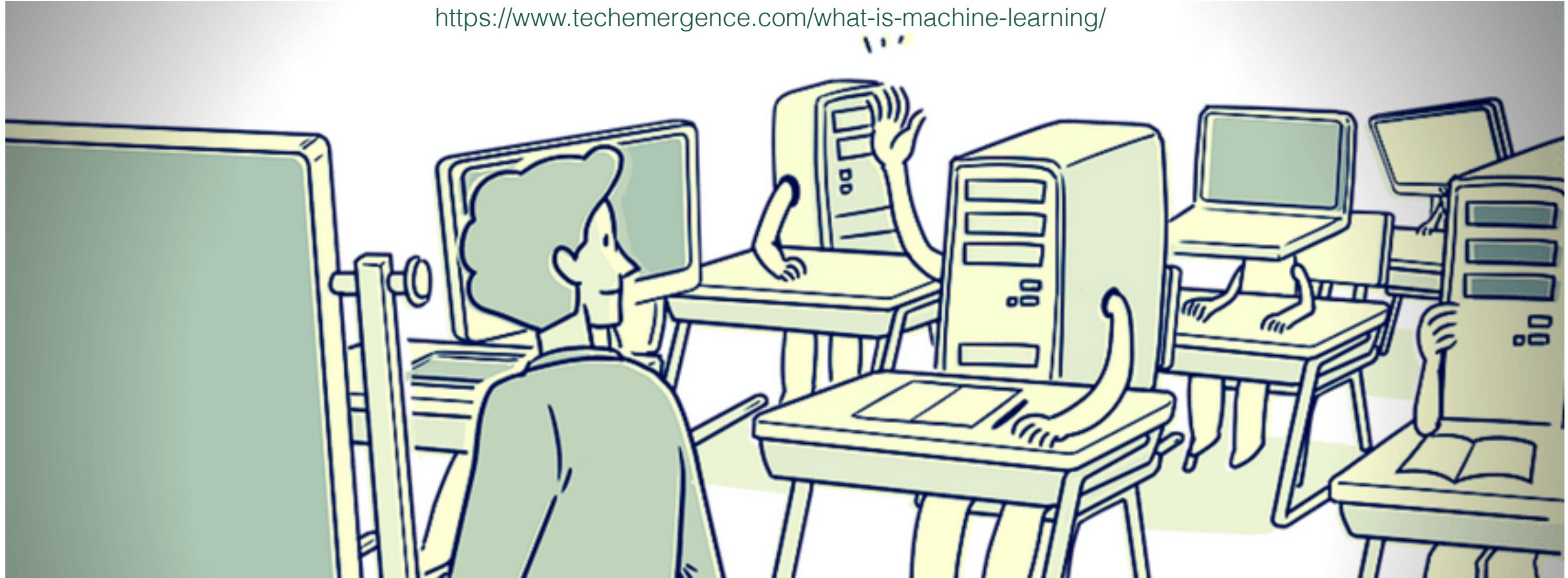


Introduction to Machine Learning

Lecture 3

<https://www.techemergence.com/what-is-machine-learning/>

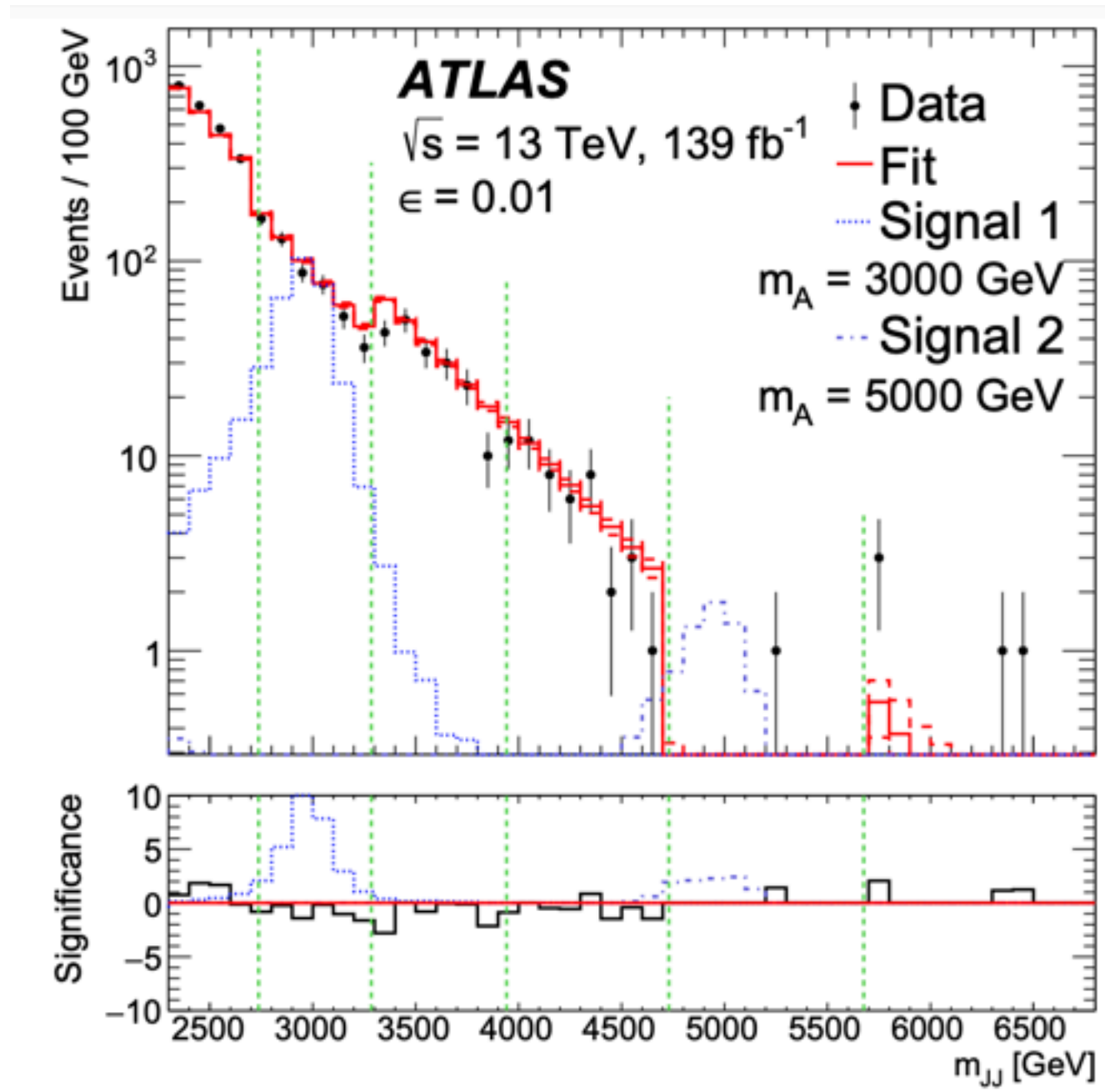


QUC Winter School 2021
KIAS
December 2021

Questions from last lecture?

ATLAS search using weak supervision

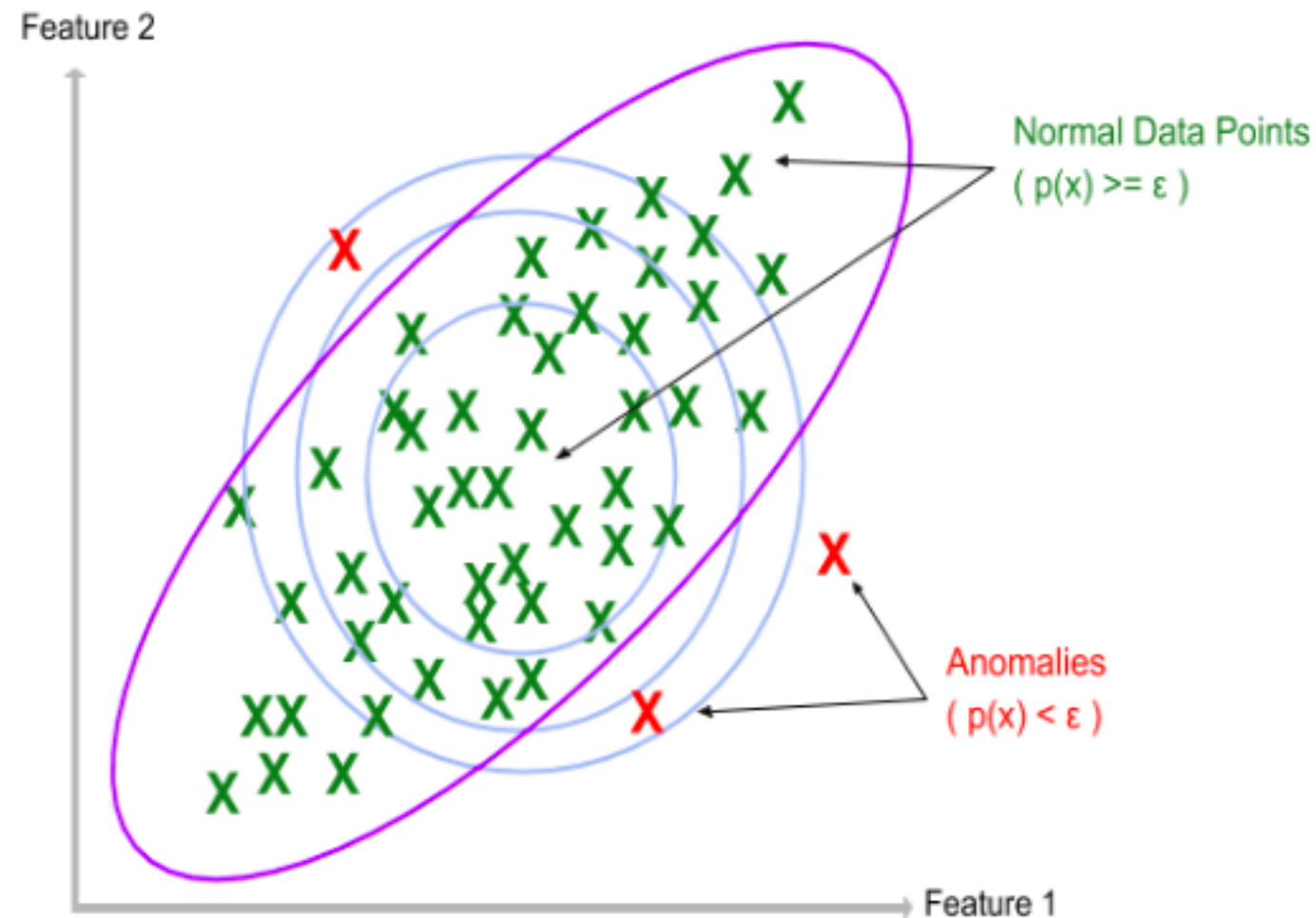
[2005.02983]



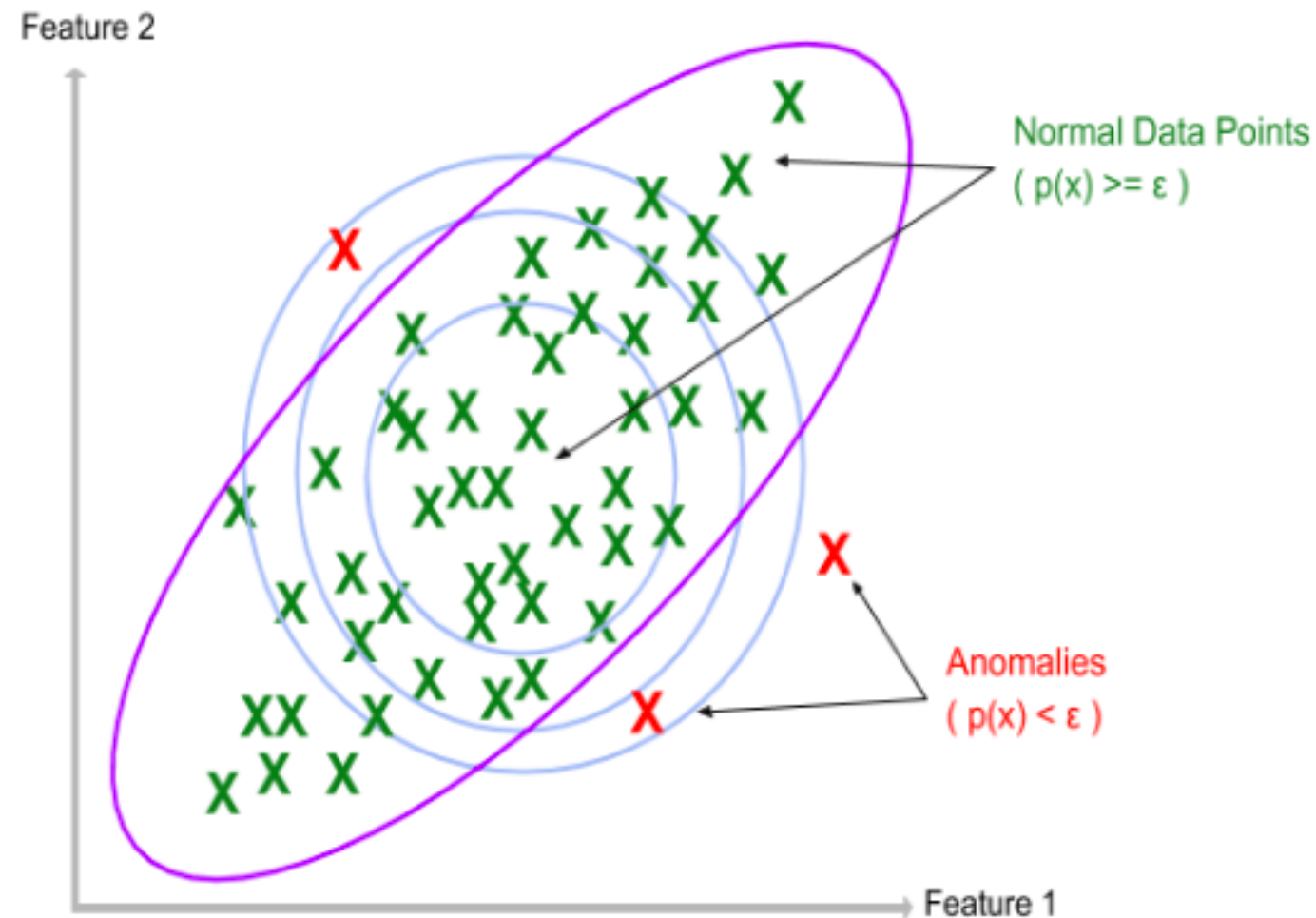
Start of lecture 3

- Anomaly Detection
- Generative Networks
- Flows for sampling/inference

What is Anomaly Detection?

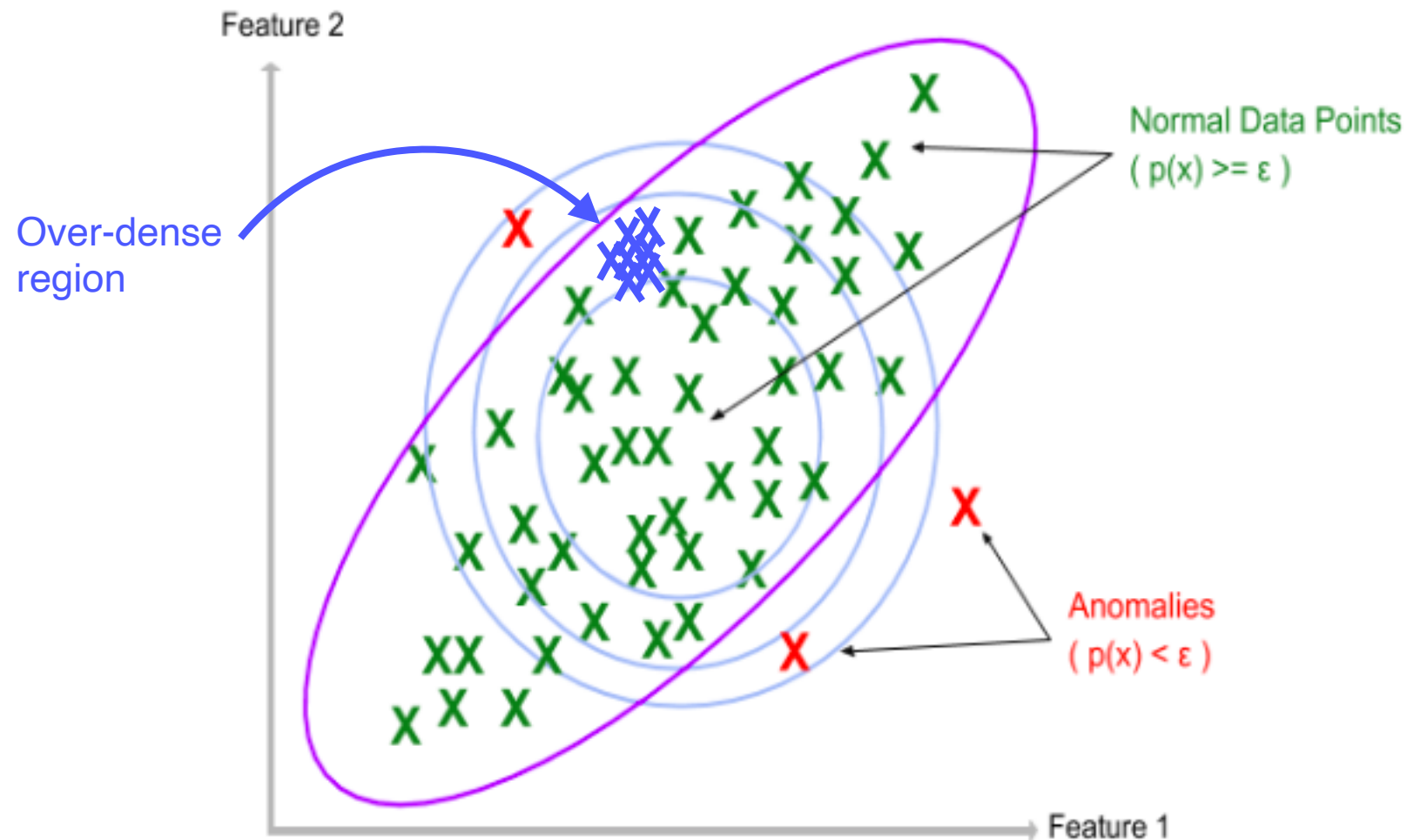


What is Anomaly Detection?



1) Events that occur with low probability (look different in some feature space)

What is Anomaly Detection?



- 1) Events that occur with low probability (look different in some feature space)
- 2) Events that don't look very different, but occur at a higher rate than expected

Why Anomaly Detection?

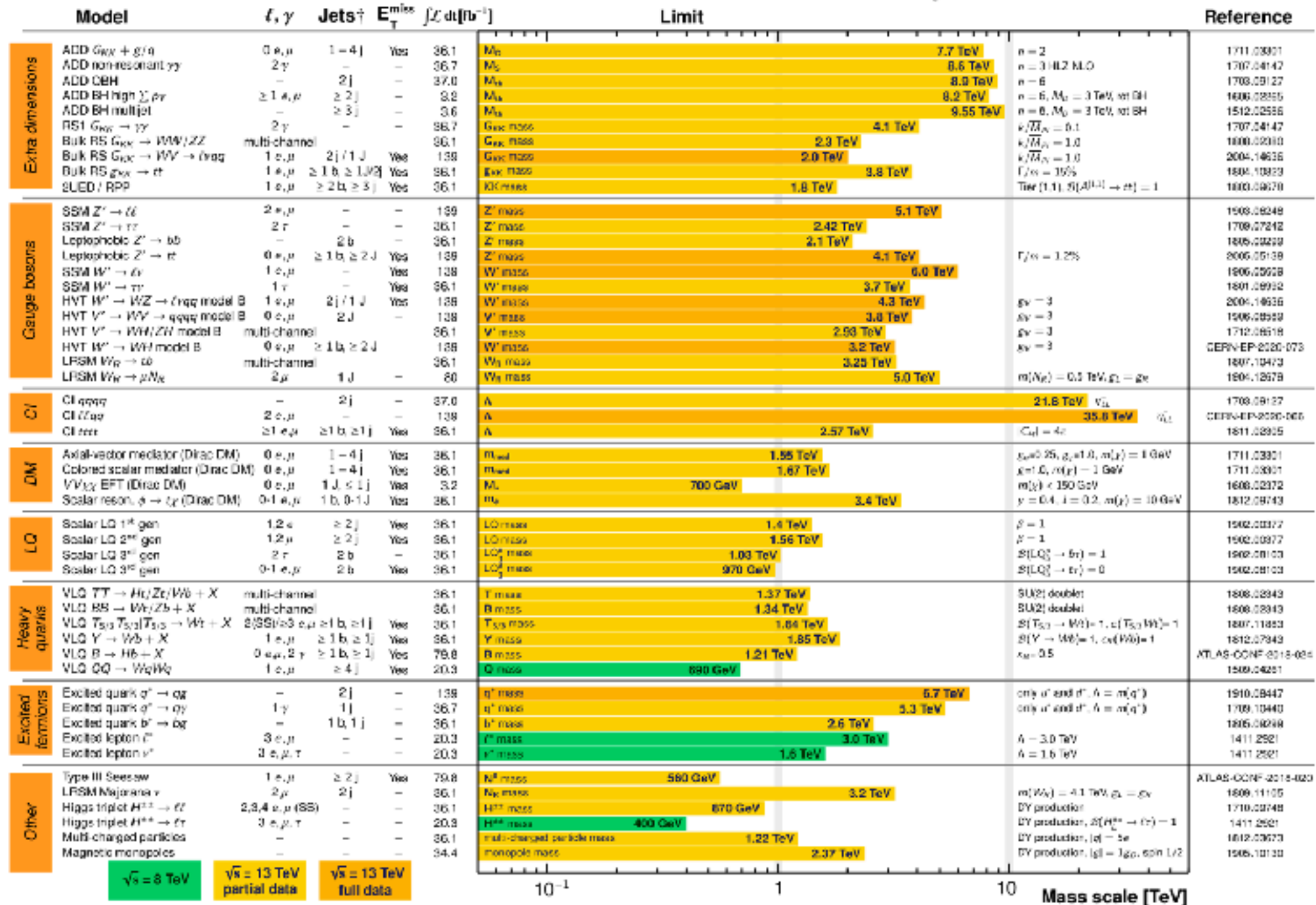
ATLAS Exotics Searches* - 95% CL Upper Exclusion Limits

Status: May 2020

ATLAS Preliminary

$\int \mathcal{L} dt = (3.2 - 139) \text{ fb}^{-1}$

$\sqrt{s} = 8, 13 \text{ TeV}$



*Only a selection of the available mass limits on new states or phenomena is shown.

† Small-radius (large-radius) jets are denoted by the letter j (J).

Why Anomaly Detection?

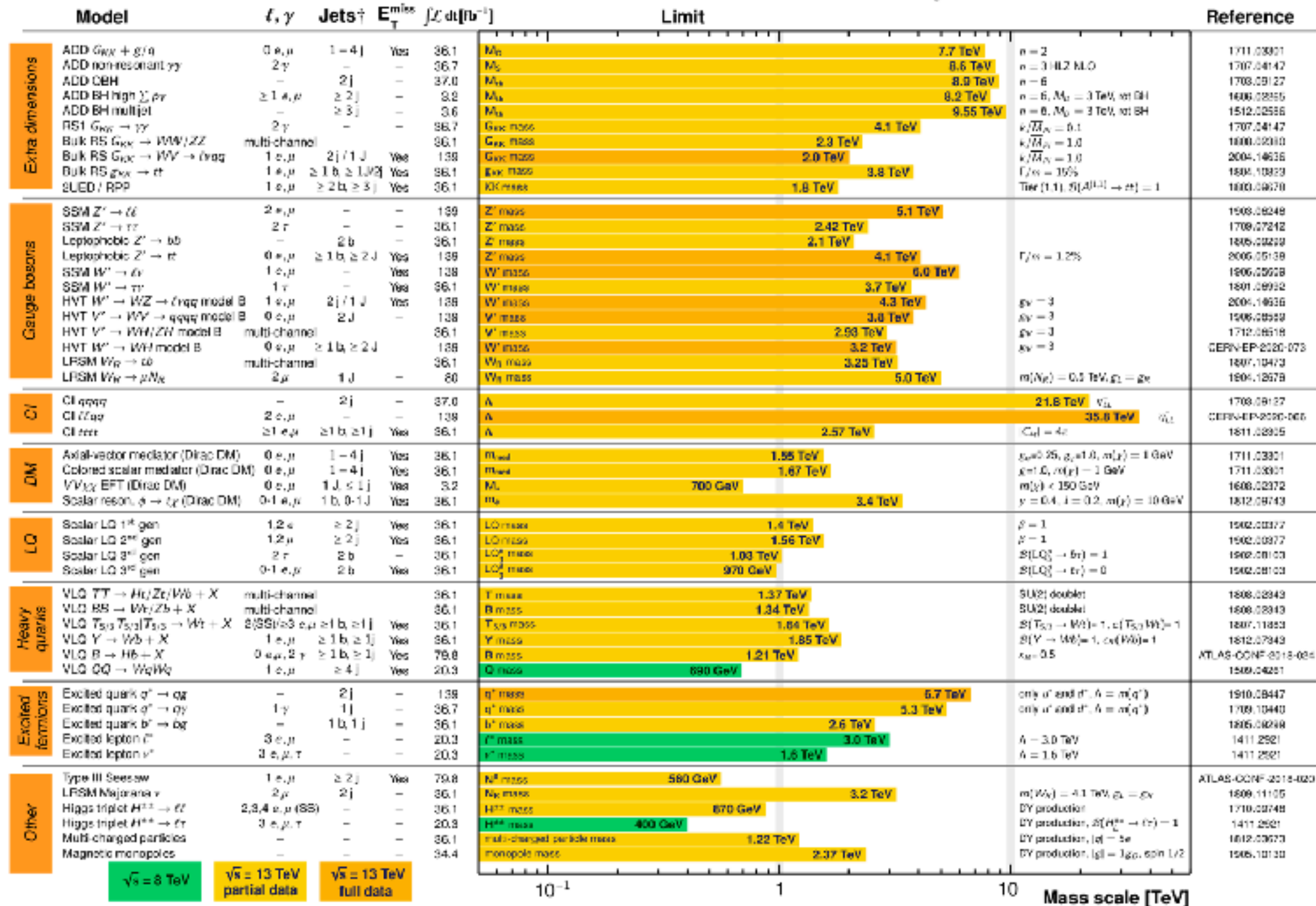
ATLAS Exotics Searches* - 95% CL Upper Exclusion Limits

Status: May 2020

ATLAS Preliminary

$\int \mathcal{L} dt = (3.2 - 139) \text{ fb}^{-1}$

$\sqrt{s} = 8, 13 \text{ TeV}$



*Only a selection of the available mass limits on new states or phenomena is shown.

† Small-radius (large-radius) jets are denoted by the letter j (J).

Higher rates, but hard to distinguish from SM

Look more different, but more low rates

Why Anomaly Detection?

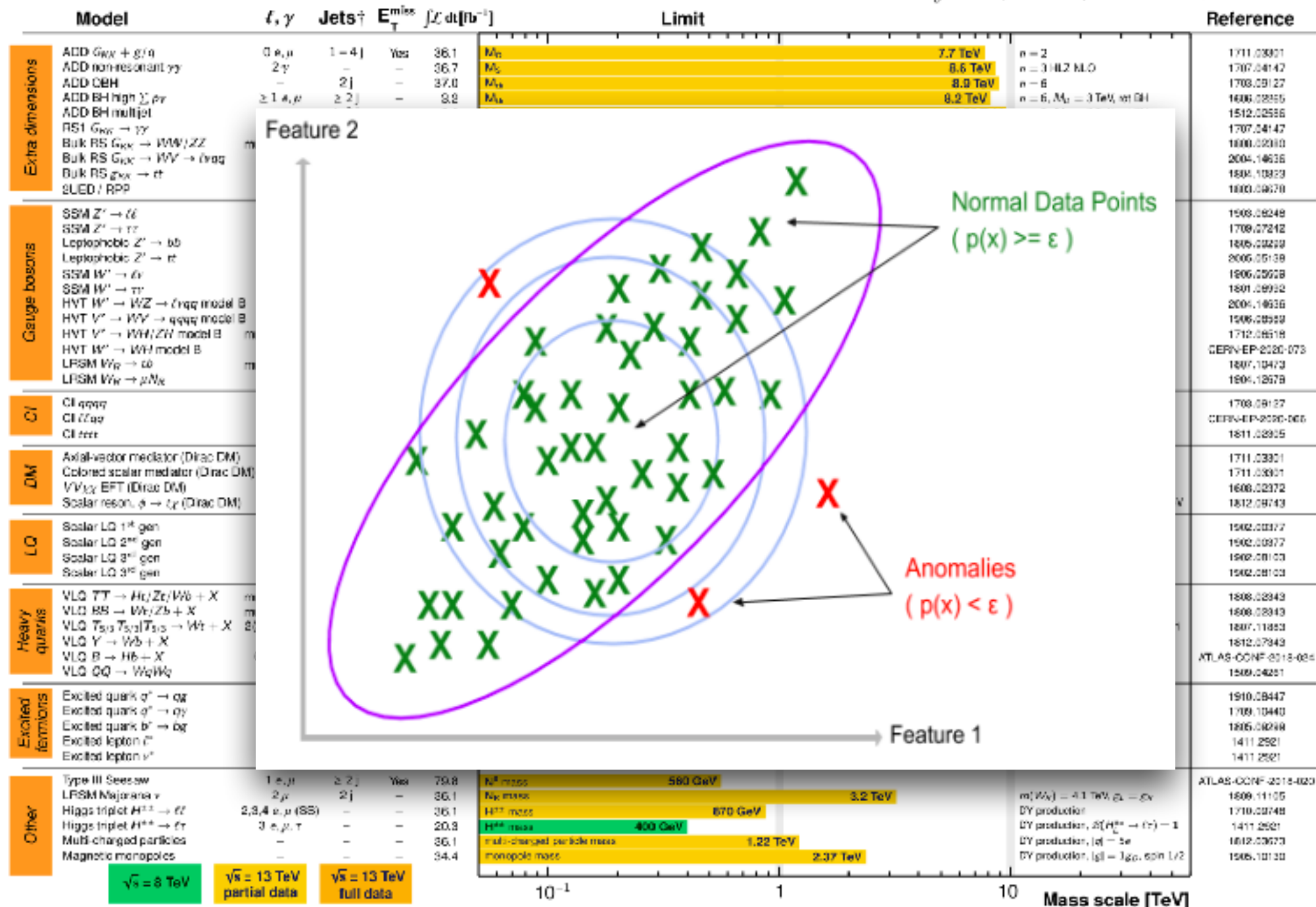
ATLAS Exotics Searches* - 95% CL Upper Exclusion Limits

Status: May 2020

ATLAS Preliminary

$$\int \mathcal{L} dt = (3.2 - 139) \text{ fb}^{-1}$$

$$\sqrt{s} = 8, 13 \text{ TeV}$$



*Only a selection of the available mass limits on new states or phenomena is shown.

† Small-radius (large-radius) jets are denoted by the letter j (J).

Higher rates, but hard to distinguish from SM

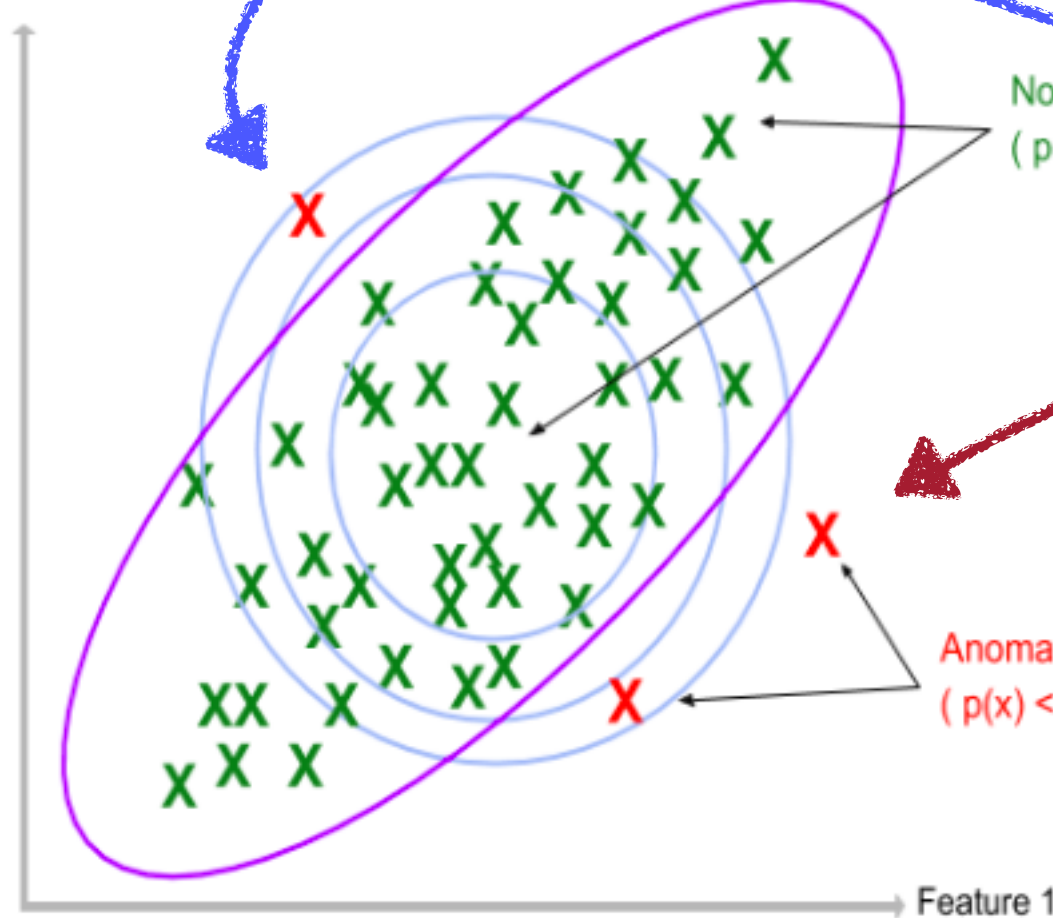
Look more different, but more low rates

Why Anomaly Detection?

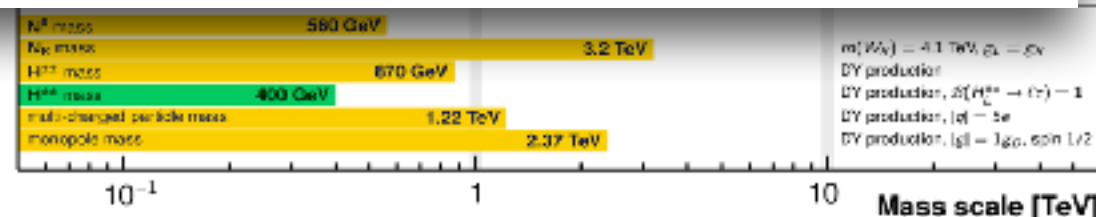
ATLAS Exotics Searches* - 95% CL Upper Excl
Status: May 2020

Model	ℓ, γ	Jets [†]	E_T^{miss}	$\int \mathcal{L} d\sqrt{s}$ [fb ⁻¹]
Extra dimensions				
ADD $G_{KK} + g/\eta$	$0, \mu$	$1-4$	Yes	36.1
ADD nonresonant $\gamma\gamma$	2γ	—	—	36.7
ADD CBH	—	2	—	37.0
ADD BH high $\sum p_T$	$\geq 1, \mu$	≥ 2	—	3.8
ADD BH multijet	—	—	—	—
RS1 $G_{KK} \rightarrow \gamma\gamma$	—	—	—	—
Bulk RS $G_{KK} \rightarrow WW/ZZ$	—	—	—	—
Bulk RS $G_{KK} \rightarrow W\gamma \rightarrow \ell\nu q$	—	—	—	—
Bulk RS $G_{KK} \rightarrow t\bar{t}$	—	—	—	—
2UED / RPP	—	—	—	—
Gauge bosons				
SSM $Z' \rightarrow \ell\ell$	—	—	—	—
SSM $Z' \rightarrow \tau\tau$	—	—	—	—
Leptophobic $Z' \rightarrow b\bar{b}$	—	—	—	—
Leptophobic $Z' \rightarrow t\bar{t}$	—	—	—	—
SSM $W' \rightarrow \ell\nu$	—	—	—	—
SSM $W' \rightarrow \tau\tau$	—	—	—	—
HVT $W' \rightarrow WZ \rightarrow \ell\nu q q$ model B	—	—	—	—
HVT $V' \rightarrow W\gamma \rightarrow q\bar{q} q q$ model B	—	—	—	—
HVT $V' \rightarrow WW/ZH$ model B	—	—	—	—
HVT $W' \rightarrow WH$ model B	—	—	—	—
LFSM $W' \rightarrow b\bar{b}$	—	—	—	—
LFSM $W' \rightarrow \mu N_k$	—	—	—	—
CI				
CI $q\bar{q}q$	—	—	—	—
CI $\ell\bar{\ell}q$	—	—	—	—
CI $t\bar{t}t$	—	—	—	—
DM				
Axial-vector mediator (Dirac DM)	—	—	—	—
Colored scalar mediator (Dirac DM)	—	—	—	—
VV_{EFT} (Dirac DM)	—	—	—	—
Scalar reson. $\phi \rightarrow \ell\ell$ (Dirac DM)	—	—	—	—
LQ				
Scalar LQ 1 st gen	—	—	—	—
Scalar LQ 2 nd gen	—	—	—	—
Scalar LQ 3 rd gen	—	—	—	—
Scalar LQ 3 rd gen	—	—	—	—
Heavy quarks				
VLO $TT \rightarrow H/Z/\gamma + X$	—	—	—	—
VLO $BB \rightarrow W/Z/\gamma + X$	—	—	—	—
VLO $T_{3/2} T_{3/2} \rightarrow Wt + X$	—	—	—	—
VLO $Y \rightarrow Wb + X$	—	—	—	—
VLO $B \rightarrow Hb + X$	—	—	—	—
VLO $QQ \rightarrow WqWq$	—	—	—	—
Excited fermions				
Excited quark $q^* \rightarrow q\gamma$	—	—	—	—
Excited quark $q^* \rightarrow q\gamma$	—	—	—	—
Excited quark $b^* \rightarrow b\gamma$	—	—	—	—
Excited lepton $\ell^* \rightarrow \ell\gamma$	—	—	—	—
Excited lepton ν^*	—	—	—	—
Other				
Type III Seesaw	$1, \mu$	≥ 2	Yes	79.8
LFSM Majorana ν	$2, \mu$	2	—	36.1
Higgs triplet $H^{\pm\pm} \rightarrow \ell\ell$	$2, 3, 4, \mu, \tau$ (SS)	—	—	36.1
Higgs triplet $H^{\pm\pm} \rightarrow \ell\tau$	—	—	—	20.9
Multi-charged particles	—	—	—	36.1
Magnetic monopoles	—	—	—	34.4

Feature 2



What if we have been looking for new physics here, rather than here



*Only a selection of the available mass limits on new states or phenomena is shown.

† Small-radius (large-radius) jets are denoted by the letter j (J).

Higher rates, but hard to distinguish from SM

Look more different, but more low rates

Why Anomaly Detection?

- 1) Model Agnostic: We are good at looking for models we know of, but what if we don't know what we should be looking for?
- 2) Simulation Independent: With no signal model, it is possible to use methods directly on data from the LHC without Monte Carlo simulations.

How to detect anomalies?

Over 50 papers in HEP anomaly detection

<https://iml-wg.github.io/HEPML-LivingReview/>

LHC Olympics [[2101.08320](#)] focuses on finding over densities in all-hadronic events

- Black Box 1: Similar to example data: 4 methods found resonance
- Black Box 2: SM only. 4 methods claimed a resonance, 1 claimed lack-of-resonance
- Black Box 3: Correct resonance not detected by any group

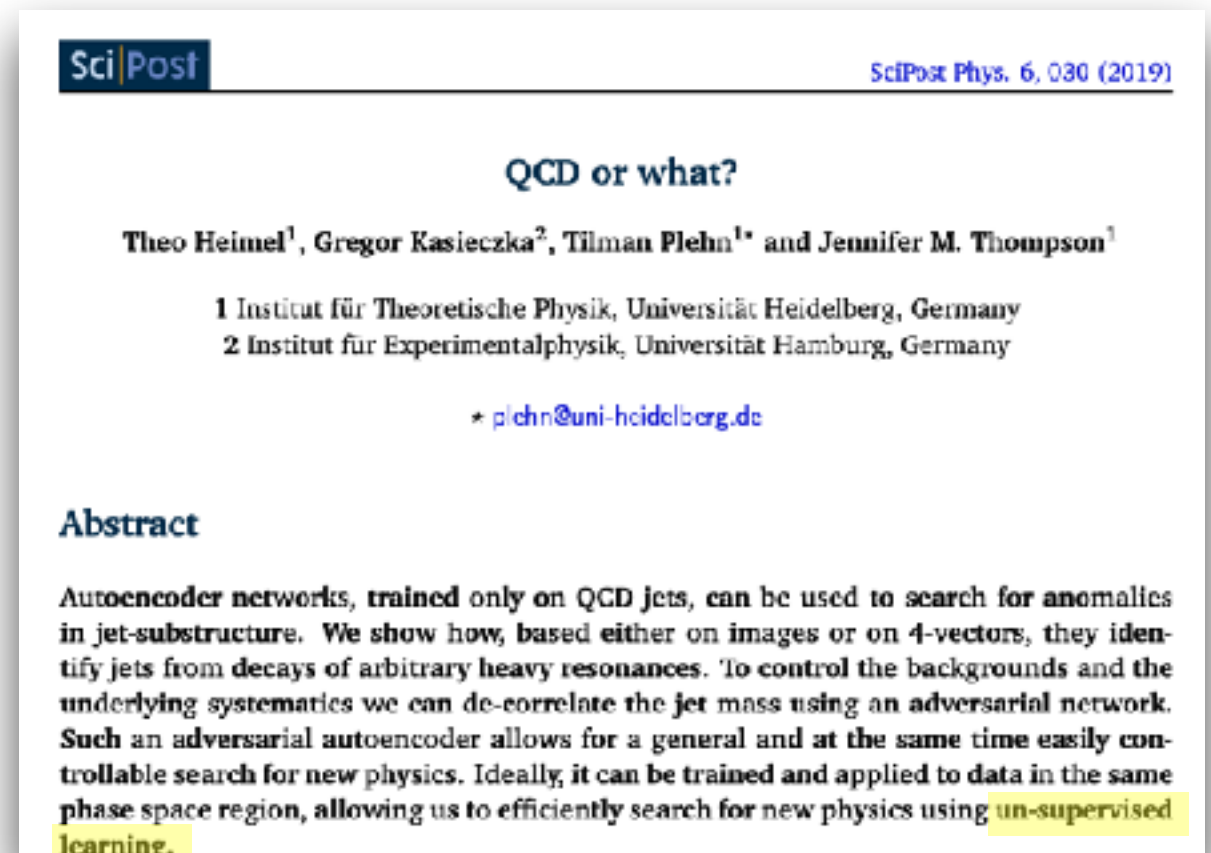
Dark Machines Challenge [[2105.14027](#)] focuses on finding individual events which look different

- Train on SM-only events, apply to many different signals
- Find methods which work best for most signals
- No method finds every new physics signal

How to detect anomalies?



[[1808.08992](#)]



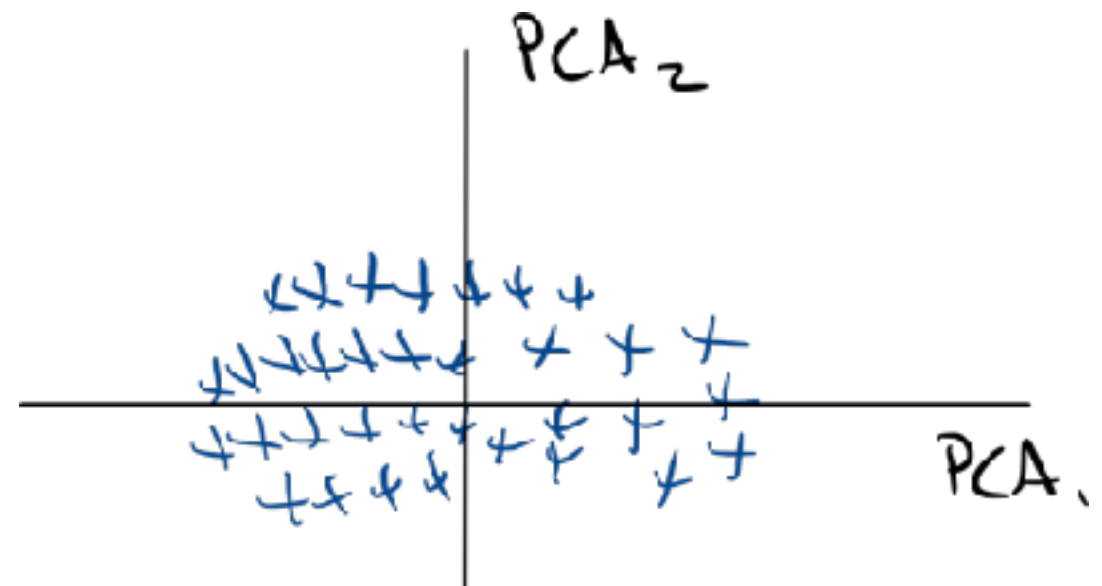
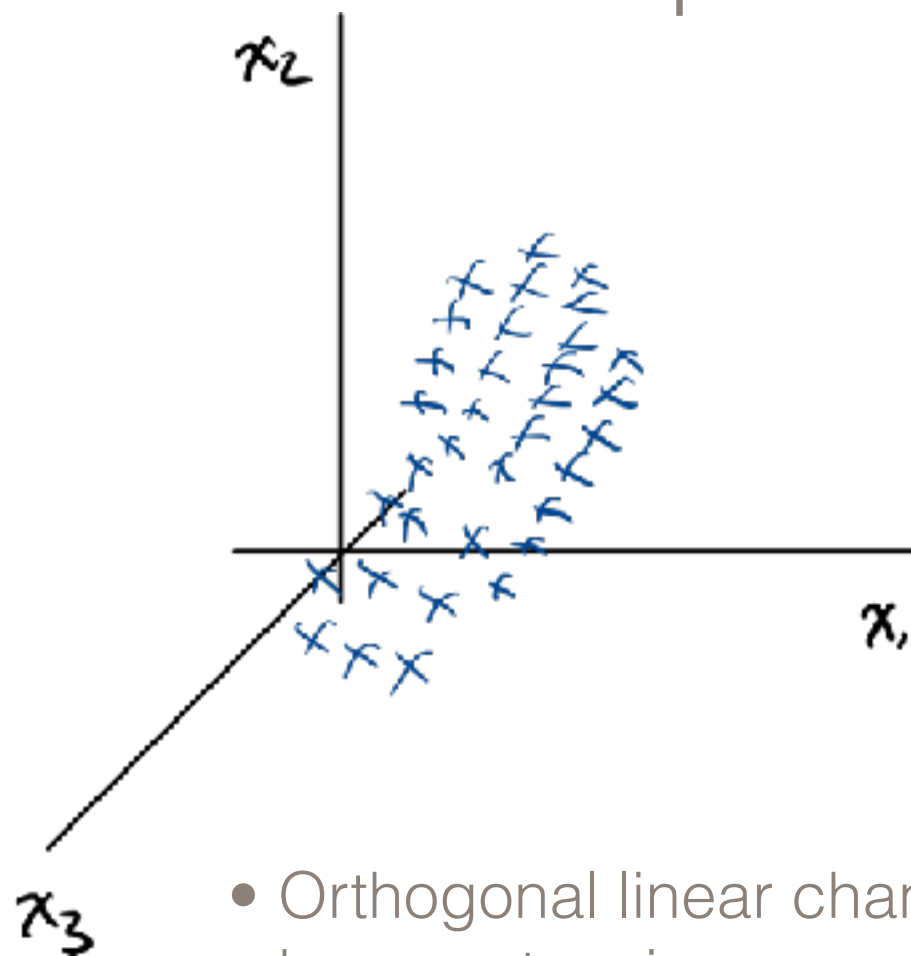
[[1808.08979](#)]

1. Use Autoecoders for anomaly detection
2. Emphasize training directly on data, how much signal can be in training and still work?

Introduction to Autoencoders

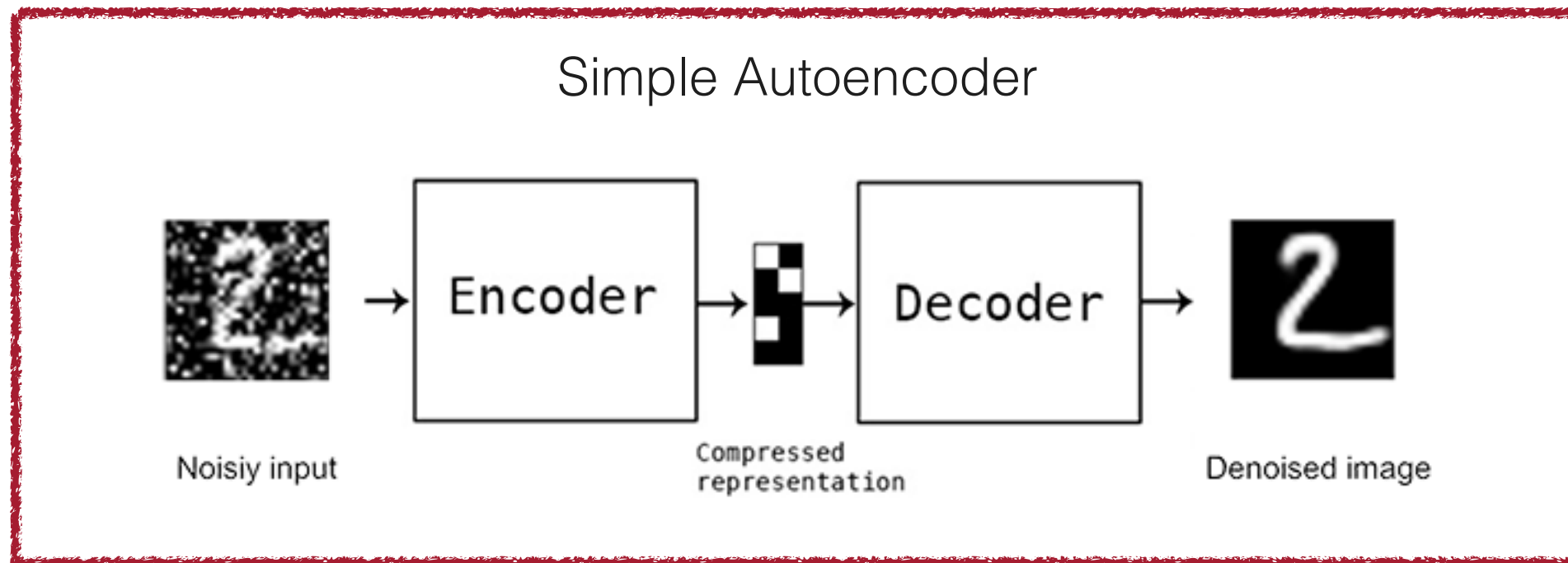
- Method for dimensionality reduction
- Focuses on important pieces of information and ignores noise

Principal Component Analysis



- Orthogonal linear change of variables such that first axis has most variance, second has second most, etc
- Which variables describe most of the variation of the data?
- Invertible (lossy if not using full basis)

Introduction to Autoencoders



- Encoding and Decoding can be non-linear
- Encoder learns what is important in the data and what is not
- Size of compressed representation chosen before training
- Compressed representation changed each training of the networks

Introduction to Autoencoders

Introduction to Autoencoders

- Unlike PCA, autoencoders need to be trained
- Need input data, output predictions, and a target
 - Target data is the same as input data

Introduction to Autoencoders


- Unlike PCA, autoencoders need to be trained
- Need input data, output predictions, and a target
 - Target data is the same as input data
- How to compare output predictions and target?
 - Use Mean Squared Error

$$d_{\text{ME},2} = \text{MSE} \equiv \frac{1}{n} \sum_{i=1}^n \left| D(E(x))_i - y_i \right|^2$$

Introduction to Autoencoders

- Unlike PCA, autoencoders need to be trained
- Need input data, output predictions, and a target
 - Target data is the same as input data
- How to compare output predictions and target?
 - Use Mean Squared Error

$$d_{\text{ME},2} = \text{MSE} \equiv \frac{1}{n} \sum_{i=1}^n \left| D(E(x))_i - y_i \right|^2$$



$y = \text{target } (=x)$
(or the input in this case)

Introduction to Autoencoders

- Unlike PCA, autoencoders need to be trained
- Need input data, output predictions, and a target
 - Target data is the same as input data
- How to compare output predictions and target?
 - Use Mean Squared Error

$$d_{\text{ME},2} = \text{MSE} \equiv \frac{1}{n} \sum_{i=1}^n \left| D(E(x))_i - y_i \right|^2$$

$D(E(x))$ = output of decoded, encoded data
use $f(x)$ for rest of talk



y = target (=x)
(or the input in this case)



Introduction to Autoencoders

- Unlike PCA, autoencoders need to be trained
- Need input data, output predictions, and a target
 - Target data is the same as input data
- How to compare output predictions and target?
 - Use Mean Squared Error

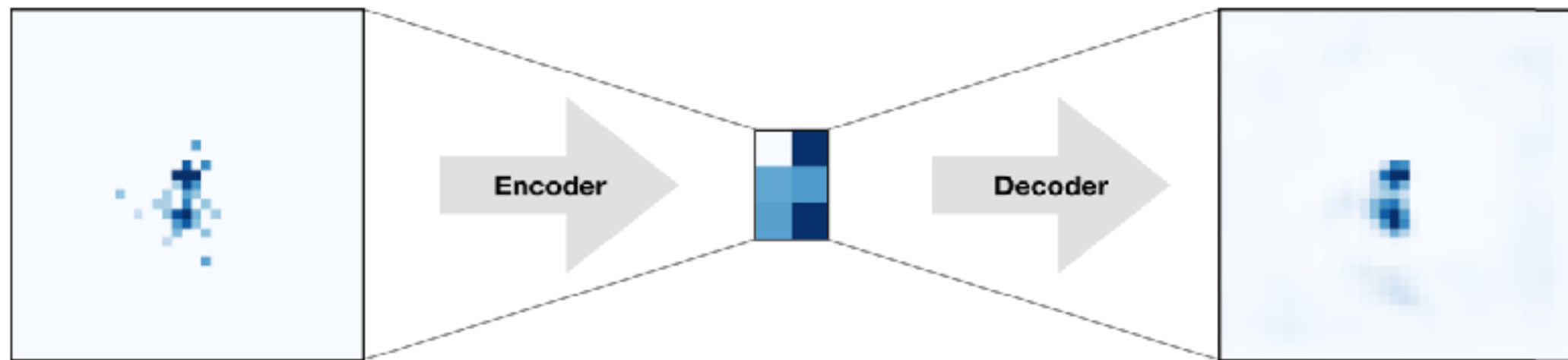
$$d_{\text{ME},2} = \text{MSE} \equiv \frac{1}{n} \sum_{i=1}^n \left| D(E(x))_i - y_i \right|^2$$

Take the average over all dimensions of the data vector

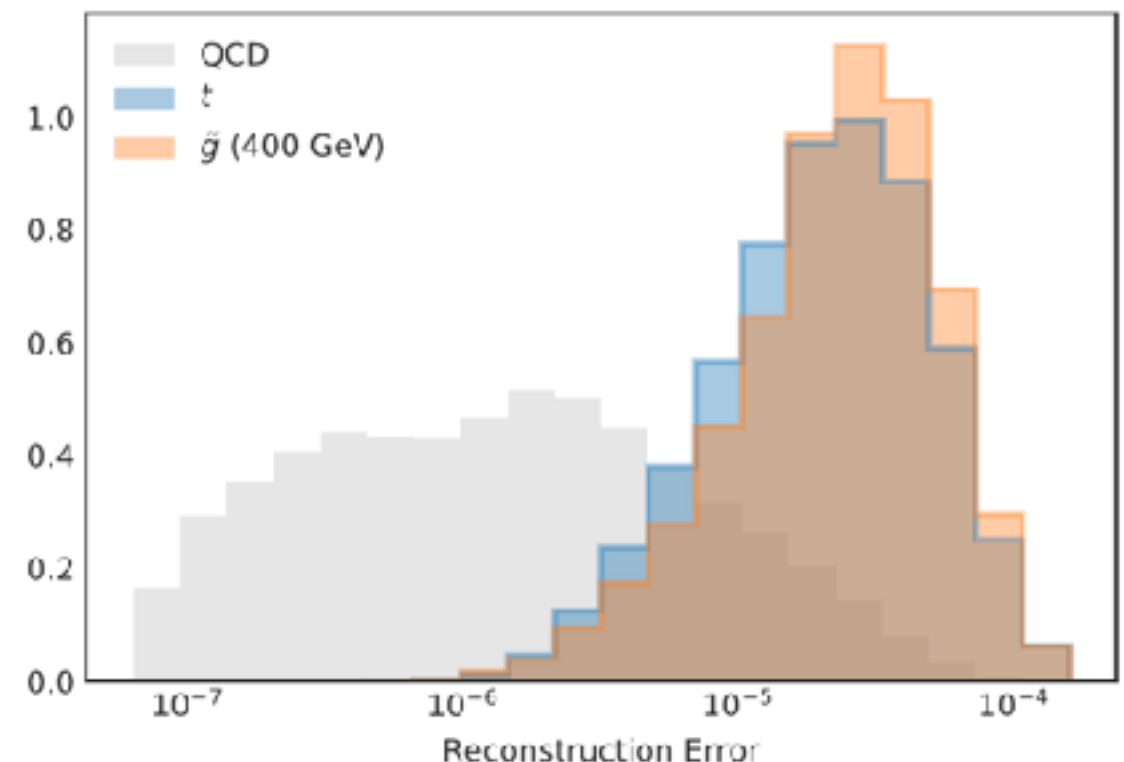
$D(E(x))$ = output of decoded, encoded data
use $f(x)$ for rest of talk

y = target (=x)
(or the input in this case)

Autoencoders for Anomaly Detection

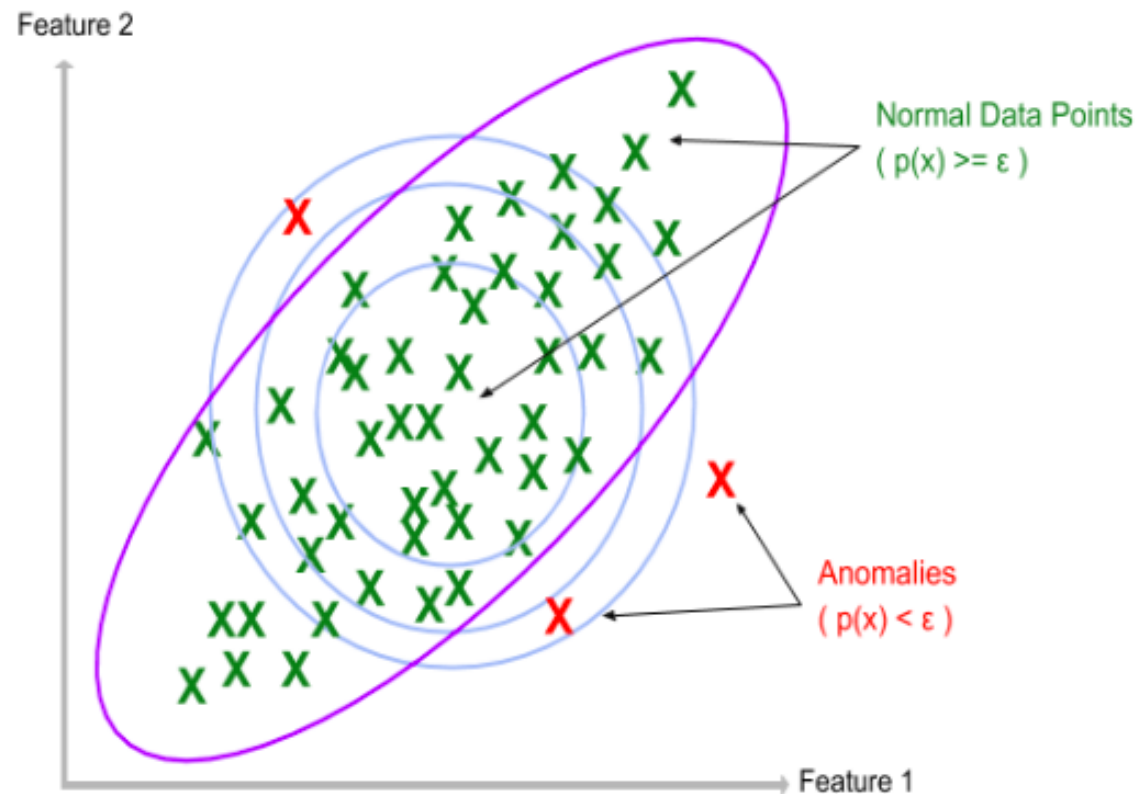


- Networks are trained to minimize the reconstruction error of events from SM background
- The encoding-decoding of BSM events will have larger reconstruction errors



[[1808.08992](#)]

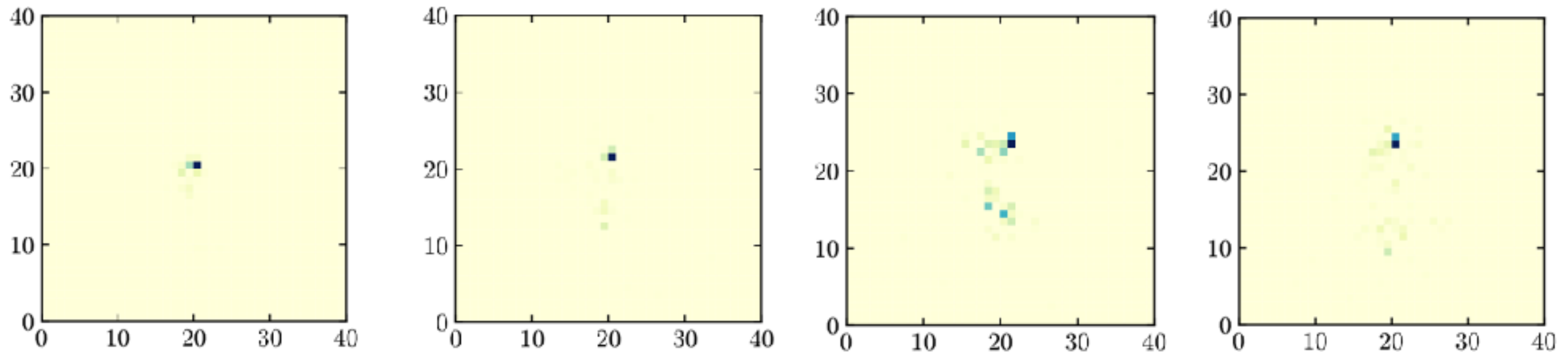
Challenges for Autoencoders



Do autoencoders fit with this picture of anomaly detection?

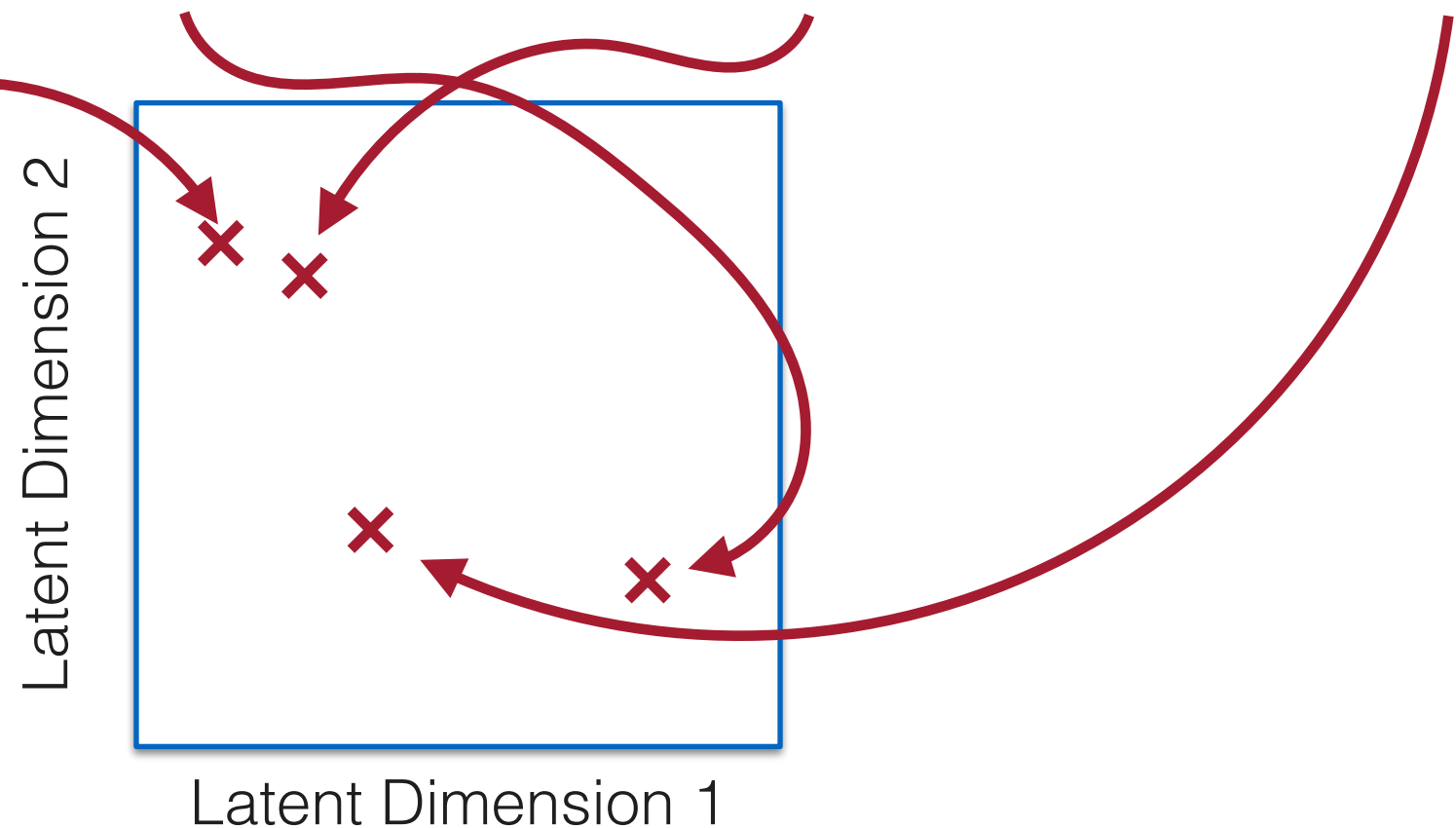
- Delicate balance between reconstructing SM well, but anomalous data poorly
- Latent space doesn't have structure, can't assign probabilistic interpretation
- Is MSE the right metric to use for reconstruction?

Plain Autoencoders

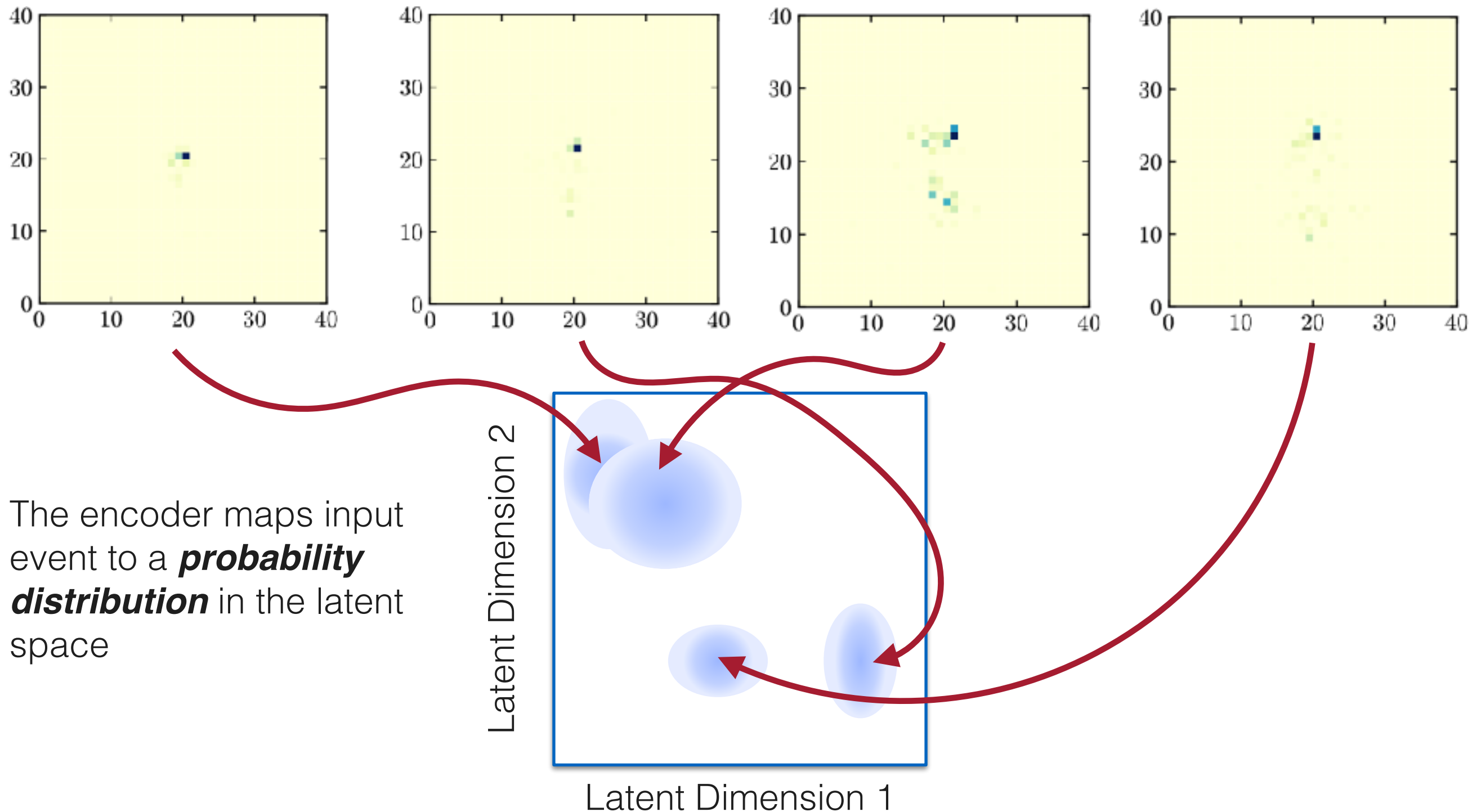


The encoder maps input event to a specific **point** in the latent space

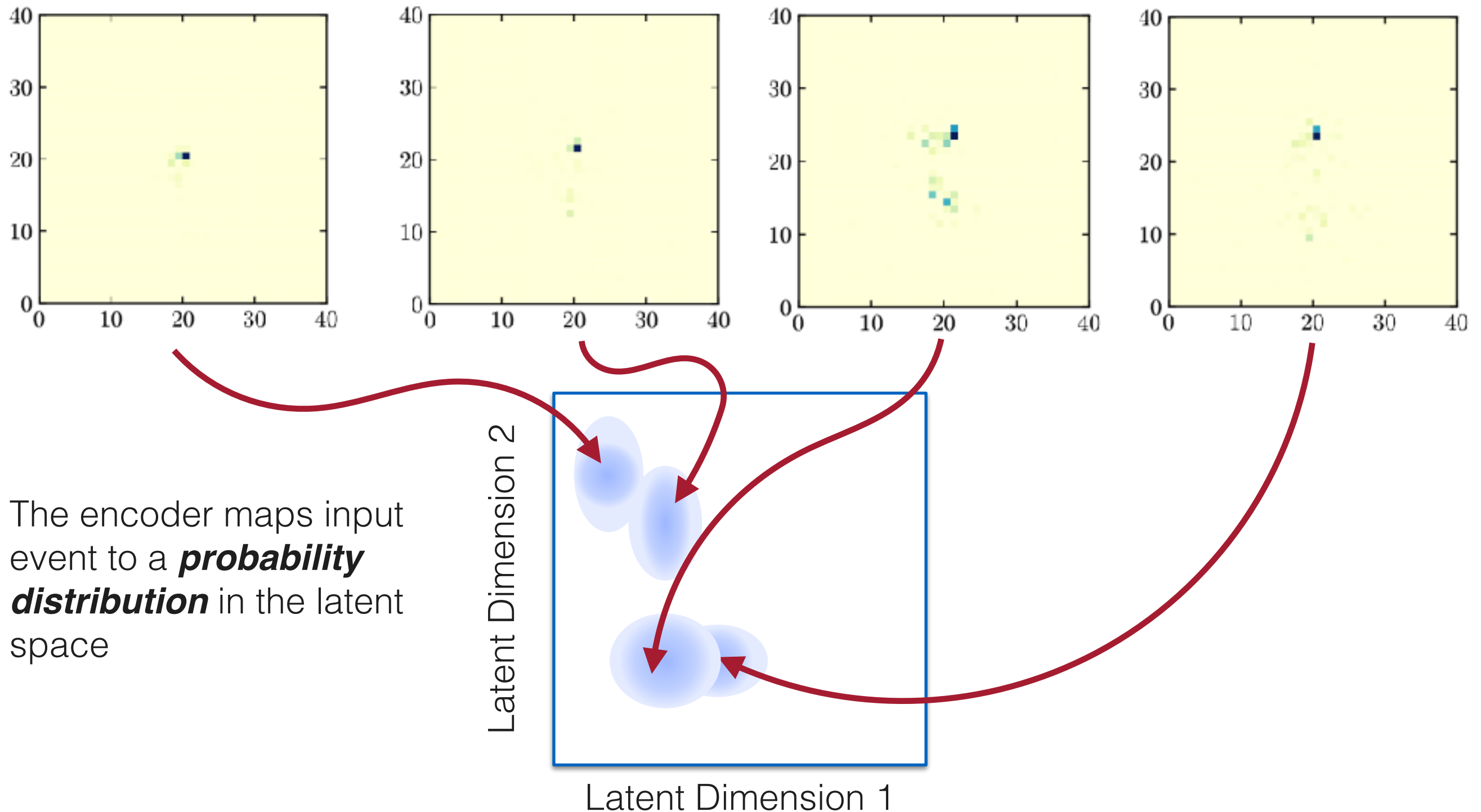
There is no intrinsic structure in the latent space of an Autoencoder



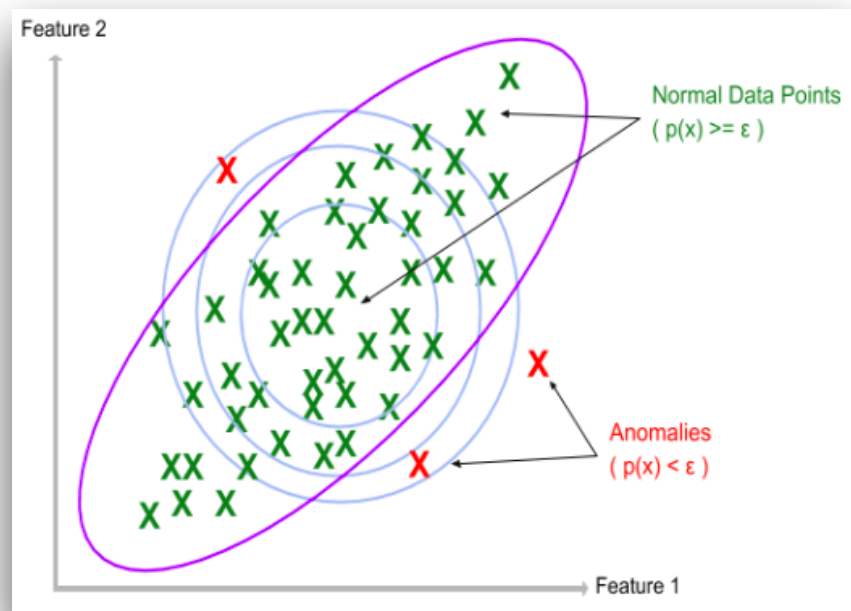
Variational Autoencoders



Variational Autoencoders



Variational Autoencoders



Add a probabilistic interpretation

How likely is a detector event,
given a point in latent space?

Model each detector pixel as a Gaussian:

$$p(E|z) \propto \exp\left(\frac{-(E - D(z))^2}{2\sigma^2}\right)$$

$$x = \{E_1, E_2, \dots, E_N\}$$

$$p(x|z) = \prod_{i=1}^N p(E_i|z)$$

Probability of the observing an event

$$p(x) = \int p(x|z)p(z)dz \equiv \mathbb{E}_{p(z)}[p(x|z)]$$

What is the latent space to
integrate over?

Variational Autoencoders

Probability of the observing an event

$$p(x) = \int p(x|z)p(z)dz \equiv \mathbb{E}_{p(z)} [p(x|z)]$$

What is the latent space to integrate over?

Use the encoder to learn/sample the latent space

$$p(x) = \int p(x|z)p(z) \frac{q(z|x)}{q(z|x)} dz = \mathbb{E}_{q(z|x)} \left[p(x|z) \frac{p(z)}{q(z|x)} \right]$$

Variational Autoencoders

Probability of the observing an event

$$p(x) = \int p(x|z)p(z)dz \equiv \mathbb{E}_{p(z)} [p(x|z)]$$

What is the latent space to integrate over?

Use the encoder to learn/sample the latent space

$$p(x) = \int p(x|z)p(z) \frac{q(z|x)}{q(z|x)} dz = \mathbb{E}_{q(z|x)} \left[p(x|z) \frac{p(z)}{q(z|x)} \right]$$

Take log likelihood and use Jensen's inequality to move log inside expectation value

$$\log p(x) \geq \mathbb{E}_{q(z|x)} \left[\log (p(x|z)) + \log \left(\frac{p(z)}{q(z|x)} \right) \right]$$

Variational Autoencoders

Probability of the observing an event

$$p(x) = \int p(x|z)p(z)dz \equiv \mathbb{E}_{p(z)} [p(x|z)]$$

What is the latent space to integrate over?

Use the encoder to learn/sample the latent space

$$p(x) = \int p(x|z)p(z) \frac{q(z|x)}{q(z|x)} dz = \mathbb{E}_{q(z|x)} \left[p(x|z) \frac{p(z)}{q(z|x)} \right]$$

Take log likelihood and use Jensen's inequality to move log inside expectation value

$$\log p(x) \geq \mathbb{E}_{q(z|x)} \left[\log (p(x|z)) + \log \left(\frac{p(z)}{q(z|x)} \right) \right]$$

~ MSE of encoded-decoded event



Variational Autoencoders

Probability of the observing an event

$$p(x) = \int p(x|z)p(z)dz \equiv \mathbb{E}_{p(z)} [p(x|z)]$$

What is the latent space to integrate over?

Use the encoder to learn/sample the latent space

$$p(x) = \int p(x|z)p(z) \frac{q(z|x)}{q(z|x)} dz = \mathbb{E}_{q(z|x)} \left[p(x|z) \frac{p(z)}{q(z|x)} \right]$$

Take log likelihood and use Jensen's inequality to move log inside expectation value

$$\log p(x) \geq \mathbb{E}_{q(z|x)} \left[\log(p(x|z)) + \log \left(\frac{p(z)}{q(z|x)} \right) \right]$$

~ MSE of encoded-decoded event

Kullback-Leibler Divergence (KLD)
between encoded representation
and latent prior

Summary of Autoencoders

Plain autoencoders: no sampling, $L = \text{MSE}$

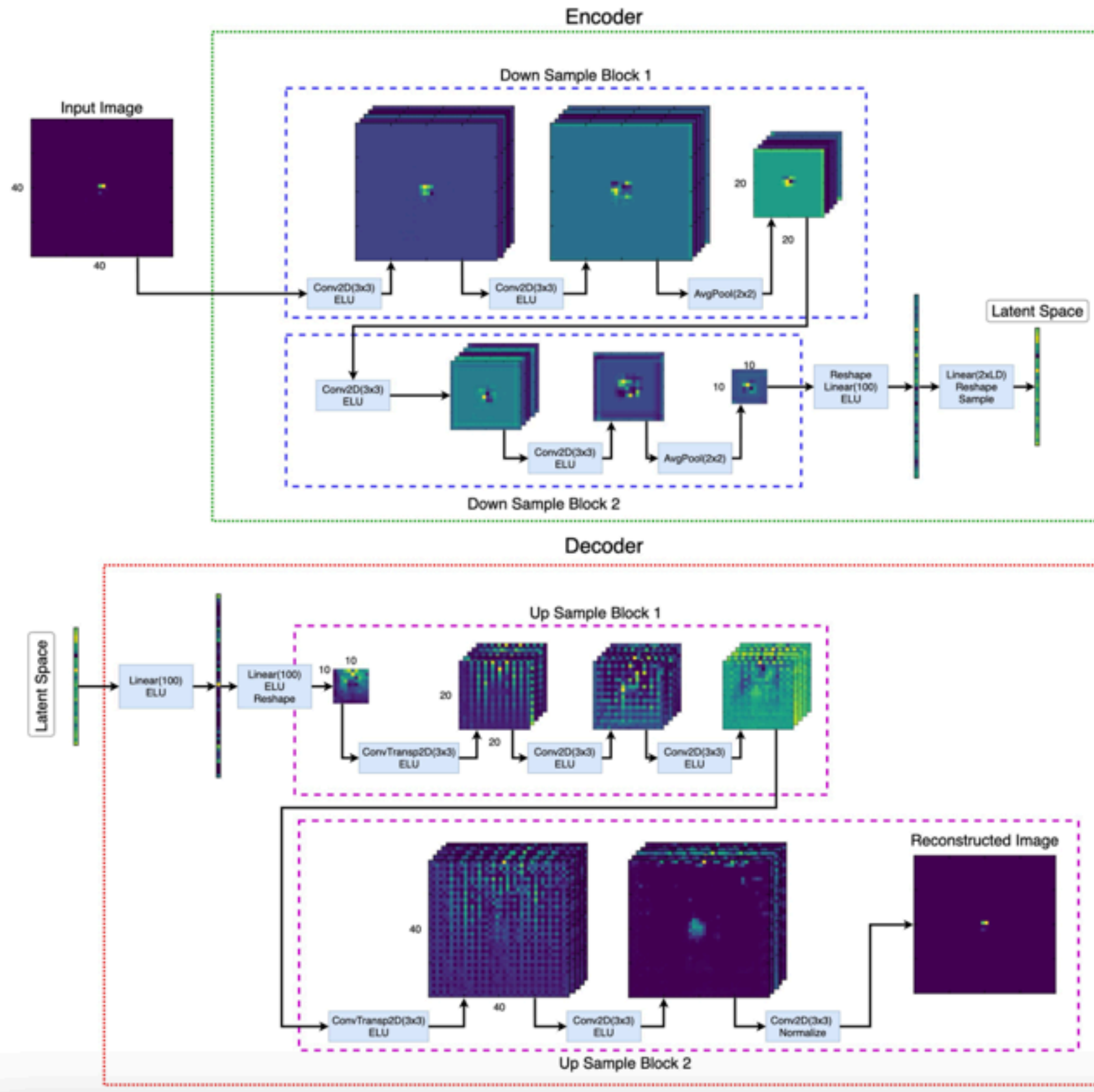
Variation autoencoders: sampling in latent space, $L = \text{MSE} + \text{KLD}$

Sampling adds/forces structure on the latent space.
KLD term helps to regularize and adds to
probabilistic interpretation

In practice, these terms may be too far apart,
introduce a scaling between them

$$L = (1 - \beta) \text{MSE} + \beta \text{KLD}$$

Generative Networks



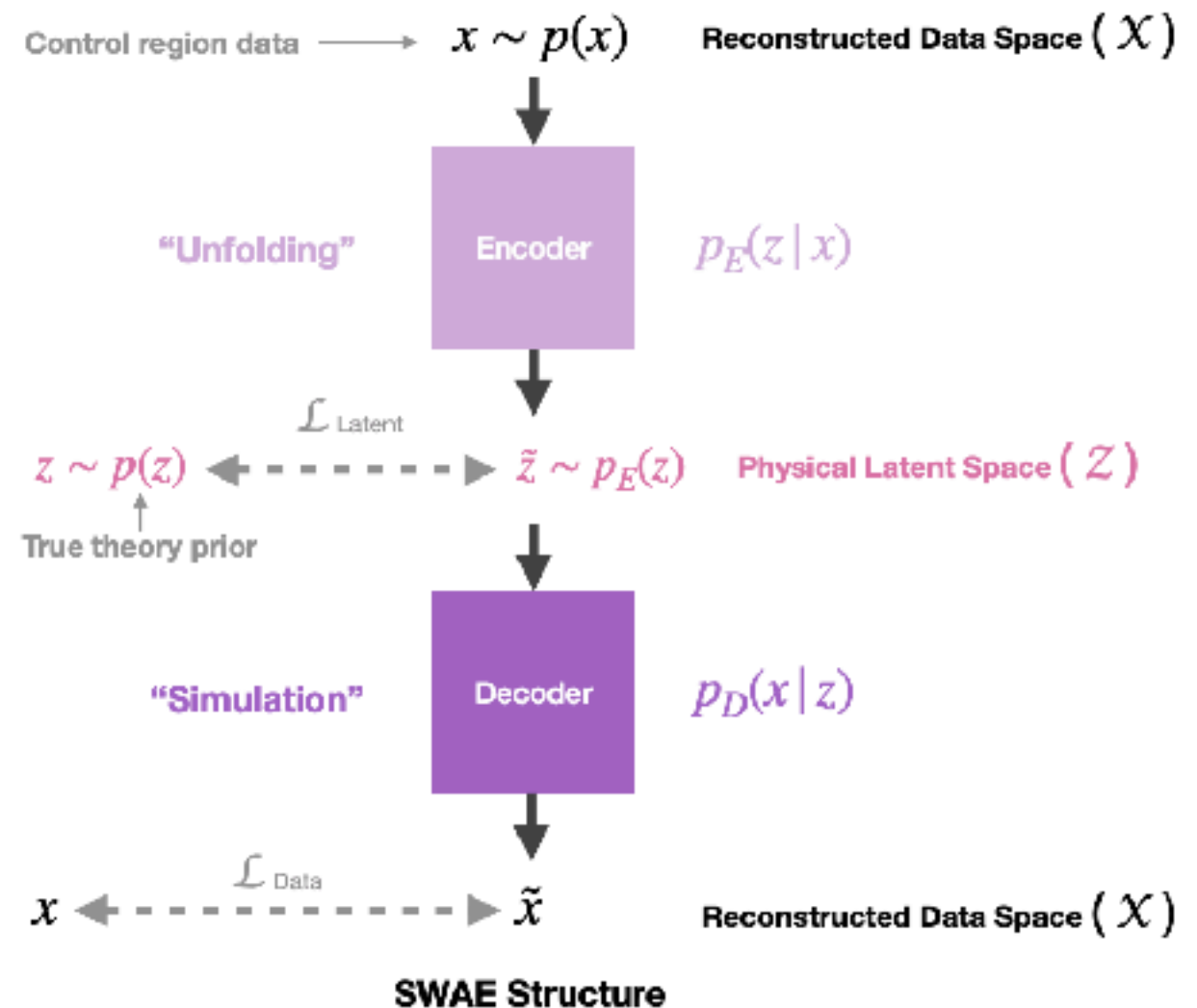
Can randomly sample in the latent space to get “new” data

Alter the latent space

- Latent loss ($\mathcal{L}_{\text{Latent}}$)
 - SW distance for finite samples
- Data loss ($\mathcal{L}_{\text{Data}}$):
 - Mean Squared Error (MSE)
- Total SWAE loss function:

$$\mathcal{L}_{\text{SWAE}} = \mathcal{L}_{\text{Data}} + \lambda \mathcal{L}_{\text{Latent}}$$
- Easy to add additional physically-motivated constraints

$$\mathcal{L} = \mathcal{L}_{\text{SWAE}} + \lambda'_i \mathcal{L}_i$$

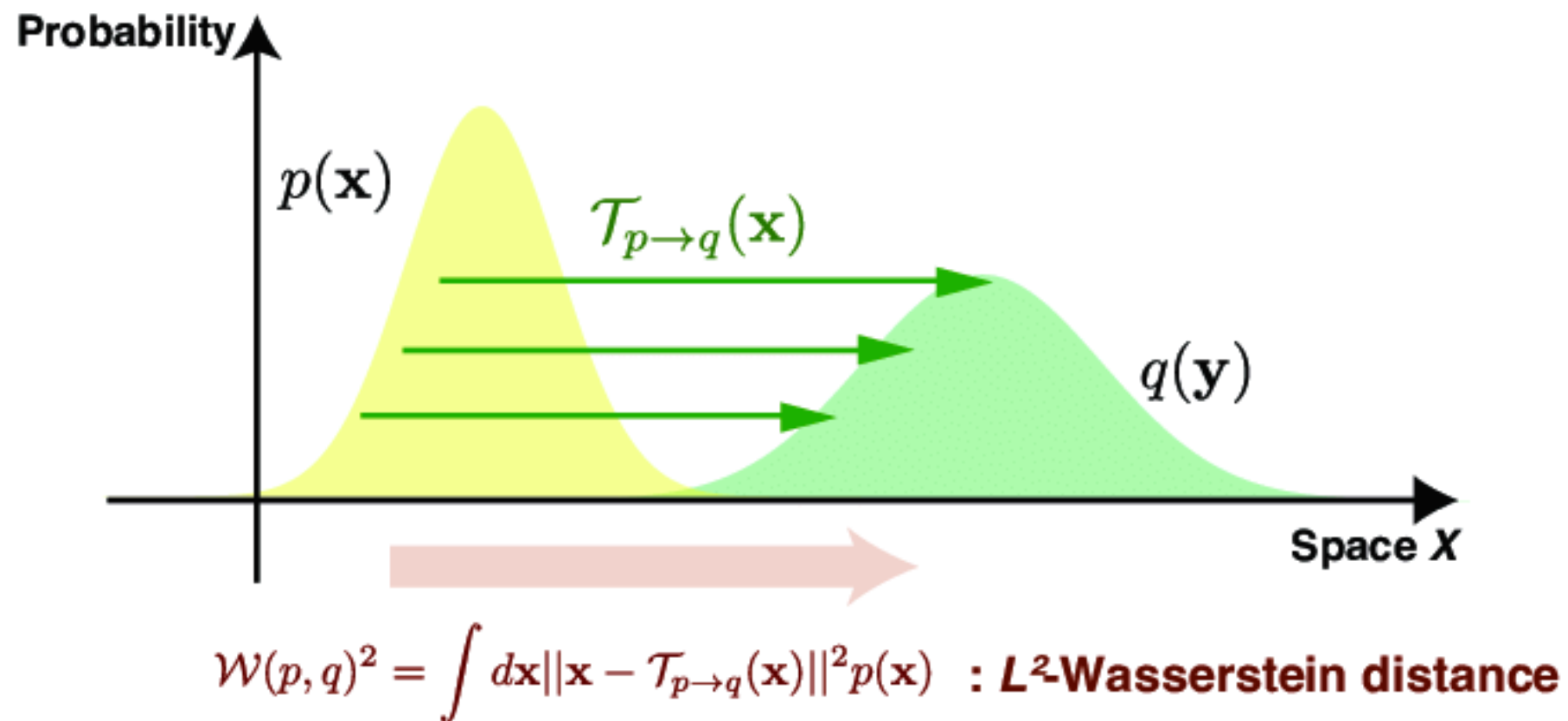


In the previous discussion of VAEs, we used a Gaussian for the prior of the latent space. However, it is also possible to use different/more motivated priors. For instance, use the parton-level distributions.

[2101.08944]

Alter the latent space

Use the Earth mover's distance (Wasserstein distance) to compare the latent prior with the learned latent distributions



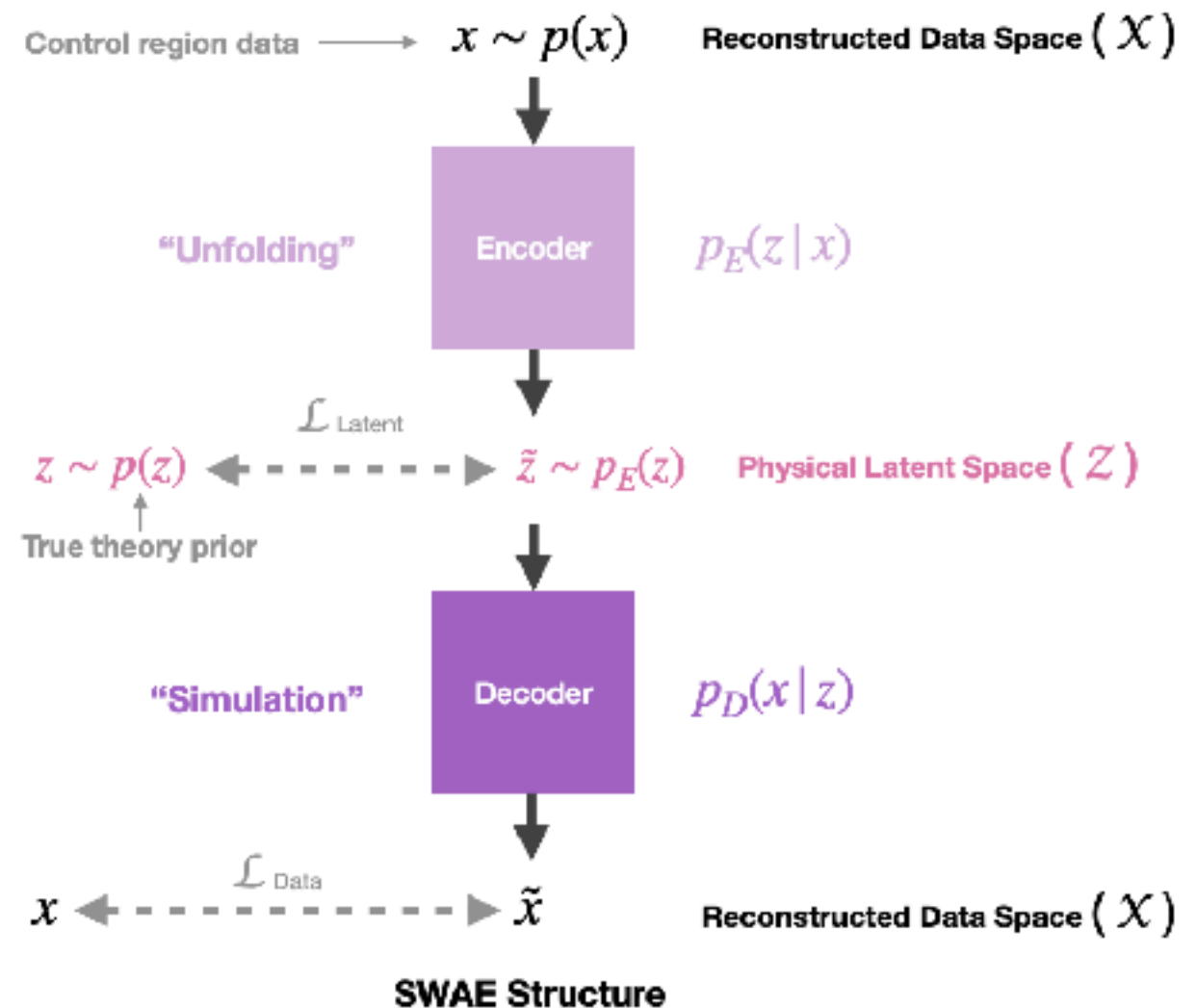
Side note: Wasserstein distances are also an active area of research for the other areas of machine learning in HEP, such as classification and anomaly detection.

Alter the latent space

- Latent loss ($\mathcal{L}_{\text{Latent}}$)
 - SW distance for finite samples
- Data loss ($\mathcal{L}_{\text{Data}}$):
 - Mean Squared Error (MSE)
- Total SWAE loss function:

$$\mathcal{L}_{\text{SWAE}} = \mathcal{L}_{\text{Data}} + \lambda \mathcal{L}_{\text{Latent}}$$
- Easy to add additional physically-motivated constraints

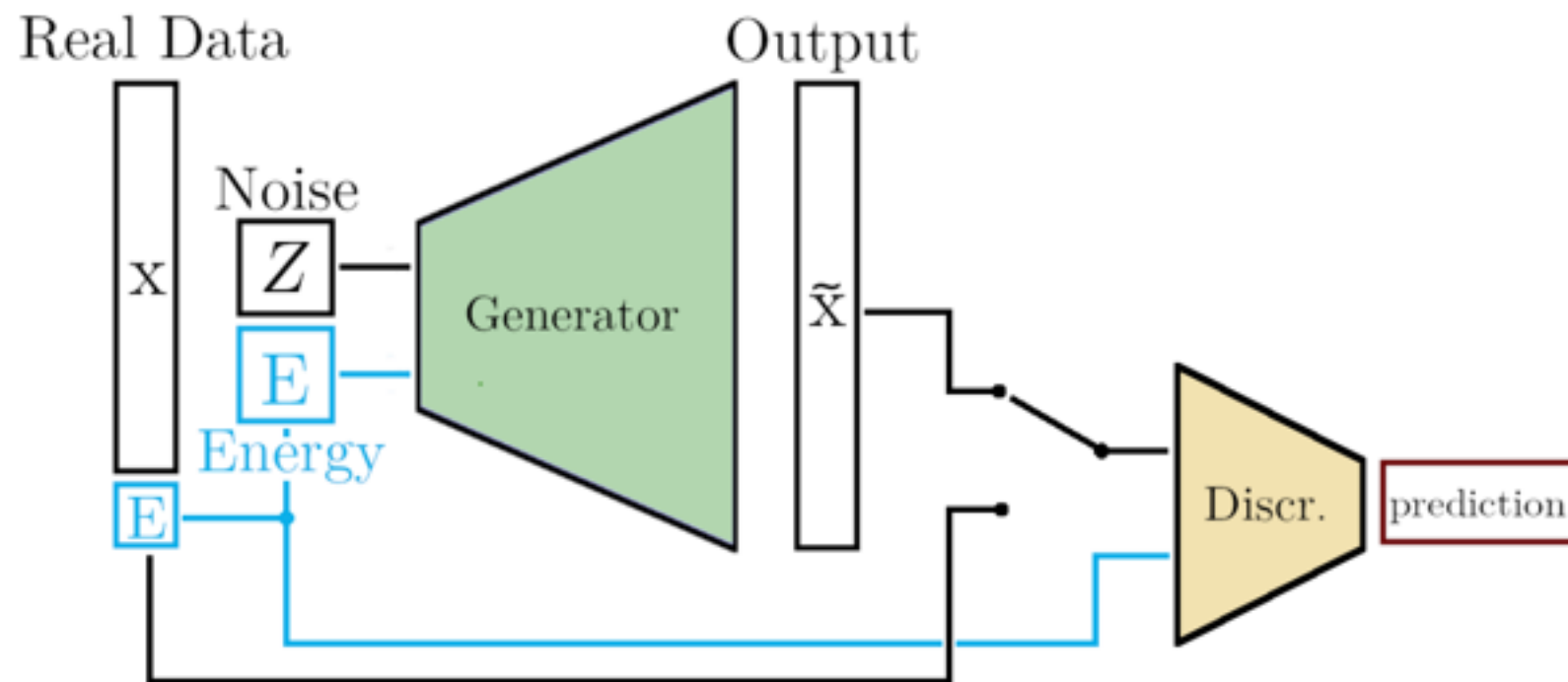
$$\mathcal{L} = \mathcal{L}_{\text{SWAE}} + \lambda'_i \mathcal{L}_i$$



Example in paper: Examine $p p \rightarrow t \bar{t}$ events with semi-leptonic decays,
 X is the detector level information (lepton 4-vector, MET, and jet four vectors)
 Z is the parton level information (lepton, neutrino, and quark four vectors)
 Many challenges, active research

[[2101.08944](#)]

Other generative models



Generative Adversarial Networks (GAN) work by taking in random noise and generating an event (image, observables, etc). There are two components to this, first is the network which generates the event, second is the discriminator network. This network is trained to tell the difference between real and generated data. The generator is trained to trick the discriminator.

Over 55 papers in the HEP literature using GANs

Other generative models

Ideas seen for generating events:

GANs -> Noise to data

VAES -> Latent representation to data

In both of these, we are trying to capture something about the underlying distributions of the true data, why not try to learn that directly?

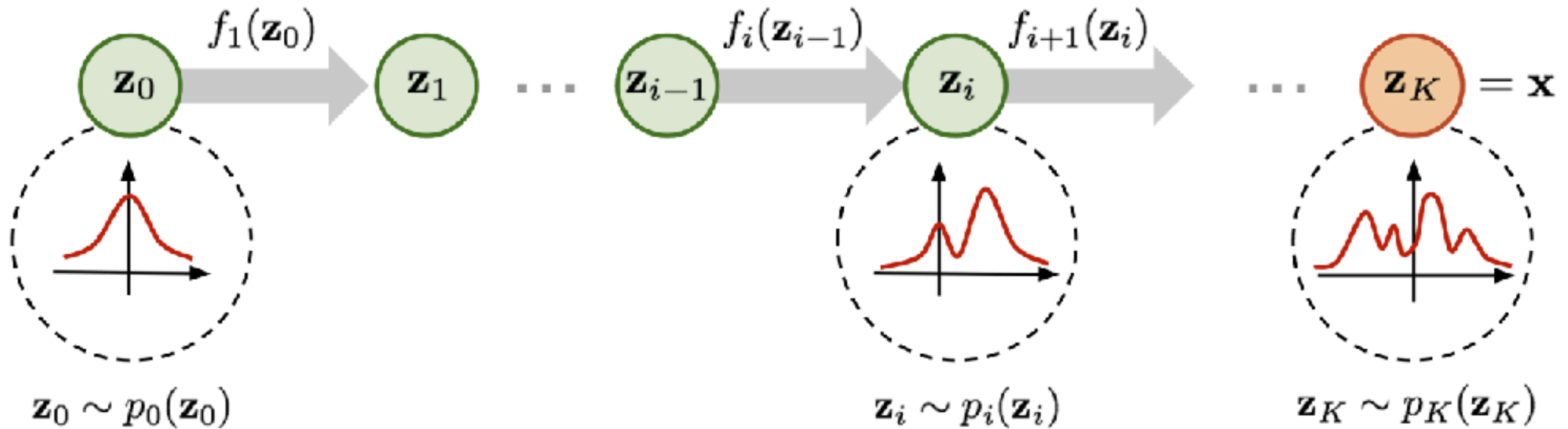
Normalizing Flows

Normalizing Flows

- Normalizing flows implement a change of variables from a simple distribution to the real distribution.
- Often use special layers which are easily invertible and have a simple Jacobian
- Once trained, can draw a random sample from the easy distribution and obtain a sample from a complex distribution
- Invertibility then also allows for taking a sample from the complex distribution and computing the likelihood

Nice review of the methods [[1912.02762.pdf](#)] (not HEP)

Normalizing Flows



These models allow for easy density estimation (learn the PDF of the data)

Allows for: Generating events, Generating field configurations for Lattice Gauge Theory, or Anomaly Detection

Conclusions

- There are many types of machine learning which do not simply use labeled data for supervision
- We saw how AEs for anomaly detection do not build in a probabilistic interpretations, but a VAE does
- VAEs can also be used to generate data (and the latent space can be made more intuitive)
- There are many other ways to make generative networks
- We can use ML to estimate probabilities which allows for inference as well

Ideas for Project (Tutorial 3)

- 1) Compare raw-4vectors, n-subjettiness basis, and jet-images for top-tagging dataset
 - a) Which does best?
 - b) Are the number of parameters the similar?
 - c) Speed of training?
- 2) Use **anomaly detection** techniques to train only on the background data. Then apply to the test set and see how the performance compares to supervised classification.

References

<https://iml-wg.github.io/HEPML-LivingReview/>

Particle Data Group has a new review chapter on ML which covers all of these techniques and more <https://pdg.lbl.gov/>