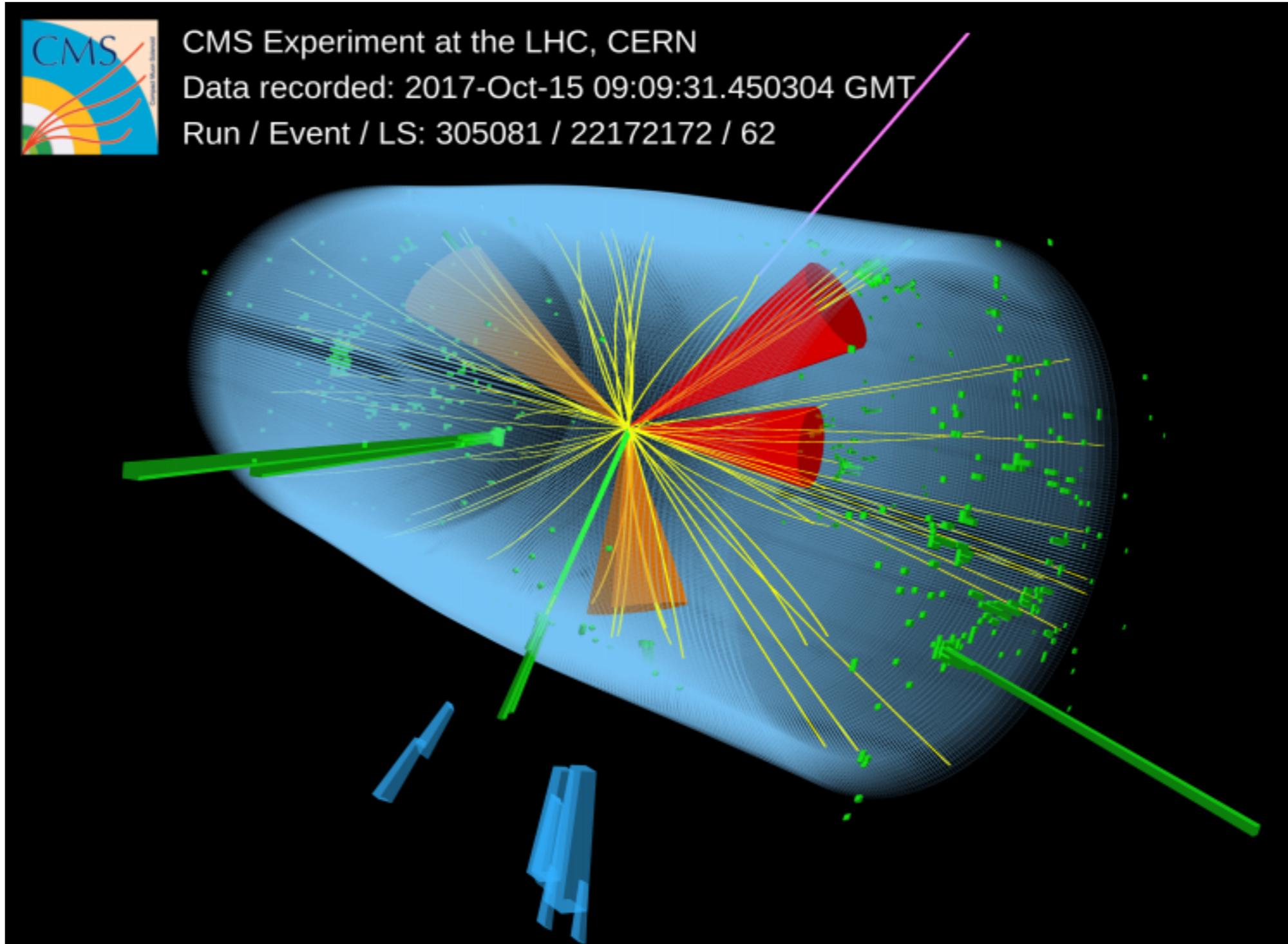


Autoencoders for anomaly detection in particle physics

Michael Krämer

(Institute for Theoretical Particle Physics and
Cosmology, RWTH Aachen University)

Particle collisions at the LHC



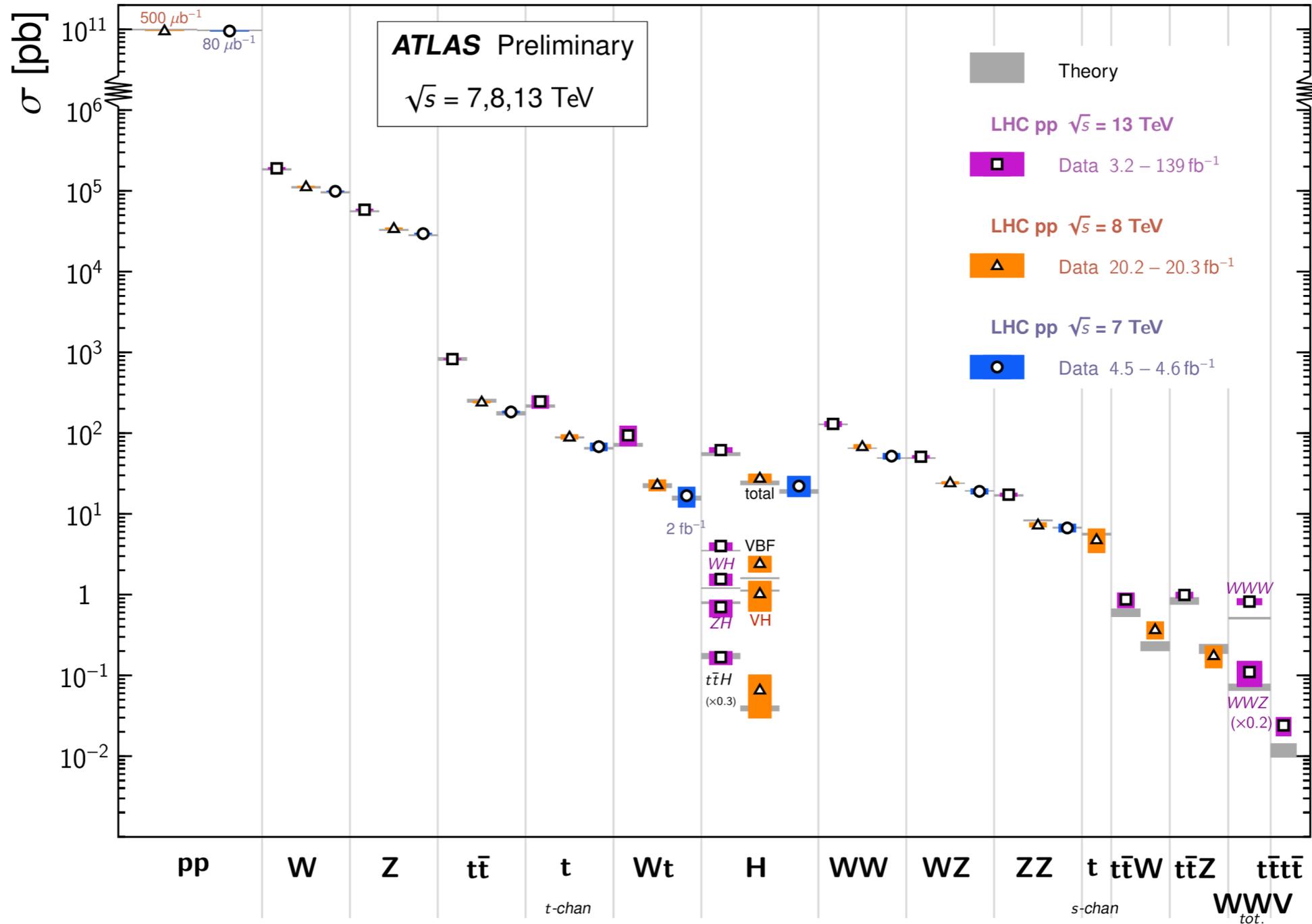
The Standard Model of Particle Physics

$$\begin{aligned}\mathcal{L} = & -\frac{1}{4} F_{\mu\nu} F^{\mu\nu} \\ & + i\bar{\psi} \not{D} \psi + h.c. \\ & + \bar{\psi}_i Y_{ij} \psi_j \phi + h.c. \\ & + |D_\mu \phi|^2 - V(\phi)\end{aligned}$$

The Standard Model of Particle Physics

Standard Model Total Production Cross Section Measurements

Status: February 2022

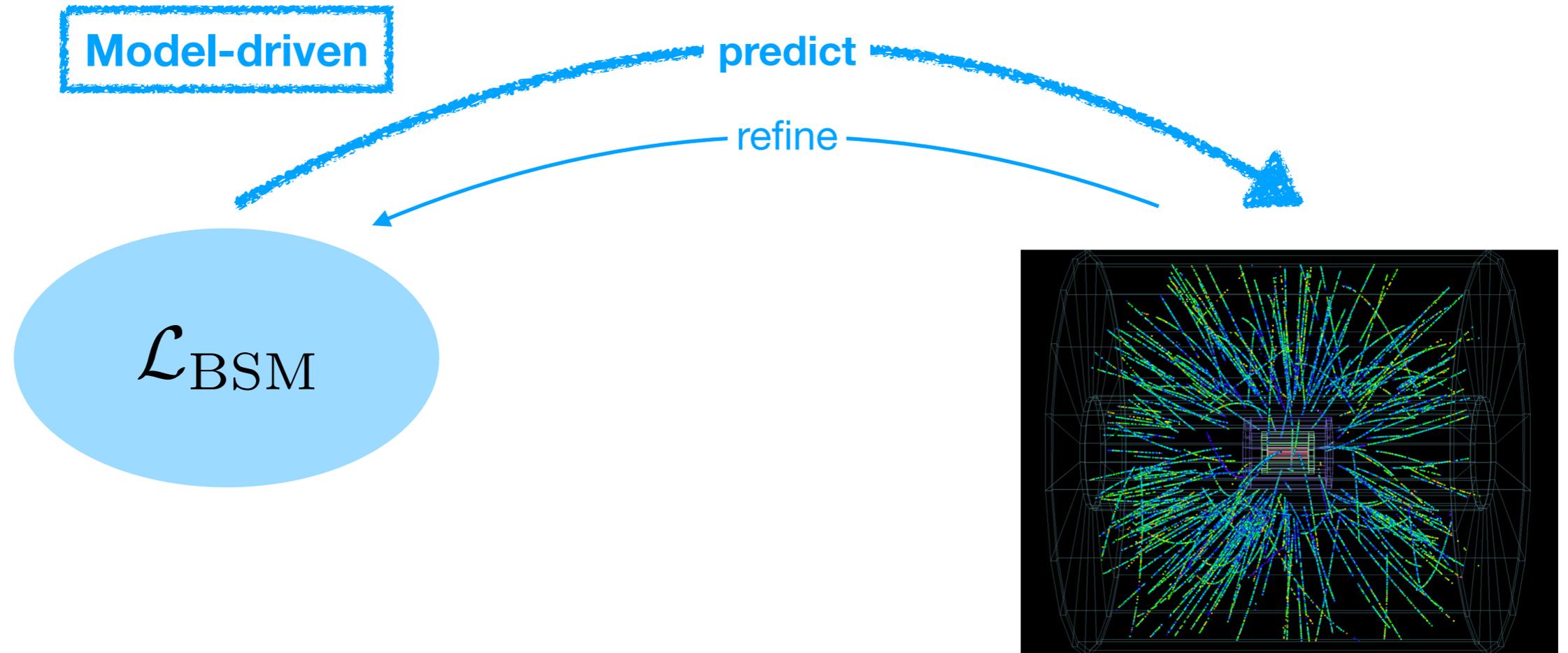


Physics beyond the Standard Model?

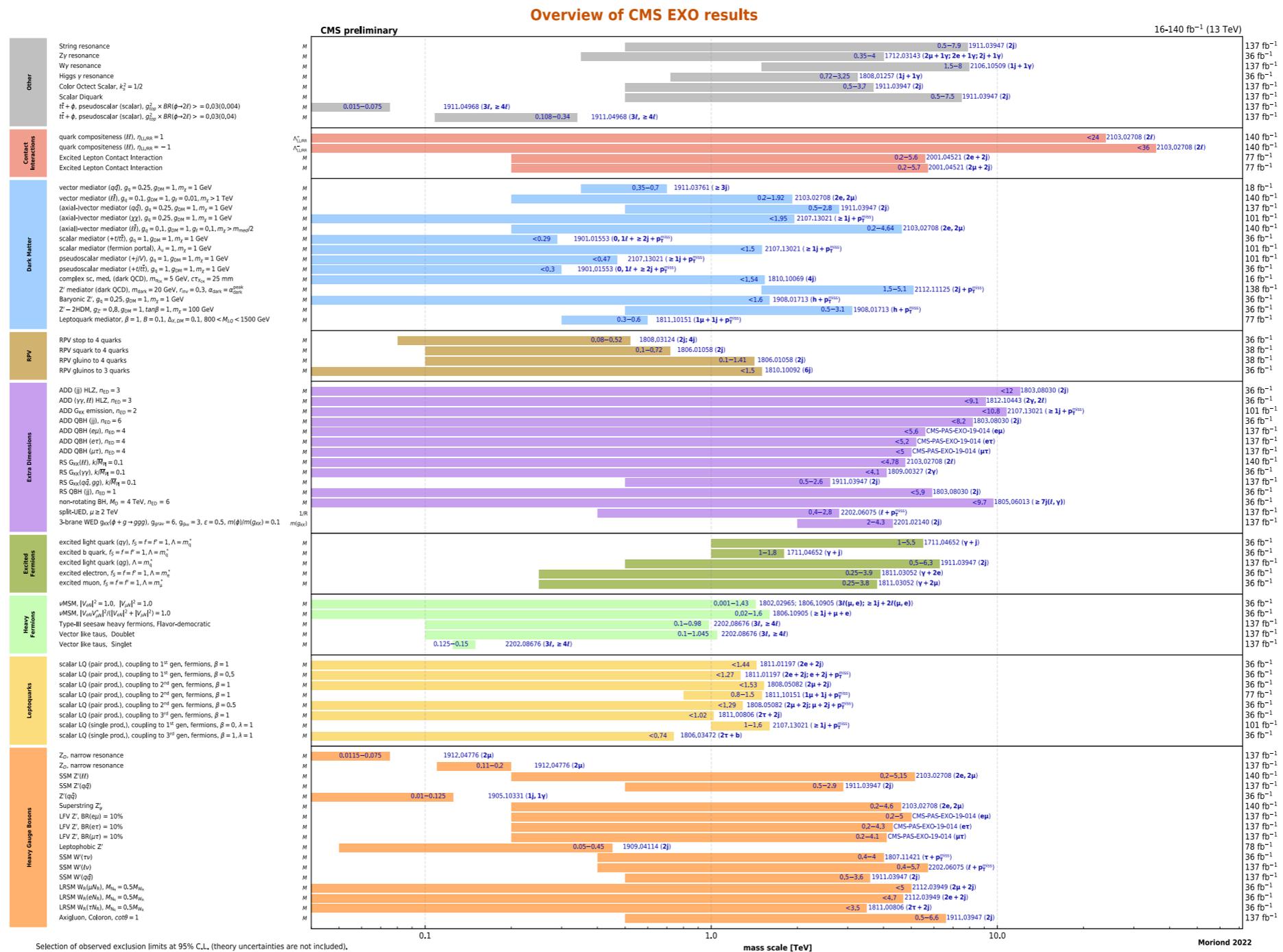
There are potential anomalies and conceptual shortcomings of the SM

- **Particle physics anomalies:** LFV in B-meson decays, $(g-2)_\mu$, ...
- **Cosmic enigmas:** dark matter, matter-antimatter asymmetry, ...
- **Conceptual questions:** origin of EWSB, mass hierarchies, unification,...

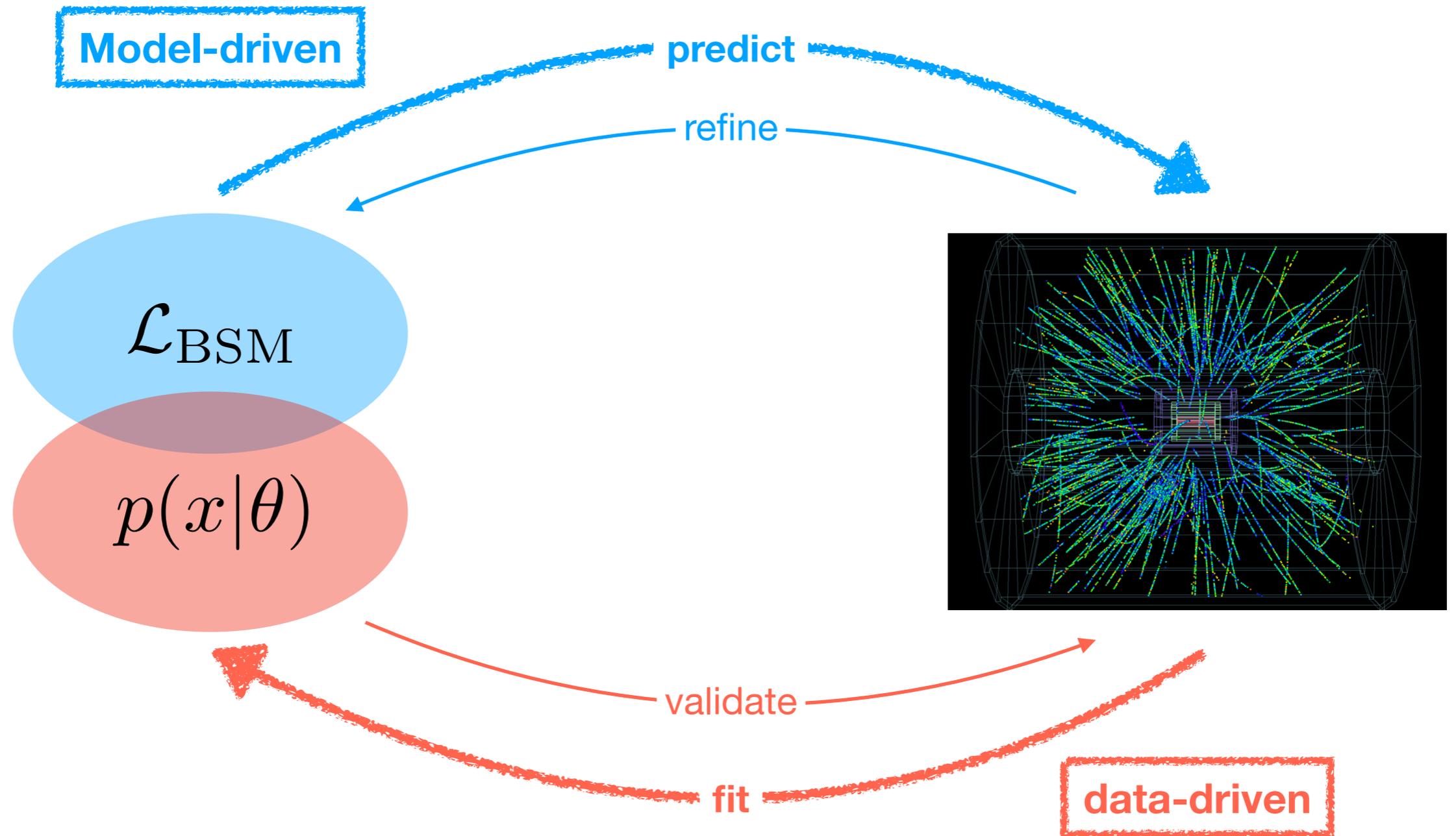
Physics beyond the Standard Model?



Physics beyond the Standard Model?



Physics beyond the Standard Model?

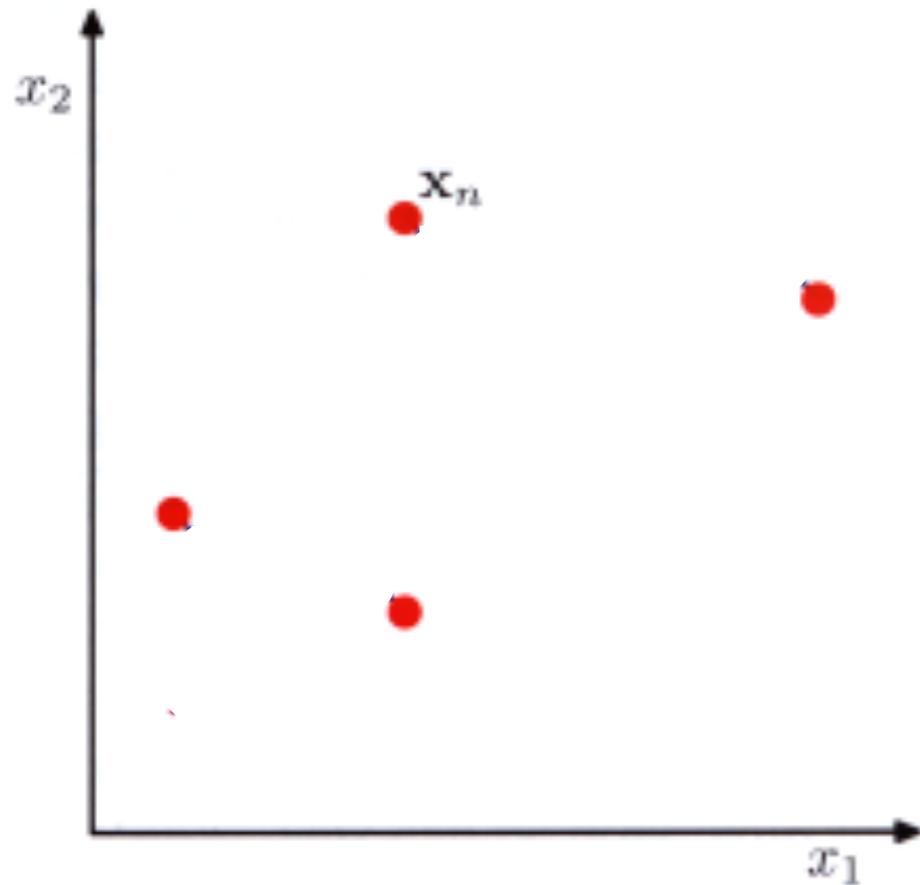


Outline

- Unsupervised learning
- Autoencoders
- Autoencoders for anomaly detection
- Anomaly searches in particle physics

The principle component analysis

- ▶ Data reduction
- ▶ Decorrelation of features



The principle component analysis

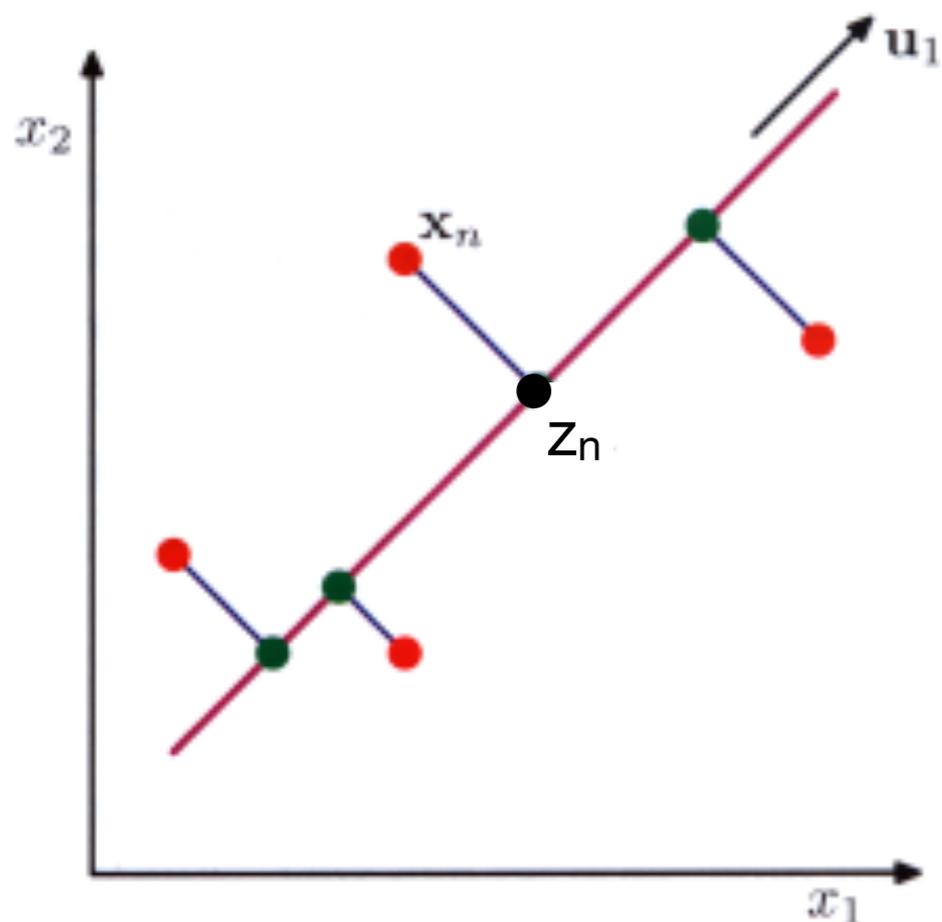
- ▶ Data reduction
- ▶ Decorrelation of features

Search a mapping

$\{\mathbf{x}_n\} \in \mathbb{R}^D \rightarrow \{\mathbf{z}_n\} \in \mathbb{R}^M$ with $M < D$
with minimal loss of information.

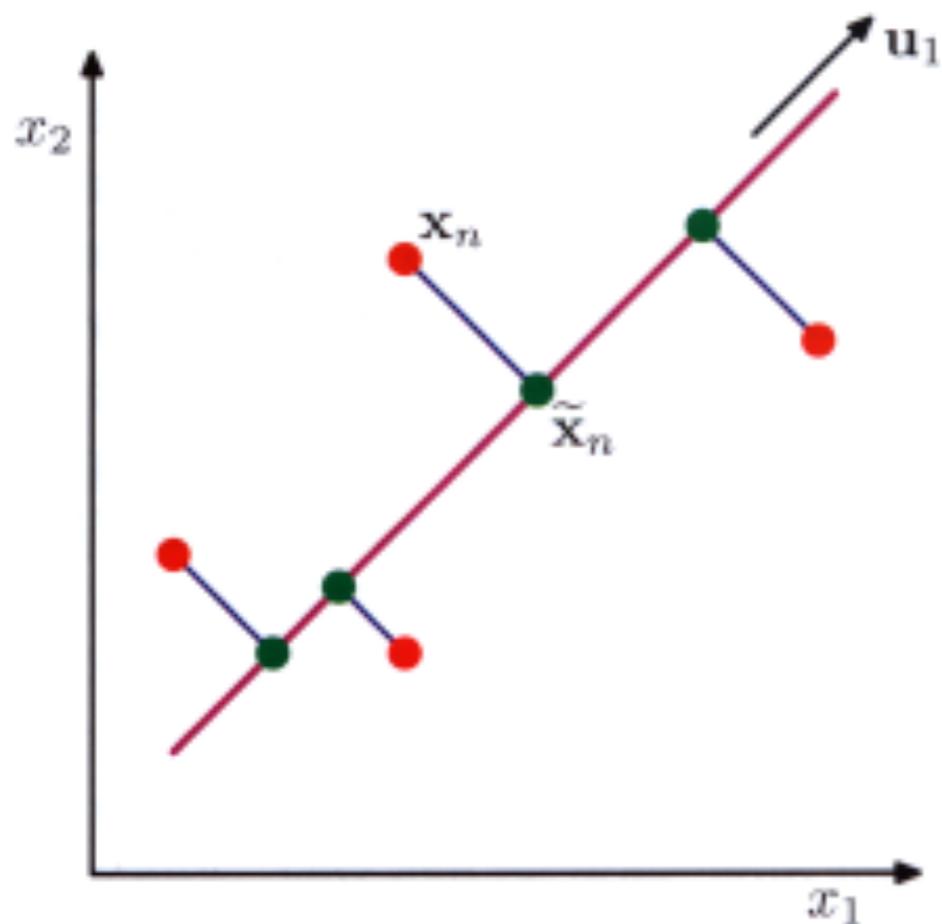
Want to find an **encoding function** $f(\mathbf{x}) = \mathbf{z}$ and a **decoding function** $g(\mathbf{z}) = \tilde{\mathbf{x}}$ such that the reconstruction error $\|\mathbf{x} - \tilde{\mathbf{x}}\|^2 = \|\mathbf{x} - g(\mathbf{z})\|^2$ is minimal.

PCA: choose $g(\mathbf{z}) = D\mathbf{z} \rightarrow f(\mathbf{x}) = D^T\mathbf{x}$



The principle component analysis

- ▶ Data reduction
- ▶ Decorrelation of features



Compute data covariance matrix

$$\mathbf{S} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \bar{\mathbf{x}})(\mathbf{x}_n - \bar{\mathbf{x}})^T$$

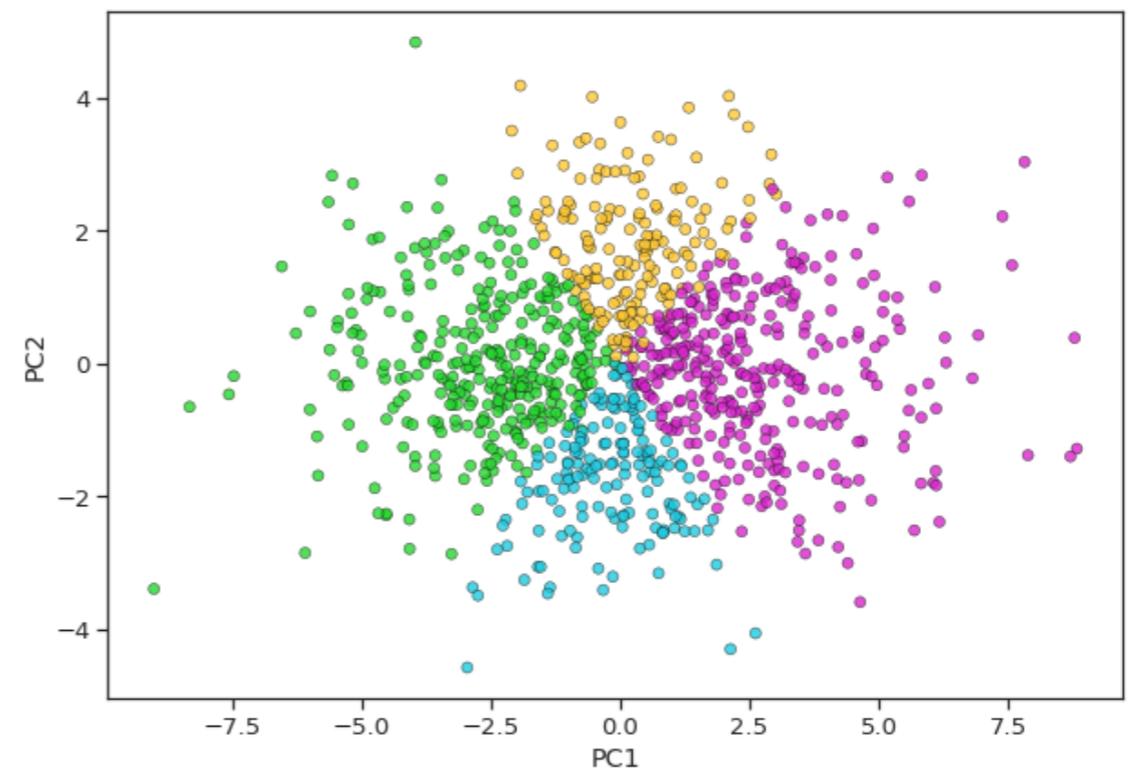
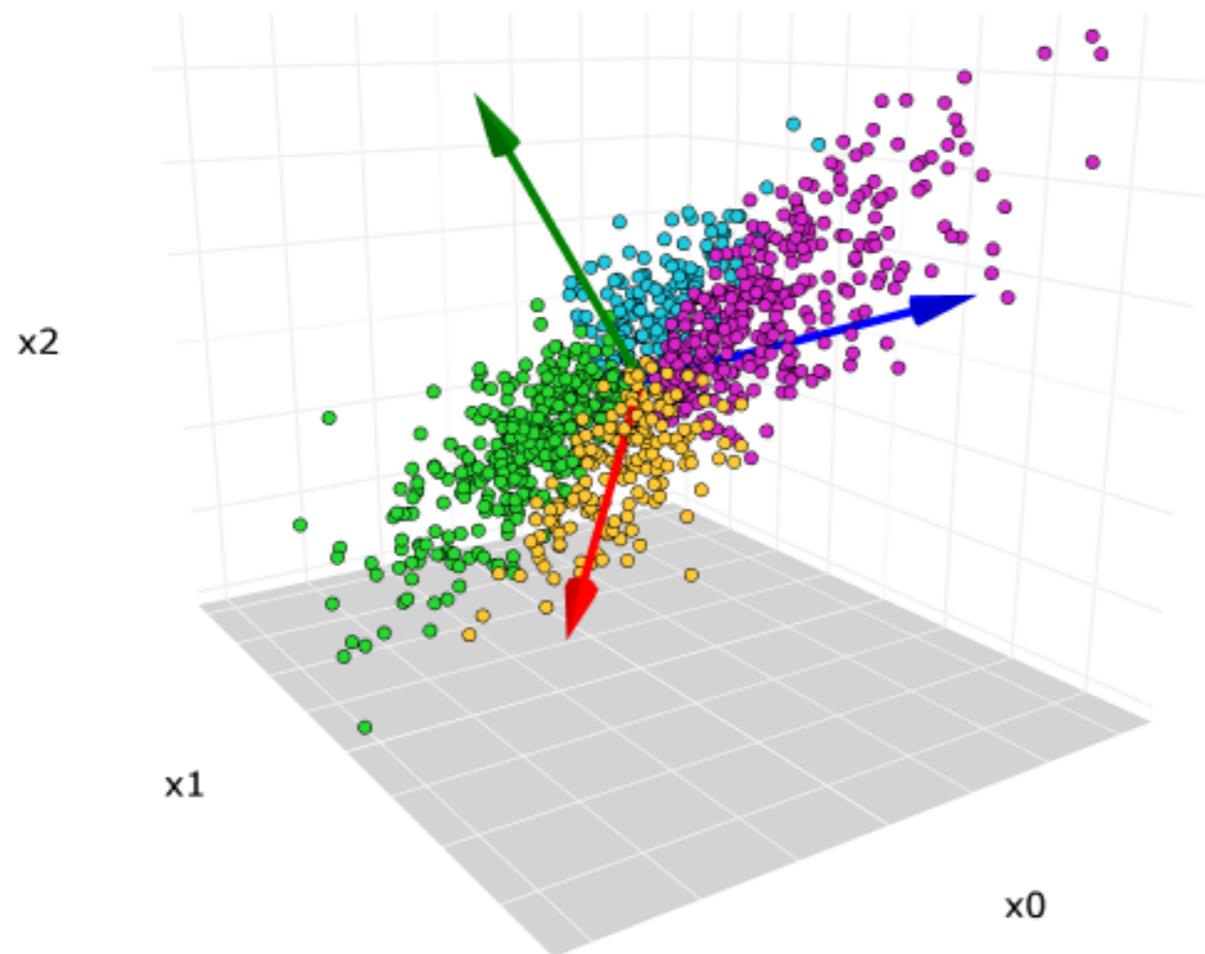
and its eigenvector decomposition.

Vector \mathbf{u}_1 is eigenvector with largest eigenvalue,

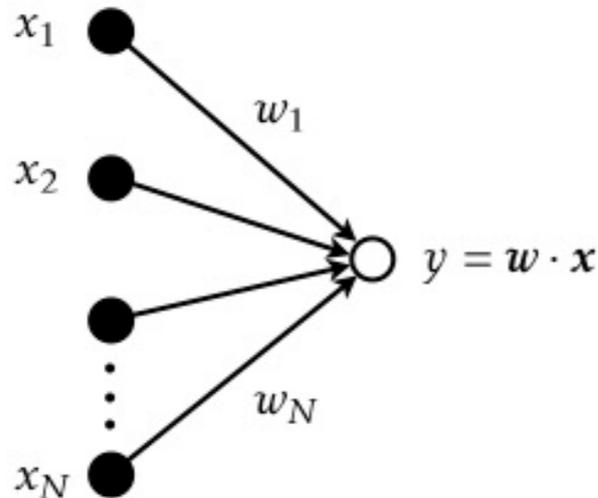
$$\mathbf{u}_1^T \mathbf{S} \mathbf{u}_1 = \lambda_1$$

The principle component analysis

- ▶ Data reduction
- ▶ Decorrelation of features



Hebbian learning



Weight update: $\mathbf{w}' = \mathbf{w} + \delta \mathbf{w}$ with $\delta \mathbf{w} = \eta y \mathbf{x}$.

The output $|y|$ becomes the larger, the more often an input feature occurs in the data.

B. Mehlig: <https://arxiv.org/abs/1901.05639v4>

Can write this as DGL: $\tau \frac{d\mathbf{w}}{dt} = \langle y \mathbf{x} \rangle = \langle \mathbf{x} \cdot \mathbf{x} \rangle \mathbf{w}$ with $\tau \propto 1/\eta$

Writing \mathbf{w} in terms of the eigenvectors of the covariance matrix, $\mathbf{w} = \sum_i c_i(t) \mathbf{u}_i$, we get:

$$\mathbf{w} = \sum_i c_i(0) e^{\lambda_i t / \tau} \mathbf{u}_i, \text{ and thus } \mathbf{w} \propto \mathbf{u}_1 \text{ for large } t \gg \tau.$$

Hebbian learning implements the principal component analysis.

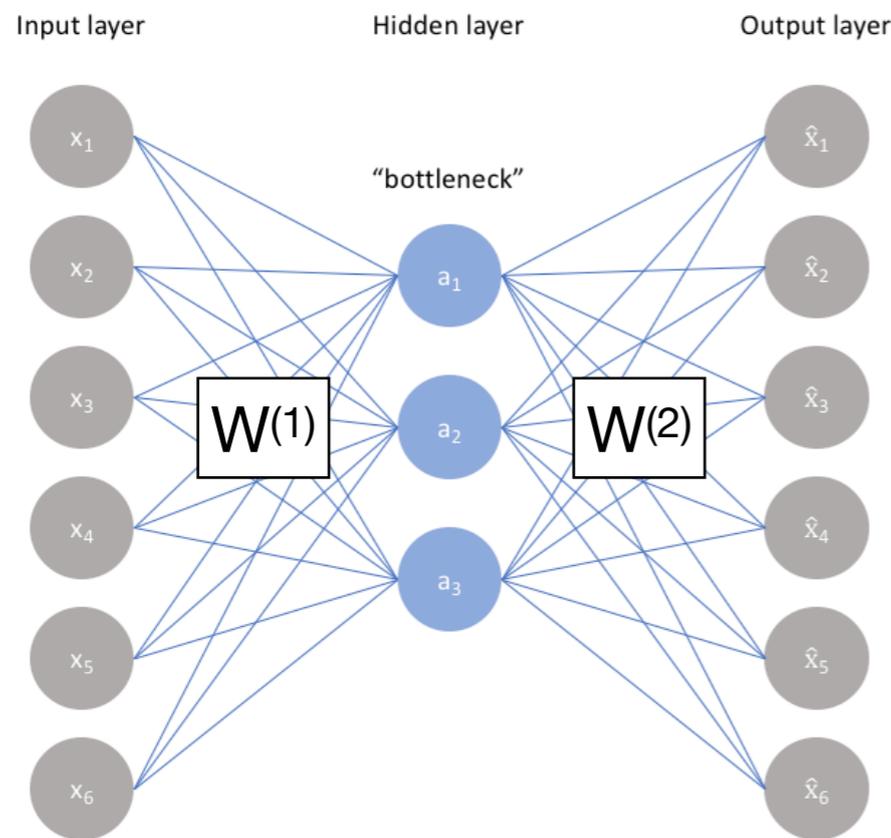
Outline

- Unsupervised learning
- **Autoencoders**
- Autoencoders for anomaly detection
- Anomaly searches in particle physics

Autoencoders

Let us try to implement a principal component analysis with a neural network.

Recall: we need an **encoding function** $f(\mathbf{x}) = D^T \mathbf{x}$ and a **decoding function** $g(\mathbf{z}) = D \mathbf{z}$ such that $\|\mathbf{x} - g(\mathbf{z})\|^2$ is minimal.



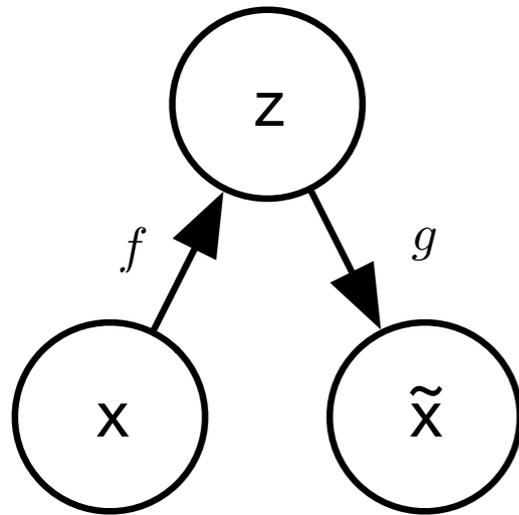
$$\mathbf{z} = f(\mathbf{x}) = W^{(1)} \mathbf{x}$$

and

$$\tilde{\mathbf{x}} = g(\mathbf{z}) = W^{(2)} \mathbf{z} = W^{(2)} W^{(1)} \mathbf{x}$$

Training the weights $W^{(1)}$ and $W^{(2)}$ to minimise the mean square error between input and output, the linear neural network (nearly) implements a principal component analysis.

Autoencoders

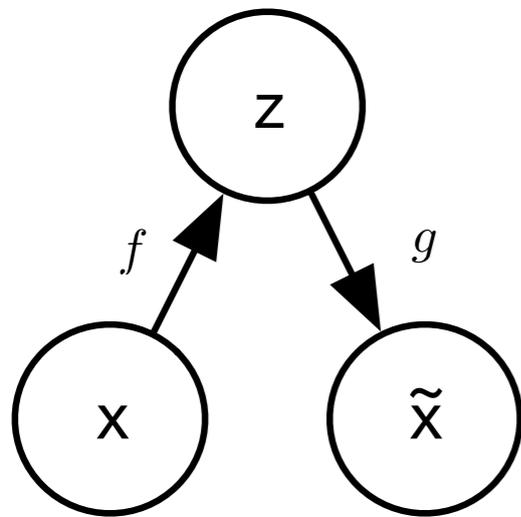


An autoencoder is a neural network that tries to learn an approximation to the identity function $\tilde{\mathbf{x}} = g(f(\mathbf{x})) \approx \mathbf{x}$

Learning the identity function itself is not very useful, but by **placing constraints** on the network, such as by limiting the number of hidden units, one can discover interesting structures about the data:

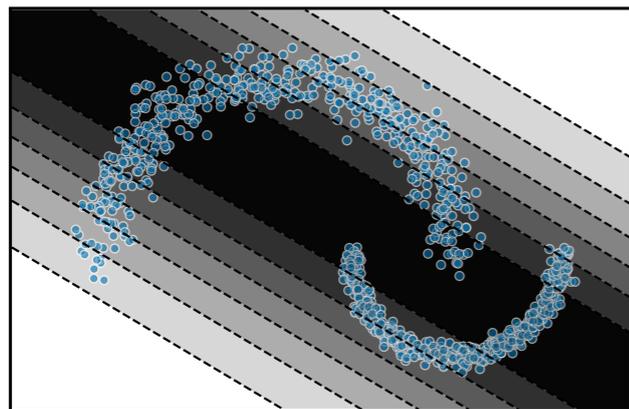
- latent space z has lower dimension than x ;
- f or g have low capacity (e.g. linear g);
- introduce regularisation, e.g. sparse autoencoders.

Autoencoders

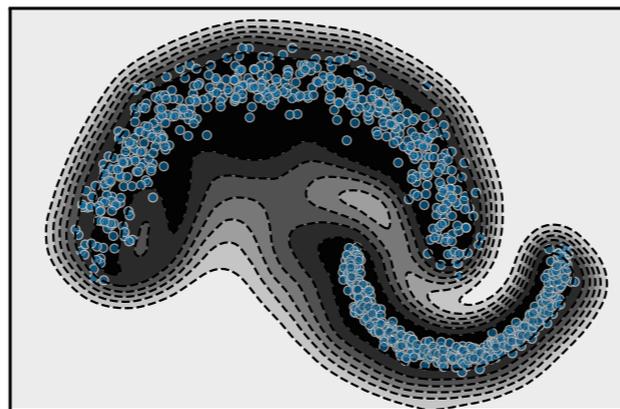


An autoencoder is a neural network that tries to learn an approximation to the identity function $\tilde{\mathbf{x}} = g(f(\mathbf{x})) \approx \mathbf{x}$

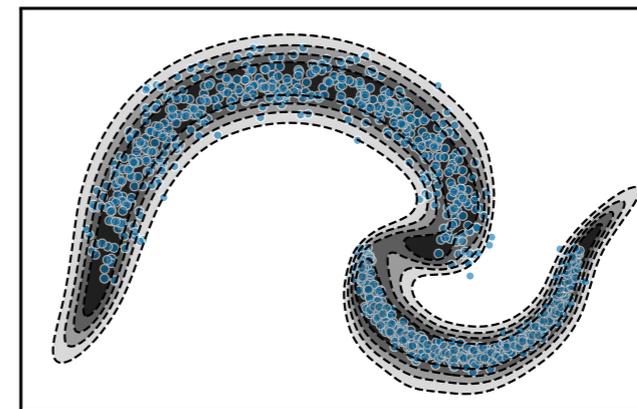
PCA (AUC=66.8)



kPCA (AUC=94.0)



AE (AUC=97.9)



Ruff et al., <https://arxiv.org/abs/2009.11732>

Variational autoencoders

One can combine the idea of an autoencoder with the concept of generative modeling.

Bayesian view: $p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x}|\mathbf{z})d\mathbf{z}$

The diagram shows the equation $p(\mathbf{x}) = \int p(\mathbf{z})p(\mathbf{x}|\mathbf{z})d\mathbf{z}$. Three arrows point from labels below to terms in the equation: 'evidence' points to $p(\mathbf{x})$, 'latent prior' points to $p(\mathbf{z})$, and 'likelihood' points to $p(\mathbf{x}|\mathbf{z})$.

Goal: maximise $p_\theta(\mathbf{x})$ by learning $p_\theta(\mathbf{z})$ and $p_\theta(\mathbf{x}|\mathbf{z})$: $\theta^* = \operatorname{argmin}_\theta \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} (-\ln p_\theta(\mathbf{x}))$

Difficult to evaluate in practice \rightarrow introduce **recognition model** $q_\theta(\mathbf{z}|\mathbf{x})$ as an approximation to true posterior $p(\mathbf{z}|\mathbf{x})$

AE terminology: $q_\theta(\mathbf{z}|\mathbf{x}) \rightarrow$ encoder
 $p_\theta(\mathbf{x}|\mathbf{z}) \rightarrow$ decoder

In a **variational autoencoder**, $q_\theta(\mathbf{z}|\mathbf{x})$ is a multivariate Gaussian, parametrised by a neural network.

Variational autoencoders

How can we use the recognition model $q(\vec{z}|\vec{x})$ to maximize likelihood?

$$\begin{aligned}\ln p(\vec{x}) &= \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \ln \left(p(\vec{x}) \frac{p(\vec{z}|\vec{x})}{p(\vec{z}|\vec{x})} \right) \\ &= \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \left(\ln \left(\frac{p(\vec{z}, \vec{x})}{q(\vec{z}, \vec{x})} \right) + \ln \left(\frac{q(\vec{z}|\vec{x})}{p(\vec{z}|\vec{x})} \right) \right) \\ &= \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \ln \left(\frac{p(\vec{z}, \vec{x})}{q(\vec{z}, \vec{x})} \right) + \underbrace{\text{KL} \left(q(\vec{z}|\vec{x}) \parallel p(\vec{z}, \vec{x}) \right)}_{\geq 0} \\ &\geq \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \ln \left(\frac{p(\vec{z}, \vec{x})}{q(\vec{z}, \vec{x})} \right)\end{aligned}$$

\leadsto evidence lower bound (ELBO)

Variational autoencoders

Rather than maximizing log-likelihood, minimize negative log-likelihood:

$$-\ln p(\vec{x}) \leq \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \ln \left(\frac{q(\vec{z}|\vec{x})}{p(\vec{z}, \vec{x})} \right)$$

$$= \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \ln \left(\frac{q(\vec{z}|\vec{x})}{p(\vec{z}) p(\vec{x}|\vec{z})} \right)$$

$$= \mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} \left(\ln \left(\frac{q(\vec{z}|\vec{x})}{p(\vec{z})} \right) - \ln p(\vec{x}|\vec{z}) \right)$$

$$= \text{KL}(q(\vec{z}|\vec{x}) \| p(\vec{z})) - \underbrace{\mathbb{E}_{\vec{z} \sim q(\vec{z}|\vec{x})} (-\ln p(\vec{x}|\vec{z}))}_{\sim \text{reconstruction error}}$$

as autoencoder structure: $\vec{x} \xrightarrow{q} \vec{z} \xrightarrow{p} \vec{x}$

Variational autoencoders

Variational autoencoder:

Minimize bound to negative log likelihood,

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \mathbb{E}_{\vec{x} \sim p_{\text{data}}} (-\ln p_{\theta}(\vec{x}))$$

$$= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N (-\ln p_{\theta}(\vec{x}^{(i)}))$$

$$\approx \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \left[\text{KL}(q_{\theta}(\vec{z}|\vec{x}^{(i)}) \| p(\vec{z})) + \mathbb{E}_{\vec{z} \sim q_{\theta}(\vec{z}|\vec{x}^{(i)})} (-\ln p_{\theta}(\vec{x}^{(i)}|\vec{z})) \right]$$

In VAE, $q_{\theta}(\vec{z}|\vec{x})$ is a multivariate Gaussian, parametrized by neural network:

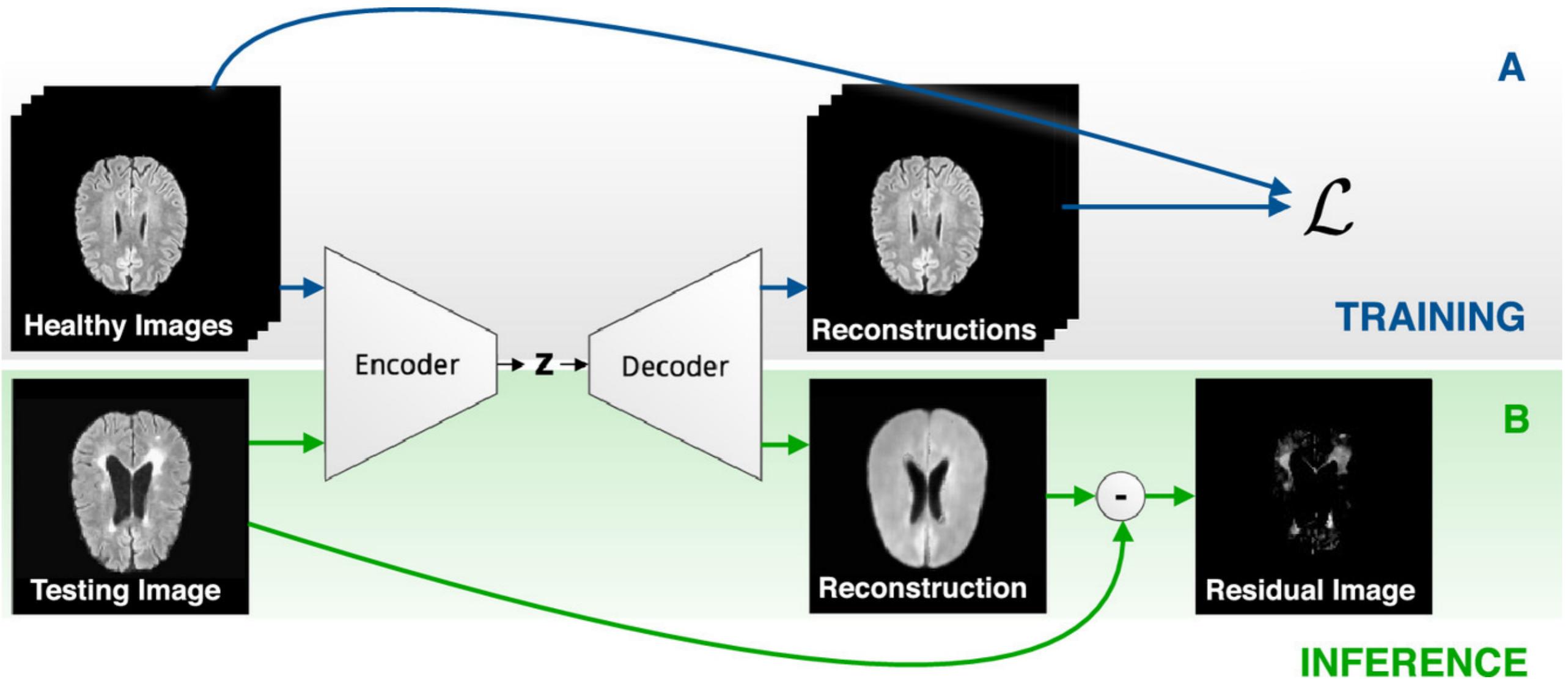
$$q_{\theta}(\vec{z}|\vec{x}) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_{\theta}(\vec{x})|^{1/2}} \exp\left(-\frac{1}{2} (\vec{z} - \vec{\mu}_{\theta}(\vec{x}))^{\top} \Sigma_{\theta}^{-1}(\vec{x}) (\vec{z} - \vec{\mu}_{\theta}(\vec{x}))\right)$$

→ mean $\vec{\mu}$ and covariance Σ_i are functions of the data and determined by a neural network.

Outline

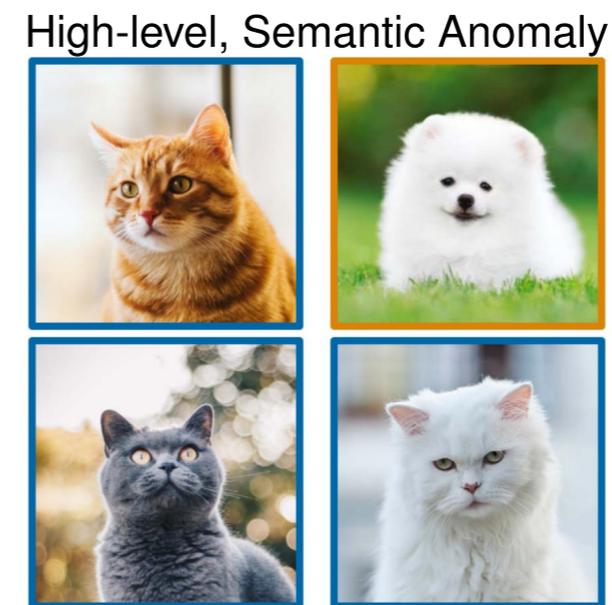
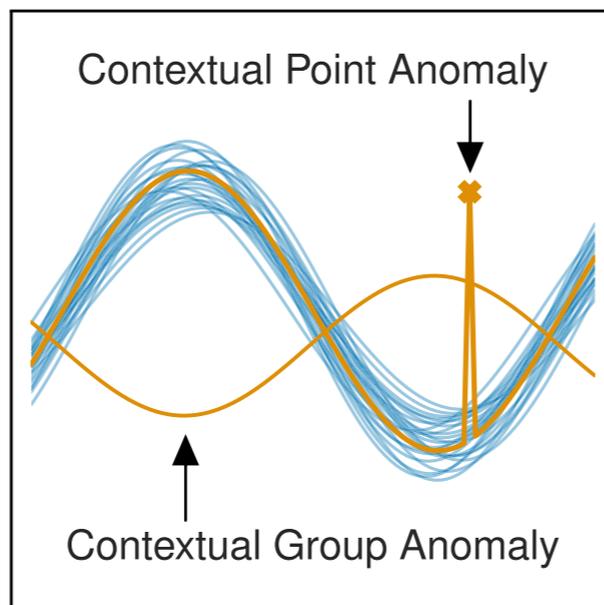
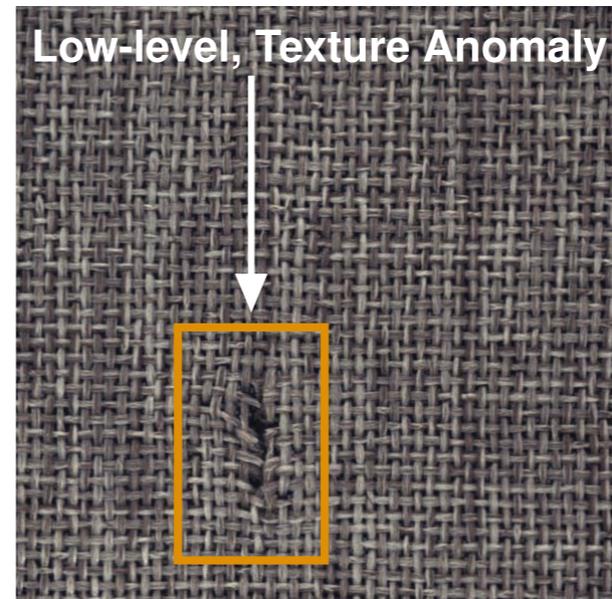
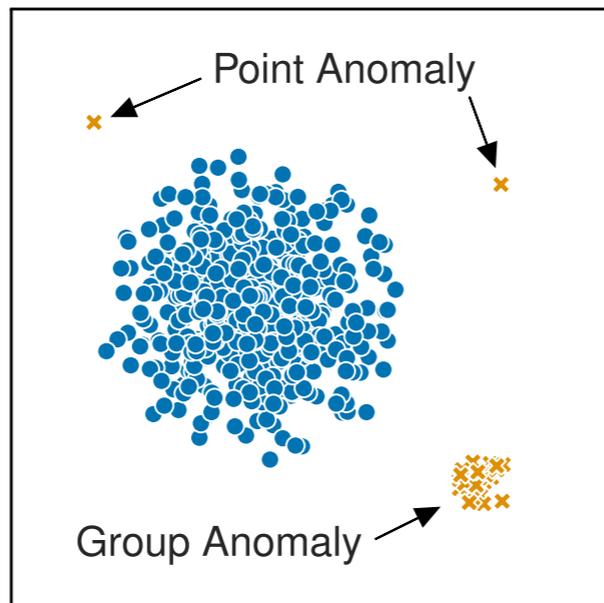
- Unsupervised learning
- Autoencoders
- **Autoencoders for anomaly detection**
- Anomaly searches in particle physics

Anomaly detection with autoencoders



Baur et al., <https://arxiv.org/abs/2004.03271>

Anomaly detection

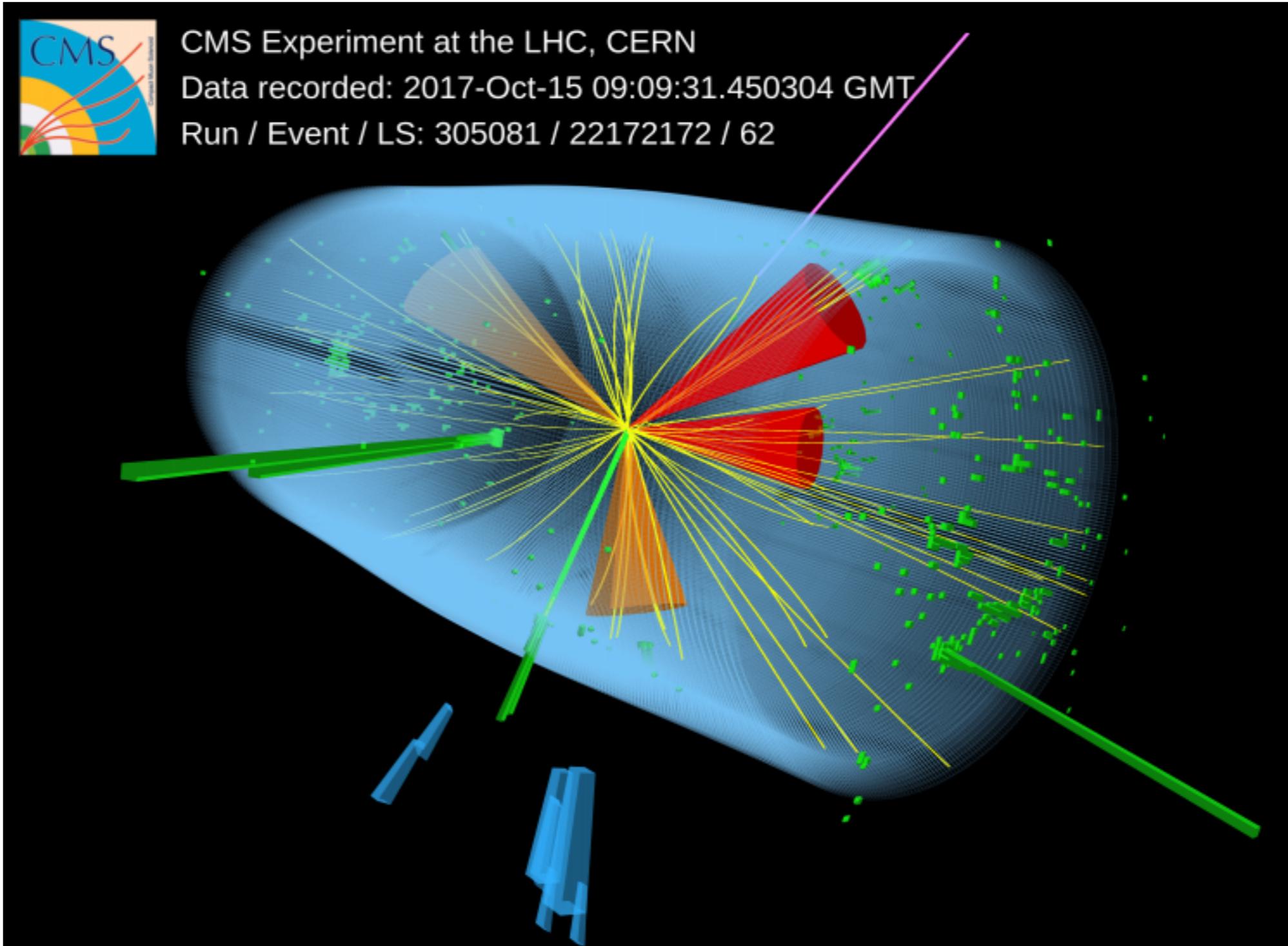


Ruff et al., <https://arxiv.org/abs/2009.11732>

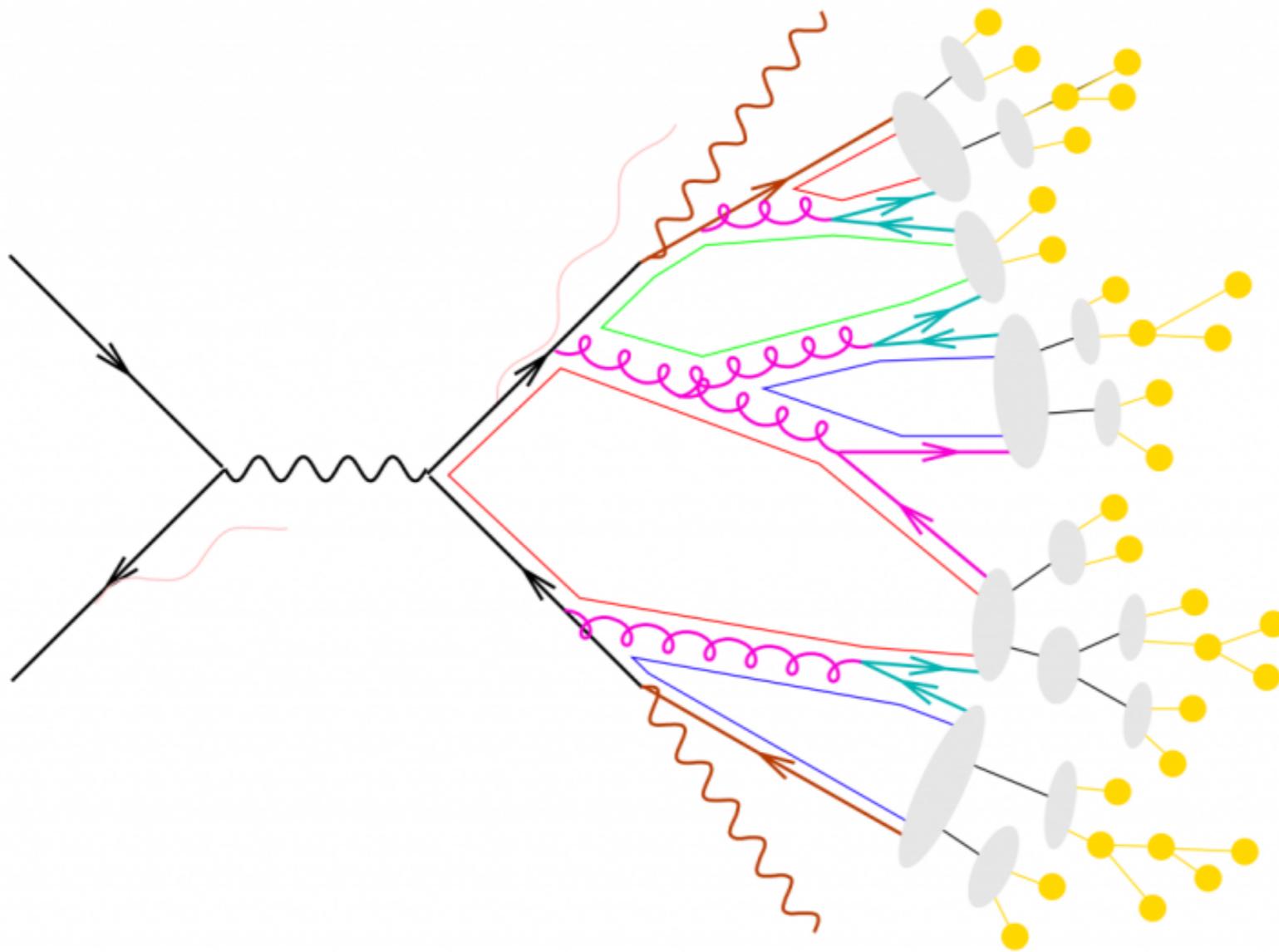
Outline

- Unsupervised learning
- Autoencoders
- Autoencoders for anomaly detection
- **Anomaly searches in particle physics**

Particle collisions at the LHC



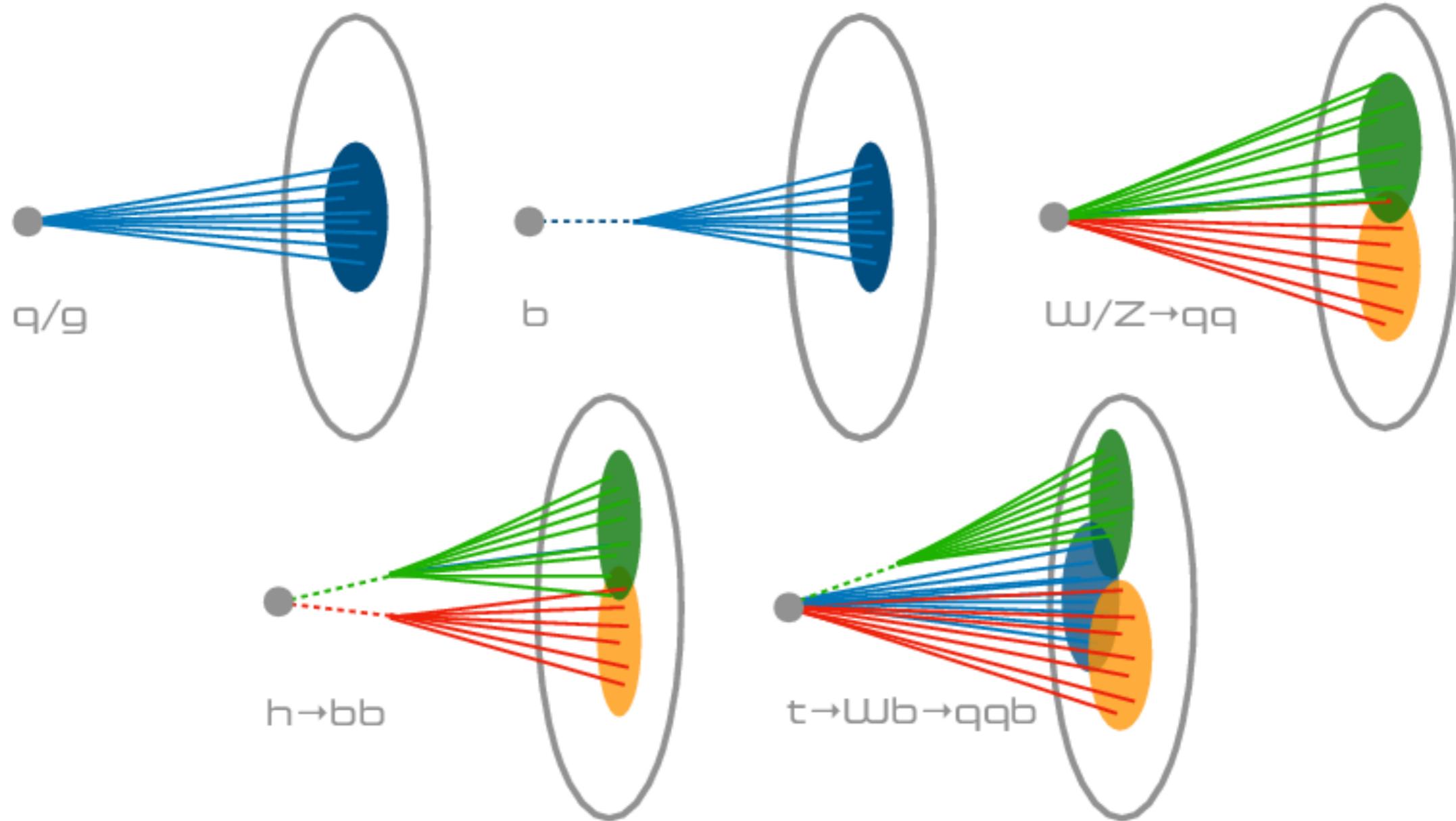
Particle collisions at the LHC



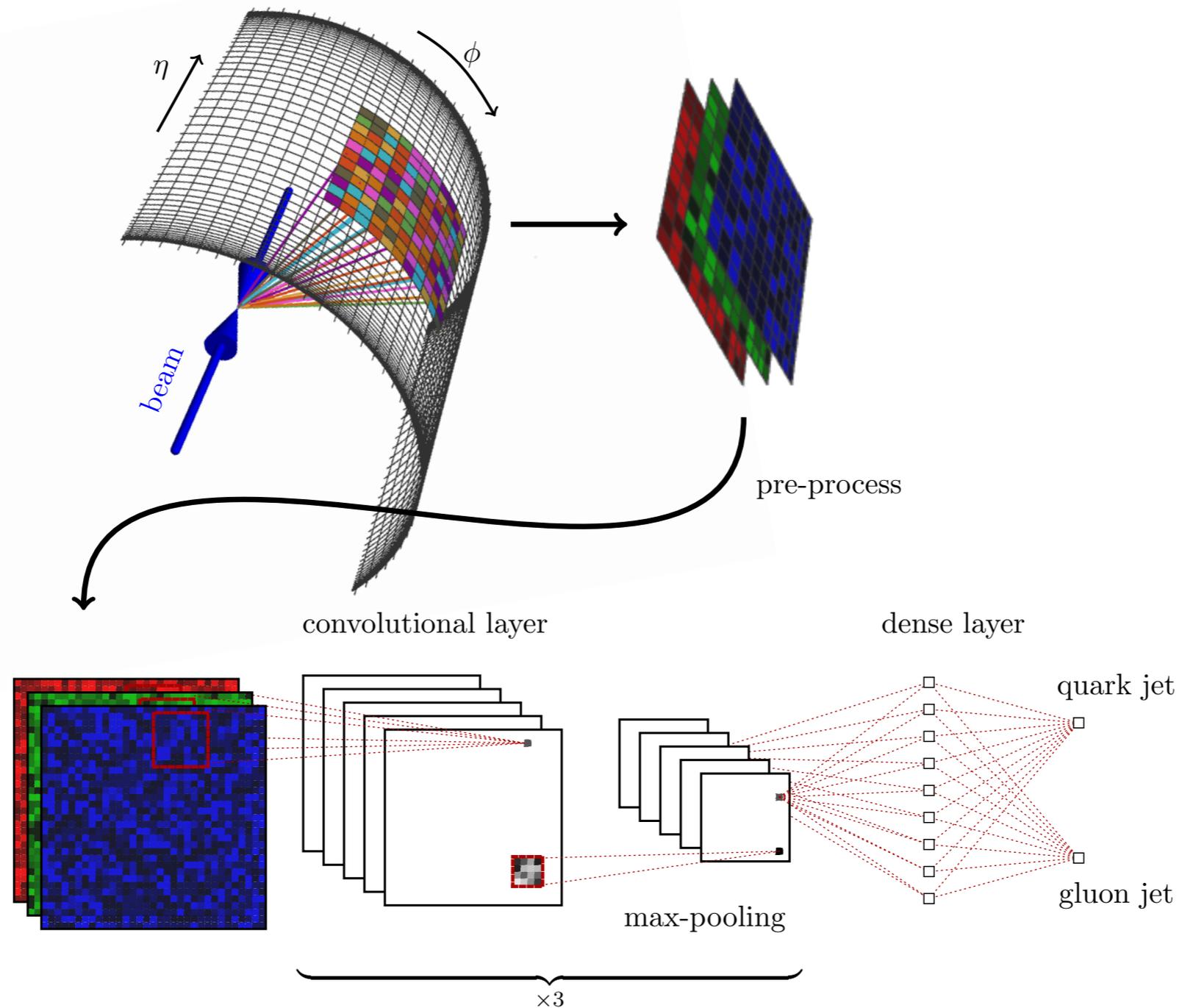
- hard scattering
- (QED) initial/final state radiation
- partonic decays, e.g. $t \rightarrow bW$
- parton shower evolution
- nonperturbative gluon splitting
- colour singlets
- colourless clusters
- cluster fission
- cluster \rightarrow hadrons
- hadronic decays

Figure from Dieter Zeppenfeld

Jets at the LHC

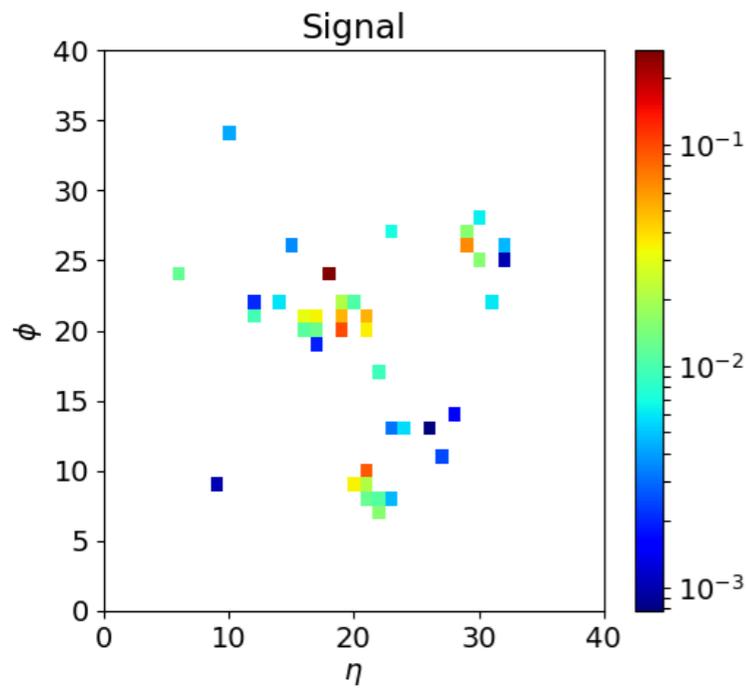


Convolutional neural networks for jet images

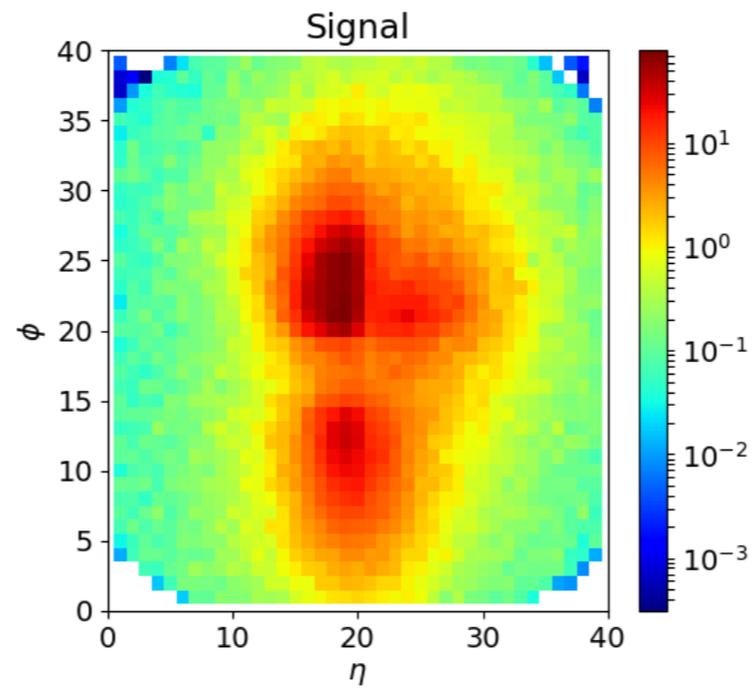


Jet images

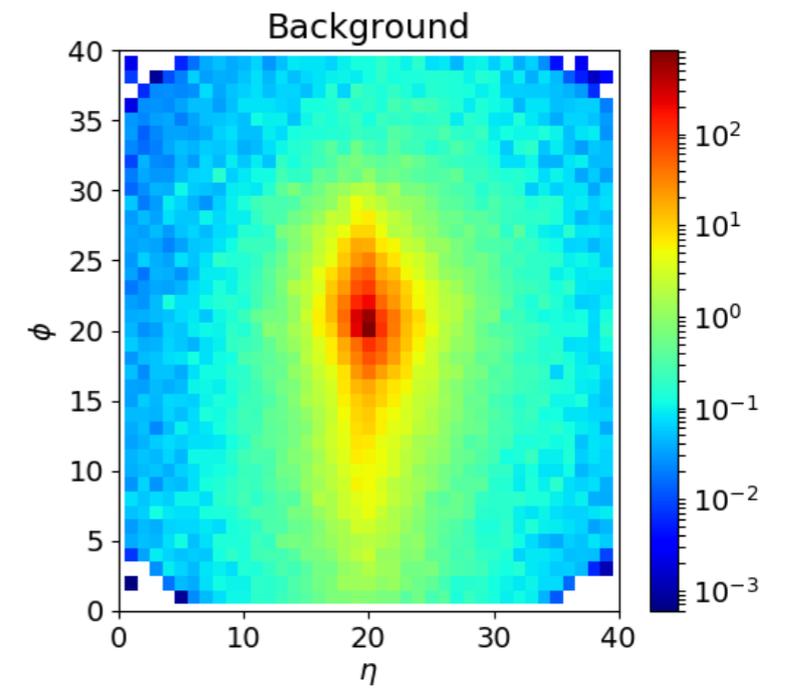
Typical single top-jet



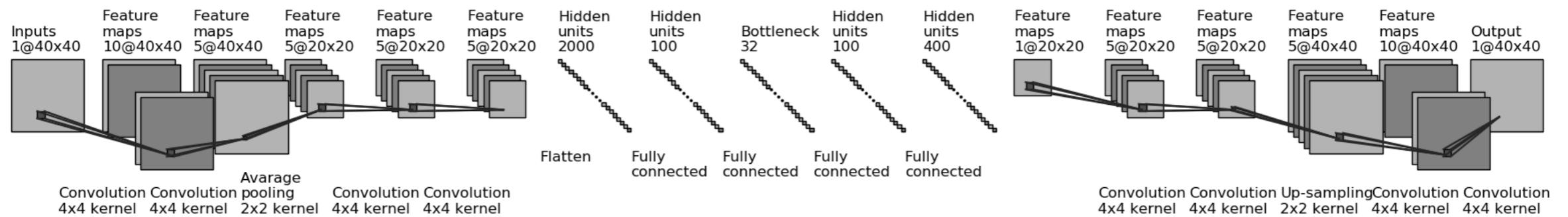
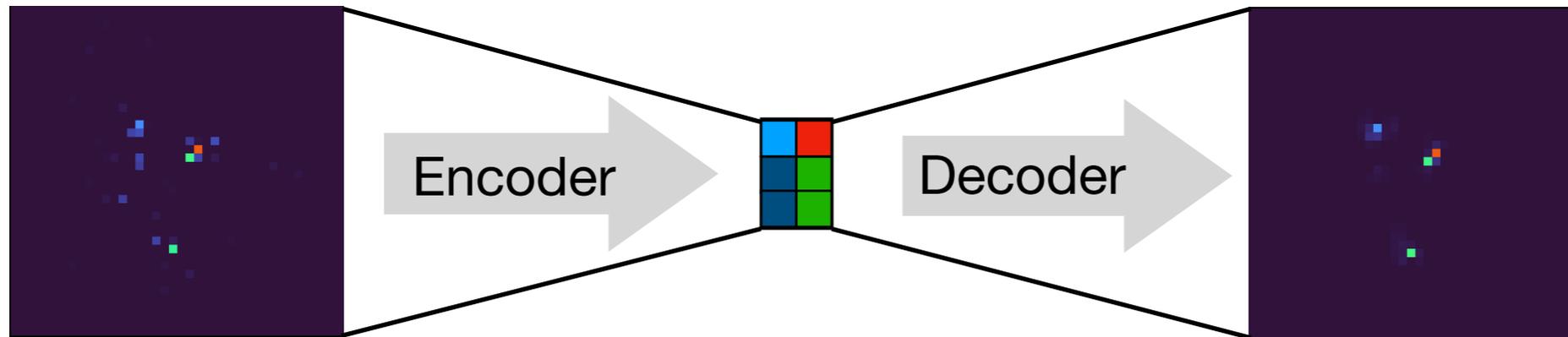
Average top-jet



Average QCD-jet

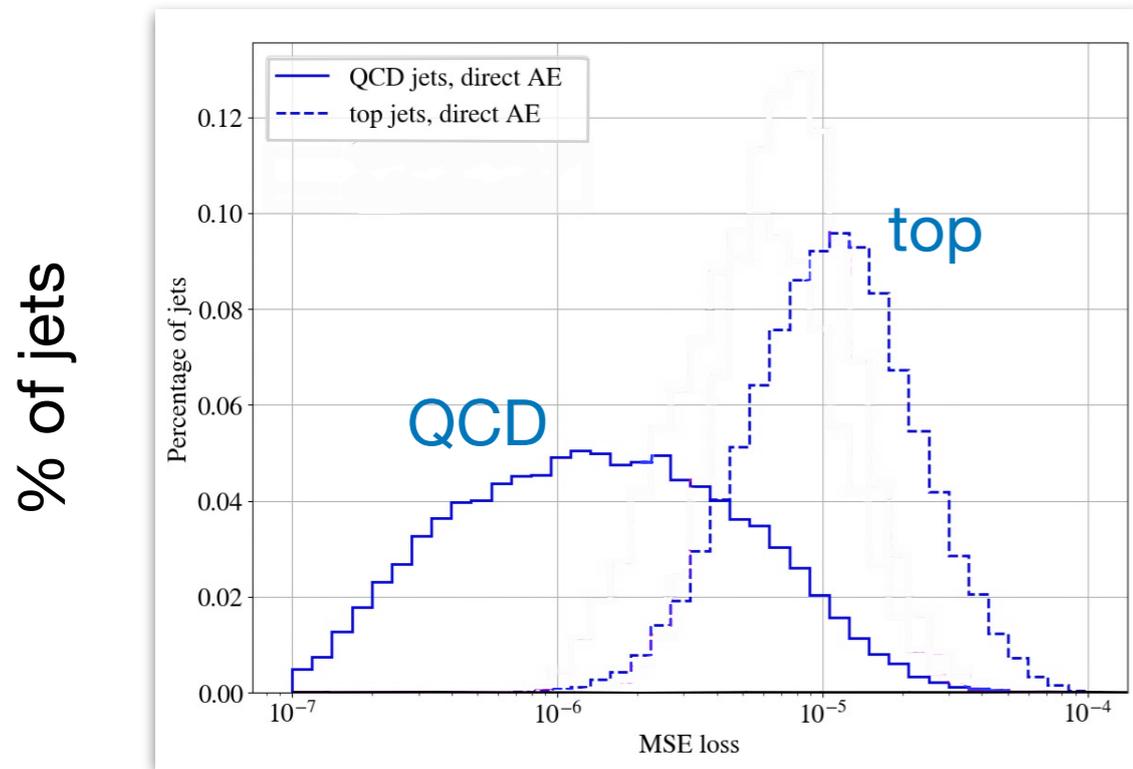
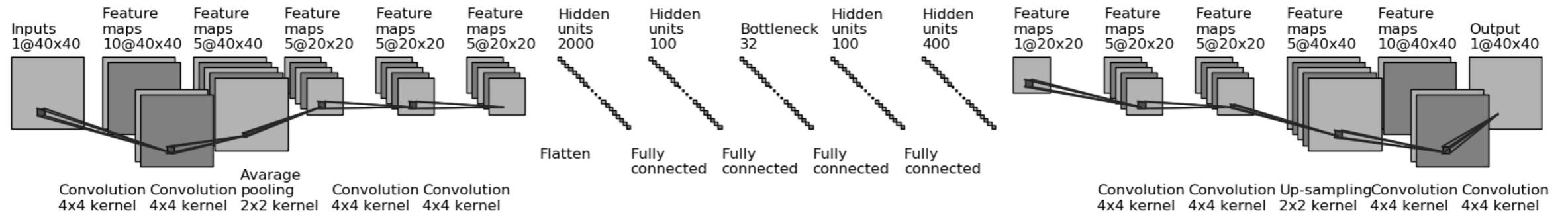


Anomaly detection with autoencoders

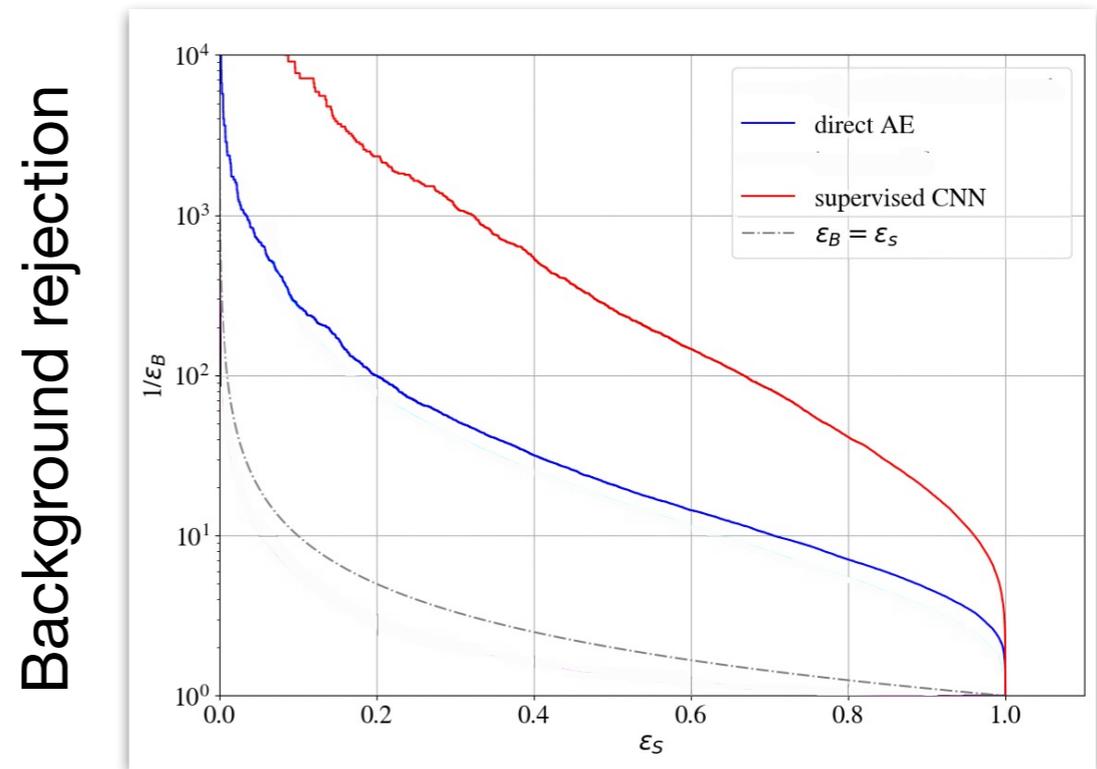


cf Heimgel, Kasieczka, Plehn, Thompson, SciPost Phys. 6, 030 (2019); Farina, Nakai, Shih, PRD 101 (2020)

Anomaly detection with autoencoders



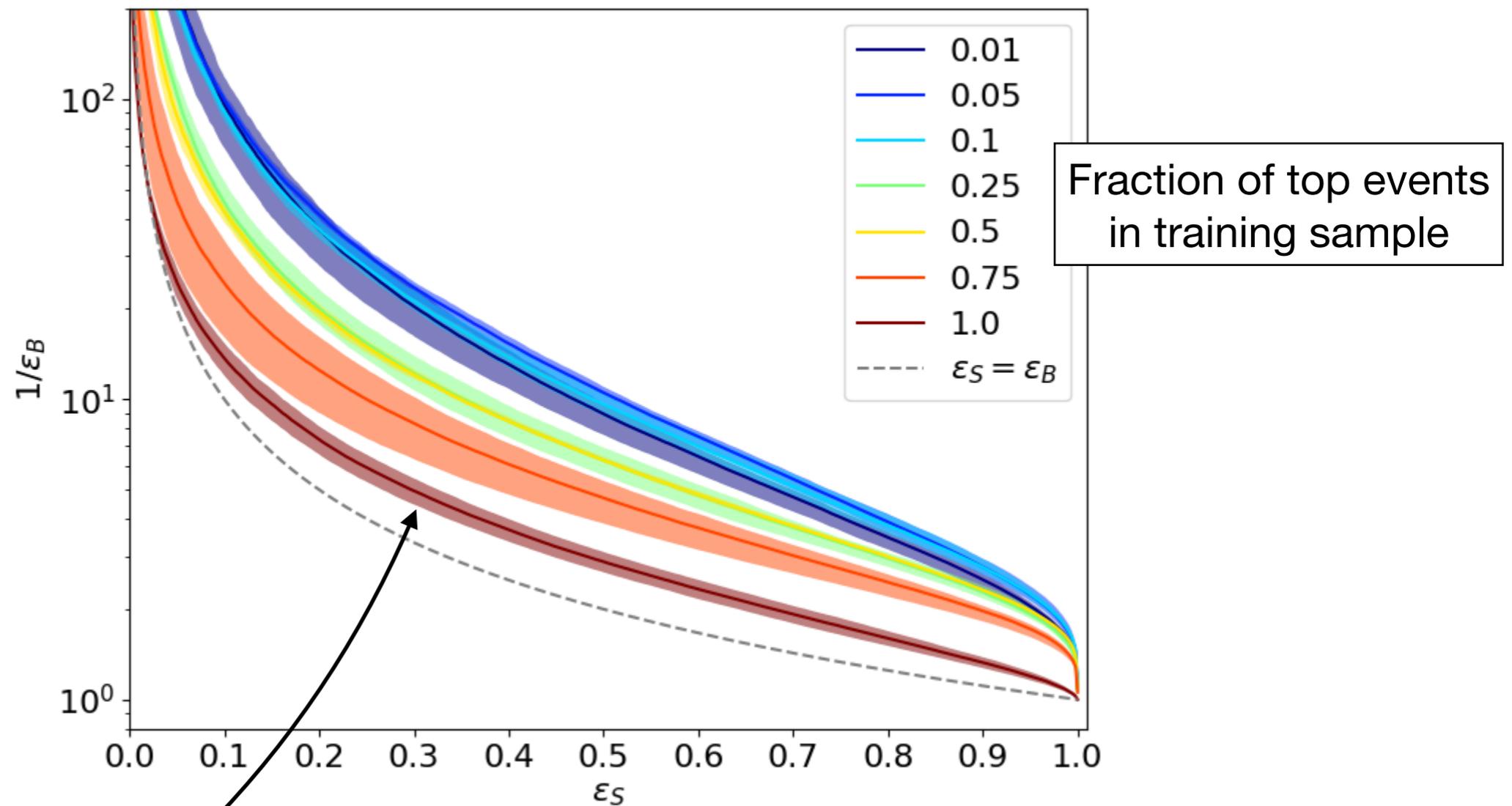
MSE loss



Top tagging efficiency

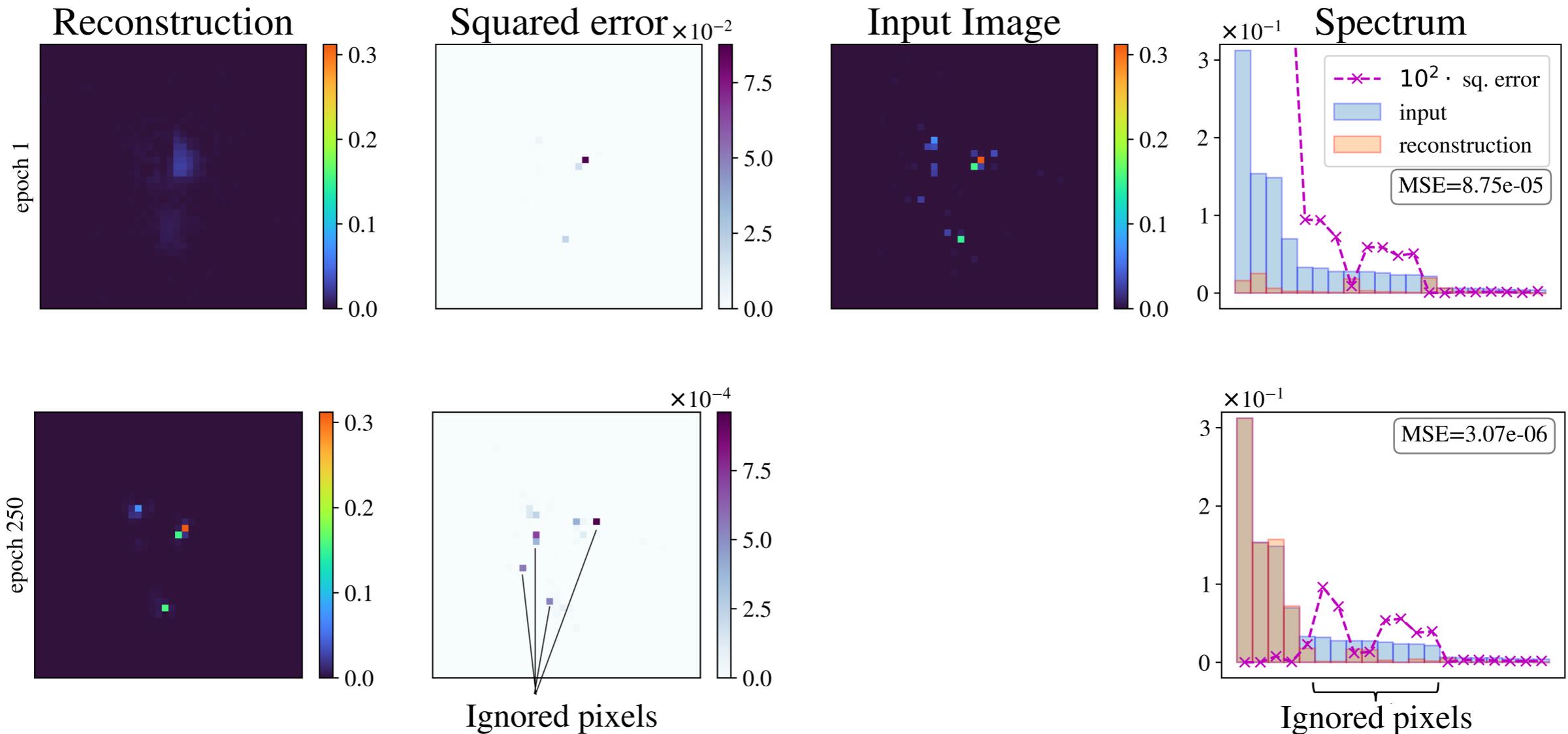
What does the autoencoder learn?

Top tagging: increase the fraction of top events (anomalies) in the training sample:



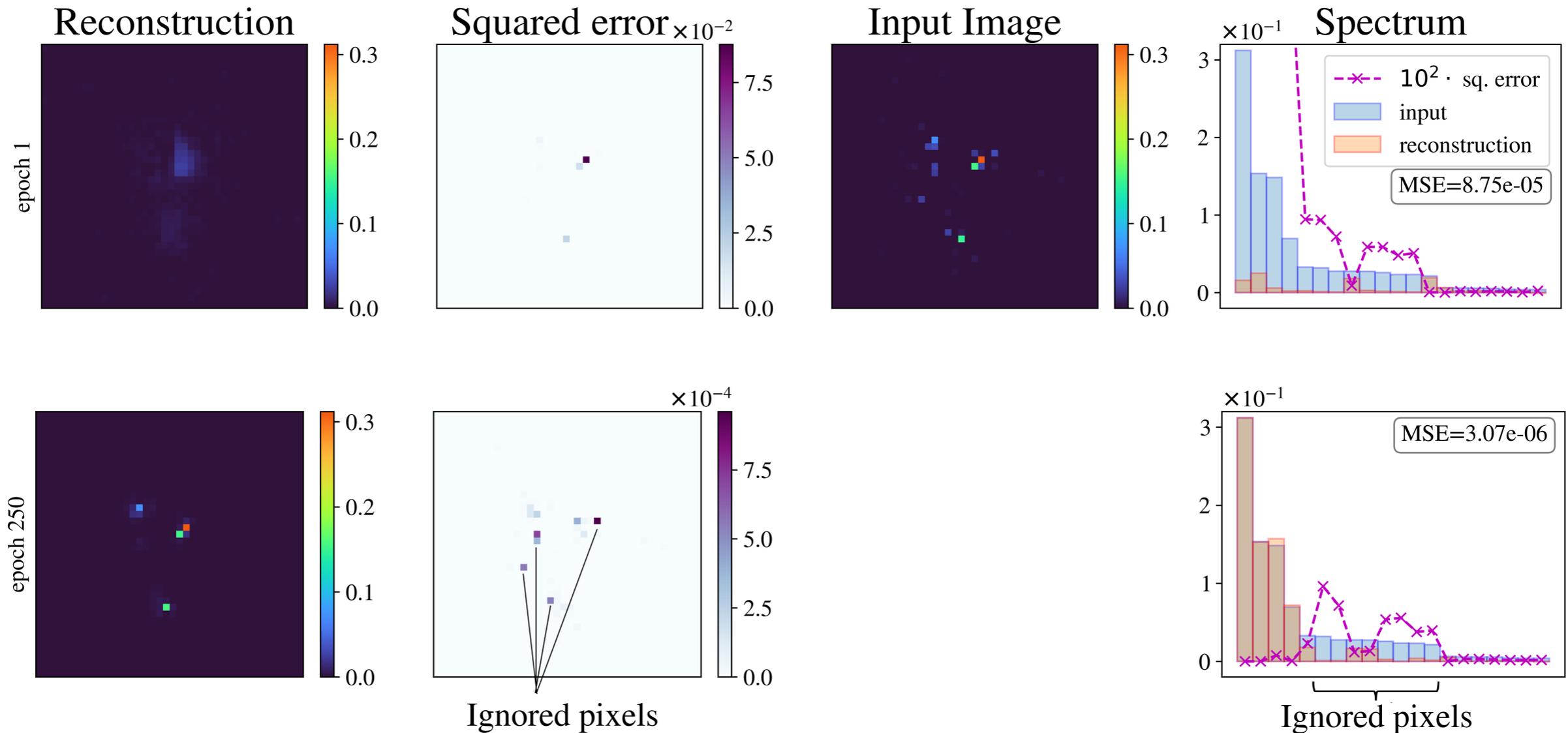
Training on top-jets only, the AE still identifies top-jets as anomalous

What does the autoencoder learn?



Vanilla autoencoder shows a very limited reconstruction capability.

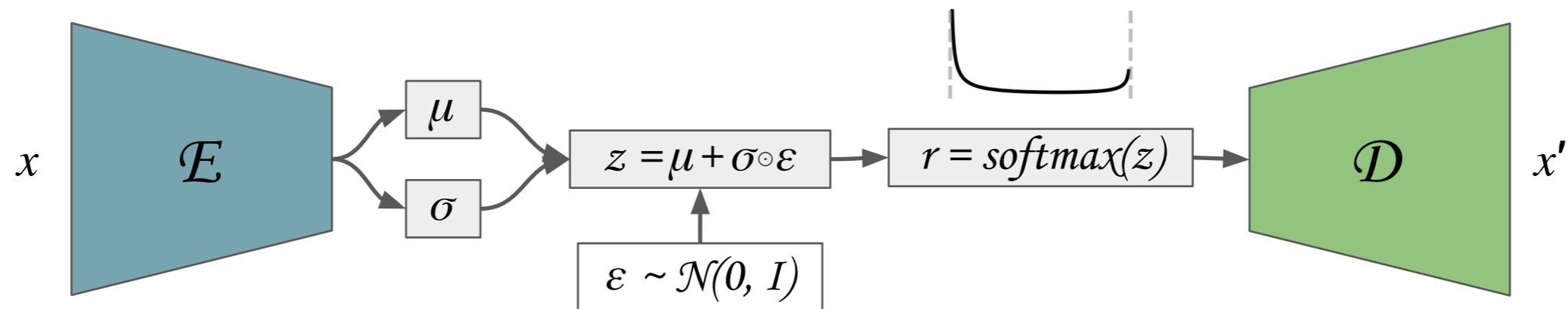
What does the autoencoder learn?



Vanilla autoencoder shows a complexity bias, it tends to better reconstruct "simpler" images.

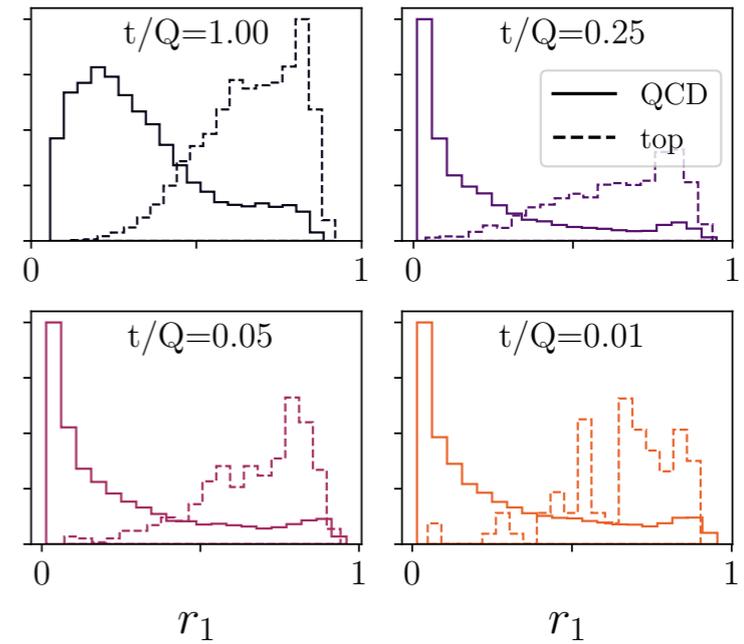
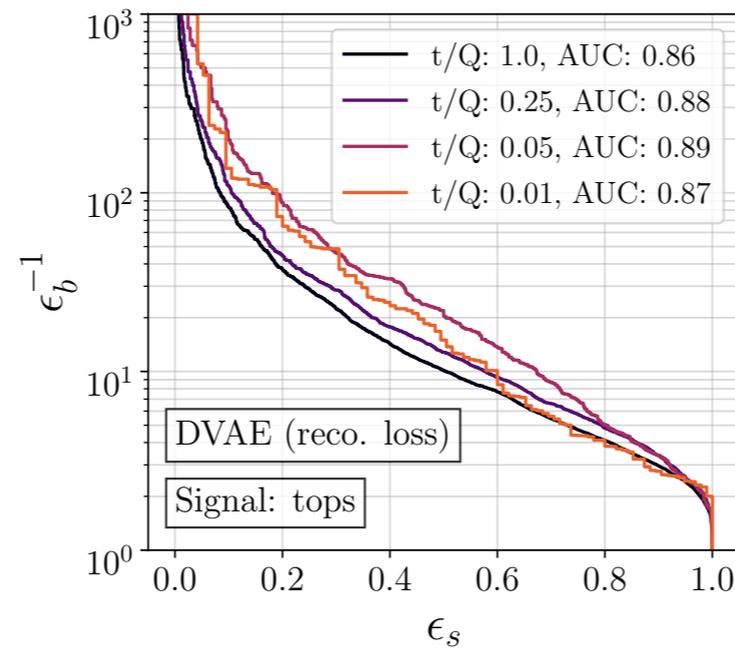
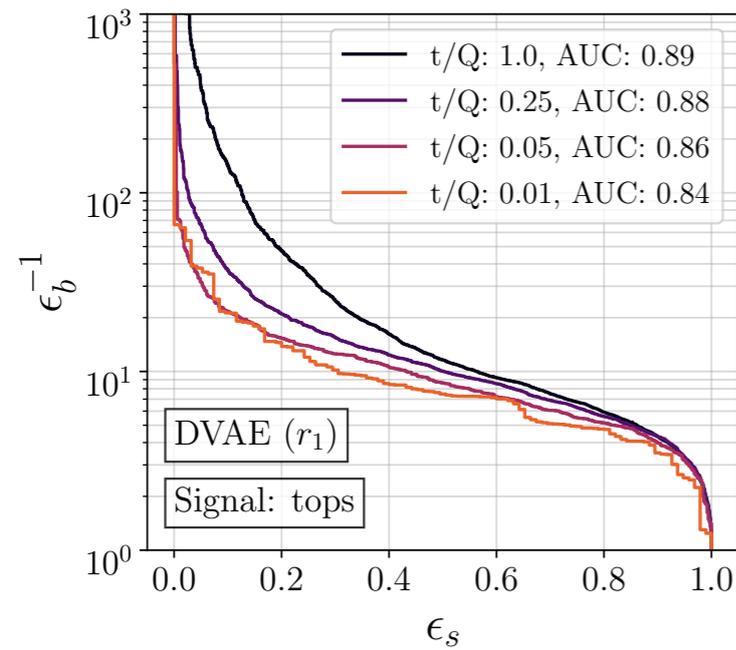
Anomaly detection with autoencoders: outlook

- **Regularise the latent space?** [Cerri et al., JHEP 05 (2019) 036, Cheng et al., e-Print: 2007.01850 [hep-ph], Dillon et al., SciPost Phys. 11 (2021) 061, ...]



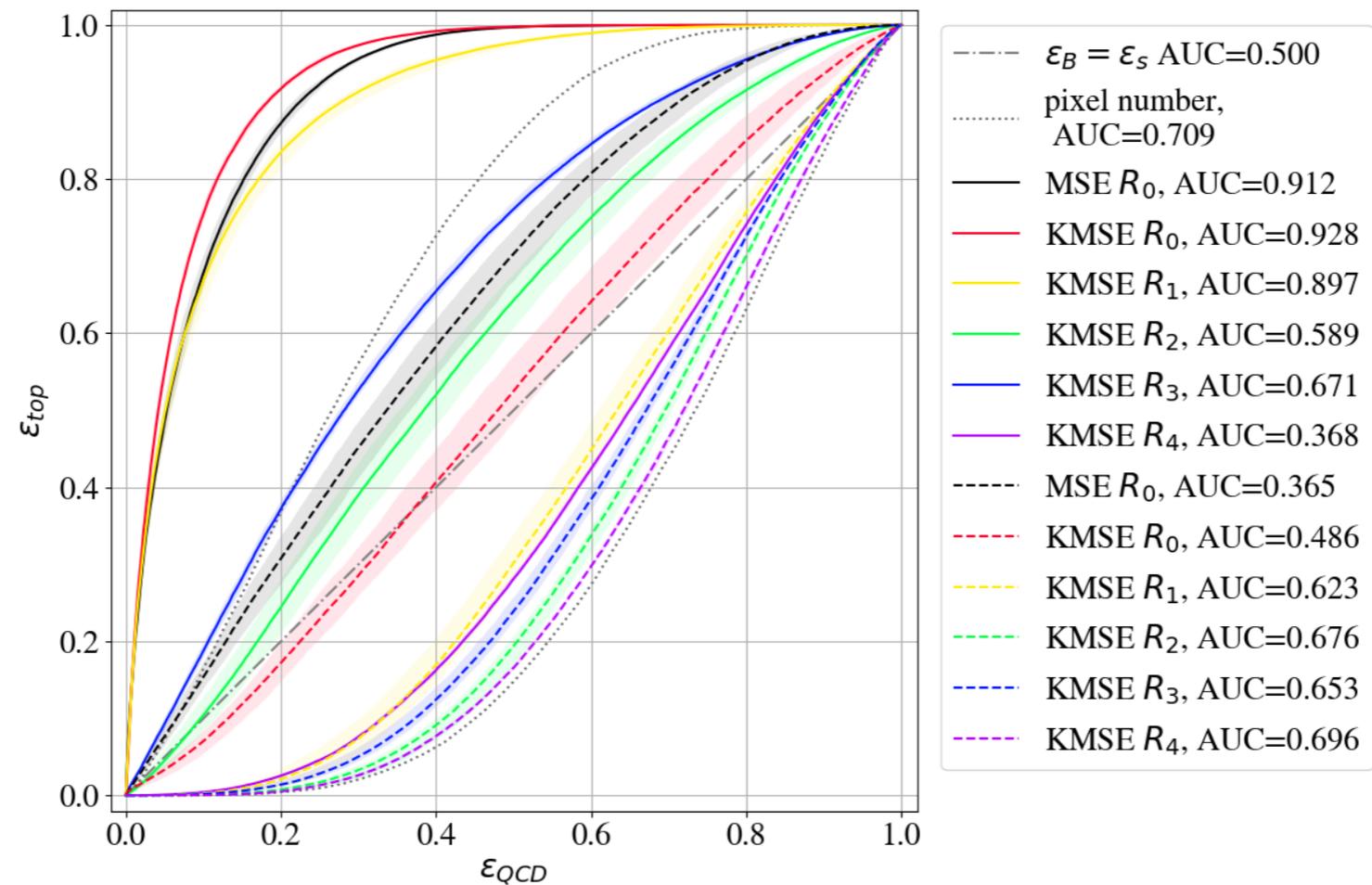
Anomaly detection with autoencoders: outlook

- **Regularise the latent space?** [Cerri et al., JHEP 05 (2019) 036, Cheng et al., e-Print: 2007.01850 [hep-ph], Dillon et al., SciPost Phys. 11 (2021) 061, ...]



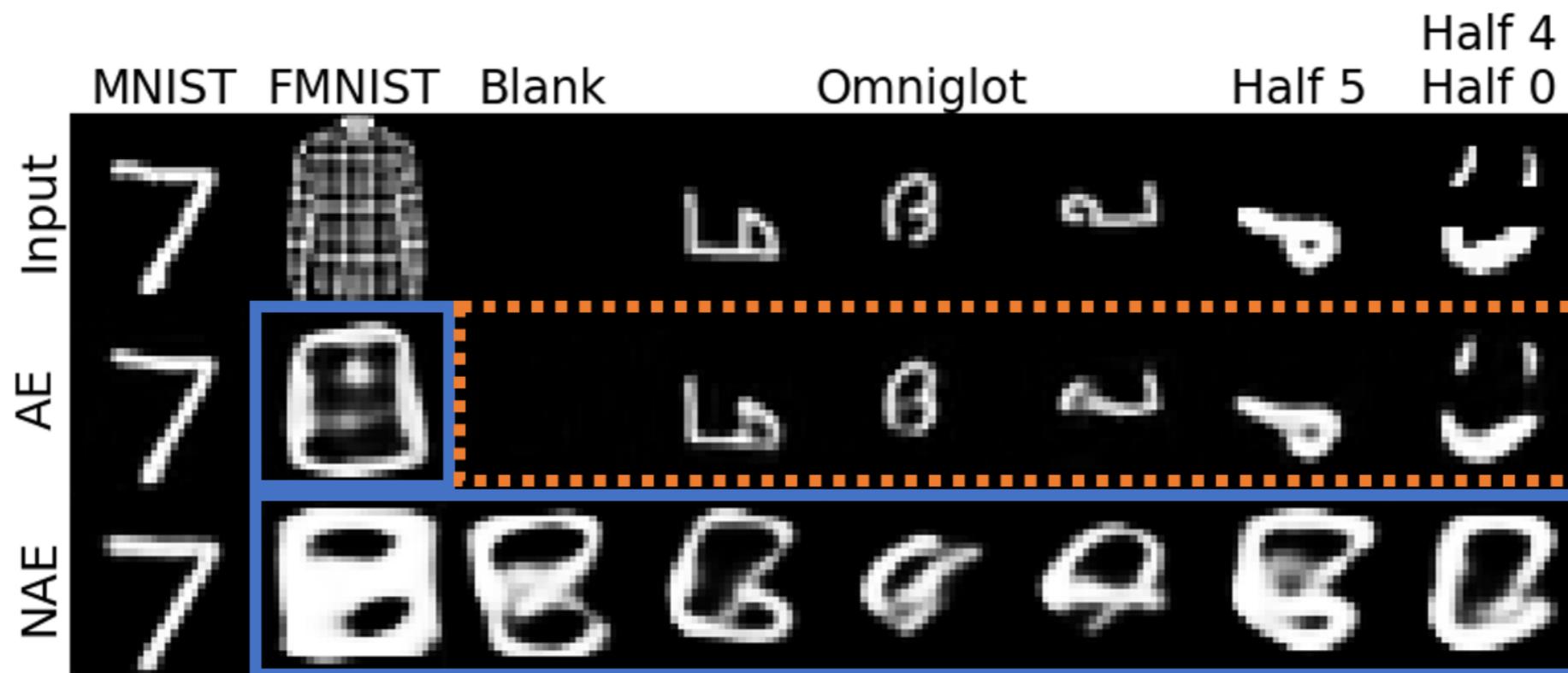
Anomaly detection with autoencoders: outlook

- Improve the performance of the AE through preprocessing: smearing and re-weighting of pixels? [Finke et al., JHEP 06 (2021) 161, Buss et al., e-Print: 2202.00686 [hep-ph], ...]



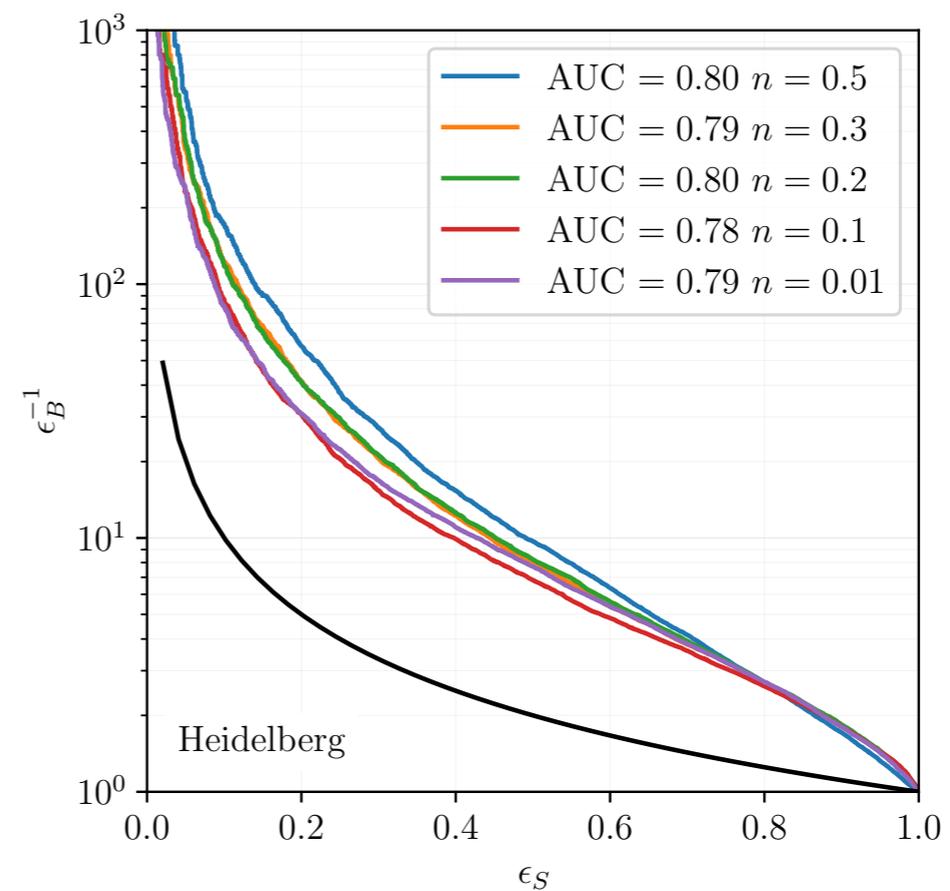
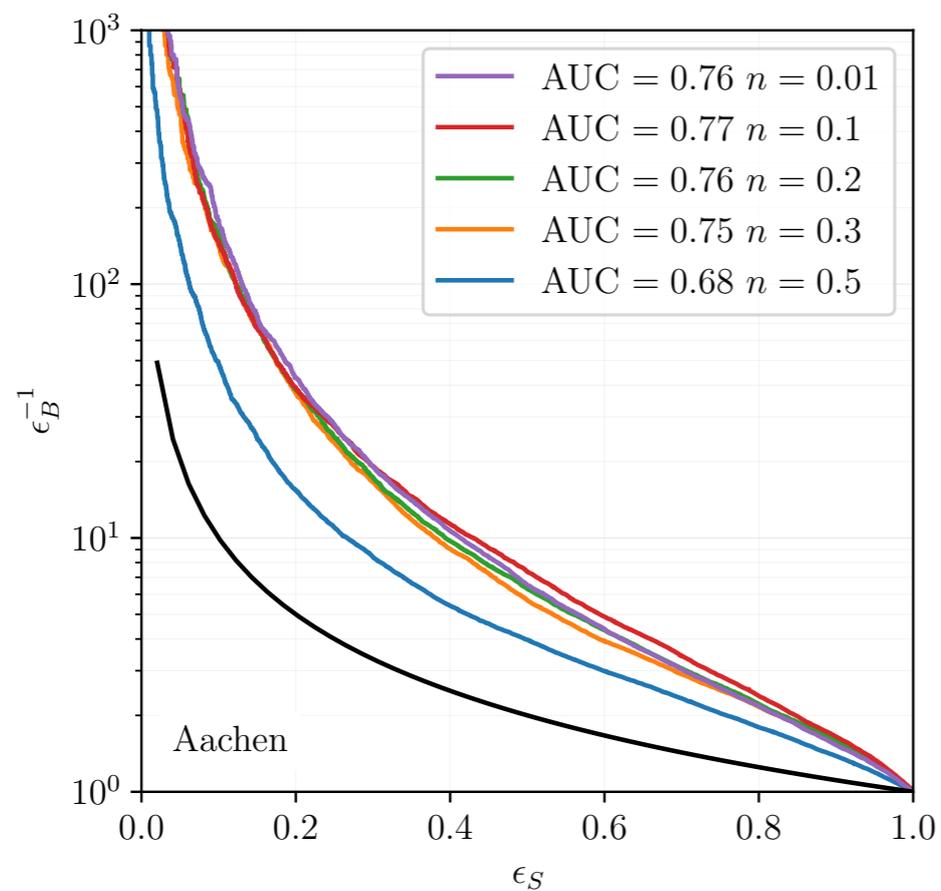
Anomaly detection with autoencoders: outlook

- **Introduce normalising condition to prevent outlier reconstruction?** [Yoon et al., arXiv:2105.05735 [cs.LG], Dillon et al., arXiv:2206.14225 [hep-ph], ...]



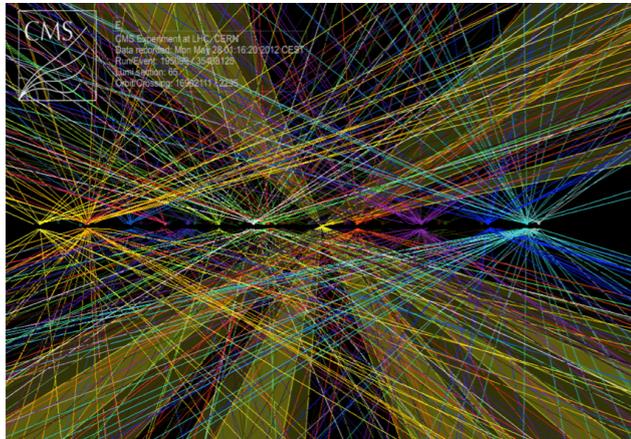
Anomaly detection with autoencoders: outlook

- Introduce normalising condition to prevent outlier reconstruction? [Yoon et al., arXiv:2105.05735 [cs.LG], Dillon et al., arXiv:2206.14225 [hep-ph], ...]



Model independence: searching for anomalies at the LHC

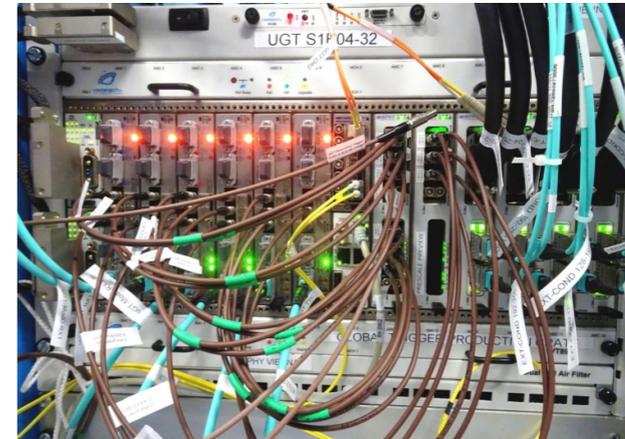
Collisions



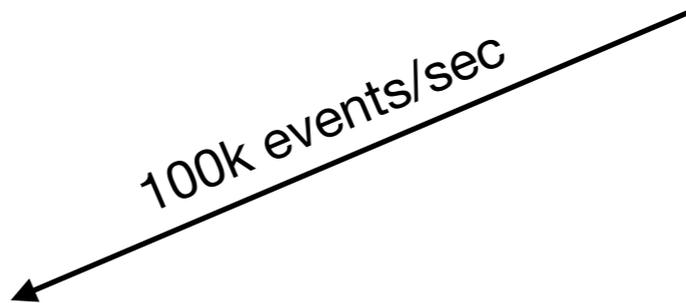
40M events/sec



Level 1 trigger



100k events/sec



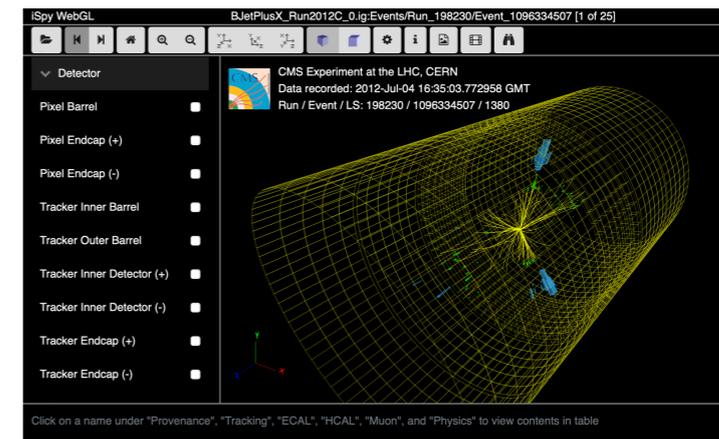
High-level trigger



1k events/sec

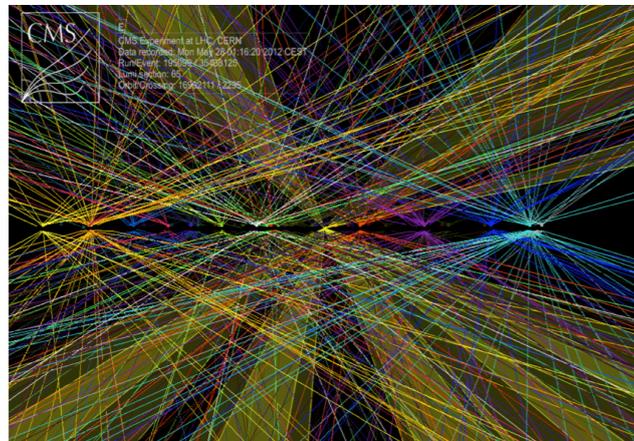


Data analysis

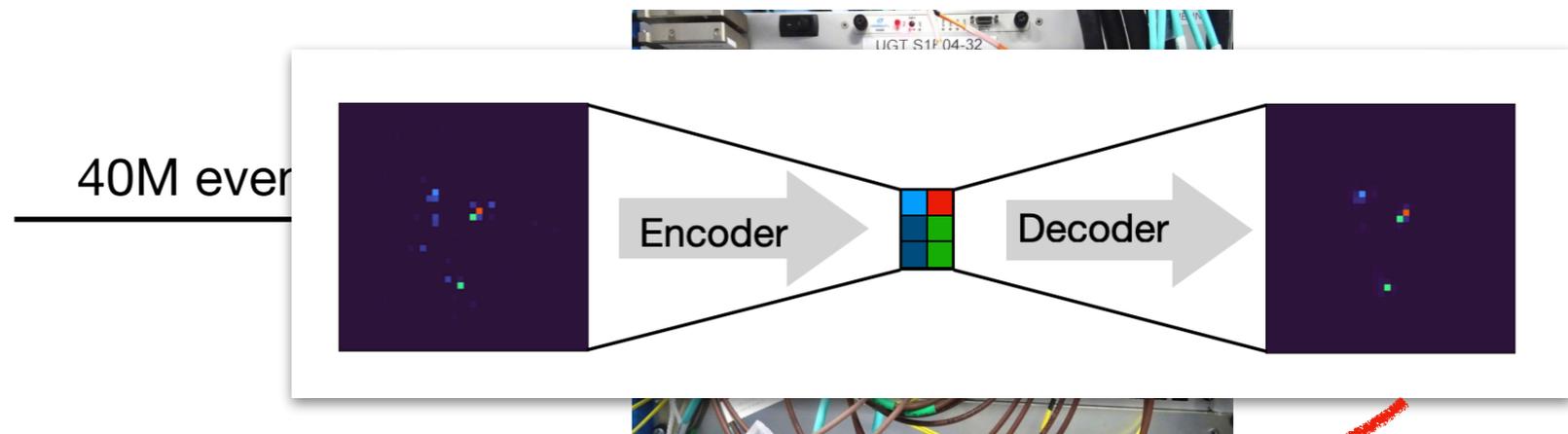


Model independence: searching for anomalies at the LHC

Collisions



Level 1 trigger



40M events/sec

High-level trigger

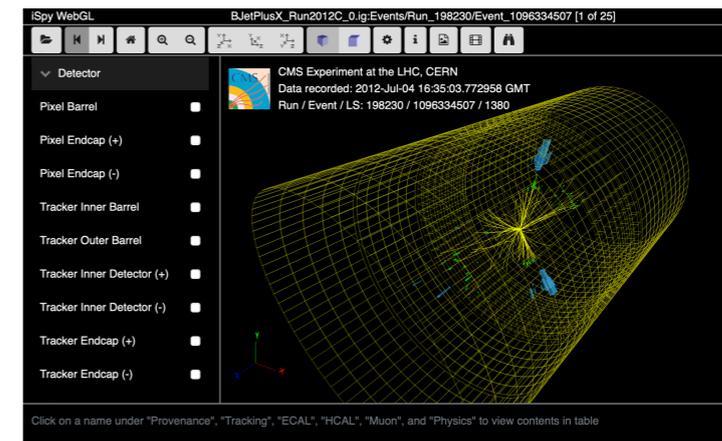


100k events/sec

New physics?

Data analysis

1k events/sec



Anomaly detection data challenge: <https://mpp-hep.github.io/ADC2021/>

Anomaly detection with autoencoders: outlook

- **Regularise the latent space?** [Cerri et al., JHEP 05 (2019) 036, Cheng et al., e-Print: 2007.01850 [hep-ph], Dillon et al., SciPost Phys. 11 (2021) 061, ...]
- **Improve the performance of the AE through preprocessing: smearing and re-weighting of pixels?** [Finke et al., JHEP 06 (2021) 161, Buss et al., e-Print: 2202.00686 [hep-ph], ...]
- **Introduce normalising condition to prevent outlier reconstruction?** [Yoon et al., arXiv:2105.05735 [cs.LG], Dillon et al., arXiv:2206.14225 [hep-ph], ...]
- **Can we construct a fast autoencoder to detect anomalies in real time?**

Thank you