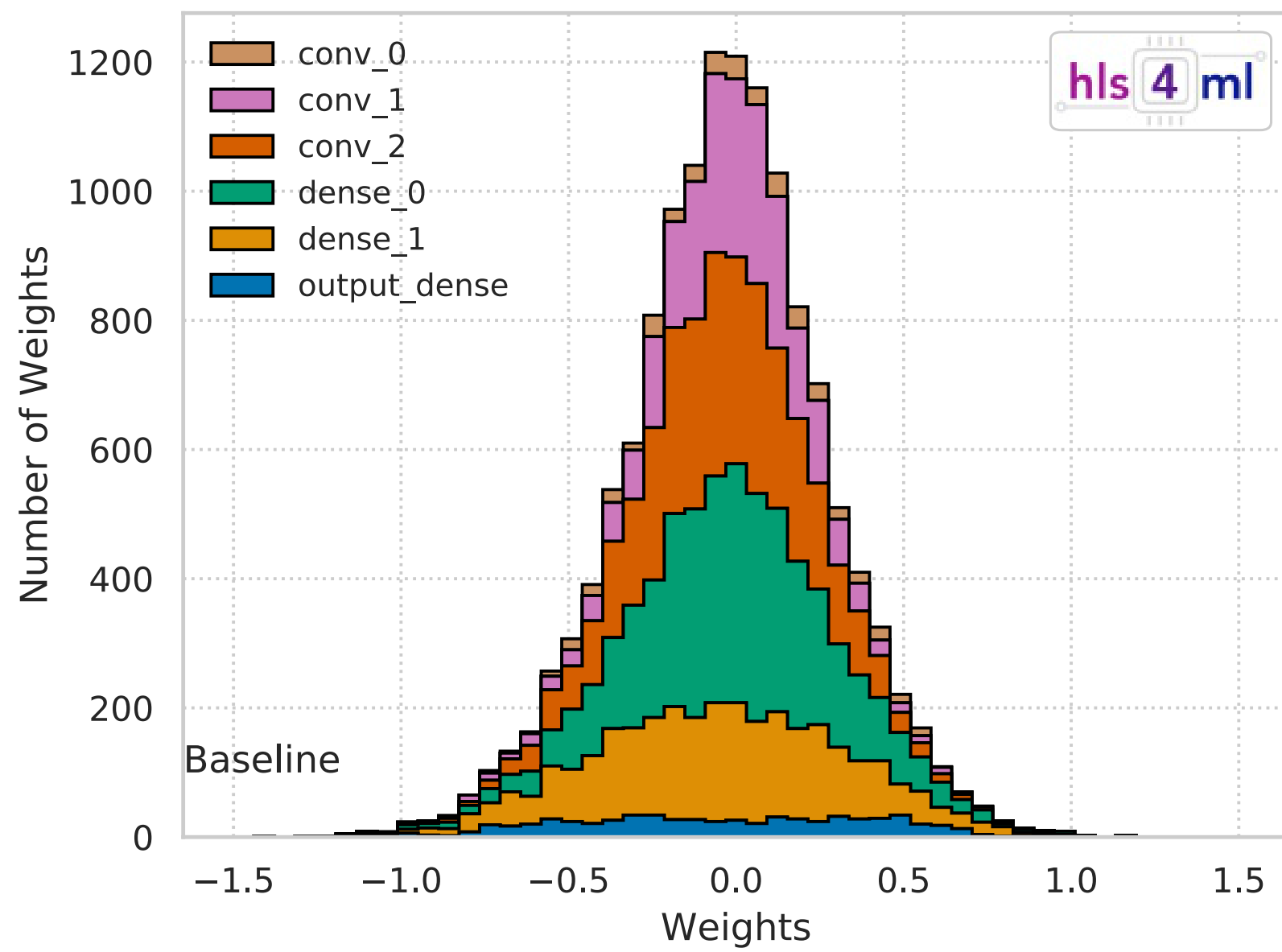
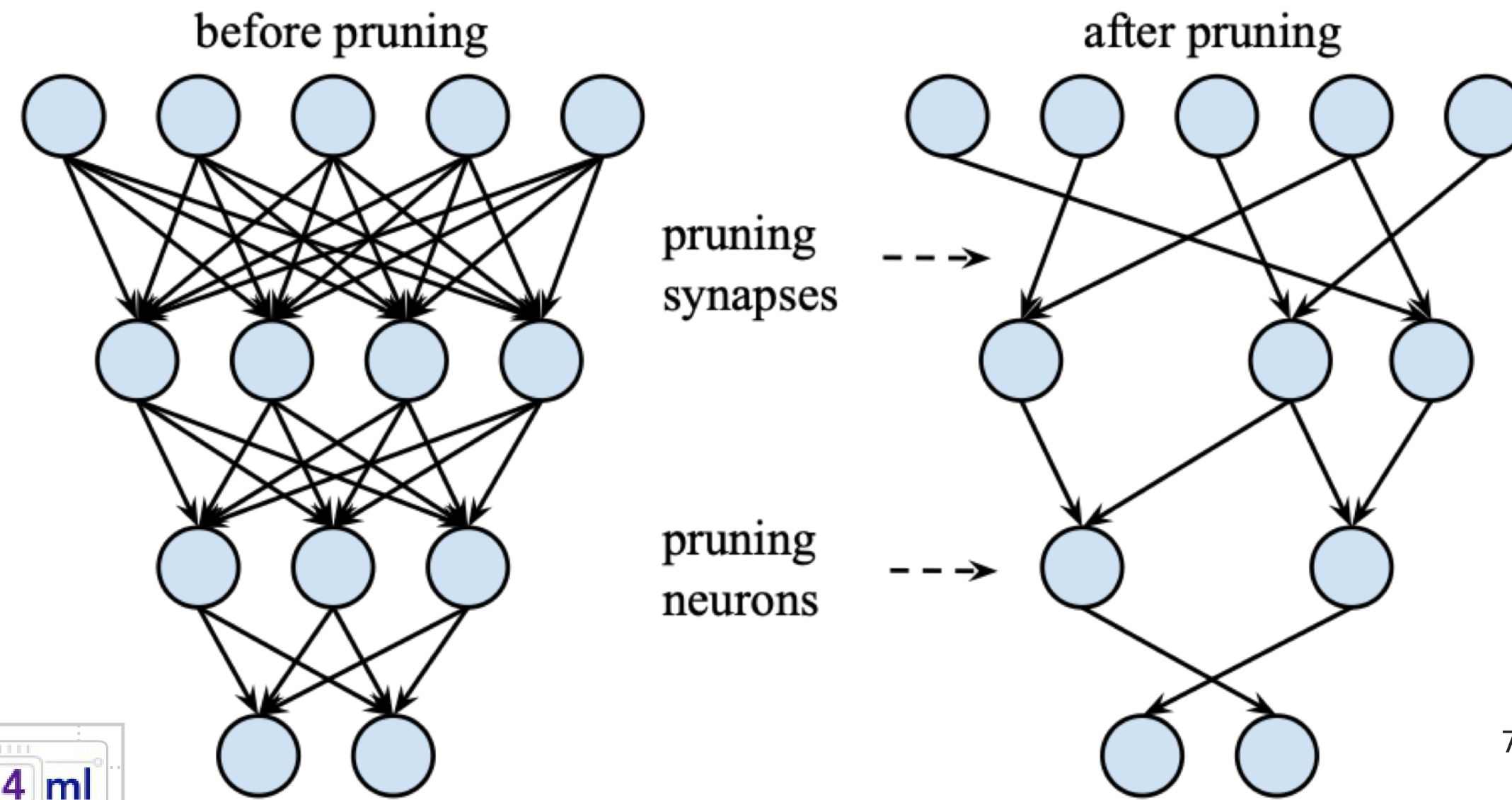


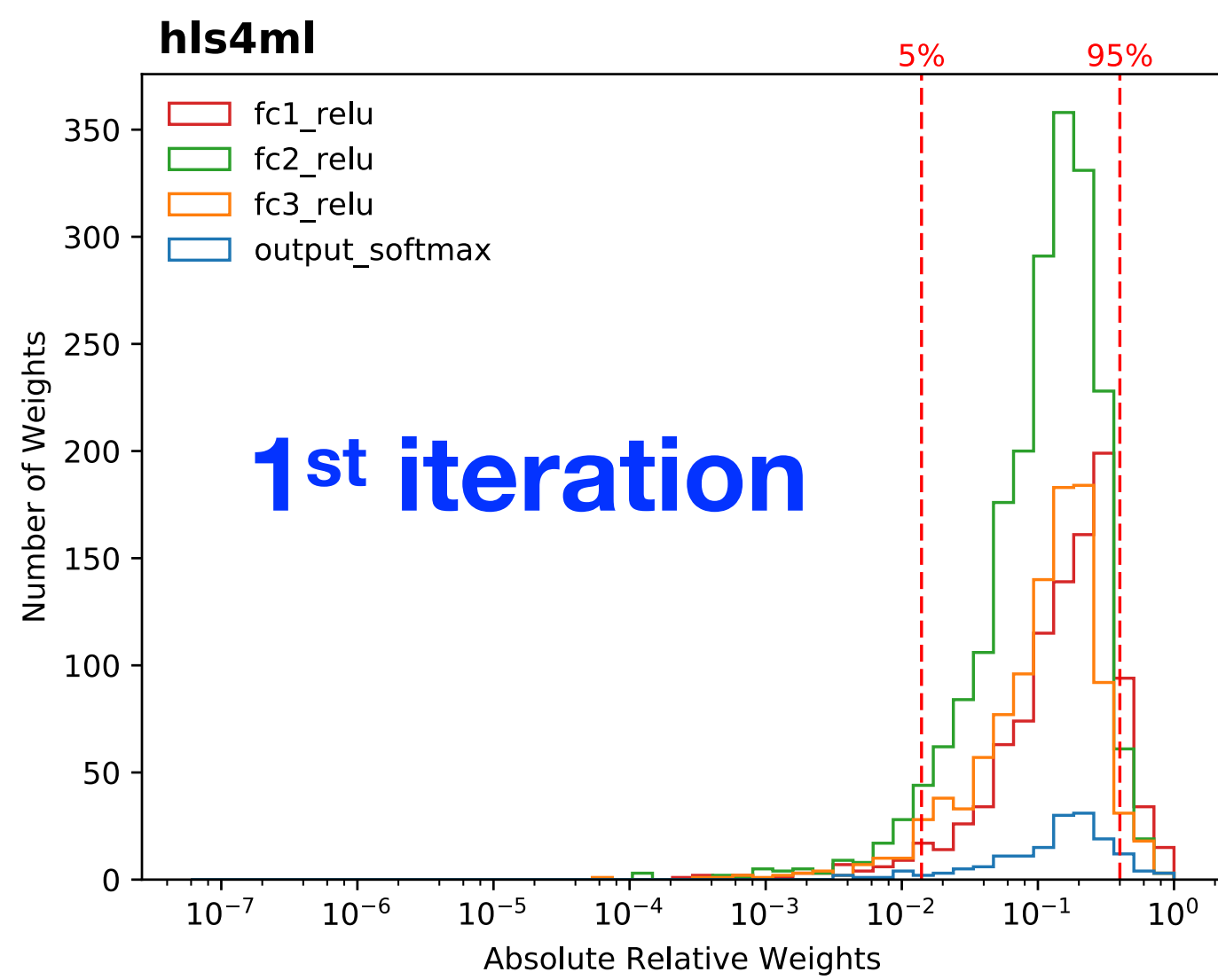
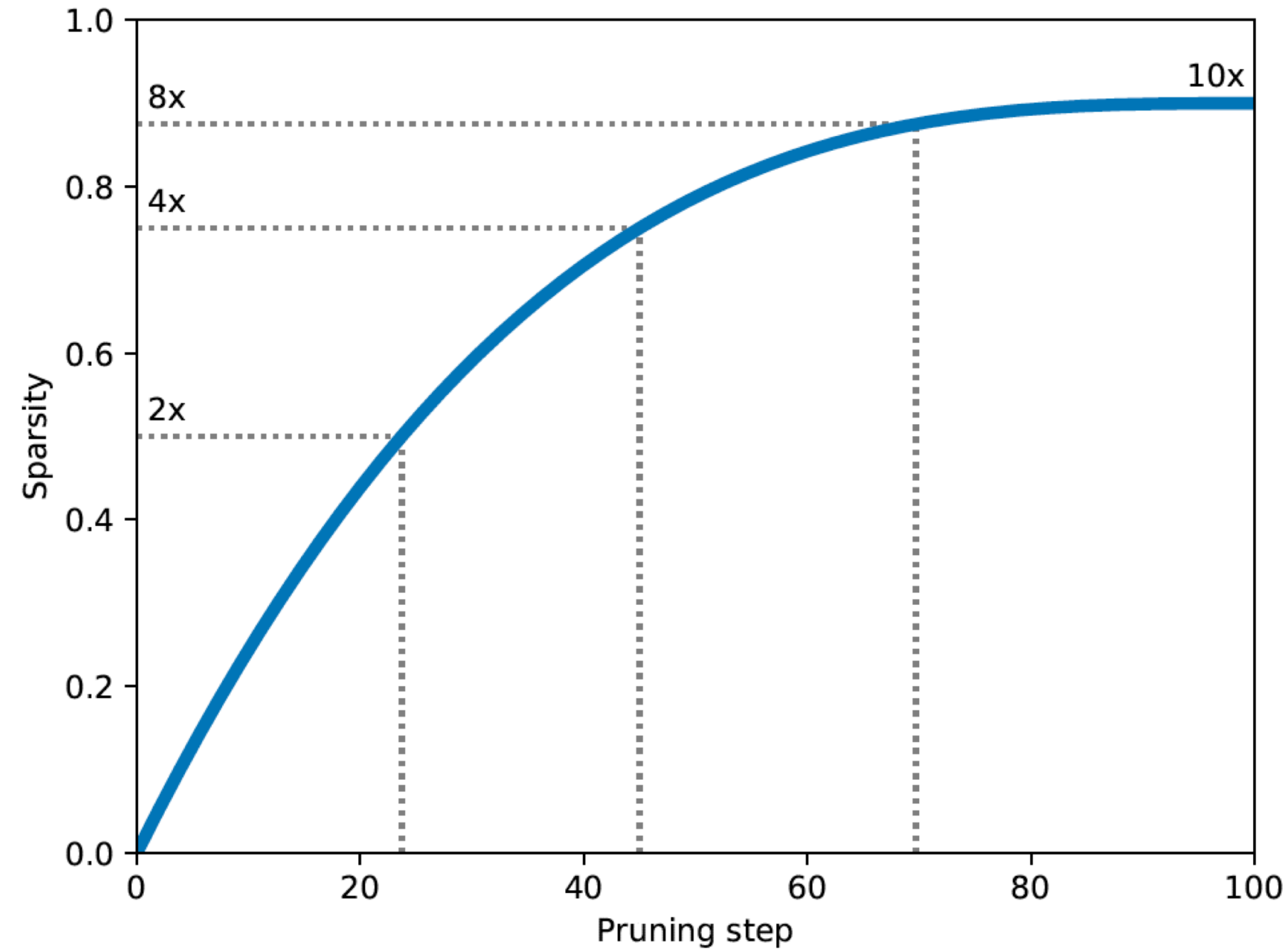
Trillions of
quantized
operations per
second

Weights can
stay **entirely**
on-chip

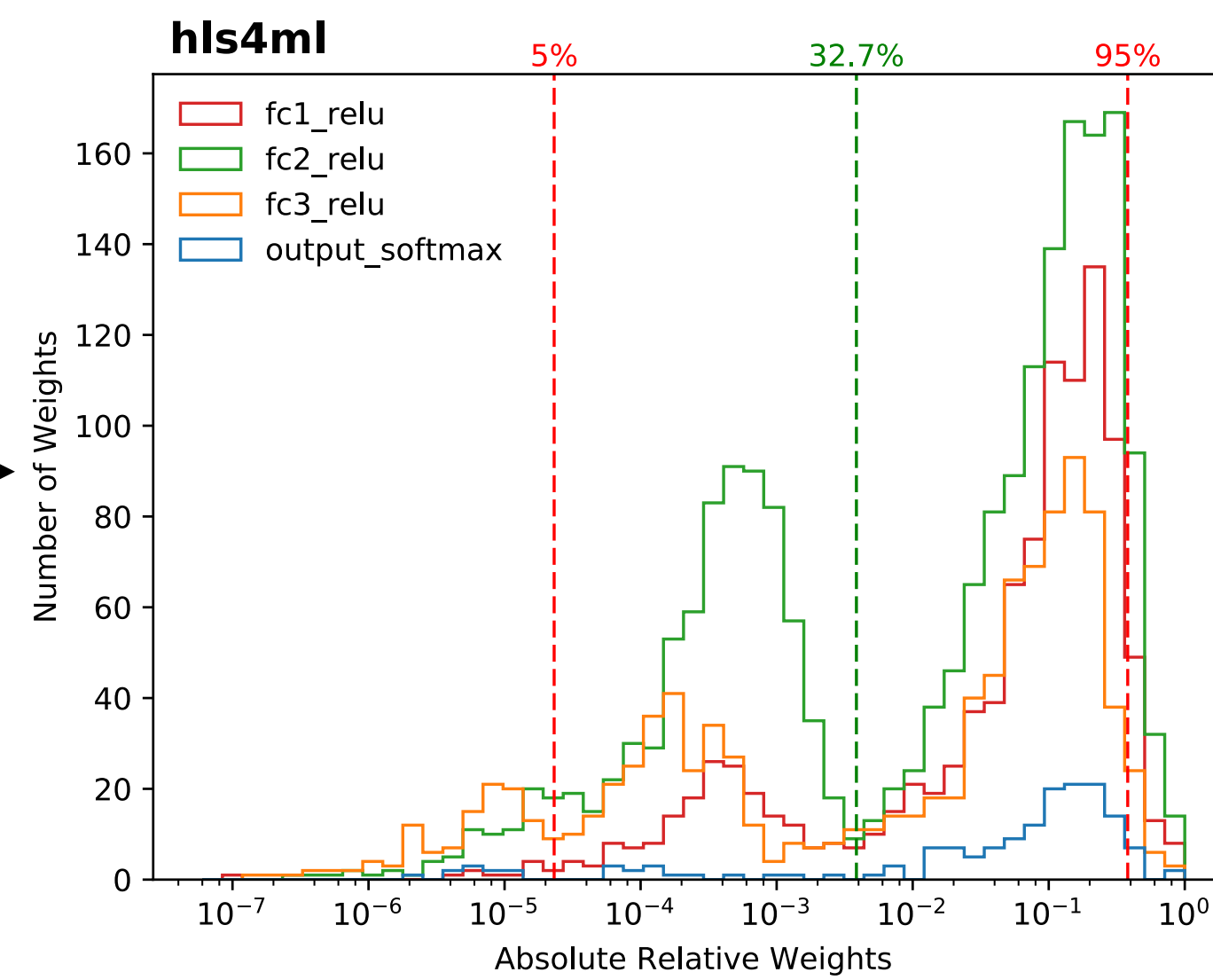
Pruning



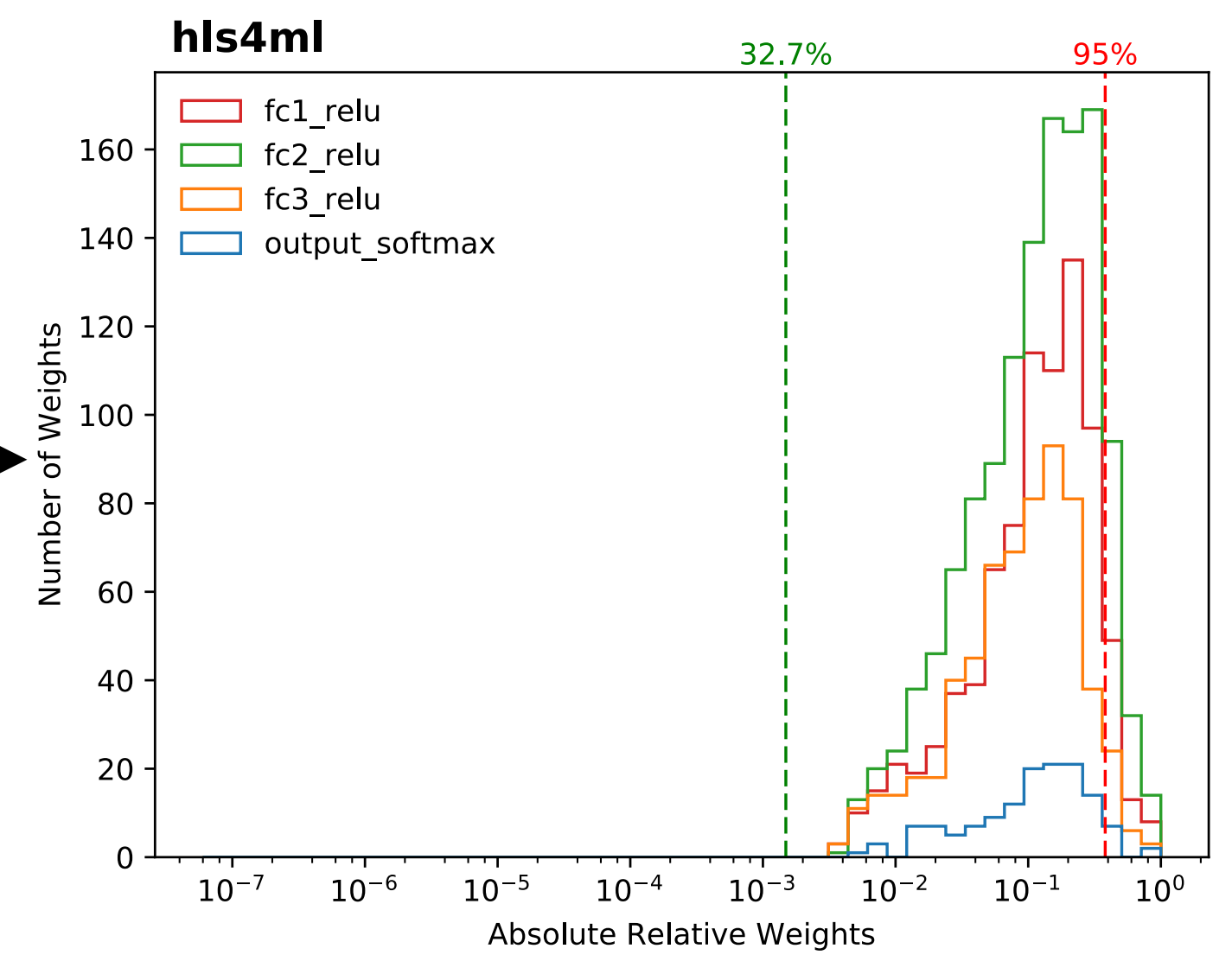




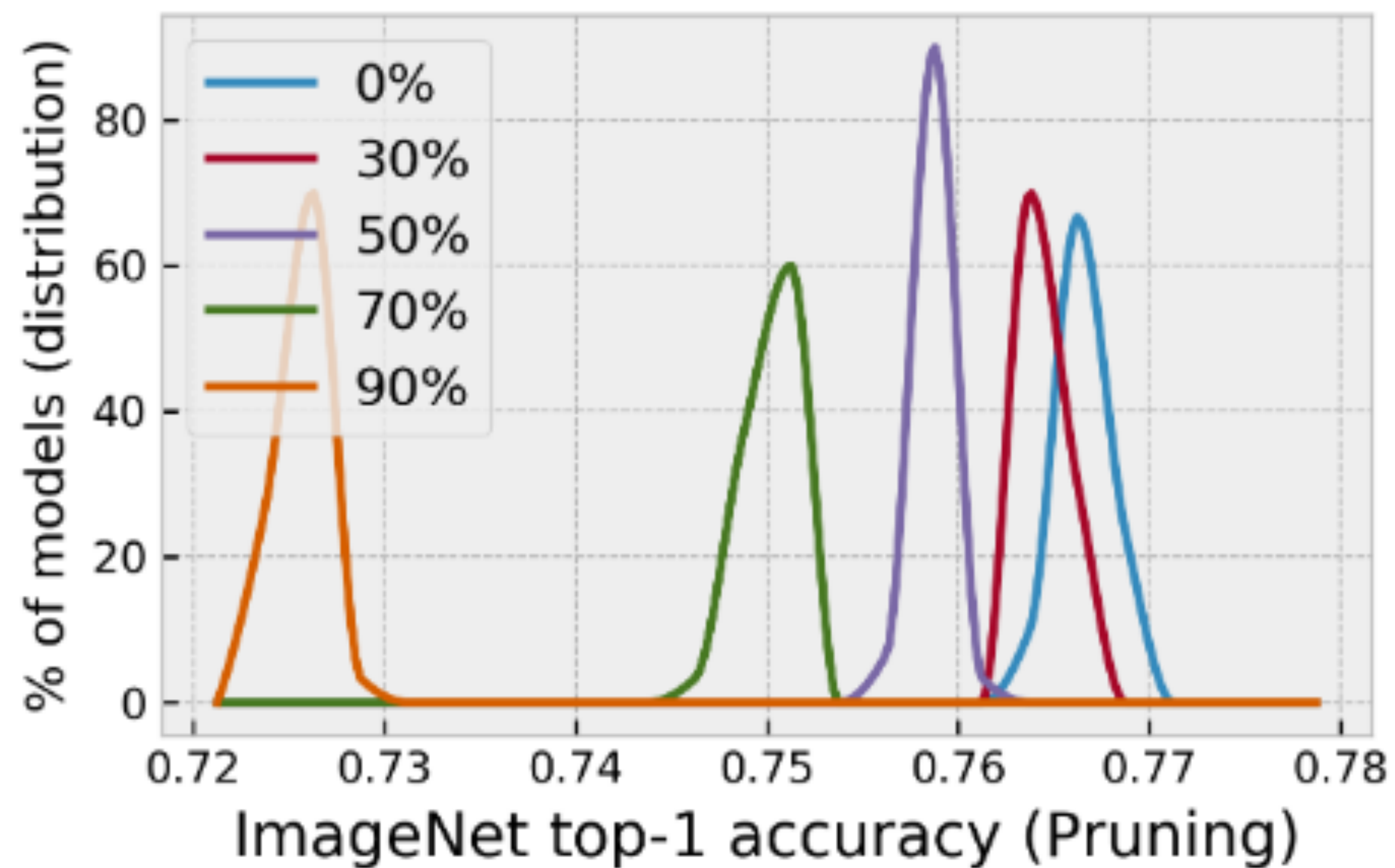
Train with L₁



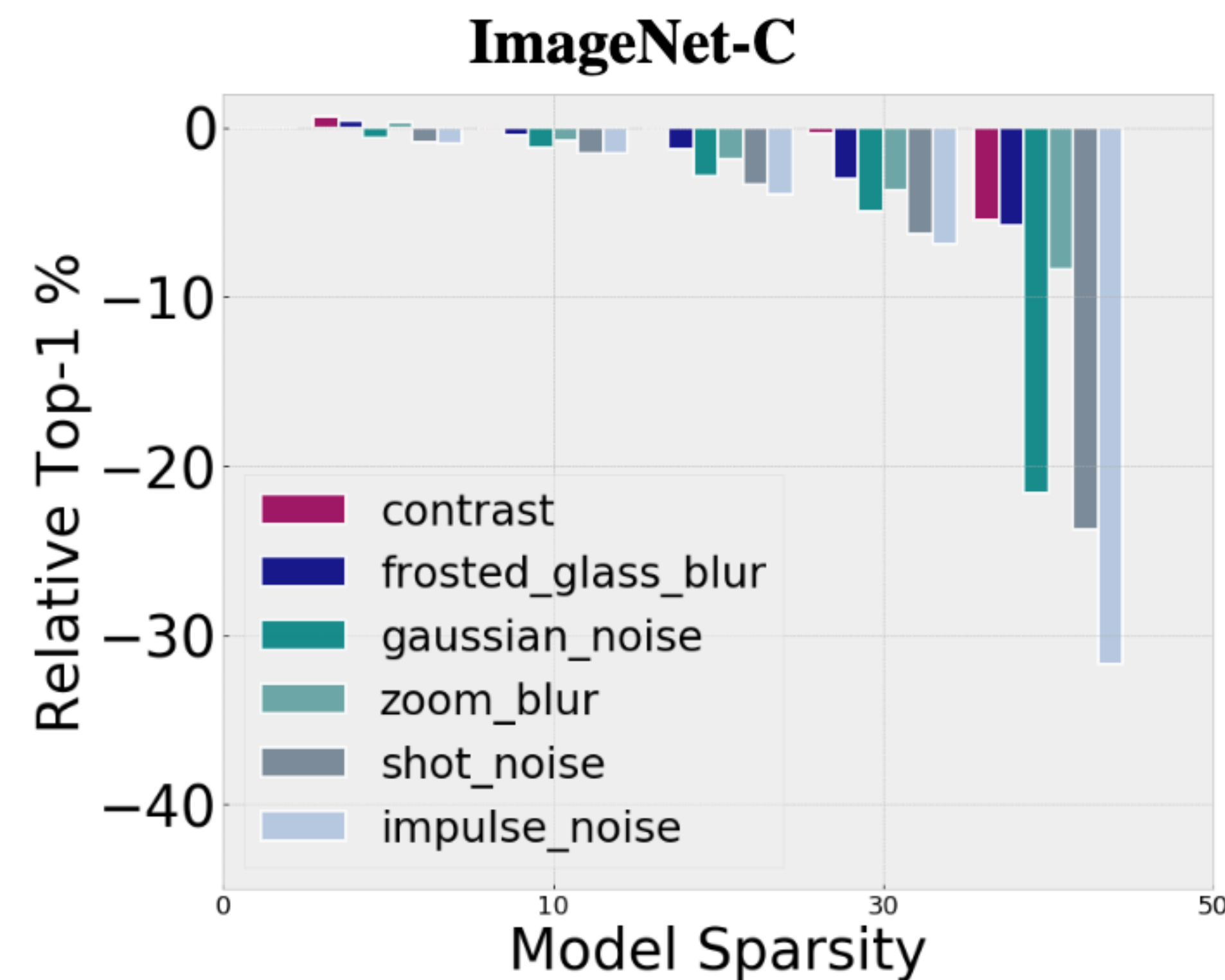
Prune



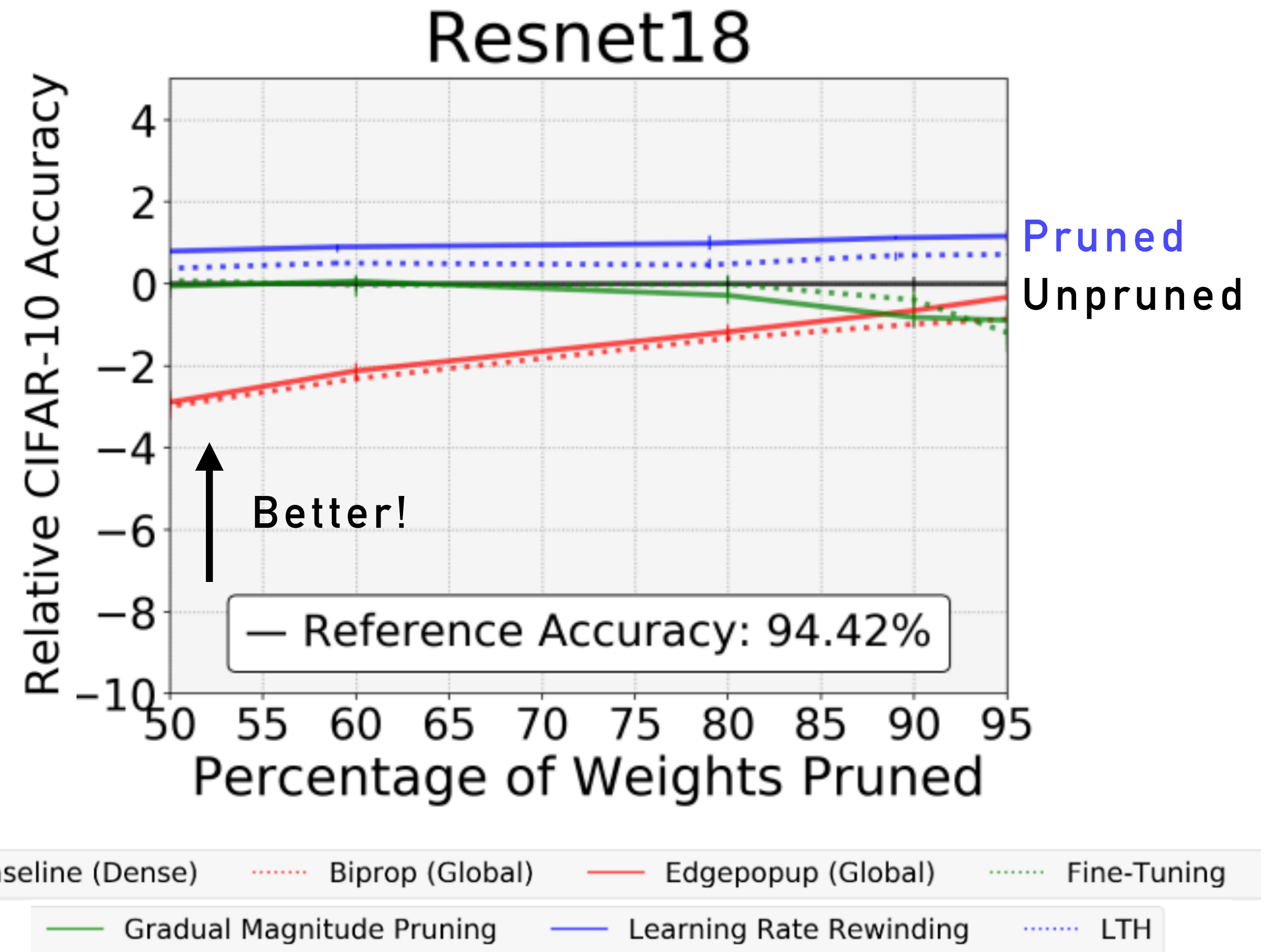
accuracy:



robustness:



Lottery ticket hypothesis




There exists a optimal network WITHIN each network (lottery ticket)
Uncover it through pruning!



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

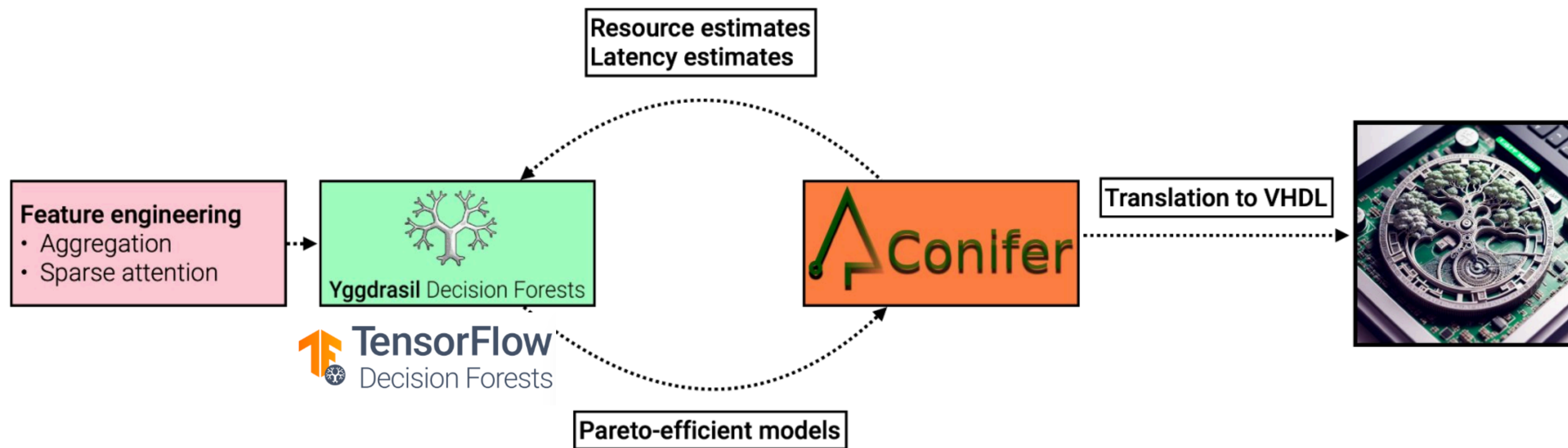
From: Çağlar Aytakin [[view email](#)]

[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%



Very high quality data can help smaller models compete

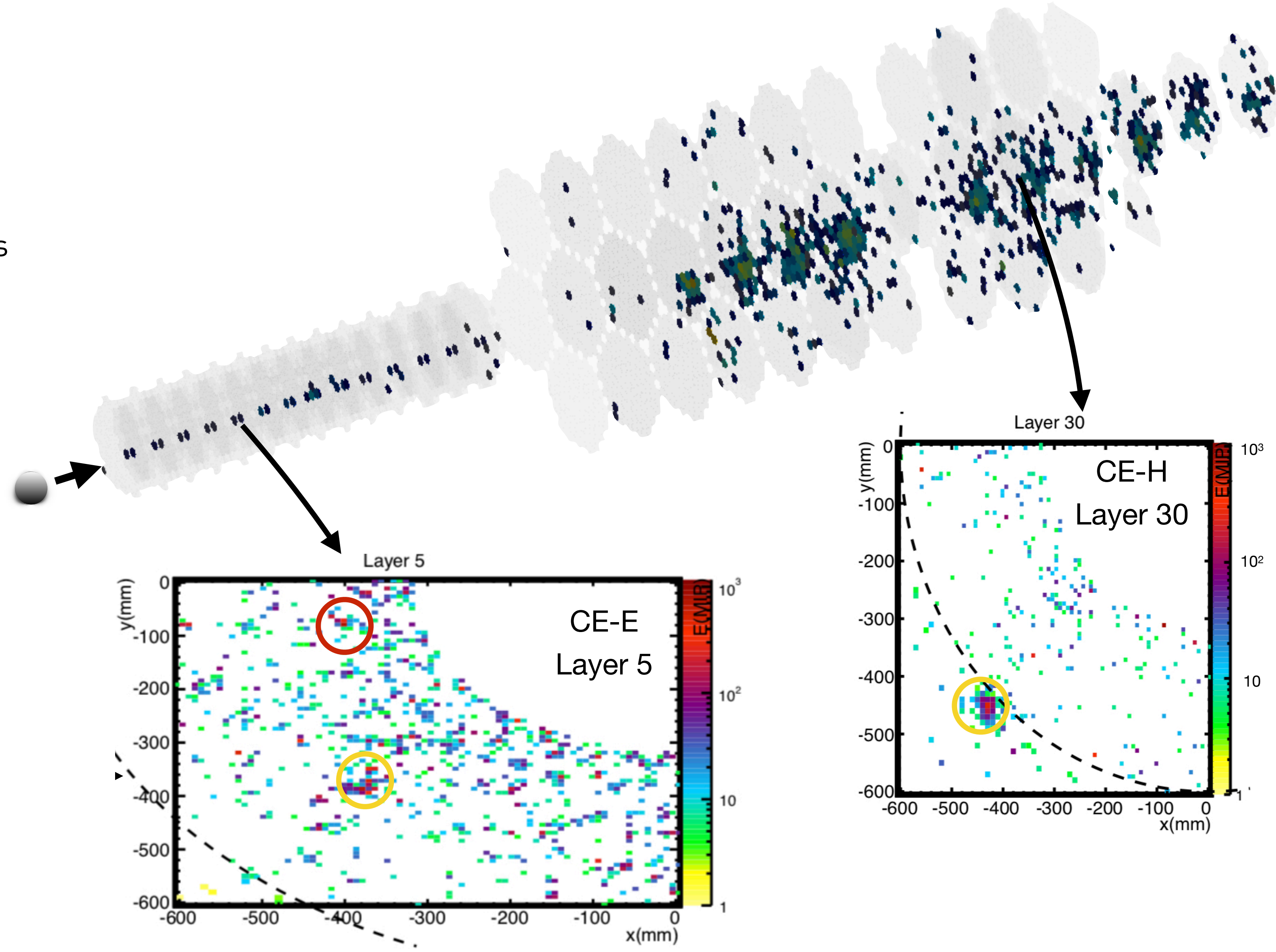
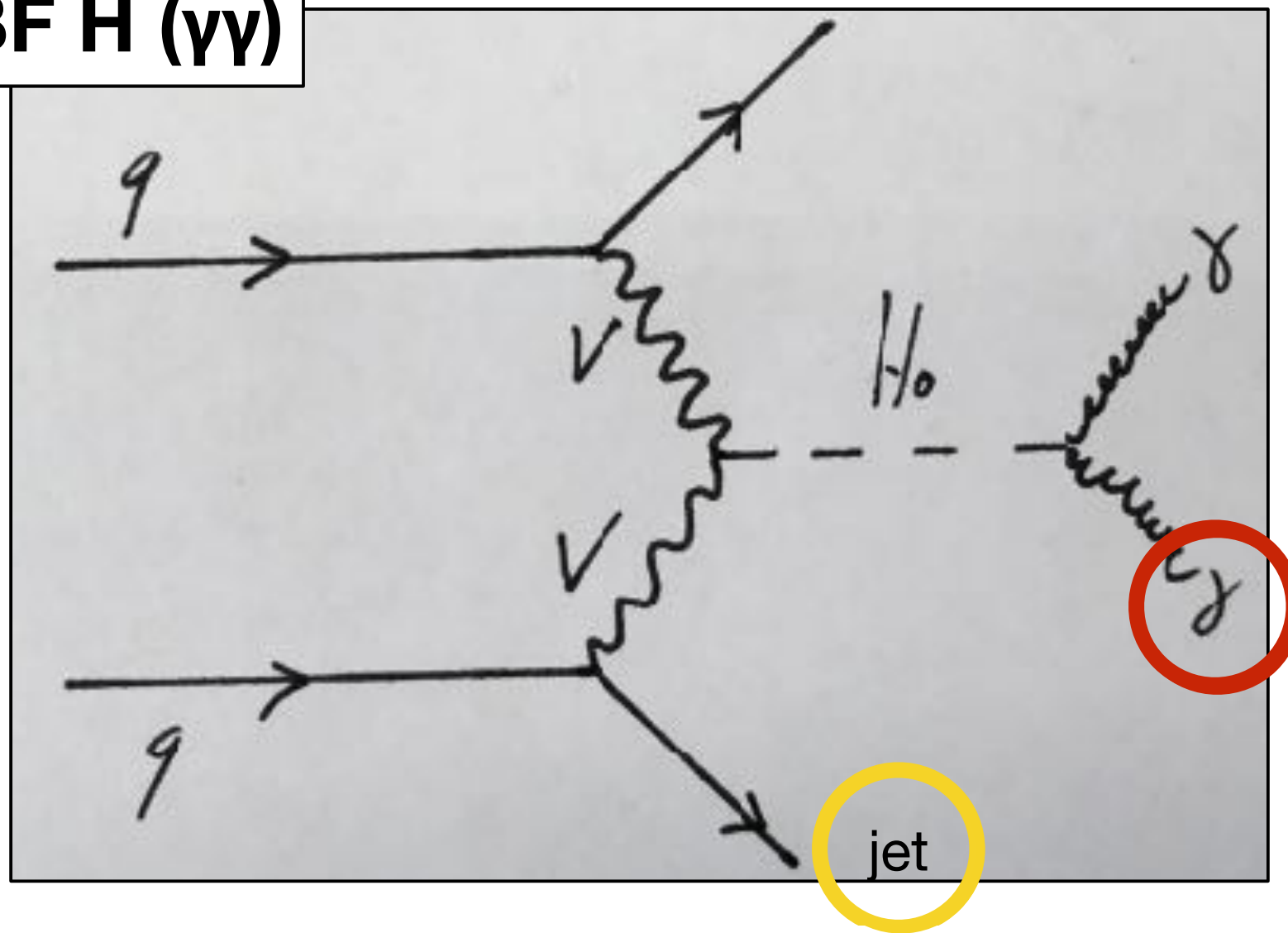
Date	Model	Model size	Dataset size (Tokens)	HumanEval (Pass@1)	MBPP (Pass@1)
2021 Jul	Codex-300M [CTJ ⁺ 21]	300M	100B	13.2%	-
2021 Jul	Codex-12B [CTJ ⁺ 21]	12B	100B	28.8%	-
2022 Mar	CodeGen-Mono-350M [NPH ⁺ 23]	350M	577B	12.8%	-
2022 Mar	CodeGen-Mono-16.1B [NPH ⁺ 23]	16.1B	577B	29.3%	35.3%
2022 Apr	PaLM-Coder [CND ⁺ 22]	540B	780B	35.9%	47.0%
2022 Sep	CodeGeeX [ZXZ ⁺ 23]	13B	850B	22.9%	24.4%
2022 Nov	GPT-3.5 [Ope23]	175B	N.A.	47%	-
2022 Dec	SantaCoder [ALK ⁺ 23]	1.1B	236B	14.0%	35.0%
2023 Mar	GPT-4 [Ope23]	N.A.	N.A.	67%	-
2023 Apr	Replit [Rep23]	2.7B	525B	21.9%	-
2023 Apr	Replit-Finetuned [Rep23]	2.7B	525B	30.5%	-
2023 May	CodeGen2-1B [NHX ⁺ 23]	1B	N.A.	10.3%	-
2023 May	CodeGen2-7B [NHX ⁺ 23]	7B	N.A.	19.1%	-
2023 May	StarCoder [LAZ ⁺ 23]	15.5B	1T	33.6%	52.7%
2023 May	StarCoder-Prompted [LAZ ⁺ 23]	15.5B	1T	40.8%	49.5%
2023 May	PaLM 2-S [ADF ⁺ 23]	N.A.	N.A.	37.6%	50.0%
2023 May	CodeT5+ [WLG ⁺ 23]	2B	52B	24.2%	-
2023 May	CodeT5+ [WLG ⁺ 23]	16B	52B	30.9%	-
2023 May	InstructCodeT5+ [WLG ⁺ 23]	16B	52B	35.0%	-
2023 Jun	WizardCoder [LXZ ⁺ 23]	16B	1T	57.3%	51.8%
2023 Jun	phi-1	1.3B	7B	50.6%	55.5%

Competitive despite being 10x smaller in model size and 100x smaller in dataset size!

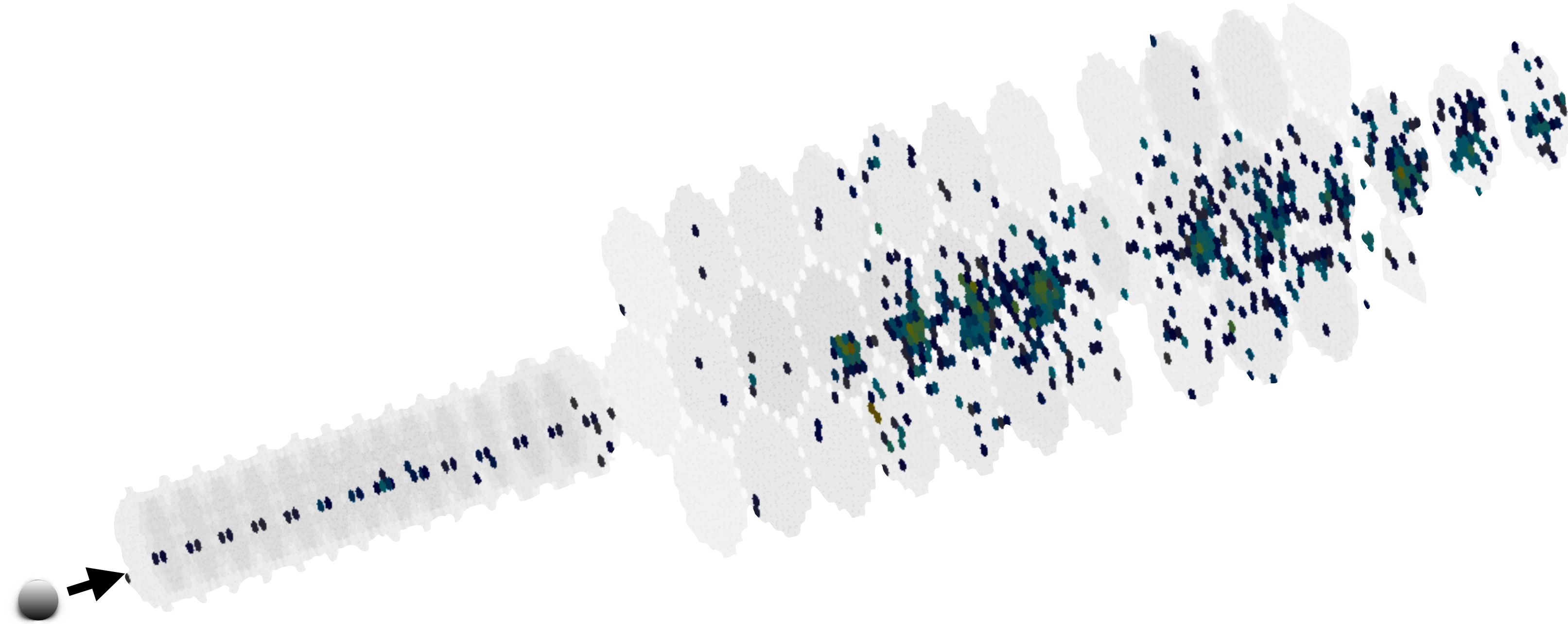
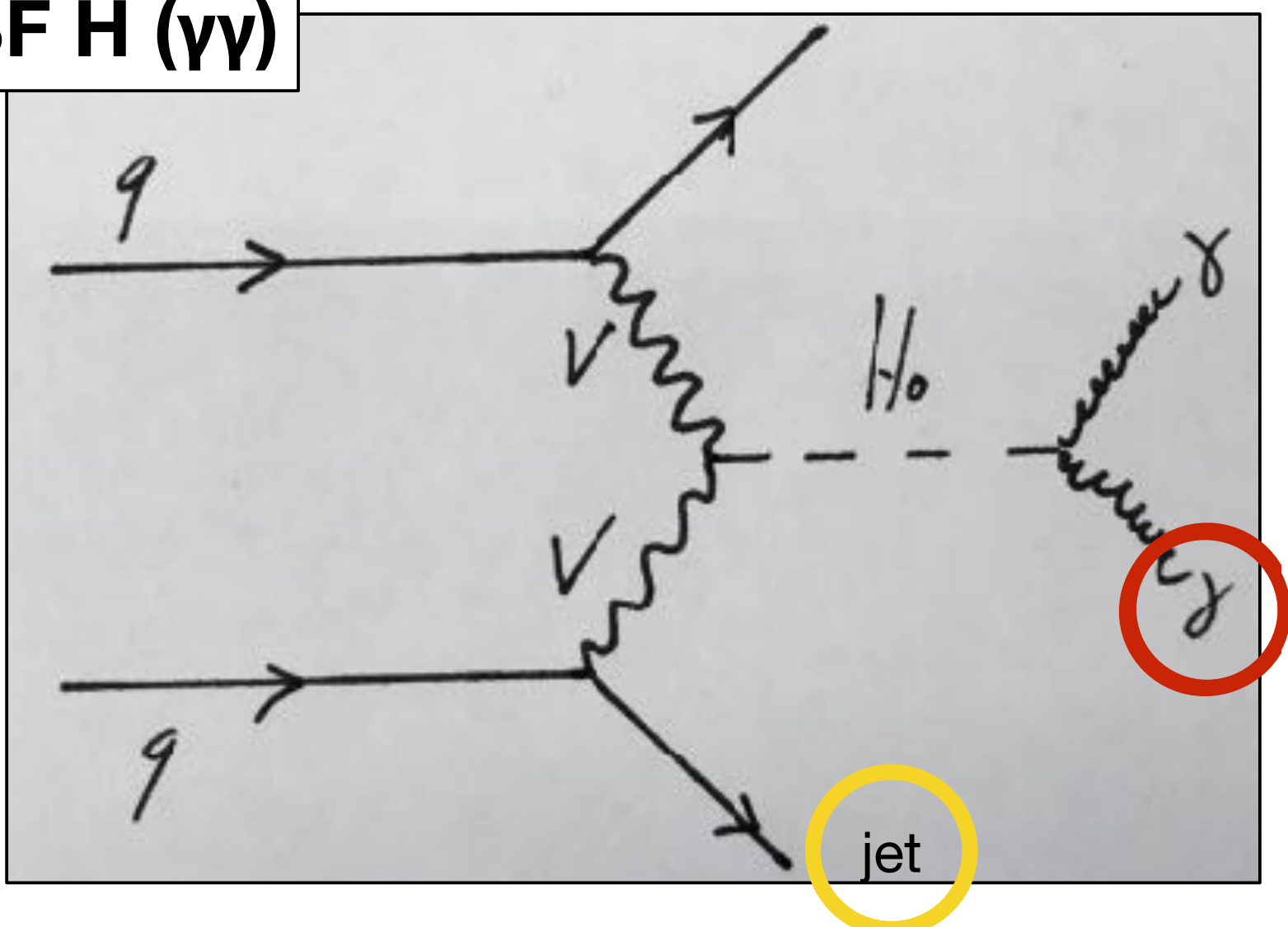
CMS High Granularity calorimeter

- 6.5 million readout channels, 50 layers

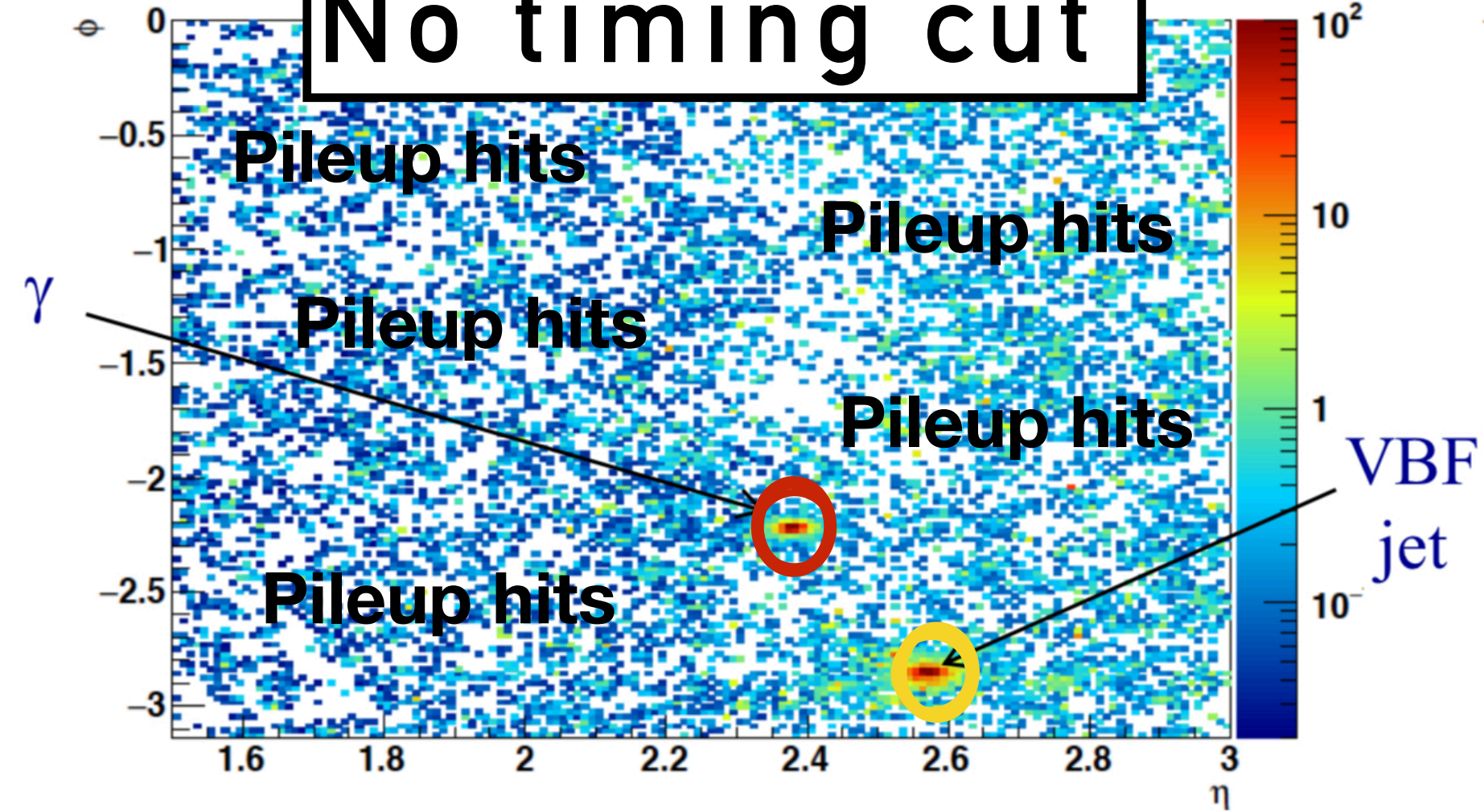
VBF H ($\gamma\gamma$)



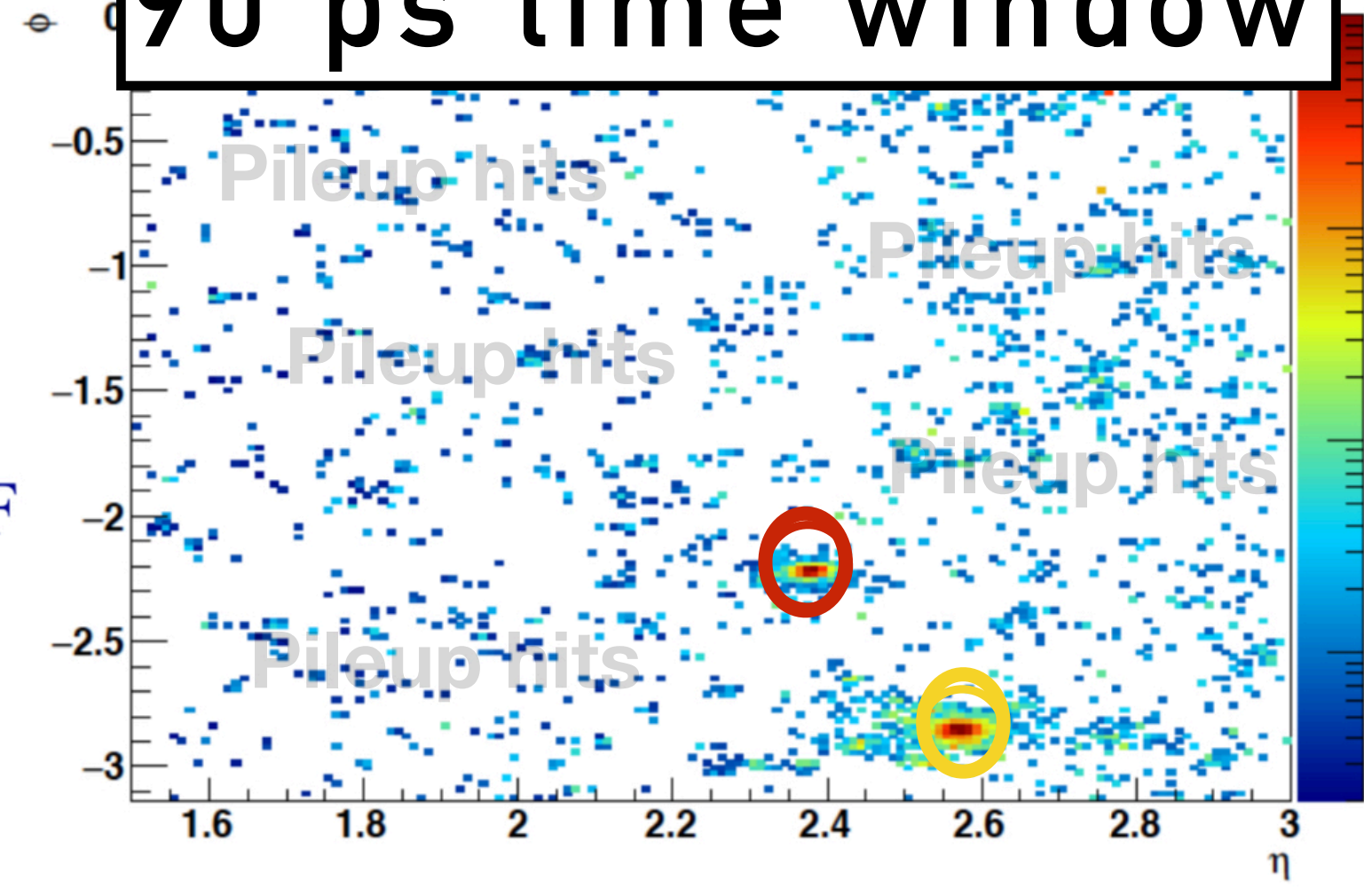
VBF H ($\gamma\gamma$)



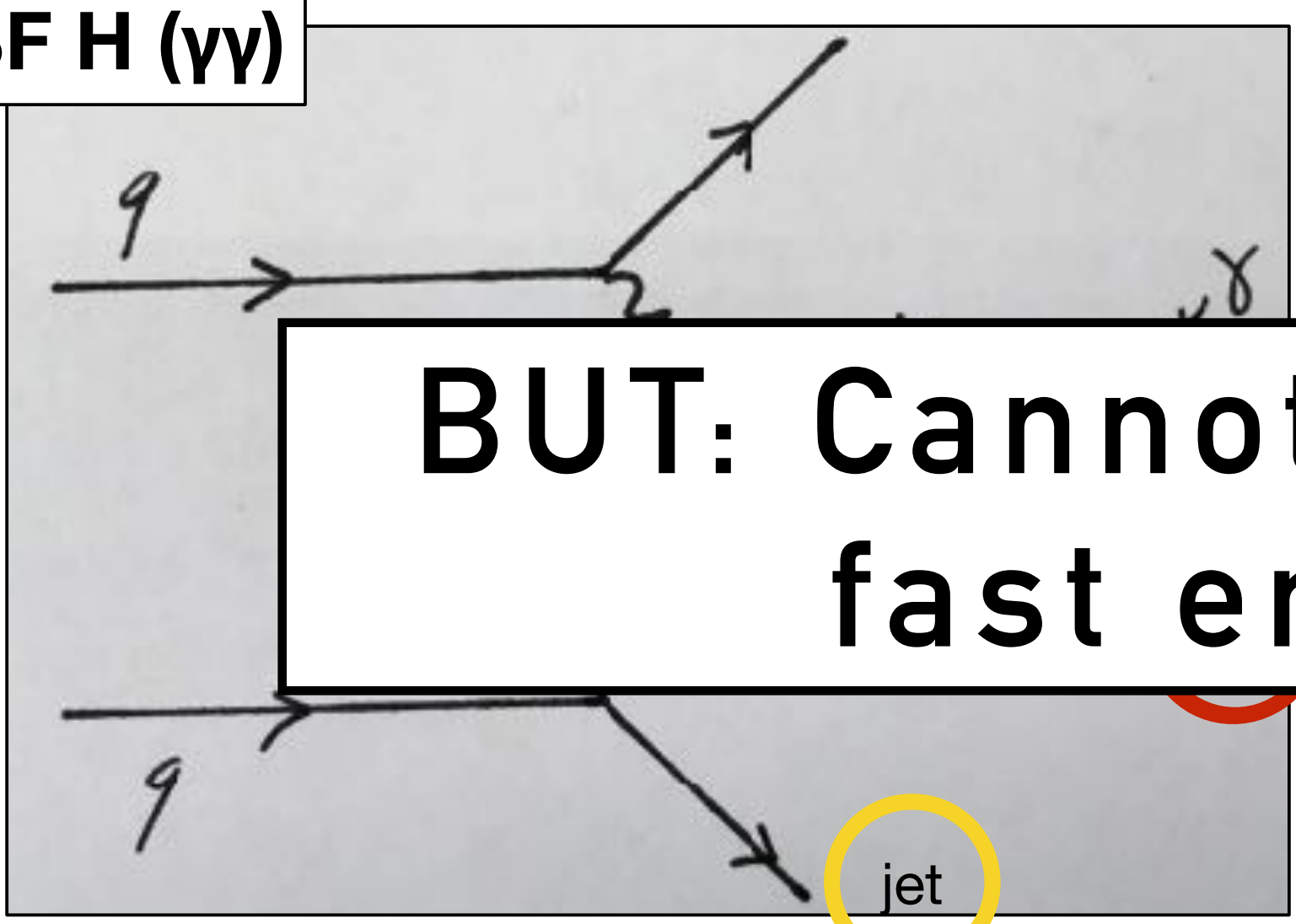
No timing cut



90 ps time window

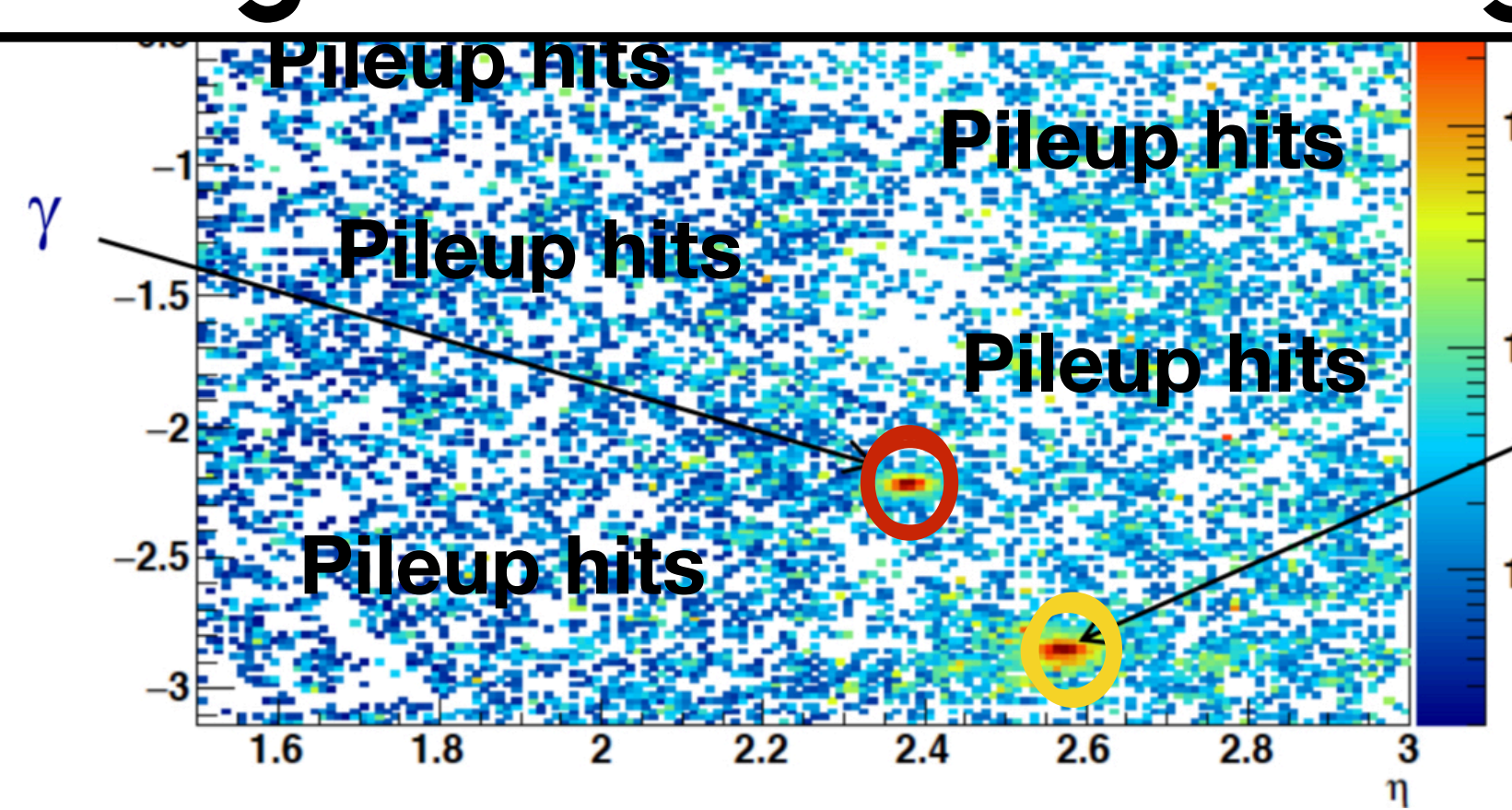
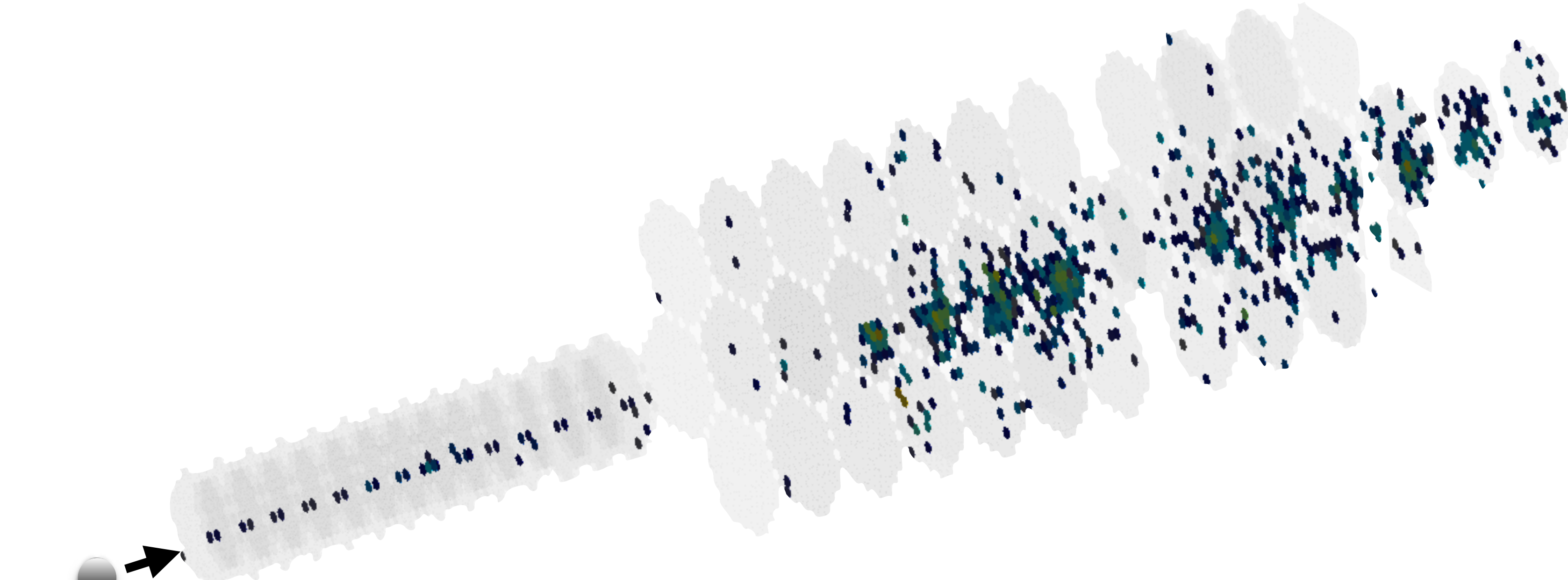


VBF H ($\gamma\gamma$)

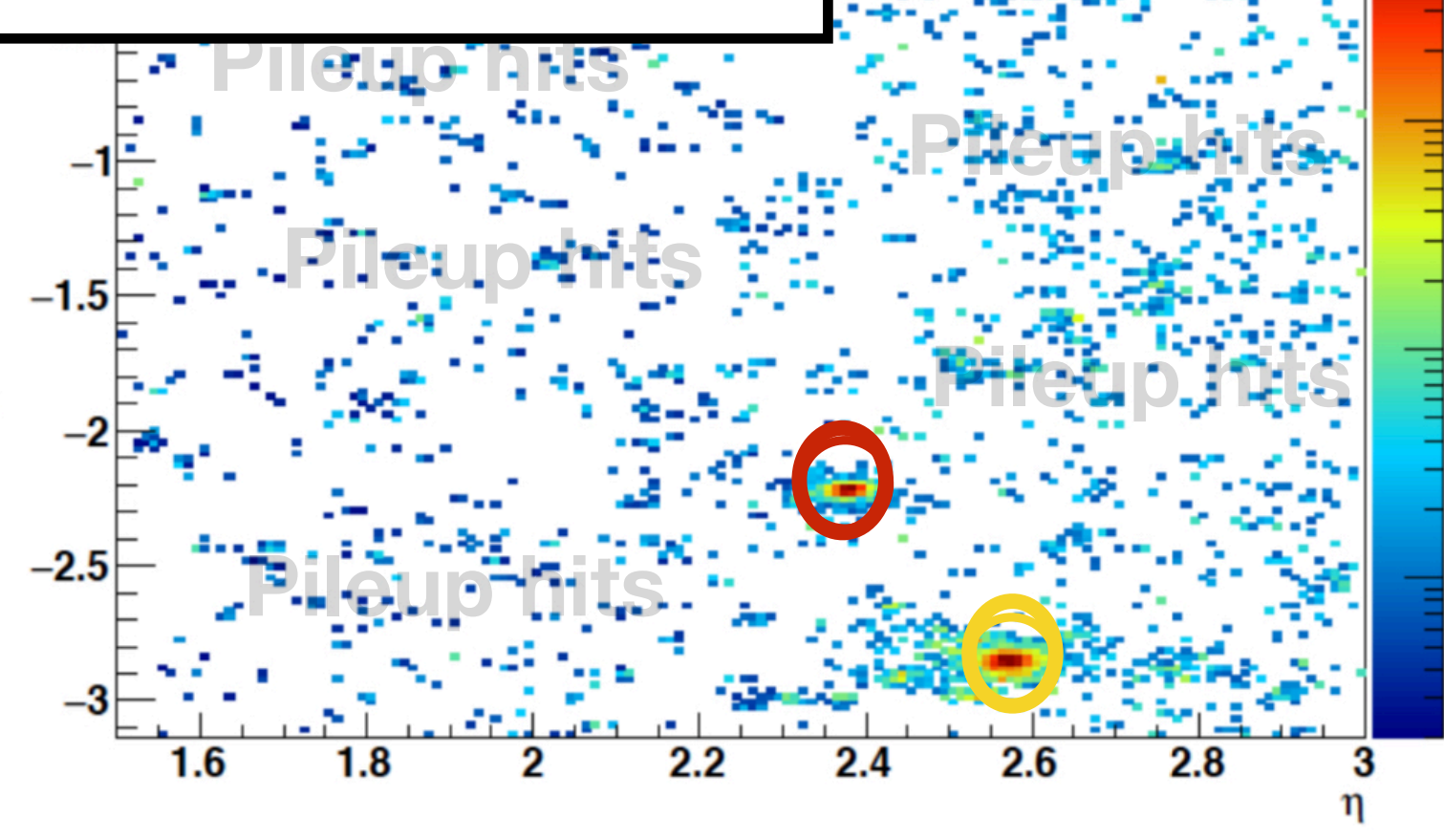


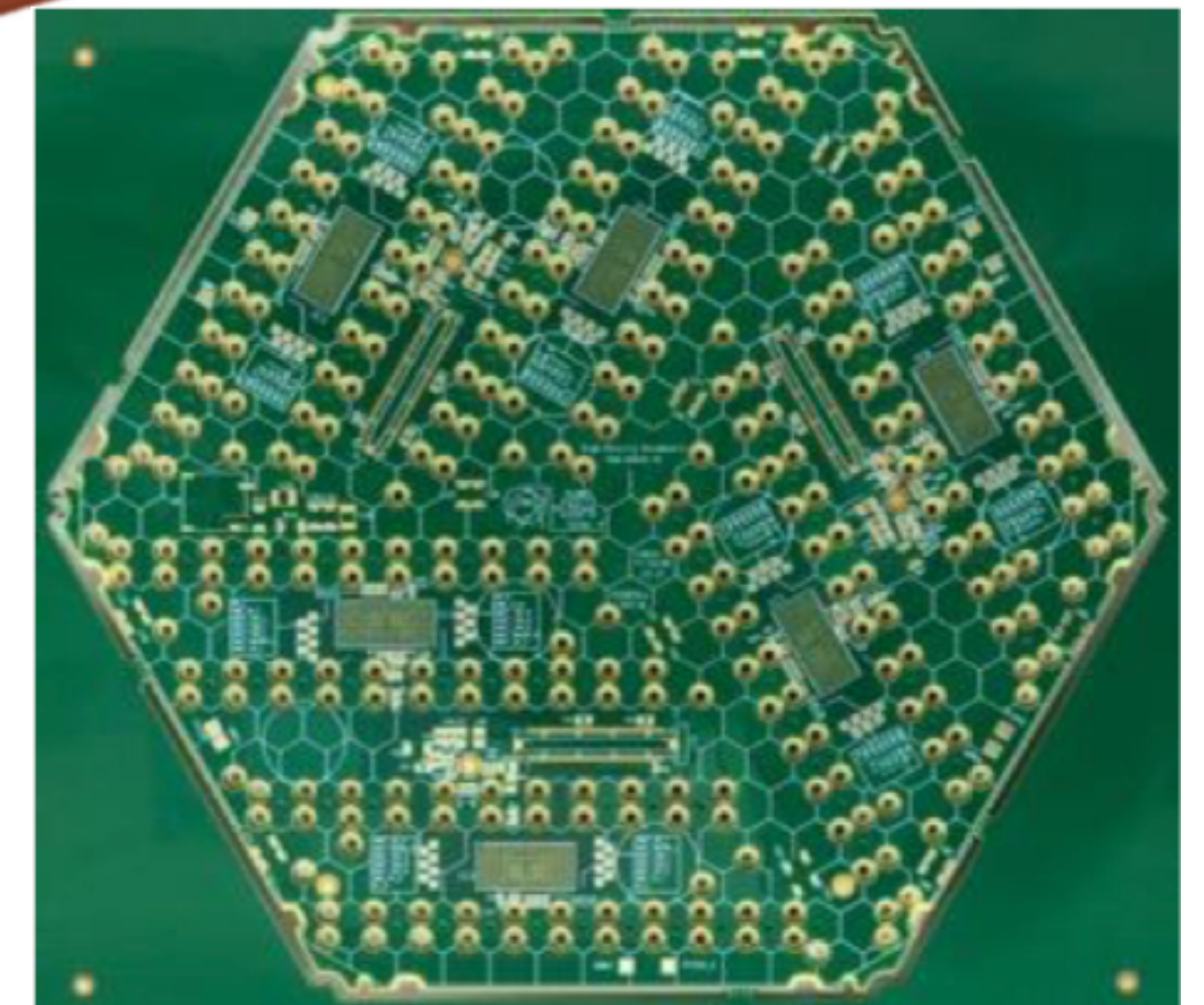
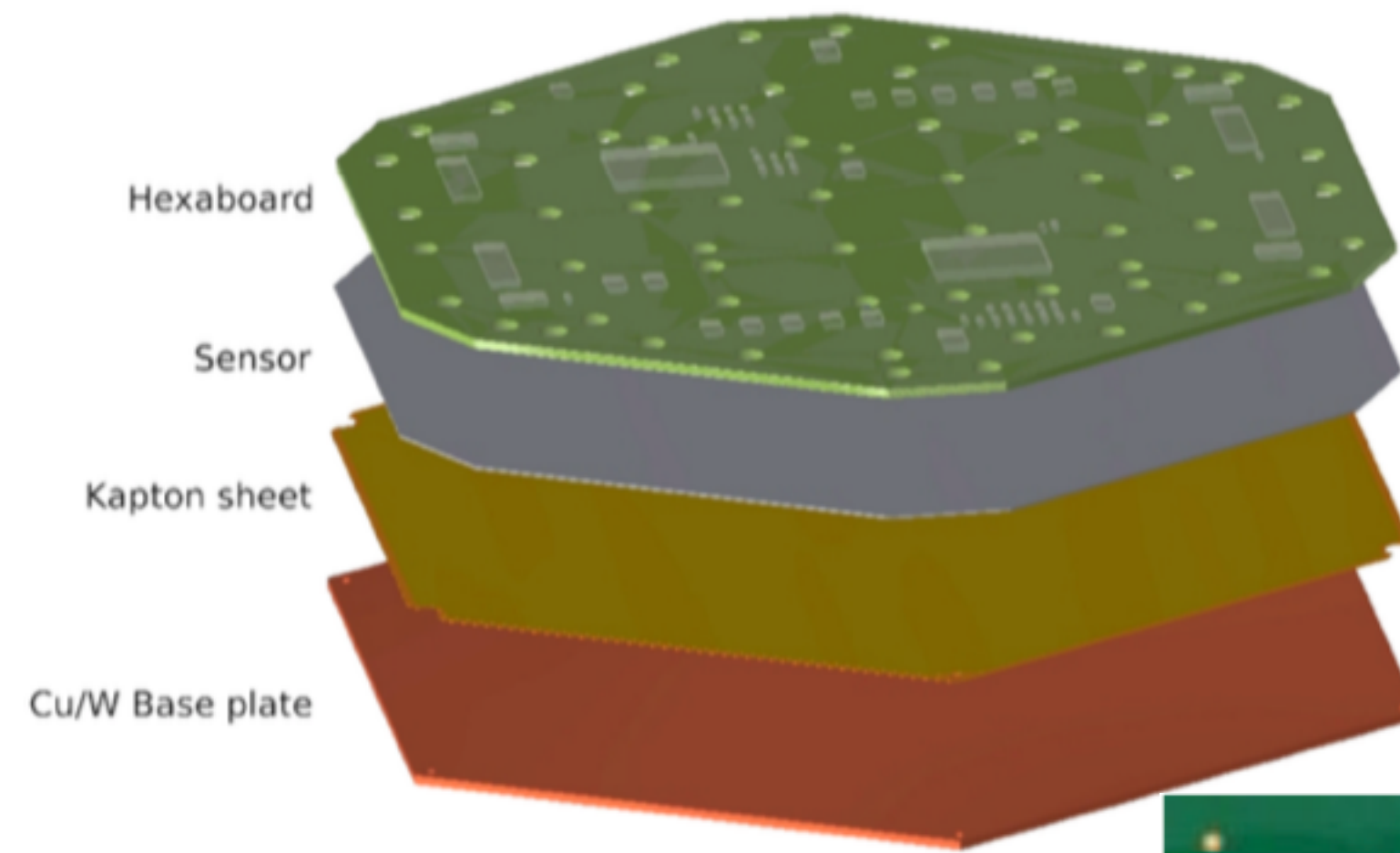
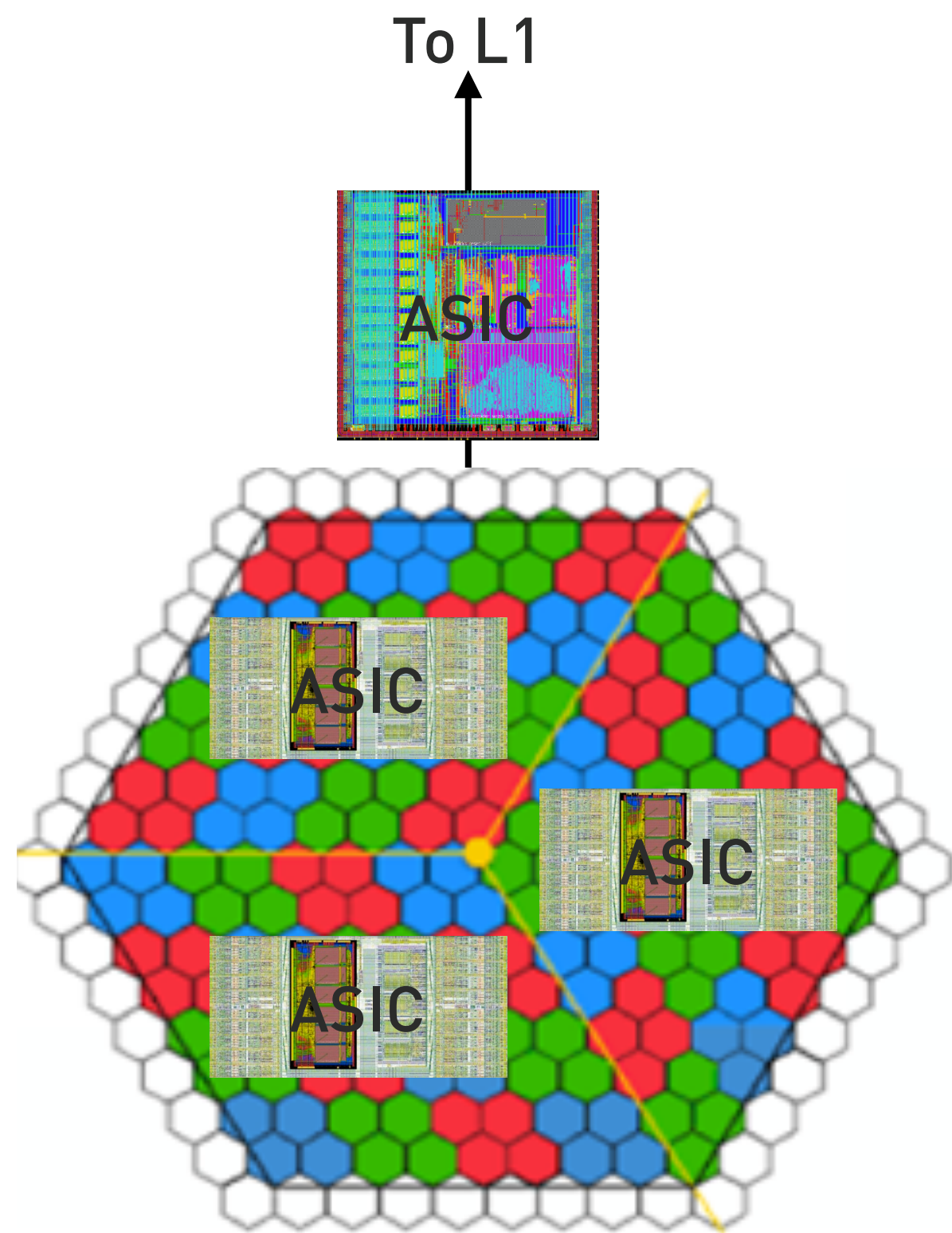
BUT: Cannot read out all these channels fast enough for L1 to trigger!

window



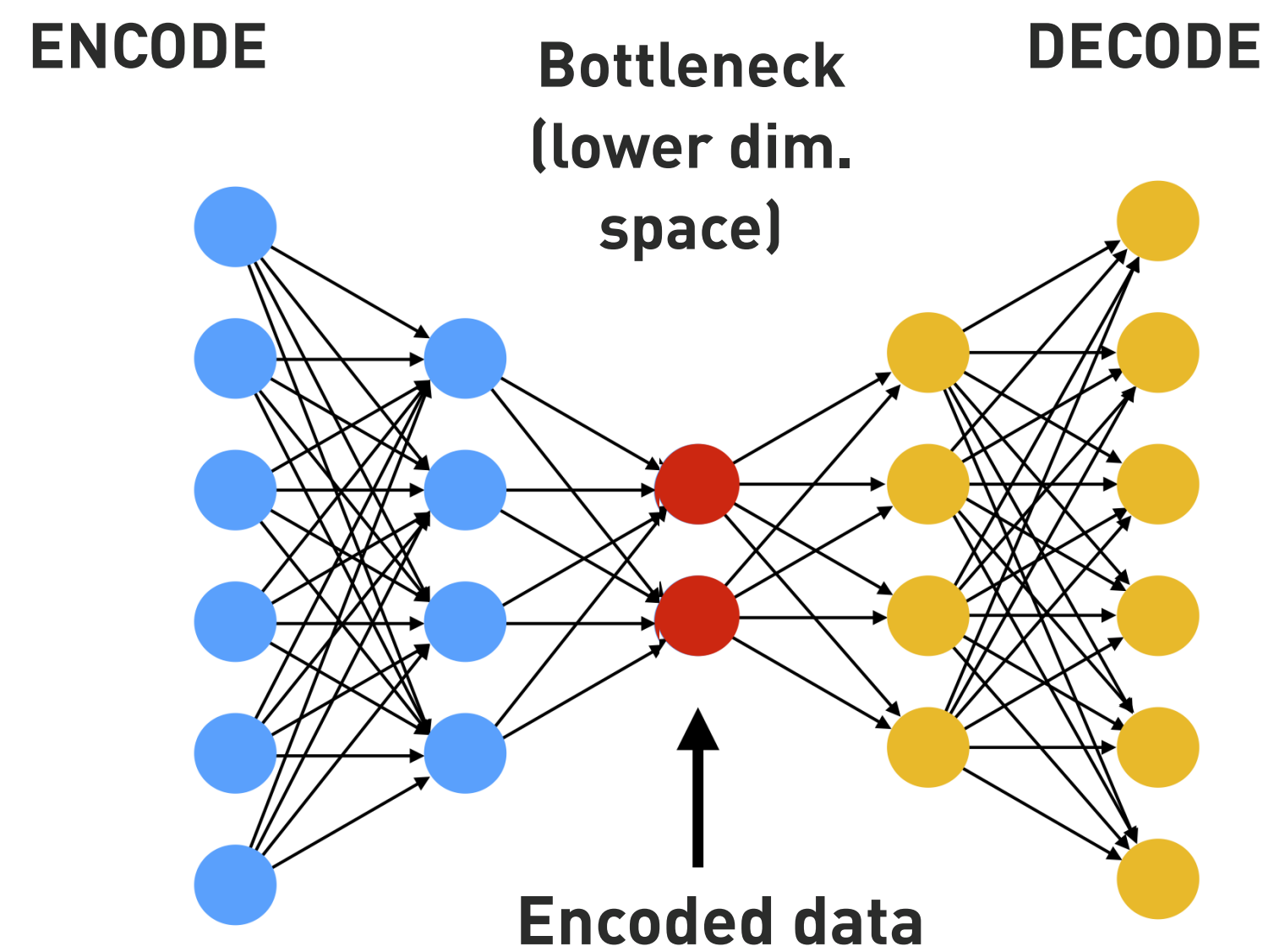
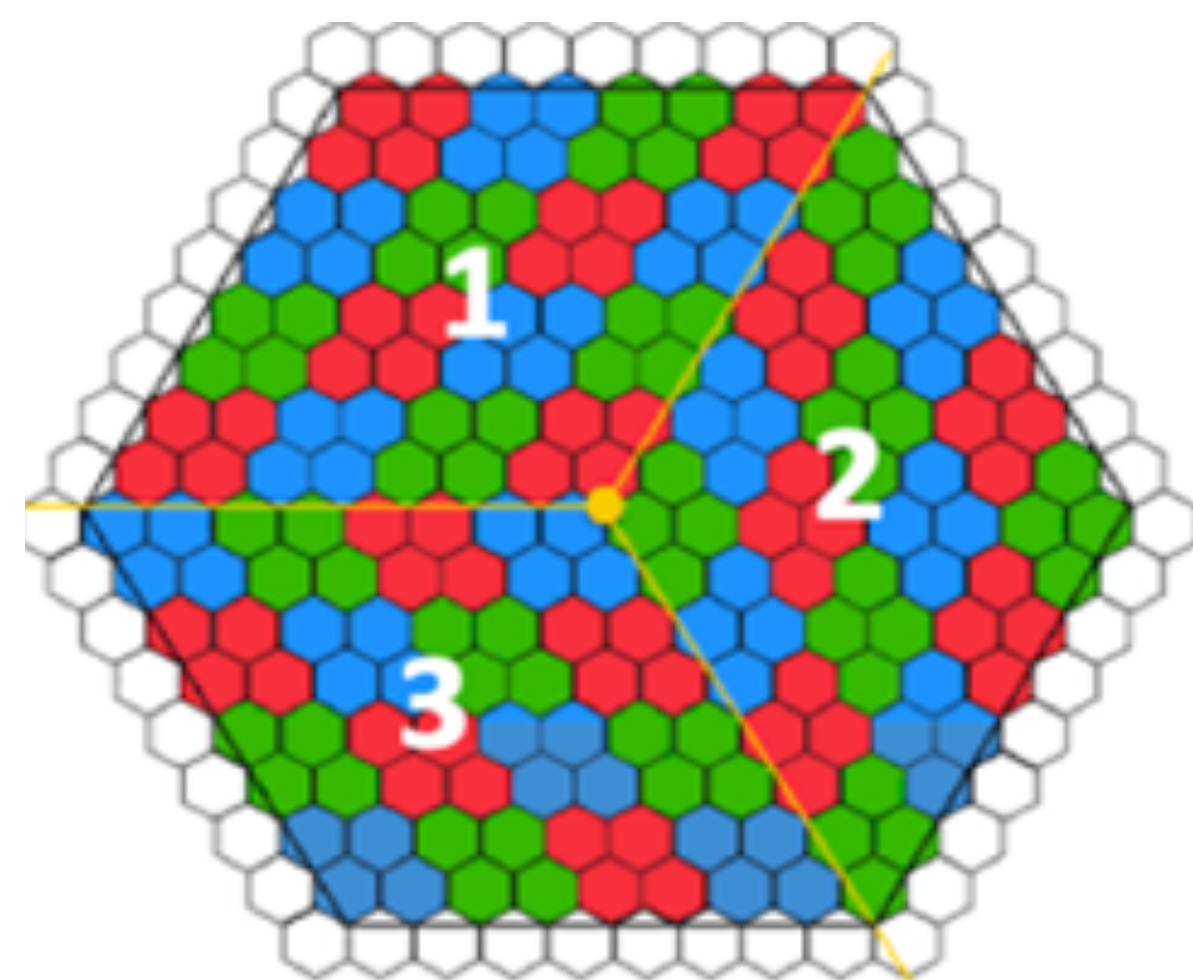
VBF
jet





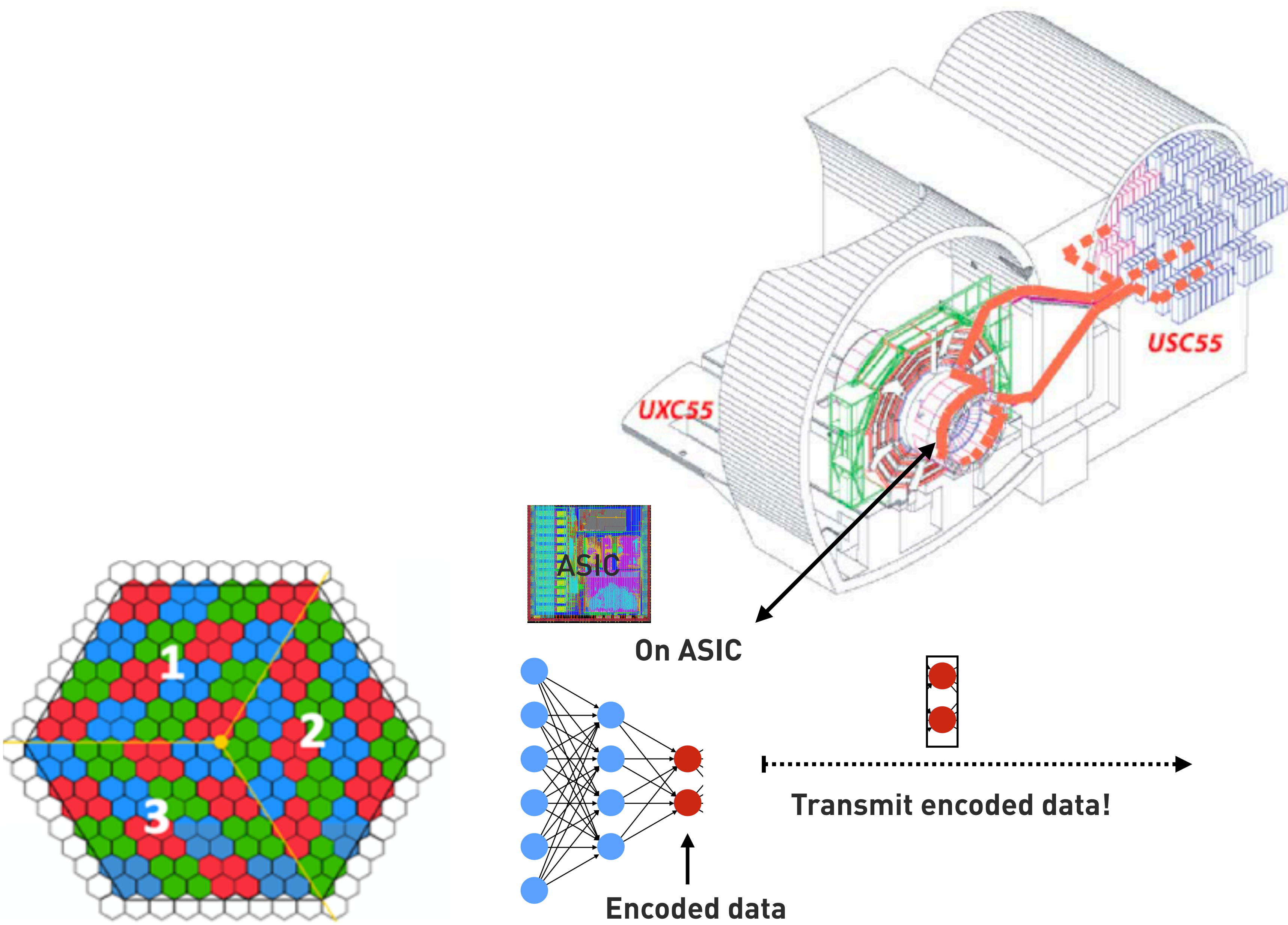
Compress data ON the detector

- ASICs reduce + transmit data
 - 40 MHz trigger data
 - 750 kHz DAQ data
- High radiation
- Cooled to -30 → low power (Max 500 mW total)
- 1.5 μ s latency

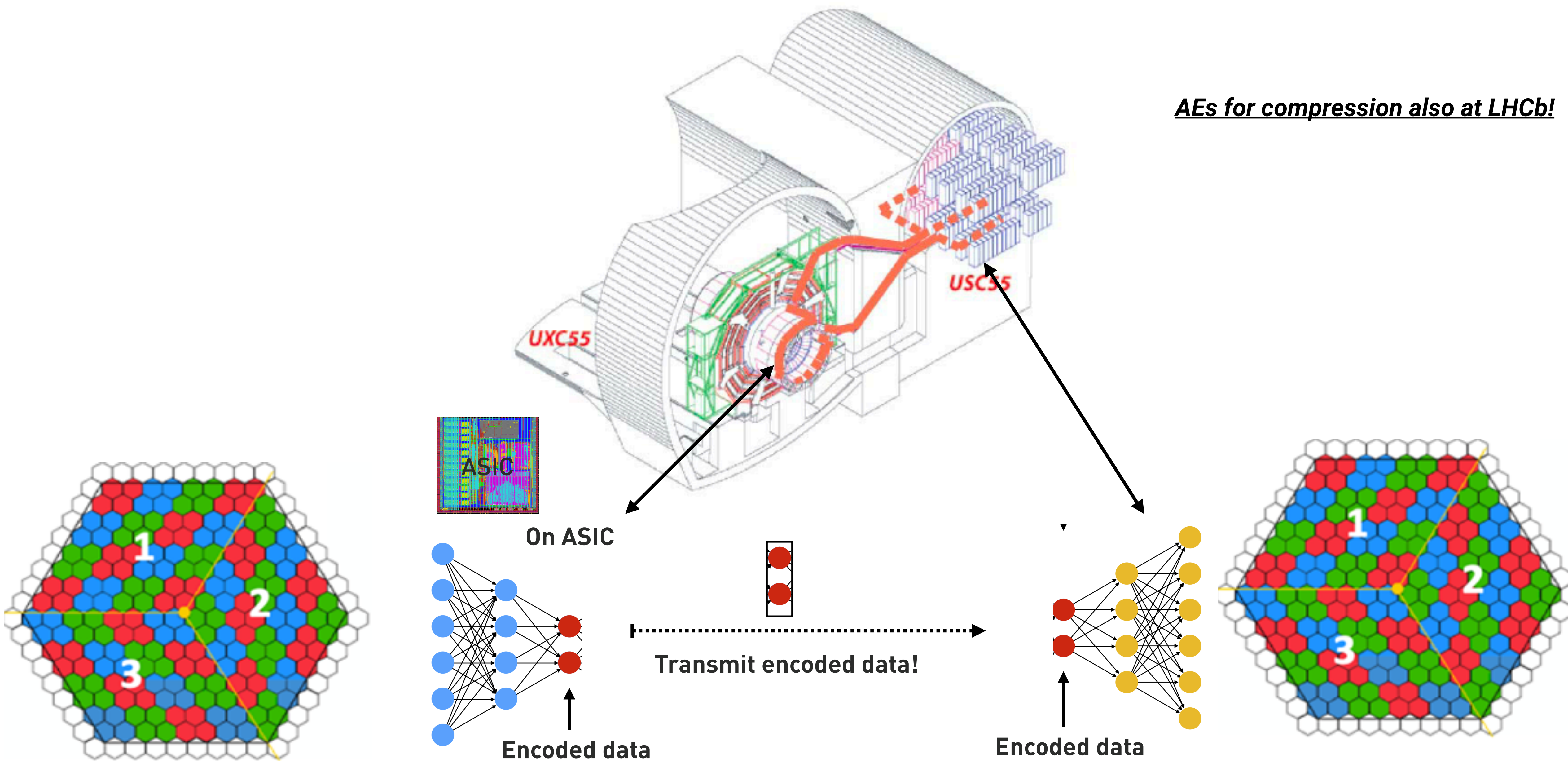


Variational Autoencoder

AEs for compression also at LHCb!

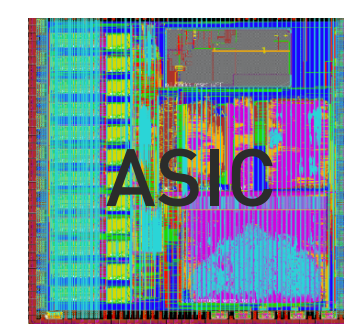
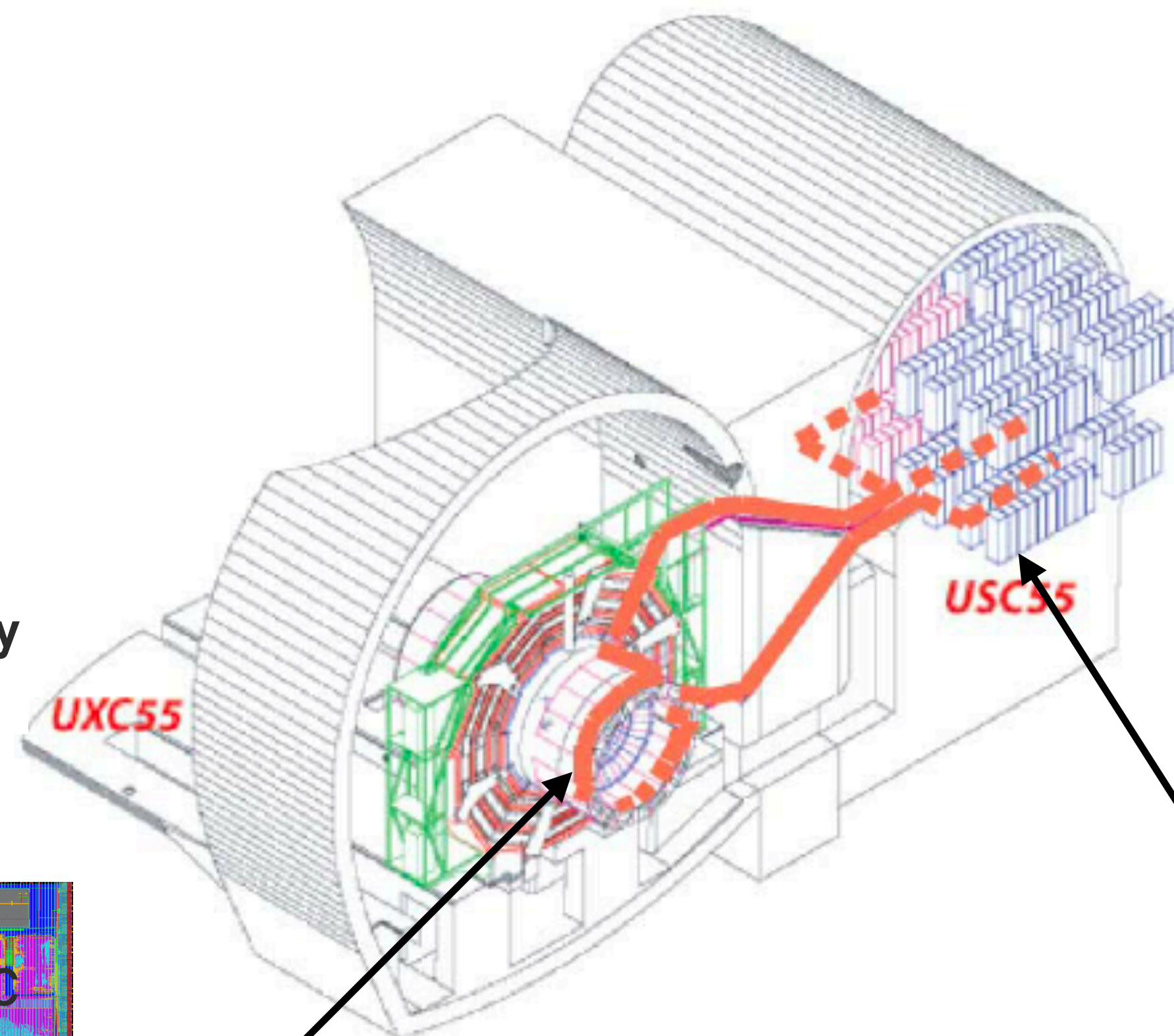


AEs for compression also at LHCb!

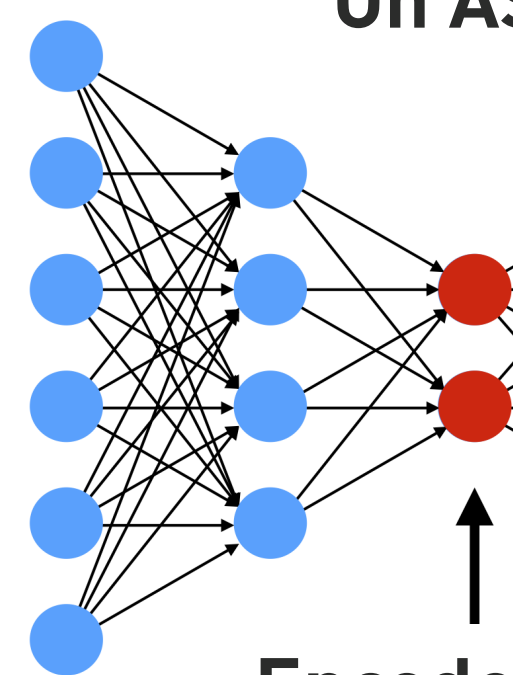


- 75-100 mW
- Triplicated w/b for radiation safety
- Reprogrammable w/b over IC2!

AEs for compression also at LHCb!



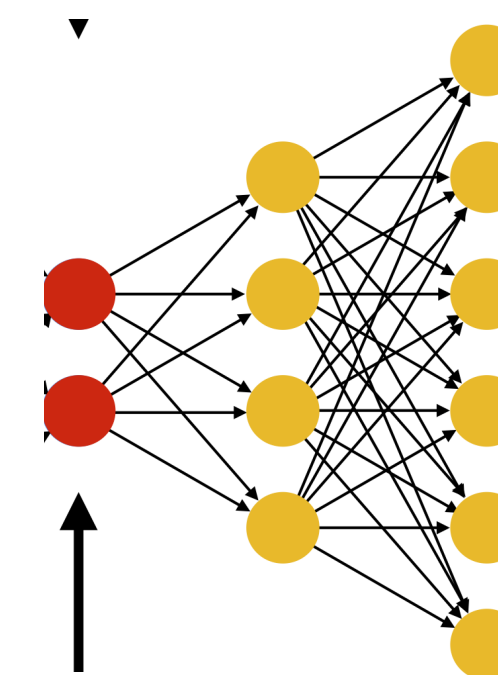
On ASIC



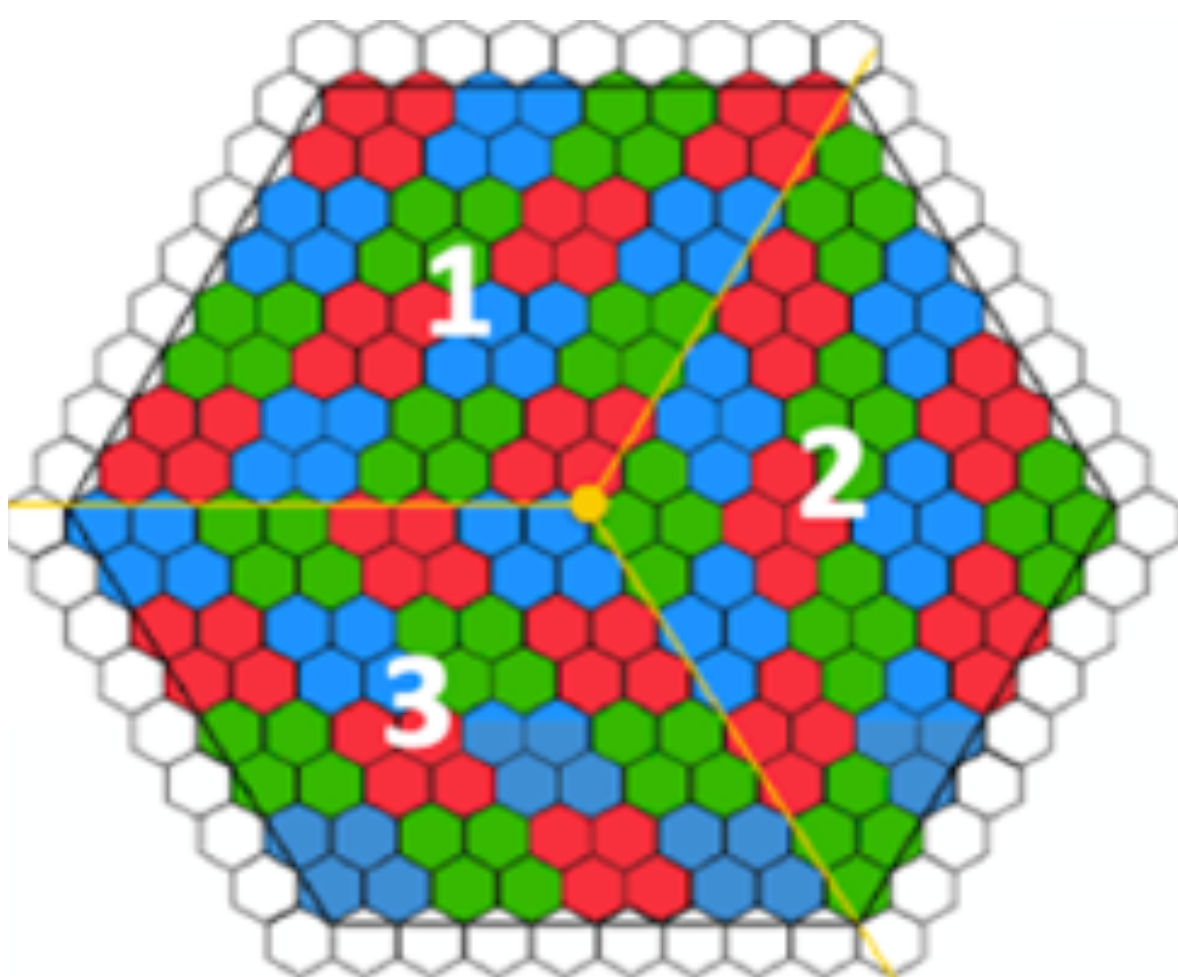
Encoded data



Transmit encoded data!

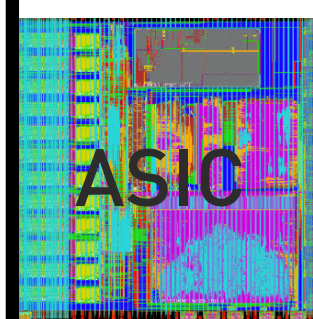
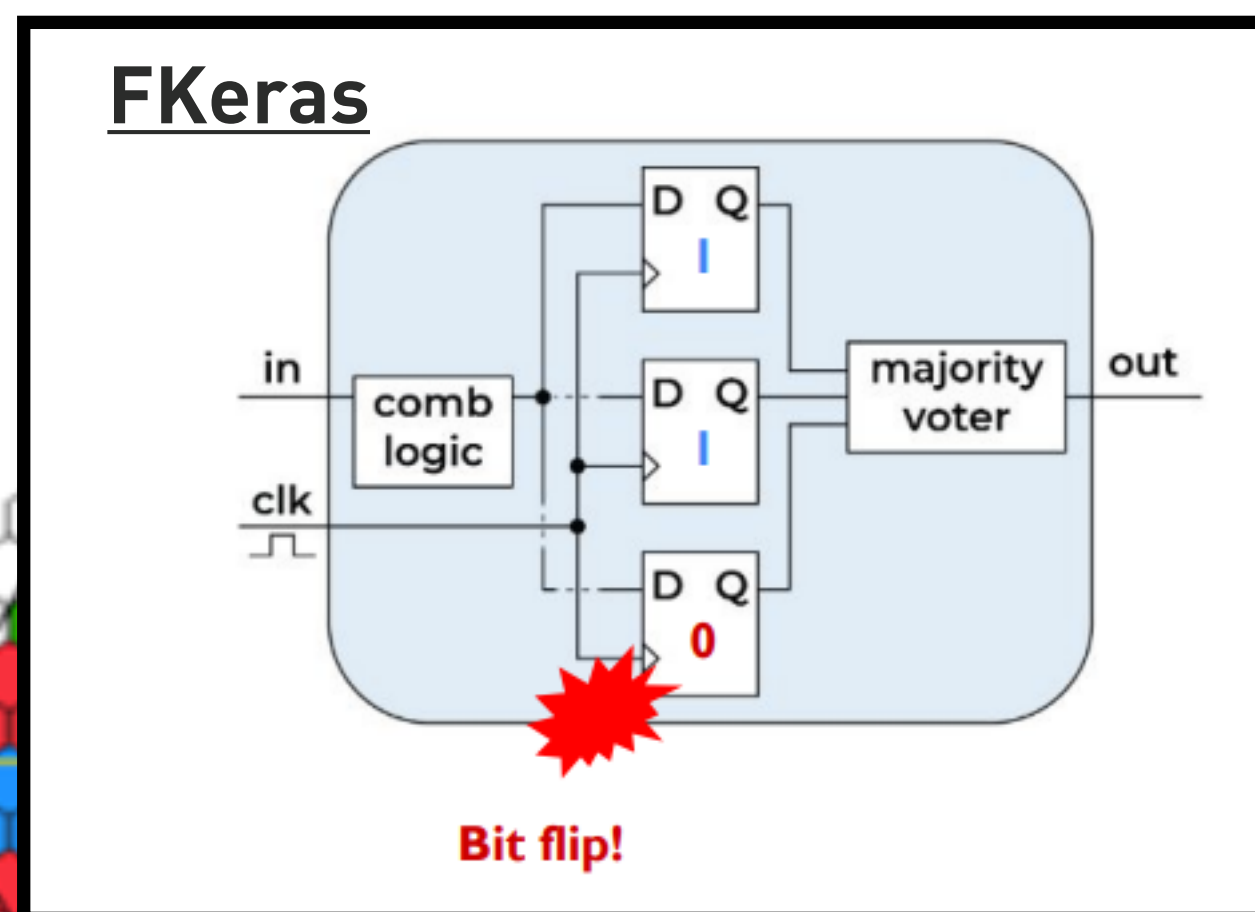
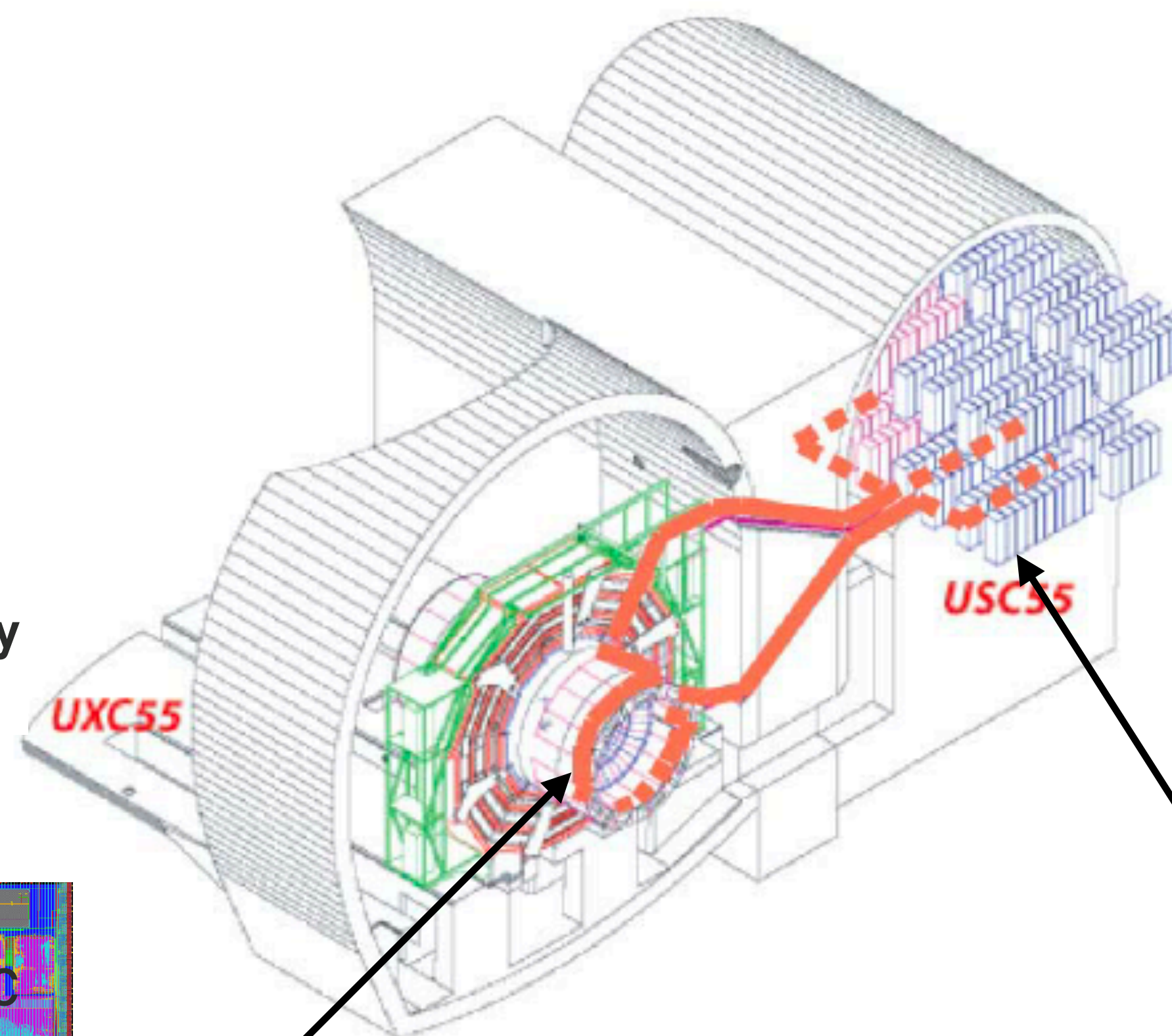


Encoded data

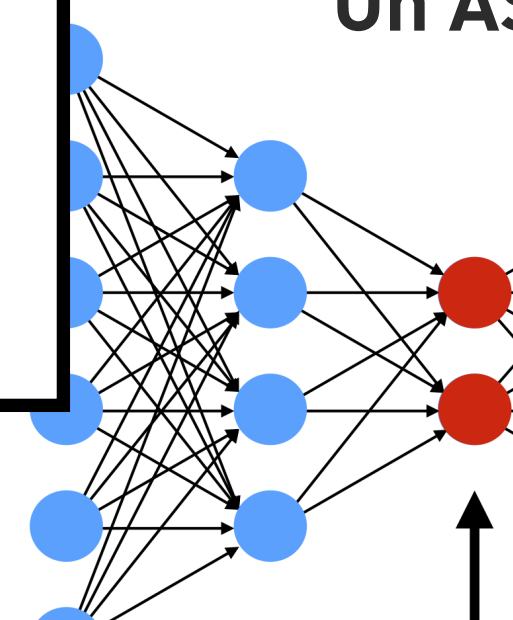


- 75-100 mW
- Triplicated w/b for radiation safety
- Reprogrammable w/b over IC2!

AEs for compression also at LHCb!



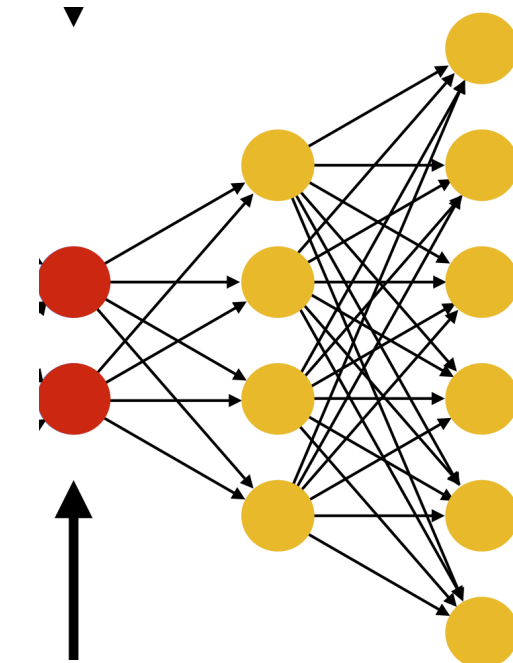
On ASIC



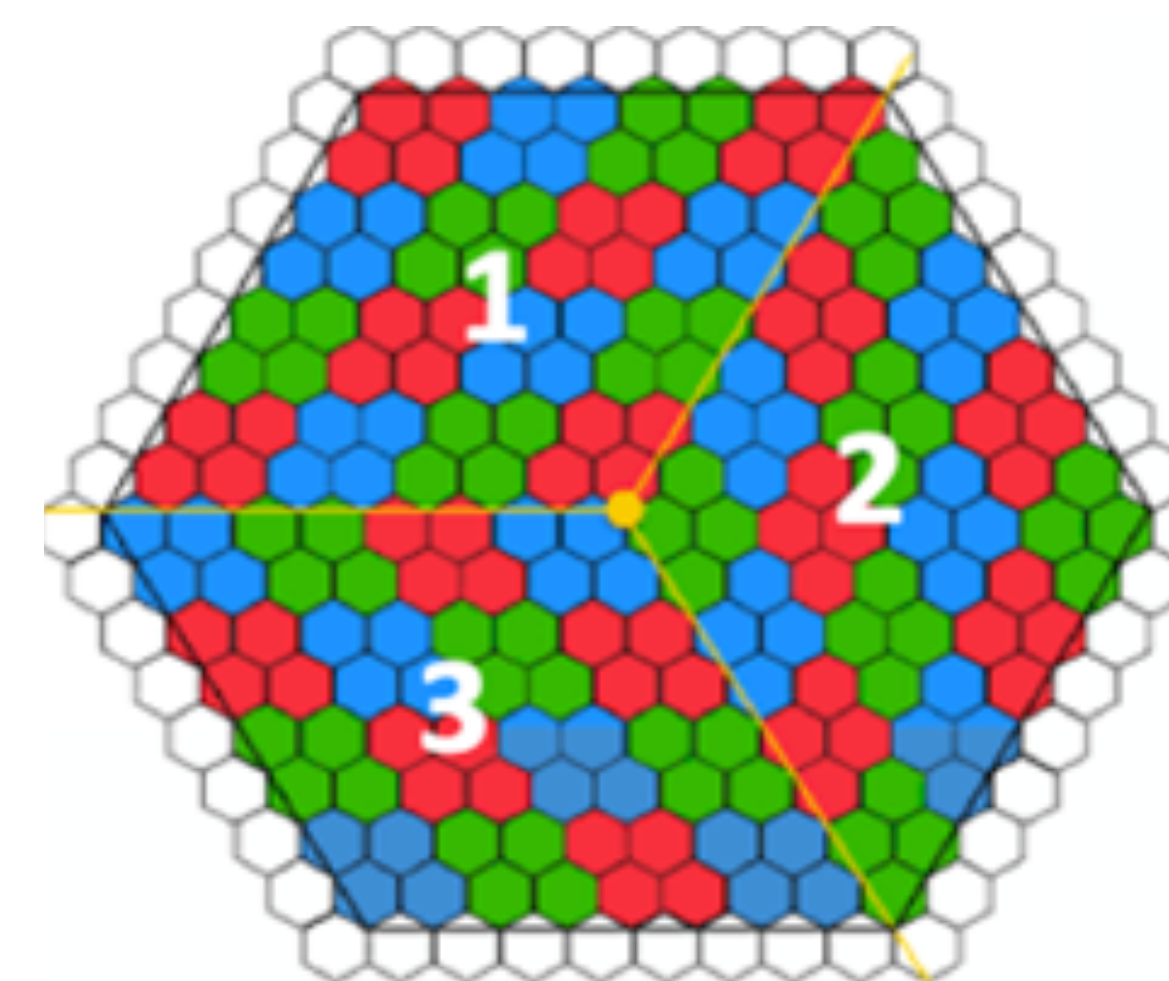
Encoded data



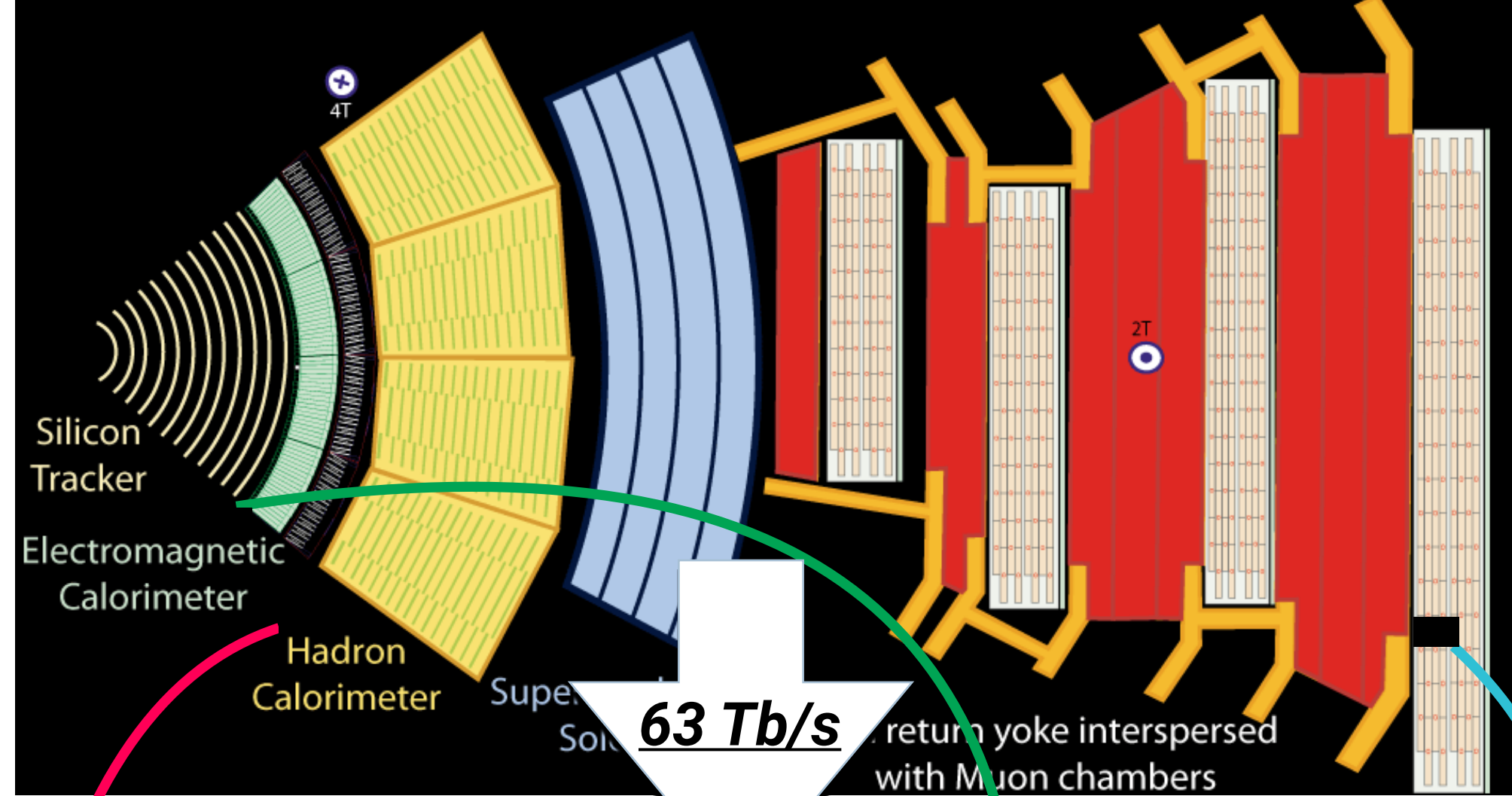
Transmit encoded data!



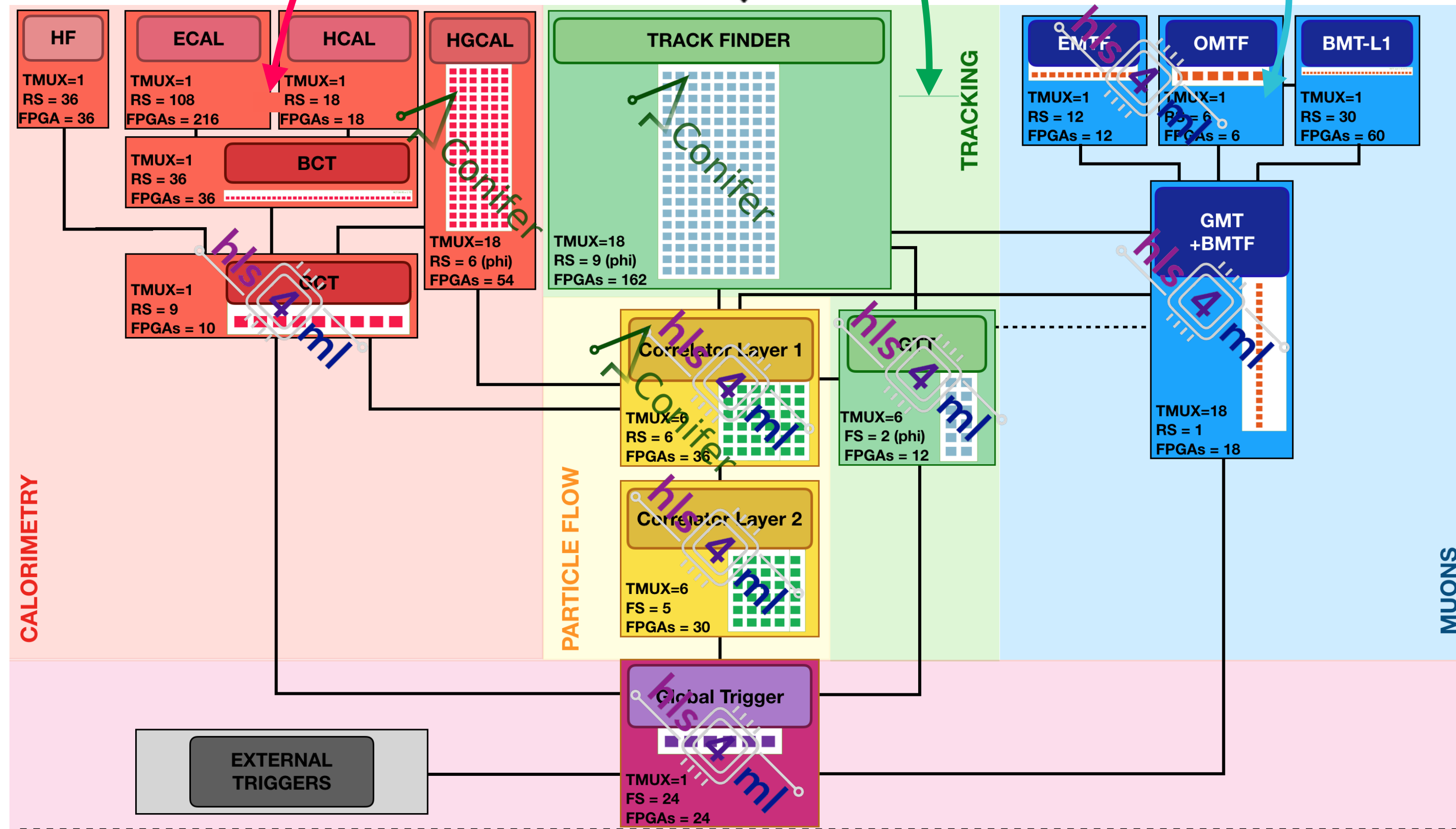
Encoded data



25 billion inferences per second!



63 Tb/s



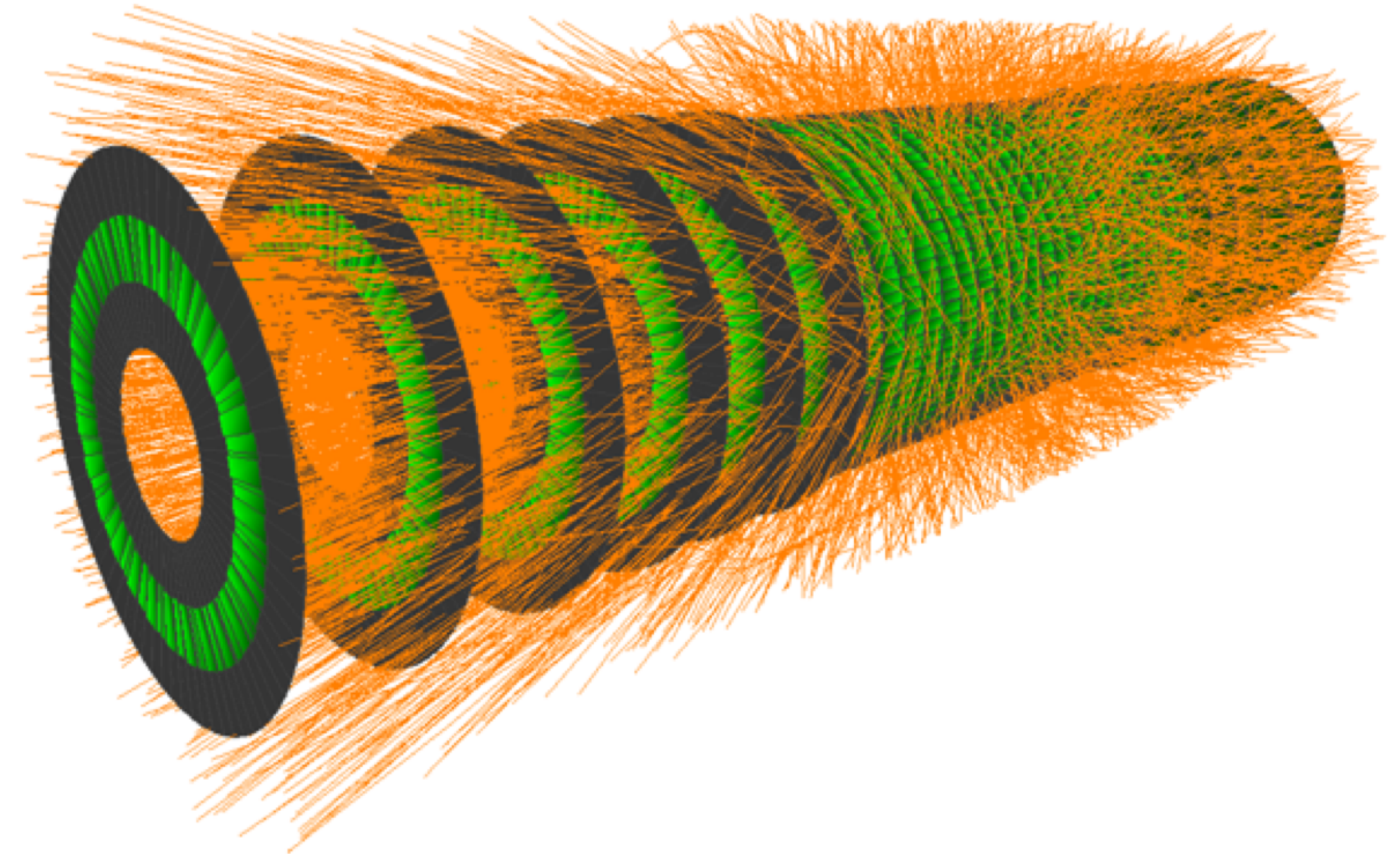
ML for tracking

In HL-LHC, will need to do track finding at L1

- $O(1000)$ hits, $O(100)$ tracks, 40 MHz rate, $\sim 5 \mu\text{s}$ latency

Graph Neural Networks for fast charged particle tracking

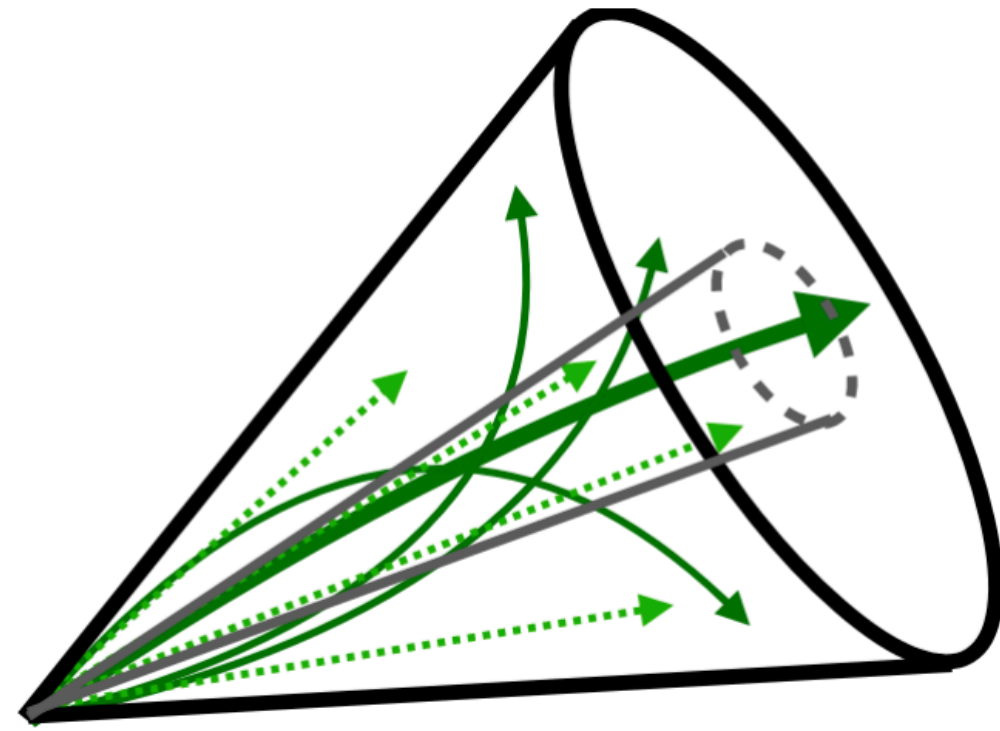
Throughput-optimized for L1 applications,
resource-optimised for co-processing



Design	$(n_{\text{nodes}}, n_{\text{edges}})$	RF	Precision	Latency [cycles]	II [cycles]	DSP [%]	LUT [%]	FF [%]	BRAM [%]
Throughput-opt.	(28, 56)	1	ap_fixed<14, 7>	59 295 ns	1	99.9	66.0	11.7	0.7
Resource-opt.	(28, 56)	1	ap_fixed<14, 7>	79 395 ns	28	56.6	17.6	3.9	13.1

Fast jet tagging

... best representation?



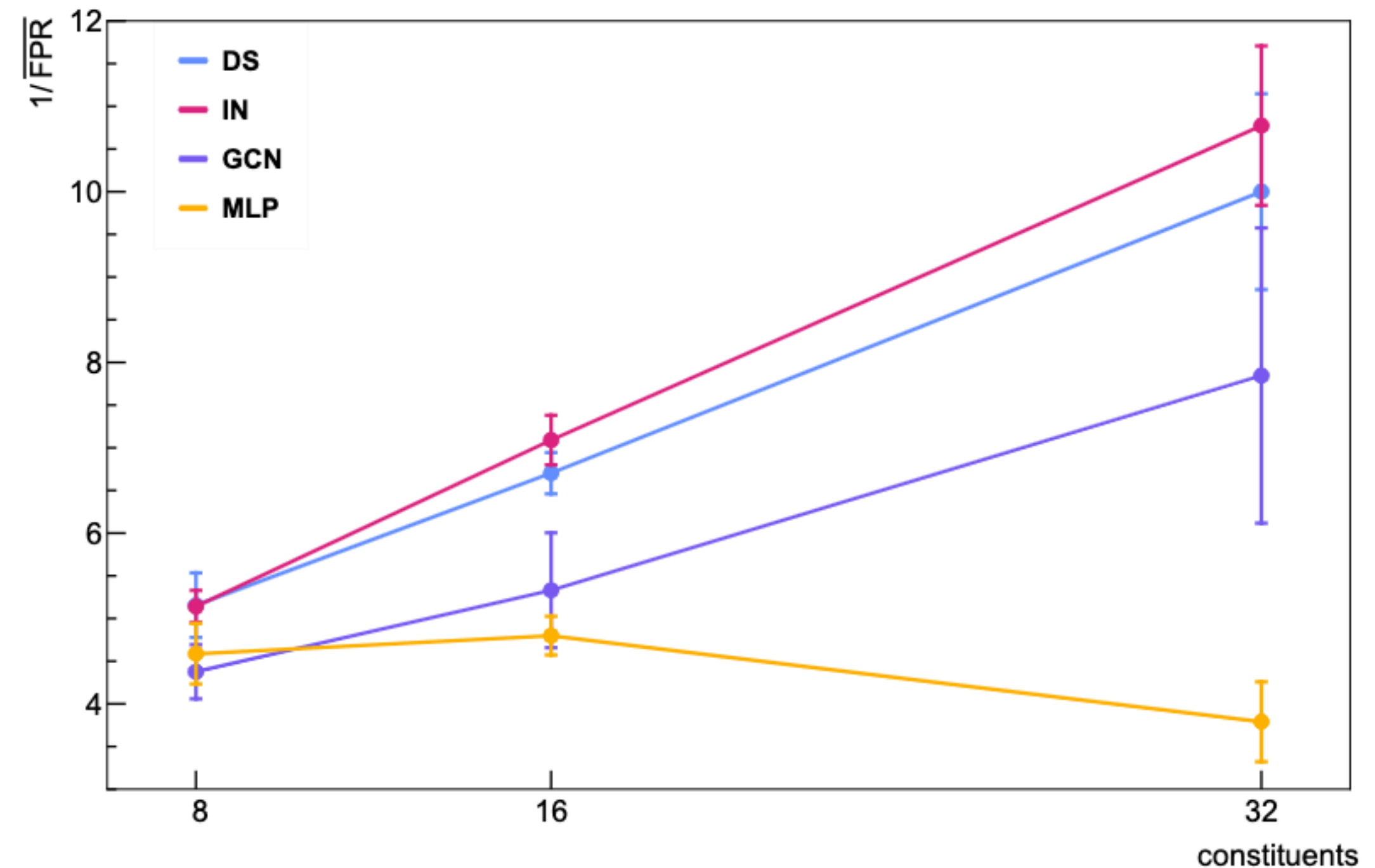
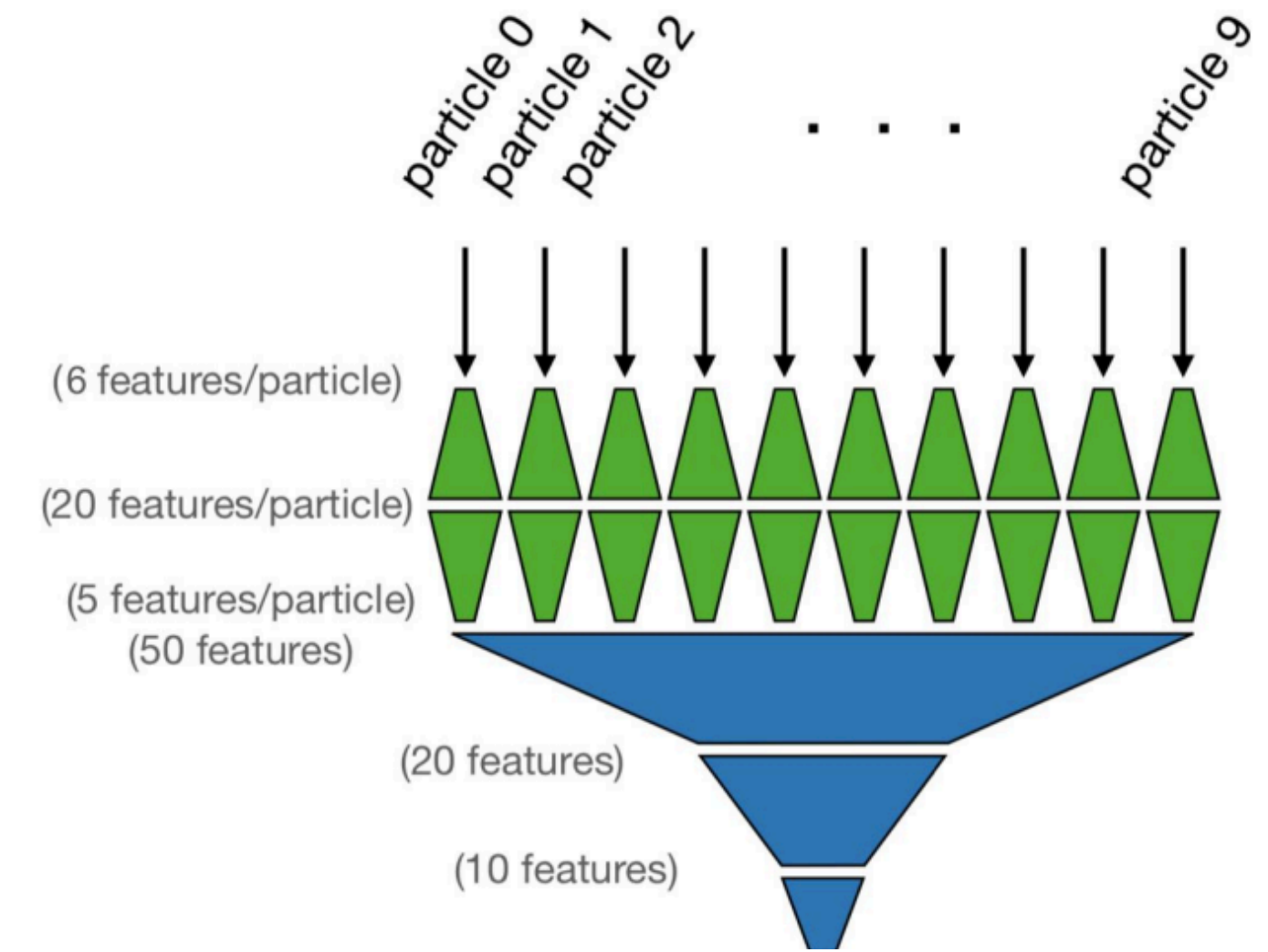
Sets: Information is only assigned to individual nodes.



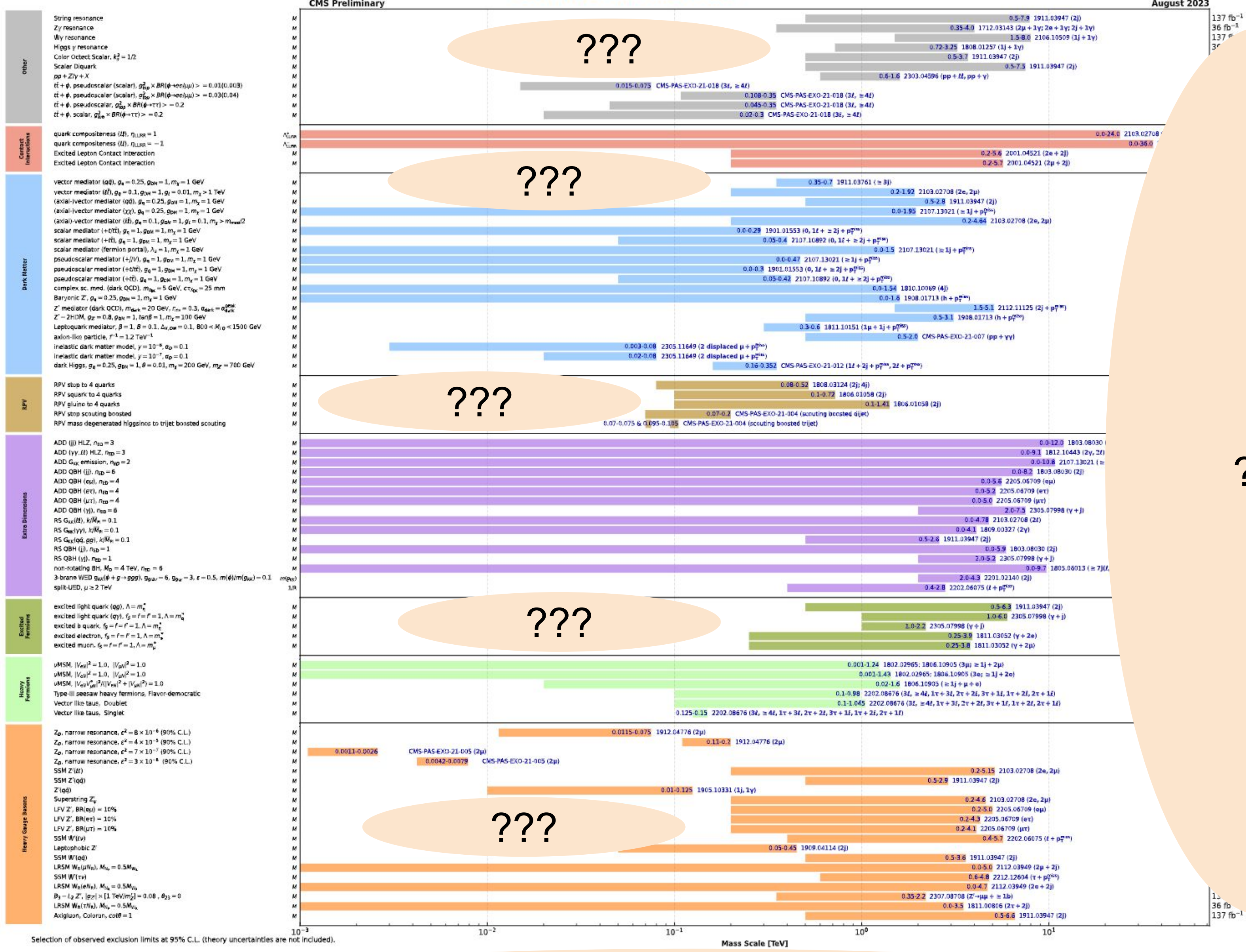
Graphs: Information is assigned to edges, i.e., pairs of nodes.

- Nodes
- Edges
- Features

(Can also do 90 ns transformers for jet tagging!)



Overview of CMS EXO results



Selection of observed exclusion limits at 95% C.L. (theory uncertainties are not included).

???

???

From A. Rizzi



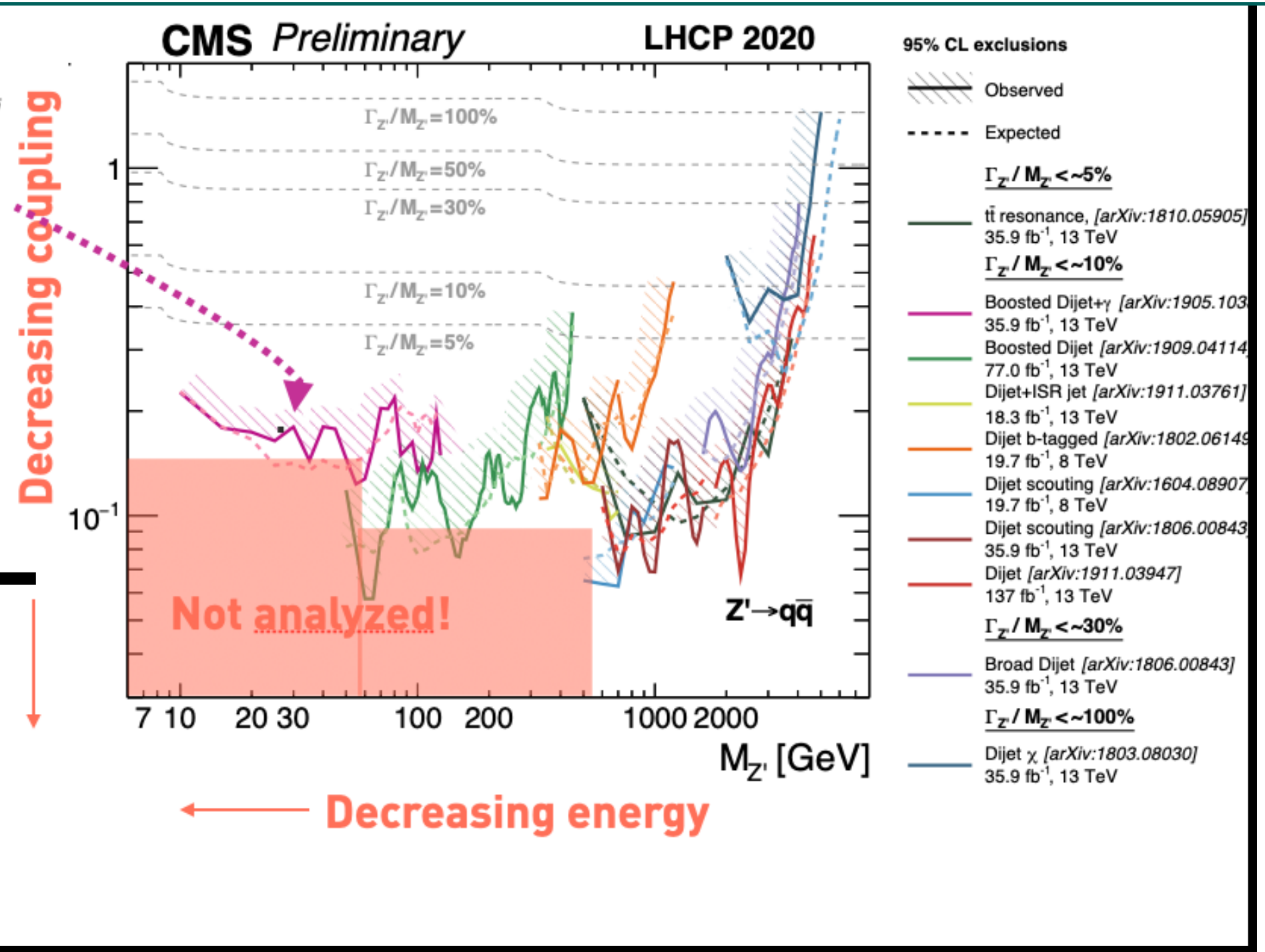
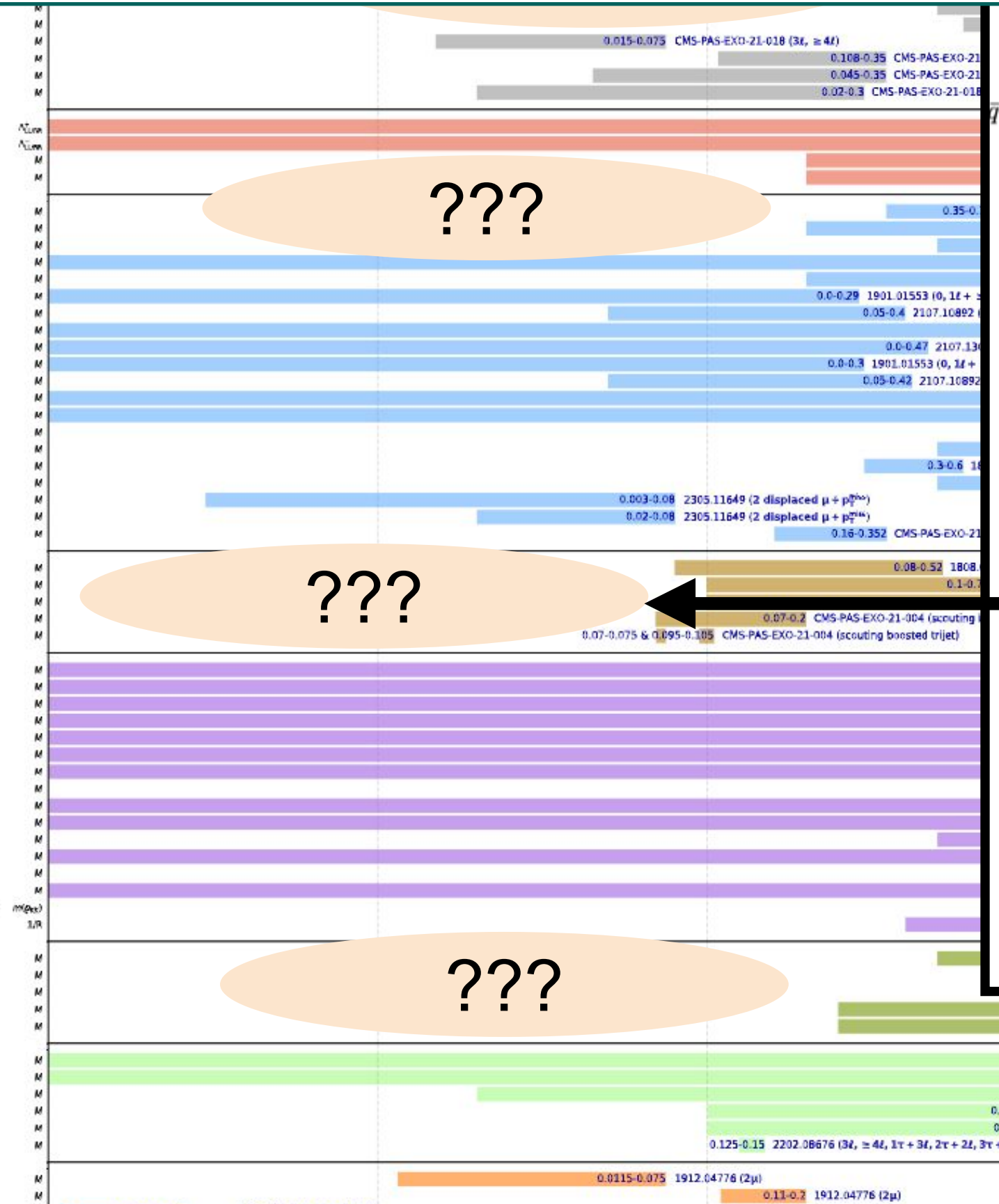
Selection of observed exclusion limits at 95% C.L. (theory uncertainties are not included).

EXPLOIT THE FULL CAPABILITIES OF THE LHC AND BE MORE GENERIC!

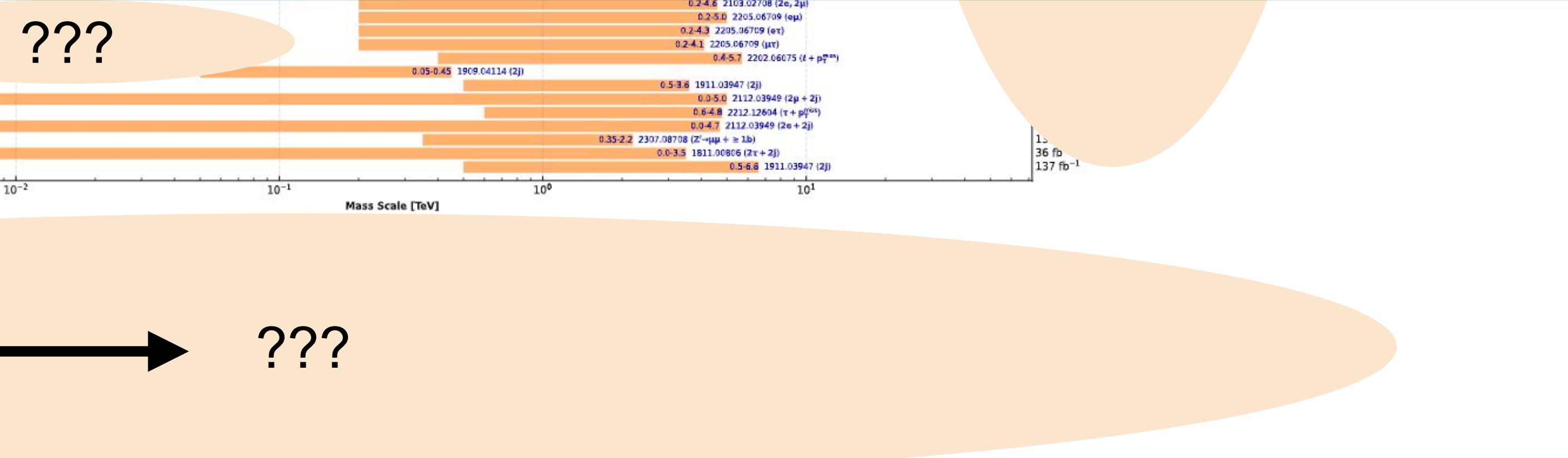
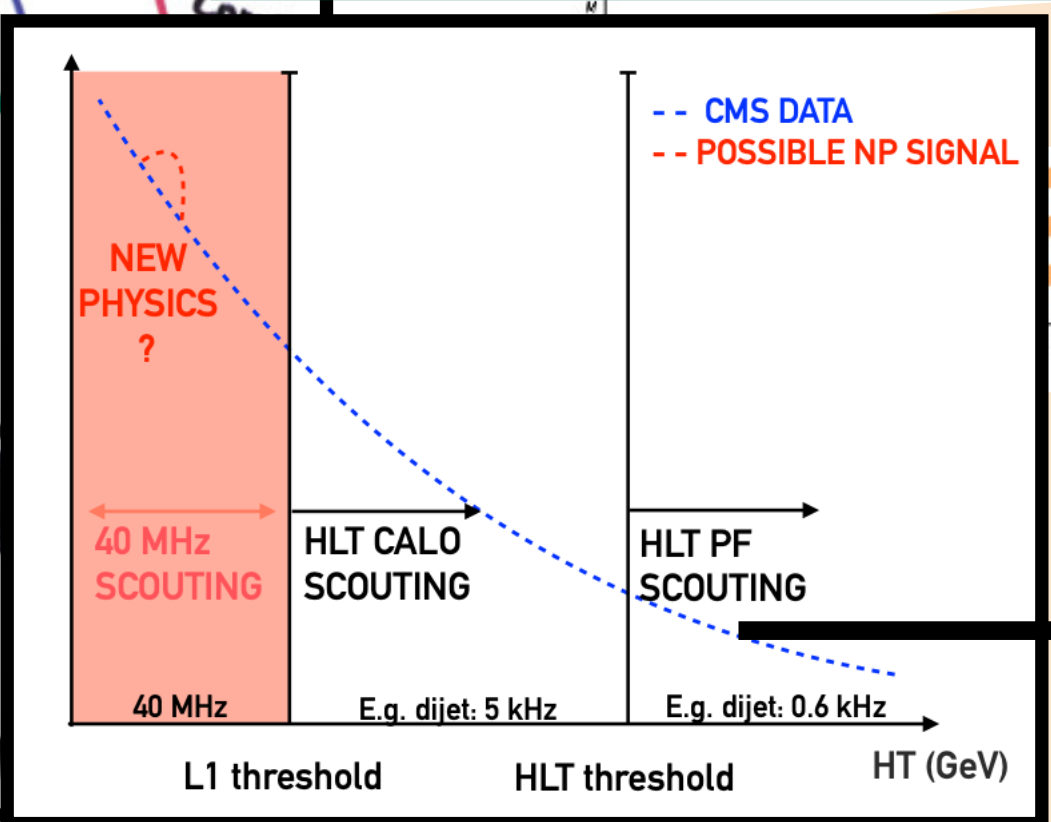
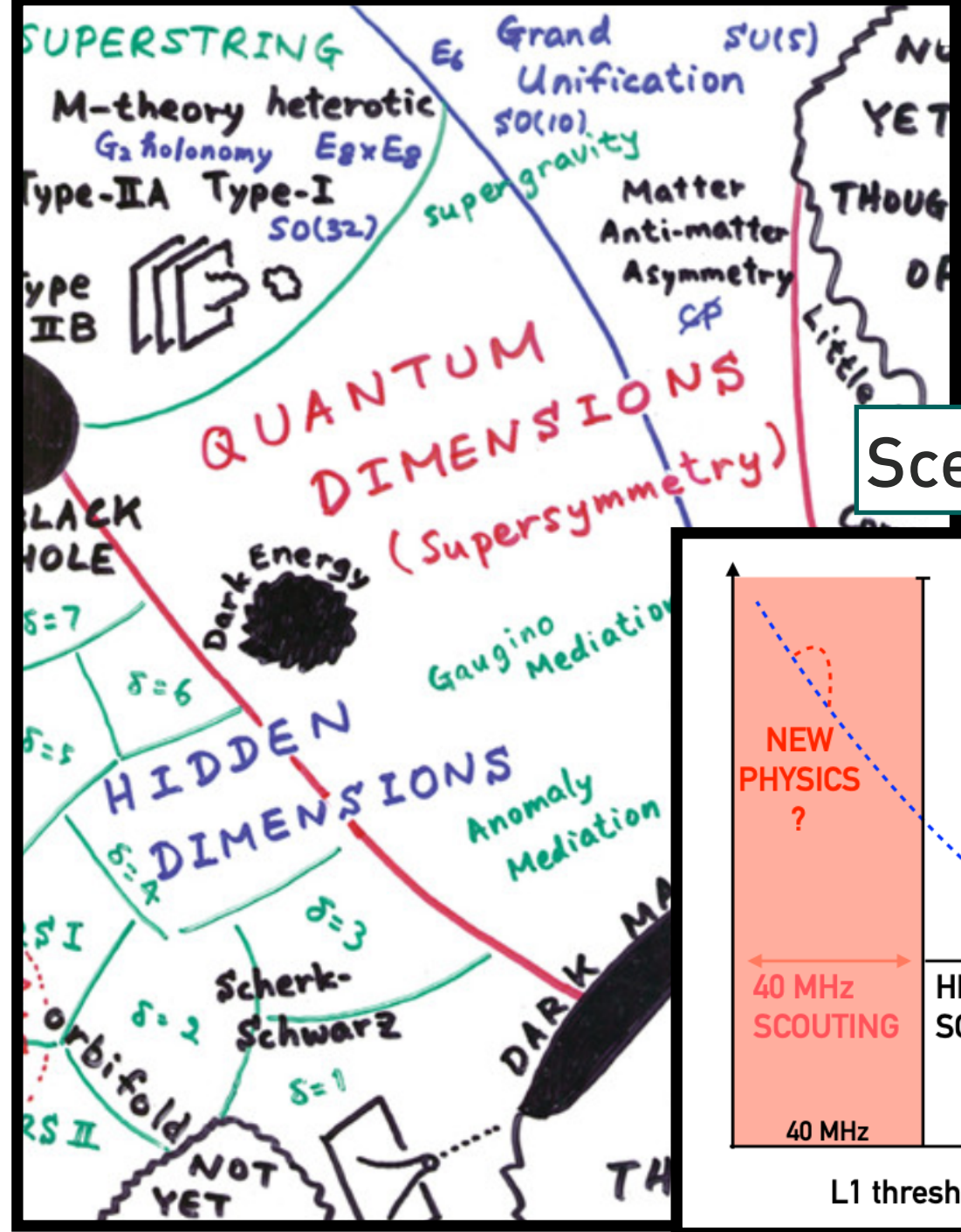
???

Scenario 2) Even for signatures we already look for, some regions out of reach due to L1 trigger

- Other
 - String resonance
 - Z' resonance
 - W' resonance
 - Higgs γ resonance
 - Color Octet Scalar, $k^2 = 1/2$
 - Scalar Diquark
 - $pp \rightarrow Z\gamma + X$
 - $t\bar{t} + \phi$, pseudoscalar (scalar), $g_{\phi q}^2 \times BR(\phi \rightarrow \tau\tau) > = 0.01(0.003)$
 - $t\bar{t} + \phi$, pseudoscalar (scalar), $g_{\phi q}^2 \times BR(\phi \rightarrow \tau\tau) > = 0.03(0.04)$
 - $t\bar{t} + \phi$, pseudoscalar, $g_{\phi q}^2 \times BR(\phi \rightarrow \tau\tau) > = 0.2$
 - $t\bar{t} + \phi$, scalar, $g_{\phi q}^2 \times BR(\phi \rightarrow \tau\tau) > = 0.2$
- Contact Interactions
 - quark compositeness (U), $n_{U,qq} = 1$
 - quark compositeness (U), $n_{U,qq} = -1$
 - Excited Lepton Contact Interaction
 - Excited Lepton Contact Interaction
- Dark Matter
 - vector mediator (v), $g_v = 0.25, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - vector mediator (v), $g_v = 0.1, g_{\text{SM}} = 1, g_{\nu} = 0.01, m_{\nu} > 1 \text{ TeV}$
 - (axial-)vector mediator (v), $g_v = 0.25, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - (axial-)vector mediator (v), $g_v = 0.25, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - (axial-)vector mediator (v), $g_v = 0.1, g_{\text{SM}} = 1, g_{\nu} = 0.1, m_{\nu} > m_{\text{WIMZ}}$
 - scalar mediator (s), $g_s = 1, g_{\text{SM}} = 1, m_{\nu} > m_{\text{WIMZ}}$
 - scalar mediator (s), $g_s = 1, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - scalar mediator (fermion portal), $\lambda_s = 1, m_{\nu} = 1 \text{ GeV}$
 - pseudoscalar mediator (p), $g_p = 1, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - pseudoscalar mediator (p), $g_p = 1, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - pseudoscalar mediator (p), $g_p = 1, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - complex sc. med. (dark OCD), $m_{\nu} = 5 \text{ GeV}, c_{\text{SM}} = 25 \text{ mm}$
 - Baryonic Z', $g_{\nu} = 0.25, g_{\text{SM}} = 1, m_{\nu} = 1 \text{ GeV}$
 - Z' mediator (dark OCD), $m_{\text{WIMZ}} = 20 \text{ GeV}, r_{\text{SM}} = 0.3, g_{\text{SM}} = g_{\text{SM}}^{\text{SM}}$
 - Z' - 2HDM, $g_{\nu} = 0.8, g_{\text{SM}} = 1, \tan\beta = 1, m_{\nu} = 100 \text{ GeV}$
 - Leptoquark mediator, $\beta = 1, \delta = 0.1, \Delta_{\text{SM}} = 0.1, 800 < M_{\text{LQ}} < 1500 \text{ GeV}$
 - axion-like particle, $f^{-1} = 1.2 \text{ TeV}^{-1}$
 - inelastic dark matter model, $\gamma = 10^{-4}, a_{\text{D}} = 0.1$
 - inelastic dark matter model, $\gamma = 10^{-7}, a_{\text{D}} = 0.1$
 - dark Higgs, $g_{\nu} = 0.25, g_{\text{SM}} = 1, \theta = 0.01, m_{\nu} = 200 \text{ GeV}, m_{\nu} = 700 \text{ GeV}$
- RPV
 - RPV stop to 4 quarks
 - RPV squark to 4 quarks
 - RPV gluino to 4 quarks
 - RPV stop scouping boosted
 - RPV mass-degenerated higgsinos to trilepton boosted scouping
- Extra Dimensions
 - ADD (d) HLZ, $n_{\text{D}} = 3$
 - ADD (y, t) HLZ, $n_{\text{D}} = 3$
 - ADD G_{μ} emission, $n_{\text{D}} = 2$
 - ADD QBH (b), $n_{\text{D}} = 6$
 - ADD QBH (e), $n_{\text{D}} = 4$
 - ADD QBH (t), $n_{\text{D}} = 4$
 - ADD QBH (u), $n_{\text{D}} = 4$
 - ADD QBH (d), $n_{\text{D}} = 4$
 - ADD QBH (s), $n_{\text{D}} = 4$
 - RS G_{μ} (t), $k/M_{\text{Pl}} = 0.1$
 - RS G_{μ} (y), $k/M_{\text{Pl}} = 0.1$
 - RS G_{μ} (q), $k/M_{\text{Pl}} = 0.1$
 - RS QBH (t), $n_{\text{D}} = 1$
 - RS QBH (y), $n_{\text{D}} = 1$



Scenario 1) There is some NP signature we haven't thought of and we do not trigger on



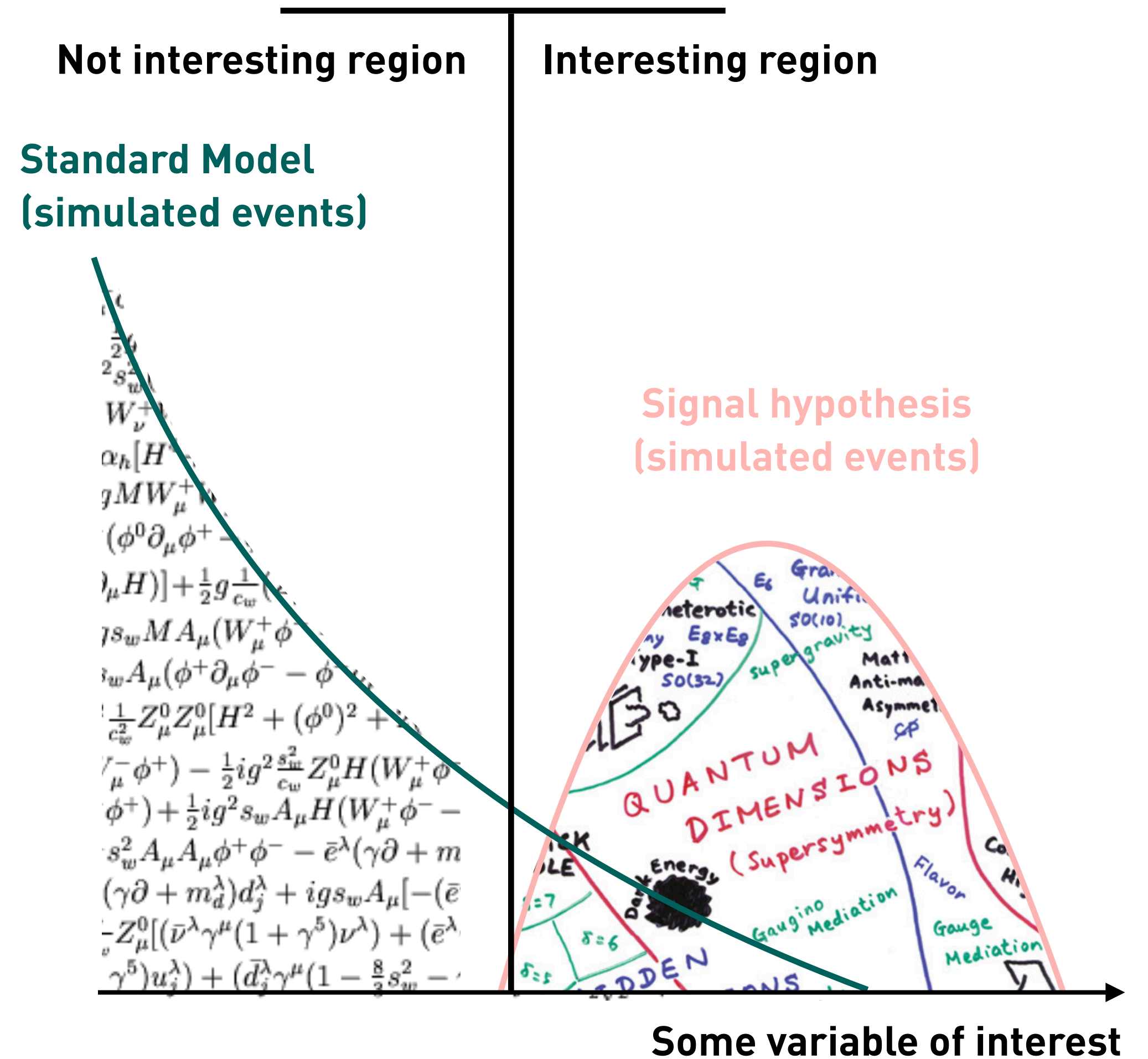
Bias in particle physics

Searches at LHC always start by

- assuming **Standard Model** and some **signal hypothesis**

This is fine when we know what “signal” is (like Higgs)

- **Tailor search** to a given theory
- Powerful, but **limited to model of choice**



Bias in particle physics

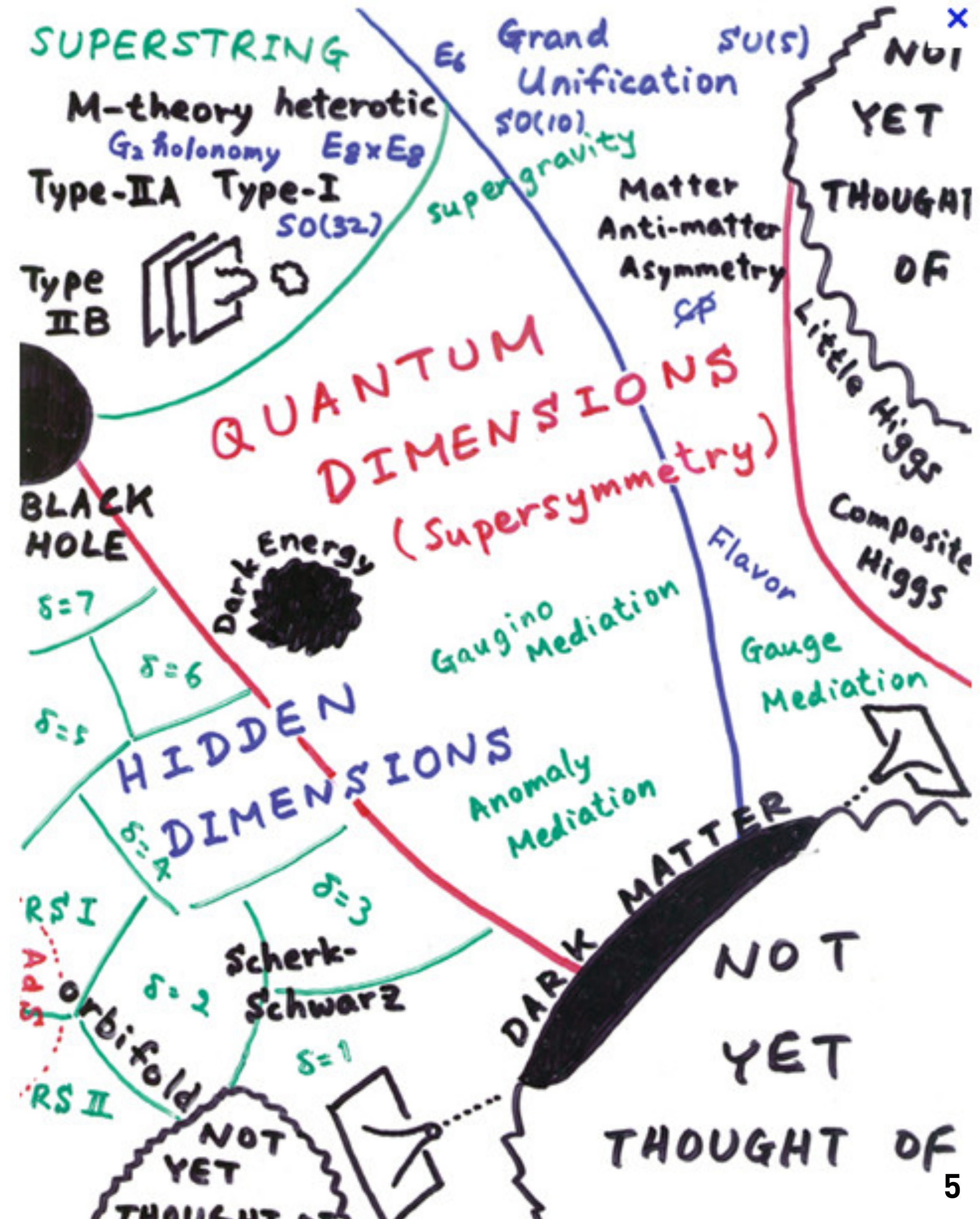
Searches at LHC always start by

- assuming **Standard Model** and some **signal hypothesis**

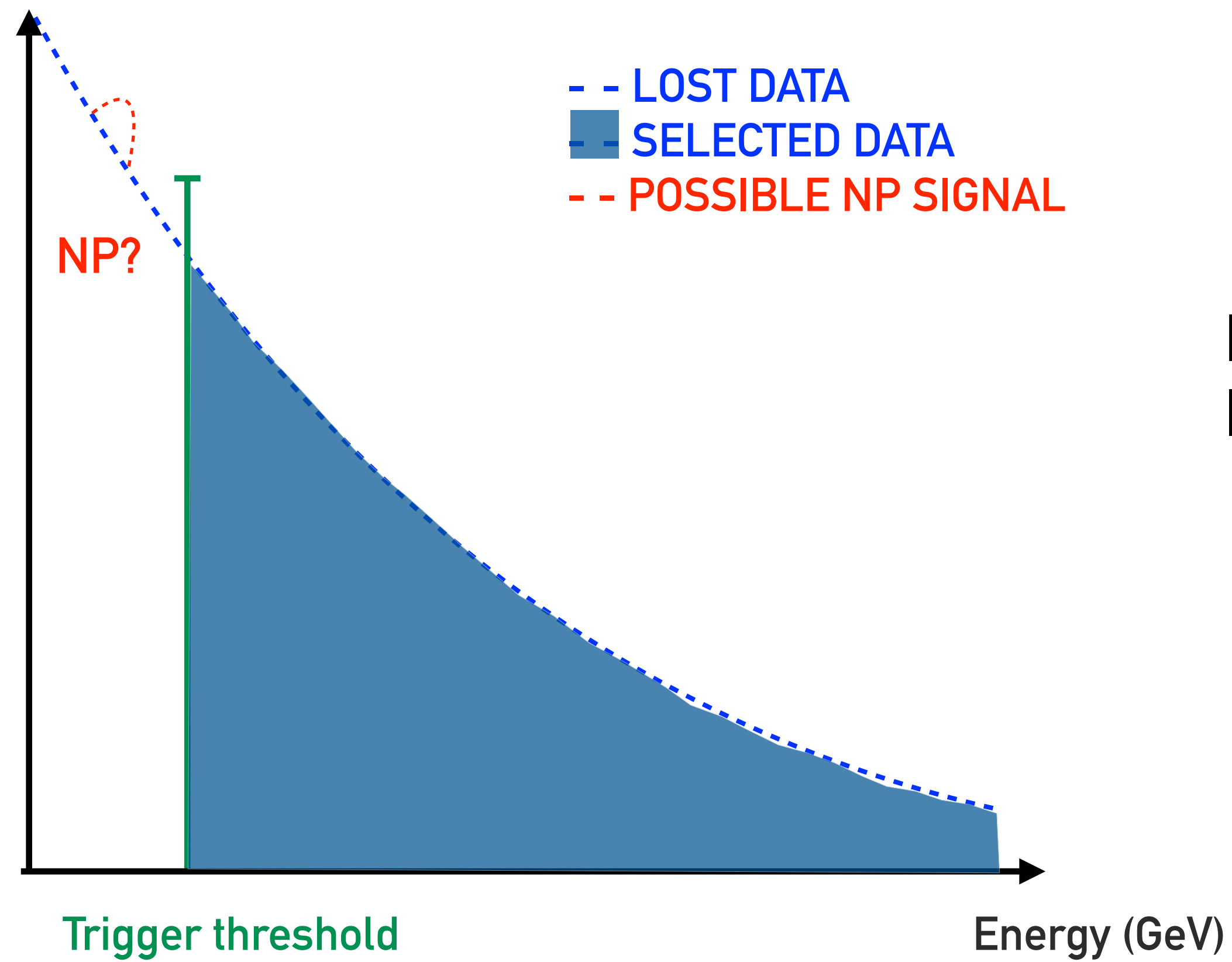
This is fine when we know what “signal” is (like Higgs)

- **Tailor search** to a given theory
- Powerful, but **limited to model of choice**

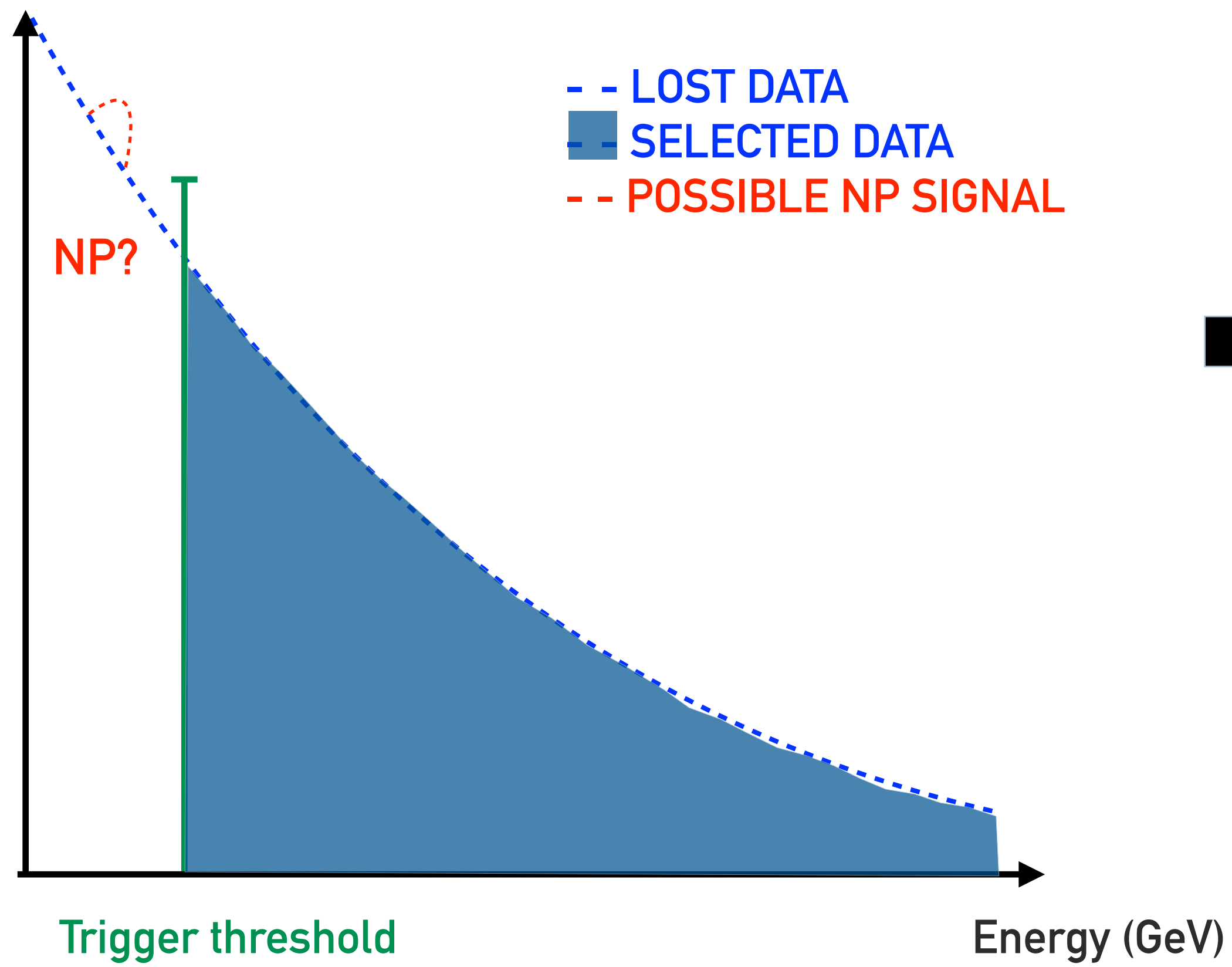
- How do we know we are looking for the right thing in the enormous New Physics model landscape?



Limitations of current trigger

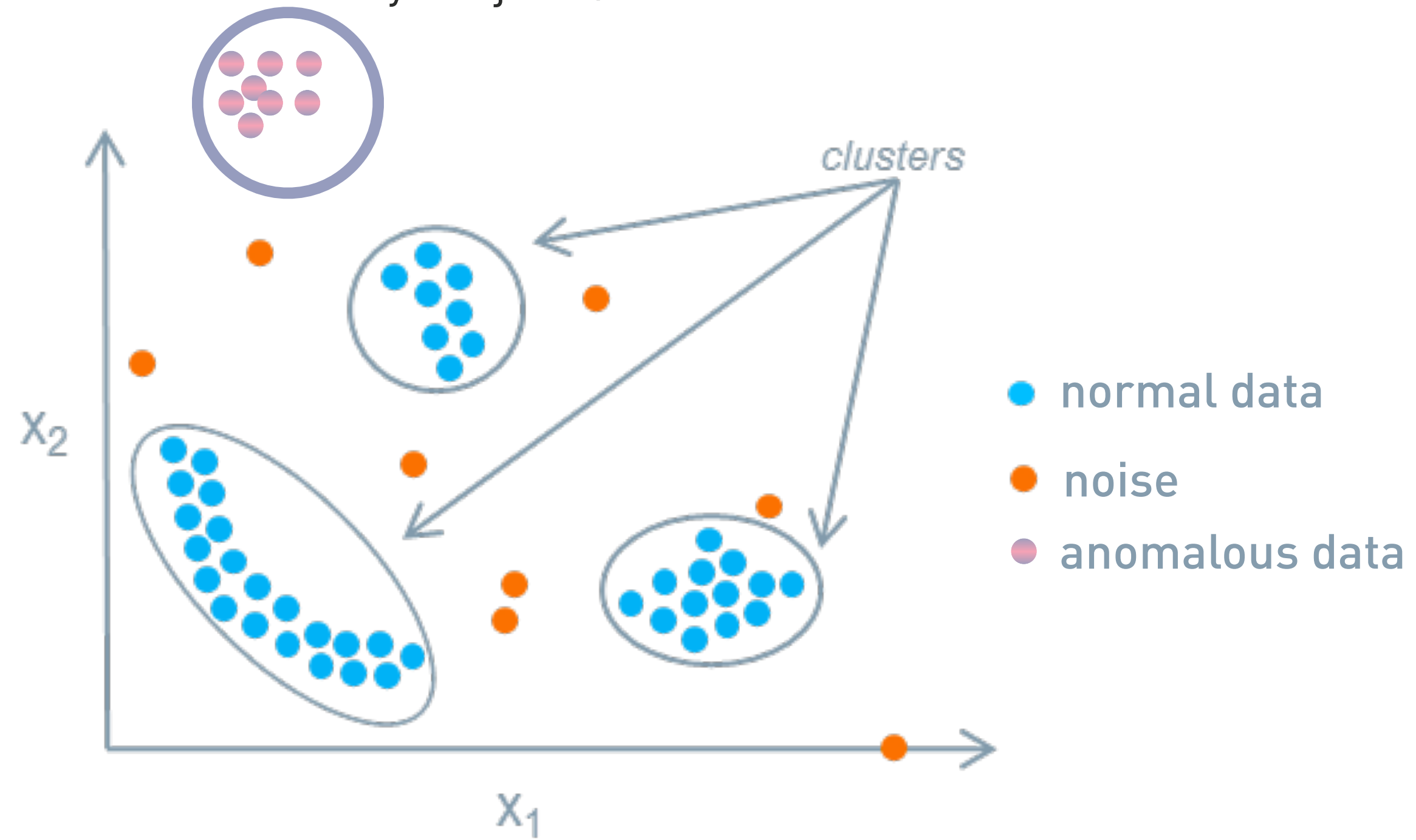


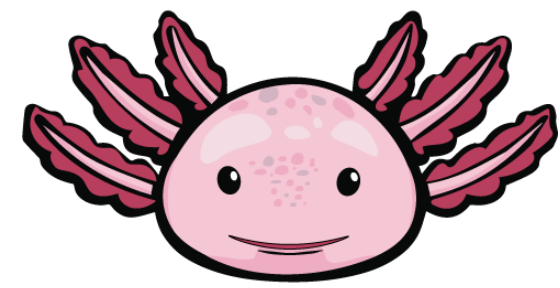
Level-1 rejects >99% of events!
Is there a smarter way to select?



Look at **data** rather than defining signal hypothesis a priori

- Can we "classify" objects/events?





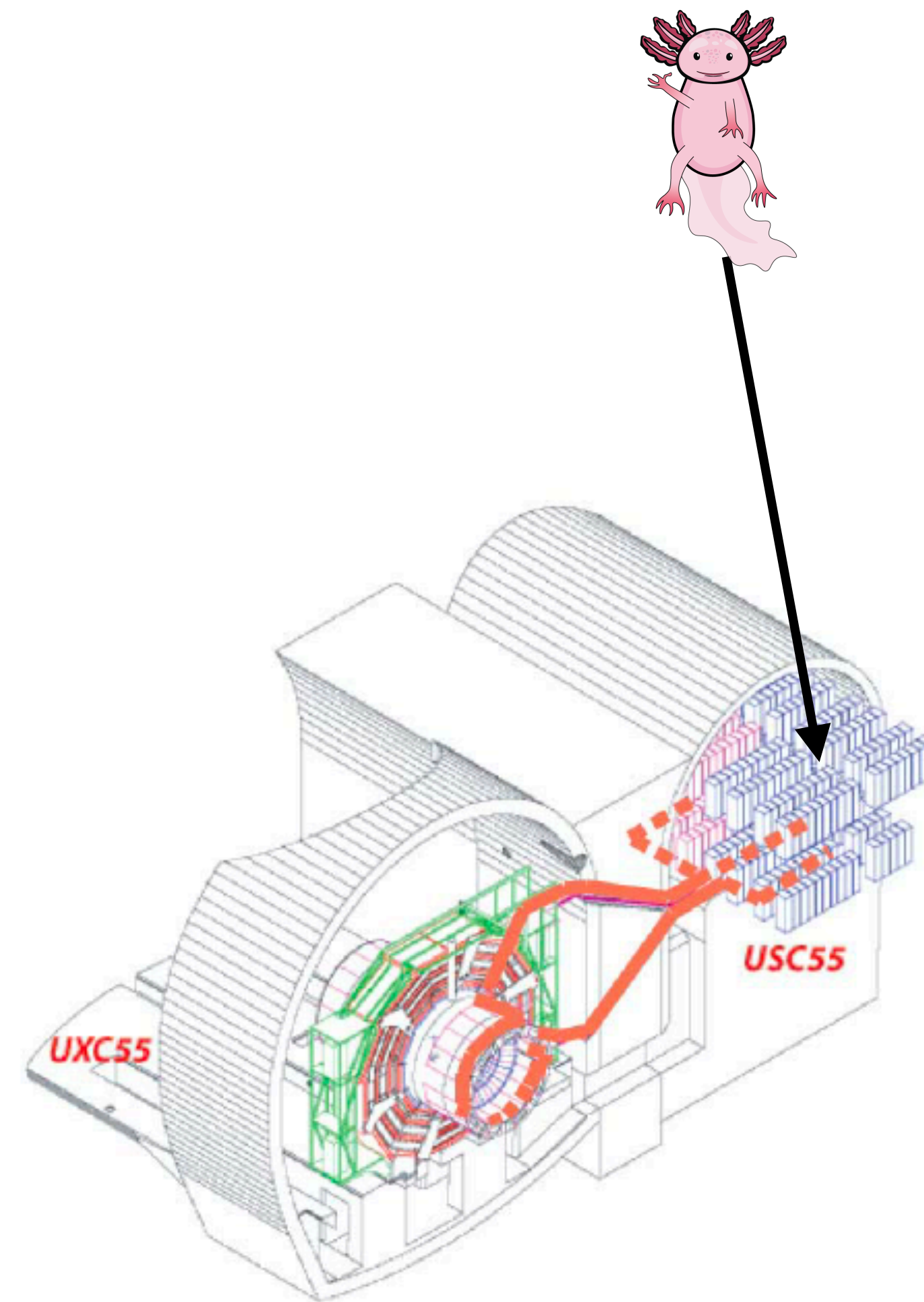
AXOLITL

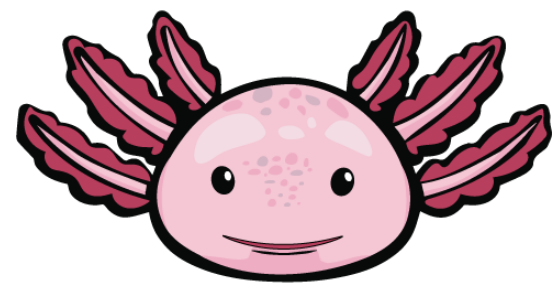
Anomaly Detection in the hardware trigger

Input from available trigger objects quantities:

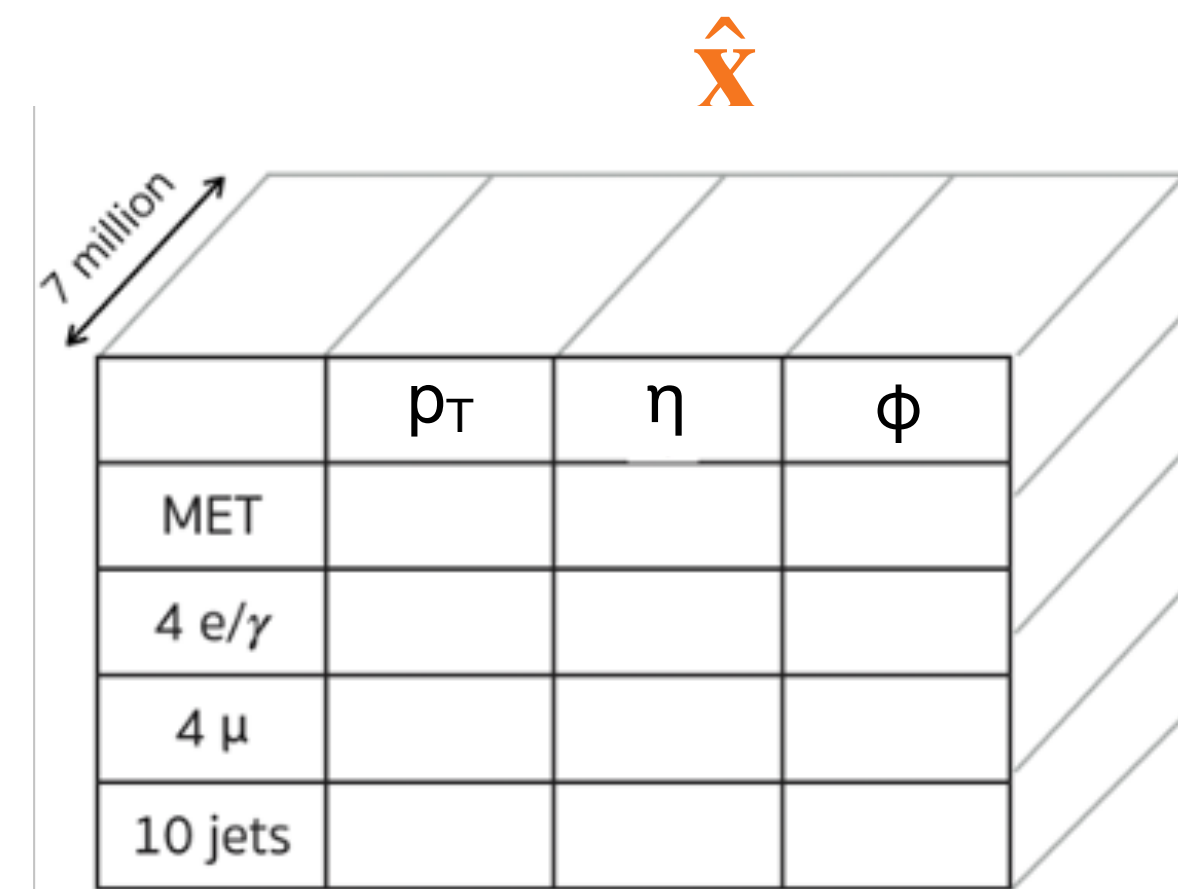
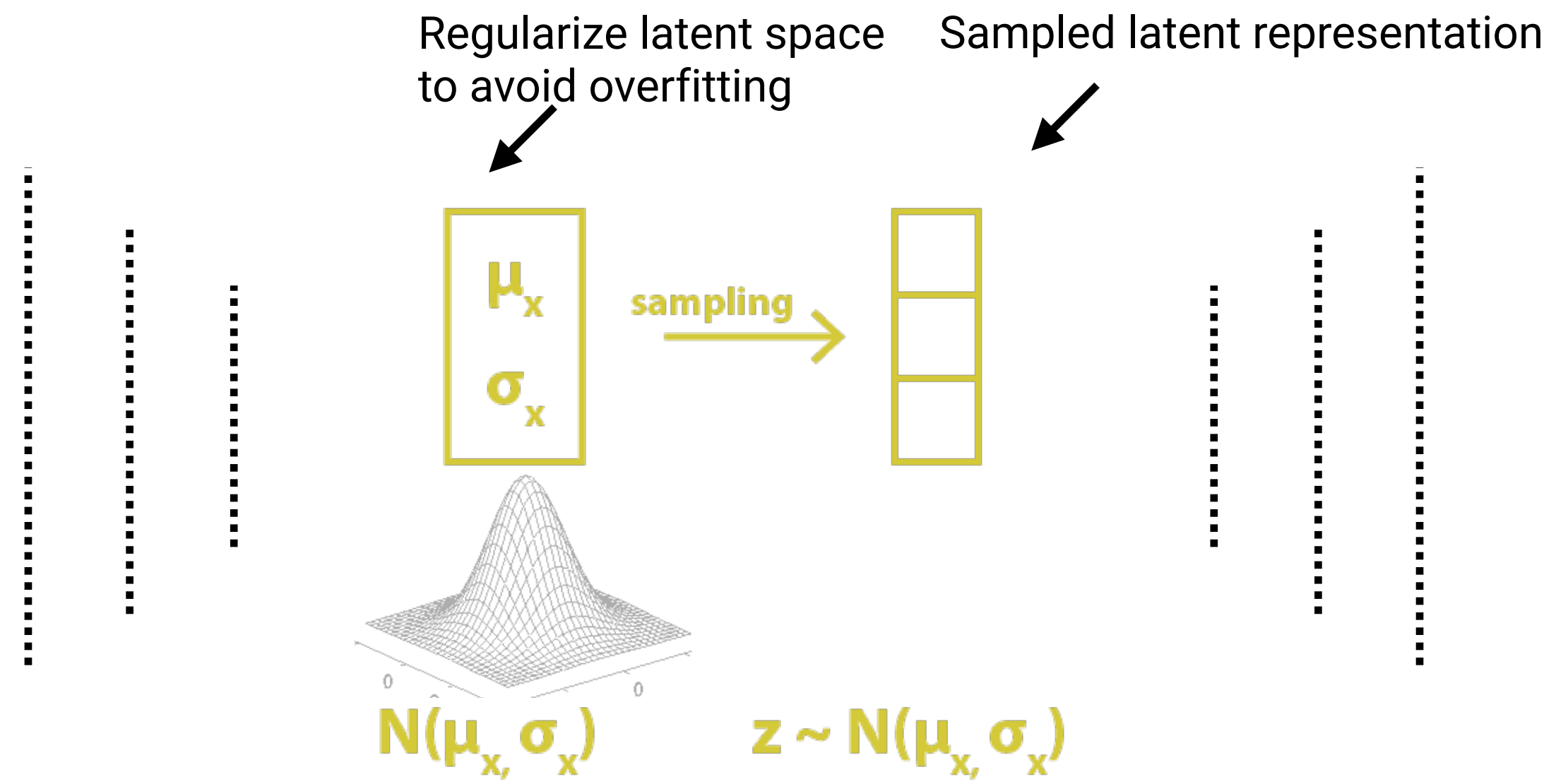
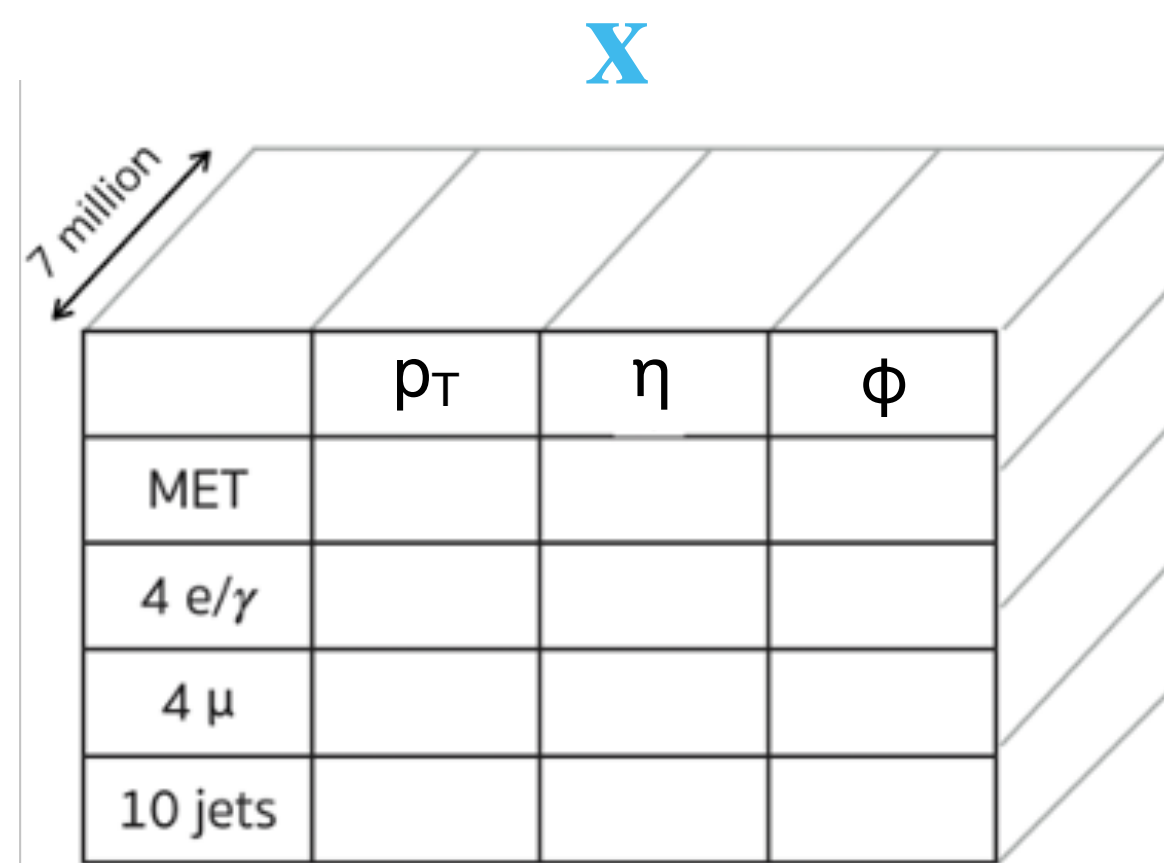
X

7 million				
	ρ_T	η	ϕ	
	MET			
	4 e/γ			
	4 μ			
	10 jets			

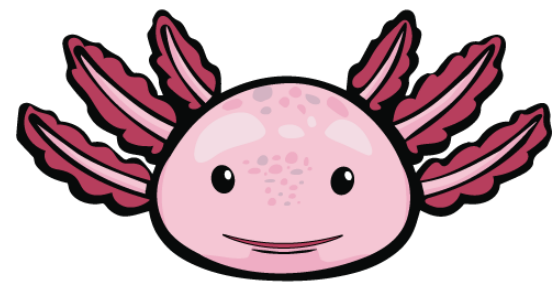




AXOLITL

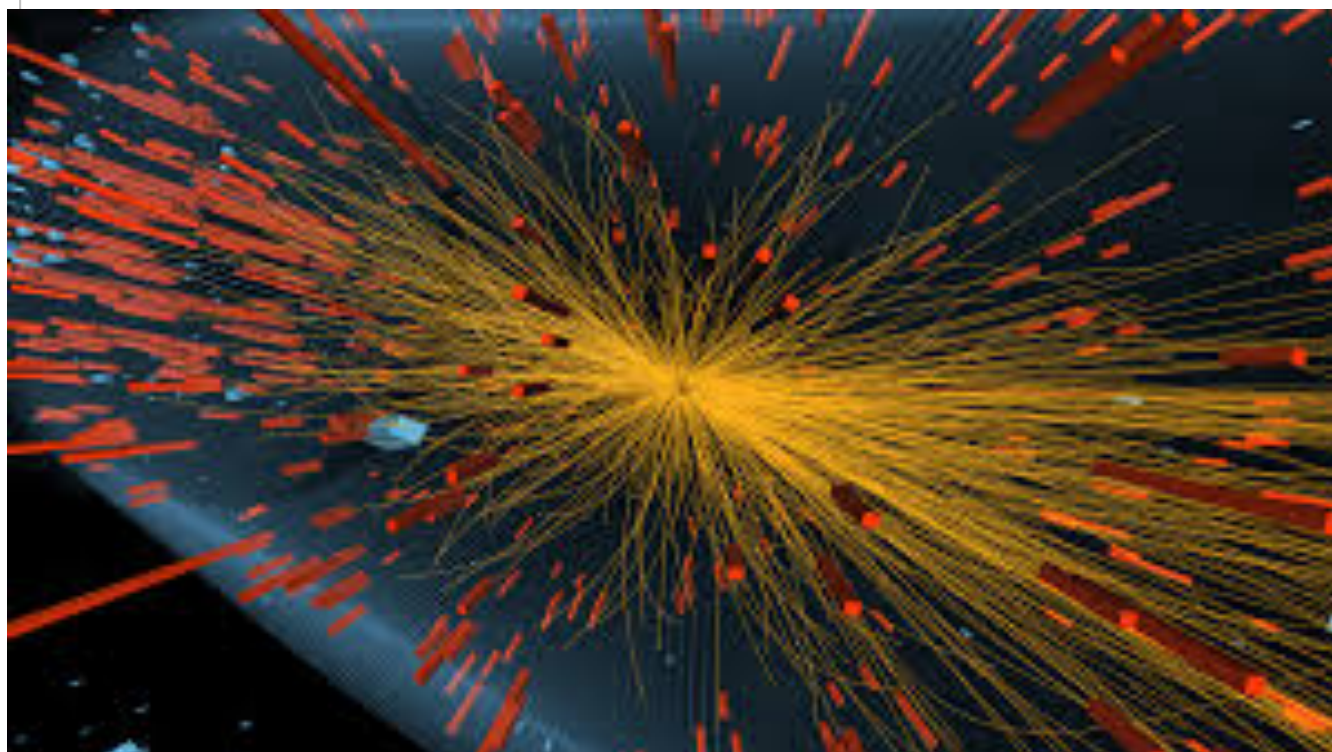


$$\text{loss} = \| \mathbf{x} - \hat{\mathbf{x}} \|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$



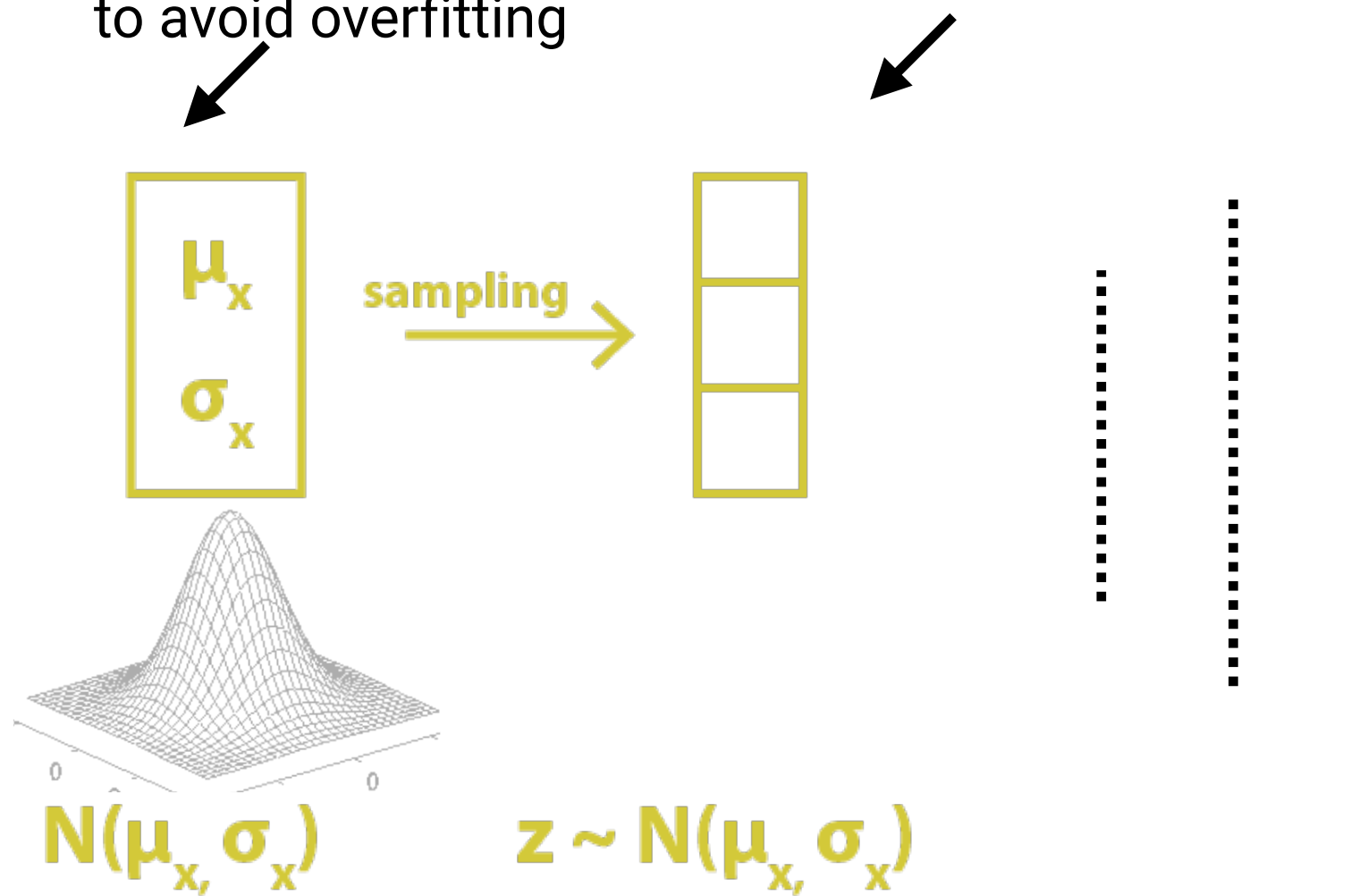
AXOLITL

x

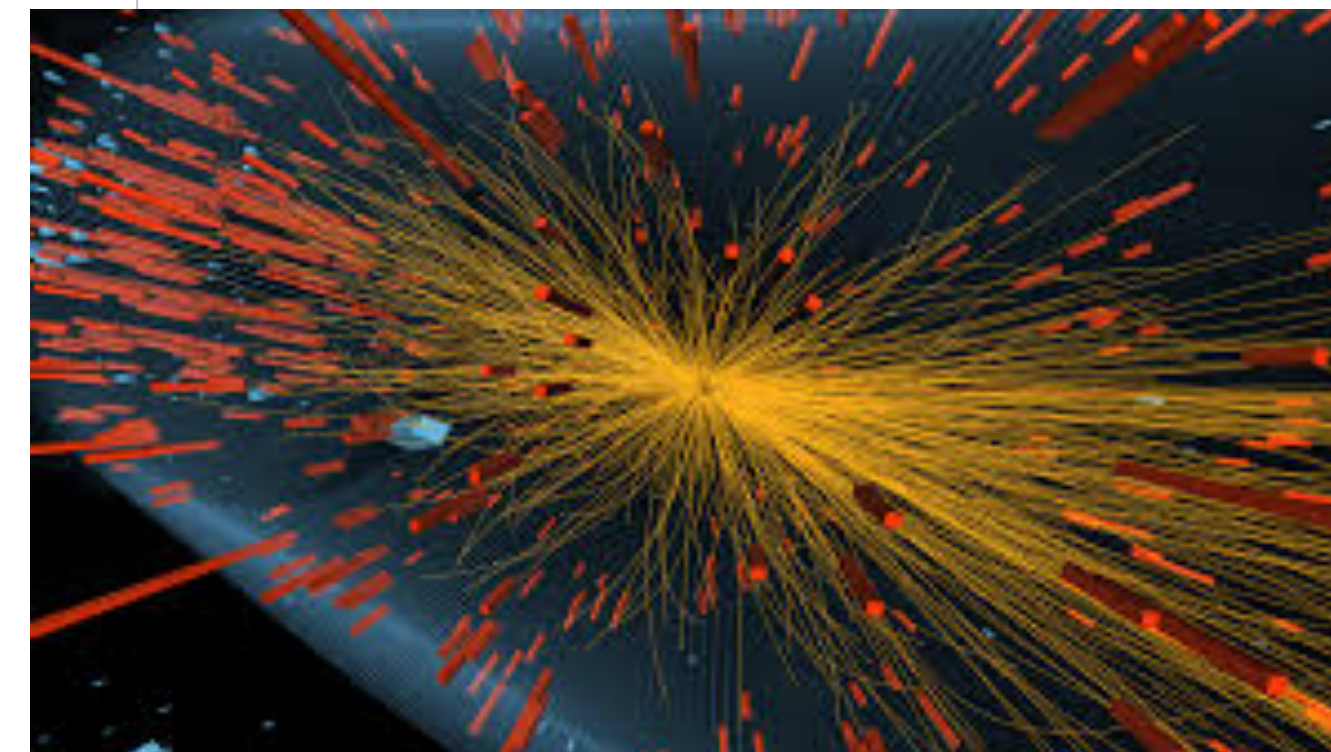


Regularize latent space to avoid overfitting

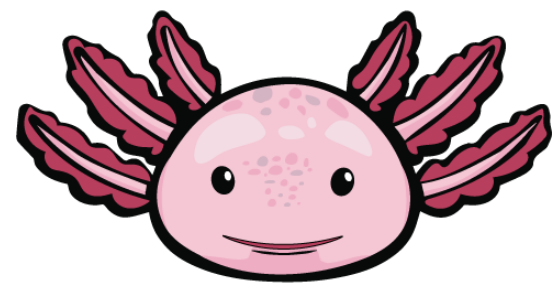
Sampled latent representation



\hat{x}

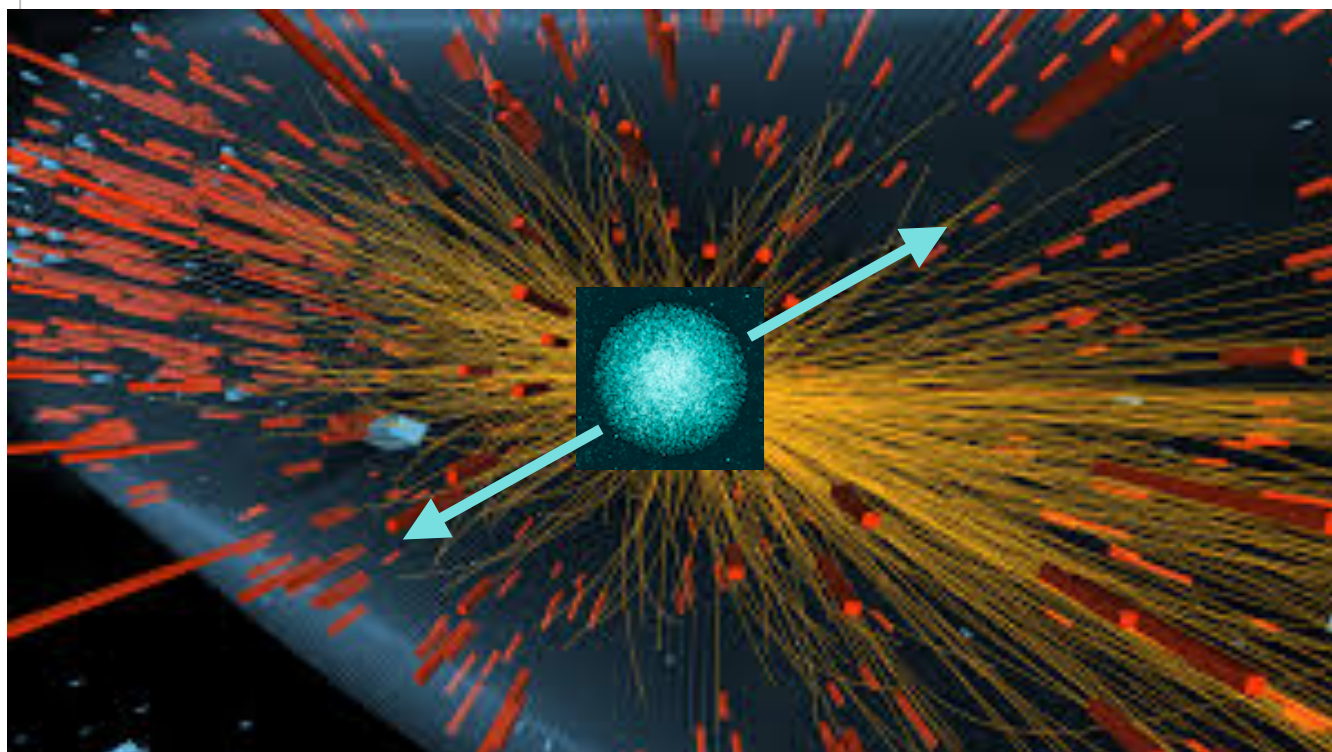


$$\text{loss} = \|x - \hat{x}\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$



AXOLITL

x



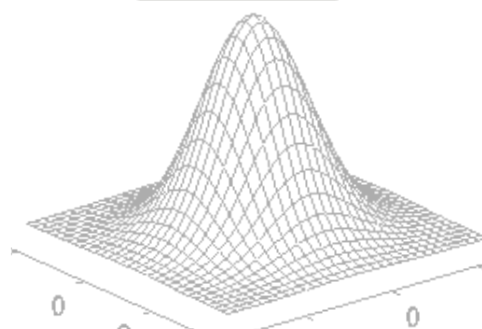
Regularize latent space to avoid overfitting

Sampled latent representation

μ_x
 σ_x

sampling

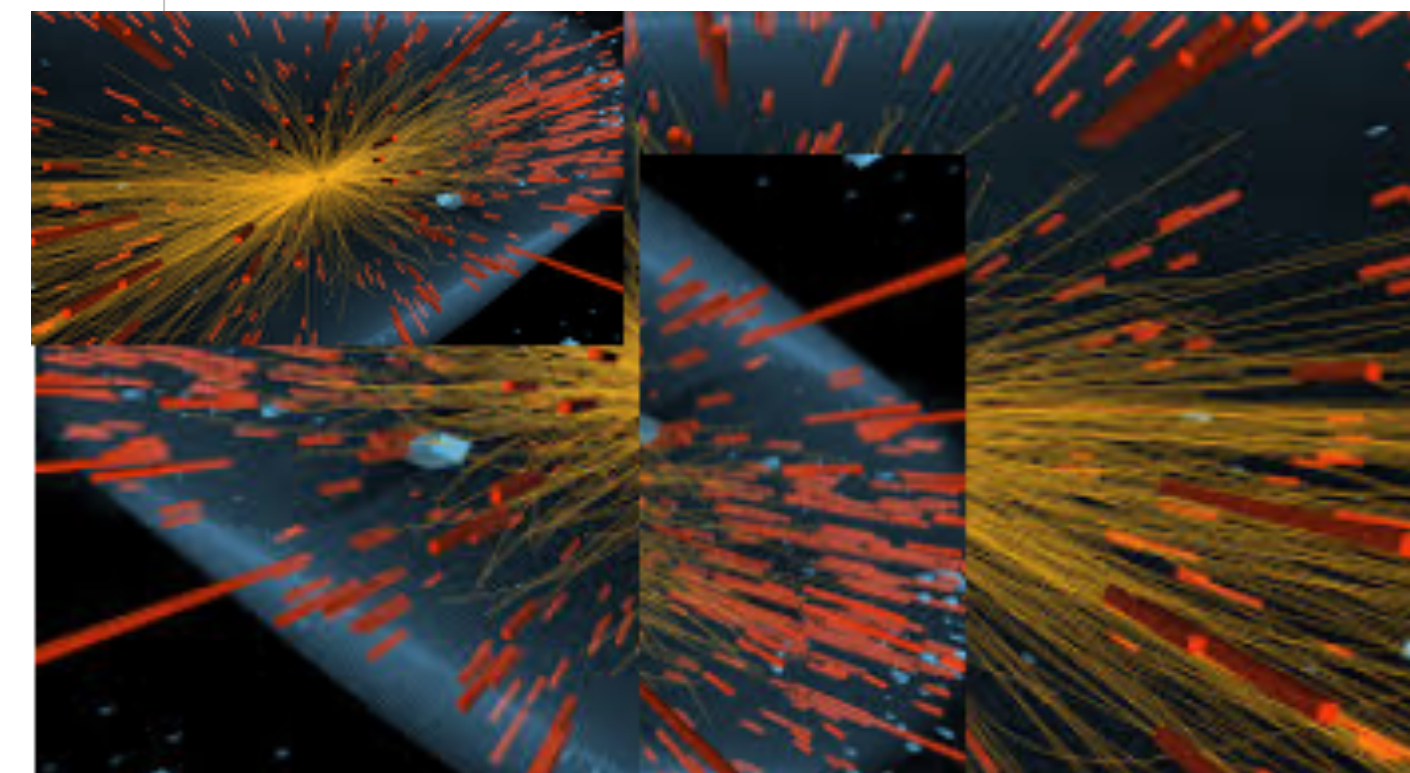
Sampled latent representation



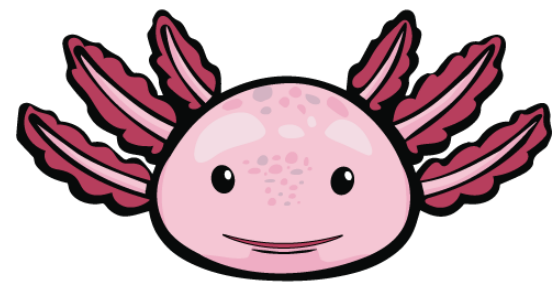
$N(\mu_x, \sigma_x)$

$z \sim N(\mu_x, \sigma_x)$

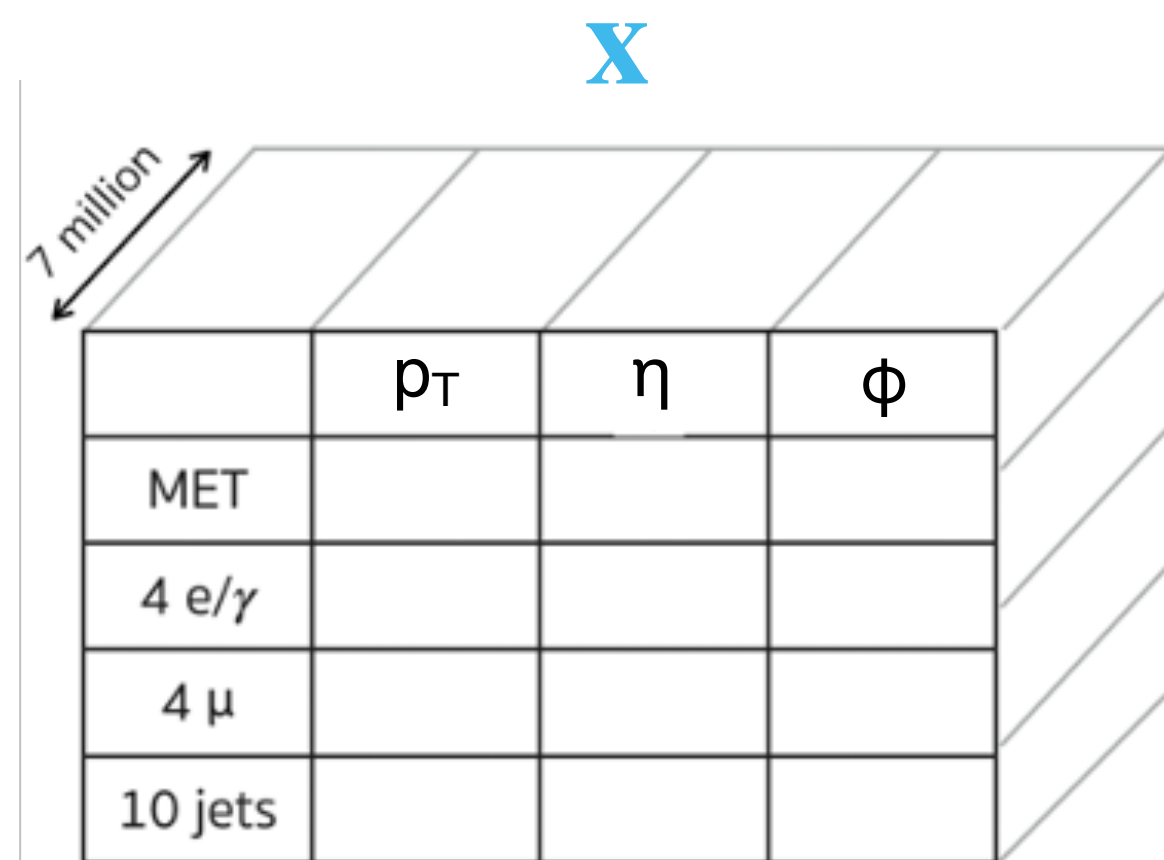
\hat{x}



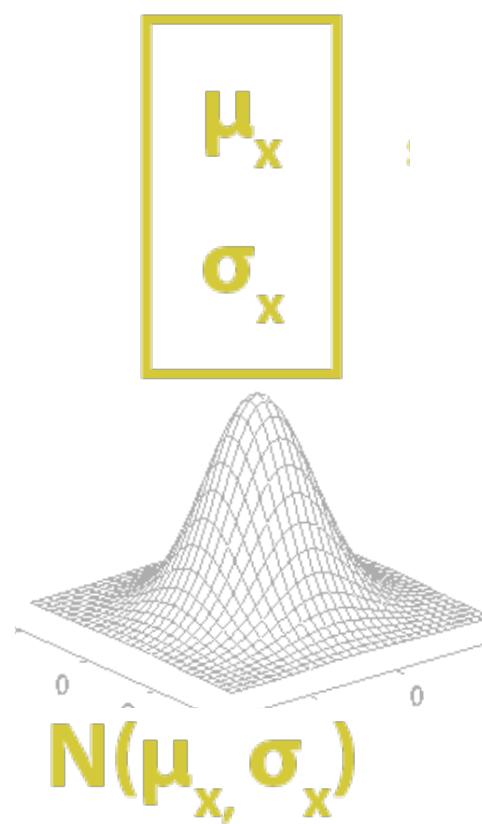
$$\text{loss} = \|x - \hat{x}\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$



AXOLITL



Anomaly score!



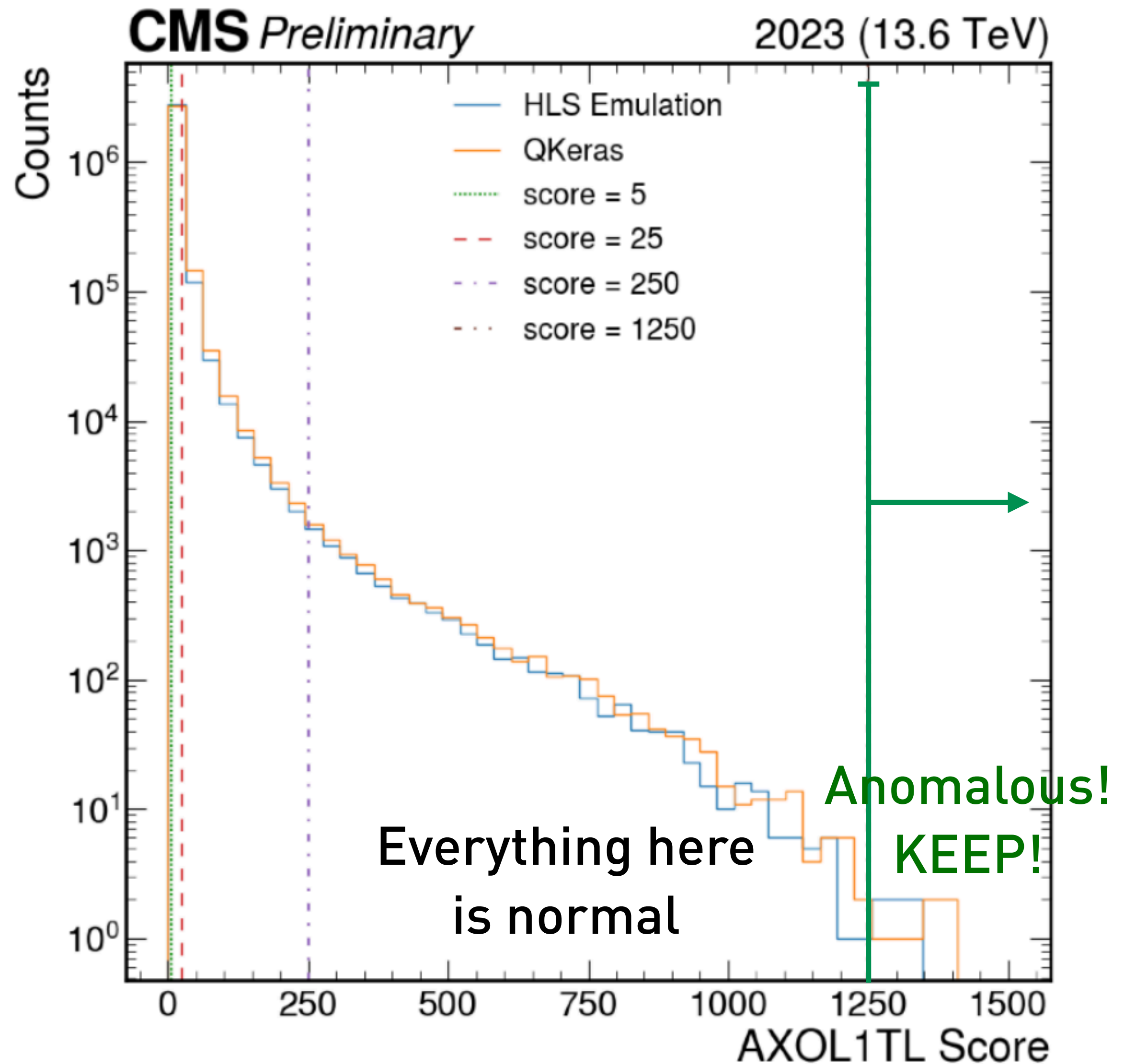
$$\text{loss} = \|x - \hat{x}\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$

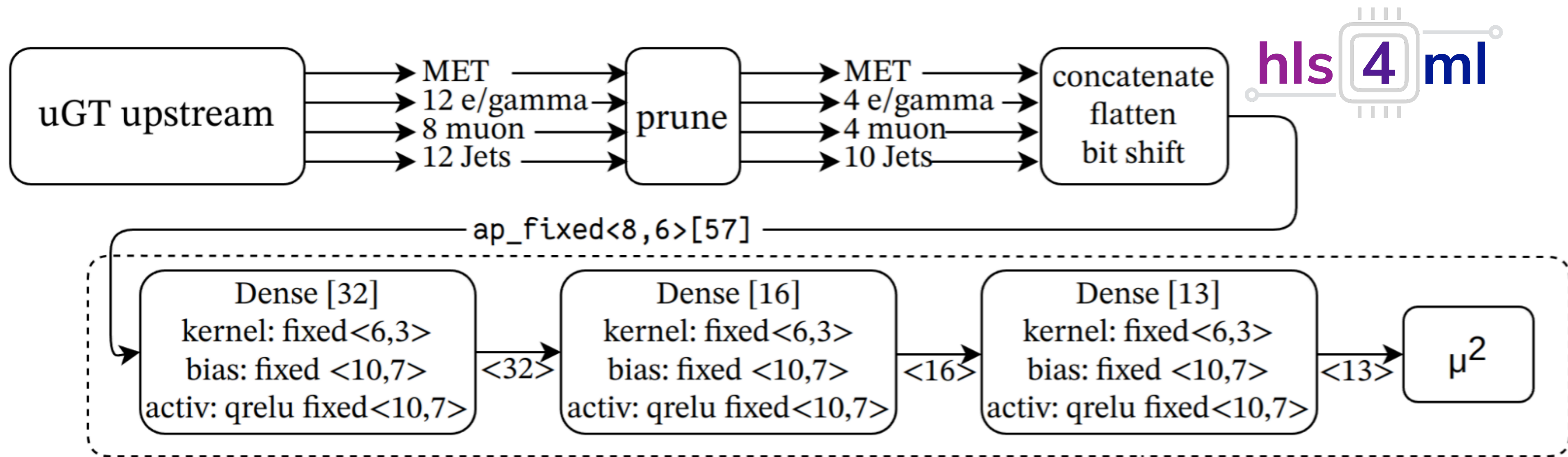
Anomaly score

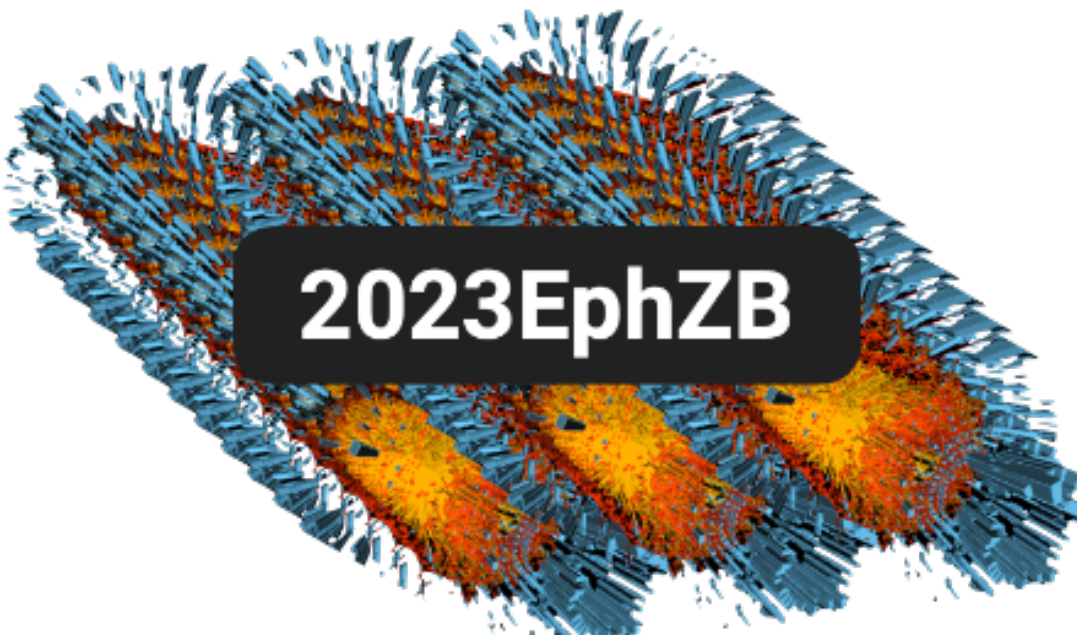
After training, network is evaluated on orthogonal test set (offline)

- Find cut on anomaly score resulting in pre-defined background rate

This model, and thresholds evaluated on test set deployed in the trigger!

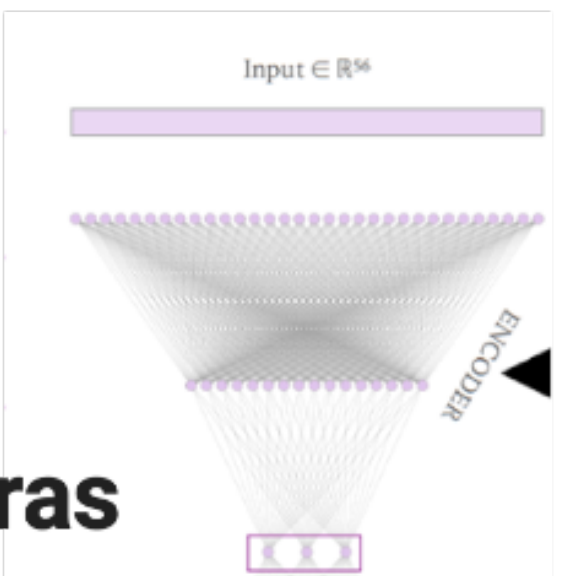






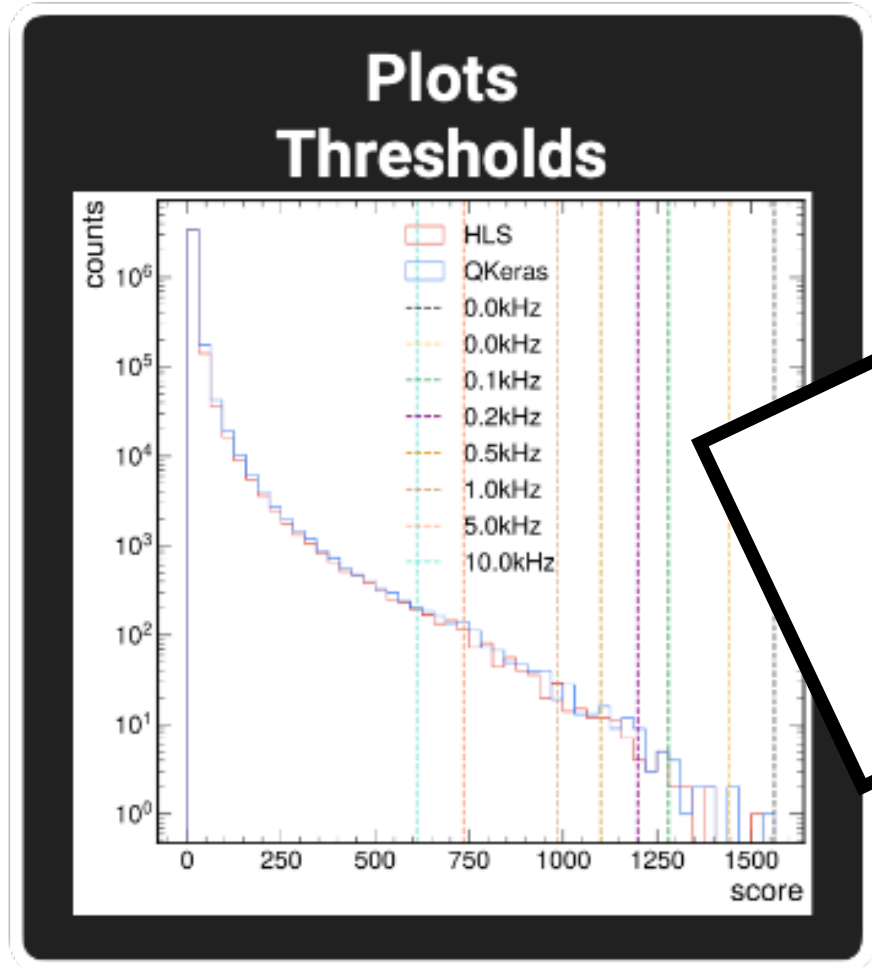
2023EphZB

Train QKeras

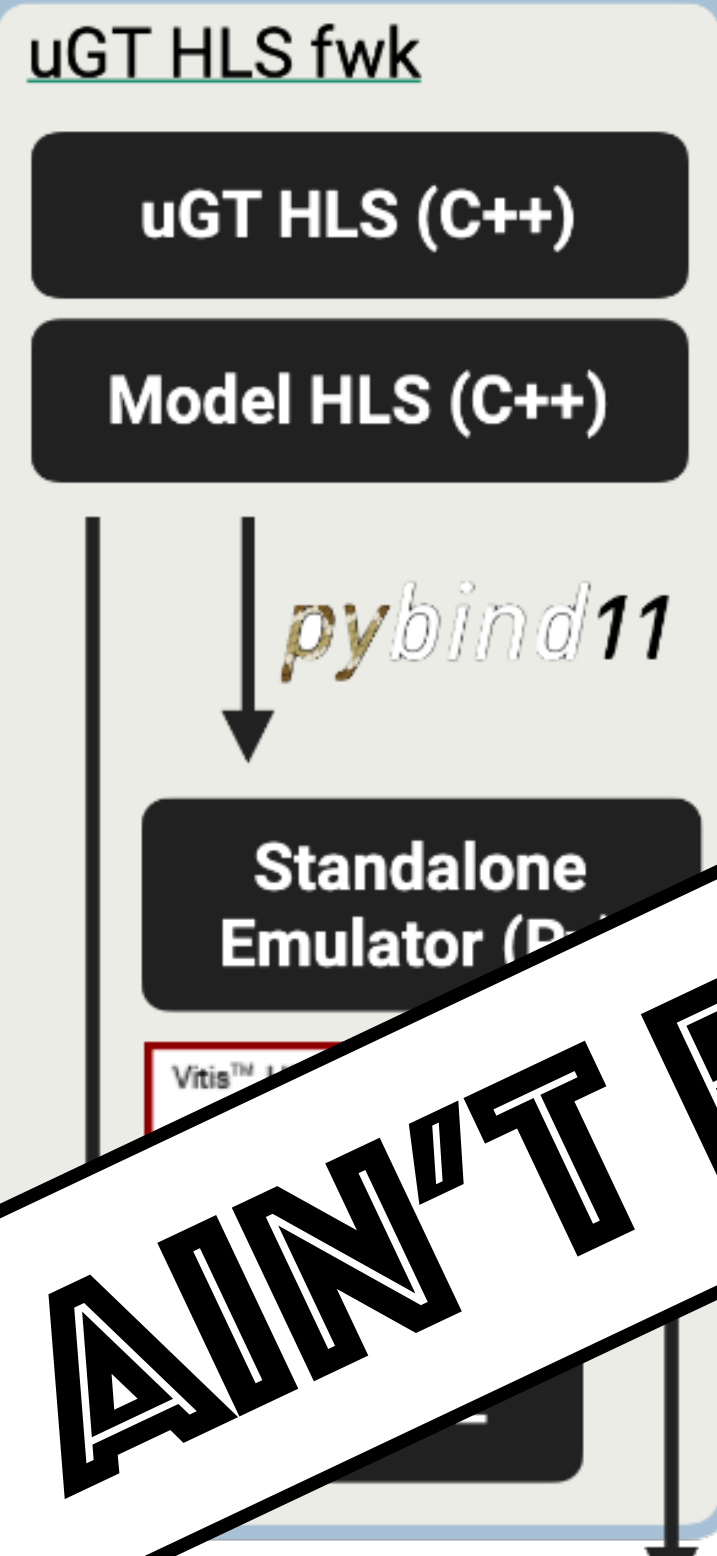


Model .h5

Evaluate QKeras

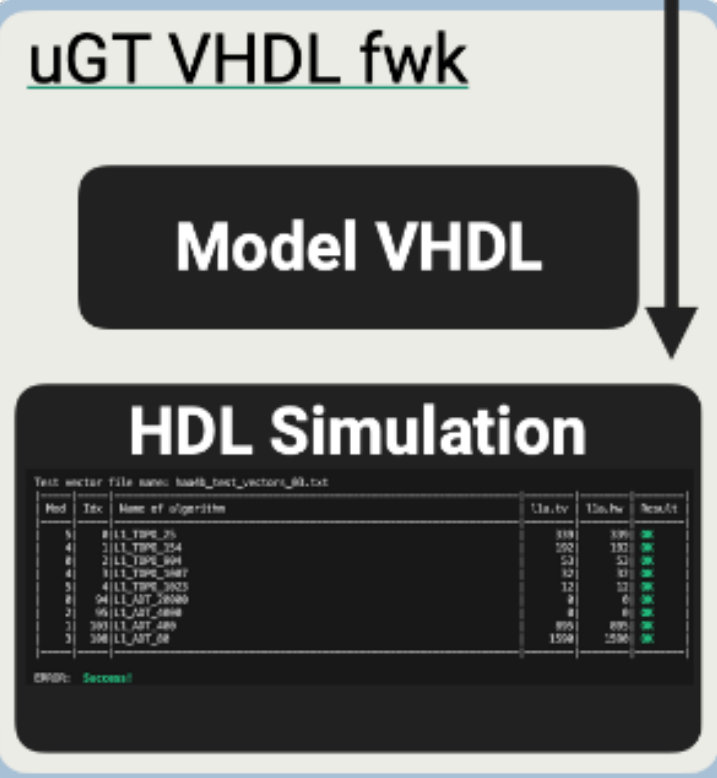


emu xml (example)



Standalone Emulator (D)

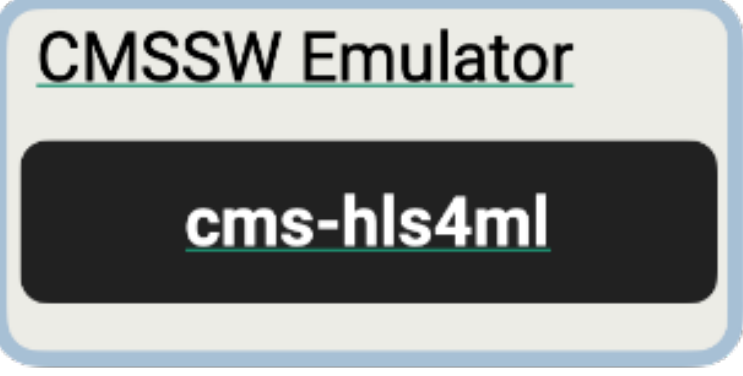
Test Vectors



Model VHDL

HDL Simulation

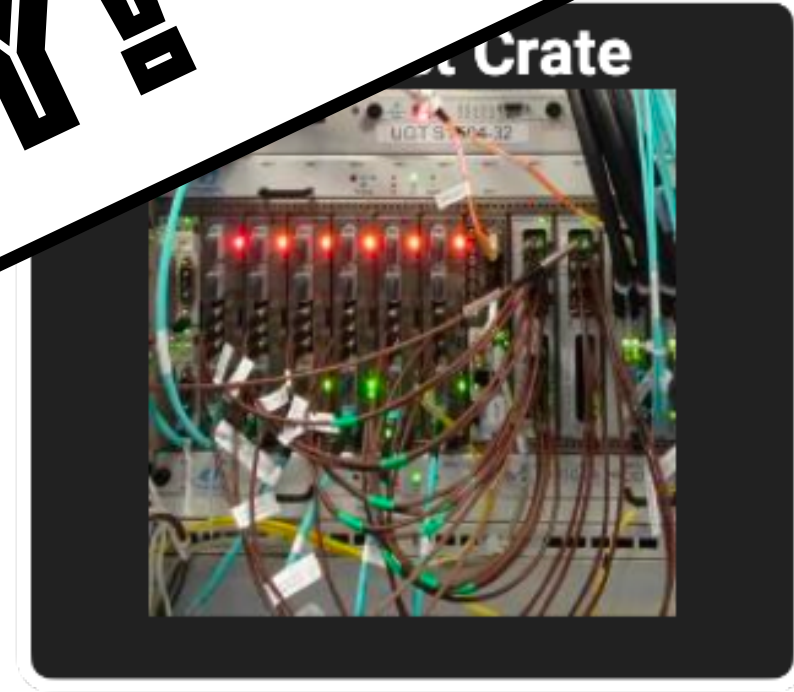
Table with columns: Test vector file name, Name of algorithm, Status, Time, Results



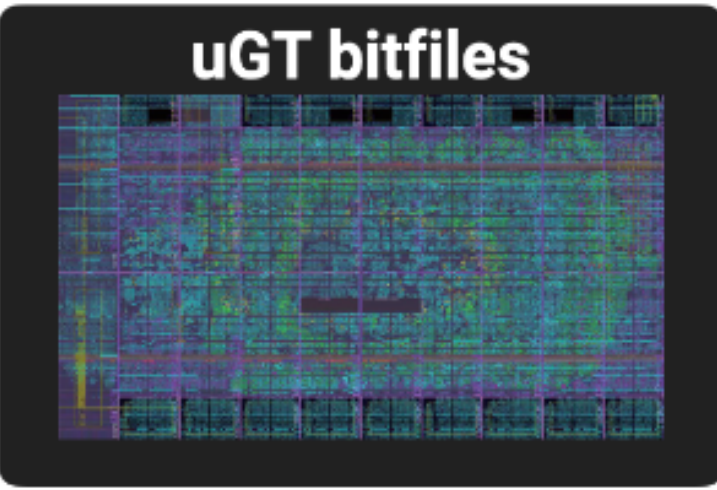
cms-hls4ml



L1 NTuples (w/TC)

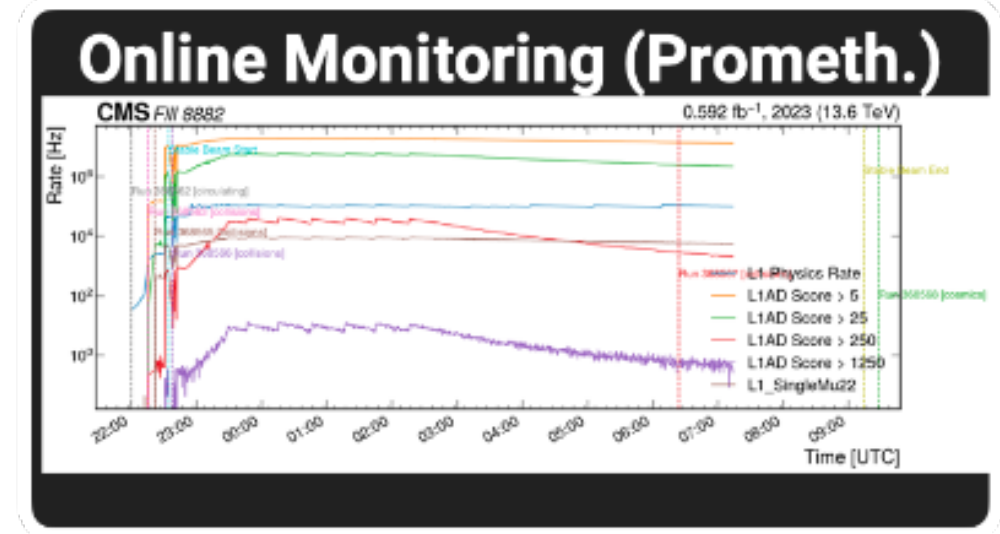
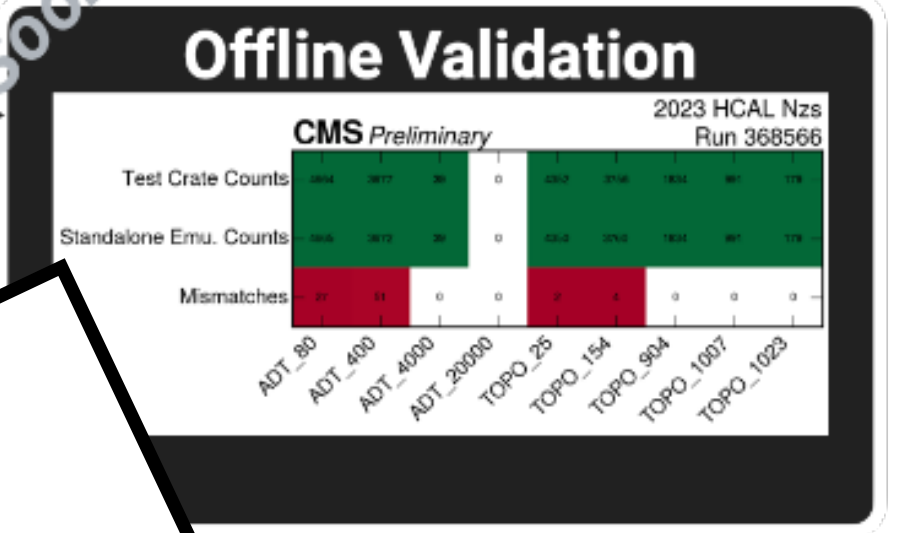


L1 keys

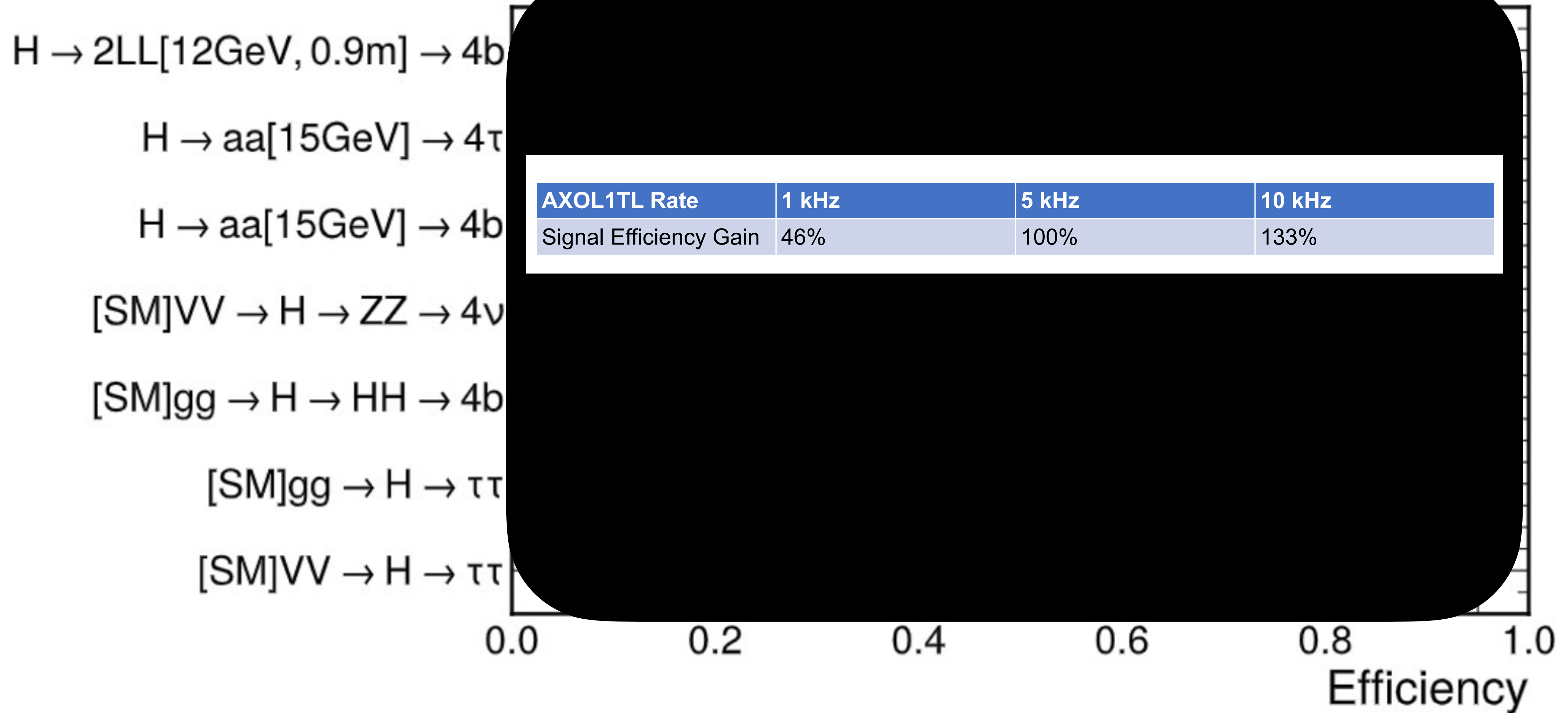


uGT bitfiles

Coming Soon



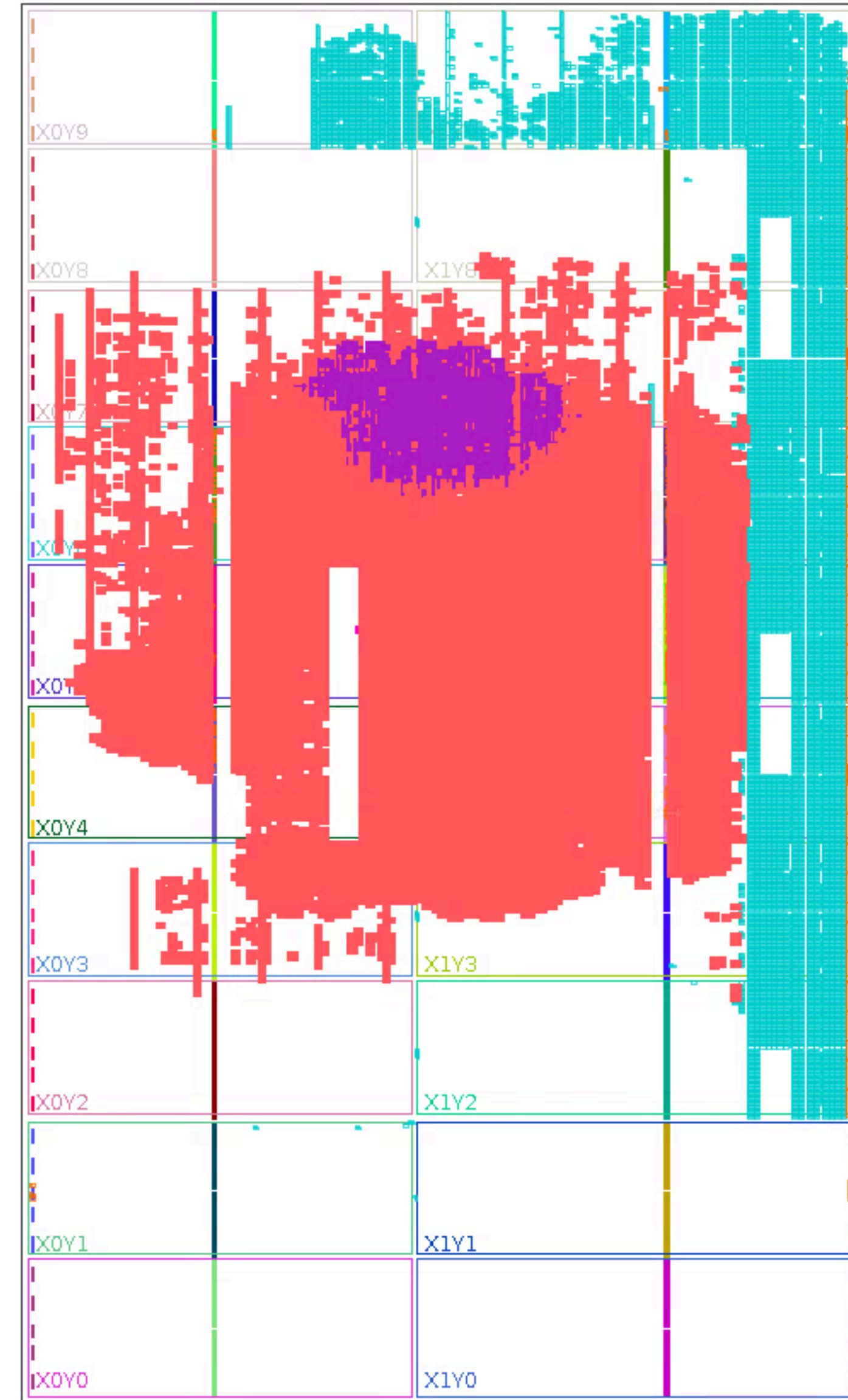
IT AIN'T EASY!



For various BSM and SM signals, significantly increase signal efficiency adding AD!

	Latency	LUTs	FFs	DSPs	BRAMs
AXOLITL	2 ticks 50 ns	2.1%	~0	0	0

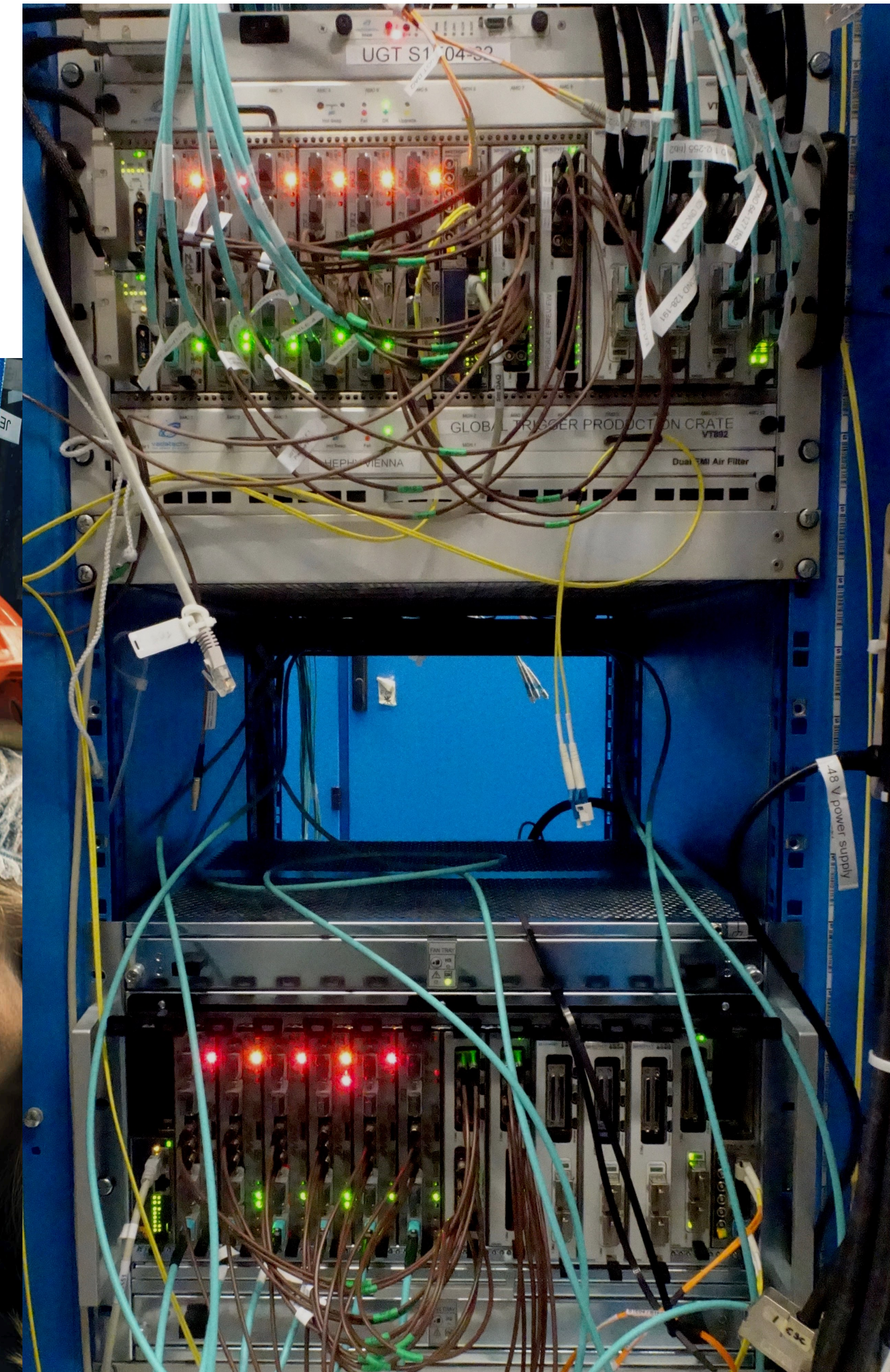
50 nanosecond anomaly detection!



uGT test crate

CMS Global Trigger test crate:

- Copy of main GT system, receiving the same input data, but not used to trigger CMS
- Excellent test bench for future ML algorithms targeting L1T FPGAs
- AXOL1TL integrated since late 2023

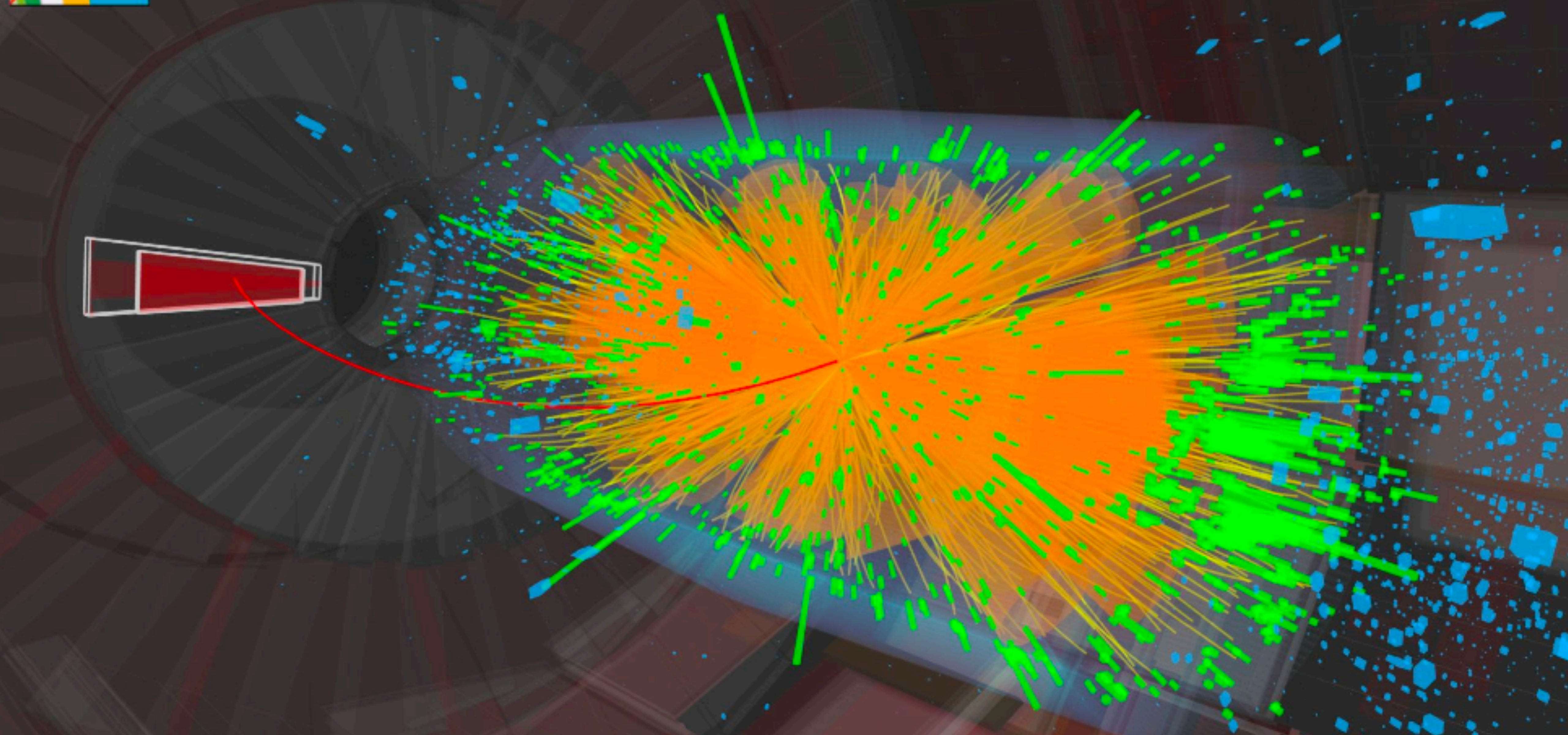




CMS Experiment at the LHC, CERN

Data recorded: 2023-May-24 01:42:17.826112 GMT

Run / Event / LS: 367883 / 374187302 / 159

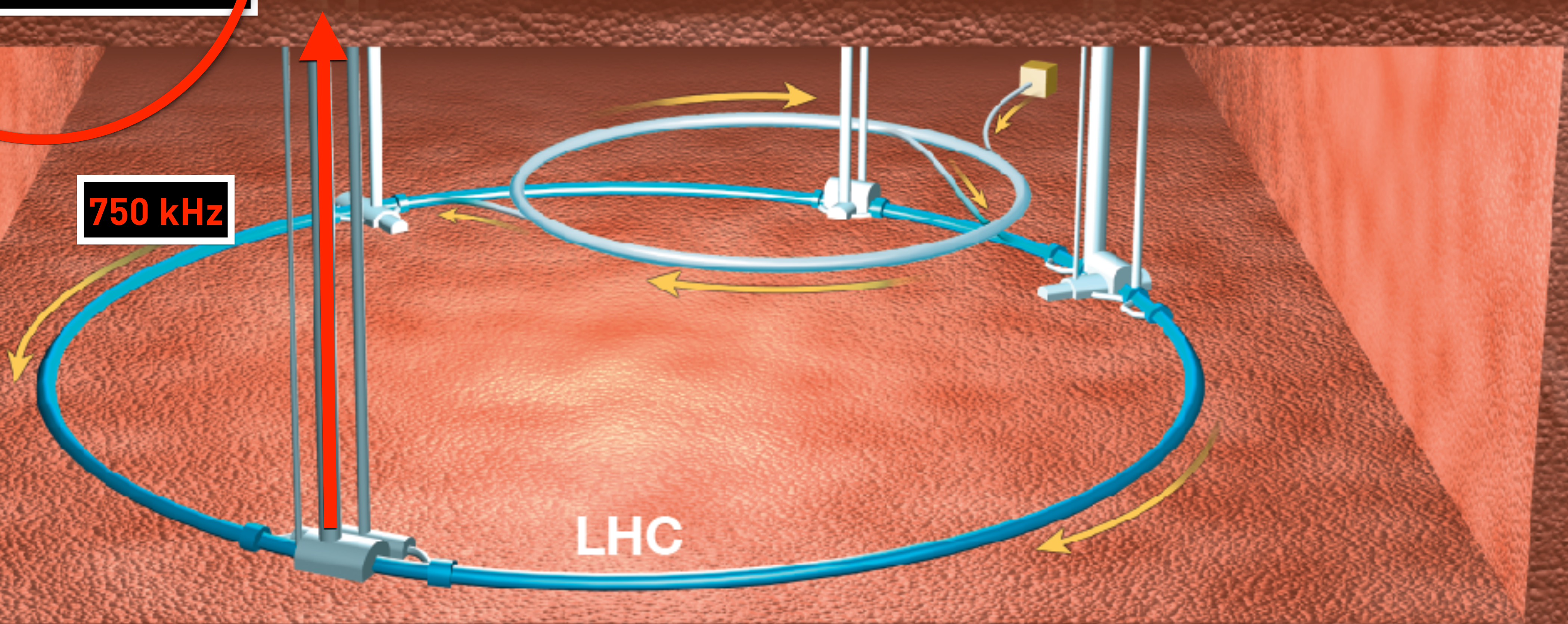




**High Level Trigger:
25'600 CPUs / 400 GPUs
Latency: 3-400 ms**



7.5 kHz



750 kHz

LHC

Events from L1 @ 750 kHz



30k cores, single-threaded

CPU tasks

Raw data

Pixel tracks
~ 48 ms

Pixel tracks

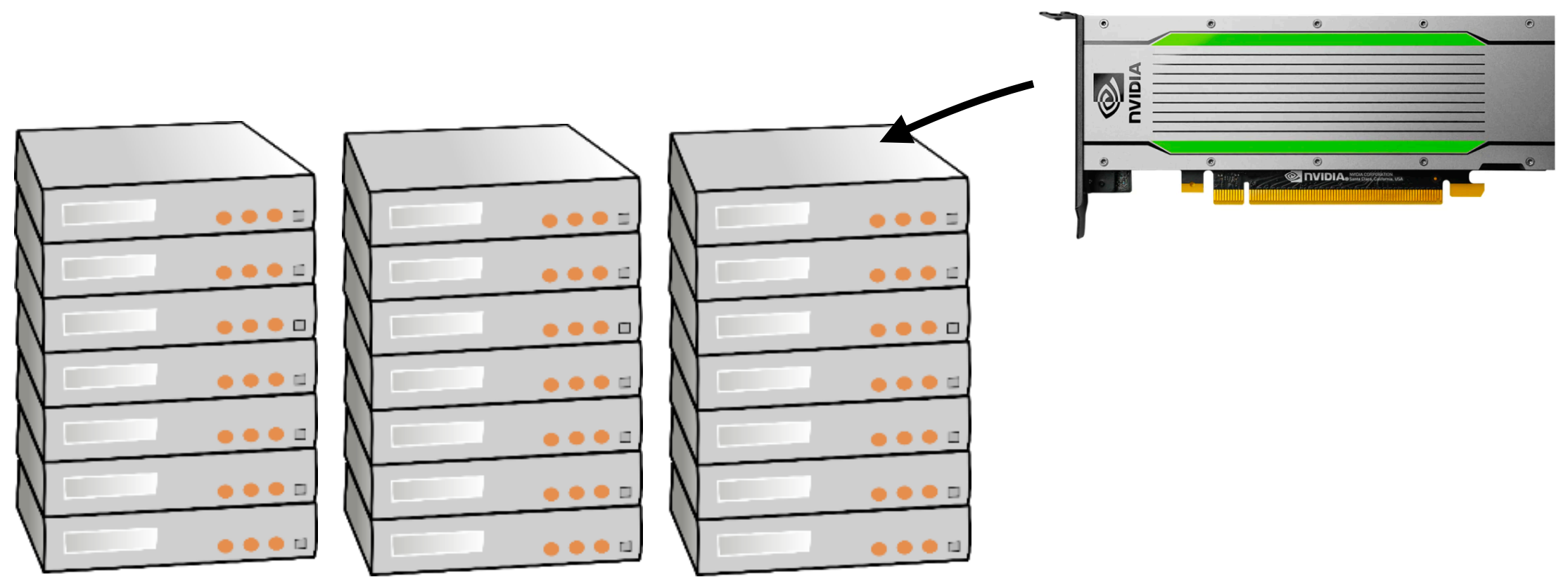
Particle Flow

300 ms

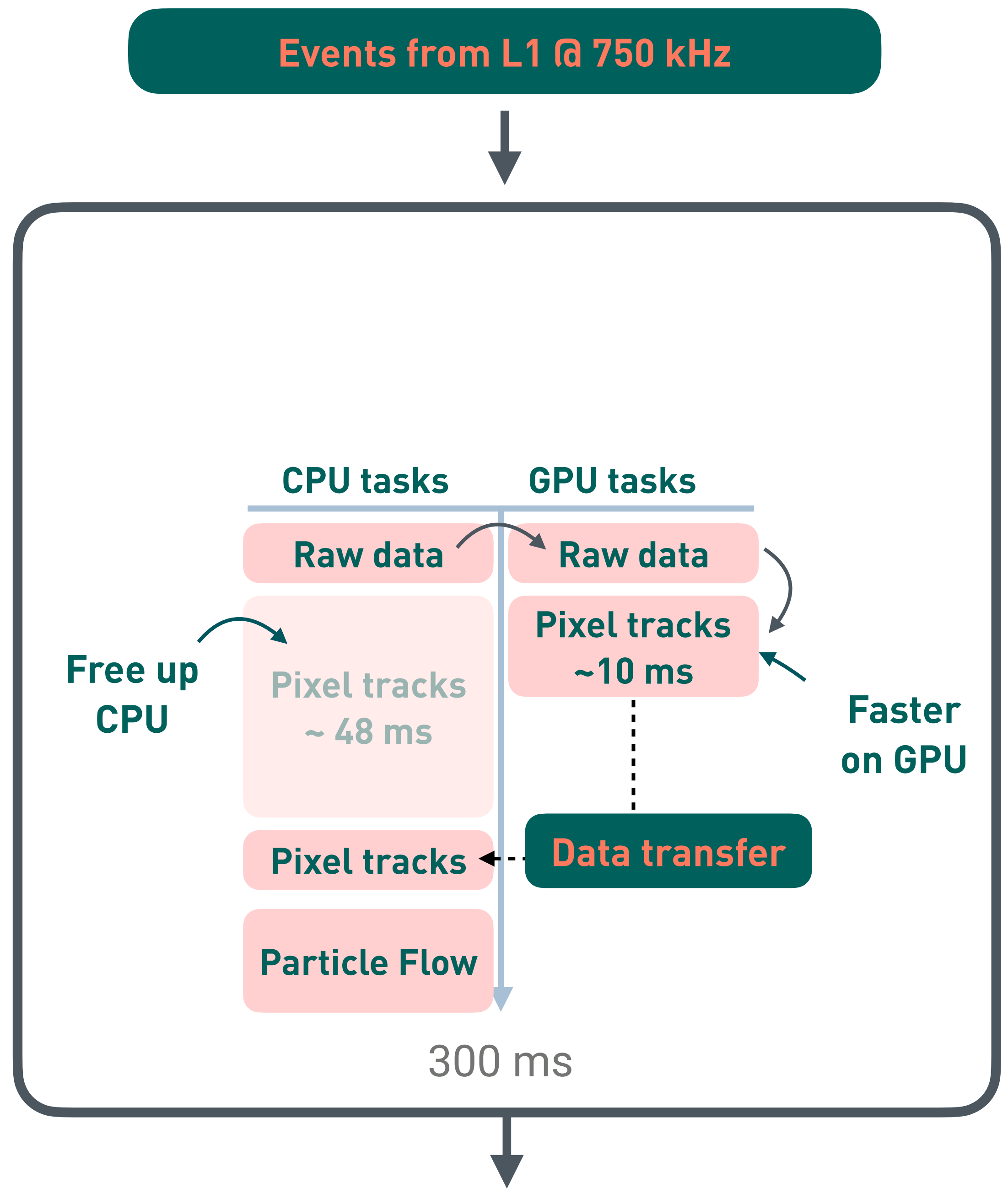


CPU node (128 cores)





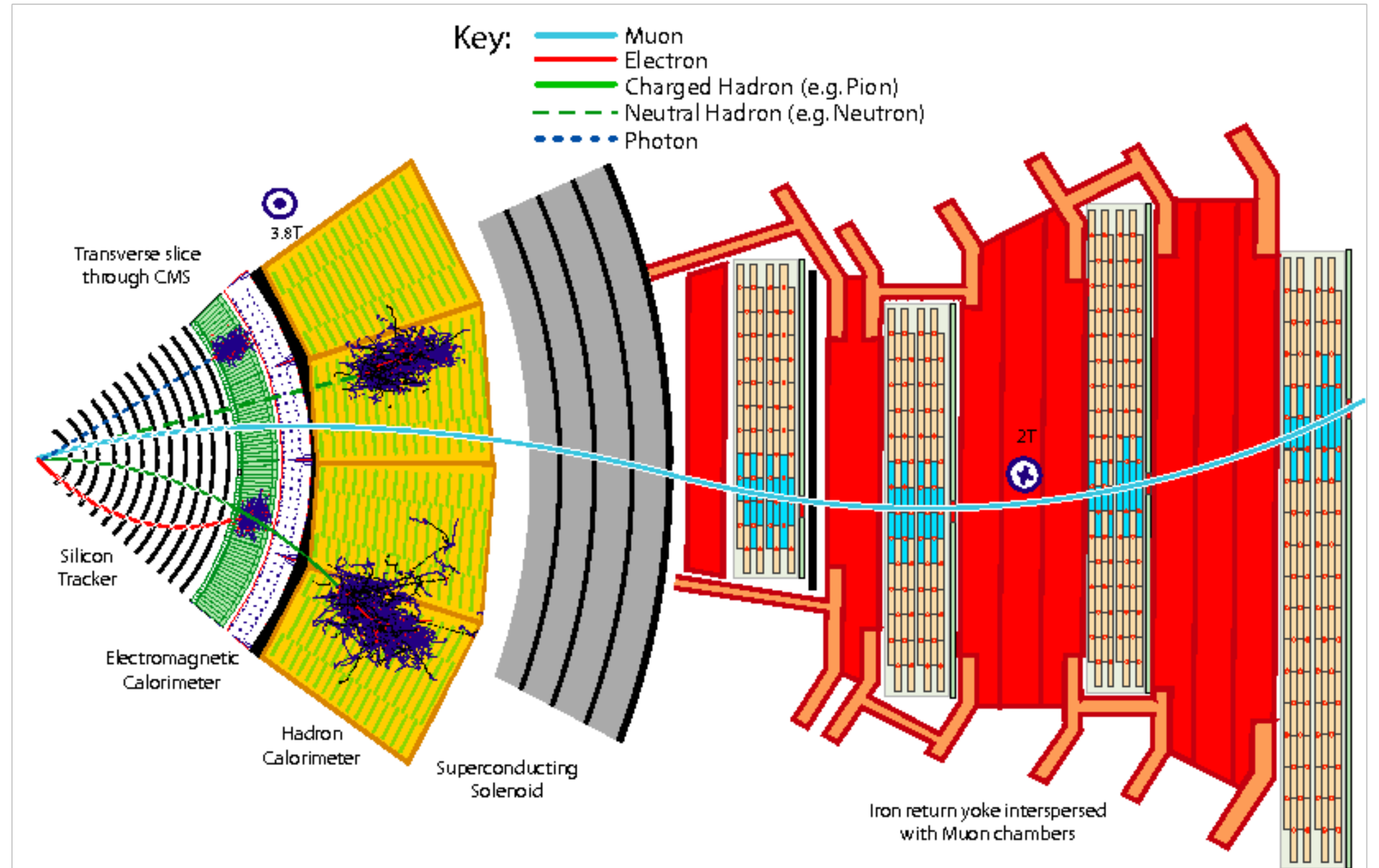
- Offload resource-intensive computations to GPU
- Increase throughput by x1.8, reduce power by 30%
 - Will offload 50-80% to computing accelerators!



ML for fast reconstruction

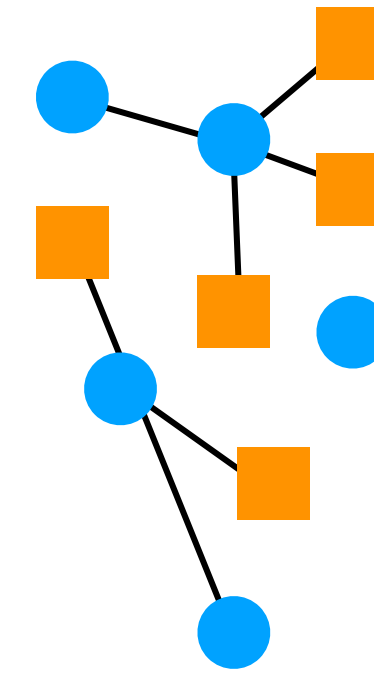
Particle Flow: Best reconstruction at HLT

- Slow, cannot run on all events (currently 17%)

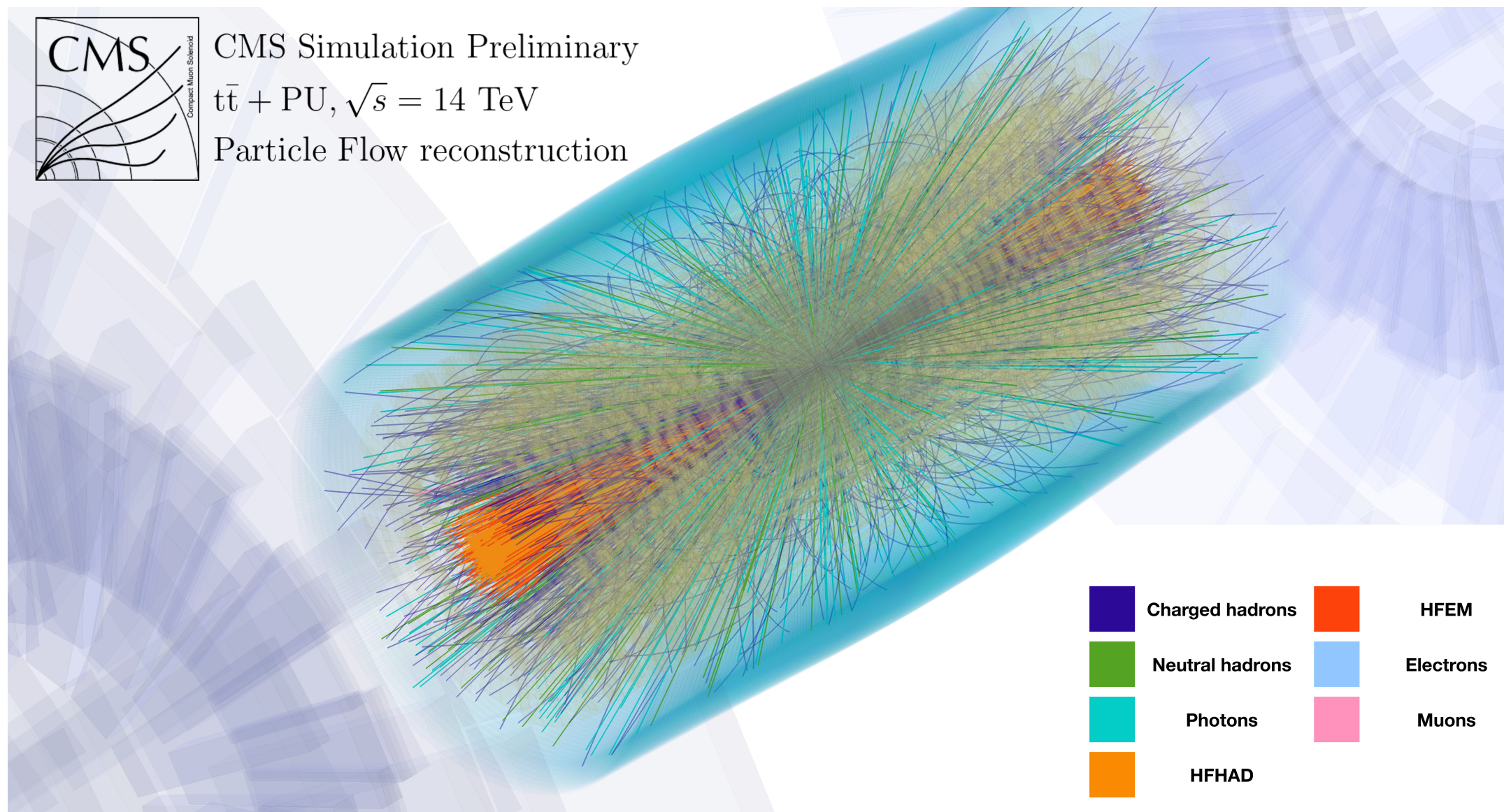


ML for fast reconstruction

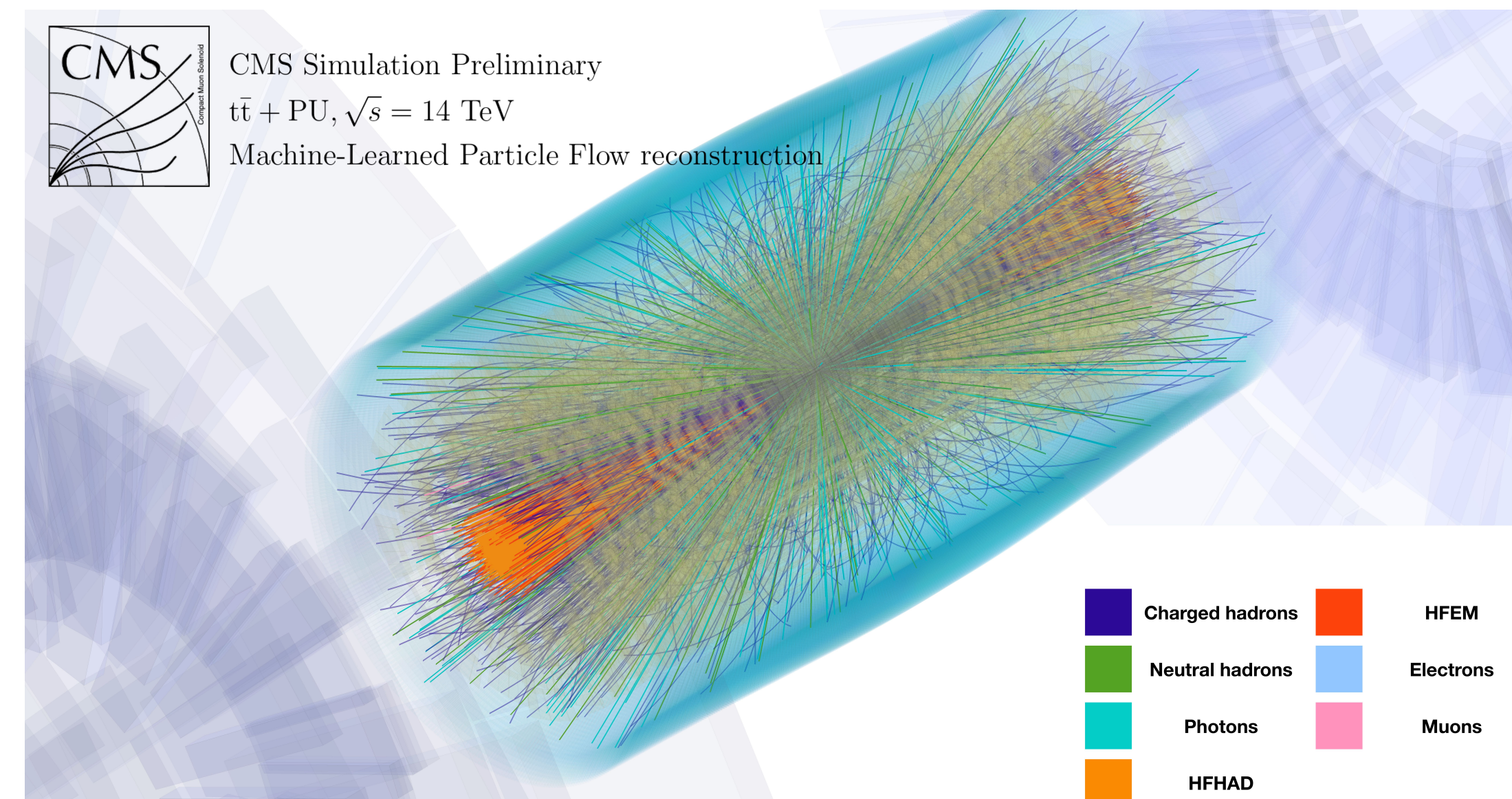
Graph Deep Neural Networks: “fast” approximations of ParticleFlow

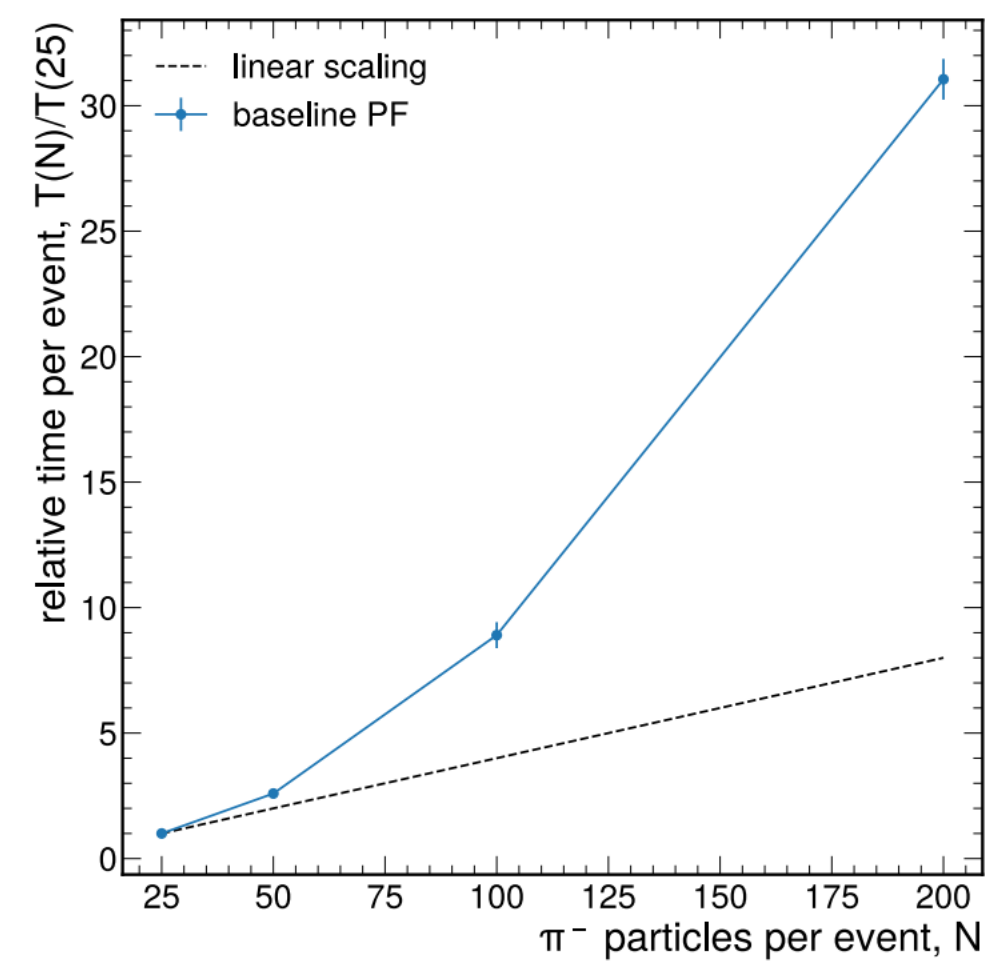


Classical Particle Flow

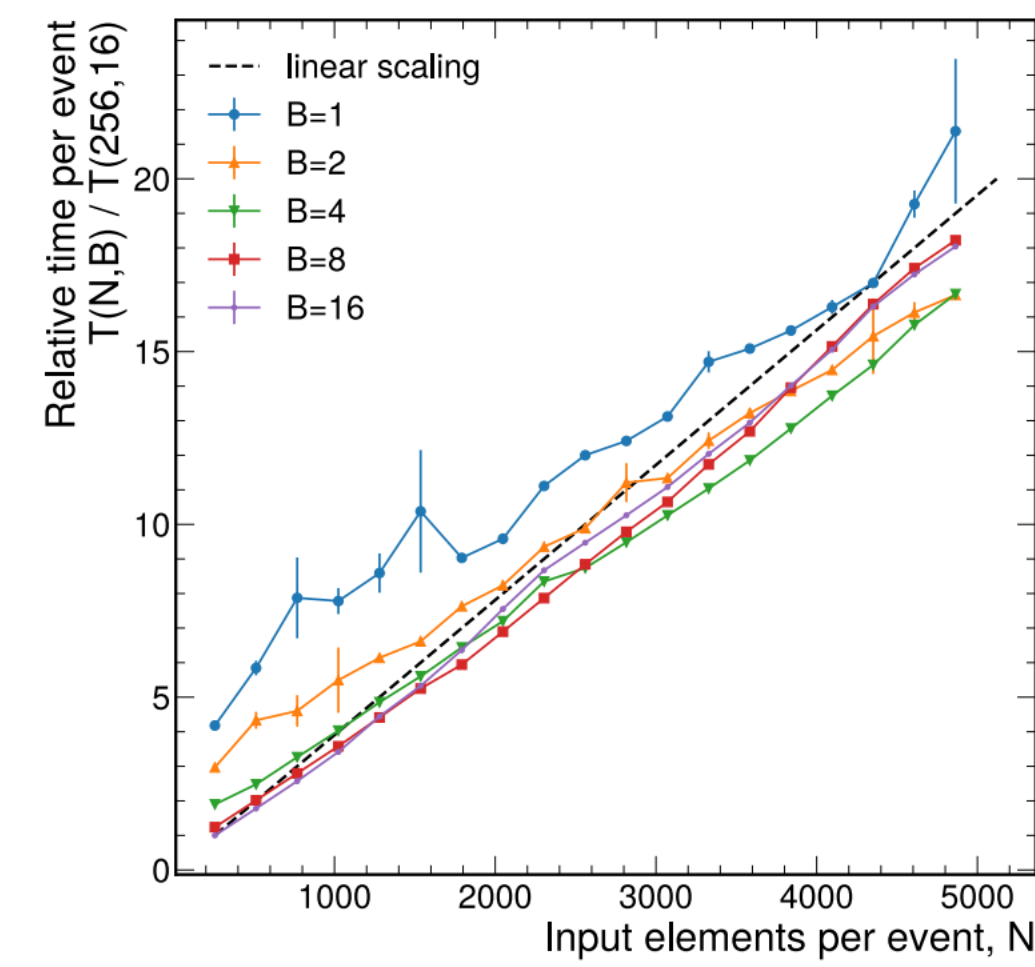


Graph Neural Network





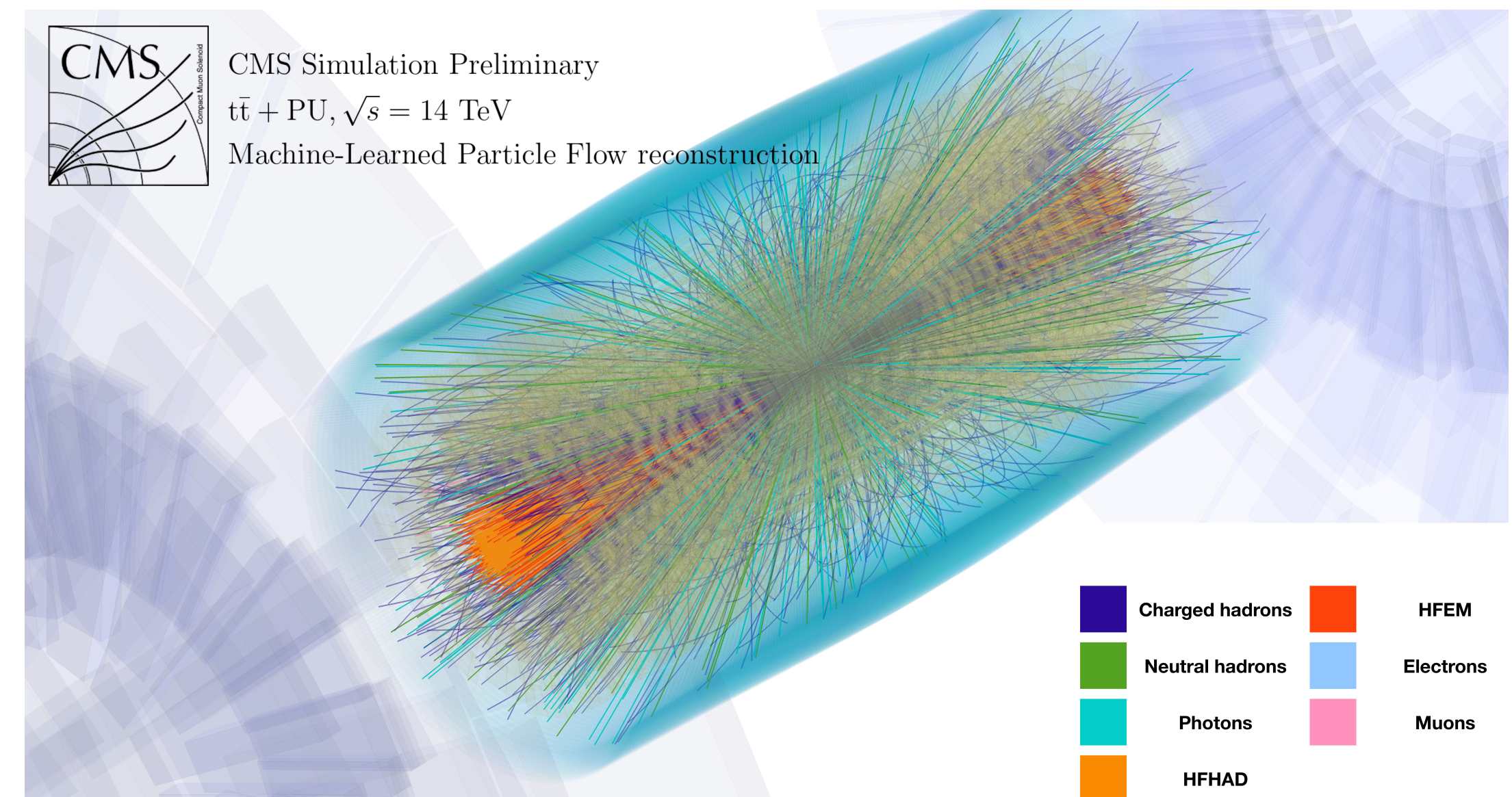
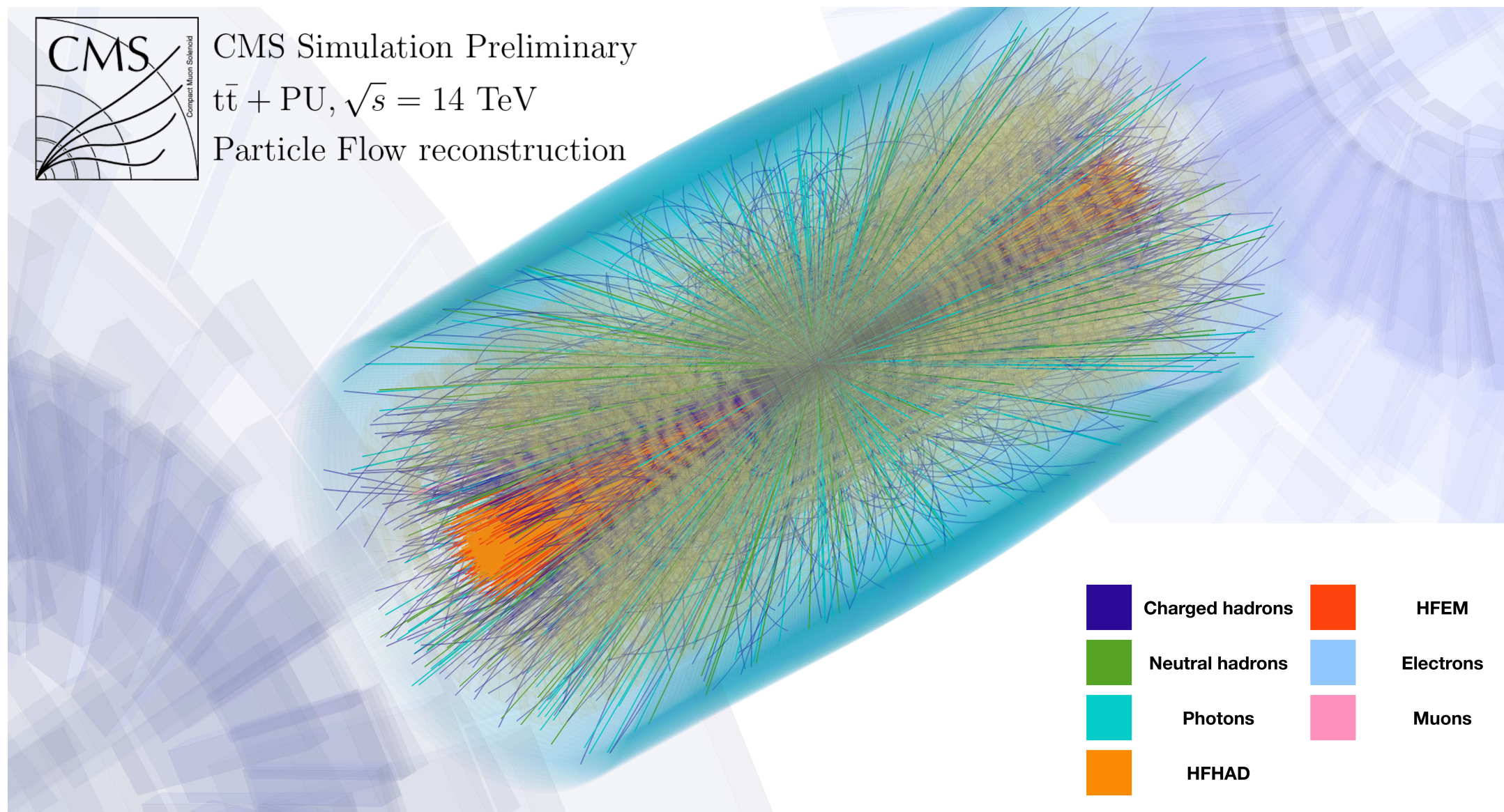
PF baseline scales non-linearly with increasing input size



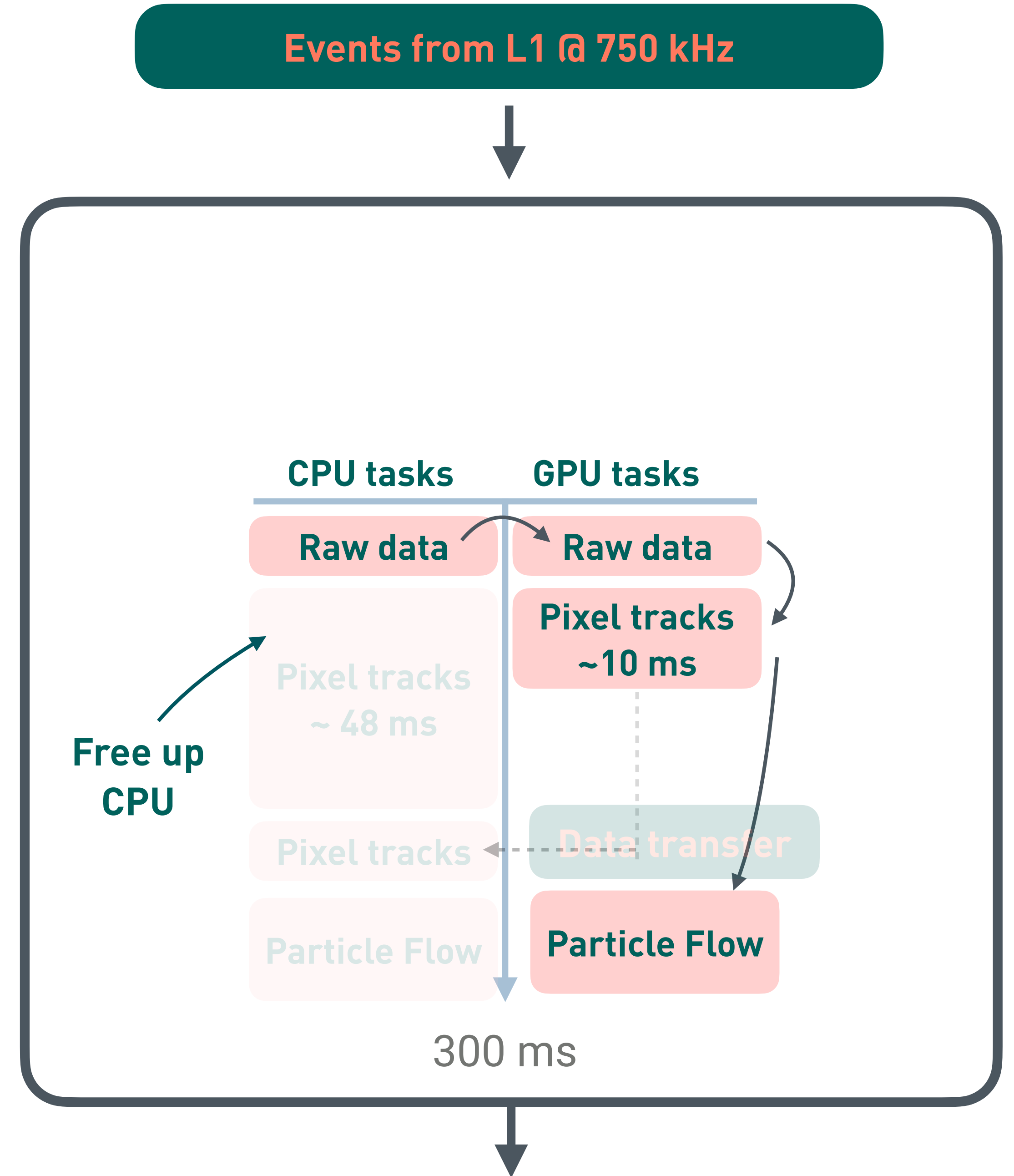
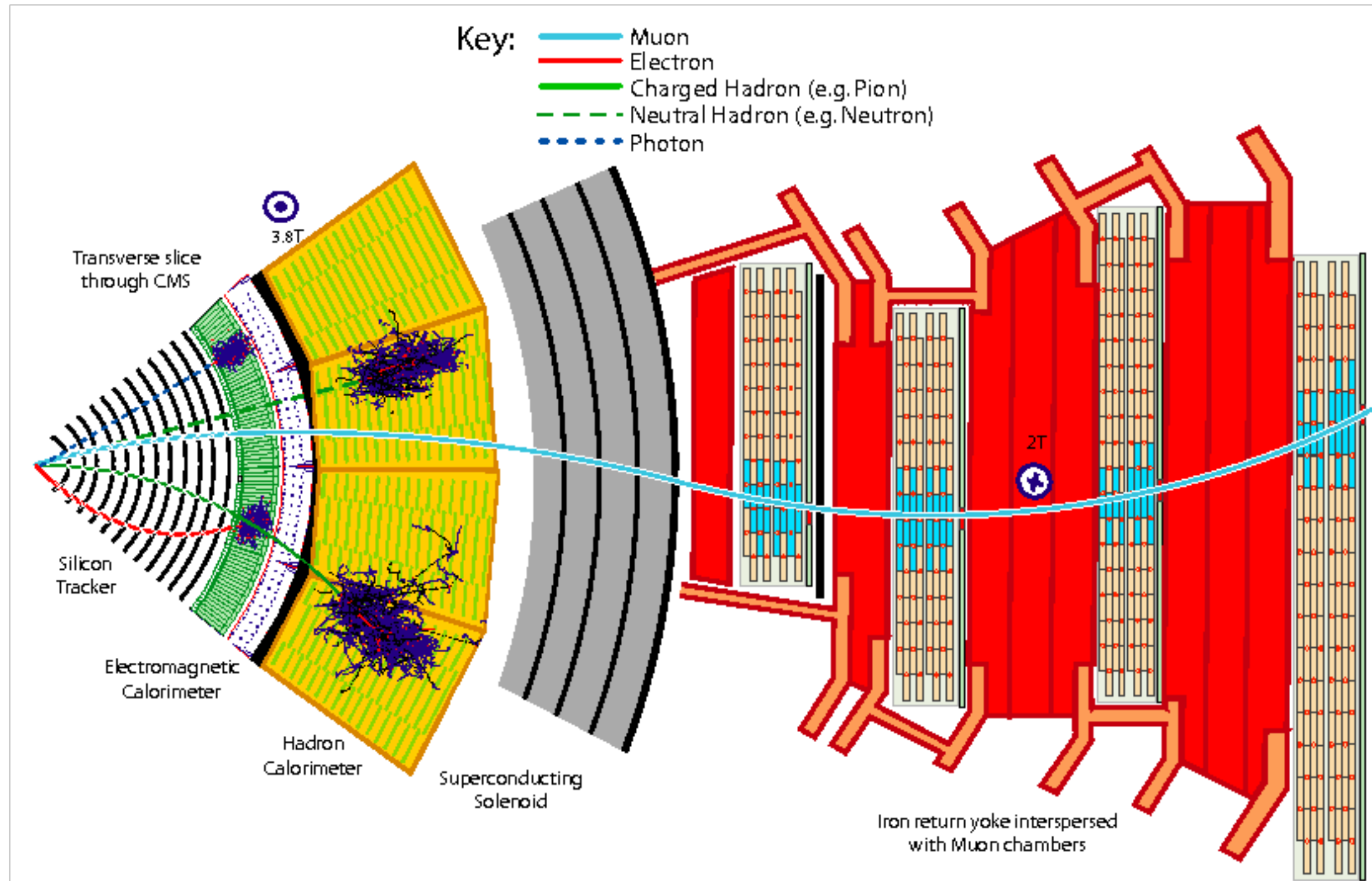
GNN-based model inference time scales approximately linearly with increasing input size

Classical Particle Flow

Graph Neural Network




ML for fast reconstruction



AI hardware

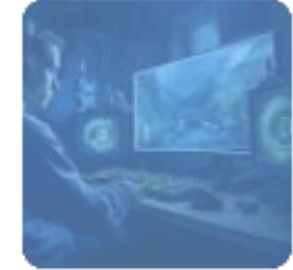
About 121'000 results (0.26 seconds)

 Game Is Hard

AMD Aims to Challenge Nvidia in the AI Hardware Market

AMD recently announced its optimistic projections for the upcoming fiscal year, with a focus on its new AI chip platform.

17 hours ago




 Tech Xplore

The future of AI hardware: Scientists unveil all-analog photoelectronic chip

Researchers from Tsinghua University, China, have developed an all-analog photoelectronic chip that combines optical and electronic...

21 hours ago



 The Information

An AI Chip Armageddon is Coming; Biden Punts on Open-Source LLMs

When I asked David Bennett, the chief customer officer of AI hardware developer Tenstorrent, about the future of startups like his,...

17 hours ago



 BBVA Openmind

Green Artificial Intelligence

As the prominence of AI continues to grow, so too does the need to address its environmental impact, particularly in terms of carbon...



AI hardware

GNNs with Versal AI, P. Schwaebig

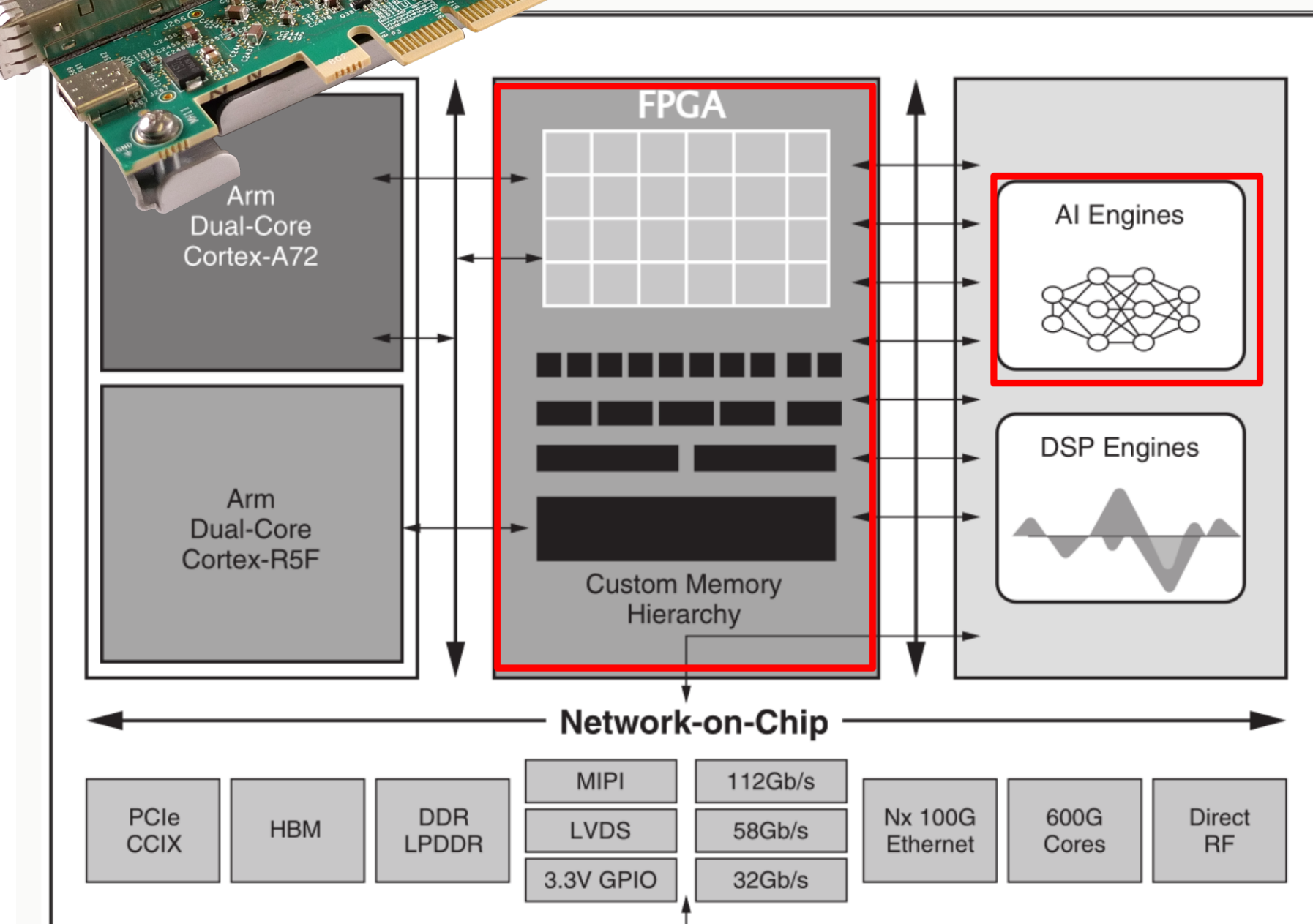
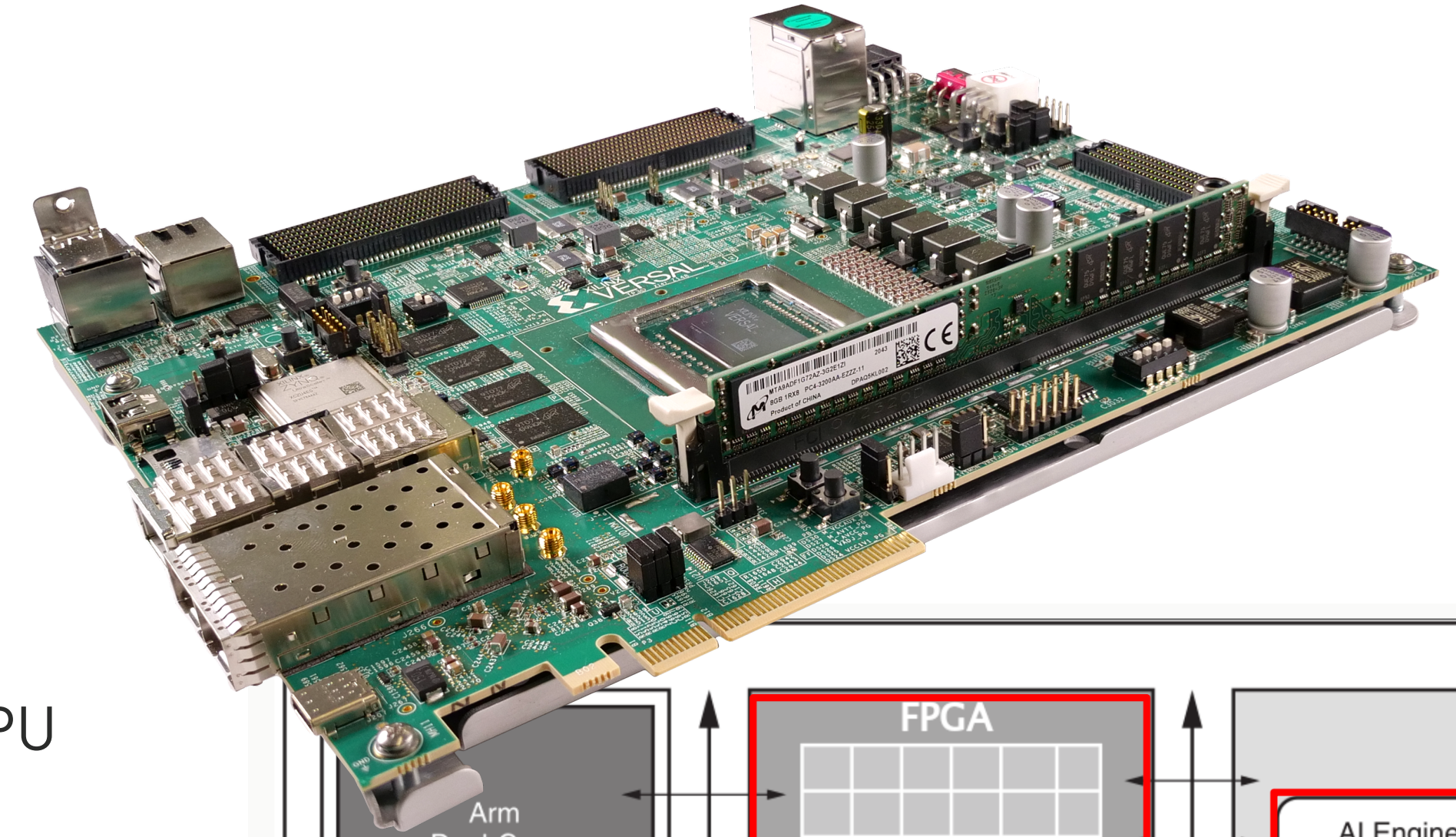
More and more dedicated AI processors on the market

Xilinx Versal AI processors

- Programmed in C/C++
- 400 AI processors, ~2M logic cells (FPGA), 2k DSPs, Arm CPU and RPU
- Data move back and forth between AI Engines and FPGA

Currently explored for real-time tracking in trigger application

- Interaction Network for pattern recognition (similar to DeZoort et al)
- Deployed on Xilinx Versal VC1902 ACAP



Source: Xilinx, WP505,

Subscribe to newsletters

Forbes

FORBES > INNOVATION > CLOUD

Groq's Record-Breaking Language Processor Hits 100 Tokens Per Second On A Massive AI Model

Paul Smith-Goodson Contributor
Moor Insights and Strategy Contributor Group

Follow

0

Aug 11, 2023, 03:29pm EDT



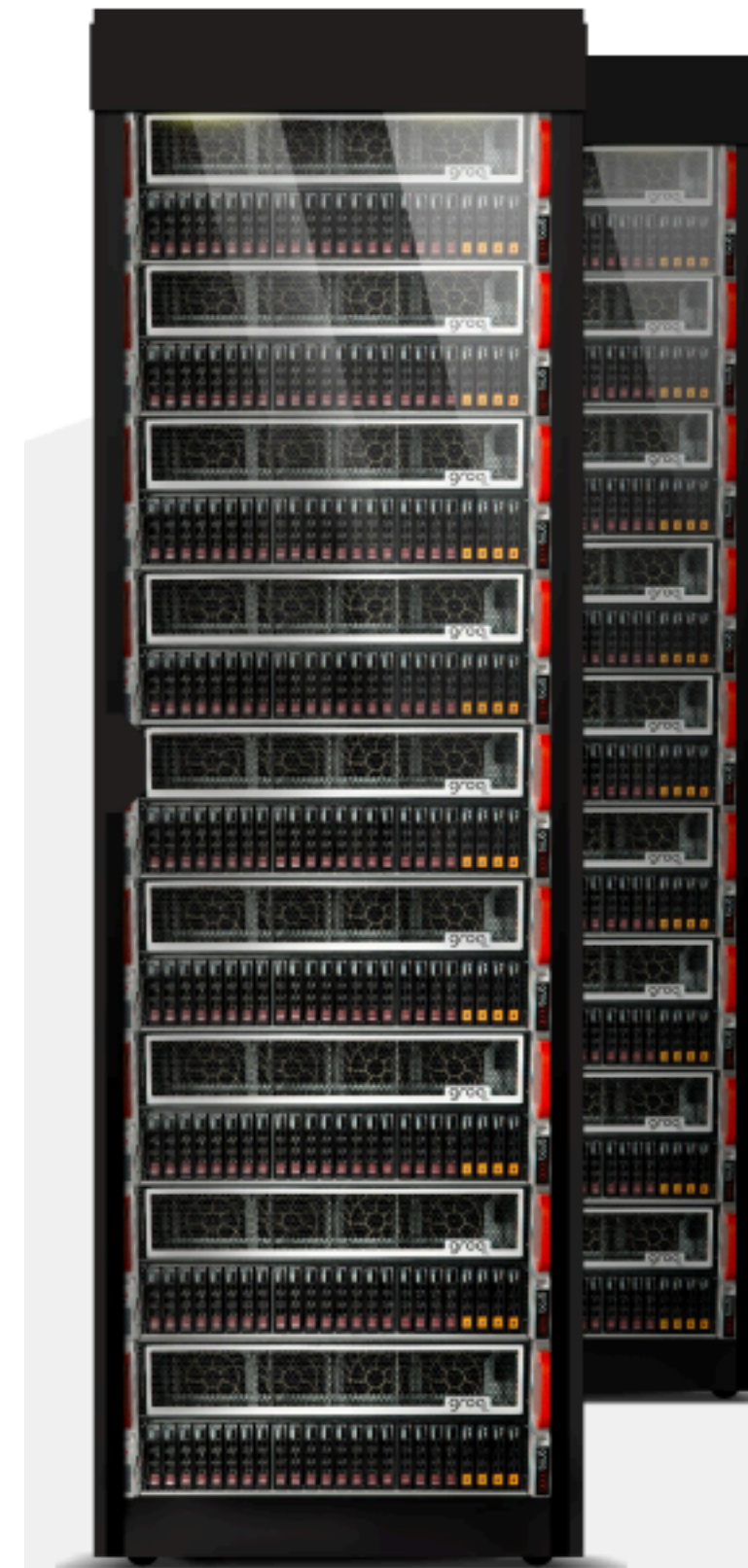
0:14

Ad 1 of 2

0:00

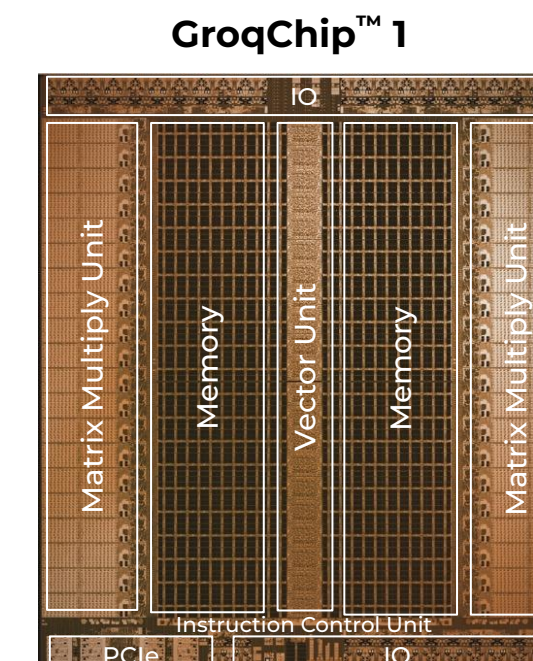


FIRST
TO REACH
> 100
TOKENS
PER SECOND
Running Llama-2 70B
on a Groq LPU™ System



Groq: ultra-low latency dedicated language processor dedicated language processors

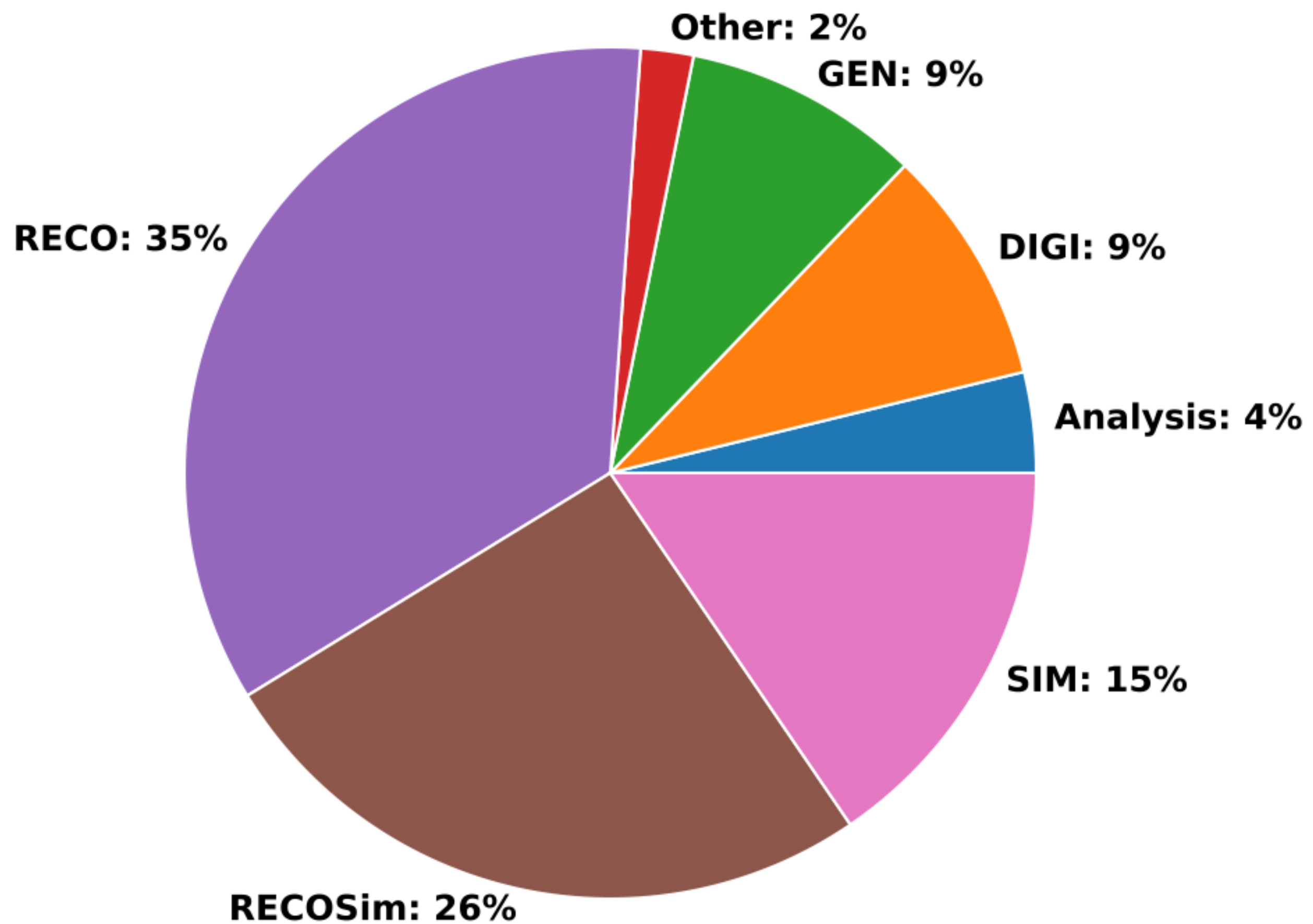
- Optimised for sequential data
- First ever 100 tokens/s (usually, ~10)



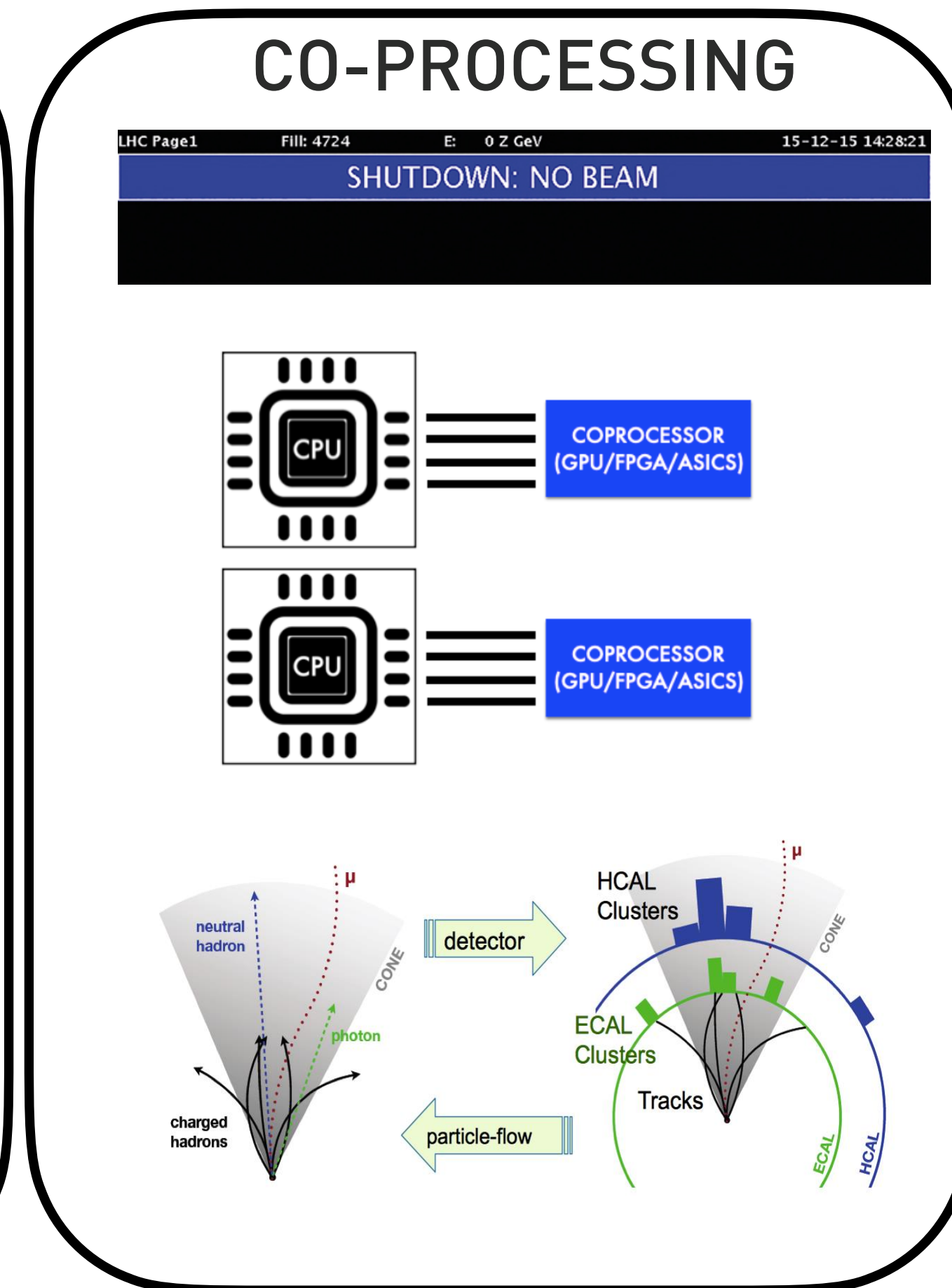
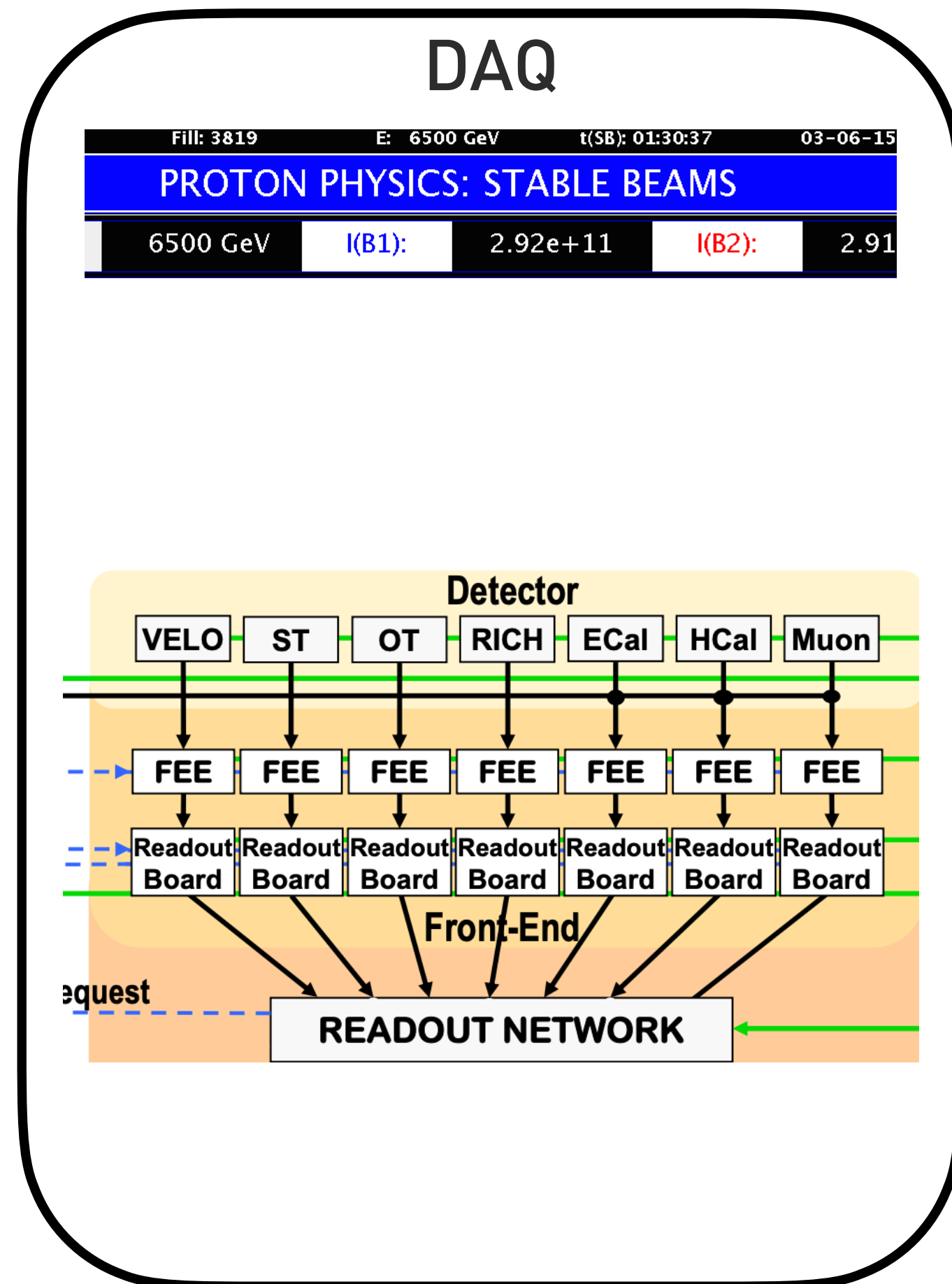
GroqRack™

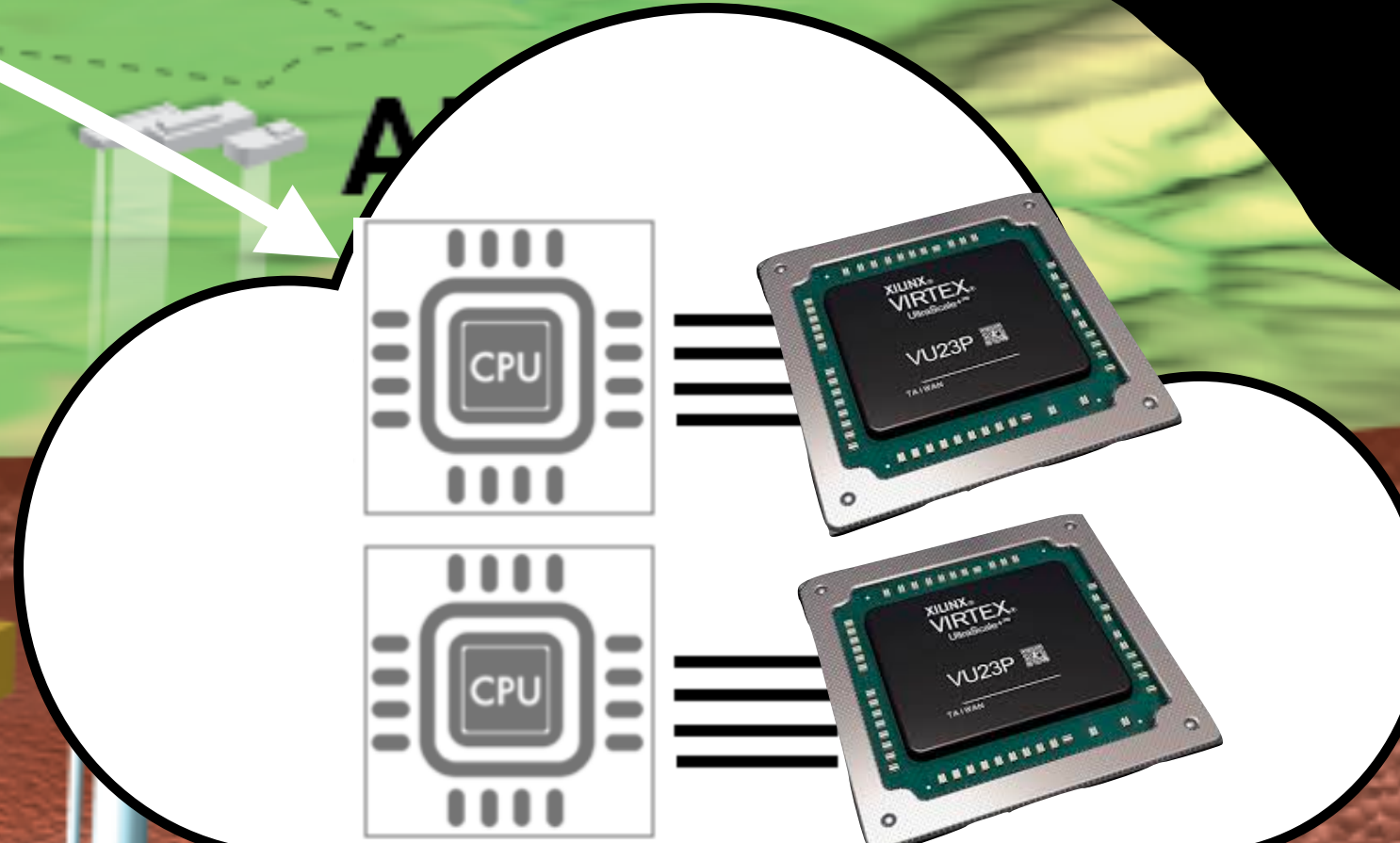
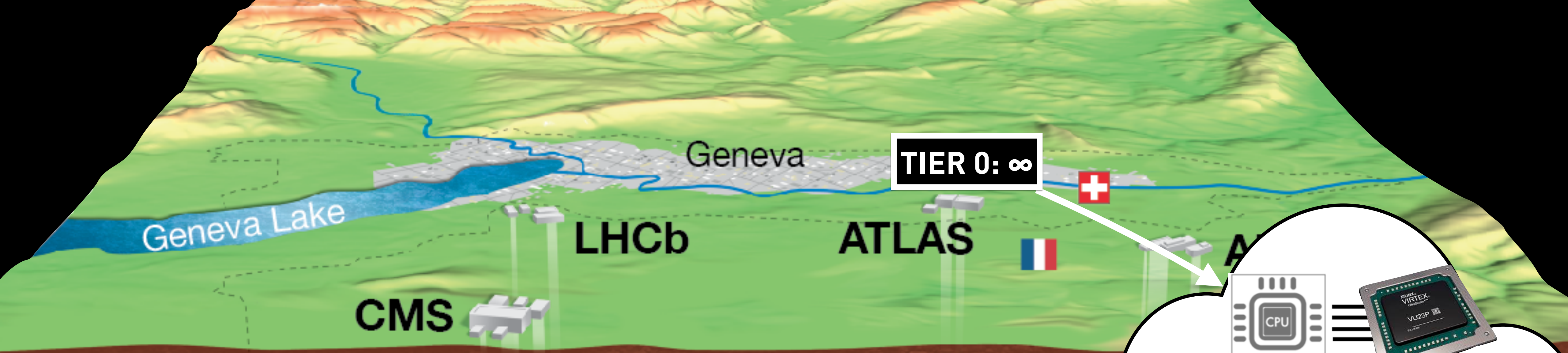
CMSPublic

Total CPU HL-LHC (2031/No R&D Improvements) fractions
2022 Estimates

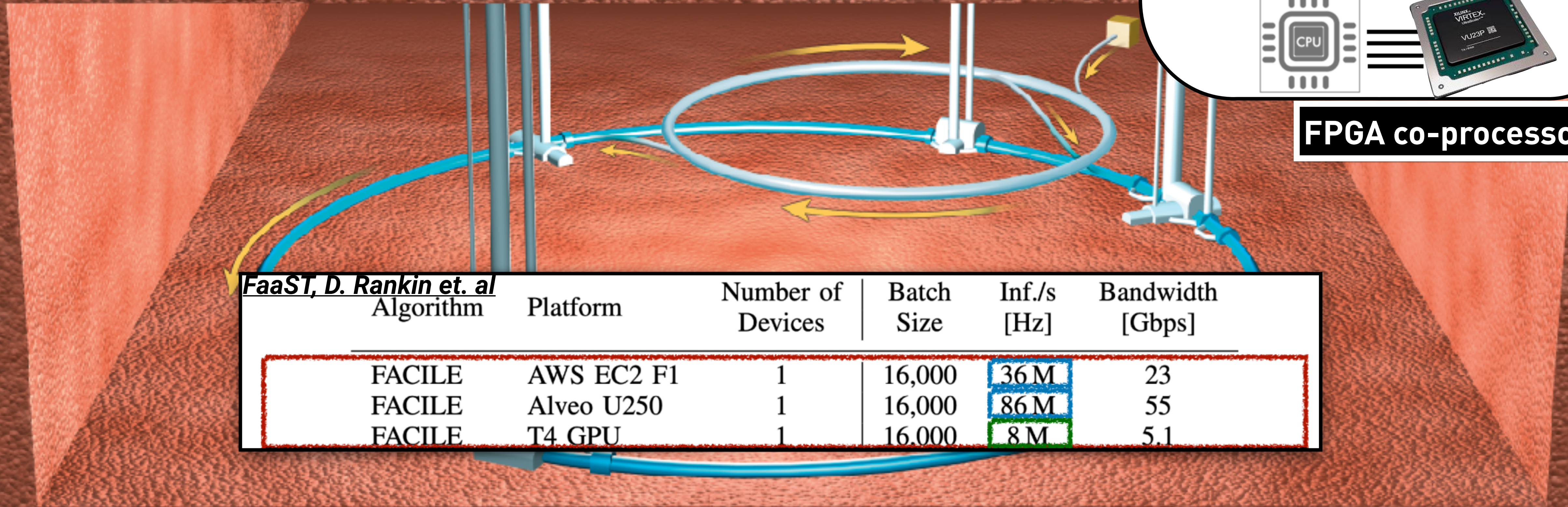


FPGAs as accelerators





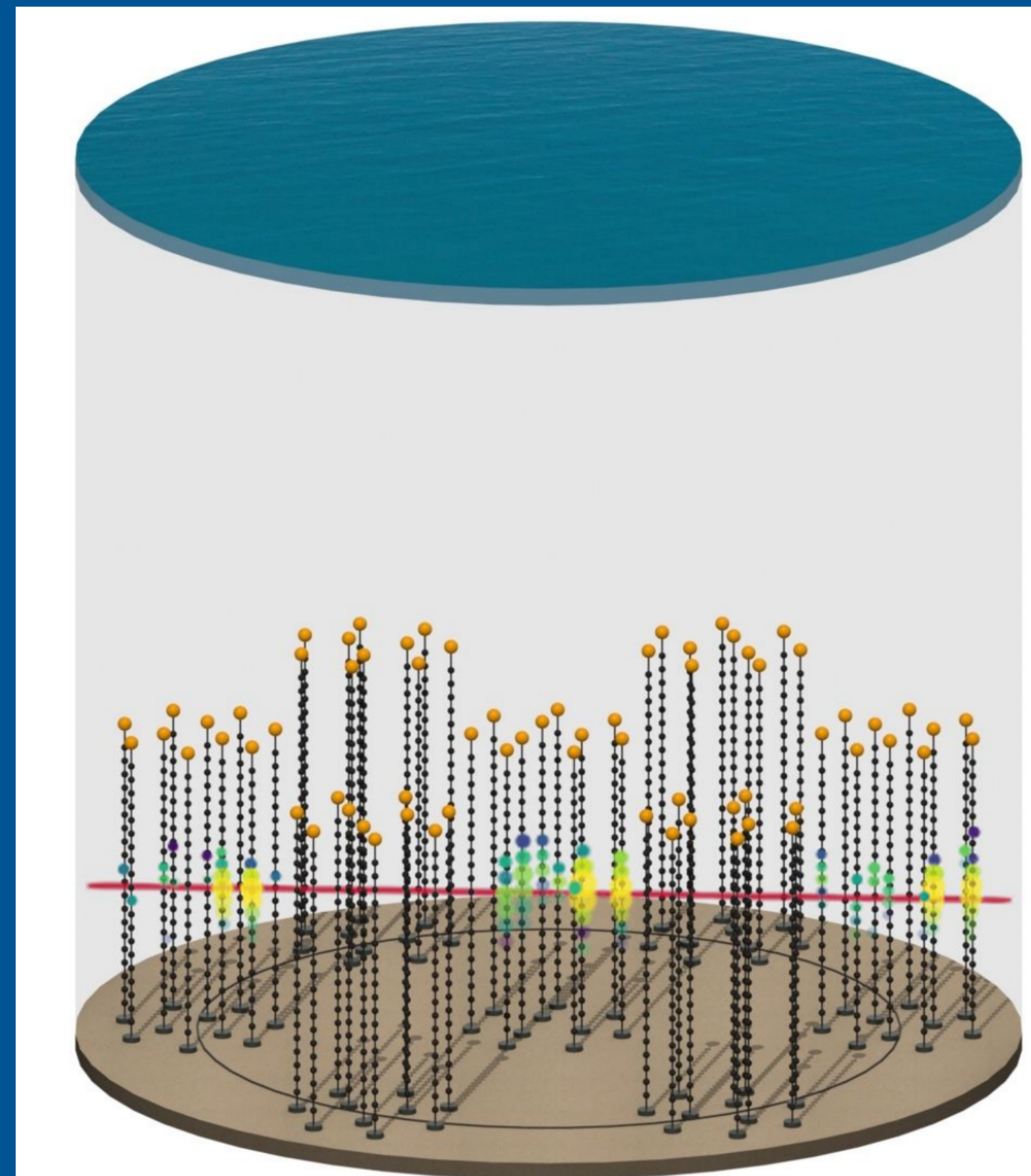
FPGA co-processors

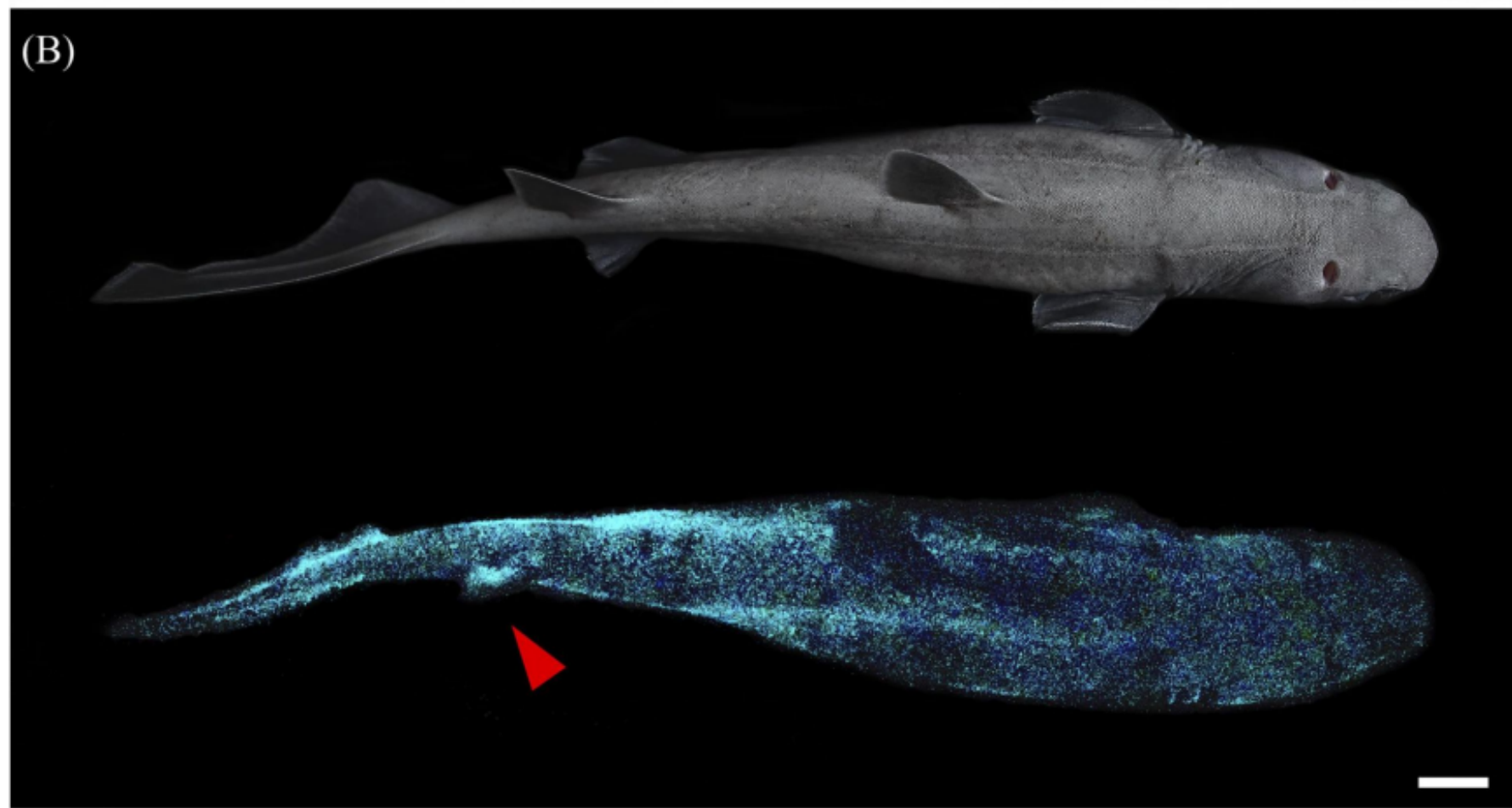
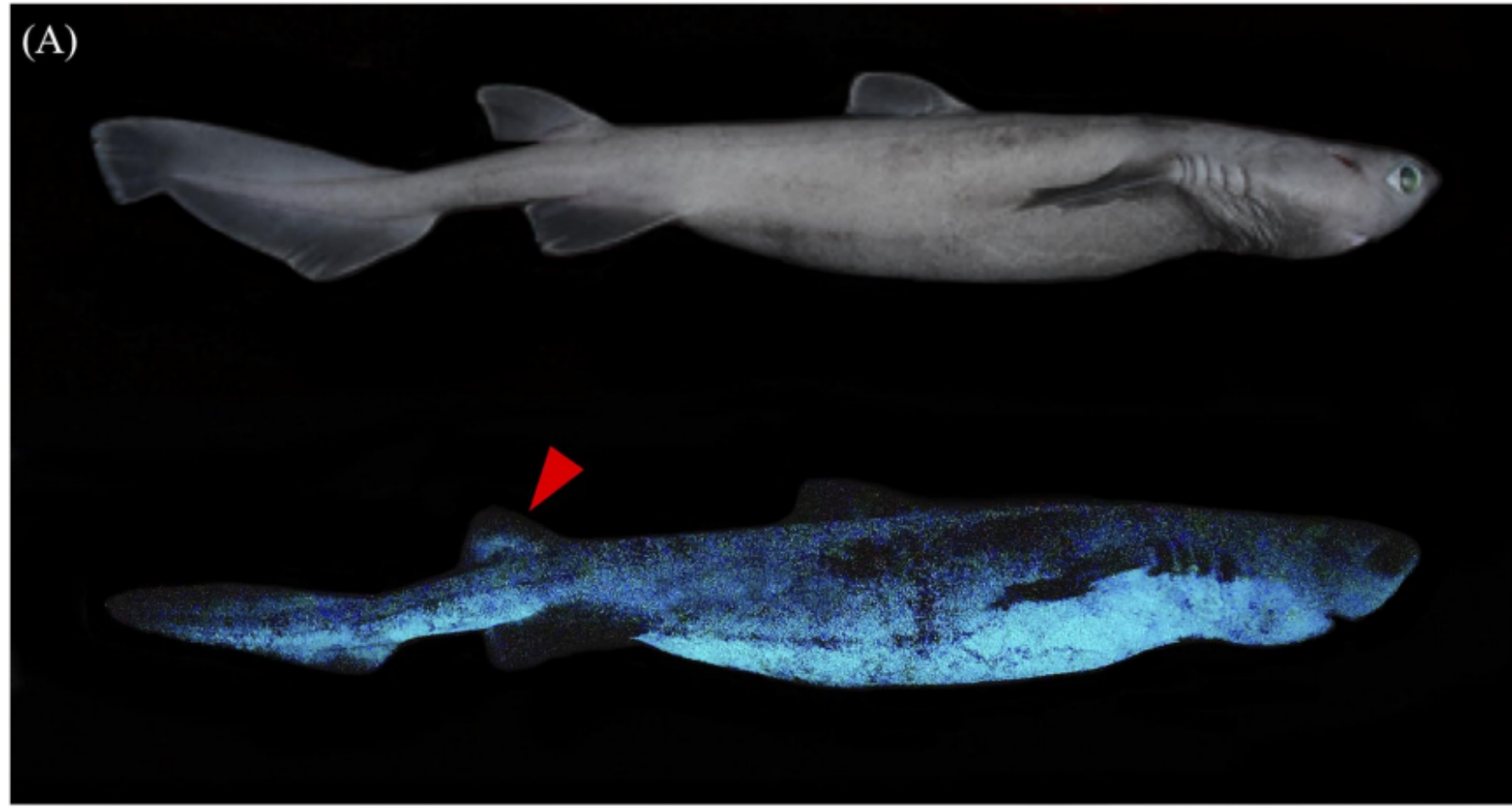


FaaST, D. Rankin et. al

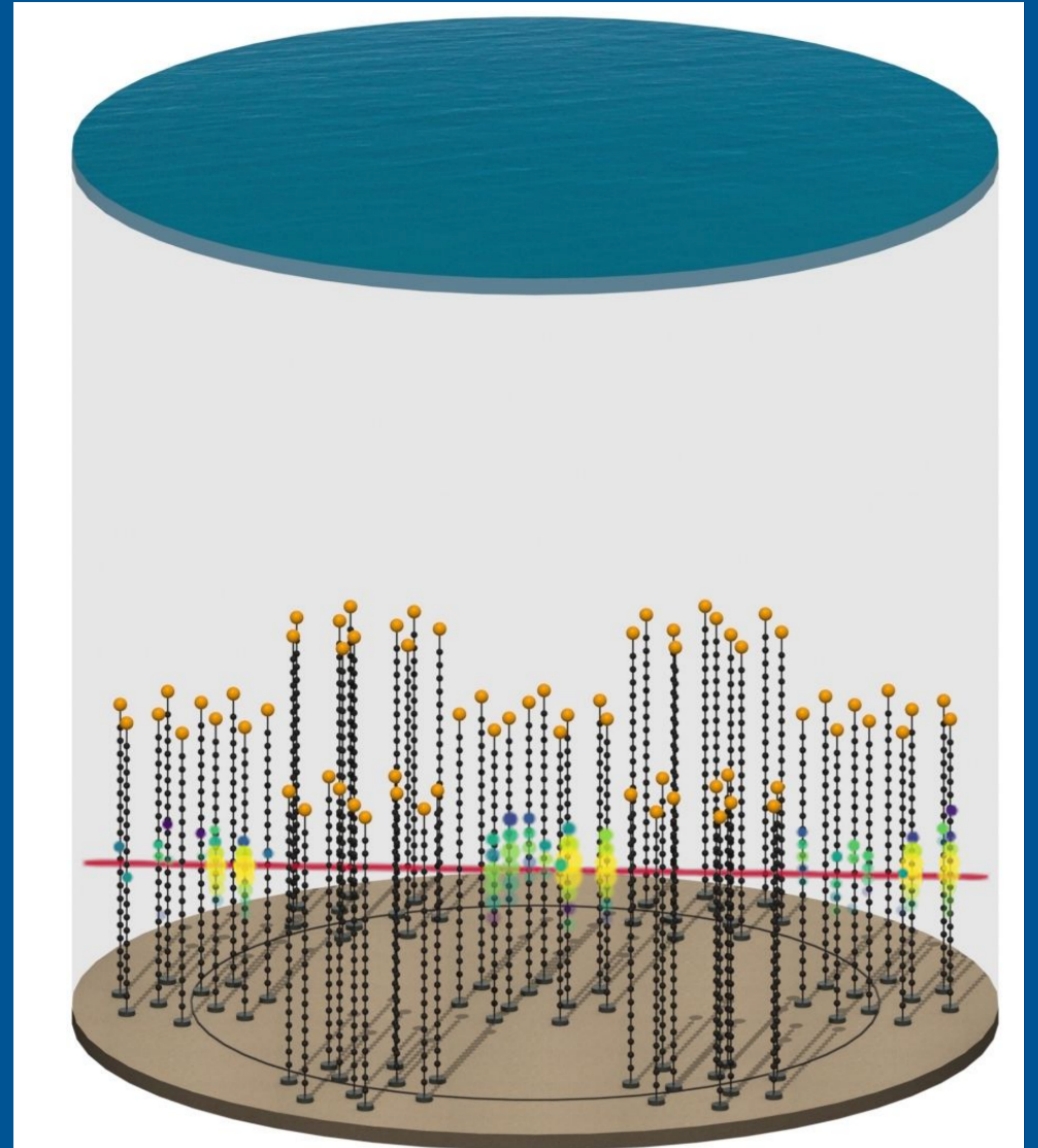
Algorithm	Platform	Number of Devices	Batch Size	Inf./s [Hz]	Bandwidth [Gbps]
FACILE	AWS EC2 F1	1	16,000	36 M	23
FACILE	Alveo U250	1	16,000	86 M	55
FACILE	T4 GPU	1	16,000	8 M	5.1

Triggering in other experiments

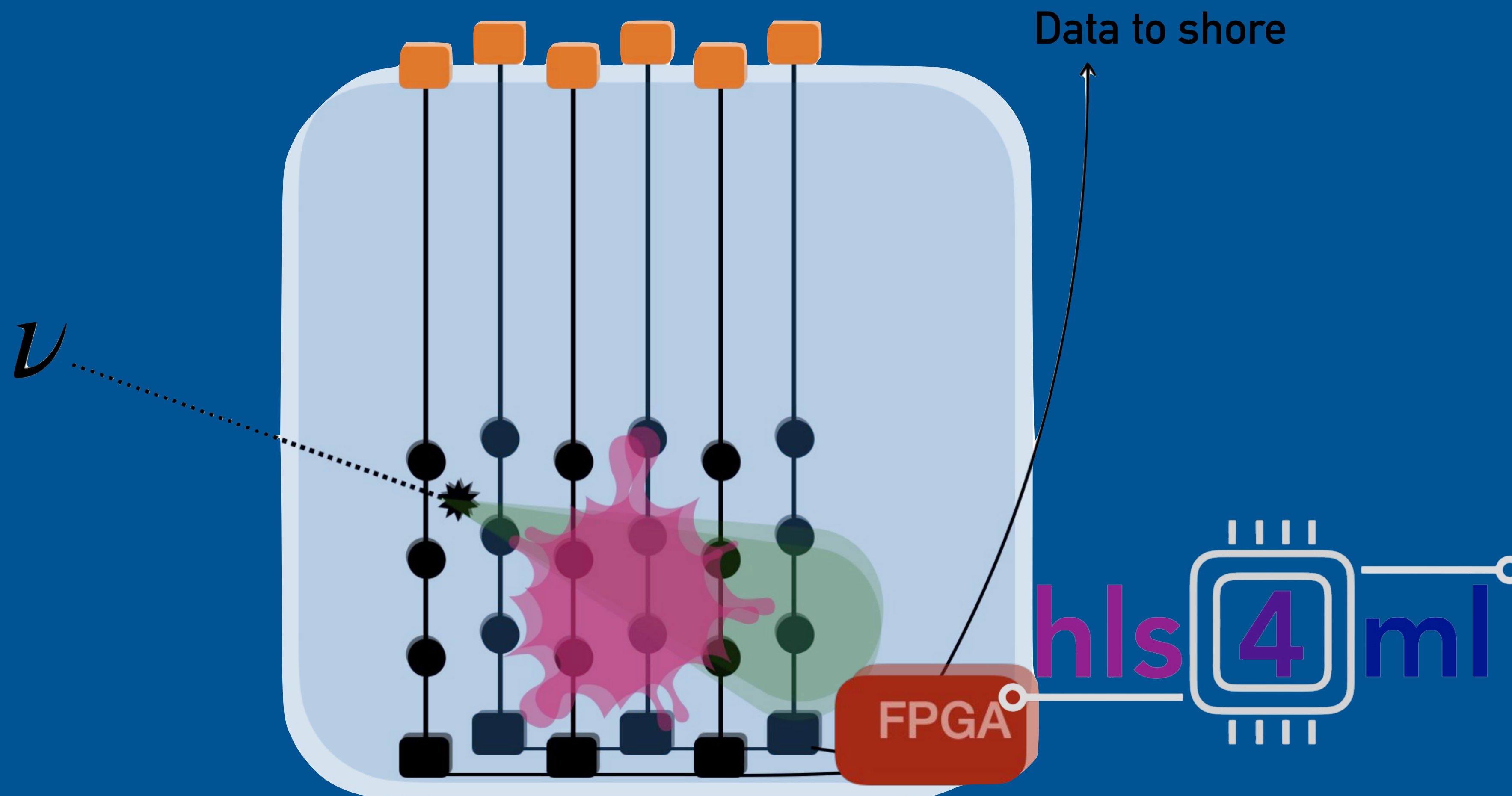




Bioluminescence bursts up to few MHz!



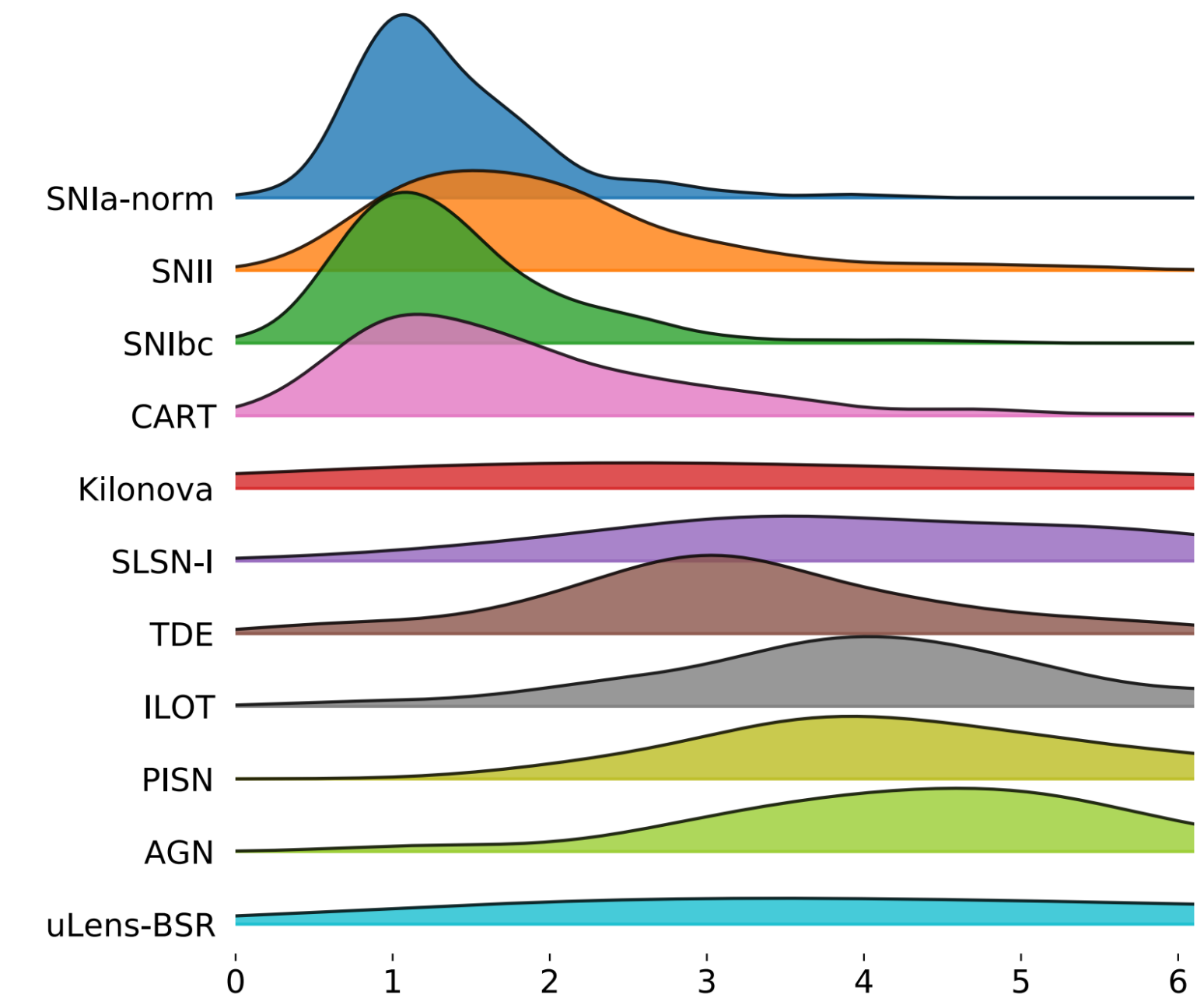
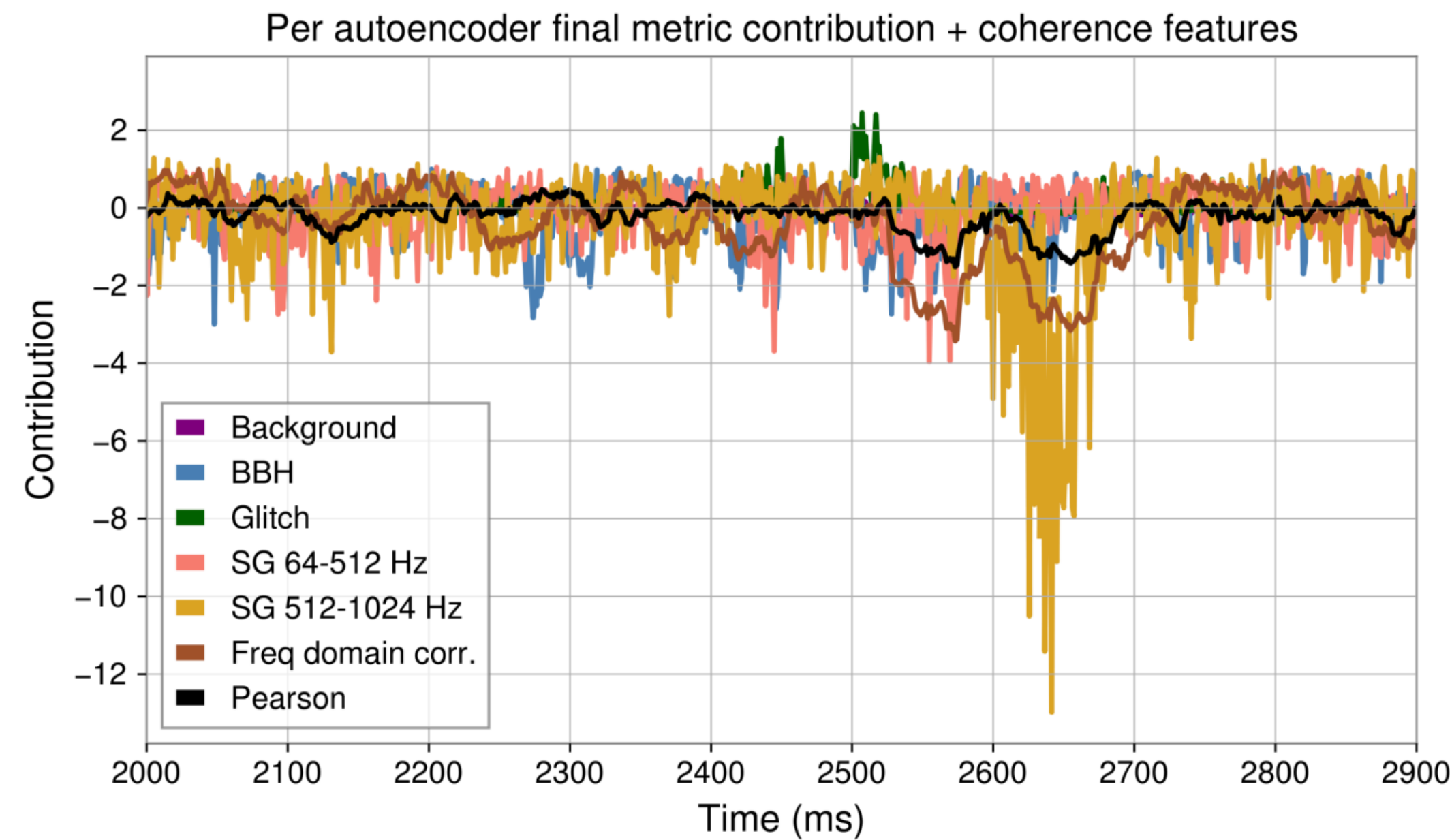
Signals and backgrounds



Gravitational wave

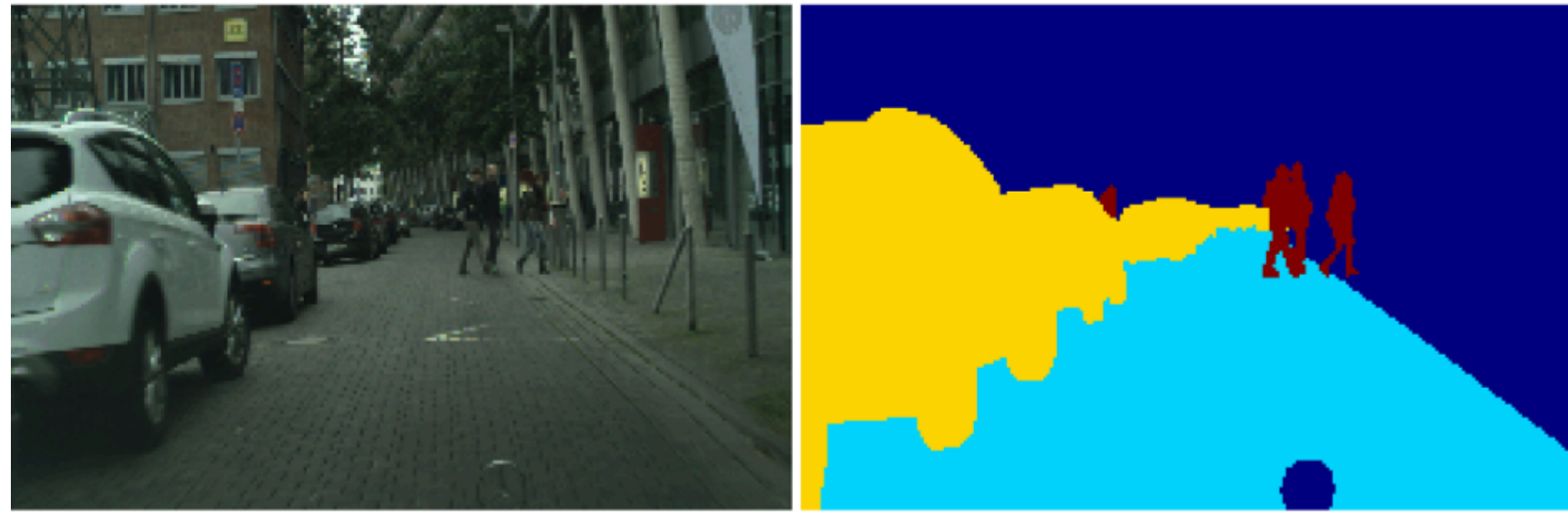
Multi-messenger astronomy

Real-time Multi-Messenger Anomaly Detection



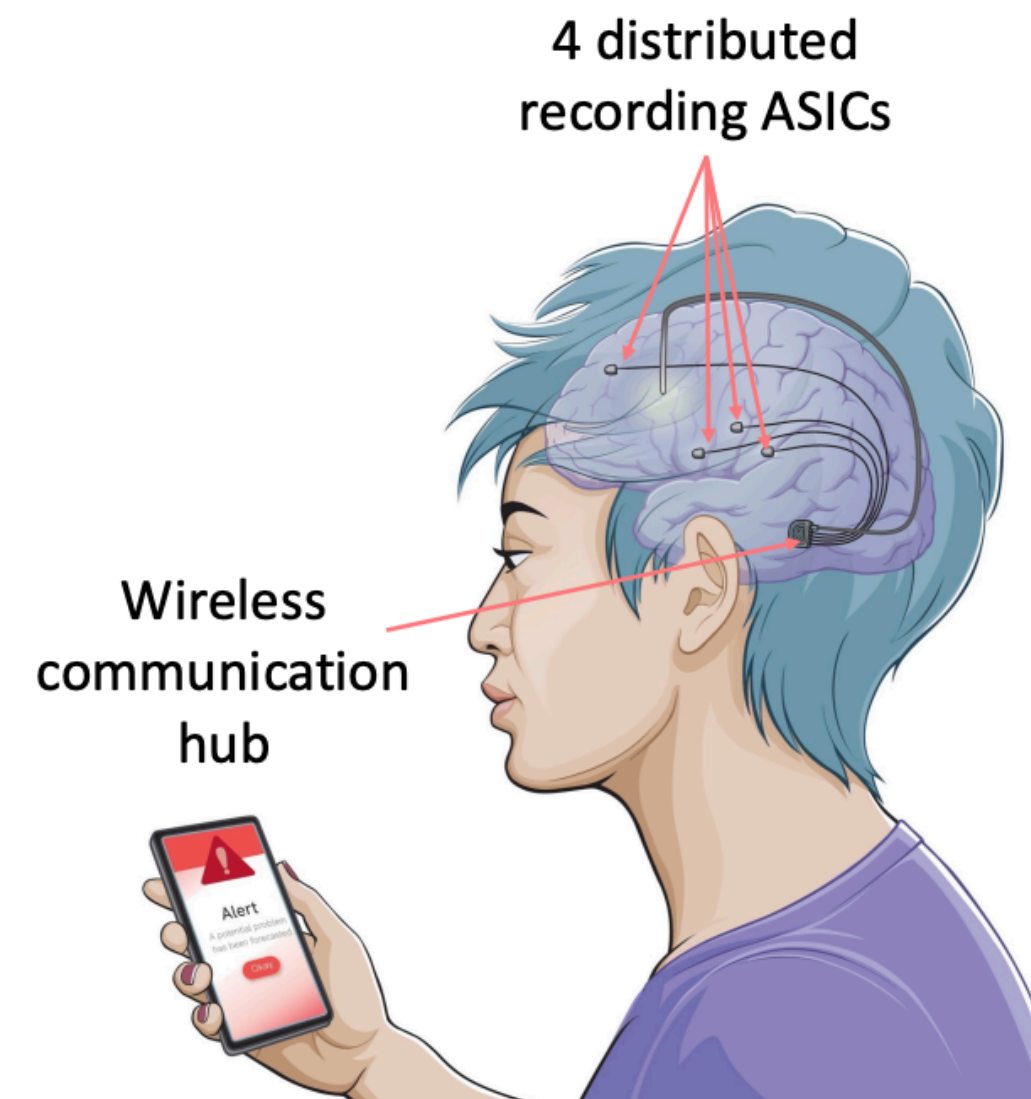
...and outside of particle physics

Semantic segmentation for autonomous vehicles



N. Ghielmetti et al.

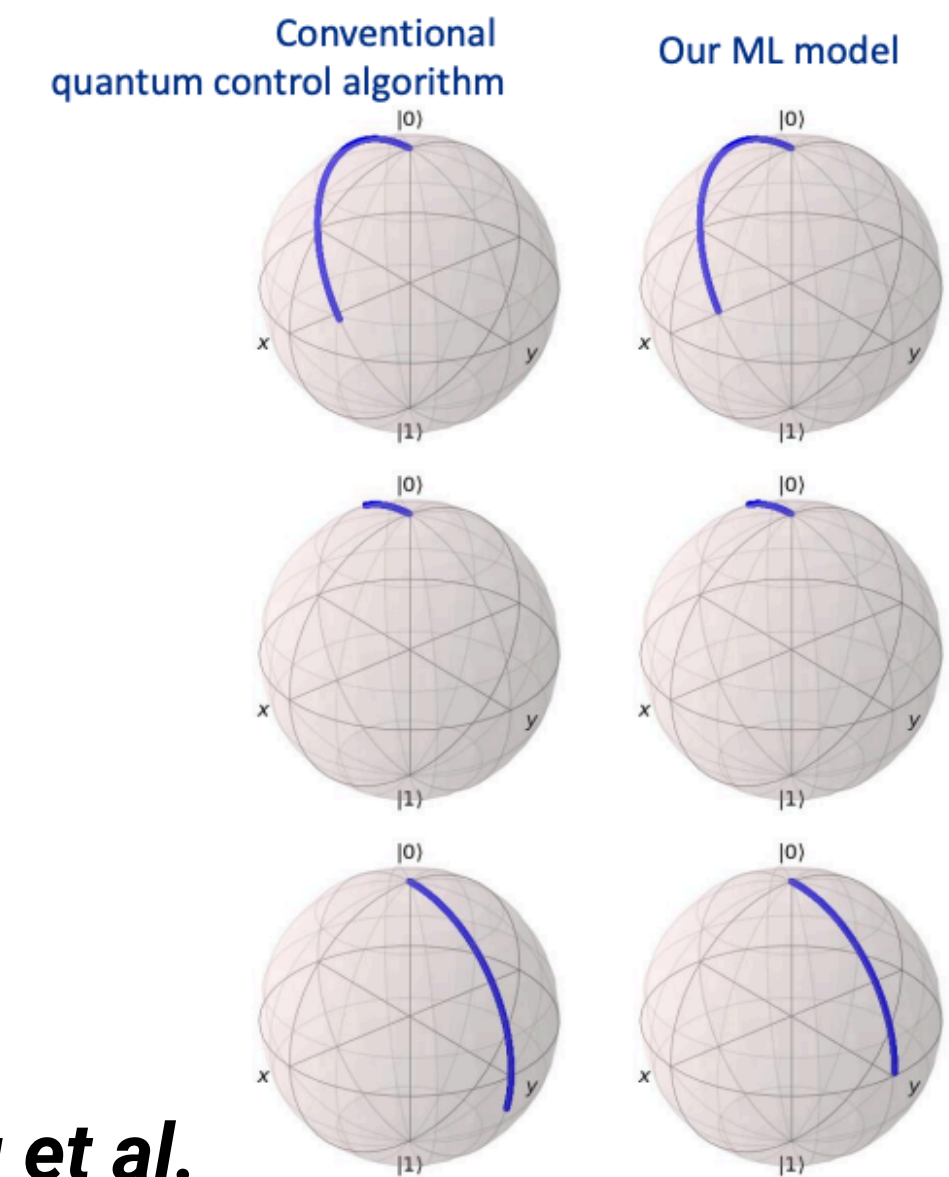
Seizure Predicting Brain Implant



W. Lemaire et al.

NN accelerator for quantum control

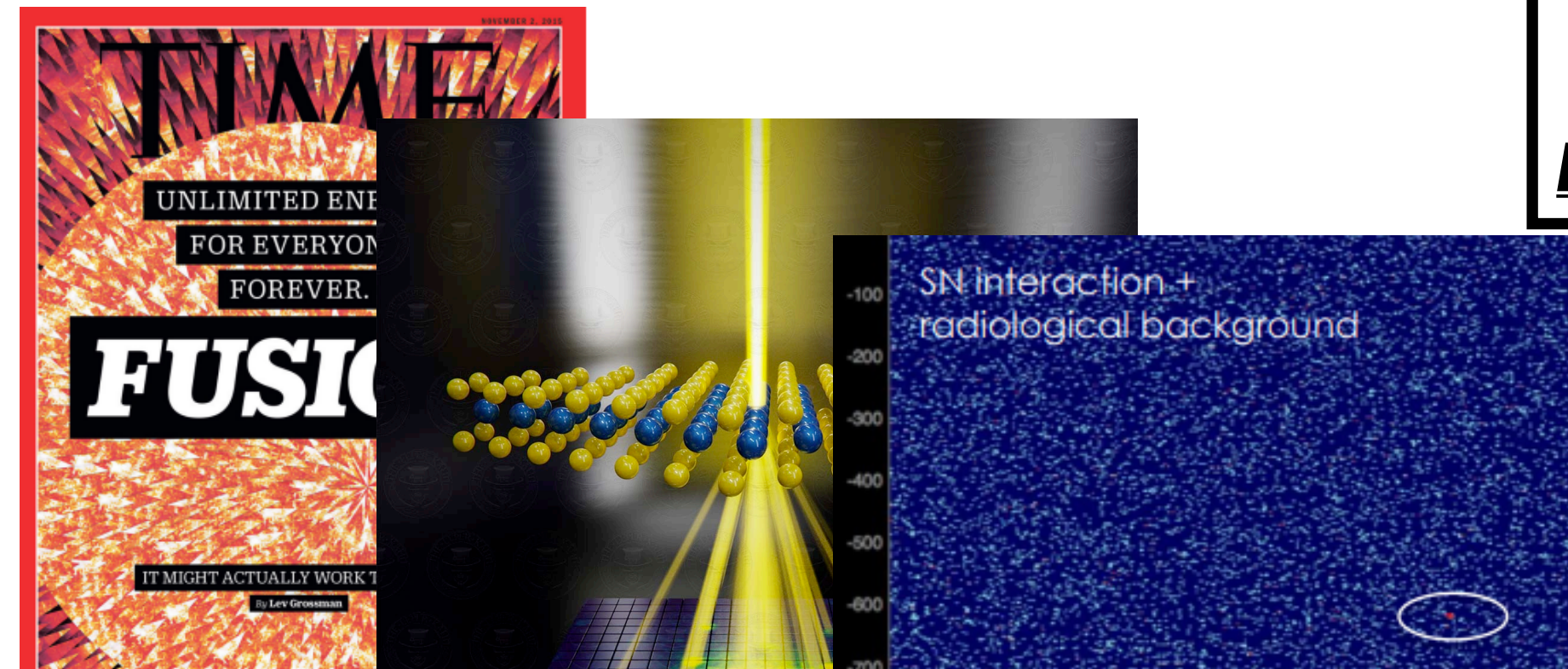
- Putting control in cryostat (e.g optimal pulse parameters)



D Xu et al.

Other examples

- ***For fusion science phase/mode monitoring***
- ***Crystal structure detection***
- ***Triggering in DUNE***
- ***Accelerator control***
- ***Magnet Quench Detection***
- ***MLPerf tinyML benchmarking***
- ***Food contamination detection***
- etc....





Fast Machine Learning Lab

FastML Lab

Real-time and accelerated ML for fundamental sciences

Fast ML Lab is a research collective of physicists, engineers, and computer scientists interested in deploying machine learning algorithms for unique and challenging scientific applications. Our projects range from real-time, on-detector and low latency machine learning applications to high-throughput heterogeneous computing big data challenges. We are interested in deploying sophisticated machine learning algorithms to advance the exploration of fundamental physics from the world's biggest colliders to the most intense particle beams to the cosmos.

<https://fastmachinelearning.org/>



Hardware and Algorithm Co-development

Developing AI methods to encode non-lattice-structured data is one main challenge in current AI systems.

[Read More >](#)



High Energy Physics

Build tools to process LHC collisions occurring 40 million times per second data in real-time using AI.

[Read More >](#)



Systems Neuroscience

Discover the computations that brain-wide neural networks perform to process sensory and motor information during behavior by using high-throughput and low-latency AI algorithms to process.

[Read More >](#)

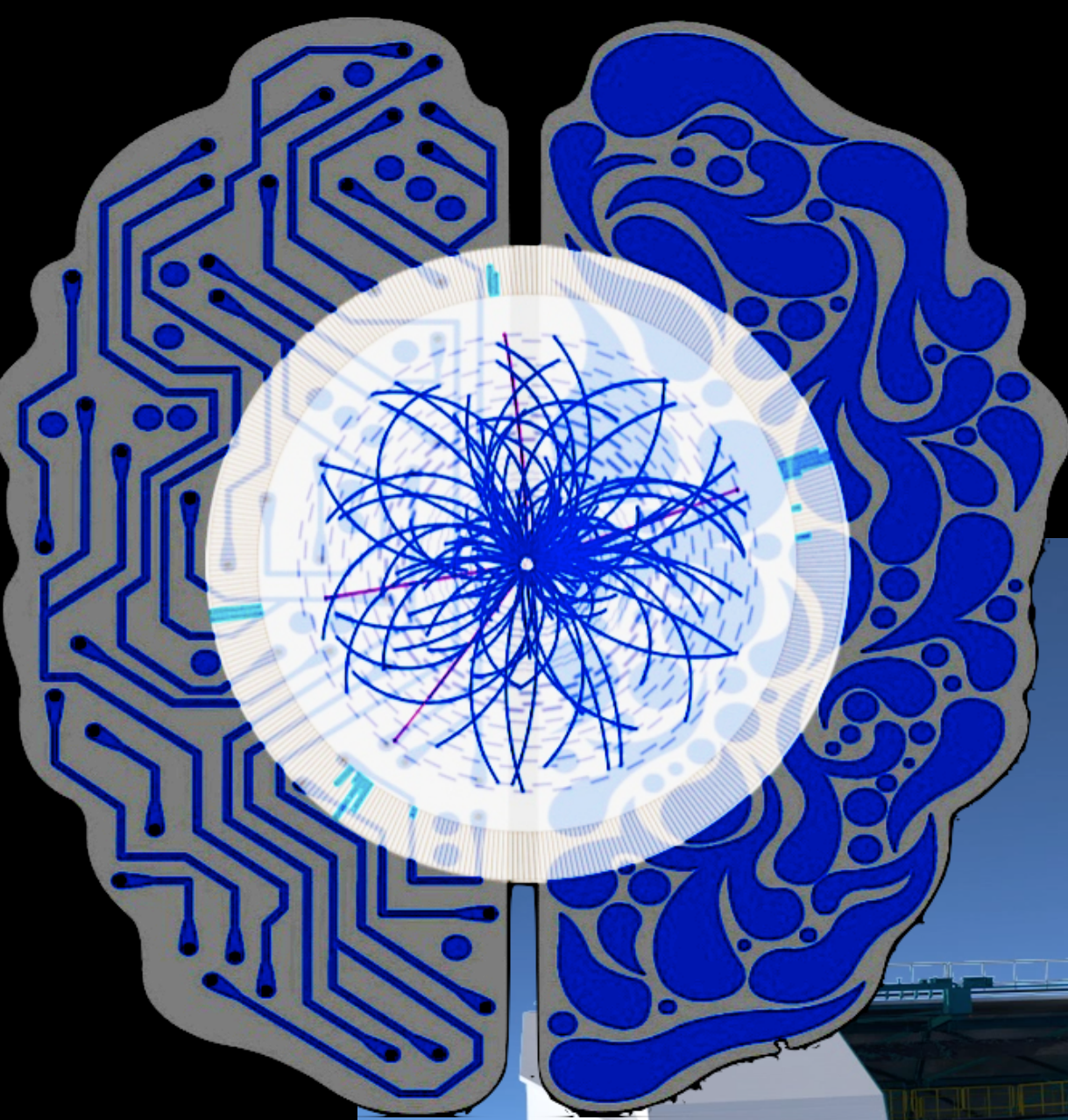


Multi-messenger Astrophysics

Process the data from telescopes, neutrino detectors, and gravitational-wave detectors to identify astronomical events corresponding to the most violent phenomena in the Cosmos.

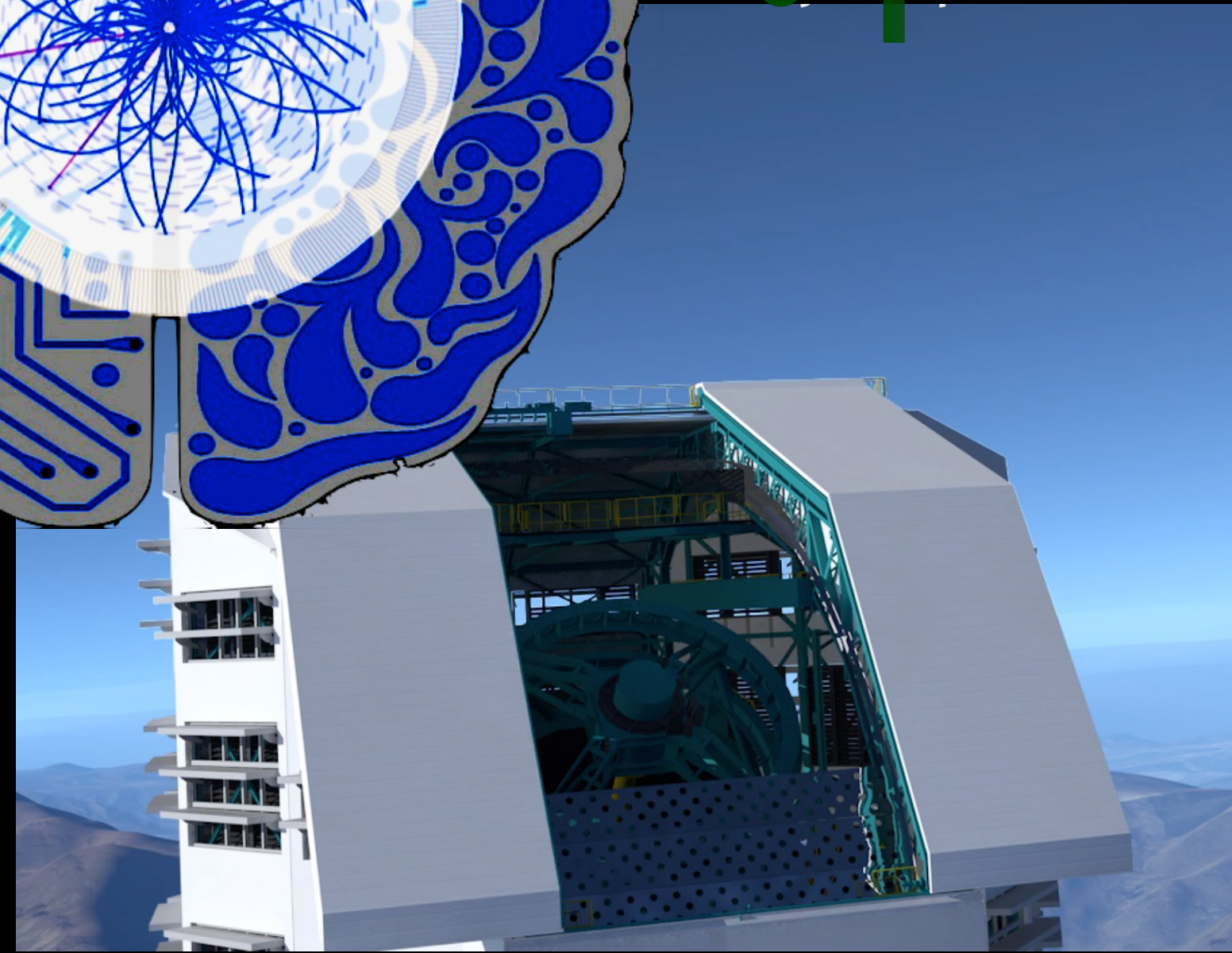
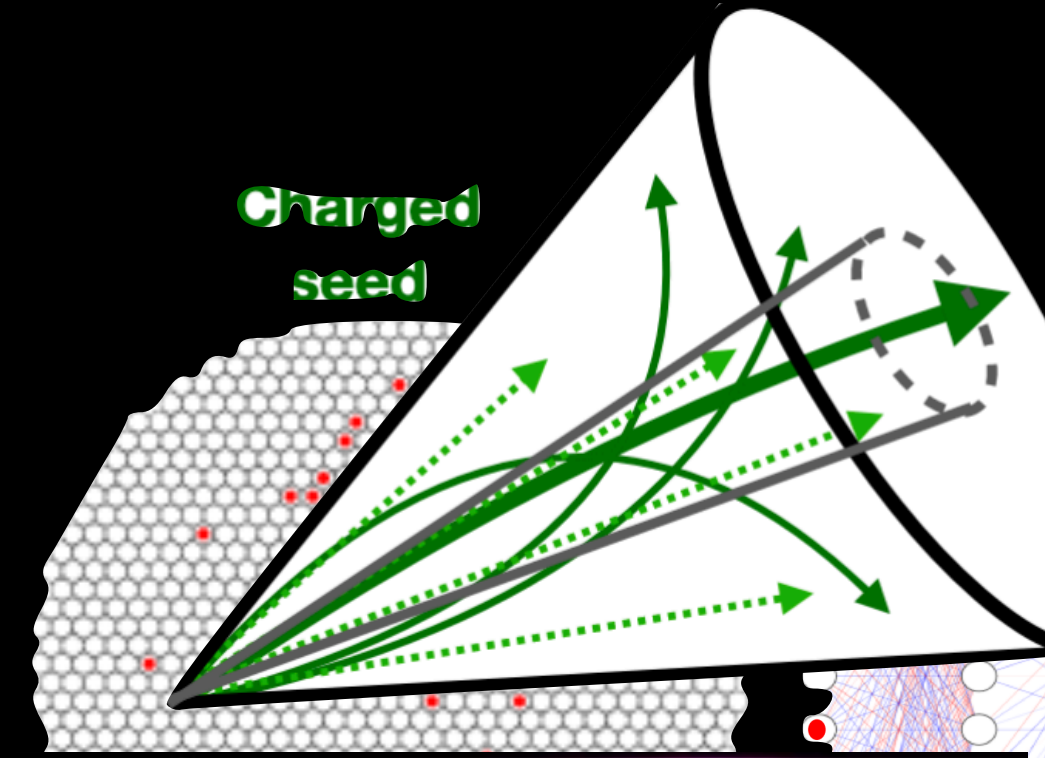
[Read More >](#)

<https://a3d3.ai/>

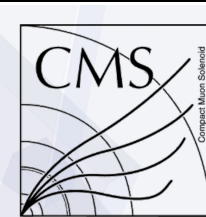
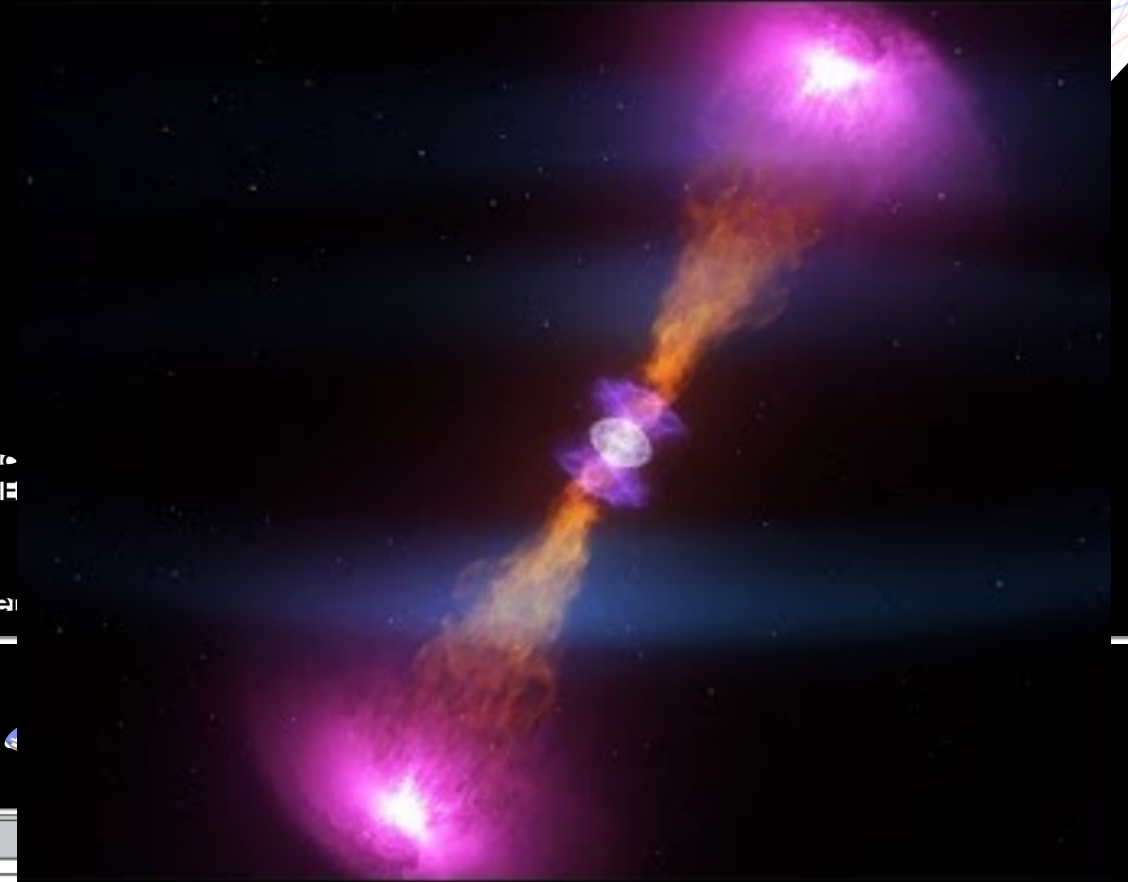


Conifer

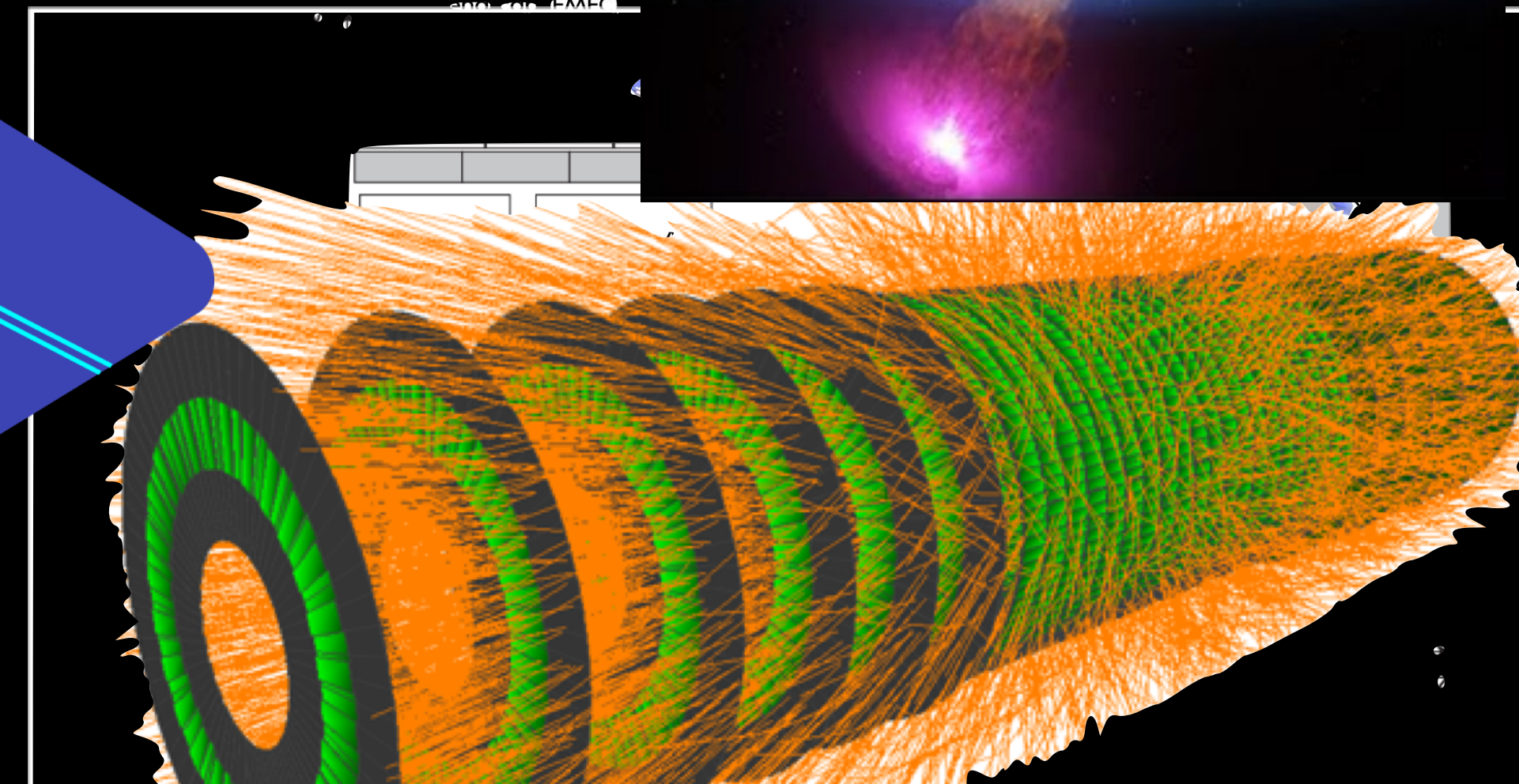
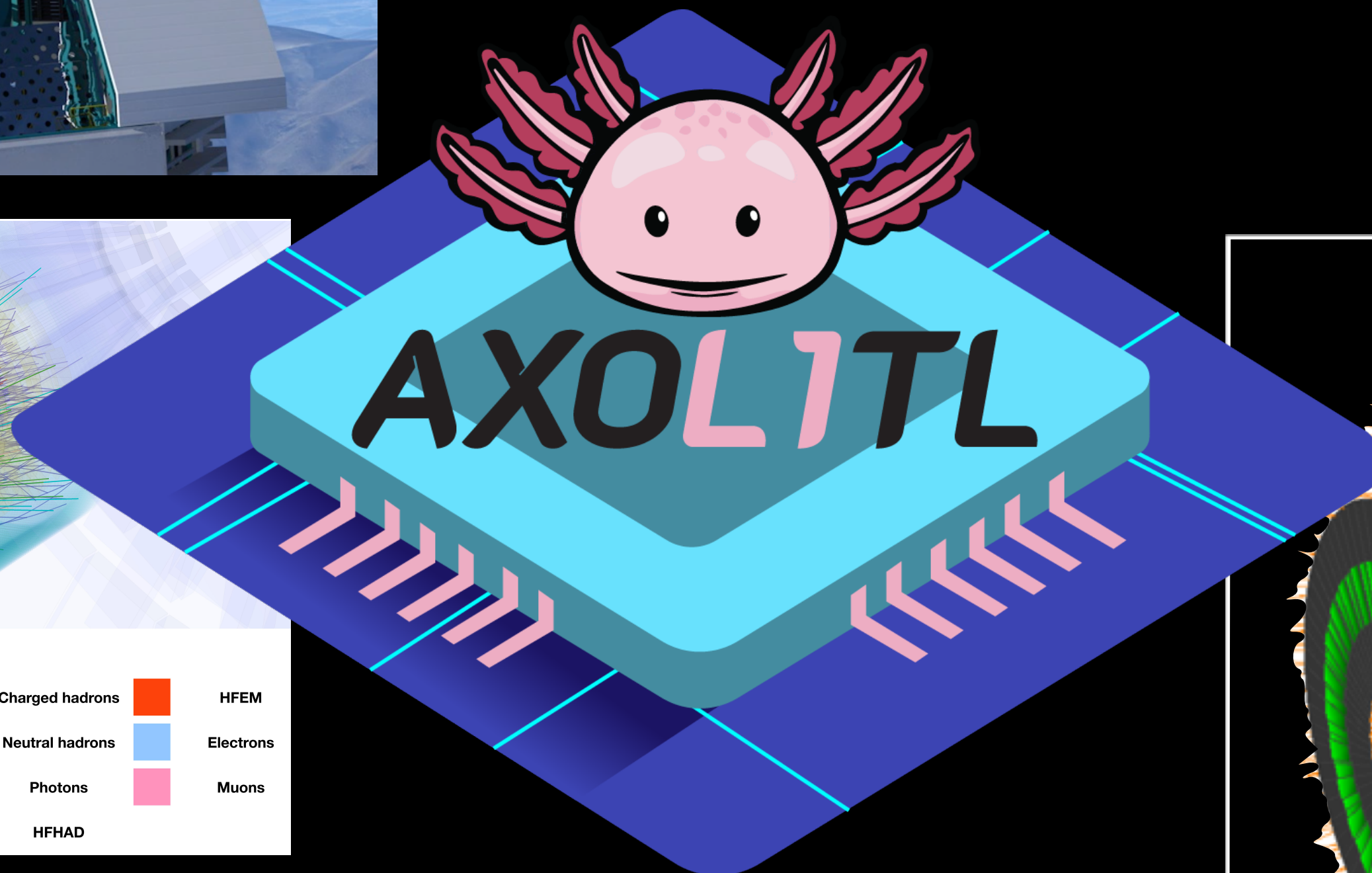
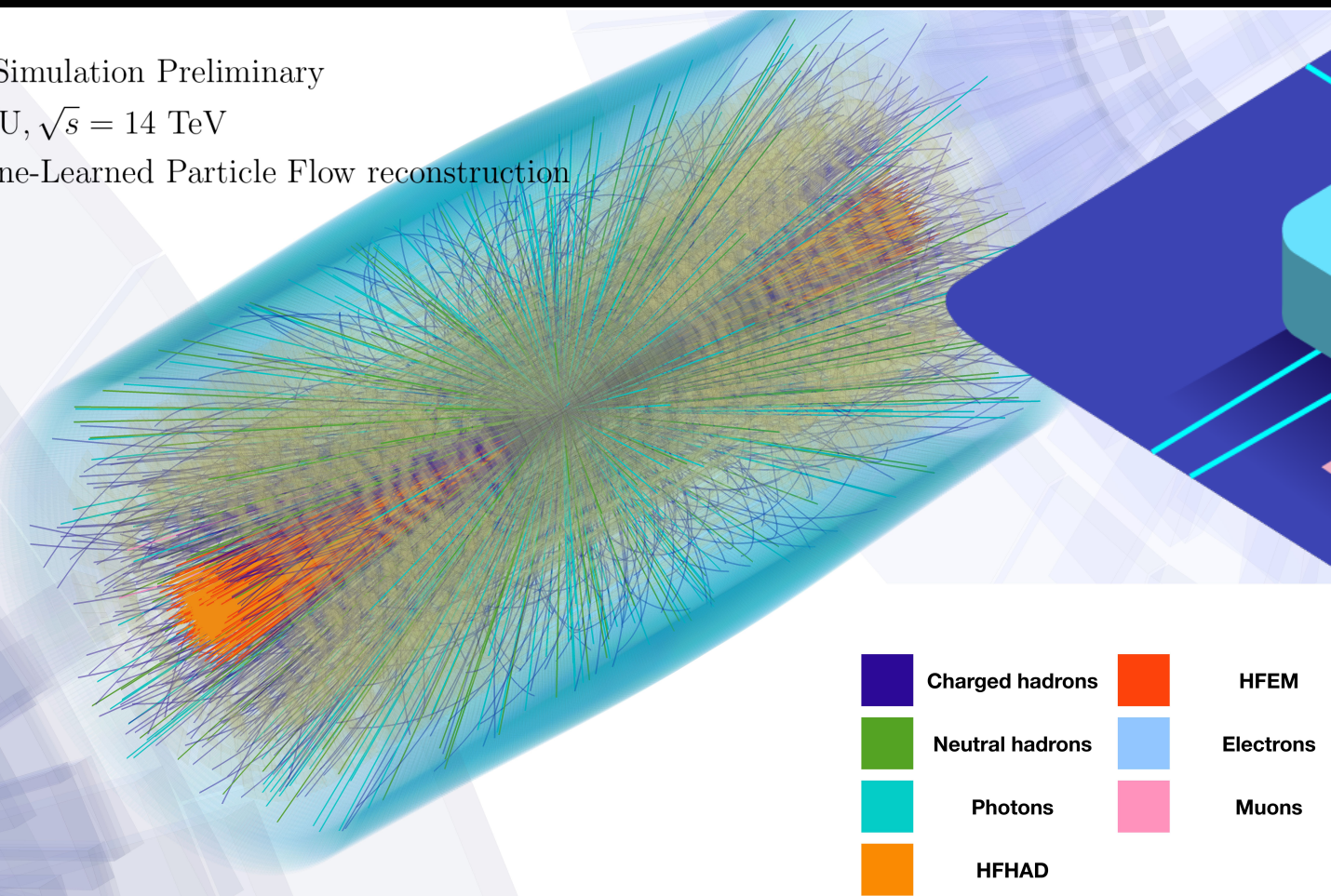
hls 4 ml



Join us at the
FastML Lab!



CMS Simulation Preliminary
 $t\bar{t} + \text{PU}, \sqrt{s} = 14 \text{ TeV}$
Machine-Learned Particle Flow reconstruction



Backup

Benchmarking

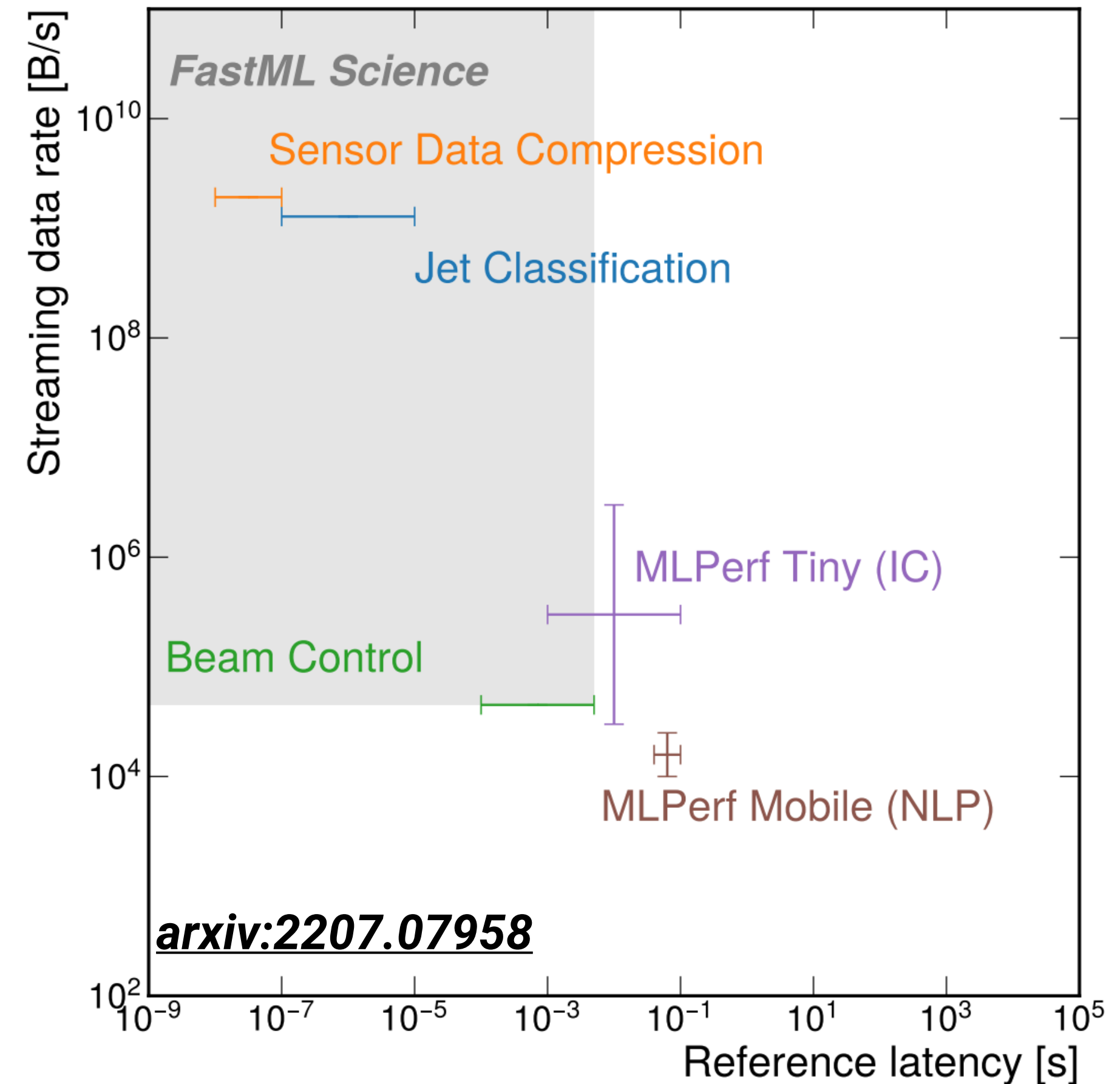
Datasets: Common FastML Science Benchmarking datasets

- guide design of edge ML hardware and software for sub-microsecond inference!

Algorithms: hls4ml-FINN benchmarked in MLPerf™

- how fast systems can process inputs and produce results
- efficient and low-latency FPGA solutions with hls4ml and FINN

Consistently competitive (QKeras+hls4ml, semantic segmentation, MLPerf)



<https://mlcommons.org/en/inference-tiny-07/>

Model	LUT		LUTRAM		FF		BRAM [36 kb]		DSP		Latency [ms]	Energy/inf. [μ J]
Pynq-Z2												
Available	53 200		17 400		106 400		140		220		–	–
IC (hls4ml)	28 544	53.7%	3 756	21.6%	49 215	46.3%	42	30.0%	4	1.8%	27.3	44 330
IC (FINN)	24 502	46.1%	2 086	12.0%	34 354	32.3%	100	71.4%	0	0.0%	1.5	2 535
AD	40 658	76.4%	3 659	21.0%	51 879	48.8%	14.5	10.4%	205	93.2%	0.019	30.1
KWS	33 732	63.4%	1 033	5.9%	34 405	32.3%	37	26.4%	1	0.5%	0.017	30.9

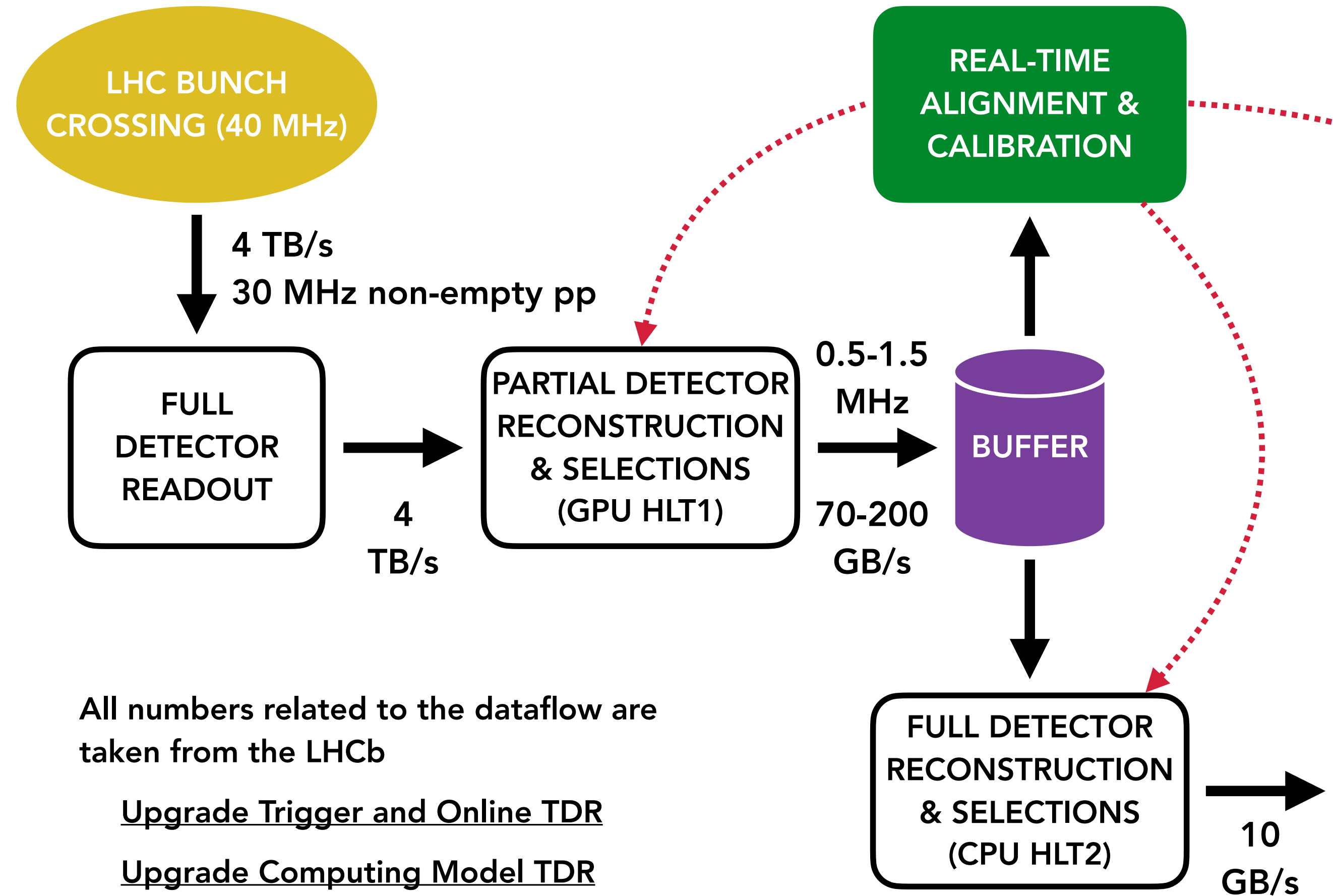
arxiv:2103.05579

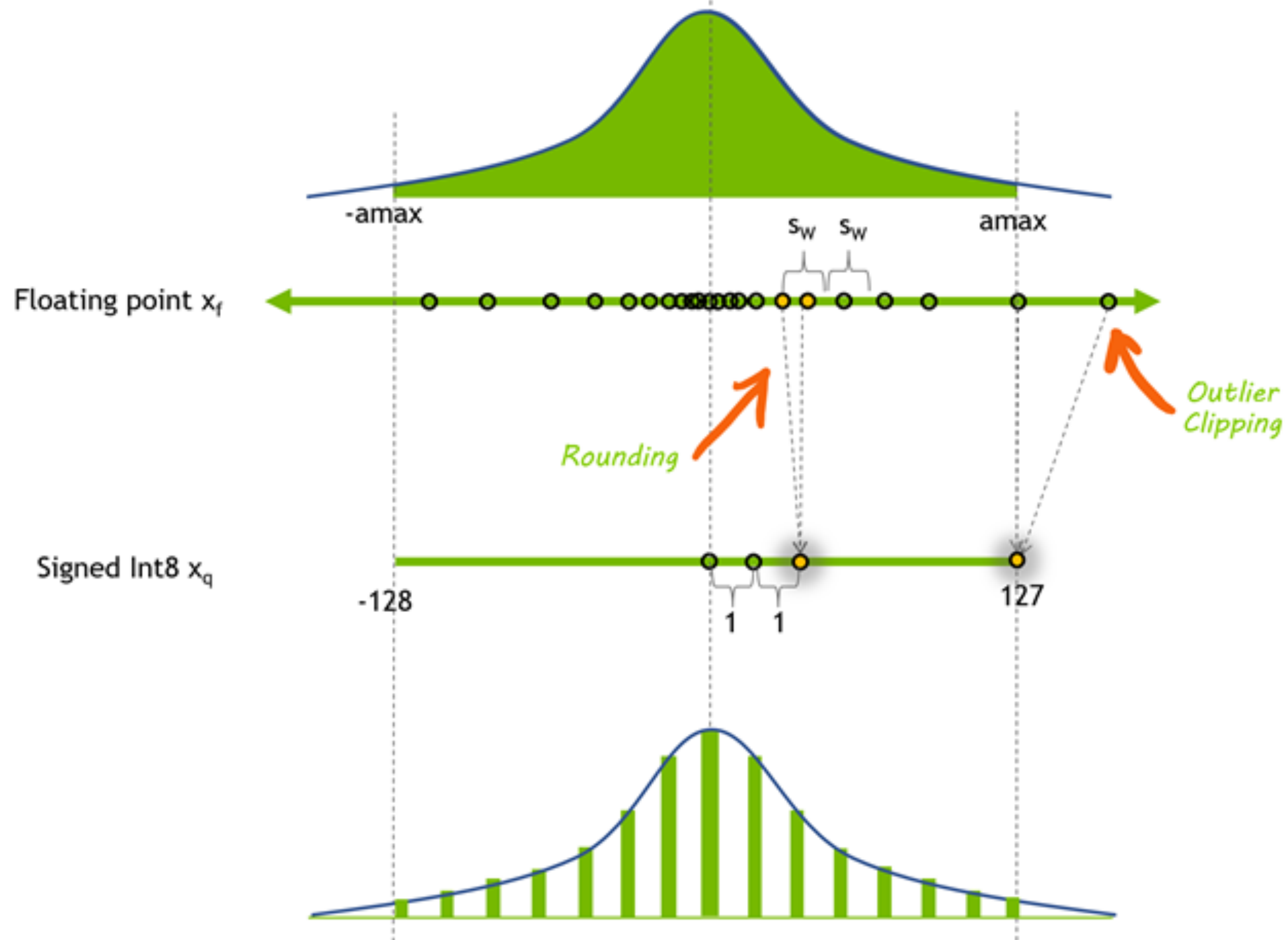
LHCb GPU trigger

Full GPU reconstruction @ 4% internet traffic

- 326 GPUs, 60 kHz per GPU

Characteristics of LHCb HLT1	Characteristics of GPUs
Intrinsically parallel problem: - Run events in parallel - Reconstruct tracks in parallel	Good for - Data-intensive parallelizable applications - High throughput applications
Huge compute load	Many TFLOPS
Full data stream from all detectors is read out → no stringent latency requirements	Higher latency than CPUs, not as predictable as FPGAs
Small raw event data (~100 kB)	Connection via PCIe → limited I/O bandwidth
Small event raw data (~100 kB)	Thousands of events fit into O(10) GB of memory





QPYTORCH ?

<https://github.com/Xilinx/brevitas>

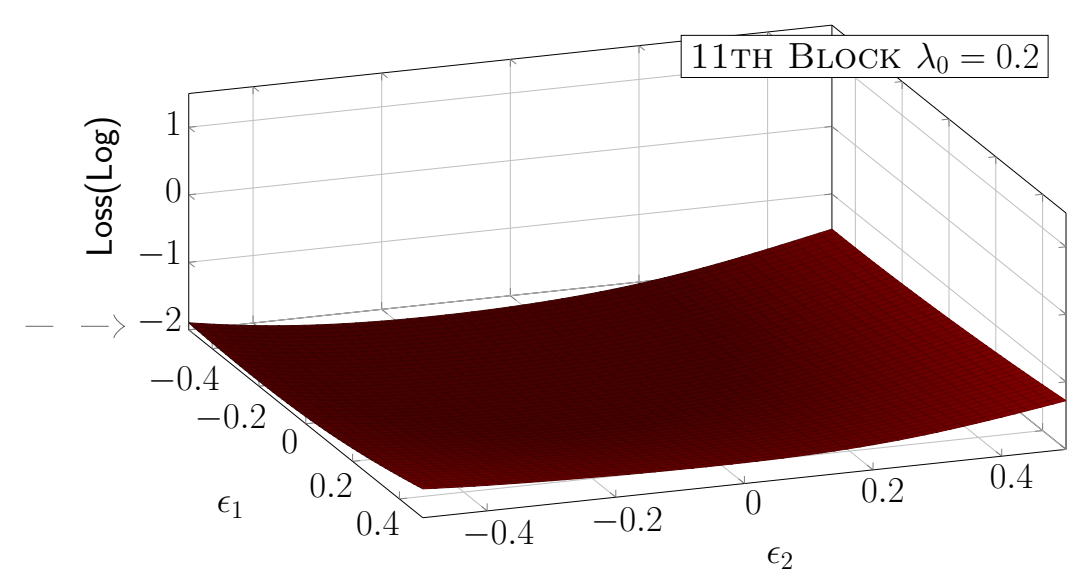
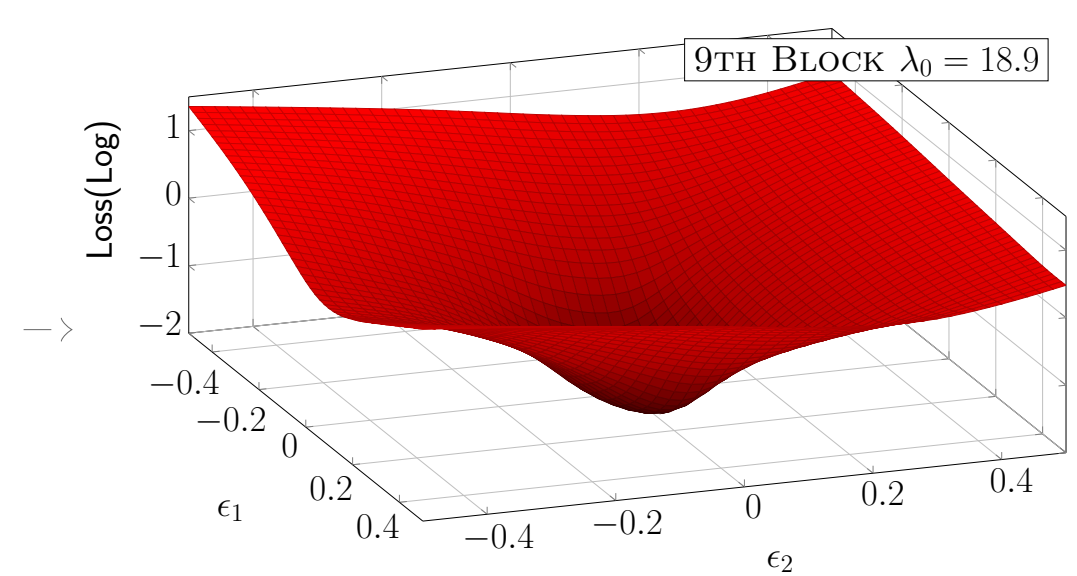
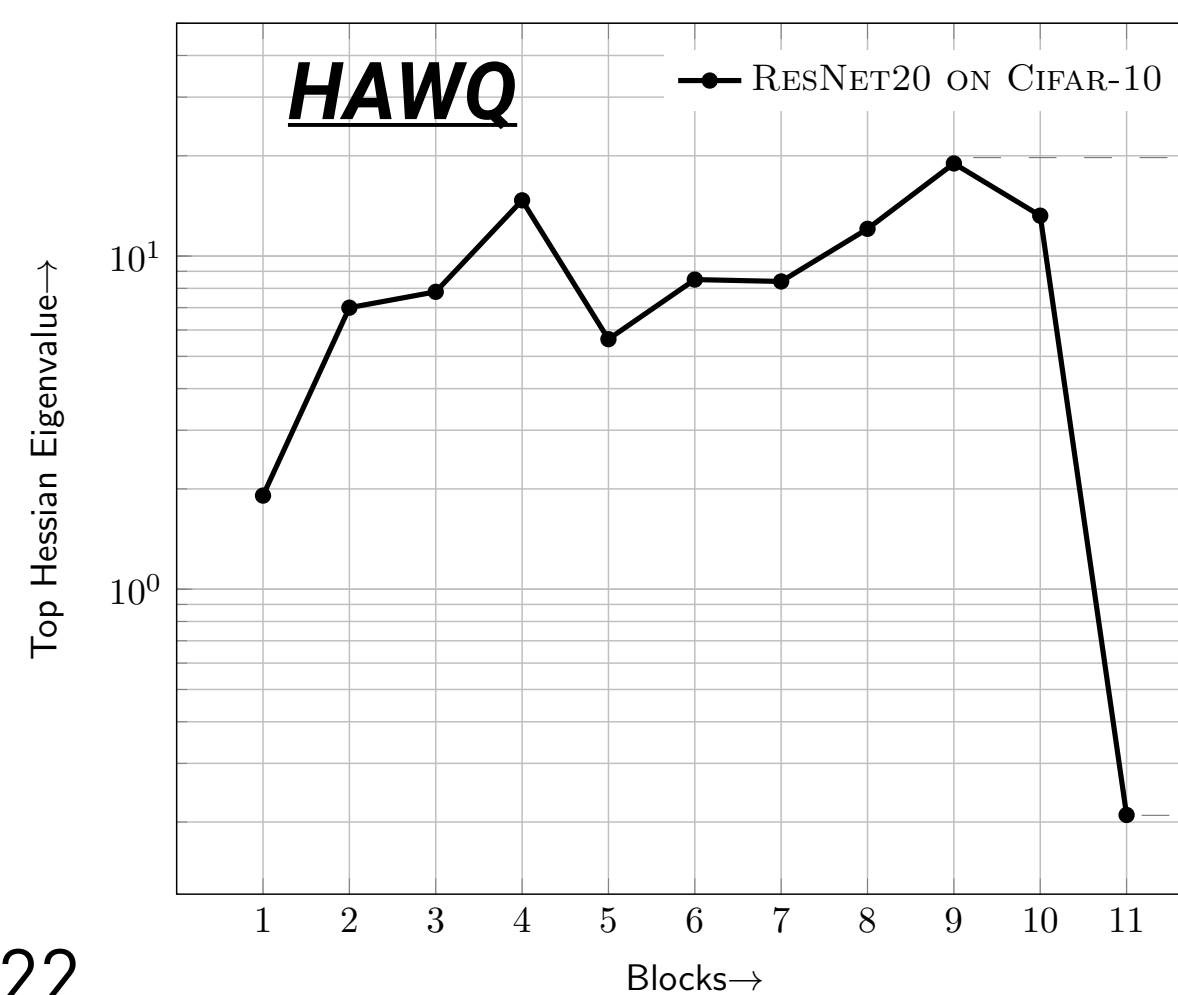
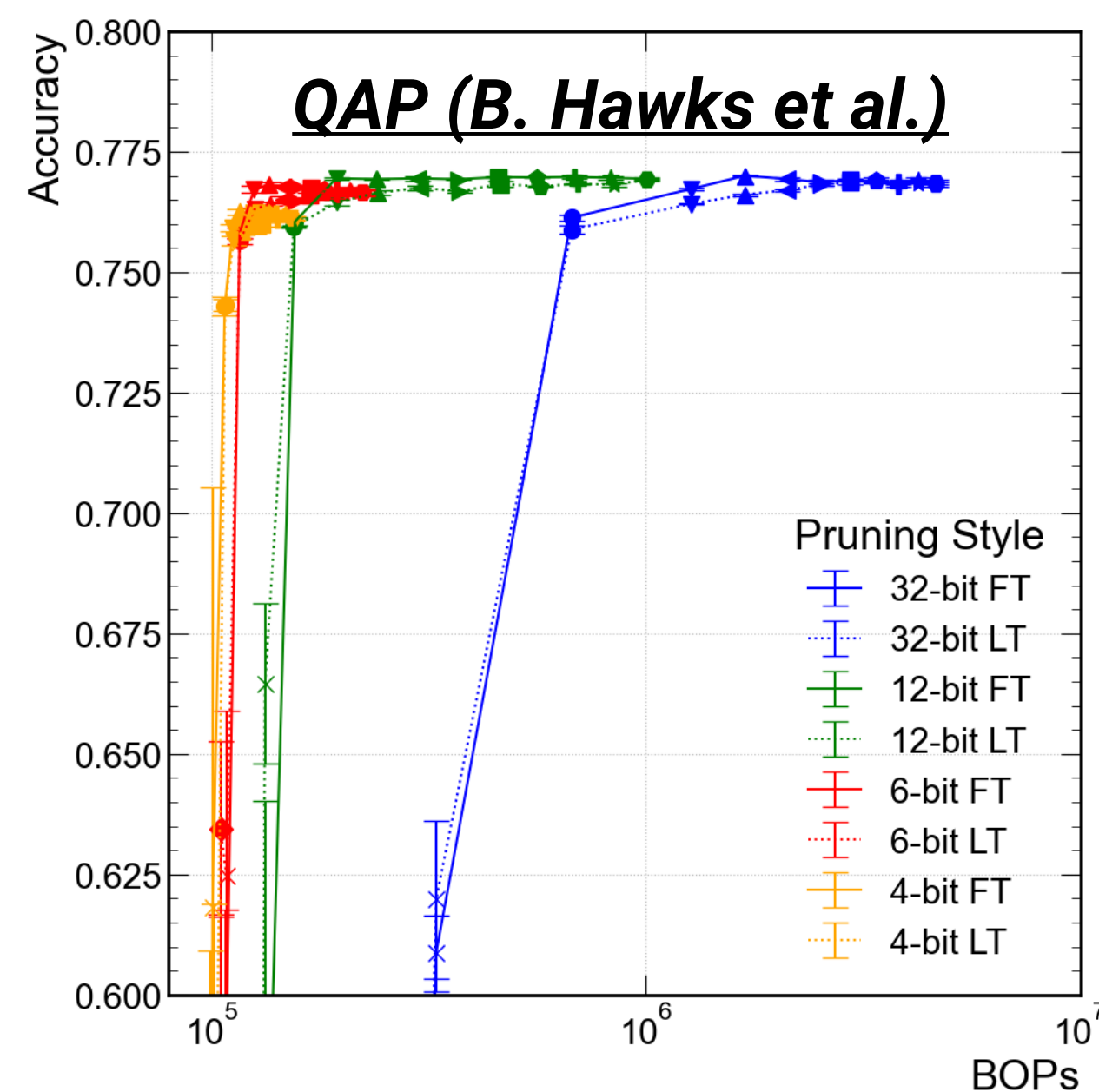
Brevitas like QKeras, but for PyTorch

- QAT library
- Support most common layers and activation functions

Other quantization techniques:

- **HAWQ: Hessian AWare Quantization**
- **Quantization Aware Pruning (B. Hawks et al.)**

```
import brevitas.nn as qnn
qnn.|
├── quant_bn (brevitas.nn)
├── QuantCat (brevitas.nn.quant_eltwise)
├── QuantTanh (brevitas.nn.quant_activation)
├── ScaleBias (brevitas.nn.quant_scale_bias)
├── quant_conv (brevitas.nn)
├── hadamard_classifier (brevitas.nn)
├── quant_accumulator (brevitas.nn)
├── quant_activation (brevitas.nn)
├── quant_avg_pool (brevitas.nn)
├── quant_convtranspose (brevitas.nn)
├── quant_dropout (brevitas.nn)
├── quant_eltwise (brevitas.nn)
├── quant_linear (brevitas.nn)
├── quant_max_pool (brevitas.nn)
├── quant_scale_bias (brevitas.nn)
├── quant_upsample (brevitas.nn)
├── BatchNorm1dToQuantScaleBias (brevitas.nn.quant_bn)
├── BatchNorm2dToQuantScaleBias (brevitas.nn.quant_bn)
├── ClampQuantAccumulator (brevitas.nn.quant_accumulator)
```

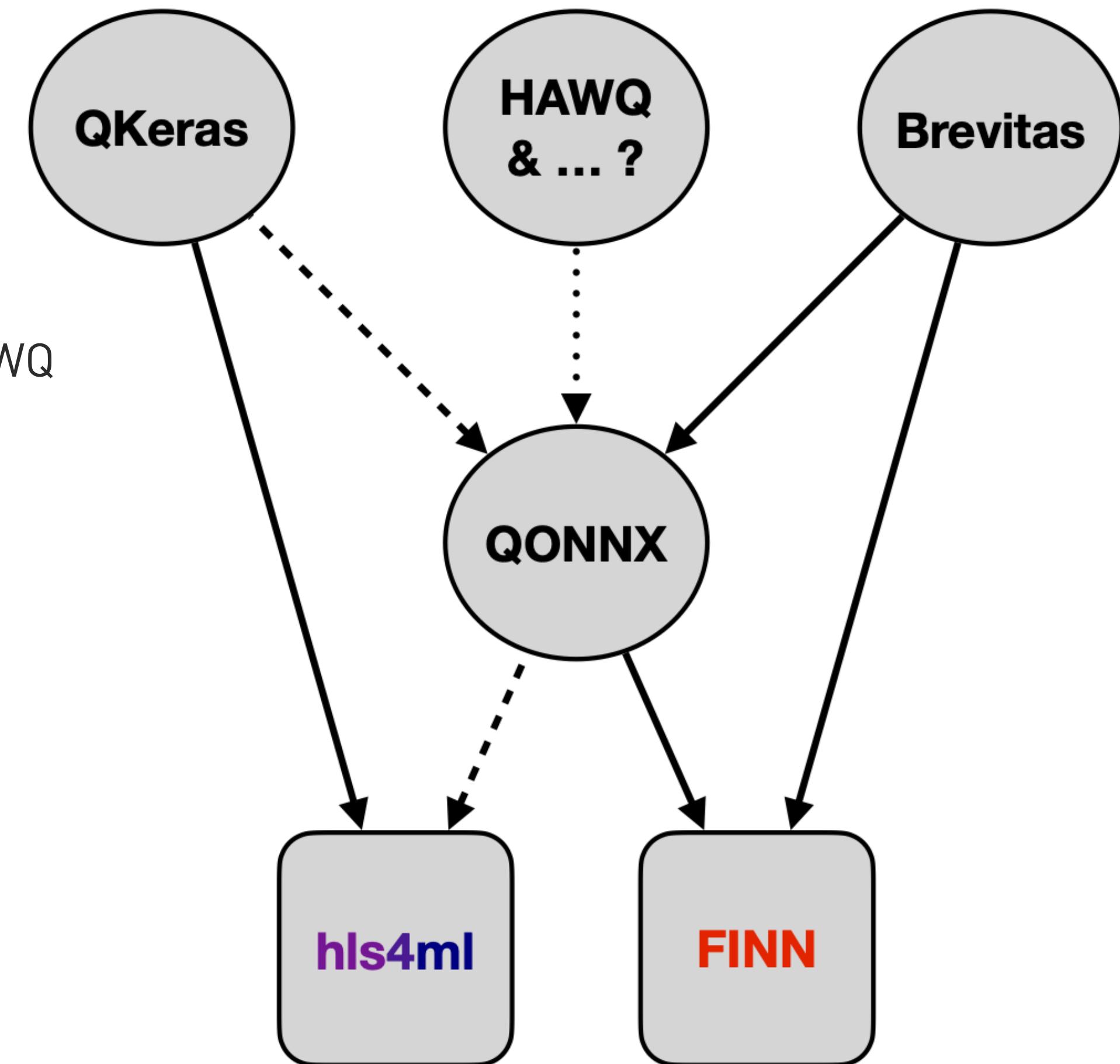


QPYTORCH ?

Quantized ONNX (QONNX), J. Mitrevski et. al

hls4ml collaborate with Xilinx Research Labs to develop QOONX

- Introducing 'Quant' node to ONNX graph
- Brevitas (PyTorch) and QKeras (Keras) can export QONNX (HAWQ export in progress): then hls4ml and FINN can import QONNX

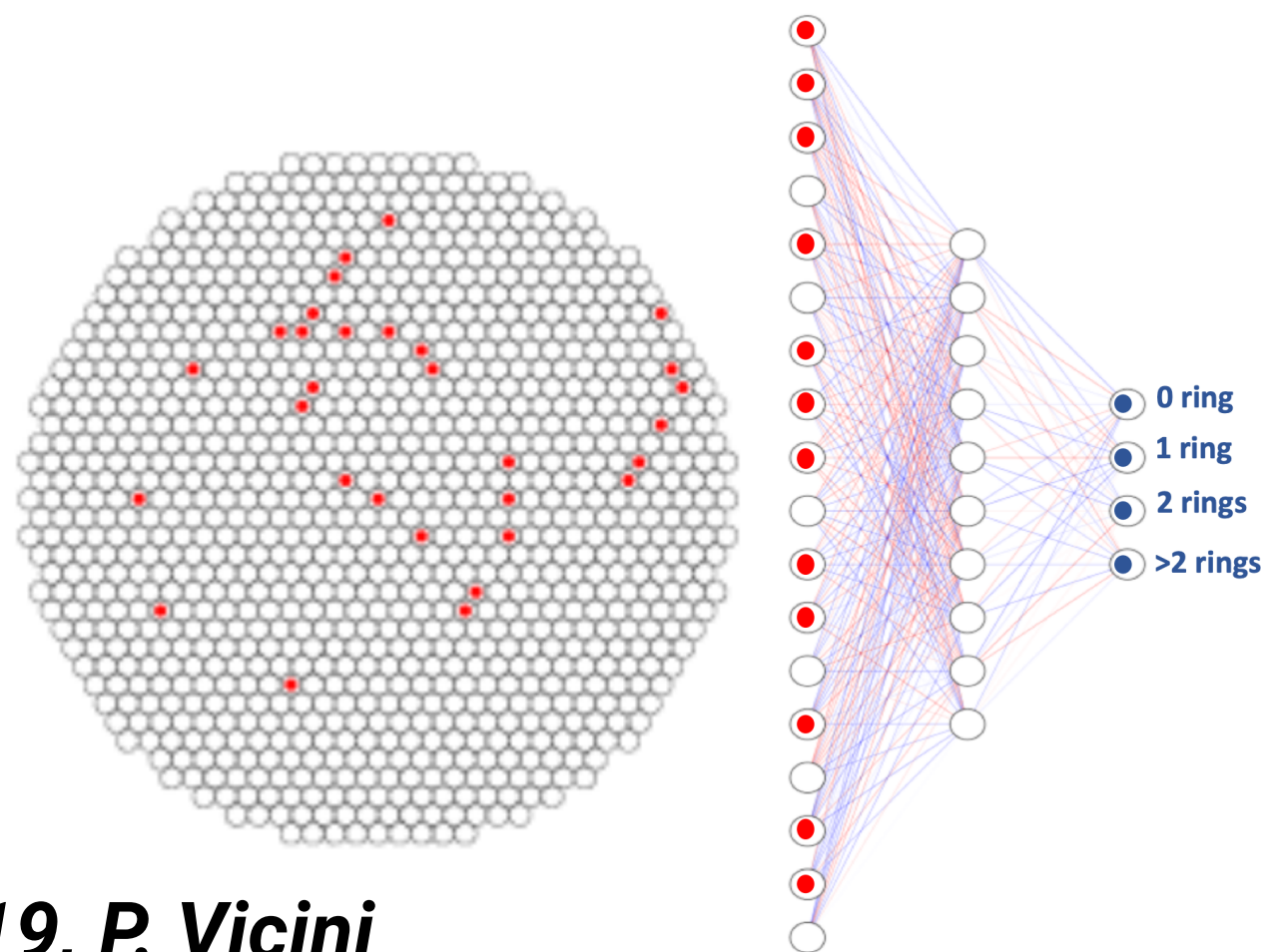


- Done
- > In progress
- > Planned

hls4ml in other CERN experiments

NA62: Measuring $BR(K^+ \rightarrow \pi^+ \nu \bar{\nu}) = O(10^{-11})$

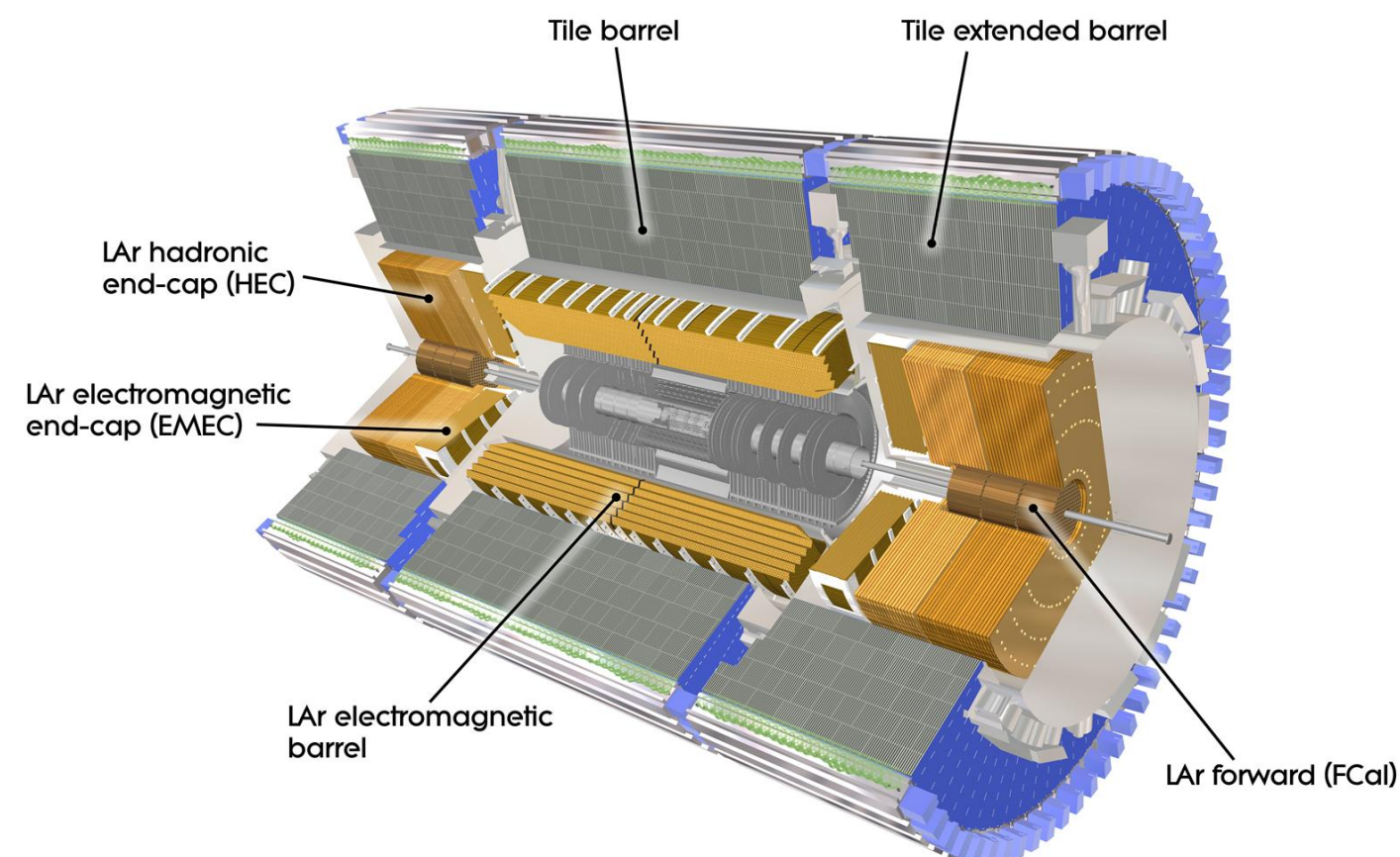
- FPGA trigger 800 MHz \rightarrow 1 MHz



CHEP 2019, P. Vicini

ATLAS Liquid Argon Calorimeter (R&D)

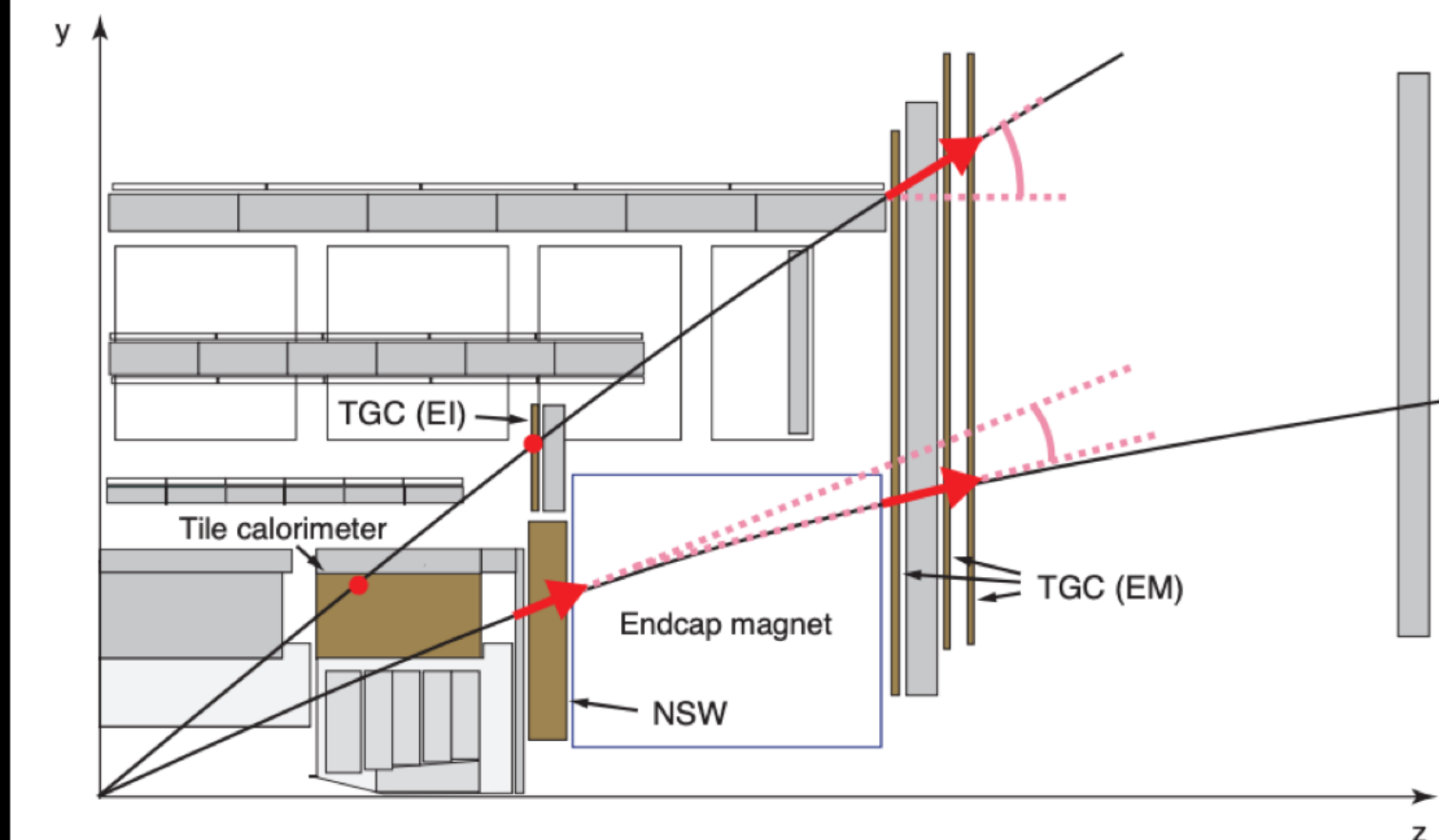
- RNN for real-time energy reconstruction
- ~200 ns on Intel Stratix-10 FPGA



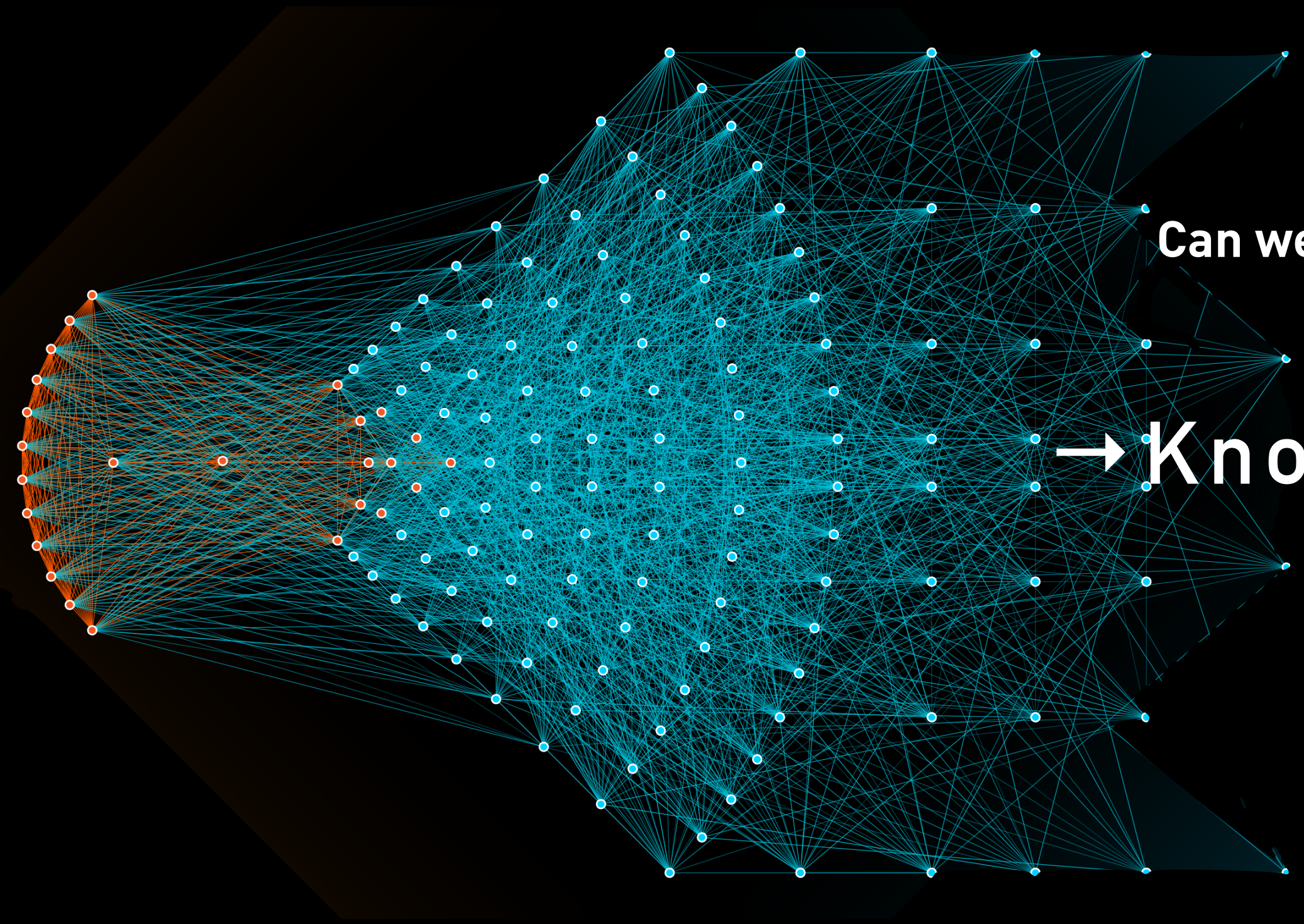
[DOI:10.1007/s41781-021-00066-y](https://doi.org/10.1007/s41781-021-00066-y)

ATLAS small wheel muon segment finding and reconstruction (R&D)

- Regression of muon position and angle
- 400 ns budget



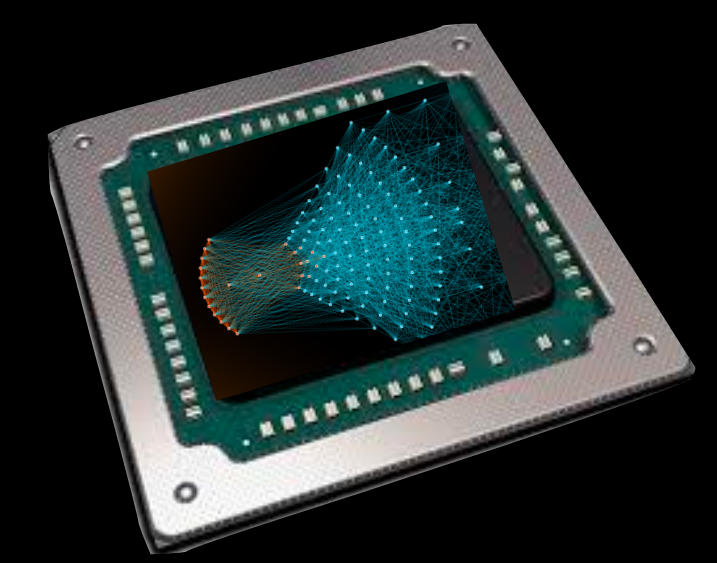
R. Teixeira de Lima, R Rojas Caballero et al.



Train

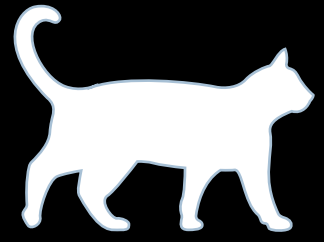
Can we have the best of both worlds?

→ Knowledge Distillation

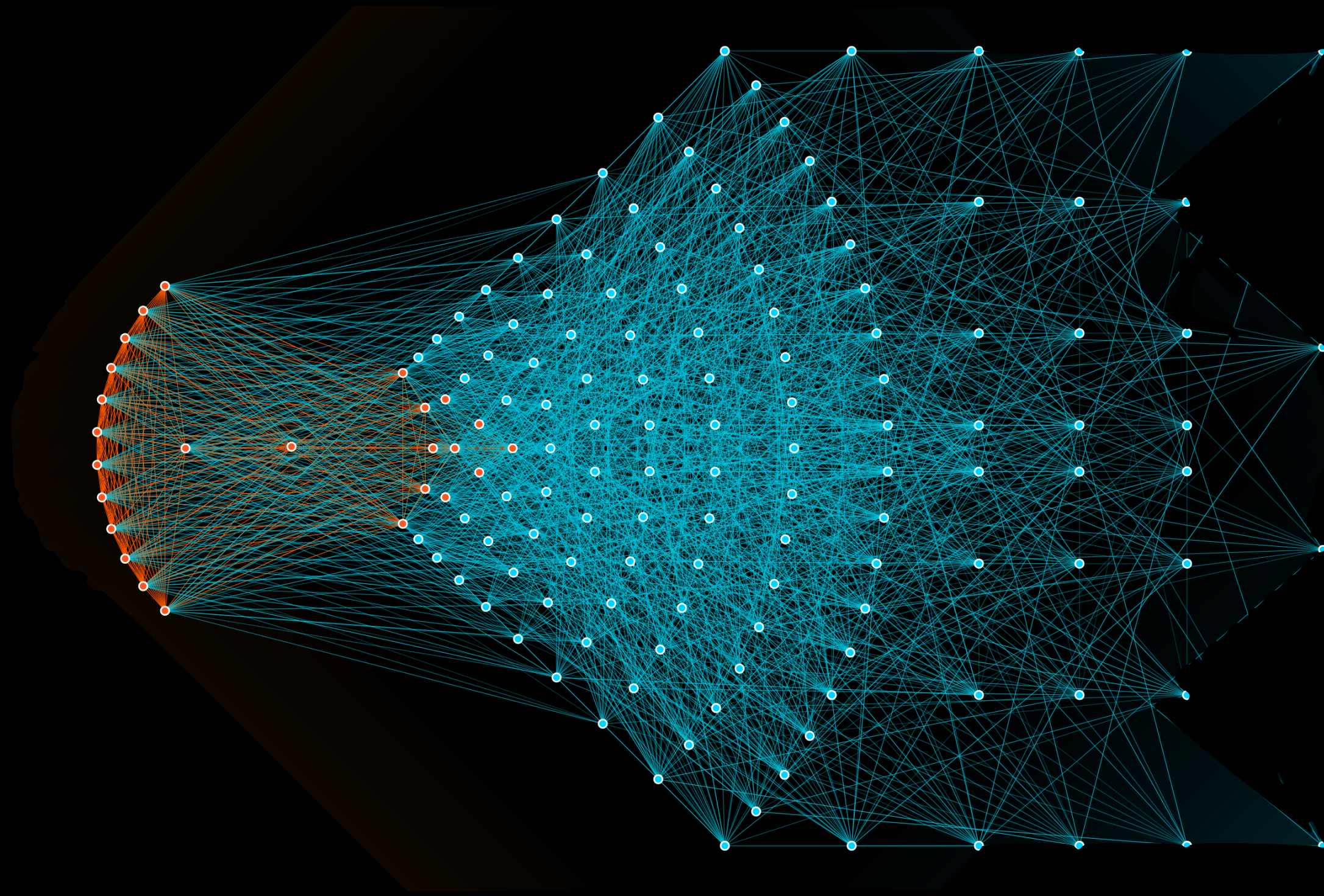
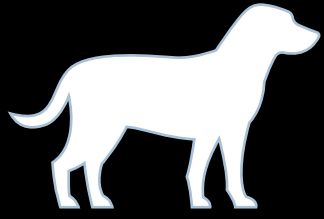


Inference

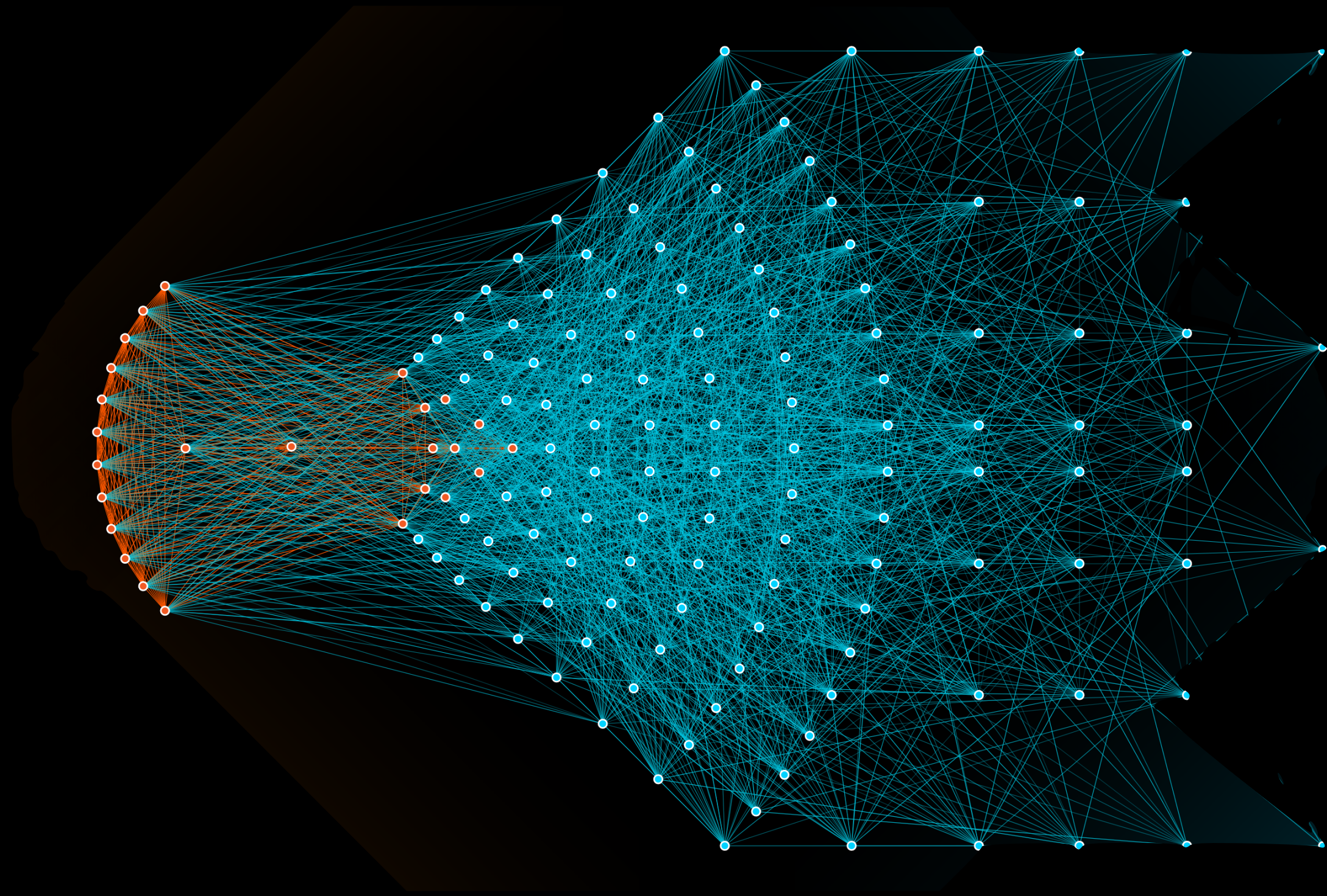
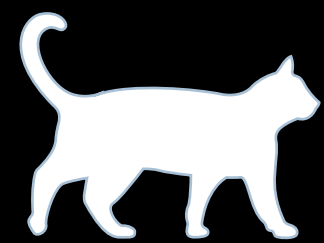
Cat



Dog



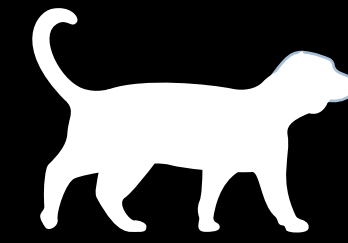
Cat



is cat

is dog

Cat

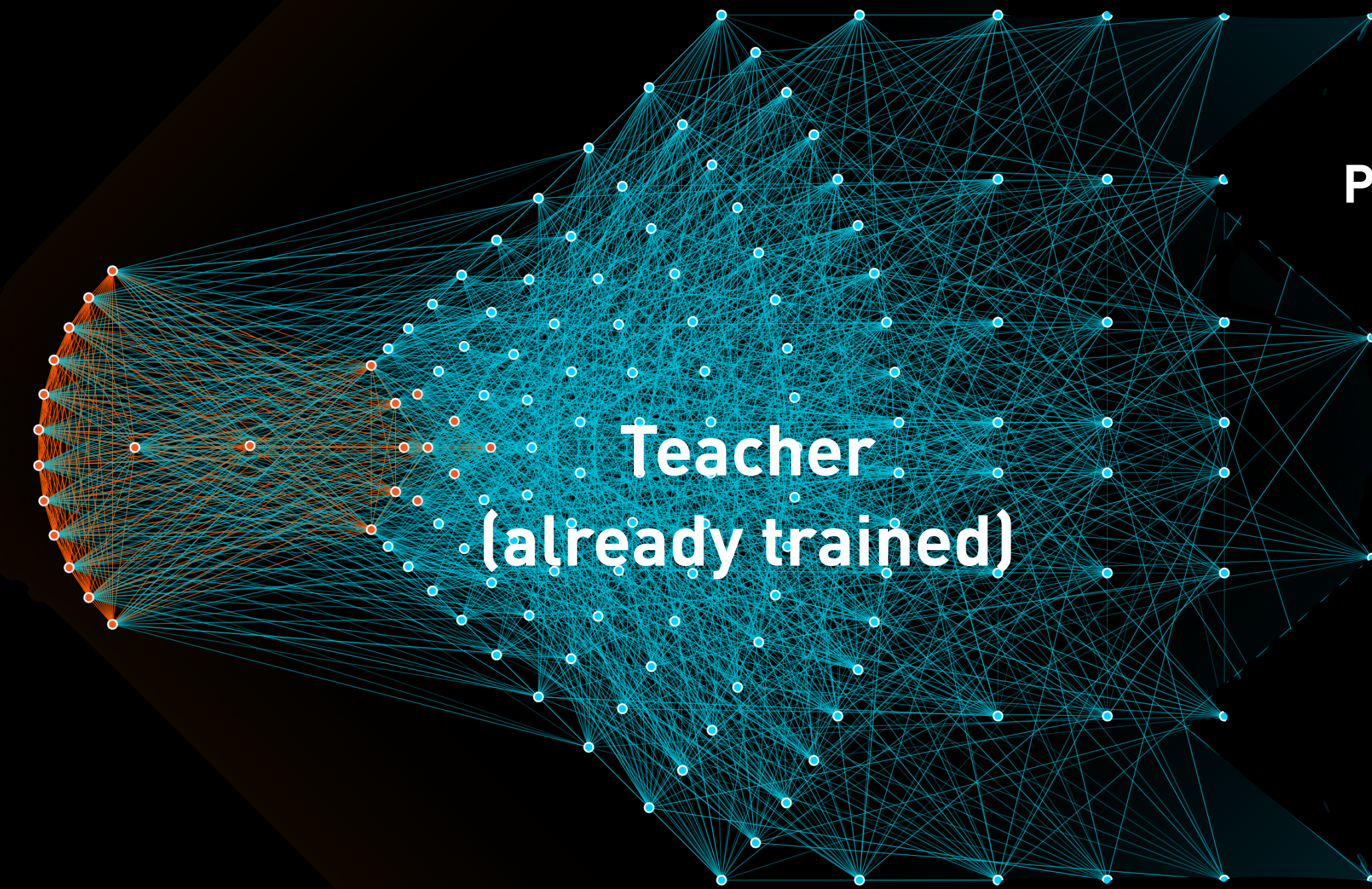


Predicted labels

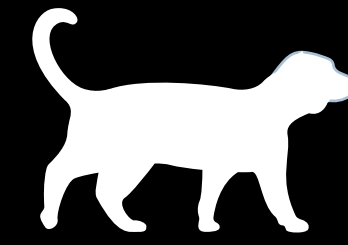
is cat = 0.89

is dog = 0.11

Teacher
(already trained)



Cat



True labels

is cat = 1

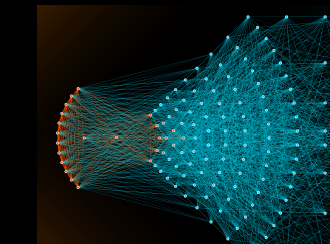
is dog = 0

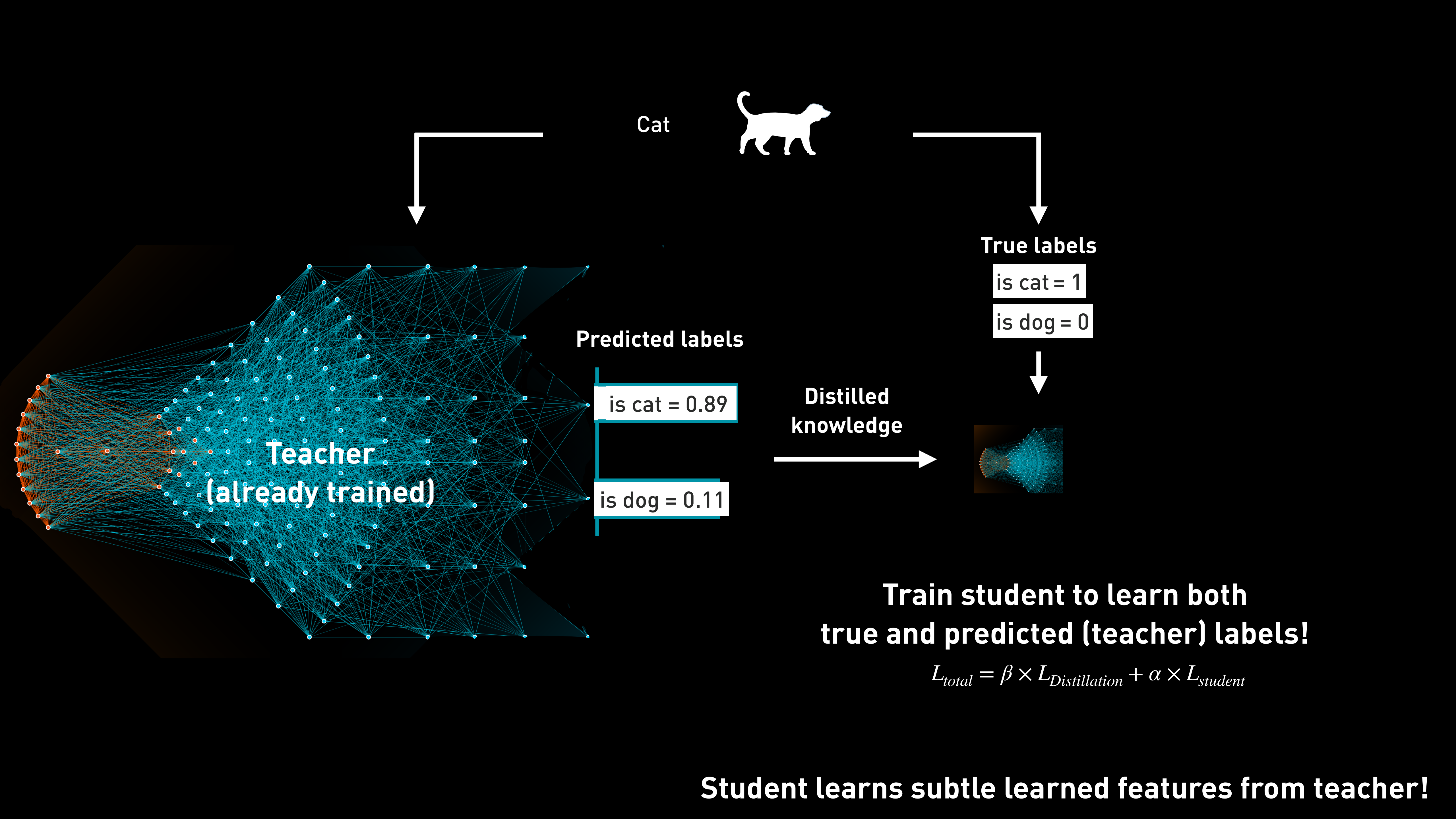
Predicted labels

is cat = 0.89

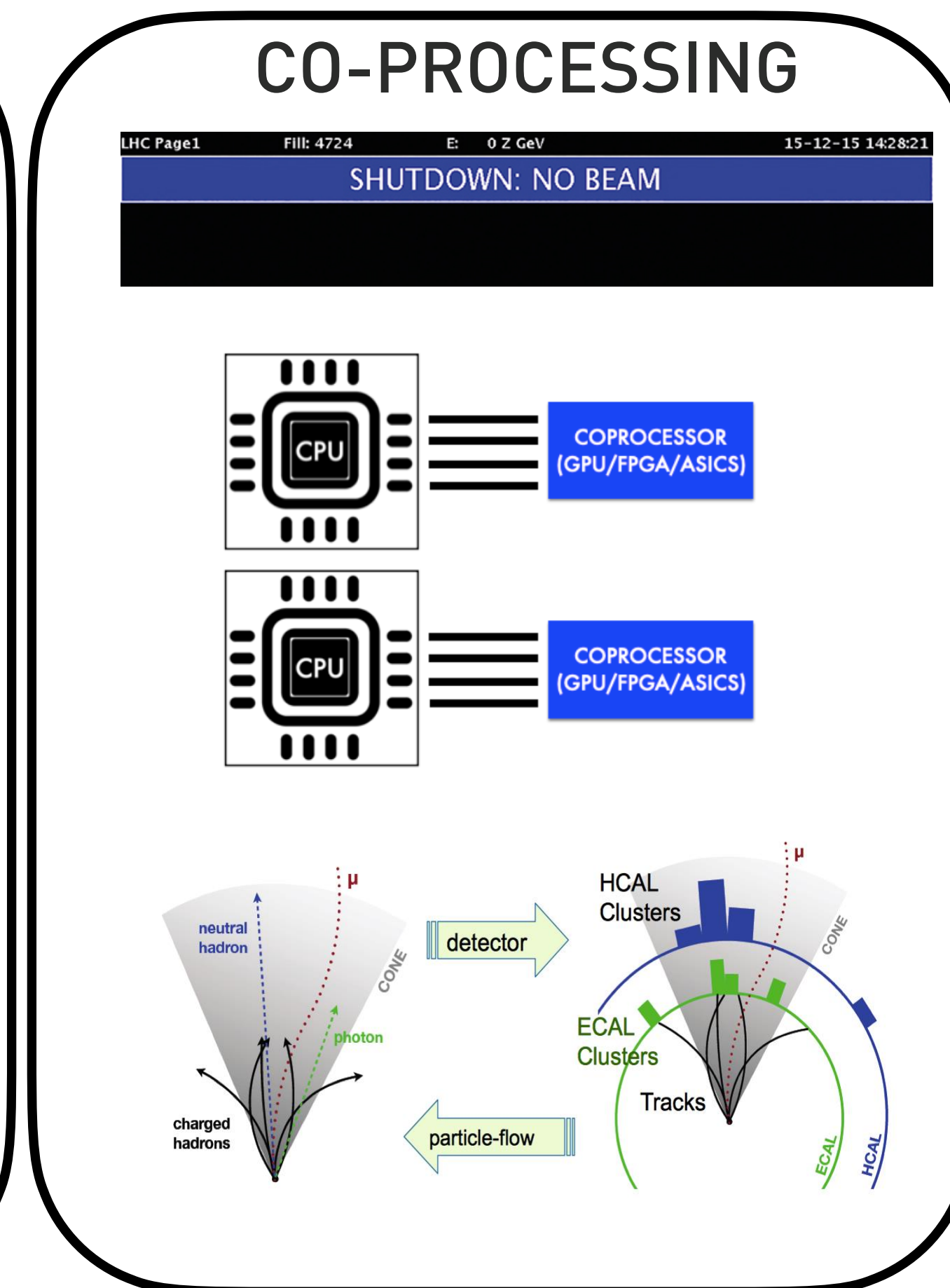
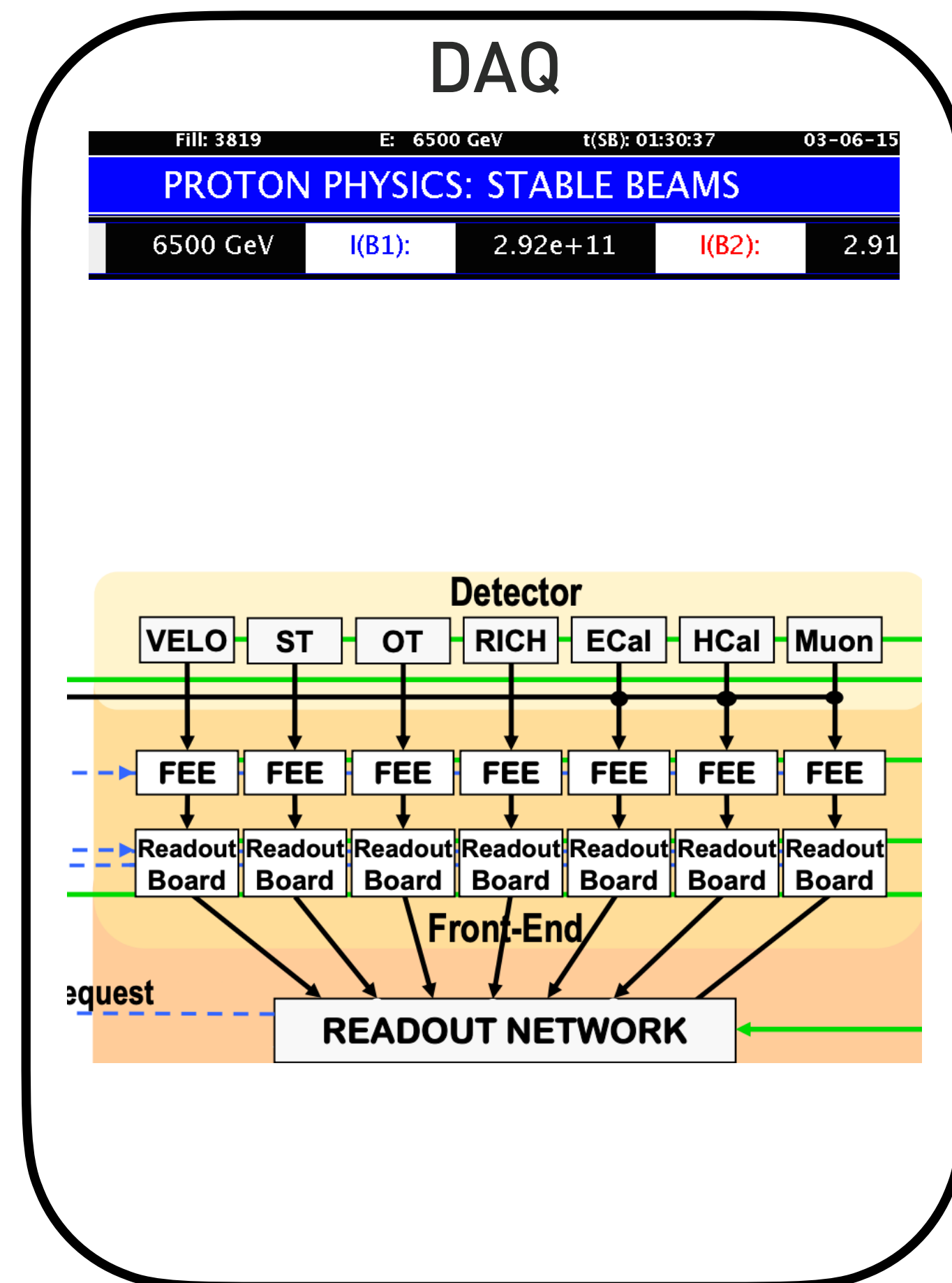
is dog = 0.11

Teacher
(already trained)



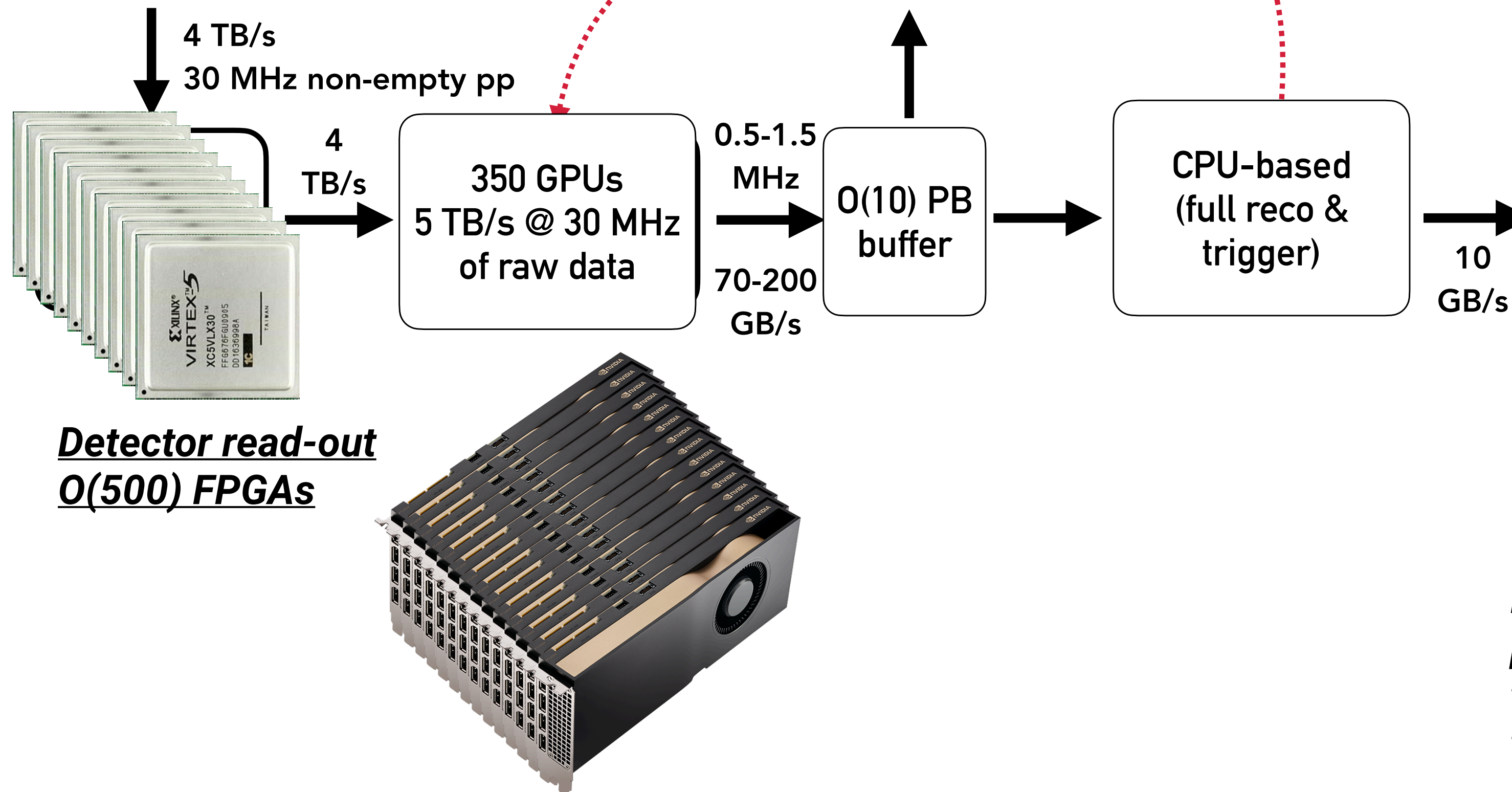
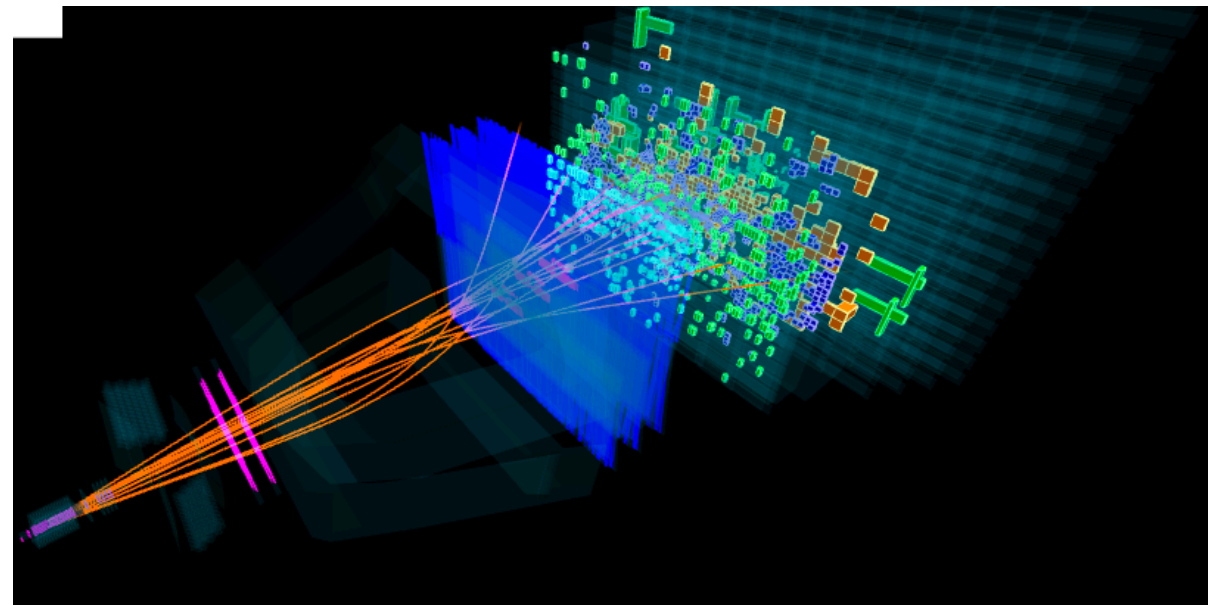


FPGAs as accelerators

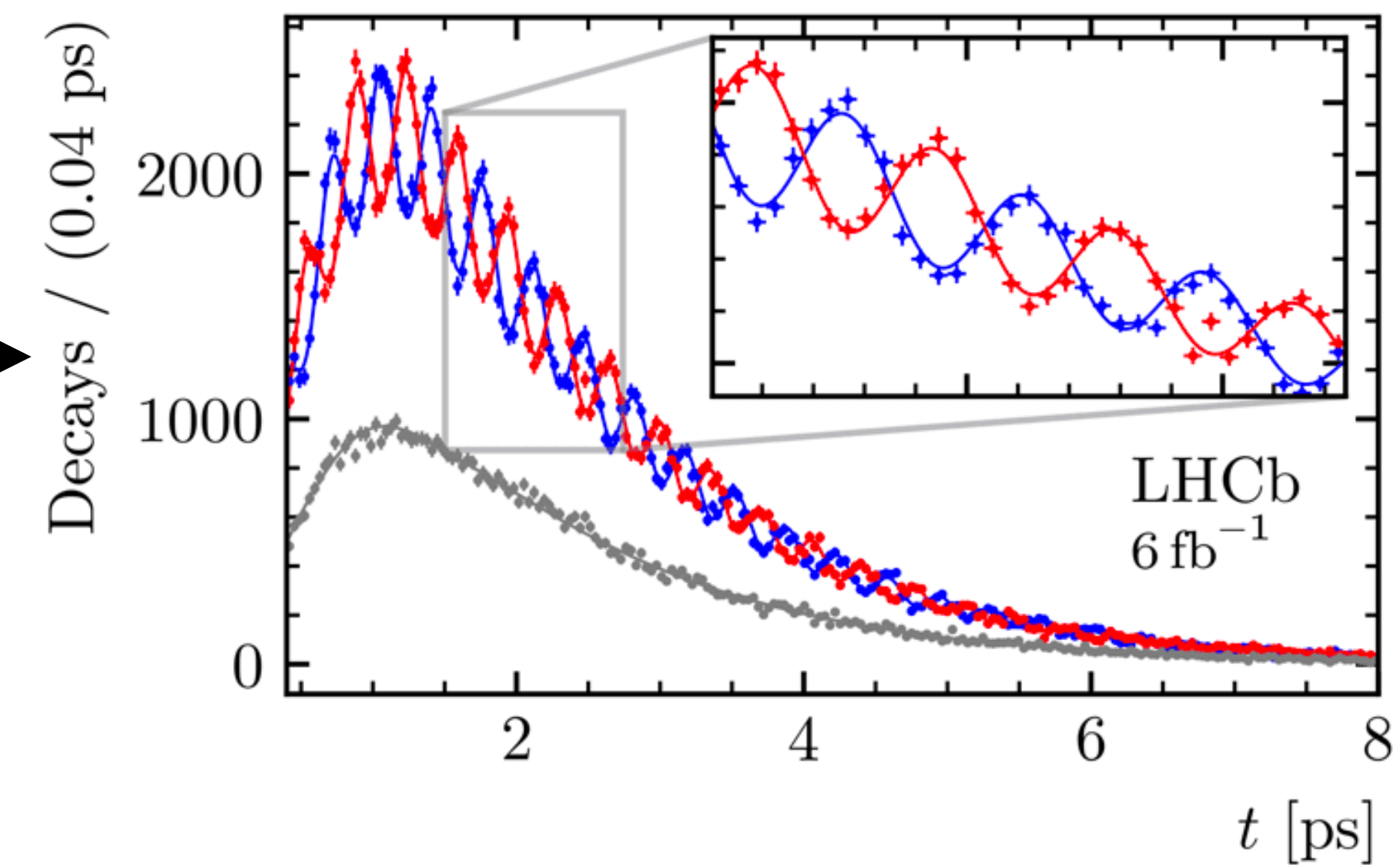


High throughput GPU triggers

40 MHz pp collisions



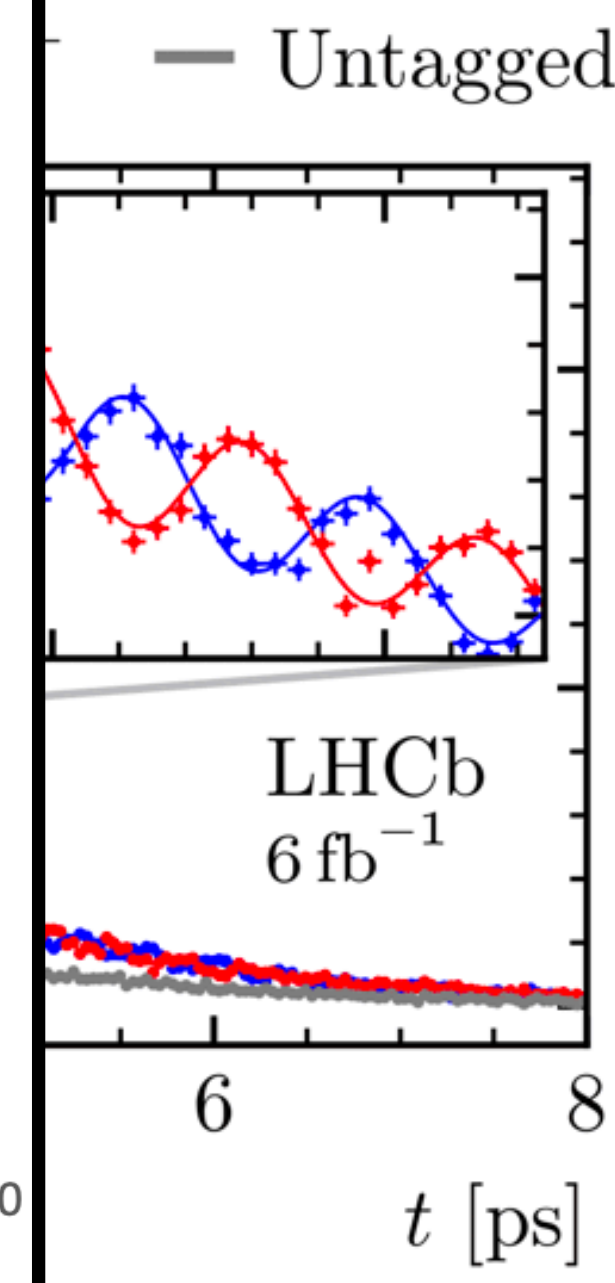
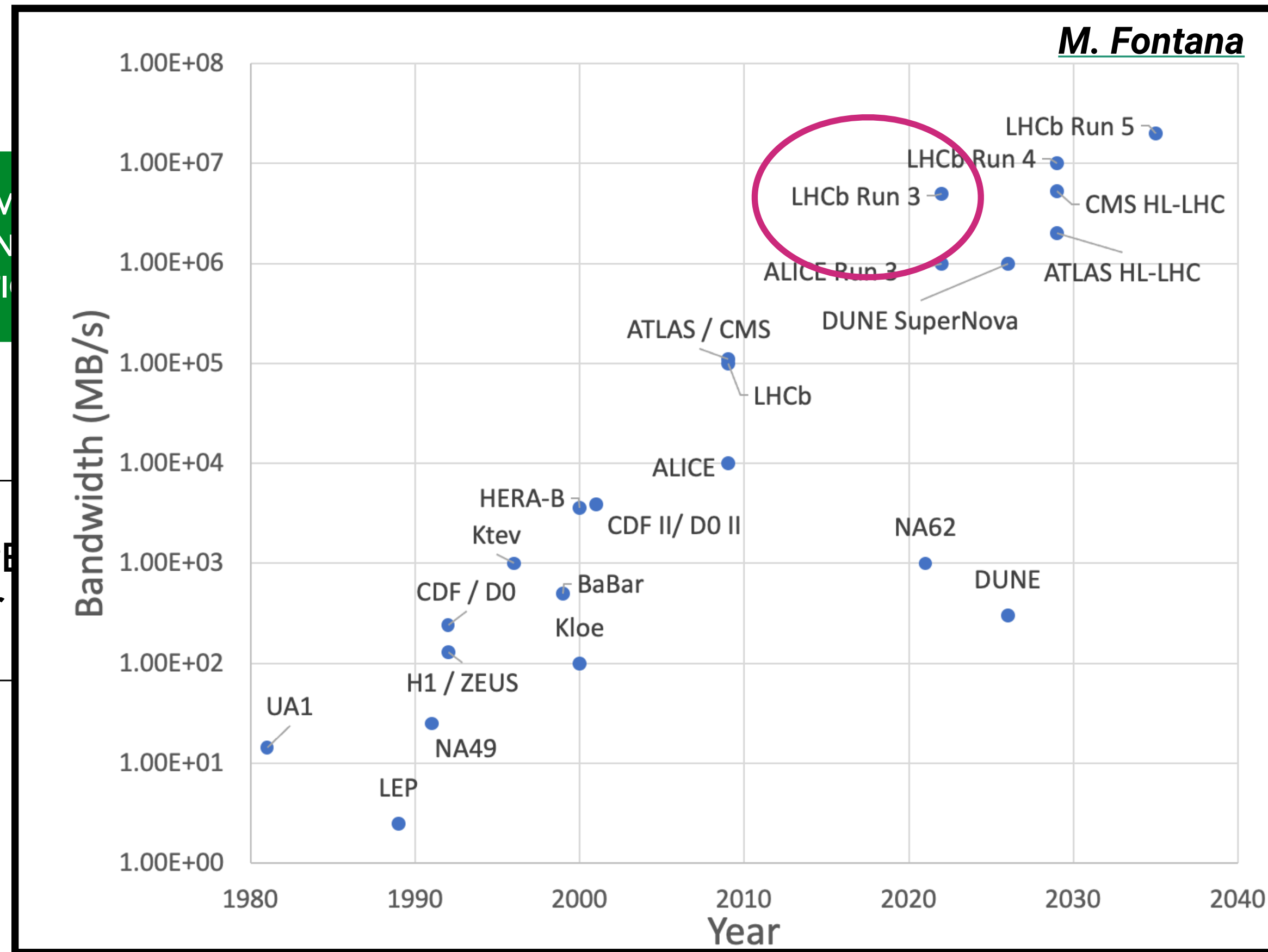
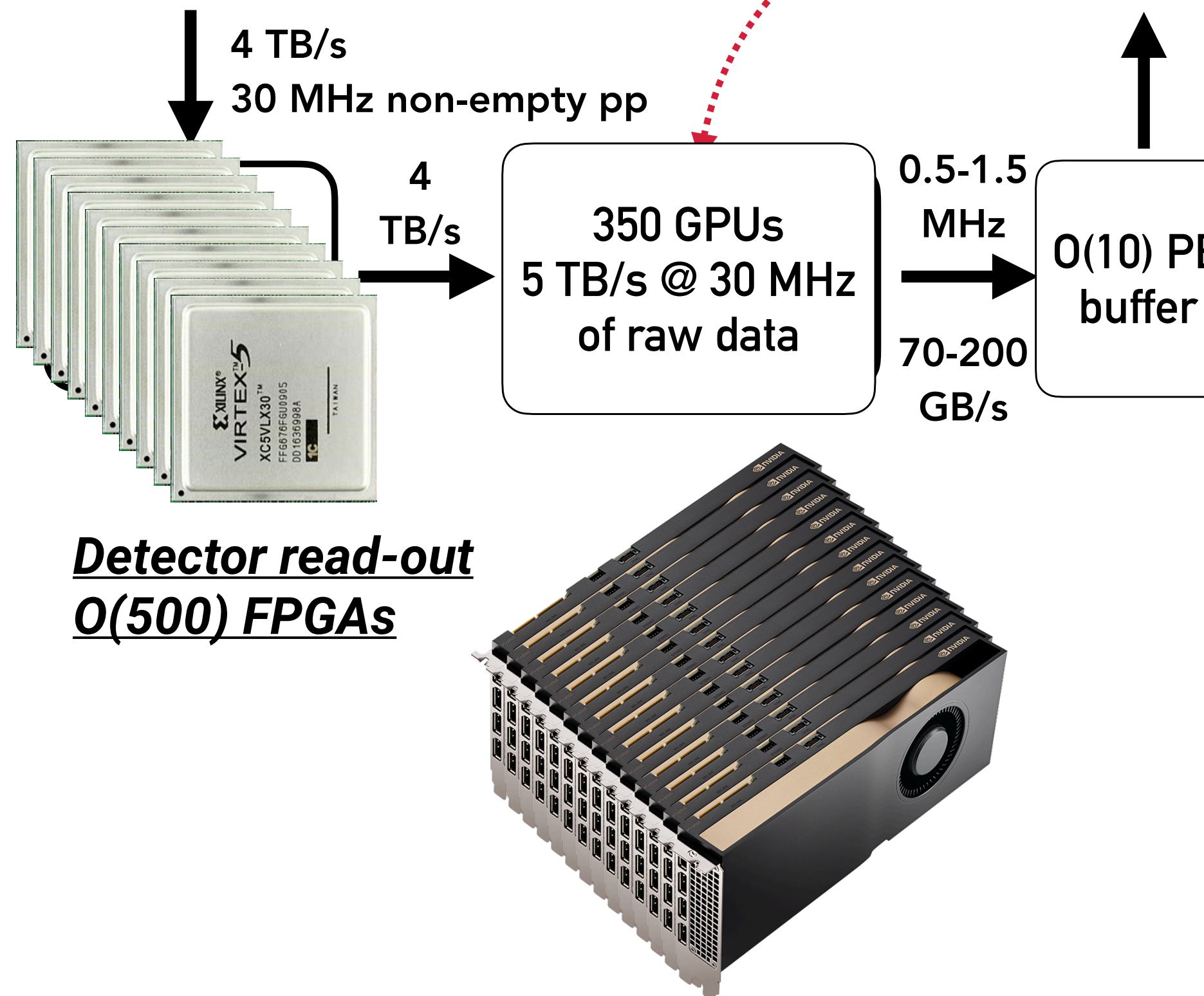
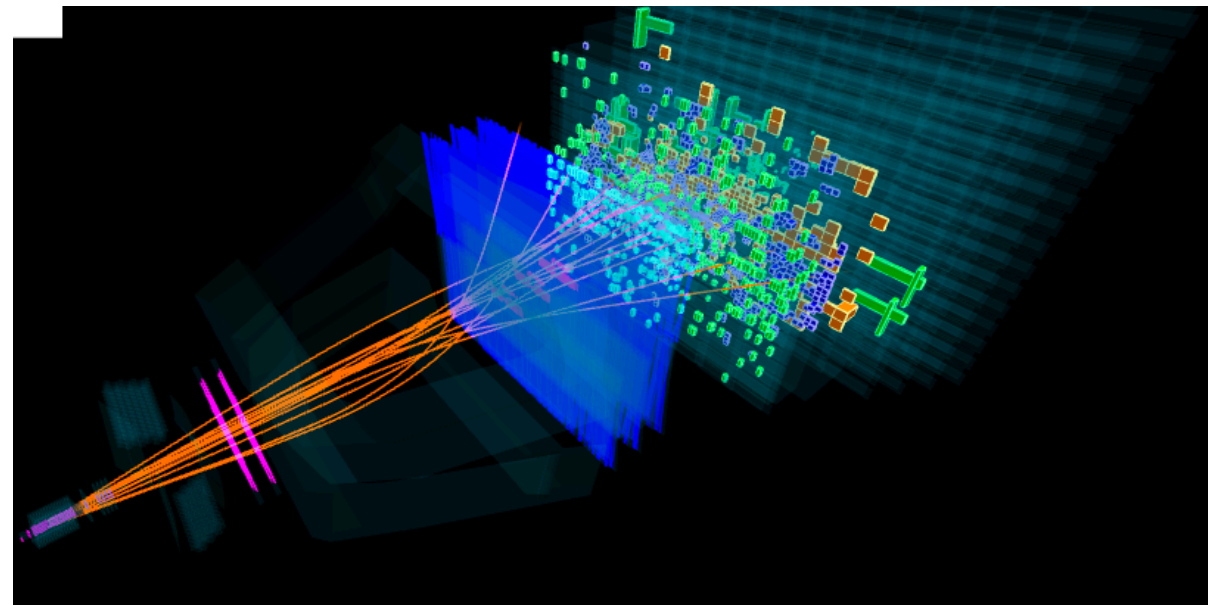
— $B_s^0 \rightarrow D_s^- \pi^+$ — $\bar{B}_s^0 \rightarrow D_s^- \pi^+$ — Untagged



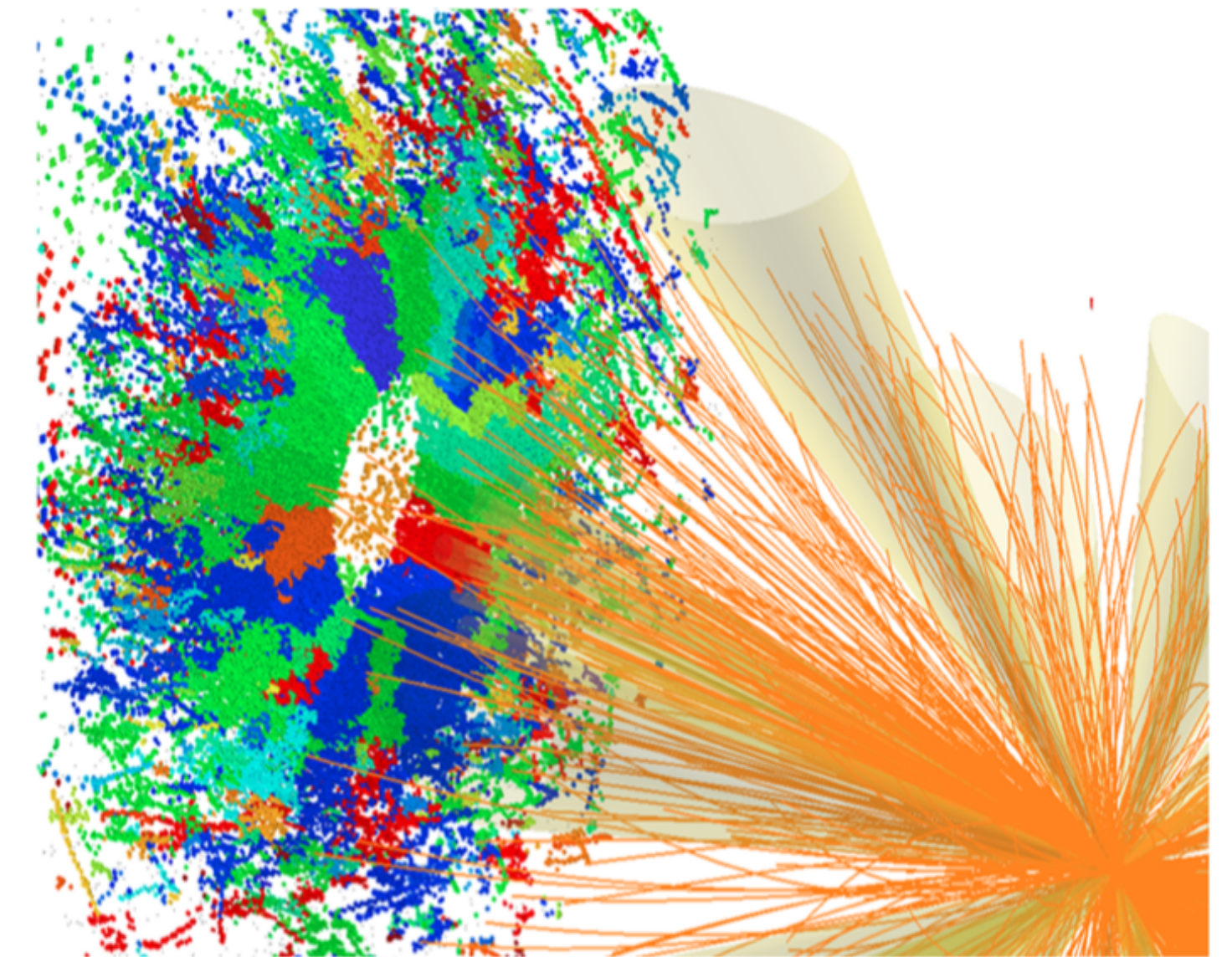
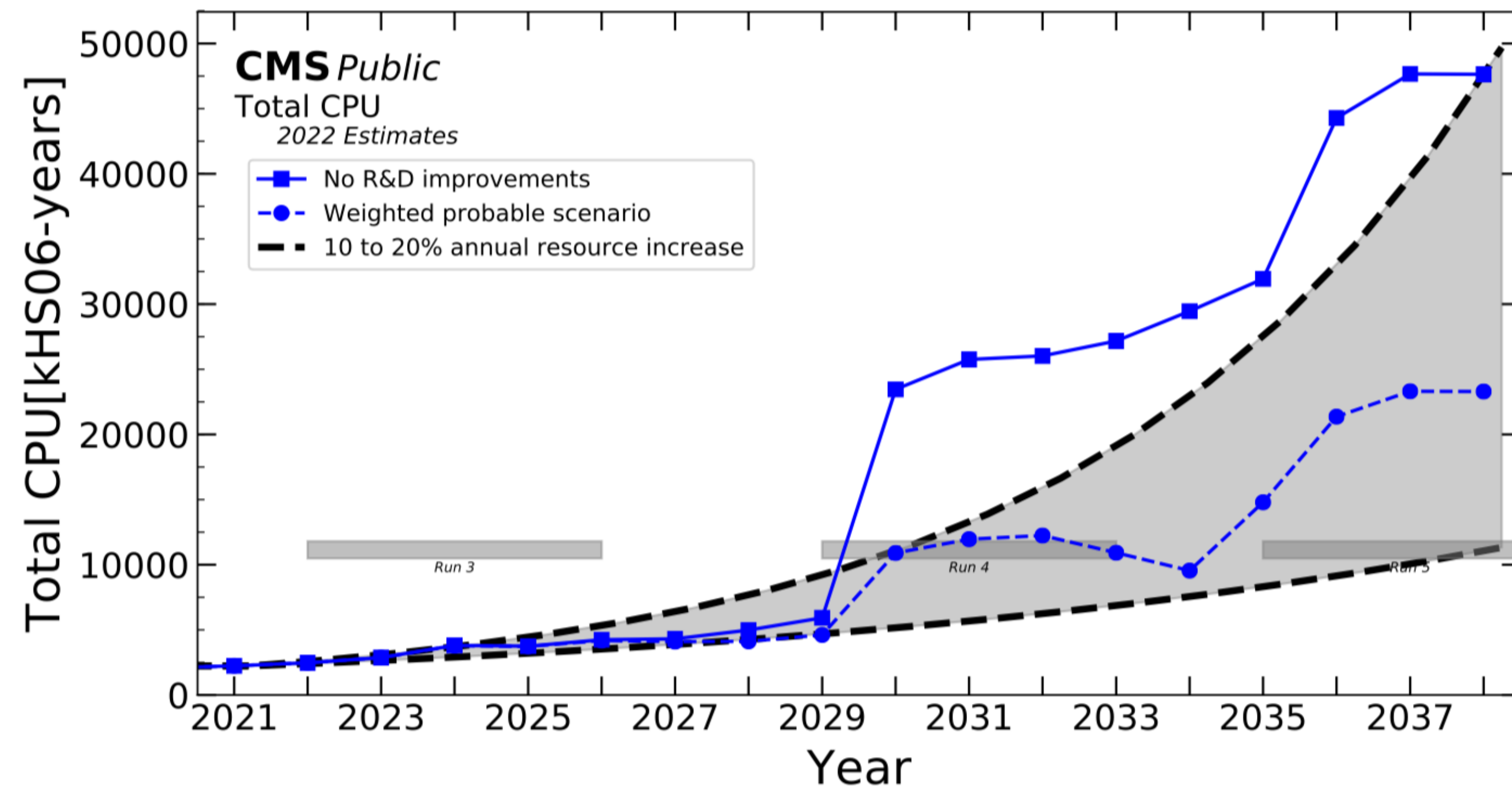
ML already used in LHCb GPU trigger for particle ID, track reconstruction, trigger decisions ... more underway!

High throughput GPU triggers

40 MHz pp collisions



CMS Offline Computing Results

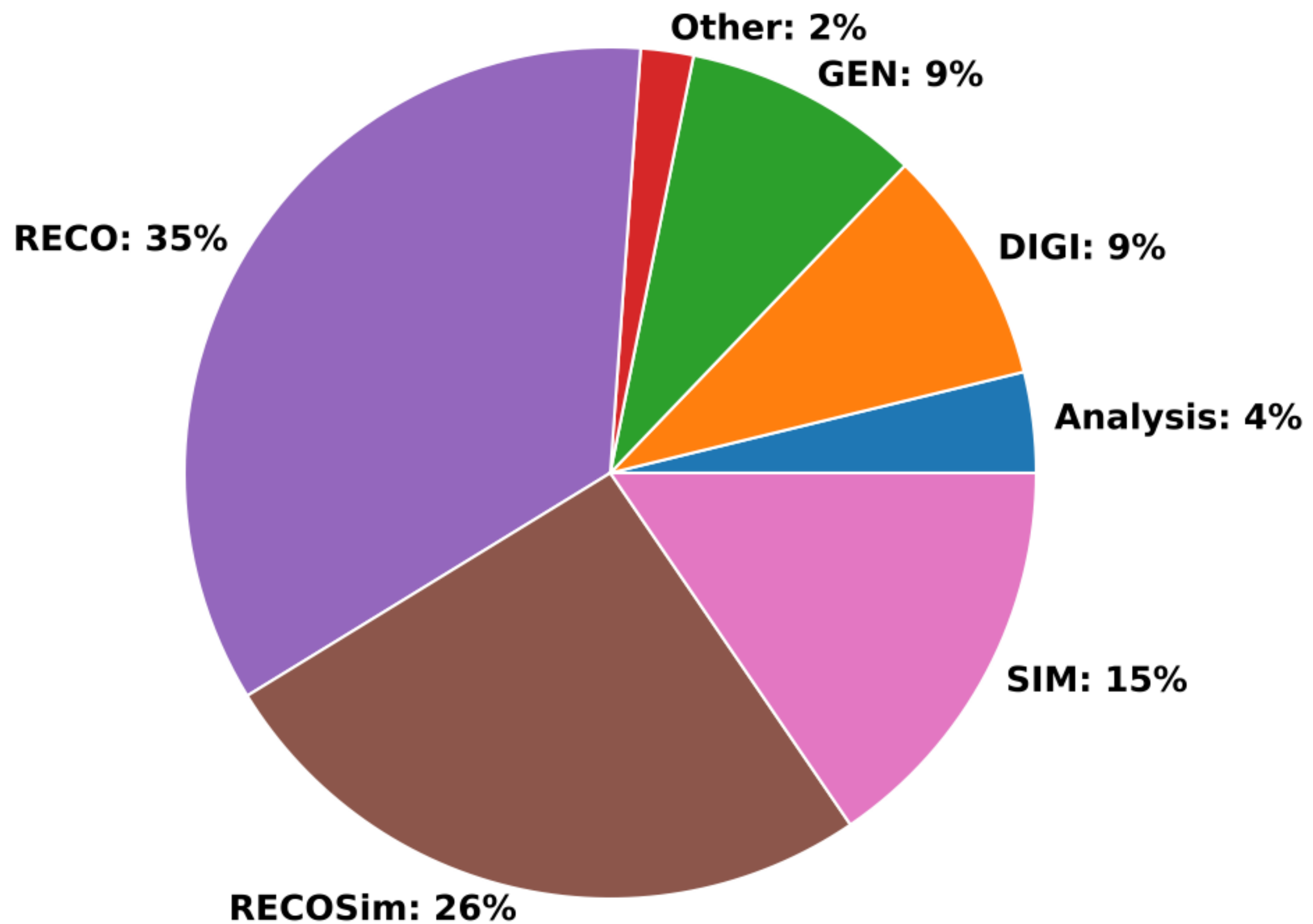


HL-LHC, Simulation of CMS HGCal with 140 PU

CMSPublic

Total CPU HL-LHC (2031/No R&D Improvements) fractions

2022 Estimates



$O(10)$

$O(10^3)$

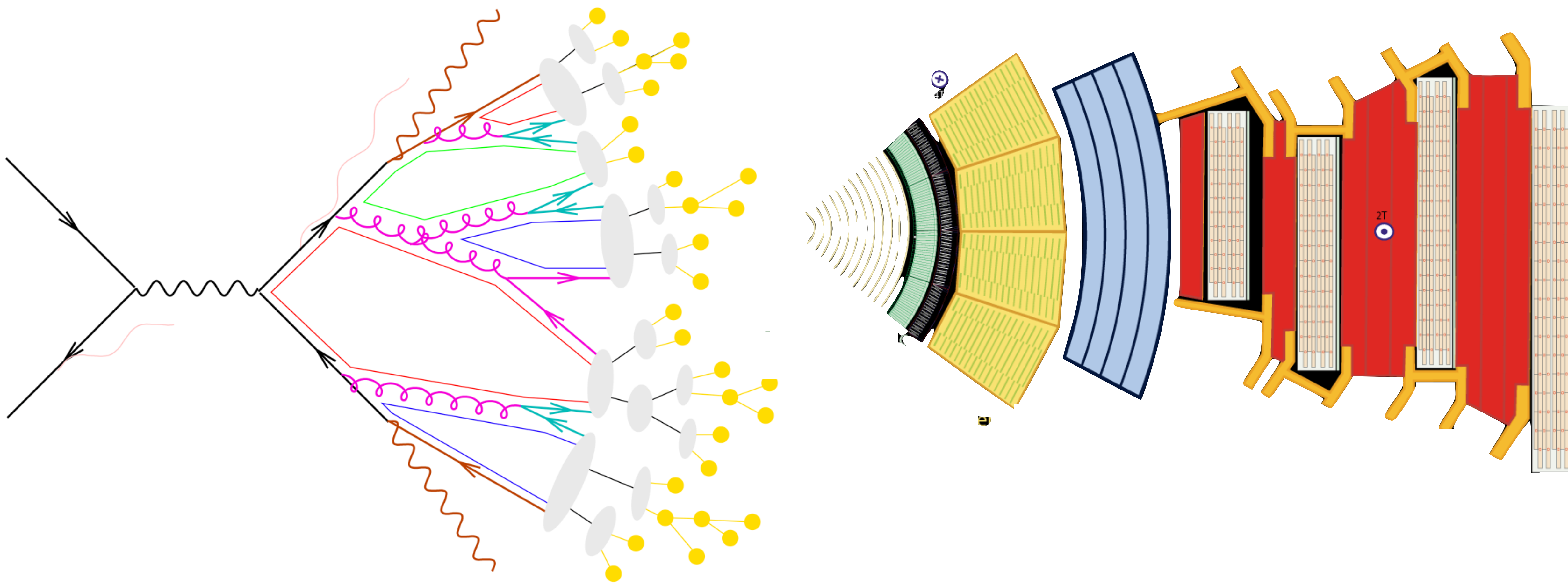
$O(10^{10})$

10^{-18}m

10^{-15}m

10^{-6}m

100m



$O(10)$

$O(10^3)$

$O(10^{10})$

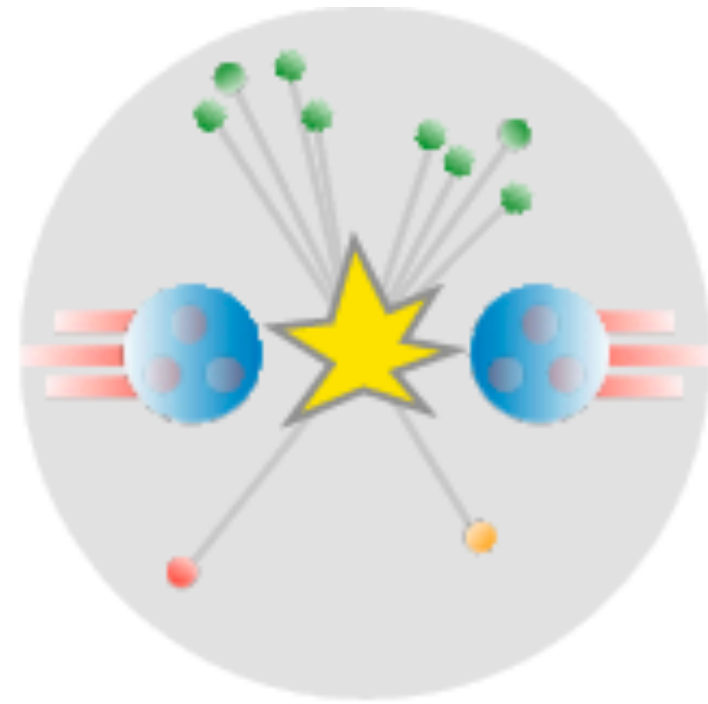
10^{-18}m

10^{-15}m

10^{-6}m

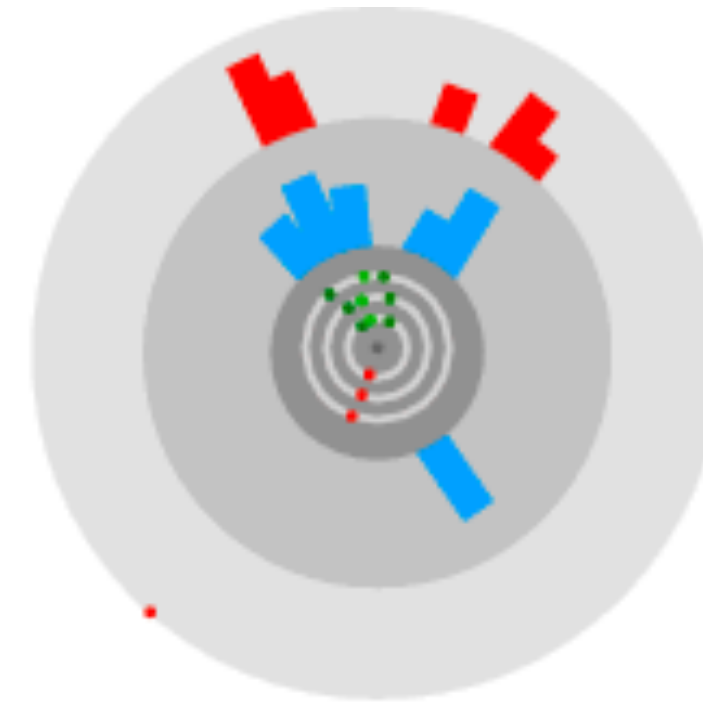
100m

GEN



pp collisions up to
production of stable
particles [Easy & Fast]

SIM



detector response
simulation [Hard & Slow]



DIGI+RECO



Energy deposits → digital
signals → reconstructed by
the reconstruction software
[Hard & Slow]

$O(10)$

$O(10^3)$

$O(10^{10})$

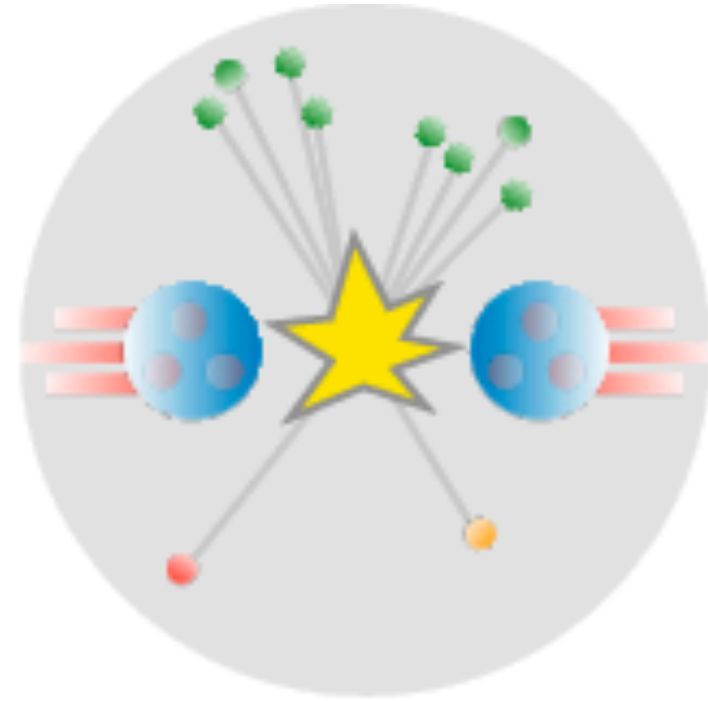
10^{-18}m

10^{-15}m

10^{-6}m

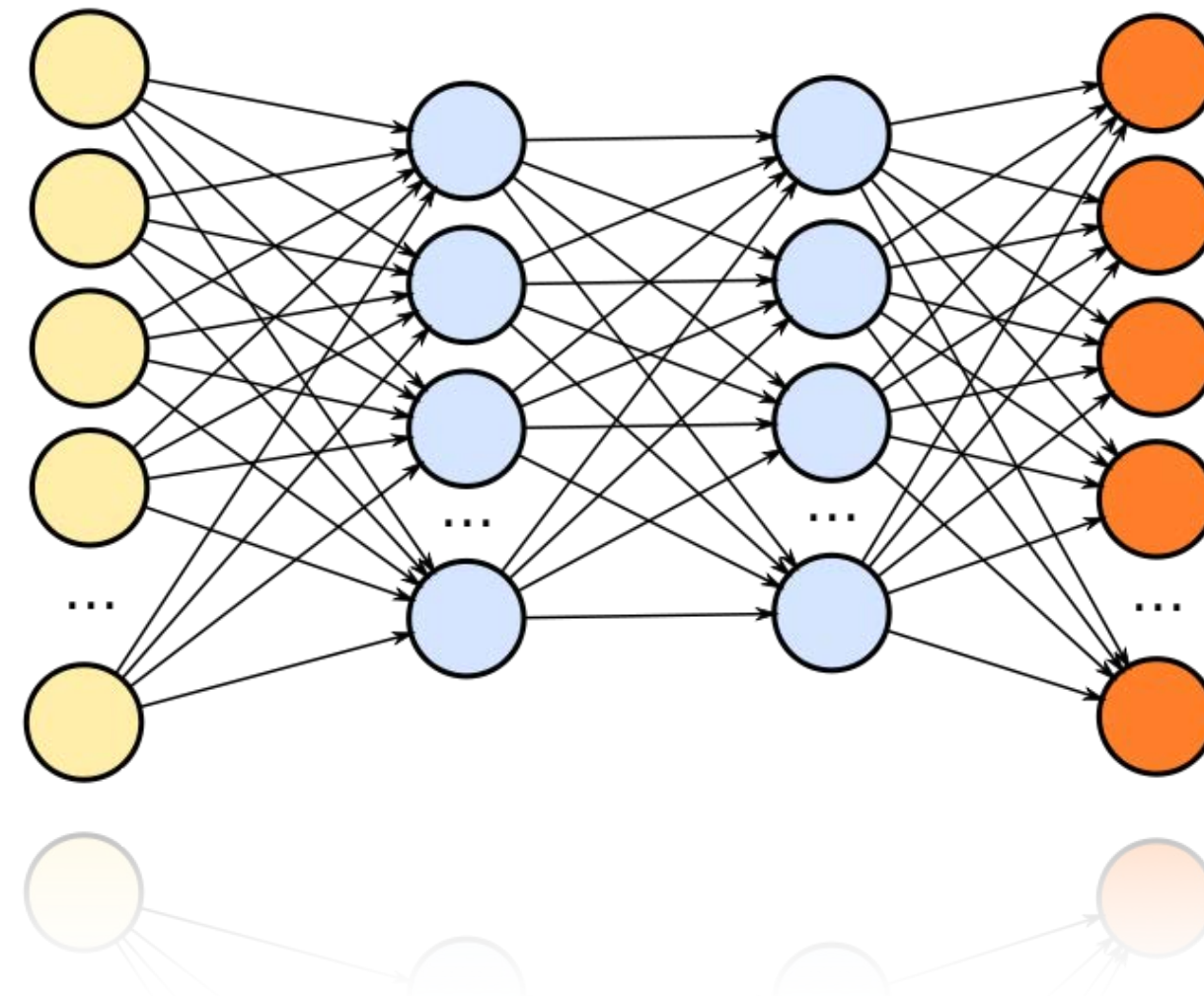
100m

GEN



pp collisions up to
production of stable
particles [Easy & Fast]

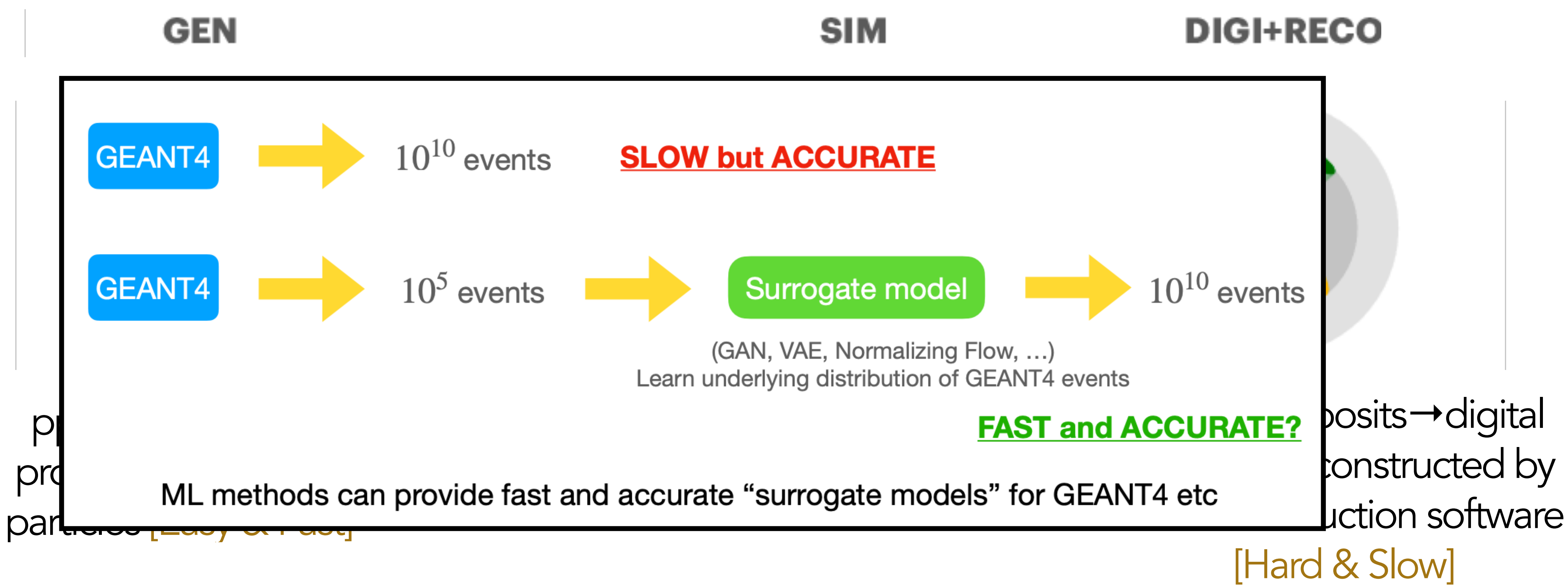
SIM



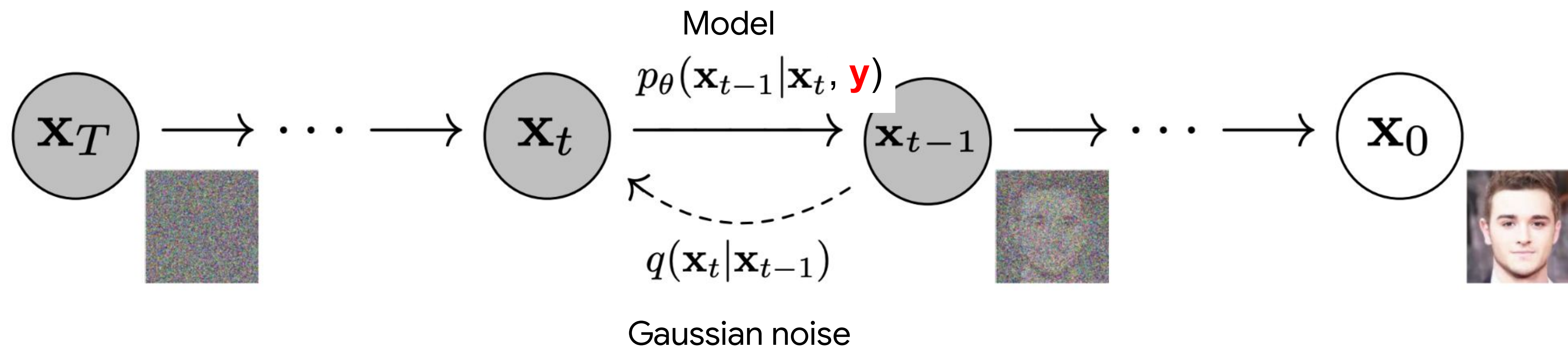
DIGI+RECO



Energy deposits → digital
signals → reconstructed by
the reconstruction software
[Hard & Slow]



Diffusion models

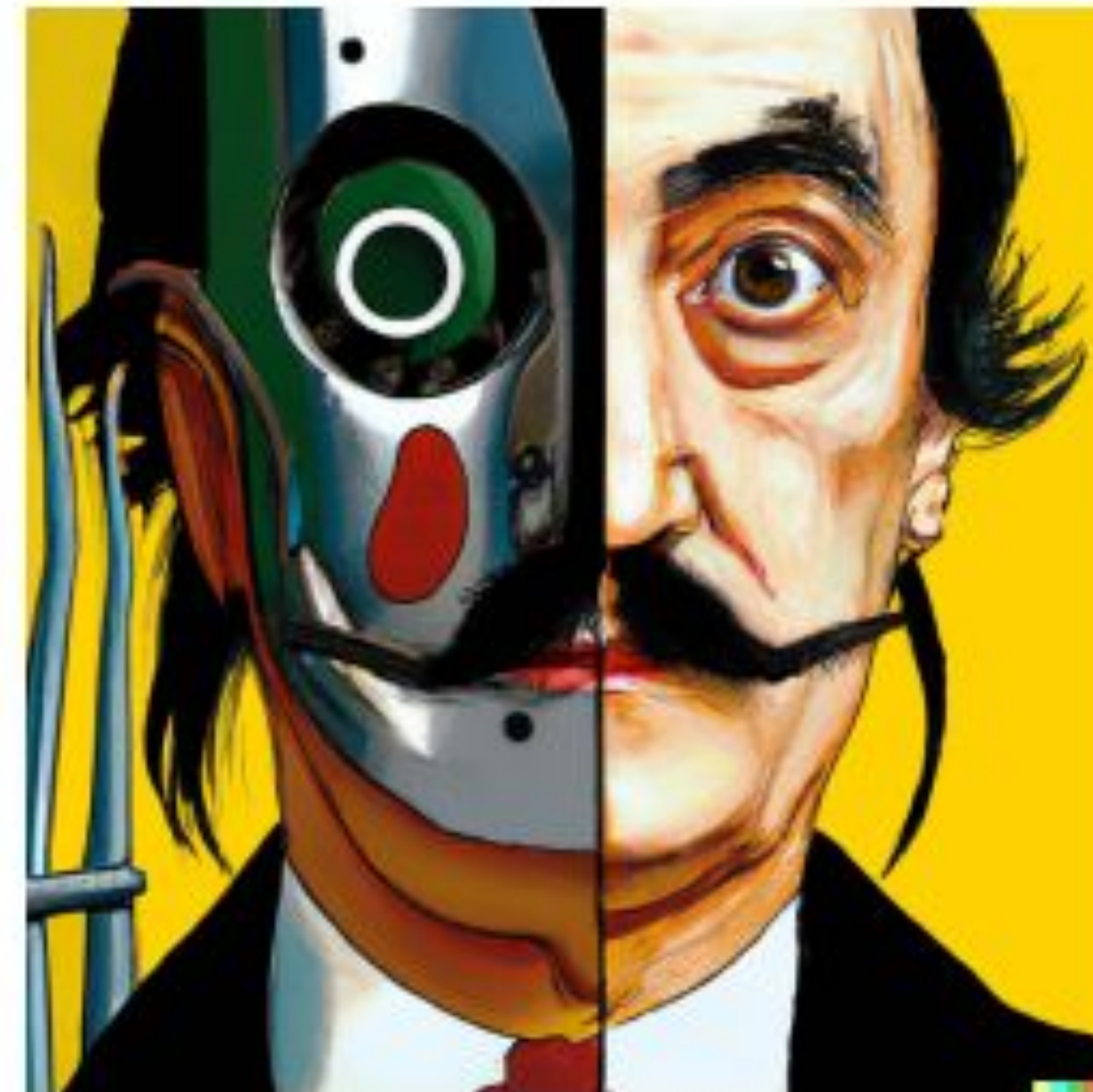


Dall-e 2



an espresso machine that makes coffee from human souls, artstation 2

Dall-e 2



vibrant portrait painting of Salvador Dalí with a robotic half face

