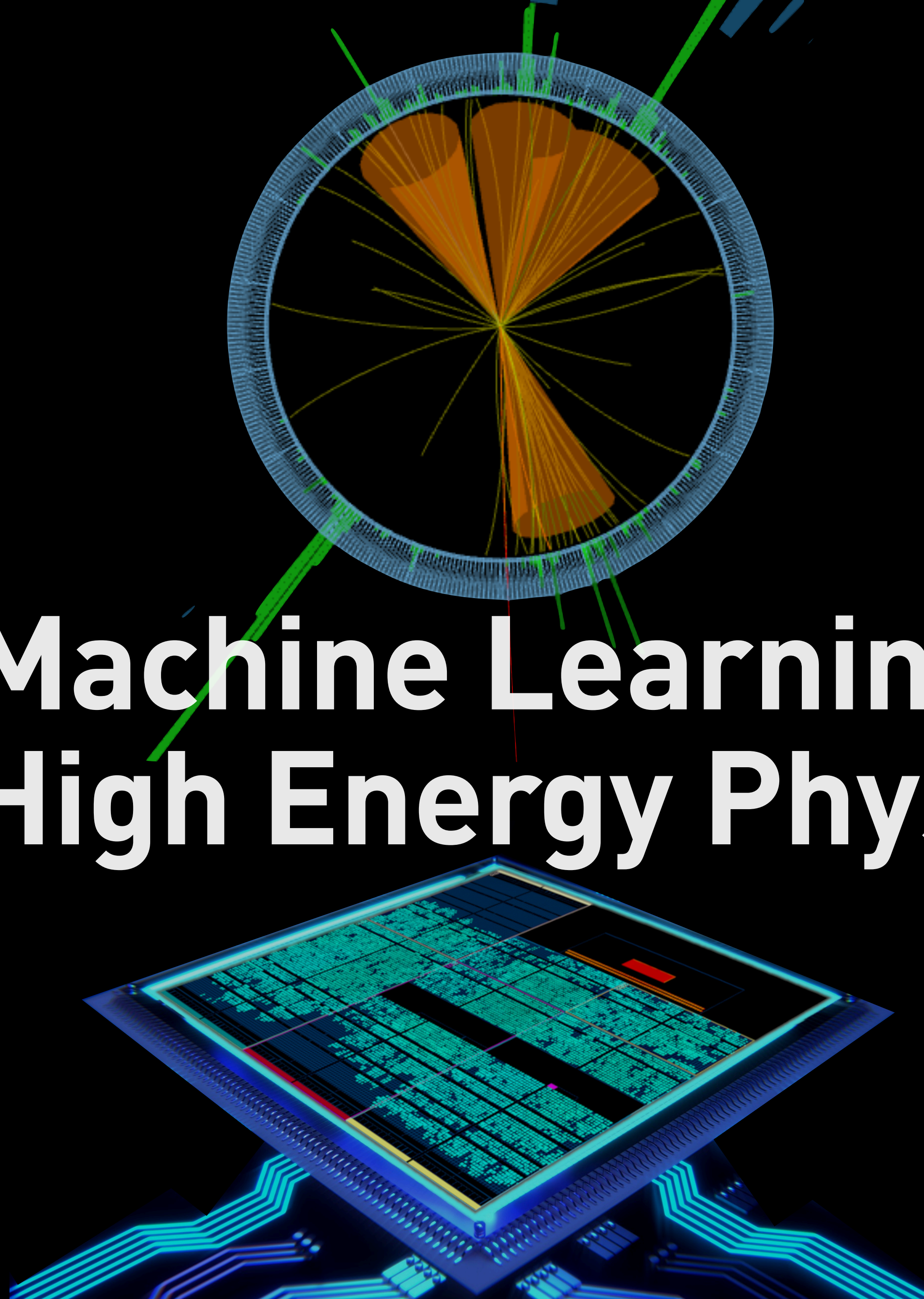# Machine Learning
# in High Energy Physics

Thea K. Årrestad (ETH Zürich)
thea.aarrestad@cern.ch
thaarres.github.io

QCD School 2024

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

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

arxiv:1407.0558

5.7σ

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

T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)

T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction

AlphaFold nature cover

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

AlphaFold nature cover

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

# GPT-3

# GPT-4 (MoE)

**175,000,000,000**
(0.16% of neurons in your brain)

**1,800,000,000,000**
(1.6% of neurons in your brain)

ChatGPT 4o ⌄

are mosquitoes in Finland known to be more aggressive than Swiss mosquitos?

Mosquitoes in Finland are often perceived as more aggressive compared to those in Switzerland. There are several reasons behind this perception:

1. **Species and Density**: Finland, especially in its northern parts, is home to a high density of mosquitoes during the summer months. The country's extensive wetlands and numerous lakes create ideal breeding grounds for mosquitoes. In contrast, Switzerland has fewer mosquito-prone areas.

2. **Climate and Activity**: The short but intense summer season in

**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 <u>carefully mapped onto hardware</u>**

AI and Memory Wall

**Computer vision: 10–100M trainable parameters ($10^{18}$ –$10^{19}$ floating point operations for training)**
**LLMs: 100M to 100Bs trainable parameters ($10^{20}$–$10^{23}$ floating point operations for training)**



Kaplan et al. (2020)

**What is deep learning?**
- innovations in network structures
- strategies to train them
- dedicated hardware

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

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

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

*Test loss*

*Compute (peta-FLOP/s-day)*

*Data set* ... *meters*

**FlashAttention**

**SRAM:** 19 TB/s (20 MB)

GPU SRAM

**HBM:** 1.5 TB/s (40 GB)

GPU HBM

**DRAM:** 12.8 GB/s (>1 TB)

Main Memory (CPU DRAM)

**Memory Hierarchy with Bandwidth & Memory Size**

Kaplan et al. (2020)

|              | Accuracy  | # params |
| ------------ | --------- | -------- |
| PFN          | 0.772     | 86.1 k   |
| P-CNN        | 0.809     | 354 k    |
| ParticleNet  | 0.844     | 370 k    |
| **ParT**     | **0.861** | 2.14 M   |
| ParT (plain) | 0.849     | 2.13 M   |

# What has changed?

Krizhevsky et al. [2012]:

      Artificial Neural Network with a **simple structure**

      (known for >20 years [LeCun et al., 1989]),

      Beat complex SOTA image recognition methods by huge margin

      **How? x100 larger and trained on a data set x100 larger**

TPU          GPU

# What has changed?

Krizhevsky et al. [2012]:

    Artificial Neural Network with a **simple structure**

    (known for >20 years [LeCun et al., 1989]),

    Beat complex SOTA image recognition methods by huge margin

    **How? x100 larger and trained on a data set x100 larger**

Made possible due to

    Graphical Processing Units (GPUs)

    Data, data and data!

TPU         GPU

# What has changed?

**Krizhevsky et al. [2012]:**
 Artificial Neural Network with a **simple structure**
 (known for >20 years [LeCun et al., 1989]),
 Beat complex SOTA image recognition methods by huge margin
 **How? x100 larger and trained on a data set x100 larger**

**Made possible due to**
 Graphical Processing Units (GPUs)
 Data, data and data!

**Deep Learning:**
 innovations in network structures,
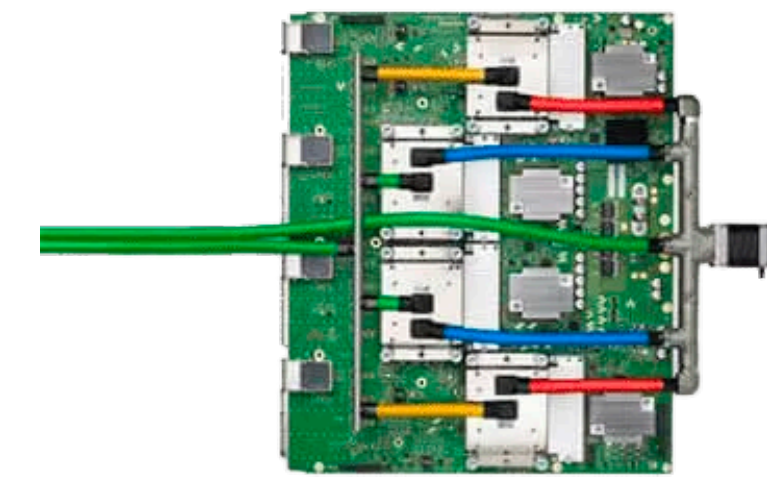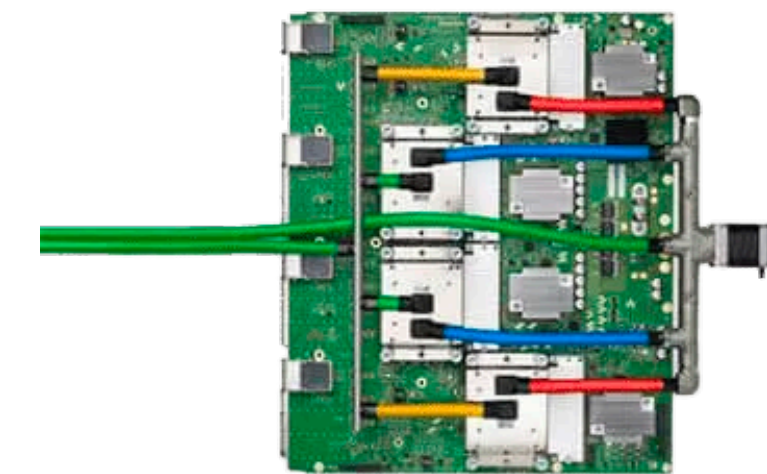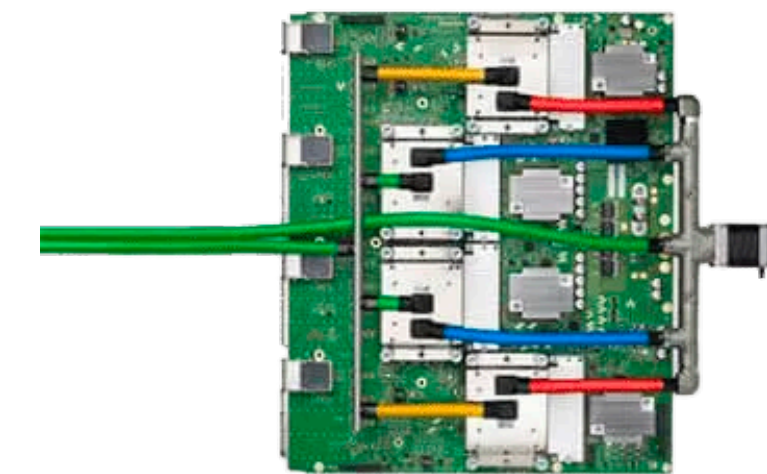 strategies to train them,
 and dedicated hardware

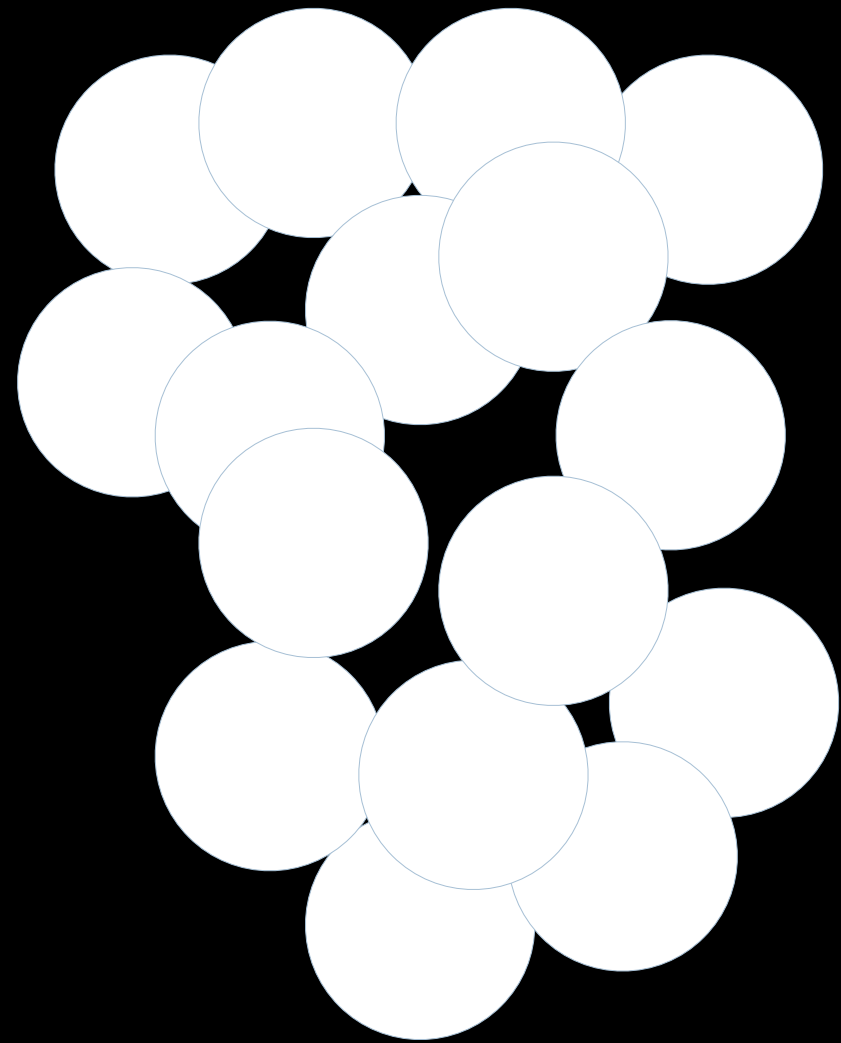**Exponential increase in size and quantity of training data [Sevilla et al., 2022]!**
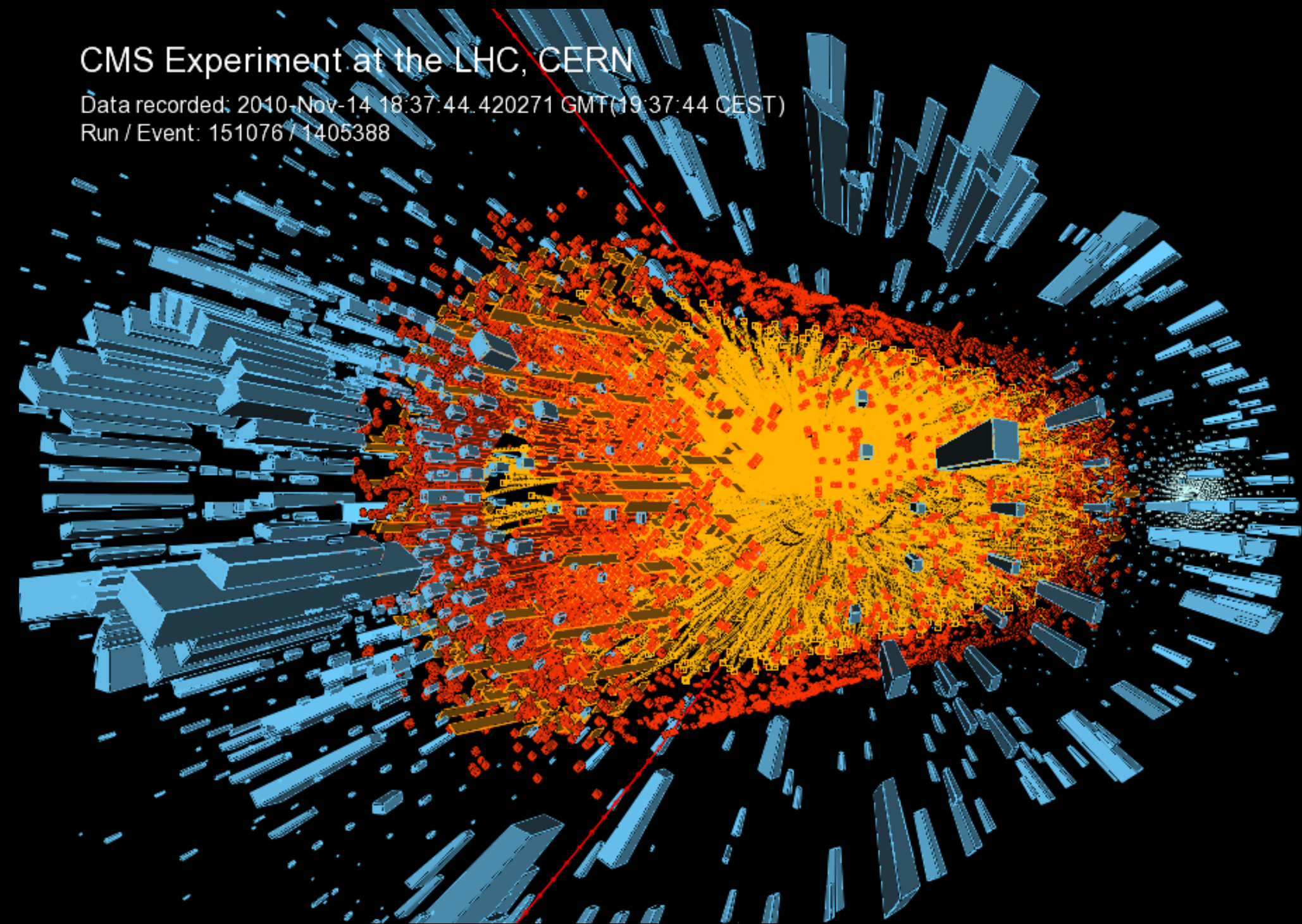


TPU  GPU

GPT-4(omni?)

?

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

**Where would you say Machine Learning has the potential to make the largest contributions within high energy physics?**

Machine learning has the potential to make significant contributions within high-energy physics in various areas. Here are some of the key areas where machine learning can have a substantial impact:

1. **Data Analysis and Event Classification:** High-energy physics experiments generate massive amounts of data, and sifting through this data to identify interesting events and classify them is a challenging task. Machine learning algorithms can help automate this process by efficiently identifying and categorizing particle collisions, reducing the workload on physicists.

2. **Anomaly Detection:** Machine learning can be used to detect anomalies or rare events in experimental data. This is crucial for identifying unexpected or new physics phenomena that may deviate from existing theories. Algorithms like autoencoders and anomaly detection models can help in this regard.

3. **Particle Tracking and Reconstruction:** Accurate reconstruction of particle trajectories and interactions is essential in high-energy physics experiments. Machine learning can improve the precision and efficiency of particle tracking, helping physicists better understand the paths and properties of particles.

4. **Simulation and Generative Models:** Generative models like Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs) can be used to simulate high-energy physics events and particle interactions. These models can generate synthetic data that can be compared with experimental data, aiding in background estimation and calibration.

5. **Detector Calibration and Optimization:** Machine learning can assist in calibrating and optimizing particle detectors. Algorithms can automatically adjust detector parameters to improve the accuracy of measurements and reduce systematic errors.

6. **Rare Signal Searches:** High-energy physics experiments often involve searching for rare signals amidst a vast amount of background noise. Machine learning techniques, including advanced statistical analysis and deep learning, can enhance the sensitivity of these searches.
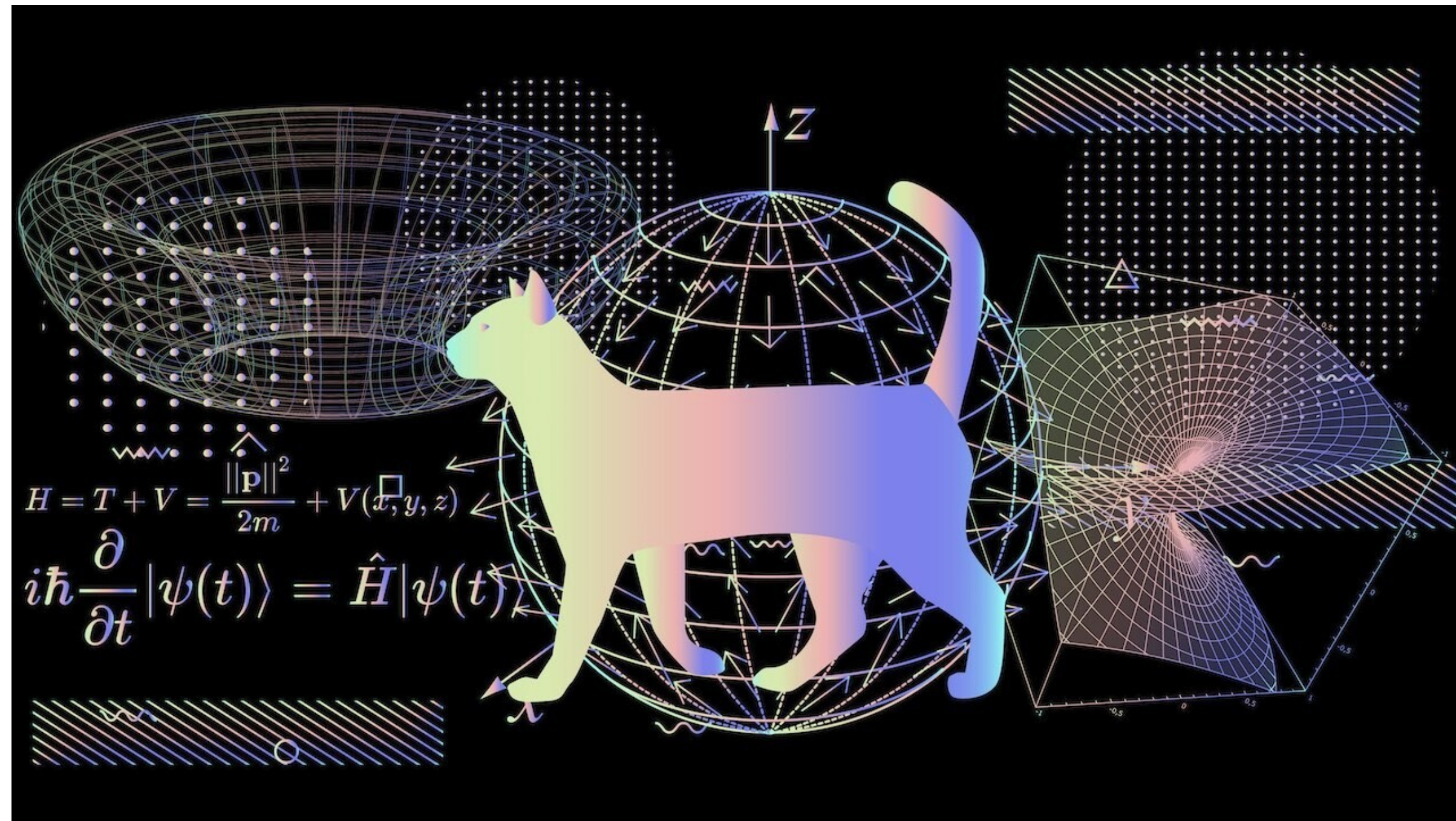
7. **Data Compression:** Storing and transmitting large datasets from particle collisions can be resource-intensive. Machine learning can be used for data compression techniques, allowing for more efficient storage and transmission of experimental data.

8. **Dark Matter and Beyond the Standard Model Physics:** Machine learning can help in the search for dark matter and other particles beyond the Standard Model. Advanced algorithms can assist in designing experiments, analyzing data, and identifying potential signatures of new physics.

9. **Optimizing Experimental Design:** Machine learning can assist in optimizing the design of high-energy physics experiments. It can help determine the most informative configurations and setups for detectors, accelerators, and data collection strategies.

10. **Reducing Human Bias:** Machine learning can help reduce human bias in data analysis by providing automated and unbiased methods for data interpretation.

It's important to note that the application of machine learning in high-energy physics often requires collaboration between physicists and machine learning experts to develop specialized algorithms and models tailored to the unique challenges of the field. Additionally, the interpretation of results and the validation of machine learning techniques remain crucial aspects of ensuring the reliability of the findings in high-energy physics experiments.
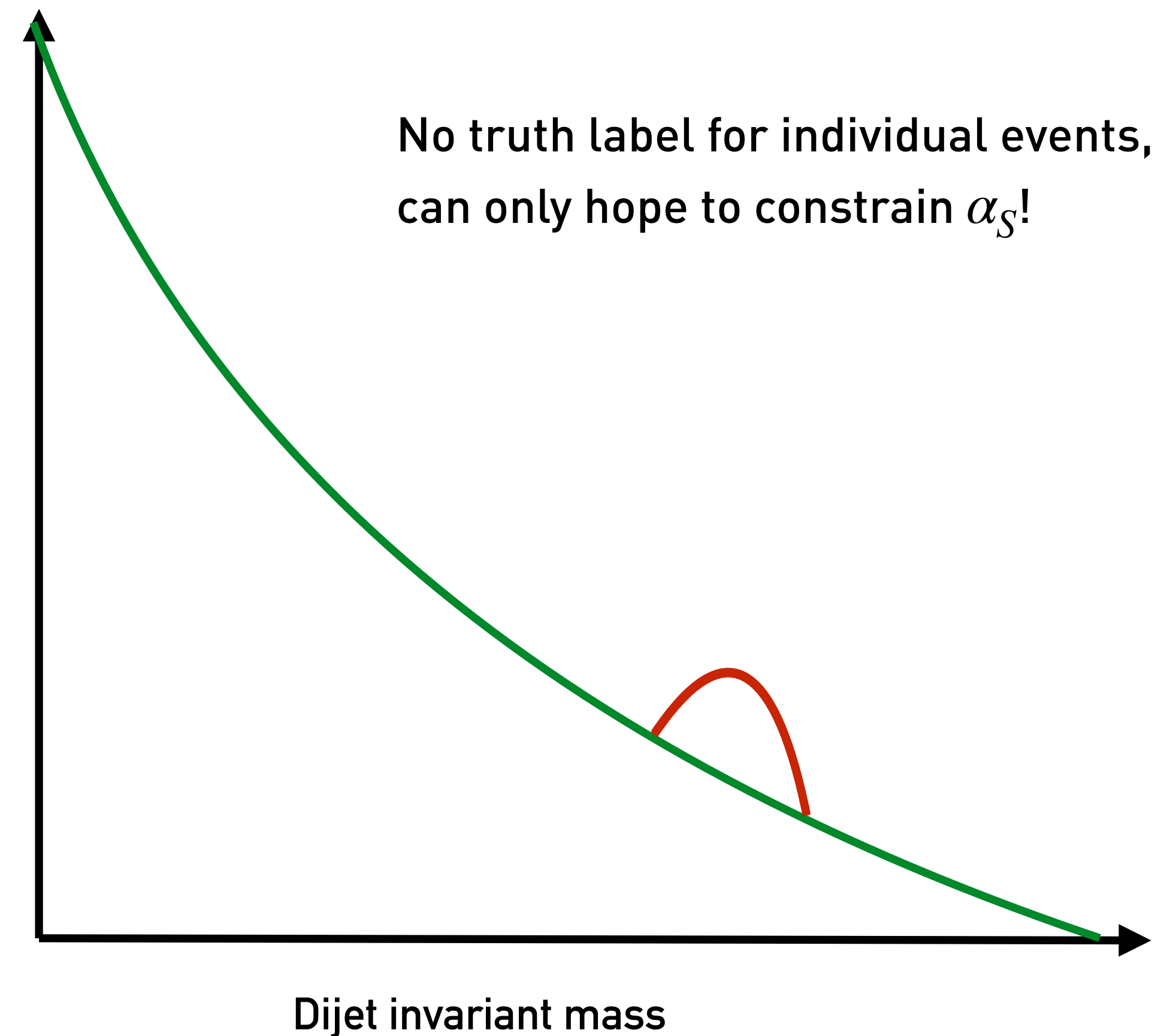
# What makes particle physics special?

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

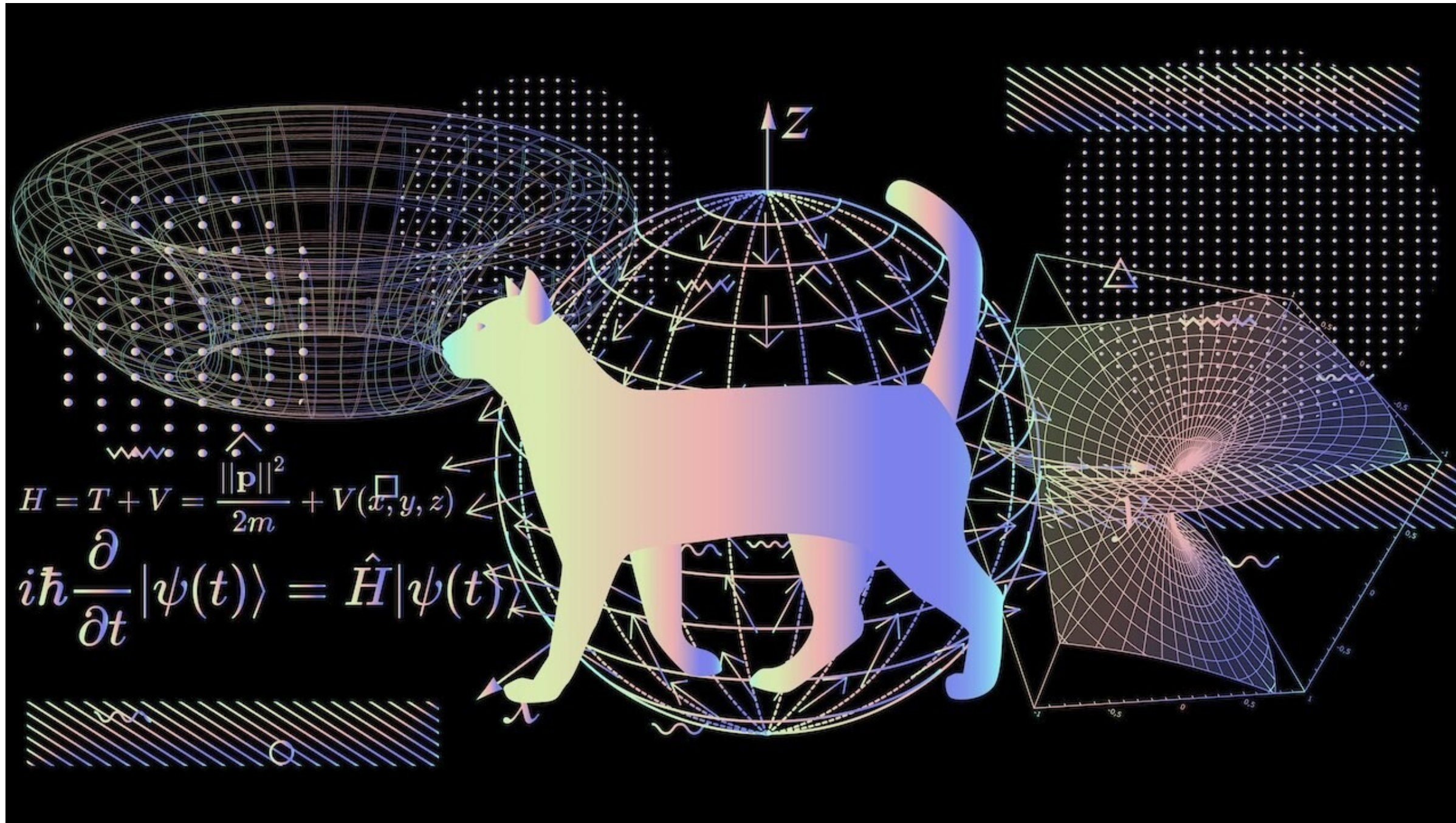$$P_{data} = \alpha_S P_S + \alpha_B P_B$$



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

$$M_S M_B {}^* + M_B M_S {}^*$$
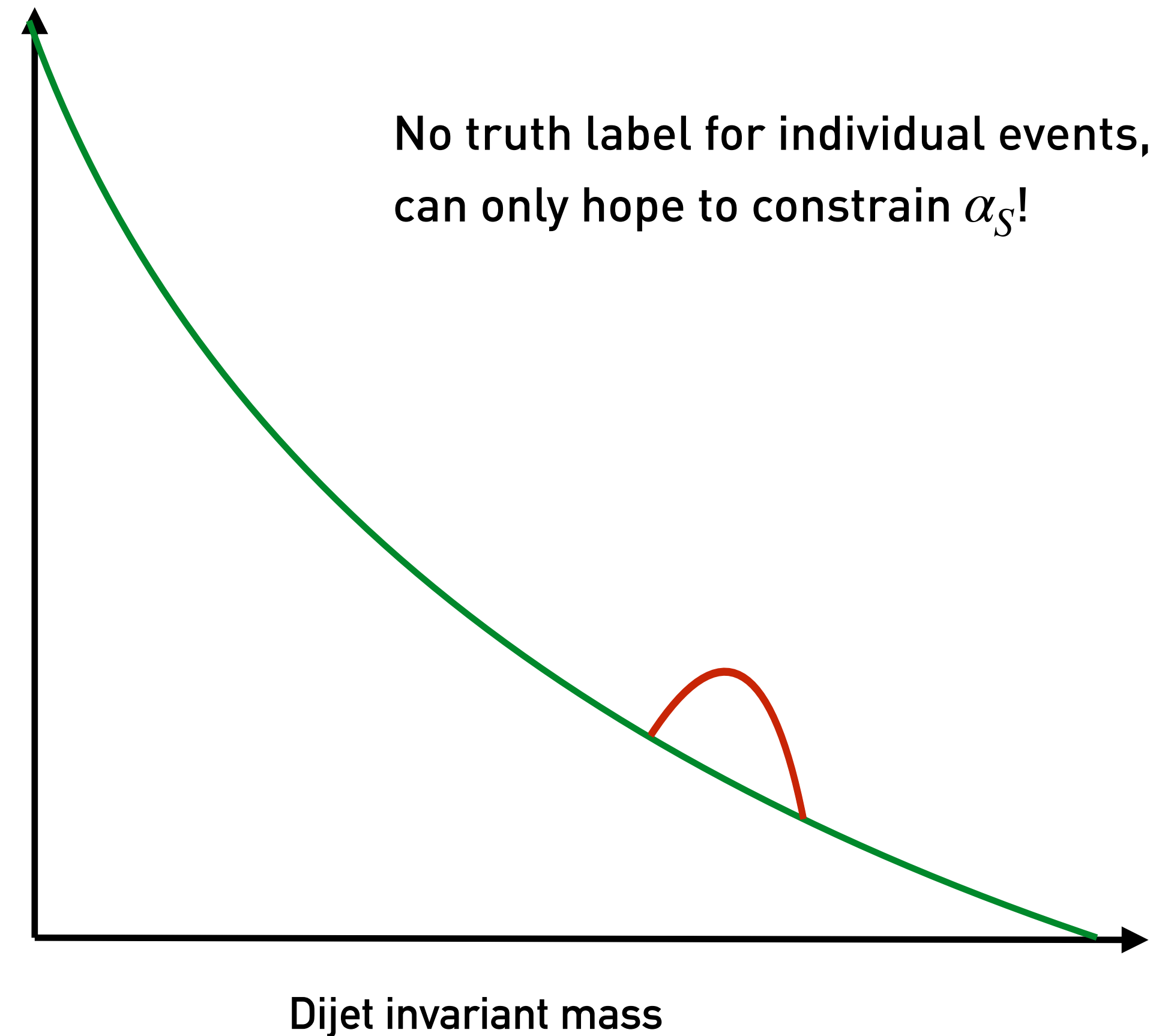
Dijet invariant mass
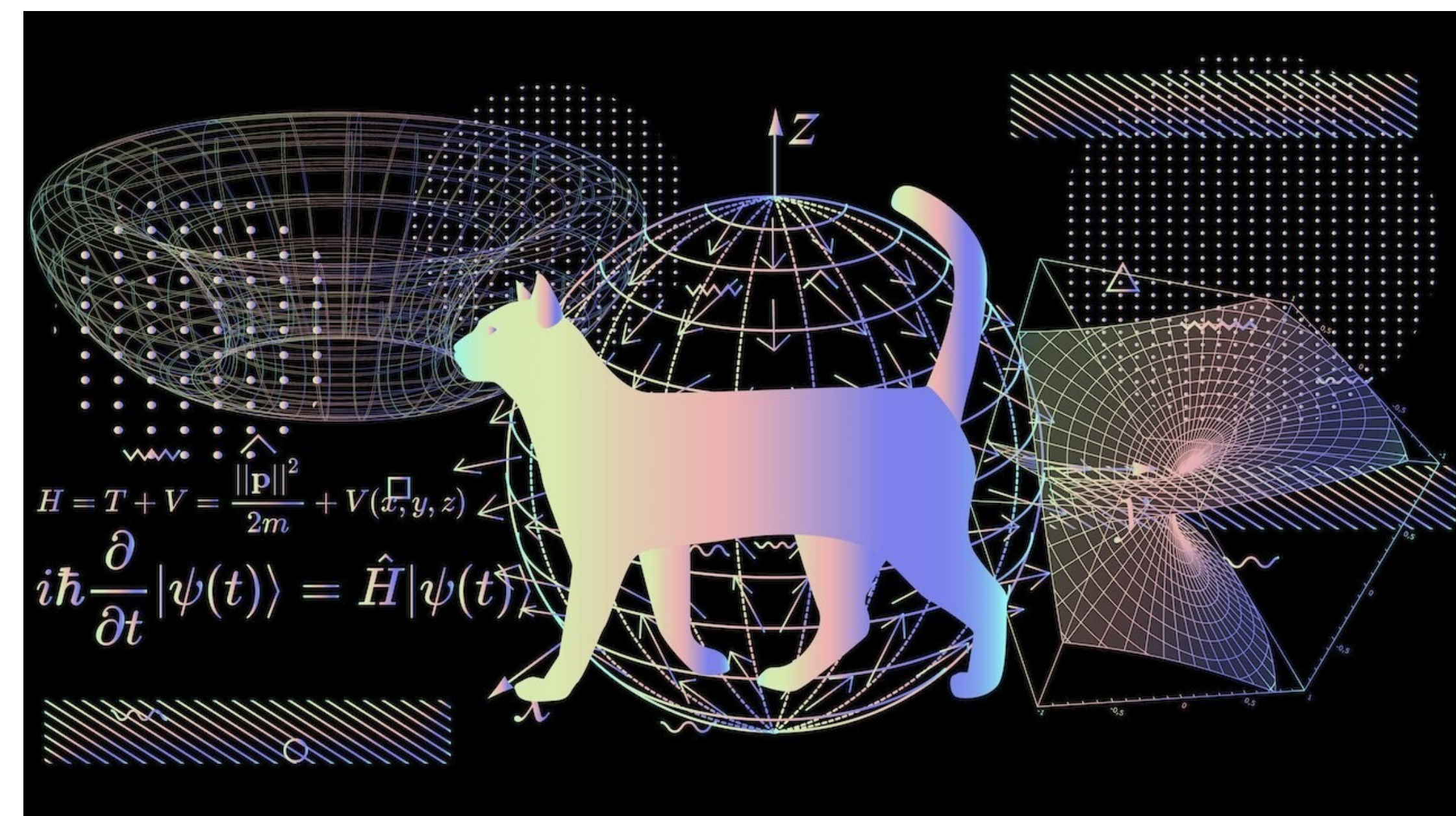
# It's against physical law to annotate our data!

$$dP_{data}^n = |M_S + M_B|^2 dp_1 dp_2 \ldots dp_n$$
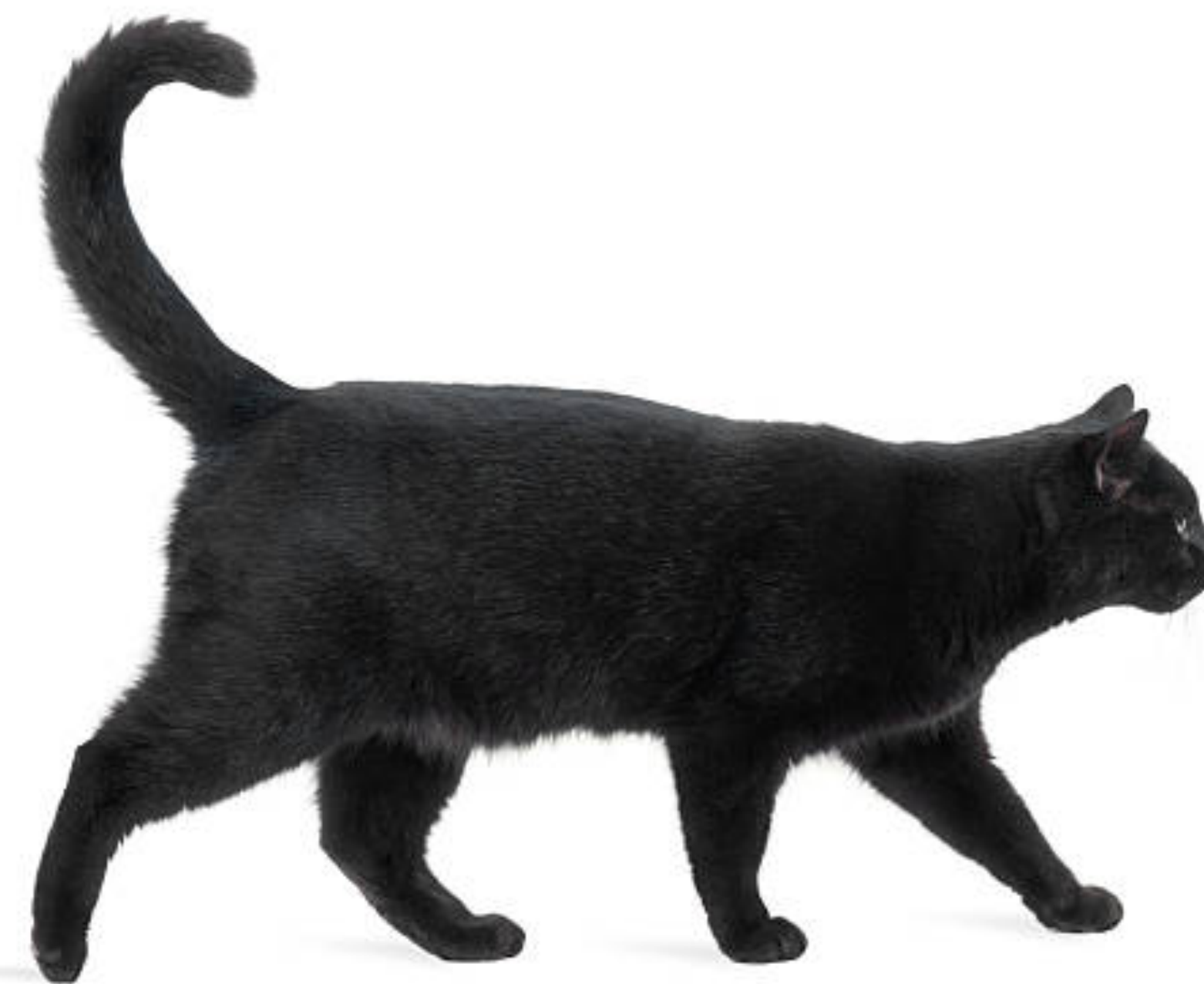
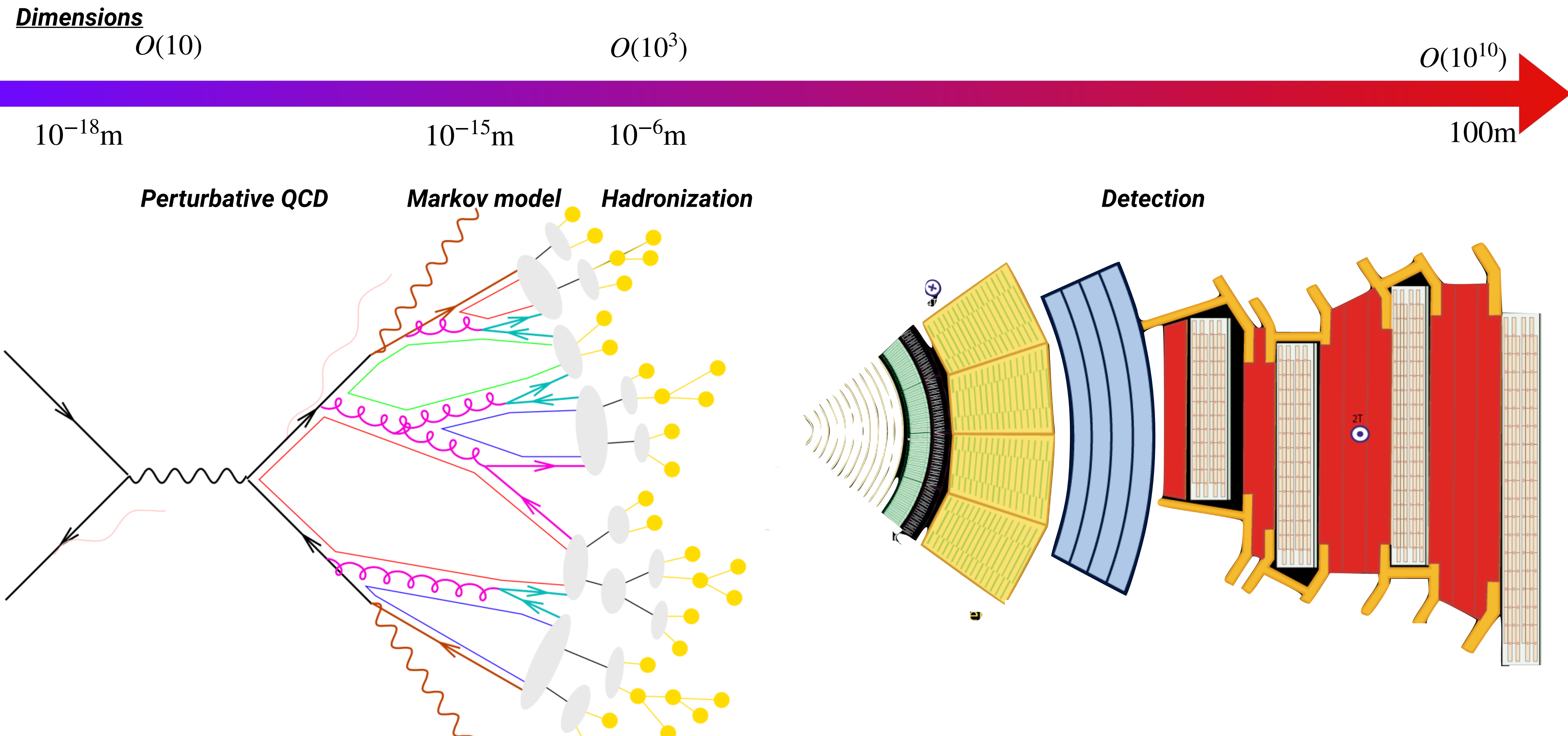$$P_{data} = \alpha_S P_S + \alpha_B P_B$$



$$M_S M_B{}^* + M_B M_S{}^*$$

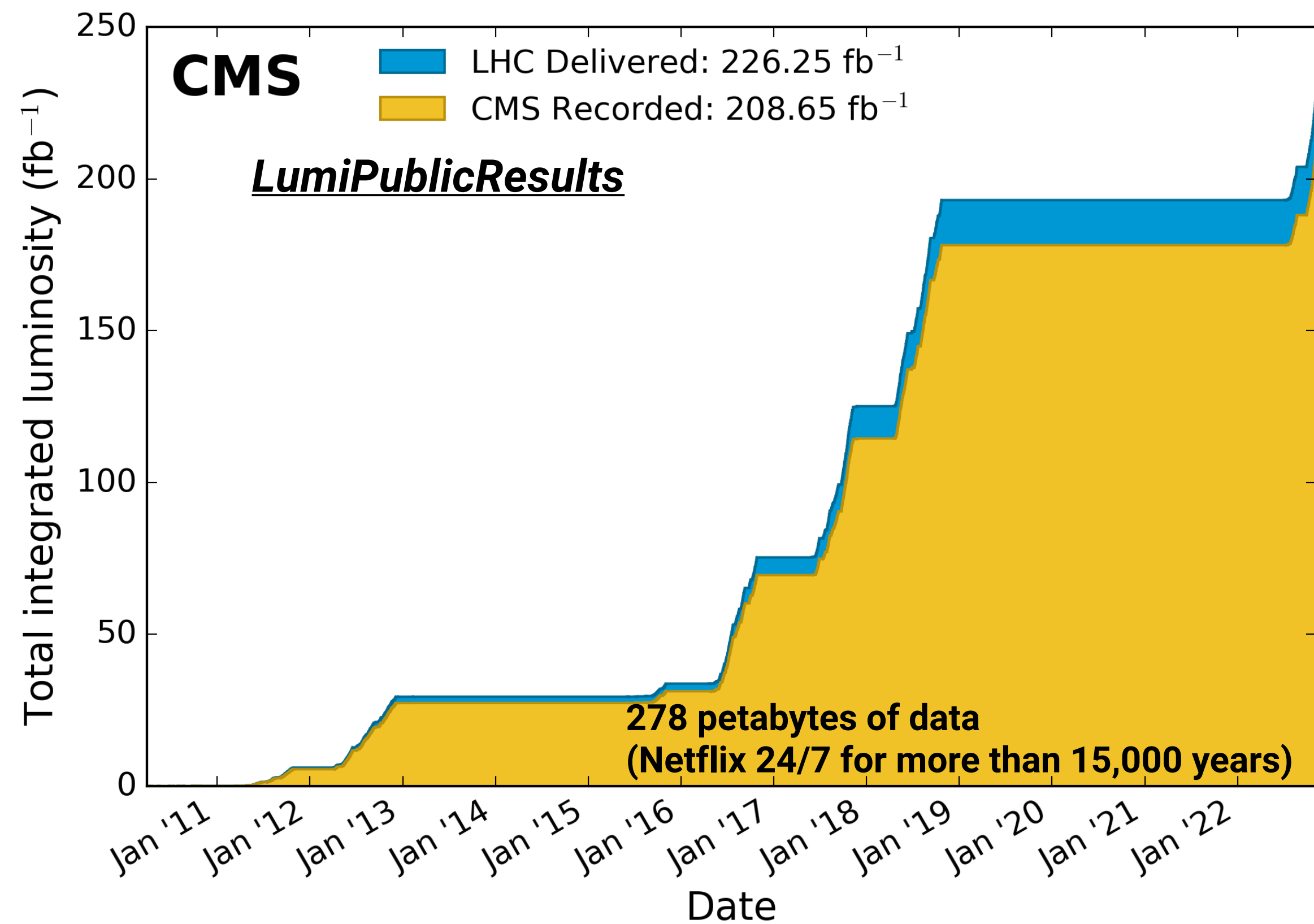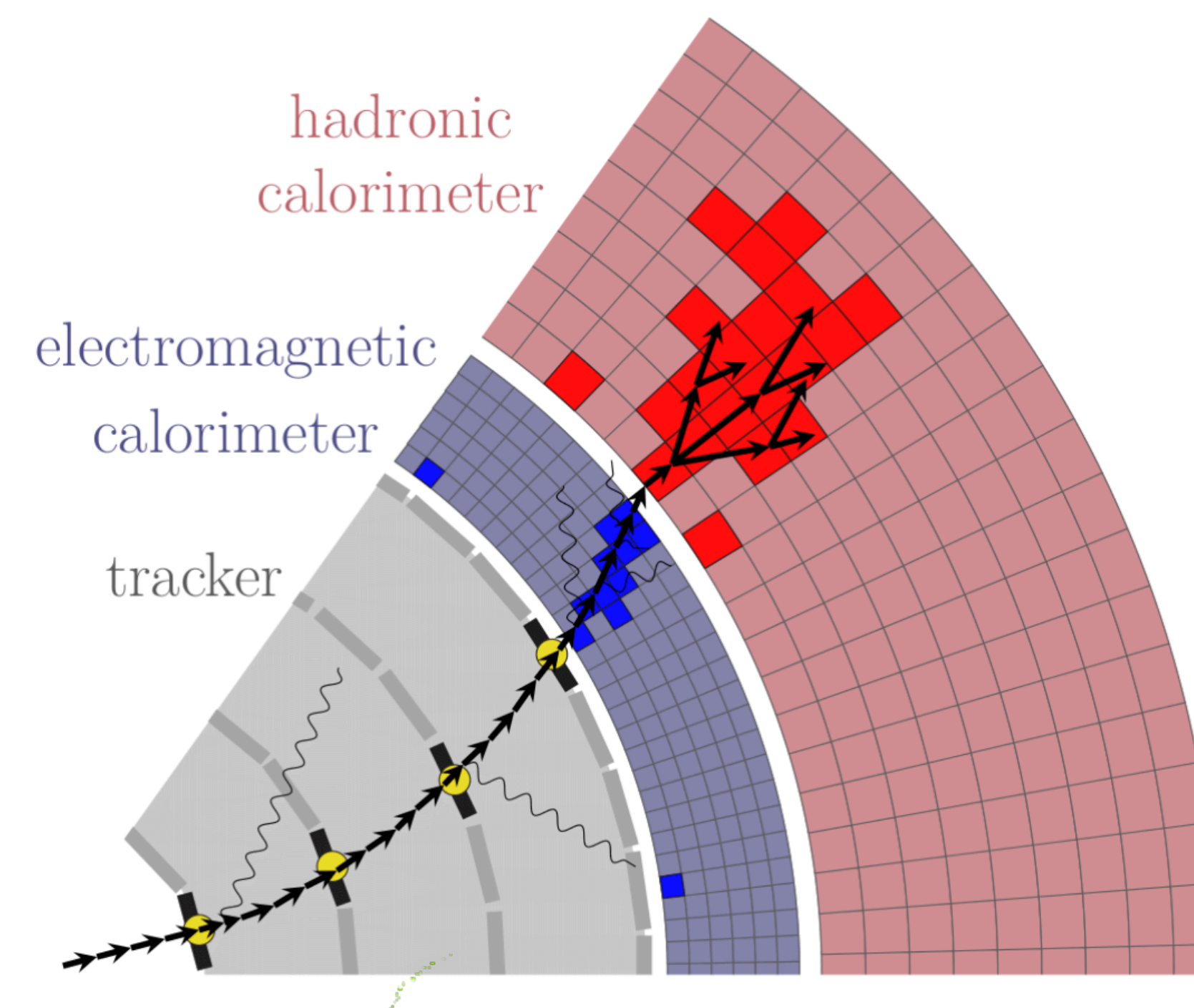No truth label for individual events, can only hope to constrain $\alpha_S$!



Dijet invariant mass

# Monte Carlo Simulation

**Dimensions**

$O(10)$

$O(10^3)$

$O(10^{10})$

$10^{-18}$m

$10^{-15}$m

$10^{-6}$m

$100$m

*Perturbative QCD*

*Markov model*

*Hadronization*

*Detection*

**~40 quadrillion collisions recorded at LHC**
**(1 fb⁻¹ ~ 100 trillion collisions)**

**CMS**

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

*LumiPublicResults*

**278 petabytes of data**
**(Netflix 24/7 for more than 15,000 years)**

**O(1) trillion simulated events**

**GEN**        **SIM**        **DIG**

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**
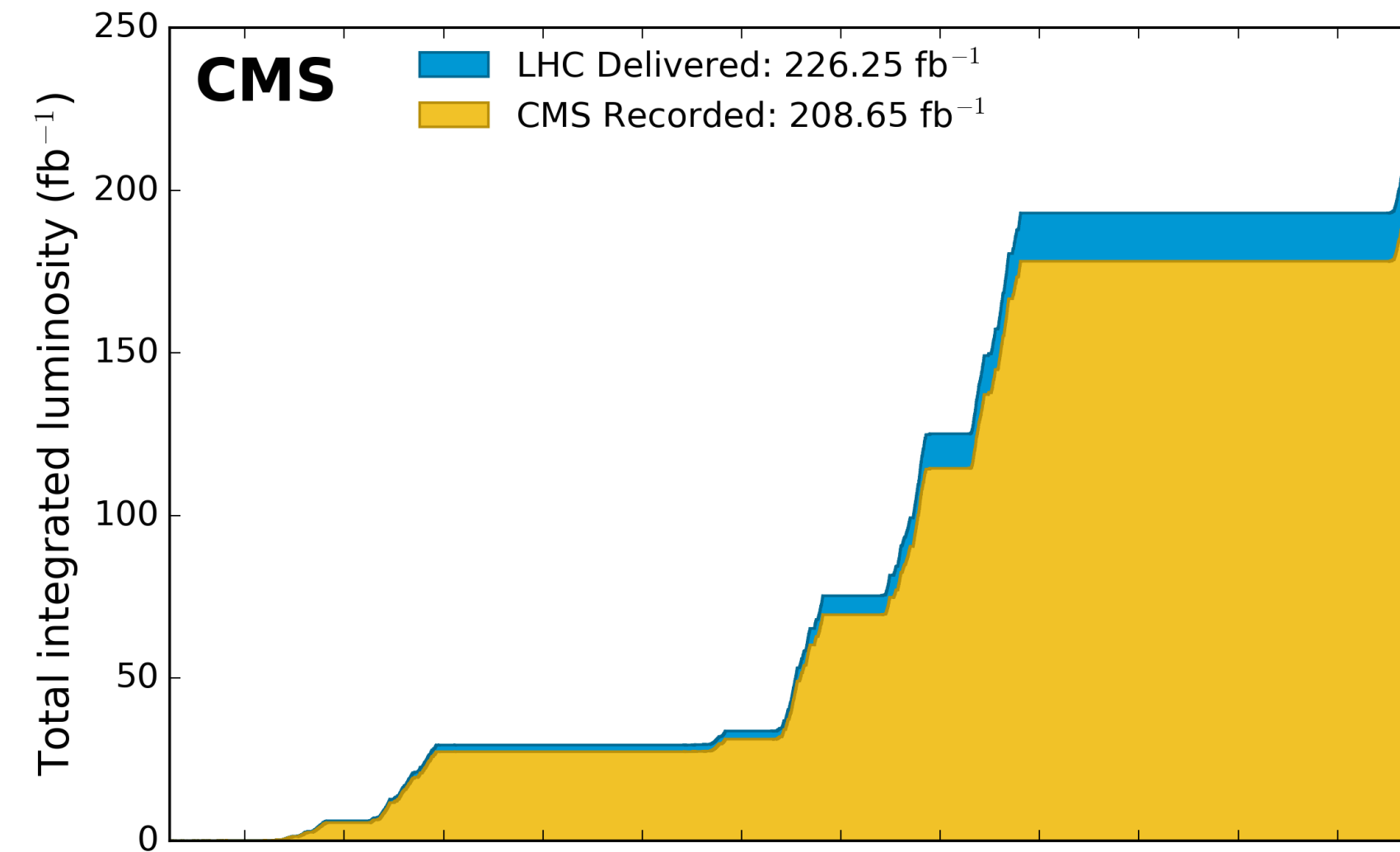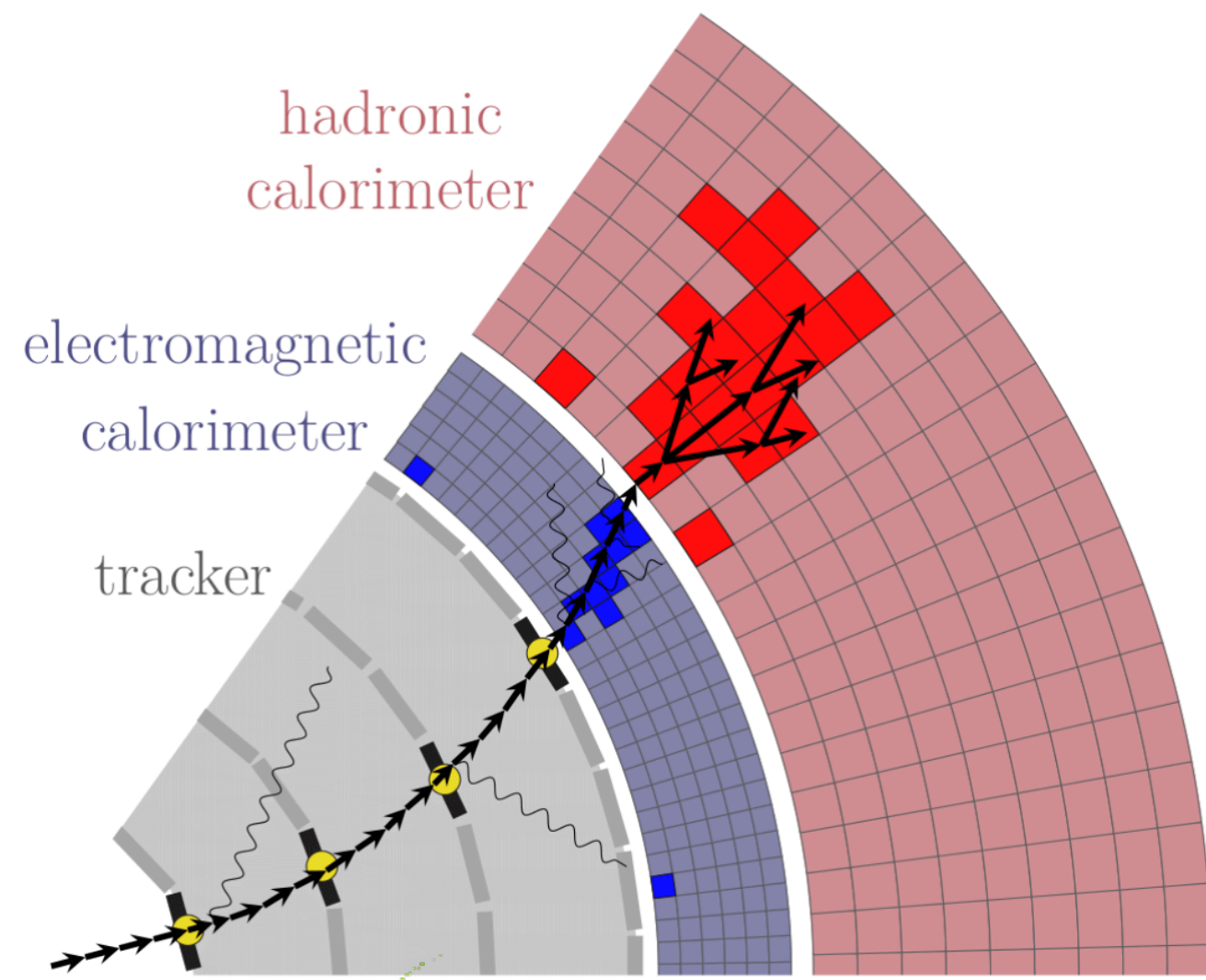- Requires truth labels
- Only possible using simulation

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

We have a lot of high quality simulated data that we want to use

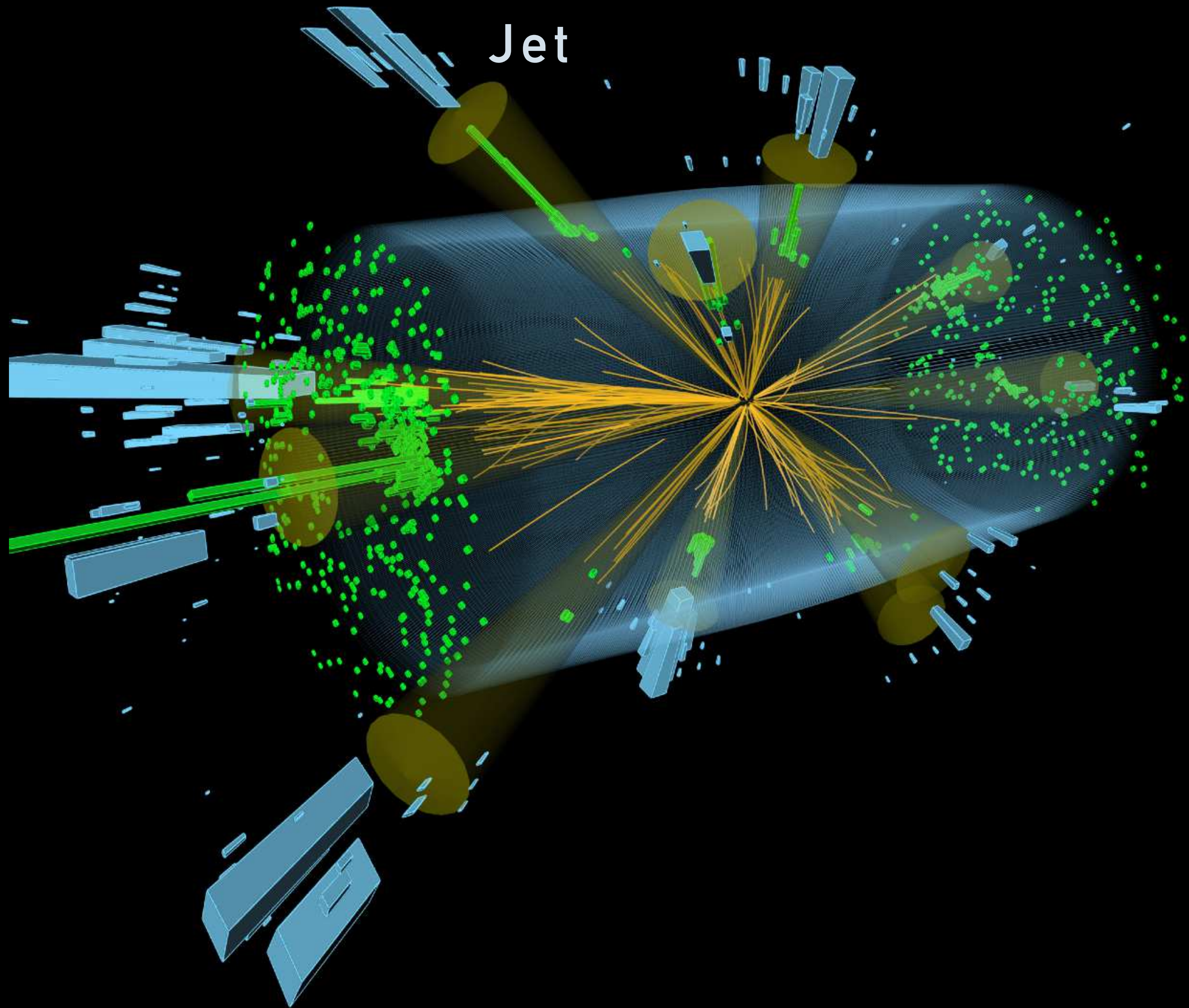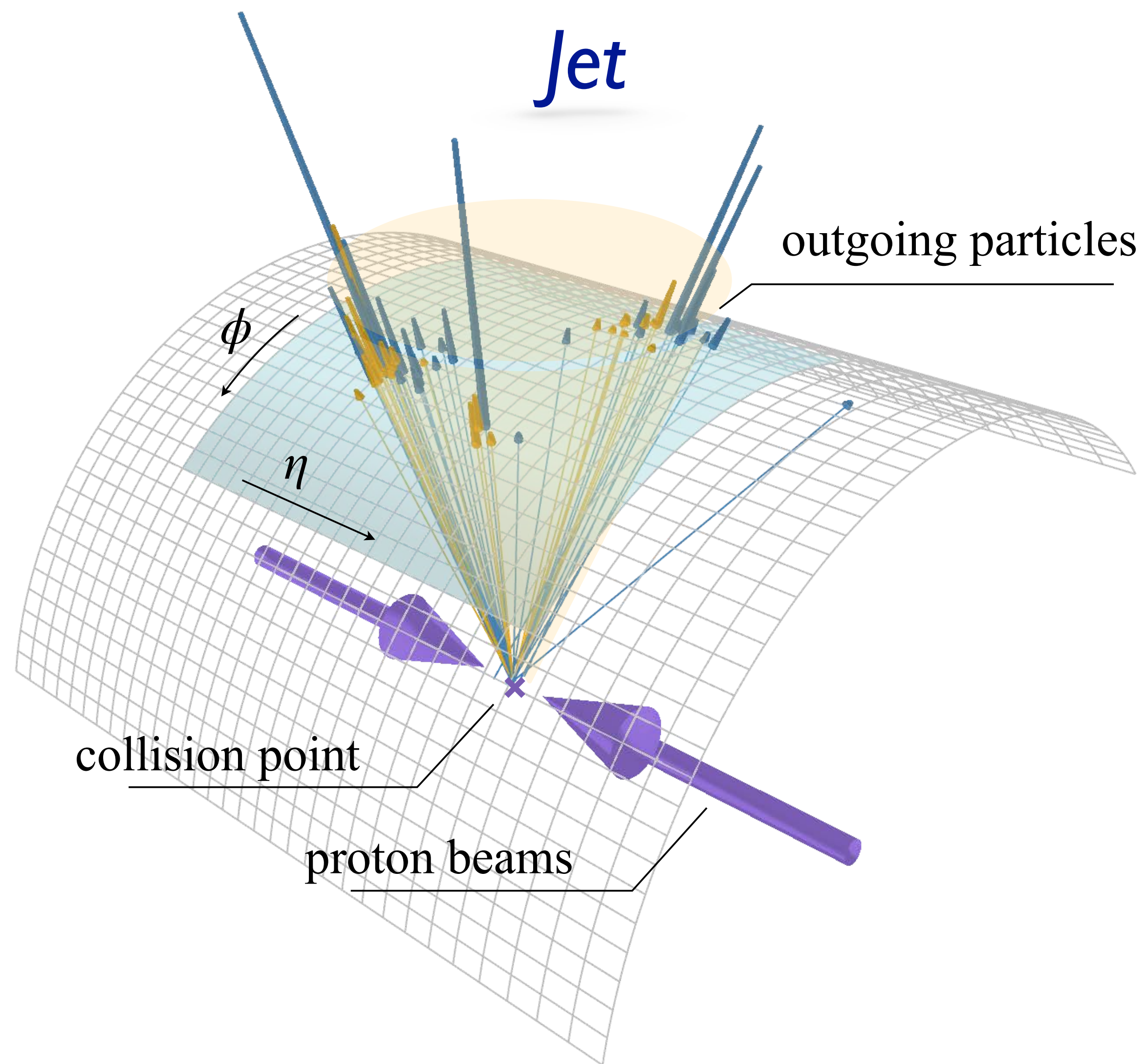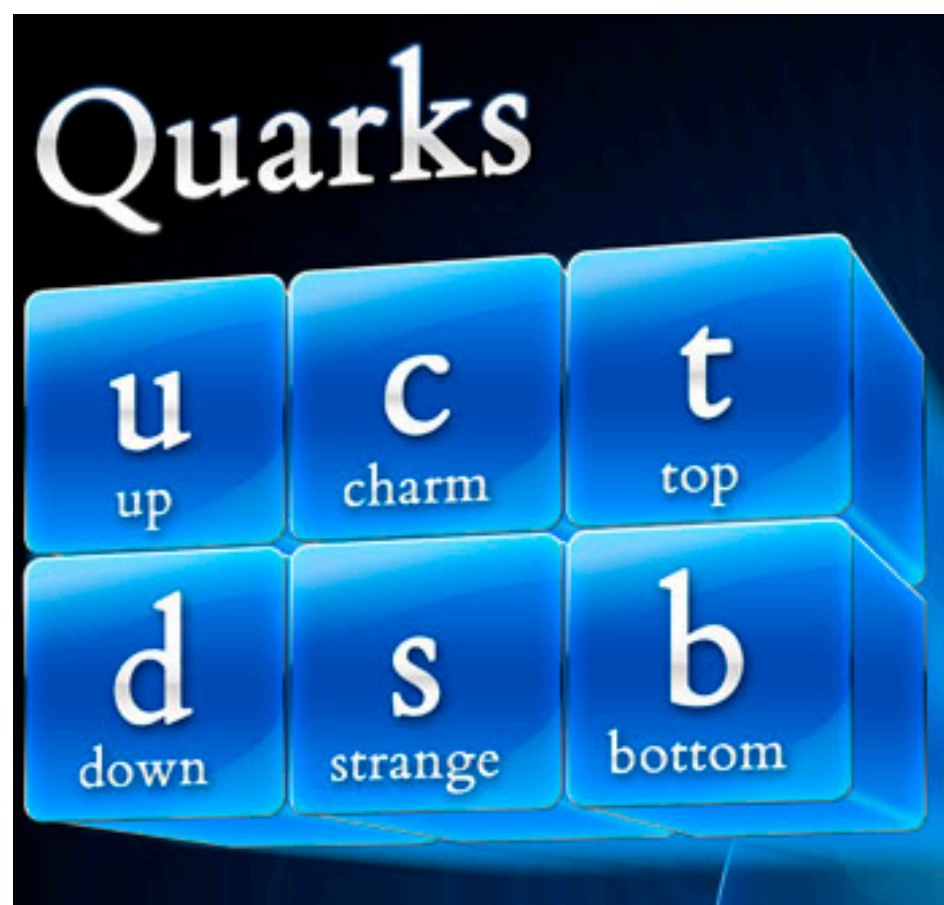We are also very keen on using this!

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!

GEN
SIM
DIGI+RECO+MINIAOD

hadronic
calorimeter

electromagnetic
calorimeter

tracker

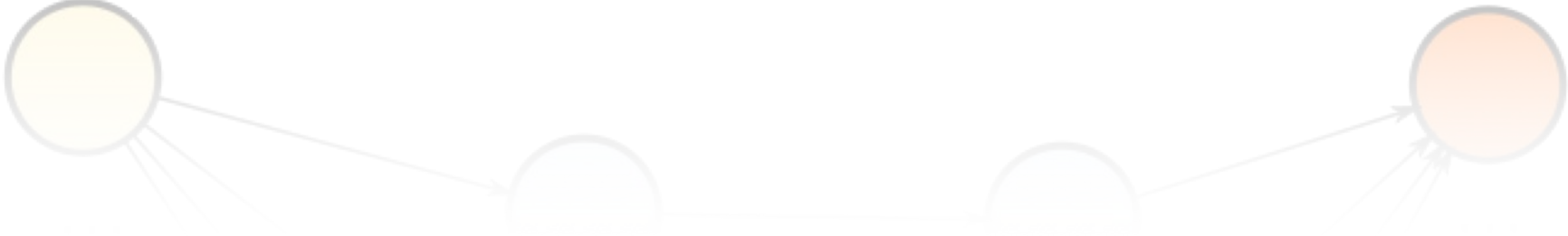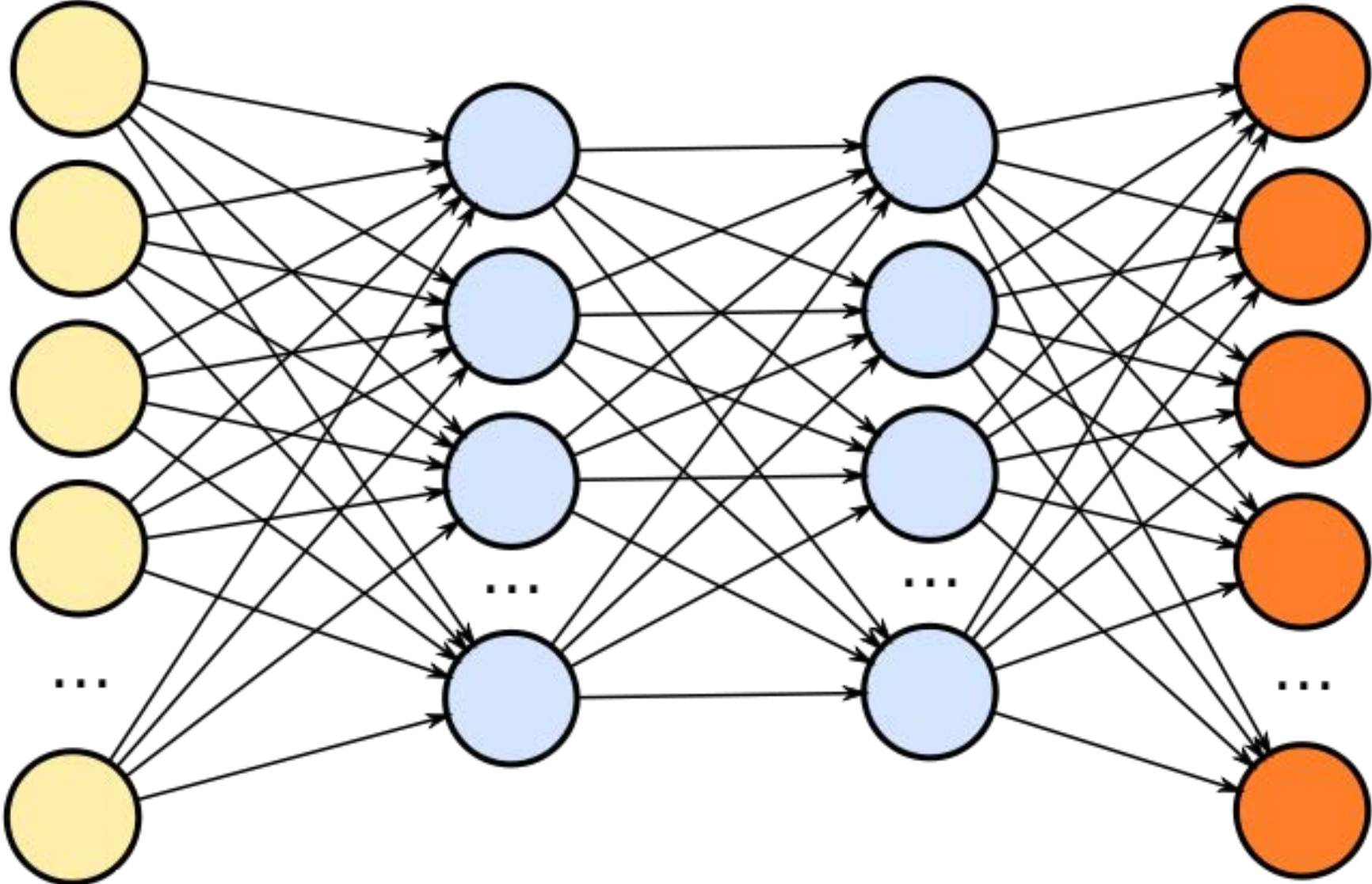Disk

10%

GEN
SIM
MINIAOD

AOD

81%

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

We have a lot of high quality
simulated data that we want to use

**?**

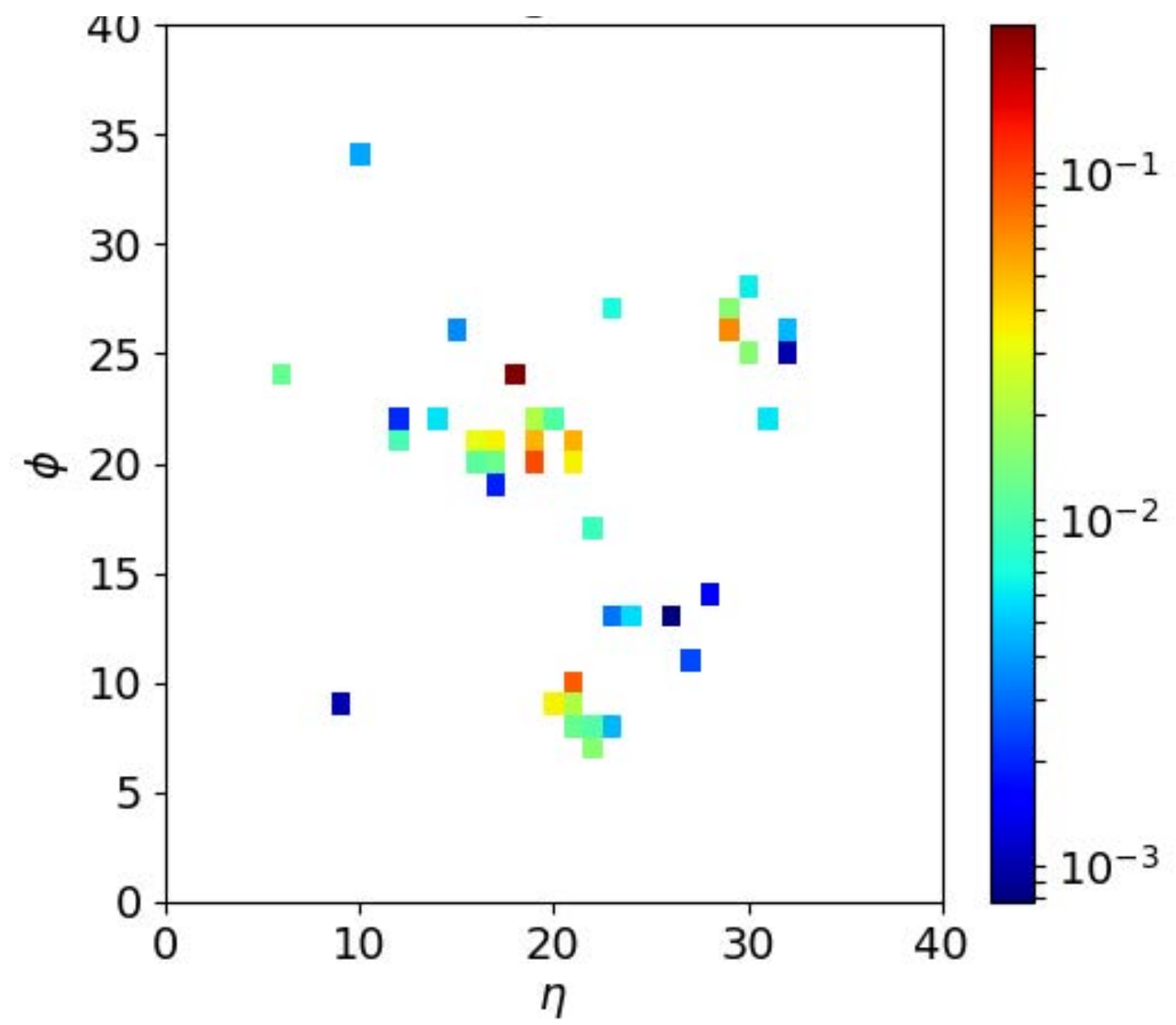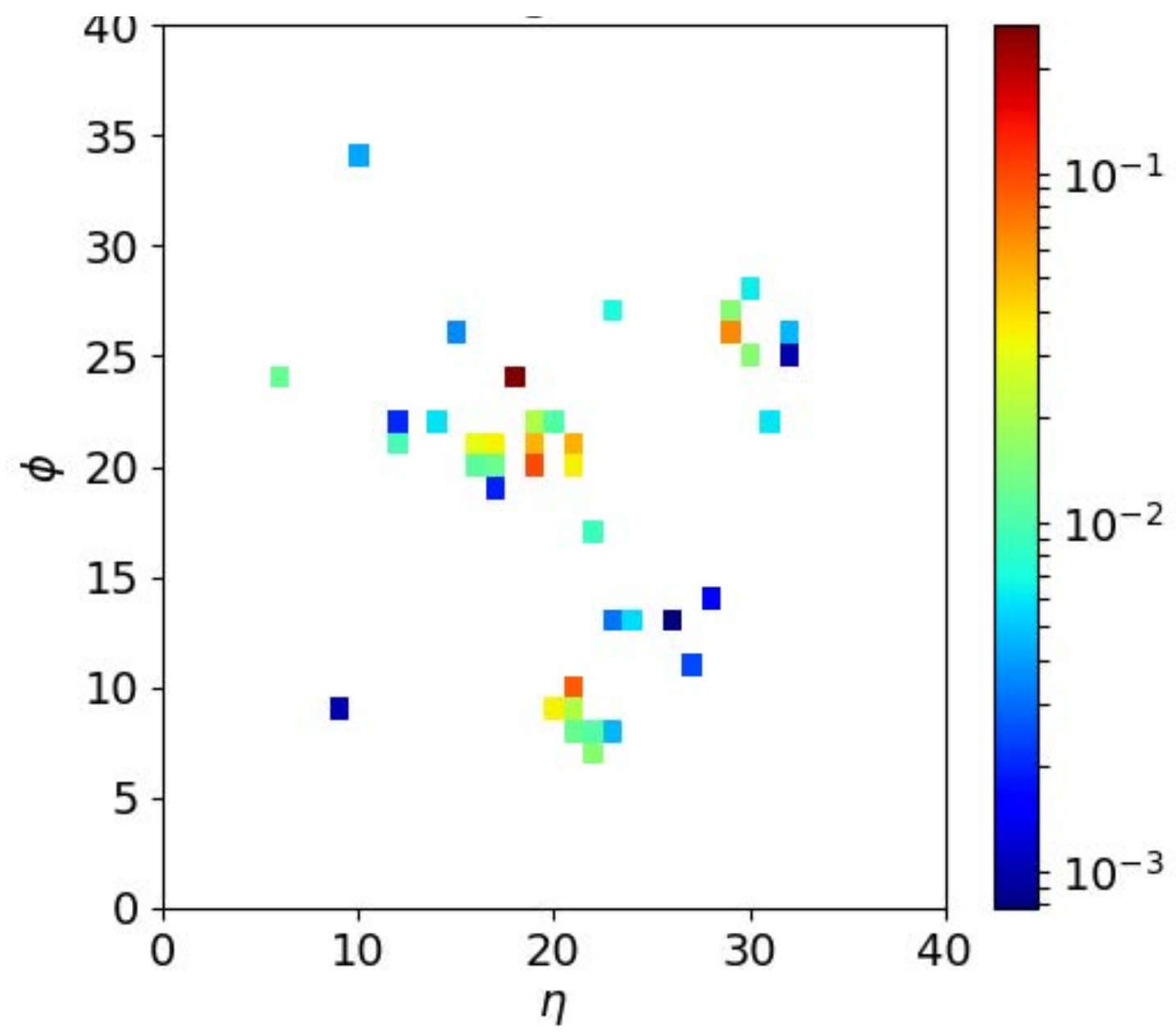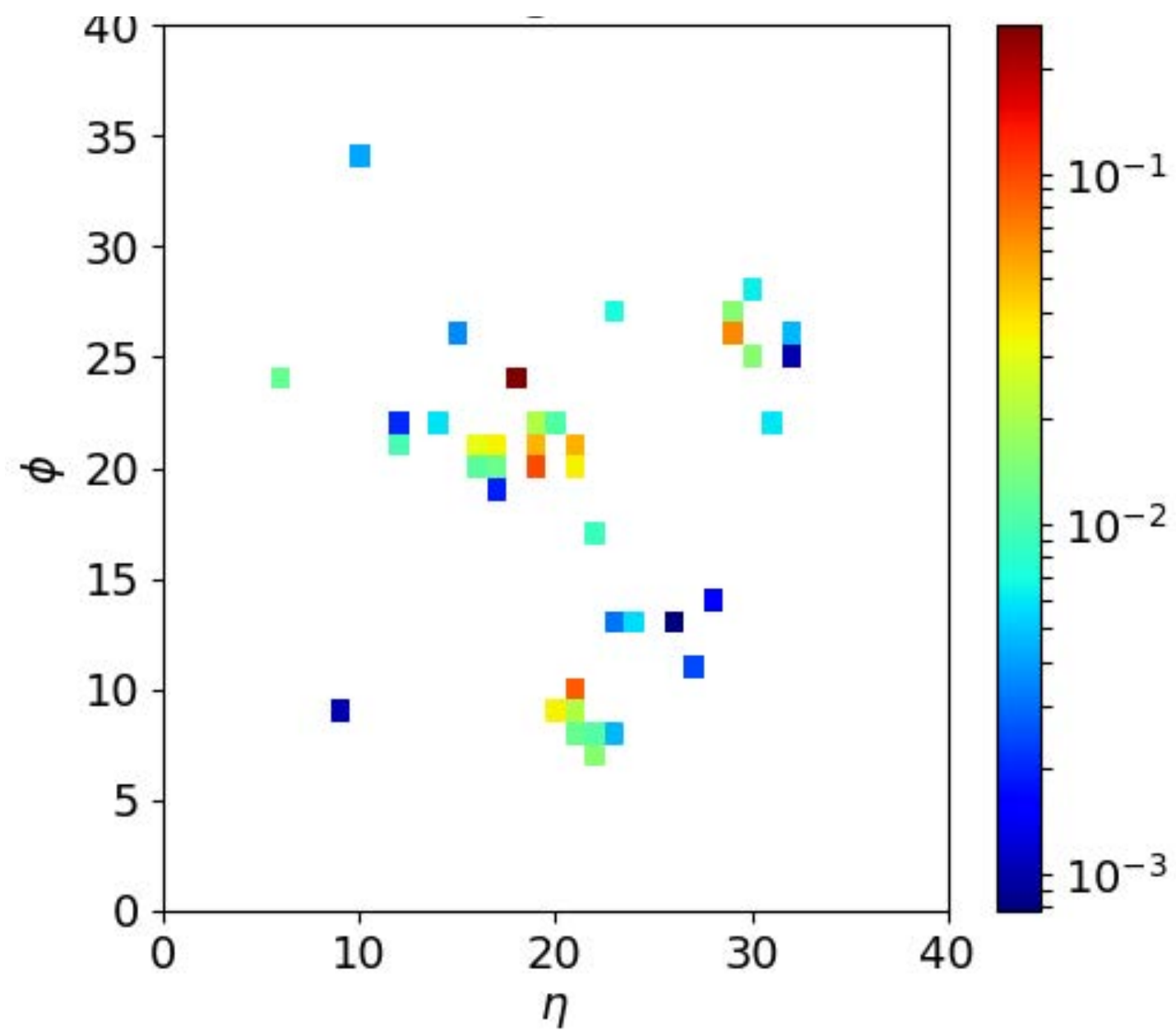Quarks
| u up | c charm | t top |
| d down | s strange | b bottom |

*Jet*

outgoing particles

$\phi$

$\eta$

collision point

proton beams

0.4

0.2

0.0

*rel*

0.2

0.4

outgoing particles

# Image

# Image

But… inhomogeneous geometry,
high sparsity

# Image

# Sequence

outgoing particles

*arXiv:1511.05190*

| 1 | 2 | 3 |
|---|---|---|
| How | are | you |

Output

MLP

LSTM States

Input Sequence

$S_1$ $S_2$ $\cdots$ $S_n$

$I_1$ $I_2$ $\cdots$ $I_n$

*arXiv:1607.08633*

# Image

# Sequence



outgoing particles

*arXiv:1511.05190*

*arXiv:1607.08633*

But... permutation-invariance

# Jet



*Permutation invariance*

=

# Sequence

| 1 | 2 | 3 |
|---|---|---|
| How | are | you |

≠

| 1 | 2 | 3 |
|---|---|---|
| are | How | you |

# Image

# Sequence

outgoing particles

outg

$\phi$

$\eta$

collision point

proton beams

**Output**

**MLP**

$S_n$ LSTM States

$I_n$ Input Sequence

# Point Cloud

$\phi$

$\eta$

*arXiv:1511.05190*
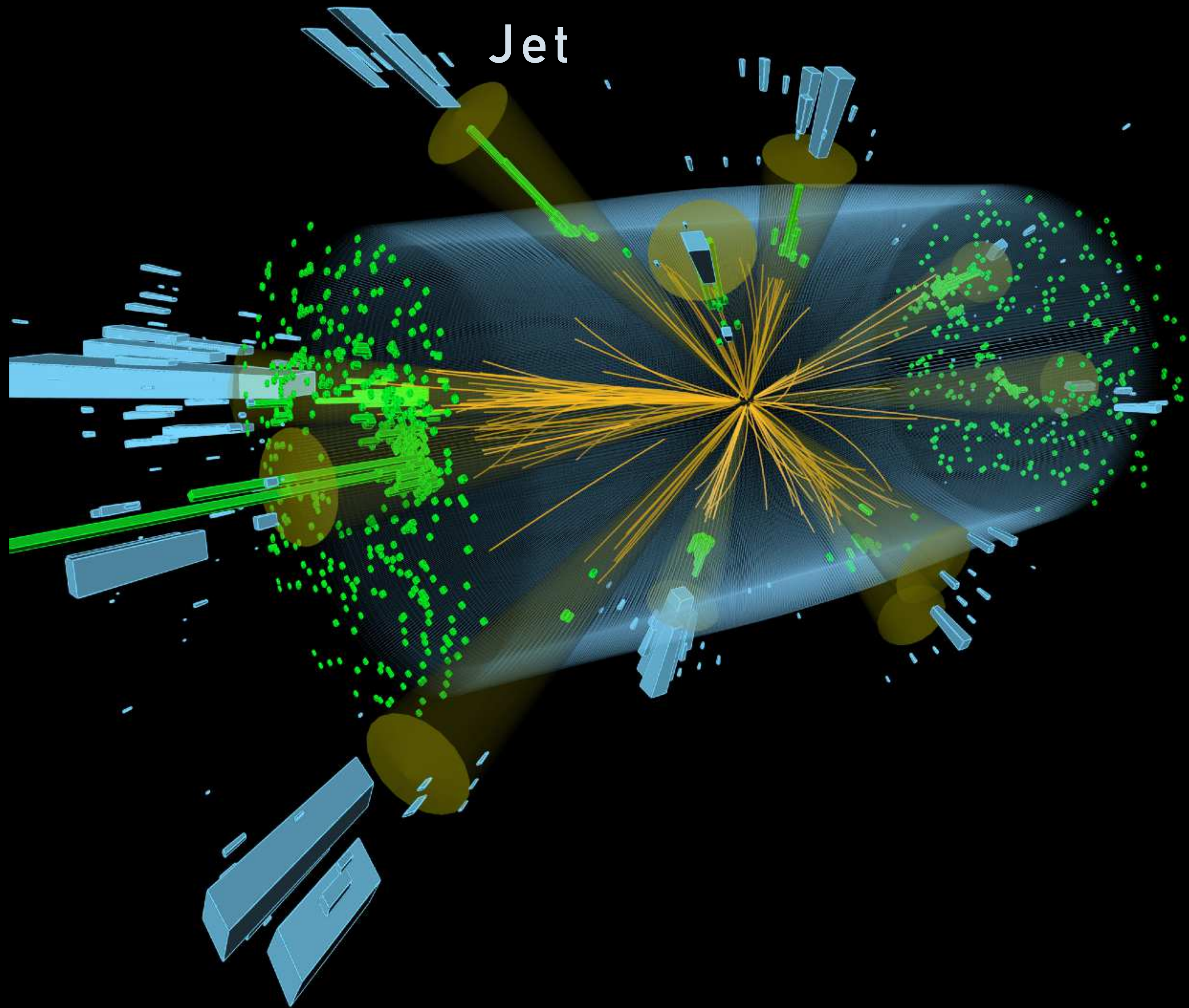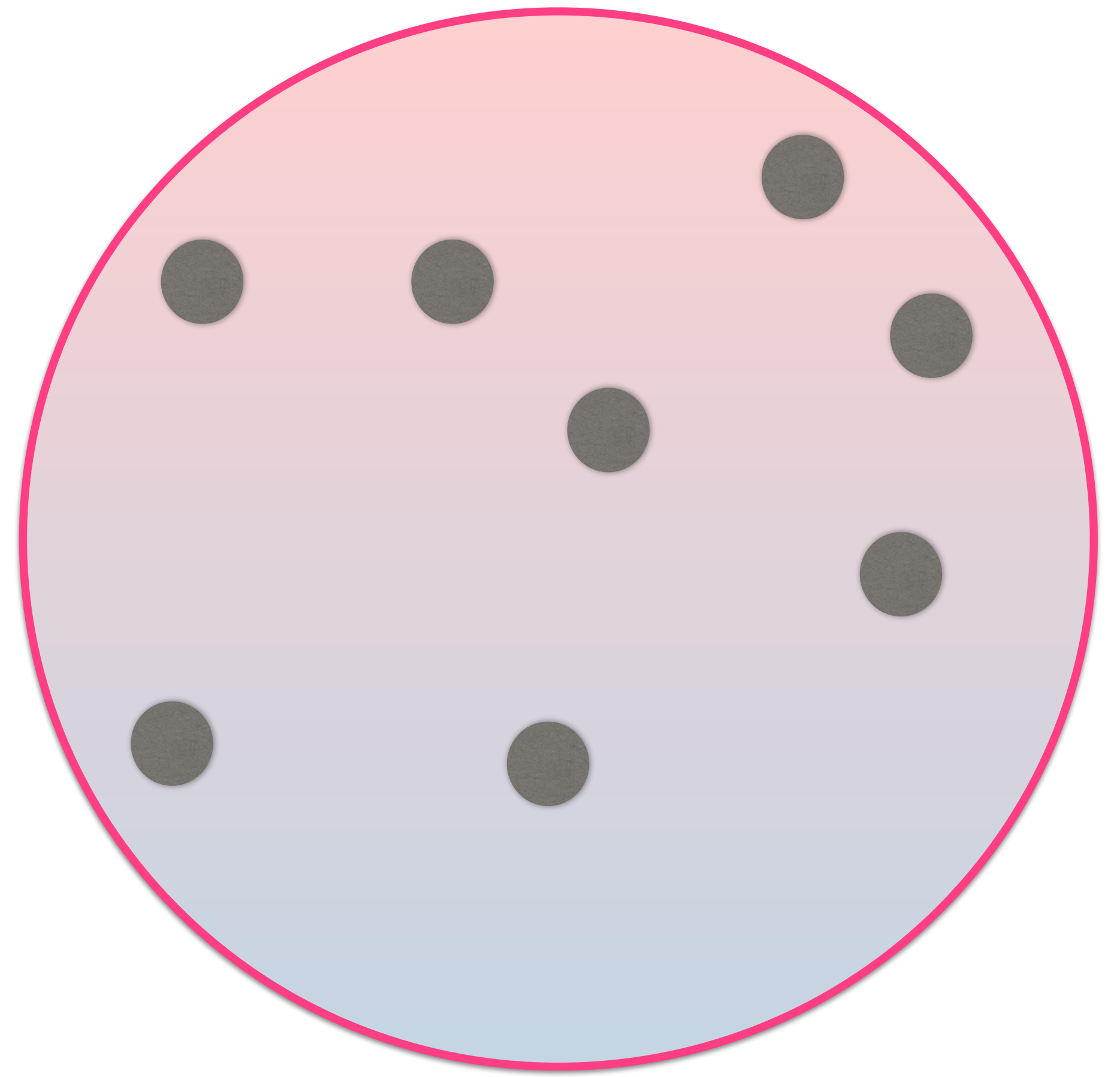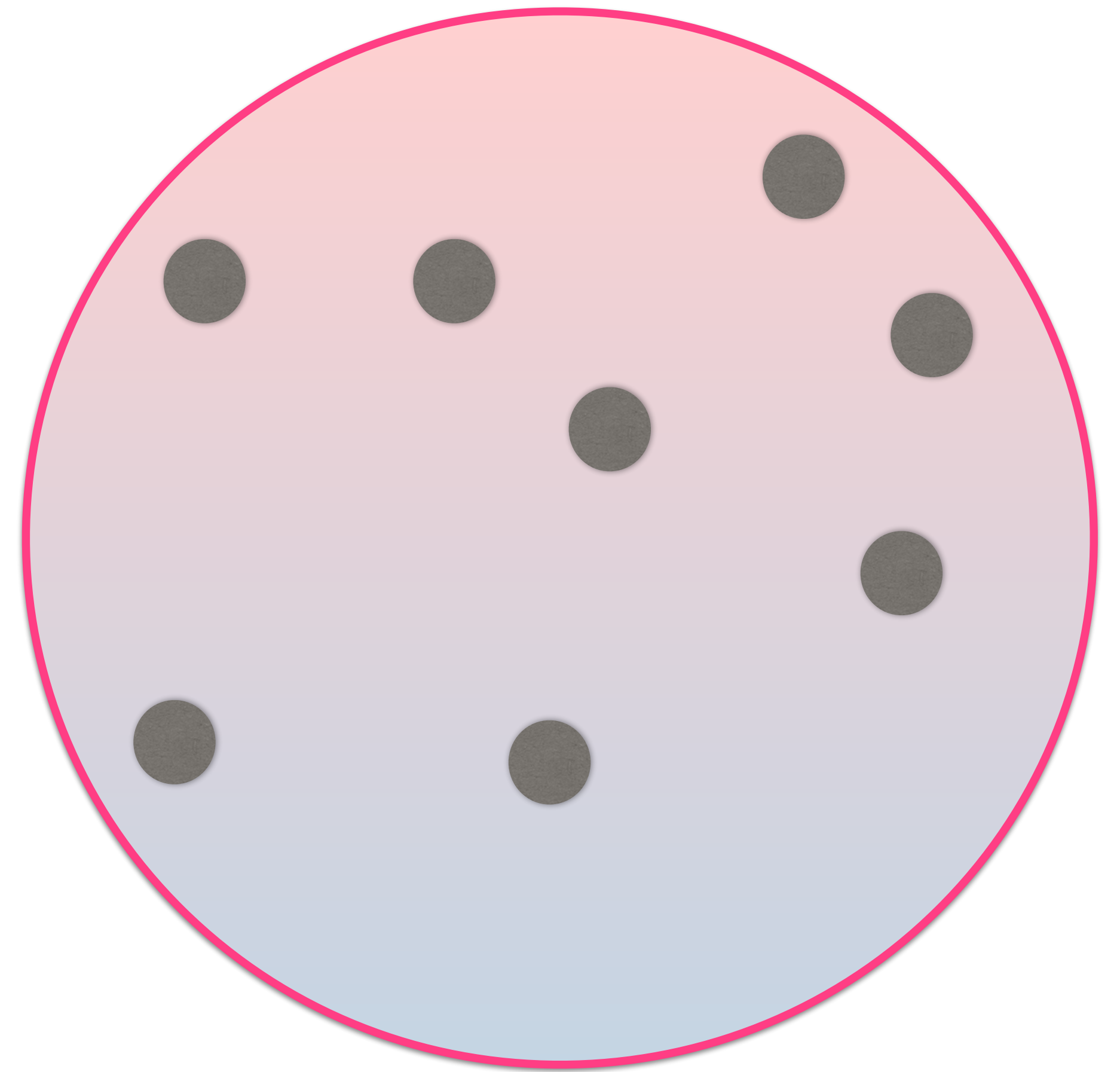
*arXiv:1607.08633*

*PRD:101.056019*

# Graph Neural Networks

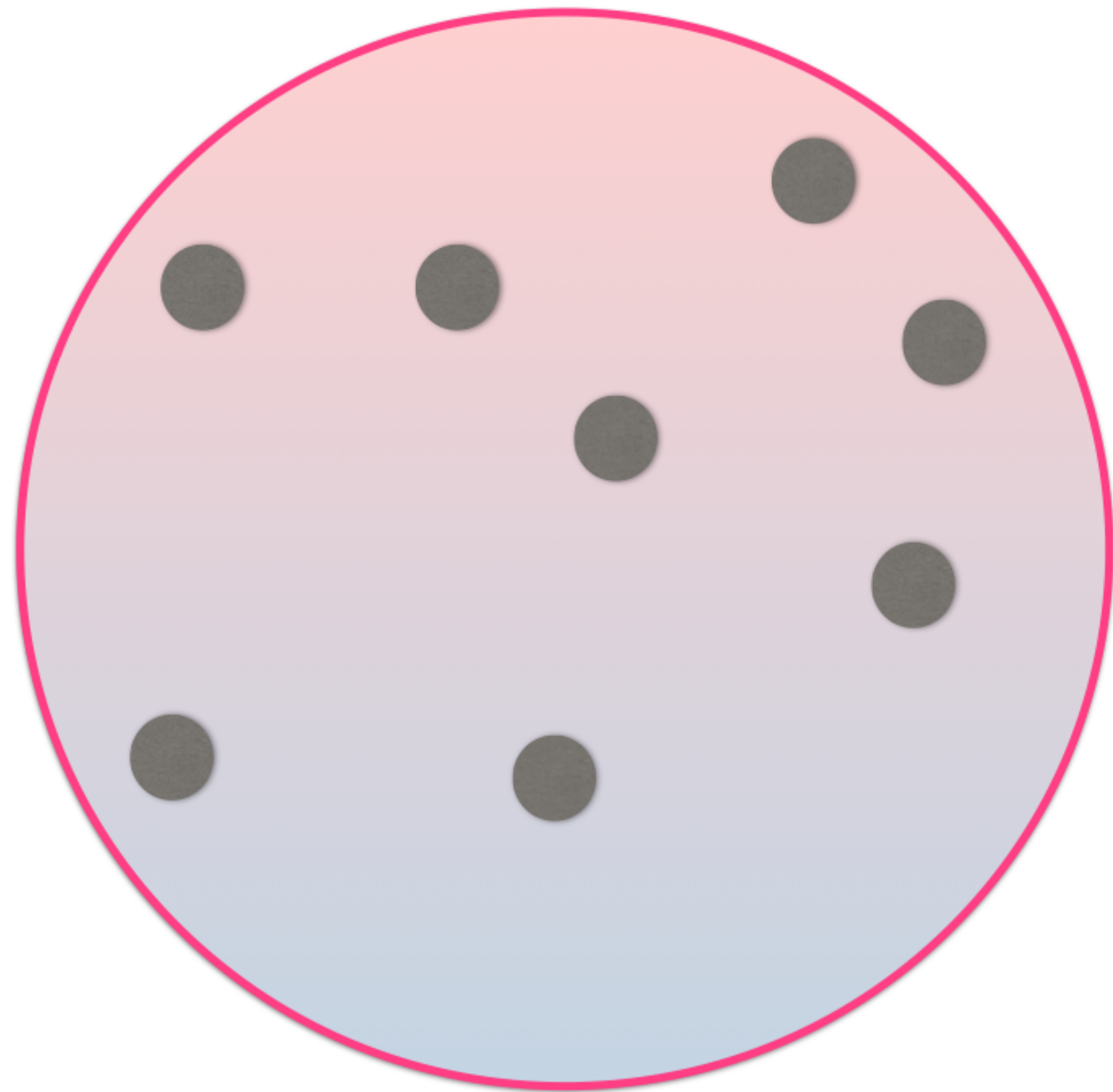Properties of physics data:
- Measurements distributed in space (and time) **irregularly**
- **Sparse** (most detector channels are empty), but pockets of density
- Complex **interdependencies** between measurements
- Physics "objects" composed of **multiple measurements**
- Inherent **symmetries** (Lorentz boosts, rotational)

Graph (or point cloud) embedding of the data can handle these properties!

**Graph (global) features $u$: jet mass**

**Node features $v_i$:**

$$p = [E, p_x, p_y, p_z] \equiv [p_{\mathrm{T}}, \eta, \phi, m]$$

**Edge features** $e_{i,j}$

$$\Delta R = \sqrt{\Delta \eta^2 + \Delta \phi^2}$$

**u**

$$m = \sqrt{\sum_{i \in \mathrm{jet}} E_i^2 - p_{x,i}^2 - p_{y,i}^2 - p_{z,i}^2}$$

# Graph Neural Networks

# Graph Neural Networks

$v_1' =$

$\mathbf{e}'_{1\to5} = \mathbf{MLP}(\vec{\mathbf{v}}_1, \vec{\mathbf{v}}_5)$

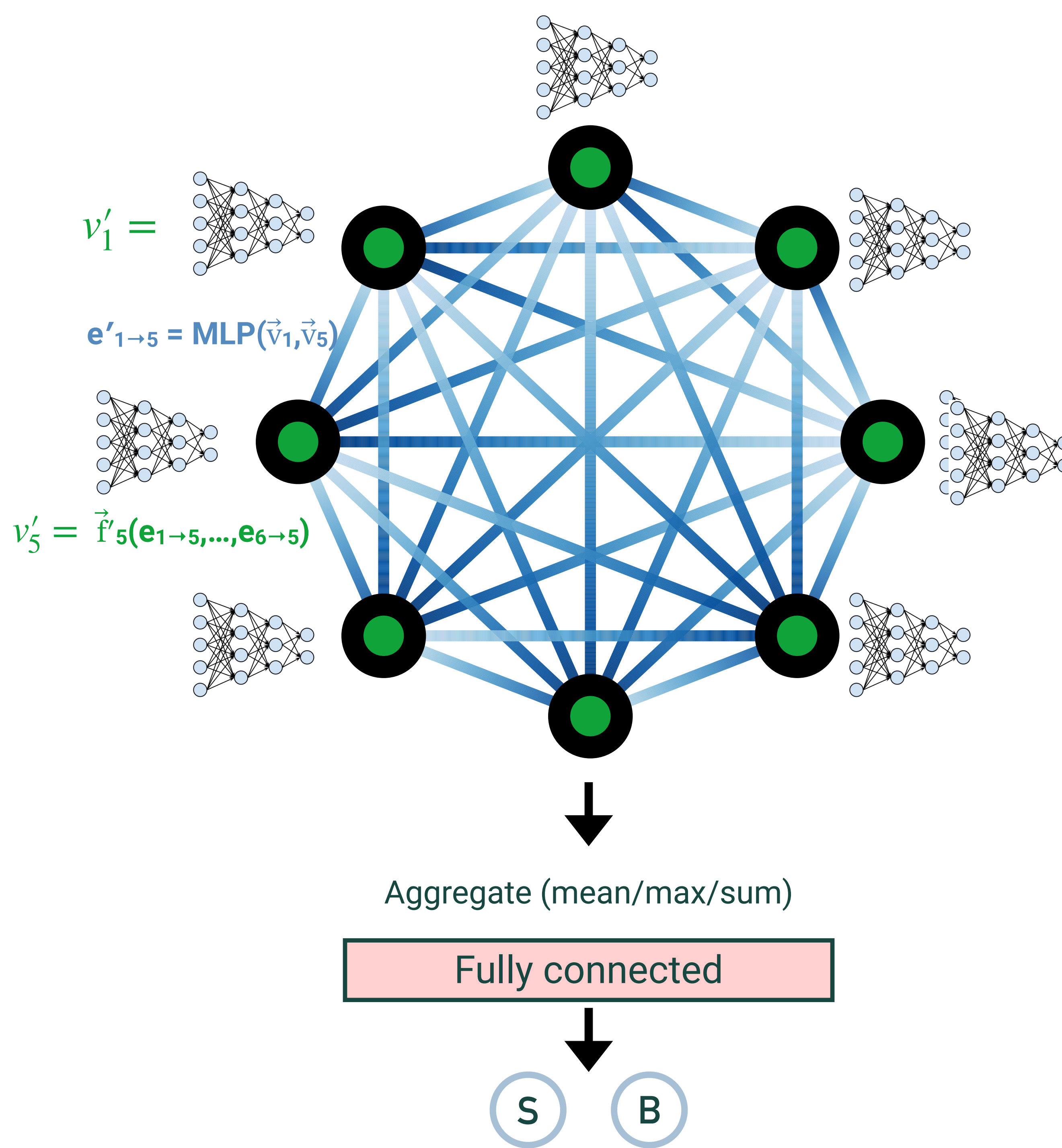$v_5' = \vec{f}'_5(\mathbf{e}_{1\to5},...,\mathbf{e}_{6\to5})$

**Want to create "new features" on the nodes,
edges, or the full graph with multiple iterations:**

- Create a new representation for each part of
  the graph
- These "updates" are usually DNNs!

$$\mathbf{e}'_k = \phi^e\left(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}, \mathbf{u}\right)$$
$$\mathbf{v}'_i = \phi^v\left(\bar{\mathbf{e}}'_i, \mathbf{v}_i, \mathbf{u}\right)$$
$$\mathbf{u}' = \phi^u\left(\bar{\mathbf{e}}', \bar{\mathbf{v}}', \mathbf{u}\right)$$

DNNs to be trained!

$$v'_1 =$$

$$e'_{1\to5} = MLP(\vec{v}_1, \vec{v}_5)$$

$$v'_5 = \vec{f}'_5(e_{1\to5}, ..., e_{6\to5})$$

Aggregate (mean/max/sum)

Fully connected

S    B

$v'_1 =$

$\mathbf{e'_{1 \to 5} = MLP(\vec{v}_1, \vec{v}_5)}$

$v'_5 = \vec{f'}_5(\mathbf{e_{1 \to 5}, ..., e_{6 \to 5}})$

Aggregate (mean/max/sum)

Fully connected

$e_{ij} = MLP(v_i, v_j, v_{ij})$

$\mathbf{v_i}$

$\mathbf{v_{ij}}$

$\mathbf{v_j}$

Lund-like features

S    B

$$v_1' =$$

$$\mathbf{e'}_{1\to5} = \mathbf{MLP}(\vec{v}_1, \vec{v}_5)$$

$$v_5' = \vec{f'}_5(\mathbf{e}_{1\to5}, ..., \mathbf{e}_{6\to5})$$
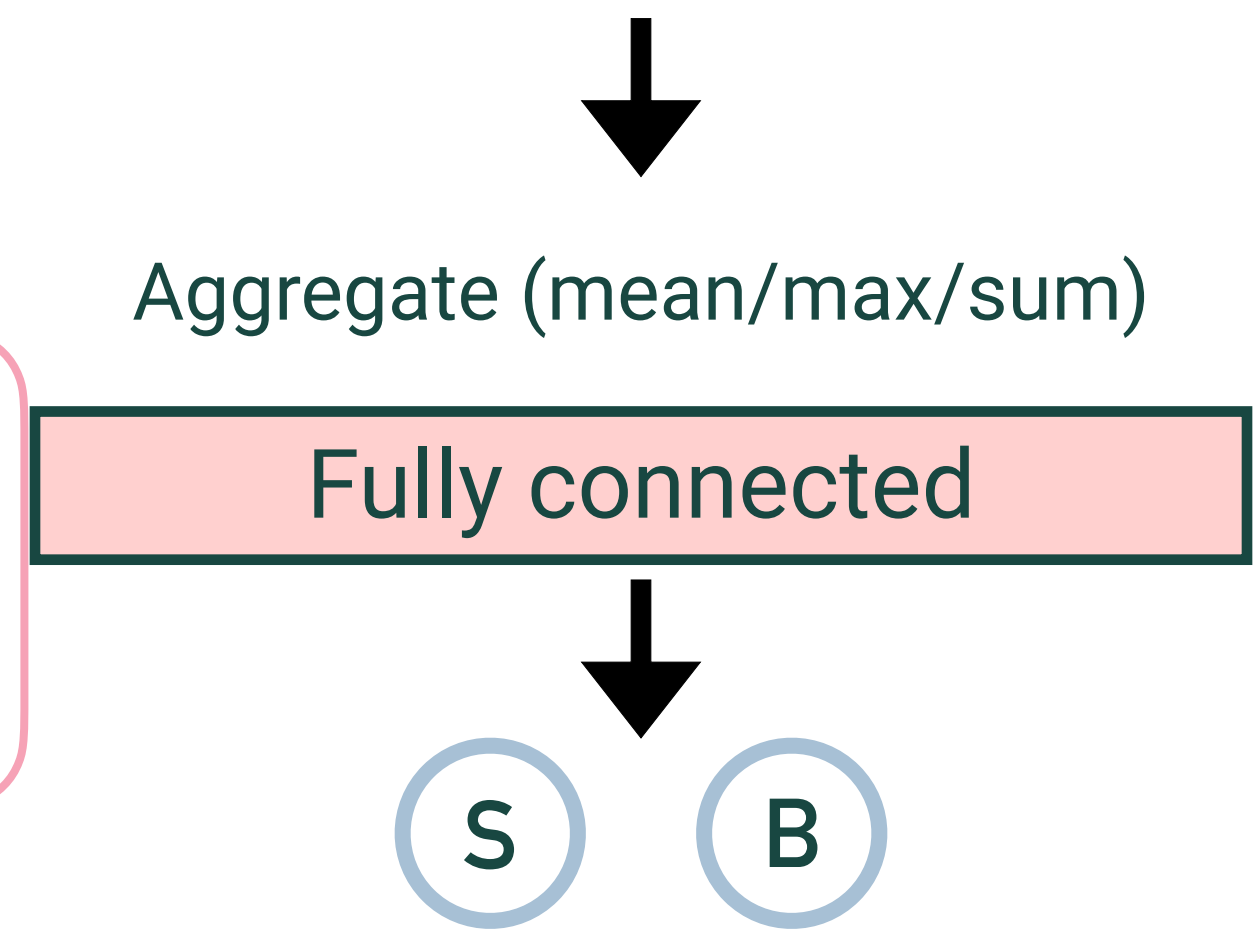
SOTA: GNNs acting on point cloud data
- ParticleNet (GNN on point cloud)
  LundNet (GNN, Lund plane)
  ABCNet (GNN, attention)
  Point Cloud Transformers (transformer, attention)
  ParticleNeXt (GNN, attention, Lund)
  ParT (transformer, attention)
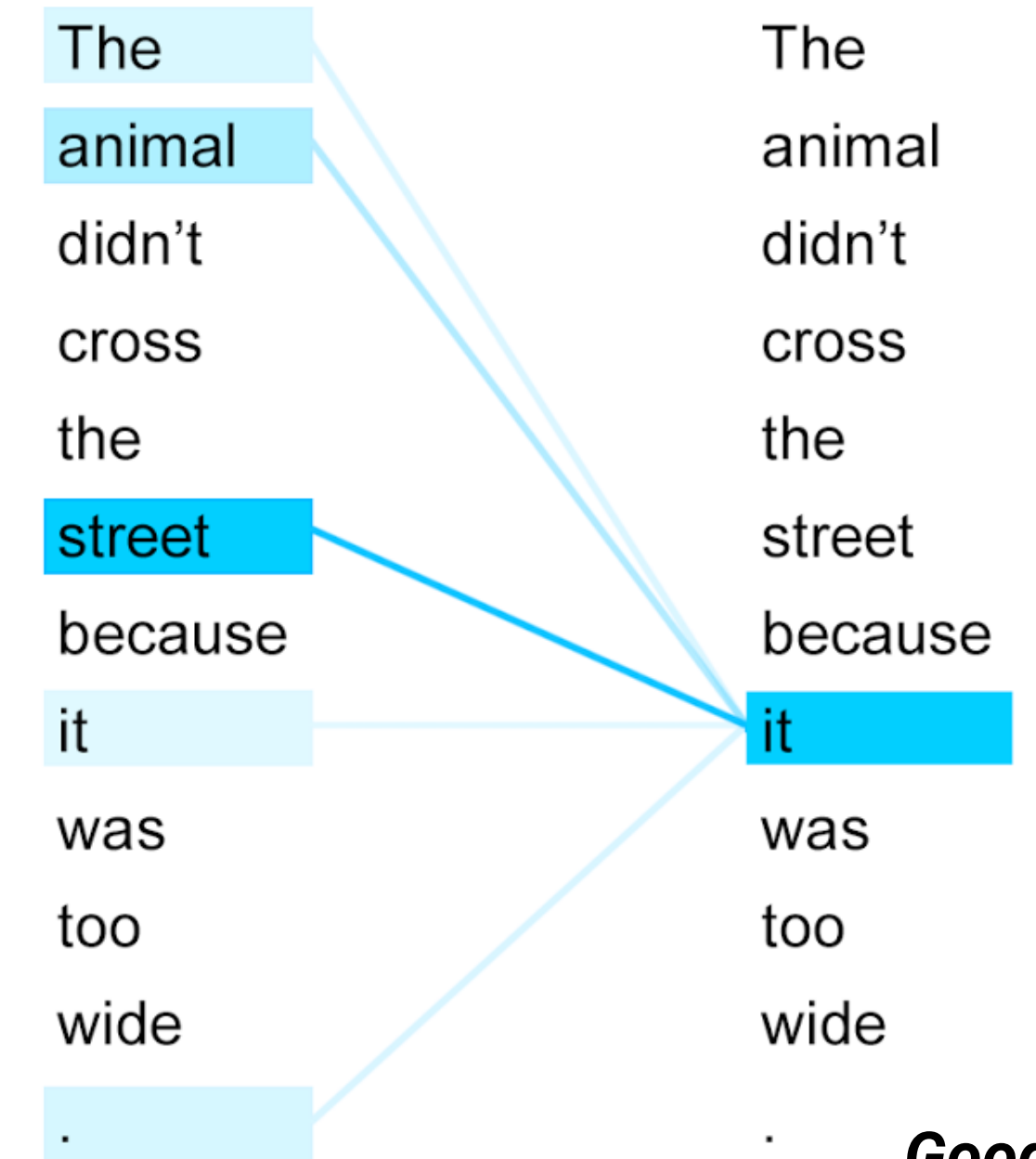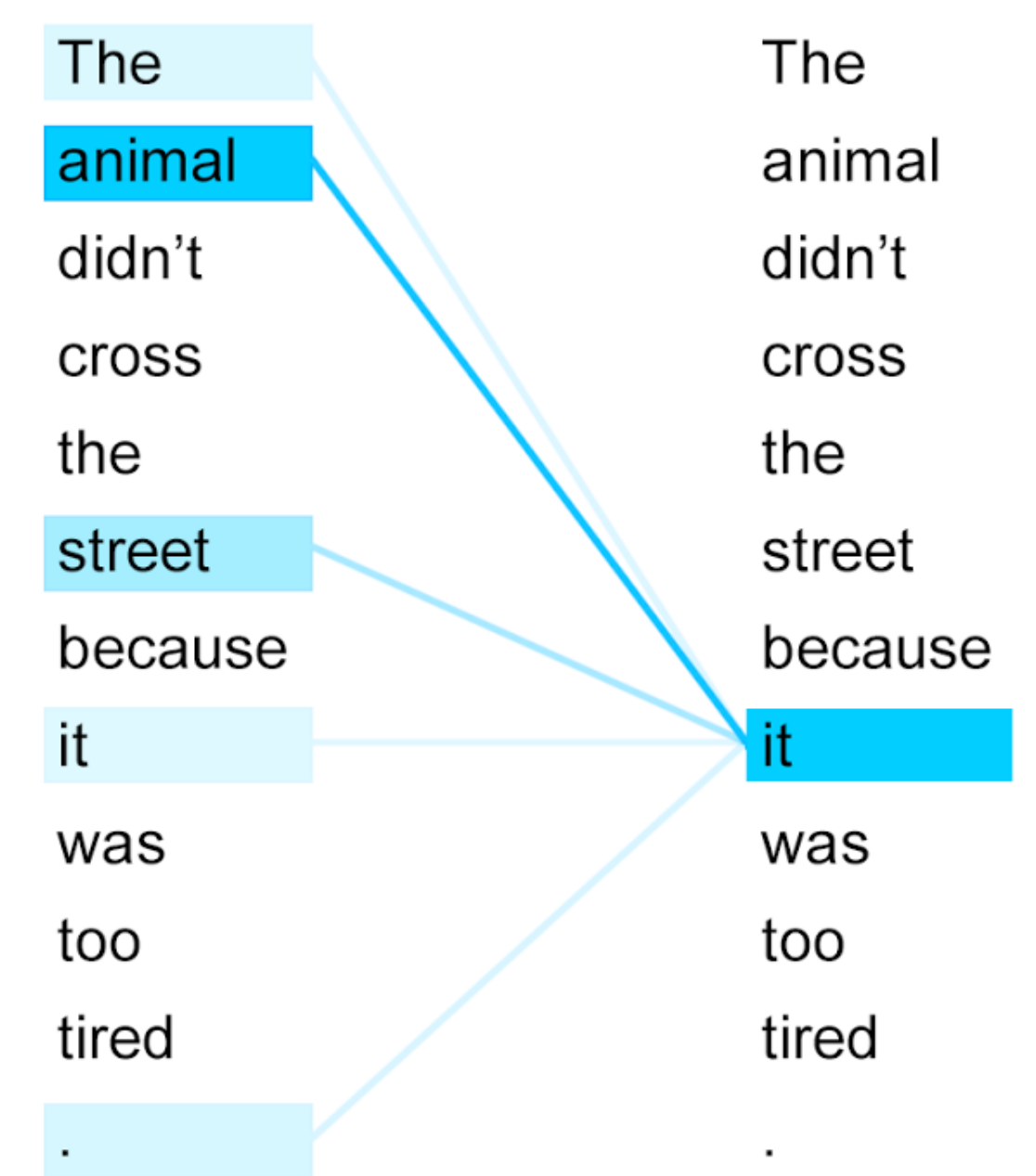
Aggregate (mean/max/sum)

Fully connected

S    B

# Transformers and (self-)attention

**(Self-)Attention**
- Allows inputs to interact with each other ("self") and find out who they should pay more attention to ("attention").
- Outputs: aggregates of interactions and attention scores

# Transformers and (self-)attention

**(Self-)Attention**
- Allows inputs to interact with each other ("self") and find out who they should pay more attention to ("attention").
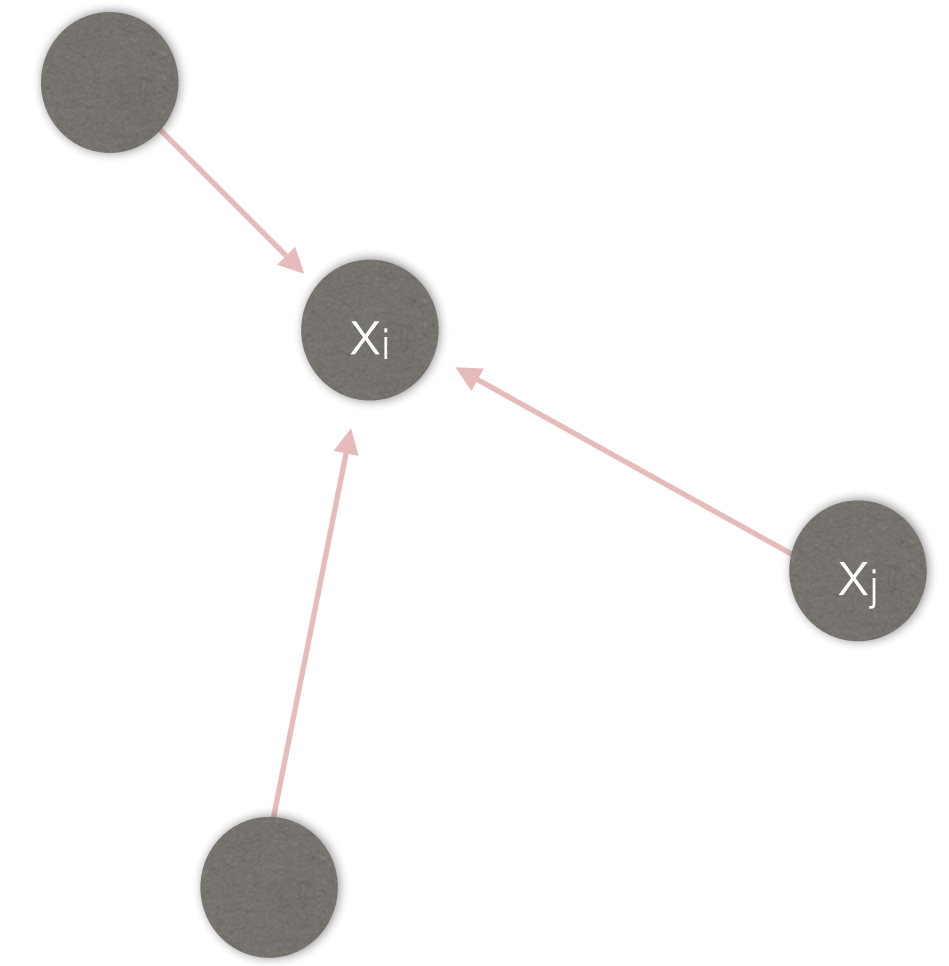- Outputs: aggregates of interactions and attention scores

Weighted sum over all input vectors:

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

Weight (how related inputs are):
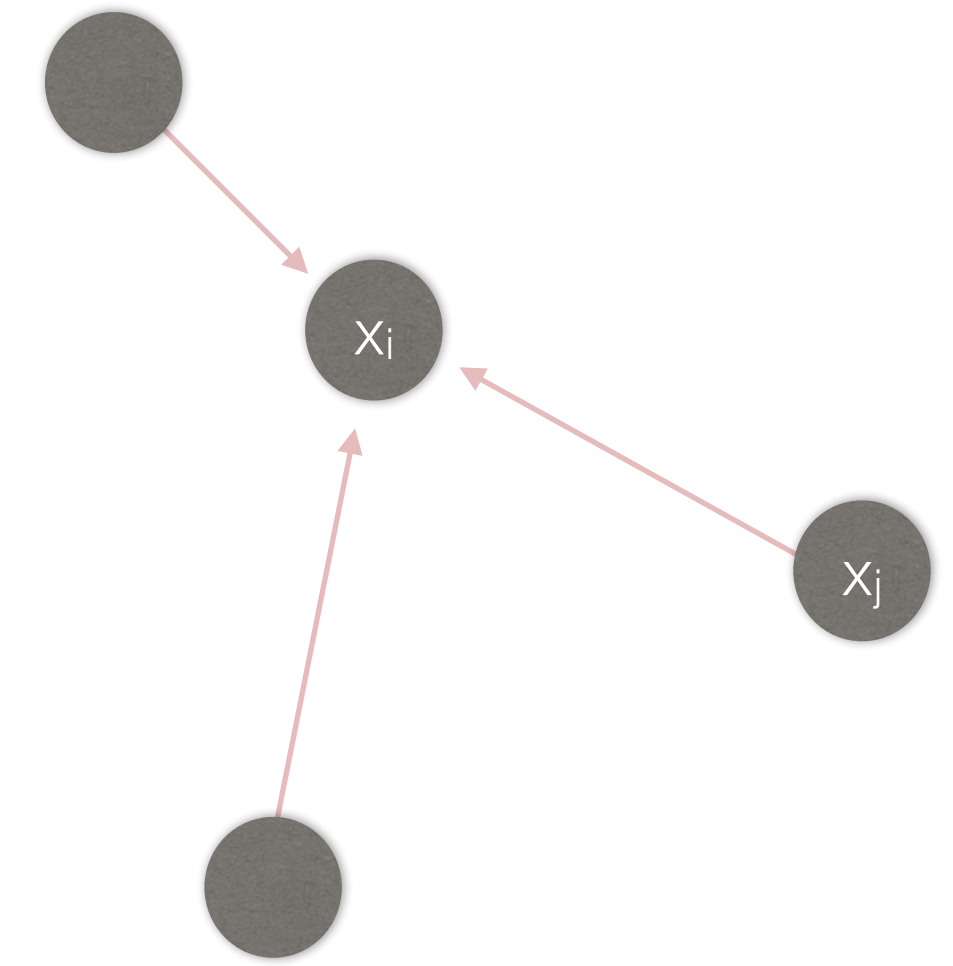
$$w'_{ij} = \mathbf{x}_i{}^\mathsf{T} \mathbf{x}_j$$

Map to [0,1]:

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

# Transformers and (self-)attention

(Self-)Attention
- Allows inputs to interact with each other ("self") and find out who they should pay more attention to ("attention").
- Outputs: aggregates of interactions and attention scores

Weighted sum over all input vectors:

$$\mathbf{y}_i = \sum_j w_{ij}\mathbf{x}_j$$

Weight (how related inputs are):

$x_j \rightarrow MLP(x_J)$
$x_i \rightarrow MLP(x_i)$

$$w'_{ij} = \mathbf{x}_i{}^{\mathsf{T}}\mathbf{x}_j$$

Map to [0,1]:

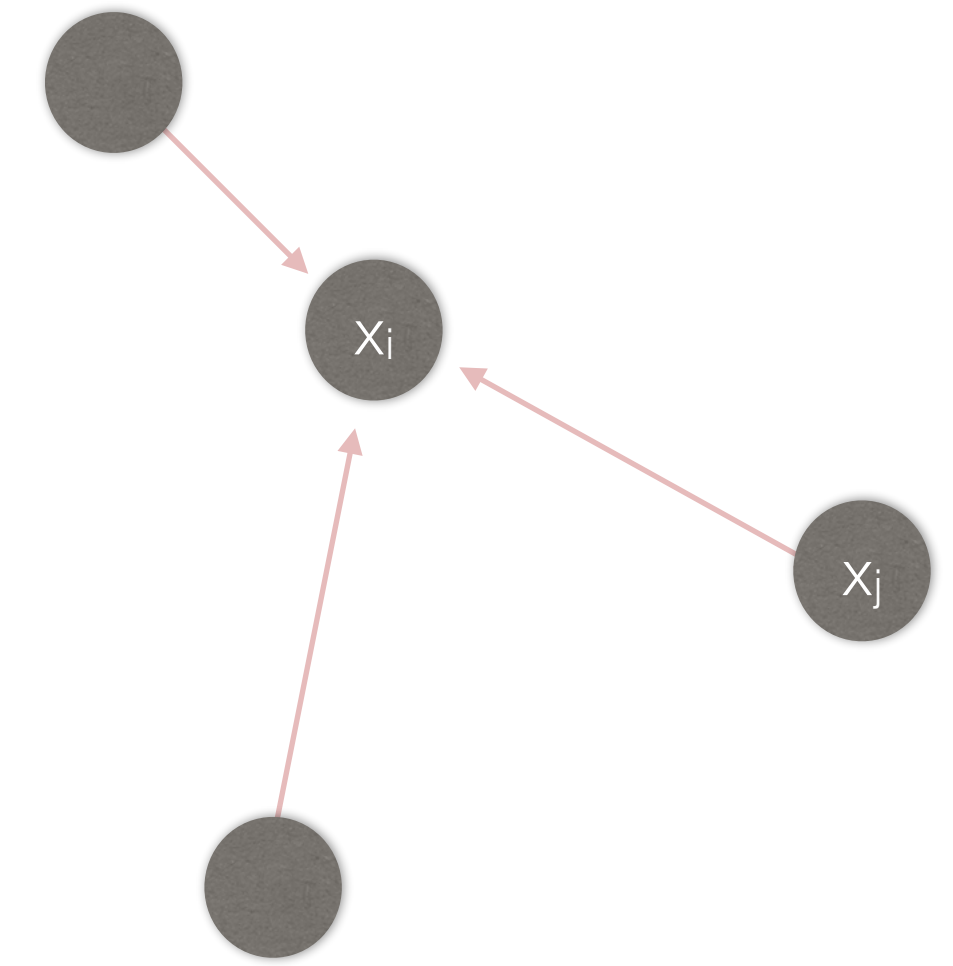$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

# Transformers and (self-)attention
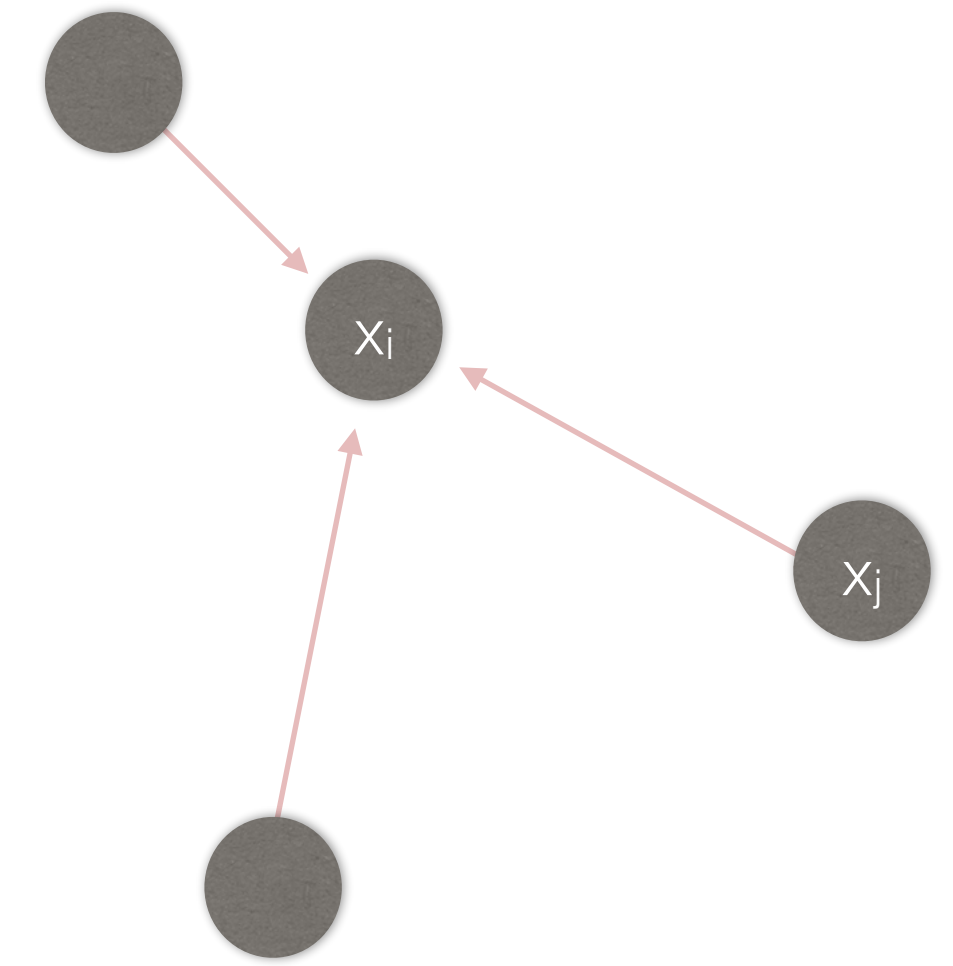


**(Self-)Attention**
- Allows inputs to interact with each other ("self") and find out who they should pay more attention to ("attention").
- Outputs: aggregates of interactions and attention scores

**Attention weights: weighted importance between each pair of particles**
- Determine relationship between all particles of point cloud
- Jet features become parameters of the model
- Several attention layers → different important features (multi-head attention)

Weighted sum over all input vectors:

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

Weight (how related inputs are):

$x_j \rightarrow MLP(x_J)$
$x_i \rightarrow MLP(x_i)$

$$w'_{ij} = \mathbf{x}_i{}^\mathsf{T} \mathbf{x}_j$$

Map to [0,1]:

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

# Transformers and (self-)attention



(Self-)Attention
- Allows inputs to interact with each other ("self") and find out who they should pay more attention to ("attention").
- Outputs: aggregates of interactions and attention scores

Attention weights: weighted importance between each pair of particles
- Determine relationship between all particles of point cloud
- Jet features become parameters of the model
- Several attention layers → different important features (multi-head attention)

Transformer:
- Only set of interaction between units is self-attention!

Weighted sum over all input vectors:

$$\mathbf{y}_i = \sum_j w_{ij} \mathbf{x}_j$$

Weight (how related inputs are):

$x_j \rightarrow MLP(x_J)$
$x_i \rightarrow MLP(x_i)$

$$w'_{ij} = \mathbf{x}_i^{\mathsf{T}} \mathbf{x}_j$$

Map to [0,1]:

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

# Transformers and (self-)attention

(Self-)Attention
- Allows inputs to interact with each other ("self") and find out who they should pay more attention to ("attention").
- Outputs: aggregates of interactions and attention scores

Attention weights: weighted importance between each pair of particles
- Determine relationship between all particles of point cloud
- Jet features become parameters of the model
- Several attention layers → different important features (multi-head attention)

Transformer:
- Only set of interaction between units is self-attention!

**Example prompt**

```
Rigor [adj.]
Something for scientists to aspire to, a state of mind
that would not be required if scientists could be trusted
to do their job.
```

[View next definition]

**GPT-3's output: 1 of 10**

```
The Literature [noun]
A name given to other people's published papers, referred
to by scientists without actually reading them.
```

Gwern.net

Weight (how related inputs are):

$x_j \rightarrow MLP(x_J)$
$x_i \rightarrow MLP(x_i)$

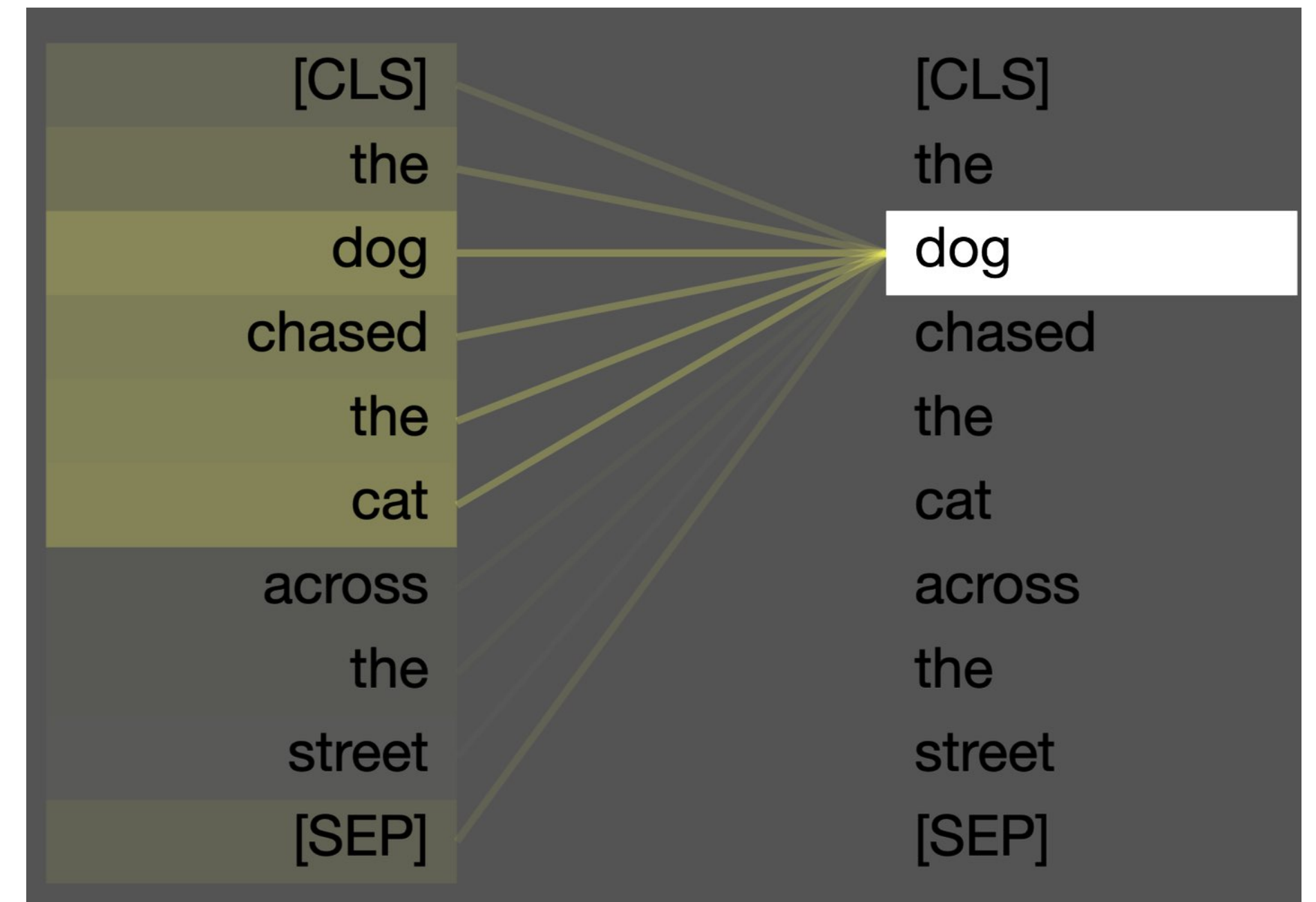$$w'_{ij} = \mathbf{x}_i{}^{\mathsf{T}} \mathbf{x}_j$$

Map to [0,1]:

$$w_{ij} = \frac{\exp w'_{ij}}{\sum_j \exp w'_{ij}}$$

# Transformers

Query, Key and Value: How self-attention is implemented
- **The query**: dog
  I'm looking for verbs, adjectives related to me
- **The key**: every word in the sentence!
  I'm a noun, an adjective or a verb
  (What am I? What features do I posses in relation to the sentence?)
- **The value:** the meaning of this word in general not specifically for this sentence (What are my embeddings? What's the semantic information I posses?)
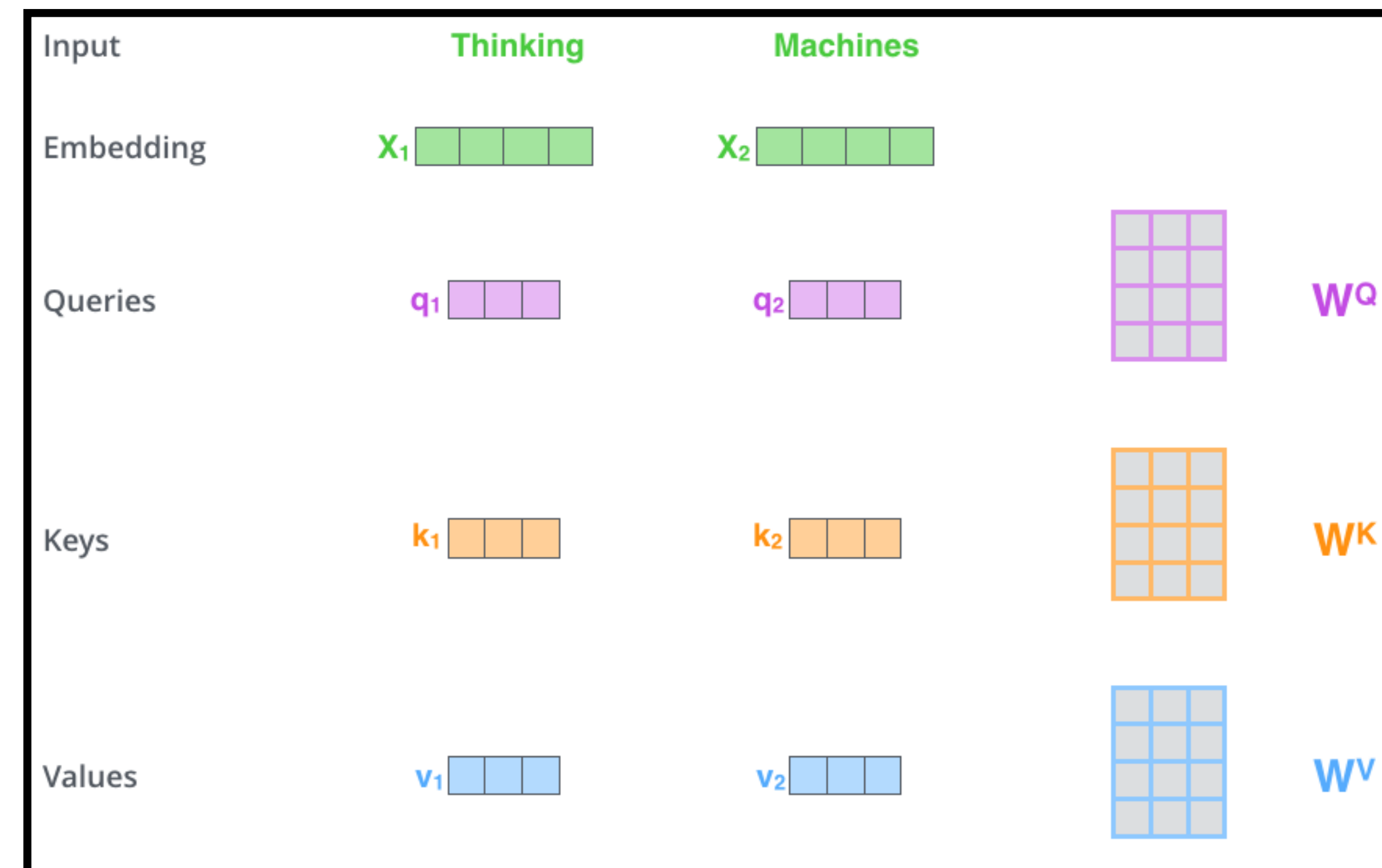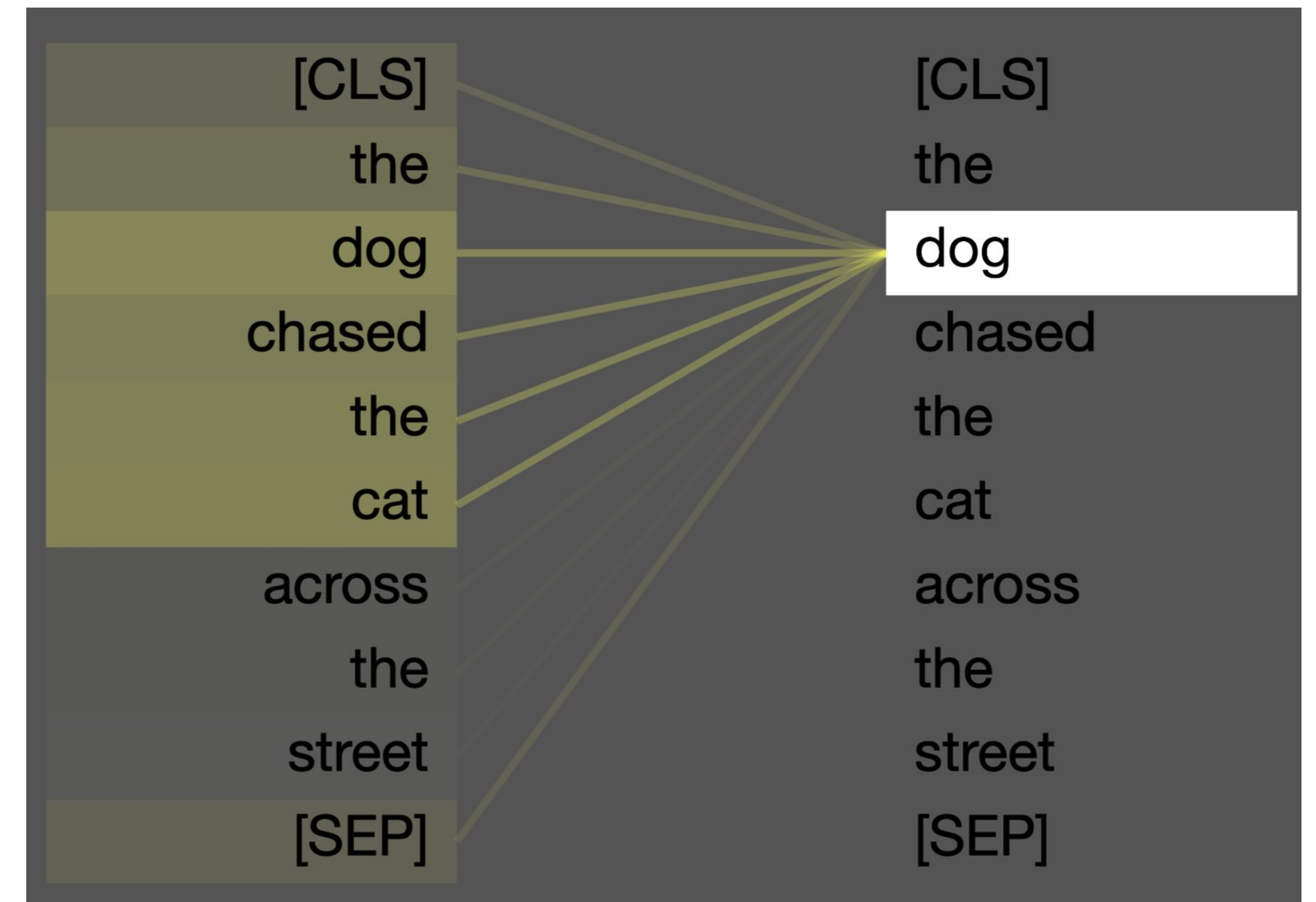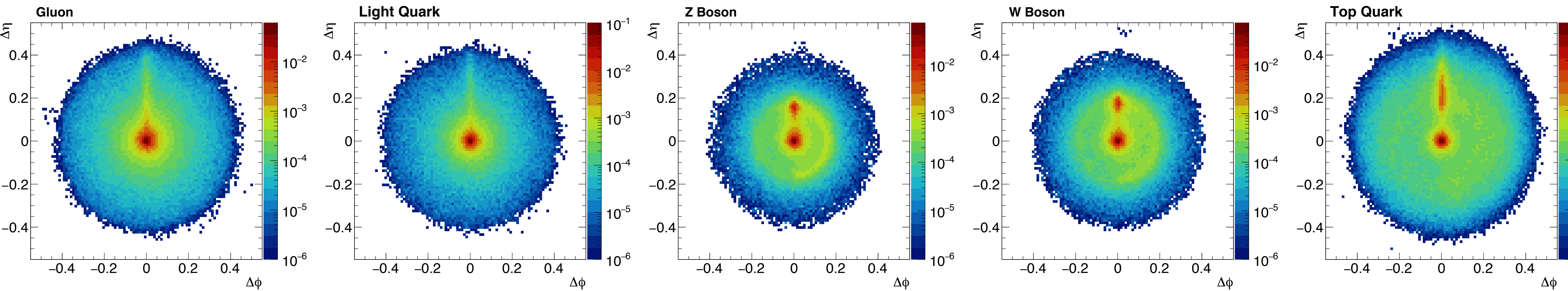
# Transformers

Query, Key and Value: How self-attention is implemented
- **The query**: dog
  I'm looking for verbs, adjectives related to me
- **The key**: every word in the sentence!
  I'm a noun, an adjective or a verb
  (What am I? What features do I posses in relation to the sentence?)
- **The value:** the meaning of this word in general not specifically for this sentence (What're my embeddings? What's the semantic information I posses?)

Self-Attention for word dog:
- **Dog (Query Vector):** Multiplied with all other words (Keys) to get Attention Map
- **Attention Map:** represent importance of every other word related to Dog
- This attention map will be multiplied by the Embeddings of the Sentence words (Values), and produce a weighted sum of the embeddings based on the relevancy of the words

ABCNet:
Pixel intensity = particle importance w.r.t most energetic particle in jet, from attention weights
No substructure information given, learned through attention layers!

# Symmetries

**Symmetries is an extremely important concept, also in Machine Learning.**

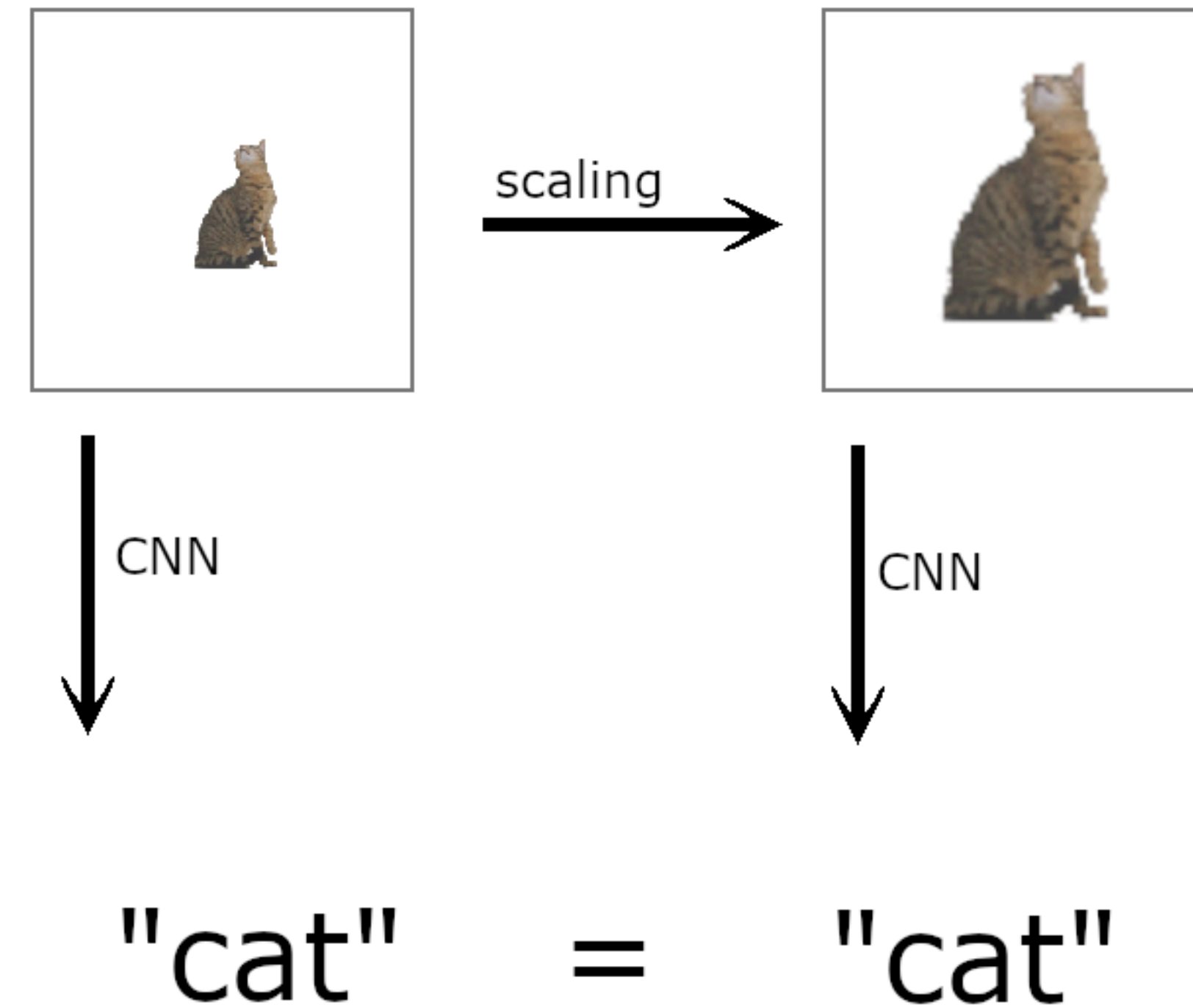• If there is a symmetry in your system, integrate it into your model, and it can do more with less!

| Limited Generalizability | Limited Data | Lack Physical Consistency | Inconsistency & Lack of Guarantees |

*Incorporate symmetry as inductive bias into models*

| Generalize across symmetry | Automatic Data Augmentation | Symmetry and Conservation | Provable Performance |

# Symmetries

**_Result changes in "the same way" as the input_**  **_Result doesn't change when you change the input_**



Equivariance      Invariance

**_E.g CNNs & GNNs (see later)! Invariances are usually obtained through weight sharing_**

# Symmetries

More and more work in HPE try to utilise symmetries when designing DNNs, e .g invariant under Lorentz symmetries!
- Other than observed data, we know the invariances that govern physical phenomena
- Conservation of mass, momentum, energy
- In many cases, we also have approximate models that can predict the system behaviour



*The label of a jet should be invariant under any transformation of the input jet, right?*

*Physics-informed networks respecting Lorentz group symmetries!*

2400 — PELICAN

LorentzNet

2000

**2022 - Transformers**

1600 — ParticleNet — ParT

*BETTER*

$R_{30}$

Disco-FFS on EFPs

1200 — ParticleNet-lite

ResNeXt

**2019 - Message passing graphs**

TreeNiN

PFN

DNN EFPs

CNN

800 — 8 Nsub

6 Nsub

LoLa

LBN

**2018 - CNNs**

P-CNN

EFN

Linear EFPs

400 — TopoDNN

LDA

**~pre-deep learning**

0

LLF taggers
HLF taggers

$10^3$  $10^4$  $10^5$  $10^6$

Parameters

*arXiv:2212:00046*

# Train on simulation, test on data



Eur. Phys. J. C 74 (2014) 3023

**ATLAS** Simulation
Discriminant for MC-Based Tagger
Pythia MC11, $\sqrt{s}$ = 7 TeV
anti-$k_t$ R=0.4, $|\eta|$ < 0.8
160 GeV<$p_T$<210 GeV

quark vs gluon
jets in **simulation**

$L = q/(q+g)$

**ATLAS**
Discriminant for Data-Driven Tagger
$\int L\, dt$ = 4.7 fb$^{-1}$, $\sqrt{s}$ = 7 TeV
anti-$k_t$ R=0.4, $|\eta|$ < 0.8
160 GeV<$p_T$<210 GeV

quark vs gluon
jets in **data**

BOOST 2018, Nachman et al.

$L = q/(q+g)$

Track Width

Track Width

$n_{trk}$

If data and simulation differ, this is sub-optimal!

**CMS**

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

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

Simulation != test data

Mostly (SM )background
samples, small signal datasets

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

We are also very keen
on using this!

# The scientific method

```
┌─────────────────┐
│  Do objective   │
│   observation   │
└────────┬────────┘
         │
┌────────▼────────┐
│  Ask questions  │
└────────┬────────┘
         │
┌────────▼────────┐
│Gather information│
└────────┬────────┘
         │
┌────────▼────────┐      ┌──────────────────┐
│ Form hypothesis │◄─────│ Ask new questions │
└────────┬────────┘      └──────────────────┘
         │                        ▲
┌────────▼────────┐               │
│ Test prediction │               │
└────────┬────────┘               │
         │                        │
┌────────▼────────┐               │
│   Do analysis   │               │
└────────┬────────┘               │
         │                        │
┌────────▼────────┐               │
│Arrive at conclusion│            │
└────────┬────────┘               │
         │                        │
┌────────▼────────┐      ┌──────────────────┐
│  Are results as │─────►│Use experimental  │
│    predicted?   │      │data as hypothesis│
└─────────────────┘      └──────────────────┘
```

# Searches at LHC



**Replaced by:**

**Standard Model (MC)**    **Signal hypothesis (MC)**

Searches at LHC (almost) always start with by
- assuming Standard Model
- and some signal hypothesis

No longer learn from observation
- Blind analysis only way we perform searches

Do objective observation

Ask questions

Gather information

Form hypothesis ← Ask new questions

Test prediction

Do analysis

Arrive at conclusion

Are results as predicted? → Use experimental data as hypothesis

# Searches at LHC

This is fine when you know what you are looking for
- Tailor search to a given theory
- Motivated by belief/disbelief
- Powerful, but limited to model of choice

Everything here is "signal"

Everything here is background

**CMS**
*Simulation*

QCD

W

W

QCD

N-subjettiness ratio $\tau_{21}$

ARE WE LOOKING IN THE WRONG WAY?

# Learning from data

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



clusters

- normal data
- noise
- anomalous data

$x_2$

$x_1$

What are "normal" data and what are "outliers" (and what is noise)?

Let's get back here!



Do objective observation

Ask questions

Gather information

Form hypothesis ⟷ Ask new questions

Test prediction

Do analysis

Arrive at conclusion

Are results as predicted? ⟶ Use experimental data as hypothesis

# Anomaly detection for New Physics searches

LEARN THIS FROM DATA

LOOK FOR ANYTING THAT DOESNT LOOK LIKE THIS

# Types of anomaly detection

## Outlier detection

Find (non-resonant) out-of-distribution datapoints



## Detecting overdensities

Find (resonant) overdensities in distributions

# Types of anomaly detection

## Outlier detection



**Non-resonant, tail of distributions**
- Often (variational) auto-encoders
- Useful for triggering/"selecting"!

## Detecting overdensities



$p_{bg}(x|m_{jj})$     $p_{bg}(x|m_{jj})$

$p_{sig+bg}(x|$

**Resonant, similar to a bump hunt**
- Density estimation methods
- Useful for offline analysis

# Outlier detection

$\mathbf{x}$

$\hat{\mathbf{x}}$

$q$

$q$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

E, $p_x$, $p_y$, $p_z$

$n \times m$

$n \times m$

$\mathfrak{R}^k$

Compressed representation of x.

Latent space $\mathfrak{R}^k$, k < m×n

prevents memorisation of input, must learn

# Outlier detection



E.g 3-prong gluino fat jet

$\mathbf{x}$

$n \times m$

$\mathfrak{R}^k$

$\hat{\mathbf{x}}$

$n \times m$

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

# Outlier detection

**x**

E.g 3-prong gluino fat jet

$n \times m$

$\mathfrak{R}^k$

arXiv:1808.08992

Large error for
abnormal data

Encoder

Decoder

# Outlier detection in analysis

E.g

# Outlier detection in ana

E.g



## Careful! Cut on score can sculpt spectrum



$q_{99}$

$q_{50}$

10% bkg eff

30% bkg eff

60% bkg eff

80% bkg eff

inclusive

$m_{jj}$

$M_{jj}$

## Can fix using quantile regression



10% bkg eff

30% bkg eff

60% bkg eff

80% bkg eff

inclusive

$m_{jj}$

60% bkg eff

80% bkg eff

inclusive

$m_{jj}$

10% bkg eff

30% bkg eff

60% bkg eff

# Outlier detection in analysis

# Example for semi-visible jets

Normalized autoencoders        : Lund Graph autoencoders

# Finding overdensities



FETA

# Weak classification without labels (CWoLa)

Mostly background

Signal + background

Mostly background

Dijet invariant mass

# Weak classification without labels (CWoLa)

## Classification Without Labels

- Lemma: "Given mixed S+B samples SB and SR, optimal classifier trained to distinguish SB and SR is also optimal for distinguishing S from B"



Dijet invariant mass

LABEL = SIGNAL

LABEL = BKG

Mostly background

Signal + background

Mostly background

Dijet invariant mass

S

B

CWola hunting in ATLAS

E.g

# DNN likelihood

Alternative approach: End-to-end DNN search
- How do we get around defining a signal hypothesis?
- What is alternate hypothesis to test reference?

Idea: Assume alternate model n(x|w) can be
parametrised in terms of reference model n(x|R)

$$n(x \,|\, \vec{w}) = n(x \,|\, R)e^{f(x;\vec{w})} \longleftarrow \text{Set of real functions}$$

- Let DNN parametrise alternative model

$$f(x; \vec{w}) = NN$$

# DNN likelihood

Alternative approach: End-to-end DNN search
  • How do we get around defining a signal hypothesis?
  • What is alternate hypothesis to test reference?

Idea: Assume alternate model n(x|w) can be
parametrised in terms of reference model n(x|R)

$$n(x \mid \overrightarrow{w}) = n(x \mid R)e^{f(x;\overrightarrow{w})} \quad \longleftarrow \text{ Set of real functions}$$

  • Let DNN parametrise alternative model

$$f(x; \overrightarrow{w}) = NN$$

  • Formulate loss as log likelihood.
    → Trained DNN **is** the maximum likelihood fit
    to data and reference log-ratio
    → best approximate of true data distribution

$$f(x, \widehat{\mathbf{w}}) \simeq \log \left[ \frac{n(x|\mathrm{T})}{n(x|\mathrm{R})} \right] \begin{array}{l} \longleftarrow \text{True underlying data distribution} \\ \longleftarrow \text{MC distribution} \end{array}$$

**INPUTS**
- any high level features

**OUTPUTS**
-t$^{obs}$ and f(x; ŵ)

QCD MC R



g. (1, 4, 1) network
g. (1, 4, 1) network

1) Best fit log ratio of data
and MC PDFs



$$f(x, \hat{w}) \simeq \log\left[\frac{n(x\,|\,T)}{n(x\,|\,R)}\right]$$

$x$ — Neural Network $\mathbf{w}$ — $f(x; \mathbf{w})$

Train $\mathcal{D}$ vs. $\mathcal{R}$

CMS DATA D

$x$ — Neural Network $\widehat{\mathbf{w}}$ — $f(x; \widehat{\mathbf{w}})$



2) test-statistic on data
sample t$_{obs}$

$$t(\mathcal{D}) = -2 \operatorname*{Min}_{\{\mathbf{w}\}} L[f]$$  ← **DNN loss function!**

Can be used to build
hypothesis test + p-value
Data → toys under R,
repeat

$$f(x, \widehat{\mathbf{w}}) \simeq \log\left[\frac{n(x\,|\,\mathrm{T})}{n(x\,|\,\mathrm{R})}\right]$$  ← True underlying data distribution
← MC distribution

# Hybrid approaches - NoVa

# Hybrid approaches - NoVa

*Aurisano et al*
*K. Sachdev*



Efficiency of selecting electron neutrinos improved by 40%

# Hybrid approaches - NoVa

*Aurisano et al*
*K. Sachdev*

Neutrino
from
Fermilab

Efficiency of selecting electron neutrinos improved by 40%



(a) A candidate $\nu_\mu$ CC interaction in ND data



(b) The muon removed or MRCC version of the event



*Similar techniques used for H → ττ by ATLAS and CMS!*

(c) A simulated electron is inserted in place of the muon to make an MRE event.

Microsoft

Machine Learning
"how"
learning algorithms

Deep Learning
features
architectures

Foundation Models
functionalities
models

Emergence of...
Homogenization of...

The New York Times

A.I. and Chatbots ›   How the A.I. Race Began   One Year of ChatGPT   Key Figures in A.I.   How A.I. Could Be Regulated

THE SHIFT

*Maybe We Will Finally Learn More About How A.I. Works*

Stanford researchers have ranked 10 major A.I. models on how openly they operate.

| | | | |
|---|---|---|---|
| BigScience | BLOOM | 176B | July 2022 |
| | T0pp | 11B | October 2021 |
| EleutherAI | GPT-J | 6B | July 2021 |
| | GPT-NeoX | 20B | February 2022 |
| Tsinghua University | GLM | 130B | August 2022 |
| Google Research | UL2 | 20B | October 2022 |
| | T5 | 11B | February 2020 |
| | OPT | 175B | June 2022 |

# AI Explainer: Foundation models and the next era of AI

Published March 23, 2023

**Next-gen (existing) applications**

Product & customer interaction / management
viable  chatdesk  Quickchat
Nevermaps  ActiveChat  exceed by GENESYS
Stateset  Sapling

Personal productivity
personal.ai
mem  O

Search engine
YOU  Google
algolia

Oogway

**Emerging net-new applications**

Application synthesis
Adept  CODEGEN

Data analyst productivity
veezoo  AI 2sql*  cogram

Developer productivity
warp  tabnine
GitHub Copilot  ASK JARVIS
repl.it  M  KD

New media generation
FABLE  DALL·E 2  MidJourney  alethea.  R

Writing assistant/text generation
AI21labs  Jasper  Snazzy AI
PR Guy  copy.ai  W  Scalenut
LAVENDER  YOUWrite
anyword  Simplified
copysmith  copymatic
LONGSHOT  Rytr  Writesonic (previously MagicFlow)

**Infrastructure**

Model /builders providers  - Big Tech
Microsoft
Google  DeepMind
Meta  NVIDIA

Model providers/builders  - Startups
OpenAI  cohere
Hugging Face  BigScience
AI21labs  Lighton
ANTHROP\C

Accessible specialized AI chips
NVIDIA  GRAPHCORE
λ  Google  Lighton

Other tooling
Humanloop  anyscale

Figure 5: Representative sample of companies that have publicly stated that they are using, building, or enabling

Foundation Models: An Explainer for Non-Experts

Share

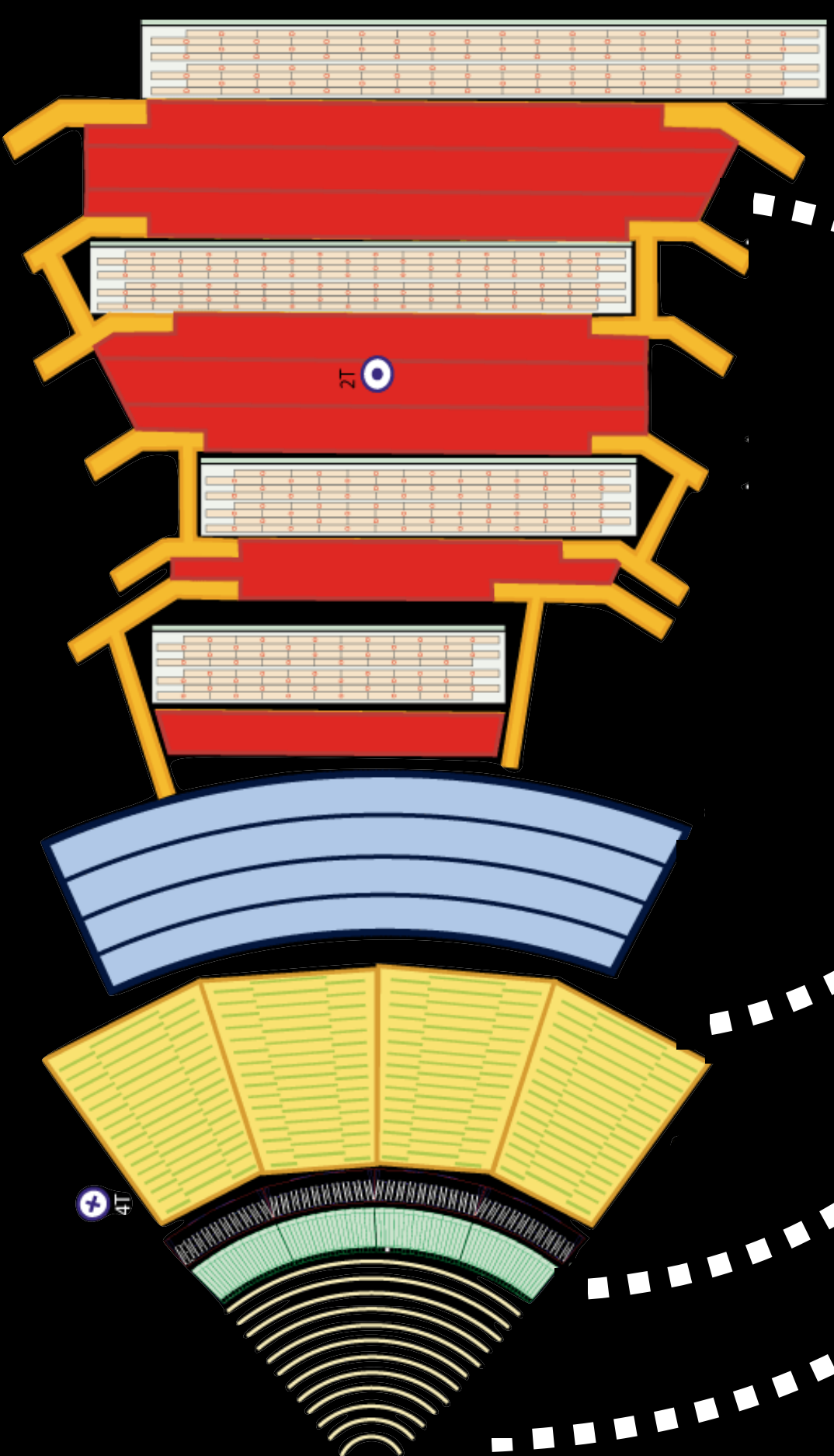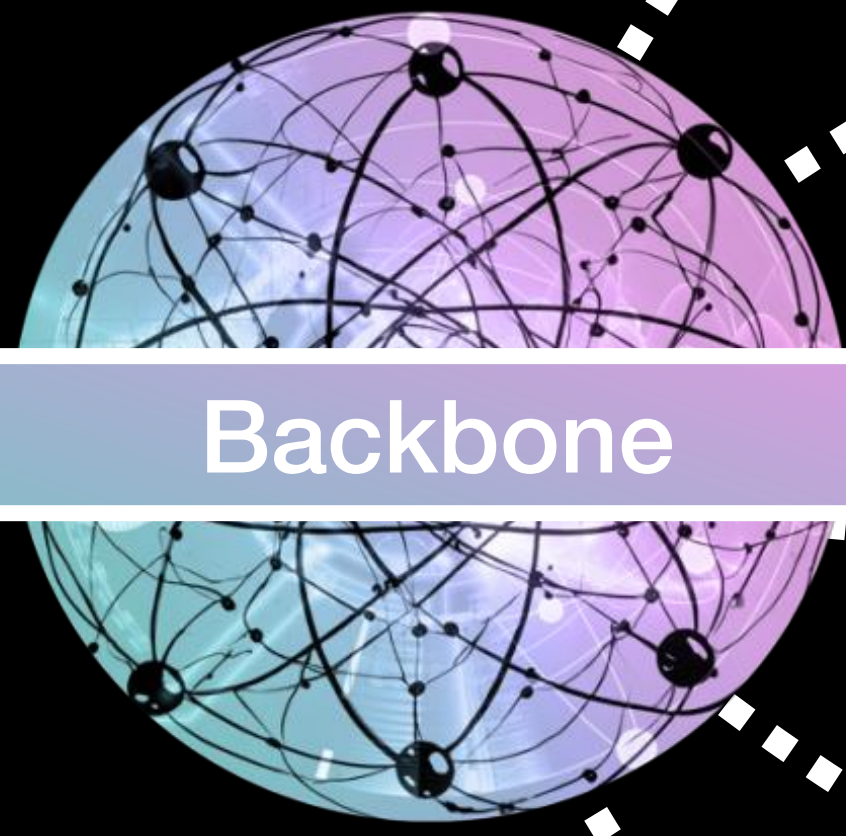TENS OR EVEN HUNDREDS OF MILLIONS

$

MORE VIDEOS

1:12 / 2:09

YouTube

# Foundation Models

Heterogeneous detector
Multi-modal input!

Pre-training → Backbone

Fine-tuning — Jet reconstruction

Fine-tuning — Electron ... on

Fine-tuning

Fine-tuning — Missing energy computation

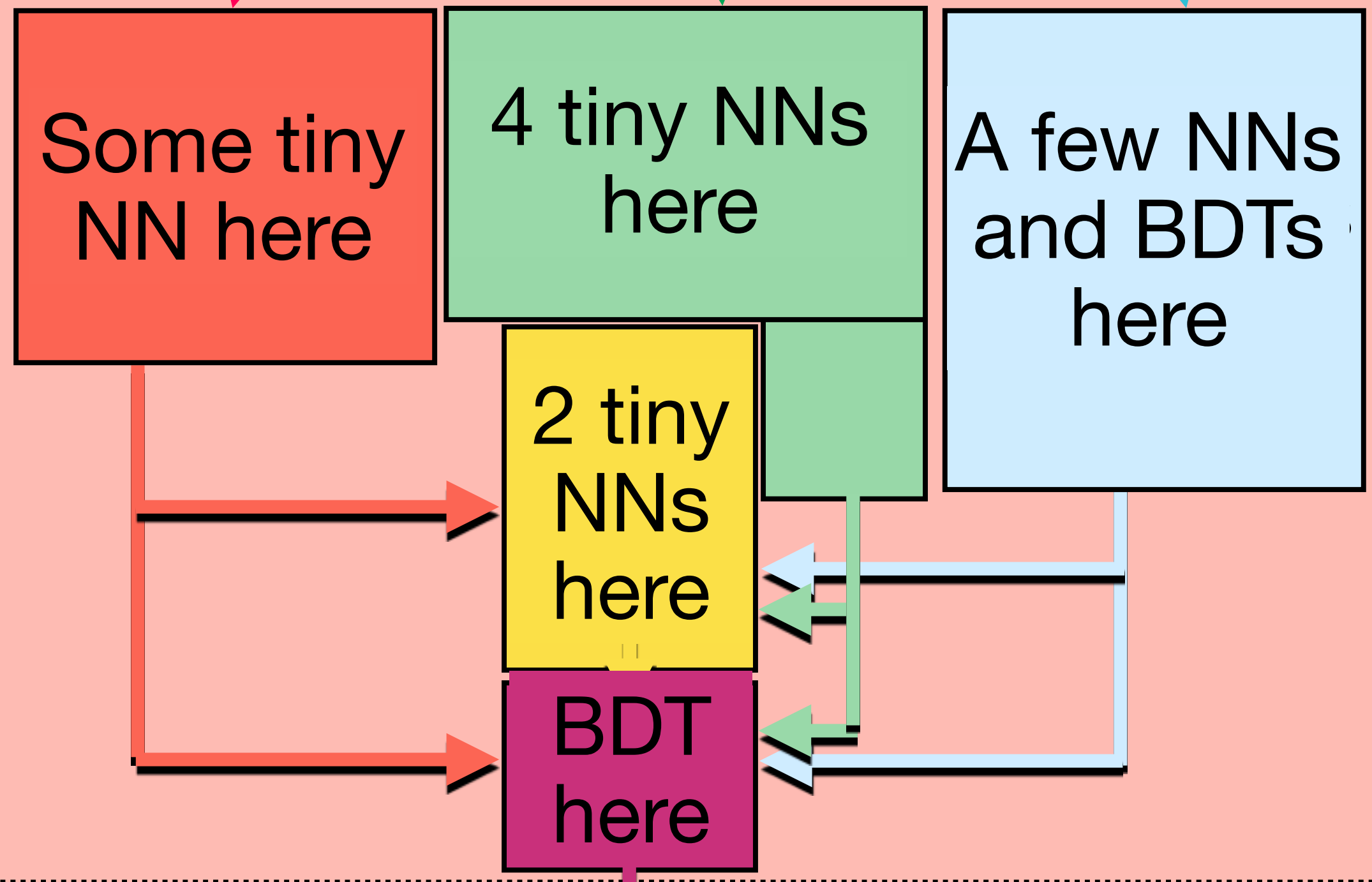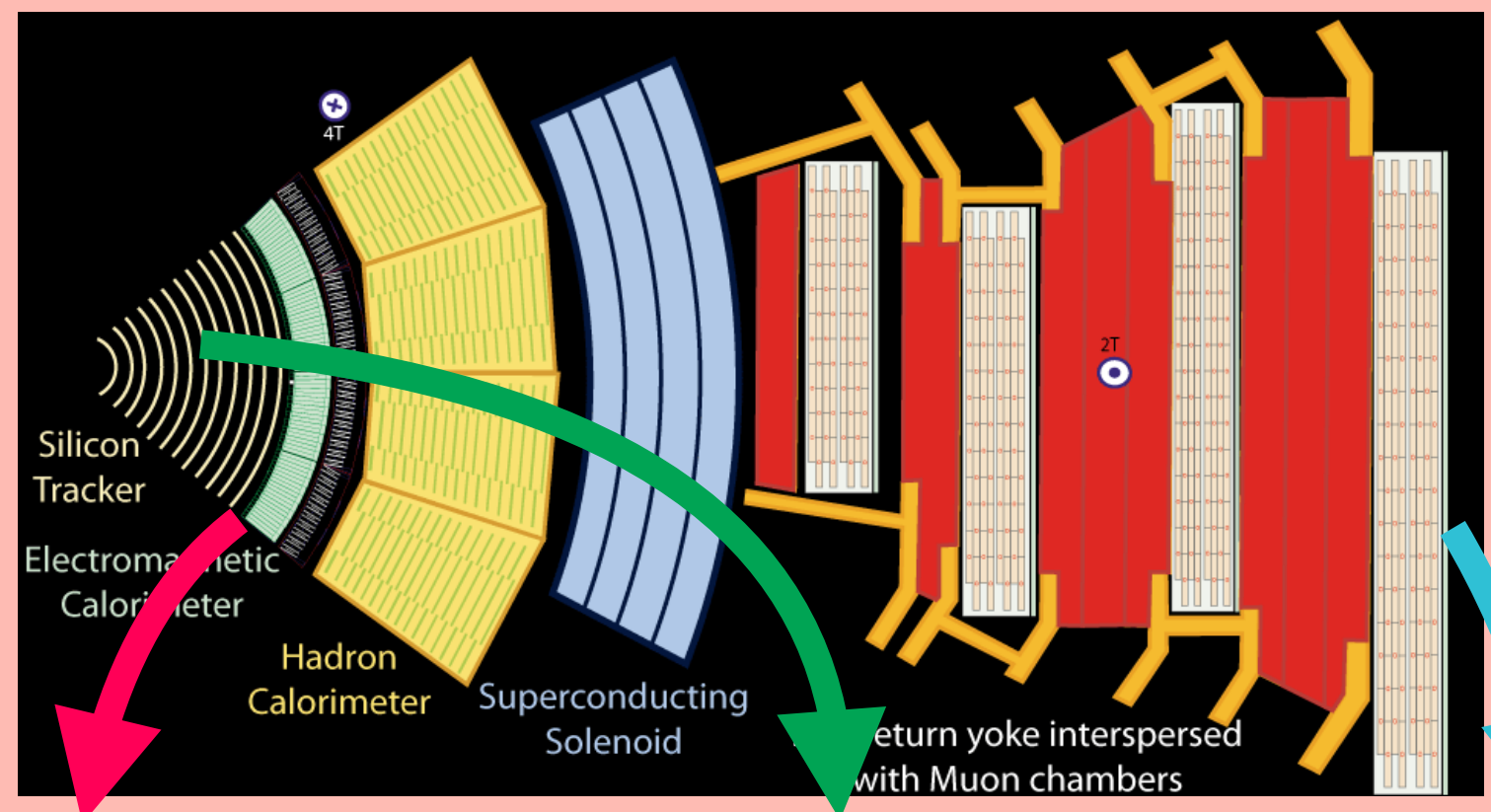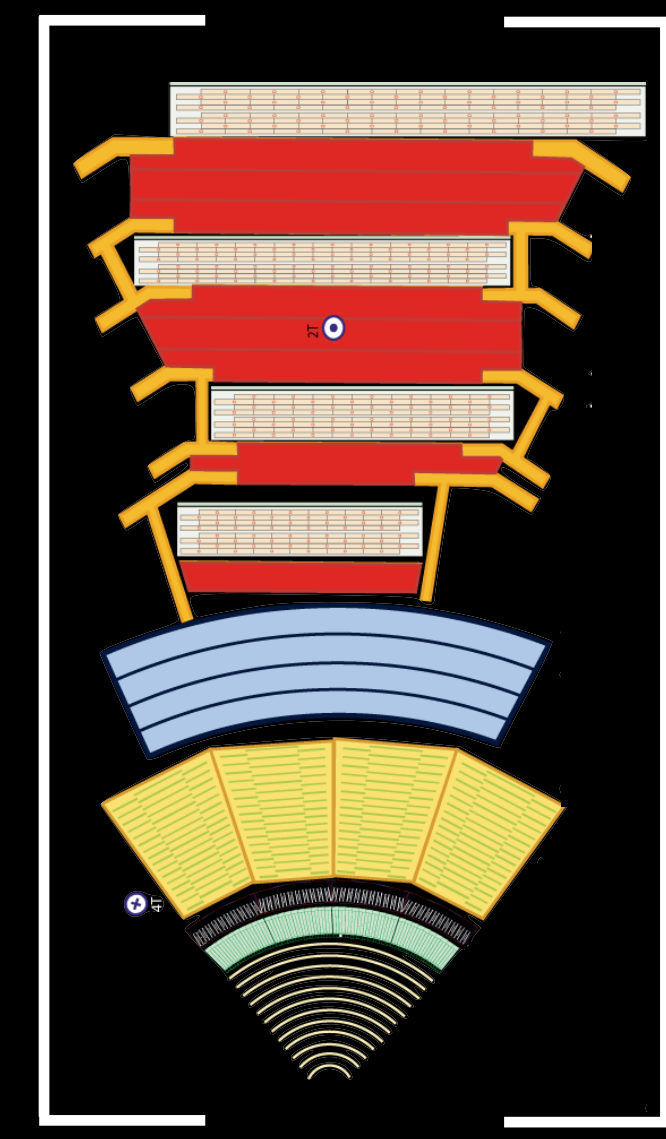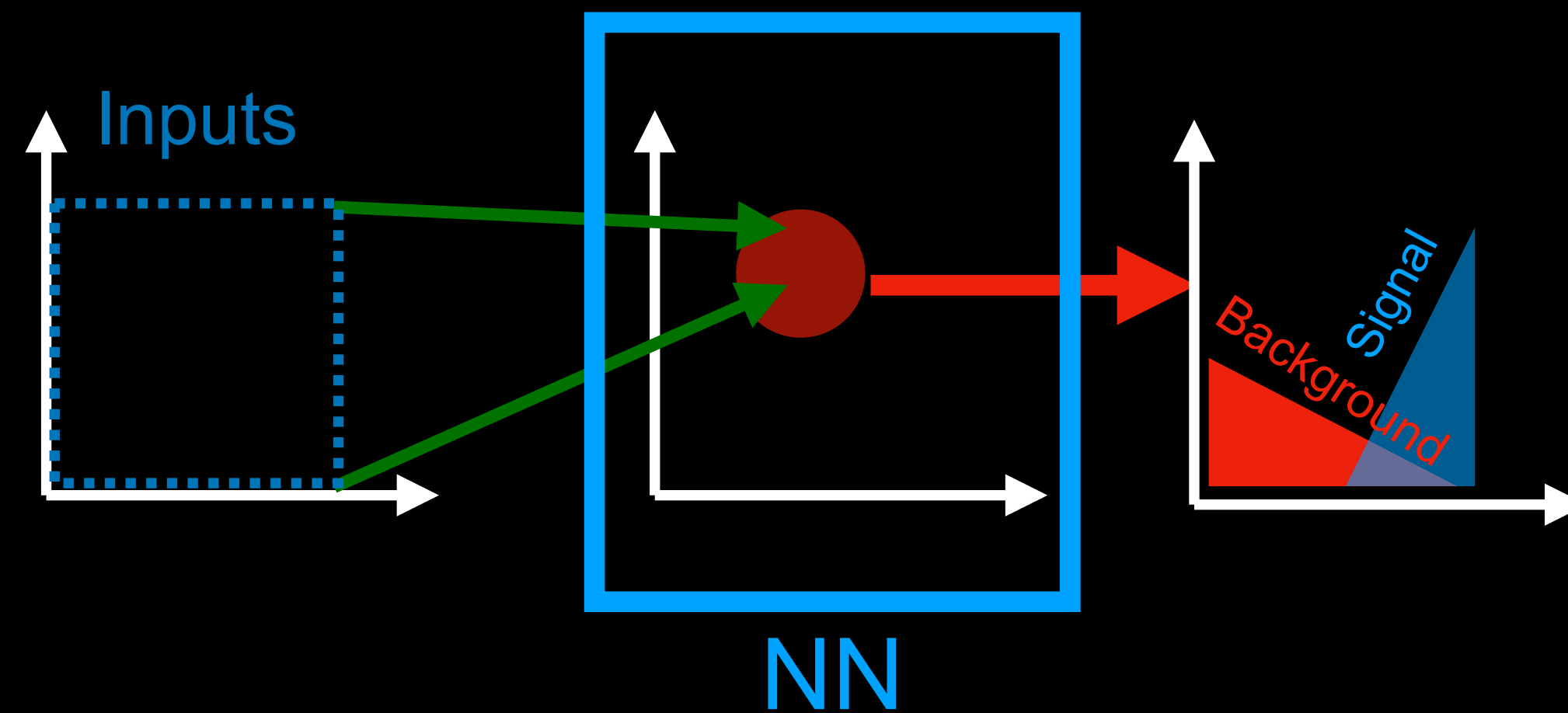Fine-tuning — Anomaly Detection

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

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

$f(x; w*)$ $\hat{y}$

# Too many models, too little learning?



## Discrimination

Metric Learning

What if we really try to focus on this space

Inputs

NN

Something New

Signal

Background

Other Signal

Background

Signal

**Neural embedding**

What if we really try to focus on this space

Inputs

NN

Something
New

Signal

Background

**Neural embedding**

# 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
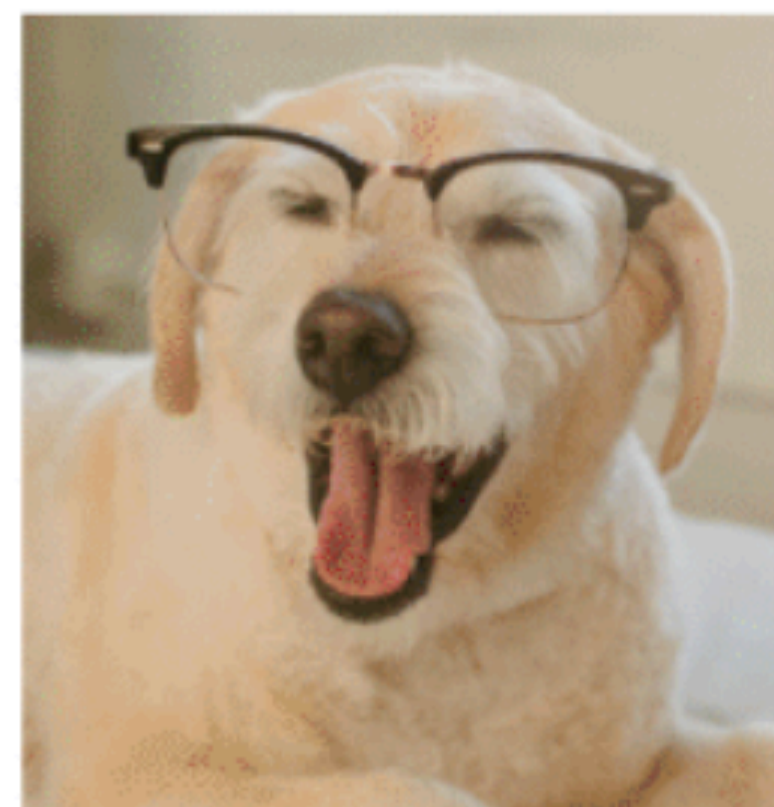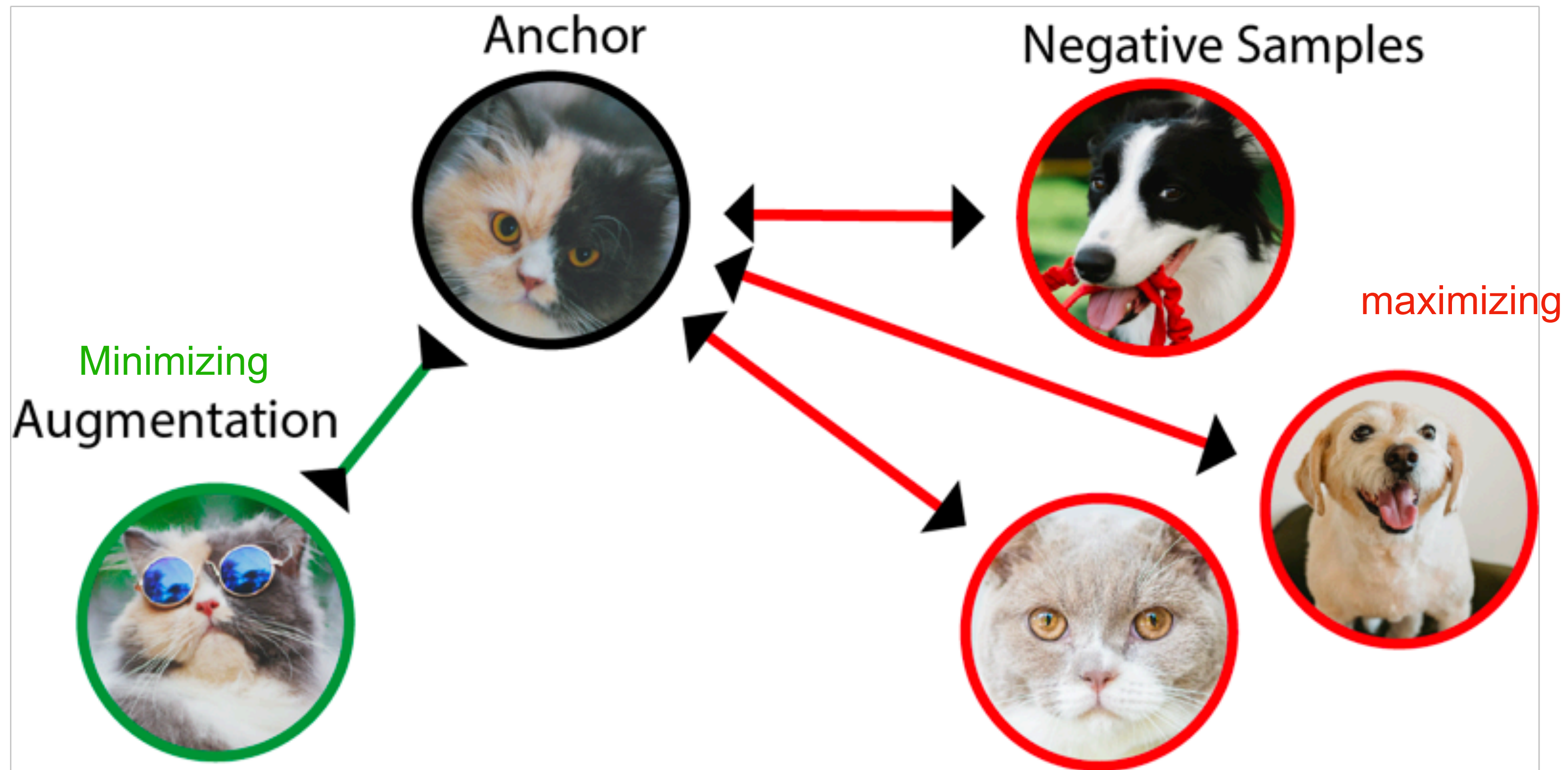
Cat A



Dog A

Cat A

Augmented Cat A

Dog A

Augmented Dog A

# Minimize

# Maximize
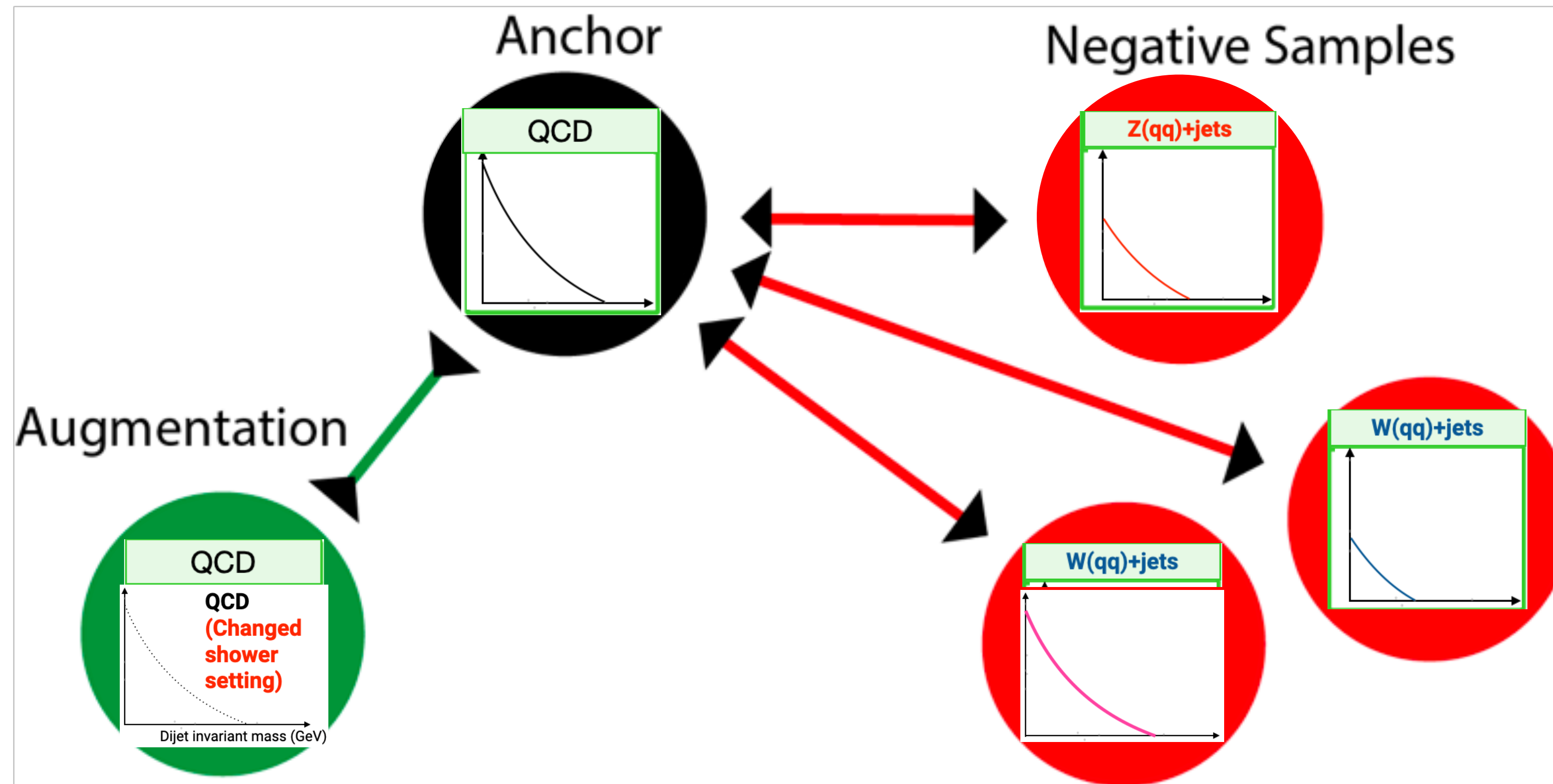
# Physically motivated augmentations?



- Minimizing and maximizing distances learns a space

Augmented Cat A

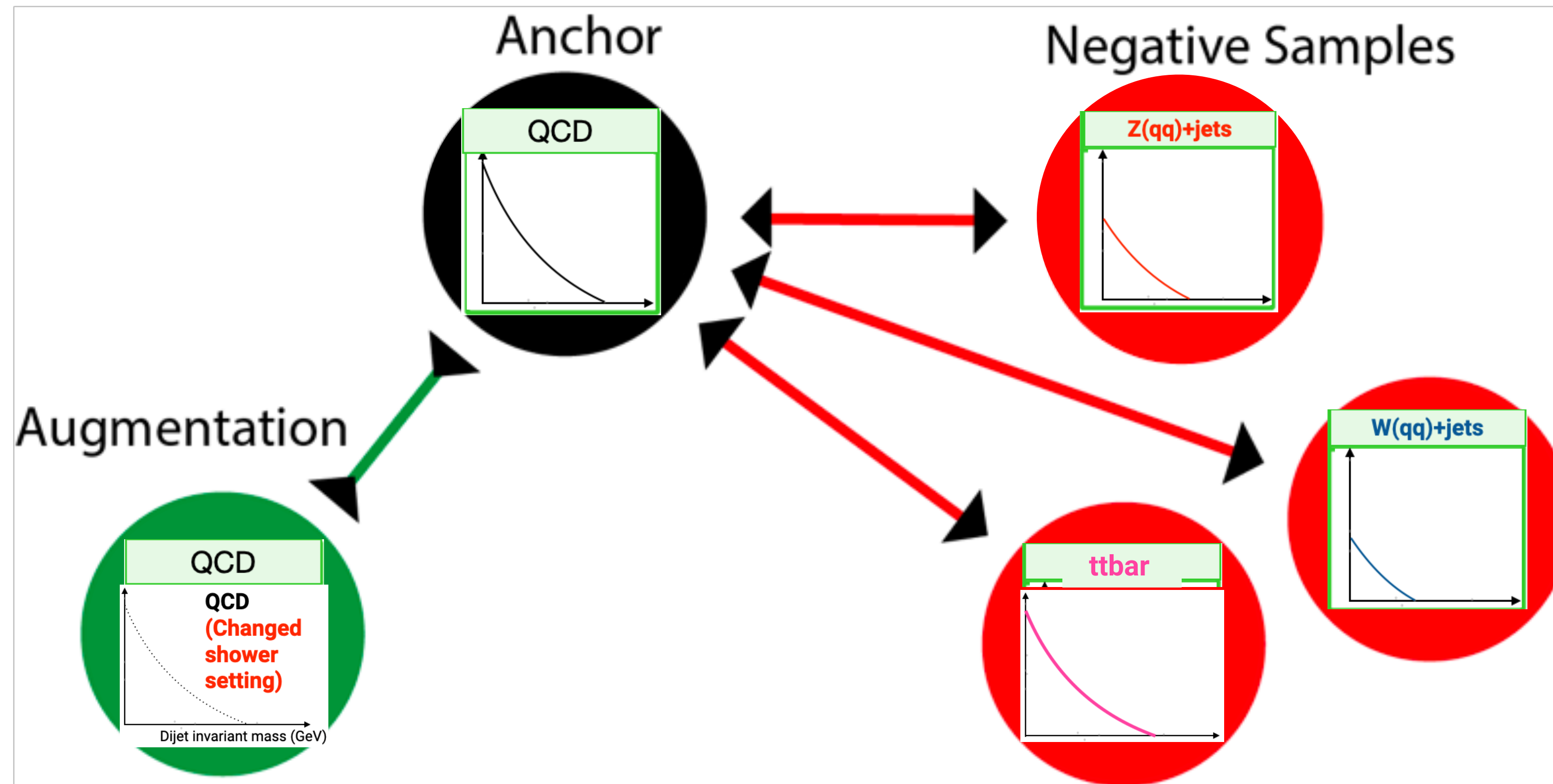Dog B

Cat A

Dog A

Cat B

Augmented Dog A

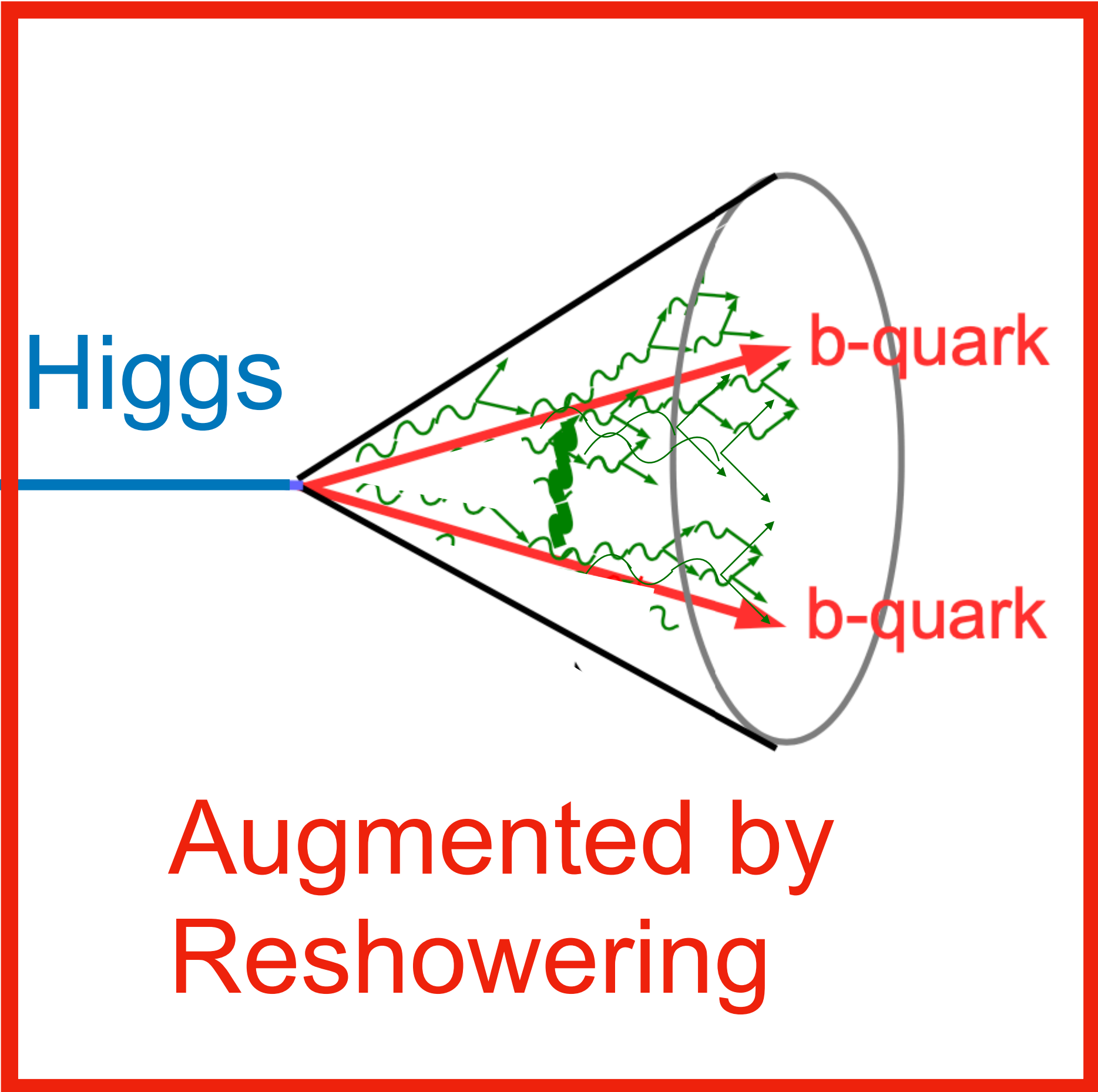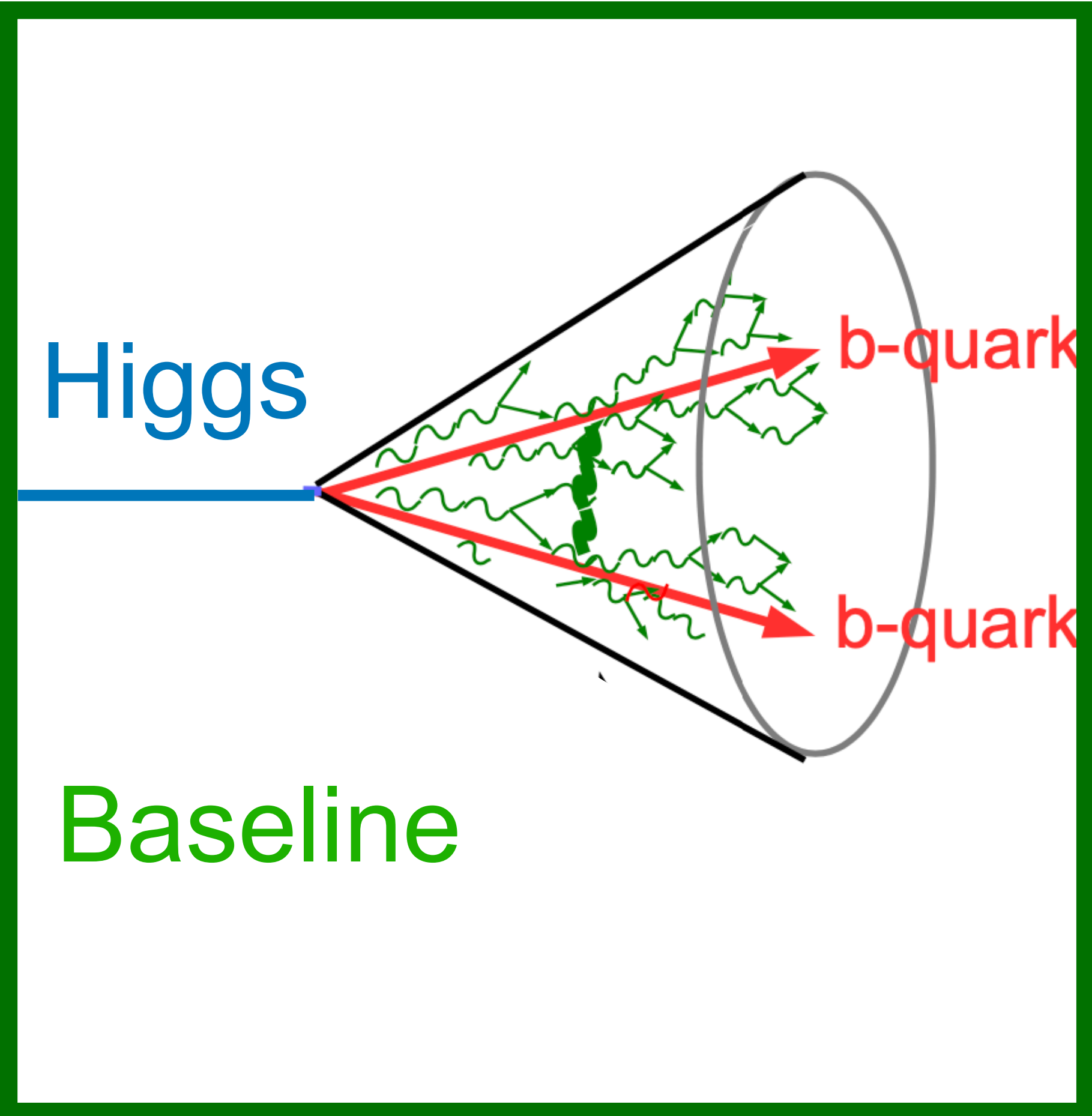# Physically motivated augmentations?



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?
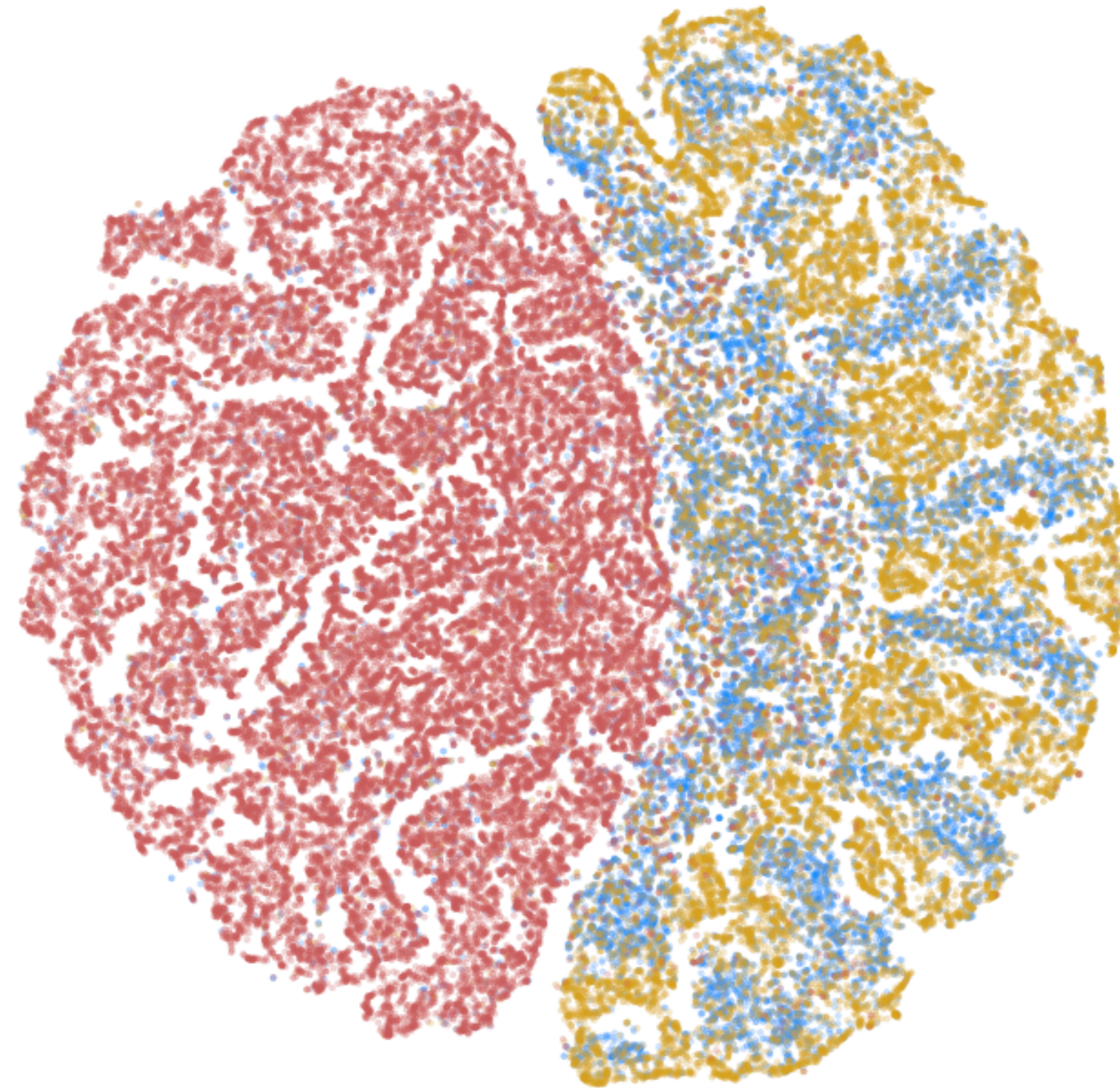
Augmentation

Higgs

Baseline

b-quark

b-quark
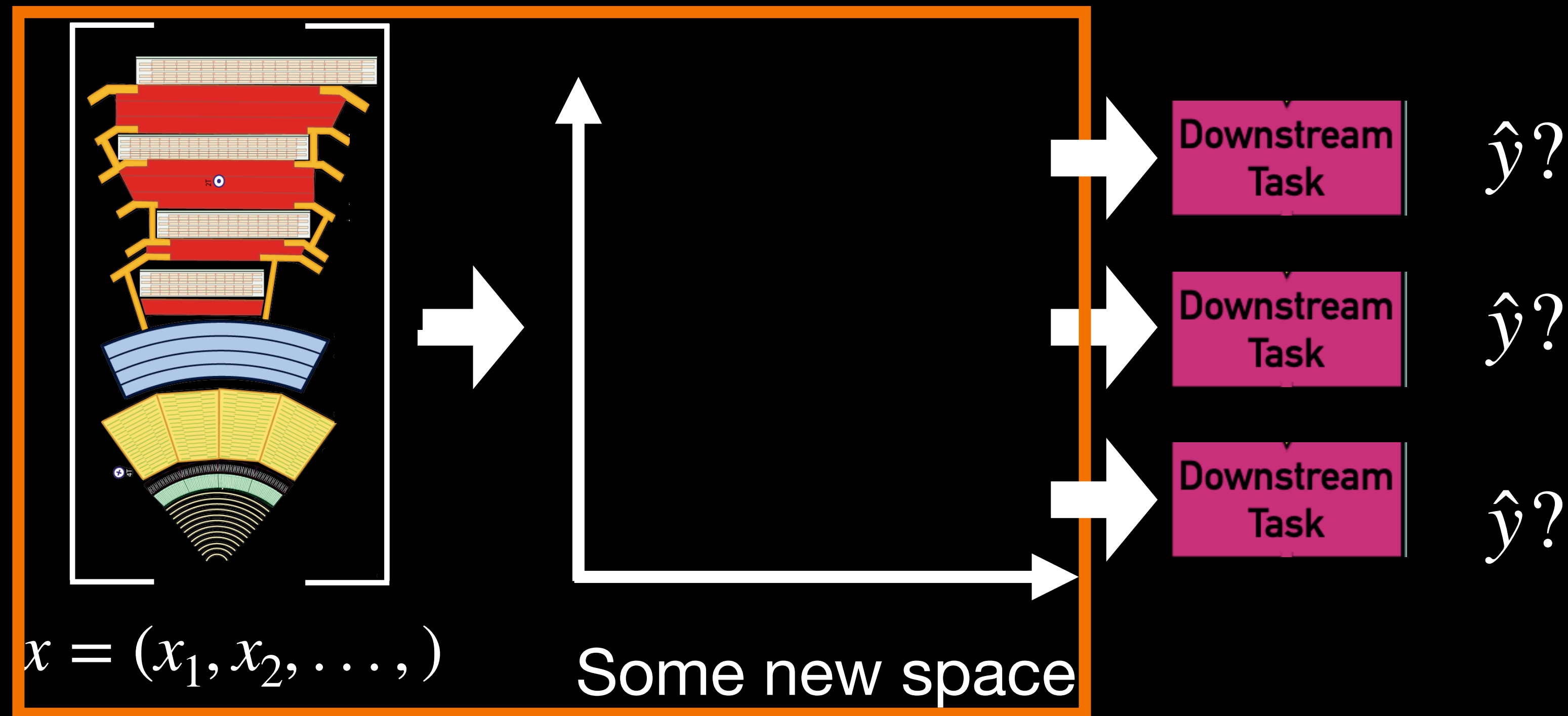
Higgs

Augmented by Reshowering

b-quark

b-quark

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 $\hat{y}$?

Downstream Task $\hat{y}$?

Downstream Task $\hat{y}$?

**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)**

Theorists

N-D Space

Signal

Other Signal

Background

Capture Physics

Signal

Background

We can replace the QCD theorist with a NN
(And it works better)

(Graph) NN

N-D Space

Signal

Background

Other
Signal

NN

Capture
Physics

Background

Signal

# Masked language modelling



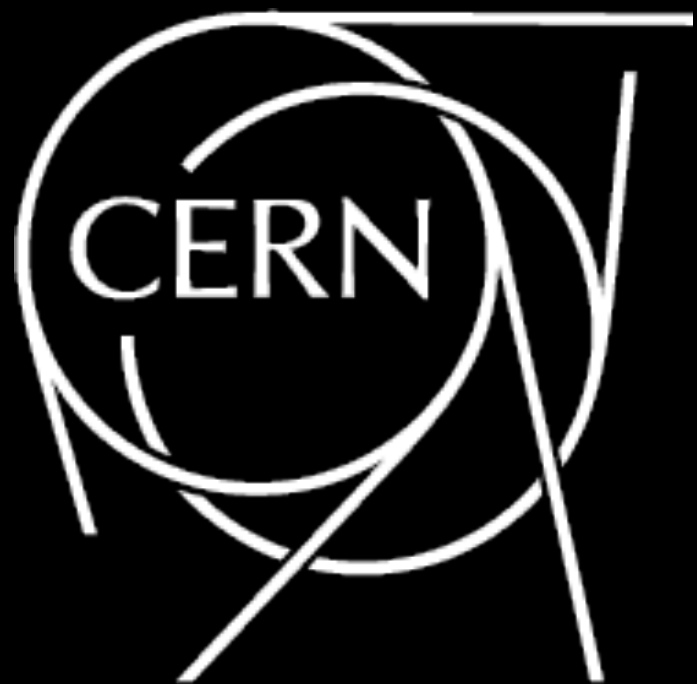| Next-token-prediction | Masked-language-modeling |
|---|---|
| The model is given a sequence of words with the goal of predicting the next word. | The model is given a sequence of words with the goal of predicting a 'masked' word in the middle. |
| Example:<br>Hannah is a ____ | Example<br>Jacob [mask] reading |
| Hannah is a *sister*<br>Hannah is a *friend*<br>Hannah is a *marketer*<br>Hannah is a *comedian* | Jacob *fears* reading<br>Jacob *loves* reading<br>Jacob *enjoys* reading<br>Jacob *hates* reading |

# Self-supervised pre-training

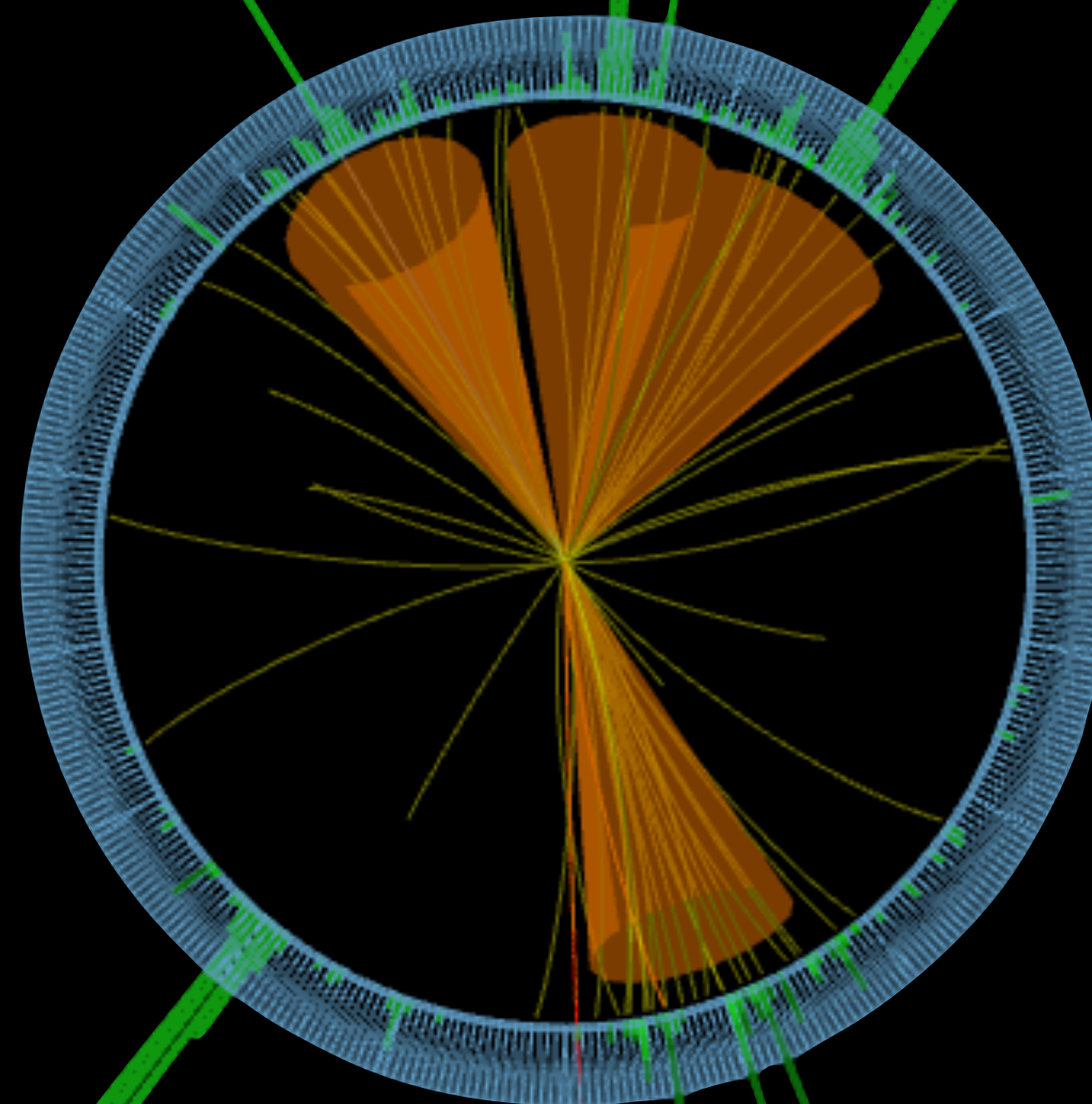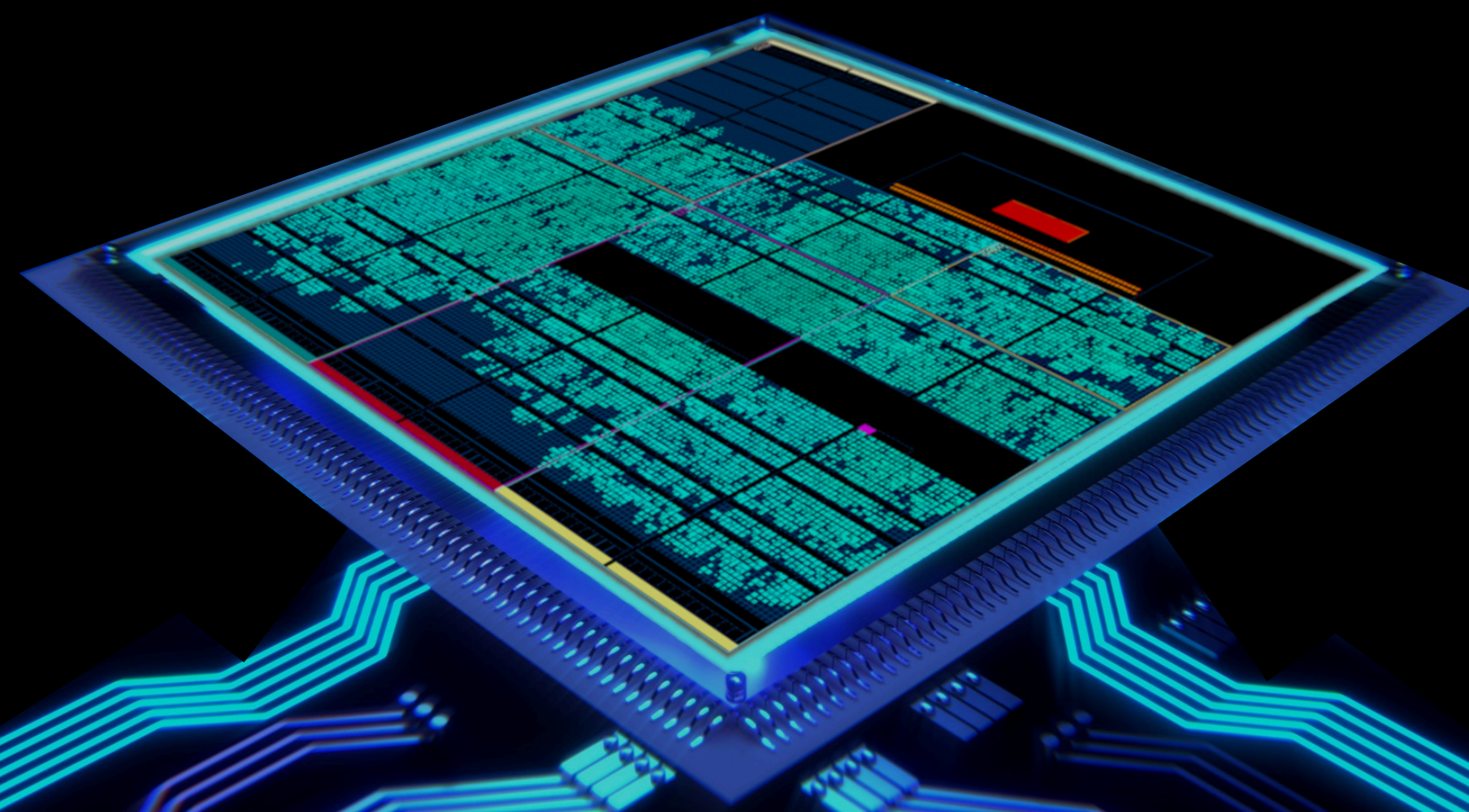# Masked particle modelling



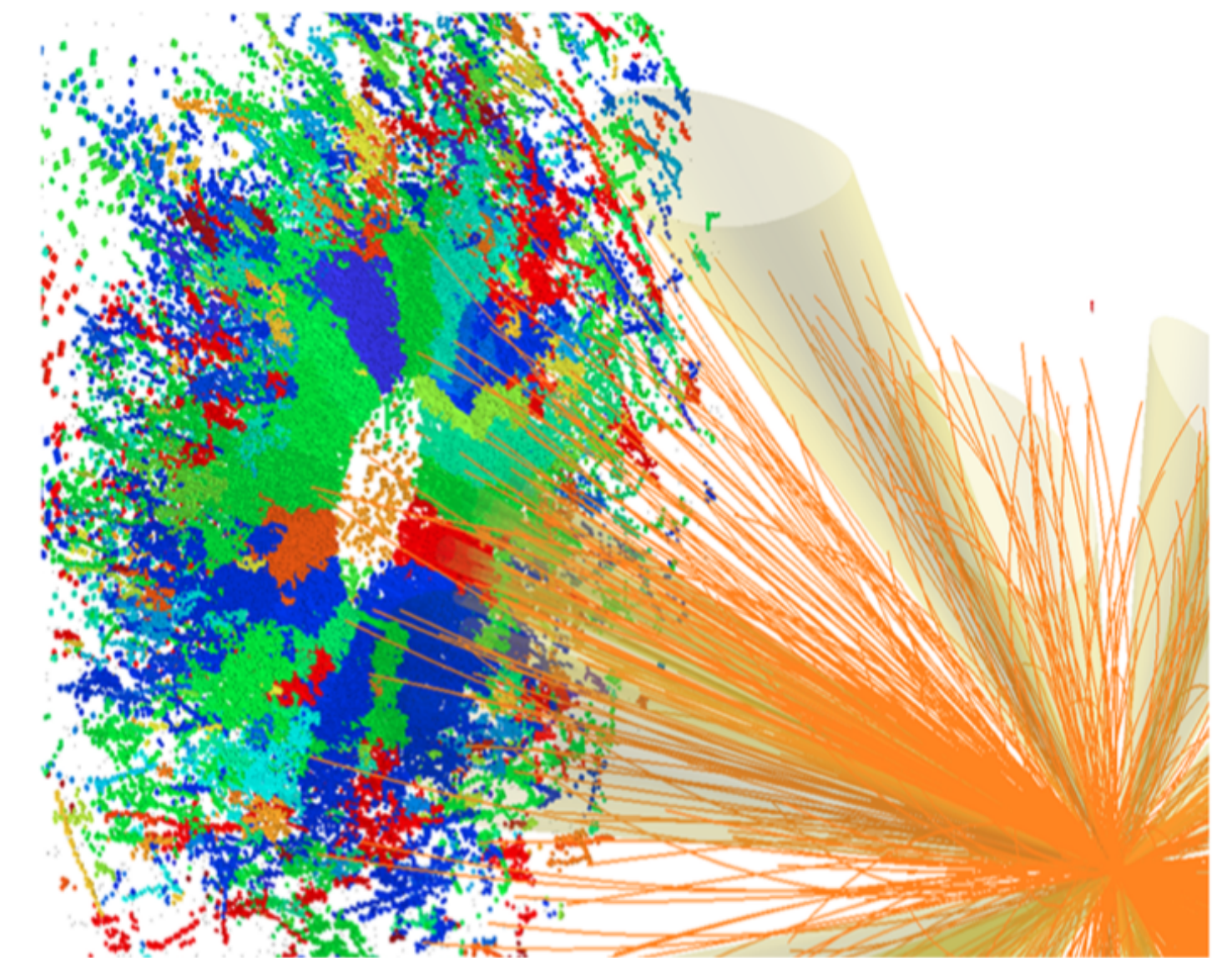# Masked calorimeter pre-training?

# Part2:
# ML in HEP
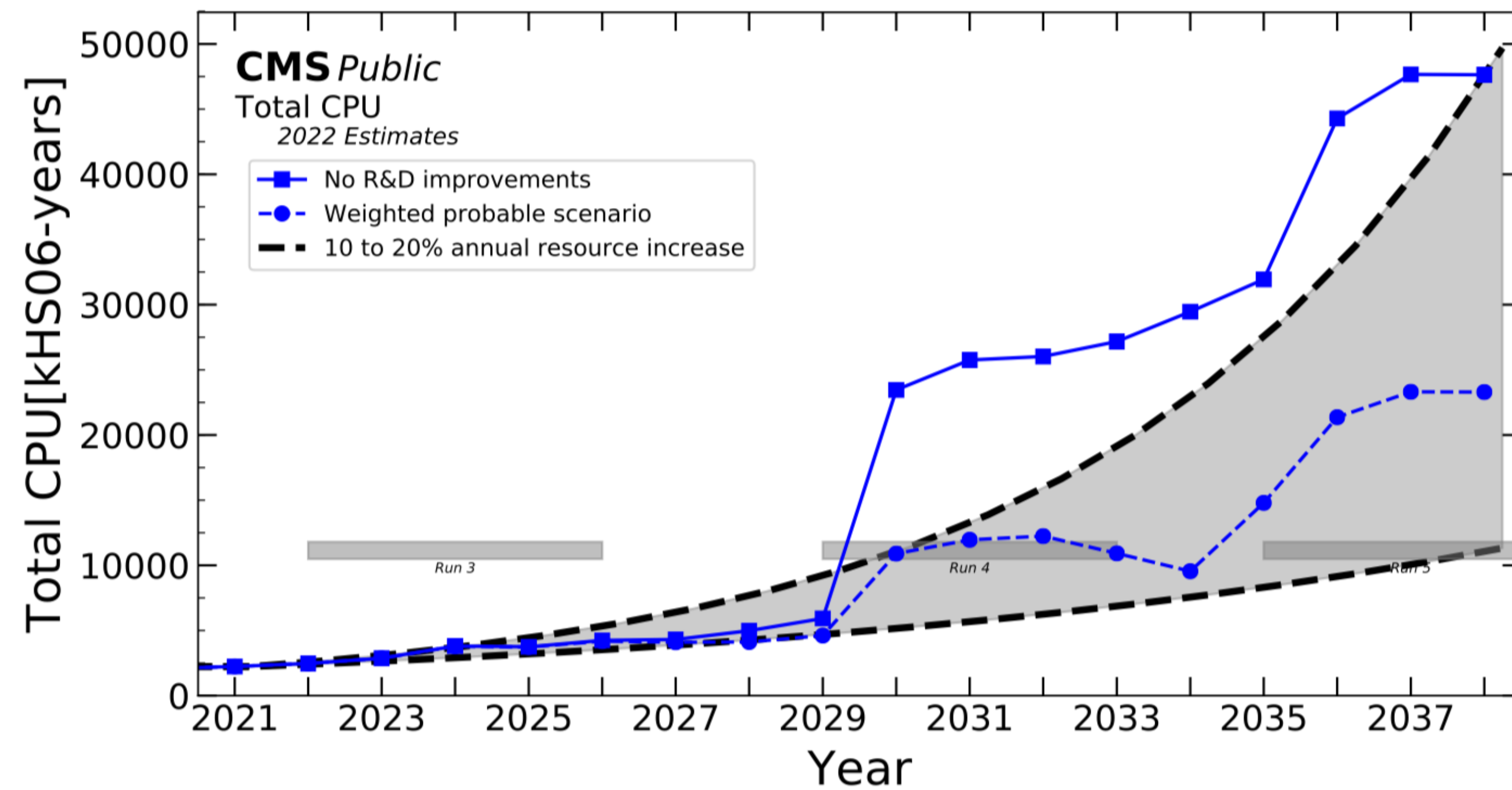
Thea K. Årrestad (ETH Zürich)
thea.aarrestad@cern.ch
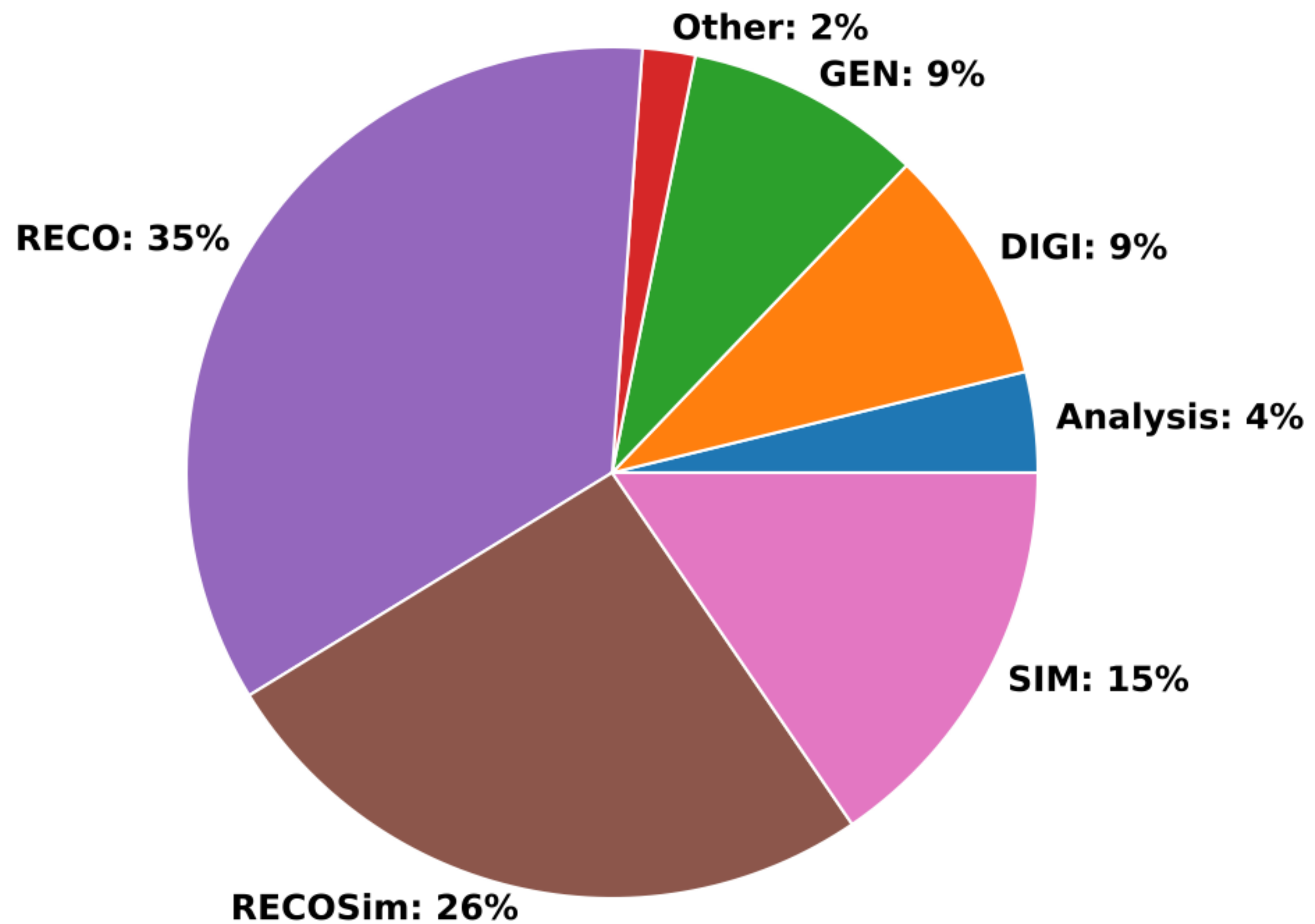thaarres.github.io

QCD School 2024

# ML for simulation

HL-LHC, Simulation of CMS HGCAL with 140 PU

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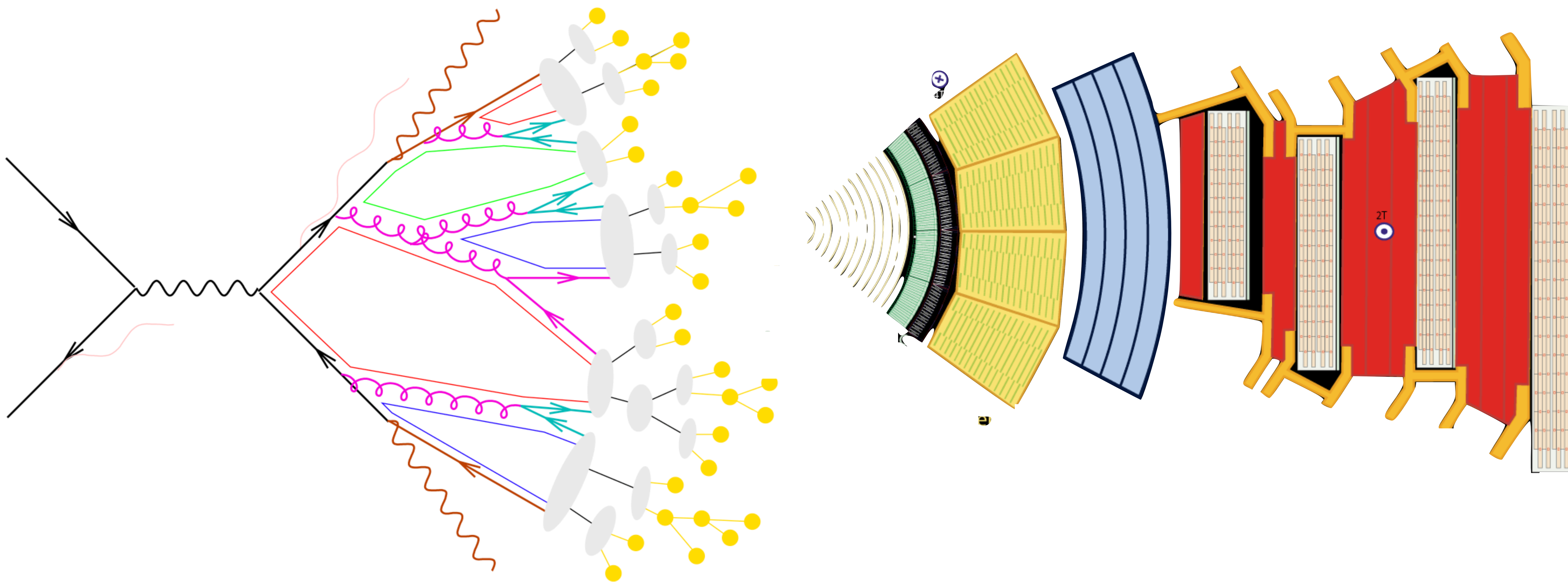*2022 Estimates*

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

$10^{-18}$m  $10^{-15}$m  $10^{-6}$m  $100$m
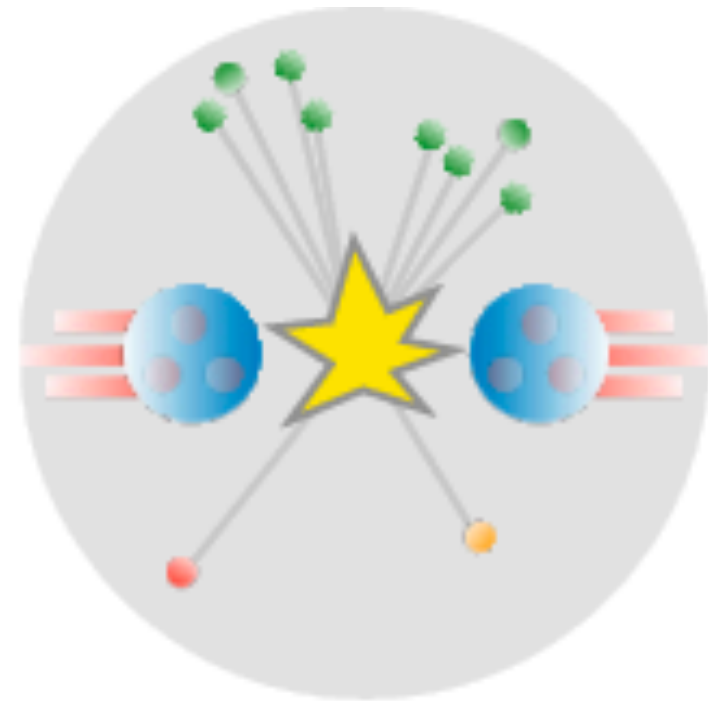
$O(10)$    $O(10^3)$    $O(10^{10})$

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

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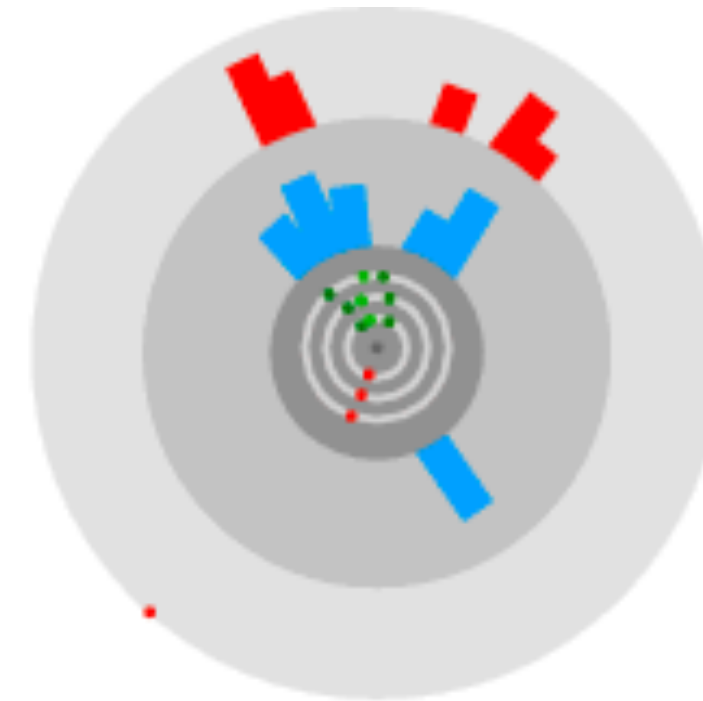
pp collisions up to production of stable particles [Easy & Fast]

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%

GEANT4
A SIMULATION TOOLKIT

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

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

GEN    SIM    DIGI+RECO+MINIAOD

collision point
proton beams
$\phi$
$\eta$

1.1%    CPU    0.1%
16.8%    24.4%
GEN
SIM
DIGI
RECO
MINIAOD
57.6%

9%    Disk    10%
SIM
MINIAOD
81%

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

pp co
produ
particle

57.6%

particle
DIG
REC
MIN
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
pro
particle

pp co
produ
1
particle

ML me

1.1%  CPU  0.1%
16.8%  24.4%

- GEN
- SIM
- DIGI
- RECO
- MINIAOD

57.6%

Disk
9%  10%
Energy deposits→digital
signals→reconstructed by
SIM
the reconstruction software
MINIAOD
[Hard & Slow]

81%

57.6%

81%

# Diffusion models



Model

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

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

Gaussian noise

**_Dall-e 2_**

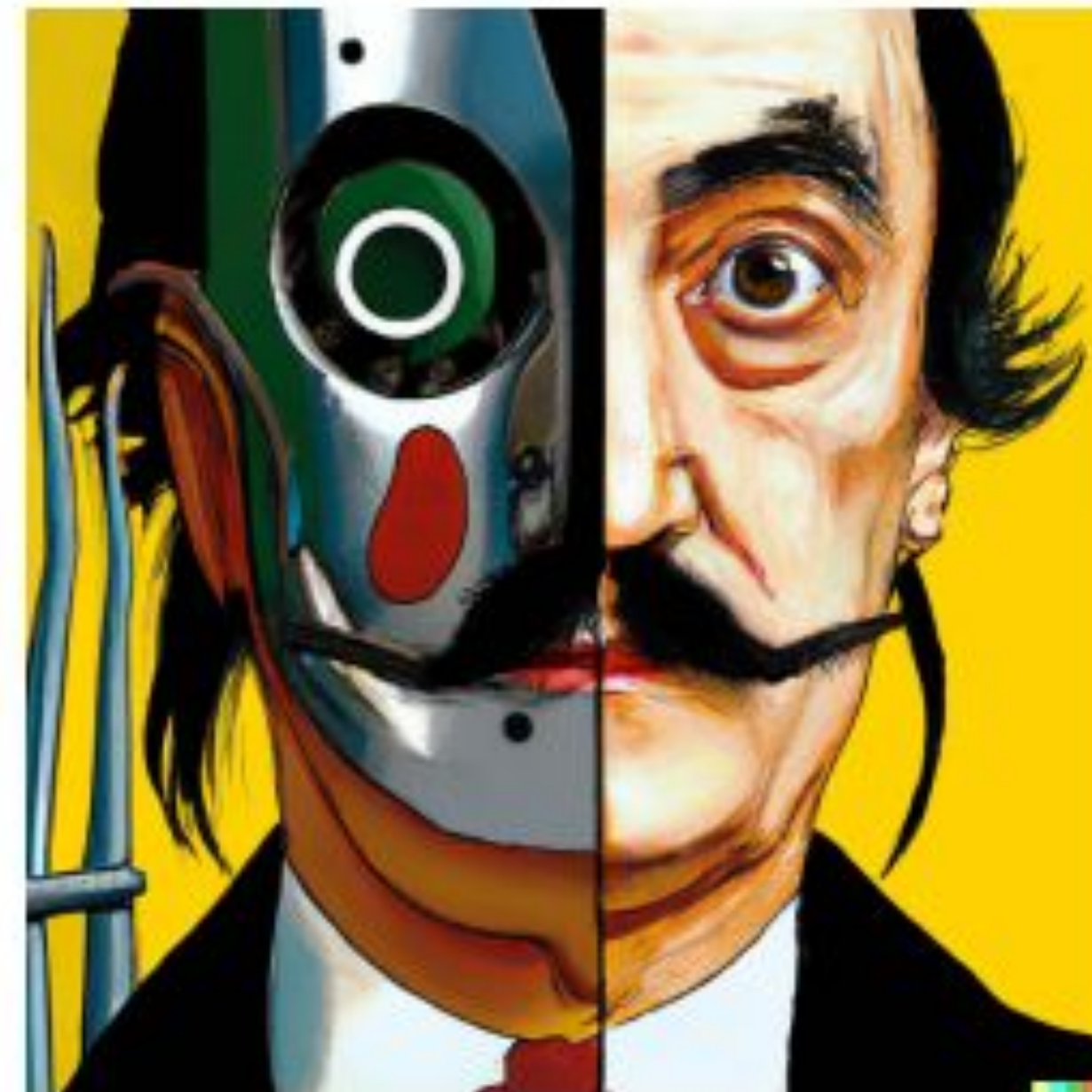**_Dall-e 2_**



an espresso machine that makes coffee from human souls, artstation



vibrant portrait painting of Salvador Dalí with a robotic half face

text
encoder

prior

decoder