# Foundation models at the edge for particle physics
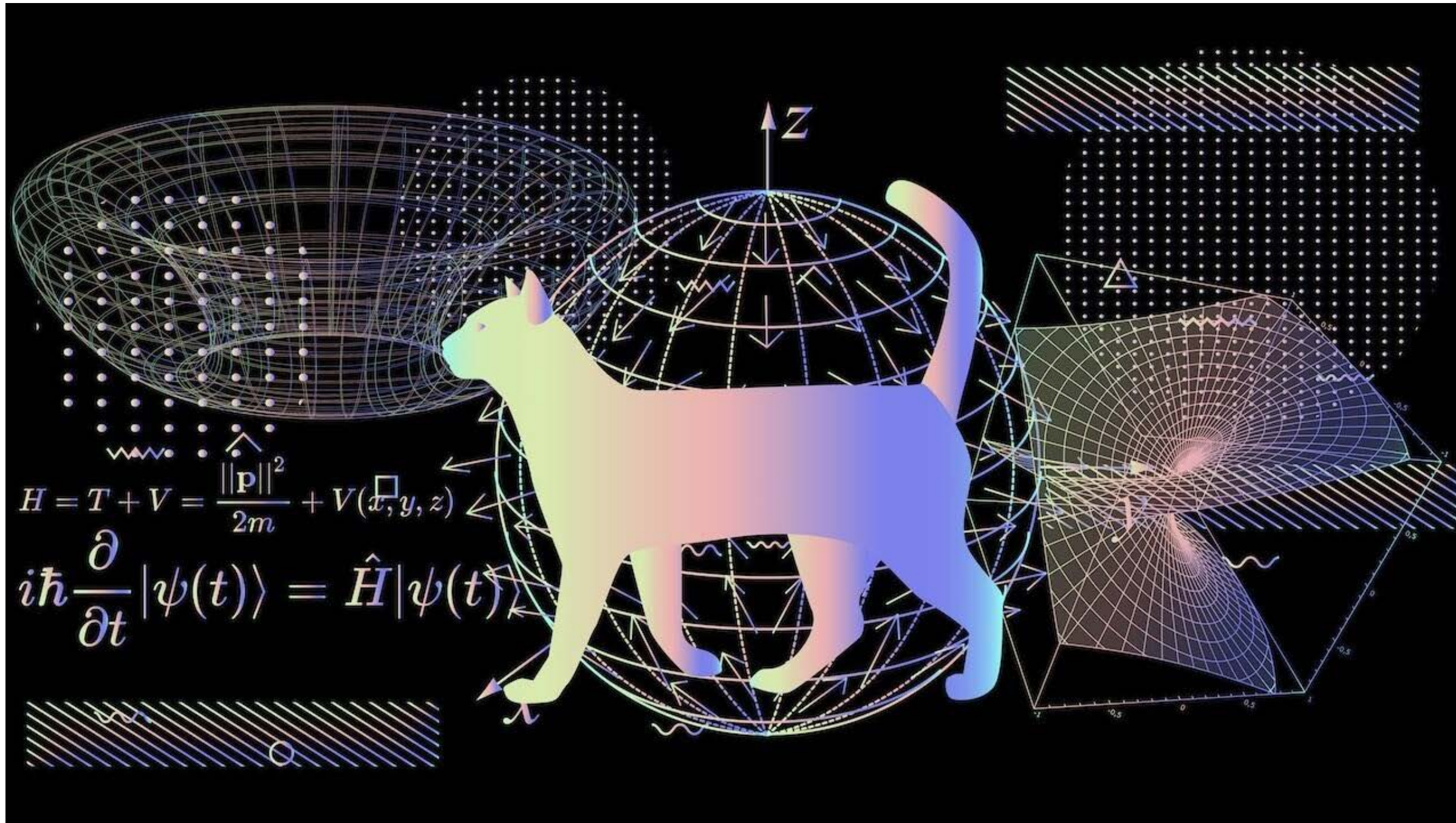
**IAIFI Symposium on Generative AI in the Physical Sciences**

Klaeboe Aarrestad (ETH Zurich)

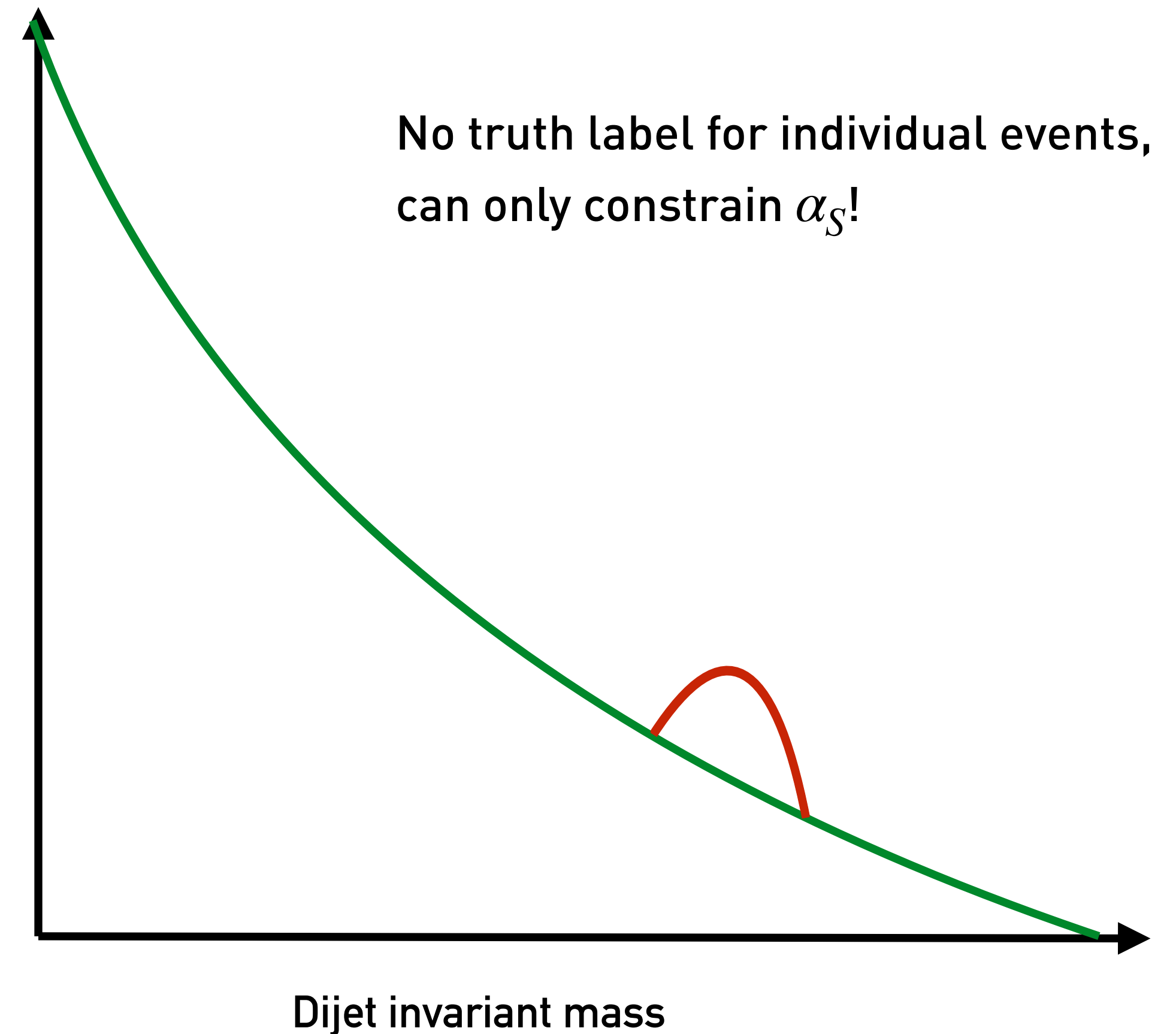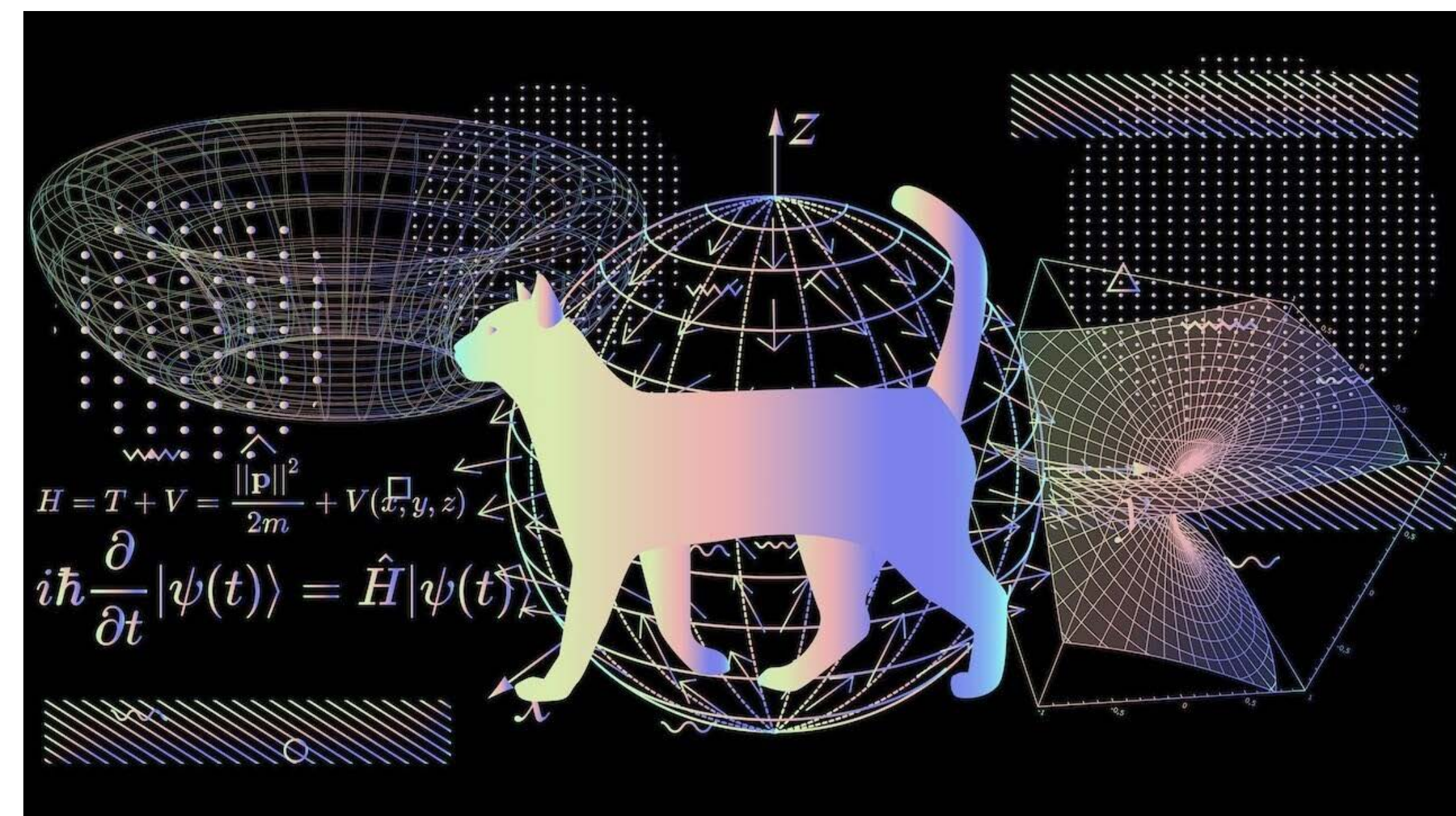# It's against physical law to annotate our data!

$$dP^n_{data} = |M_S + M_B|^2 dp_1 dp_2 \ldots dp_n$$

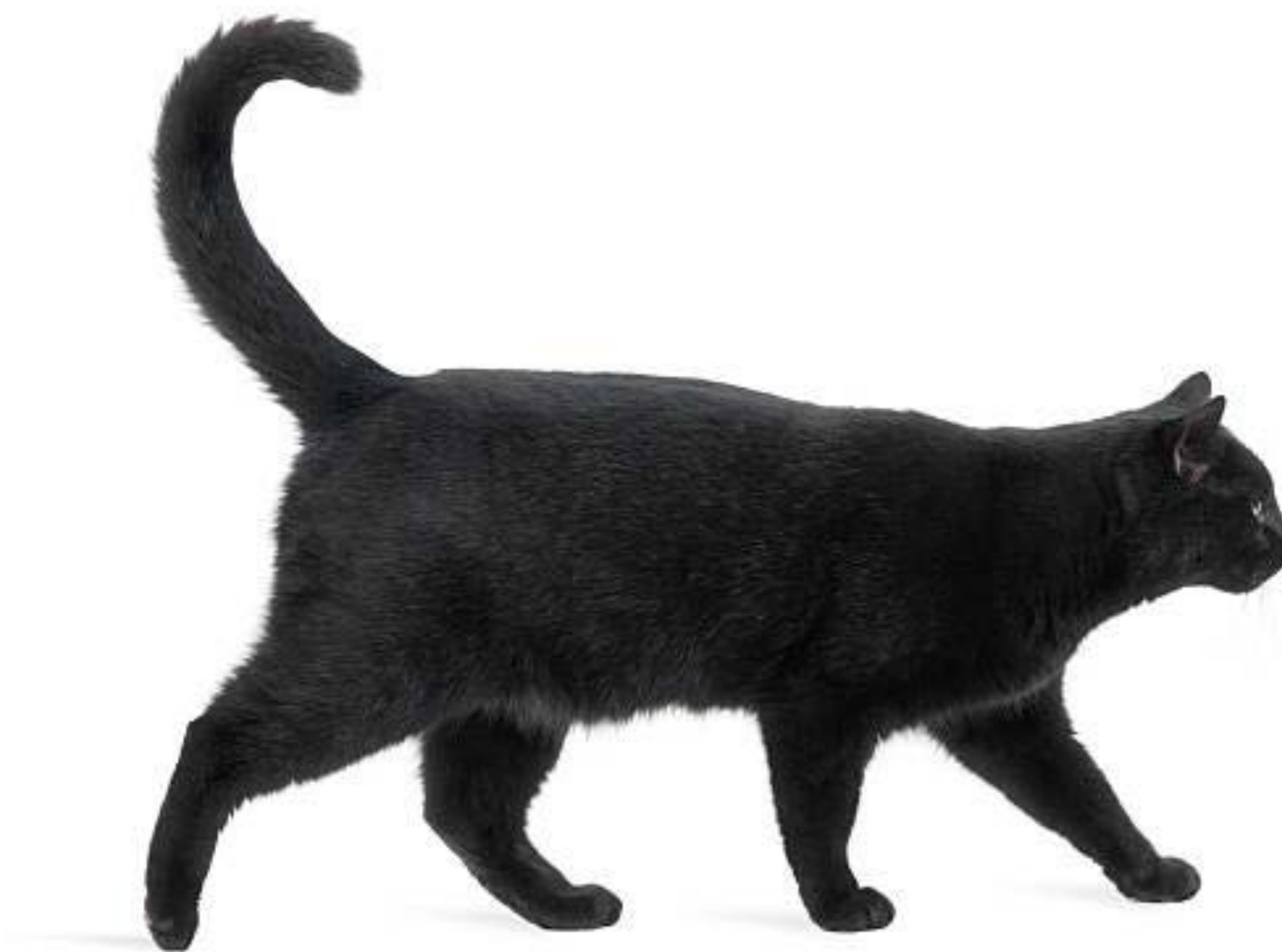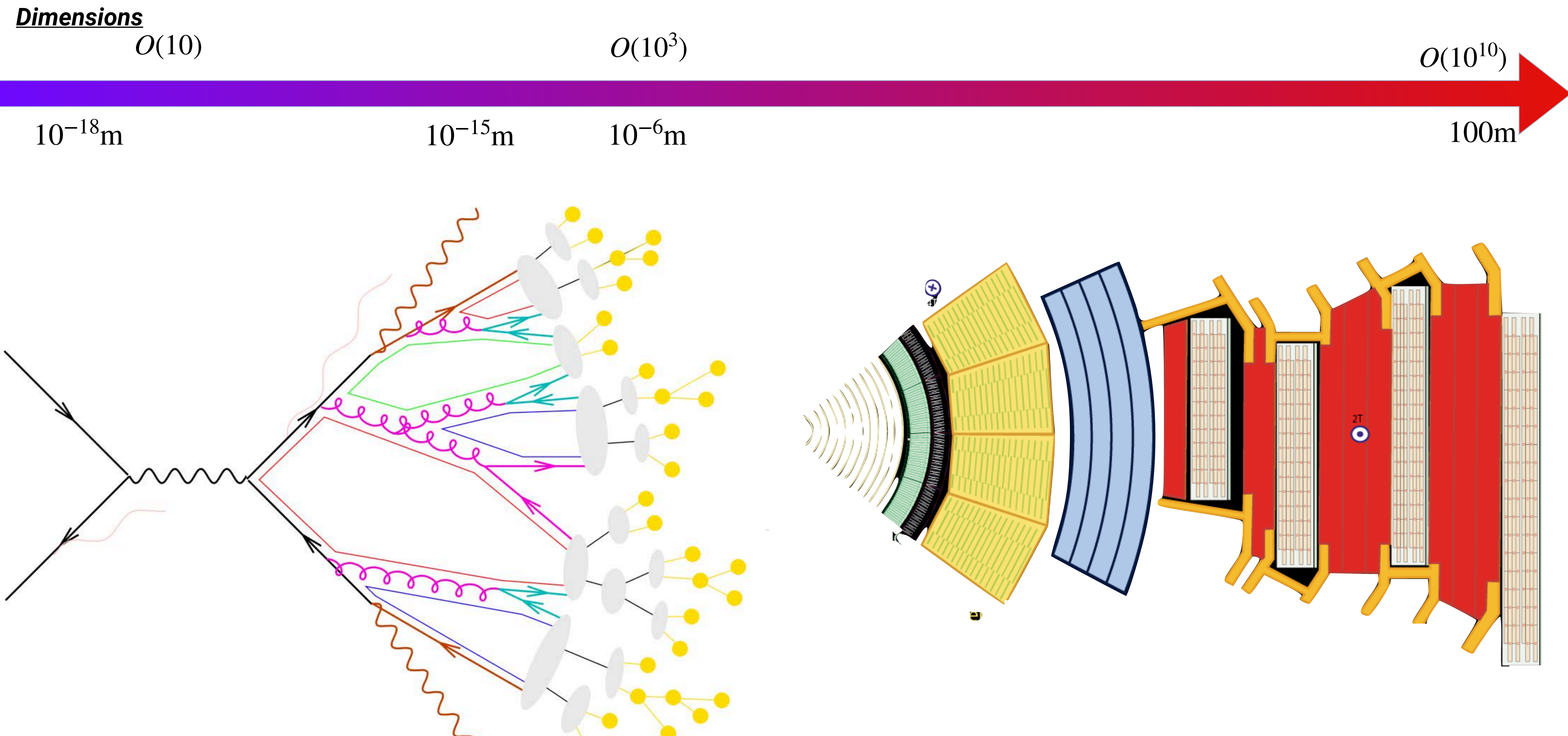$$P_{data} = \alpha_S P_S + \alpha_B P_B$$

No truth label for individual events, can only constrain $\alpha_S$!

$$M_S M_B * + M_B M_S *$$

Dijet invariant mass

$$!=$$

# Monte Carlo Simulation

**Dimensions**

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

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

# ~40 quadrillion collisions recorded at LHC



CMS

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

*LumiPublicResults*

# O(1) trillion simulated events



GEN

SIM

DIGI

hadronic calorimeter

electromagnetic calorimeter

tracker

1.1%

16.8%

GEN
SIM
DIGI
RECO
MINIAOD

57.6%

Disk

81%

**GEN**

**SIM**

**DIGI+RECO+MINIAOD**

hadronic
calorimeter

electromagnetic
calorimeter

tracker

**Disk** 10%

GEN
SIM
MINIAOD

81%

**CMS**

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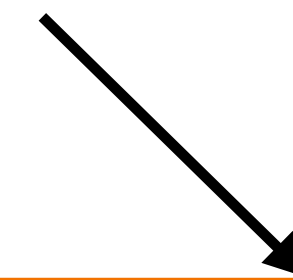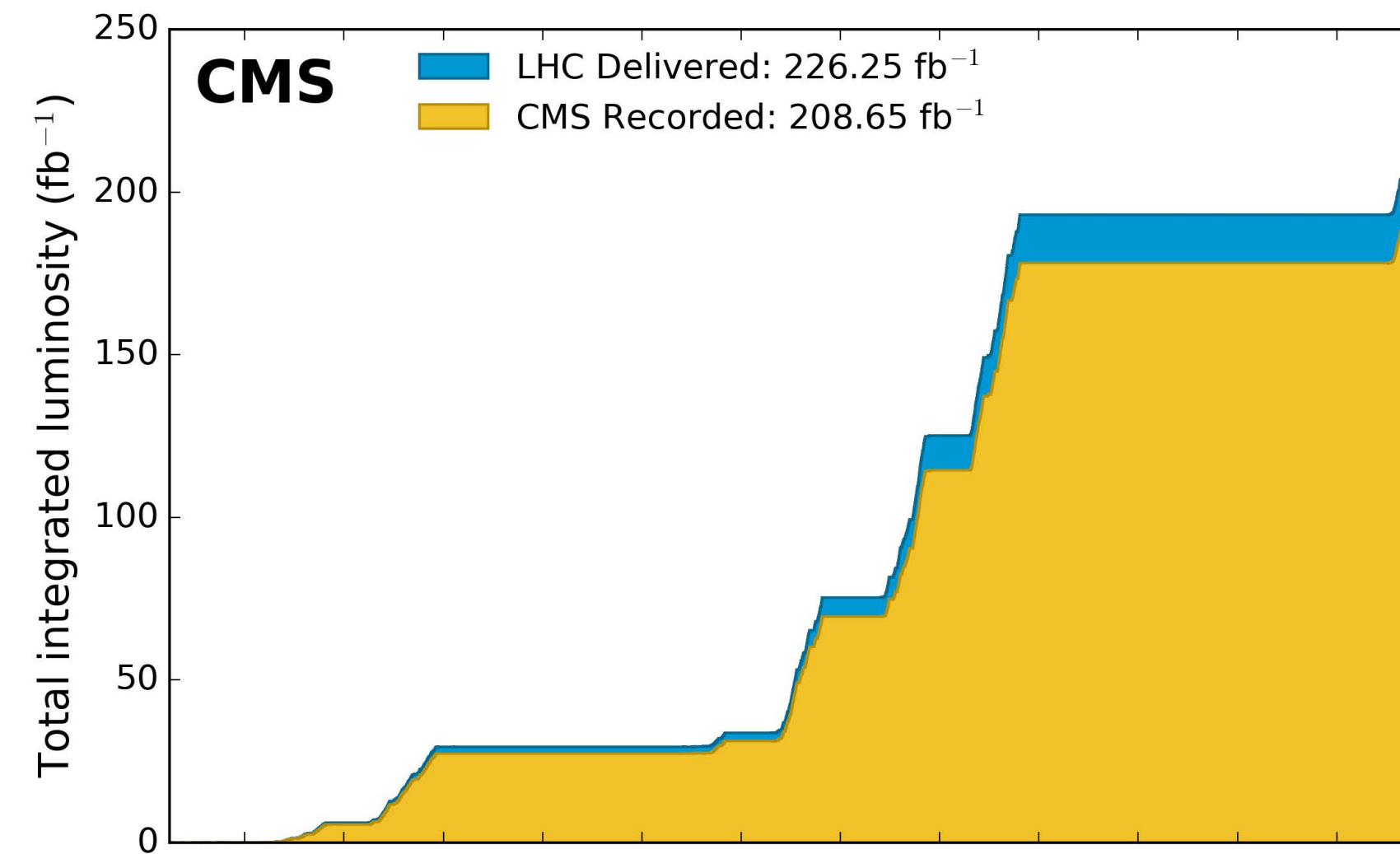CMS Recorded: 208.65 fb$^{-1}$

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

**Fully supervised**
- Requires truth labels
- Only possible using simulation

**Unsupervised/SSL**
No labels, completely data driven
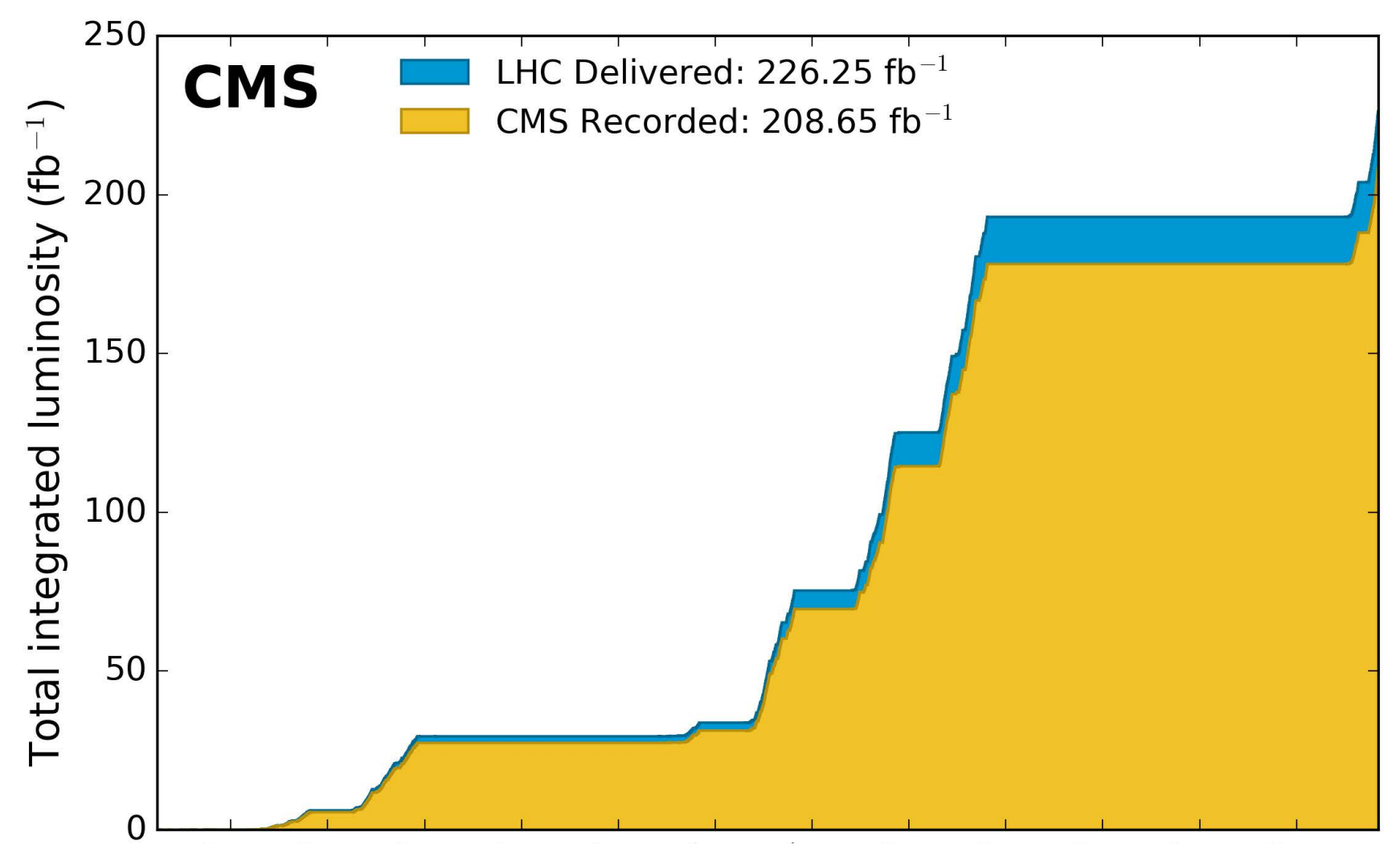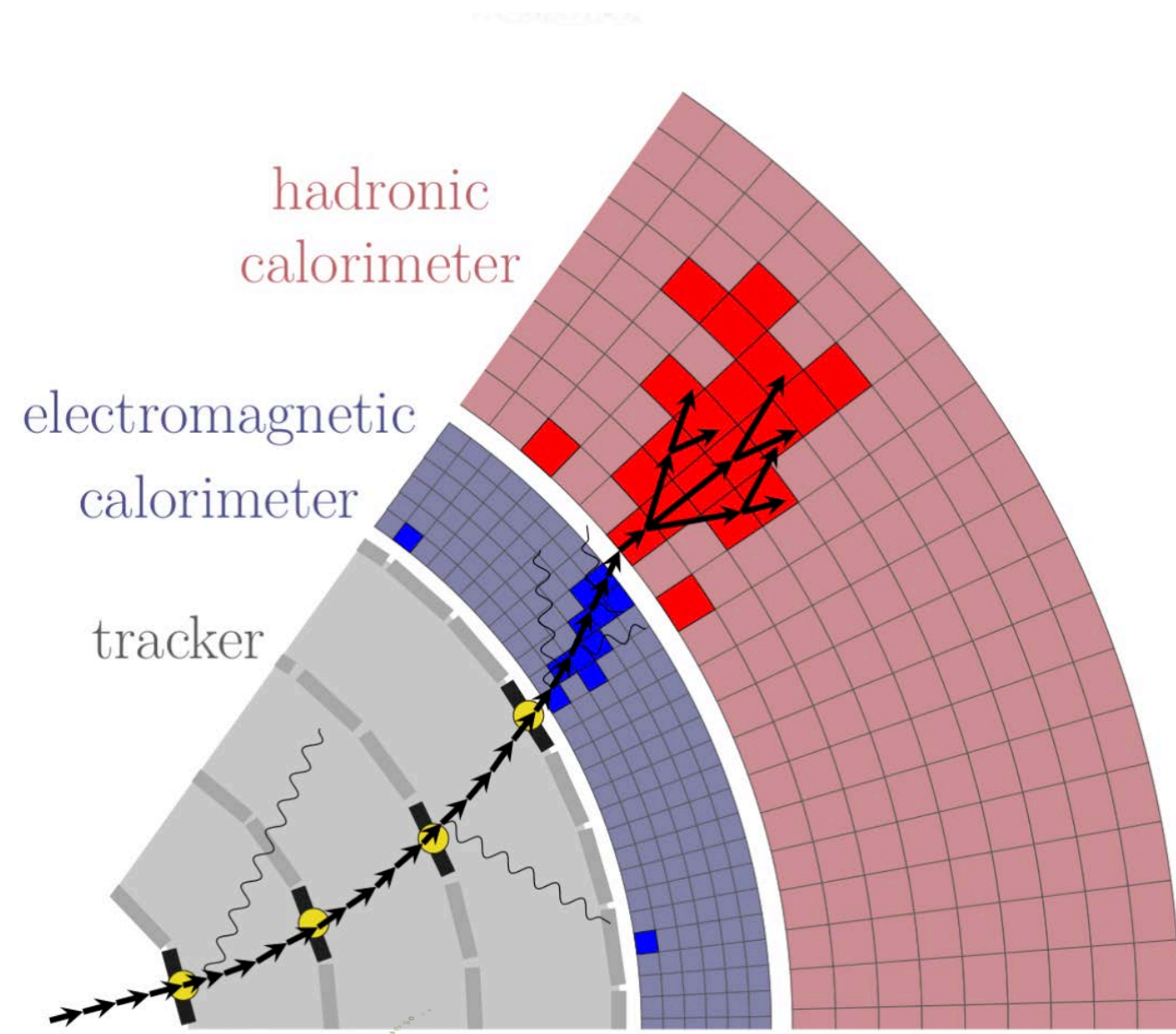
We have a lot of high quality
simulated data that we want to use

We are also very keen
on using this!

GEN · SIM · DIGI+RECO+MINIAOD

hadronic calorimeter

electromagnetic calorimeter

tracker

Disk — 10%

GEN
SIM
MINIAOD

81%

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

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

250
200
150
100
50
0

Simulation != test data

**Fully supervised**
- Requires truth labels
- Only possible using simulation

Mostly (SM )background samples, small signal datasets

**Unsupervised/SSL**
No labels, completely data driven

We have a lot of high quality simulated data that we want to use
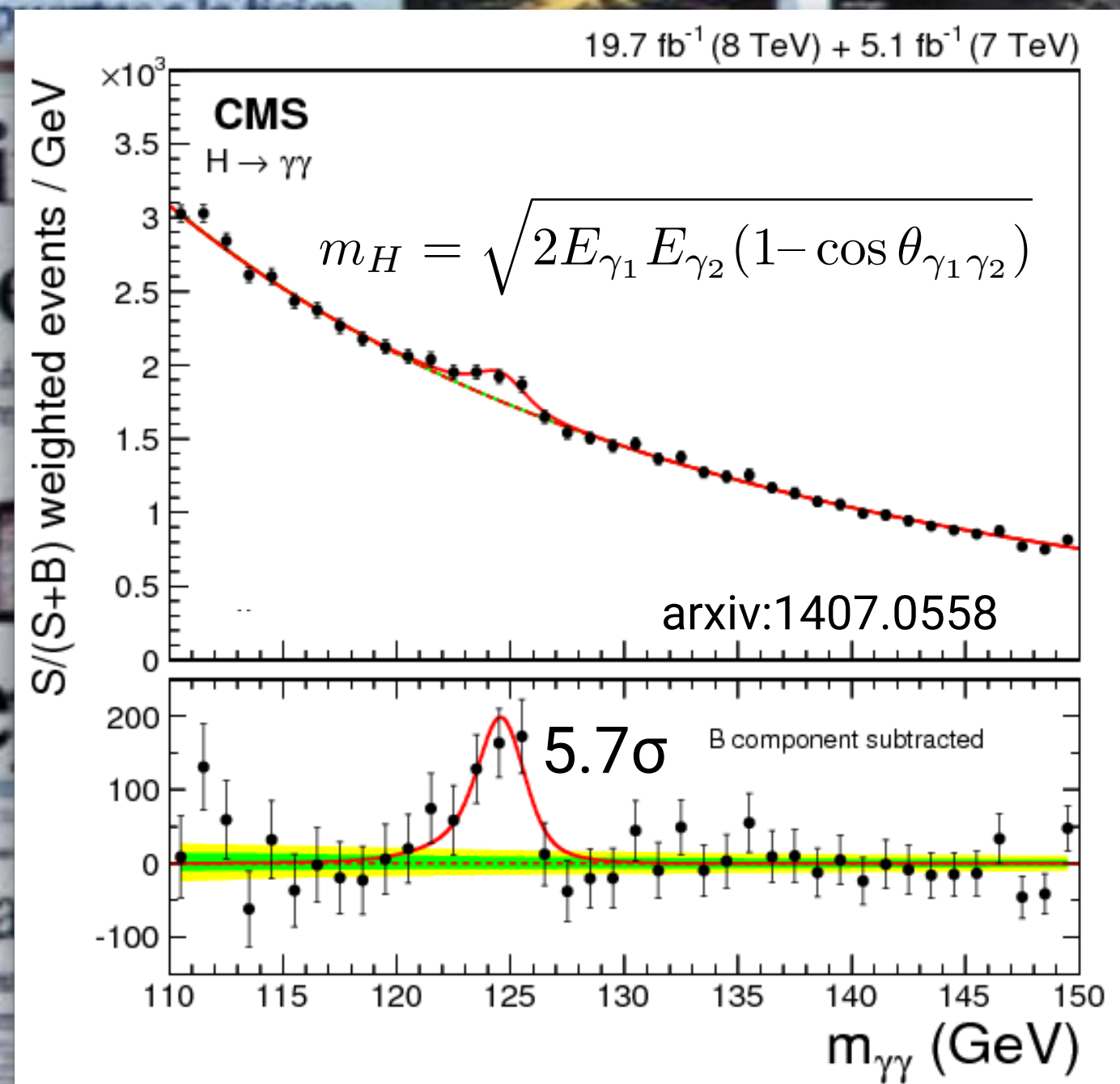
We are also very keen on using this!

Inspire:
("machine learning" or "deep learning" or neural) and (hep-ex or hep-ph or hep-th)

Selected Papers: 457
Total Papers: 457
Year: 2023

Date of paper

Selected Papers: 100
Total Papers: 100
Year: 2024

1992

2024

$$m_H = \sqrt{2E_{\gamma_1}E_{\gamma_2}(1-\cos\theta_{\gamma_1\gamma_2})}$$

CMS
H → γγ

19.7 fb⁻¹ (8 TeV) + 5.1 fb⁻¹ (7 TeV)

arxiv:1407.0558

5.7σ

B component subtracted

CERN Summer student 2012

**Nature Review**

| Analysis | Years of data collection | Sensitivity without machine learning | Sensitivity with machine learning | Ratio of $P$ values | Additional data required |
|---|---|---|---|---|---|
| CMS[24] $H \to \gamma\gamma$ | 2011–2012 | $2.2\sigma$, $P = 0.014$ | $2.7\sigma$, $P = 0.0035$ | 4.0 | 51% |
| ATLAS[43] $H \to \tau^+\tau^-$ | 2011–2012 | $2.5\sigma$, $P = 0.0062$ | $3.4\sigma$, $P = 0.00034$ | 18 | 85% |
| ATLAS[99] $VH \to bb$ | 2011–2012 | $1.9\sigma$, $P = 0.029$ | $2.5\sigma$, $P = 0.0062$ | 4.7 | 73% |
| ATLAS[41] $VH \to bb$ | 2015–2016 | $2.8\sigma$, $P = 0.0026$ | $3.0\sigma$, $P = 0.00135$ | 1.9 | 15% |
| CMS[100] $VH \to bb$ | 2011–2012 | $1.4\sigma$, $P = 0.081$ | $2.1\sigma$, $P = 0.018$ | 4.5 | 125% |

**We were using ML for discovery very early on**

CERN Summer student 2012

11–15 Mar 2024
Charles B. Wang Center, Stony Brook University
US/Eastern timezone

**Now happening:** Exascale Infrastructures for Science (Theatre) *08:45 - 09:15*

Enter your search term

Overview

Scientific Programme

Info for presenters

Timetable

Contribution List

Registration

Accommodations

Travel Information

- About Stony Brook and Long Island
- Important dates
- Getting Around and Parking, Internet access, Venue and Registration
- Food and Drinks
- Things to Do near SBU
- What to do in New York City

ACAT Organization

## 22nd International Workshop on Advanced Computing and Analysis Techniques in Physics Research

The 22nd International Workshop on Advanced Computing and Analysis Techniques in Physics Research (ACAT 2024 will take place between Monday 11th and Friday, 15th March, 2024 at the Stony Brook University, Stony Brook, Long Island NY, USA.

The 22nd edition of ACAT will — once again — bring together computational experts from a wide range of disciplines, including particle-, nuclear-, astro-, and accelerator-physics as well as high performance computing. Through this unique forum, we will explore the areas where these disciplines overlap with computer science, fostering the exchange of ideas related to cutting-edge computing, data-analysis, and theoretical-calculation technologies.

Our Theme will be **Foundation Models for Physics - Nexus of Computation and Physics through Embracing the Era of Foundation Models:** The 2024 ACAT workshop invites the vanguard of computational and physics experts to delve into the transformative potential of foundation models. As the intersection between physics and computational realms deepens, these advanced models, underpinned by colossal datasets and capable of generating nuanced outputs, are redefining the research spectrum and increasingly reshaping the way researchers approach complex problems, simulations, and data analyses. As we chart this new territory, we'll address challenges and opportunities encompassing integration into computational ecosystems, innovative data practices, training nuances, infrastructure evolution, uncertainty metrics, ethical dimensions, and collaborative vistas across disciplines.

Selected Papers: 457
Total Papers: 457
Year: 2023

Selected Papers: 100
Total Papers: 100
Year: 2024

Re-Simulation-based Self-Supervised Learning for Pre-Training Foundation Models  #3
Philip Harris, Michael Kagan, Jeffrey Krupa, Benedikt Maier, Nathaniel Woodward (Mar 11, 2024)
e-Print: 2403.07066 [hep-ph]

pdf    cite    claim    reference search    0 citations

OmniJet-$\alpha$: The first cross-task foundation model for particle physics  #5
Joschka Birk, Anna Hallin, Gregor Kasieczka (Mar 8, 2024)
e-Print: 2403.05618 [hep-ph]

pdf    cite    claim    reference search    0 citations
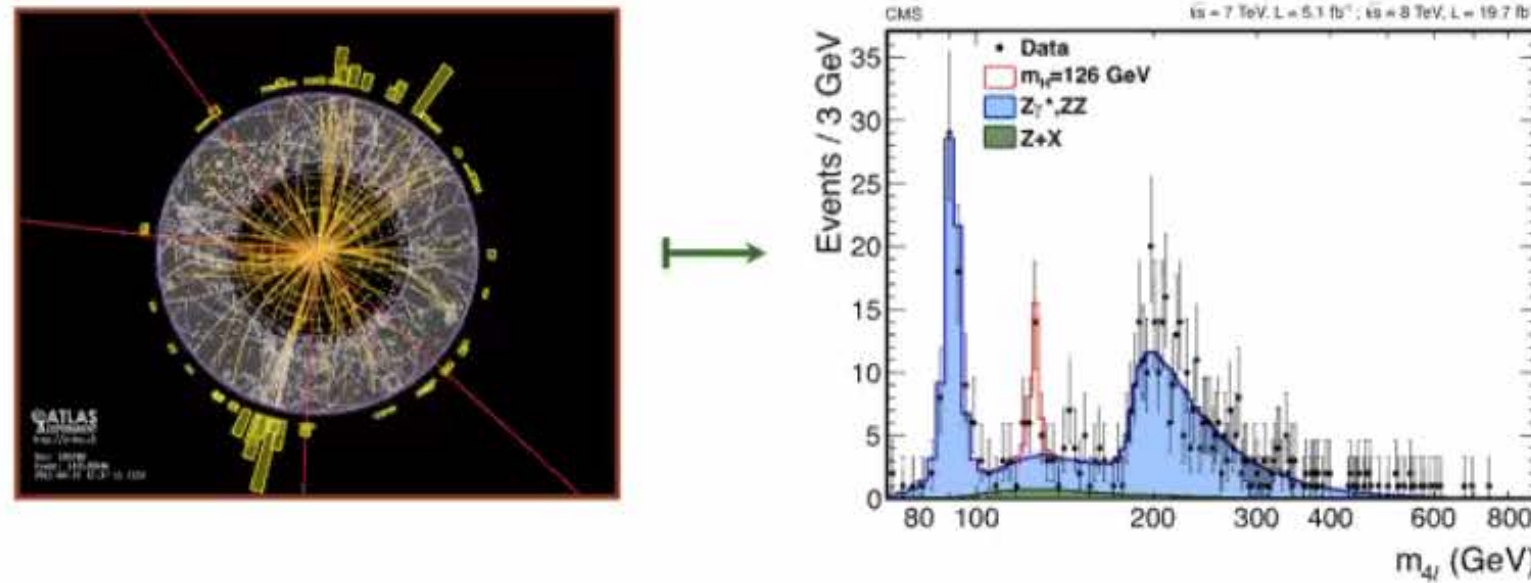
# AI + Physics: *A new frontier?*

Framing: Kyle Cranmer

Many fields within AI4Science are pushing the frontiers of AI… what about physics?

**Reliable inference with complex forward models**

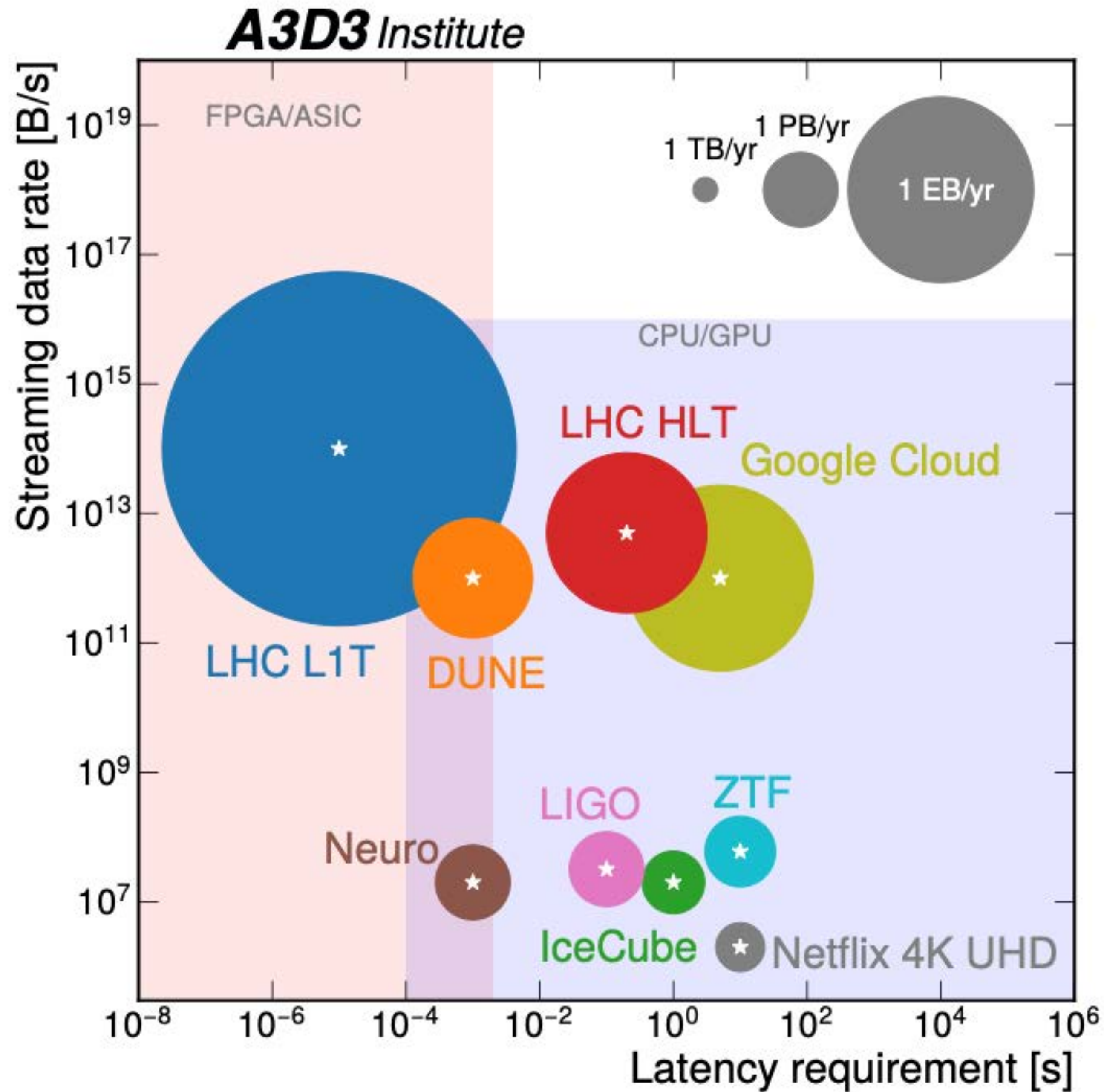**Extremely fast real-time inference**

**High Energy Physics**

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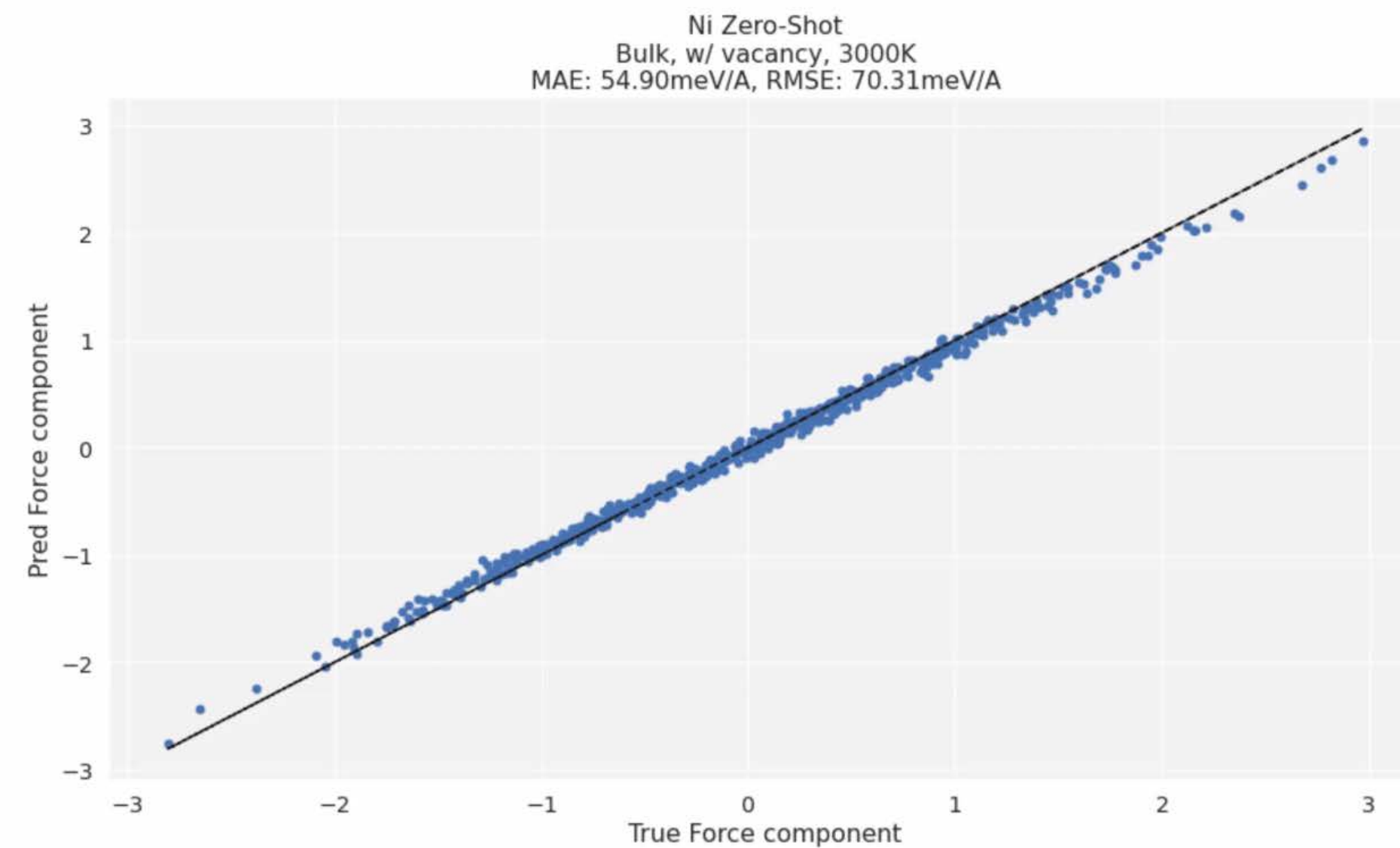
Read More >

(From A3D3 website)

- Sampling under complex symmetries and exactness guarantees (e.g., in lattice QFT)
- Statistical anomaly detection
- Highly structured models/data-generating processes
- …

FastML:

Pioneering
AI in the
physical
sciences

From Simon. 60 million parameter model



Can we combine 12μs latency and O(100M) parameter models?

chat.openai.com

**ChatGPT**

Explore GPTs

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

Today

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

2023

Upgrade plan
Collaborate on a Team plan

Thea Aarrestad

## How can I help you today?

| Design a database schema | Create a personal webpage for me |
| for an online merch store | after asking me three questions |

| Recommend a dish | Tell me a fun fact |
| to bring to a potluck | about the Roman Empire |

Message ChatGPT...

ChatGPT can make mistakes. Consider checking important information.

AlphaFold nature cover

T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

sequence—the structure prediction component of the 'protein folding problem'[8]—has been an important open research problem for more than 50 years[9]. Despite recent

T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

The international journal of science / 30 July 2021

nature

**ALPHAFOLD MANIA**
The number of research papers and preprints citing the AlphaFold2 AI software has shot up since its source code was released in July 2021*.

■ Journal article   ■ Preprint

AlphaFold2 announced as winner of protein-folding software contest.

Paper describing AlphaFold2 released, with source code.

Number of research articles

125
100
75
50
25
0

Dec Jan Feb Mar Apr May Jun Jul Aug Sep Oct Nov Dec Jan Feb Mar
2021▶                                                    2022▶

*Nature analysis using Dimensions database; removing duplicate preprints and papers/R. Van Noorden, E. Callaway.

©nature

sequence—the structure prediction component of the 'protein folding problem'[8]—has been an important open research problem for more than 50 years[9]. Despite recent

**Train (GPT-4):**
- **$2.15^{25}$ floating point operations**
- **25,000 A100 GPUs**
- **90-100 days**
- **$63 million**
- **Trained on 13 trillion tokens**

**Train (GPT-4):**
- **$2.15^{25}$ floating point operations**
- **25,000 A100 GPUs**
- **90-100 days**
- **$63 million**
- **Trained on 13 trillion tokens**

**Inference (GPT-4):**
- **Multiple clusters of 128 GPUs**
- **Model carefully mapped onto hardware**

**Transformer Size: 240x / 2 yrs**

**AI HW Memory:    2x / 2 yrs**

10TB Baidu RecSys

2TB Baidu RecSys

GShard

GPT-3

Microsoft T-NLG

A100-80 (80GB)

Megatron LM

V100 (32GB)          TPUv3 (32GB)

A100 (40GB)

TPUv2 (16GB)

GPT-2

P100 (12GB)

BERT

ALBERT

Inception V4          ResNext101

GPT-1

ELECTRA

Transformer

ResNet50          DenseNet

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

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

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

Kaplan et al. (2020)

**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**

**Resources:** 128 interconnected GPUs
**Latency :** 10 seconds

**Resources:** O(10) single chips
**Latency:** 1 millionth of a second
5% of internet traffic

GPT-4

?

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

Geneva

Geneva Lake

CMS

LHCb

ATLAS

ALICE

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

$2.2 \cdot 10^{11}$ protons

25 ns

7.5 m

$$E = mc^2$$

Quarks

Leptons

Force Carriers

H Higgs boson

Higgs: 125 GeV

Masses span 9 orders of magnitude!

TeV

GeV

MeV

keV

eV

meV

**Quarks**

| | | |
|---|---|---|
| u up | c charm | t top |
| d down | s strange | b bottom |

H Higgs boson

**Force Carriers**

| | |
|---|---|
| Z Z boson | γ photon |
| W W boson | g gluon |

**Leptons**

| | | |
|---|---|---|
| e electron | μ muon | τ tau |
| $\nu_e$ electron neutrino | $\nu_\mu$ muon neutrino | $\nu_\tau$ tau neutrino |

4T

2T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with Muon chambers

Quarks

| u up | c charm | t top |
| d down | s strange | b bottom |

Force Carriers

| Z z boson | γ photon |
| W w boson | g gluon |

H Higgs boson

Leptons

| e electron | μ muon | τ tau |
| νe electron neutrino | νμ muon neutrino | ντ tau neutrino |

4T

2T

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Iron return yoke interspersed with Muon chambers

**CMS** *Supplementary*

$\sqrt{s} = 13$ TeV: L = 4.3 fb$^{-1}$ (2016)

- H(125)
- gg/q$\bar{q} \to$ ZZ,Z$\gamma^*$
- EW
- Z+X
- Data

Events/4 GeV

$m_{4l}$ [GeV]

$118 < m_{4l} < 130$ GeV

Events/1 GeV

$m_{4l}$ [GeV]

*We had to collide billions of protons,*
*only around 10 signal events were needed to claim discovery!*

We had to collide billions of protons,
only around 10 signal events were needed to claim discovery!

# The Standard Model

$$-\tfrac{1}{2}\partial_\nu g^a_\mu \partial_\nu g^a_\mu - g_s f^{abc}\partial_\mu g^a_\nu g^b_\mu g^c_\nu - \tfrac{1}{4}g_s^2 f^{abc}f^{ade}g^b_\mu g^c_\nu g^d_\mu g^e_\nu +$$
$$\tfrac{1}{2}ig_s^2(\bar{q}^\sigma_i \gamma^\mu q^\sigma_j)g^a_\mu + \bar{G}^a \partial^2 G^a + g_s f^{abc}\partial_\mu \bar{G}^a G^b g^c_\mu - \partial_\nu W^+_\mu \partial_\nu W^-_\mu -$$
$$M^2 W^+_\mu W^-_\mu - \tfrac{1}{2}\partial_\nu Z^0_\mu \partial_\nu Z^0_\mu - \tfrac{1}{2c_w^2}M^2 Z^0_\mu Z^0_\mu - \tfrac{1}{2}\partial_\mu A_\nu \partial_\mu A_\nu - \tfrac{1}{2}\partial_\mu H \partial_\mu H -$$
$$\tfrac{1}{2}m_h^2 H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - M^2 \phi^+ \phi^- - \tfrac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \tfrac{1}{2c_w^2}M\phi^0\phi^0 - \beta_h[\tfrac{2M^2}{g^2}+$$
$$\tfrac{2M}{g}H + \tfrac{1}{2}(H^2 + \phi^0\phi^0 + 2\phi^+\phi^-)] + \tfrac{2M^4}{g^2}\alpha_h - igc_w[\partial_\nu Z^0_\mu(W^+_\mu W^-_\nu -$$
$$W^+_\nu W^-_\mu) - Z^0_\nu(W^+_\mu \partial_\nu W^-_\mu - W^-_\mu \partial_\nu W^+_\mu) + Z^0_\mu(W^+_\nu \partial_\nu W^-_\mu -$$
$$W^-_\nu \partial_\nu W^+_\mu)] - igs_w[\partial_\nu A_\mu(W^+_\mu W^-_\nu - W^+_\nu W^-_\mu) - A_\nu(W^+_\mu \partial_\nu W^-_\mu -$$
$$W^-_\mu \partial_\nu W^+_\mu) + A_\mu(W^+_\nu \partial_\nu W^-_\mu - W^-_\nu \partial_\nu W^+_\mu)] - \tfrac{1}{2}g^2 W^+_\mu W^-_\mu W^+_\nu W^-_\nu +$$
$$\tfrac{1}{2}g^2 W^+_\mu W^-_\nu W^+_\mu W^-_\nu + g^2 c_w^2(Z^0_\mu W^+_\mu Z^0_\nu W^-_\nu - Z^0_\mu Z^0_\mu W^+_\nu W^-_\nu) +$$
$$g^2 s_w^2(A_\mu W^+_\mu A_\nu W^-_\nu - A_\mu A_\mu W^+_\nu W^-_\nu) + g^2 s_w c_w[A_\mu Z^0_\nu(W^+_\mu W^-_\nu -$$
$$W^+_\nu W^-_\mu) - 2A_\mu Z^0_\mu W^+_\nu W^-_\nu] - g\alpha[H^3 + H\phi^0\phi^0 + 2H\phi^+\phi^-] -$$
$$\tfrac{1}{8}g^2\alpha_h[H^4+(\phi^0)^4+4(\phi^+\phi^-)^2+4(\phi^0)^2\phi^+\phi^-+4H^2\phi^+\phi^-+2(\phi^0)^2 H^2]-$$
$$gMW^+_\mu W^-_\mu H - \tfrac{1}{2}g\tfrac{M}{c_w^2}Z^0_\mu Z^0_\mu H - \tfrac{1}{2}ig[W^+_\mu(\phi^0\partial_\mu\phi^- - \phi^-\partial_\mu\phi^0) -$$
$$W^-_\mu(\phi^0\partial_\mu\phi^+ - \phi^+\partial_\mu\phi^0)]+\tfrac{1}{2}g[W^+_\mu(H\partial_\mu\phi^- - \phi^-\partial_\mu H) - W^-_\mu(H\partial_\mu\phi^+ -$$
$$\phi^+\partial_\mu H)]+\tfrac{1}{2}g\tfrac{1}{c_w}(Z^0_\mu(H\partial_\mu\phi^0 - \phi^0\partial_\mu H)-ig\tfrac{s_w^2}{c_w}MZ^0_\mu(W^+_\mu\phi^- - W^-_\mu\phi^+)+$$
$$igs_w MA_\mu(W^+_\mu\phi^- - W^-_\mu\phi^+) - ig\tfrac{1-2c_w^2}{2c_w}Z^0_\mu(\phi^+\partial_\mu\phi^- - \phi^-\partial_\mu\phi^+) +$$
$$igs_w A_\mu(\phi^+\partial_\mu\phi^- - \phi^-\partial_\mu\phi^+) - \tfrac{1}{4}g^2 W^+_\mu W^-_\mu[H^2 + (\phi^0)^2 + 2\phi^+\phi^-] -$$
$$\tfrac{1}{4}g^2\tfrac{1}{c_w^2}Z^0_\mu Z^0_\mu[H^2 + (\phi^0)^2 + 2(2s_w^2-1)^2\phi^+\phi^-] - \tfrac{1}{2}g^2\tfrac{s_w^2}{c_w}Z^0_\mu\phi^0(W^+_\mu\phi^- +$$
$$W^-_\mu\phi^+) - \tfrac{1}{2}ig^2\tfrac{s_w^2}{c_w}Z^0_\mu H(W^+_\mu\phi^- - W^-_\mu\phi^+) + \tfrac{1}{2}g^2 s_w A_\mu\phi^0(W^+_\mu\phi^- +$$
$$W^-_\mu\phi^+) + \tfrac{1}{2}ig^2 s_w A_\mu H(W^+_\mu\phi^- - W^-_\mu\phi^+) - g^2\tfrac{s_w}{c_w}(2c_w^2-1)Z^0_\mu A_\mu\phi^+\phi^- -$$
$$g^1 s_w^2 A_\mu A_\mu\phi^+\phi^- - \bar{e}^\lambda(\gamma\partial + m_e^\lambda)e^\lambda - \bar{\nu}^\lambda\gamma\partial\nu^\lambda - \bar{u}^\lambda_j(\gamma\partial + m_u^\lambda)u^\lambda_j -$$
$$\bar{d}^\lambda_j(\gamma\partial + m_d^\lambda)d^\lambda_j + igs_w A_\mu[-(\bar{e}^\lambda\gamma^\mu e^\lambda) + \tfrac{2}{3}(\bar{u}^\lambda_j\gamma^\mu u^\lambda_j) - \tfrac{1}{3}(\bar{d}^\lambda_j\gamma^\mu d^\lambda_j)] +$$
$$\tfrac{ig}{4c_w}Z^0_\mu[(\bar{\nu}^\lambda\gamma^\mu(1 + \gamma^5)\nu^\lambda) + (\bar{e}^\lambda\gamma^\mu(4s_w^2 - 1 - \gamma^5)e^\lambda) + (\bar{u}^\lambda_j\gamma^\mu(\tfrac{4}{3}s_w^2 -$$
$$1 - \gamma^5)u^\lambda_j) + (\bar{d}^\lambda_j\gamma^\mu(1 - \tfrac{8}{3}s_w^2 - \gamma^5)d^\lambda_j)] + \tfrac{ig}{2\sqrt{2}}W^+_\mu[(\bar{\nu}^\lambda\gamma^\mu(1 + \gamma^5)e^\lambda) +$$
$$(\bar{u}^\lambda_j\gamma^\mu(1 + \gamma^5)C_{\lambda\kappa}d^\kappa_j)] + \tfrac{ig}{2\sqrt{2}}W^-_\mu[(\bar{e}^\lambda\gamma^\mu(1 + \gamma^5)\nu^\lambda) + (\bar{d}^\kappa_j C^\dagger_{\lambda\kappa}\gamma^\mu(1 +$$
$$\gamma^5)u^\lambda_j)] + \tfrac{ig}{2\sqrt{2}}\tfrac{m^\lambda_e}{M}[-\phi^+(\bar{\nu}^\lambda(1 - \gamma^5)e^\lambda) + \phi^-(\bar{e}^\lambda(1 + \gamma^5)\nu^\lambda)] -$$
$$\tfrac{g}{2}\tfrac{m^\lambda_e}{M}[H(\bar{e}^\lambda e^\lambda) + i\phi^0(\bar{e}^\lambda\gamma^5 e^\lambda)] + \tfrac{ig}{2M\sqrt{2}}\phi^+[-m^\kappa_d(\bar{u}^\lambda_j C_{\lambda\kappa}(1 - \gamma^5)d^\kappa_j) +$$
$$m^\lambda_u(\bar{u}^\lambda_j C_{\lambda\kappa}(1+\gamma^5)d^\kappa_j] + \tfrac{ig}{2M\sqrt{2}}\phi^-[m^\lambda_d(\bar{d}^\lambda_j C^\dagger_{\lambda\kappa}(1+\gamma^5)u^\kappa_j) - m^\kappa_u(\bar{d}^\lambda_j C^\dagger_{\lambda\kappa}(1-$$
$$\gamma^5)u^\kappa_j] - \tfrac{g}{2}\tfrac{m^\lambda_u}{M}H(\bar{u}^\lambda_j u^\lambda_j) - \tfrac{g}{2}\tfrac{m^\lambda_d}{M}H(\bar{d}^\lambda_j d^\lambda_j) + \tfrac{ig}{2}\tfrac{m^\lambda_u}{M}\phi^0(\bar{u}^\lambda_j\gamma^5 u^\lambda_j) -$$
$$\tfrac{ig}{2}\tfrac{m^\lambda_d}{M}\phi^0(\bar{d}^\lambda_j\gamma^5 d^\lambda_j) + \bar{X}^+(\partial^2 - M^2)X^+ + \bar{X}^-(\partial^2 - M^2)X^- + \bar{X}^0(\partial^2 -$$
$$\tfrac{M^2}{c_w^2})X^0+\bar{Y}\partial^2 Y+igc_w W^+_\mu(\partial_\mu\bar{X}^0 X^- -\partial_\mu\bar{X}^+ X^0)+igs_w W^+_\mu(\partial_\mu\bar{Y}X^- -$$
$$\partial_\mu\bar{X}^+ Y) + igc_w W^-_\mu(\partial_\mu\bar{X}^- X^0 - \partial_\mu\bar{X}^0 X^+) + igs_w W^-_\mu(\partial_\mu\bar{X}^- Y -$$
$$\partial_\mu\bar{Y}X^+) + igc_w Z^0_\mu(\partial_\mu\bar{X}^+ X^+ - \partial_\mu\bar{X}^- X^-) + igs_w A_\mu(\partial_\mu\bar{X}^+ X^+ -$$
$$\partial_\mu\bar{X}^- X^-) - \tfrac{1}{2}gM[\bar{X}^+ X^+ H + \bar{X}^- X^- H + \tfrac{1}{c_w^2}\bar{X}^0 X^0 H] +$$
$$\tfrac{1-2c_w^2}{2c_w}igM[\bar{X}^+ X^0\phi^+ - \bar{X}^- X^0\phi^-] + \tfrac{1}{2c_w}igM[\bar{X}^0 X^-\phi^+ - \bar{X}^0 X^+\phi^-] +$$
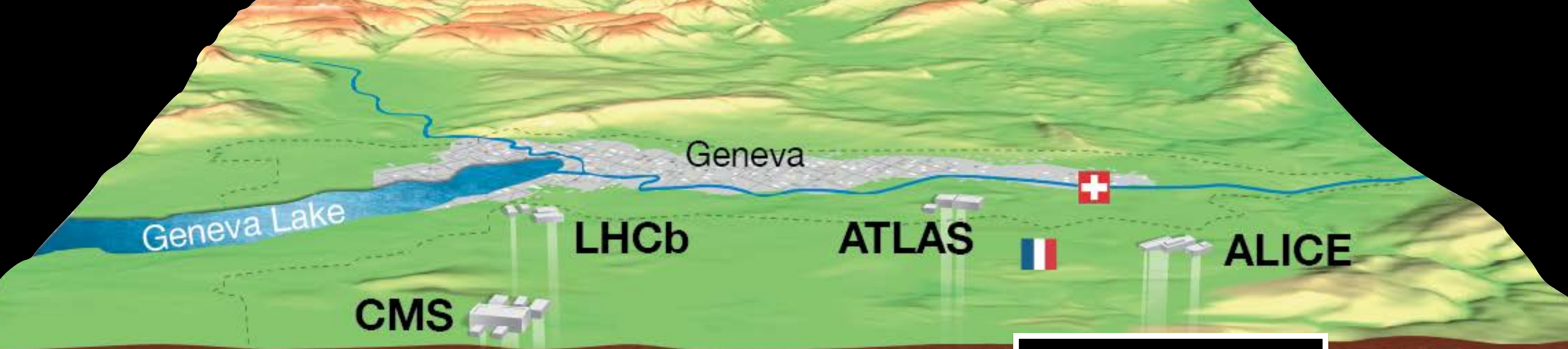$$igMs_w[\bar{X}^0 X^-\phi^+ - \bar{X}^0 X^+\phi^-] + \tfrac{1}{2}igM[\bar{X}^+ X^+\phi^0 - \bar{X}^- X^-\phi^0]$$

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

CMS

LHCb

ATLAS

ALICE

**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)**

LHC

Geneva

Geneva Lake

LHCb

AT

CMS

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

**5% of internet traffic to L1**

**L1 trigger: ~1000 FPGAs**

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

**Reject >99% of collisions!**

LHC

Geneva

Geneva Lake

LHCb

CMS

AT

LHC

L1 bit:
Accept = 1
Reject = 0

Geneva

Geneva Lake

LHCb

ATLAS

ALICE

CMS

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

LHC

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

Reject further 99%!

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

LHCb    ATLAS    ALICE

LHC

TIER 0: ∞

High Level Trigger:
Latency O(100) ms

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

Geneva

LHCb

ATLAS

CMS

LHC

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

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

**TIER 0: ∞**

Geneva

LHCb

ATLAS

CMS

LHC

High Level Trigger: Latency O(100) ms

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

TIER 0: ∞

0.0025% of collision events remaining

Geneva

LHCb

ATLAS

CMS

LHC

To make sure we select "the right" 0.0025%, algorithms must be
- Fast (get more data through)
- Accurate (select the right data)

"Probability" of producing "anything"

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

Need **more** data!

New Physics?

# 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 |

# High Luminosity LHC

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

200 vertices
(average 140)



CMS: event with 78 reconstructed vertices

6 cm

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

Run 4+5

Maintain physics acceptance → better detectors

CMS High Granularity (endcap) calorimeter
- X20 times more readout channels (6.5 million!!)

More collisions
More readout channels



silicon
scintillator
~10λ
~5λ
CE-E    CE-H

CMS HGCAL TDR

HL-LHC Level-1:

Complete re-design of Level-1

# HL-LHC Level-1:

**Complete re-design of Level-1**
- Charged particle tracks (6.4 Tb/s, 200 FPGAs)



*Simulated event display with average pileup of 140*

**em**

*ics*

*, no tracking information*



$\langle \mu \rangle = 32$

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

# HL-LHC Level-1:

**Complete re-design of Level-1**
- Charged particle tracks
- Particle Flow (40 FPGAs)

# HL-LHC Level-1:

**Complete re-design of Level-1**
- Charged particle tracks
- Particle Flow (40 FPGAs)

# HL-LHC Level-1:

**Complete re-design of Level-1**
- Charged particle tracks
- Particle Flow
- HGCal (4 Tb/s, 200 FPGAs)

# HL-LHC Level-1:

**Complete re-design of Level-1**
- Charged particle tracks
- Particle Flow
- HGCal

**Input data**
- 2 Tb/s $\rightarrow$ **63 Tb/s**

**Latency**
- 4 μs $\rightarrow$ **12 μs**

Extremely high data complexity,

Extremely little time

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Super... Sol...

4T

2T

... 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

CALORIMETRY

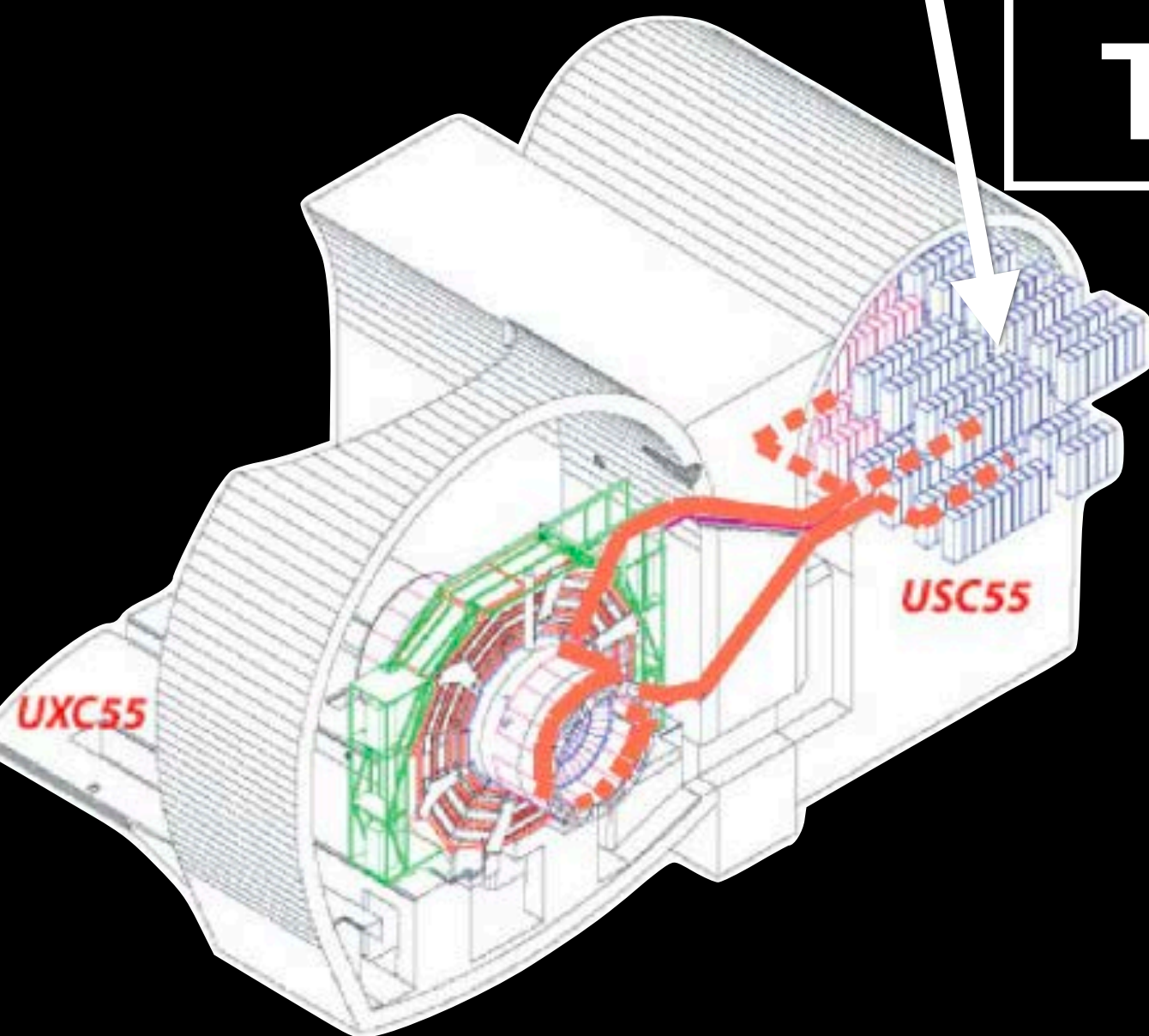PARTICLE FLOW

MUONS

**PARTICLE FLOW: 66 FPGAs**

**GLOBAL TRIGGER: 12 FPGAs**

EXTERNAL TRIGGERS

**Trigger accept/reject**

**12.5 µs**

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

**63 Tb/s**

return yoke interspersed with Muon chambers

CALORIMETRY: 370 FPGAs

*54 for HGCAL only!

Xilinx Ultrascale+ FPGAs

TRACKING 174 FPGAs

MUONS:

*Simulated event display with average pileup of 140*

PARTICLE FLOW

UXC55

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

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

$<\mu> = 32$

**63 Tb/s**

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Xilinx Ultrascale+ FPGAs

return yoke interspersed with Muon chambers

**BRAND NEW SYSTEM, NOT YET BUILT!**
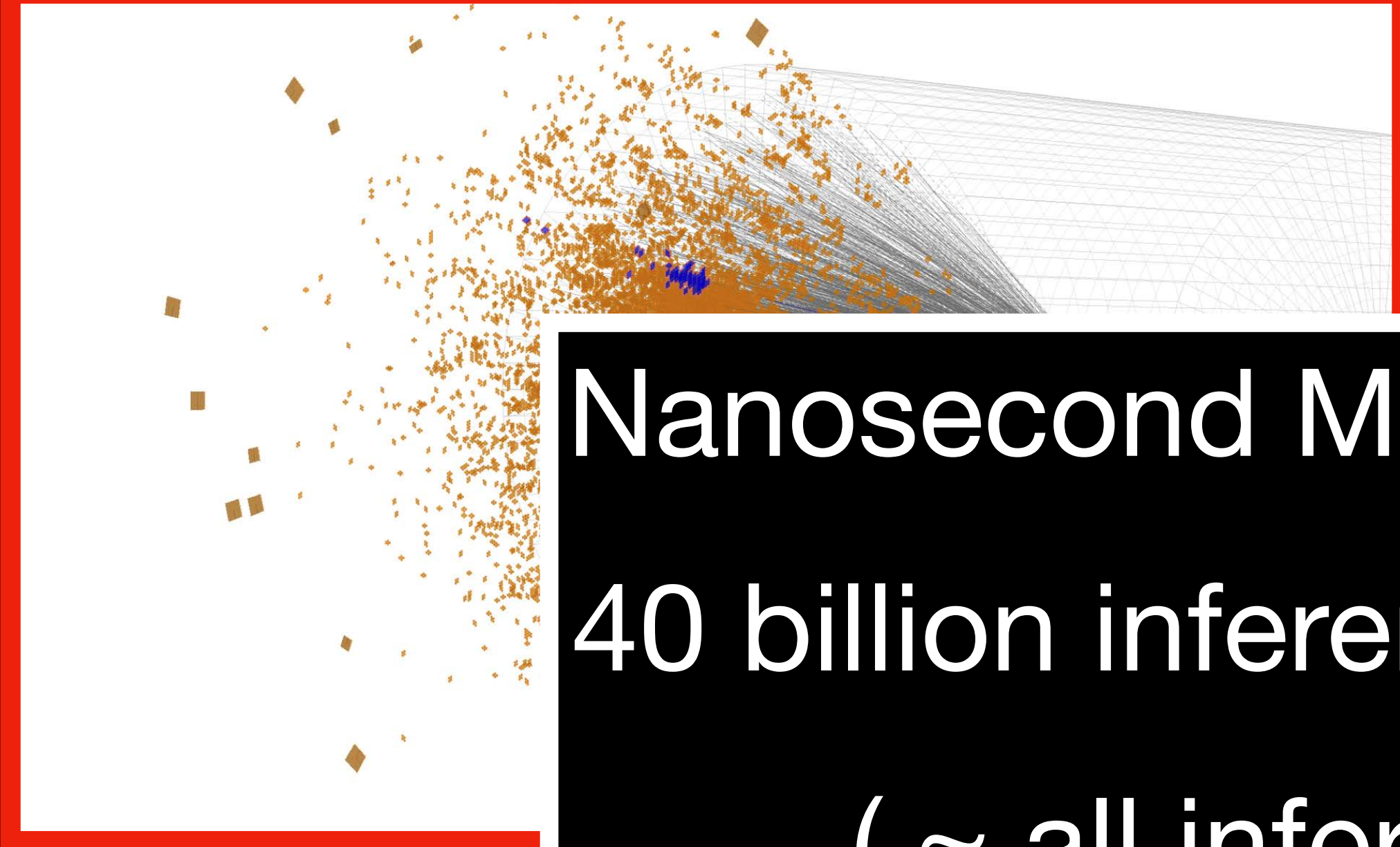**The time to design algorithms is now**

event display with average pileup of 140

PARTICLE FLOW

UXC55

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

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

$<\mu> = 32$

12 microseconds latency

Processing 5% of internet traffic

*Simulated event display with average pileup of 140*

# System

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

## **Journey to HL-LHC**

*s*

*LT*

**run:**

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

**TLAS, 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**     **40 MHz**

**L1 output: 75 kHz**   **L1 output: 100 kHz**

Detectors

Front end pipelines

Readout buffers

Switching network

Processor farms

**~3 kHz**

**200 Hz**

*HLT output:   ~1 kHz*

**40 MHz**

**L1 output: 100 kHz**

**750 kHz**

**100 Hz**

**7.5 kHz**

Lvl-1

Lvl-2

Lvl-3

LvI-1

HLT

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

MUONS

Trigger

ept/reject

**12.5 µs**

$\langle\mu\rangle = 32$

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

Recorded Luminosity (pb⁻¹/1.00)

Mean number of interactions per crossing

Nanosecond ML inference on FPGAs!

40 billion inferences/s during HL-LHC

( ≈ all inferences at Google)

**Journey to HL-LHC**

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

**run:**

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

**TLAS, CMS operating:**

ccept ≤ 100 kHz,

ncy ≤ 2.5 (AT), 4 μsec (CM)

Accept ≤ 1 kHz

**LAS & CMS will be:**

**5 x 10³⁴**

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

Trigger
ept/reject

**12.5 μs**

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**

Conifer  hls 4 ml

**run:**

**7 x 10$^{33}$, PU = 30, E = 7 TeV, 50 nsec bunch spacing**

**TLAS, CMS operating:**

ccept ≤ 100 kHz,

ncy ≤ 2.5 (AT), 4 µsec (CM)

Accept ≤ 1 kHz

**LAS & CMS will be:**

**5 x 10$^{34}$**

$<\mu> = 32$

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

Mean number of interactions per crossing

**LHC**

40 MHz → Detectors → 40 MHz

Front end pipelines

Lvl-1

L1 output: 75 kHz

L1 output: 100 kHz

Lvl-2

Readout buffers

~3 kHz

Switching network

Lvl-3

Processor farms

200 Hz

HLT output: ~1 kHz

**HL-LHC**

40 MHz → Detectors

Front end pipelines

Lvl-1

L1 output: 100 kHz

Readout buffers

Switching network

HLT

Processor farms

100 Hz

750 kHz

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

MUONS

Trigger ept/reject

**12.5 µs**

**Current HL-LHC design**

**Foundation-model based trigger**

CALORIMETRY
(370 FPGAs)

CHARGED PARTICLE
TRACKING
(174 FPGAs)

MUON CHAMBERS
(96 FPGAs)

PARTICLE
FLOW
(66 FPGAs)

GLOBAL
TRIGGER
(13 FPGAs)

Accept / Reject

63 Tb/s

12.5 μs

CALORIMETER
PRE-PROCESSING

TRACKING
PRE-PROCESSING

MUON CHAMBERS
PRE-PROCESSING

End-to-end reconstruction model

Latent representation

Downstream
Task

Downstream
Task

Downstream
Task

Downstream
Task

Accept / Reject

# Why FPGAs?

# Why FPGAs?

- Latency (resource parallelism)



resource parallelism

# Why FPGAs?

- Throughput (pipeline parallelism)



pipeline

parallelism

**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 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

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
```

Ideally

Reality

Ideally

- **Quantization**
- **Pruning**
- **Parallelisation**
- **Knowledge distillation**

Reality

# Quantization



**Floating point 32:**
**4B numbers in [-3.4e38, +3.4e38]**

# Quantization



**Quantising:**
**int8 $2^8$=256 numbers in [-128,127]**

$$x_q = Clip(Round(\frac{x_f}{scale}))$$

Weights Layer 1

FP 32

Weights Layer 2

FP 32

Fixed point

Weights Layer 1

< 4,0 >

Weights Layer 2

< 4,0 >

Fixed point

0101.1011101010

integer    fractional
        width

Weights Layer 1

< 4,0 >

Weights Layer 2

< 4,0 >

hls4ml + Google

Quantization-aware training

Forward pass →

← Back propagation

Nature Machine Intelligence 3 (2021)

**hls4ml** + **Google**

# Quantization-aware training

```
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)
```
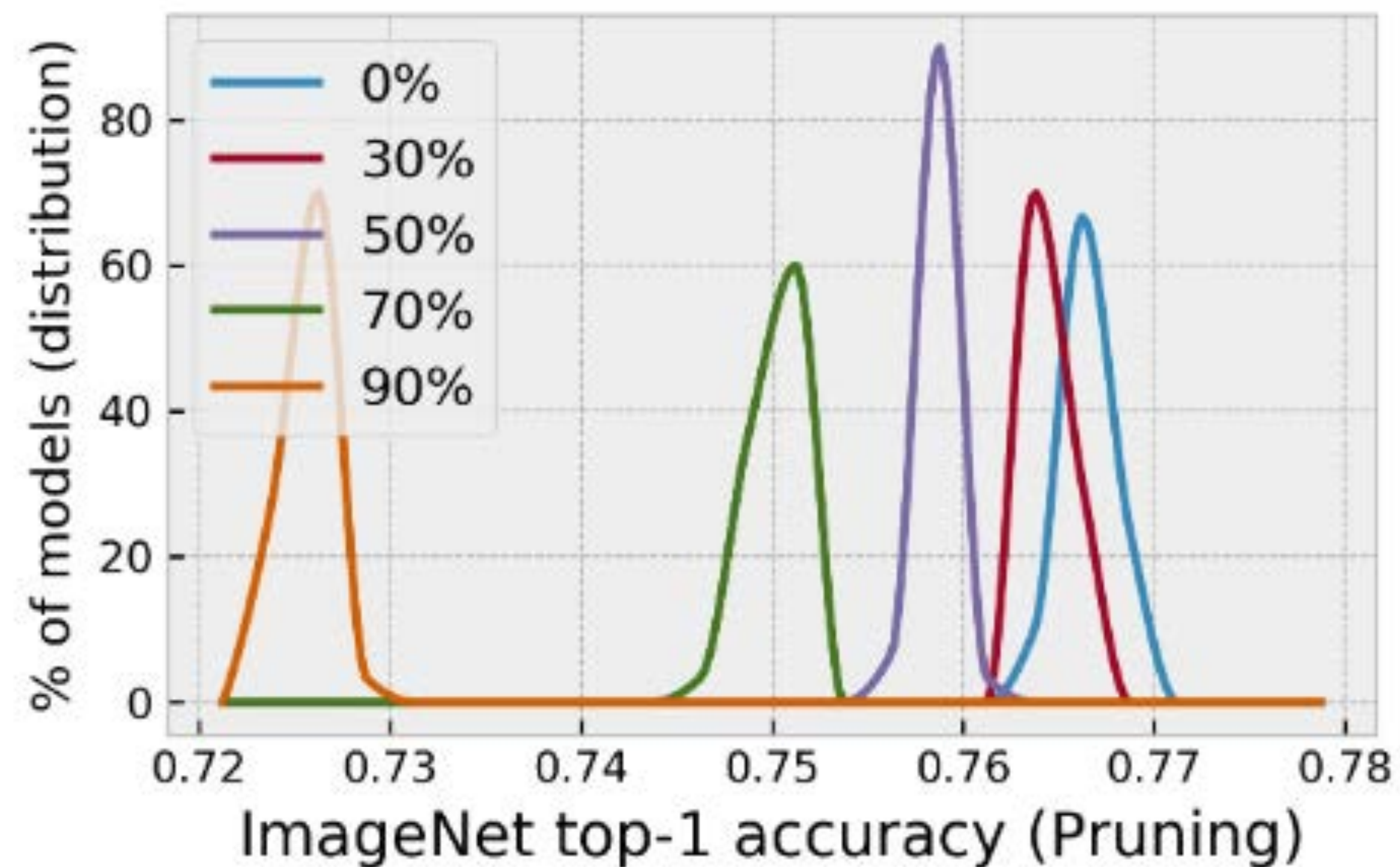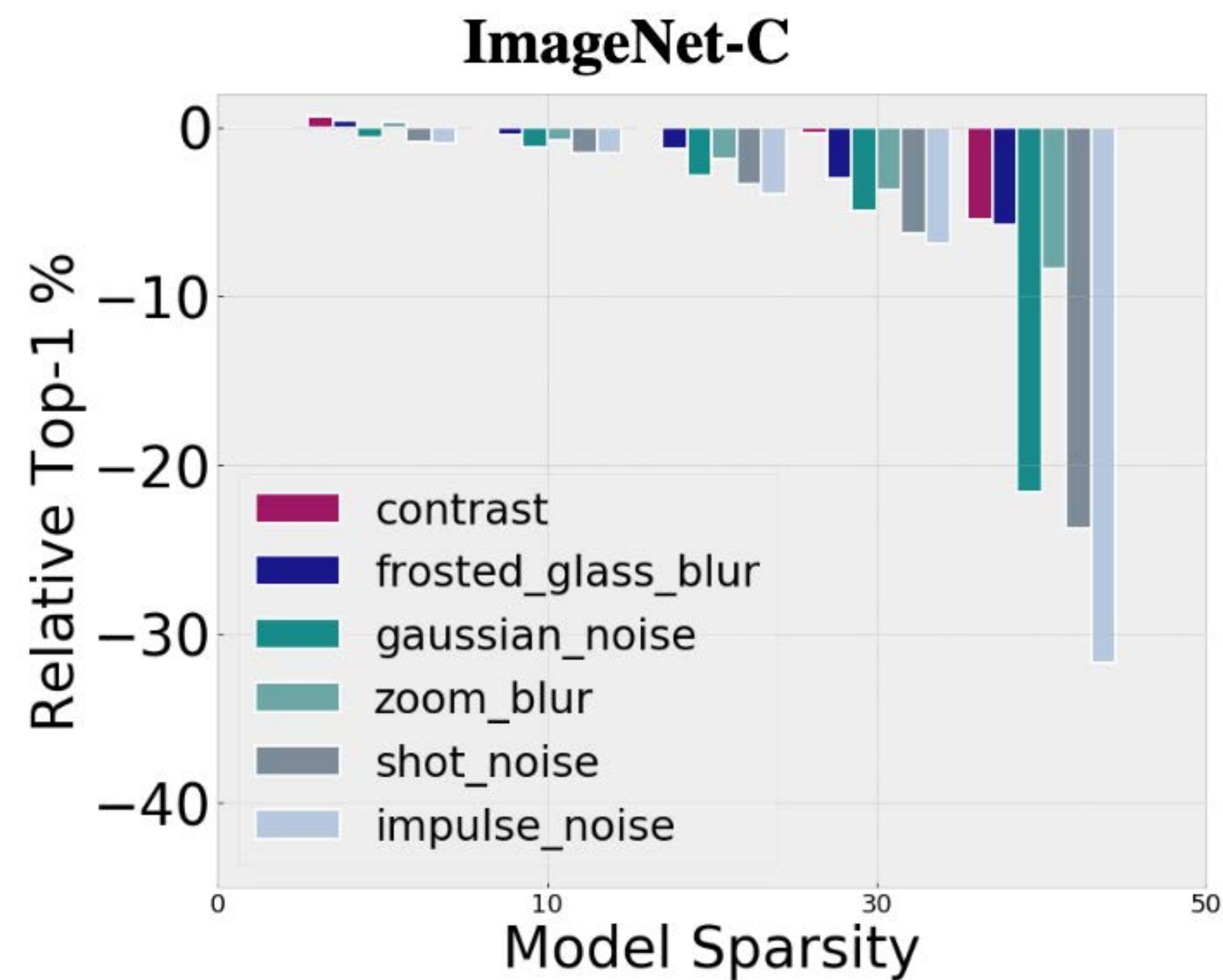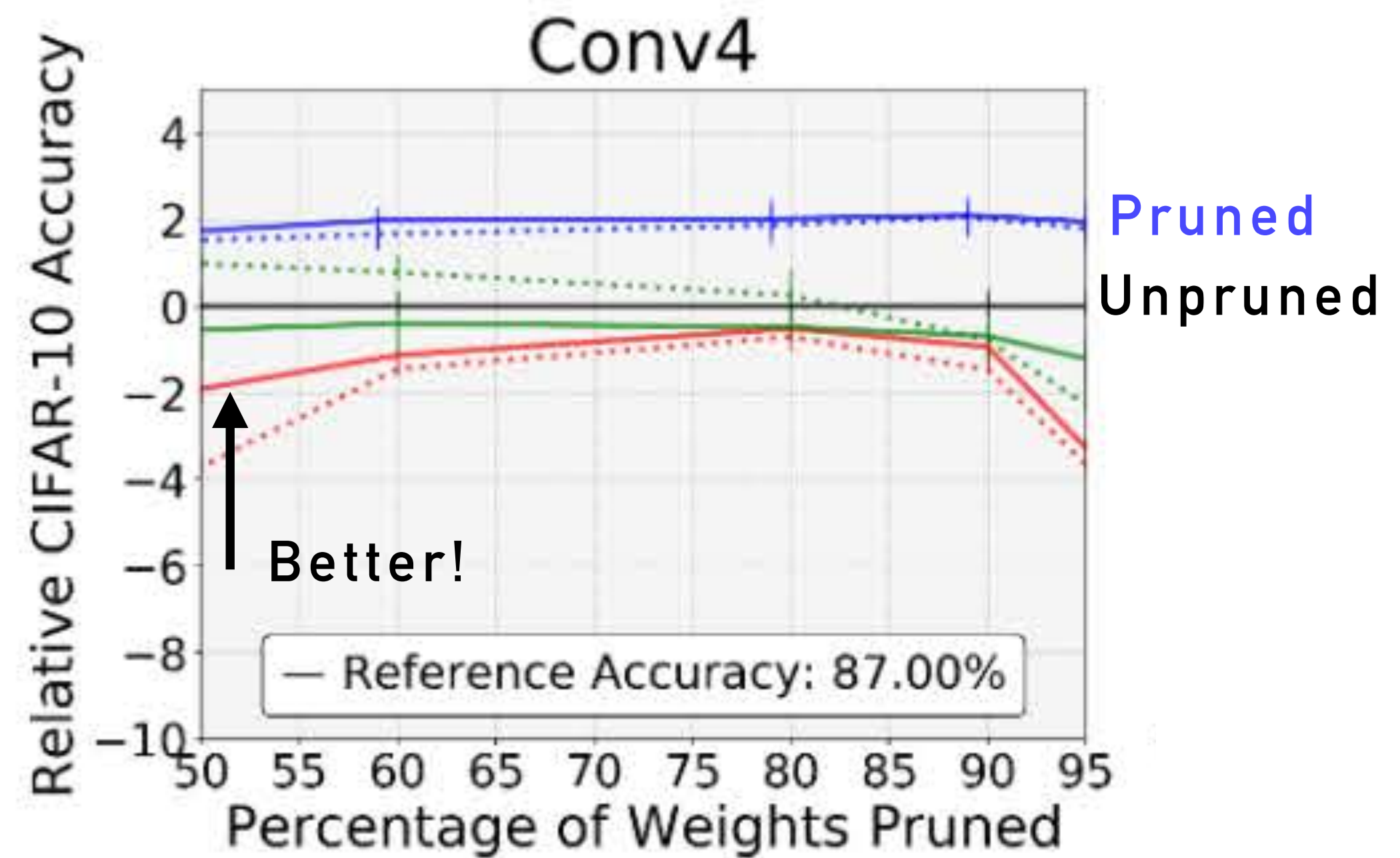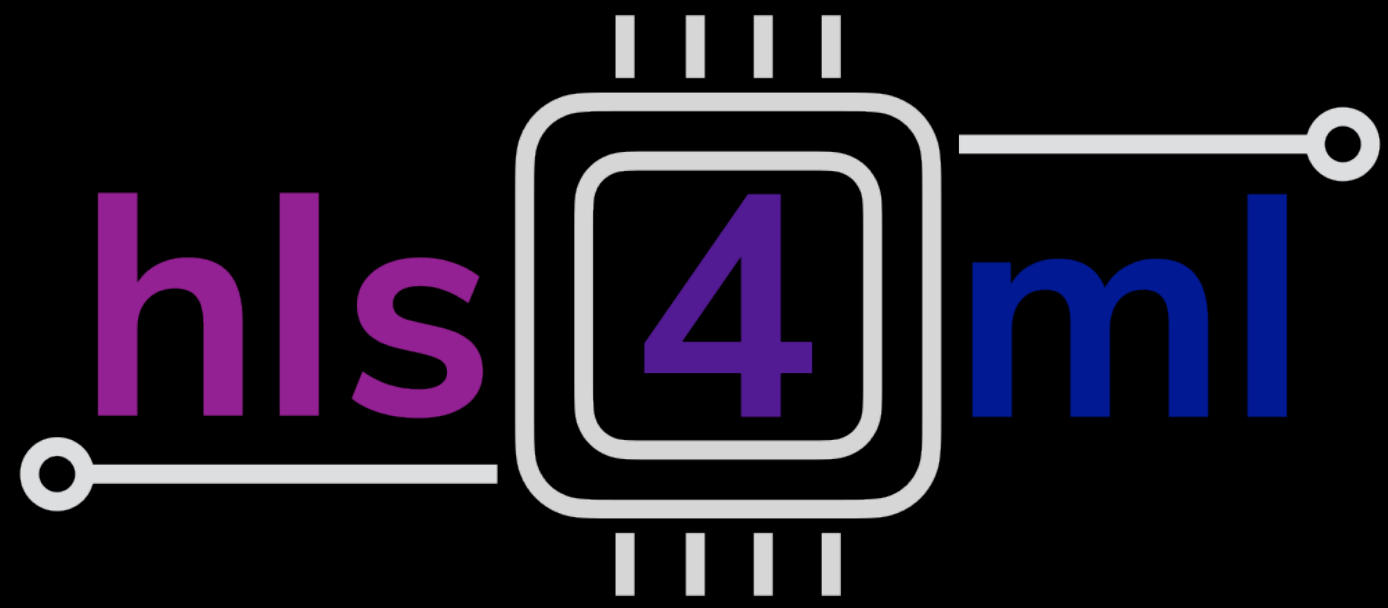
**Nature Machine Intelligence 3 (2021)**

# Pruning



before pruning

# Pruning

*From Brian Bartoldson*



Original image | Gaussian Noise | Shot Noise | Impulse Noise

**ImageNet-C**



contrast
frosted_glass_blur
gaussian_noise
zoom_blur
shot_noise
impulse_noise

Relative Top-1 %

Model Sparsity



% of models (distribution)

0%
30%
50%
70%
90%

ImageNet top-1 accuracy (Pruning)

Conv4

Pruned
Unpruned

Better!

Reference Accuracy: 87.00%

Relative CIFAR-10 Accuracy
Percentage of Weights Pruned

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)

hls4ml tutorial

Nanosecond ML inference on FPGAs!

40 billion inferences/s during HL-LHC

( ≈ all inferences at Google)

HEP developed libraries for fast ML on FPGAs

VBF H (γγ)

200 vertices
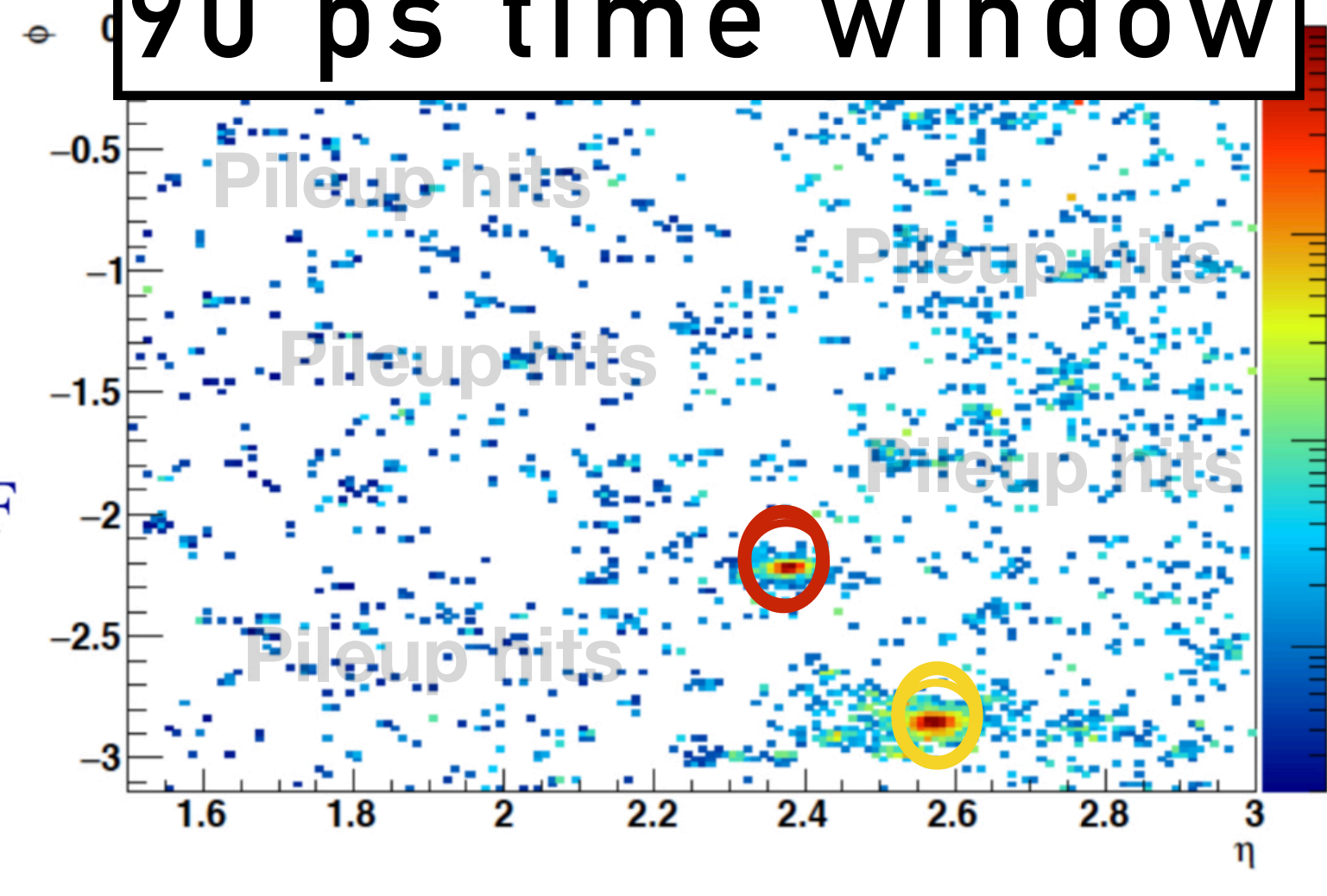
$10^2$

$10$

+

Layer 30

CE-H
Layer 30

$10^3$

$10^3$

$10^3$

$10^2$

y(mm)

0

-100

-200

Layer 5

$10^2$

d one VBF jet in the same quadrant,

jet

Cut $\Delta t < 90ps$ (3σ at 30ps)

Layers projected onto one plane
-require hits within 90ps time window-

0

-0.5

-1

VBF H (γγ)

**No timing cut**

**90 ps time window**

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

$$\text{Cut } \Delta t < 90\text{ps} \quad (3\sigma \text{ at } 30\text{ps})$$

Layers projected onto one plane

-require hits within 90ps time window-
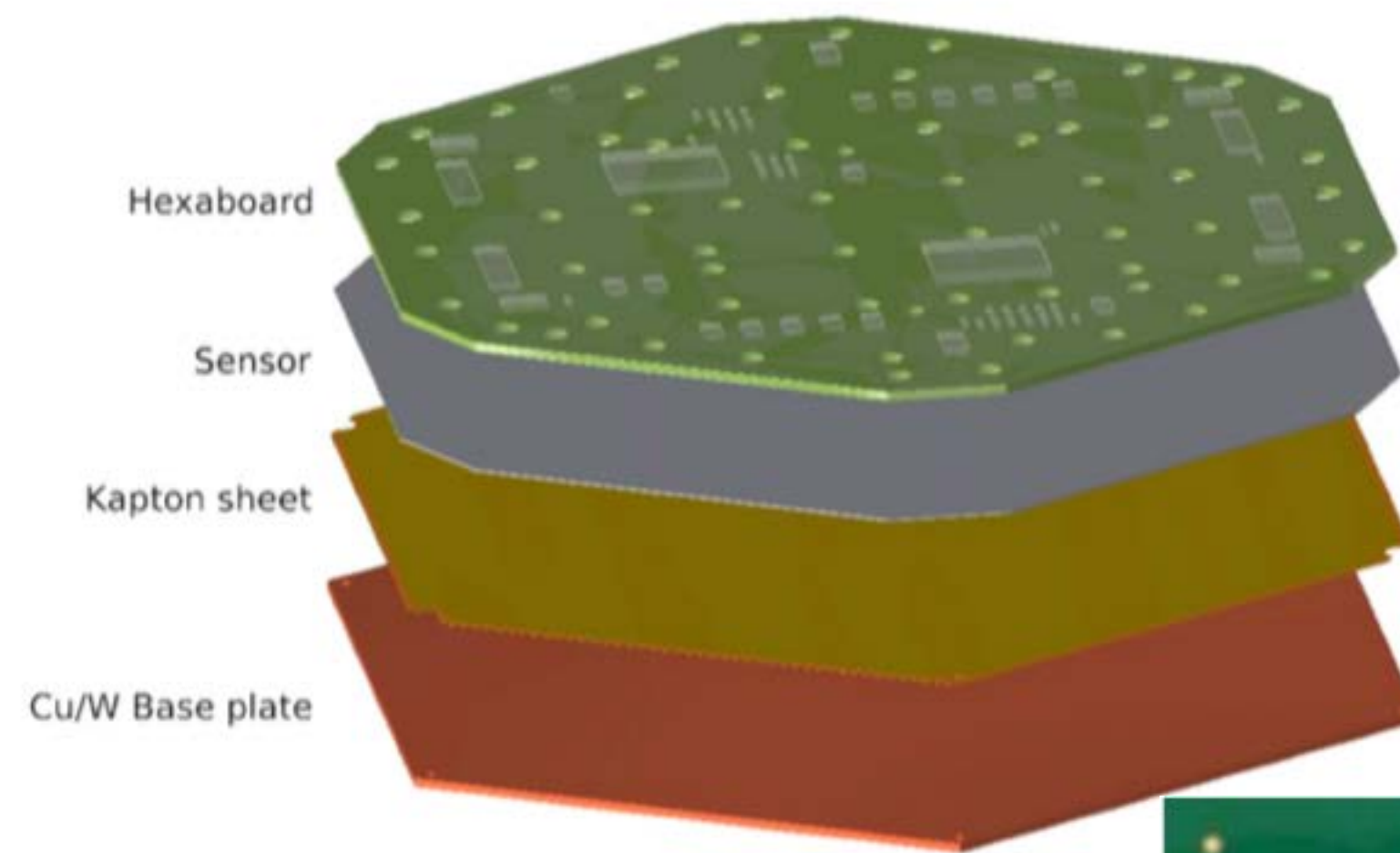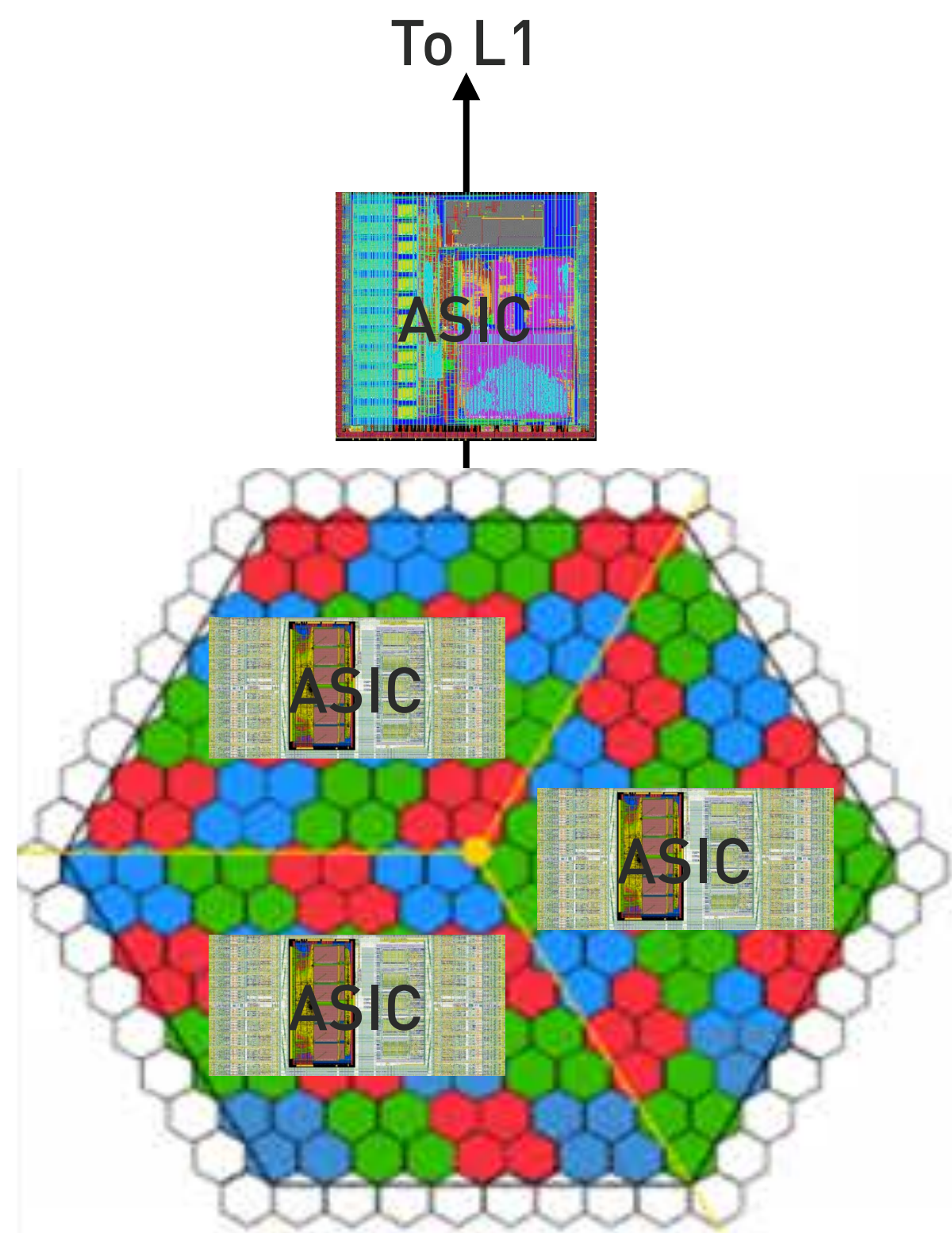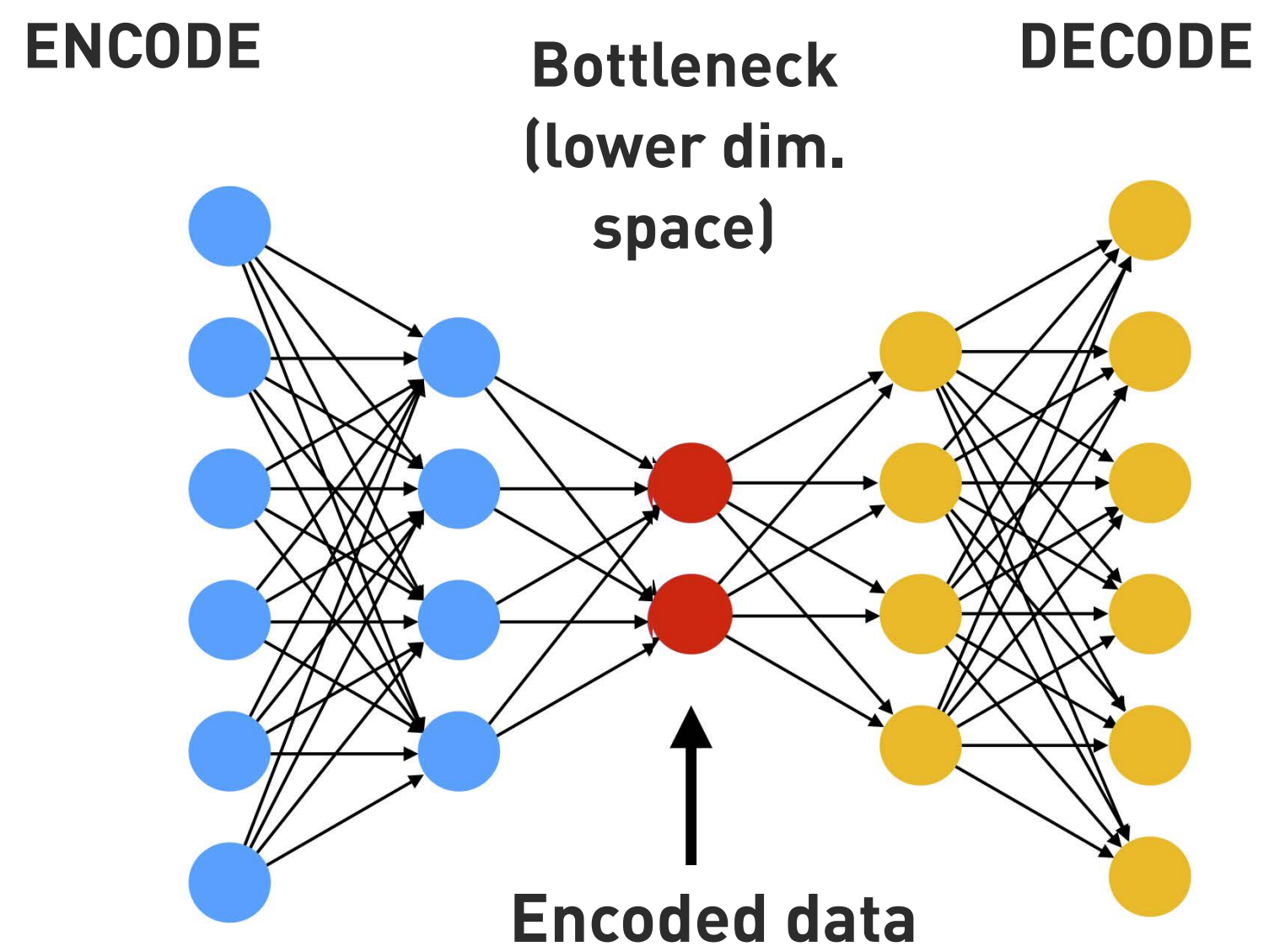
jet

200 vertices

$10^2$

$10$

+

VBF H (γγ)

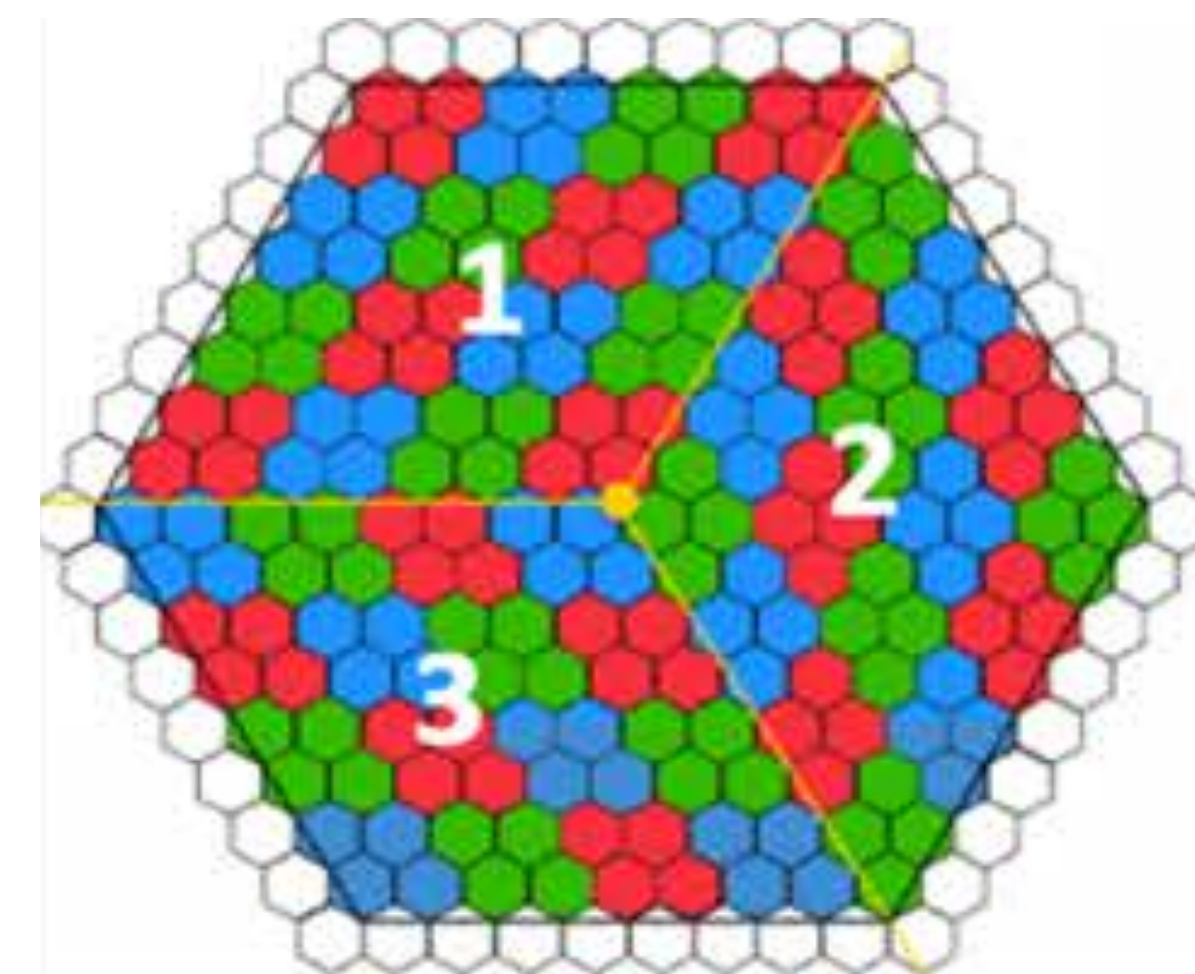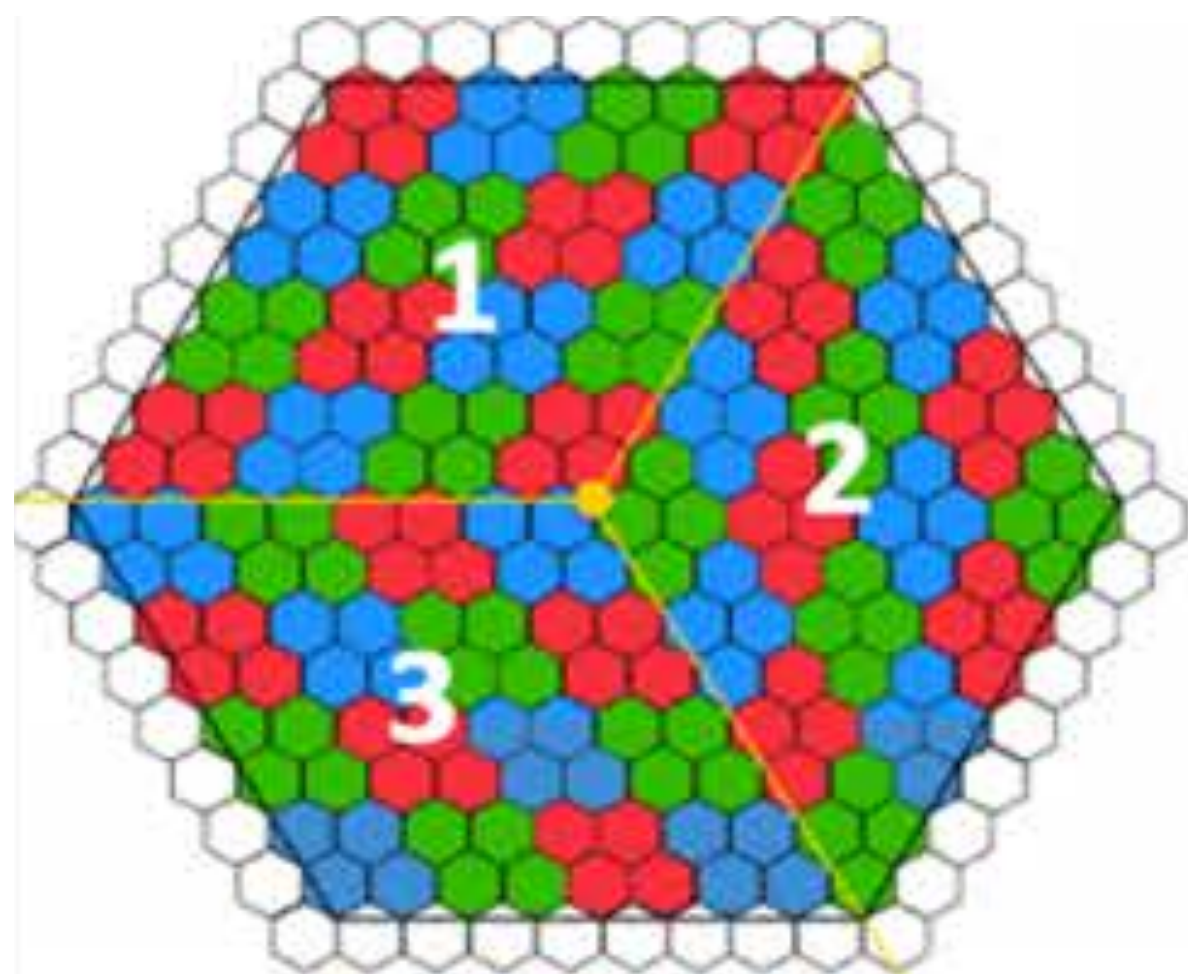BUT: Cannot read out all these channels fast enough for L1 to trigger!
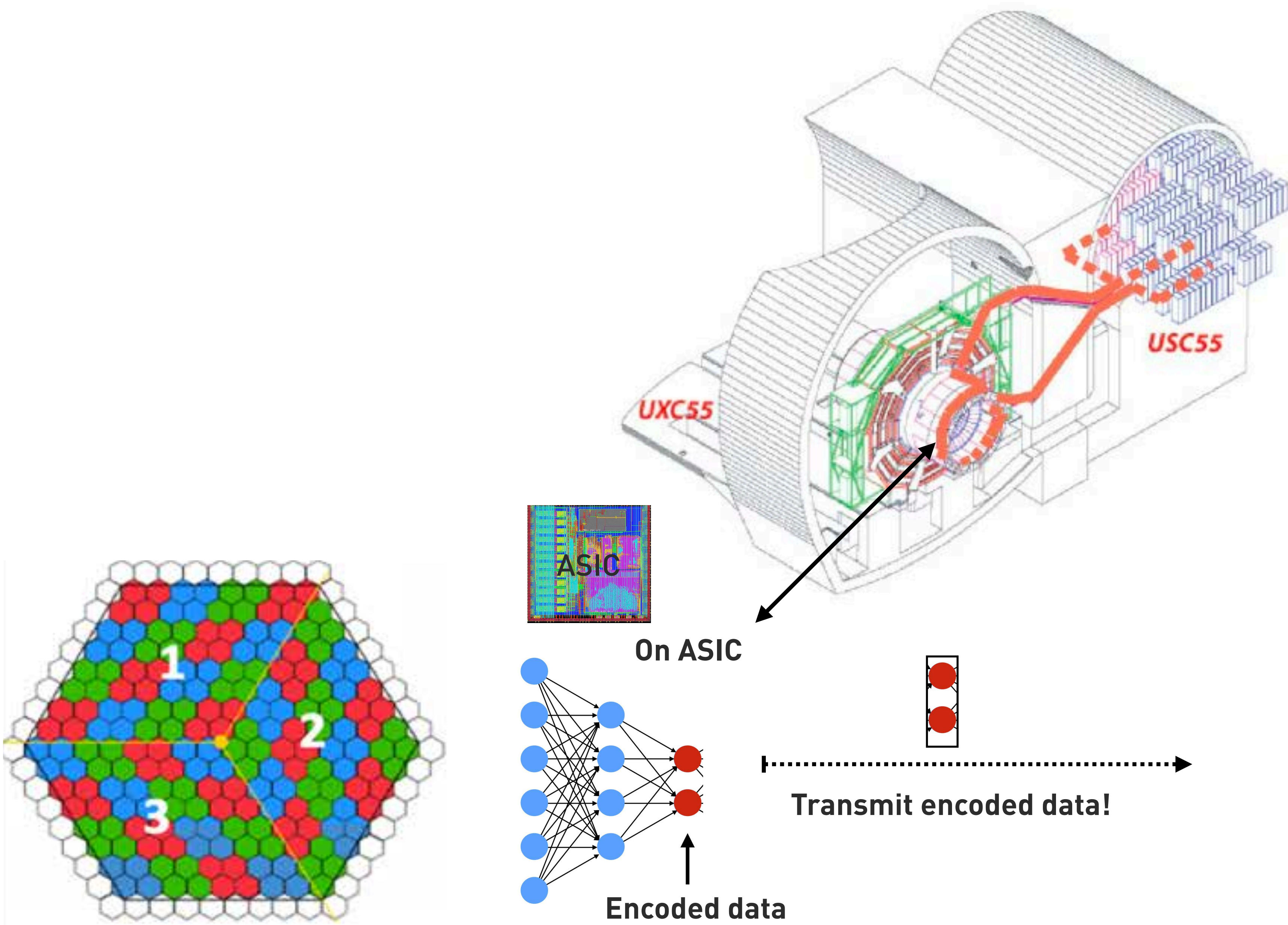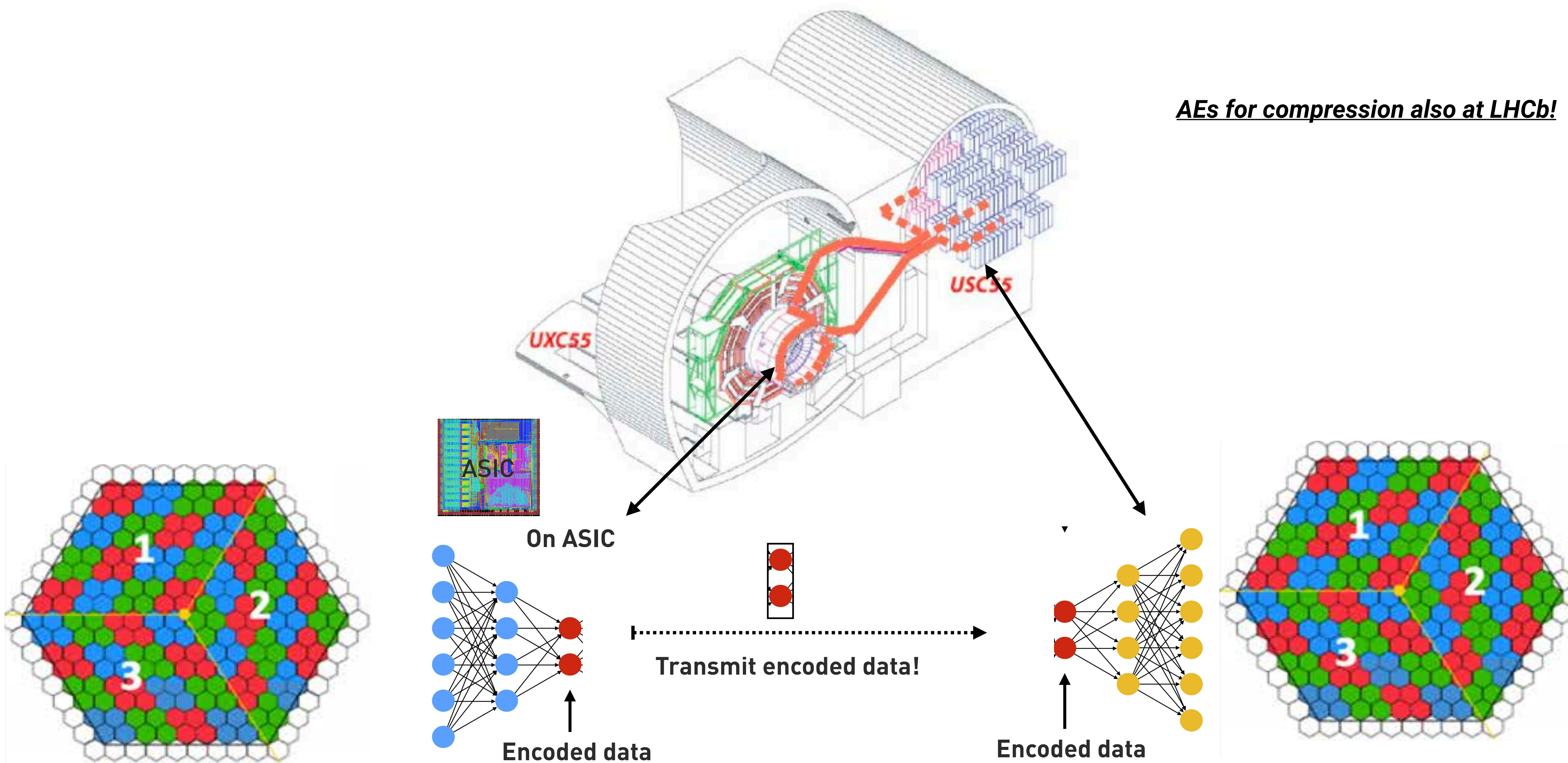
window

jet

200 vertices

$10^2$

10

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

Layers projected onto one plane

-require hits within 90ps time window-

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
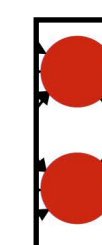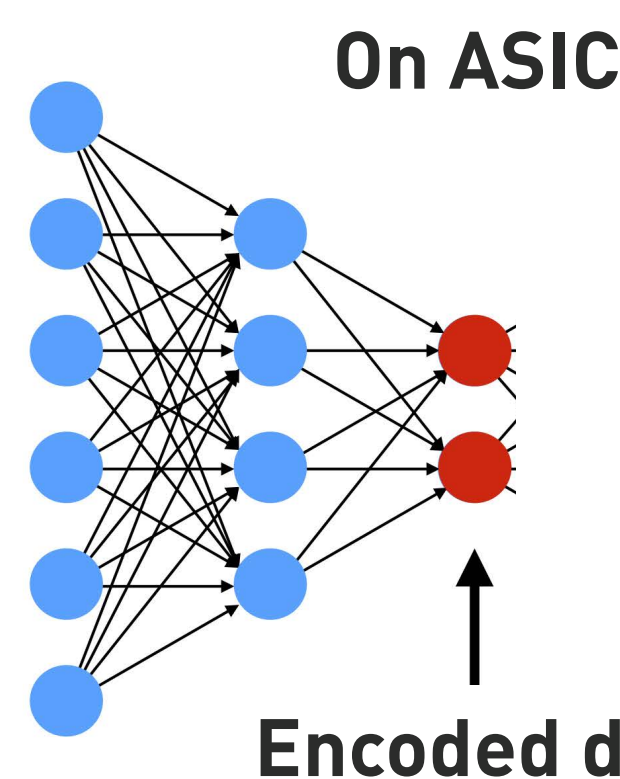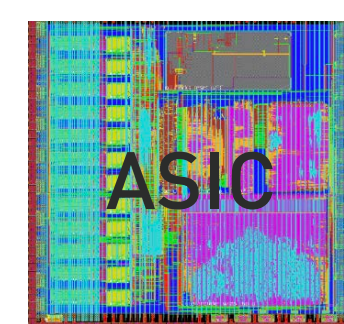
ENCODE     Bottleneck (lower dim. space)     DECODE

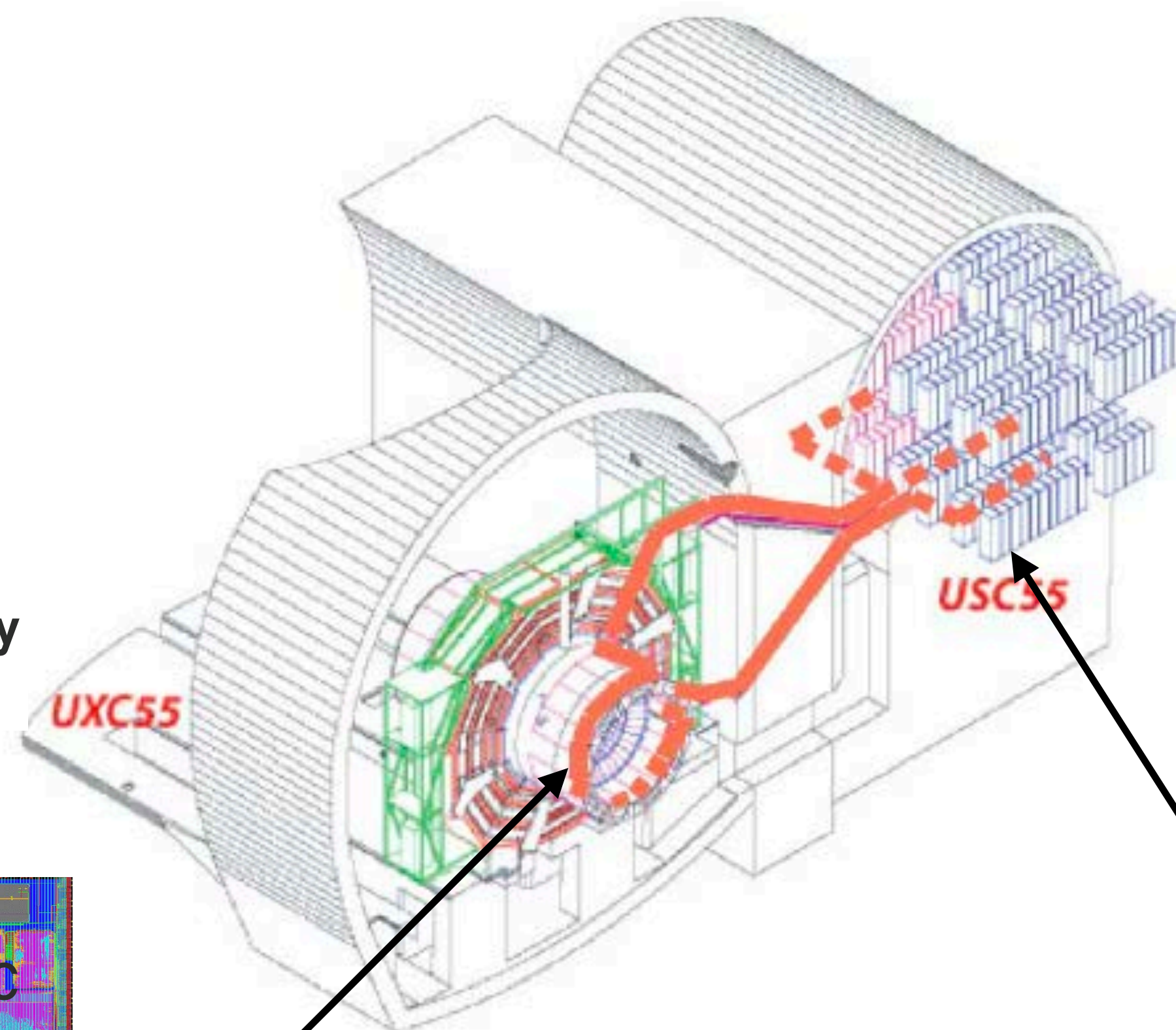Encoded data

# 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)

Signal

TOT

TOA

Time

signal pulse

Interpretation of recorded TOA amplitudes

Sensor PCBs

On detector    Off detector

FE ASICs    FE ASICs    FE ASICs

~160-320Gbit/s

DAQ

Concentrator ASIC    Optical link

About 8k pairs at 10Gbit

Trigger

About 8k pairs at 10Gbit

Panel PCBs

DAQ, Trigger data

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

5-6mm

Power, Clock, Trigger & Slow-control

3

**Encoded data**

**Transmit encoded data!**

*AEs for compression also at LHCb!*



On ASIC

Transmit encoded data!

Encoded data

Encoded data

*AEs for compression also at LHCb!*
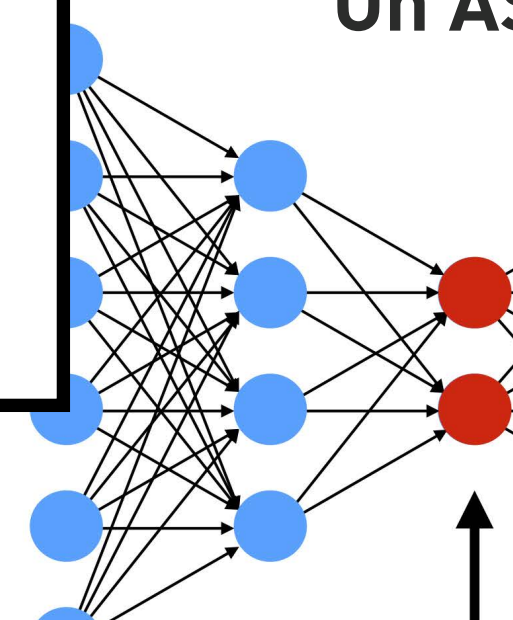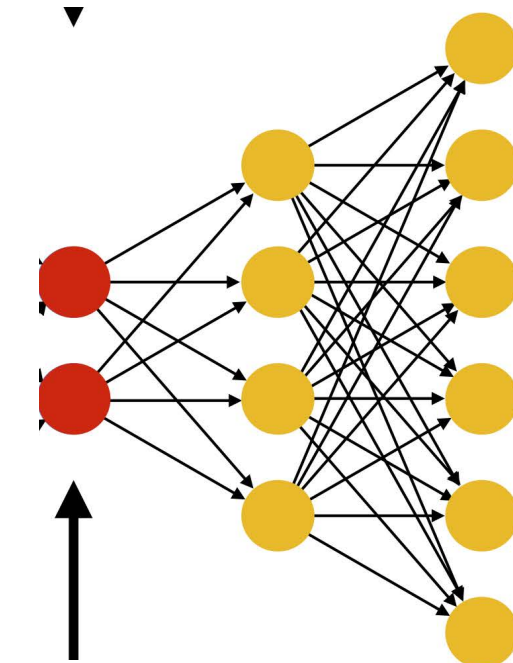
- **75-100 mW**
- **Triplicated w/b for radiation safety**
  **Reprogrammable w/b over IC2!**



**On ASIC**

**Transmit encoded data!**

**Encoded data**

**Encoded data**

*AEs for compression also at LHCb!*

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

**FKeras**
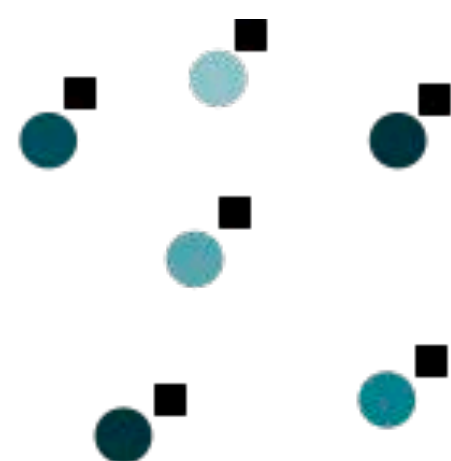
**On ASIC**

**Bit flip!**
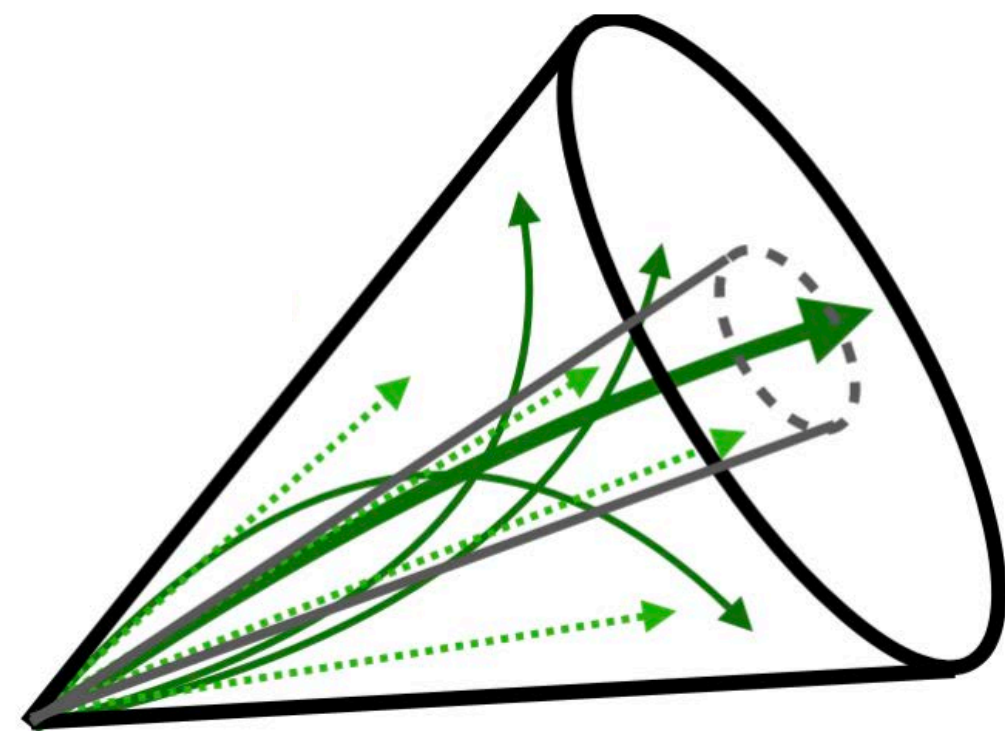
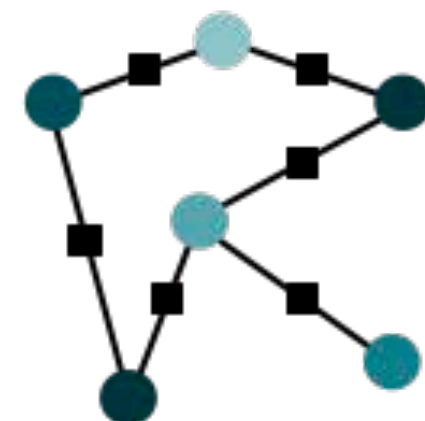**Transmit encoded data!**

**Encoded data**

**Encoded data**

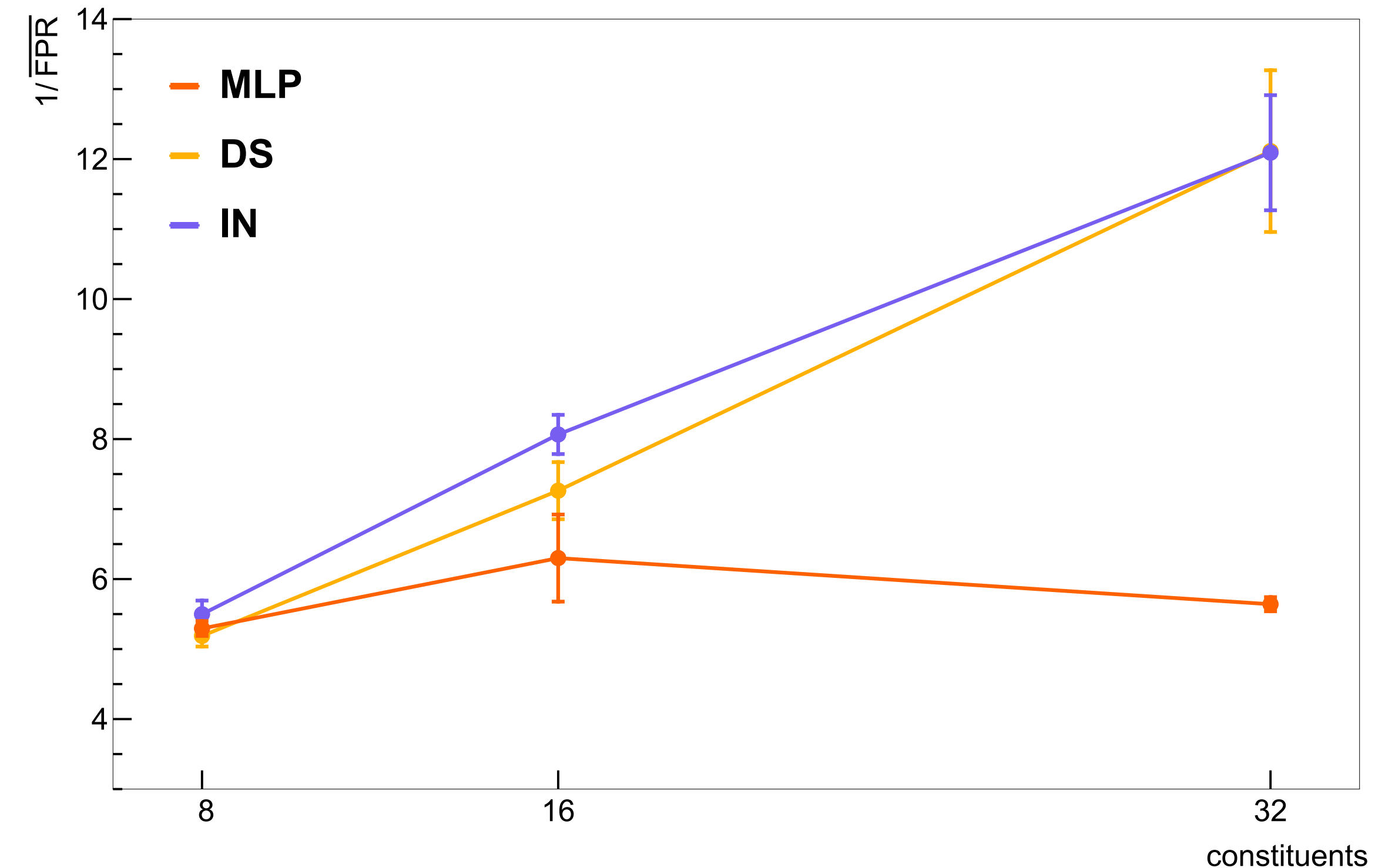# Invariance vs equivariance, sets vs graphs for smaller models?



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

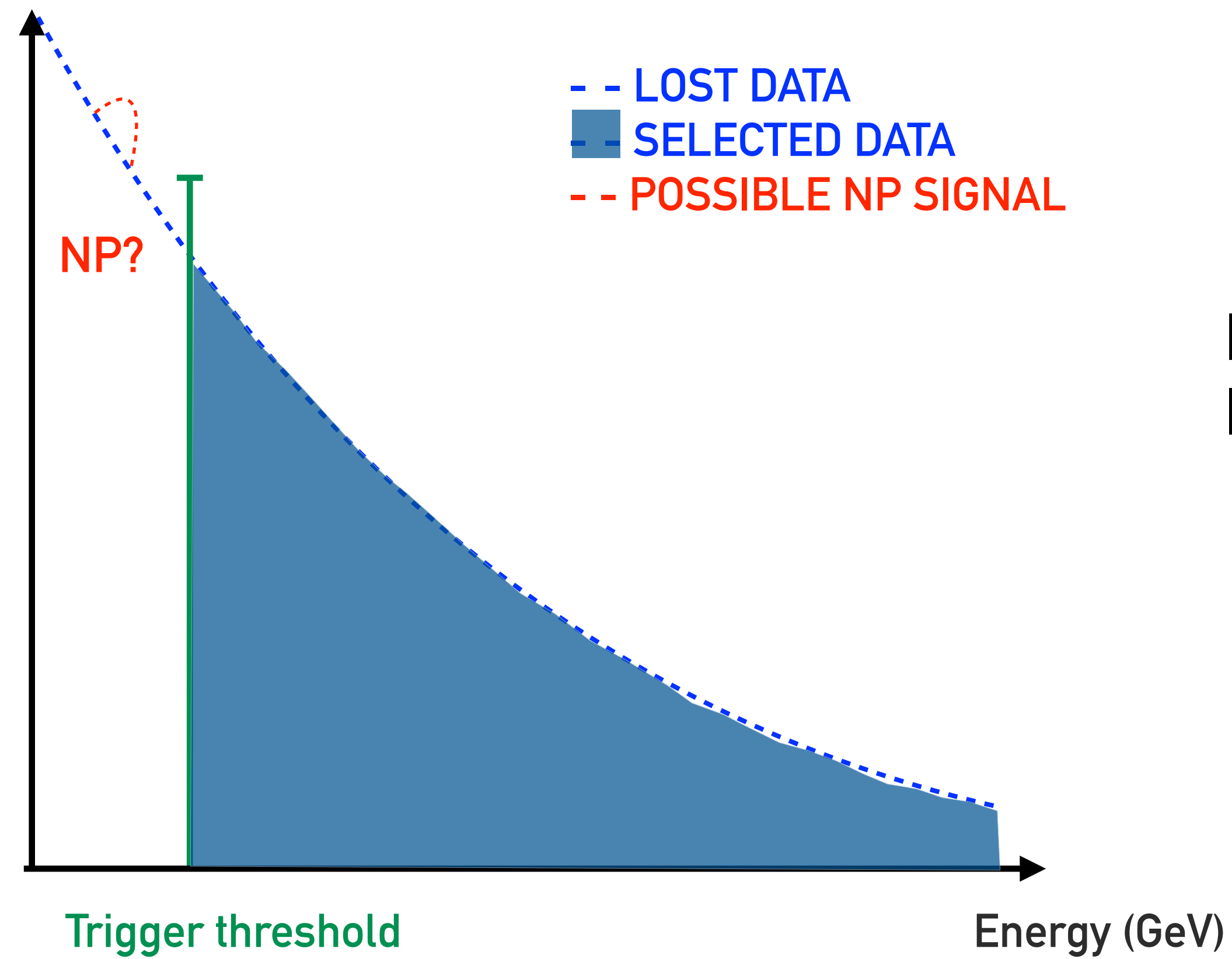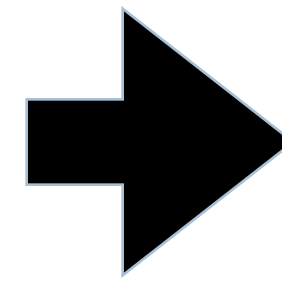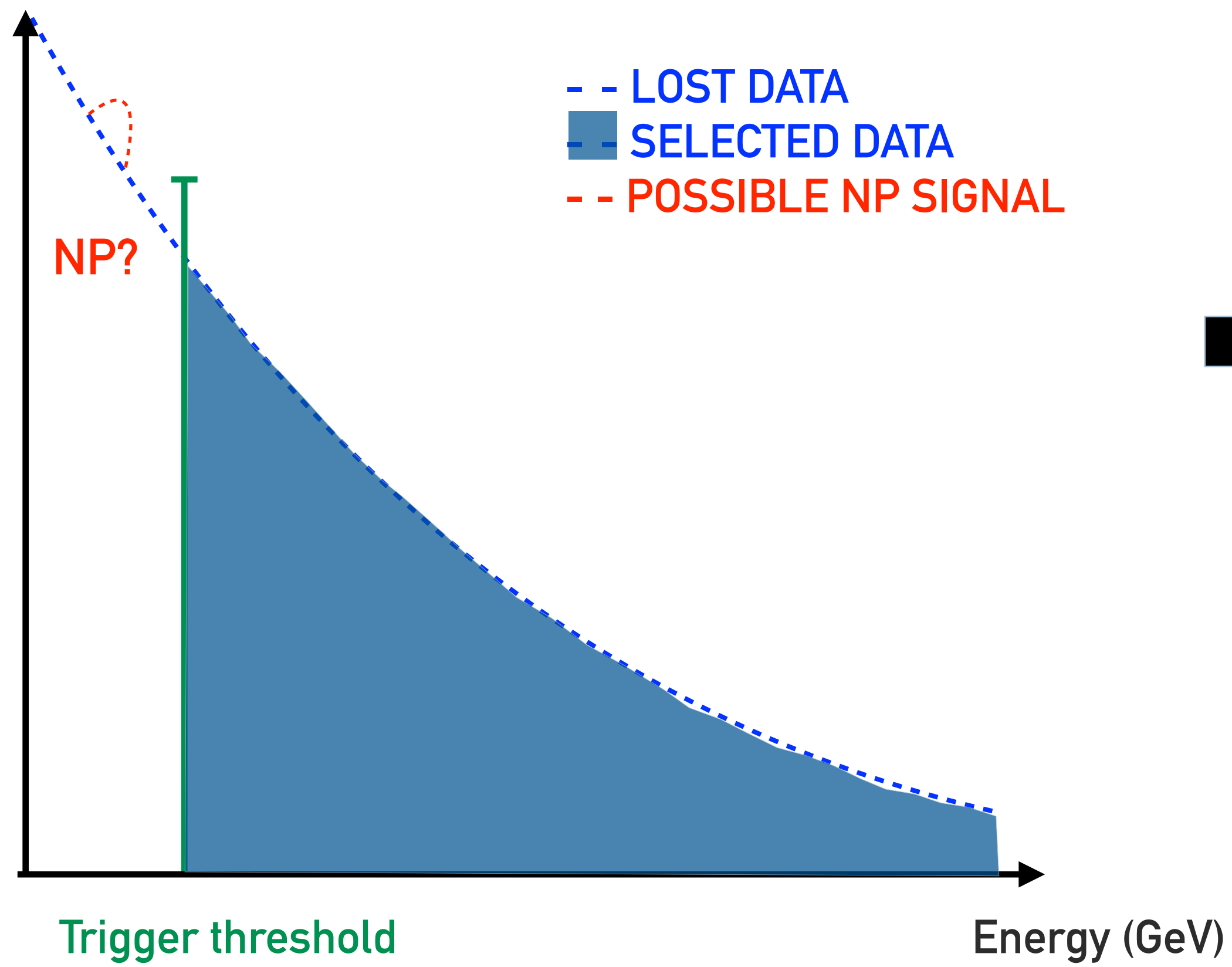**Graphs:** Information is assigned to edges, i.e., pairs of nodes.

● **Nodes**
— **Edges**
■ **Features**



FPGA: Xilinx Virtex UltraScale+ VU13P

| Architecture | Constituents | RF | Latency [ns] (cc) | II [ns] (cc) | DSP | LUT |
|---|---|---|---|---|---|---|
| MLP | 8 | 1 | 105 (21) | 5 (1) | 262 (2.1%) | 155,080 (9.0%) |
| | 16 | 1 | 100 (20) | 5 (1) | 226 (1.8%) | 146,515 (8.5%) |
| | 32[a] | 1 | 105 (21) | 5 (1) | 262 (2.1%) | 155,080 (7.2%) |
| DS | 8 | 2 | 95 (19) | 15 (3) | 626 (5.1%) | 386,294 (22.3%) |
| | 16 | 4 | 115 (23) | 15 (3) | 555 (4.5%) | 747,374 (43.2%) |
| | 32[a] | 8 | 130 (26) | 10 (2) | 434 (3.5%) | 903,284 (52.3%) |
| IN | 8 | 2 | 160 (32) | 15 (3) | 2,191 (17.8%) | 472,140 (27.3%) |
| | 16 | 4 | 180 (36) | 15 (3) | 5,362 (43.6%) | 1,387,923 (80.3%) |
| | 32[a] | 8 | 205 (41) | 15 (3) | 2,120 (17.3%) | 1,162,104 (67.3%) |

# Limitations of current trigger



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

**LOST DATA**
**SELECTED DATA**
**POSSIBLE NP SIGNAL**

NP?

Trigger threshold

Energy (GeV)

Look at **data** rather than defining signal hypothesis a priori
• Can we "classify" objects/events?

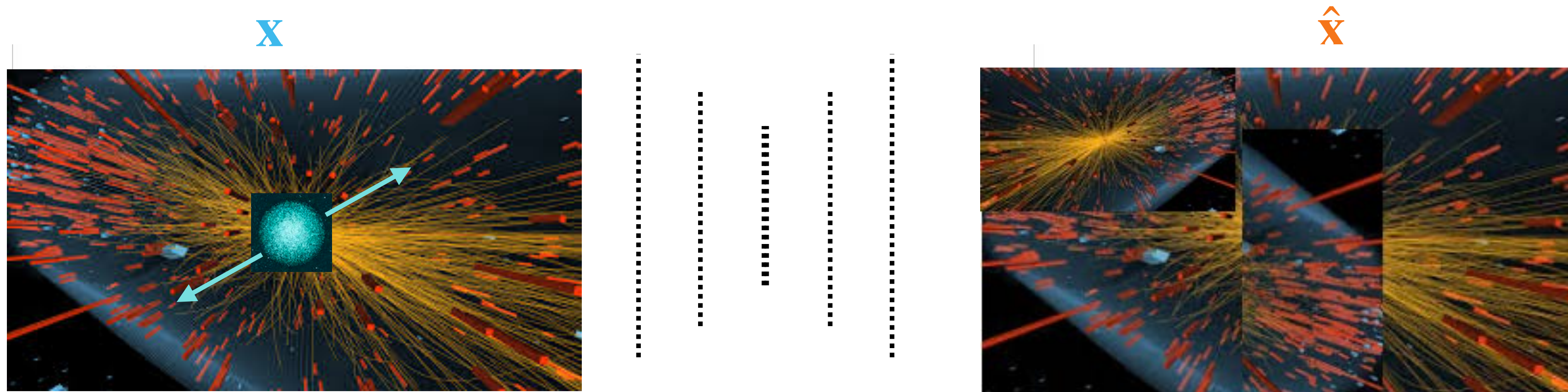clusters

$x_2$

$x_1$

● normal data
● noise
● anomalous data

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

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

LOST DATA
SELECTED DATA
POSSIBLE NP SIGNAL

Everything here is normal

Everything here is abnormal

NP?

Reconstruction error

AD threshold

....in 50 nanoseconds!

Semantic segmentation for autonomous vehicles

**_N. Ghielmetti 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....

Seizure Predicting Brain Implant

4 distributed recording ASICs

Wireless communication hub

**_W. Lemaire et al._**

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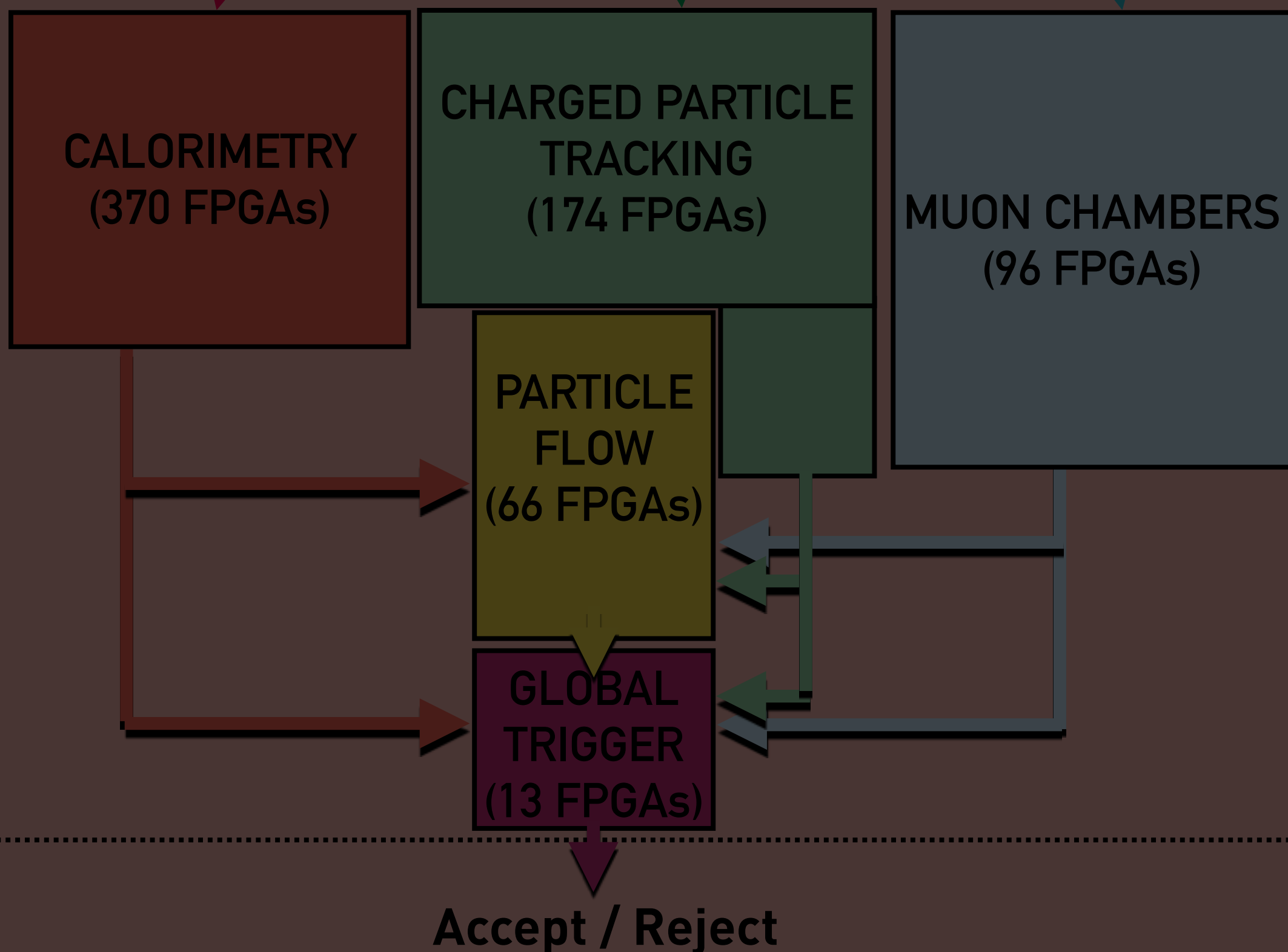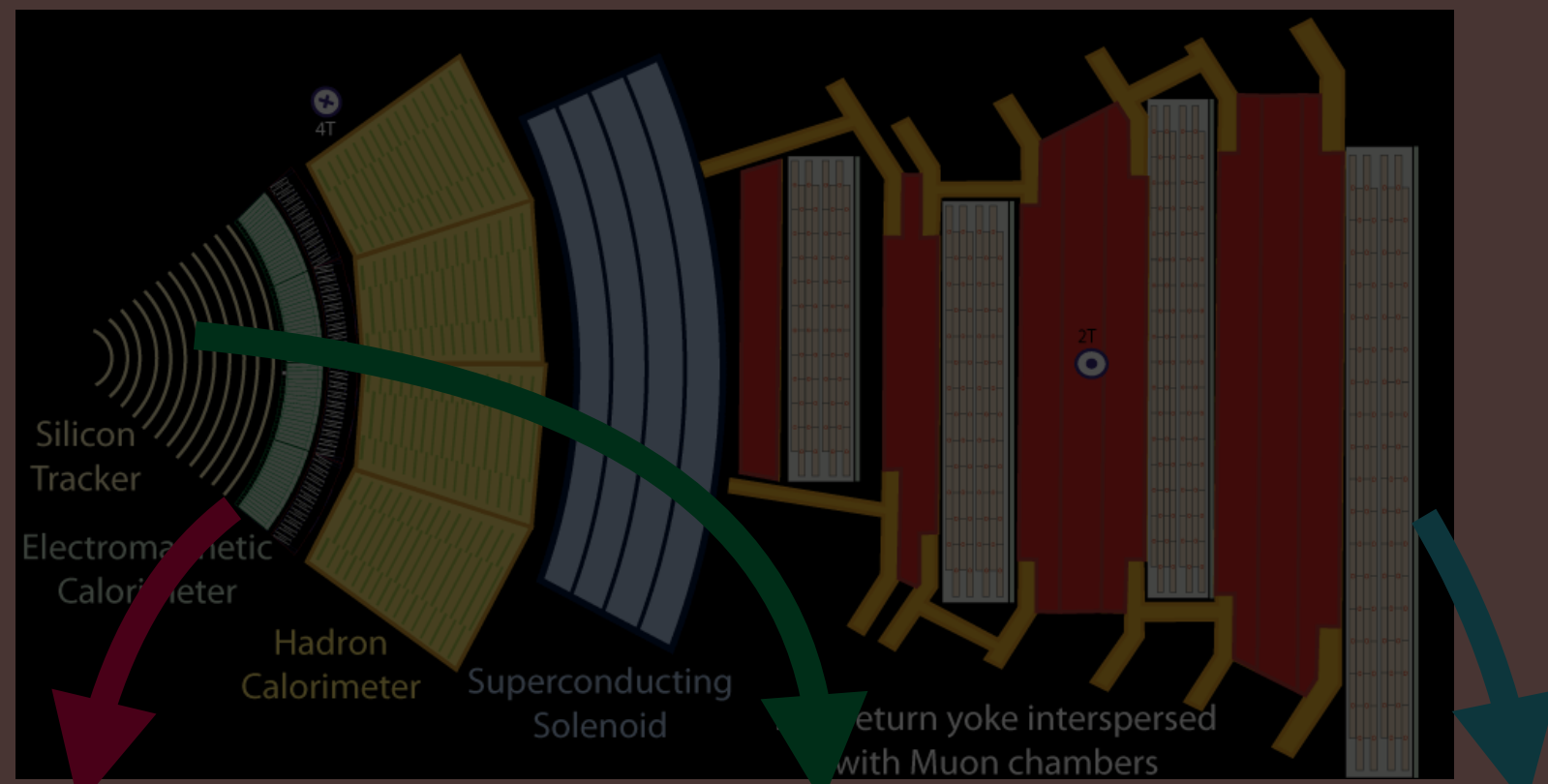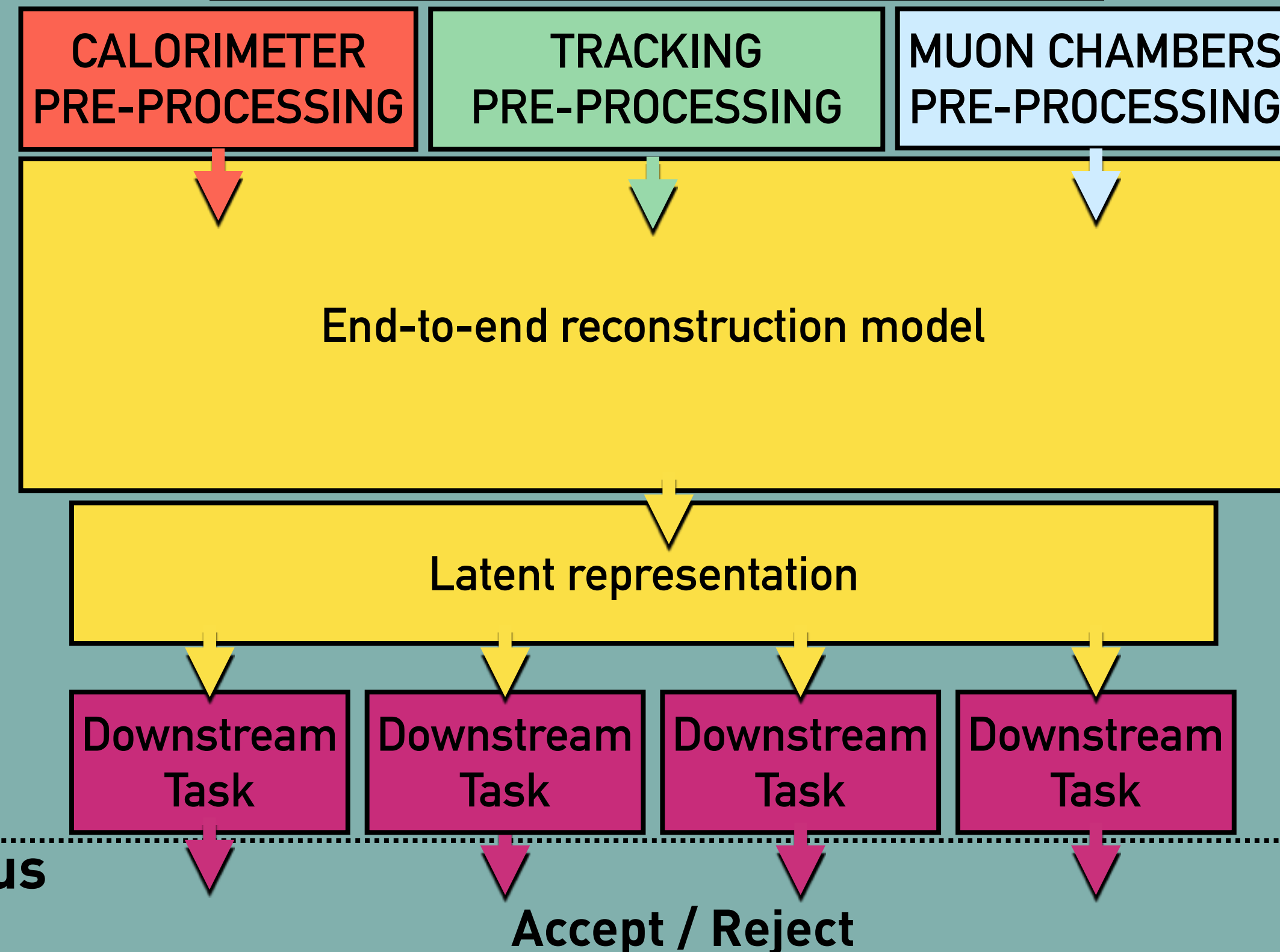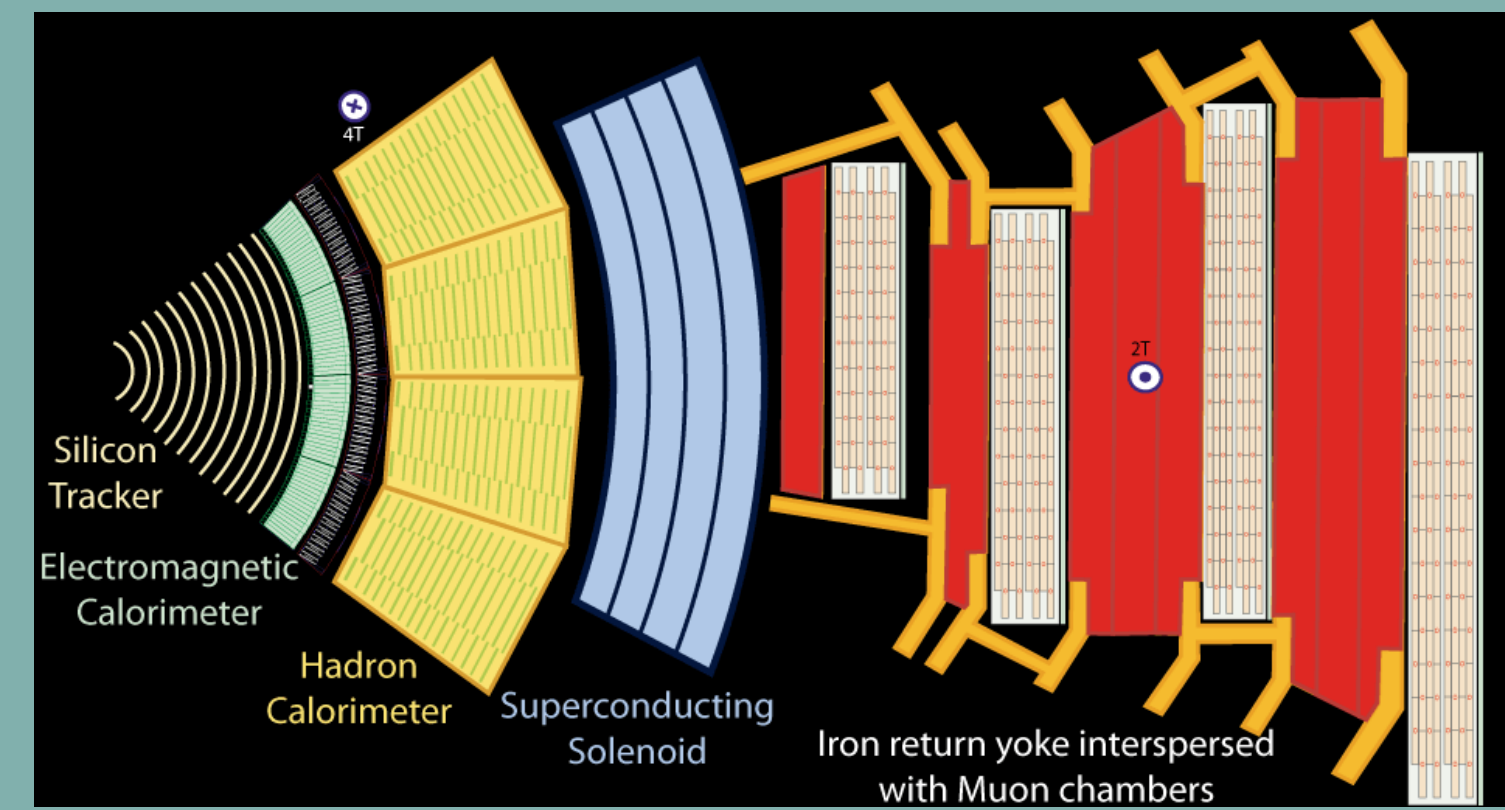- Putting control in cryostat (e.g optimal pulse parameters)

Conventional quantum control algorithm

Our ML model

**_D Xu et al._**

SN interaction + radiological background

**63 Tb/s**

**12.5 μs**

CALORIMETRY
(370 FPGAs)

CHARGED PARTICLE
TRACKING
(174 FPGAs)

MUON CHAMBERS
(96 FPGAs)

PARTICLE
FLOW
(66 FPGAs)

GLOBAL
TRIGGER
(13 FPGAs)

Accept / Reject

**Current HL-LHC design**

CALORIMETER
PRE-PROCESSING

TRACKING
PRE-PROCESSING

MUON CHAMBERS
PRE-PROCESSING

End-to-end reconstruction model

Latent representation

Downstream
Task

Downstream
Task

Downstream
Task

Downstream
Task

Accept / Reject

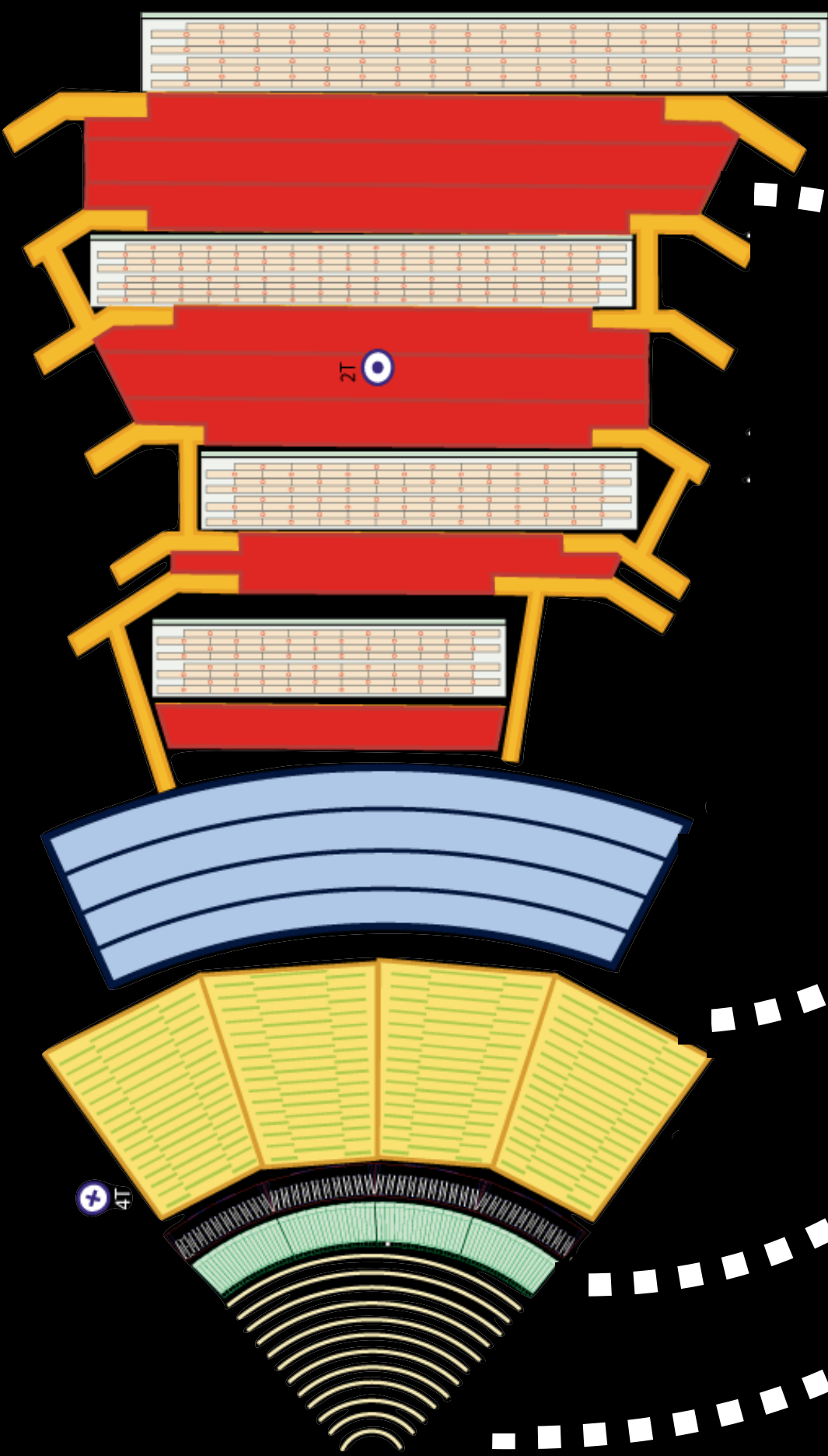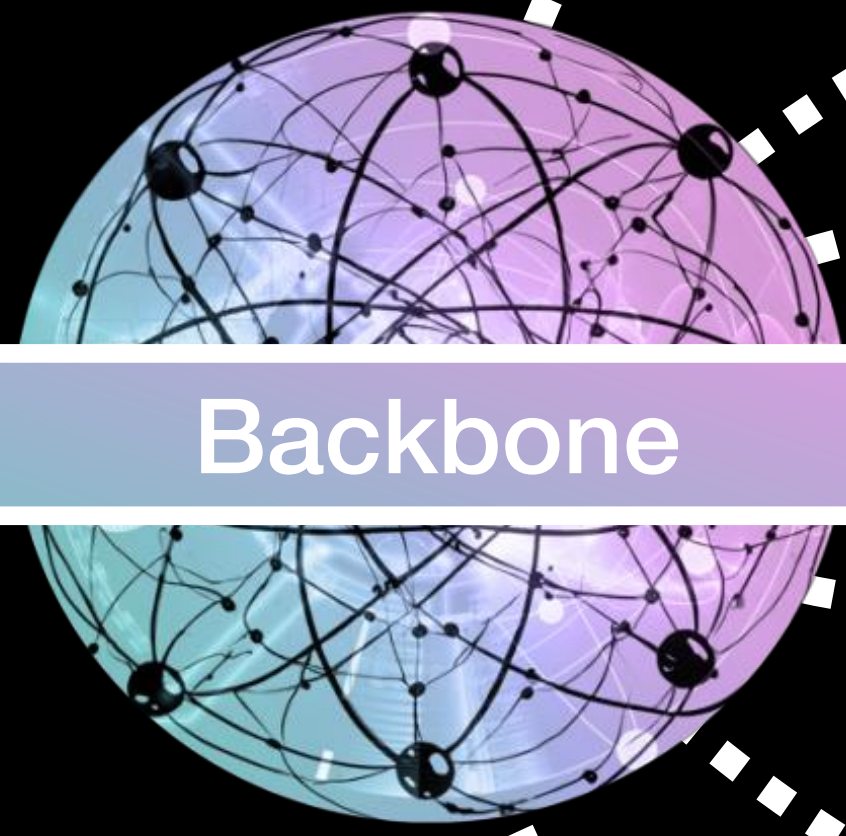**This project**

Heterogeneous detector
Multi-modal input!

Pre-training → Backbone

Fine-tuning → Jet reconstruction

Fine-tuning → Electron on
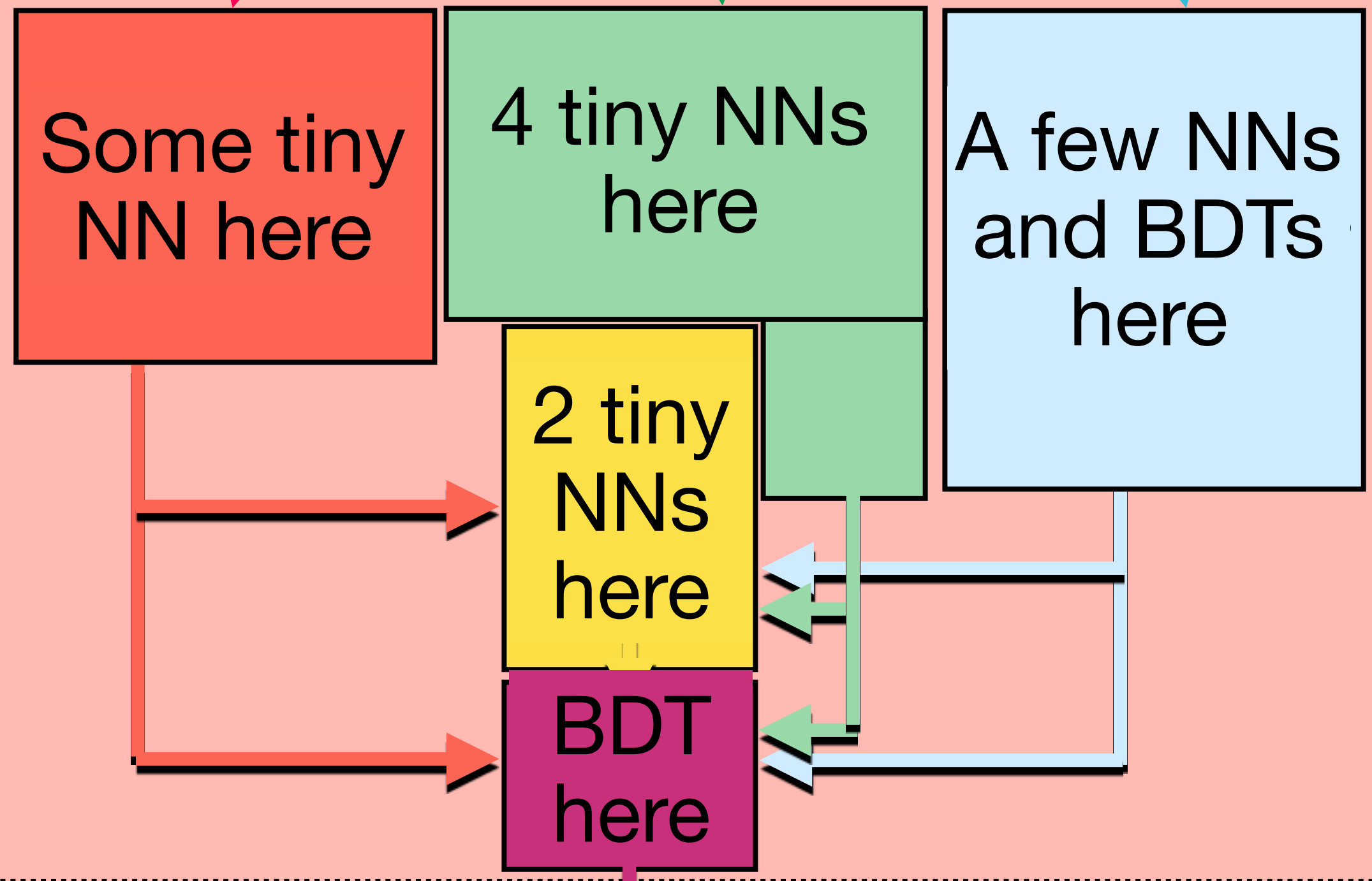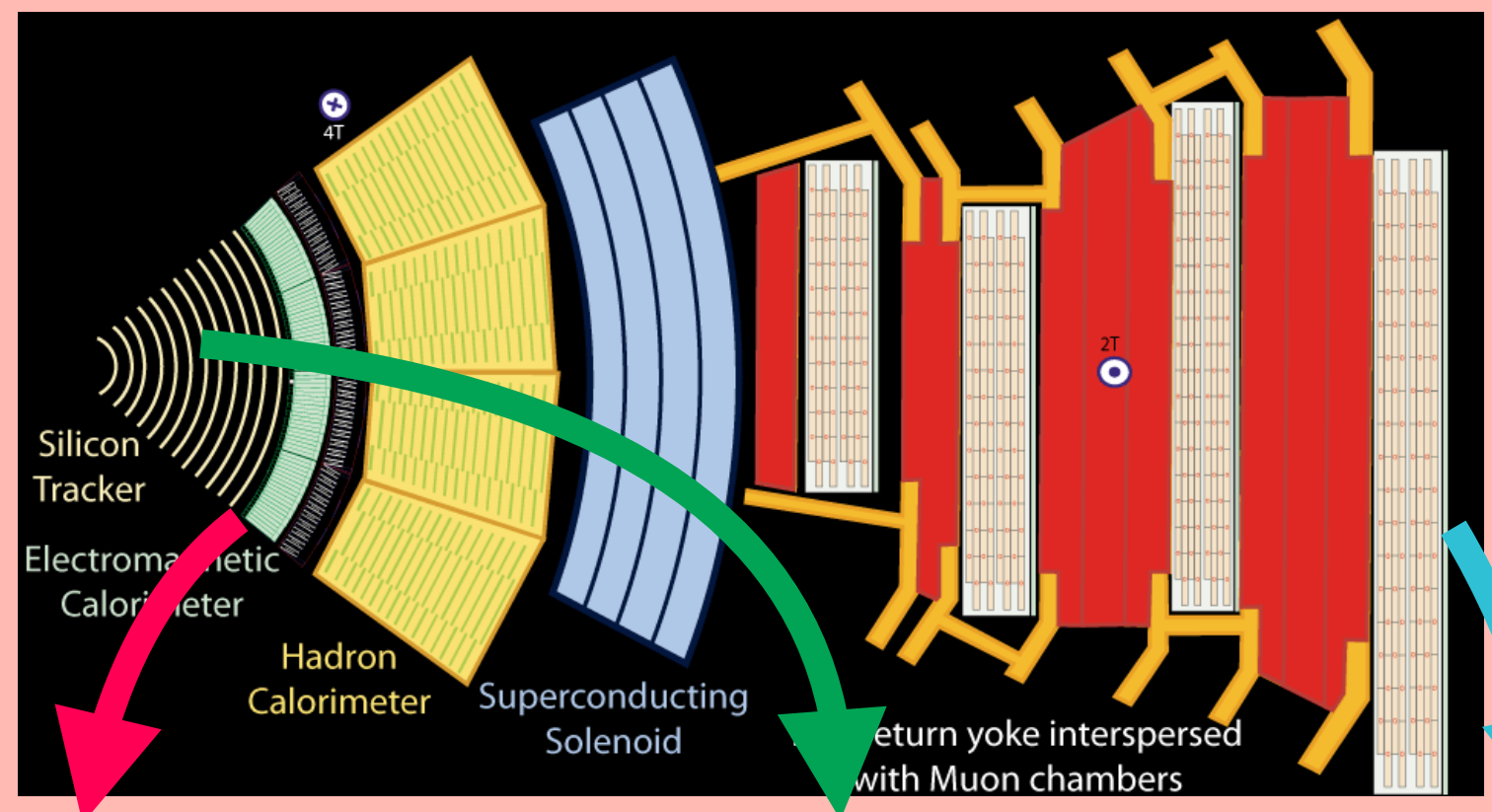
Fine-tuning

Fine-tuning → Missing energy computation

Fine-tuning → Anomaly Detection

Fine-tuning → 0/1?
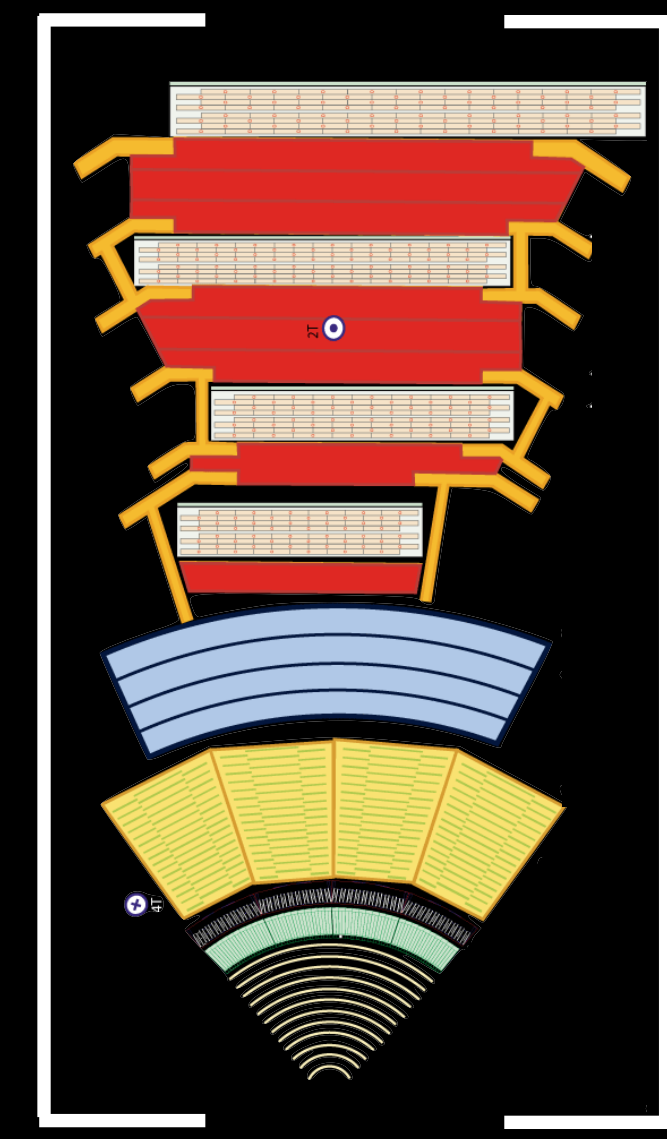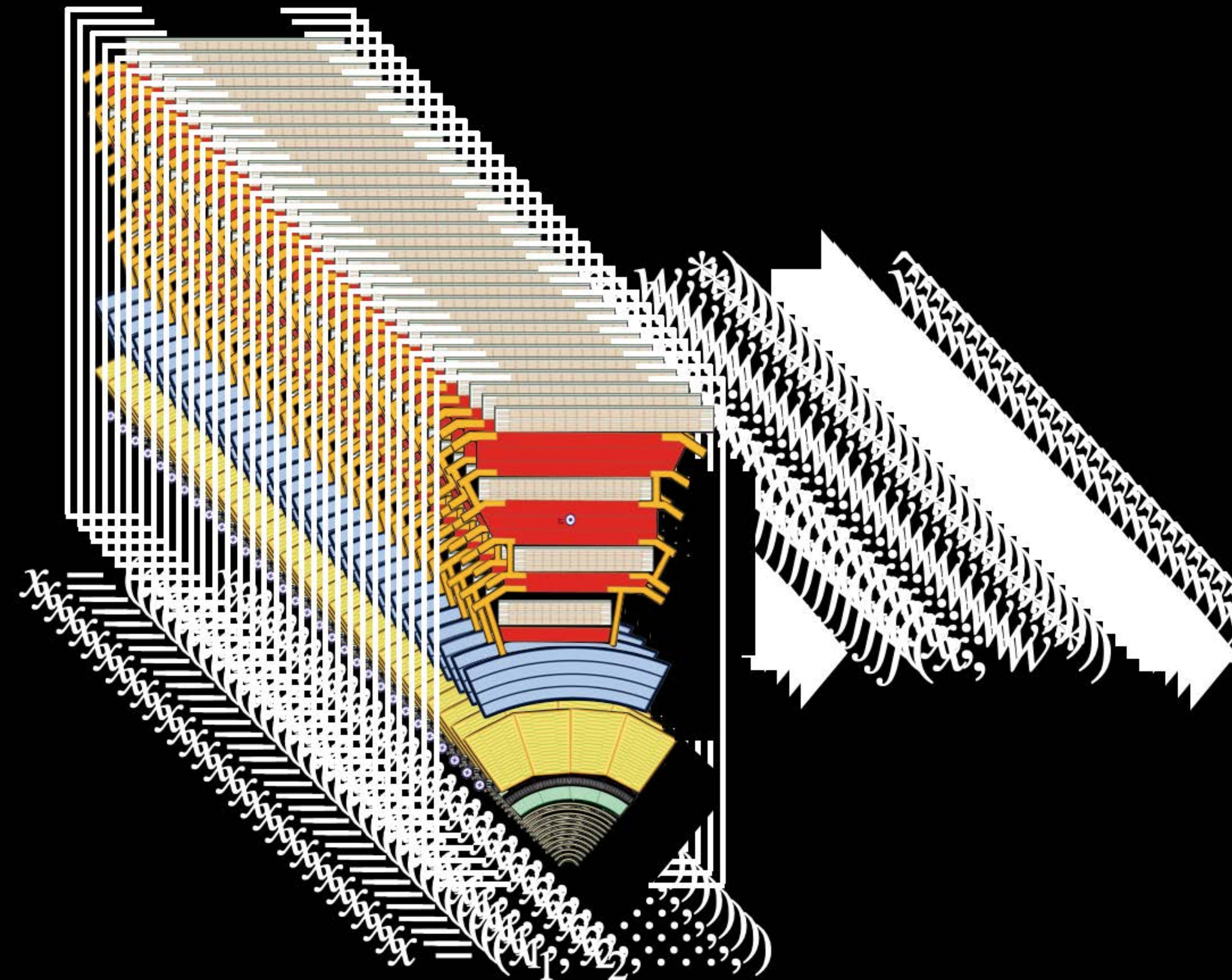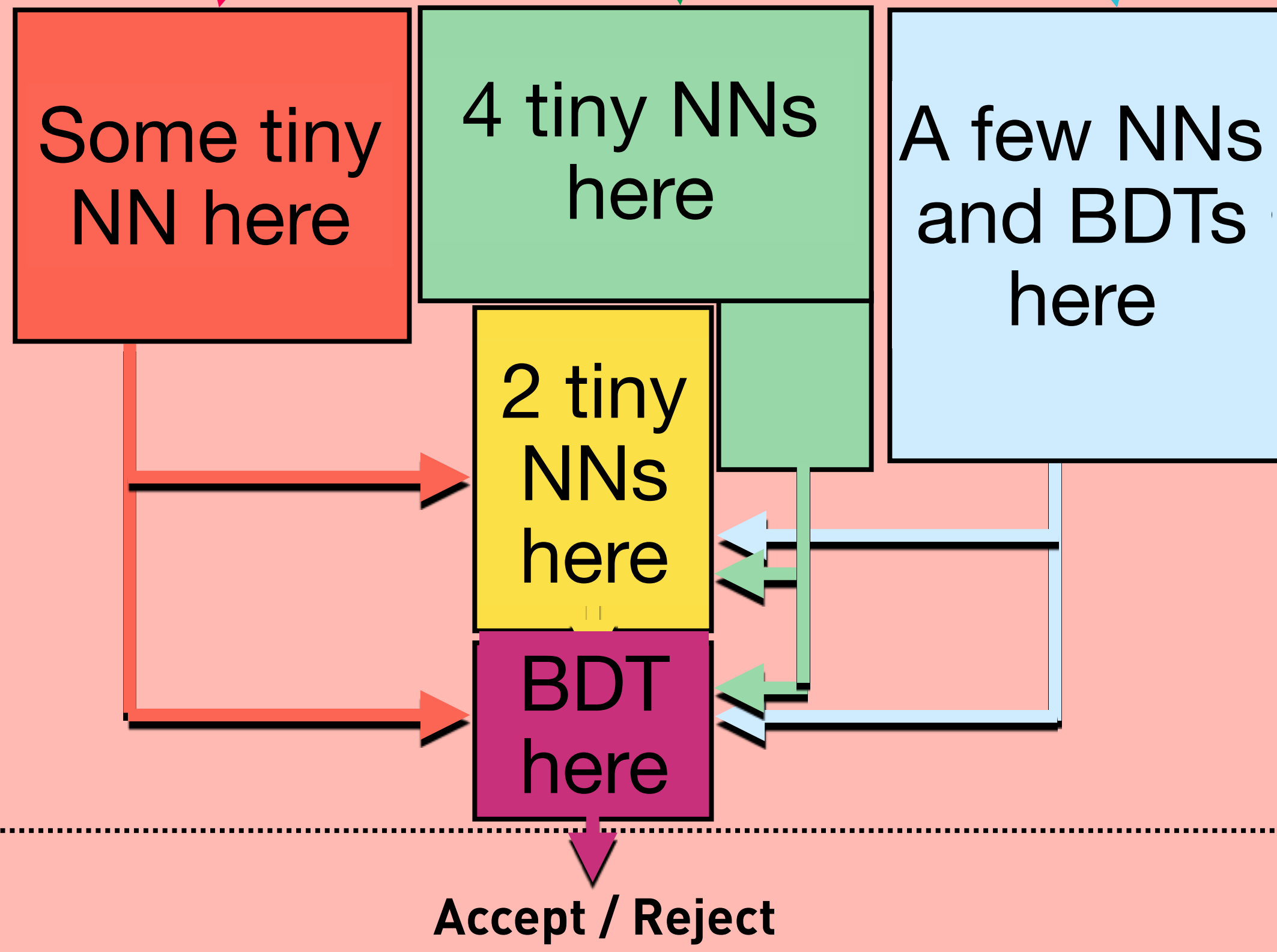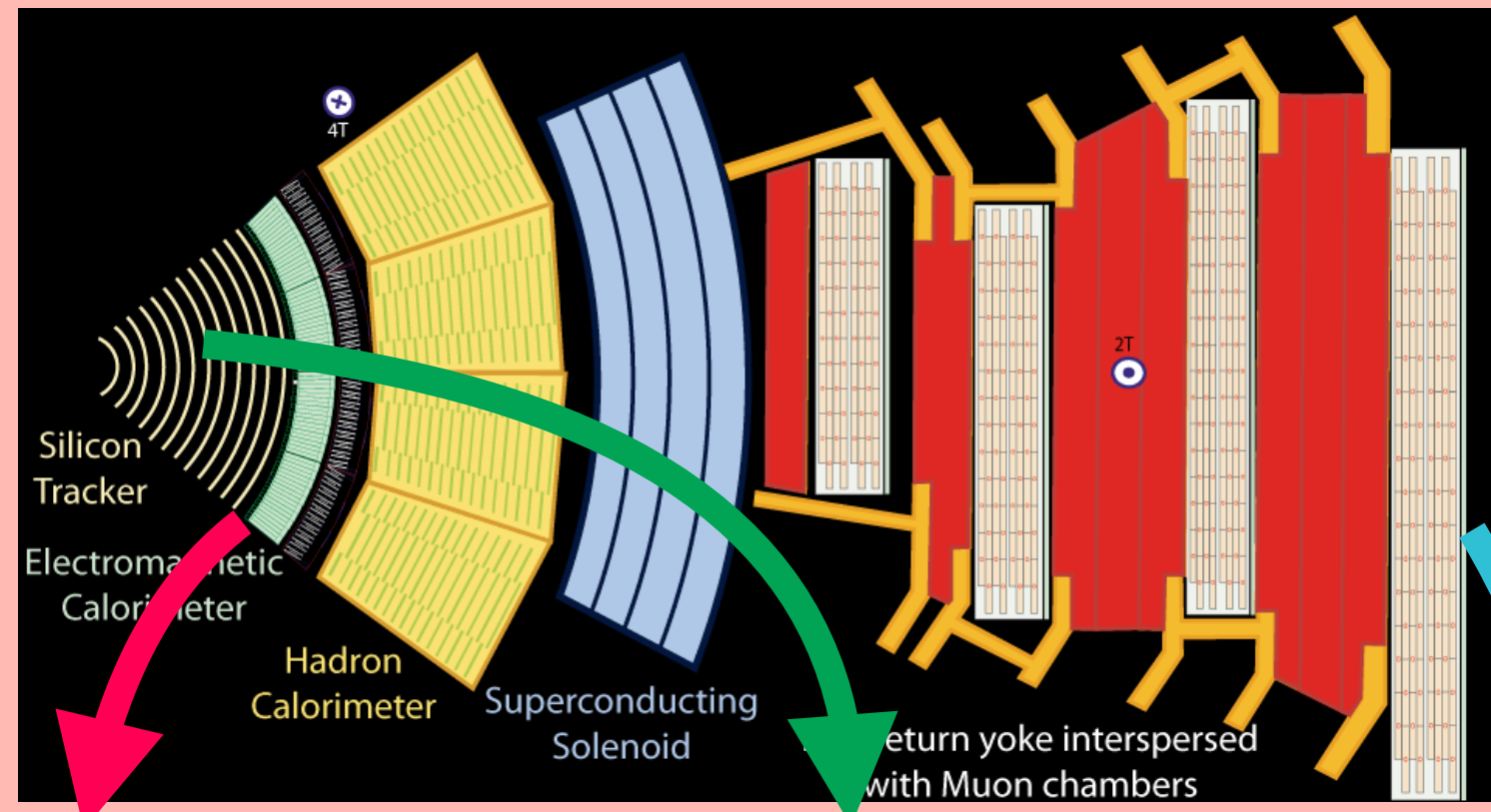
( Generate simulation? )

Some tiny NN here

4 tiny NNs here

A few NNs and BDTs here

2 tiny NNs here

BDT here

Accept / Reject

$x = (x_1, x_2, \ldots, )$

$f(x; w^*)$ → $\hat{y}$

Some tiny NN here

4 tiny NNs here

A few NNs and BDTs here

2 tiny NNs here

BDT here

Accept / Reject

Too many models, too little learning?

# One model, learn neural embedding?



CALORIMETER PRE-PROCESSING

TRACKING PRE-PROCESSING

MUON CHAMBERS PRE-PROCESSING

End-to-end reconstruction model

Latent representation

Downstream Task

Downstream Task

Downstream Task

Downstream Task

μs

Accept / Reject

$x = (x_1, x_2, \ldots, )$

Some new space

Downstream Task

Downstream Task

Downstream Task

# One model, learn neural embedding?

**CALORIMETER PRE-PROCESSING**

**TRACKING PRE-PROCESSING**

**MUON CHAMBERS PRE-PROCESSING**

End-to-end reconstruction model

Latent representation

Downstream Task

Downstream Task

Downstream Task

Downstream Task

μs

**Accept / Reject**

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

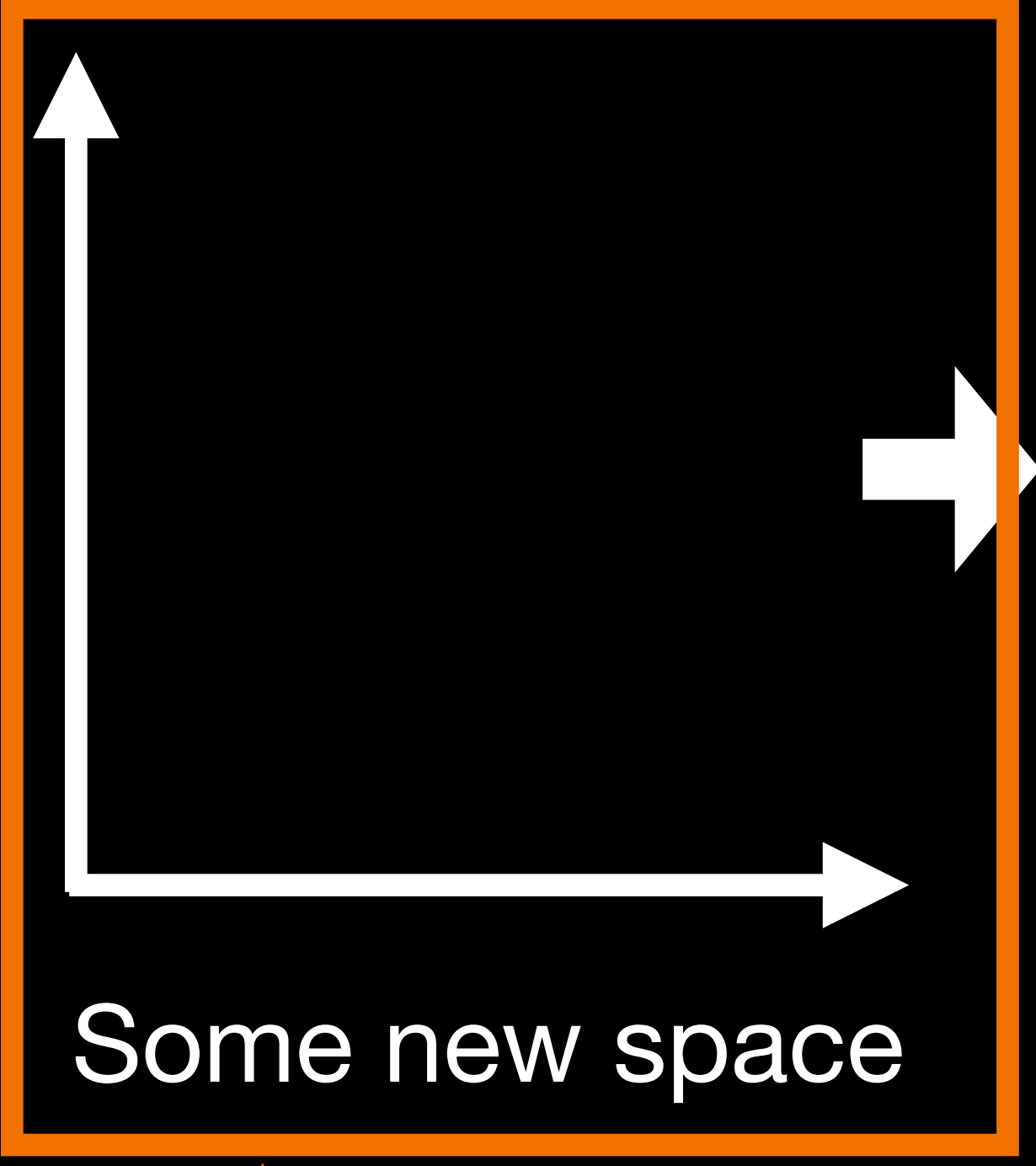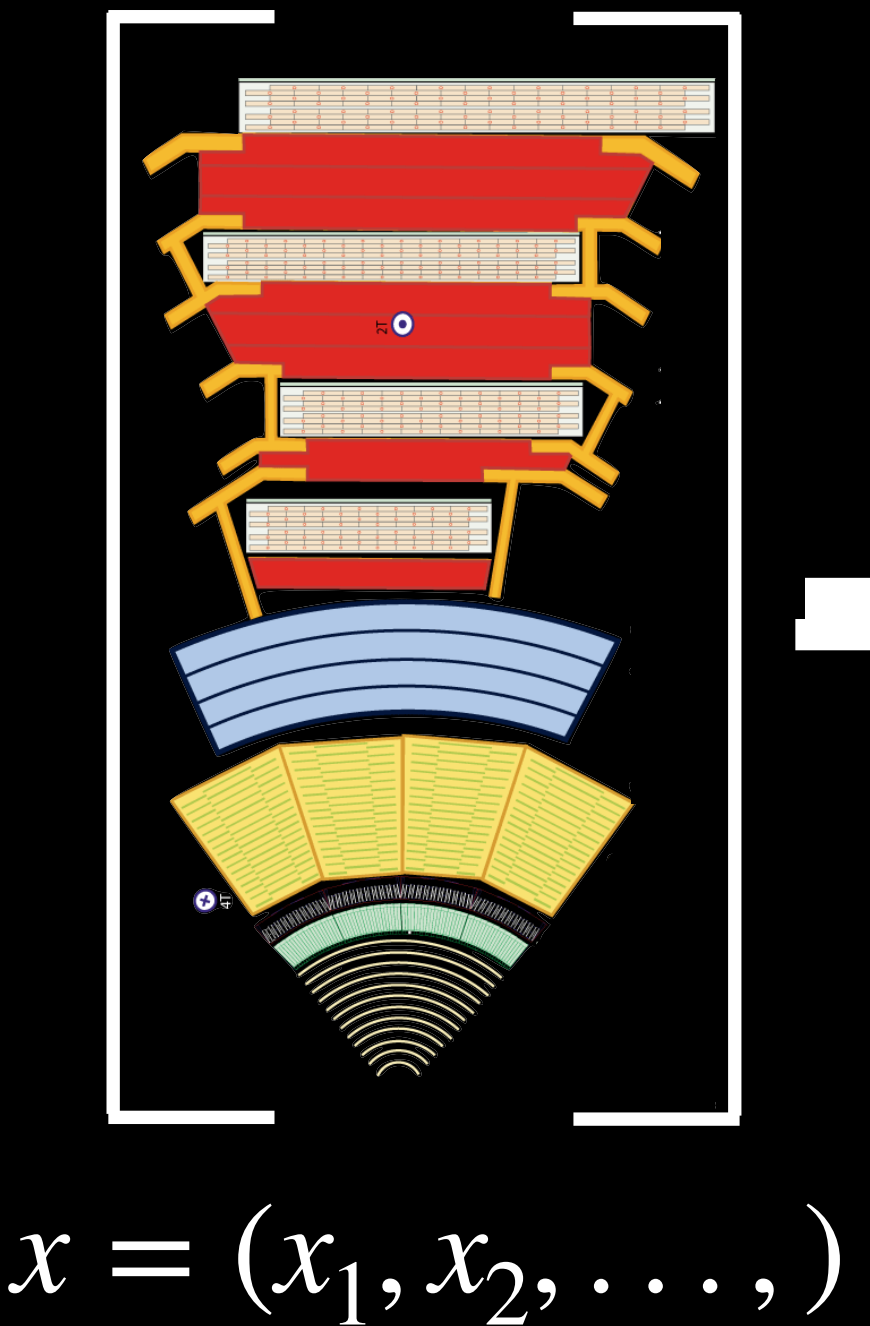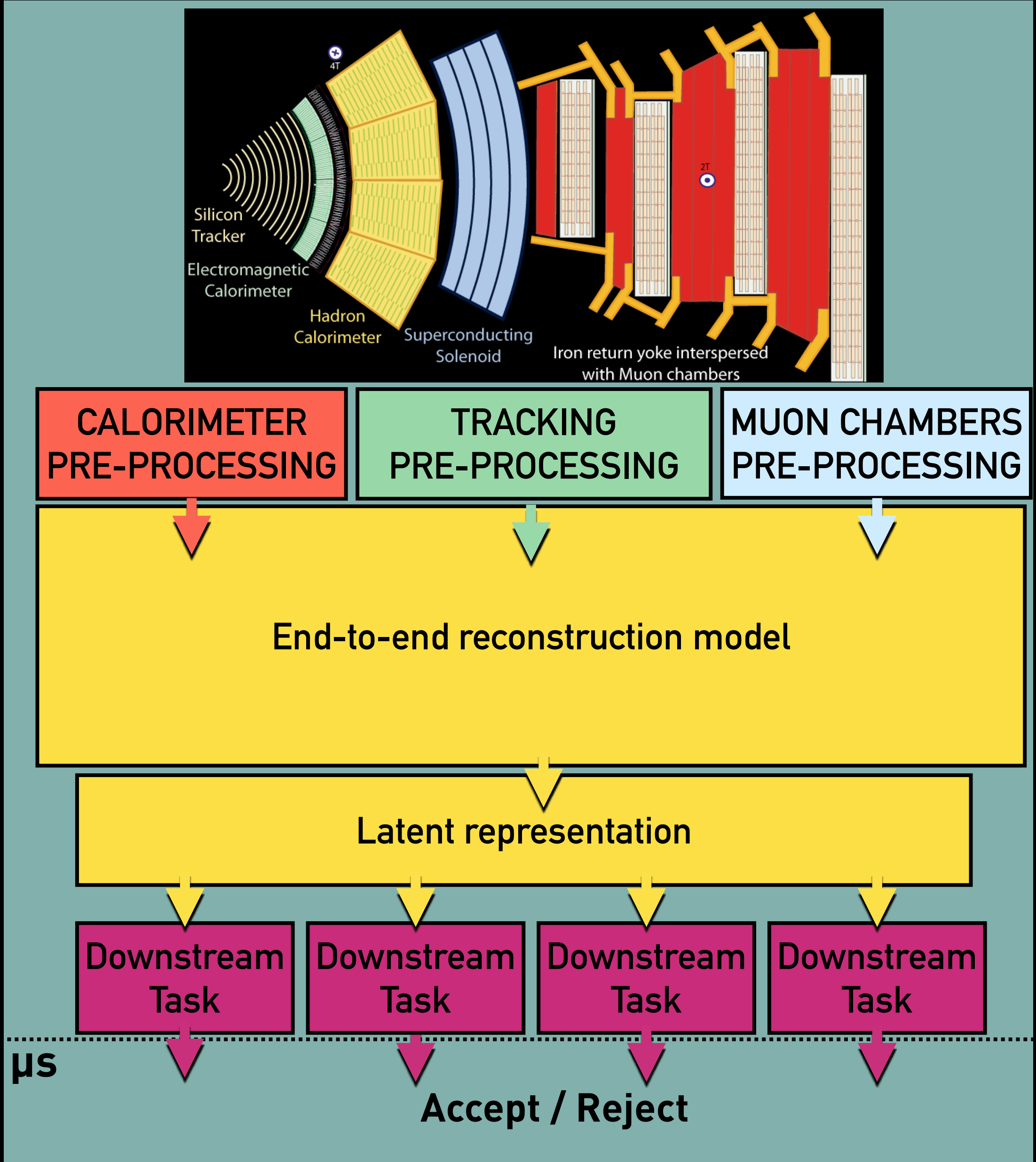Iron return yoke interspersed with Muon chambers

$x = (x_1, x_2, \ldots, )$
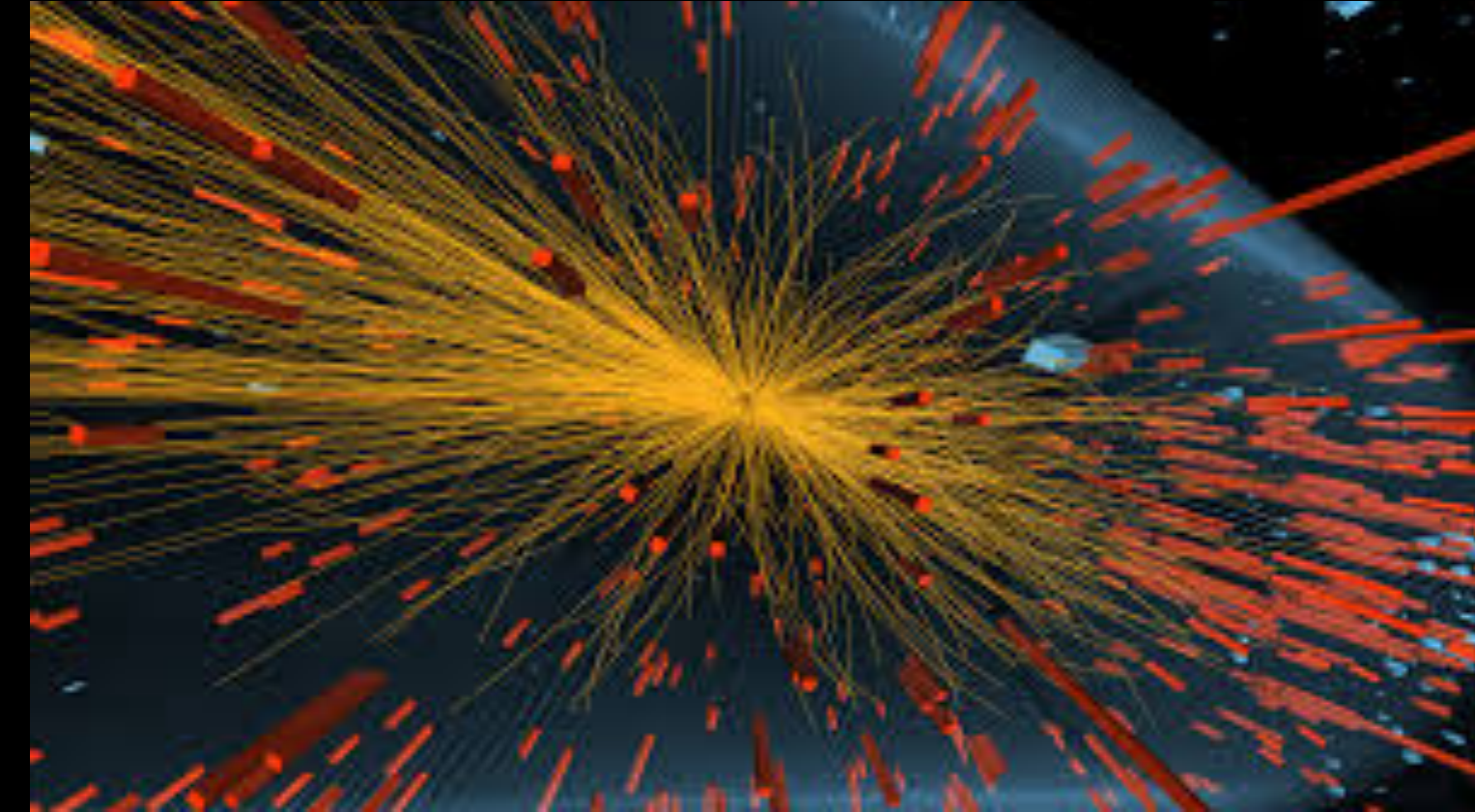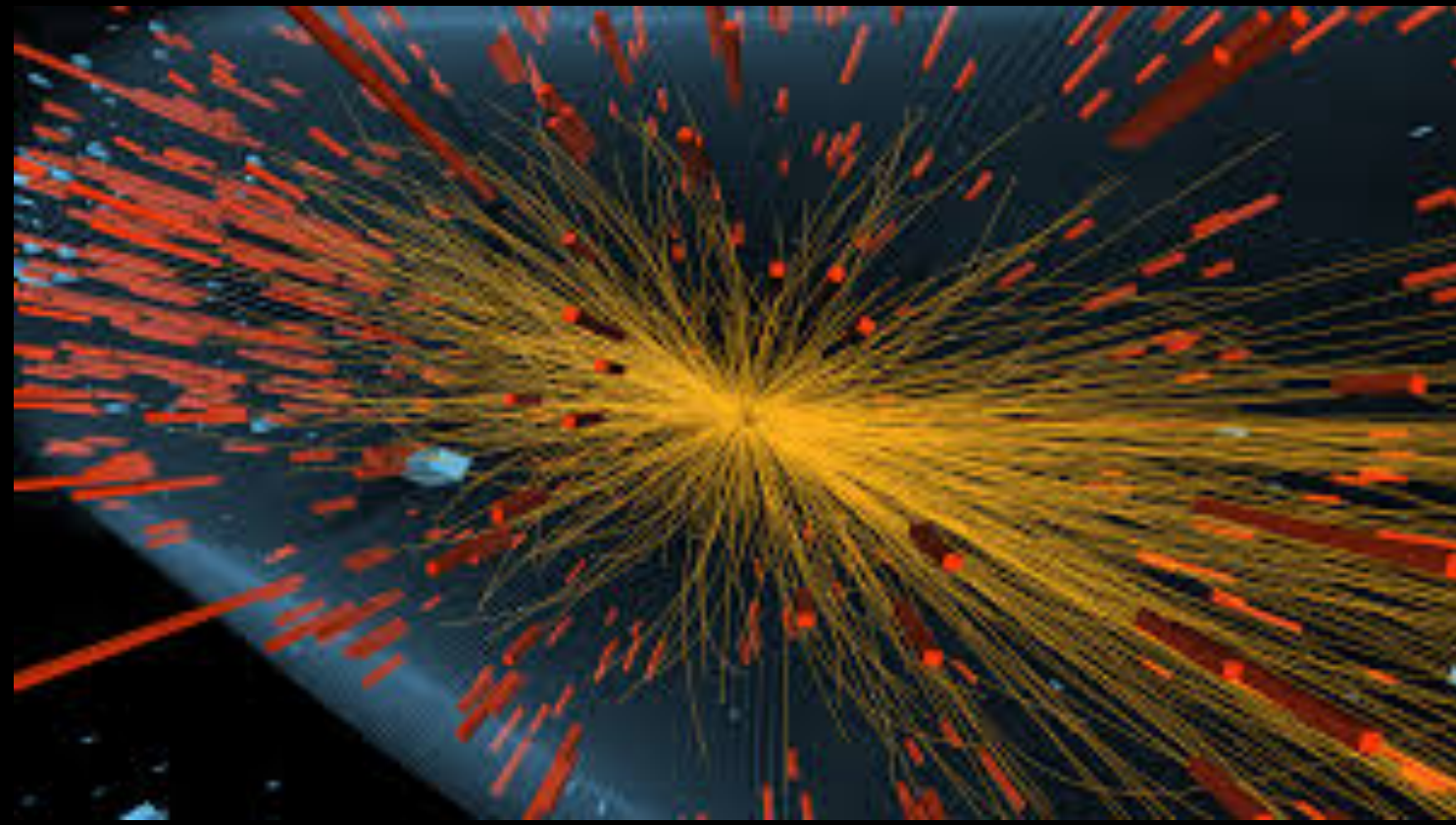
Some new space

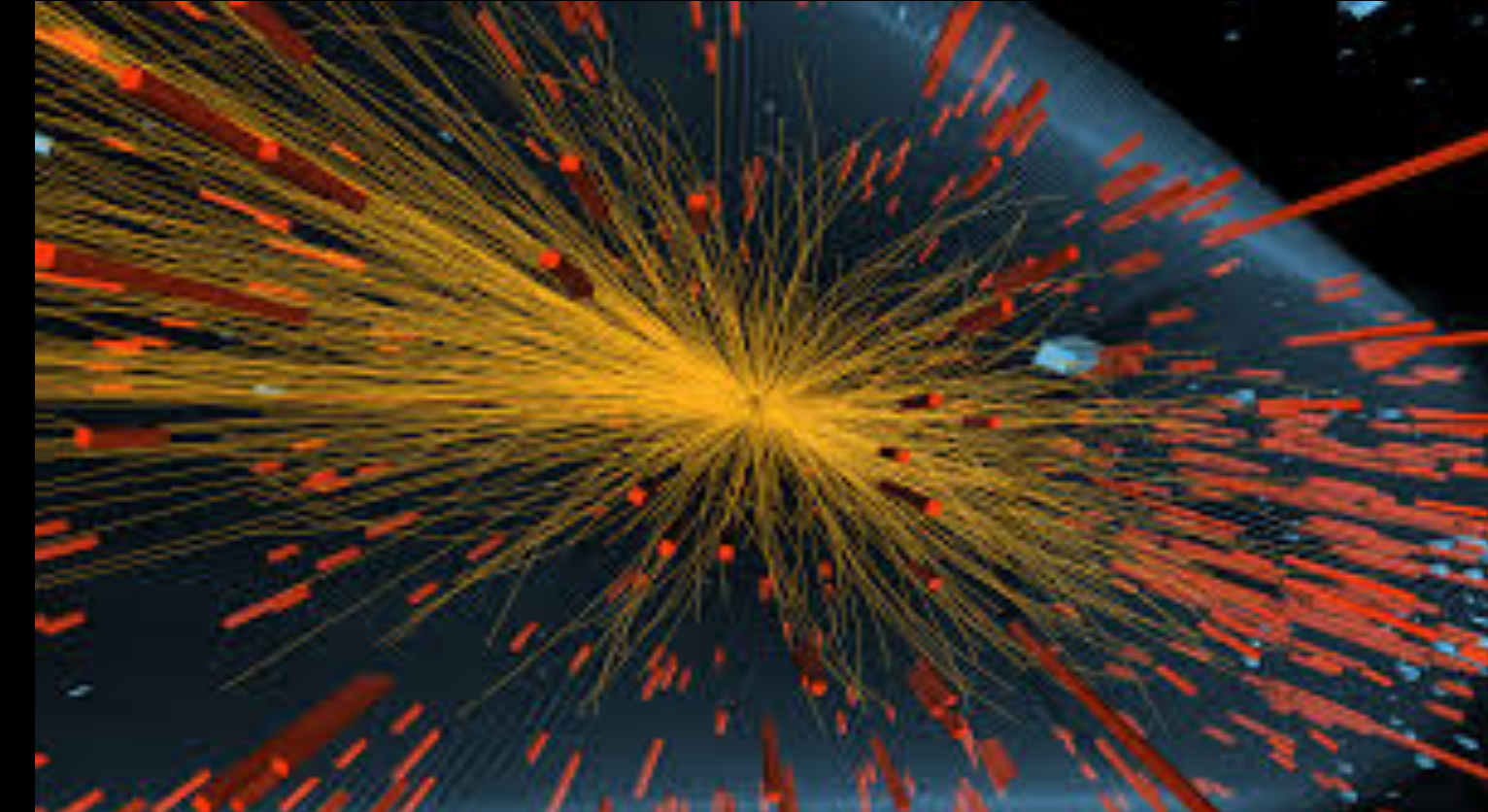Downstream Task

Downstream Task
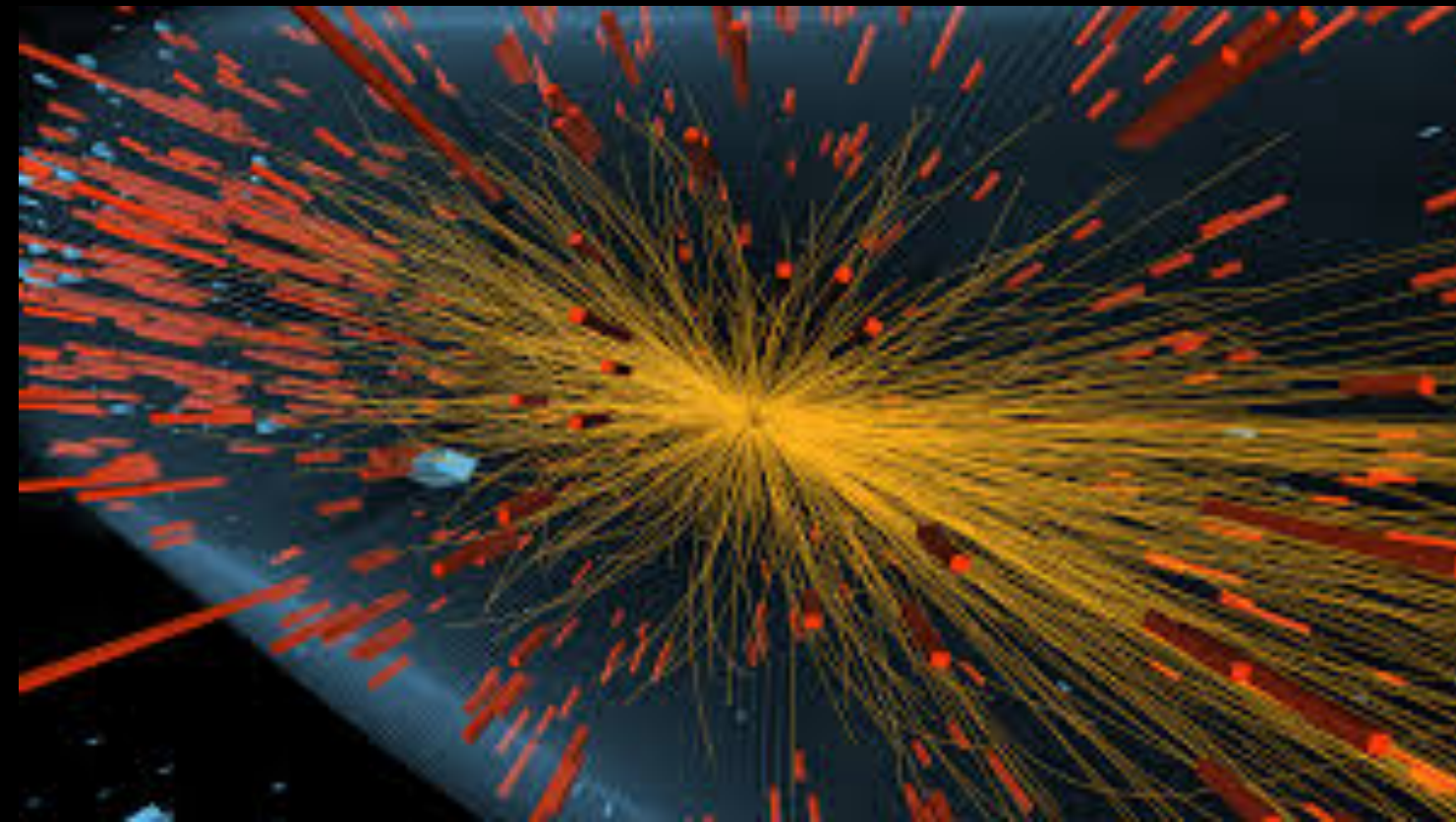
Downstream Task

**Let's build this space!**

# Learning the space
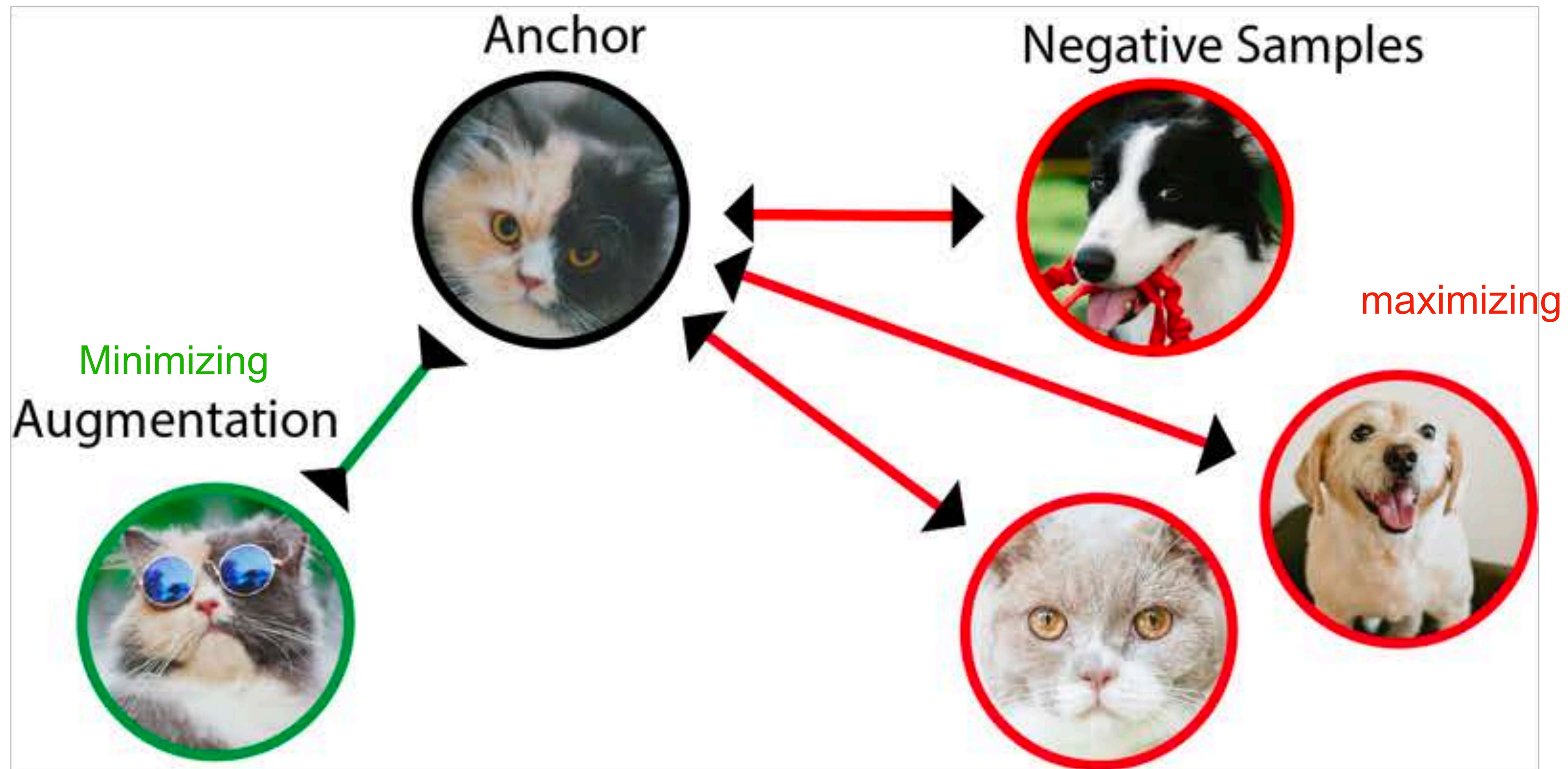
# Learning the space

- By looking at data, we can learn a lot
  - Go over input  piece by piece
  - Analyze every aspect
  - Compare every feature
- Find distinctive style of the input
  - can be done e.g by looking for a deviation

# Physically motivated augmentations?



- Minimizing and maximizing distances learns a space

Augmented Cat A

Dog B

Cat A

Dog A

Cat B

Augmented Dog A

# Physically motivated augmentations?



No class labels used in training! How do we augment detector data?

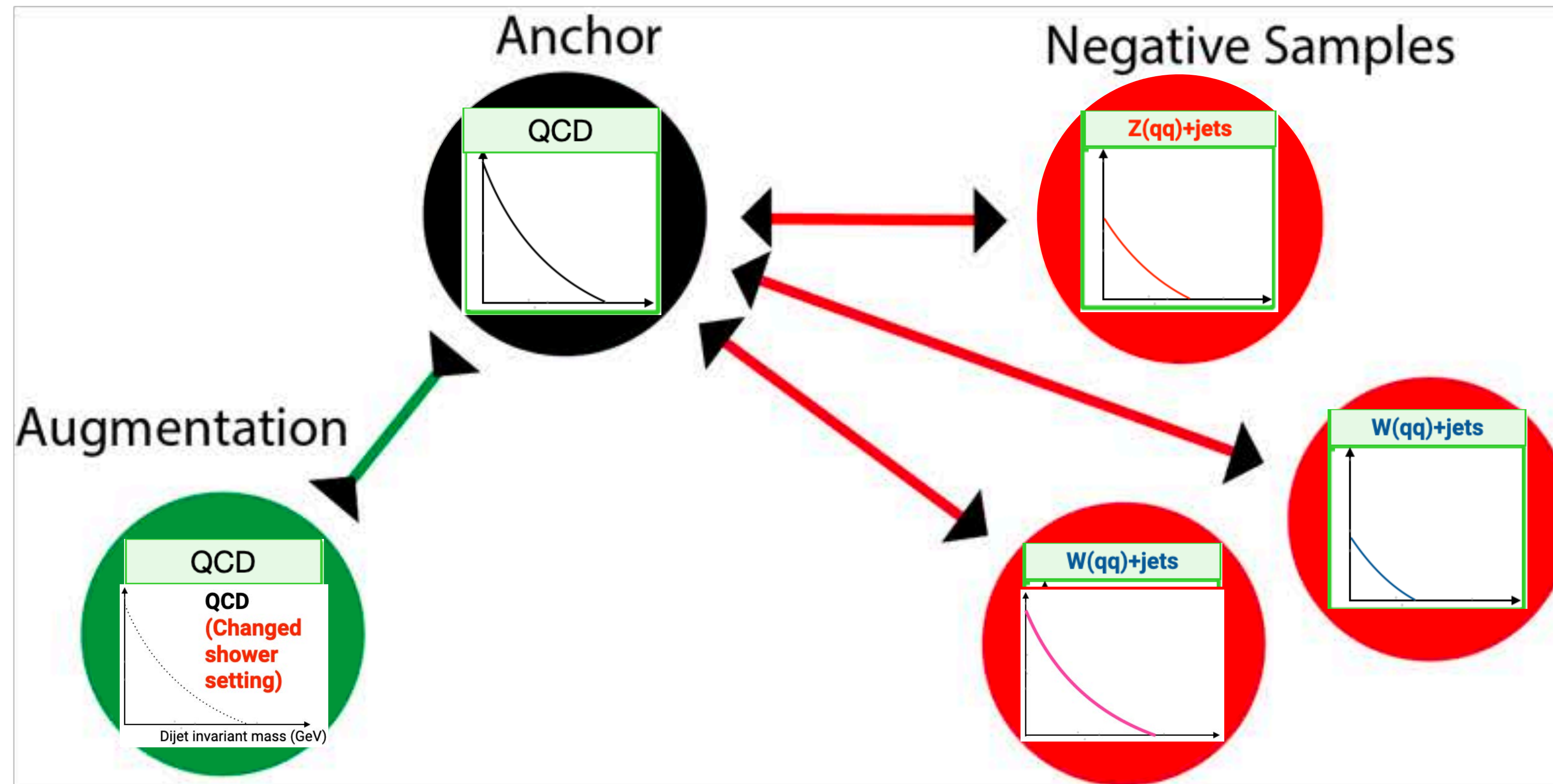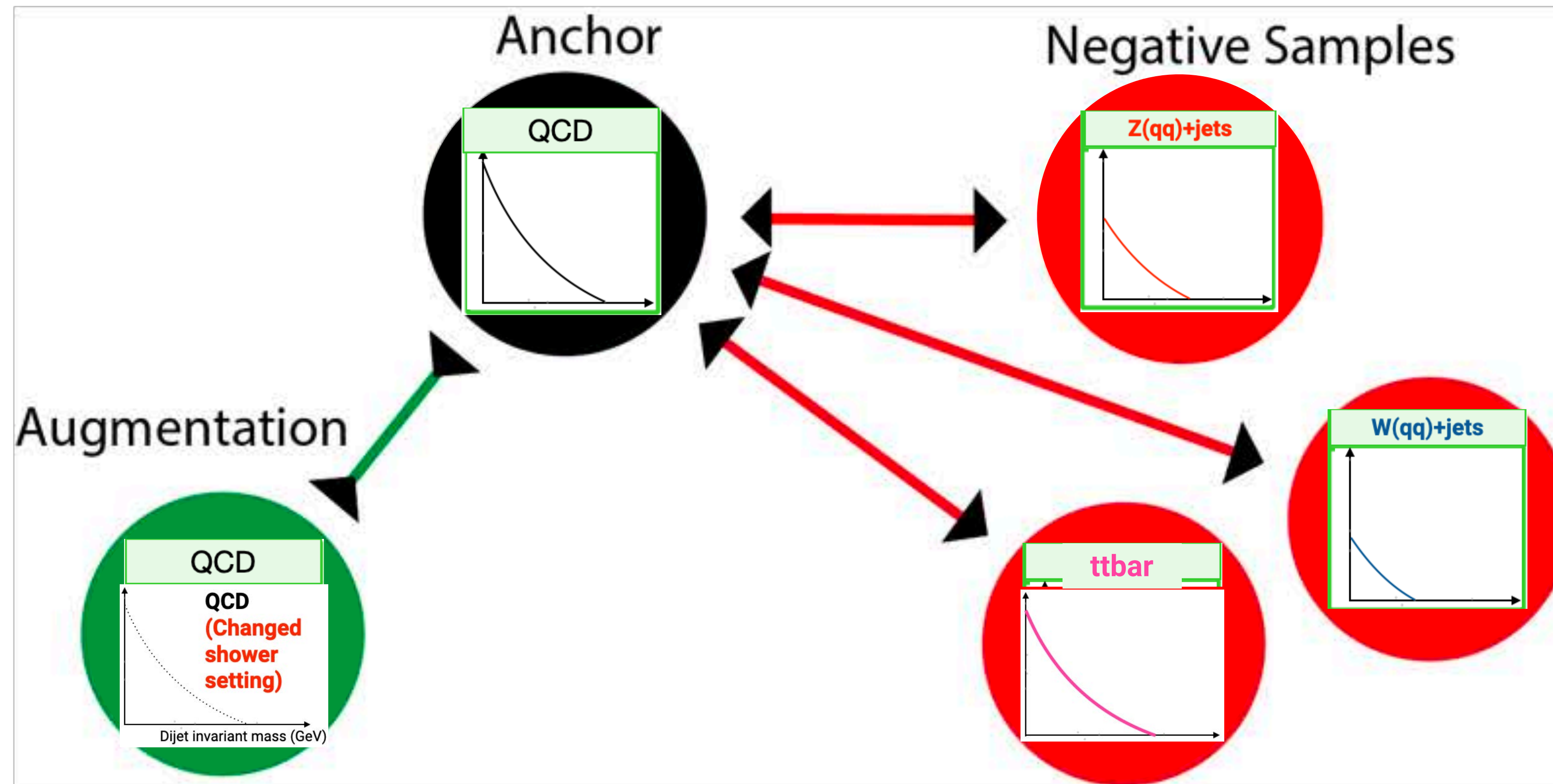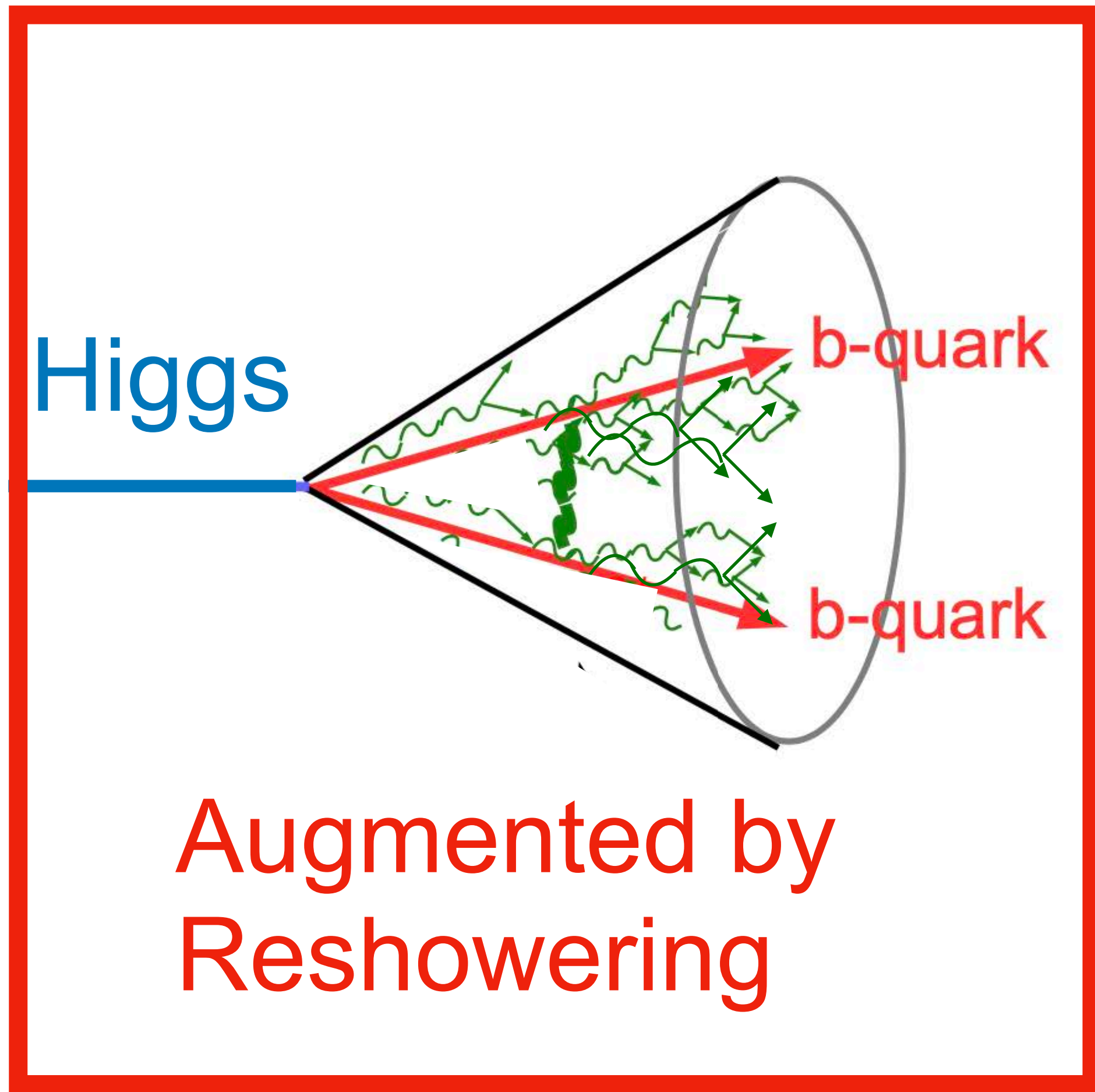# Physically motivated augmentations?



No class labels used in training! How do we augment detector data?

Augmentation

Higgs

Baseline

Higgs

Augmented by Reshowering

Embedded Space can use any NN to embed

# QM foundation models



→ embedding quantum mechanics into AI algorithm

$x = (x_1, x_2, \ldots, )$

Some new space

Downstream Task $\quad \hat{y}?$

Downstream Task $\quad \hat{y}?$

Downstream Task $\quad \hat{y}?$

**Training 1: Learn neural embedding (on a lot of data, for a long time) On simulation? On data?**

$$x = (x_1, x_2, \ldots, )$$

Some new space

$\hat{y}?$

$\hat{y}?$

$\hat{y}?$

**Training 2: Fine tune for specific task (fast, small dataset, simulation)**

# Foundation model of the Level-1 trigger



63 Tb/s

CALORIMETER PRE-PROCESSING

TRACKING PRE-PROCESSING

MUON CHAMBERS PRE-PROCESSING

Foundation model

Latent

Downstream Task

Downstream Task

Downstream Task

Downstream Task

μs

Accept / Reject

# Do I really think this will be possible?



**Foundation model**

**Latent**

**Foundation model**

**Latent**

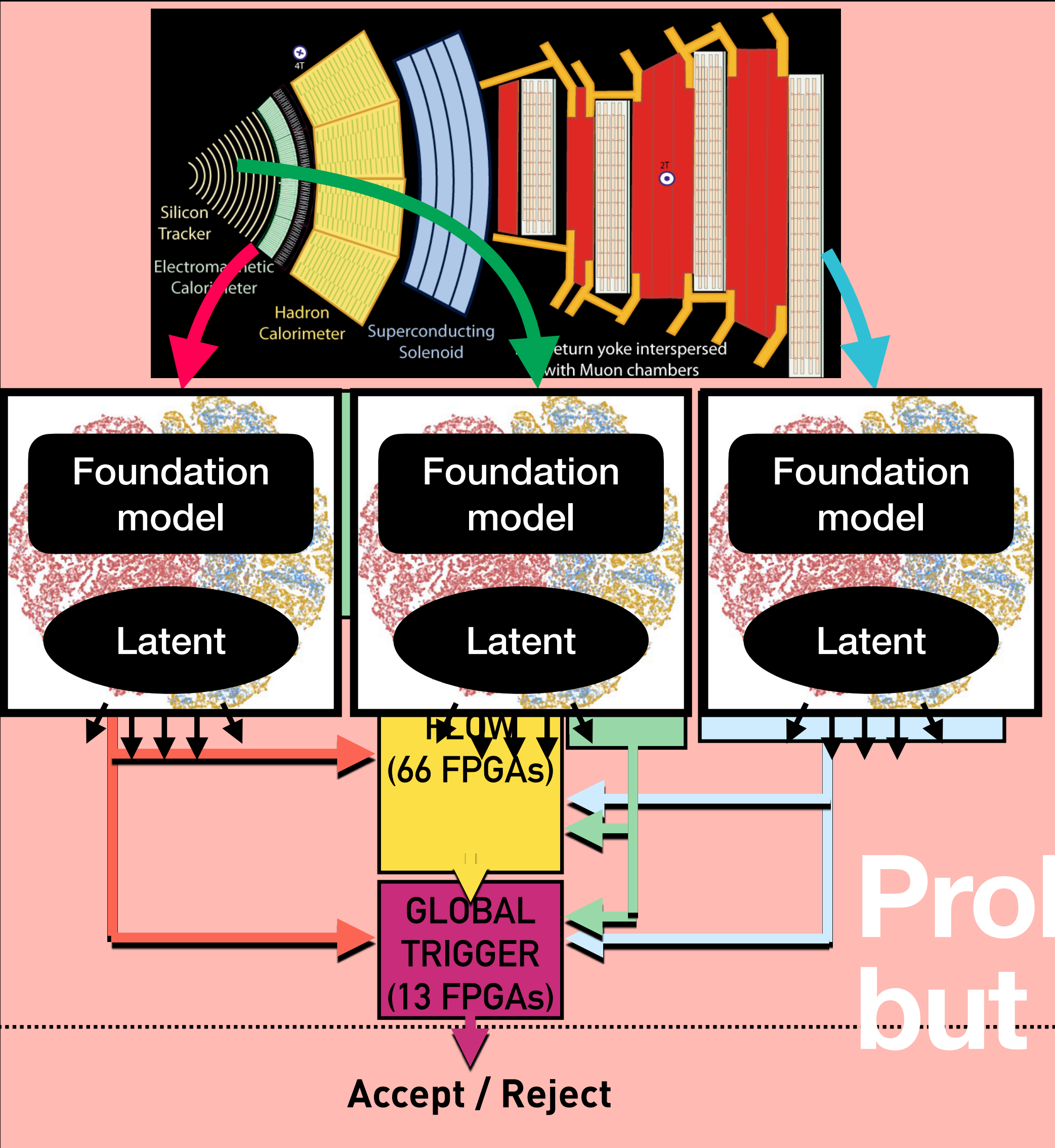**Foundation model**

**Latent**

(66 FPGAs)

GLOBAL TRIGGER
(13 FPGAs)

Accept / Reject

**Probably not,
but at some scale**

Silicon Tracker

Electromagnetic Calorimeter

Hadron Calorimeter

Superconducting Solenoid

Key:
Muon — Electron
Neutral Hadron (e.g. Neutron)

# Careful software-hardware co-design

**O(1M) parameter model on 1000 FPGAs and do inference in O(1)μs?**



Layer 1 (FPGA 1)

Layer 1 (FPGA 2)

Layer 1 (FPGA 3)

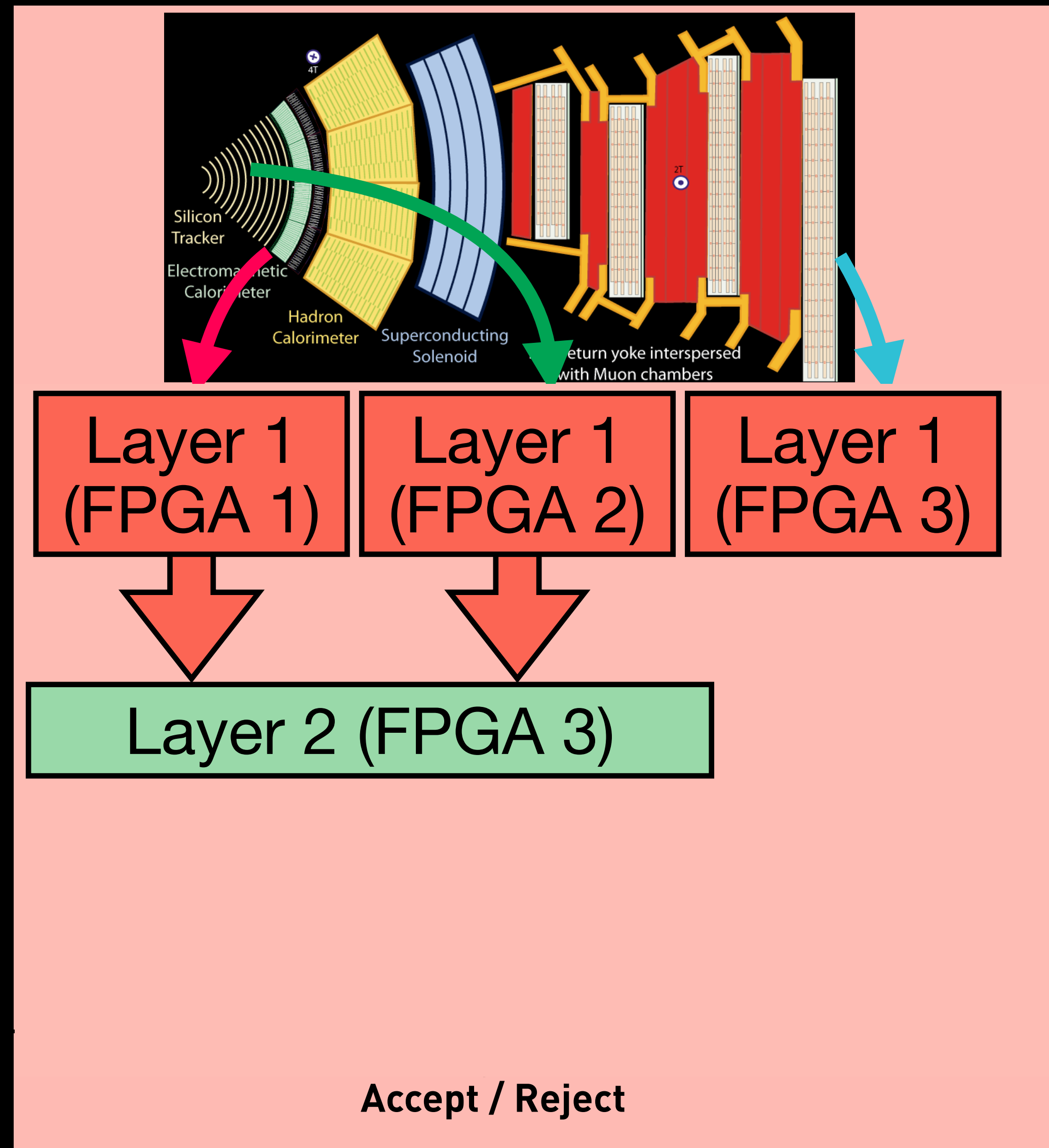Layer 2 (FPGA 3)

**Accept / Reject**

**Similar for GPT-4, layers carefully map onto hardware**

<u>Algean</u>

# Careful software-hardware co-design

**Designed our own protocol to make boards talk to each other fast enough**

**(25 Gbs to transfer data LHC-synchronously between boards)**



Layer 1 (FPGA 1)

Layer 1 (FPGA 2)

Layer 1 (FPGA 3)

Layer 2 (FPGA 3)

**Accept / Reject**

# Masked language modelling



**Next-token-prediction**

The model is given a sequence of words with the goal of predicting the next word.

Example:
Hannah is a ____

Hannah is a *sister*
Hannah is a *friend*
Hannah is a *marketer*
Hannah is a *comedian*

**Masked-language-modeling**

The model is given a sequence of words with the goal of predicting a 'masked' word in the middle.
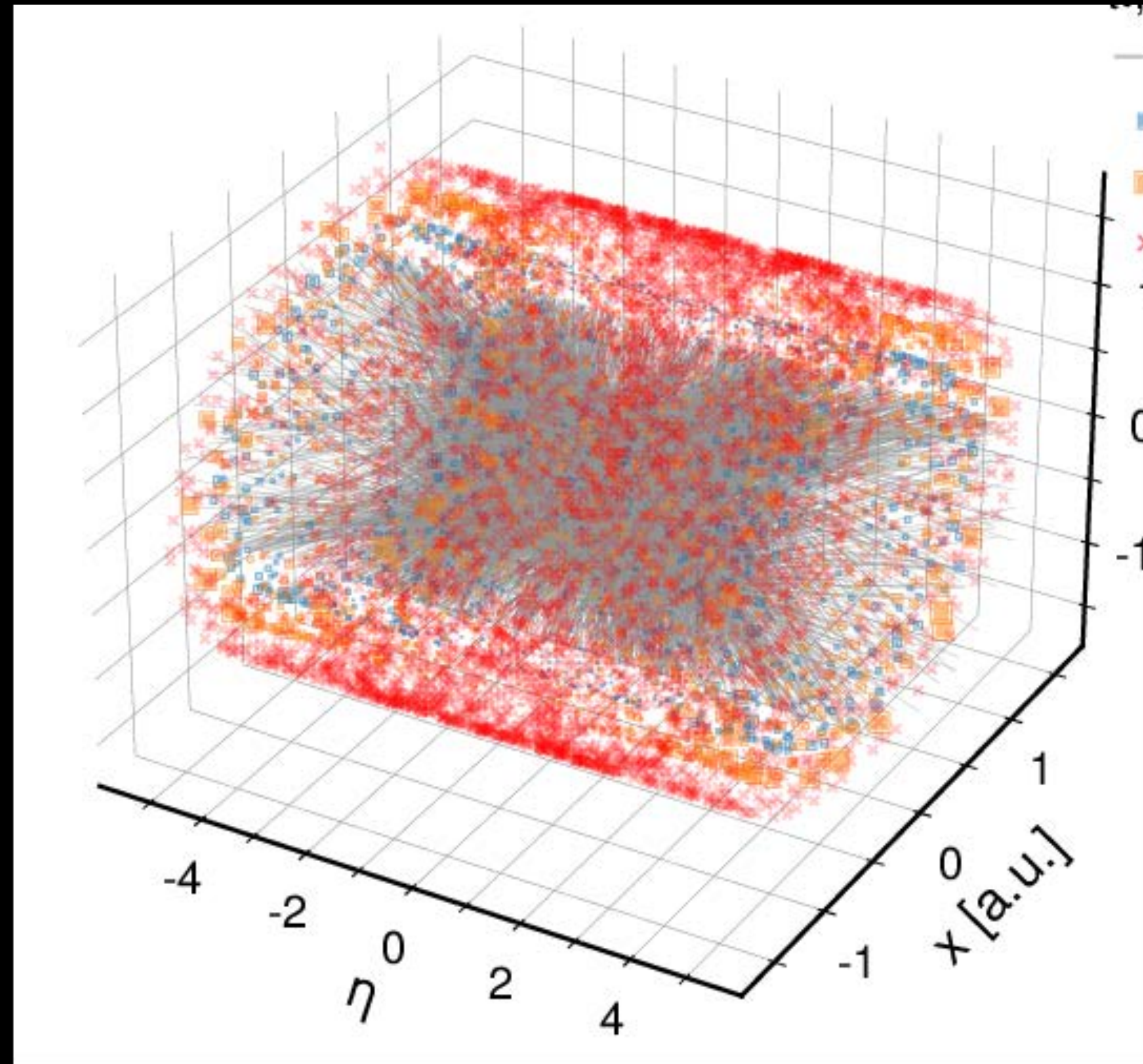
Example
Jacob [mask] reading

Jacob *fears* reading
Jacob *loves* reading
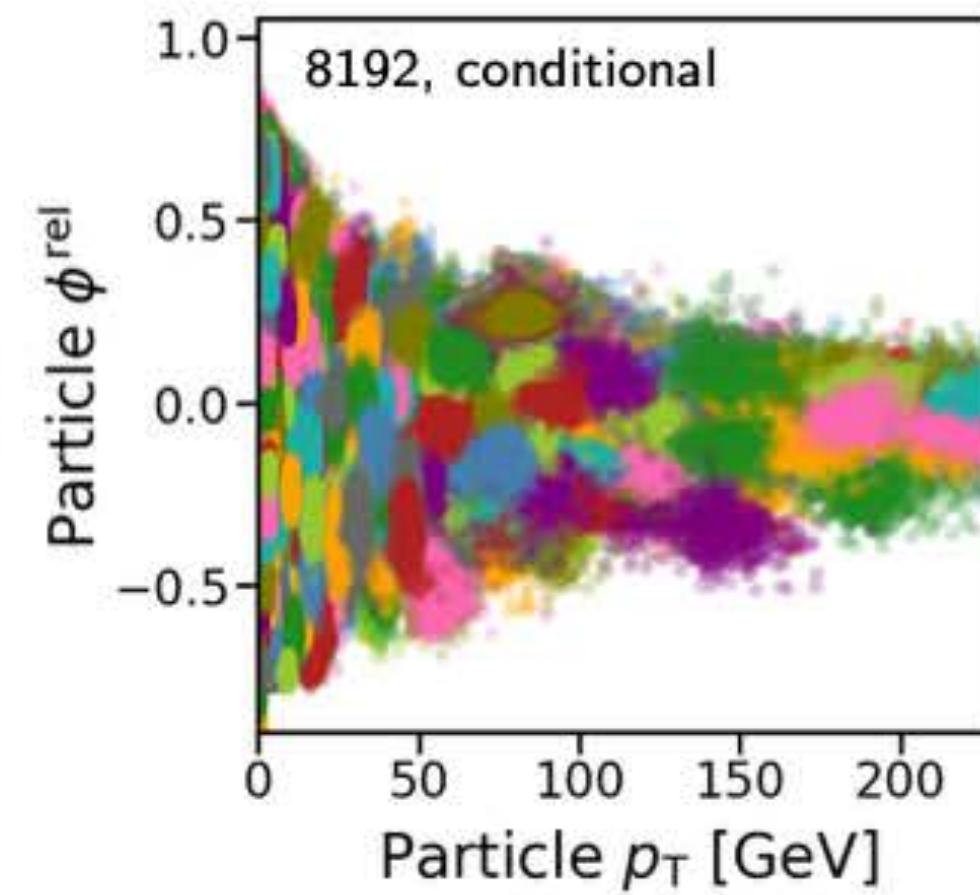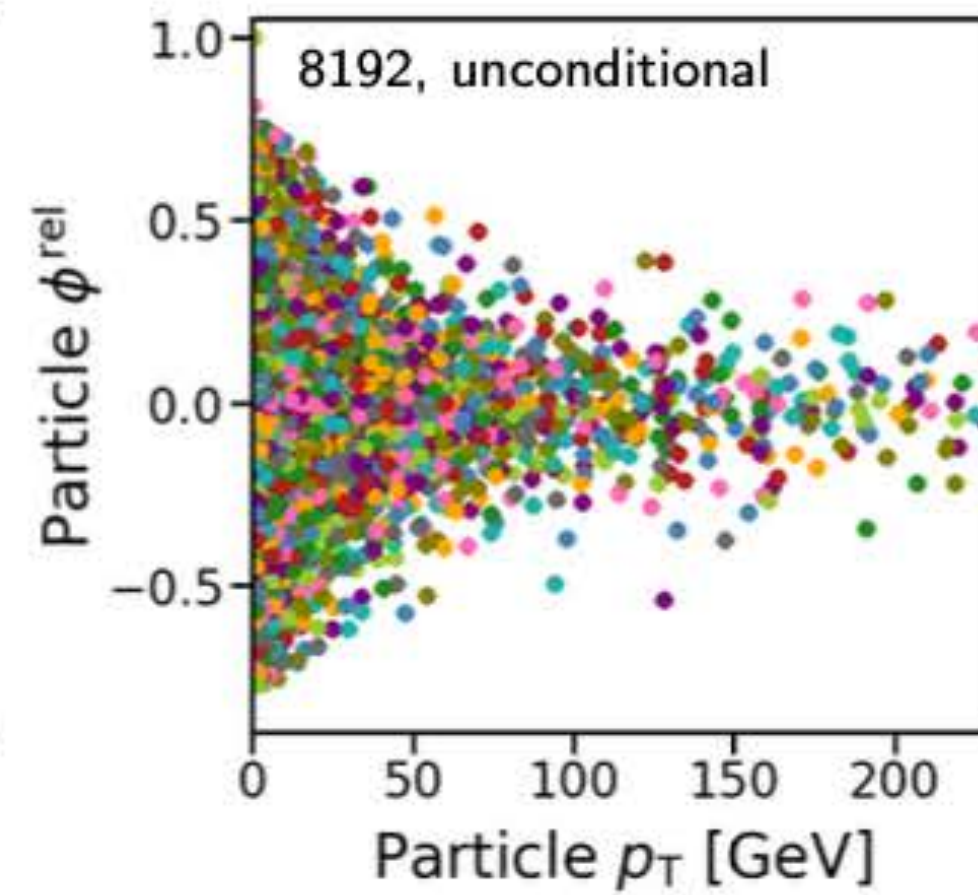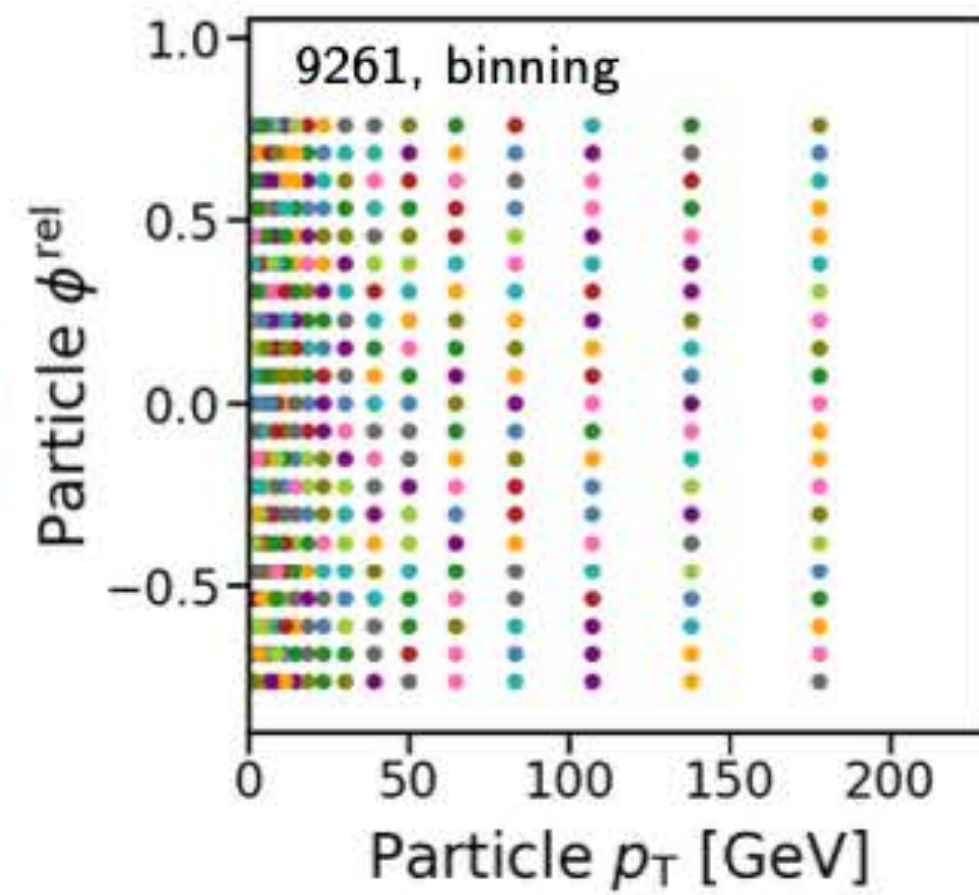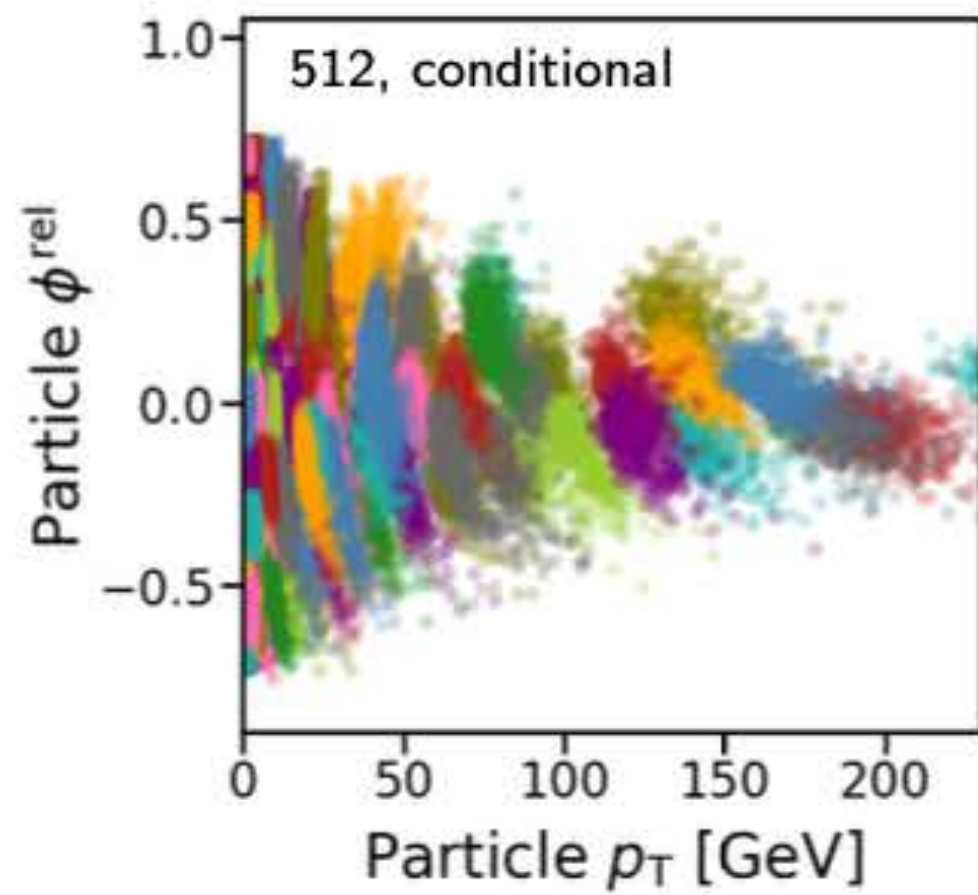Jacob *enjoys* reading
Jacob *hates* reading

# Self-supervised pre-training

# Masked particle modelling



# Masked calorimeter pre-training?

# Tokenisation?

# Hardware?



Groq: ultra-low latency dedicated language processor dedicated language processors
- Optimised for sequential data
- First ever 100 tokens/s (usually, ~10)

GPT-4

? 

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

# Backup

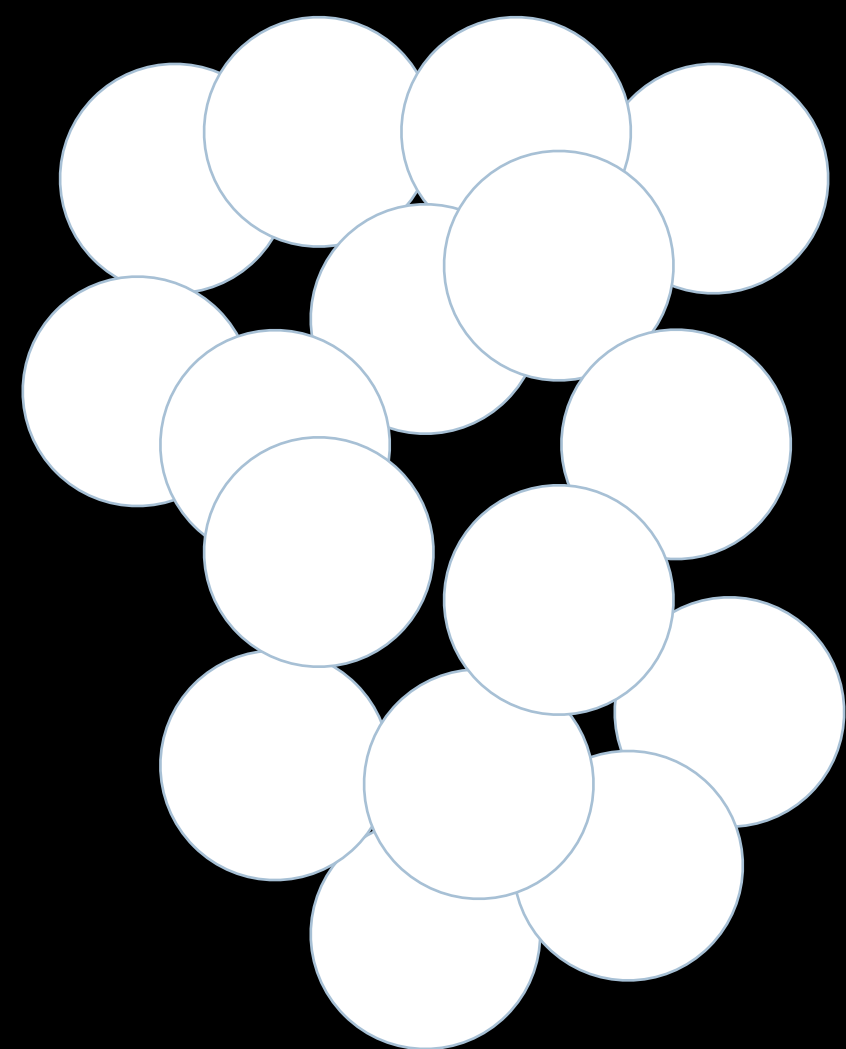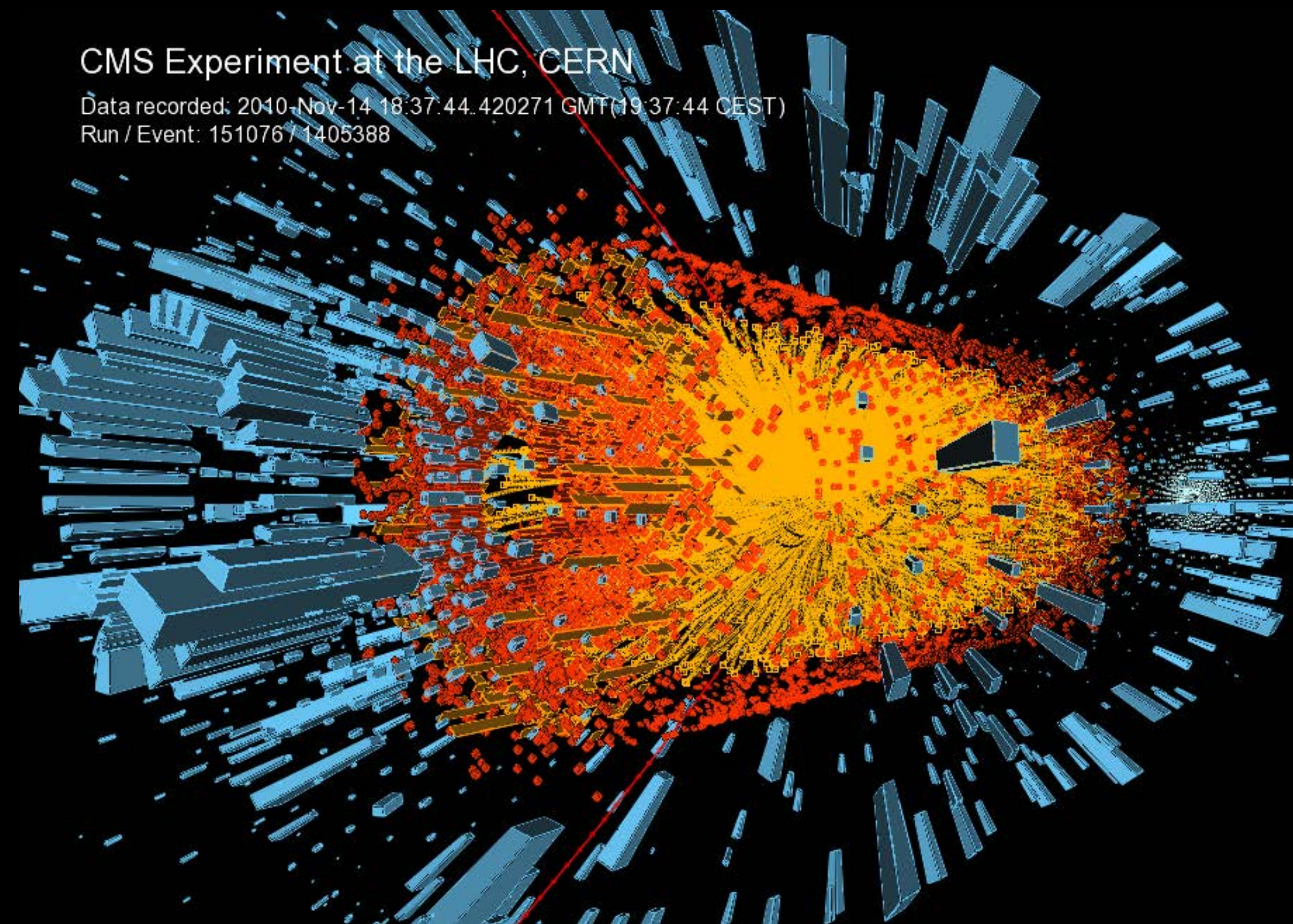# Why FPGAs?

# Why FPGAs?

- Latency (resource parallelism)



resource parallelism

# Why FPGAs?

- Throughput (pipeline parallelism)



**pipeline**

**parallelism**

Latency (resource parallelism)

Can work on different parts of problem, different data simultaneously

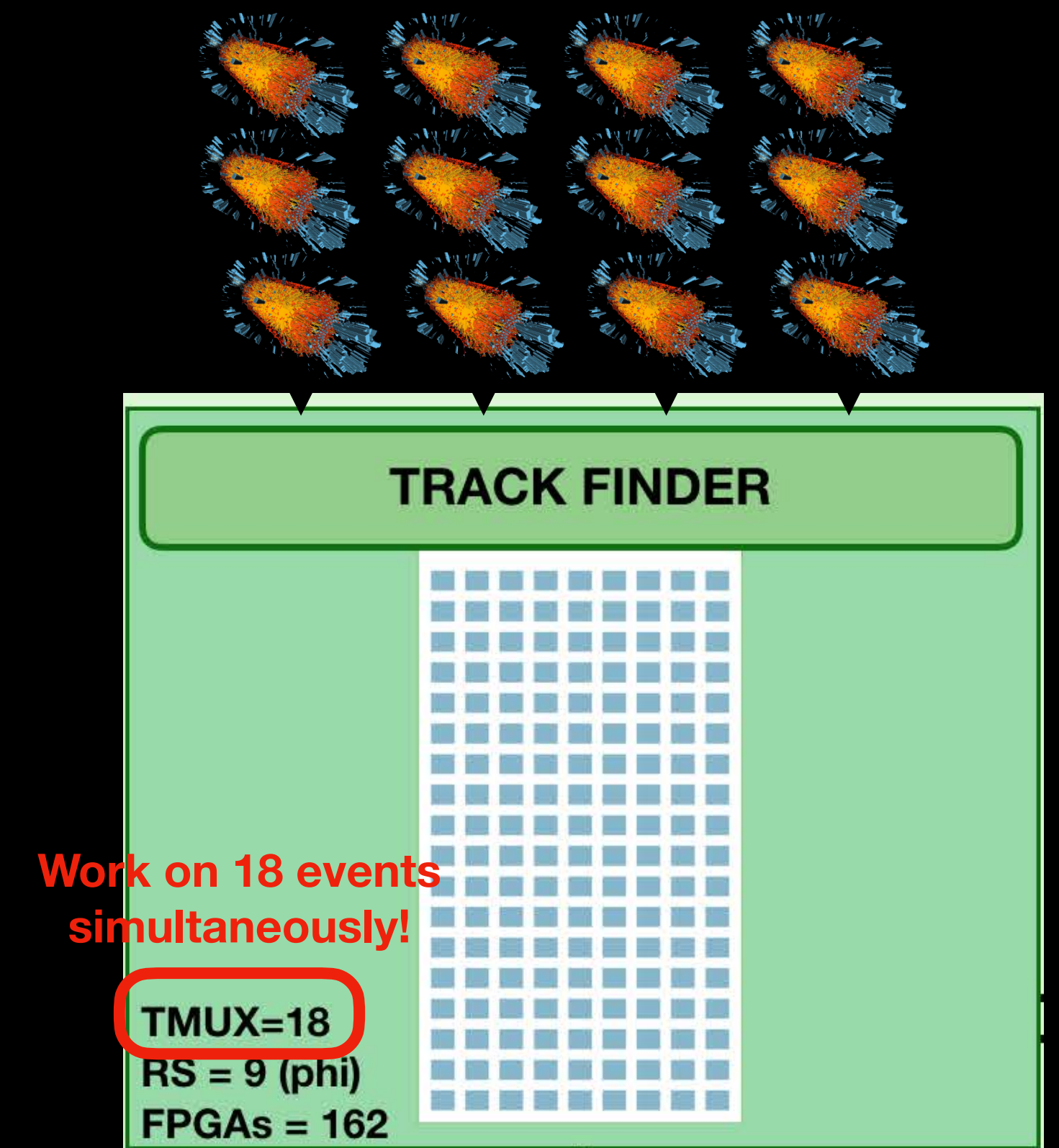Latency strictly limited by detector frontend buffer

High bandwidth (pipeline parallelism)

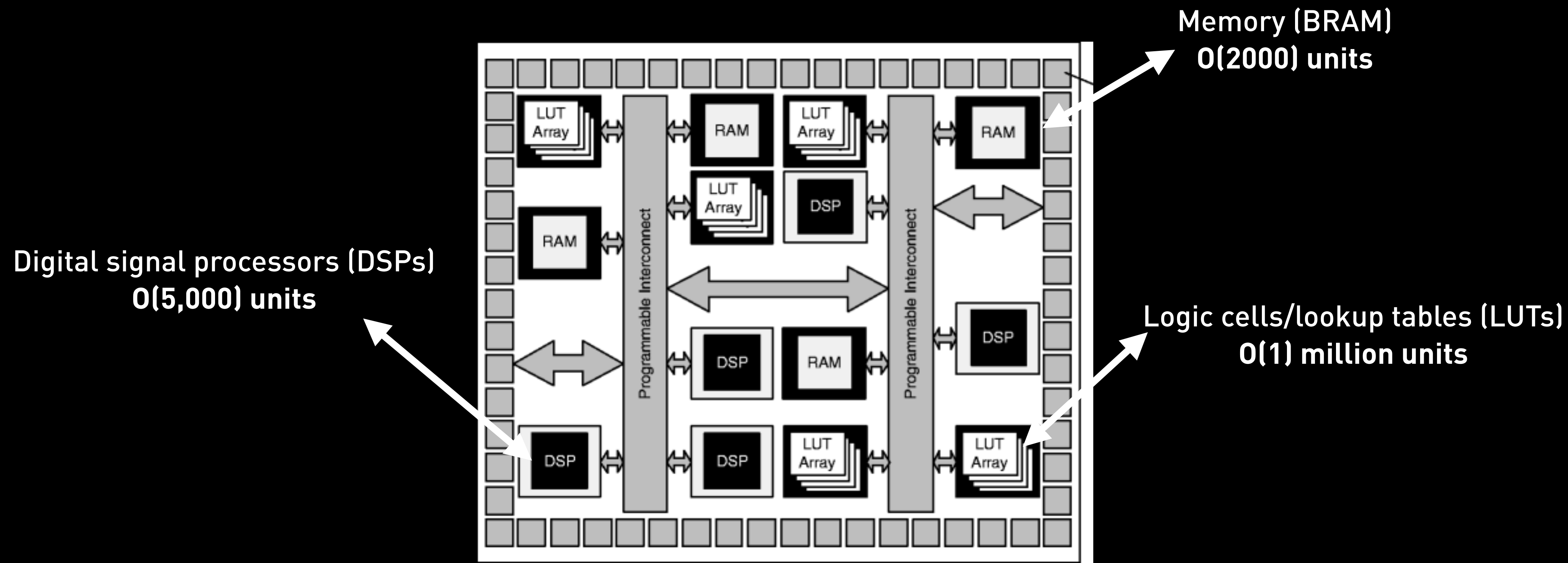Phase 2 L1T processes 5% of total internet traffic

Latency deterministic

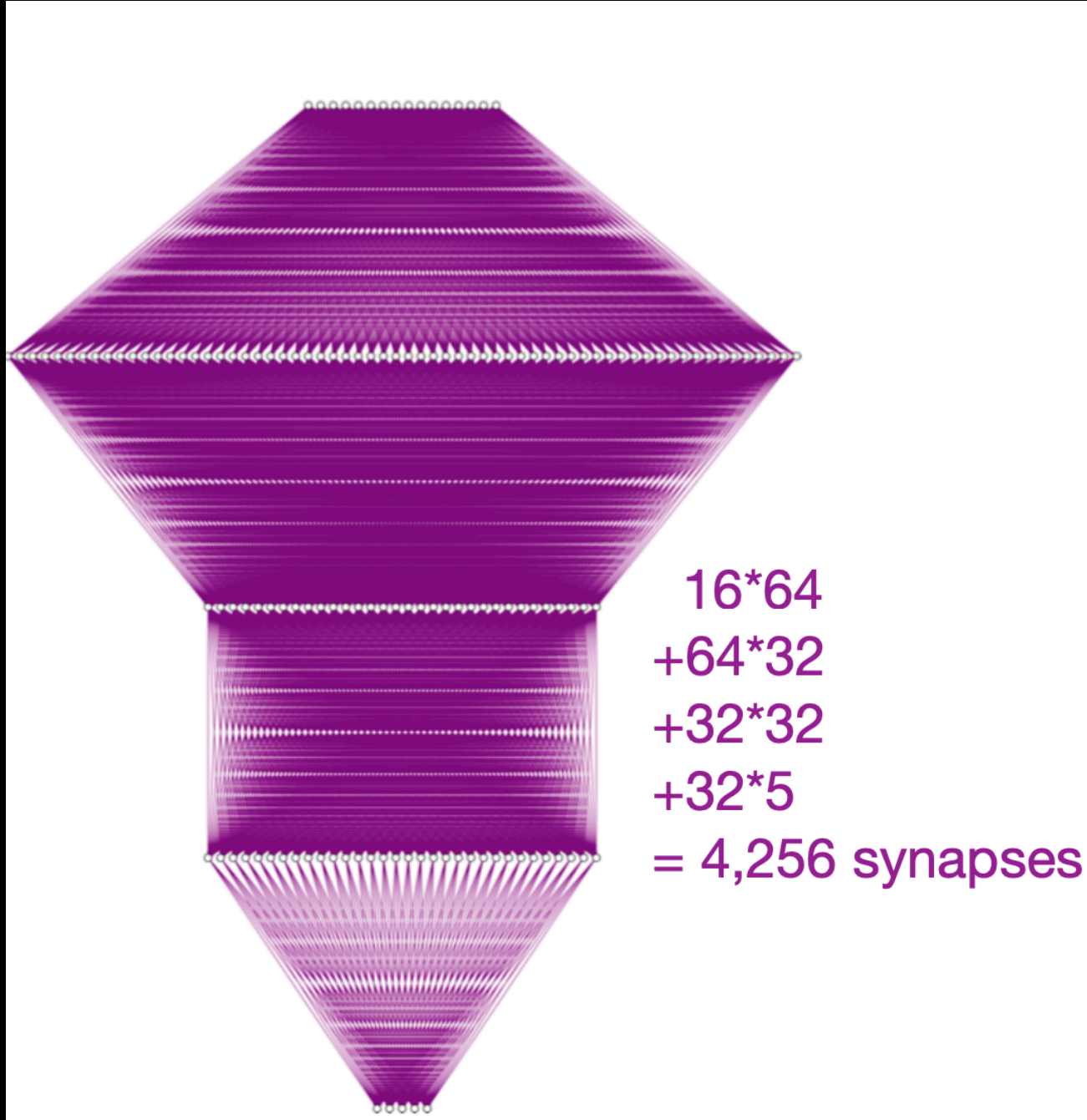CPU/GPU processing randomness,

FPGAs repeatable and predictable latency



TRACK FINDER

Work on 18 events simultaneously!

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

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

activation function      multiplication           addition

precomputed and              DSPs               logic cells
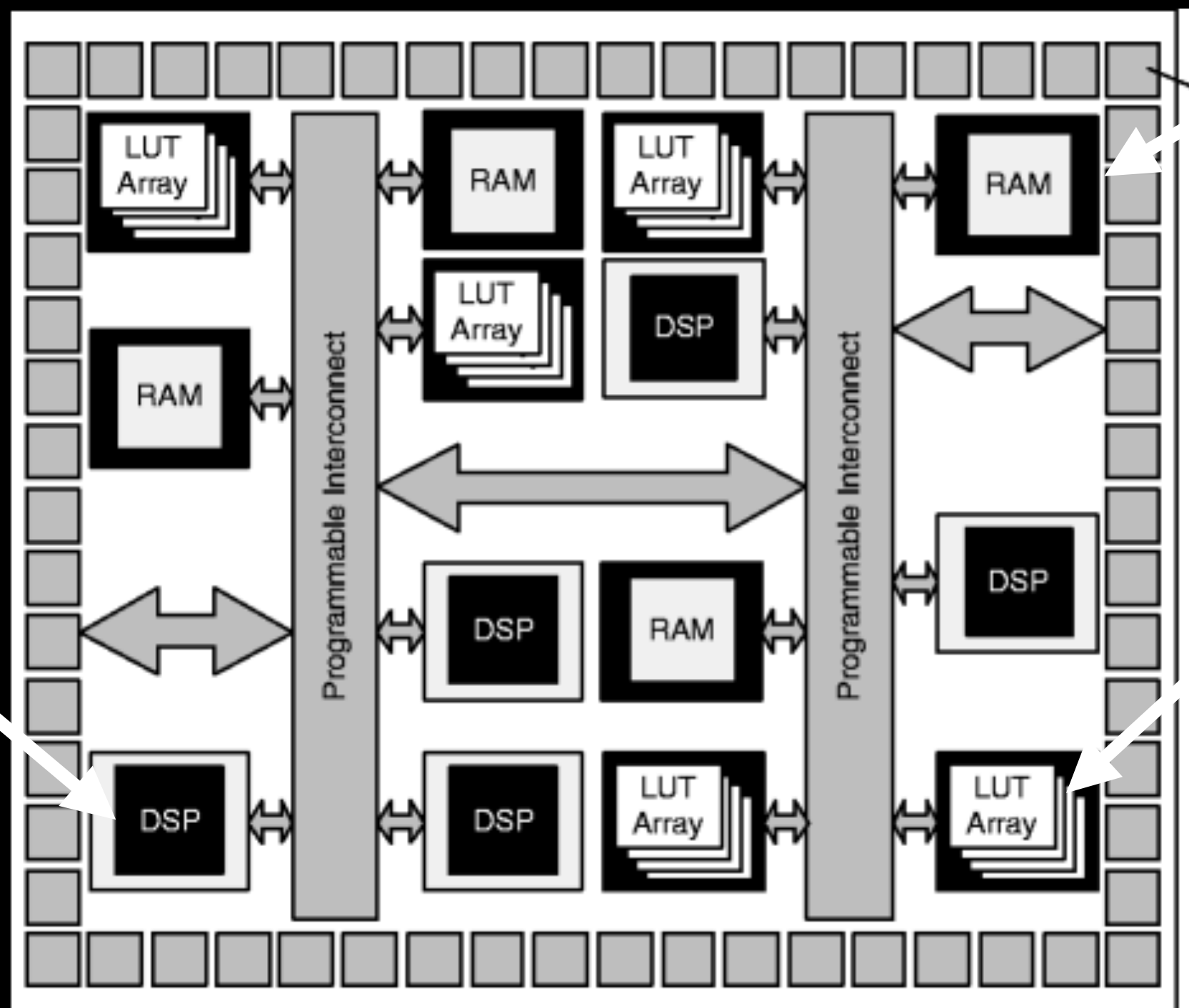stored in BRAMs

Memory (BRAM)
O(2000) units

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

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

16*64
+64*32
+32*32
+32*5
= 4,256 synapses

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

activation function
precomputed and stored in BRAMs

multiplication
DSPs

addition
logic cells

Memory (BRAM)
O(2000) units

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

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

**Ideally**

• Quantization

**Reality**

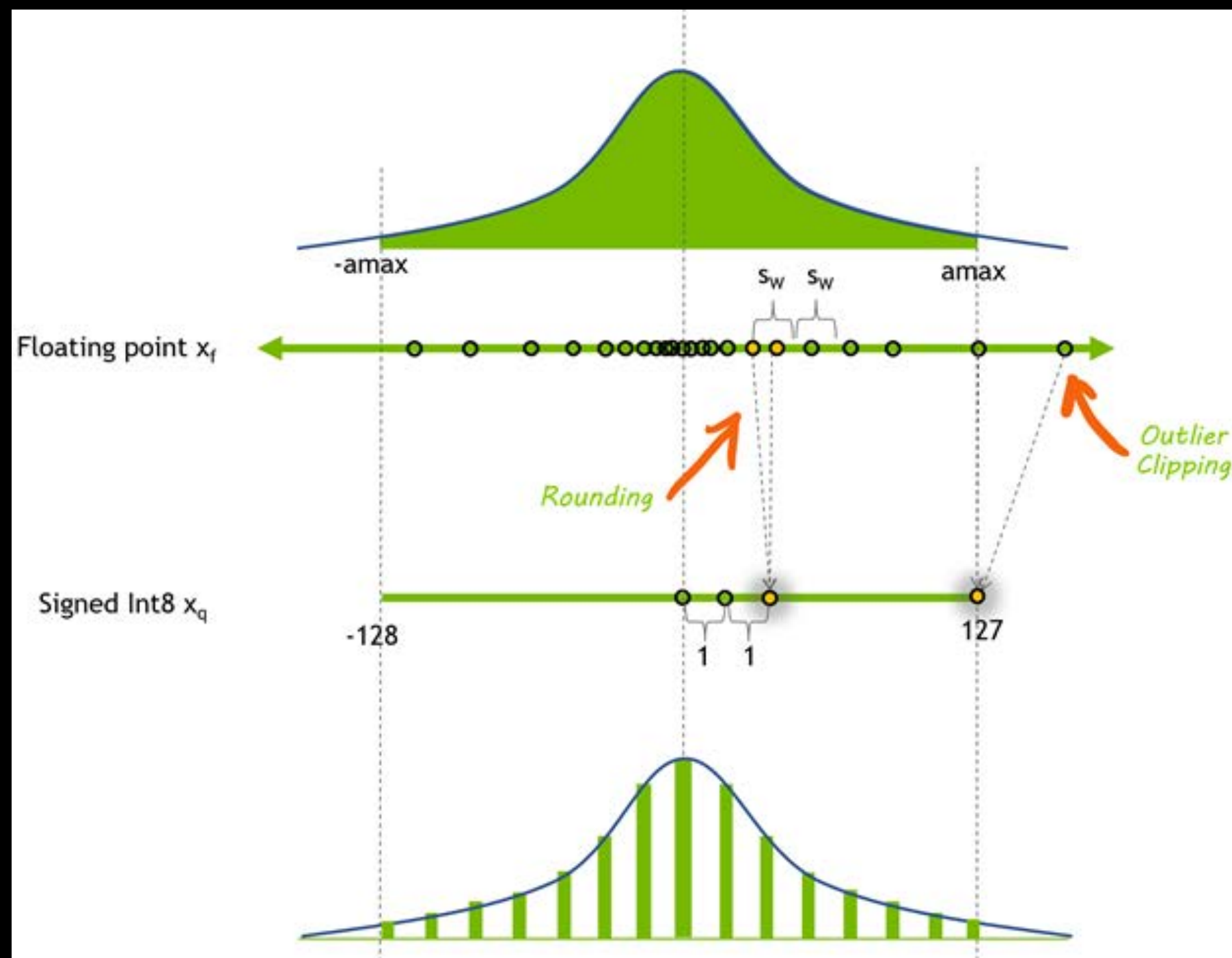# Quantization



Floating point $x_f$

**Floating point 32:**
**4B numbers in [-3.4e38, +3.4e38]**

# Quantization



**Quantising:**
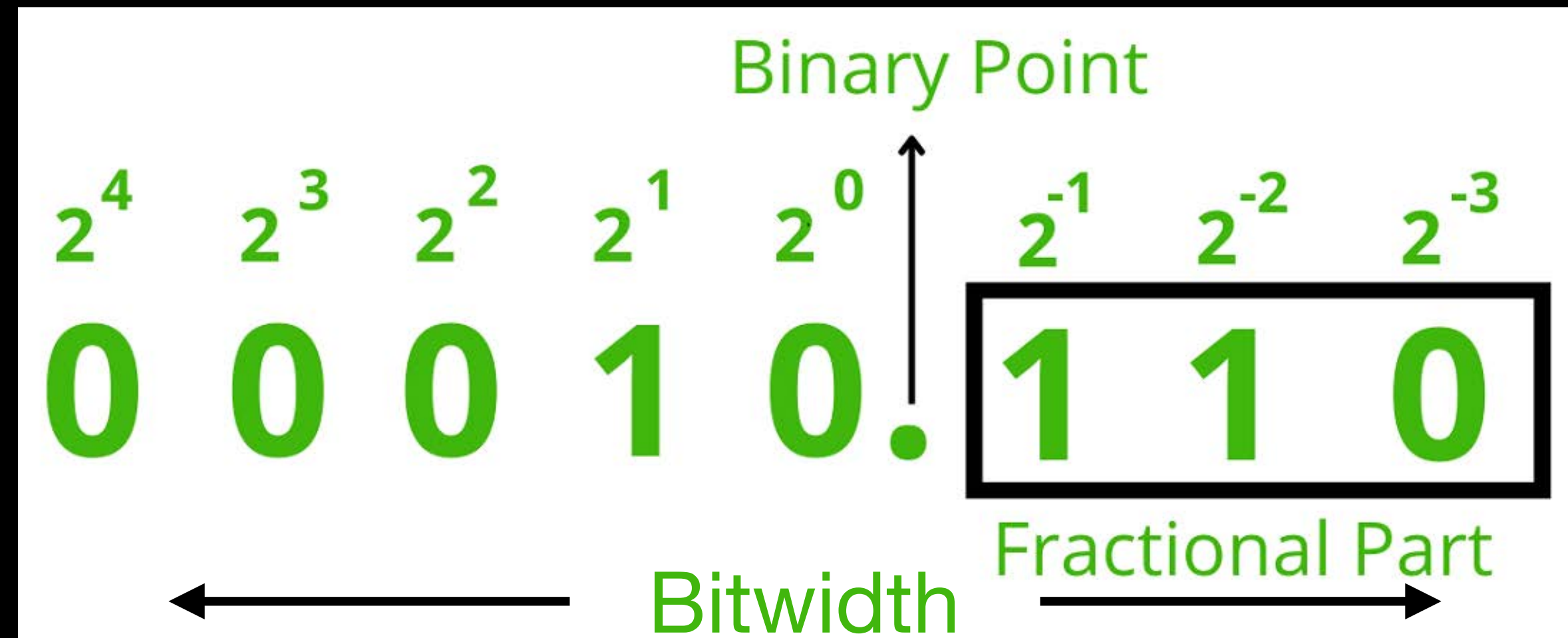**int8 $2^8$=256 numbers in [-128,127]**

$$x_q = Clip(Round(\frac{x_f}{scale}))$$

# Fixed-point $< W, I >$

a way to express fractions with integers!



$$= 2^4 \cdot 0 + 2^3 \cdot 0 + 2^2 \cdot 0 + 2^1 \cdot 1 + 2^0 \cdot 0 + 2^{-1} \cdot 1 + 2^{-2} \cdot 1 + 2^{-3} \cdot 0 = 2.75$$

# Fixed-point $< W, I >$



Trade off: range (integer bits) and precision (fractional bits). E.g $< 8,0 >$ :

$$\text{Precision} = \frac{1}{2^F} = \frac{1}{2^8} = 0.00390625$$

$$\text{Range} = [-2^0, -2^0 - 1] = [-1,0]$$

| 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**

AMD UltraScale+ MPSoC ZU19EG (conservative estimates)