# Machine Learning at the Edge: Nanosecond inference at the Large Hadron Collider
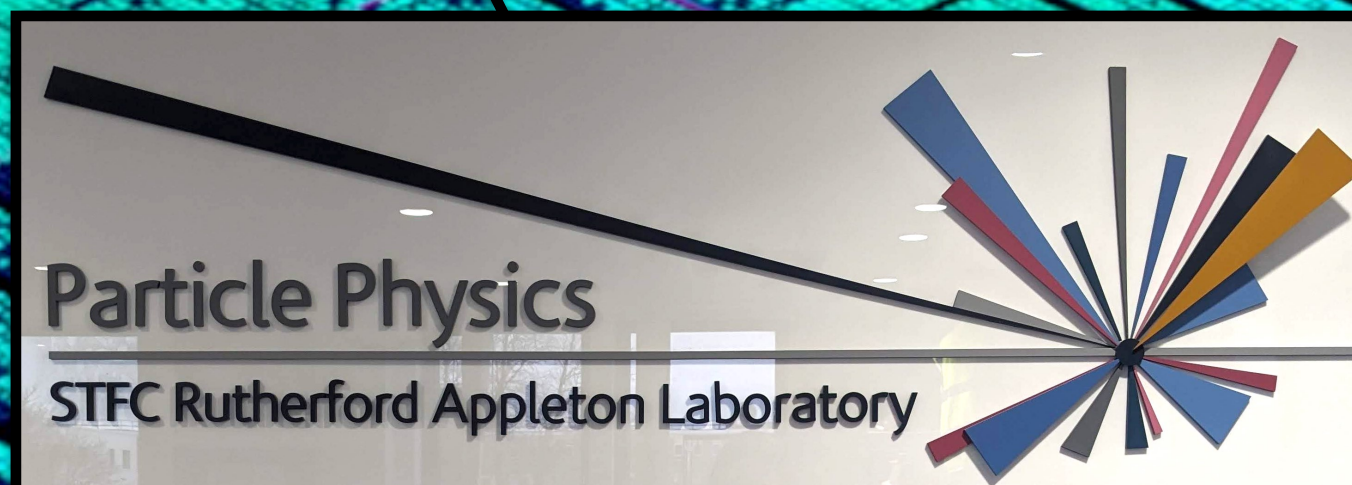
Thea Klæboe Årrestad
(ETH Zürich)

Particle Physics

STFC Rutherford Appleton Laboratory

Teaching | CMS general | Python | Papers | PSI | Useful tools | VVana_2017 | Jets | MET | b-tagging | Trigger | Jet substructure | DQM Shift Work | TTbar Info | Electrons | Private | UZH | Meetings | Talks | All Bookmarks

ChatGPT

Explore GPTs

**NEW Explore GPTs**
Now you can discover GPTs created by the community

Today

New chat

IEEE Ref Style Article Summary
IEEE Citation Style Format
Format IEEE Reference
IEEE Citation for Neuromorphic C
Advanced ML for L1T Upgrade
Cite Website Details Needed

Yesterday

CMS L1T Upgrade Tasks
IEEE Reference for Article

Previous 7 Days

New chat
IEEE Style Reference Retrieval
Anomaly Detection in Particle Phy
BibTeX Website Entry Example

Previous 30 Days

Thesis Citation in BibTeX
BibTeX for Physics Paper
ETH's CMS Trigger Development
Add Git to environment.yml
Change Hyperlinks to Black
Calculate Invariant Mass Python
Anomaly Detection Challenges
LaTeX Package Compatibility Issu
GSHPs Use Refrigerant

**Upgrade plan**
Collaborate on a Team plan

Thea Aarrestad

**You**

IEEE style reference please: @ARTICLE{9447722,
 author={Guglielmo, Giuseppe Di and Fahim, Farah and Herwig, Christian and Valentin, Manuel Blanco and Duarte, Javier and Gingu, Cristian and Harris, Philip and Hirschauer, James and Kwok, Martin and Loncar, Vladimir and Luo, Yingyi and Miranda, Llovizna and Ngadiuba, Jennifer and Noonan, Daniel and Ogrenci-Memik, Seda and Pierini, Maurizio and Summers, Sioni and Tran, Nhan},
 journal={IEEE Transactions on Nuclear Science},
 title={A Reconfigurable Neural Network ASIC for Detector Front-End Data Compression at the HL-LHC},
 year={2021},
 volume={68},
 number={8},
 pages={2179-2186},
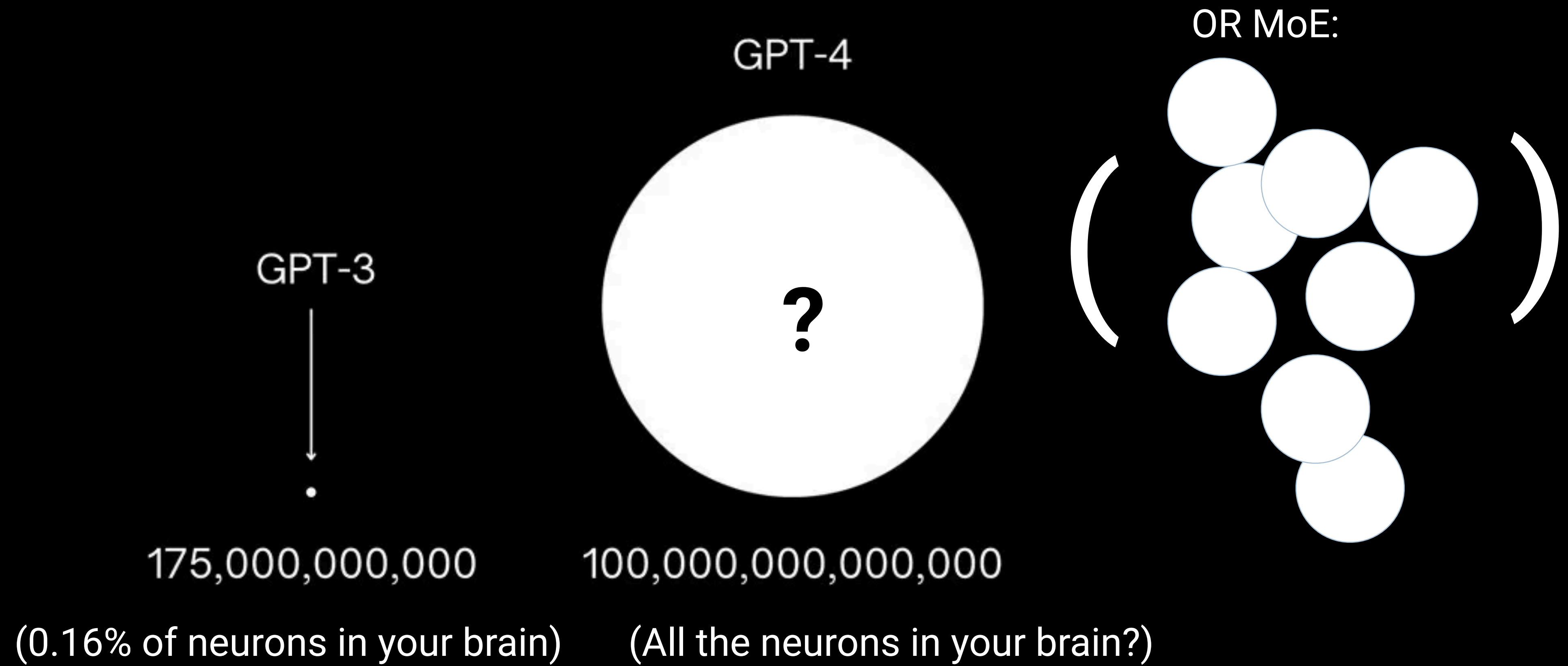 doi={10.1109/TNS.2021.3087100}}

**ChatGPT**

Here's how you would reference the specified article in IEEE style:

G. Di Guglielmo et al., "A Reconfigurable Neural Network ASIC for Detector Front-End Data Compression at the HL-LHC," in IEEE Transactions on Nuclear Science, vol. 68, no. 8, pp. 2179-2186, 2021, doi: 10.1109/TNS.2021.3087100.

Message ChatGPT...

ChatGPT can make mistakes. Consider checking important information.

**Train (GPT-3):**
- **285,000 CPU cores**
- **10,000 GPUs**
- **400 Gb/s network**
- **Several weeks**
- **Trained on ~3OOB tokens**

# NVIDIA Press Release 2020



**Train (GPT-3):**
- **285,000 CPU cores**
- **10,000 GPUs**
- **400 Gb/s network**
- **Several weeks**
- **Trained on ~300B tokens**



**You**

Who's this "Appleton" guy anyways?

**ChatGPT**

"Appleton" could refer to a variety of individuals or entities, depending on the context. It's a common surname, and it might also refer to locations or institutions. Here are a few possibilities:

1. **Individuals with the Surname Appleton**: There have been many notable people with the last name Appleton, including artists, scientists, and other professionals. Without more context, it's hard to identify a specific individual.
2. **Appleton, Wisconsin**: A city in the United States known for its role in the paper industry and its proximity to the Fox River.
3. **Appleton Estate**: A famous Jamaican rum producer.
4. **Appleton Company or Brand**: It could refer to a company or brand named Appleton, involved in various industries.
5. **Historical Figures**: There have been historical figures with the surname Appleton who played significant roles in various fields.

If you can provide more specific information or context about which "Appleton" you're referring to, I could give a more detailed and accurate response.

**You**

He might have been a friend of Rutherford

**ChatGPT**

If the "Appleton" you're referring to was a friend or contemporary of Ernest Rutherford, a prominent New Zealand-born physicist known for his pioneering work in nuclear physics, it's likely you're talking about Sir E ↓ rd Victor Appleton. Here are some key points about him:

Message ChatGPT...

**Inference (GPT-3):**
- **?**

GPT-3

175,000,000,000

# FP16 vs FP32

|   | 1 | 8 | 23 | Range | Accuracy |
|---|---|---|---|---|---|
| FP32 | S | E | M | $10^{-38} - 10^{38}$ | 0.000006% |

|   | 1 | 5 | 10 | | |
|---|---|---|---|---|---|
| FP16 | S | E | M | $6\times10^{-5} - 6\times10^{4}$ | 0.05% |

~700 GB of memory
→ (175B par × 4 bytes/par)
→ $O(10^1)$ larger than max memory in single GPU

→ ~350 GB
(175B param × 2 bytes/par)
→ 11 NVIDIA V100  ($10 000/ea)

AI and Memory Wall

Transformer Size: 240x / 2 yrs

AI HW Memory: 2x / 2 yrs

10TB Baidu RecSys

2TB Baidu RecSys

GShard

GPT-3

Microsoft T-NLG

Megatron LM

A100-80 (80GB)

V100 (32GB)  TPUv3 (32GB)

TPUv2 (16GB)

A100 (40GB)

P100 (12GB)

GPT-2

BERT

5 of these to fit one GPT-3 at inference time! Price: $17,000

Inception V4

ResNext101

ALBERT

ResNet50

DenseNet

Transformer

GPT-1

ELECTRA

Parameter Count (Billion)

YEAR

AI & Memory Wall

AI and Memory Wall

Transformer Size: 240x / 2 yrs

AI HW Memory: 2x / 2 yrs

Parameter Count (Billion)

YEAR

10TB Baidu RecSys

2TB Baidu RecSys

GShard

GPT-3

Microsoft T-NLG

Megatron LM

A100-80 (80GB)

V100 (32GB)

TPUv3 (32GB)

A100 (40GB)

TPUv2 (16GB)

P100 (12GB)

BERT

GPT-1

Transformer

Inception V4

ResNext101

ResNet50

DenseNet

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

AI & Memory Wall

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



$$L = (C_{\min}/2.3 \cdot 10^8)^{-0.050}$$

$$L = (D/5.4 \cdot 10^{13})^{-0.095}$$

$$L = (N/8.8 \cdot 10^{13})^{-0.076}$$

Test loss

Compute (peta-FLOP/s-day)

Data set size (tokens)

Number of parameters

Kaplan et al. (2020)

**_Resources:_** _11 interconnected GPUs_

**_Latency :_** _$10^1$ seconds_



**You**

Who's this "Appleton" guy anyways?

**ChatGPT**

Resources: 11 interconnected GPUs
Latency :    $10^1$ seconds

Resources: One single chip
Latency:    $10^{-9}$ seconds



CMS Experiment at the LHC, CERN
Data recorded: 2010-Nov-14 18:37:44.420271 GMT(19:37:44 CEST)
Run / Event: 151076 / 1405388

You
Who's this "Appleton" guy anyways?

ChatGPT

# EFFICIENT AI

Efficient Deep Learning: A Survey on Making Deep Learning Models Smaller, Faster, and Better, G Menghani
tinyML

Francois Lanusse

the Rubin Observatory Legacy Survey of Space and Time

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

Caltech/LIGO



Time to merger (s)

S/N ≈ 9    S/N ≈ 13    S/N ≈ 22

# Cherenkov Telescope Array (CTA)

CTA trigger: 72 Tb/s - 33 µs - FPGA based

Geneva

Geneva Lake

LHCb

ATLAS

ALICE

CMS

2,500 bunches
$10^{11}$ protons
11,000 times/s

$2.2 \cdot 10^{11}$ protons

25 ns

7.5 m

Geneva

Geneva Lake

CMS

LHCb

ATLAS

ALICE

O(1) billion collisions per second
O(1) PB of data per second

"Probability" of producing "anything"

"Probability" of producing a Higgs

Higgs produced
~1 in a billion collisions!

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

Geneva

Geneva Lake

LHCb

ATLAS

ALICE

CMS

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

LHC

Geneva

Geneva Lake

LHCb

AT[

CMS

**Data temporarily stored in detector electronics for 4 μs (frontend buffering limit)**

**2 Tb/s to L1**

**L1 trigger: (~1000 FPGAs)**

**Decide which event to keep within ~4 μs**

**Reject >99% of collisions!**

LHC

Geneva

Geneva Lake

LHCb

AT...

CMS

LHC

L1 bit:
Accept = 1
Reject = 0

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

Reject further 99%!

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

ATLAS

ALICE

LHC

High Level Trigger:
Latency O(100) ms

TIER 0: ∞

0.0025% of collision events remaining

HLT accept:
O(1) kHz
~Gb/s

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

Geneva

LHC

ATLAS

CMS

LHC

*HIG-19-001*

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 <u>more collisions</u> to observe rare processes

**High Luminosity LHC**
- ×10 data size
- ×3 collisions/s



| 2022 - 2025 | 2026 - 2028 | 2029 - 2038 |
| --- | --- | --- |
| **LHC (TODAY!)** | **MAJOR UPGRADE** | **HL-LHC** |
| Run 3 | | Run 4+5 |

structure → pile-up of ~ 60 events/x-ing
ts/x-ing)

# High Luminosity LHC

200 vertices
(average 140)

CMS: event with 78 reconstructed vertices

6 cm

ATLAS
EXPERIMENT
HL-LHC t̄t event in ATLAS ITK
at <μ>=200

Run 4+5

Maintain physics acceptance → better detectors

CMS High Granularity (endcap) calorimeter
- 85K (today) → 6M (HL-LHC) readout channels

**More** collisions
**More** readout channels



silicon

scintillator

~10λ

~5λ

CE-E    CE-H

η = 3.0

CMS HGCAL TDR

# Computing resources



... **flat** computing budget

Todays algorithms will not be sustainable in HL-LHC!

→ Need modern Machine Learning to become

faster
better
and do more

More complex architectures to deal with increased
data complexity!

Geneva

Geneva Lake

LHCb

AT

CMS

LHC

HL-LHC L1:
• Tracking
• Particle Flow
• HGCal

Input data
• 2 Tb/s → **63 Tb/s**

Latency
• 4 μs → **12 μs**

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Super Solo...

...return yoke interspersed with Muon chambers

63 Tb/s

USC55

UXC55

Xilinx Ultrascale+ FPGAs

**CALORIMETRY: 370 FPGAs**

*54 for HGCAL only!

**TRACKING 174 FPGAs**

**MUONS: 96 FPGAs**

5 μs

**PARTICLE FLOW: 66 FPGAs**

CALORIMETRY

PARTICLE FLOW

MUONS

**GLOBAL TRIGGER: 24 FPGAs**

EXTERNAL TRIGGERS

Trigger accept/reject

12.5 μs

63 Tb/s

Xilinx Ultrascale+ FPGAs

CALORIMETRY:
370 FPGAs

*54 for HGCAL only!

TRACKING
174 FPGAs

MUONS:

PARTICLE
FLOW:

Simulated event display with average pileup of 140

UXC55

*ATLAS & CMS:* **Trigger System**

- Current trigger systems
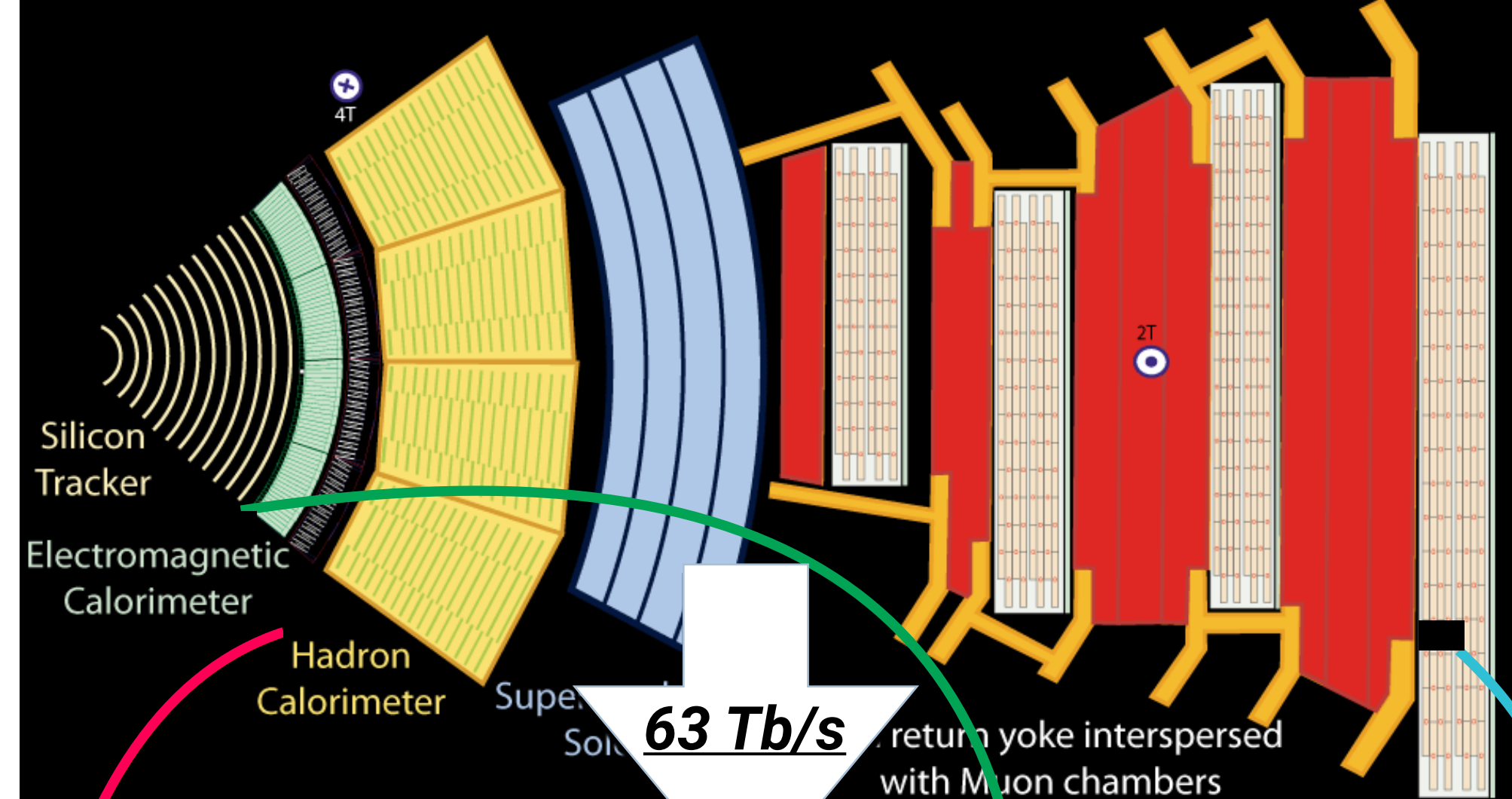  - **L1 trigger**
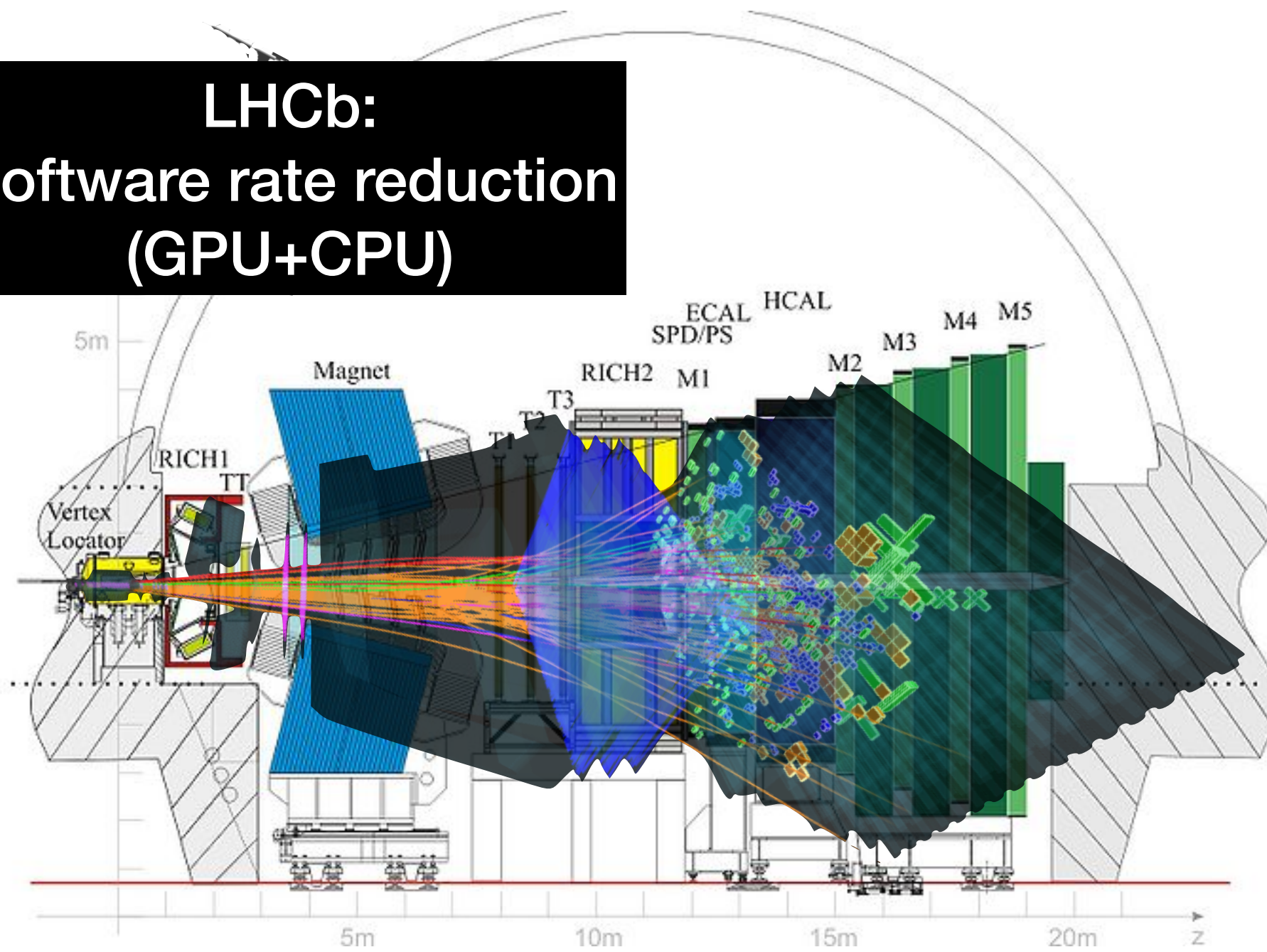    - *Hardware-based, implemented in custom-built electronics*

# 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...

LHCb:
Software rate reduction
(GPU+CPU)

→ **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)**

Full GPU reconstruction @ 4 TB/s
- 326 GPUs, 60 kHz per GPU

# Why FPGAs?

Latency, latency, latency (cannot do much on a GPU IN 4 µs)
- 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
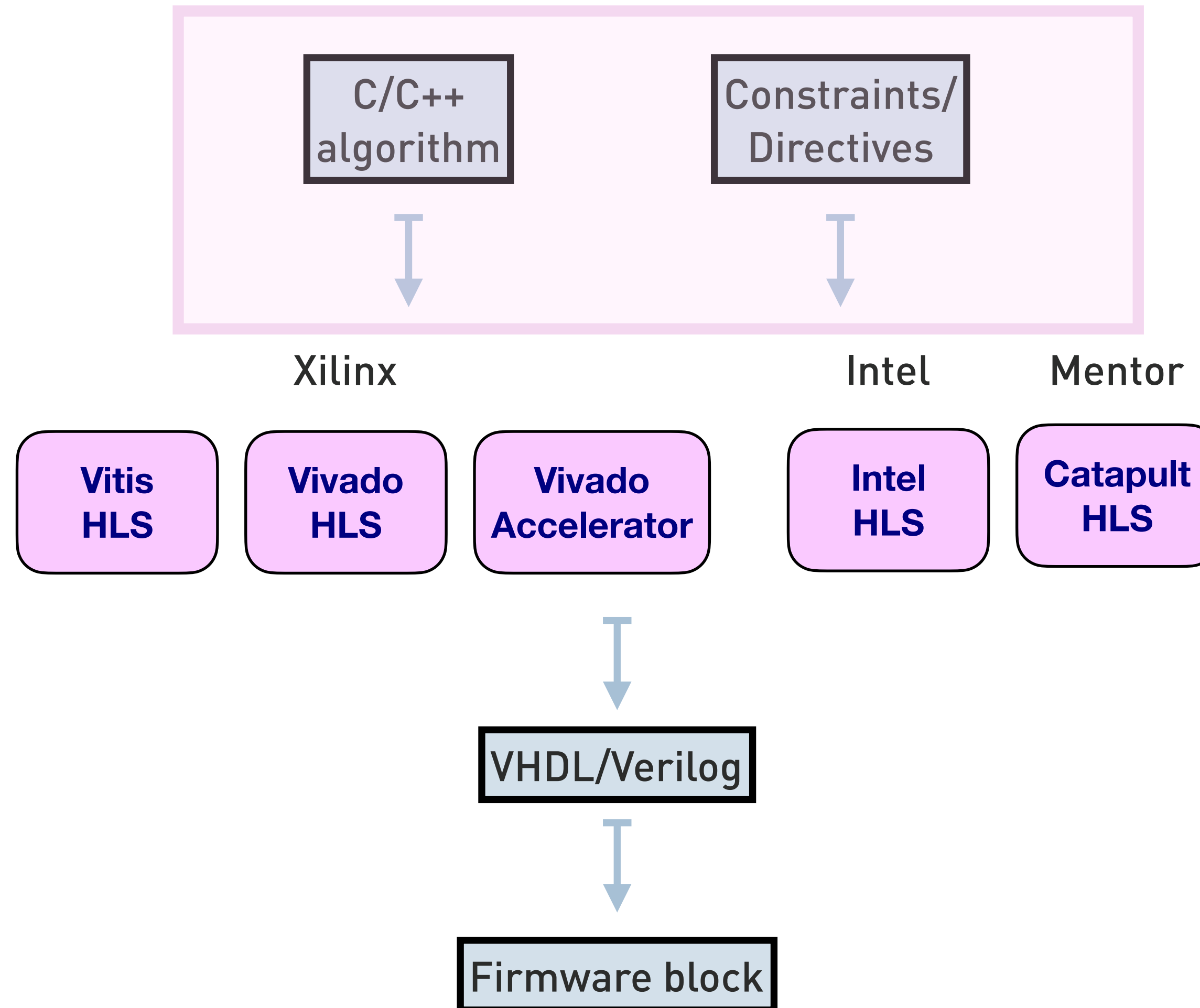
High bandwidth
- L1T processes 5% of total internet traffic, dissipate heat of ~7W/cm$^2$

TRACK FINDER

Work on 18 events simultaneously!

TMUX=18
RS = 9 (phi)
FPGAs = 162

# Programming an FPGA

# Programming an FPGA

C/C++
algorithm

Constraints/
Directives

Xilinx

Intel

Mentor

**Vitis HLS**

**Vivado HLS**

**Vivado Accelerator**

**Intel HLS**

**Catapult HLS**

VHDL/Verilog

Firmware block

```vhdl
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_MAS

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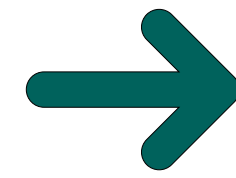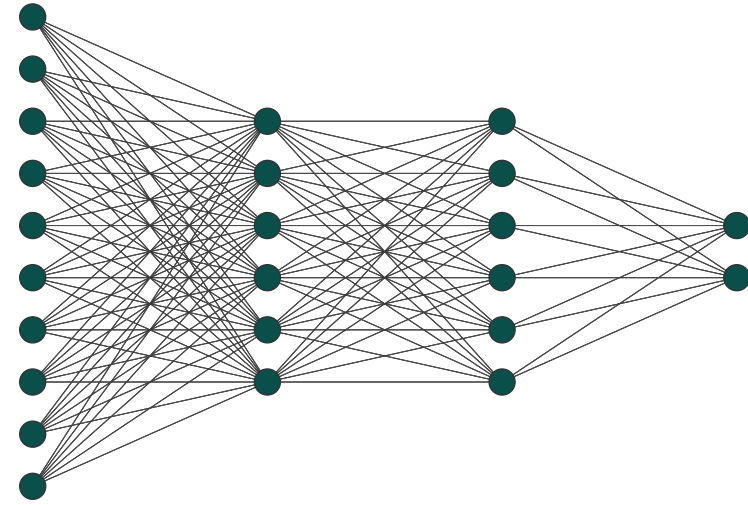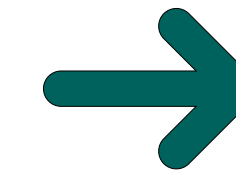
$$\mathbf{x}_n = g_n(\mathbf{W}_{n,n-1}\mathbf{x}_{n-1} + \mathbf{b}_n)$$

# Generic (superfast) HLS implementations for DNN inference?

```vhdl
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_MAS

    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;
```
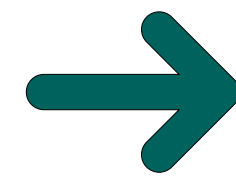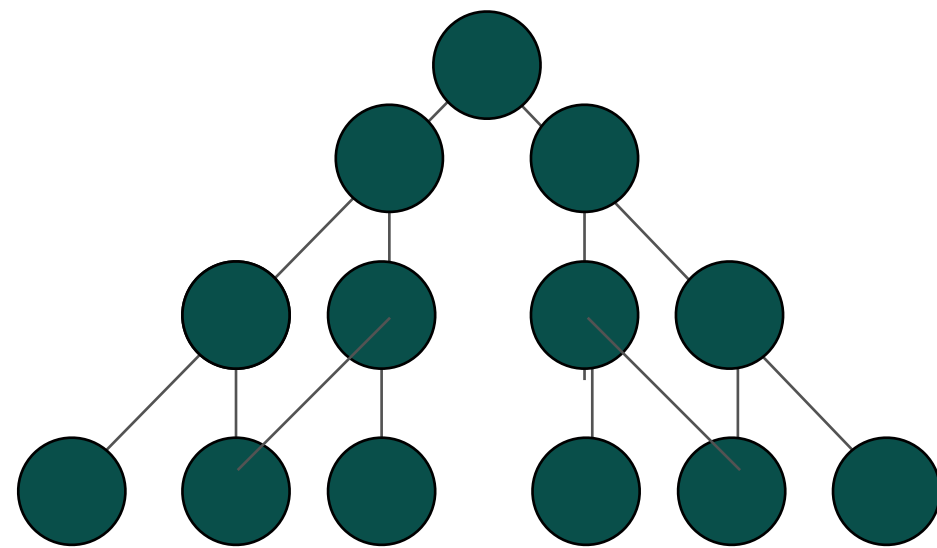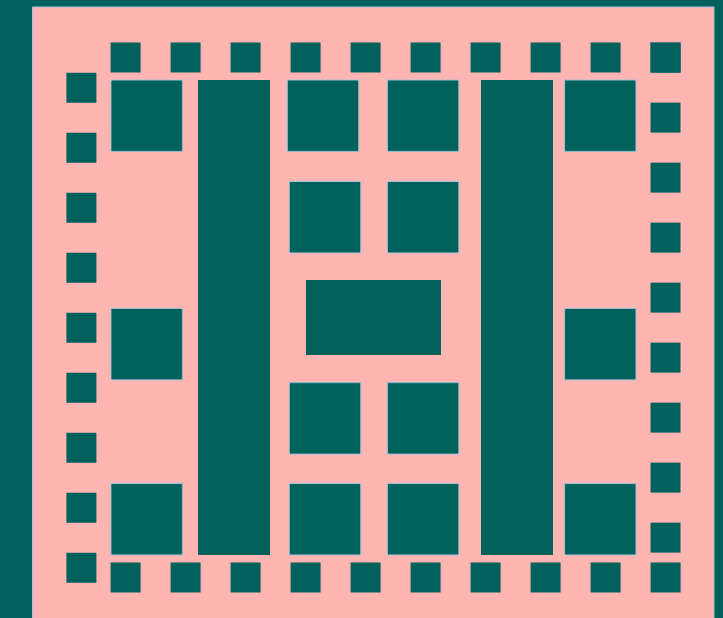
KERAS / PyTorch / ONNX

TensorFlow DF / scikit-learn / XGBoost

hls 4 ml

Conifer

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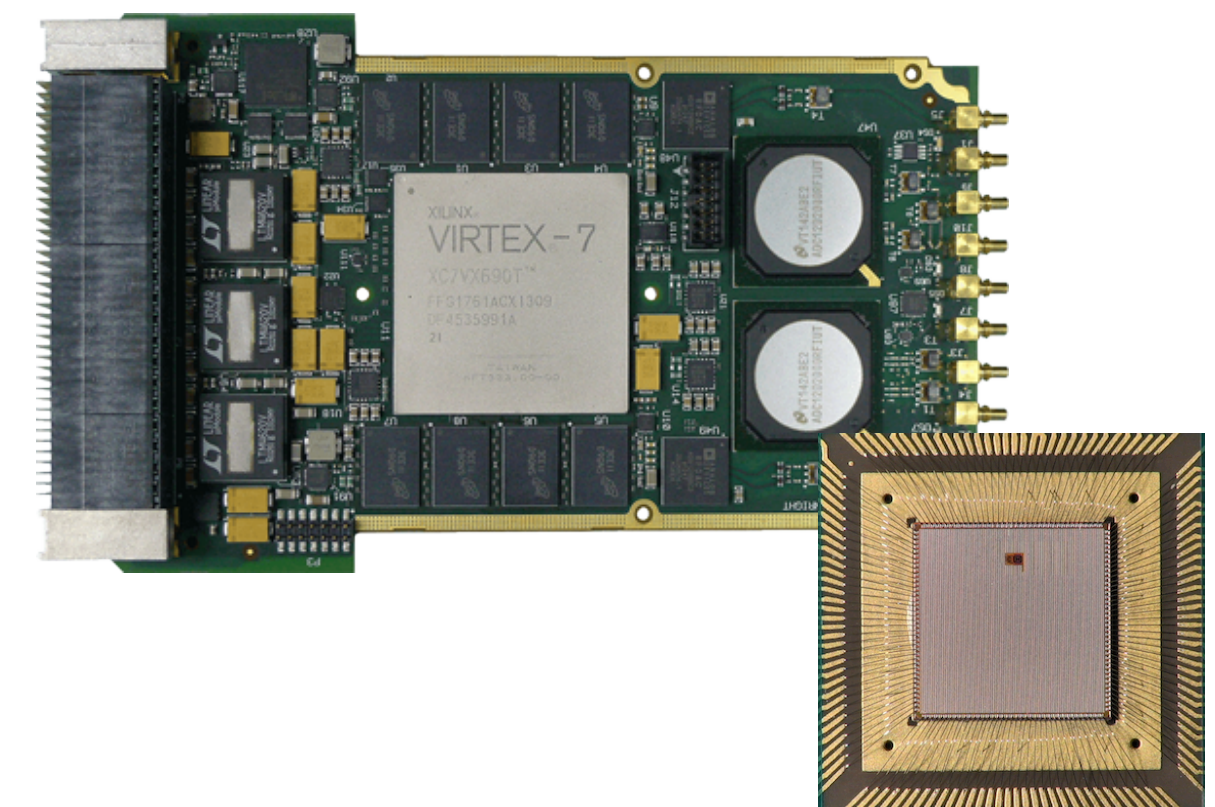Vivado / Vitis / Intel Quartus /
IntelOne API / Catapult

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

K + TensorFlow

ONNX

PYTORCH

Model
(quantized/pruned)

*Quantized:*

QONNX

QKERAS

ONNX

K + TensorFlow

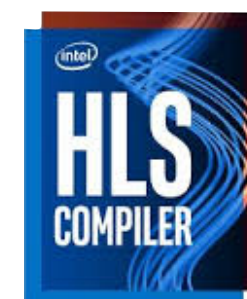PYTORCH

Model
(quantized/pruned)

Quantized:

QONNX

QKERAS

hls4ml

Convert model to internal
representation

Write HLS project targeting
specified backend

Run emulation

Run synthesis

Vivado™ HLS

Mentor
A Siemens Business

intel HLS COMPILER

Vivado/Vitis best supported
Intel Quartus
Intel One API
Mentor Catapult HLS
available soon

**hls4ml**

Convert model to internal representation

Write HLS project targeting specified backend

Run emulation

Run synthesis

**K** + TensorFlow    ONNX

**PYTORCH**

Model
(quantized/pruned)

*Quantized:*

QONNX

Q**K**ERAS

Vivado™ HLS    Mentor® A Siemens Business    intel HLS COMPILER

Vivado/Vitis best supported
Intel Quartus
Intel One API
Mentor Catapult HLS
available soon

*Co-processing kernel
(Xilinx accelerators/SoCs)*

*FPGA custom designs
(eg trigger algorithms)*

*ASICs*

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

(from Sioni S Summers)

**Ideally**

**Reality**

**Ideally**

- Quantization
- Pruning
- Parallelisation
- Knowledge distillation

**Reality**

# Quantization

4B numbers in [-3.4e38, +3.4e38]

## Floating point 32

sign exponent(8-bit)    fraction (23-bit)

0 0 1 1 1 1 1 0 0 0 0 1 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 = 0.15625

31        23                                    0

## Weights Layer 1

FP 32

## Weights Layer 2

FP 32

GPT-3

175,000,000,000

FP16 vs FP32

| | Range | Accuracy | |
|---|---|---|---|
| FP32 | $10^{-38}$ - $10^{38}$ | 0.000006% → | ~700 GB |
| FP16 | $6 \times 10^{-5}$ - $6 \times 10^{4}$ | 0.05% → | ~350 GB |

FP32: S 1 | E 8 | M 23

FP16: S 1 | E 5 | M 10

# Quantization

2⁴ numbers in $s*[-8, +7]$

## Fixed point



## Weights Layer 1

4-bit

## Weights Layer 2

4-bit

# Quantization

$2^4$ numbers in $s*[-8, +7]$

## Fixed point

0101.1011101010

integer ← → fractional
← width →

## Weights Layer 1

4-bit

## Weights Layer 2

4-bit

# hls4ml + Google
## Quantization-aware training

Forward pass →

← Back propagation

**Nature Machine Intelligence 3 (2021)**

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)
```

**Nature Machine Intelligence 3 (2021)**

# 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 $\uparrow$ ...
- ... and lowest resource consumption $\downarrow$

$\rightarrow$ 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 [$\mu$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 |

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

***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!

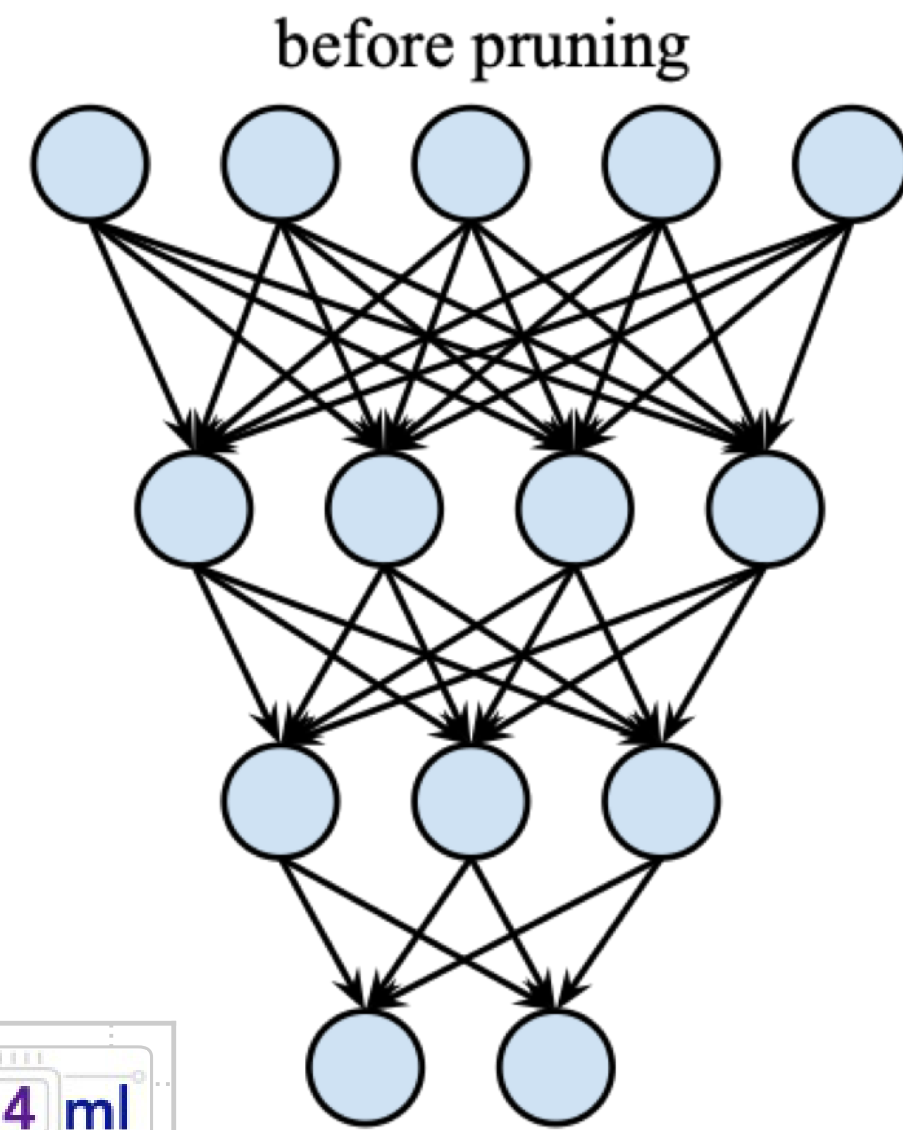AMD UltraScale+ MPSoC ZU19EG (conservative estimates)

| Precision | Approx. **Peak GOPS** |
|-----------|----------------------|
| **1b** | **64 000** |
| **4b** | **16 000** |
| **8b** | **4 000** |
| **32b** | **300** |

200x

| On-chip weights |
|-----------------|
| **~64 M** |
| **~16 M** |
| **~8 M** |
| **~2 M** |

30x

**Trillions** of quantized operations per second

Weights can stay **entirely on-chip**

# Pruning



before pruning



Baseline

*https://blog.tensorflow.org/2019/05/tf-model-optimization-toolkit-pruning-API.html*

# Lottery ticket hypothesis



Resnet18

Better!

— Reference Accuracy: 94.42%

Pruned

Unpruned

Baseline (Dense) · · · · · Biprop (Global) — Edgepopup (Global) · · · · · Fine-Tuning
Gradual Magnitude Pruning — Learning Rate Rewinding · · · · · LTH

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

Diffenderfer, Bartoldson, et al. (2021)

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

*Leo Grinsztajn, Edouard Oyallon, Gael Varoquaux*

**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 (∼10K samples) even without accounting

---

Cornell University
the Simo

arXiv > cs > arXiv:2210.05189

Computer Science > Machine Learning

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

## Neural Networks are Decision Trees

Caglar Aytekin

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.

*https://arxiv.org/abs/2210.05189*

# OpenReview.net

Search...

Help | Advanc

## Why do tree-based outperform deep le typical tabular dat

*Leo G*

*Varo*

06 Jun

Bench

Revisi

**Abstra**

and im

contri

metho

Forest

combi

domai

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 (≈10K samples) even without accounting

### Deep learning underperforms on data tables
[Grinsztajn... 2022]

Classification

Regression

Normalized test accuracy of best model (on valid set) up to this iteration

HistGradientBoostingTree
XGBoost
RandomForest
GradientBoostingTree
FT Transformer
MLP Resnet SAINT
Random search time (seconds)
1e+01  1e+02  1e+03  1e+04  1e+05

HistGradientBoostingTree
GradientBoostingTree
SAINT
RandomForest
FT Transformer
XGBoost
Resnet
MLP
Random search time (seconds)
1e+01  1e+03  1e+05

ision Trees

2022 (this version, v3)]

al network with any activation function can be
valence and not an approximation, thus

e cool

# From yesterdays IML workshop!

data + GPU solves everything

■ Trees are still very relevant

(maybe in a few months, we'll make them obsolete 😀)

■ It's <u>really</u> not all about scale
- Architecture matters
- Some physics does constrain compute
- The economics matter, and are dodgy

represent
holds bot
skip conn

Subjects: **Mac**
Cite as:    **arXi**
            (or a
            **http**

**Submissio**

From: Çağlar
**[v1]** Tue, 11 (
**[v2]** Mon, 17
**[v3]** Tue, 25 (

Training cost
a single run

$1M
$1k
$1
AlphaGo Zero
PaLM (540B)

2010    2020

Language   Multimodal
Vision     Drawing
Other      Speech
Unknown
Games

Gemini Ultra
AlphaGo Zero

$10^{24}$

$10^{18}$

1 day compute of largest supercomputer

AlexNet

# Parameters

Amount of RAM for $1000

1T
1G
1M
1k

GPT-3
Mistral 7B
AlexNet
LeNet-5

2000    2010    2020

1PB
1TB
1GB
1MB
1kB

bfloat16 footprint

**Rebound effect**

*https://arxiv.org/abs/2210.05189*

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

12 microseconds latency

Processing 5% of internet traffic

Simulated event display with average pileup of 140

# Nanosecond ML inference on FPGAs!

# 40 billion inferences/s during HL-LHC?

# ( ≈ all inferences at Google)

- **L1 trigger**
  - *Hardware-based, implemented in custom-built electronics*
  - *Muon & calorimeter information with reduced granularity, no tracking information*

## Journey to HL-LHC

**run:**

**7 x 10³³, PU = 30, E = 7 TeV, 50 nsec bunch spacing**

**ATLAS, CMS operating:**

**ccept ≤ 100 kHz,**

**ncy ≤ 2.5 (AT), 4 μsec (CM)**

**Accept ≤ 1 kHz**

**LAS & CMS will be:**

**5 x 10³⁴**

| | LHC | | HL-LHC | |
|---|---|---|---|---|
| 40 MHz | | Detectors | | Detectors |
| | Lvl-1 | Front end pipelines | Lvl-1 | Front end pipelines |
| L1 output: 75 kHz | L1 output: 100 kHz | Readout buffers | L1 output: 100 kHz | Readout buffers |
| | Lvl-2 | | | |
| ~3 kHz | | Switching network | | Switching network |
| | Lvl-3 | Processor farms | HLT | Processor farms |
| 200 Hz | HLT output: ~1 kHz | | | |

40 MHz
40 MHz

100 Hz
7.5 kHz

<μ> = 32

$\sigma_{in}^{pp} = 69.2\ mb$

**Mean number of interactions per crossing**

**L1 trigger decision in ~2.5 (4) μs for ATLAS (CMS)**

MUONS

Trigger
ept/reject

12.5 μs

VBF H (γγ)

CMS High Granularity calorimeter
• 6.5 million readout channels, 50 layers

Layer 30

CE-H
Layer 30

Layer 5

es

d one VBF jet in the same quadrant,

jet

Cut $\Delta t < 90$ps   $(3\sigma$ at 30ps$)$

Layers projected onto one plane

-require hits within 90ps time window-

200 vertices

+

VBF H (γγ)

$q$

$H_0$

$γ$

$q$

No timing cut

90 ps time window

...d one VBF jet in the same quadrant,

γ

jet

...t

200 vertices

$10^2$

$+$

10

Cut $\Delta t < 90ps$  ($3\sigma$ at 30ps)

Layers projected onto one plane

-require hits within 90ps time window-

$\phi$  0

$-0.5$

$-1$

es

VBF H (γγ)

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

window

jet

Cut $\Delta t < 90ps$   ($3\sigma$ at 30ps)

Layers projected onto one plane

-require hits within 90ps time window-

200 vertices

$10^2$

$10$

+

0

−0.5

−1

**To L1**

HGCROC ASIC both for silicon and SiPMs    ECON as concentrator ASIC

ASIC

Hexaboard

Sensor

Kapton sheet

Cu/W Base plate

- 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 µs latency

ENCODE

Bottleneck
(lower dim.
space)

DECODE

Encoded data

# Variational Autoencoder

*AEs for compression also at LHCb!*

USC55

UXC55

**On ASIC**

V3 *HGCROC* ASIC both for silicon and SiPMs

ECON as concentrator ASIC

ASIC

Time-of-arrival (TOA) & time-over-threshold (TOT)

TOT

TOA

Interpretation of recorded TOA amplitudes

Sensor PCBs

FE ASICs

FE ASICs

FE ASICs

On detector

Off detector

DAQ

~160-320Gbit/s

Concentrator ASIC

Optical link

About 4k pairs at 10Gbit

About 8k pairs at 10Gbit

Trigger

Panel PCBs

Sensors
Sensor PCB
Front-end ASICs
Concentrator & Link
Panel PCB

DAQ, Trigger data

Power, Clock, Trigger & Slow-control

3

**Encoded data**

**Transmit encoded data!**

ECON-T, D. Noonan

AEs for compression also at LHCb!

USC55

UXC55

V3 *HGCROC* ASIC both for silicon and SiPMs    ECON as concentrator ASIC

ASIC

On ASIC

Transmit encoded data!

Encoded data

Encoded data

1

2

3

**AEs for compression also at LHCb!**

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

*UXC55*

*USC55*

**V3** *HGCROC* ASIC both for silicon and SiPMs

**ECON** as concentrator ASIC

**ASIC**

**On ASIC**

Time-of-arrival (TOA) & time-over-threshold (TOT)

Interpretation of recorded TOA amplitudes

Sensor PCBs

FE ASICs

FE ASICs

FE ASICs

~160-320Gbit/s

Concentrator ASIC

Optical link

About 4k pairs at 10Gbit

About 8k pairs at 10Gbit

Panel PCBs

DAQ

Trigger

On detector | Off detector

DAQ, Trigger data

Sensors
Sensor PCB
Front-end ASICs
Concentrator & Link
Panel PCB

Power, Clock, Trigger & Slow-control

**Encoded data**

**Transmit encoded data!**

**Encoded data**

1

2

3

**AEs for compression also at LHCb!**

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

USC55

UXC55

**FKeras**

in — comb logic — D Q / D Q / D Q 0 — majority voter — out

clk

**Bit flip!**

**On ASIC**

ASIC

CON as concentrator ASIC

Off detector

DAQ

Trigger

DAQ, Trigger data

Sensors
Sensor PCB
Front-end ASICs
Concentrator & Link
Panel PCB

Power, Clock,
Trigger & Slow-control

**Encoded data**

**Transmit encoded data!**

**Encoded data**

1

2

3

# ML for tracking

In HL-LHC, will need to do track finding at L1
  • O(1000) hits, O(100) tracks, 40 MHz rate, ~5 μs latency

Graph Neural Networks for fast charged particle tracking



| Design | $(n_{nodes}, n_{edges})$ | 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



?          ?

**Sets:** Information is
only assigned to
individual nodes.

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

# Fast jet tagging



? ?

**Sets:** Information is only assigned to individual nodes.

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

(Can also do 90 ns <u>transformers for jet tagging</u>!)

**Overview of CMS EXO results**

CMS Preliminary
August 2023

*From A. Rizzi*

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

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

# Limitations of current trigger



LOST DATA
SELECTED DATA
POSSIBLE NP SIGNAL

NP?

Trigger threshold

Energy (GeV)

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?

$$\text{loss} = \| x - \hat{x} \|^2$$

AXOLITL

x

x̂

$$\text{loss} = \| x - \hat{x} \|^2$$

LOST DATA
SELECTED DATA
POSSIBLE NP SIGNAL

Everything here is normal

Everything here is abnormal

NP?

AD threshold

Reconstruction error

AXOL1TL

....in 50 nanoseconds!

# 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

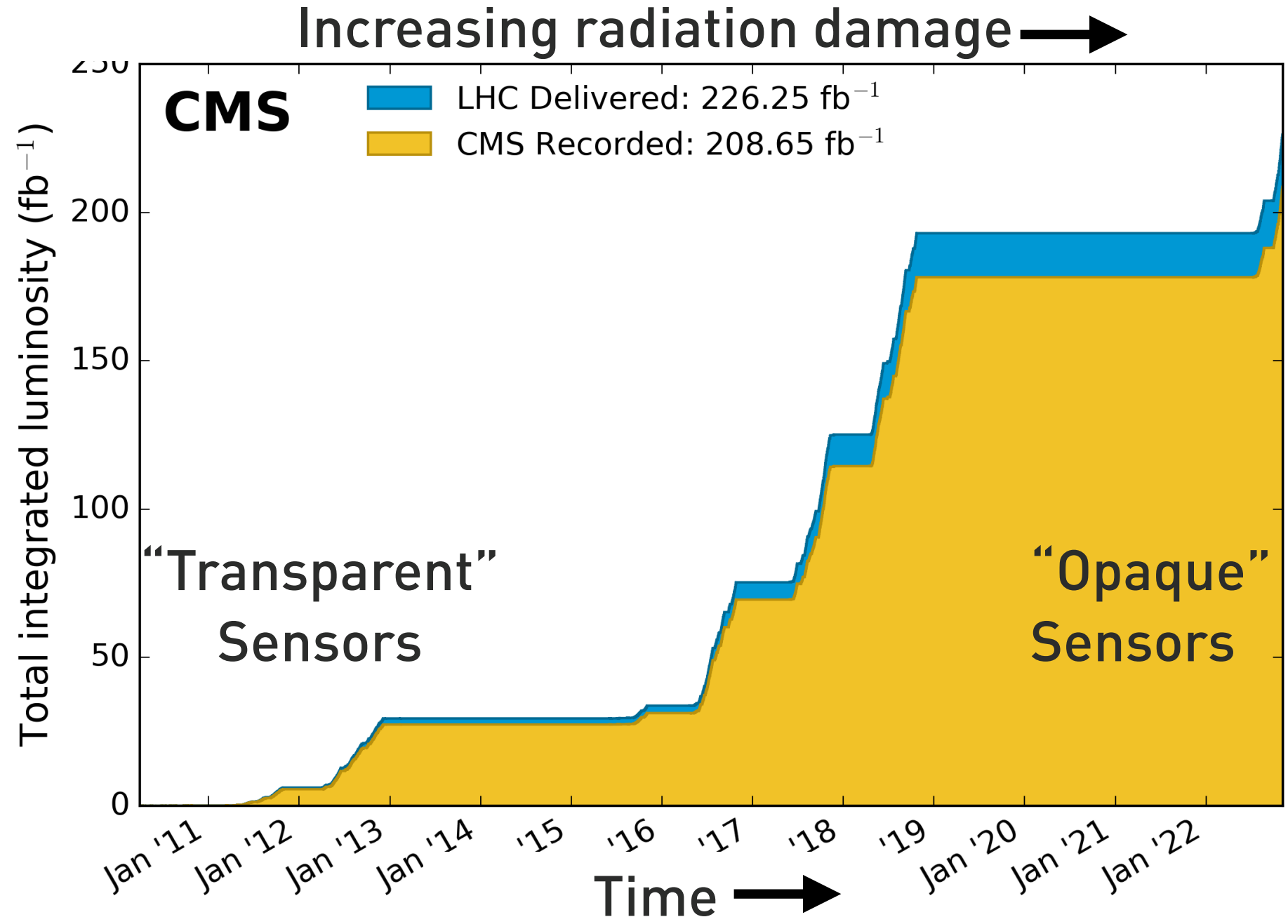- AXOL1TL integrated since late 2023
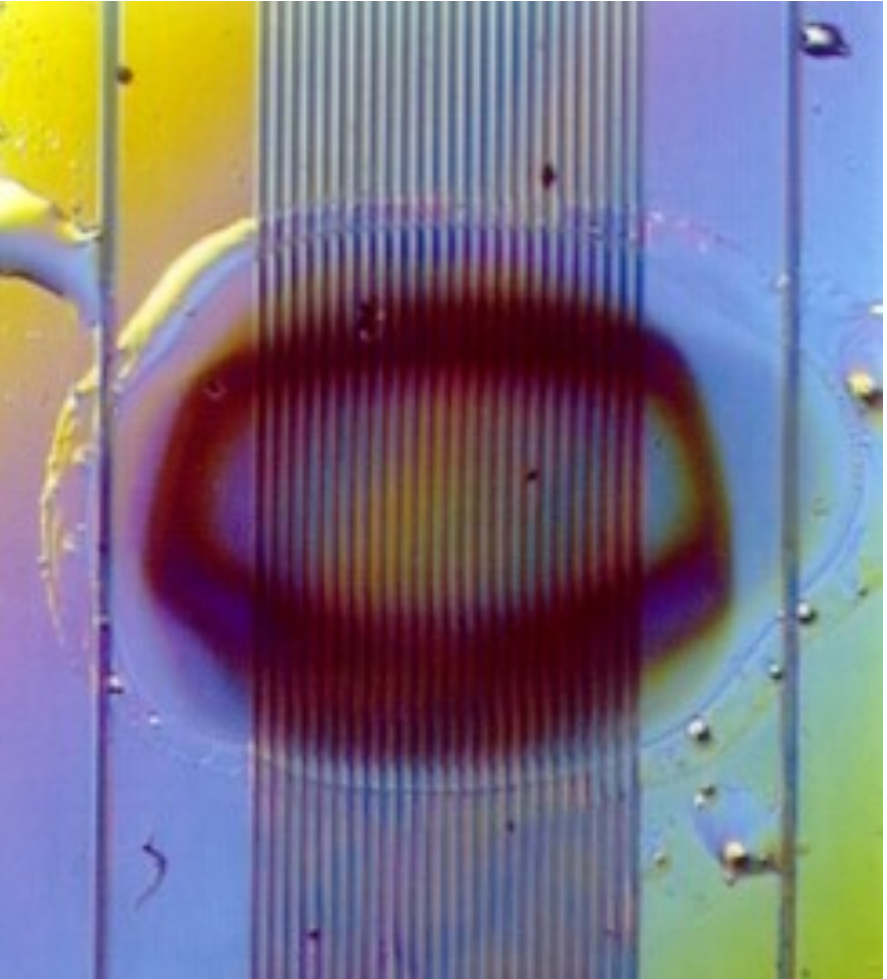
# Continual learning

Increasing radiation damage ➡



CMS

LHC Delivered: 226.25 fb$^{-1}$
CMS Recorded: 208.65 fb$^{-1}$

Total integrated luminosity (fb$^{-1}$)

"Transparent" Sensors

"Opaque" Sensors

Jan '11  Jan '12  Jan '13  Jan '14  '15  '16  '17  '18  '19  Jan '20  Jan '21  Jan '22

Time ➡

Radiation damage of silicon detector



**Many reasons for changing conditions**
- Detector position slightly changes
- Radiation damage

# Continual learning

Increasing radiation damage →



**Radiation damage of silicon detector**



**Many reasons for changing conditions**
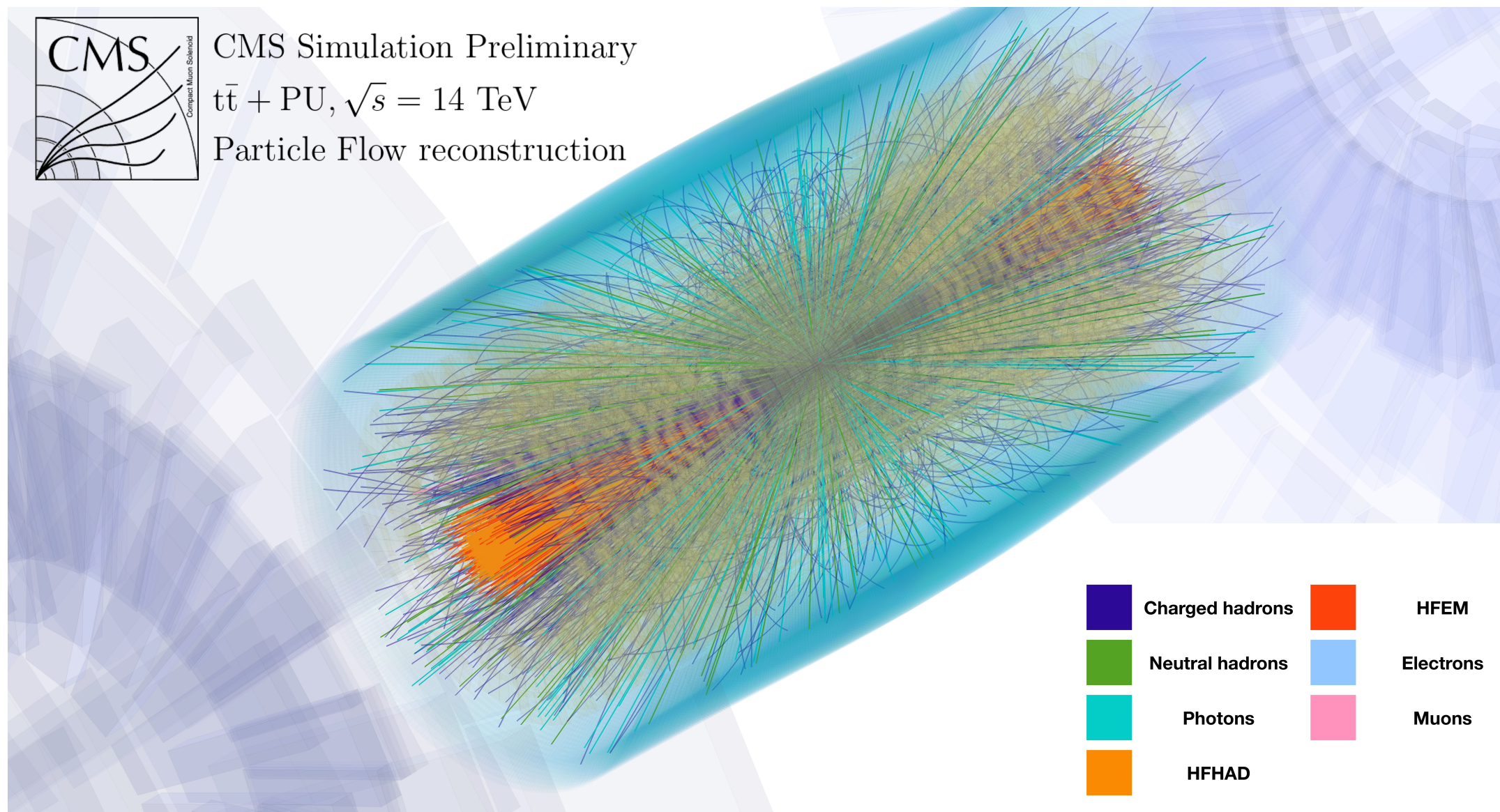- Detector position slightly changes
- Radiation damage

**Continual learning to the aid for self-supervised training?**
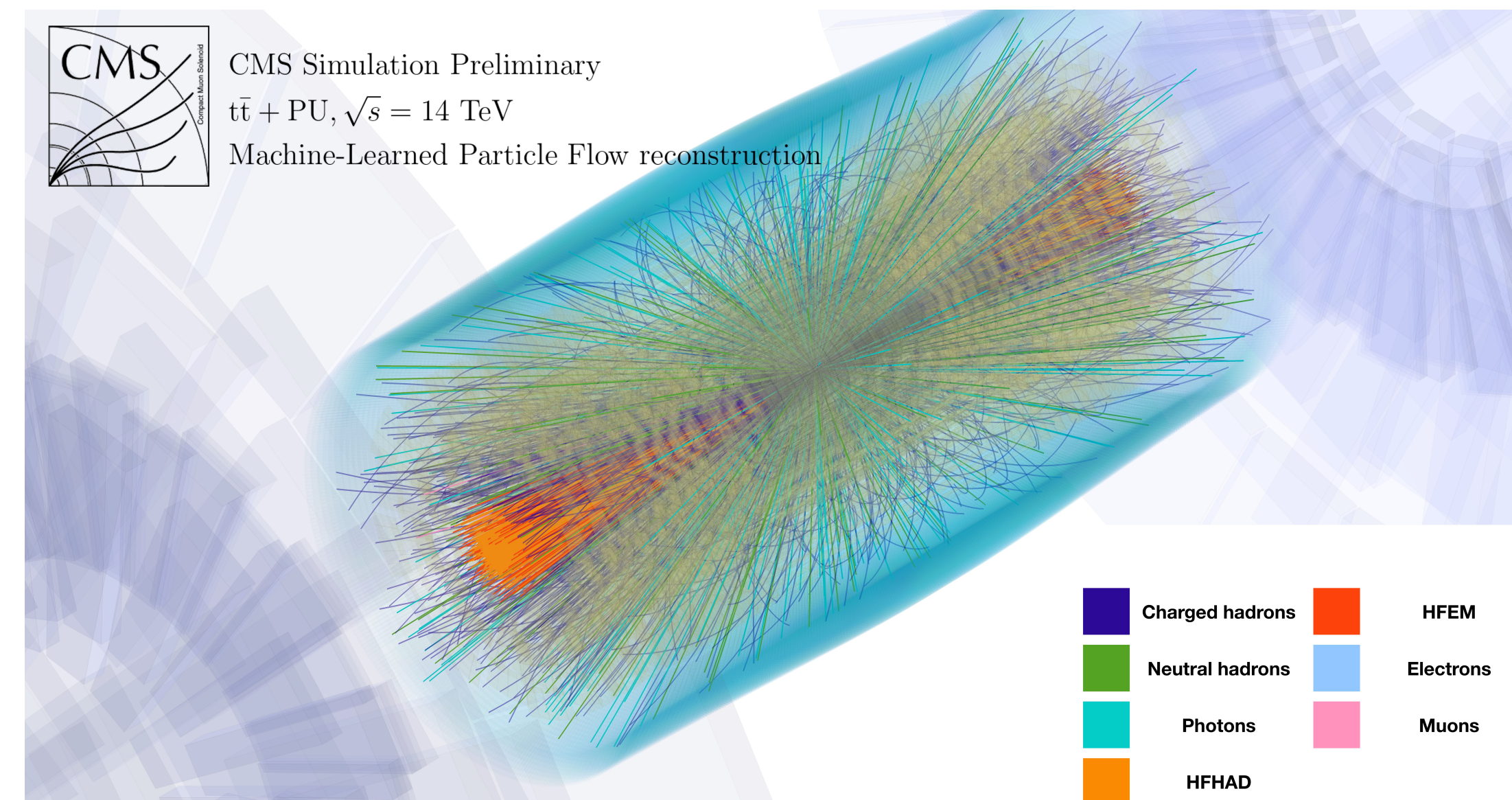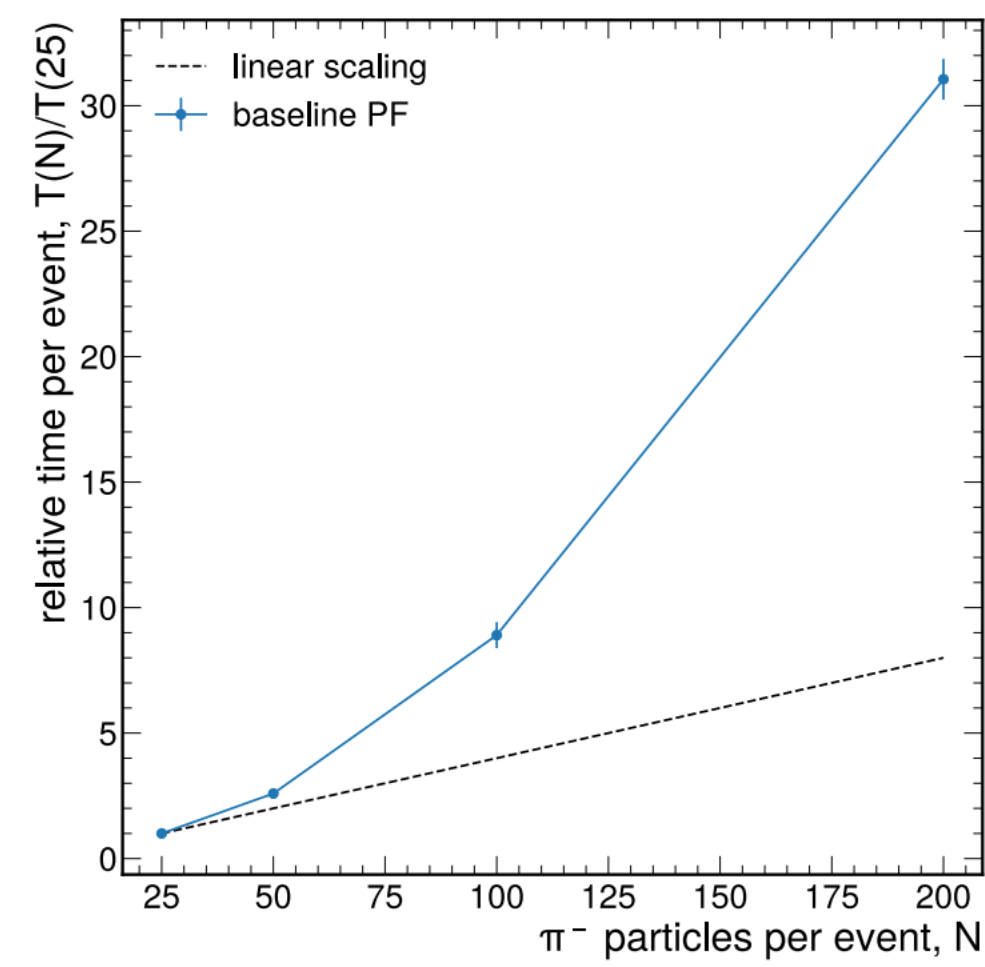- Avoid re-training on TBs of data, adapt to gradual changes!
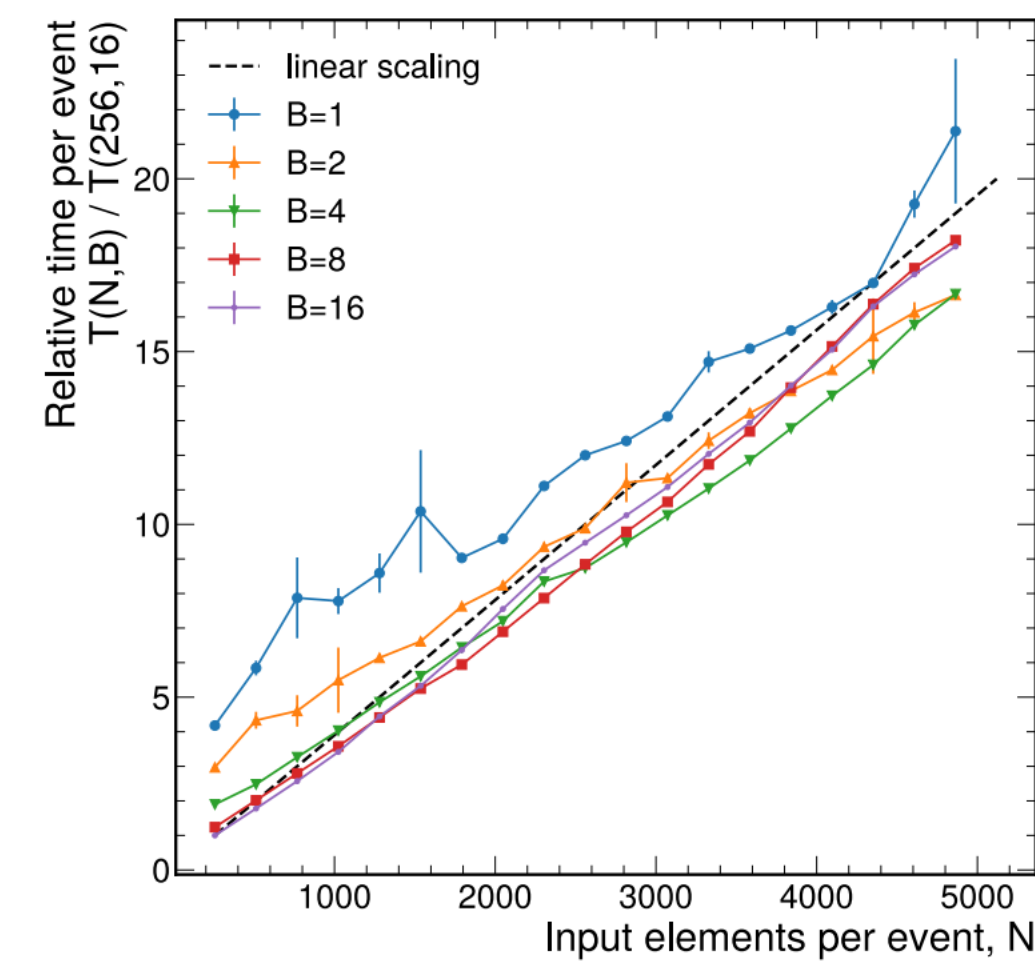
# End-to-end?

## Classical Particle Flow



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

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

## Graph Neural Network



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

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

PF baseline scales non-linearily
with increasing input size



GNN-based model inference time
scales approximately linearly with
increasing input size
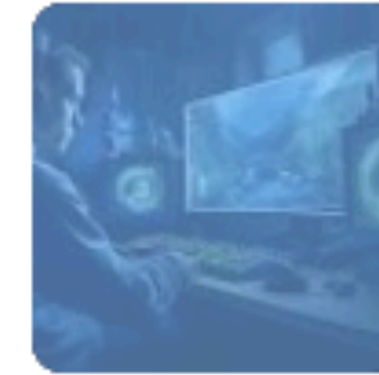
**Classical Particle Flow**



**Graph Neural Network**



arXiv:2101.08578

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...
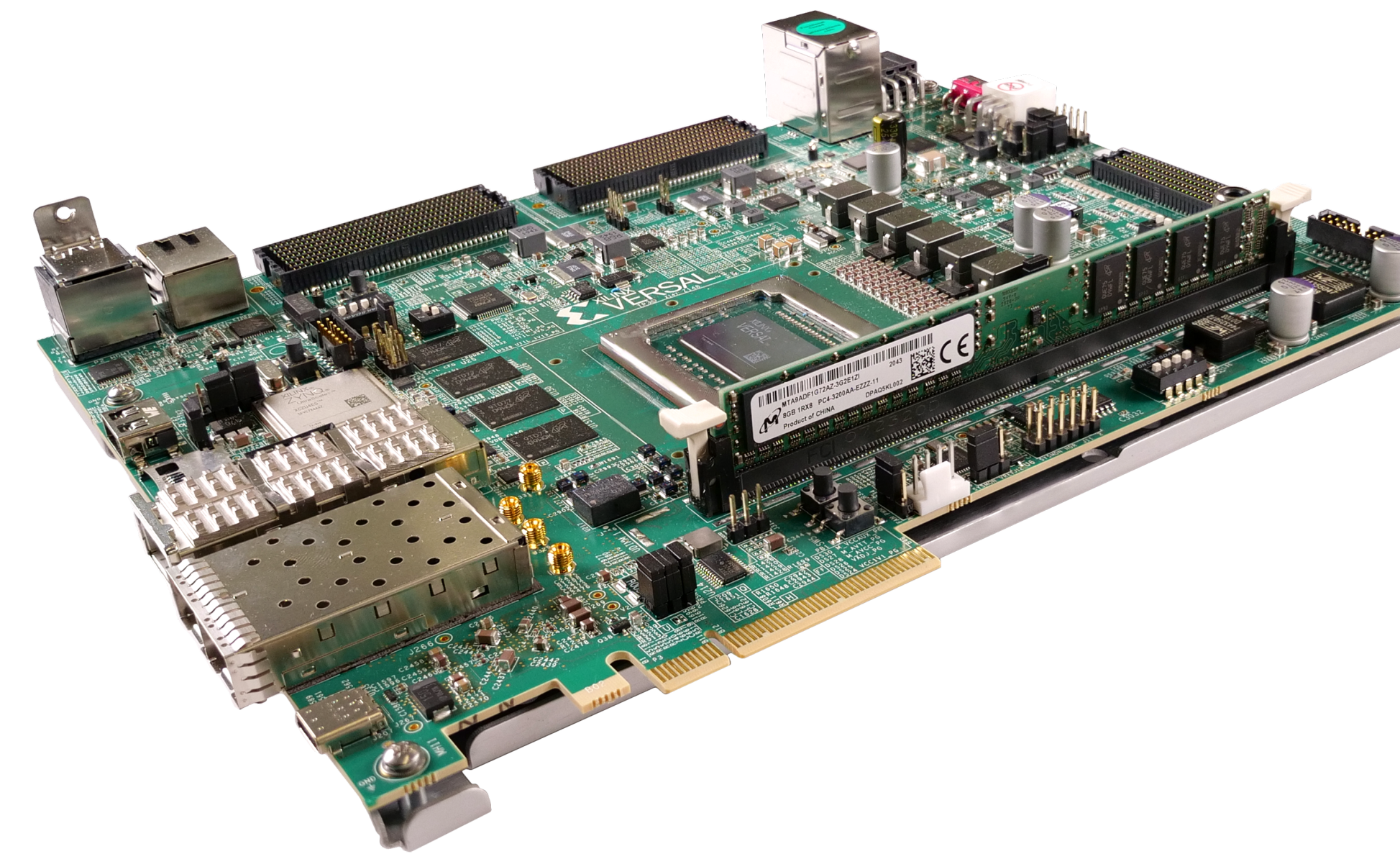
# AI hardware

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
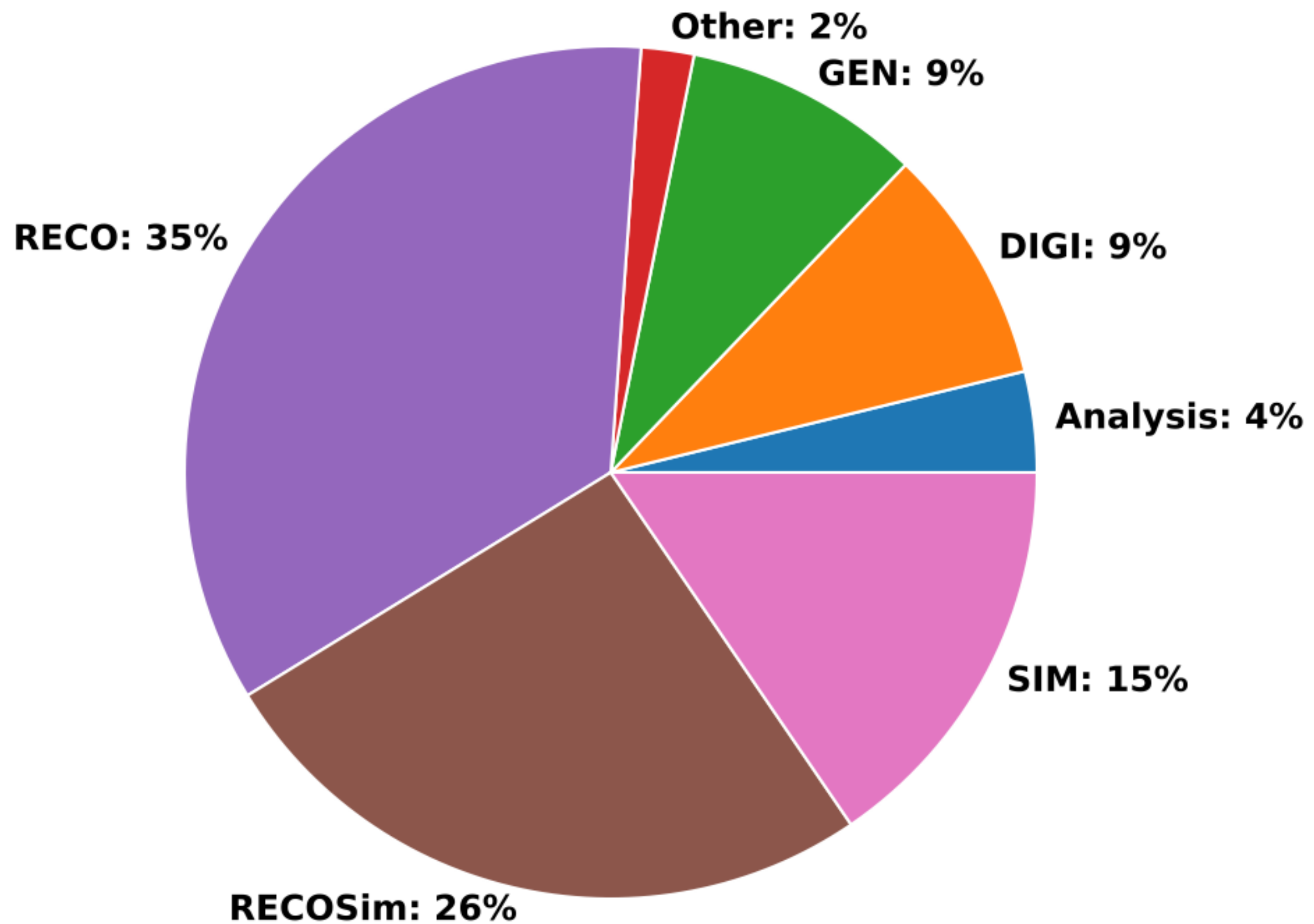- Data move back and forth between AI Engines and FPGA

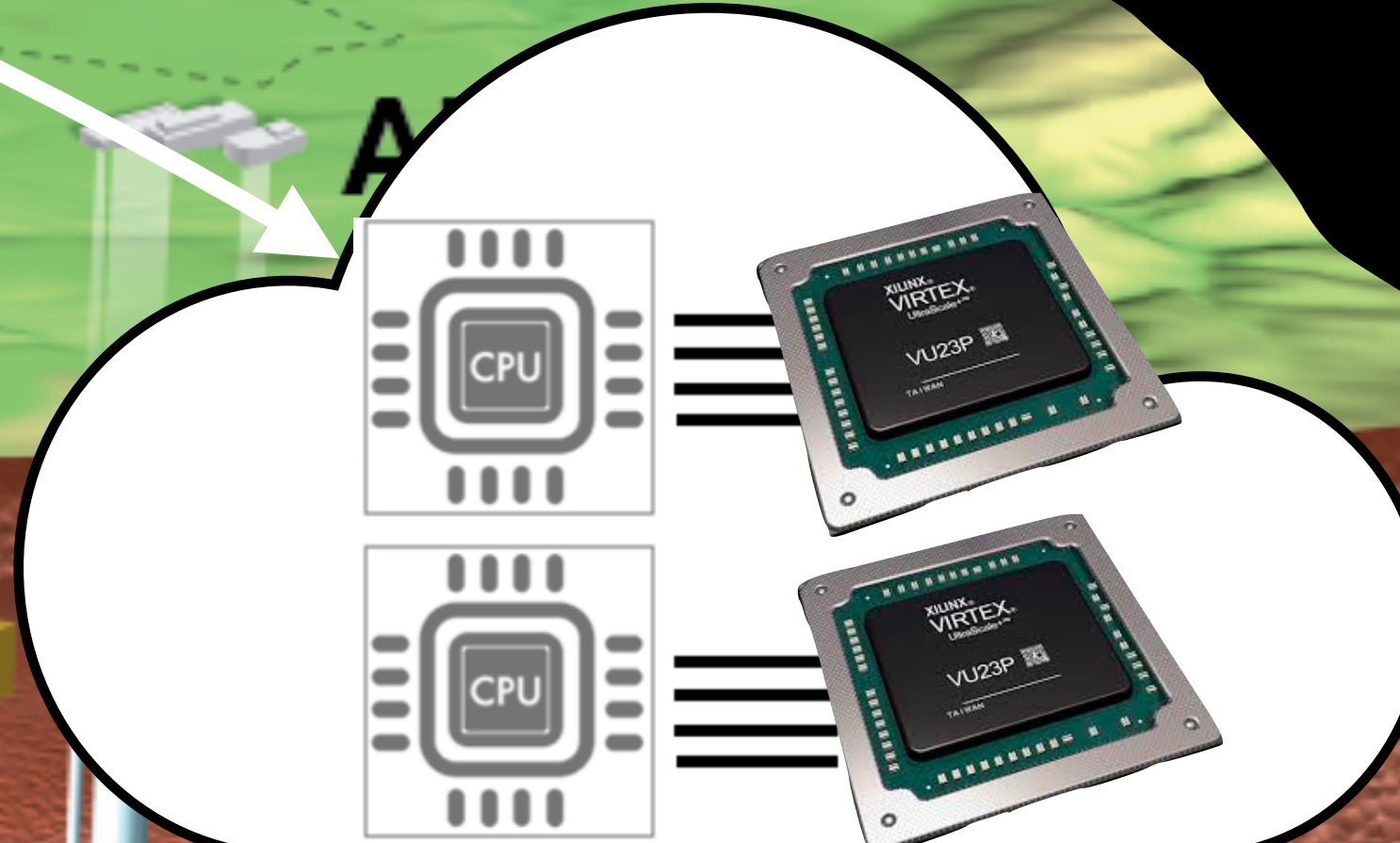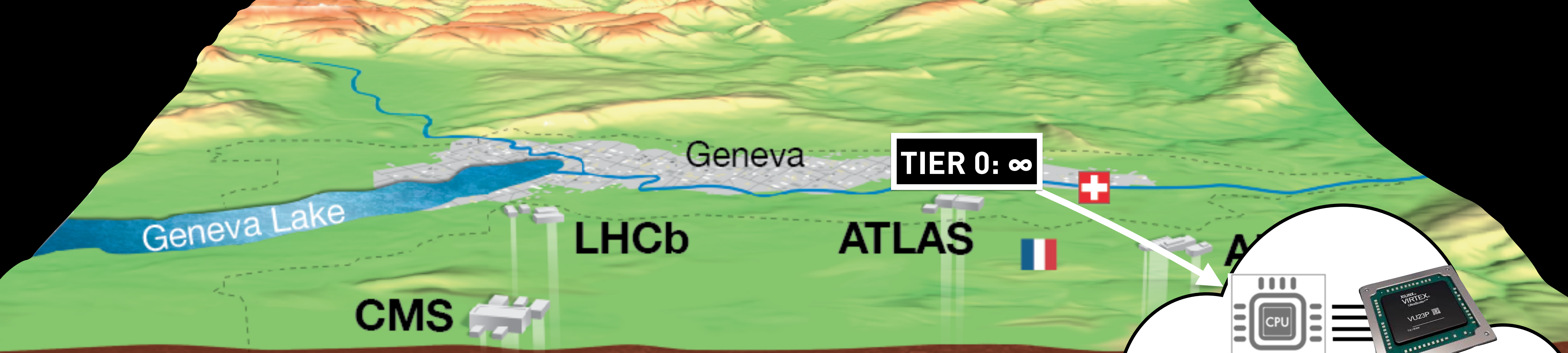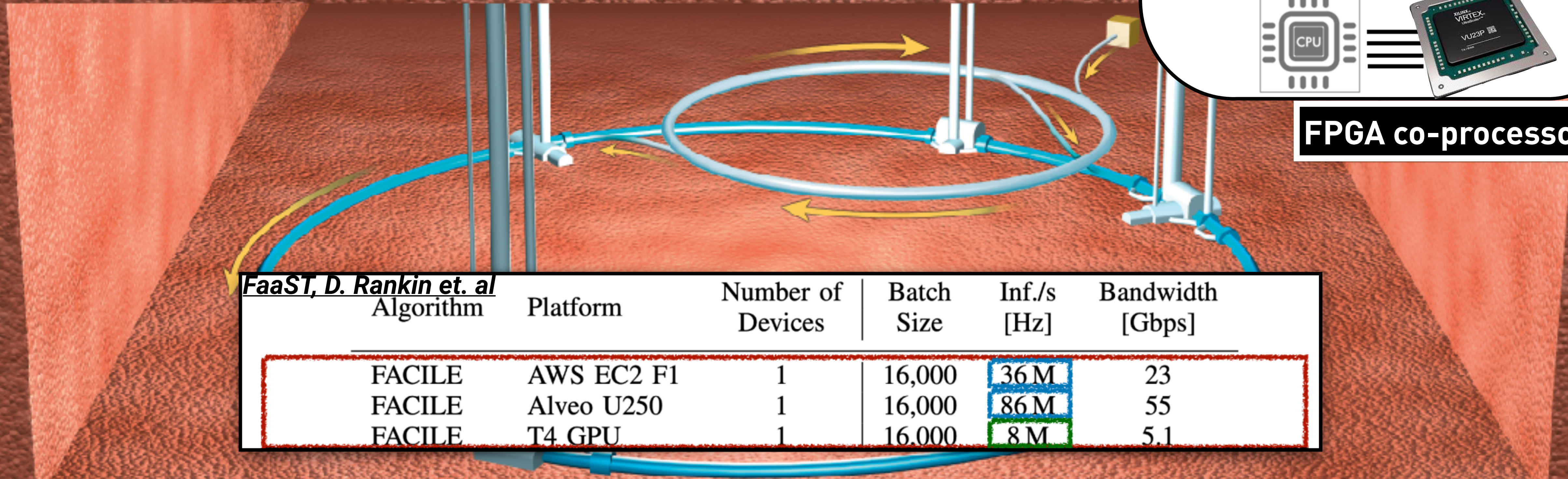Currently explored for real-time tracking in trigger application

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

TIER 0: ∞

Geneva

Geneva Lake

LHCb

ATLAS

CMS

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 |

# FPGAs as accelerators



## DAQ

| Fill: 3819 | E: 6500 GeV | t(SB): 01:30:37 | 03-06-15 |
| --- | --- | --- | --- |

**PROTON PHYSICS: STABLE BEAMS**

| 6500 GeV | I(B1): | 2.92e+11 | I(B2): | 2.91 |

**Detector**

| VELO | ST | OT | RICH | ECal | HCal | Muon |

| FEE | FEE | FEE | FEE | FEE | FEE | FEE |

| Readout Board | Readout Board | Readout Board | Readout Board | Readout Board | Readout Board | Readout Board |

**Front-End**

equest

**READOUT NETWORK**

## CO-PROCESSING

| LHC Page1 | Fill: 4724 | E: 0 Z GeV | 15-12-15 14:28:21 |
| --- | --- | --- | --- |

**SHUTDOWN: NO BEAM**

CPU — COPROCESSOR (GPU/FPGA/ASICS)

CPU — COPROCESSOR (GPU/FPGA/ASICS)

neutral hadron

photon

charged hadrons

detector

particle-flow

HCAL Clusters

ECAL Clusters

Tracks

*C. Beteta, I. Bezshyiko, N. Serra*

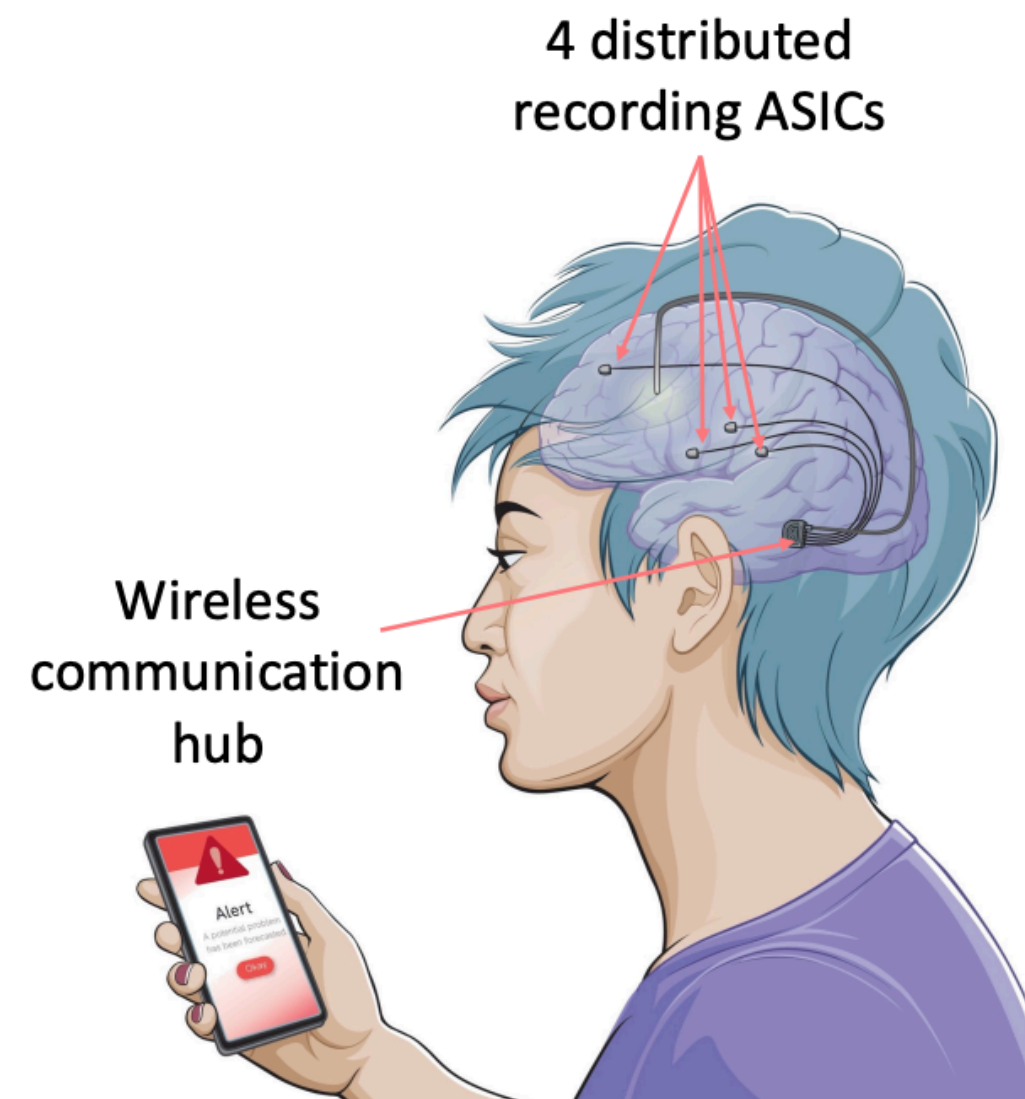*F. Capel et al.*

**Bioluminescence bursts up to few MHz!**

# ...and outside of particle physics


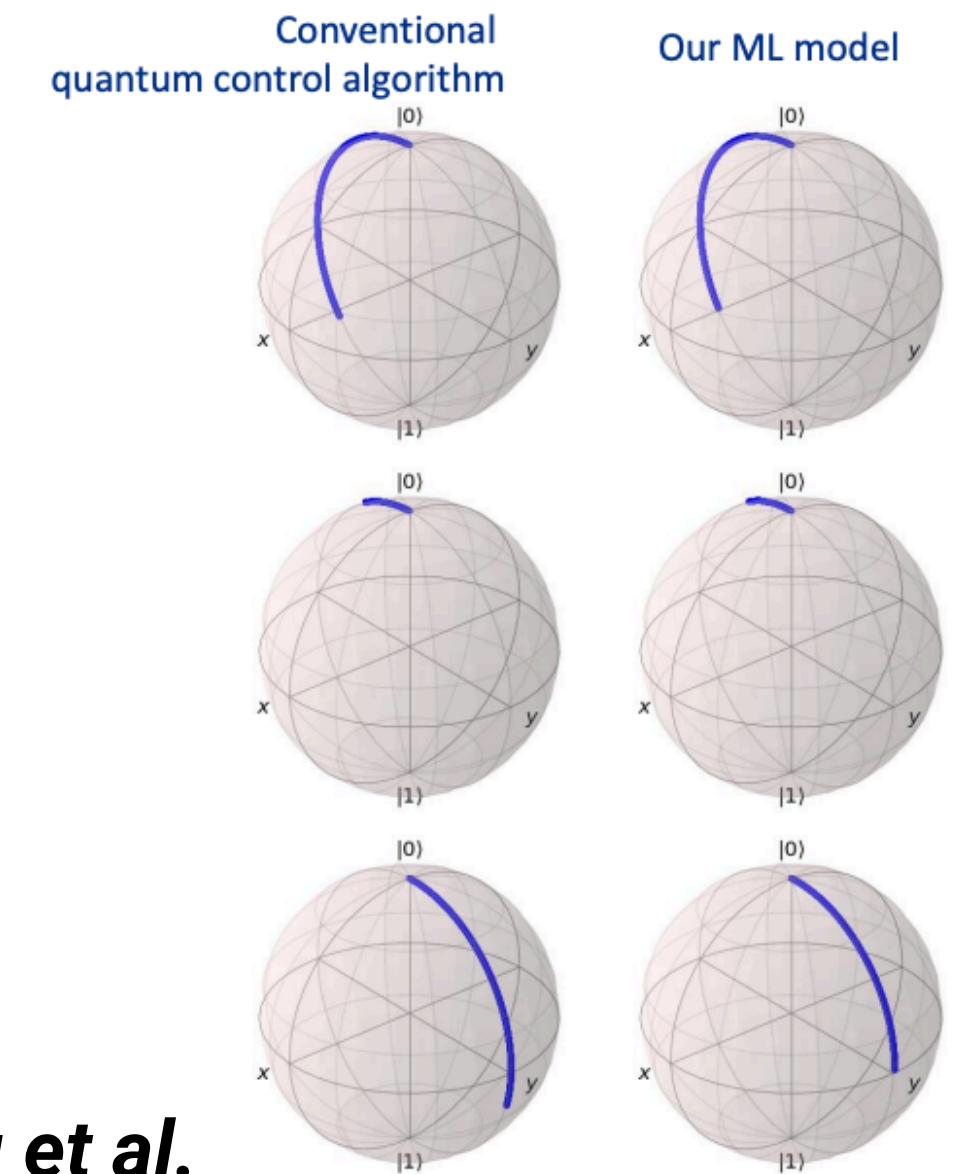Semantic segmentation for autonomous vehicles

***N. Ghielmetti et al.***


Seizure Predicting Brain Implant

4 distributed
recording ASICs

Wireless
communication
hub

***W. Lemaire et al.***


NN accelerator for **quantum control**
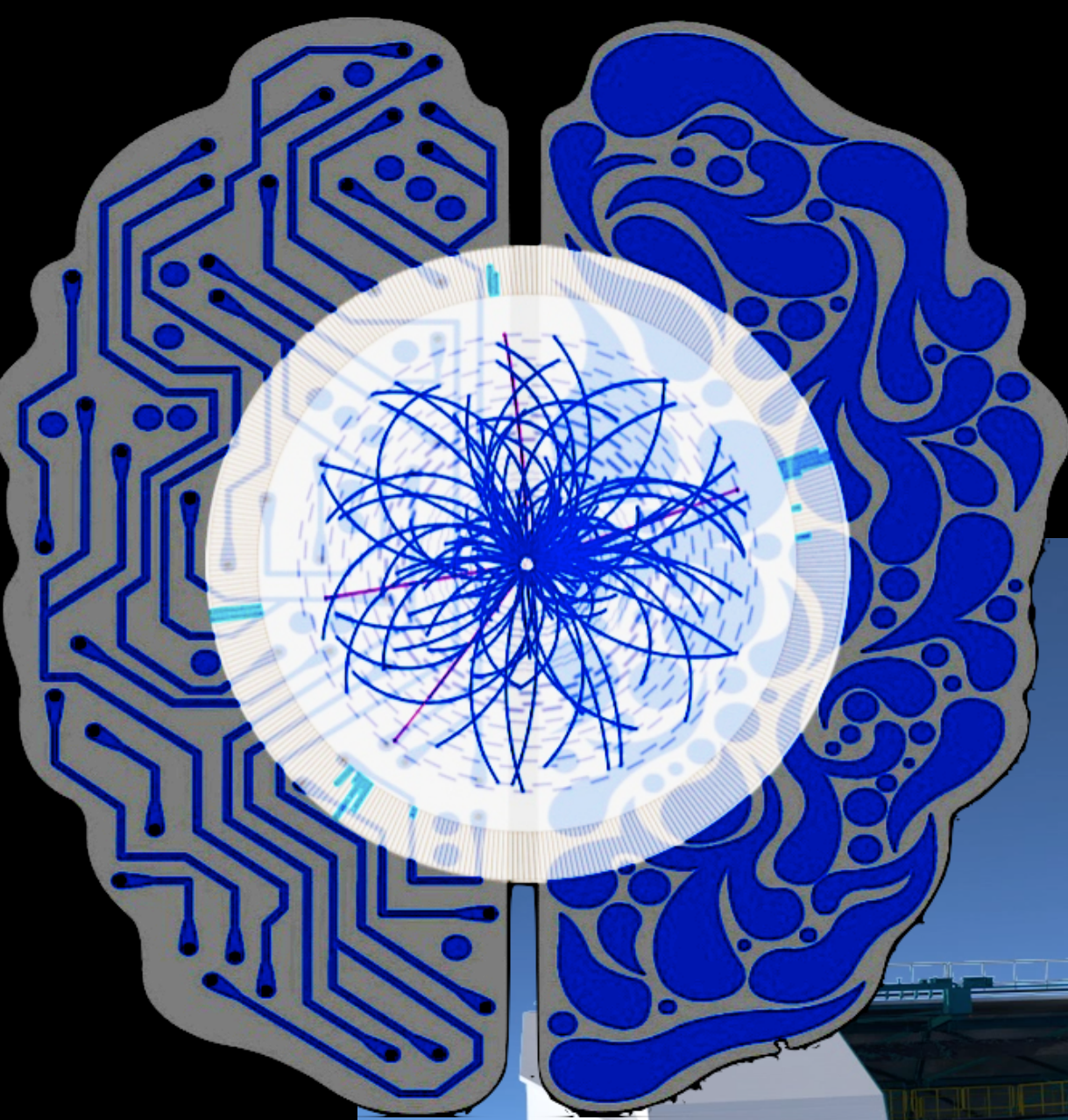- Putting control in cryostat
  (e.g optimal pulse parameters)

Conventional
quantum control algorithm    Our ML model

***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***
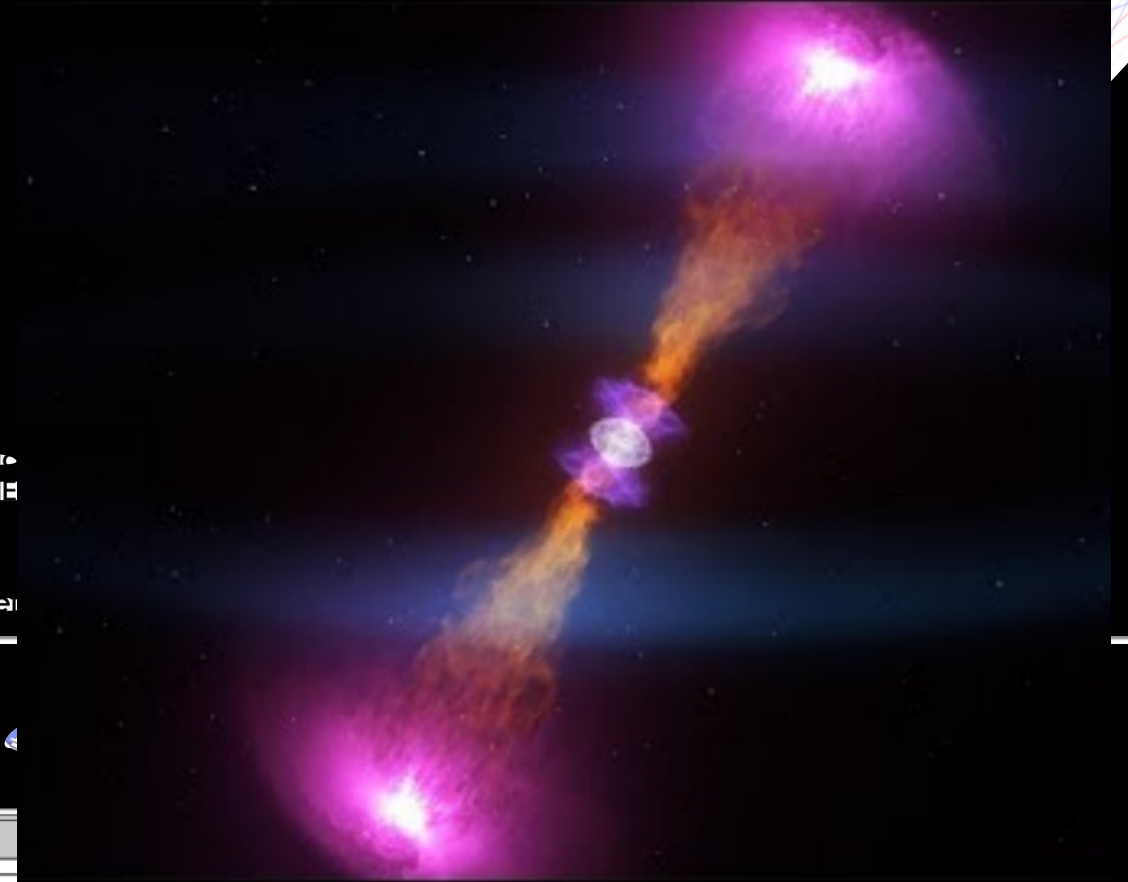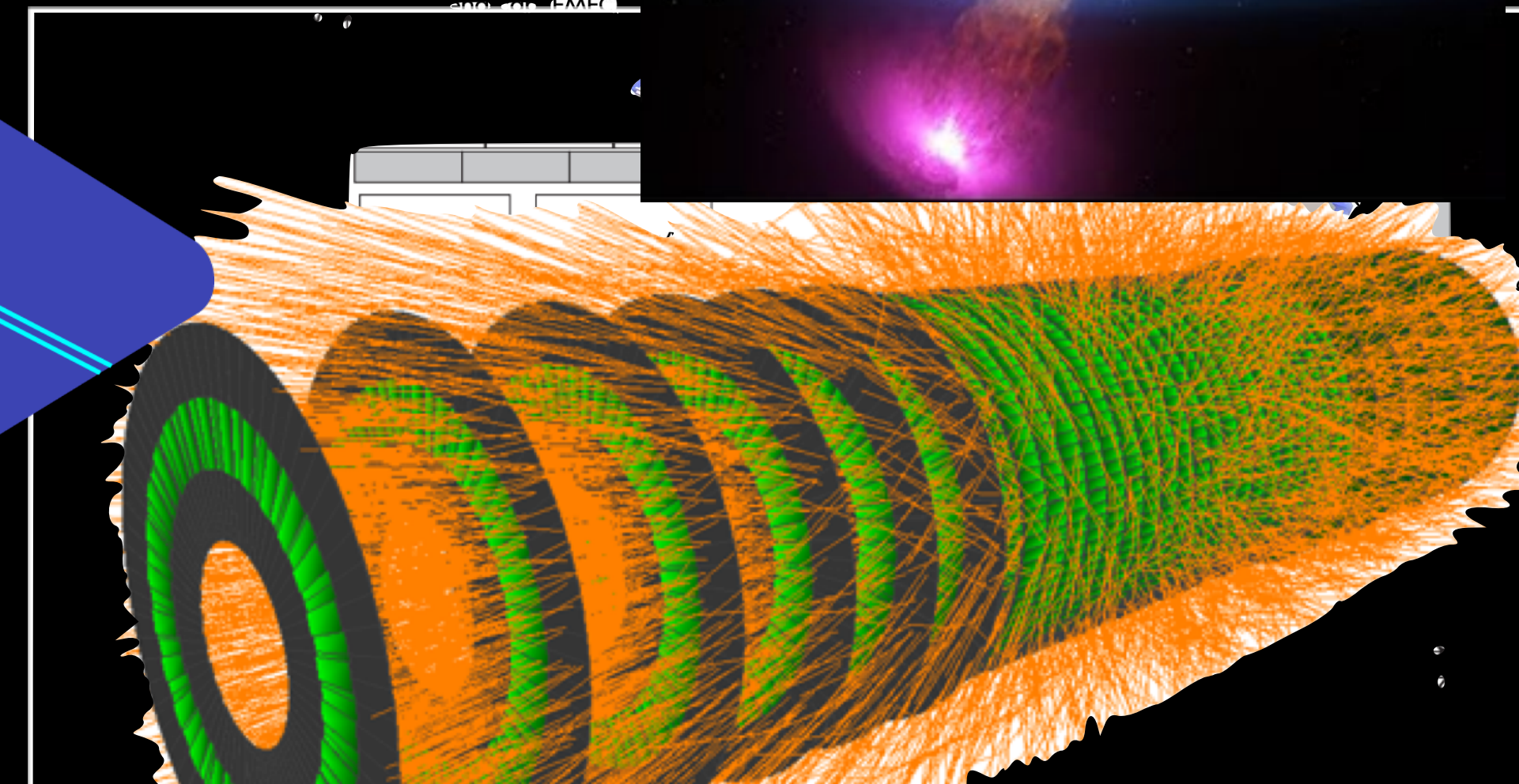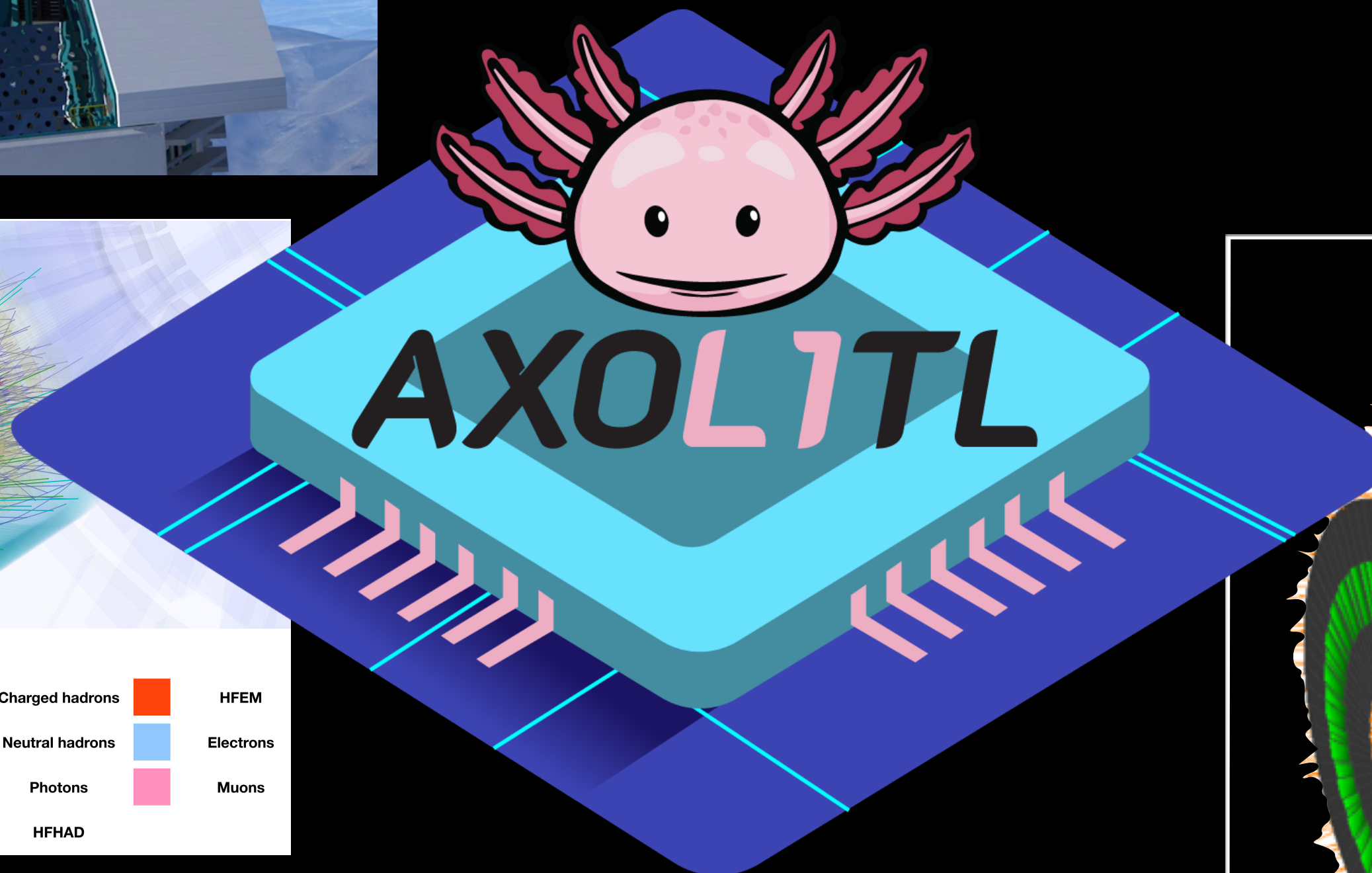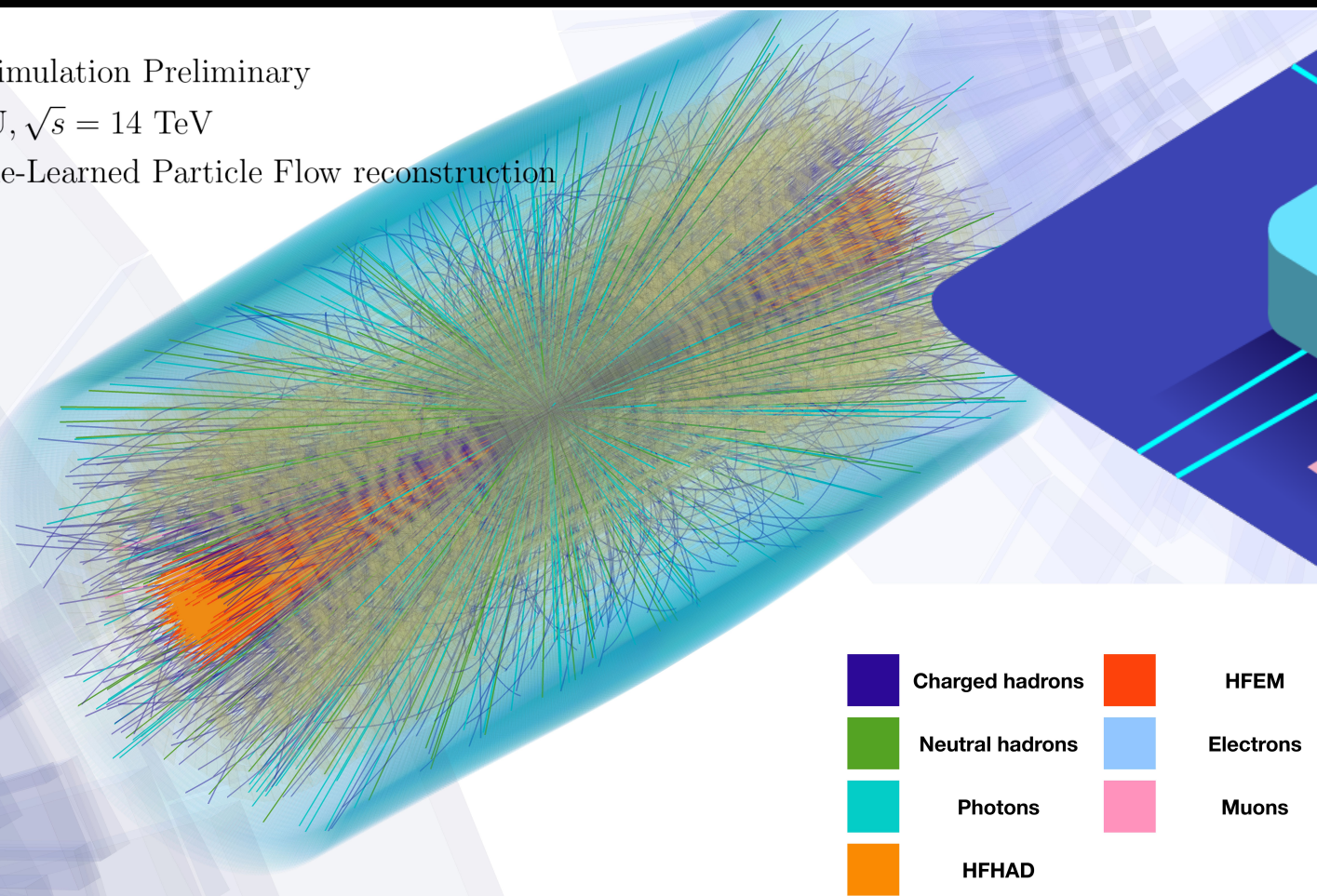- ***Food contamination detection***
- etc....

Conifer

hls4ml

Charged
seed

Join us at the

**FastML Lab!**

AXOLITL

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

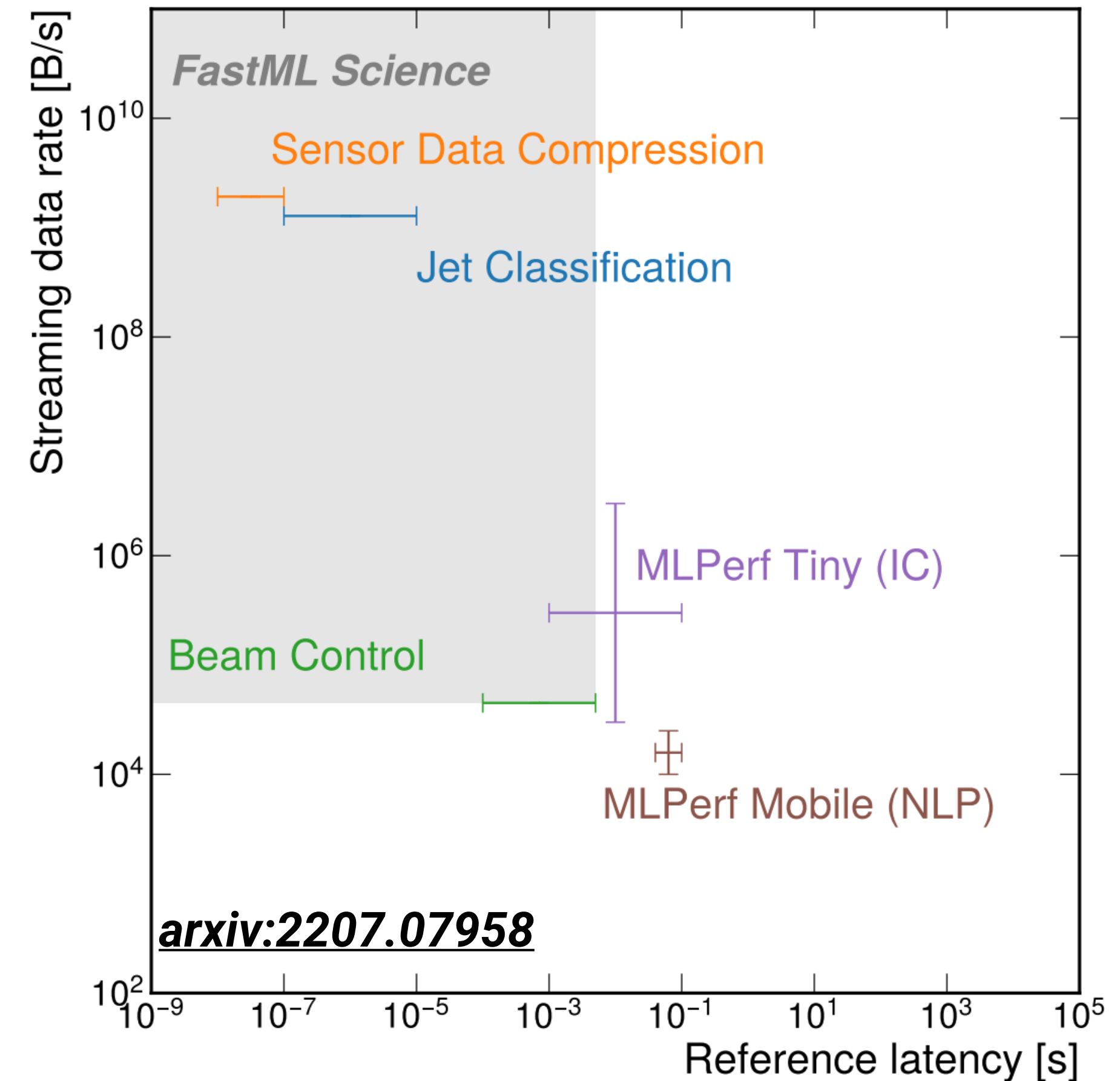| Charged hadrons | HFEM |
| Neutral hadrons | Electrons |
| Photons | Muons |
| HFHAD | |

# 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/*



*arxiv:2207.07958*

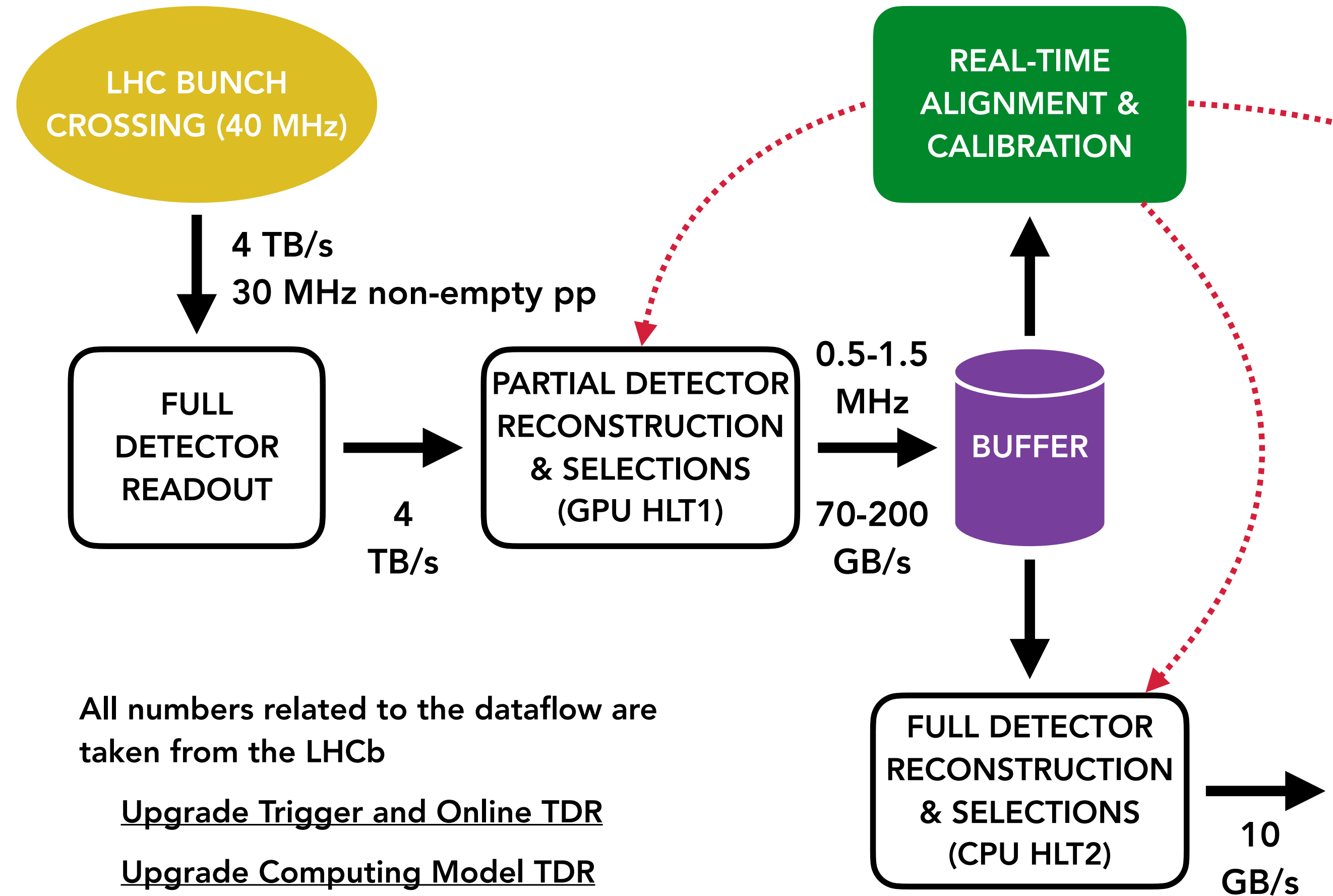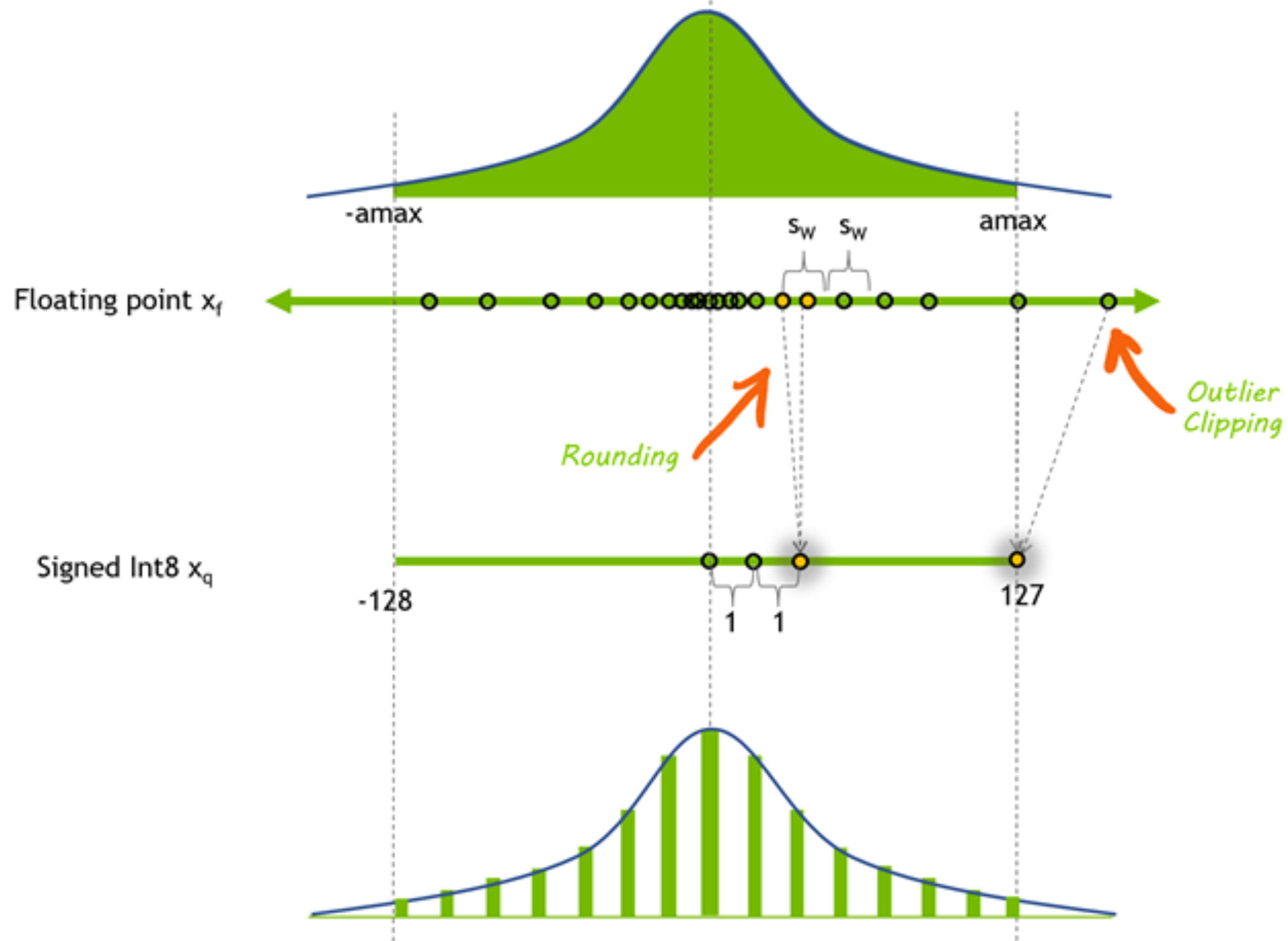| Model | LUT | | LUTRAM | | FF | | BRAM [36 kb] | | DSP | | Latency [ms] | Energy/inf. [$\mu$J] |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | Pynq-Z2 | | | | | | |
| Available | 53 200 | | 17 400 | | 106 400 | | 140 | | 220 | | – | – |
| IC (hls4ml) | 28 544 | 53.7% | 3 756 | 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:<br>  - Run events in parallel<br>  - Reconstruct tracks in parallel | Good for<br>  - Data-intensive parallelizable applications<br>  - 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 |

**LHC BUNCH CROSSING (40 MHz)**

4 TB/s
30 MHz non-empty pp

**FULL DETECTOR READOUT**

4 TB/s

**PARTIAL DETECTOR RECONSTRUCTION & SELECTIONS (GPU HLT1)**

0.5-1.5 MHz

70-200 GB/s

**BUFFER**

**REAL-TIME ALIGNMENT & CALIBRATION**

**FULL DETECTOR RECONSTRUCTION & SELECTIONS (CPU HLT2)**

10 GB/s

All numbers related to the dataflow are taken from the LHCb

Upgrade Trigger and Online TDR

Upgrade Computing Model TDR

*C. Fitzpatrick*

# Q PYT🔥RCH ?

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.)***



*QAP (B. Hawks et al.)*



*HAWQ*

# Q P Y T Ö R C H ?

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→1 MHz



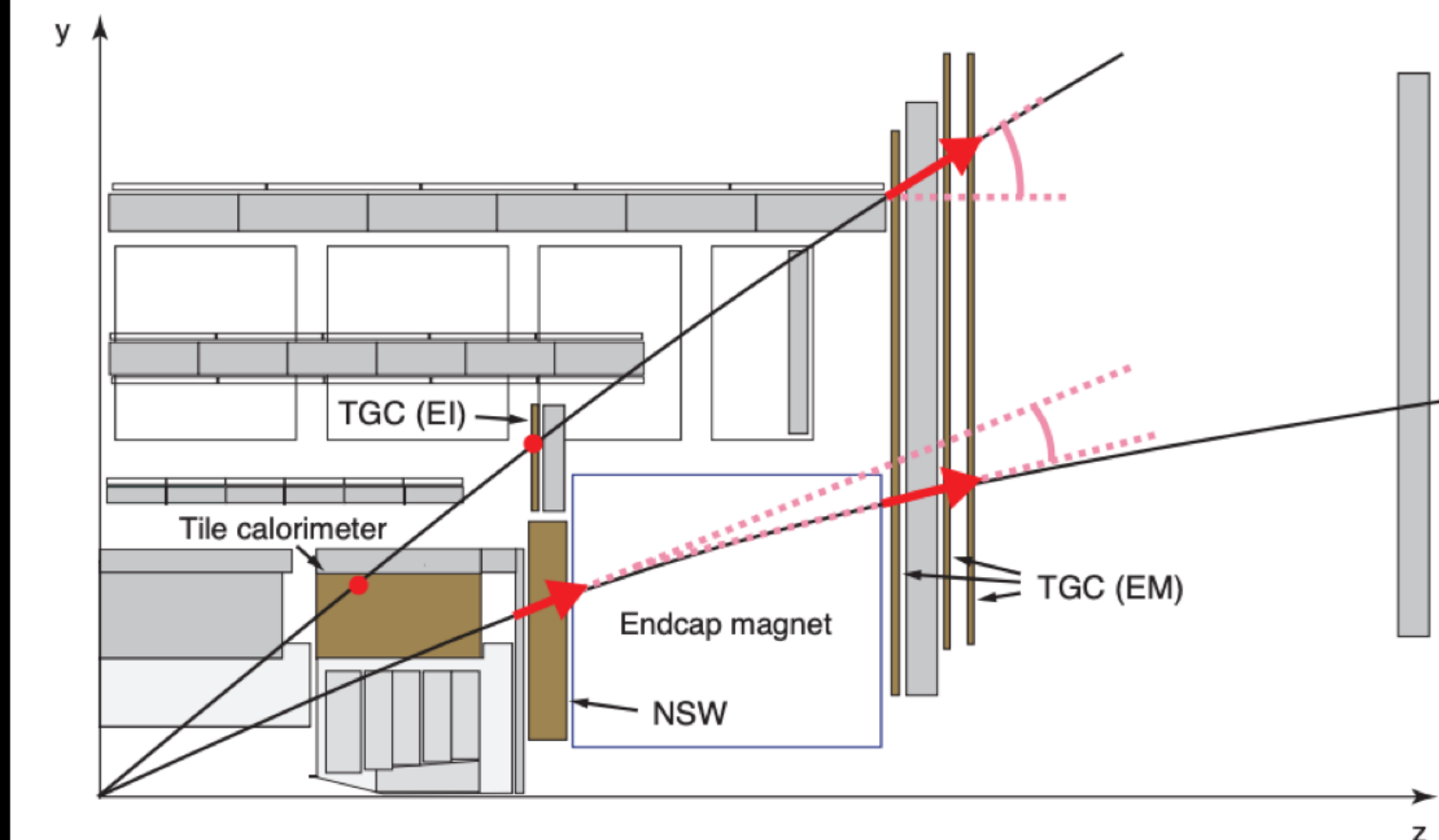*CHEP 2019, P. Vicini*

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*

ATLAS small wheel muon segment finding and reconstruction (R&D)
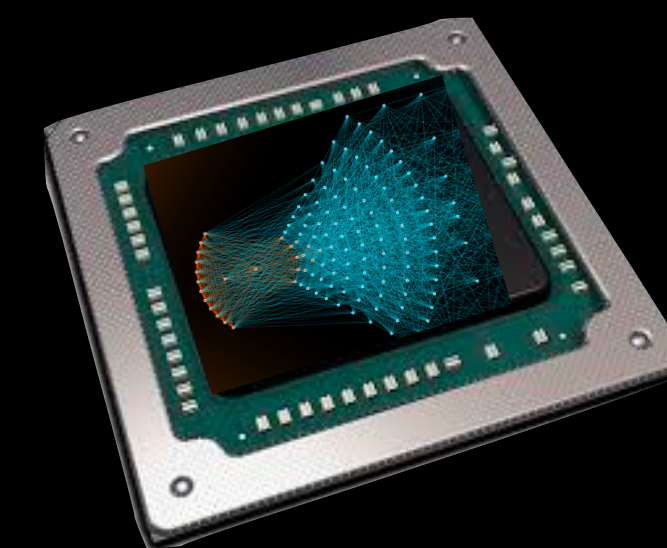- Regression of muon position and angle
- 400 ns budget
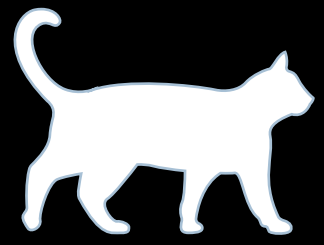


*R. Teixeira de Lima, R Rojas Caballero et al.*

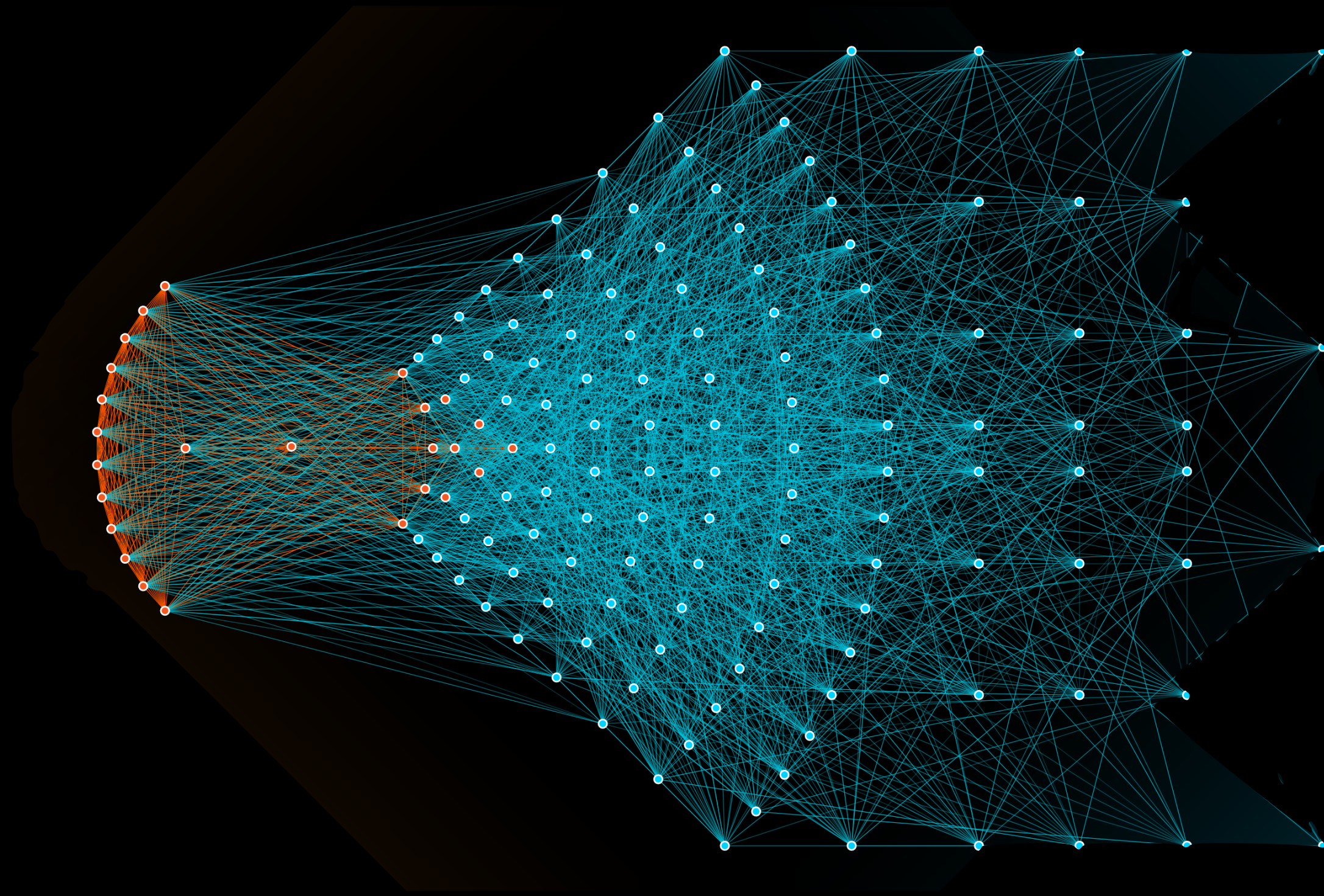Can we have the best of both worlds?

→ Knowledge Distillation

Train

Inference

Cat

Dog

Cat
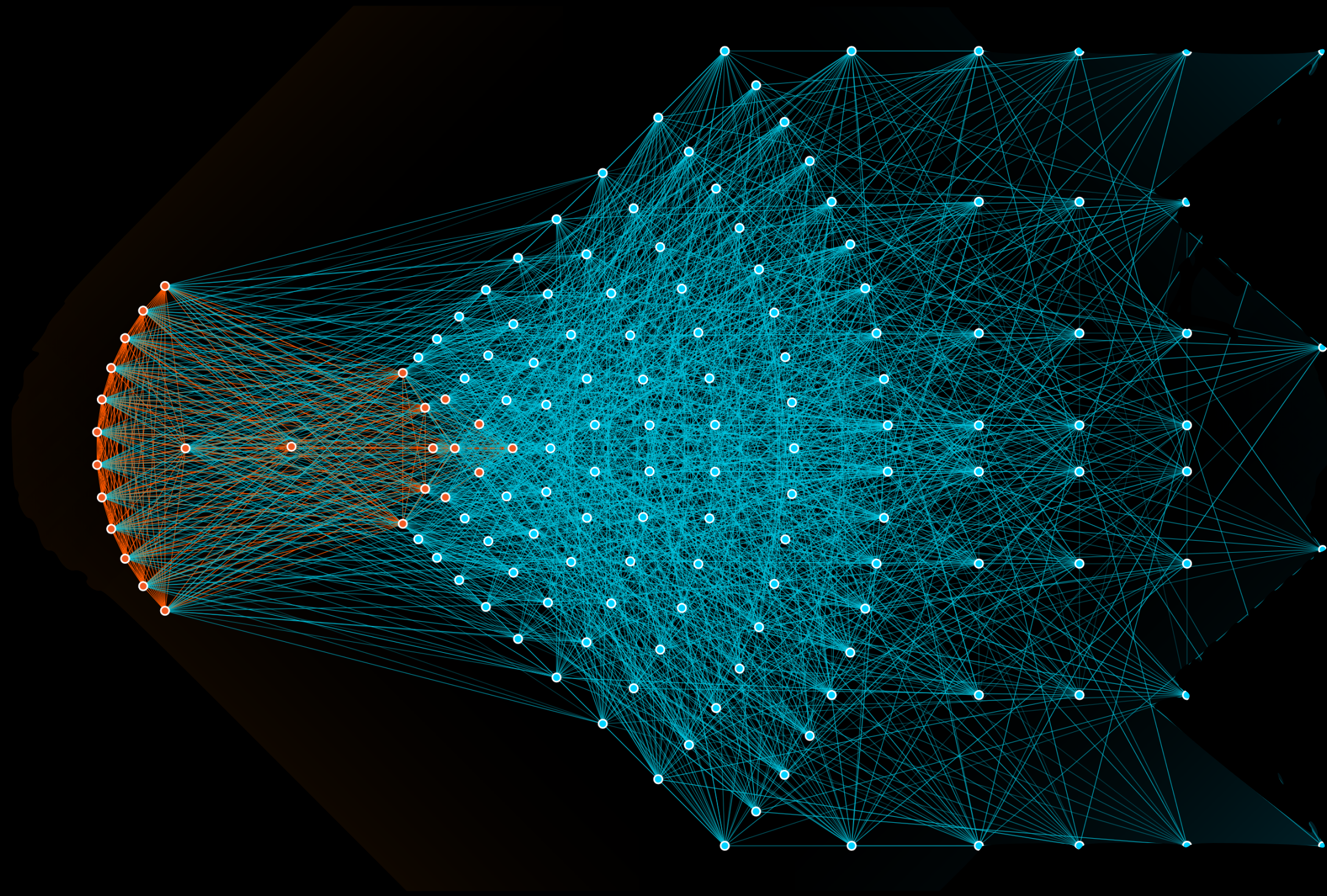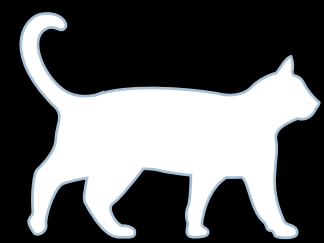
is cat

is dog

Cat

Predicted labels

Teacher
(already trained)

is cat = 0.89

is dog = 0.11

Cat

Predicted labels

Teacher
(already trained)

is cat = 0.89

is dog = 0.11

True labels

is cat = 1

is dog = 0

Cat

Teacher
(already trained)
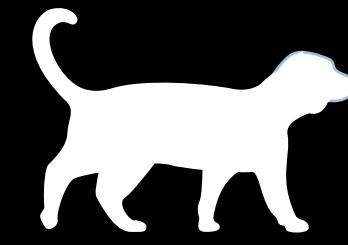
Predicted labels

is cat = 0.89

is dog = 0.11

Distilled
knowledge

True labels

is cat = 1

is dog = 0

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!

# FPGAs as accelerators

## DAQ



## CO-PROCESSING



*C. Beteta, I. Bezshyiko, N. Serra*

# High throughput GPU triggers

**40 MHz pp collisions**



4 TB/s
30 MHz non-empty pp

**_Detector read-out O(500) FPGAs_**

4 TB/s → **350 GPUs 5 TB/s @ 30 MHz of raw data** →

0.5-1.5 MHz

70-200 GB/s

**O(10) PB buffer** →

**REAL-TIME ALIGNMENT & CALIBRATION**

**CPU-based (full reco & trigger)** →

10 GB/s



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

Decays / (0.04 ps)

LHCb
6 fb$^{-1}$

$t$ [ps]

*ML already used in LHCb GPU trigger for particle ID,*
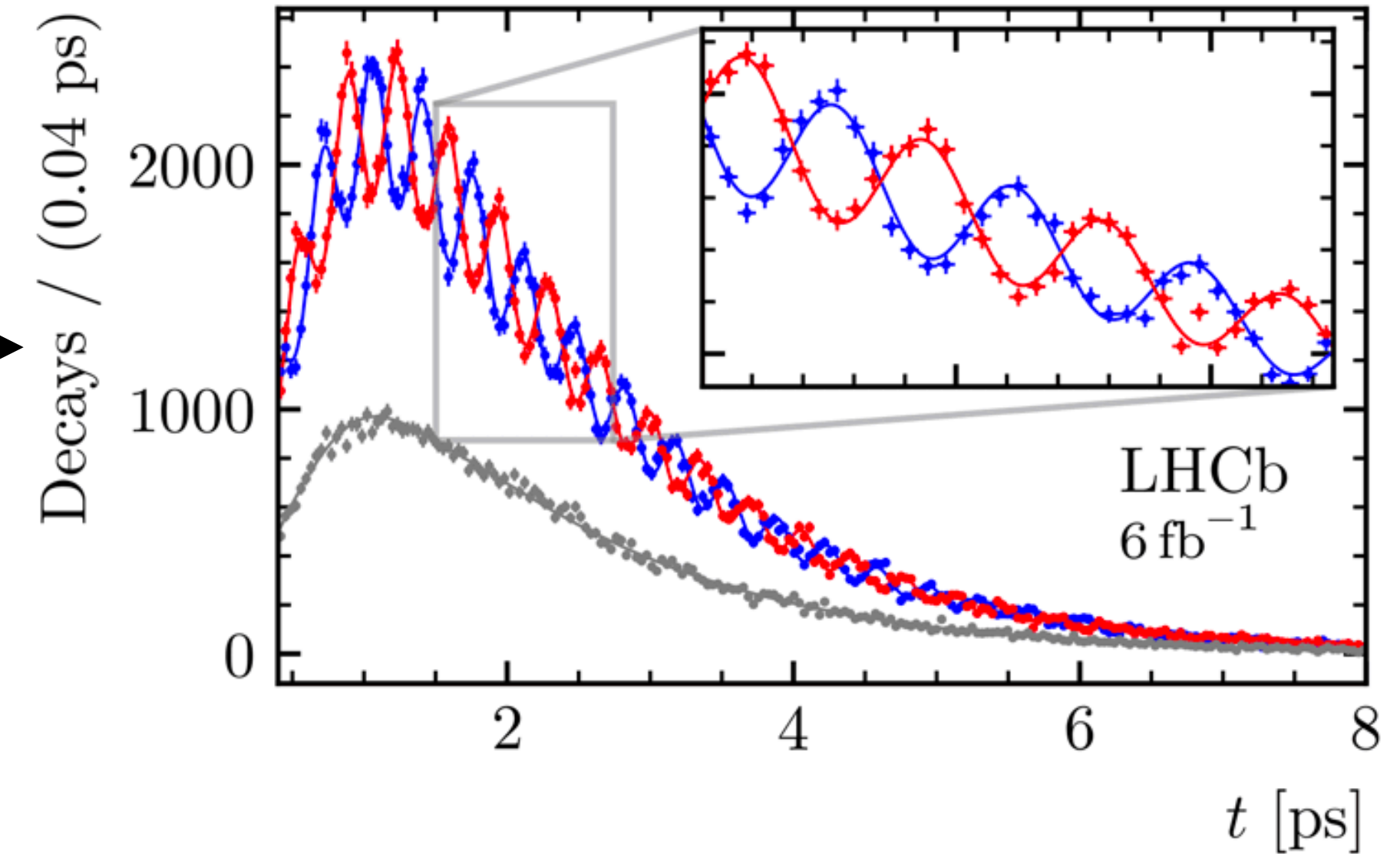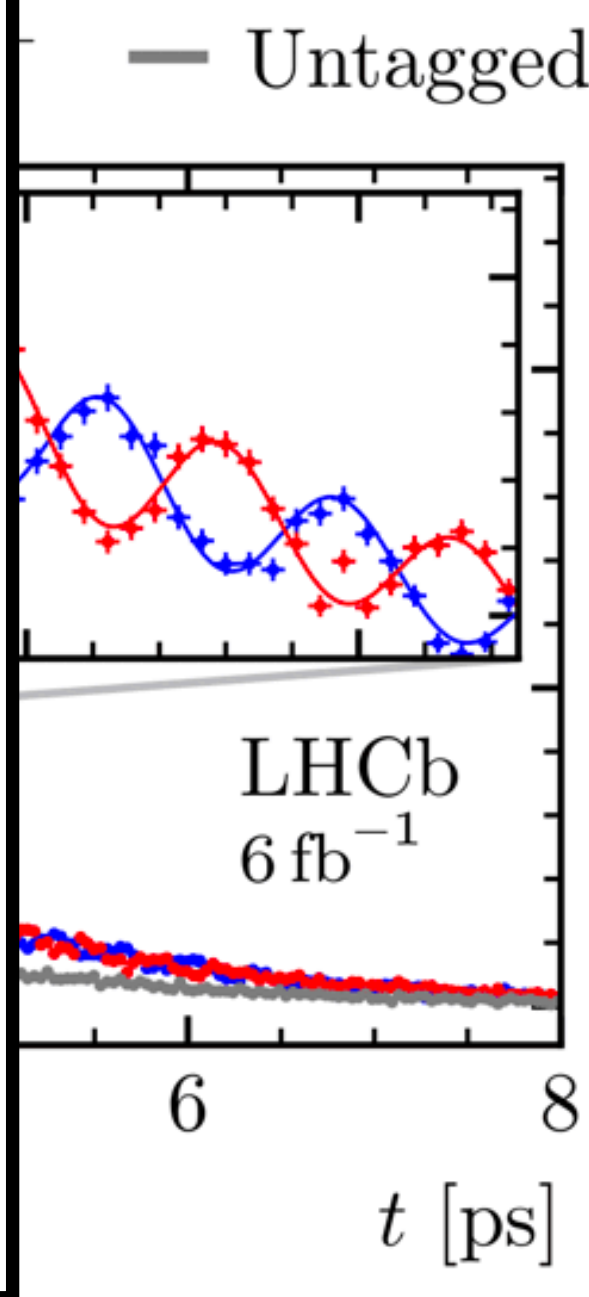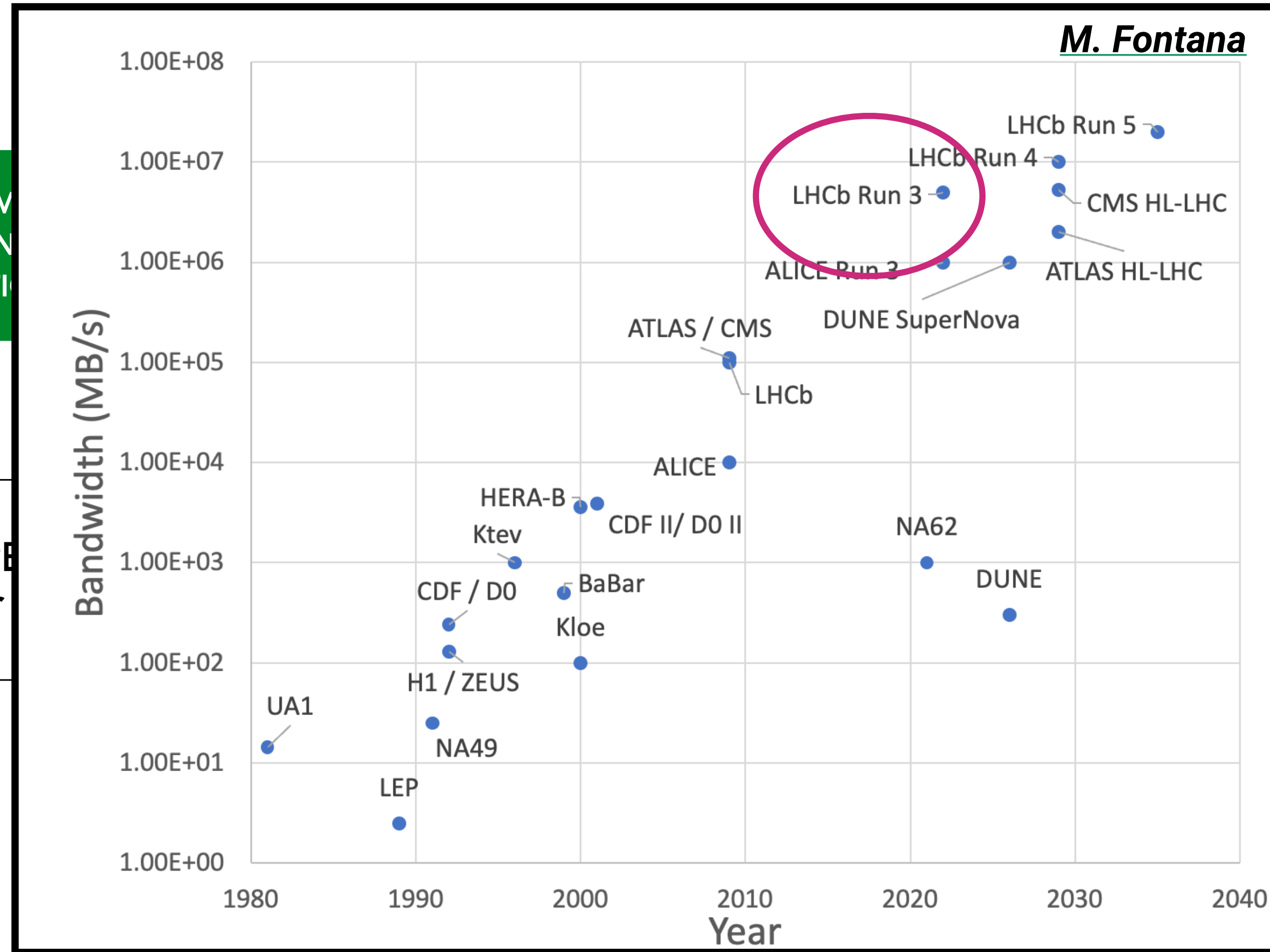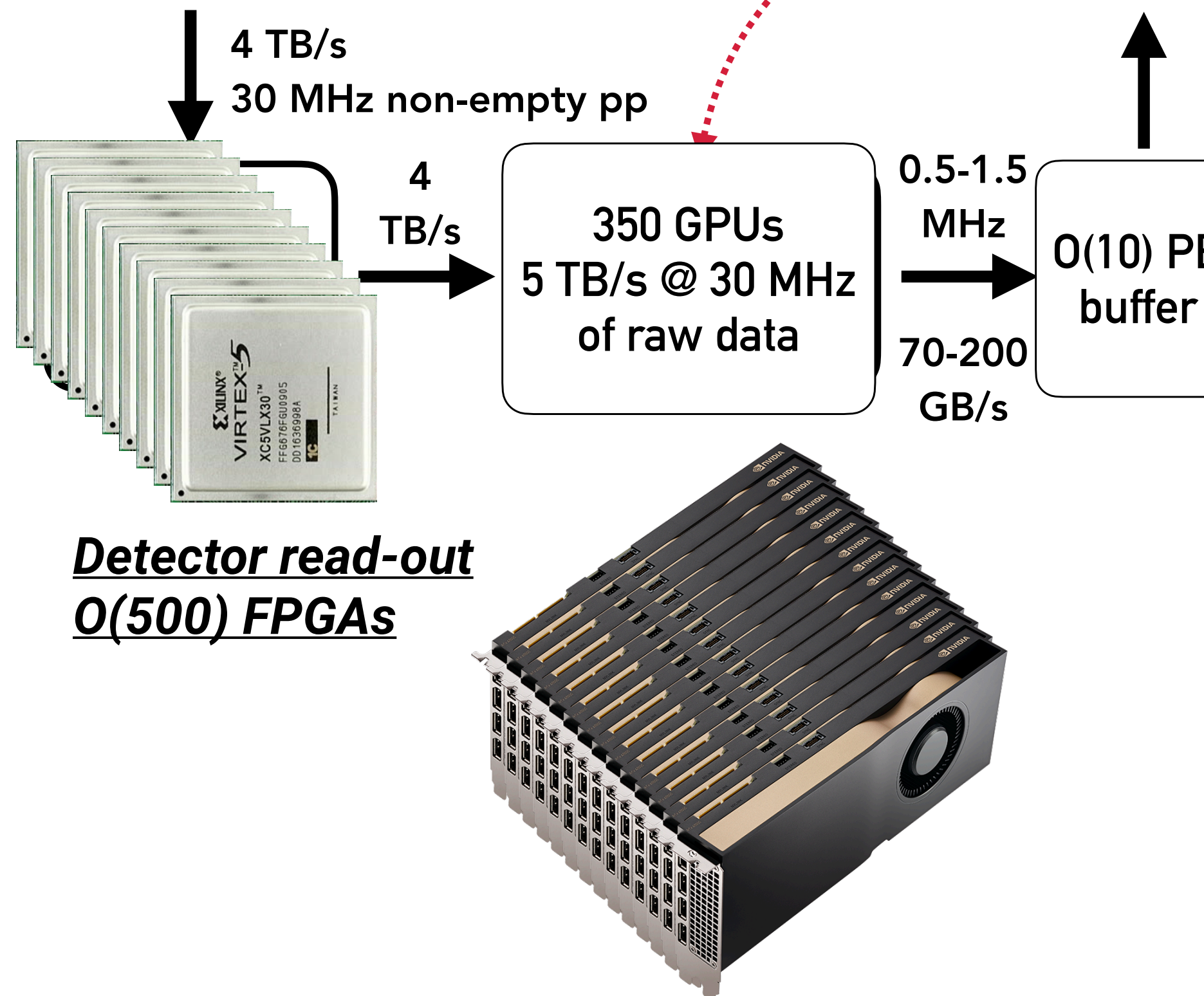*track reconstruction,*
*trigger decisions ... more underway!*
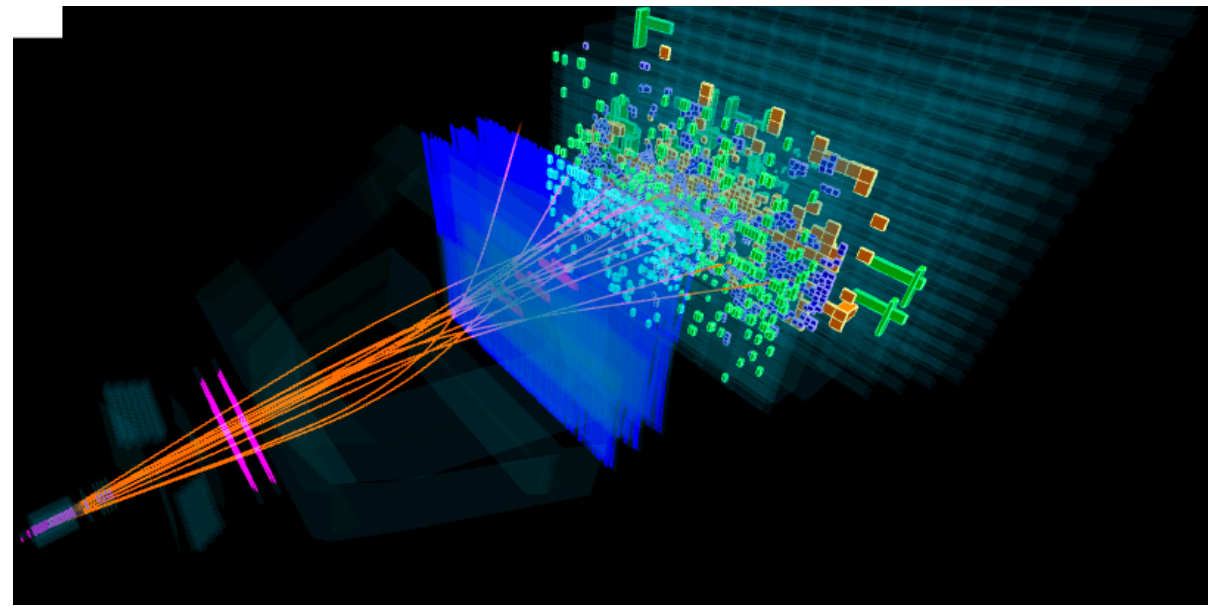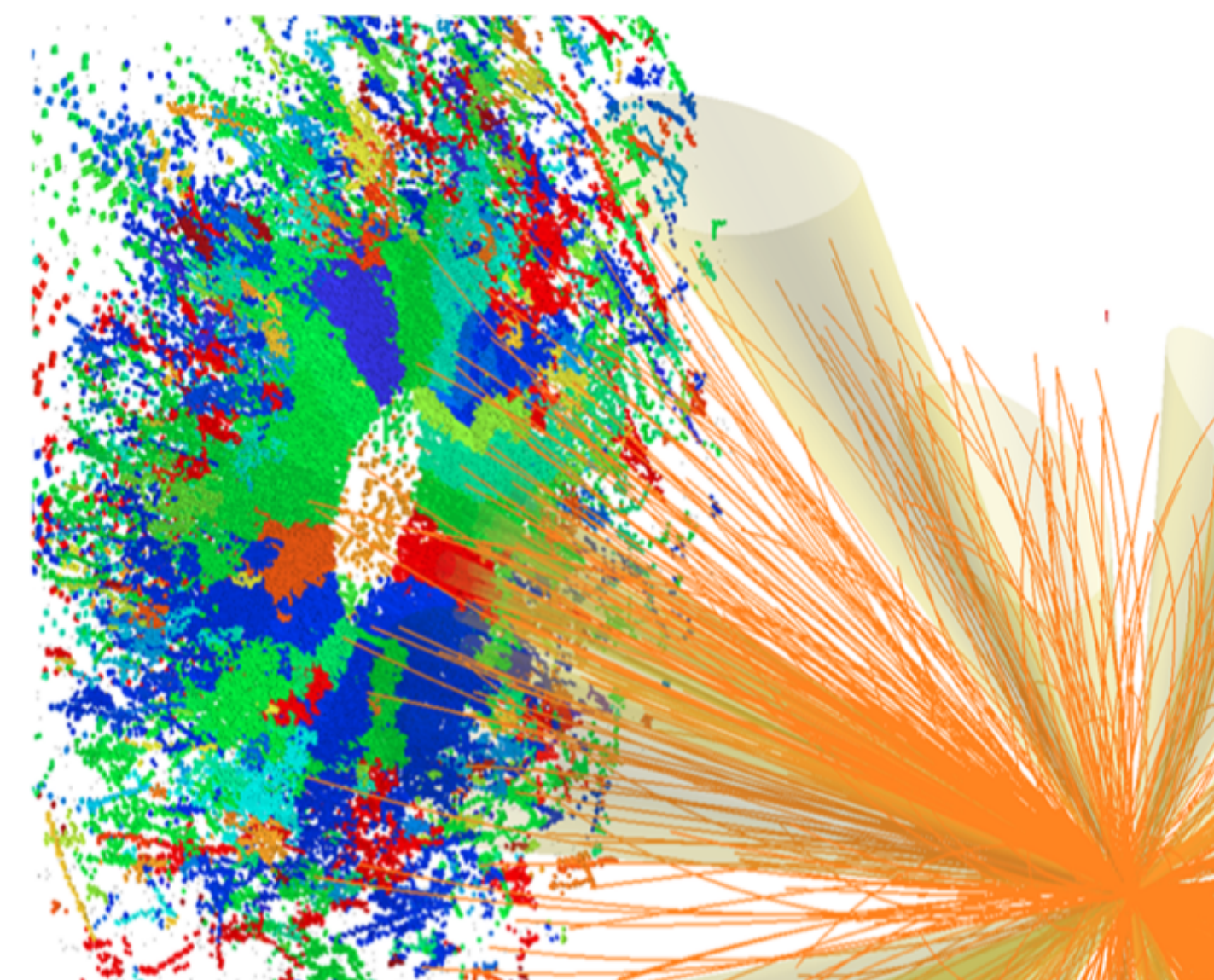
# High throughput GPU triggers

**40 MHz pp collisions**

170 KHz

LHCb 2021
Allen v1r7

Allen throughput (kHz)



*M. Fontana*

Bandwidth (MB/s)

Year

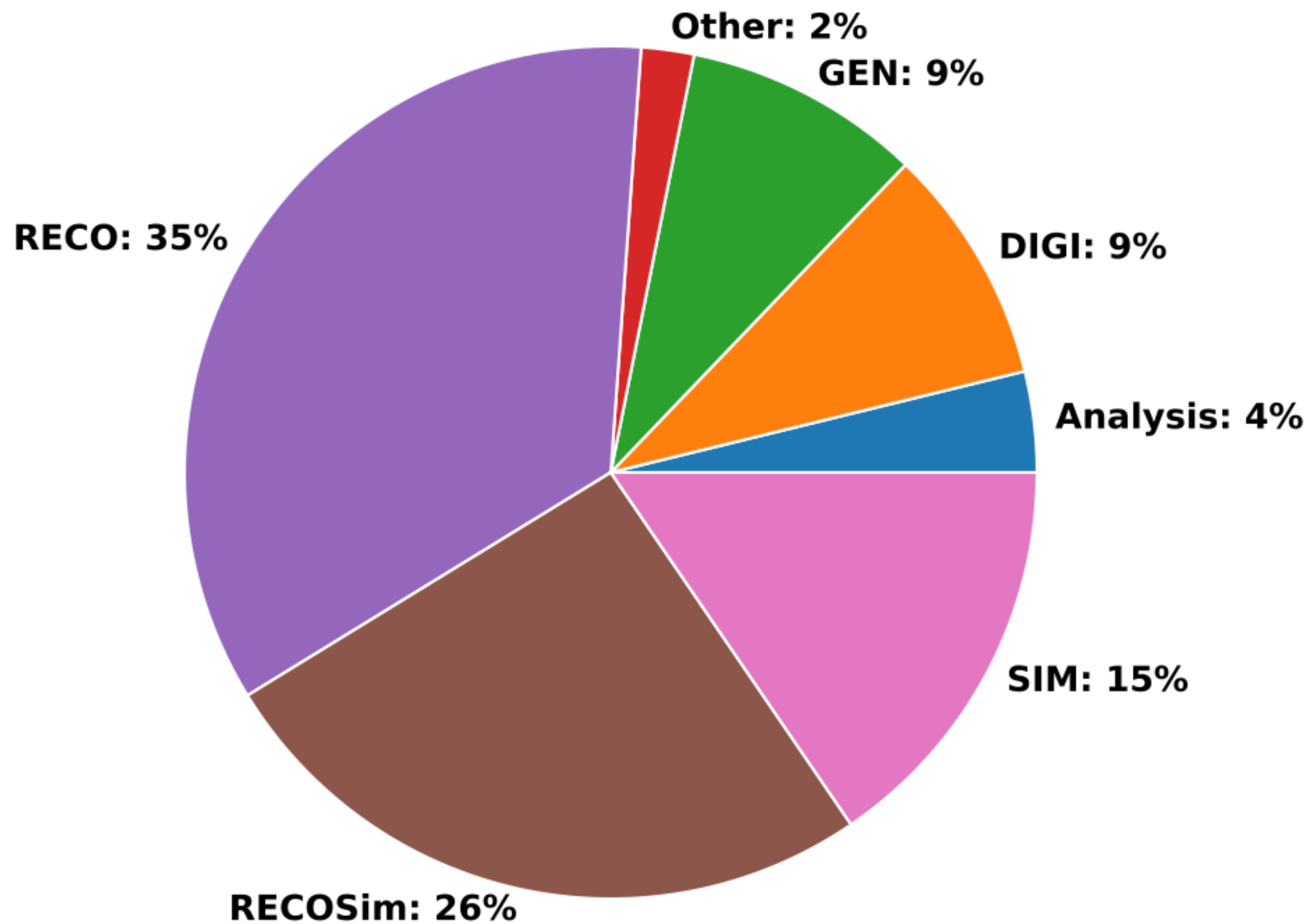Untagged

LHCb
6 fb$^{-1}$

$t$ [ps]

HL-LHC, Simulation of CMS HGCAL with 140 PU

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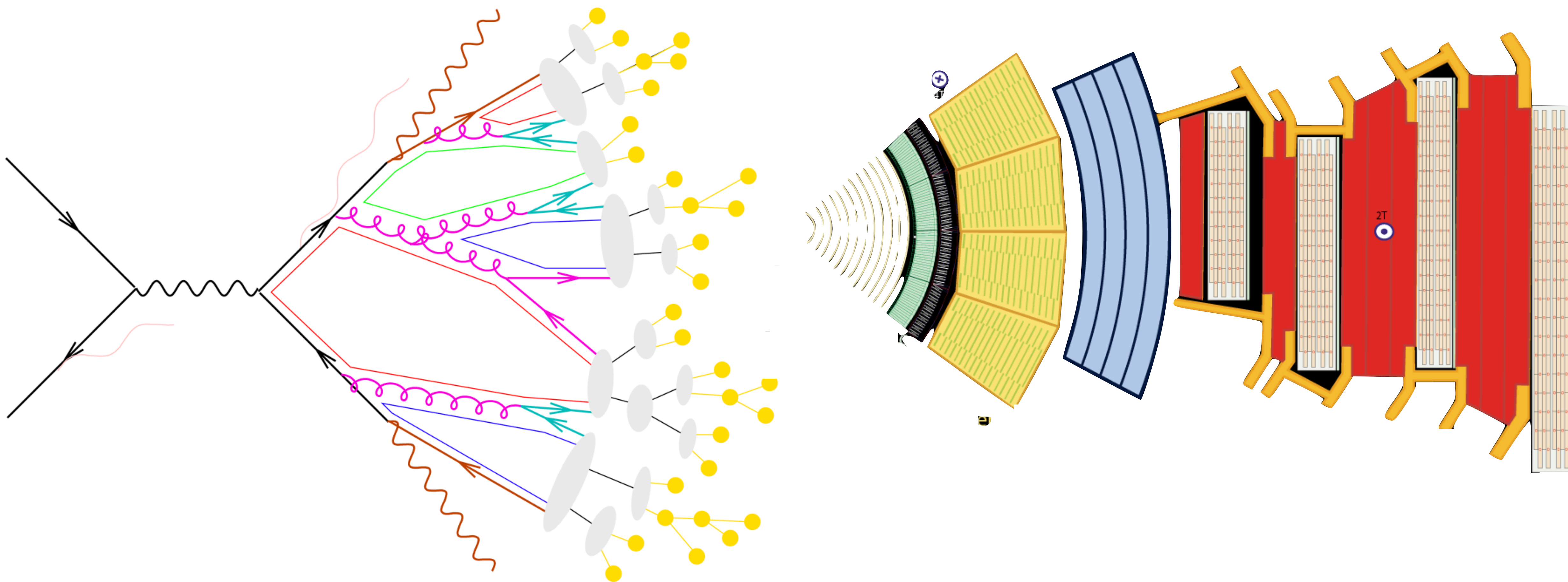*2022 Estimates*

$O(10)$           $O(10^3)$           $O(10^{10})$

$10^{-18}$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  $100$m

**GEN**  **SIM**  **DIGI+RECO**
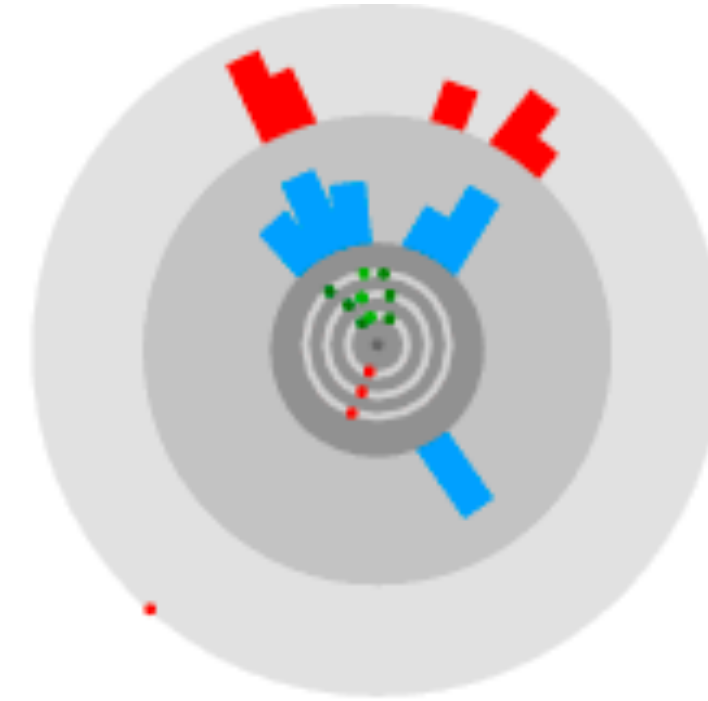
pp collisions up to production of stable particles [Easy & Fast]

detector response simulation [Hard & Slow]

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

CPU

1.1%  0.1%
16.8%  24.4%

GEN
SIM
DIGI
RECO
MINIAOD

57.6%

81%

$O(10)$

$O(10^3)$

$O(10^{10})$

$10^-$

$^{-6}$m

100m

GEN

SIM

DIGI+RECO+MINIAOD

collision point

proton beams

$\phi$

$\eta$

pp co
produ
particle

particle

DIG
REC
MIN

1.1%          CPU          0.1%

16.8%                        24.4%

GEN
SIM
DIGI
RECO
MINIAOD

57.6%

9%          Disk          10%

Energy deposits→digital signals→reconstructed by the reconstruction software

[Hard & Slow]

SIM
MINIAOD

81%

57.6%

57.6%

81%

$O(10)$        $O(10^3)$        $O(10^{10})$

$10^{-18}$m        $10^{-15}$m    $10^{-6}$m        100m

**GEN**      **SIM**      **DIGI+RECO+MINIAOD**

GEANT4

GEANT4

pp
produ
particle

pp co
produ
1
particle

ML me

**CPU**
1.1%    0.1%
16.8%    24.4%

GEN
SIM
DIGI
RECO
MINIAOD

57.6%

57.6%

57.6%

**Disk**
9%    10%

SIM
MINIAOD

81%

81%

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

# Diffusion models



Model

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y})$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

Gaussian noise

***Dall-e 2***

***Dall-e 2***

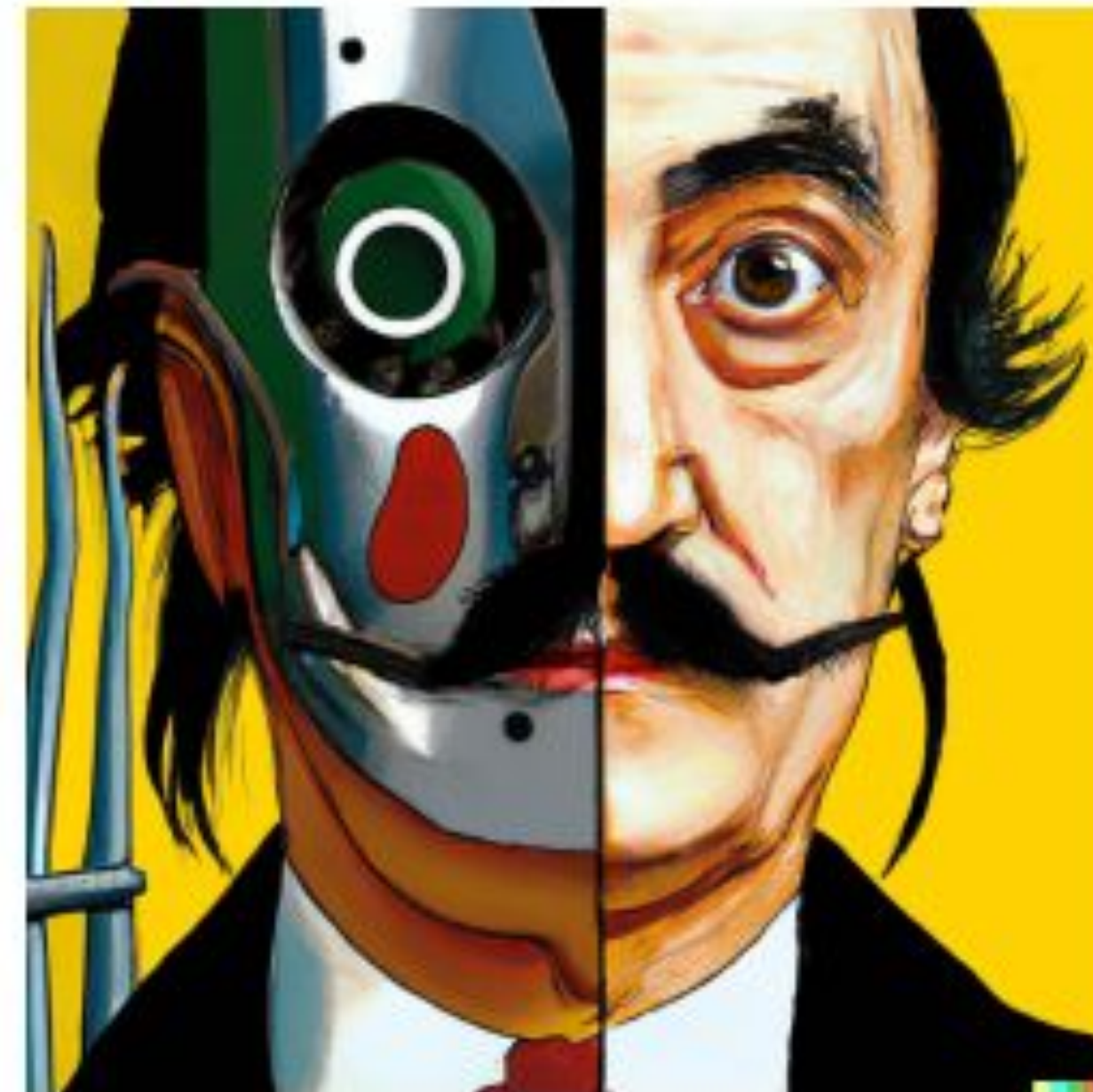an espresso machine that makes coffee from human souls, artstation

text encoder

prior          decoder

vibrant portrait painting of Salvador Dalí with a robotic half face