

Generative Adversarial Networks and their application to fast simulation in HEP



LAMBDA • HSE

Artem Maevskiy

National Research University Higher School of Economics,
Laboratory of Methods for Big Data Analysis

Introduction

- Generative Adversarial Networks (or simply GANs) are famous for their ability to create realistic images
- E.g., among the images below, can you tell which was generated and which is real?



Introduction

- Generative Adversarial Networks (or simply GANs) are famous for their ability to create realistic images
- E.g., among the images below, can you tell which was generated and which is real?



GENERATED



Introduction

- Generative Adversarial Networks (or simply GANs) are famous for their ability to create realistic images
- E.g., among the images below, can you tell which was generated and which is real?



GENERATED



GENERATED



Introduction

- Generative Adversarial Networks (or simply GANs) are famous for their ability to create realistic images
- E.g., among the images below, can you tell which was generated and which is real?



GENERATED



GENERATED



GENERATED

Introduction

- Generative Adversarial Networks (or simply GANs) are famous for their ability to create realistic images
- E.g., among the images below, can you tell which was generated and which is real?



GENERATED



GENERATED



GENERATED

<https://www.thispersondoesnotexist.com/>

GANs in HEP?

- On previous slide: examples of **human face images** that our **brain accepts as real**
- Now, replace:
 - «**human face images**» → «**vectors of numbers describing collision events**»
 - «**brain accepts as real**» → «**analysis remains unbiased with**»
- and you get a good Monte-Carlo generator
 - typically much faster than an actual generator from HEP

Simulation workflow



- One may imagine any part of this chain to be replaced by GAN

Simulation workflow



- One may imagine any part of this chain to be replaced by GAN
- Here we'll demonstrate the following approach:



Simulation workflow



- One may imagine any part of this chain to be replaced by GAN
- Here we'll demonstrate the following approach:



- Disclaimer: we're not talking about the actual total replacement of honest full simulation, but rather about new means for creating simplified fast-sim models

Outline

- A little theory
 - 'Vanilla' GAN and its problems
 - Some advancements (Wasserstein GAN)
- Using GANs for fast simulation at LHCb

Generative Models

Problem statement:

Given a finite sample of i.i.d. examples from an unknown distribution

$$\{\mathbf{x}_i\} \sim p_{\text{data}}$$

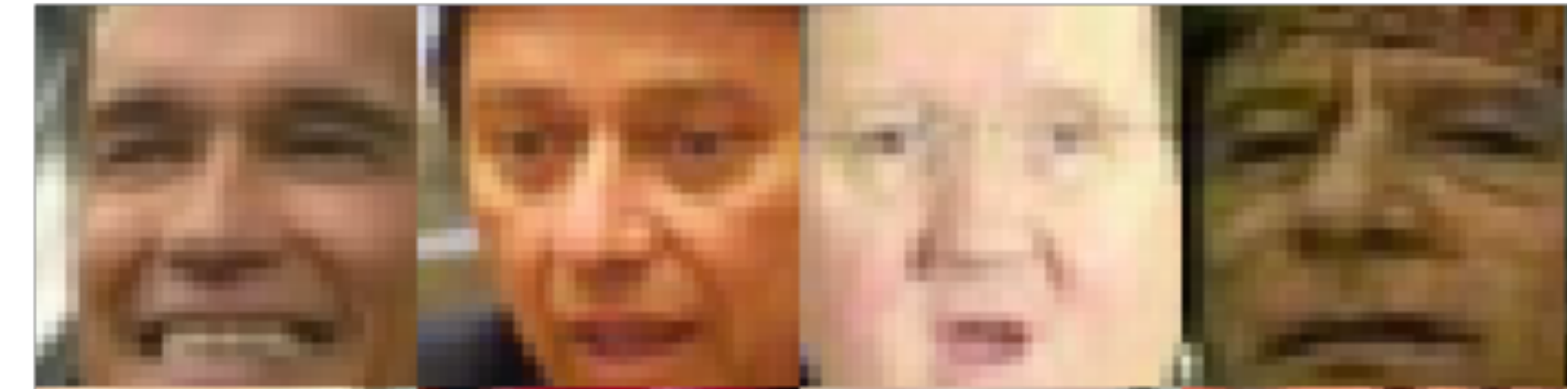


Generative Models

Problem statement:

Given a finite sample of i.i.d. examples from an unknown distribution

$$\{\mathbf{x}_i\} \sim p_{\text{data}}$$



We want to learn to draw more samples from that distribution

$$\{\mathbf{x}_j\} \sim p_{\text{data}}$$



Generative Models

Problem statement:

Given a finite sample of i.i.d. examples from an unknown distribution

$$\{x_i\} \sim p_{\text{data}}$$



We want to learn to draw more samples from ~~that distribution~~

a 'similar' distribution

~~$$\{x_j\} \sim p_{\text{data}}$$~~



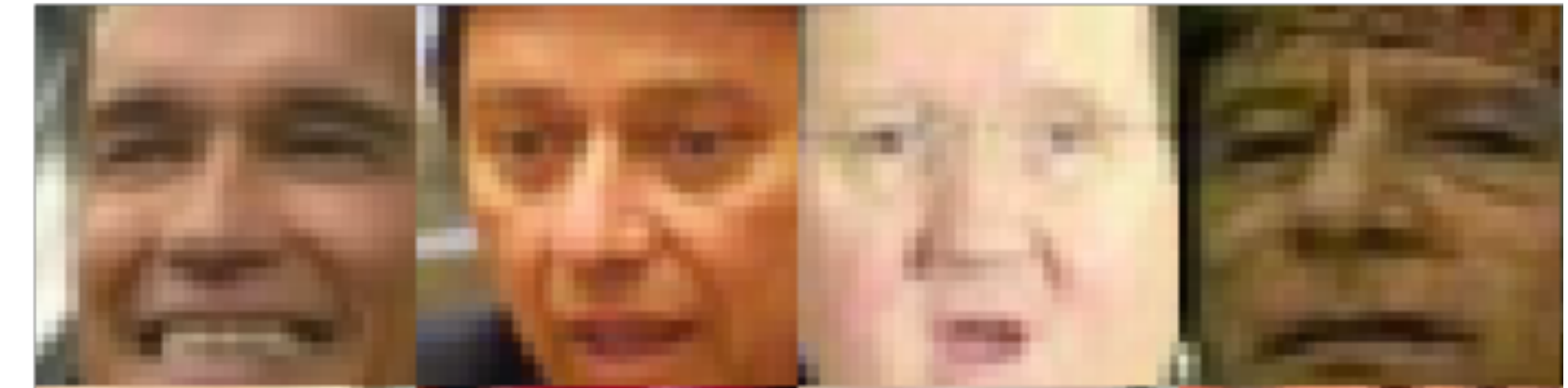
$$p_{\text{gen},\theta} \xrightarrow{\theta \rightarrow \theta^*} p_{\text{data}}$$

Generative Models

Problem statement:

Given a finite sample of i.i.d. examples from an unknown distribution

$$\{x_i\} \sim p_{\text{data}}$$



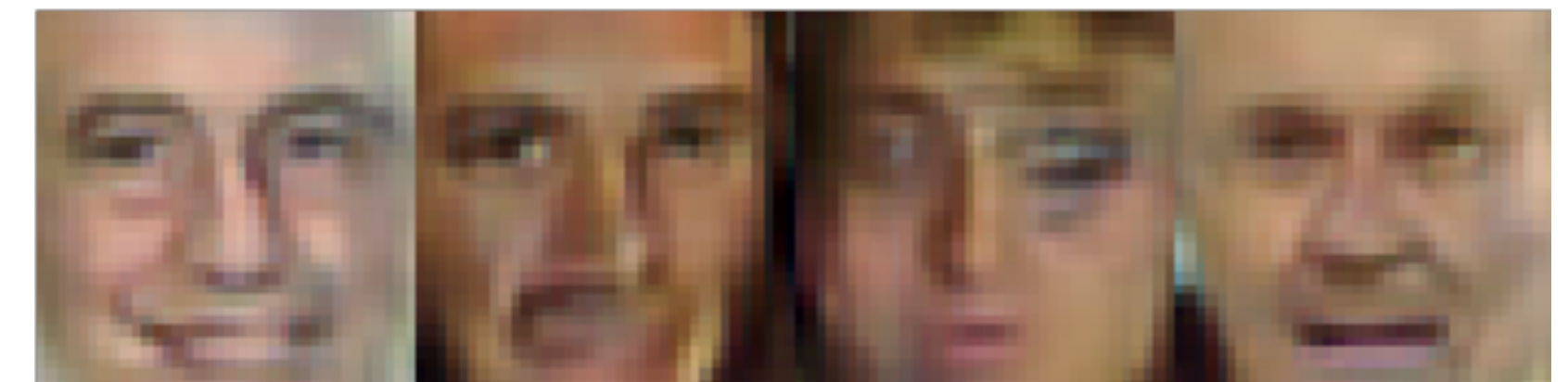
We want to learn to draw more samples from ~~that distribution~~

a 'similar' distribution

~~$$\{x_j\} \sim p_{\text{data}}$$~~



$$\{\tilde{x}_j\} \sim p_{\text{gen},\theta}$$



$$p_{\text{gen},\theta} \xrightarrow{\theta \rightarrow \theta^*} p_{\text{data}}$$

Generative Adversarial Networks

$\{\tilde{\mathbf{x}}_j\} \sim p_{gen,\theta}$ — this distribution is built in the following way:

$$\tilde{\mathbf{x}}_j = G_\theta(\mathbf{z}_j)$$

$$\mathbf{z}_j \sim N(\mathbf{0}, \mathbb{I})$$

G_θ — generator neural network

I. Goodfellow, et. al. *Generative adversarial nets*. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

Generative Adversarial Networks

$\{\tilde{\mathbf{x}}_j\} \sim p_{gen,\theta}$ — this distribution is built in the following way:

$$\tilde{\mathbf{x}}_j = G_\theta(\mathbf{z}_j)$$

$$\mathbf{z}_j \sim N(\mathbf{0}, \mathbb{I})$$

G_θ — generator neural network

We add a classifying network D_ϕ (discriminator) to distinguish between the real and generated samples and train both as follows:

$$V(\phi, \theta) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{gen,\theta}} [\log (1 - D_\phi(\tilde{\mathbf{x}}))]$$

$$V(\phi, \theta) \rightarrow \min_{\theta} \max_{\phi} V(\phi, \theta)$$

I. Goodfellow, et. al. *Generative adversarial nets*. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

Generative Adversarial Networks

$\{\tilde{\mathbf{x}}_j\} \sim p_{gen,\theta}$ — this distribution is built in the following way:

$$\tilde{\mathbf{x}}_j = G_\theta(\mathbf{z}_j)$$

$$\mathbf{z}_j \sim N(\mathbf{0}, \mathbb{I})$$

G_θ — generator neural network

Probability the sample was drawn from data

Probability the sample was generated

We add a classifying network D_ϕ (discriminator) to distinguish between the real and generated samples and train both as follows:

$$V(\phi, \theta) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{gen,\theta}} [\log (1 - D_\phi(\tilde{\mathbf{x}}))]$$

$$V(\phi, \theta) \rightarrow \min_{\theta} \max_{\phi} V(\phi, \theta)$$

I. Goodfellow, et. al. *Generative adversarial nets*. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

Generative Adversarial Networks

$\{\tilde{\mathbf{x}}_j\} \sim p_{gen,\theta}$ — this distribution is built in the following way:

$$\tilde{\mathbf{x}}_j = G_\theta(\mathbf{z}_j)$$

$$\mathbf{z}_j \sim N(\mathbf{0}, \mathbb{I})$$

G_θ — generator neural network

Probability the
sample was
drawn from data

Probability the
sample was
generated

We add a classifying network D_ϕ (discriminator) to distinguish between the real and generated samples and train both as follows:

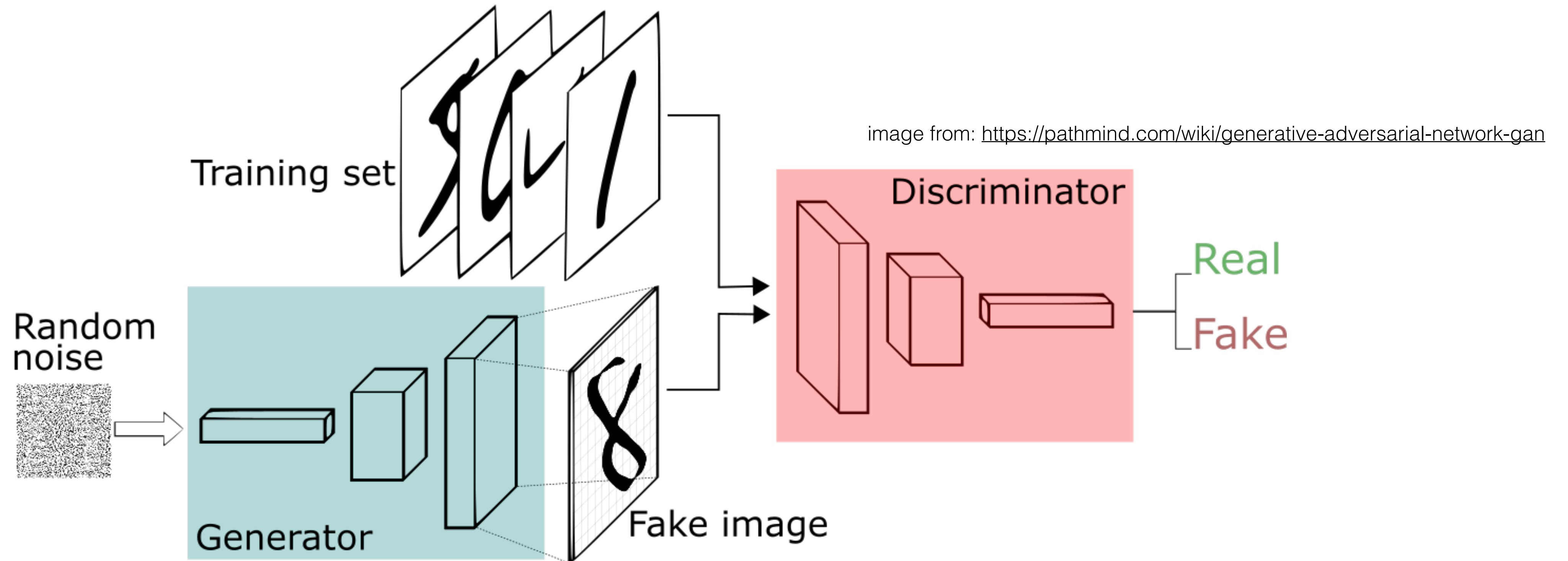
$$V(\phi, \theta) = \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_\phi(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{gen,\theta}} [\log (1 - D_\phi(\tilde{\mathbf{x}}))]$$

$$V(\phi, \theta) \rightarrow \min_{\theta} \max_{\phi} V(\phi, \theta)$$

Likelihood

I. Goodfellow, et. al. *Generative adversarial nets*. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

Game interpretation



Min-max game:

- goal of discriminator: distinguish between real and generated samples
- goal of generator: 'fool' the discriminator

Problems with GANs

(vanishing gradients)

Problems with GANs

(vanishing gradients)

Let's have another look at the objective:

$$\begin{aligned} V(\phi, \theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{gen}, \theta}} [\log (1 - D_{\phi}(\tilde{\mathbf{x}}))] \\ &= \int [p_{\text{data}}(\mathbf{x}) \log D_{\phi}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x}) \log (1 - D_{\phi}(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

Problems with GANs

(vanishing gradients)

Let's have another look at the objective:

$$\begin{aligned} V(\phi, \theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{gen}, \theta}} [\log (1 - D_{\phi}(\tilde{\mathbf{x}}))] \\ &= \int [p_{\text{data}}(\mathbf{x}) \log D_{\phi}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x}) \log (1 - D_{\phi}(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

For a given $p_{\text{gen}, \theta}$ the optimal discriminator is $D_{\phi^*}(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})}$

Problems with GANs

(vanishing gradients)

Let's have another look at the objective:

$$\begin{aligned} V(\phi, \theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{gen}, \theta}} [\log (1 - D_{\phi}(\tilde{\mathbf{x}}))] \\ &= \int [p_{\text{data}}(\mathbf{x}) \log D_{\phi}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x}) \log (1 - D_{\phi}(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

For a given $p_{\text{gen}, \theta}$ the optimal discriminator is $D_{\phi^*}(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})}$

So the objective is:

$$V(\phi^*, \theta) = \int \left[p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} + p_{\text{gen}, \theta}(\mathbf{x}) \log \frac{p_{\text{gen}, \theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} \right] d\mathbf{x}$$

Problems with GANs

(vanishing gradients)


Let's have another look at the objective:

$$\begin{aligned} V(\phi, \theta) &= \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{gen}, \theta}} [\log (1 - D_{\phi}(\tilde{\mathbf{x}}))] \\ &= \int [p_{\text{data}}(\mathbf{x}) \log D_{\phi}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x}) \log (1 - D_{\phi}(\mathbf{x}))] d\mathbf{x} \end{aligned}$$

For a given $p_{\text{gen}, \theta}$ the optimal discriminator is $D_{\phi^*}(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})}$

So the objective is: $V(\phi^*, \theta) = \int \left[p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} + p_{\text{gen}, \theta}(\mathbf{x}) \log \frac{p_{\text{gen}, \theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} \right] d\mathbf{x}$

In case p_{data} and $p_{\text{gen}, \theta}$ have non-overlapping support



Problems with GANs

(vanishing gradients)

Let's have another look at the objective:

$$V(\phi, \theta) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\tilde{\mathbf{x}} \sim p_{\text{gen}, \theta}} [\log (1 - D_{\phi}(\tilde{\mathbf{x}}))] \\ = \int [p_{\text{data}}(\mathbf{x}) \log D_{\phi}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x}) \log (1 - D_{\phi}(\mathbf{x}))] d\mathbf{x}$$

For a given $p_{\text{gen}, \theta}$ the optimal discriminator is $D_{\phi^*}(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})}$

So the objective is:

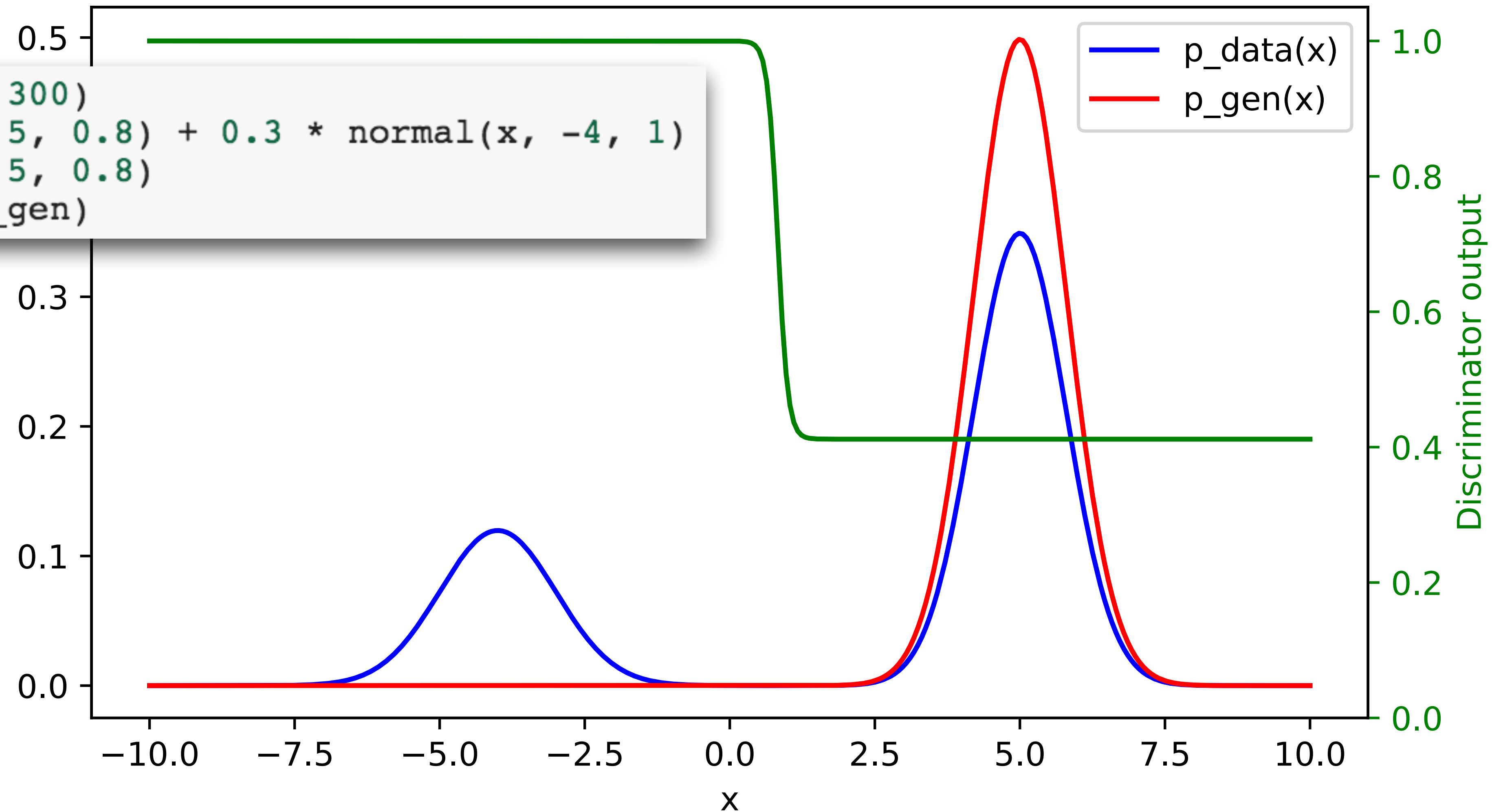
$$V(\phi^*, \theta) = \int \left[p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} + p_{\text{gen}, \theta}(\mathbf{x}) \log \frac{p_{\text{gen}, \theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} \right] d\mathbf{x} = 0$$

(= const)

In case p_{data} and $p_{\text{gen}, \theta}$ have non-overlapping support

Problems with GANs (mode collapse)

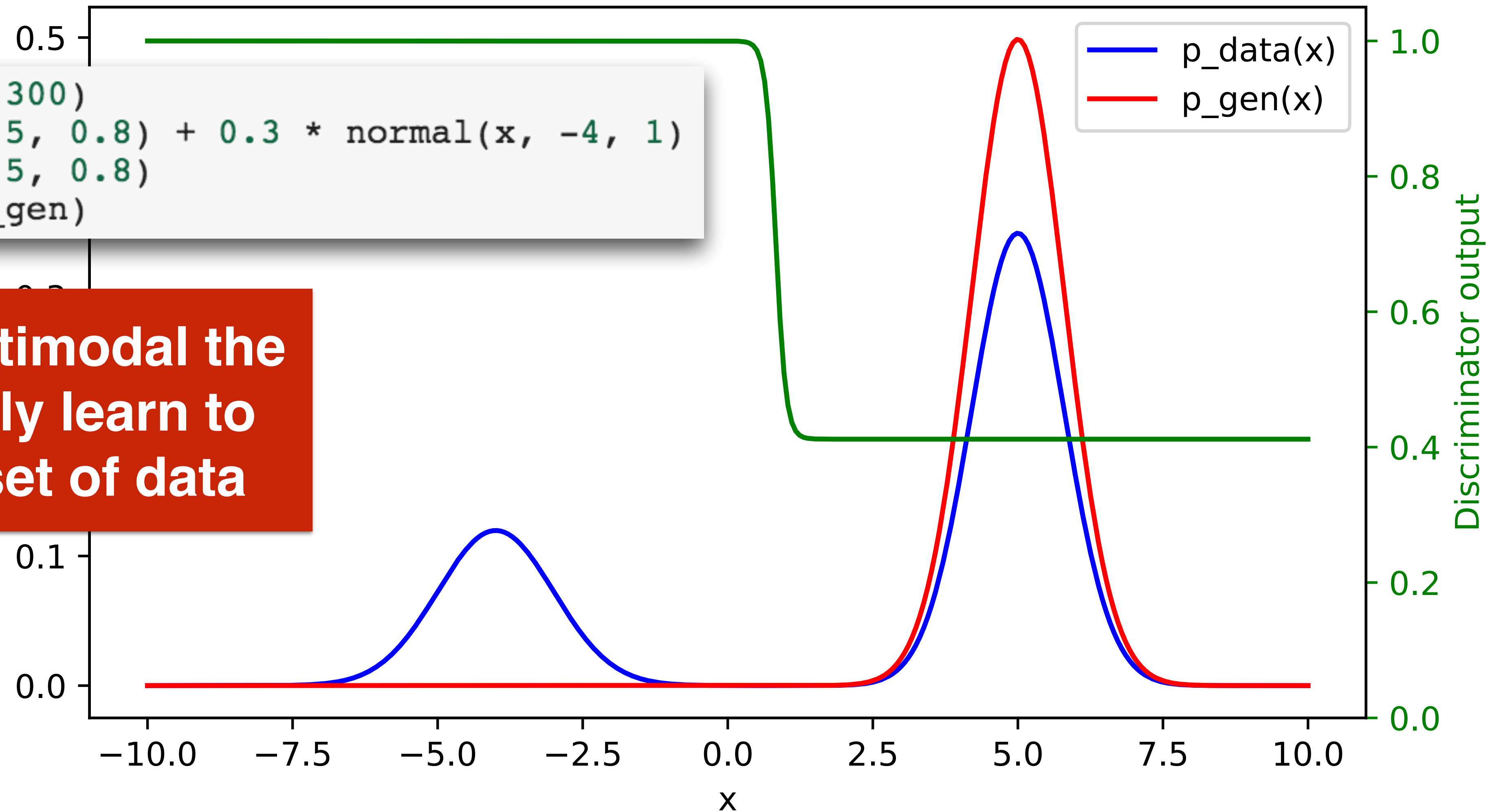
```
x = np.linspace(-10, 10, 300)
p_data = 0.7 * normal(x, 5, 0.8) + 0.3 * normal(x, -4, 1)
p_gen = 1.0 * normal(x, 5, 0.8)
D = p_data / (p_data + p_gen)
```



Problems with GANs (mode collapse)

```
x = np.linspace(-10, 10, 300)
p_data = 0.7 * normal(x, 5, 0.8) + 0.3 * normal(x, -4, 1)
p_gen = 1.0 * normal(x, 5, 0.8)
D = p_data / (p_data + p_gen)
```

In case data is multimodal the generator may only learn to reproduce a subset of data

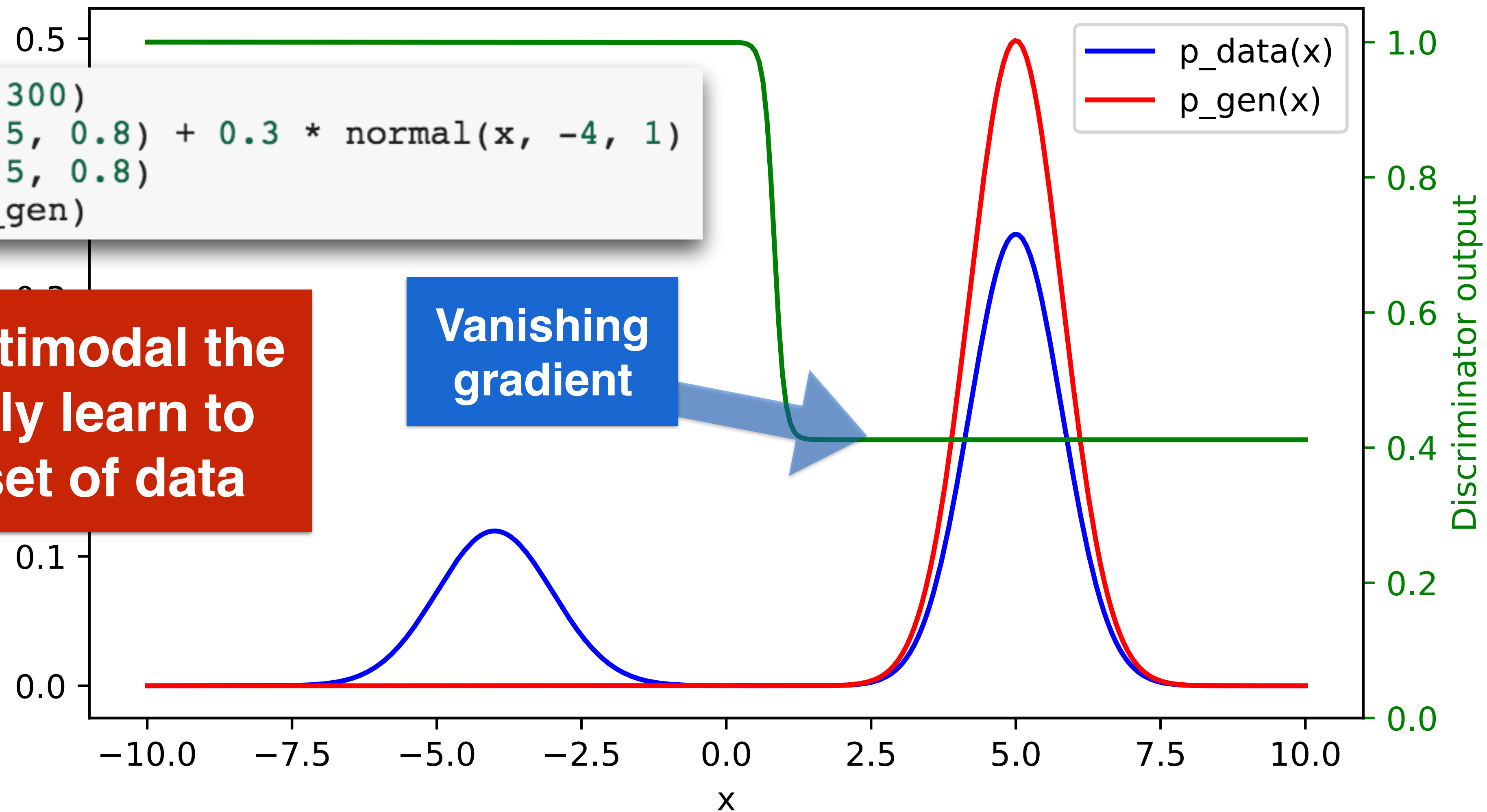


Problems with GANs (mode collapse)

```
x = np.linspace(-10, 10, 300)
p_data = 0.7 * normal(x, 5, 0.8) + 0.3 * normal(x, -4, 1)
p_gen = 1.0 * normal(x, 5, 0.8)
D = p_data / (p_data + p_gen)
```

In case data is multimodal the generator may only learn to reproduce a subset of data

Vanishing gradient



Problems with GANs

(mode collapse)

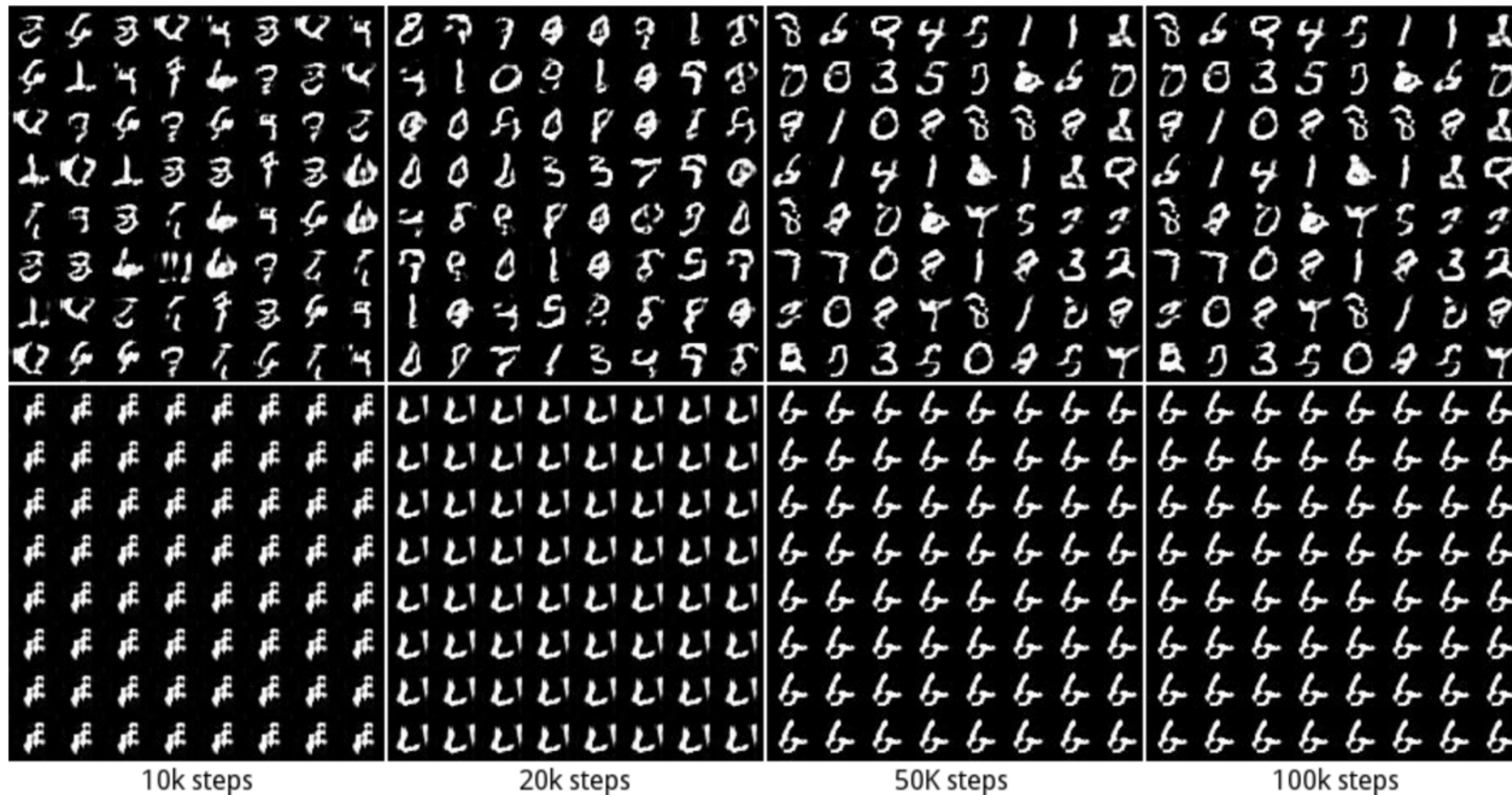


image taken from: https://medium.com/@jonathan_hui/gan-why-it-is-so-hard-to-train-generative-adversary-networks-819a86b3750b

Problems with GANs

(mode collapse)

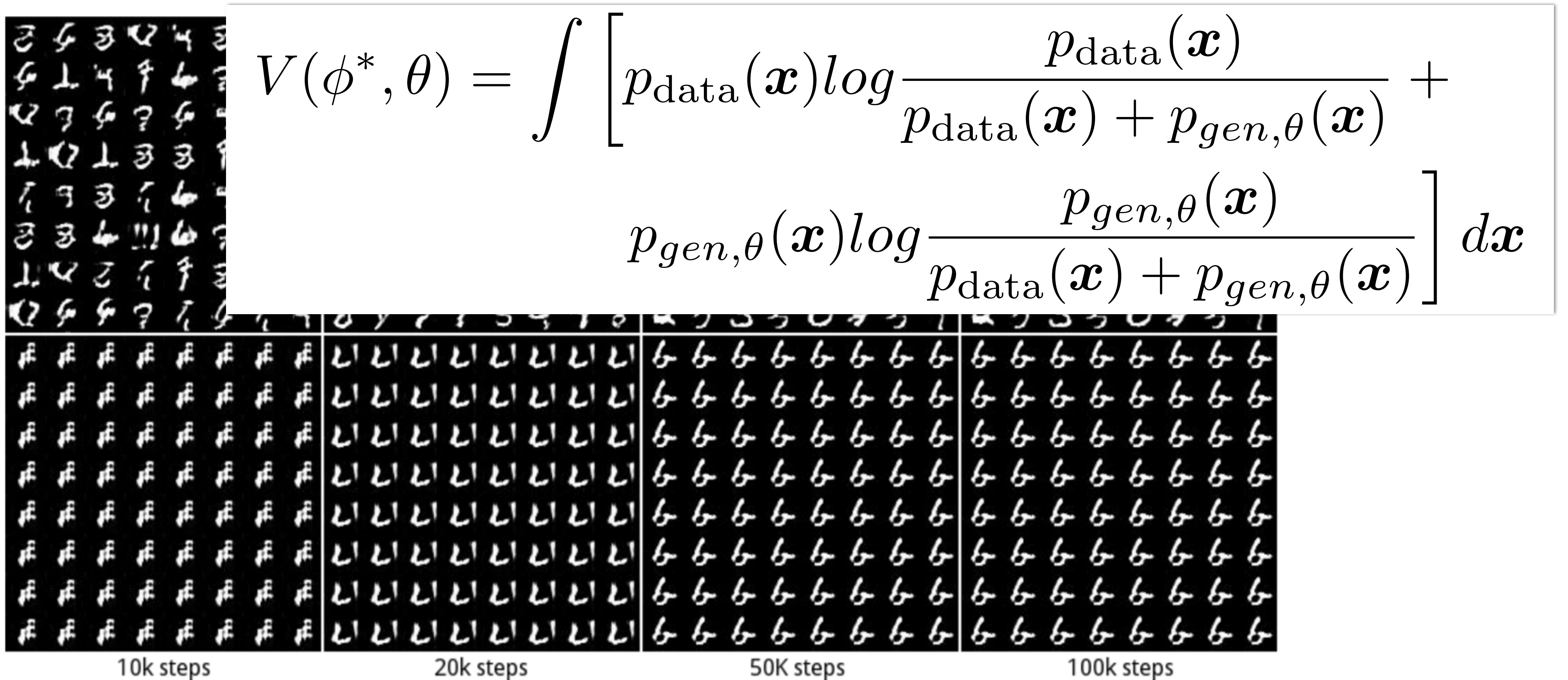


image taken from: https://medium.com/@jonathan_hui/gan-why-it-is-so-hard-to-train-generative-adversary-networks-819a86b3750b

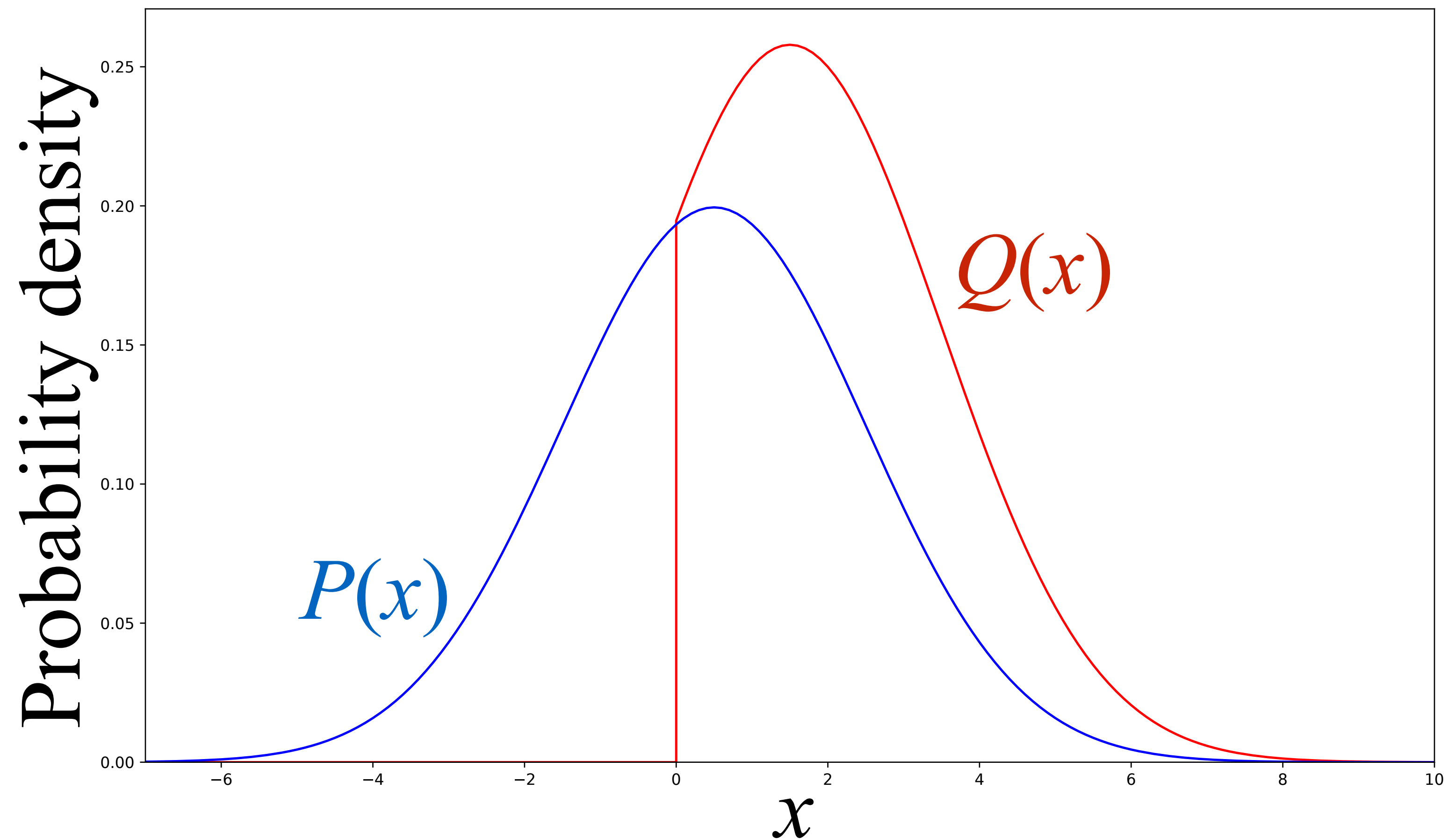
Wasserstein distance

Also called «Earth mover's distance» (EMD)

Wasserstein distance

Also called «Earth mover's distance» (EMD)

Distributions $P(x)$ and $Q(x)$ are viewed as describing the amounts of 'dirt' at point x

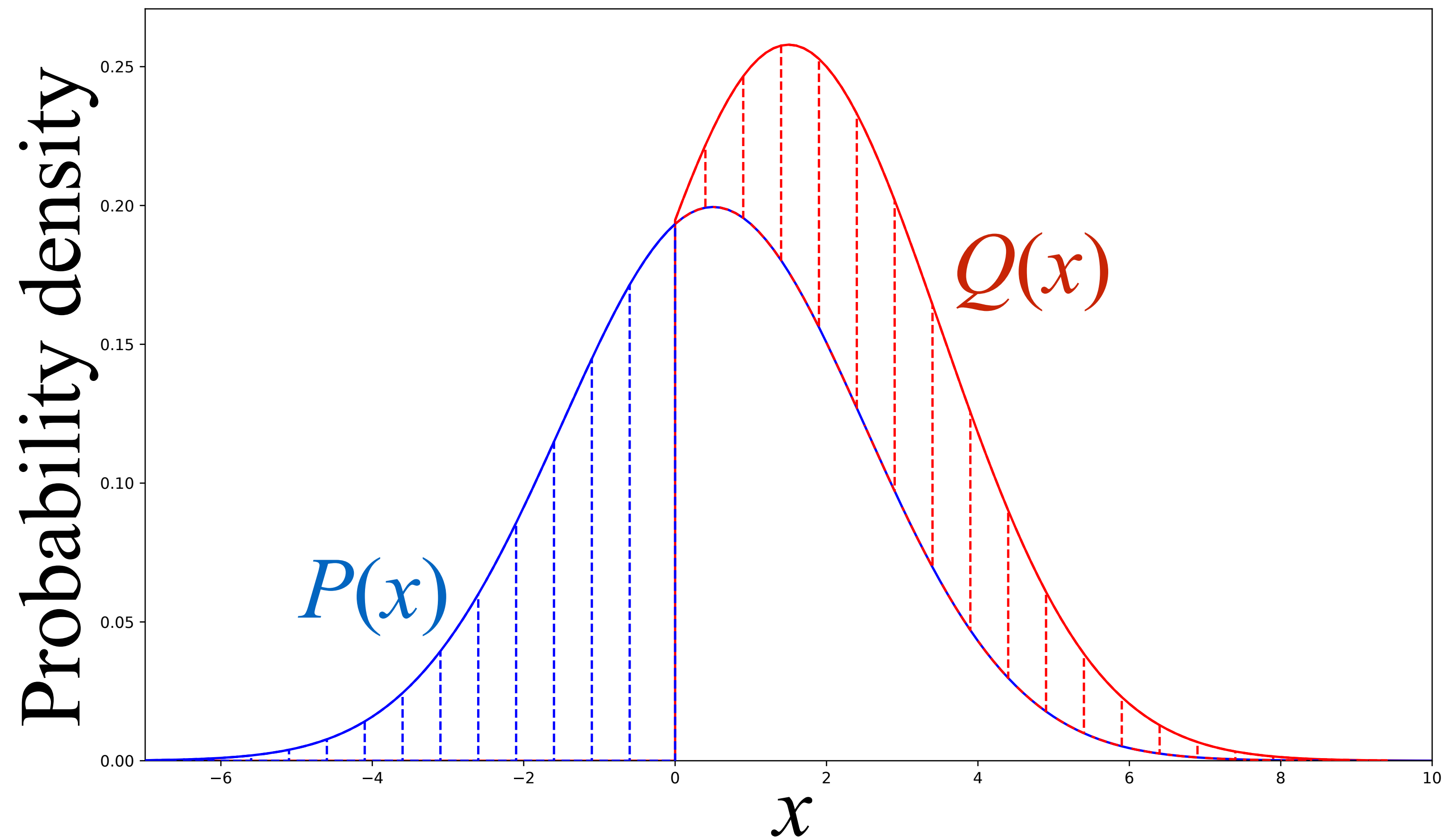


Wasserstein distance

Also called «Earth mover's distance» (EMD)

Distributions $P(x)$ and $Q(x)$ are viewed as describing the amounts of 'dirt' at point x

We want to convert one distribution into the other by moving around some amounts of dirt

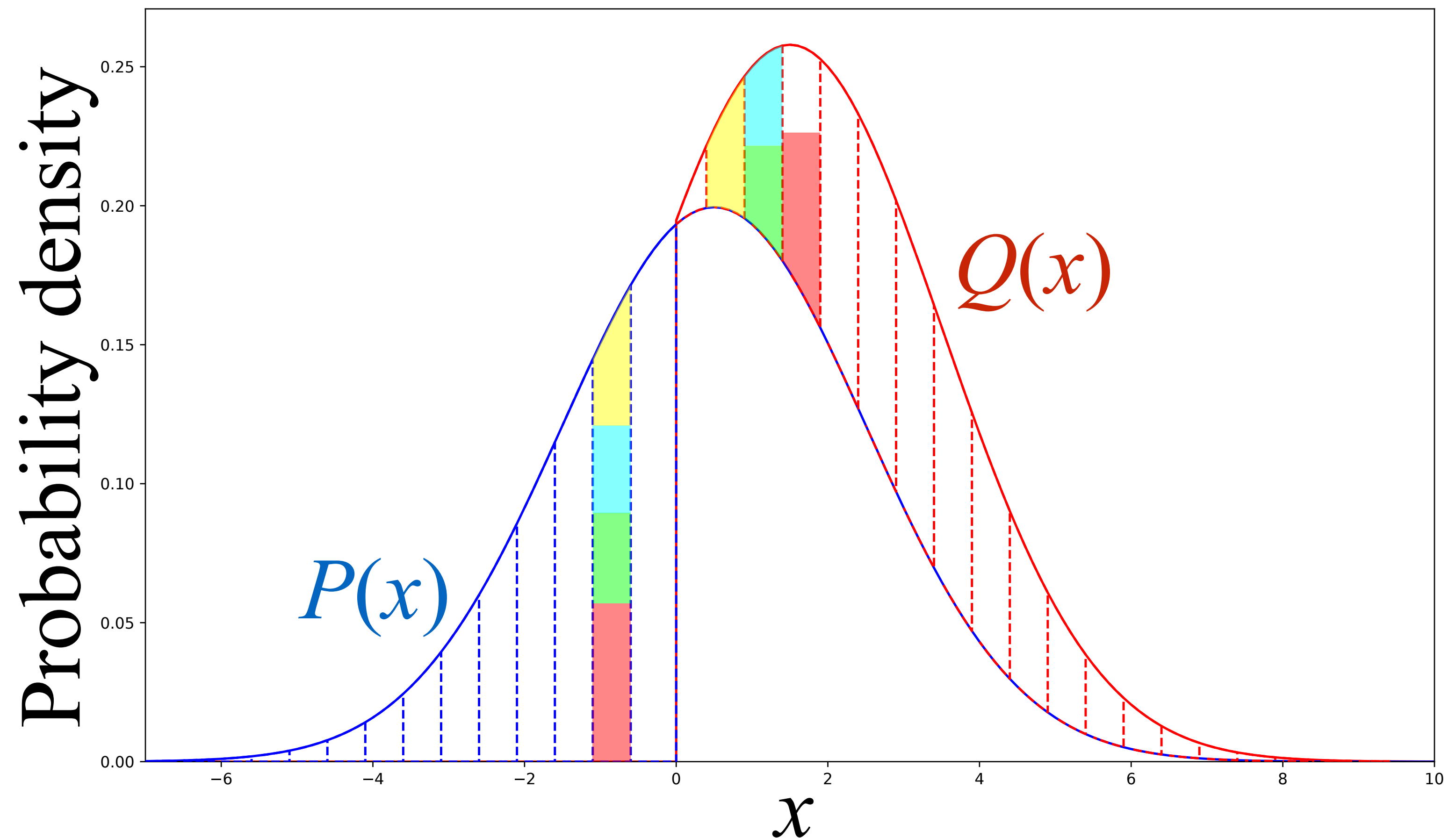


Wasserstein distance

Also called «Earth mover's distance» (EMD)

Distributions $P(x)$ and $Q(x)$ are viewed as describing the amounts of 'dirt' at point x

We want to convert one distribution into the other by moving around some amounts of dirt



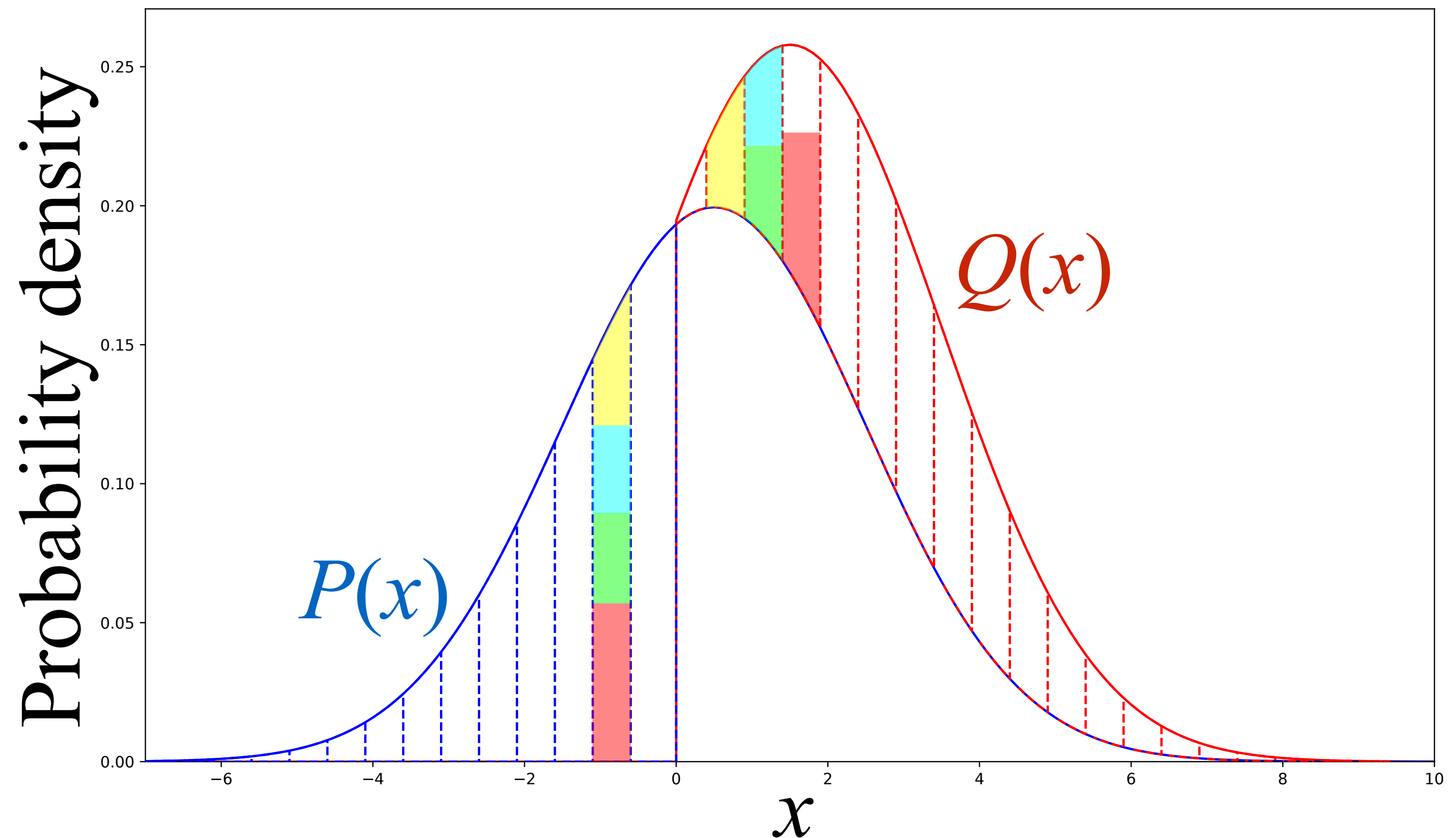
Wasserstein distance

Also called «Earth mover's distance» (EMD)

Distributions $P(x)$ and $Q(x)$ are viewed as describing the amounts of 'dirt' at point x

We want to convert one distribution into the other by moving around some amounts of dirt

The cost of moving an amount m from x_1 to x_2 is $m \times \|x_2 - x_1\|$



Wasserstein distance

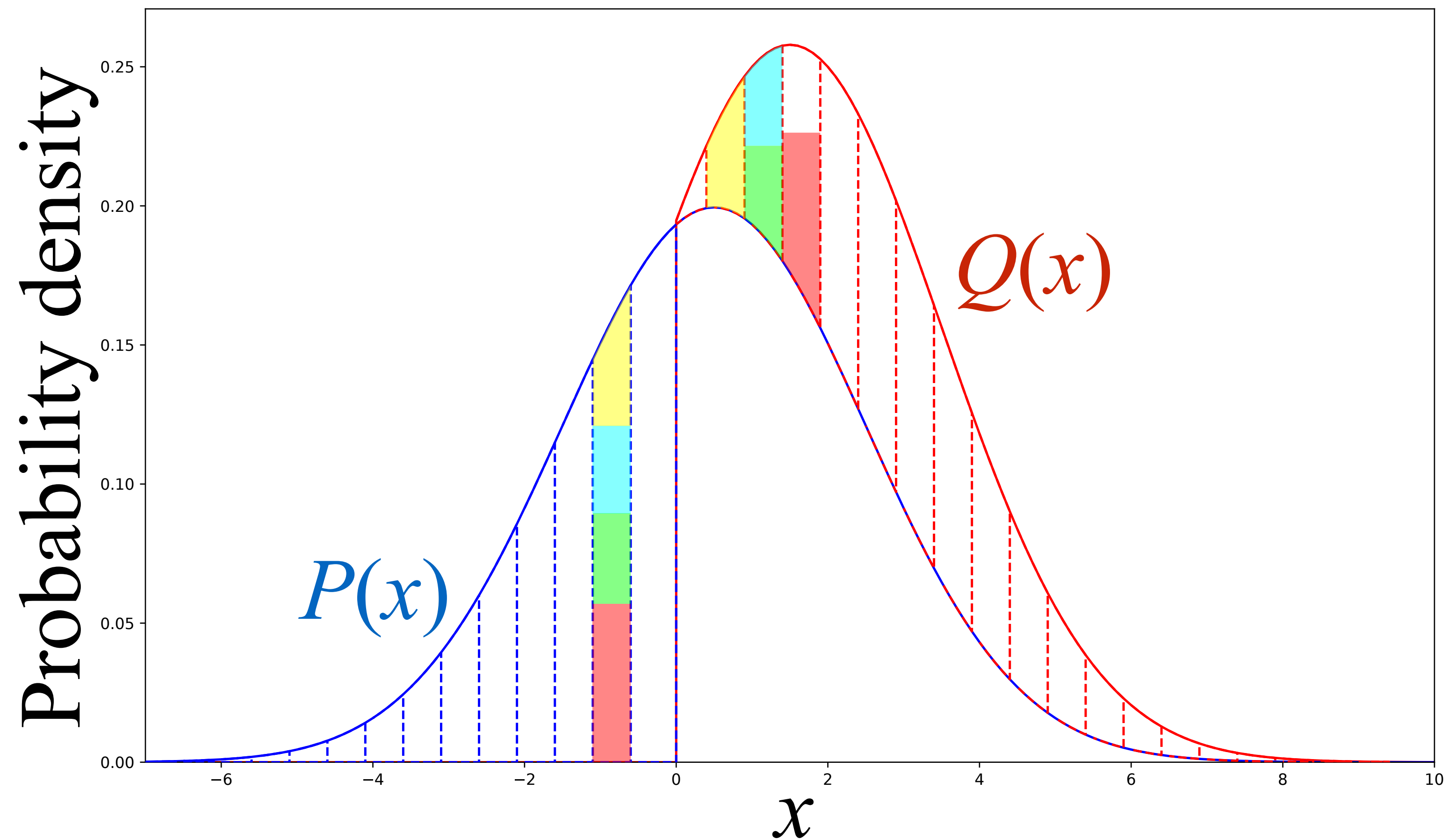
Also called «Earth mover's distance» (EMD)

Distributions $P(x)$ and $Q(x)$ are viewed as describing the amounts of 'dirt' at point x

We want to convert one distribution into the other by **moving around** some amounts of dirt

The cost of moving an amount m from x_1 to x_2 is $m \times \|x_2 - x_1\|$

$\text{EMD}(P, Q) =$ **minimum total cost** of converting P into Q



WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

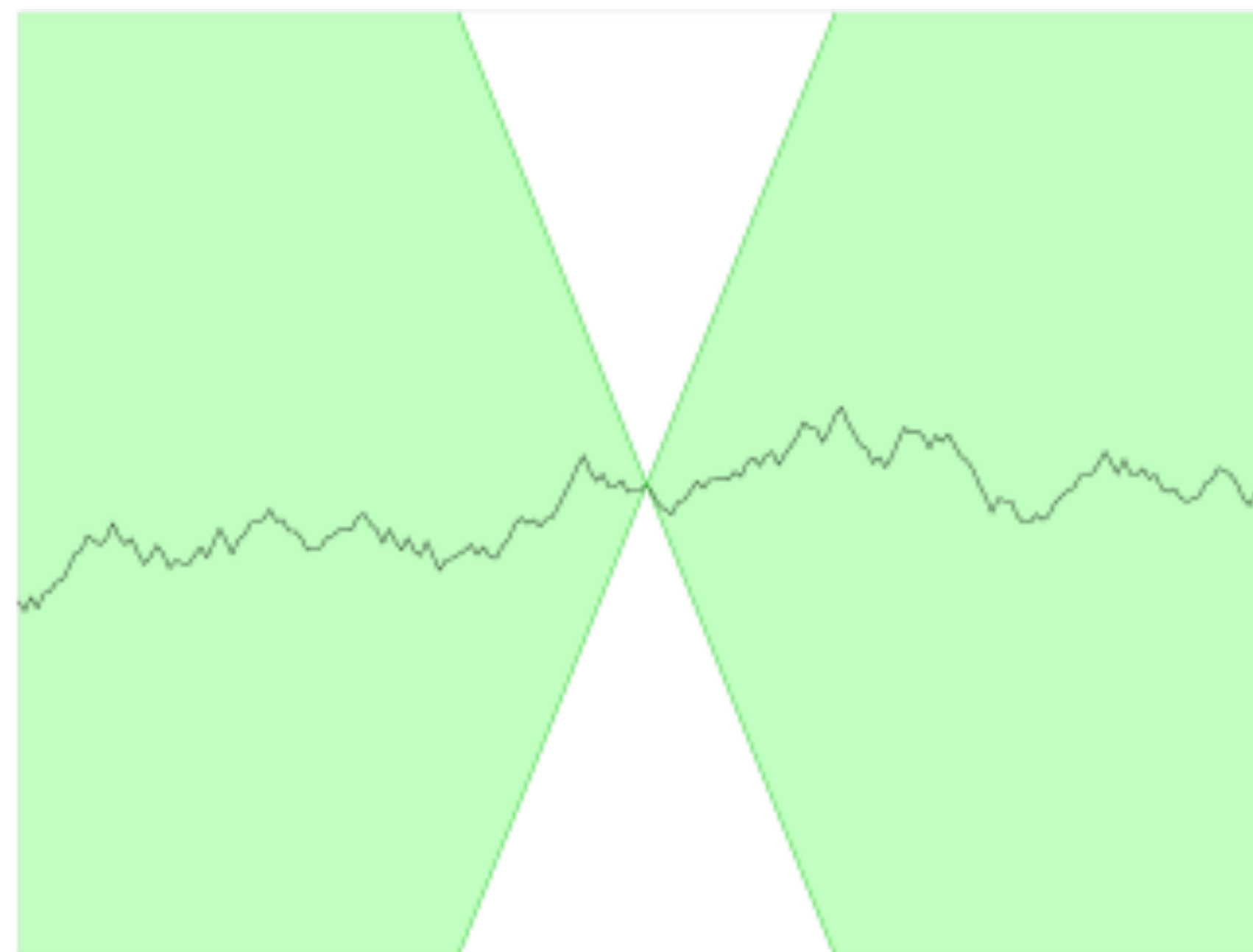
WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

$$\|f\|_L \leq 1$$



$$|f(a) - f(b)| \leq \|a - b\|, \quad \forall a, b$$



WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

<https://arxiv.org/abs/1701.07875>

WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

Expectations can be estimated as
sample mean (per batch)

<https://arxiv.org/abs/1701.07875>

WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

The function can be expressed with a neural net – discriminator ('critic' in the original paper)

Expectations can be estimated as sample mean (per batch)

<https://arxiv.org/abs/1701.07875>

WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

The function can be expressed with a neural net – discriminator ('critic' in the original paper)

Expectations can be estimated as sample mean (per batch)

This property can be replaced by a weaker one: $\|f\|_L \leq k : |f(a) - f(b)| \leq k \cdot \|a - b\|, \forall a, b$
In such case we'll estimate $k \times \text{EMD}(P, Q)$ instead of $\text{EMD}(P, Q)$

<https://arxiv.org/abs/1701.07875>

WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

The function can be expressed with a neural net – discriminator ('critic' in the original paper)

Expectations can be estimated as sample mean (per batch)

This property can be replaced by a weaker one: $\|f\|_L \leq k : |f(a) - f(b)| \leq k \cdot \|a - b\|, \forall a, b$
In such case we'll estimate $k \times \text{EMD}(P, Q)$ instead of $\text{EMD}(P, Q)$

This can be achieved by clipping critic's weights at some value: $w \rightarrow \text{clip}(w, -c, c)$

<https://arxiv.org/abs/1701.07875>

WGAN

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

The function can be expressed with a neural net – discriminator ('critic' in the original paper)

Expectations can be estimated as sample mean (per batch)

This property can be replaced by a

weaker one: $\|f\|_L \leq k : |f(a) - f(b)| \leq k \cdot \|a - b\|, \forall a, b$

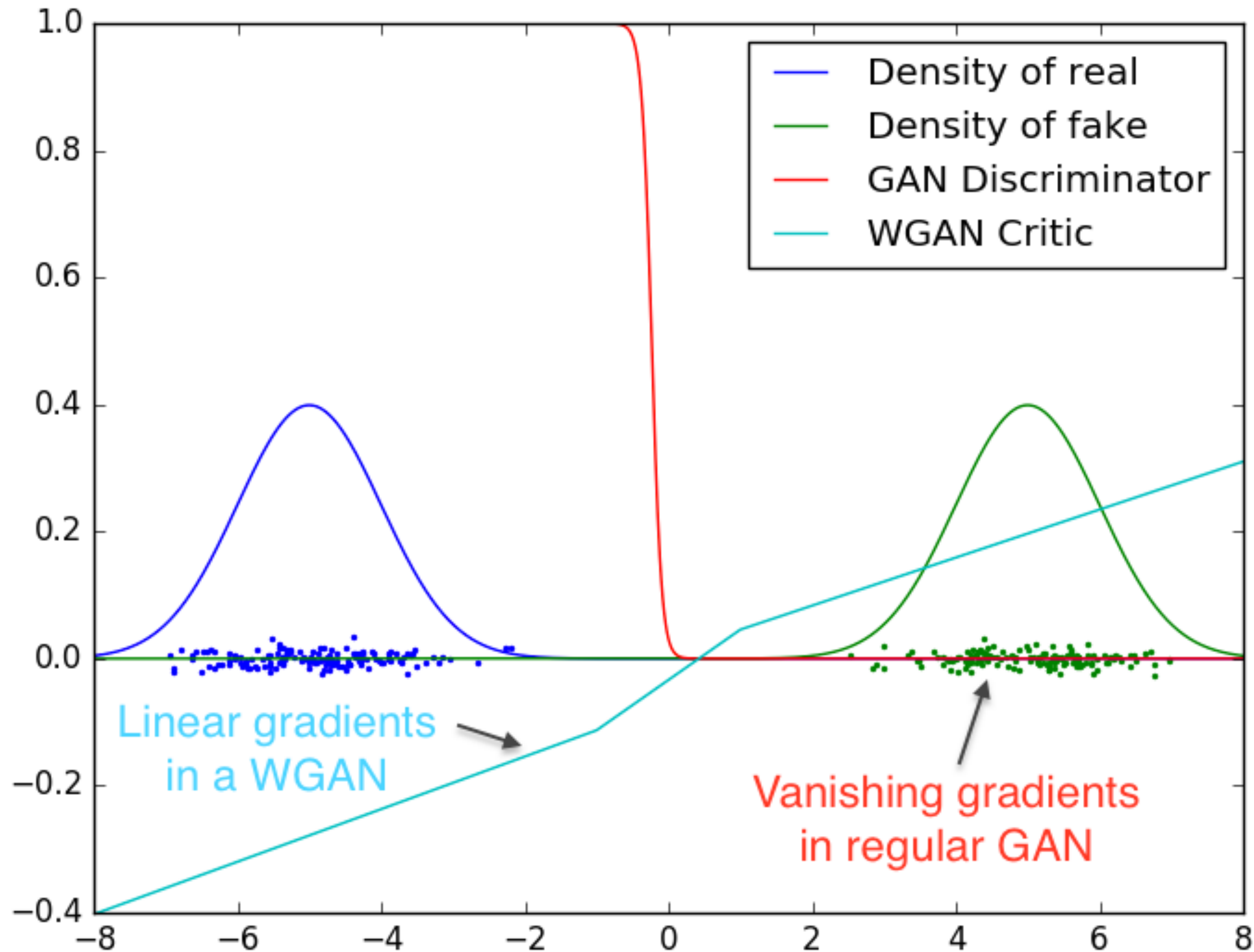
In such case we'll estimate $k \times \text{EMD}(P, Q)$ instead of $\text{EMD}(P, Q)$

This can be achieved by clipping critic's weights at some value: $w \rightarrow \text{clip}(w, -c, c)$

We wouldn't know what k is, but it doesn't matter: all we want is to **minimize** the EMD!

<https://arxiv.org/abs/1701.07875>

WGAN



<https://arxiv.org/abs/1701.07875>

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

**Weight clipping makes critic less expressive
and training process harder to converge**

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

Weight clipping makes critic less expressive
and training process harder to converge

Optimal f should satisfy $\|\nabla f\| = 1$ almost
everywhere under P and Q

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

Weight clipping makes critic less expressive
and training process harder to converge

Optimal f should satisfy $\|\nabla f\| = 1$ almost
everywhere under P and Q

$$\|f\|_L \leq 1 \quad + \quad \iff \quad \|\nabla f\| \leq 1$$

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

Weight clipping makes critic less expressive
and training process harder to converge

Optimal f should satisfy $\|\nabla f\| = 1$ almost
everywhere under P and Q

$$\|f\|_L \leq 1 \quad + \quad \iff \quad \|\nabla f\| \leq 1$$

Weight clipping can be replaced by
penalizing the gradients to be of unit
magnitude

WGAN-GP

«Gradient penalty»

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

Weight clipping makes critic less expressive
and training process harder to converge

Optimal f should satisfy $\|\nabla f\| = 1$ almost
everywhere under P and Q

$$\|f\|_L \leq 1 \quad + \quad \iff \quad \|\nabla f\| \leq 1$$

Weight clipping can be replaced by
penalizing the gradients to be of unit
magnitude

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

Weight clipping makes critic less expressive and training process harder to converge

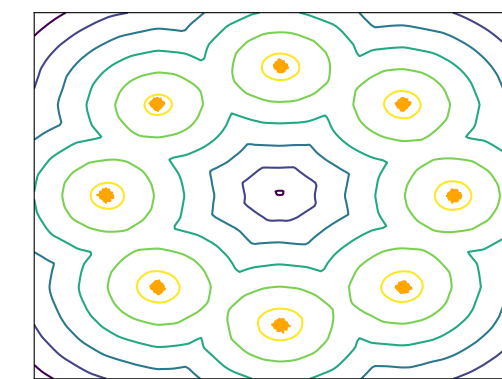
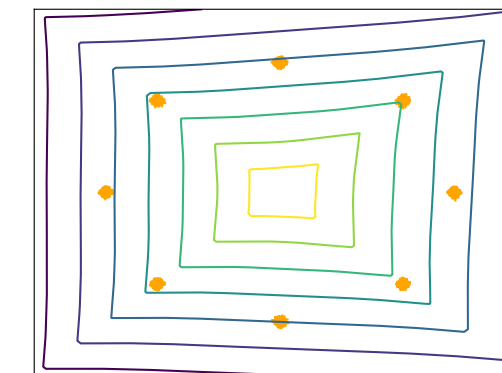
Optimal f should satisfy $\|\nabla f\| = 1$ almost everywhere under P and Q

$$\|f\|_L \leq 1 \quad + \quad \iff \quad \|\nabla f\| \leq 1$$

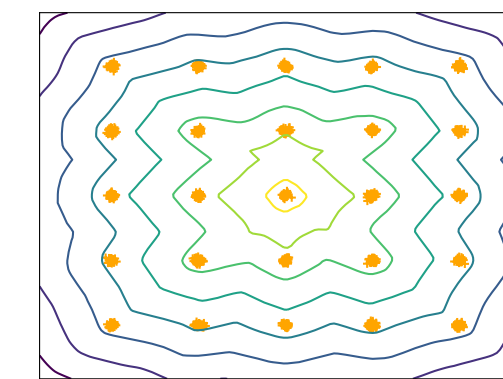
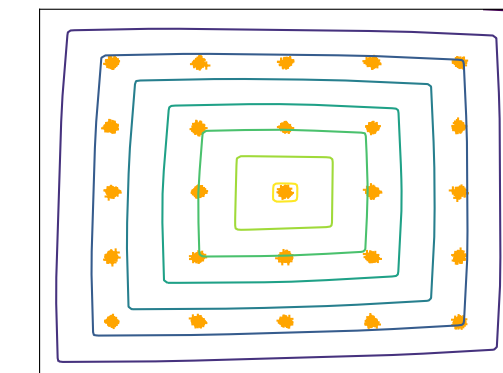
Weight clipping can be replaced by penalizing the gradients to be of unit magnitude

«Gradient penalty»

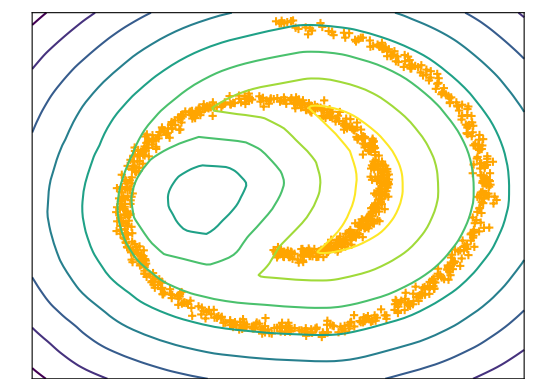
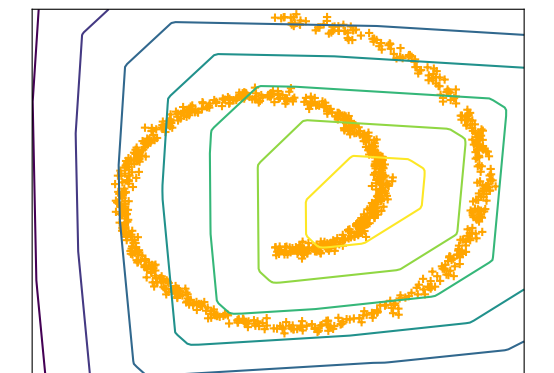
8 Gaussians



25 Gaussians



Swiss Roll



Value surfaces of f

WGAN
WGAN-GP

WGAN-GP

Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

Weight clipping makes critic less expressive and training process harder to converge

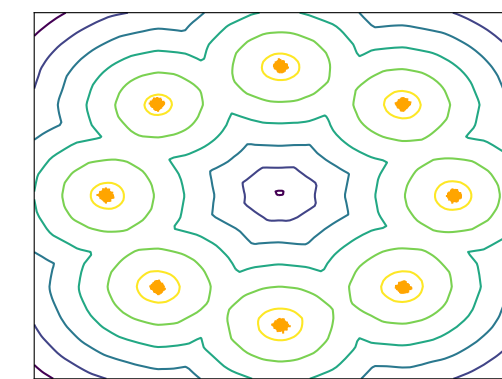
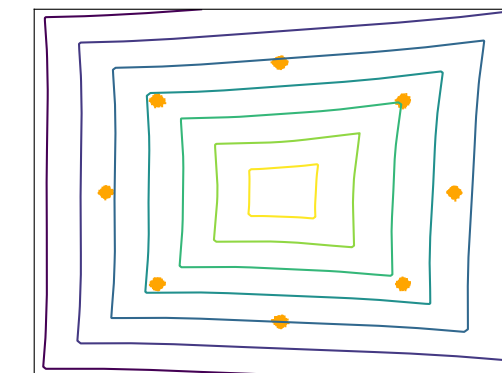
Optimal f should satisfy $\|\nabla f\| = 1$ almost everywhere under P and Q

$$\|f\|_L \leq 1 \quad + \quad \iff \quad \|\nabla f\| \leq 1$$

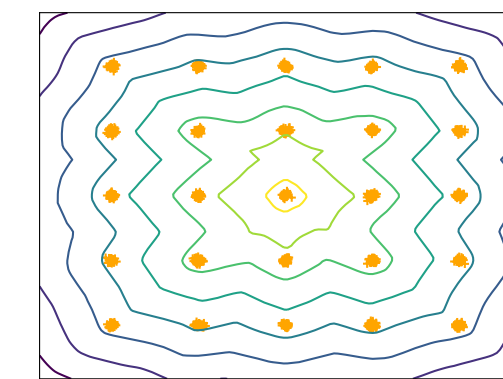
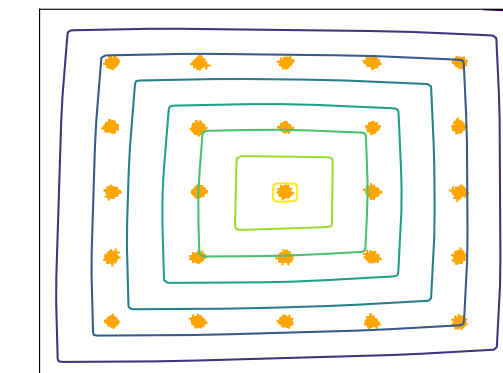
Weight clipping can be replaced by penalizing the gradients to be of unit magnitude

«Gradient penalty»

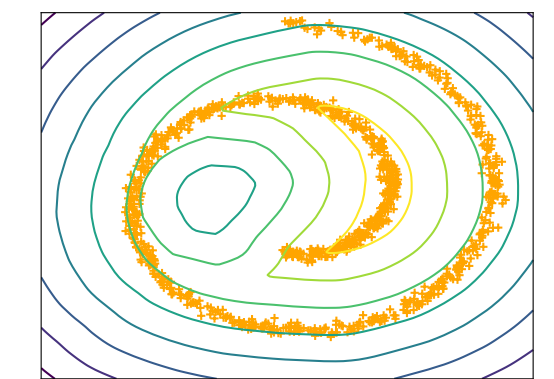
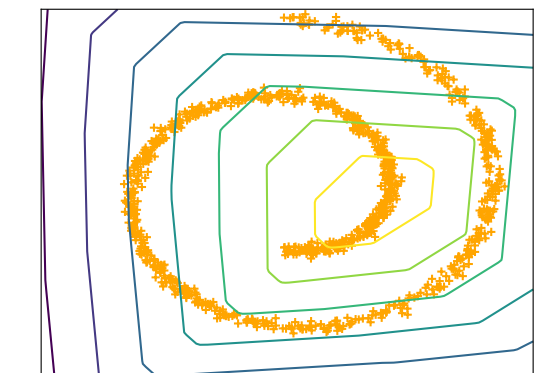
8 Gaussians



25 Gaussians



Swiss Roll



WGAN
WGAN-GP

Value surfaces of f

$$\text{GP} = \lambda \mathbb{E}_{\tilde{x} \sim \mathbb{P}_{\tilde{x}}} \left[\left(\|\nabla_{\tilde{x}} f(\tilde{x})\| - 1 \right)^2 \right]$$

$$\mathbb{P}_{\tilde{x}} : \begin{bmatrix} \tilde{x} = \alpha x_1 + (1-\alpha)x_2 \\ \alpha \sim \text{Uniform}(0,1) \\ x_1 \sim P \\ x_2 \sim Q \end{bmatrix}$$

WGAN-GP

«Gradient penalty»

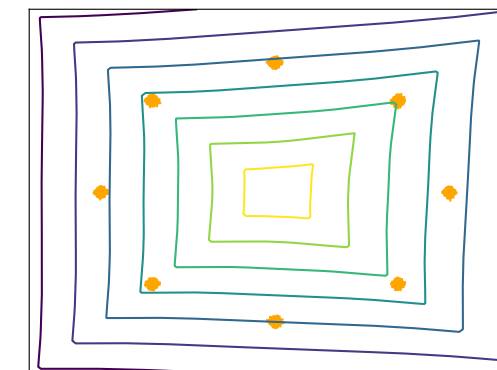
Improved Training of Wasserstein GANs

<https://arxiv.org/abs/1704.00028>

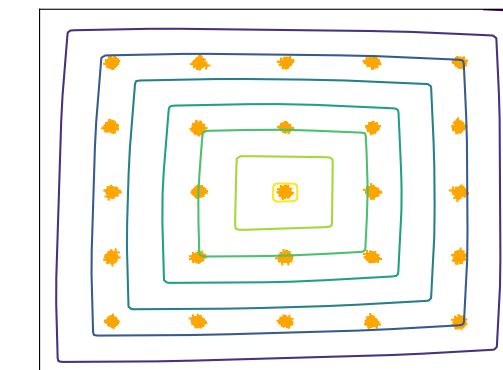
Weight clipping makes critic less expressive and training process harder to converge

Optimal f should satisfy $\|\nabla f\| = 1$ almost everywhere under P and Q

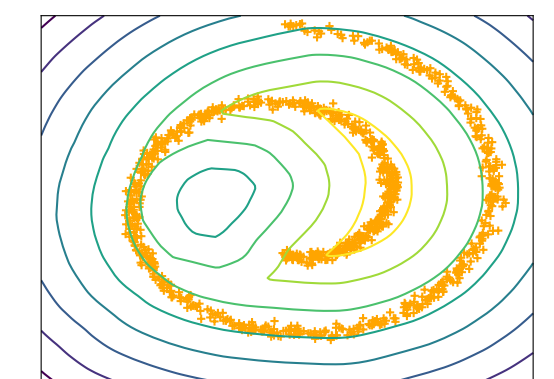
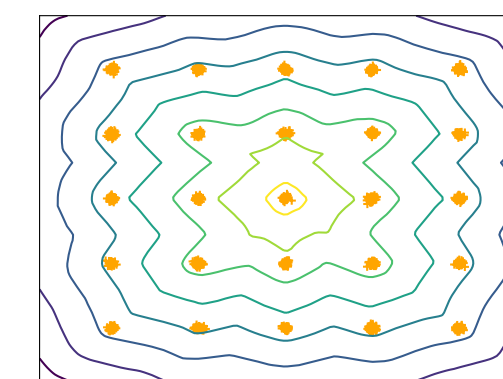
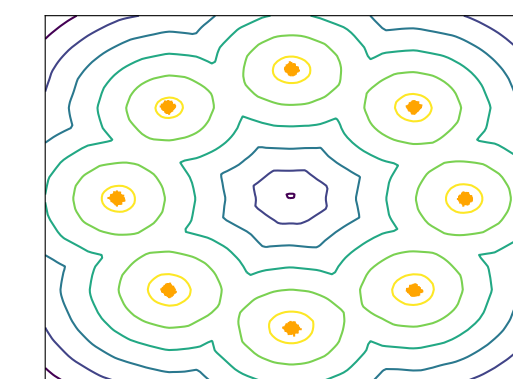
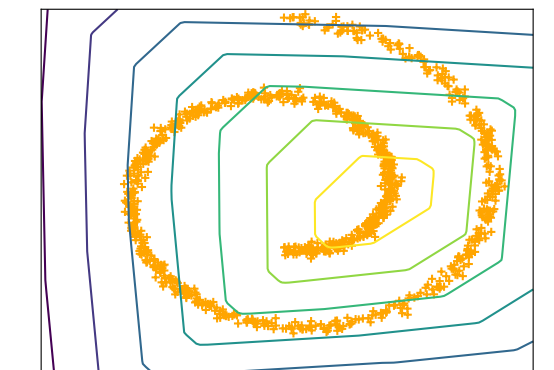
8 Gaussians



25 Gaussians



Swiss Roll



WGAN
WGAN-GP

Value surfaces of f

$$\|f\|_L \leq 1 \iff \|\nabla f\| \leq 1$$

W
pe
m

Alternative 'one-sided' penalty:

$$\text{GP} = \lambda \mathbb{E}_{\tilde{x} \sim P_{\tilde{x}}} [\max(0, \|\nabla_{\tilde{x}} f(\tilde{x})\| - 1)^2]$$

$$\text{GP} = \lambda \mathbb{E}_{\tilde{x} \sim P_{\tilde{x}}} [(\|\nabla_{\tilde{x}} f(\tilde{x})\| - 1)^2]$$

$$P_{\tilde{x}} : \begin{cases} \tilde{x} = \alpha x_1 + (1-\alpha)x_2 \\ \alpha \sim \text{Uniform}(0,1) \\ x_1 \sim P \\ x_2 \sim Q \end{cases}$$

WGAN-GP

Improved Training of Wasserstein GANs
<https://arxiv.org/abs/1704.00028>

DCGAN

LSGAN

WGAN (clipping)

WGAN-GP (ours)

Baseline (G : DCGAN, D : DCGAN)



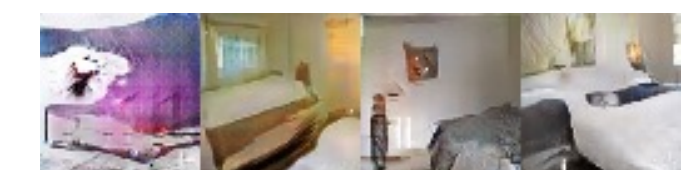
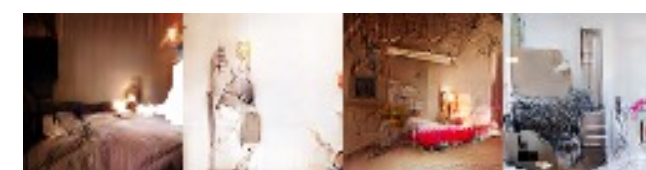
G : No BN and a constant number of filters, D : DCGAN



G : 4-layer 512-dim ReLU MLP, D : DCGAN



No normalization in either G or D



Gated multiplicative nonlinearities everywhere in G and D



tanh nonlinearities everywhere in G and D



101-layer ResNet G and D



WGAN-GP

Improved Training of Wasserstein GANs
<https://arxiv.org/abs/1704.00028>

DCGAN	LSGAN	WGAN (clipping)	WGAN-GP (ours)
Baseline (G : DCGAN, D : DCGAN)			
G : No BN and a constant number of filters, D : DCGAN			
G : 4-layer 512-dim ReLU MLP, D : DCGAN			
No normalization in either G or D			
Gated multiplicative nonlinearities everywhere in G and D			
tanh nonlinearities everywhere in G and D			
101-layer ResNet G and D			

This technique allowed for very deep networks to be used for GANs

Conditional distributions

Sometimes it's necessary to learn not just $P(\mathbf{x})$, but $P(\mathbf{x}|a)$

Conditional distributions

~~Sometimes~~ it's necessary to learn not just $P(\mathbf{x})$, but $P(\mathbf{x}|a)$
Most of the times

Conditional distributions

~~Sometimes~~ it's necessary to learn not just $P(x)$, but $P(x|a)$

Most of the times

E.g., stochastic detector output y for given particle parameters x

Conditional distributions

~~Sometimes~~ it's necessary to learn not just $P(\mathbf{x})$, but $P(\mathbf{x}|\mathbf{a})$

Most of the times

E.g., stochastic detector output y for given particle parameters x

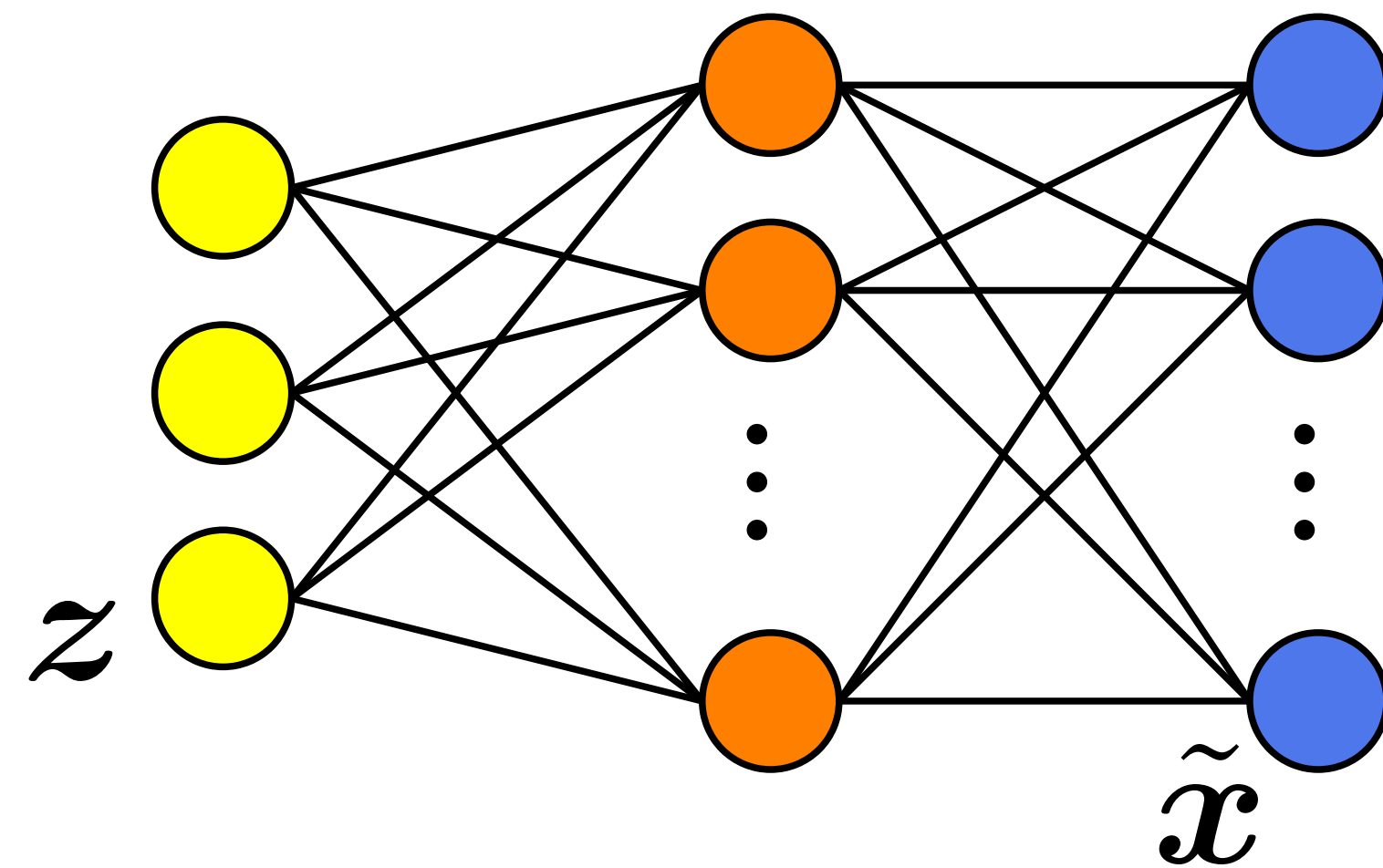
This can be achieved with any GAN architecture:

$$\begin{aligned}\tilde{\mathbf{x}}_j = G_\theta(\mathbf{z}_j) &\quad \rightarrow \quad \tilde{\mathbf{x}}_j = G_\theta(\mathbf{z}_j, \mathbf{a}_j) \\ D_\phi = D_\phi(\mathbf{x}_i) &\quad \rightarrow \quad D_\phi = D_\phi(\mathbf{x}_i, \mathbf{a}_i) \\ \tilde{D}_\phi = D_\phi(\tilde{\mathbf{x}}_j) &\quad \rightarrow \quad \tilde{D}_\phi = D_\phi(\tilde{\mathbf{x}}_j, \mathbf{a}_j)\end{aligned}$$

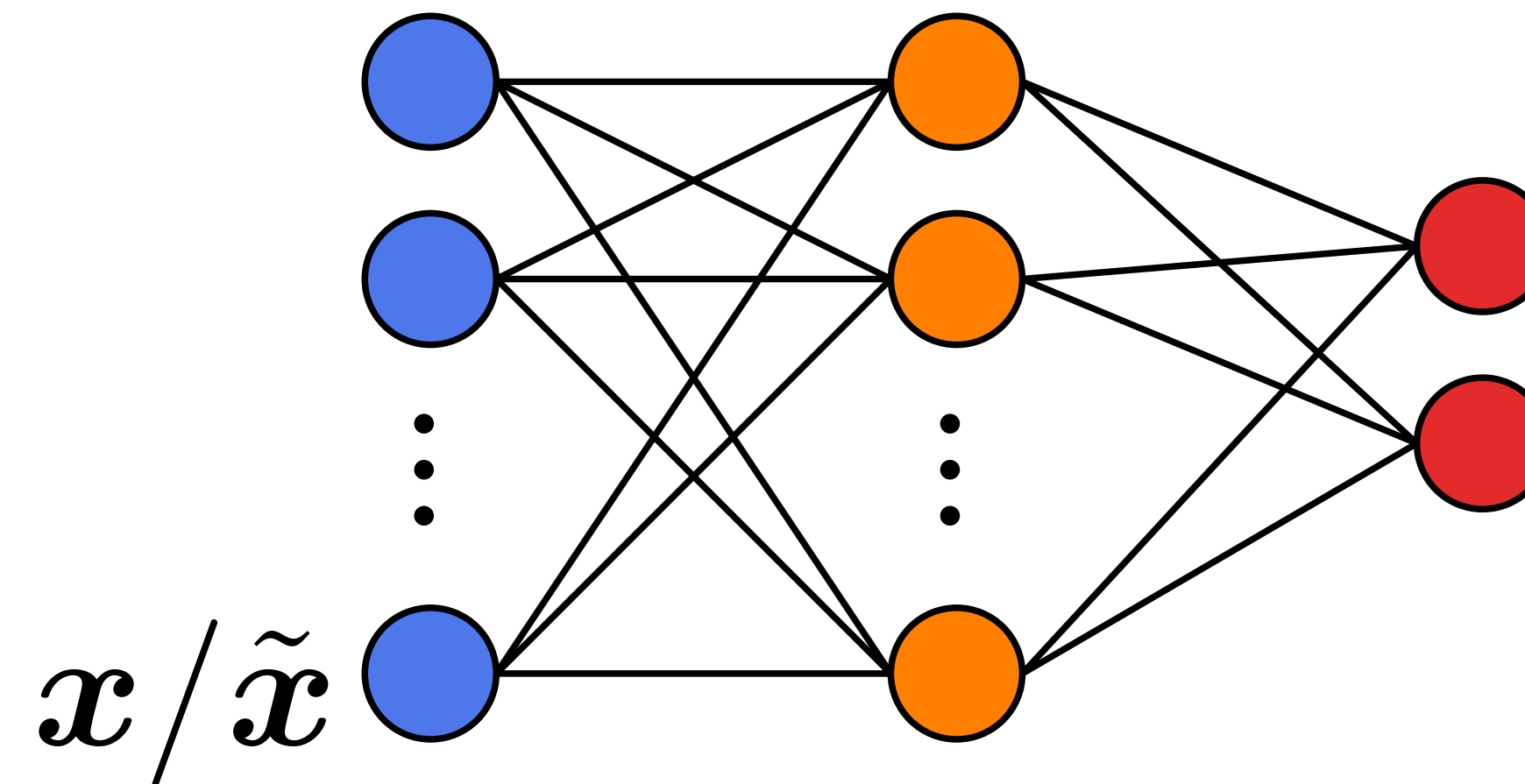
Conditional distributions

(simple case — fully-connected layers)

Generator



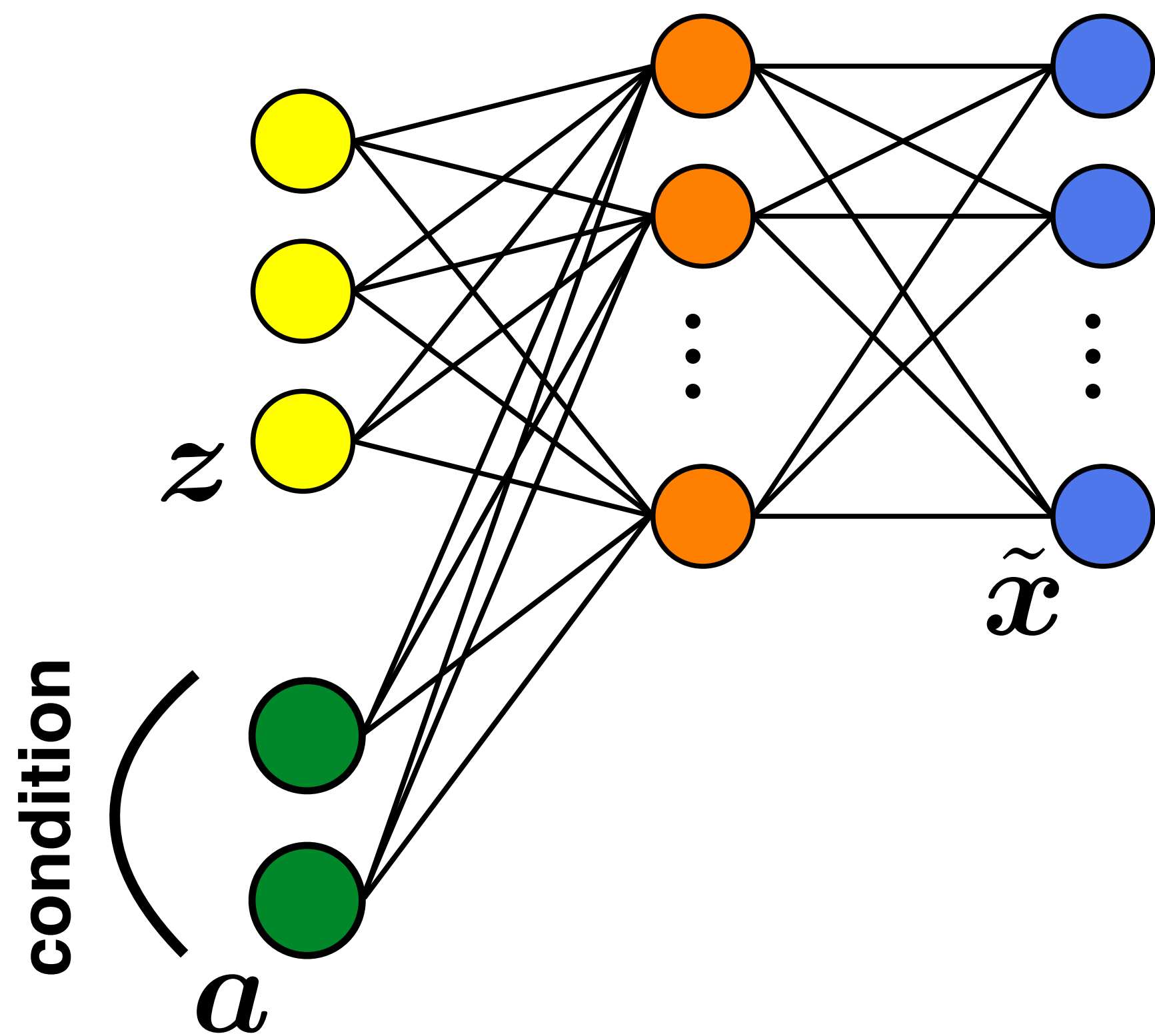
Discriminator



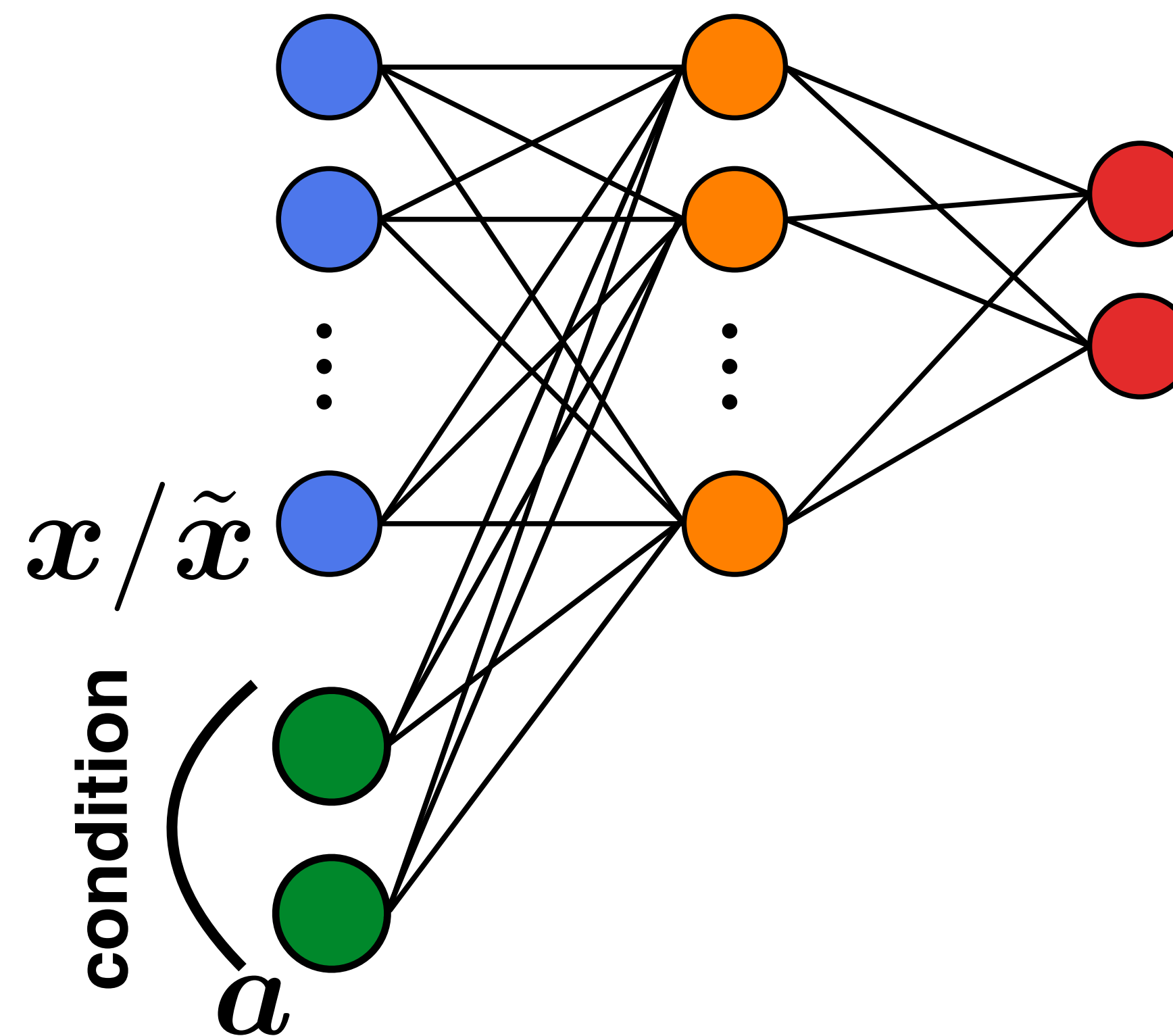
Conditional distributions

(simple case — fully-connected layers)

Generator



Discriminator



Fast Data-Driven Simulation of Cherenkov Detectors Using Generative Adversarial Networks

Artem Maevskiy, Denis Derkach, Nikita Kazeev, Andrey Ustyuzhanin, Maksim Artemev, Lucio Anderlini

(Submitted on 28 May 2019 (v1), last revised 26 Sep 2019 (this version, v2))

The increasing luminosities of future Large Hadron Collider runs and next generation of collider experiments will require an unprecedented amount of simulated events to be produced. Such large scale productions are extremely demanding in terms of computing resources. Thus new approaches to event generation and simulation of detector responses are needed. In LHCb, the accurate simulation of Cherenkov detectors takes a sizeable fraction of CPU time. An alternative approach is described here, when one generates high-level reconstructed observables using a generative neural network to bypass low level details. This network is trained to reproduce the particle species likelihood function values based on the track kinematic parameters and detector occupancy. The fast simulation is trained using real data samples collected by LHCb during run 2. We demonstrate that this approach provides high-fidelity results.

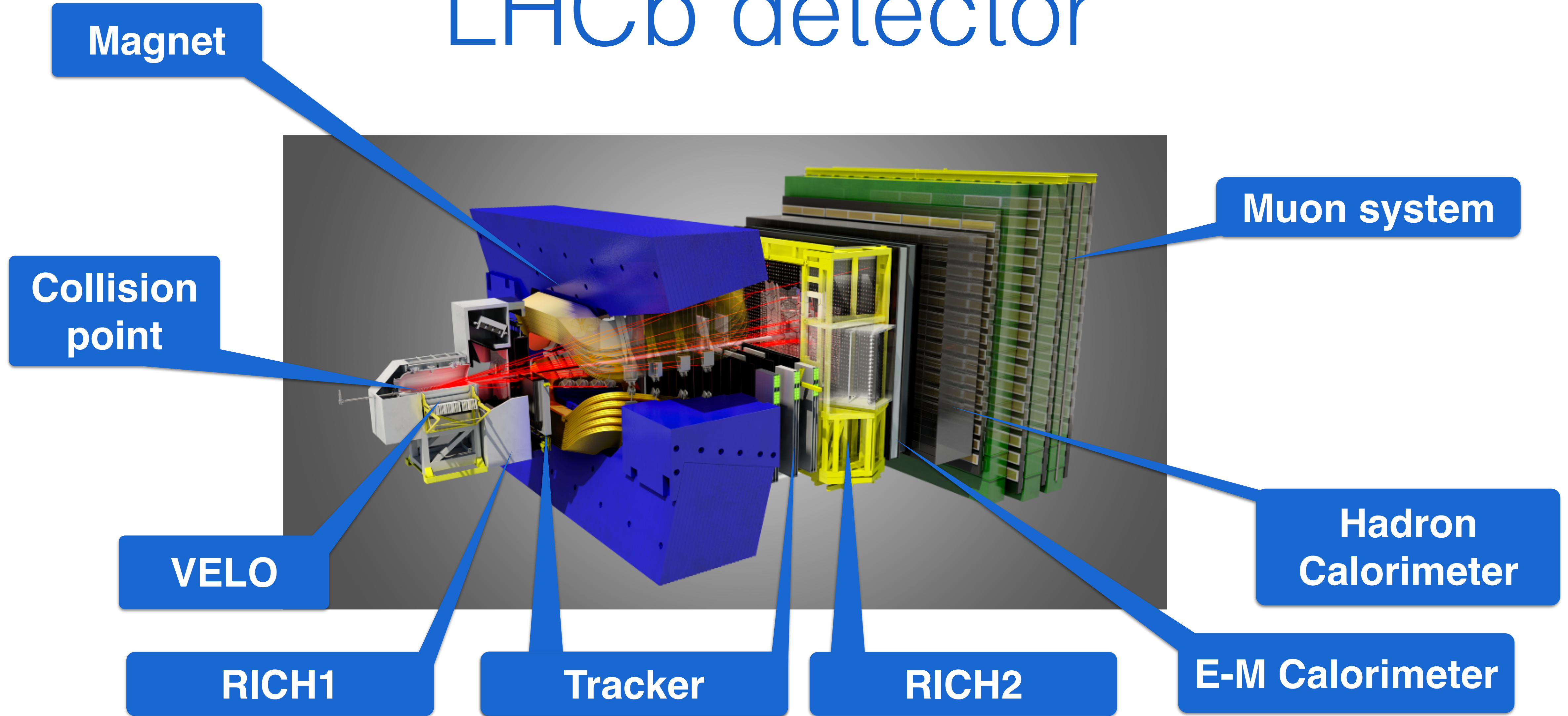
Comments: Proceedings for 19th International Workshop on Advanced Computing and Analysis Techniques in Physics Research. (Fixed typos and added one missing reference in the revised version.)

Subjects: **Instrumentation and Detectors (physics.ins-det)**; Machine Learning (cs.LG); High Energy Physics – Experiment (hep-ex)

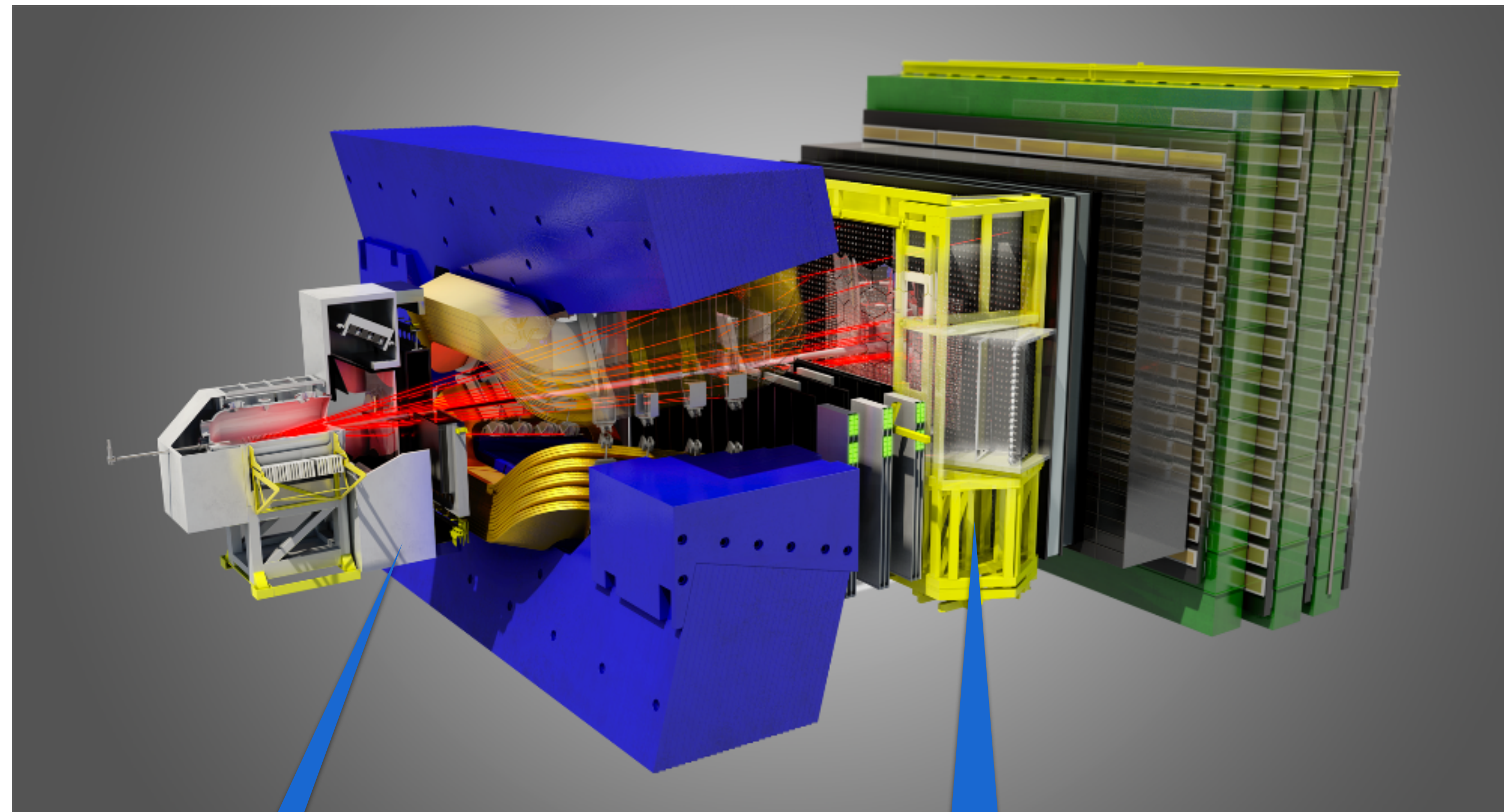
Cite as: [arXiv:1905.11825](https://arxiv.org/abs/1905.11825) [physics.ins-det]

(or [arXiv:1905.11825v2](https://arxiv.org/abs/1905.11825v2) [physics.ins-det] for this version)

LHCb detector



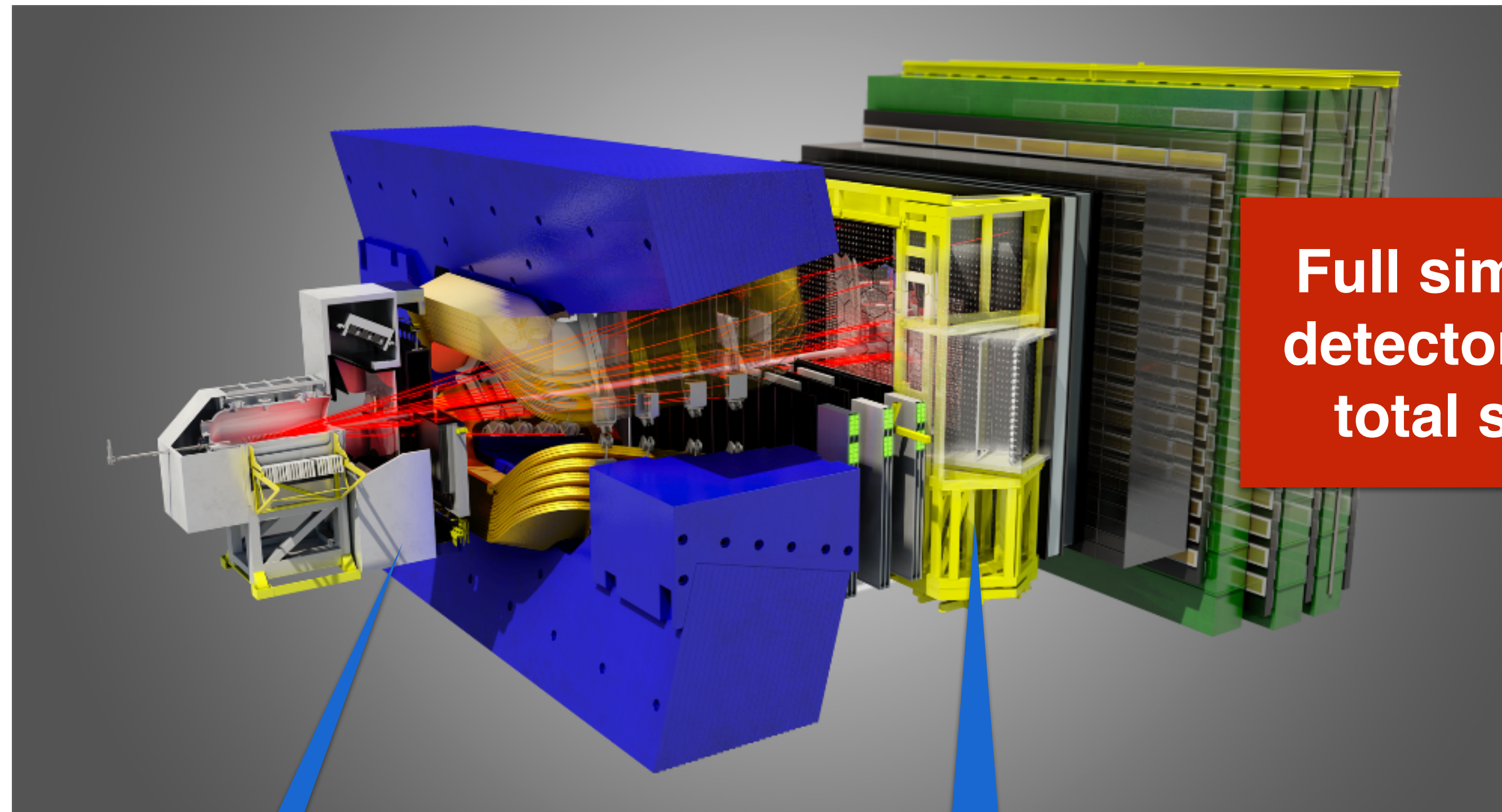
LHCb detector



RICH1

RICH2

LHCb detector

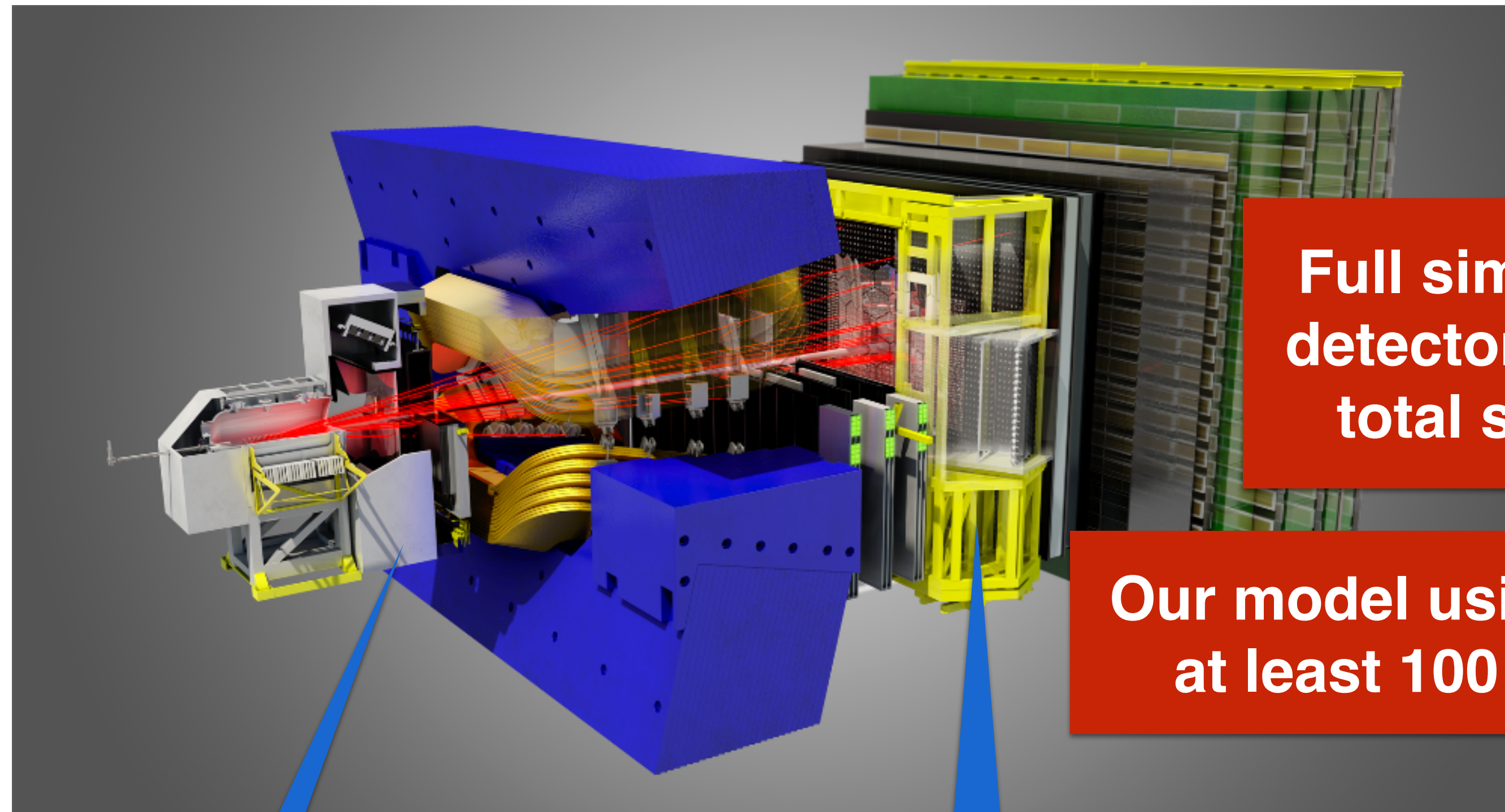


Full simulation of RICH detectors takes ~30% of total simulation time

RICH1

RICH2

LHCb detector



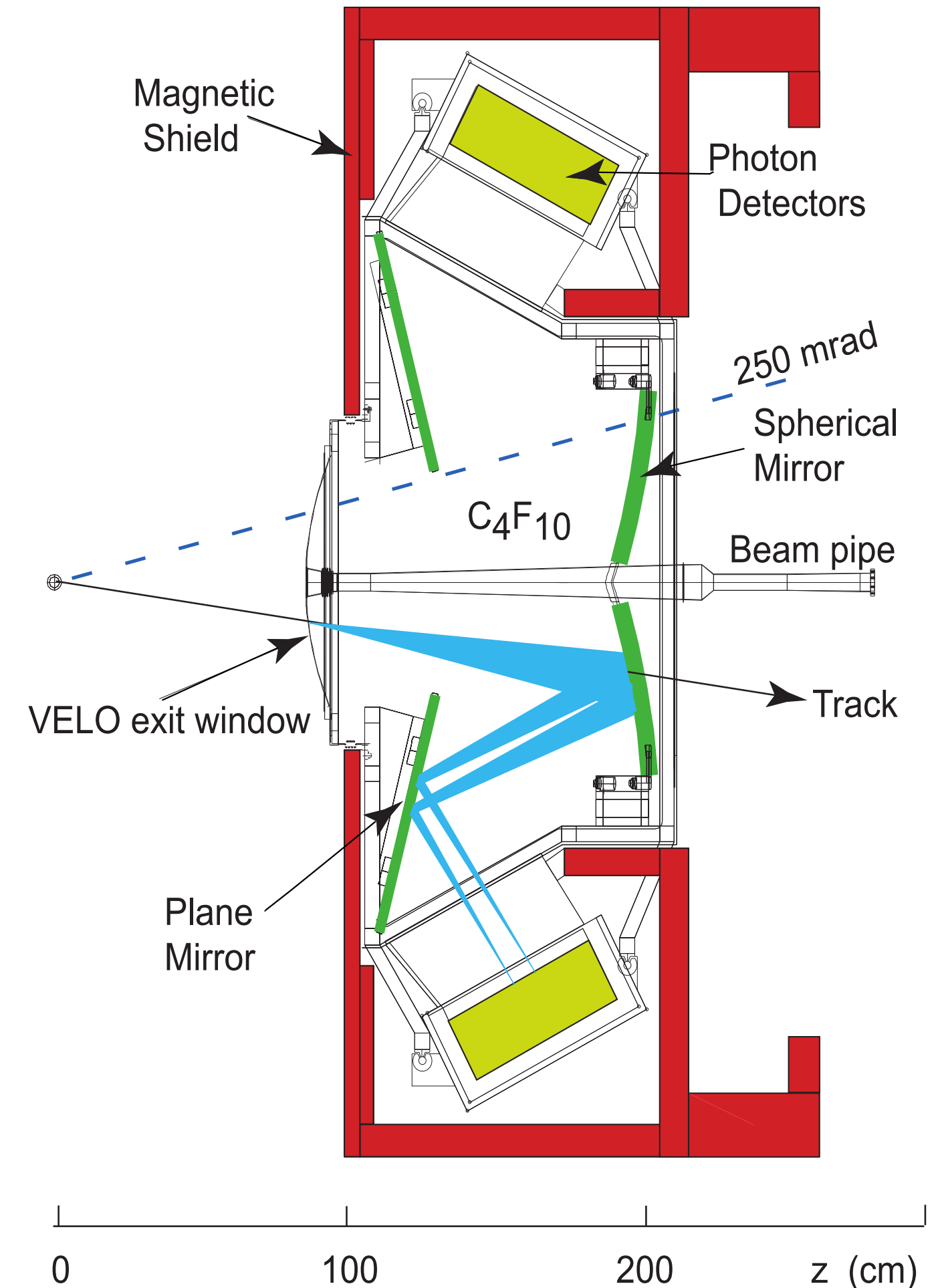
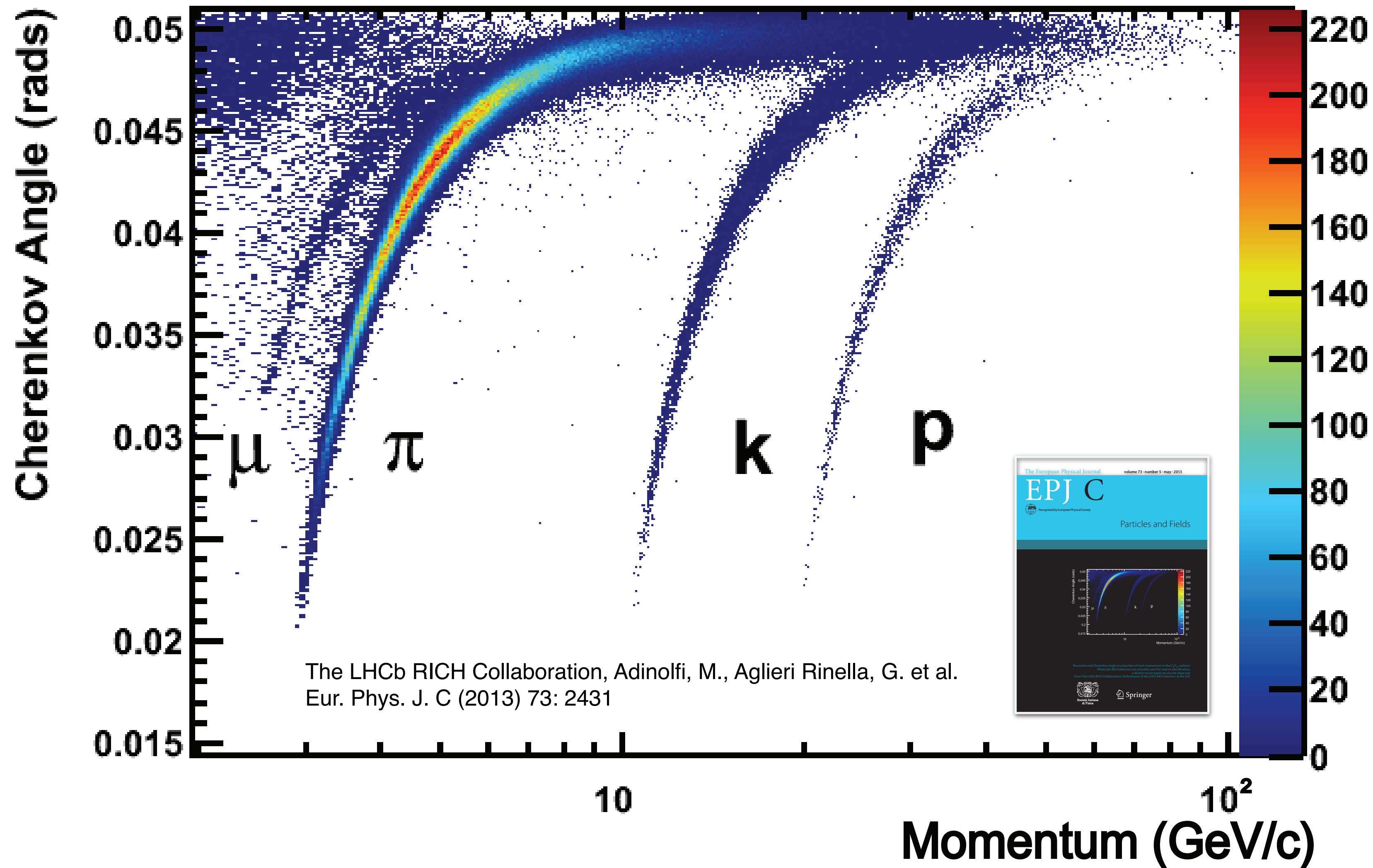
Full simulation of RICH detectors takes ~30% of total simulation time

Our model using GANs runs at least 100 times faster

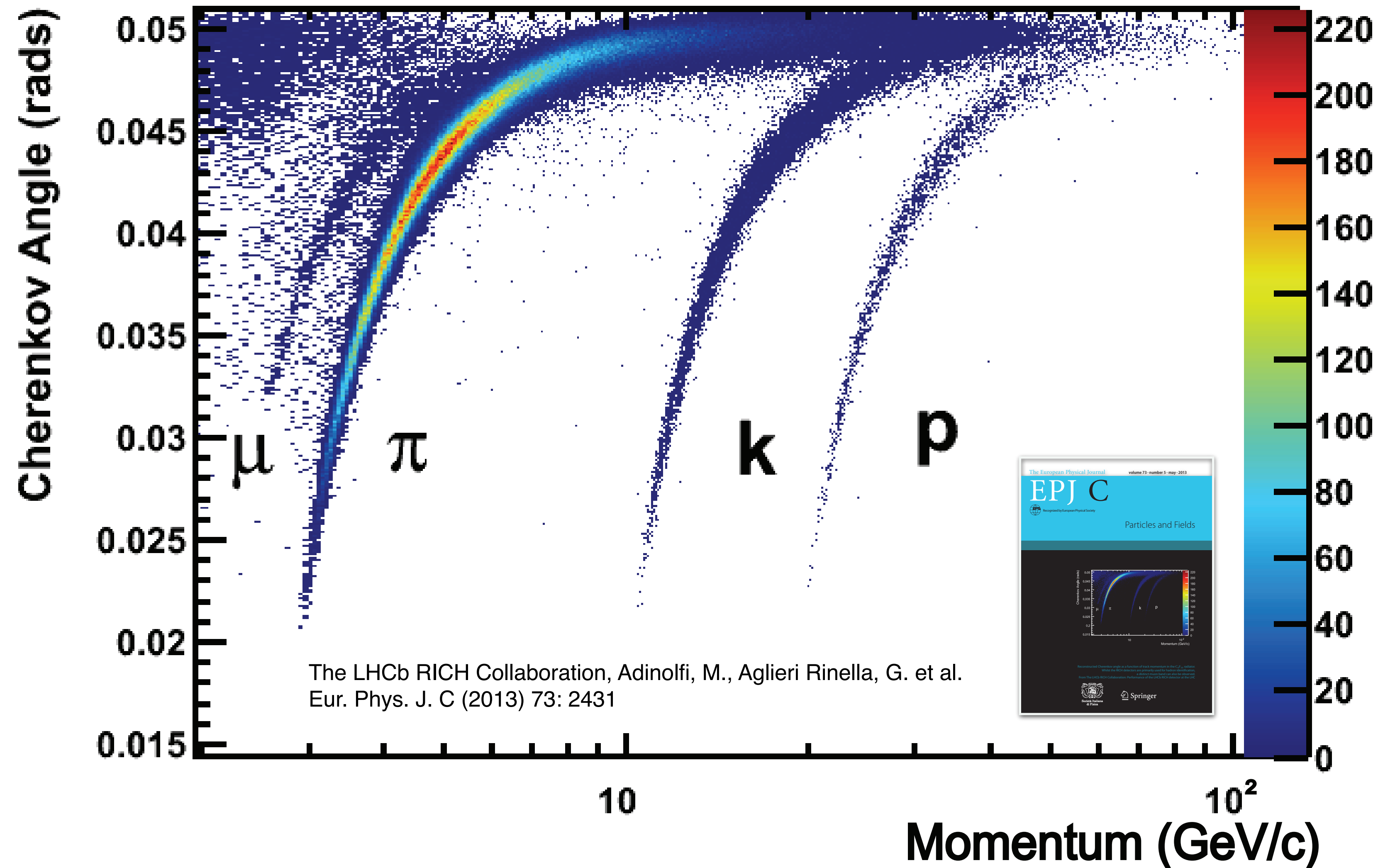
RICH1

RICH2

Ring Imaging Cherenkov Detectors (RICH)



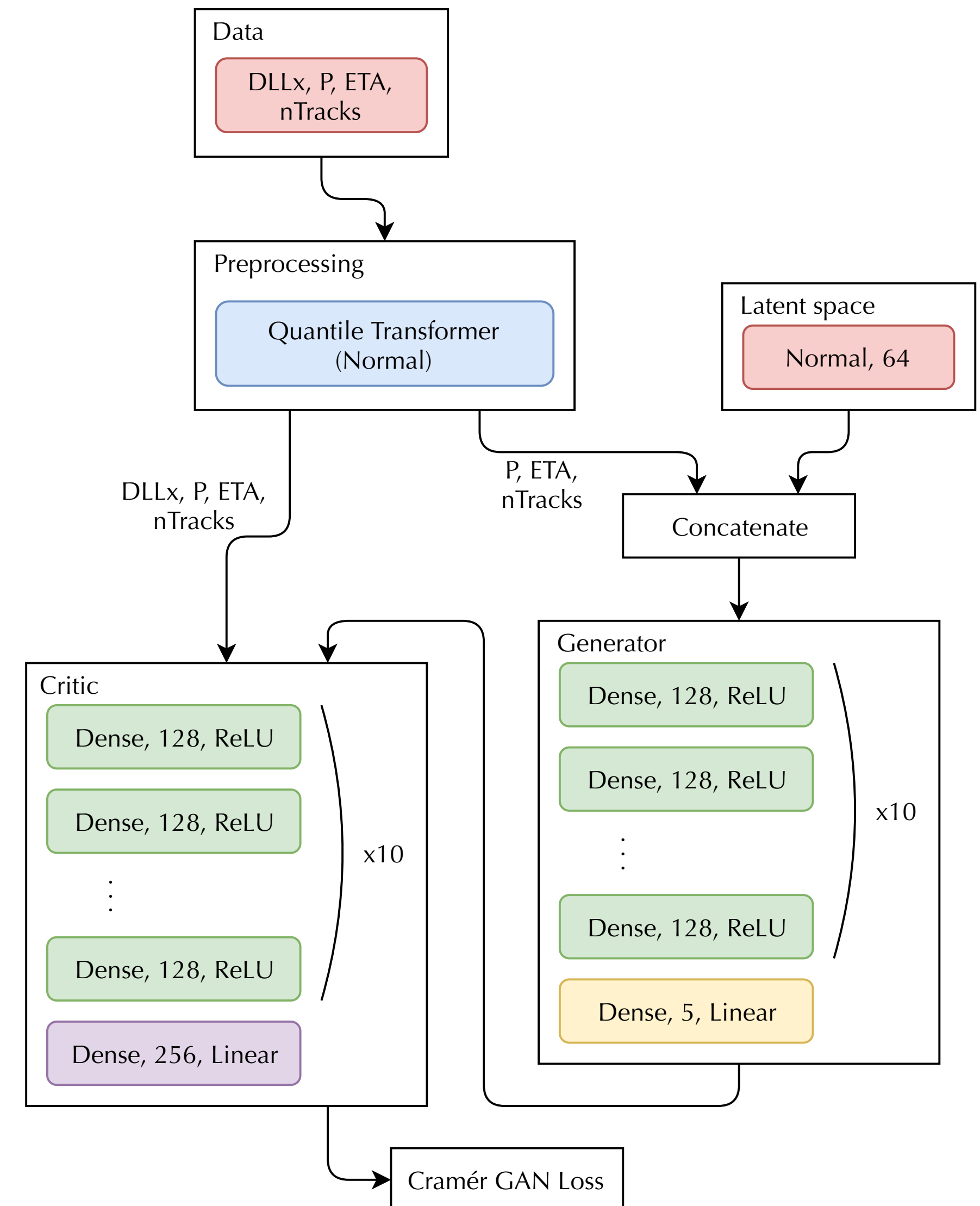
Ring Imaging Cherenkov Detectors (RICH)



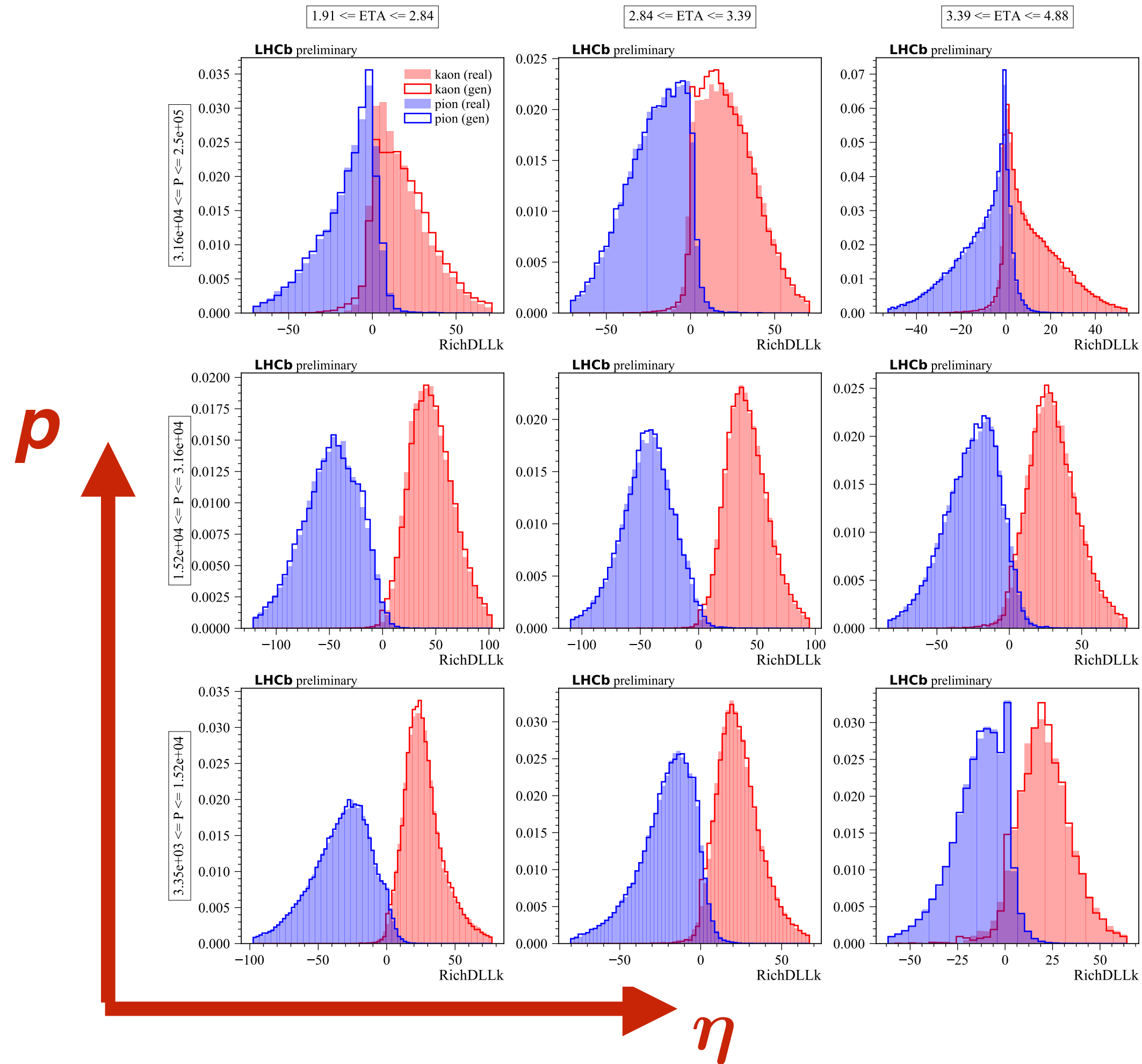
- PID with RICH is done with a **global log-likelihood method**
- PID information encoded in **log-likelihood differences (DLL)** between particle type hypotheses

Implementation details

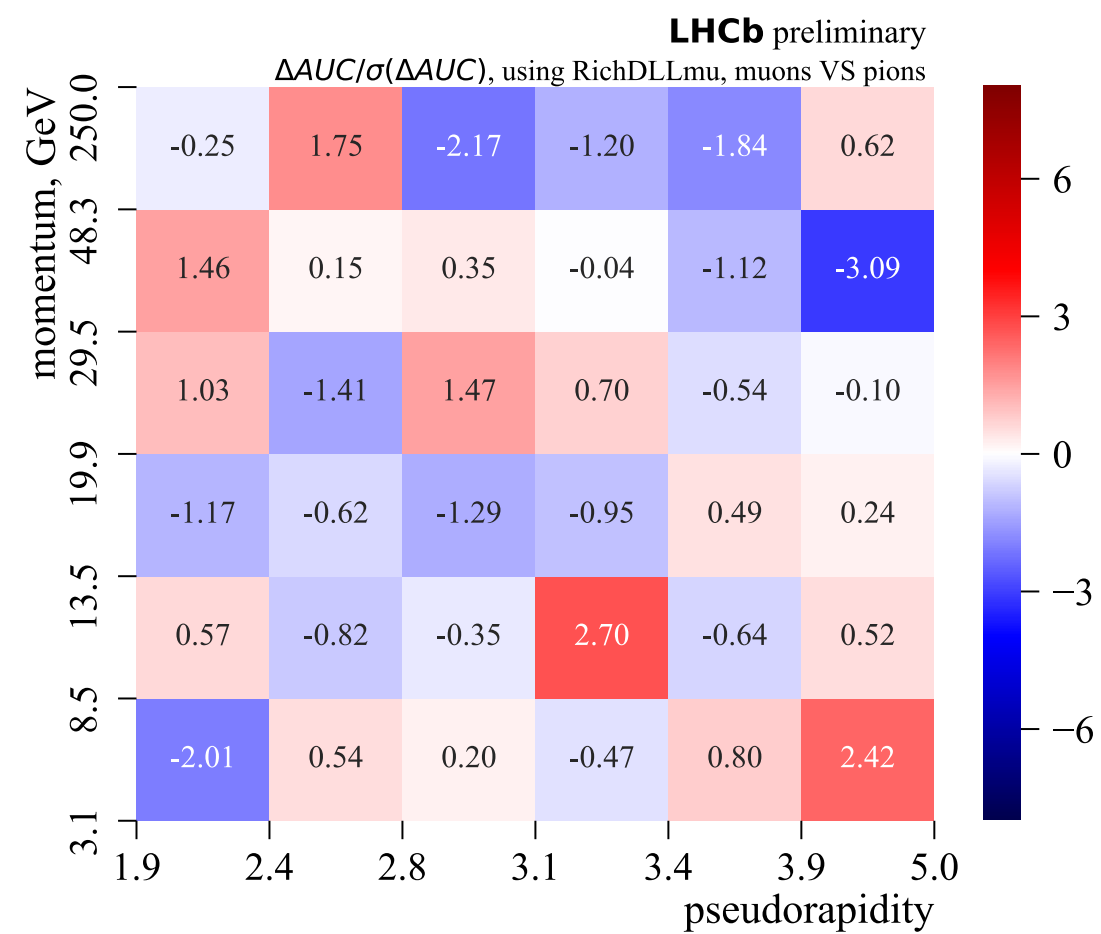
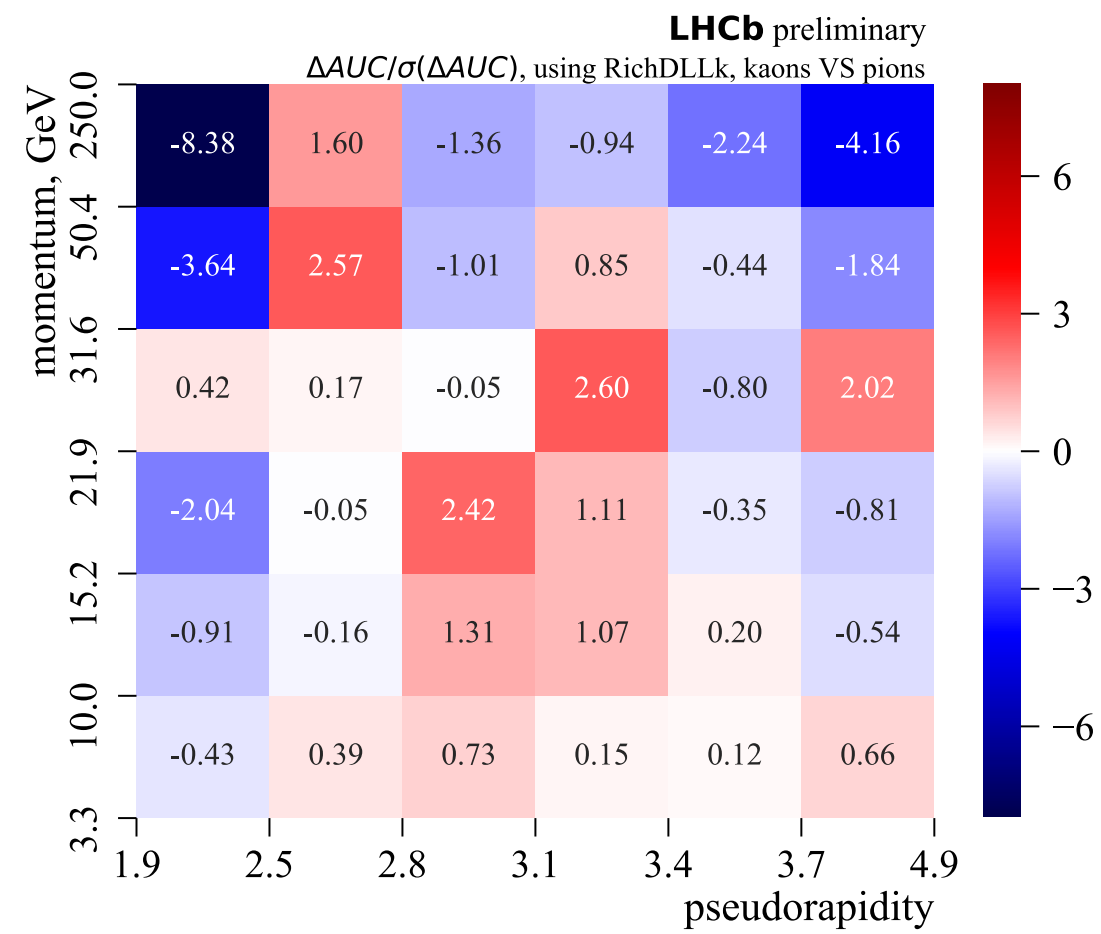
- Trained on real data!
 - 2016 calibration samples
- sPlot technique used for signal extraction
 - we apply s-weights to the loss function
- Independent models for each of the particle types
- Cramér GAN used
 - advancement of WGAN
 - more information in the backup



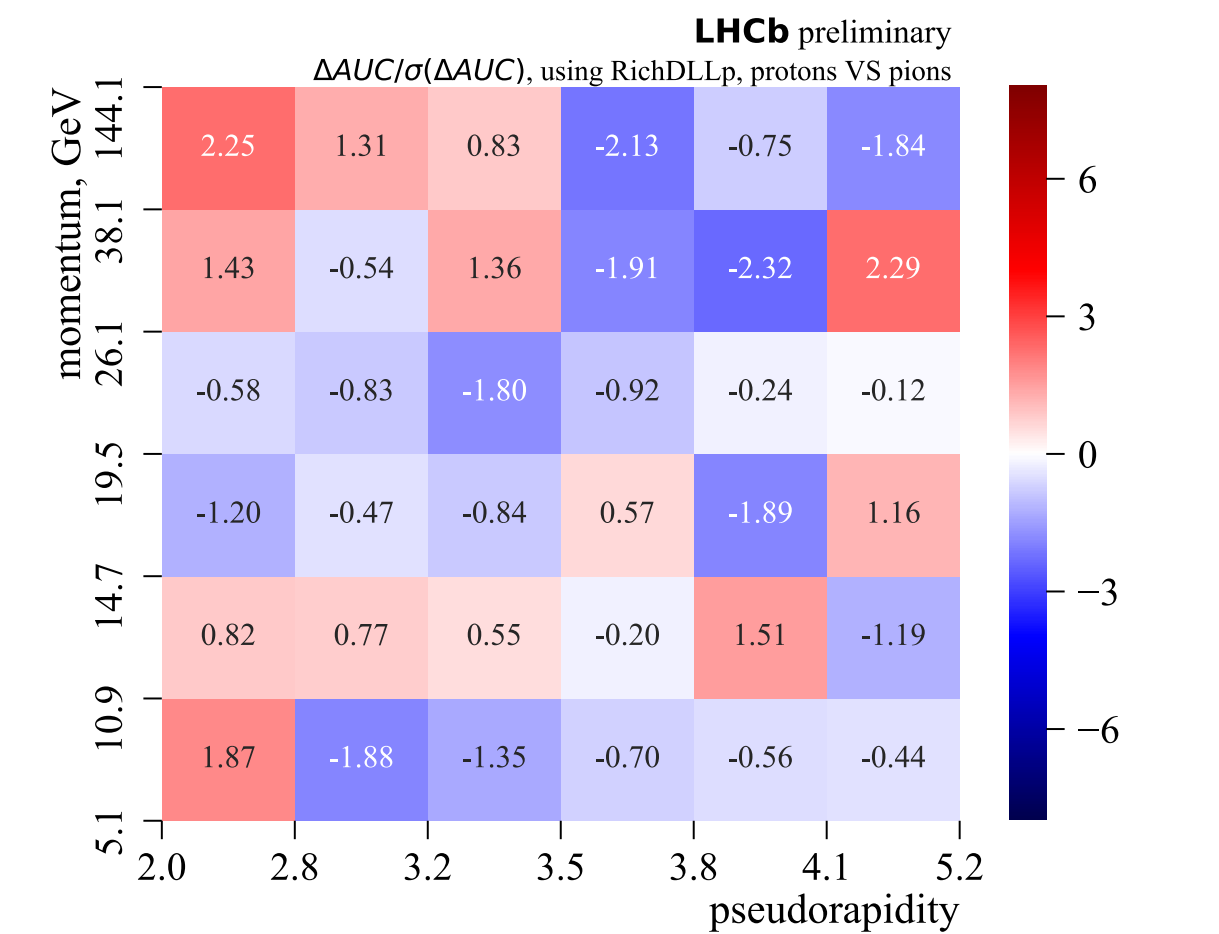
Results



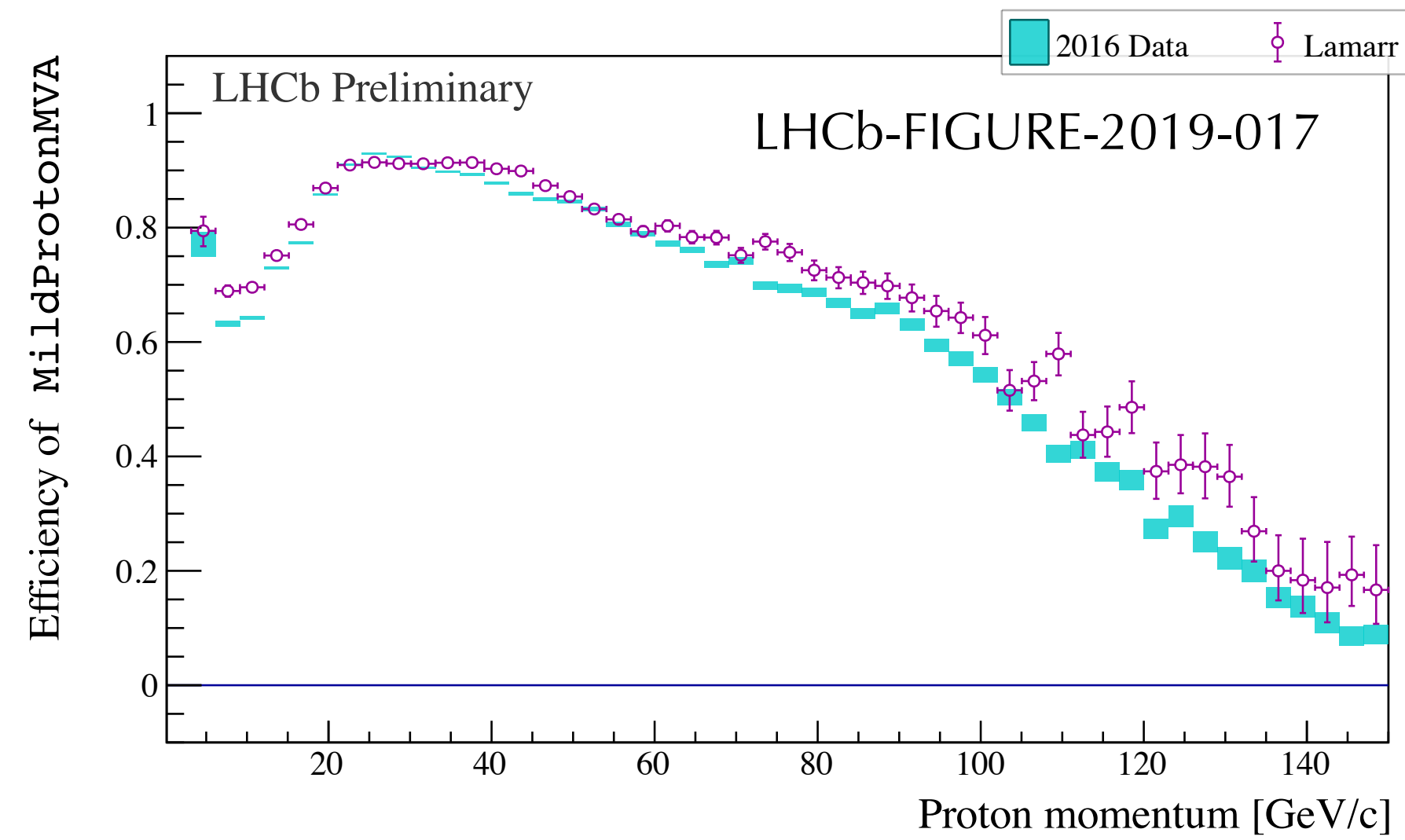
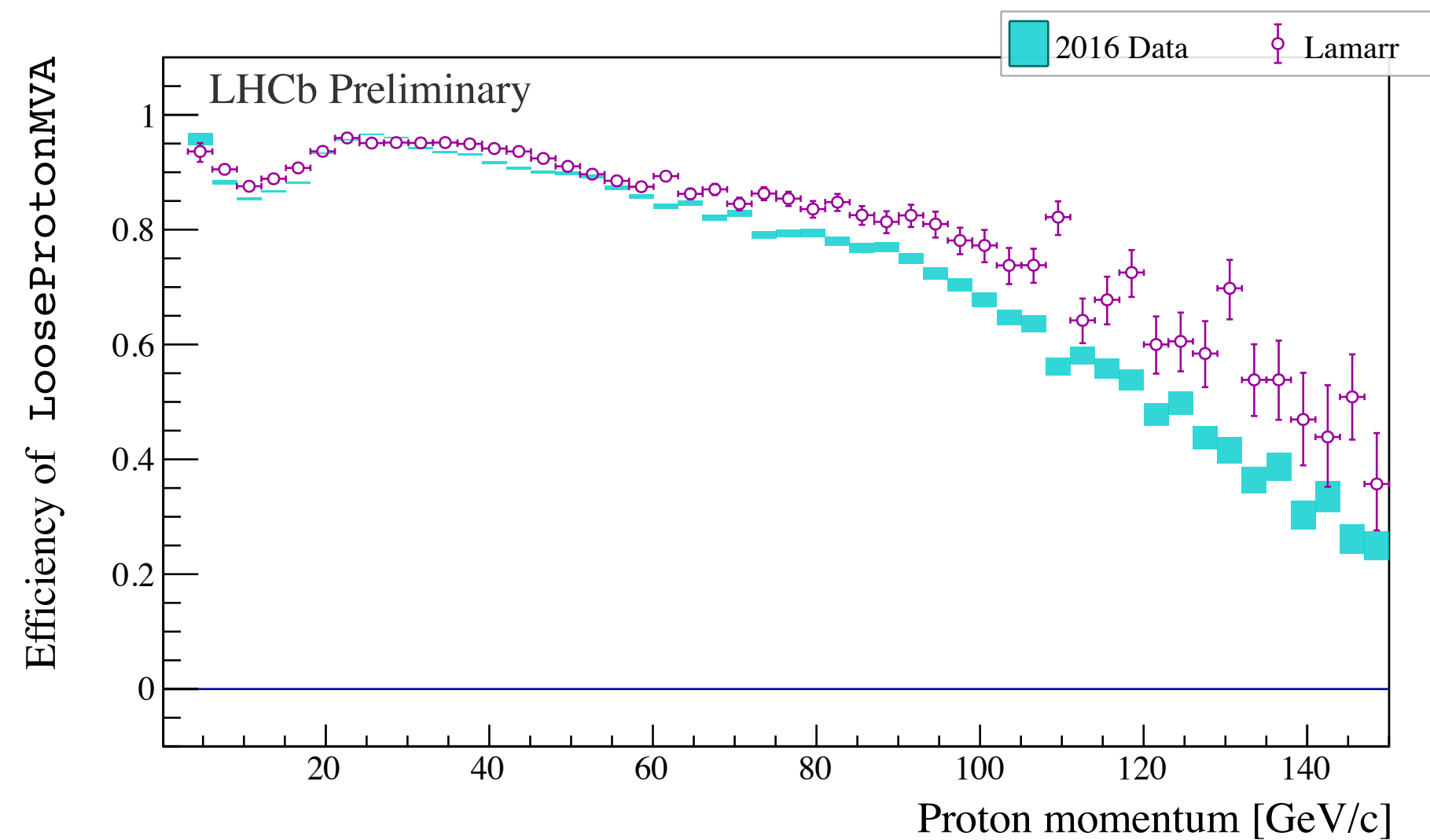
Detector response to kaons and pions
(real data vs generated with GAN)



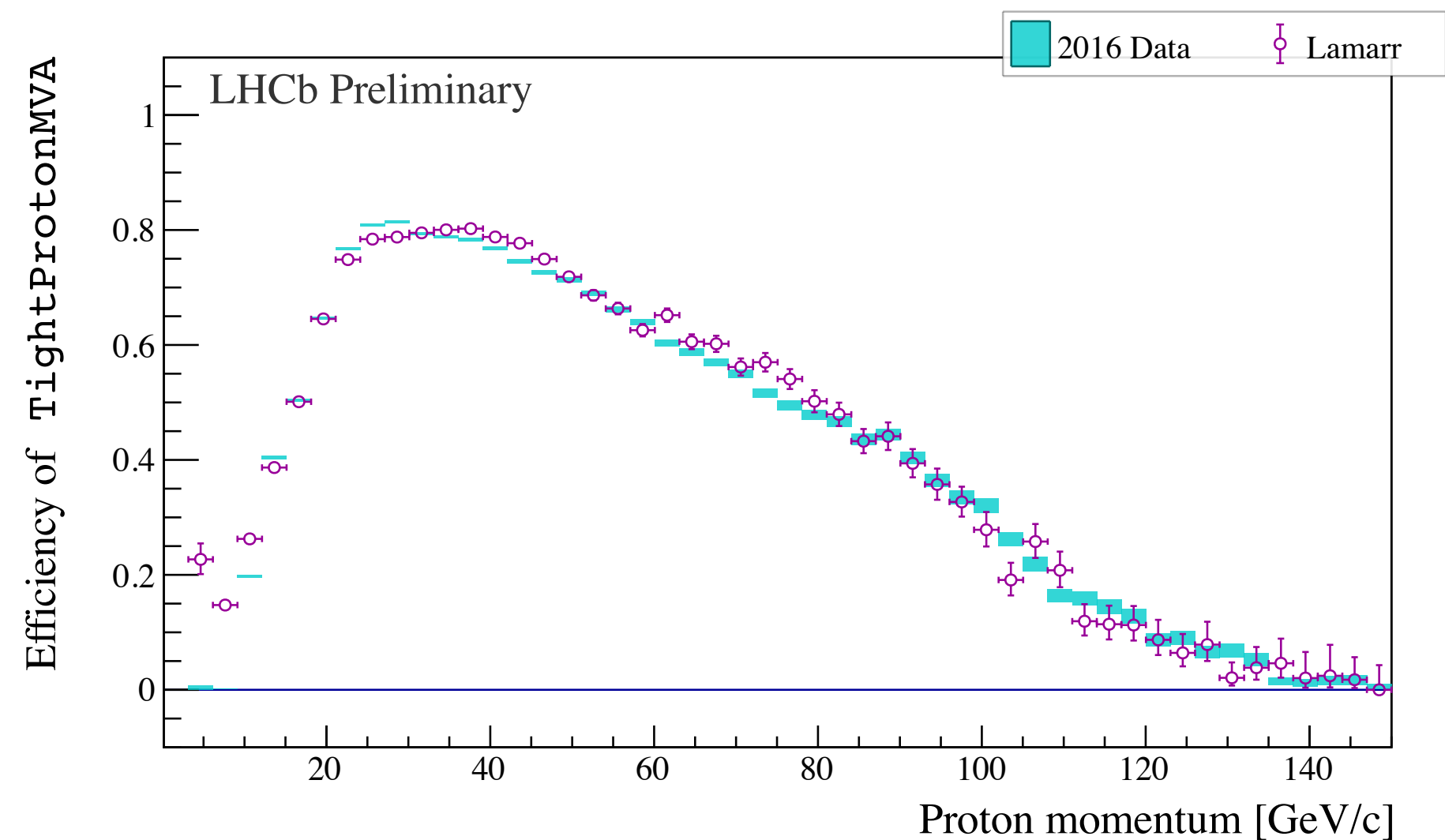
Difference between
real and generated
ROC AUCs in units of
stat uncertainty



Performance in full-stack implementation



The agreement is good and allows to use these results in further analyses



- PID performed combining information from all LHCb subsystems, part of which was simulated using GANs
- The simulated proton identification efficiency is compared with an independent sample of calibration data as a function of the proton momentum

Summary

- GANs are a really fascinating area of research in DL
- They are of great interest to experimental HEP as they provide a natural tool for fast-simulation
- Fast simulation of LHCb RICH detectors using GANs shows promising results

BACKUP

Wasserstein distance

How to define it formally?

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from \mathbf{x}_1 to \mathbf{x}_2 : $\gamma(\mathbf{x}_1, \mathbf{x}_2)$

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from \mathbf{x}_1 to \mathbf{x}_2 : $\gamma(\mathbf{x}_1, \mathbf{x}_2)$

Since we want to get from P to Q , our plan has to satisfy these conditions:

$$\int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 = Q(\mathbf{x}_2), \quad \int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = P(\mathbf{x}_1)$$

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from \mathbf{x}_1 to \mathbf{x}_2 : $\gamma(\mathbf{x}_1, \mathbf{x}_2)$

Since we want to get from P to Q , our plan has to satisfy these conditions:

$$\int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 = Q(\mathbf{x}_2), \quad \int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = P(\mathbf{x}_1)$$

The cost for some particular values \mathbf{x}_1 and \mathbf{x}_2 is: $dC = \underbrace{\|\mathbf{x}_1 - \mathbf{x}_2\|}_{\text{distance}} \underbrace{\gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2}_{\text{amount}}$

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from \mathbf{x}_1 to \mathbf{x}_2 : $\gamma(\mathbf{x}_1, \mathbf{x}_2)$

Since we want to get from P to Q , our plan has to satisfy these conditions:

$$\int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 = Q(\mathbf{x}_2), \quad \int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = P(\mathbf{x}_1)$$

The cost for some particular values \mathbf{x}_1 and \mathbf{x}_2 is: $dC = \underbrace{\|\mathbf{x}_1 - \mathbf{x}_2\|}_{\text{distance}} \underbrace{\gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2}_{\text{amount}}$

Then, total cost can simply be written as:

$$C = \iint \gamma(\mathbf{x}_1, \mathbf{x}_2) \|\mathbf{x}_1 - \mathbf{x}_2\| d\mathbf{x}_1 d\mathbf{x}_2$$

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from \mathbf{x}_1 to \mathbf{x}_2 : $\gamma(\mathbf{x}_1, \mathbf{x}_2)$

Since we want to get from P to Q , our plan has to satisfy these conditions:

$$\int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 = Q(\mathbf{x}_2), \quad \int \gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2 = P(\mathbf{x}_1)$$

The cost for some particular values \mathbf{x}_1 and \mathbf{x}_2 is: $dC = \underbrace{\|\mathbf{x}_1 - \mathbf{x}_2\|}_{\text{distance}} \underbrace{\gamma(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_1 d\mathbf{x}_2}_{\text{amount}}$

Then, total cost can simply be written as:

$$C = \iint \gamma(\mathbf{x}_1, \mathbf{x}_2) \|\mathbf{x}_1 - \mathbf{x}_2\| d\mathbf{x}_1 d\mathbf{x}_2 = \mathbb{E}_{\mathbf{x}_1, \mathbf{x}_2 \sim \gamma(\mathbf{x}_1, \mathbf{x}_2)} \|\mathbf{x}_1 - \mathbf{x}_2\|$$

Wasserstein distance

How to define it formally?

Wasserstein distance

How to define it formally?

Let π be the set of all plans that convert P into Q , i.e.:

$$\pi = \left\{ \gamma : \int \gamma(x_1, x_2) dx_1 = Q(x_2), \int \gamma(x_1, x_2) dx_2 = P(x_1) \right\}$$

Wasserstein distance

How to define it formally?

Let π be the set of all plans that convert P into Q , i.e.:

$$\pi = \left\{ \gamma : \int \gamma(x_1, x_2) dx_1 = Q(x_2), \int \gamma(x_1, x_2) dx_2 = P(x_1) \right\}$$

Then, the Wasserstein distance is defined as:

$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

WGAN

(Wasserstein GAN)

Imagine a GAN trying to minimize EMD between the real and generated distributions. You may be wondering:



$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

WGAN

(Wasserstein GAN)

Imagine a GAN trying to minimize EMD between the real and generated distributions. You may be wondering:

How would that be better than a regular GAN?

And more importantly...



$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

WGAN (Wasserstein GAN)

Imagine a GAN trying to minimize EMD between the real and generated distributions. You may be wondering:

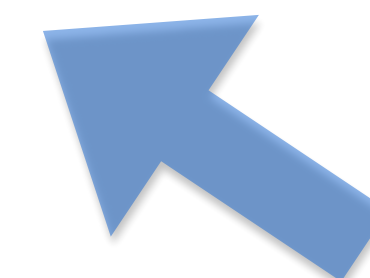
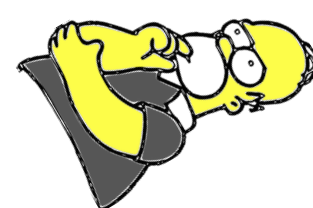
How would that be better than a regular GAN?

How does one calculate this madness?!

And more importantly...



$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$



Kantorovich-Rubinstein duality

(or how to calculate EMD)

$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

disclaimer: not a strict
mathematical derivation

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

Let's add the following term to this expression:

$$+ \inf_{\gamma} \sup_f \mathbb{E}_{x_1, x_2 \sim \gamma} [\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) - (f(x_1) - f(x_2))]$$

$f(x)$ — real-valued function

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

Let's add the following term to this expression:

$$+ \inf_{\gamma} \sup_f \mathbb{E}_{x_1, x_2 \sim \gamma} [\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) - (f(x_1) - f(x_2))]$$

$f(x)$ — real-valued function

These cancel out when $\gamma \in \pi$
otherwise supremum over $f(x)$ goes to $+\infty$

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

Let's add the following term to this expression:

$$+ \inf_{\gamma} \sup_f \mathbb{E}_{x_1, x_2 \sim \gamma} [\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) - (f(x_1) - f(x_2))]$$

$f(x)$ — real-valued function

These cancel out when $\gamma \in \pi$
otherwise supremum over $f(x)$ goes to $+\infty$

Therefore, we can remove the $\gamma \in \pi$ condition from the whole expression:

$$= \inf_{\gamma} \sup_f \mathbb{E}_{x_1, x_2 \sim \gamma} [\|x_1 - x_2\| + \mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) - (f(x_1) - f(x_2))]$$

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$\text{EMD}(P, Q) = \inf_{\gamma \in \pi} \mathbb{E}_{x_1, x_2 \sim \gamma} \|x_1 - x_2\|$$

Let's add the following term to this expression:

$$+ \inf_{\gamma} \sup_f \mathbb{E}_{x_1, x_2 \sim \gamma} [\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) - (f(x_1) - f(x_2))]$$

$f(x)$ — real-valued function

These cancel out when $\gamma \in \pi$
otherwise supremum over $f(x)$ goes to $+\infty$

Therefore, we can remove the $\gamma \in \pi$ condition from the whole expression:

$$= \inf_{\gamma} \sup_f \mathbb{E}_{x_1, x_2 \sim \gamma} [\|x_1 - x_2\| + \mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) - (f(x_1) - f(x_2))]$$

Infimum and supremum operations can be swapped under certain conditions

(satisfied here — see <https://vincentherrmann.github.io/blog/wasserstein/> for more detailed info)

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$= \sup_f \inf_{\gamma} \left[\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) + \mathbb{E}_{x_1, x_2 \sim \gamma} [\|x_1 - x_2\| - (f(x_1) - f(x_2))] \right]$$

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$= \sup_f \inf_{\gamma} \left[\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) + \mathbb{E}_{x_1, x_2 \sim \gamma} [\|x_1 - x_2\| - (f(x_1) - f(x_2))] \right]$$

Consider the following case: $|f(a) - f(b)| \leq \|a - b\|, \quad \forall a, b$

We'll denote it as: $\|f\|_L \leq 1$

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$= \sup_f \inf_{\gamma} \left[\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) + \mathbb{E}_{x_1, x_2 \sim \gamma} [\|x_1 - x_2\| - (f(x_1) - f(x_2))] \right]$$

Consider the following case: $|f(a) - f(b)| \leq \|a - b\|, \quad \forall a, b$

We'll denote it as: $\|f\|_L \leq 1$

For such case this term is 0

Otherwise the whole expression is $-\infty$

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$= \sup_f \inf_{\gamma} \left[\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) + \mathbb{E}_{x_1, x_2 \sim \gamma} [||x_1 - x_2|| - (f(x_1) - f(x_2))] \right]$$

Consider the following case: $|f(a) - f(b)| \leq ||a - b||, \quad \forall a, b$

We'll denote it as: $||f||_L \leq 1$

For such case this term is 0

Otherwise the whole expression is $-\infty$

Therefore we can finally rewrite the whole thing as:

$$\text{EMD}(P, Q) = \sup_{||f||_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

Kantorovich-Rubinstein duality

(or how to calculate EMD)

disclaimer: not a strict
mathematical derivation

$$= \sup_f \inf_{\gamma} \left[\mathbb{E}_{s \sim P} f(s) - \mathbb{E}_{t \sim Q} f(t) + \mathbb{E}_{x_1, x_2 \sim \gamma} [\|x_1 - x_2\| - (f(x_1) - f(x_2))] \right]$$

Consider the following case: $|f(a) - f(b)| \leq \|a - b\|, \quad \forall a, b$

We'll denote it as: $\|f\|_L \leq 1$

For such case this term is 0

Otherwise the whole expression is $-\infty$

Therefore we can finally rewrite the whole thing as:

$$\text{EMD}(P, Q) = \sup_{\|f\|_L \leq 1} [\mathbb{E}_{x \sim P} f(x) - \mathbb{E}_{x \sim Q} f(x)]$$

How is this
simpler?



WGAN

Algorithm 1 WGAN, our proposed algorithm. All experiments in the paper used the default values $\alpha = 0.00005$, $c = 0.01$, $m = 64$, $n_{\text{critic}} = 5$.

Require: : α , the learning rate. c , the clipping parameter. m , the batch size. n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_0 , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSPProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSPProp}(\theta, g_\theta)$ 
12: end while
```

<https://arxiv.org/abs/1701.07875>

Biased Wasserstein gradients

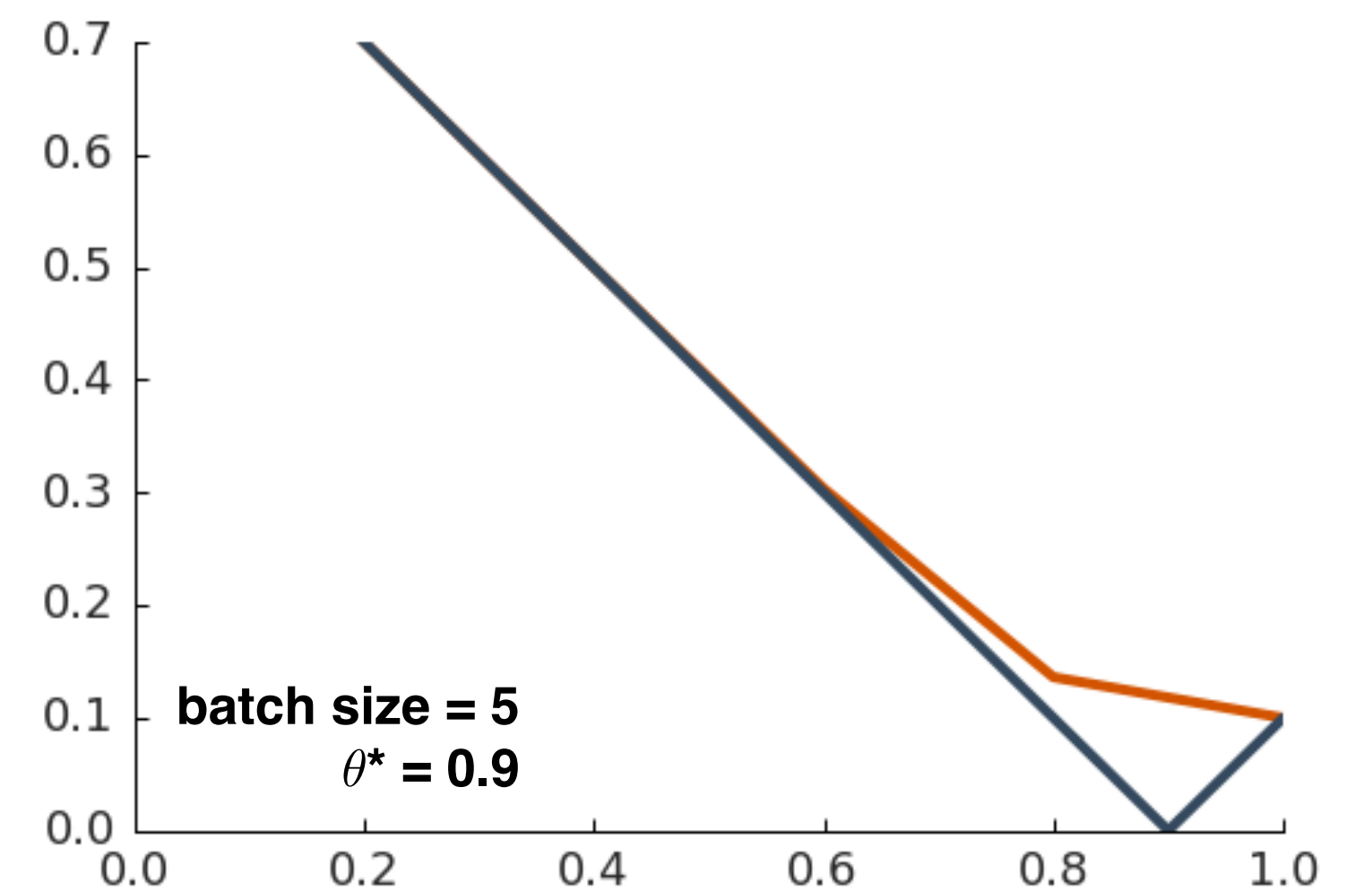
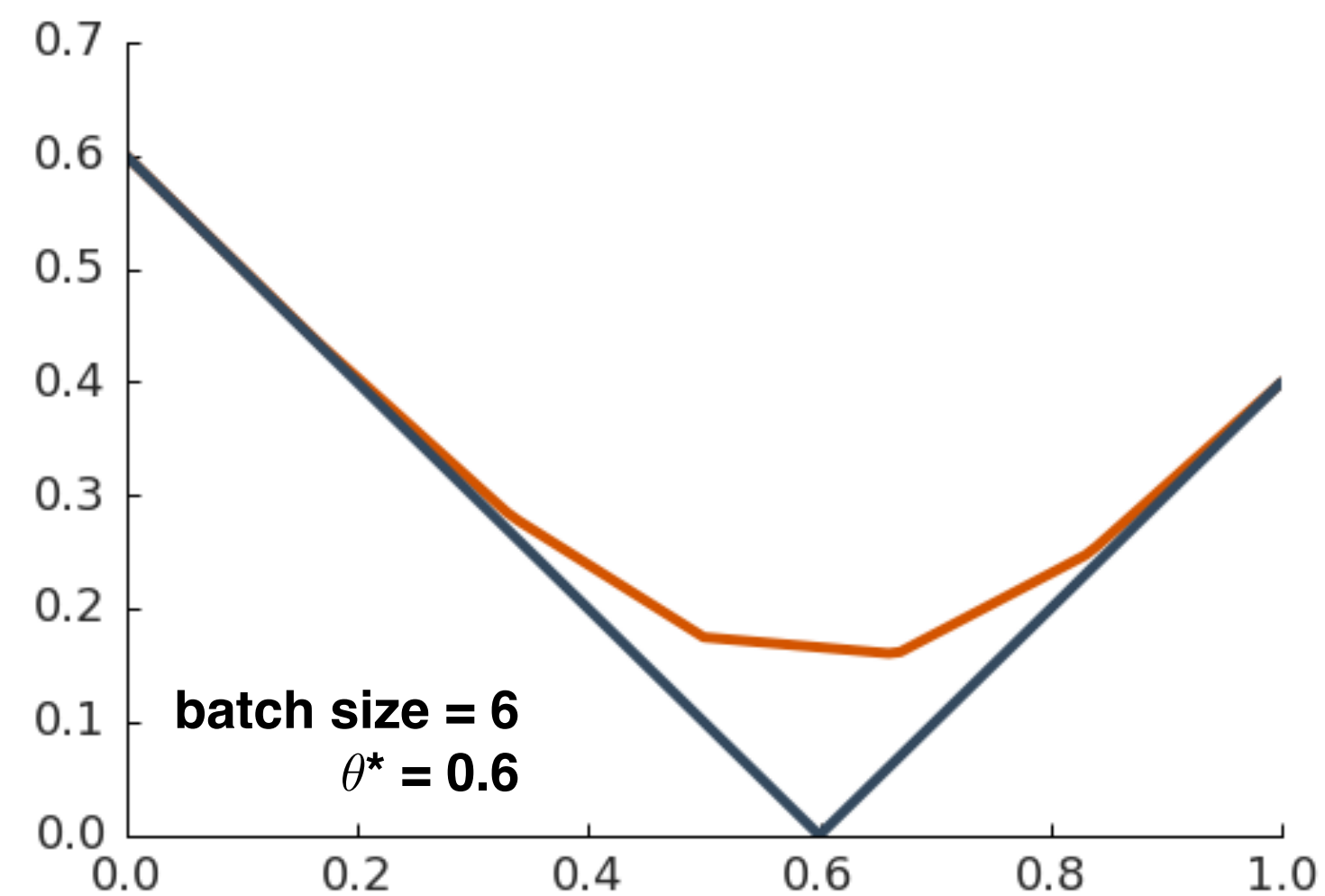
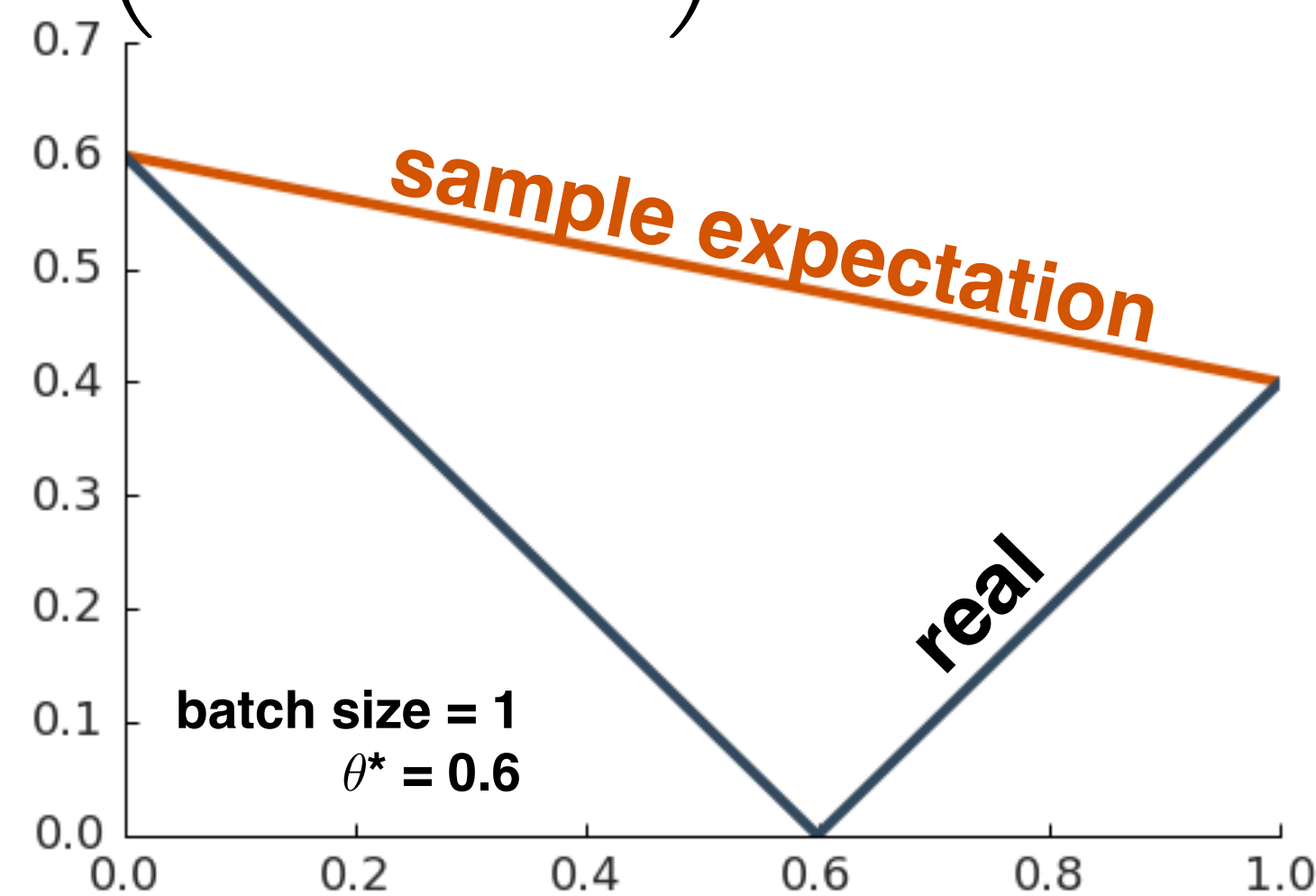
Authors of <https://arxiv.org/abs/1705.10743> showed for EMD that expected sample gradients may differ from true gradients

This may lead to a wrong minimum:

Bernoulli distribution

$$x \sim \mathcal{B}(\theta) : \begin{cases} \mathbb{P}(x = 0) = 1 - \theta \\ \mathbb{P}(x = 1) = \theta \end{cases}$$

$$\text{EMD}(\mathcal{B}(\theta^*), \mathcal{B}(\theta))$$



Parameter θ

Cramer and energy distances

Cramer distance:

$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$

CDFS for P and Q

Cramer and energy distances

Cramer distance:

$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$

CDFS for P and Q

Energy distance — a natural extension of the Cramer distance to the multivariate case:

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

Cramer and energy distances

Cramer distance:

$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$

CDFS for P and Q

Energy distance — a natural extension of the Cramer distance to the multivariate case:

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

– preserves nice properties of EMD

Cramer and energy distances

Cramer distance:

$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$

CDFS for P and Q

Energy distance — a natural extension of the Cramer distance to the multivariate case:

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

- preserves nice properties of EMD
- is proven to have unbiased sample gradients

Cramer and energy distances

Cramer distance:

$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$

CDFS for P and Q

Energy distance — a natural extension of the Cramer distance to the multivariate case:

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

- preserves nice properties of EMD
- is proven to have unbiased sample gradients

In one-dimensional case: $l_2^2(P, Q) = \frac{1}{2} \mathcal{E}(P, Q)$

Cramer GAN

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

<https://arxiv.org/abs/1705.10743>

Cramer GAN

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

Actually, this can already be used to train just the generator (regression task)

<https://arxiv.org/abs/1705.10743>

Cramer GAN

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

Actually, this can already be used to train just the generator (regression task)

However, much better quality can be achieved with a critic transforming P and Q into P^* and Q^* :

$$\mathcal{L}_{generator} = \mathcal{E}(P^*, Q^*)$$

$$\mathcal{L}_{critic} = -\mathcal{E}(P^*, Q^*) + \text{regularization}$$

$$x \sim P^* : x = f(x'), \quad x' \sim P$$

$$x \sim Q^* : x = f(x'), \quad x' \sim Q$$

<https://arxiv.org/abs/1705.10743>

Cramer GAN

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

Actually, this can already be used to train just the generator (regression task)

However, much better quality can be achieved with a critic transforming P and Q into P^* and Q^* :

$$\mathcal{L}_{generator} = \mathcal{E}(P^*, Q^*)$$

$$\mathcal{L}_{critic} = -\mathcal{E}(P^*, Q^*) + \text{regularization}$$

$$x \sim P^* : x = f(x'), \quad x' \sim P$$

$$x \sim Q^* : x = f(x'), \quad x' \sim Q$$

**Regularization (e.g. GP)
needed to prevent critic from
scaling the data to infinity**

<https://arxiv.org/abs/1705.10743>

Cramer GAN

$$\mathcal{E}(P, Q) = 2\mathbb{E}_{\substack{x_1 \sim P \\ x_2 \sim Q}} \|x_1 - x_2\| - \mathbb{E}_{x_1, x'_1 \sim P} \|x_1 - x'_1\| - \mathbb{E}_{x_2, x'_2 \sim Q} \|x_2 - x'_2\|$$

Actually, this can already be used to train just the generator (regression task)

However, much better quality can be achieved with a critic transforming P and Q into P^* and Q^* :

$$\mathcal{L}_{generator} = \mathcal{E}(P^*, Q^*)$$

$$\mathcal{L}_{critic} = -\mathcal{E}(P^*, Q^*) + \text{regularization}$$

$$x \sim P^* : x = f(x'), \quad x' \sim P$$

$$x \sim Q^* : x = f(x'), \quad x' \sim Q$$

**Regularization (e.g. GP)
needed to prevent critic from
scaling the data to infinity**

**In the paper 'surrogate' loss
($x'_1 \rightarrow 0$) is used to avoid
sampling from data twice**

<https://arxiv.org/abs/1705.10743>

Evaluating generative models

Evaluating generative models

- No single guide to follow

Evaluating generative models

- No single guide to follow
- Approaches are very problem-specific

Evaluating generative models

- No single guide to follow
- Approaches are very problem-specific
 - E.g. perceived visual quality of generated images vs. quality of a generated invariant mass distribution

Evaluating generative models

- No single guide to follow
- Approaches are very problem-specific
 - E.g. perceived visual quality of generated images vs. quality of a generated invariant mass distribution
 - Most solutions are adapted to or invented for a given particular task

Evaluating generative models

- No single guide to follow
- Approaches are very problem-specific
 - E.g. perceived visual quality of generated images vs. quality of a generated invariant mass distribution
 - Most solutions are adapted to or invented for a given particular task
- We'll mention some notable examples

Evaluating generative models

(most obvious thing to do)

- By-eye comparison
 - Compare individual objects or whole distributions (e.g. in projections) where possible
 - There might be no need to do any complicated evaluation if the model's results simply look bad

Evaluating generative models

(simple things to do)

- Compare meaningful physical characteristics
 - Integral characteristics (e.g. total energy of a calorimeter cluster)
 - Means, medians, standard deviations, etc.

Evaluating generative models

(simple things to do)

- Compare meaningful physical characteristics
 - Integral characteristics (e.g. total energy of a calorimeter cluster)
 - Means, medians, standard deviations, etc.
- Statistical tests (Chi2, Kolmogorov-Smirnov, etc.)
 - Between individual dimensions or projections
 - p -values might look insane, probably better to compare the statistics themselves

Evaluating generative models

(the 'additional classifier' way)

- Train an independent classifier (e.g. xgboost) to distinguish real and fake samples
- Evaluate your GAN by checking the classifier's score (e.g. ROC AUC)
- Pros:
 - An objective quality measure
- Cons:
 - Resource consuming
 - Requires hyper-parameter tuning
 - May get picky to things that are not important

Evaluating generative models

(Inception score)

Evaluating generative models

(Inception score)

- Introduced in <https://arxiv.org/abs/1606.03498>

Evaluating generative models

(Inception score)

- Introduced in <https://arxiv.org/abs/1606.03498>
- Apply the Inception model (pre-trained image classifier) to obtain conditional label distribution $p(y|x)$ for each image x
 - this should be low-entropy (the classifier should be certain)

Evaluating generative models

(Inception score)

- Introduced in <https://arxiv.org/abs/1606.03498>
- Apply the Inception model (pre-trained image classifier) to obtain conditional label distribution $p(y|x)$ for each image x
 - this should be low-entropy (the classifier should be certain)
- calculate marginal $p(y) = \int p(y|x = G(z))p(z)dz$
 - this should be high-entropy (diversity of samples)

Evaluating generative models

(Inception score)

- Introduced in <https://arxiv.org/abs/1606.03498>
- Apply the Inception model (pre-trained image classifier) to obtain conditional label distribution $p(y|x)$ for each image x
 - this should be low-entropy (the classifier should be certain)
- calculate marginal $p(y) = \int p(y|x = G(z))p(z)dz$
 - this should be high-entropy (diversity of samples)
- Combining these two requirements:

$$\text{IS} = \exp \left(\mathbb{E}_x \left[\text{KL}(p(y|x) || p(y)) \right] \right)$$

Evaluating generative models

(Fréchet Inception Distance – FID)

<https://arxiv.org/abs/1706.08500>

Evaluating generative models

(Fréchet Inception Distance – FID)

<https://arxiv.org/abs/1706.08500>

- One of drawbacks of IS is that it doesn't care about the true distribution

Evaluating generative models

(Fréchet Inception Distance – FID)

<https://arxiv.org/abs/1706.08500>

- One of drawbacks of IS is that it doesn't care about the true distribution
- Instead one can compare distributions of activations at some Inception layer (originally – last pooling layer)

Evaluating generative models

(Fréchet Inception Distance – FID)

<https://arxiv.org/abs/1706.08500>

- One of drawbacks of IS is that it doesn't care about the true distribution
- Instead one can compare distributions of activations at some Inception layer (originally – last pooling layer)
- The authors proposed calculating the Fréchet (aka Wasserstein-2) distance

Evaluating generative models

(Fréchet Inception Distance – FID)

<https://arxiv.org/abs/1706.08500>

- One of drawbacks of IS is that it doesn't care about the true distribution
- Instead one can compare distributions of activations at some Inception layer (originally – last pooling layer)
- The authors proposed calculating the Fréchet (aka Wasserstein-2) distance
- Distance between multivariate Gaussian approximations:

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2})$$

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Definition 1. For $\alpha, \beta \in (0, 1]$, the probability distribution Q has precision α at recall β w.r.t. P if there exist distributions μ, ν_P and ν_Q such that

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Definition 1. For $\alpha, \beta \in (0, 1]$, the probability distribution Q has precision α at recall β w.r.t. P if there exist distributions μ, ν_P and ν_Q such that

Decomposition with
a common part

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Definition 1. For $\alpha, \beta \in (0, 1]$, the probability distribution Q has precision α at recall β w.r.t. P if there exist distributions μ, ν_P and ν_Q such that

Decomposition with
a common part

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

Definition 2. The set of attainable pairs of precision and recall of a distribution Q w.r.t. a distribution P is denoted by $\text{PRD}(Q, P)$ and it consists of all (α, β) satisfying Definition 1 and the pair $(0, 0)$.

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Definition 1. For $\alpha, \beta \in (0, 1]$, the probability distribution Q has precision α at recall β w.r.t. P if there exist distributions μ, ν_P and ν_Q such that

Decomposition with a common part

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

Definition 2. The set of attainable pairs of precision and recall of a distribution Q w.r.t. a distribution P is denoted by $\text{PRD}(Q, P)$ and it consists of all (α, β) satisfying Definition 1 and the pair $(0, 0)$.

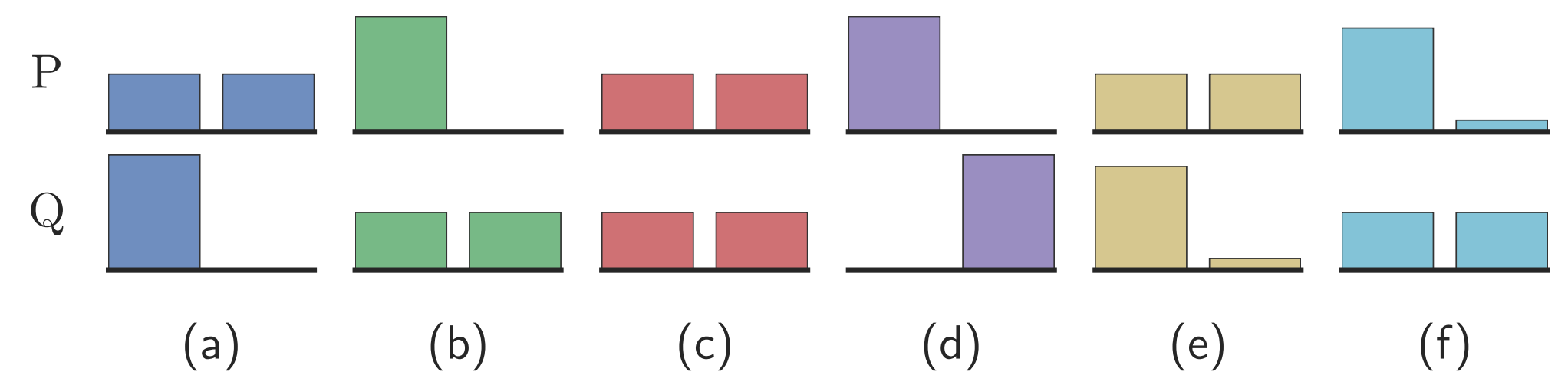


Figure 2: Intuitive examples of P and Q .

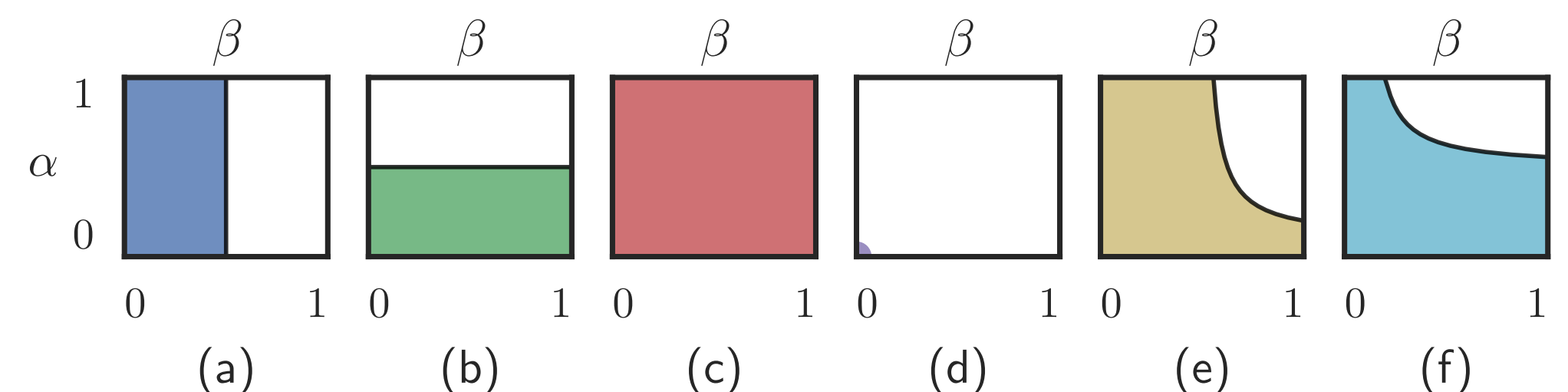


Figure 3: $\text{PRD}(Q, P)$ for the examples above.

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Definition 1. For $\alpha, \beta \in (0, 1]$, the probability distribution Q has precision α at recall β w.r.t. P if there exist distributions μ, ν_P and ν_Q such that

Decomposition with a common part

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

Definition 2. The set of attainable pairs of precision and recall of a distribution Q w.r.t. a distribution P is denoted by $\text{PRD}(Q, P)$ and it consists of all (α, β) satisfying Definition 1 and the pair $(0, 0)$.

- The authors provide an algorithm to calculate it for discrete distributions

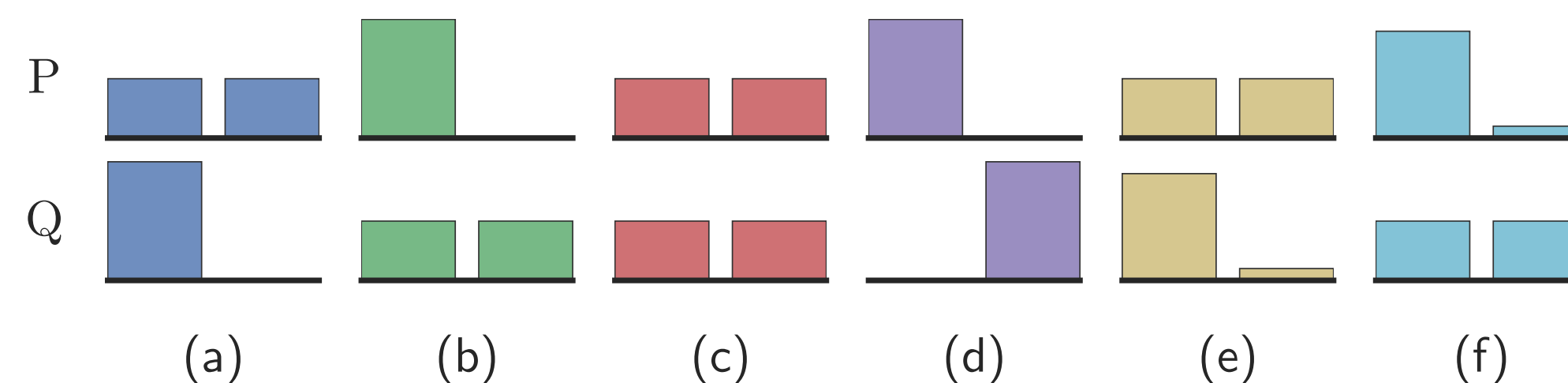


Figure 2: Intuitive examples of P and Q .

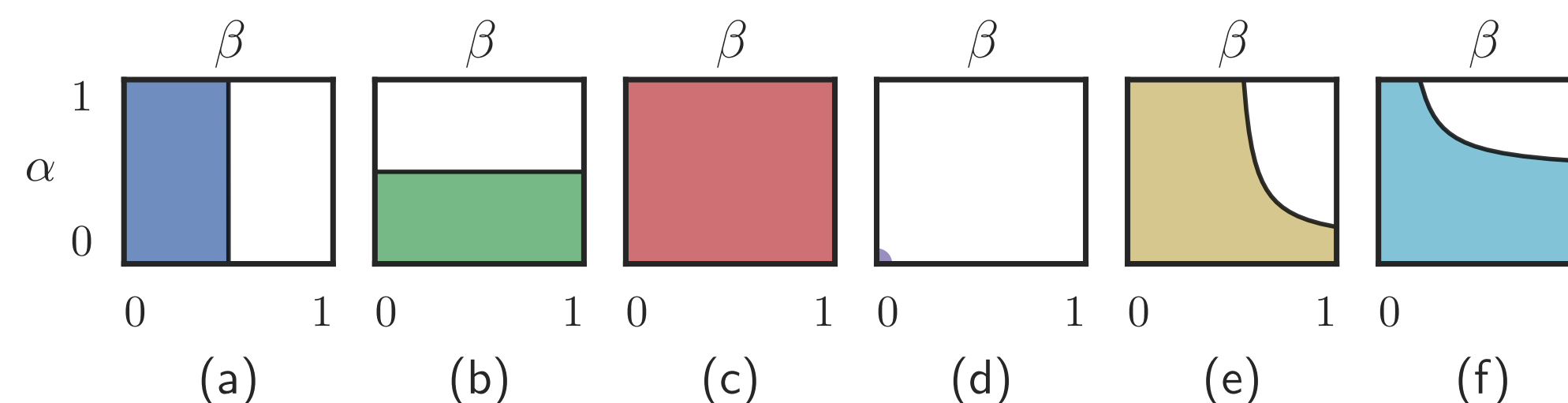


Figure 3: $\text{PRD}(Q, P)$ for the examples above.

Evaluating generative models

<https://arxiv.org/abs/1806.00035>

(Precision and Recall)

Definition 1. For $\alpha, \beta \in (0, 1]$, the probability distribution Q has precision α at recall β w.r.t. P if there exist distributions μ, ν_P and ν_Q such that

Decomposition with
a common part

$$P = \beta\mu + (1 - \beta)\nu_P \quad \text{and} \quad Q = \alpha\mu + (1 - \alpha)\nu_Q. \quad (3)$$

Definition 2. The set of attainable pairs of precision and recall of a distribution Q w.r.t. a distribution P is denoted by $\text{PRD}(Q, P)$ and it consists of all (α, β) satisfying Definition 1 and the pair $(0, 0)$.

- The authors provide an algorithm to calculate it for discrete distributions
- They convert Inception activations to discrete distribution using k -means clustering

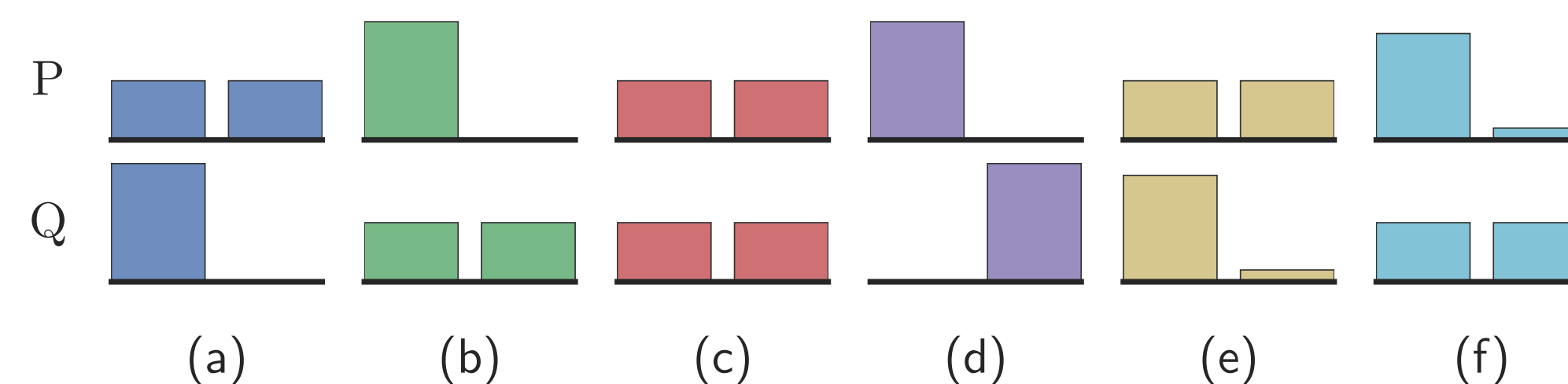


Figure 2: Intuitive examples of P and Q .

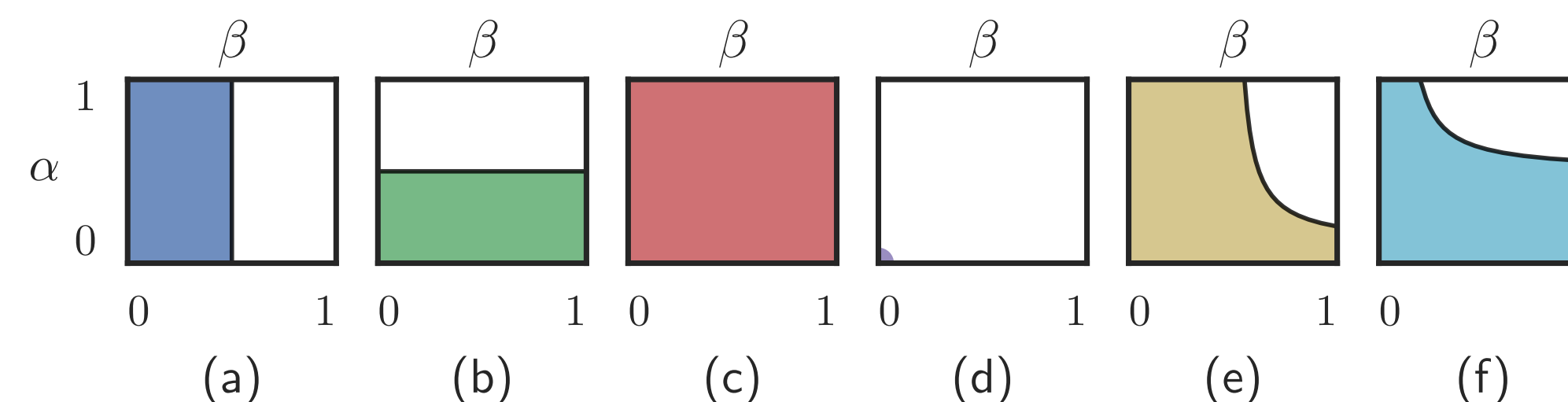


Figure 3: $\text{PRD}(Q, P)$ for the examples above.

Evaluating generative models

More on the subject
(extensive comparison of a
large variety of measures):

<https://arxiv.org/abs/1802.03446>

	Measure	Description
Quantitative	1. Average Log-likelihood [18, 22]	• Log likelihood of explaining realworld held out/test data using a density estimated from the generated data (e.g. using KDE or Parzen window estimation). $L = \frac{1}{N} \sum_i \log P_{model}(\mathbf{x}_i)$
	2. Coverage Metric [33]	• The probability mass of the true data “covered” by the model distribution $C := P_{data}(dP_{model} > t)$ with t such that $P_{model}(dP_{model} > t) = 0.95$
	3. Inception Score (IS) [3]	• KLD between conditional and marginal label distributions over generated data. $\exp(\mathbb{E}_{\mathbf{x}} [\mathbb{K}L(p(y \mathbf{x}) p(y))])$
	4. Modified Inception Score (m-IS) [34]	• Encourages diversity within images sampled from a particular category. $\exp(\mathbb{E}_{\mathbf{x}_i} [\mathbb{E}_{\mathbf{x}_j} [(\mathbb{K}L(P(y \mathbf{x}_i) P(y \mathbf{x}_j)))]])$
	5. Mode Score (MS) [35]	• Similar to IS but also takes into account the prior distribution of the labels over real data. $\exp(\mathbb{E}_{\mathbf{x}} [\mathbb{K}L(p(y \mathbf{x}) p(y^{train}))]) - \mathbb{K}L(p(y) p(y^{train}))$
	6. AM Score [36]	• Takes into account the KLD between distributions of training labels vs. predicted labels, as well as the entropy of predictions. $\mathbb{K}L(p(y^{train}) p(y)) + \mathbb{E}_{\mathbf{x}} [H(y \mathbf{x})]$
	7. Fréchet Inception Distance (FID) [37]	• Wasserstein-2 distance between multi-variate Gaussians fitted to data embedded into a feature space $FID(r, g) = \ \mu_r - \mu_g\ _2^2 + Tr(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{\frac{1}{2}})$
	8. Maximum Mean Discrepancy (MMD) [38]	• Measures the dissimilarity between two probability distributions P_r and P_g using samples drawn independently from each distribution. $M_k(P_r, P_g) = \mathbb{E}_{\mathbf{x}, \mathbf{x}' \sim P_r} [k(\mathbf{x}, \mathbf{x}')] - 2\mathbb{E}_{\mathbf{x} \sim P_r, \mathbf{y} \sim P_g} [k(\mathbf{x}, \mathbf{y})] + \mathbb{E}_{\mathbf{y}, \mathbf{y}' \sim P_g} [k(\mathbf{y}, \mathbf{y}')]$
	9. The Wasserstein Critic [39]	• The critic (e.g. an NN) is trained to produce high values at real samples and low values at generated samples $\hat{W}(\mathbf{x}_{test}, \mathbf{x}_g) = \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_{test}[i]) - \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_g[i])$
	10. Birthday Paradox Test [27]	• Measures the support size of a discrete (continuous) distribution by counting the duplicates (near duplicates)
	11. Classifier Two Sample Test (C2ST) [40]	• Answers whether two samples are drawn from the same distribution (e.g. by training a binary classifier)
	12. Classification Performance [1, 15]	• An indirect technique for evaluating the quality of unsupervised representations (e.g. feature extraction; FCN score). See also the GAN Quality Index (GQI) [41].
	13. Boundary Distortion [42]	• Measures diversity of generated samples and covariate shift using classification methods.
	14. Number of Statistically-Different Bins (NDB) [43]	• Given two sets of samples from the same distribution, the number of samples that fall into a given bin should be the same up to sampling noise
	15. Image Retrieval Performance [44]	• Measures the distributions of distances to the nearest neighbors of some query images (i.e. diversity)
	16. Generative Adversarial Metric (GAM) [31]	• Compares two GANs by having them engaged in a battle against each other by swapping discriminators or generators. $p(\mathbf{x} y=1; M_1)/p(\mathbf{x} y=1; M_2) = (p(y=1 \mathbf{x}; D_1)p(\mathbf{x}; G_2))/(p(y=1 \mathbf{x}; D_2)p(\mathbf{x}; G_1))$
	17. Tournament Win Rate and Skill Rating [45]	• Implements a tournament in which a player is either a discriminator that attempts to distinguish between real and fake data or a generator that attempts to fool the discriminators into accepting fake data as real.
	18. Normalized Relative Discriminative Score (NRDS) [32]	• Compares n GANs based on the idea that if the generated samples are closer to real ones, more epochs would be needed to distinguish them from real samples.
	19. Adversarial Accuracy and Divergence [46]	• Adversarial Accuracy. Computes the classification accuracies achieved by the two classifiers, one trained on real data and another on generated data, on a labeled validation set to approximate $P_g(y \mathbf{x})$ and $P_r(y \mathbf{x})$. Adversarial Divergence: Computes $\mathbb{K}L(P_g(y \mathbf{x}), P_r(y \mathbf{x}))$
	20. Geometry Score [47]	• Compares geometrical properties of the underlying data manifold between real and generated data.
	21. Reconstruction Error [48]	• Measures the reconstruction error (e.g. L_2 norm) between a test image and its closest generated image by optimizing for z (i.e. $\min_{\mathbf{z}} \ G(\mathbf{z}) - \mathbf{x}^{(test)}\ ^2$)
	22. Image Quality Measures [49, 50, 51]	• Evaluates the quality of generated images using measures such as SSIM, PSNR, and sharpness difference
	23. Low-level Image Statistics [52, 53]	• Evaluates how similar low-level statistics of generated images are to those of natural scenes in terms of mean power spectrum, distribution of random filter responses, contrast distribution, etc.
	24. Precision, Recall and F_1 score [23]	• These measures are used to quantify the degree of overfitting in GANs, often over toy datasets.
Qualitative	1. Nearest Neighbors	• To detect overfitting, generated samples are shown next to their nearest neighbors in the training set
	2. Rapid Scene Categorization [18]	• In these experiments, participants are asked to distinguish generated samples from real images in a short presentation time (e.g. 100 ms); i.e. real v.s fake
	3. Preference Judgment [54, 55, 56, 57]	• Participants are asked to rank models in terms of the fidelity of their generated images (e.g. pairs, triples)
	4. Mode Drop and Collapse [58, 59]	• Over datasets with known modes (e.g. a GMM or a labeled dataset), modes are computed as by measuring the distances of generated data to mode centers
	5. Network Internals [1, 60, 61, 62, 63, 64]	• Regards exploring and illustrating the internal representation and dynamics of models (e.g. space continuity) as well as visualizing learned features

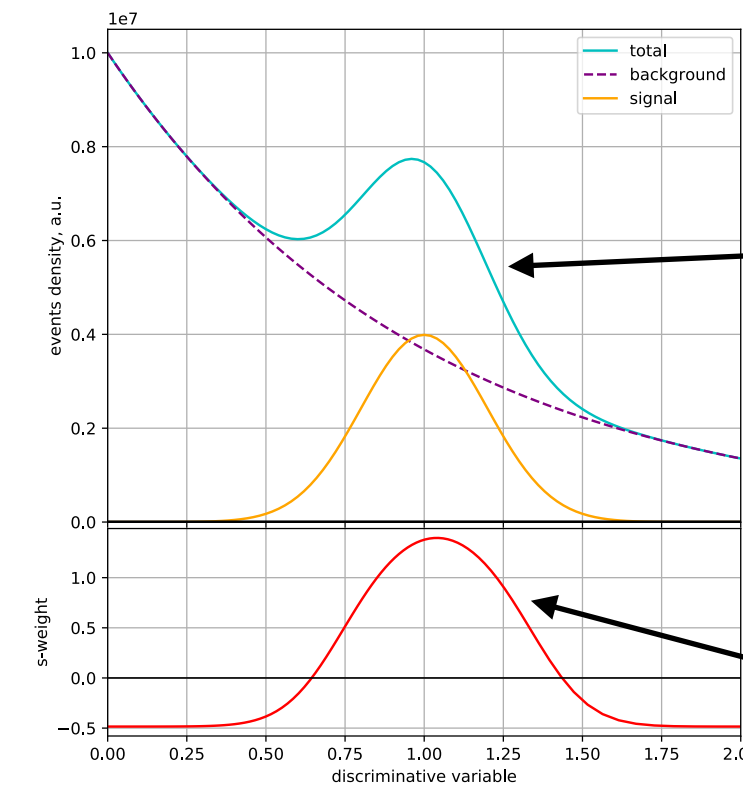
Training on "noisy" data

- Training on real data rather than simulated events allows to avoid calibration-related issues
- Special datasets were selected with a selection strategy carefully studied to avoid bias on the PID variables[†]
- The data is contaminated with background events
- Signal can be statistically extracted with a maximum likelihood fit
- *sPlot* technique^{††} was utilized to extract the signal distributions of the output variables
 - loss function was weighted with the s-weights

[†]Aaij, R., Anderlini, L., Benson, S. et al. Selection and processing of calibration samples to measure the particle identification performance of the LHCb experiment in Run 2. EPJ Techn Instrum **6**, 1 (2019)

^{††}M. Pivk and F. R. Le Diberder, *sPlot*: A statistical tool to unfold data distributions, Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment **555** (2005) 356–369

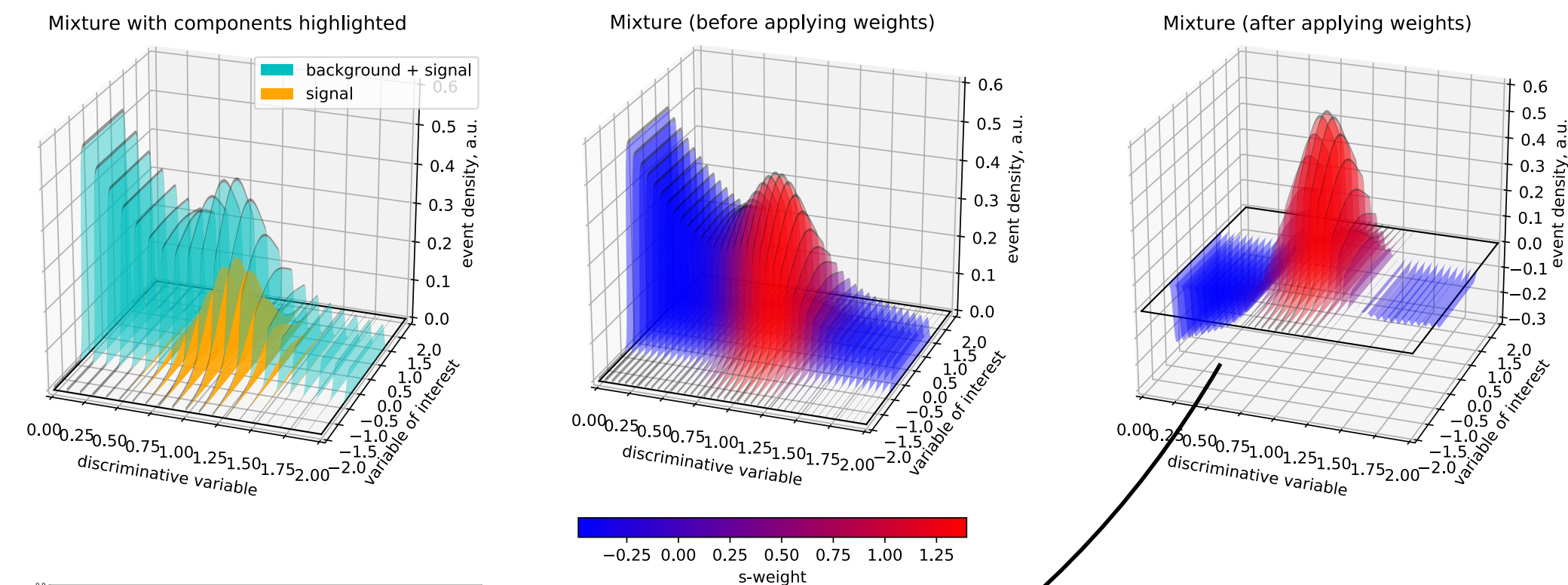
sPlot demonstrated on toy data



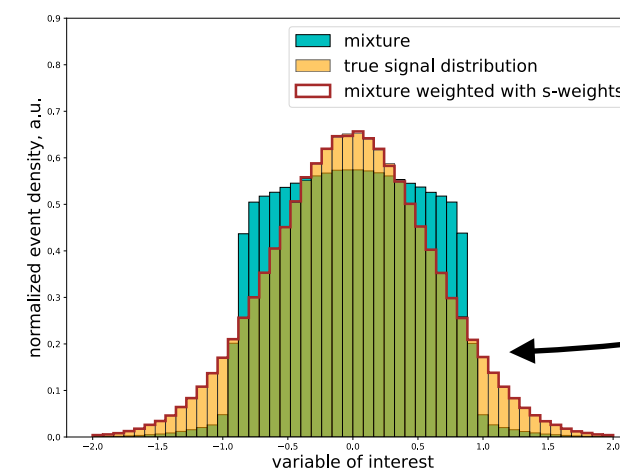
signal + background mixture that can be statistically disentangled (e.g. with a maximum likelihood fit)

sPlot technique allows one to reconstruct per-component distributions of variables that are independent from the discriminative variable

it does so using weights ("s-weights") calculated from the relative probabilities of components



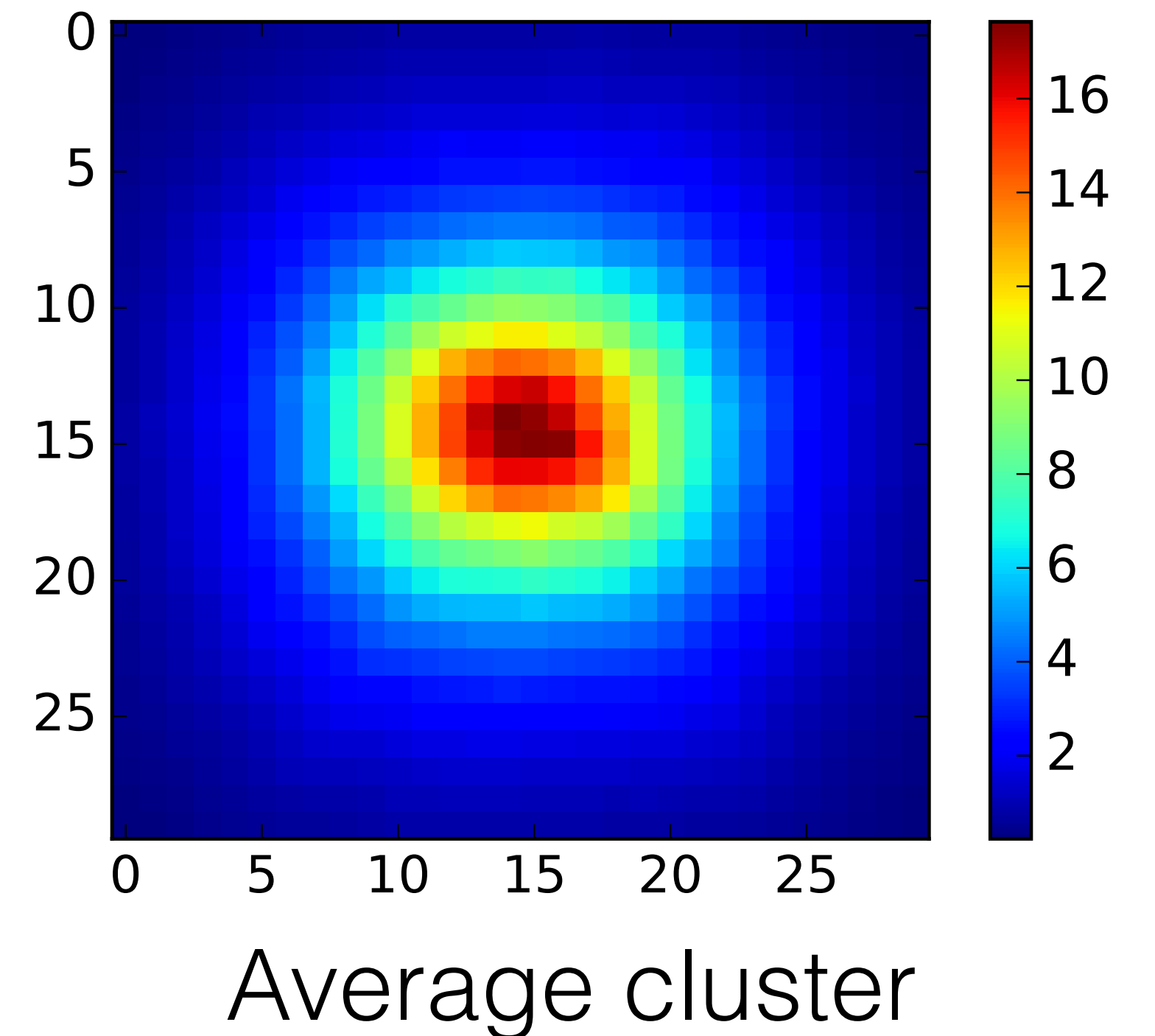
marginalizing out the discriminative variable reconstructs the distribution for the variable of interest



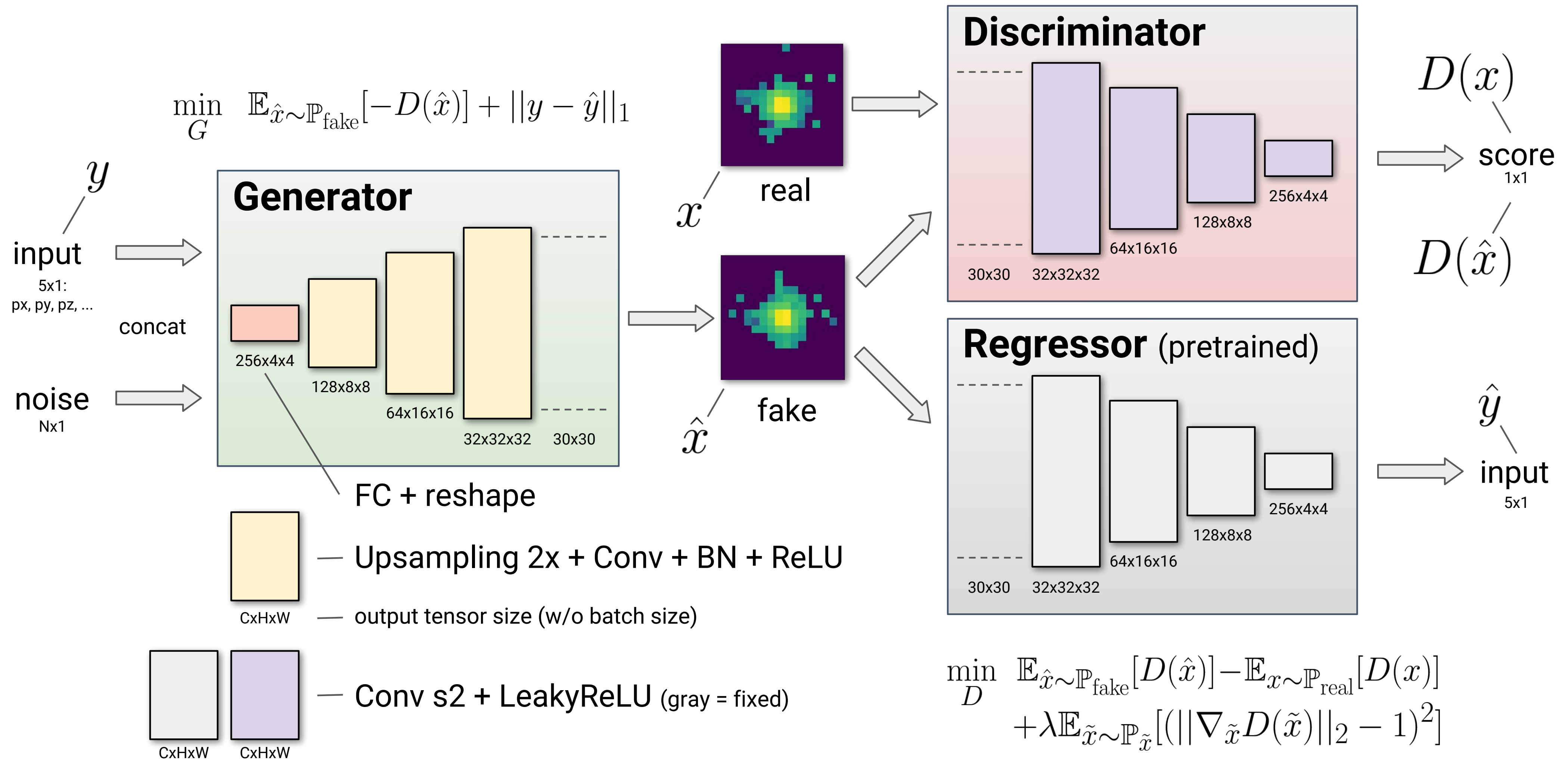
FAST LHCb CALORIMETER SIMULATION

Setup

- LHCb inspired calorimeter in GEANT4
 - 30x30 cells
- 5 conditional parameters per particle
 - 3D momentum
 - 2D coordinate
- Electrons from particle gun shot at 1x1 cm square at the center of the calorimeter face

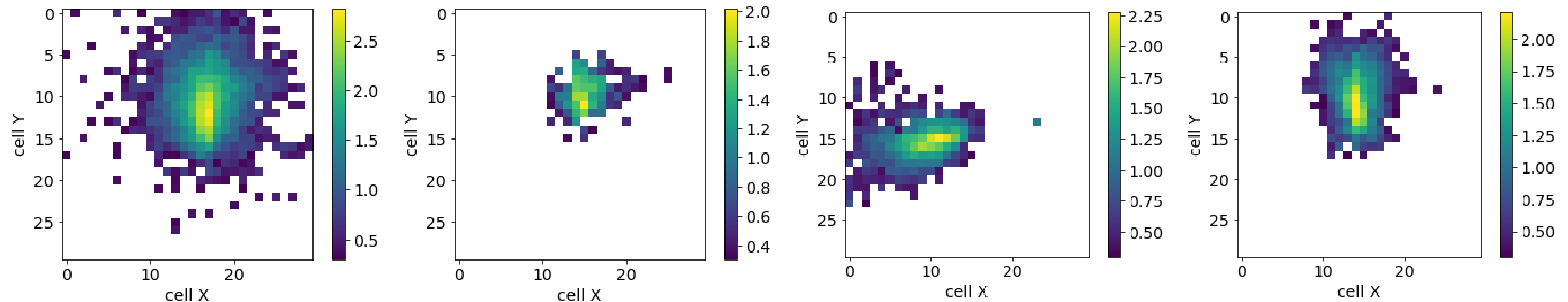
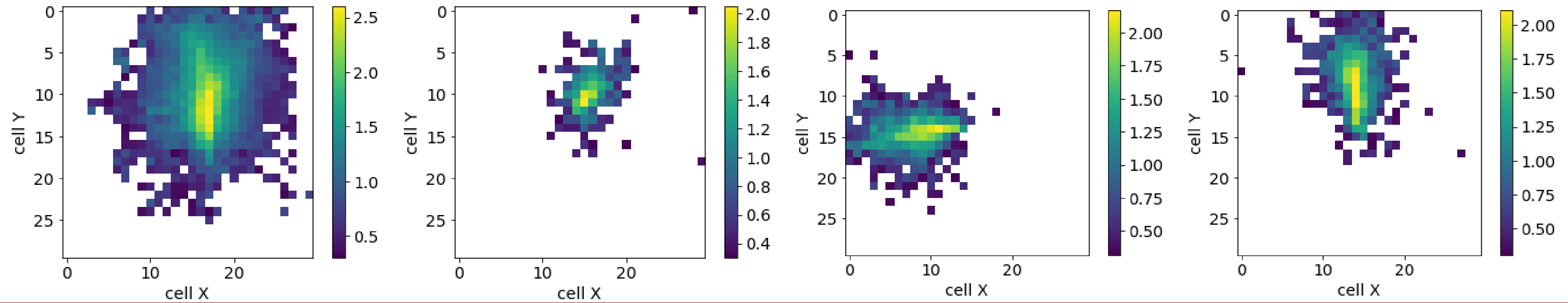


Model architecture



Energy deposits

GEANT4



GAN

(a)

$$E_0 = 63.7 \text{ GeV}$$

(b)

$$E_0 = 6.5 \text{ GeV}$$

(c)

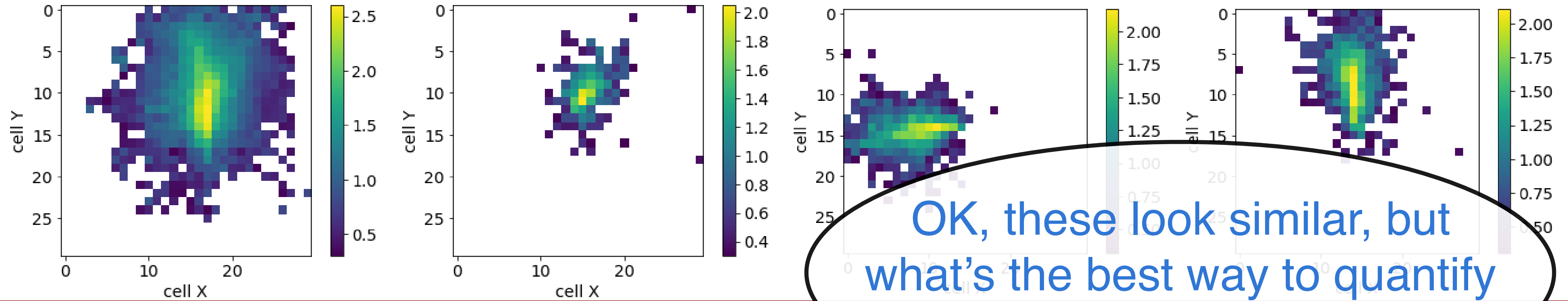
$$E_0 = 15.6 \text{ GeV}$$

(d)

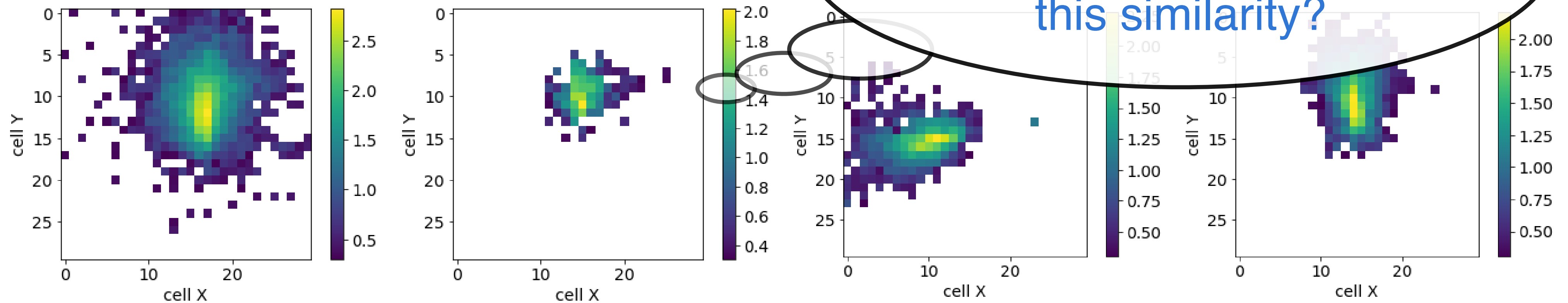
$$E_0 = 15.9 \text{ GeV}$$

Energy deposits

GEANT4



OK, these look similar, but what's the best way to quantify this similarity?



GAN

(a)

$$E_0 = 63.7 \text{ GeV}$$

(b)

$$E_0 = 6.5 \text{ GeV}$$

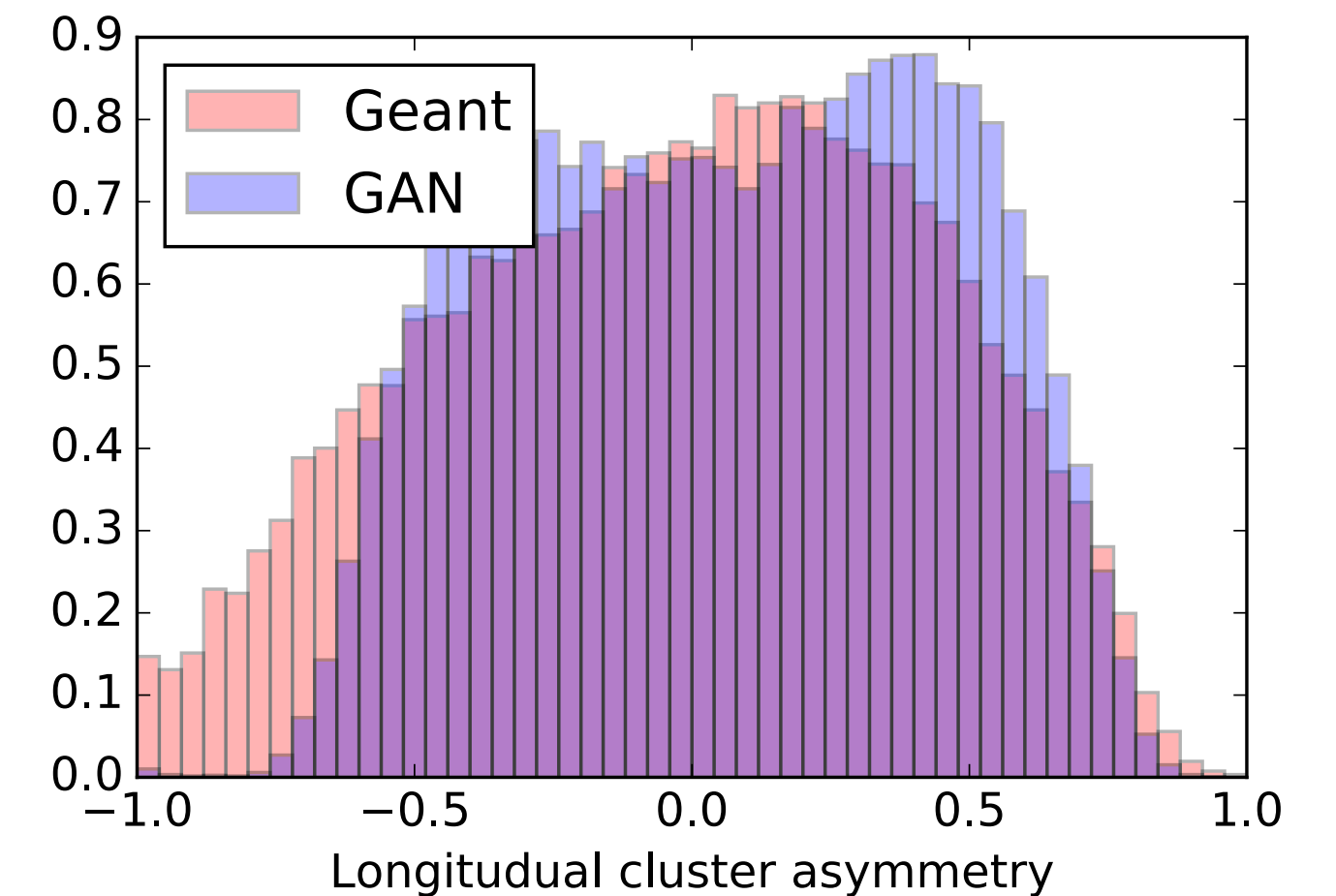
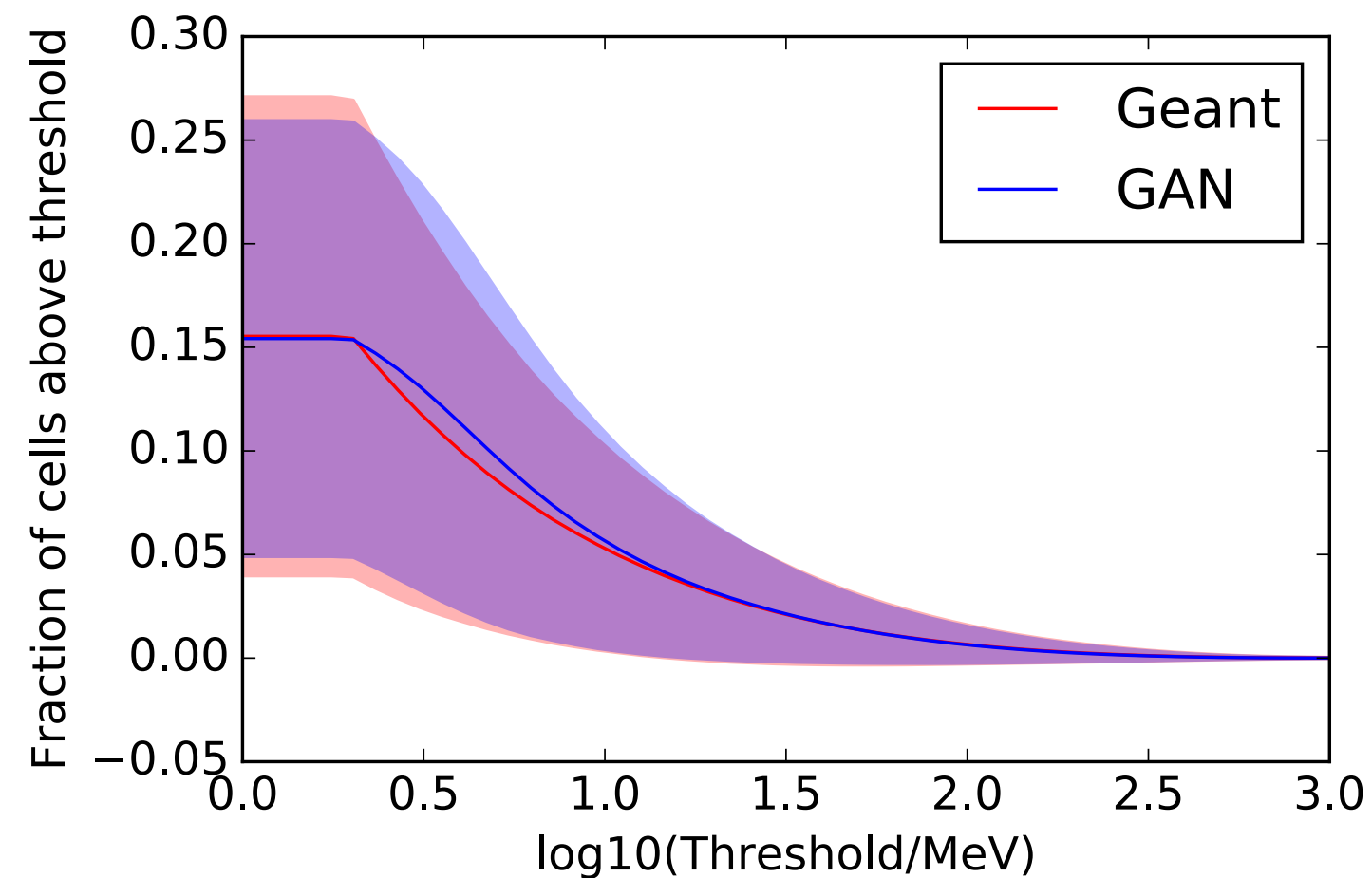
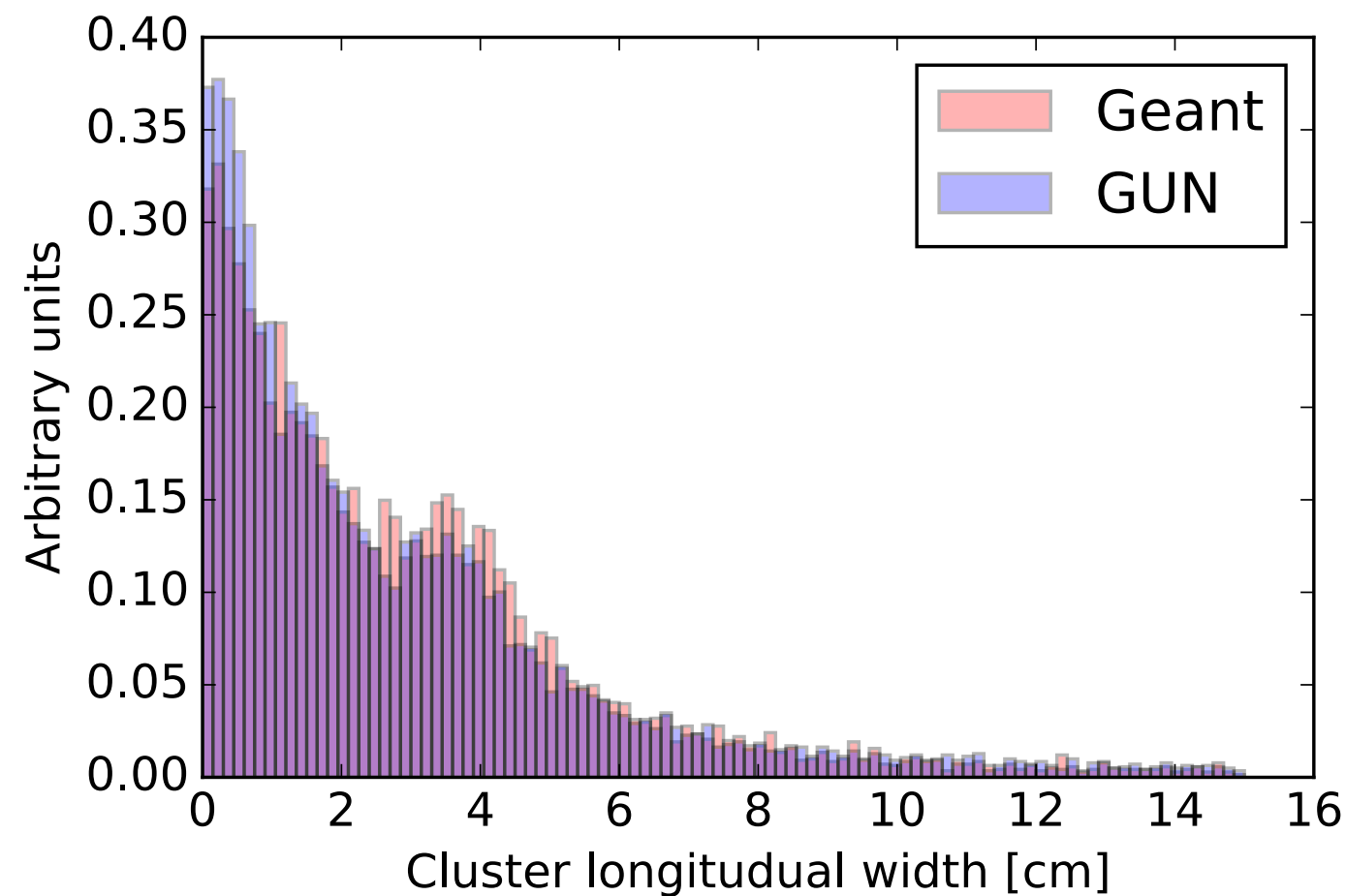
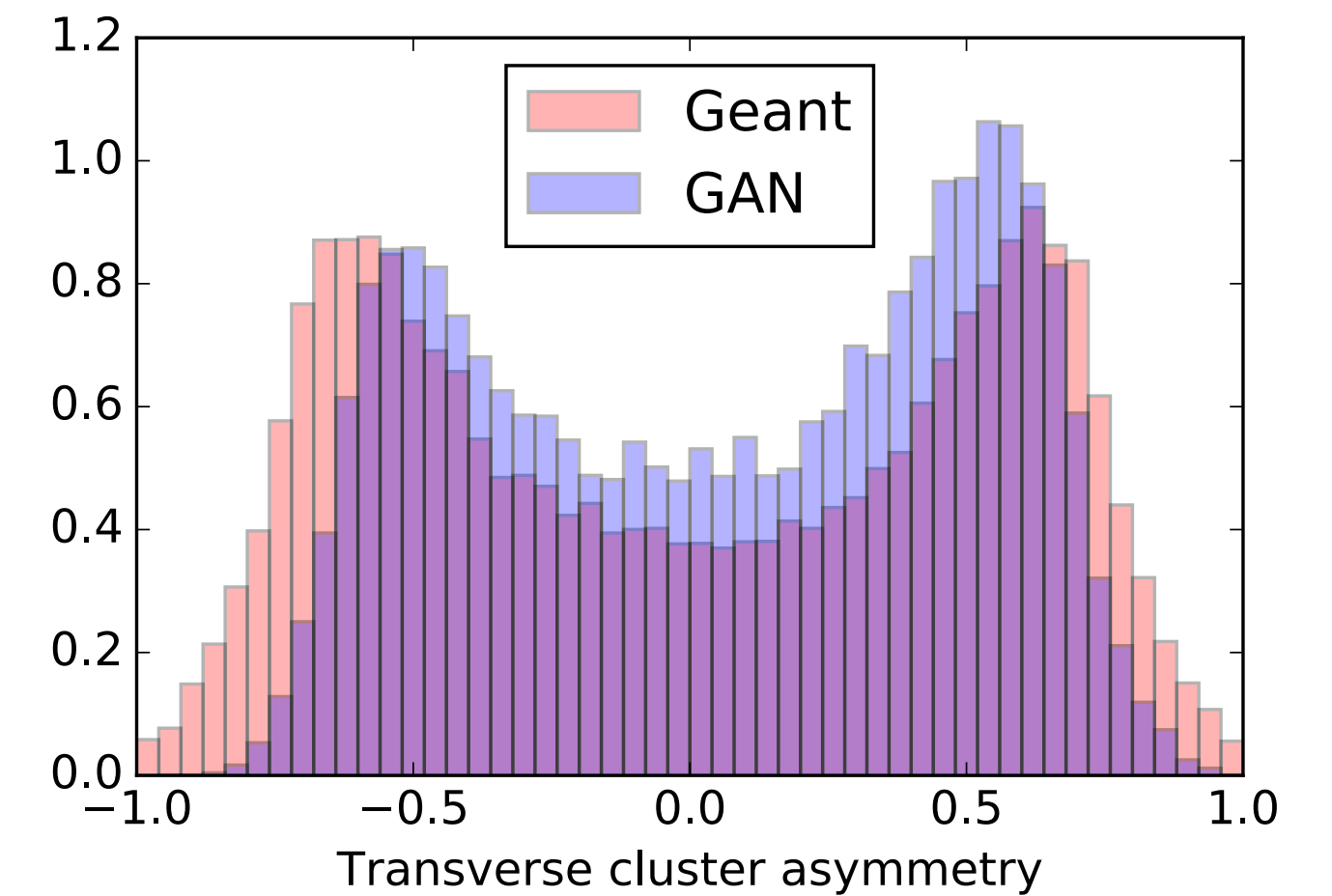
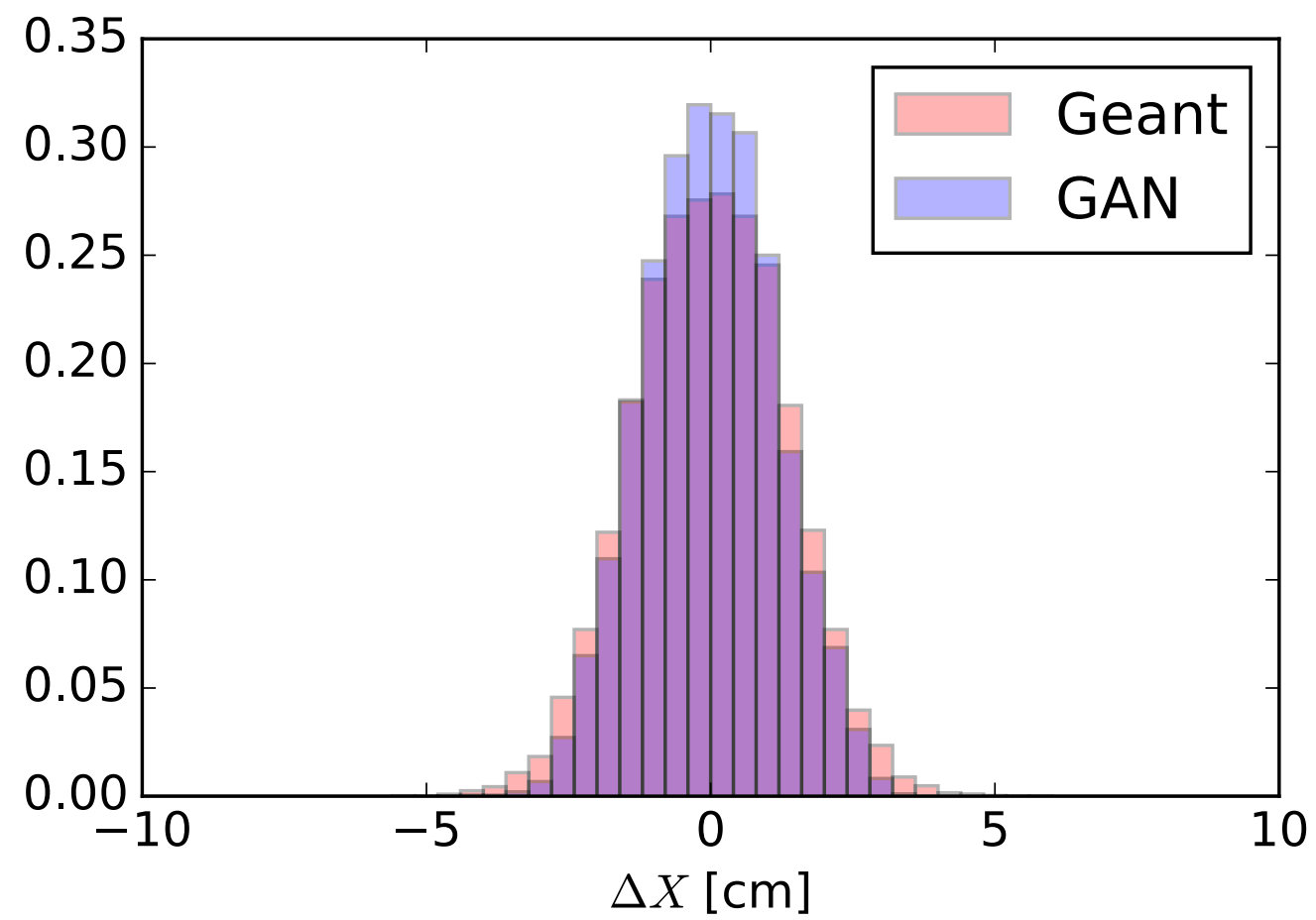
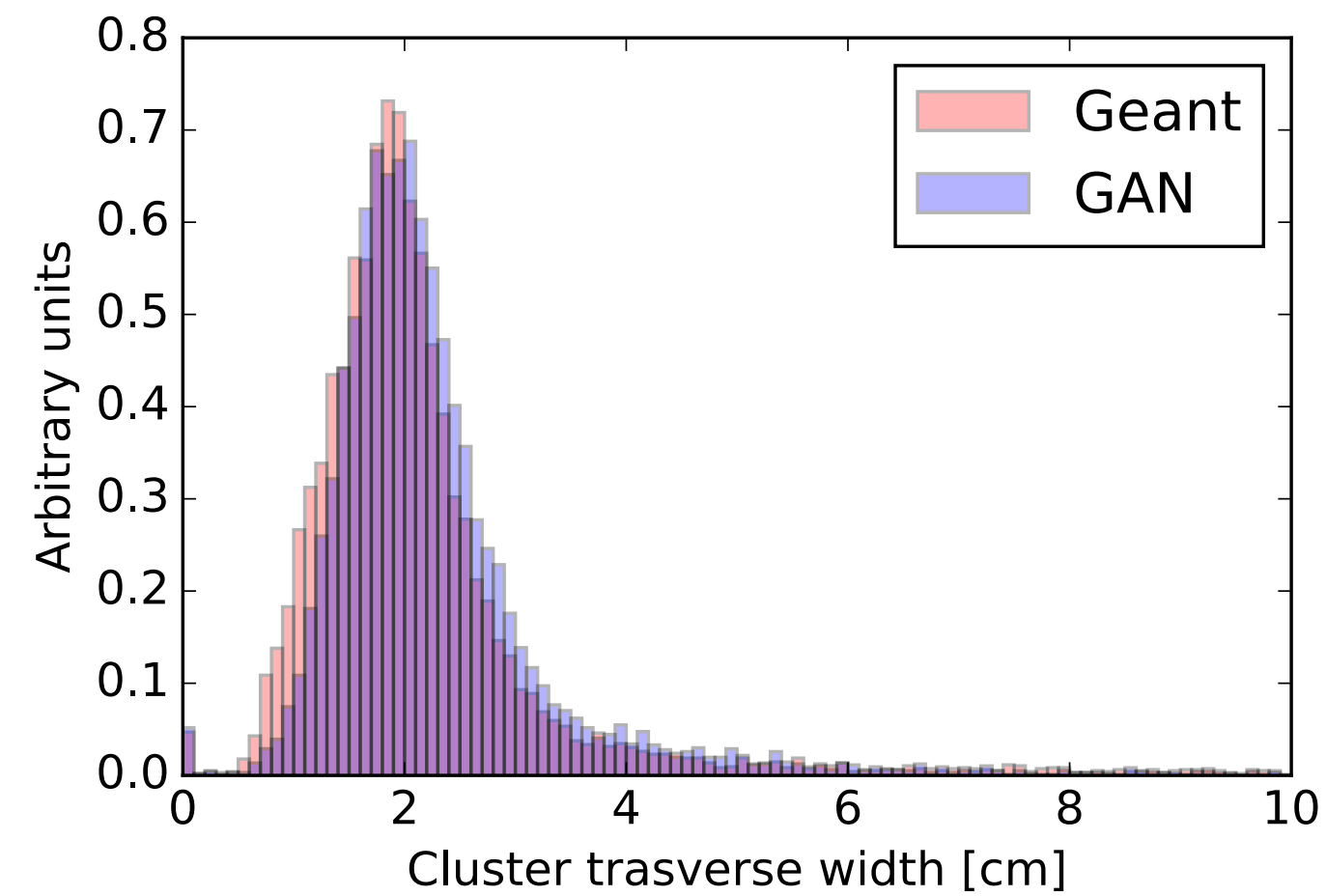
(c)

$$E_0 = 15.6 \text{ GeV}$$

(d)

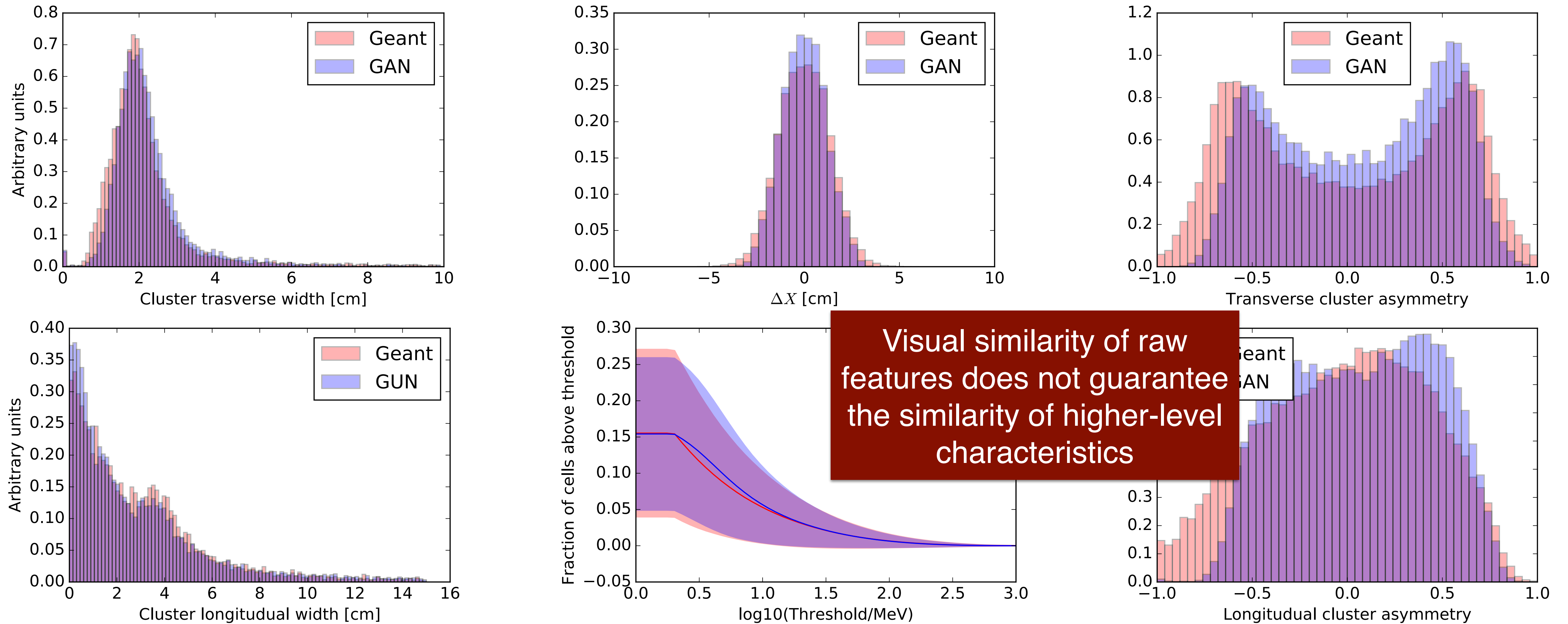
$$E_0 = 15.9 \text{ GeV}$$

Physics-motivated characteristics



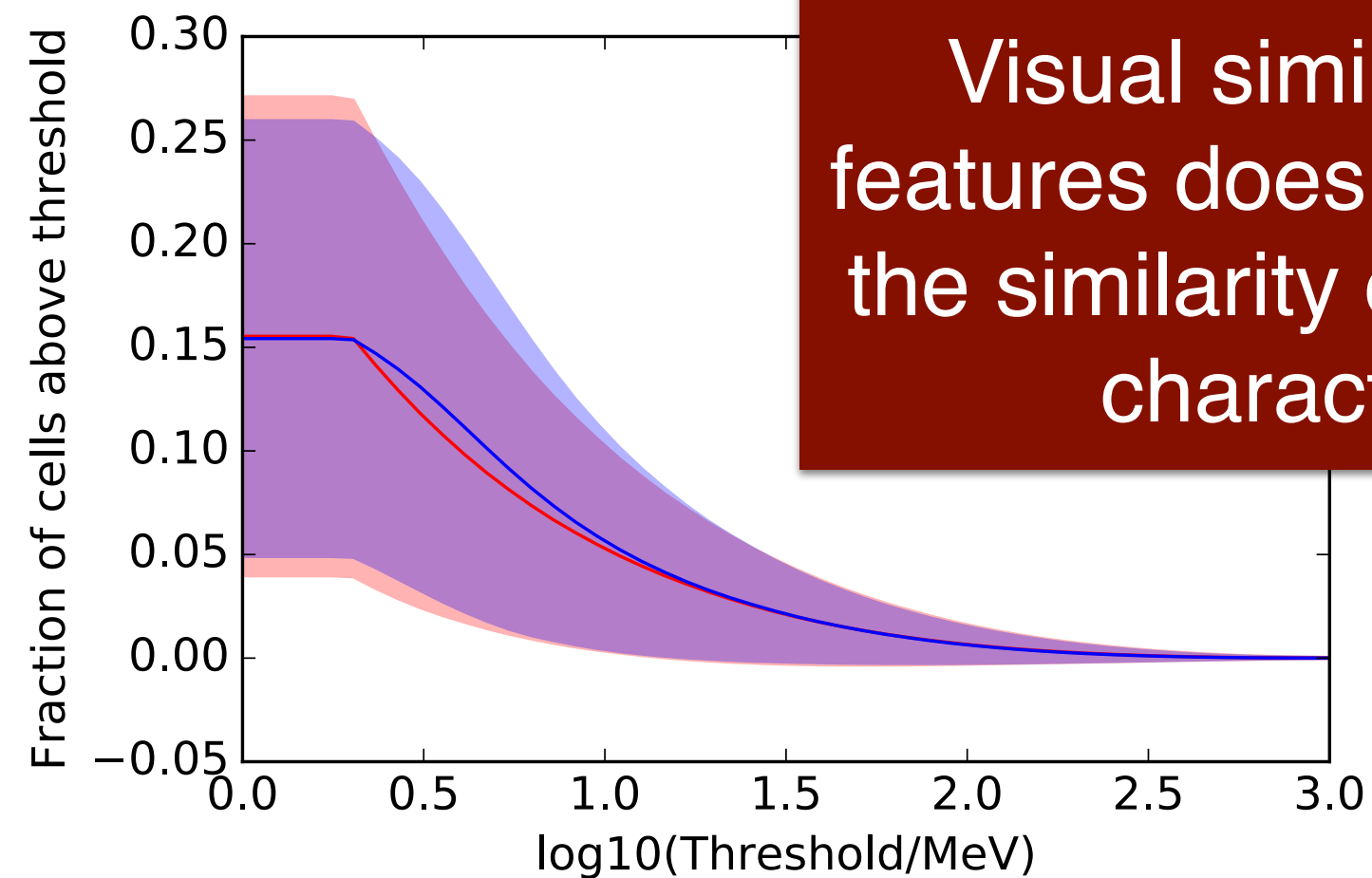
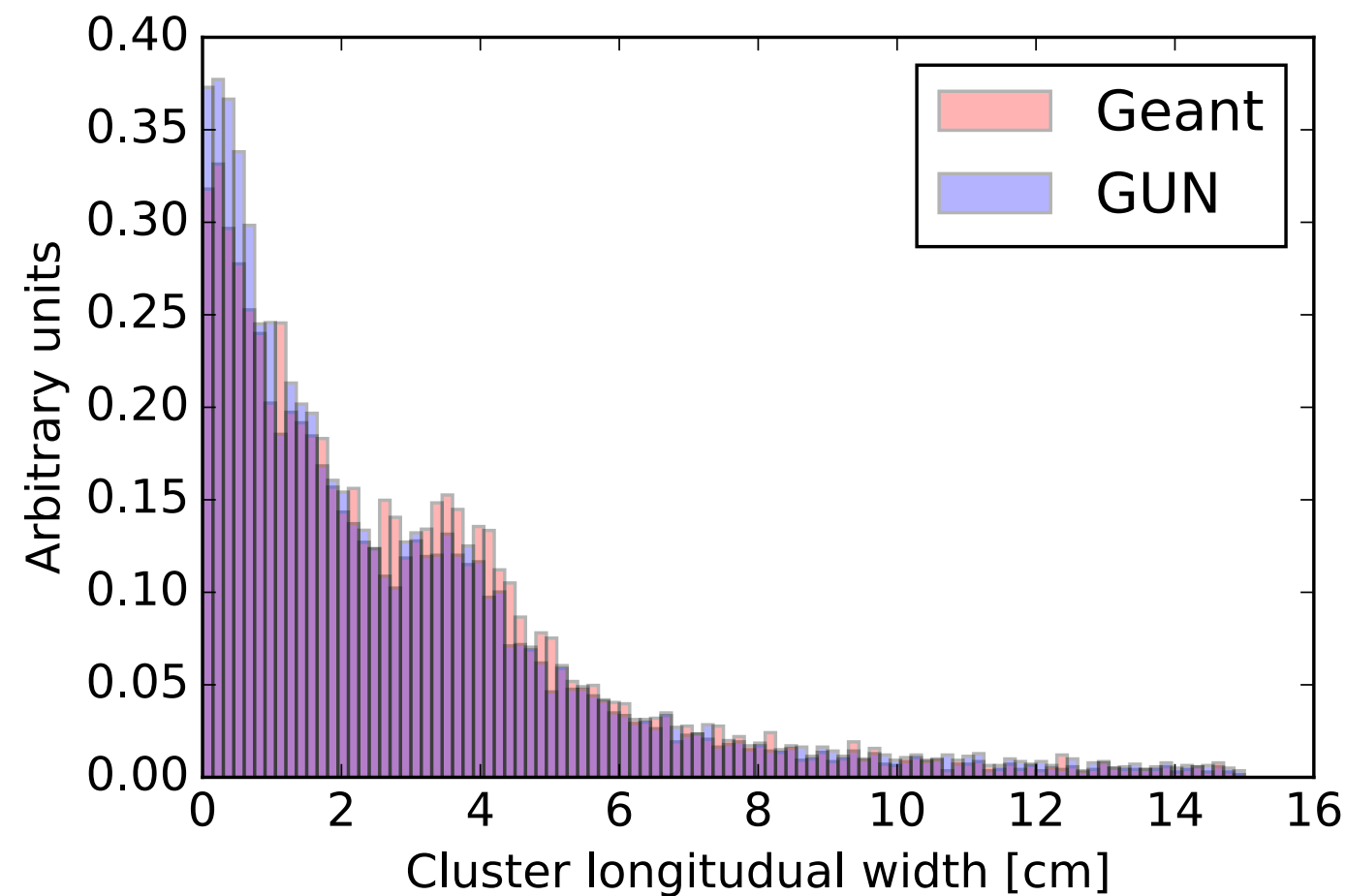
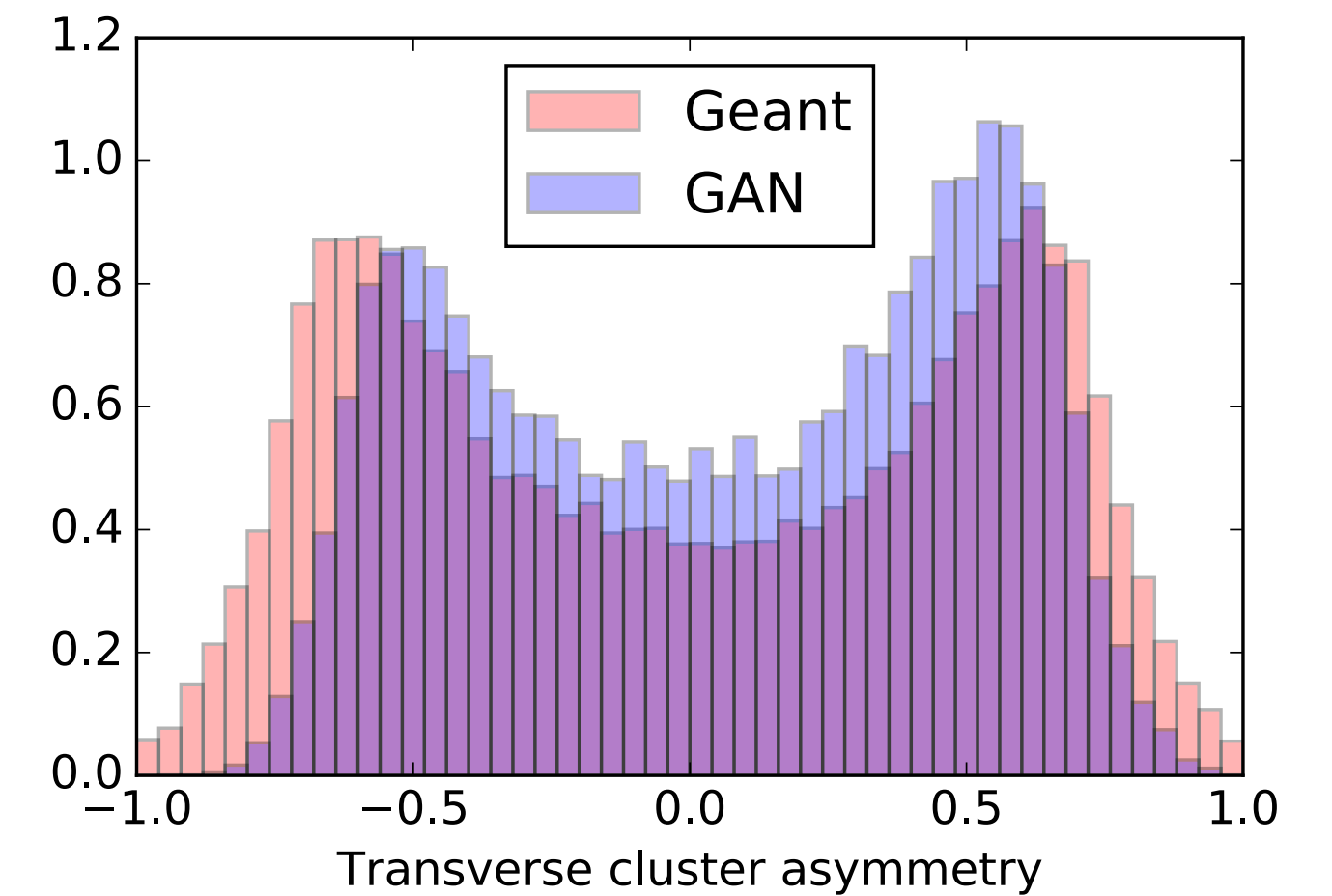
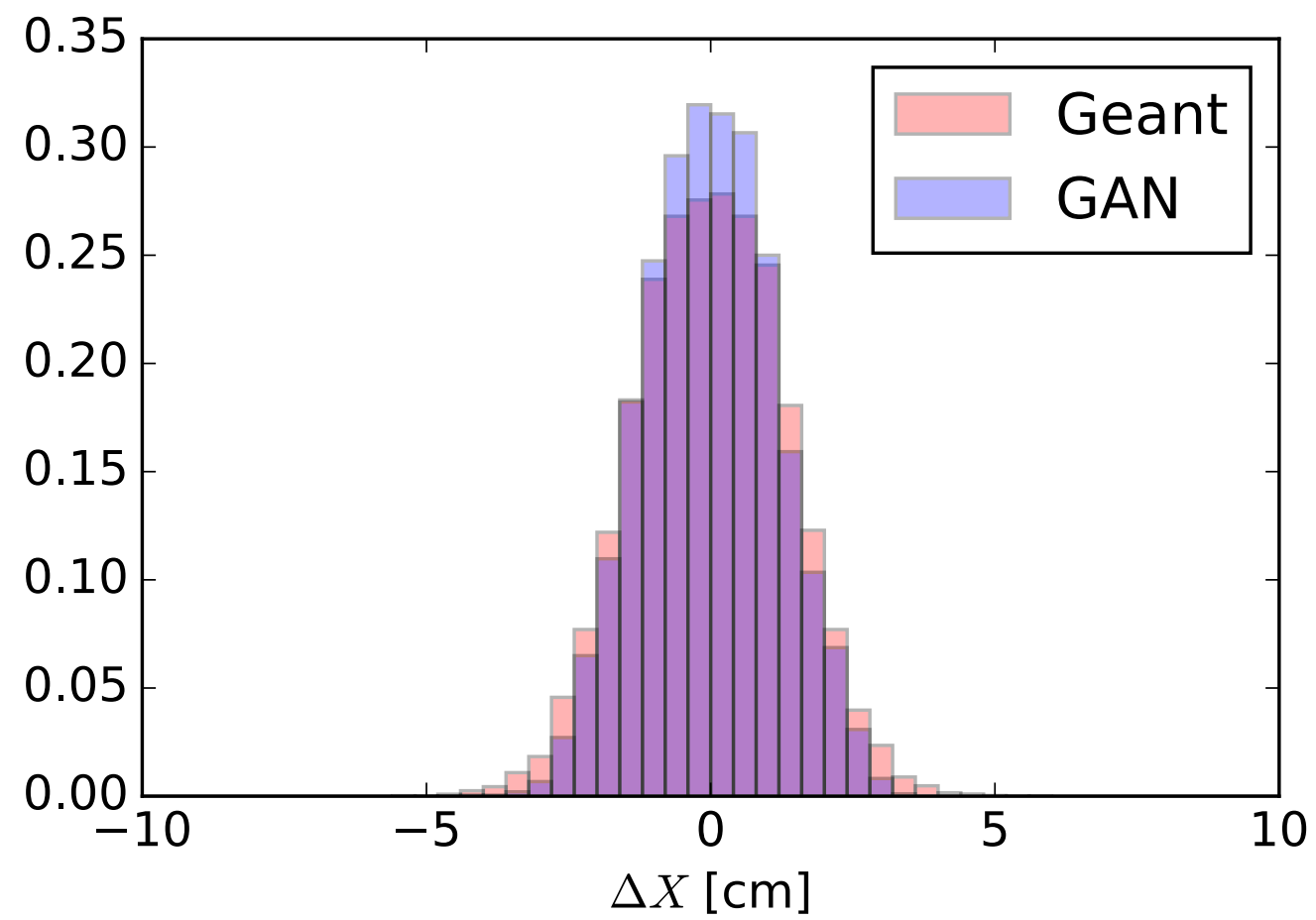
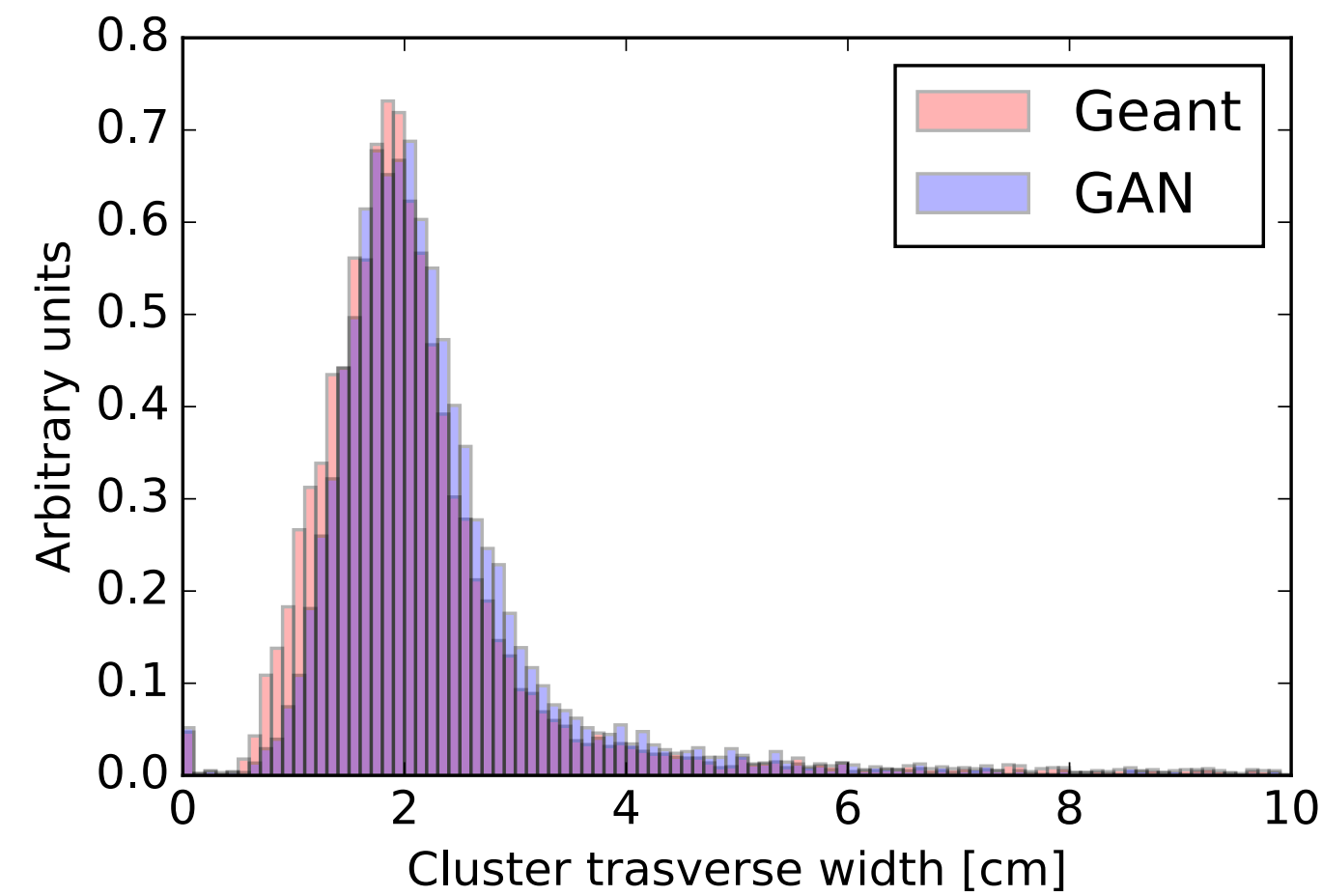
- Idea: convert each cluster to a meaningful single value
 - ▶ Then compare real vs generated distributions of such values

Physics-motivated characteristics



- Idea: convert each cluster to a meaningful single value
 - ▶ Then compare real vs generated distributions of such values

Physics-motivated characteristics



Visual similarity of features does not guarantee the similarity of higher order characteristics

Some of these are reproduced quite well though!

- Idea: convert each cluster to a meaningful single value
 - ▶ Then compare real vs generated distributions of such values

Other HEP uses of GANs

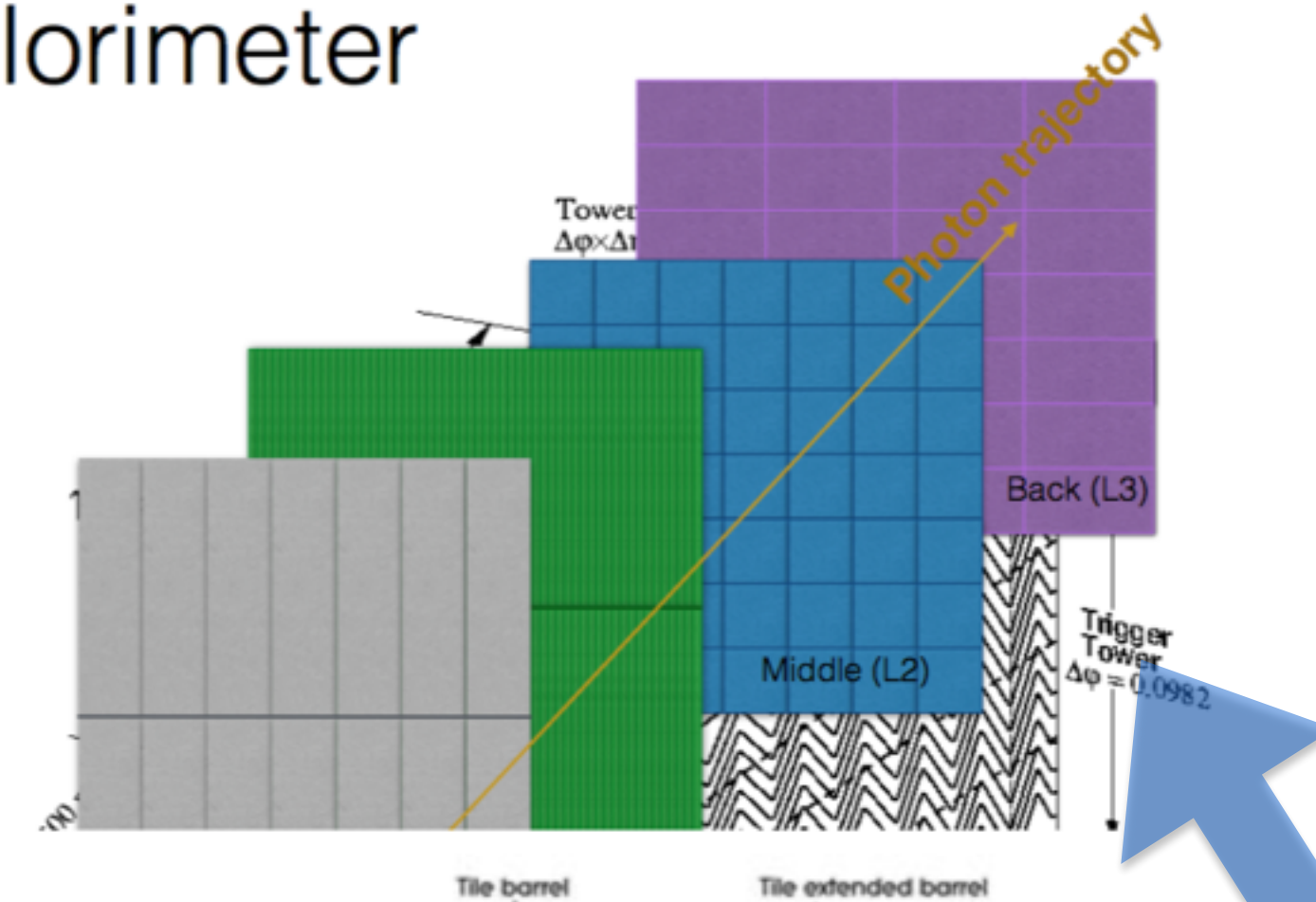
EM calorimeter shower fast simulation in ATLAS

The Calorimeter

2-D Axis: ϕ vs η

Particle goes through 4 layers in this order:

0. **Pre-Sampler** : Some energy deposit
1. **Strips**: Very granular in η ; more energy deposit
2. **Middle**: Thickest layer, maximum energy deposit
3. **Back**: Little Energy deposits



Generating output for 3D calorimeter structure (4 calorimeter layers, 266 output channels)

[https://indico.cern.ch/
event/766872/
contributions/3357991/](https://indico.cern.ch/event/766872/contributions/3357991/)

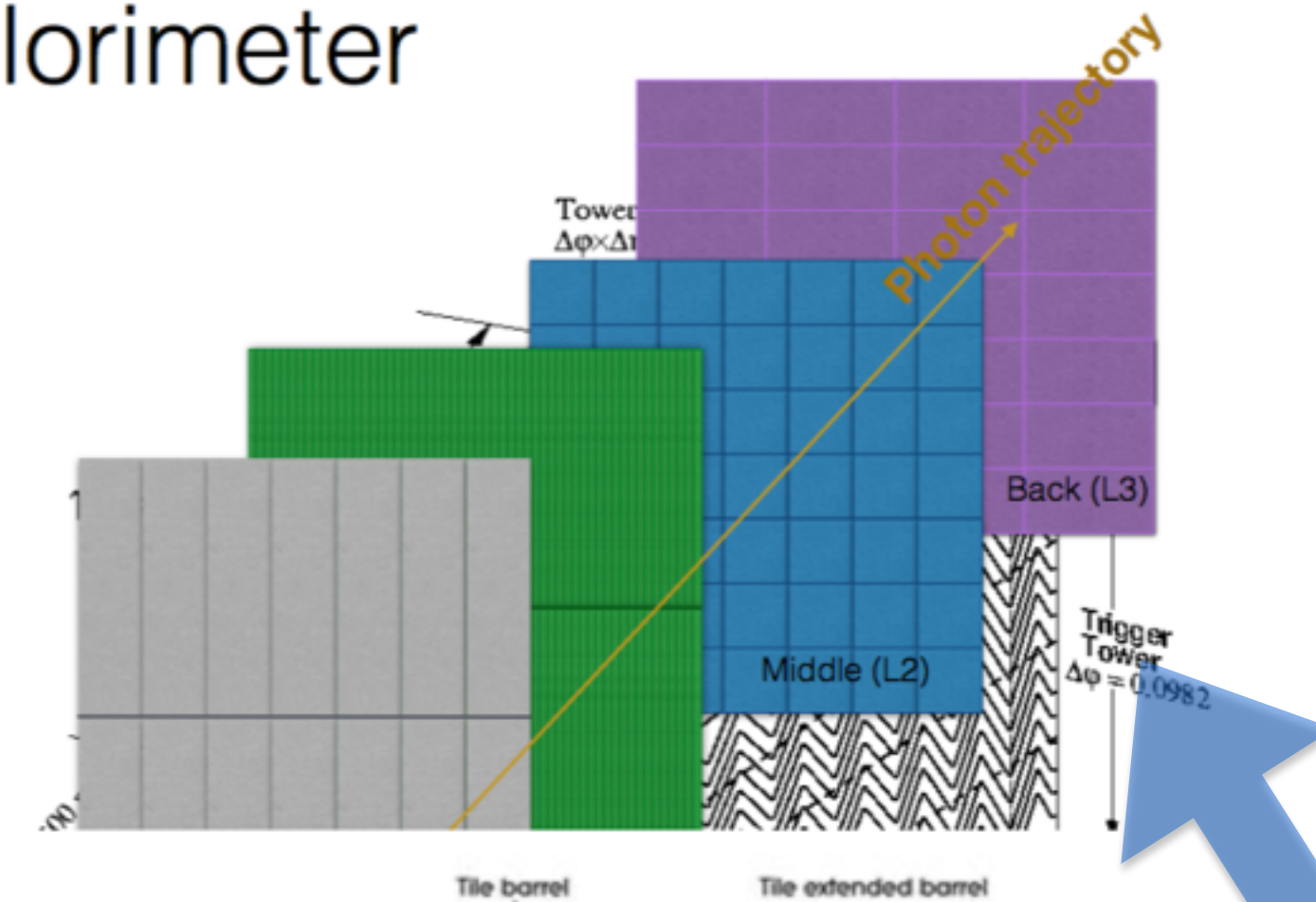
EM calorimeter shower fast simulation in ATLAS

The Calorimeter

2-D Axis: ϕ vs η

Particle goes through 4 layers in this order:

0. **Pre-Sampler** : Some energy deposit
1. **Strips**: Very granular in η ; more energy deposit
2. **Middle**: Thickest layer, maximum energy deposit
3. **Back**: Little Energy deposits



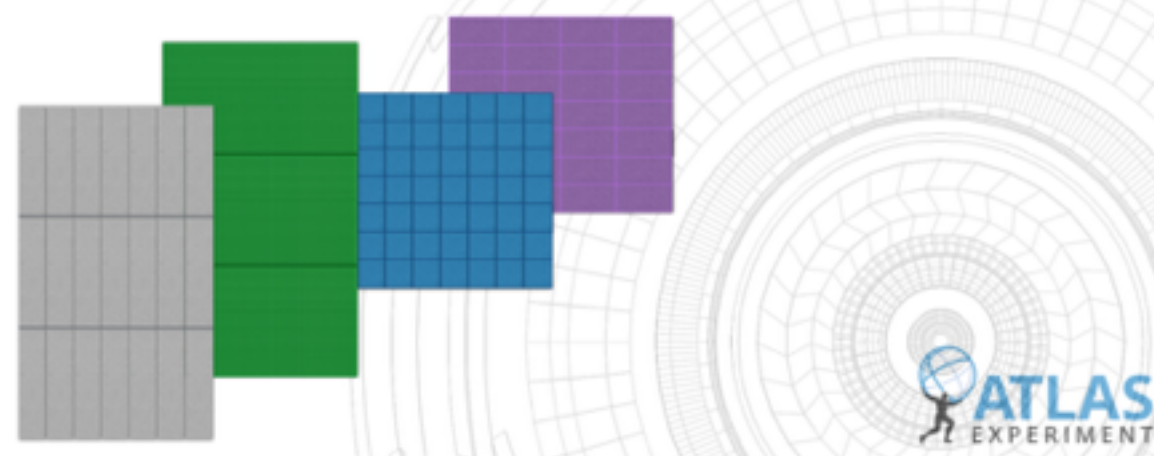
Generating output for 3D calorimeter structure (4 calorimeter layers, 266 output channels)

Deep Generative Models for Fast Simulation in ATLAS

Graeme Stewart (CERN), Aishik Ghosh, David Rousseau (LAL, Orsay), Kyle Cranmer (NYU), Michele Fucci Giannelli, Serena Palazzo (University of Edinburgh), Stefan Gadatsch, Tobias Golling, Johnny Raine, Dalia Salamani, Siava Voloshynovskiy (UniGe), Gilles Louppe (ULiège)

on behalf of ATLAS

IML Workshop
17 April 2019



Evaluating GAN performance through physics observables ($\Delta\eta$, $\Delta\phi$, $E_{\text{sim}}/E_{\text{truth}}$, etc.)

<https://indico.cern.ch/event/766872/contributions/3357991/>

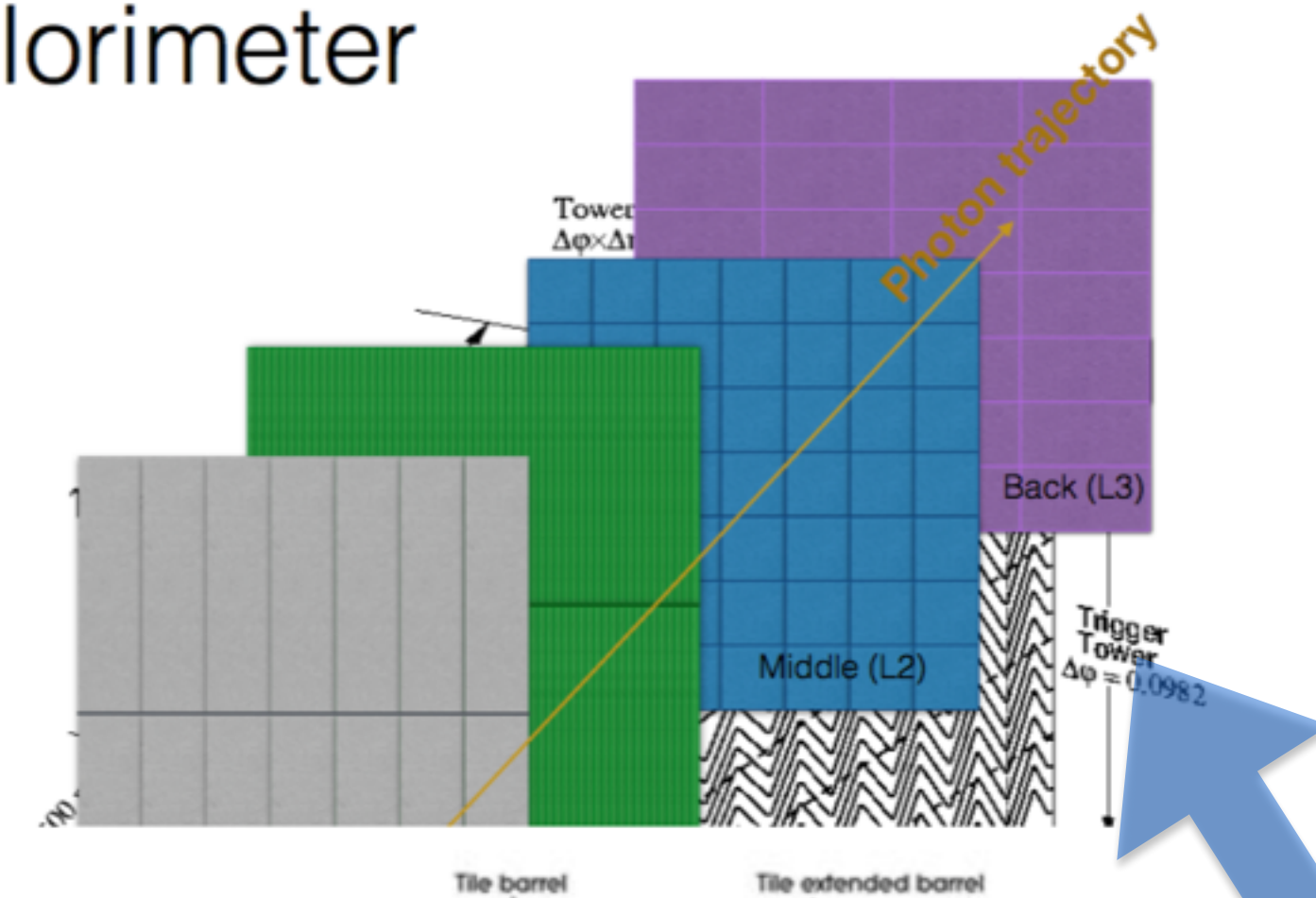
EM calorimeter shower fast simulation in ATLAS

The Calorimeter

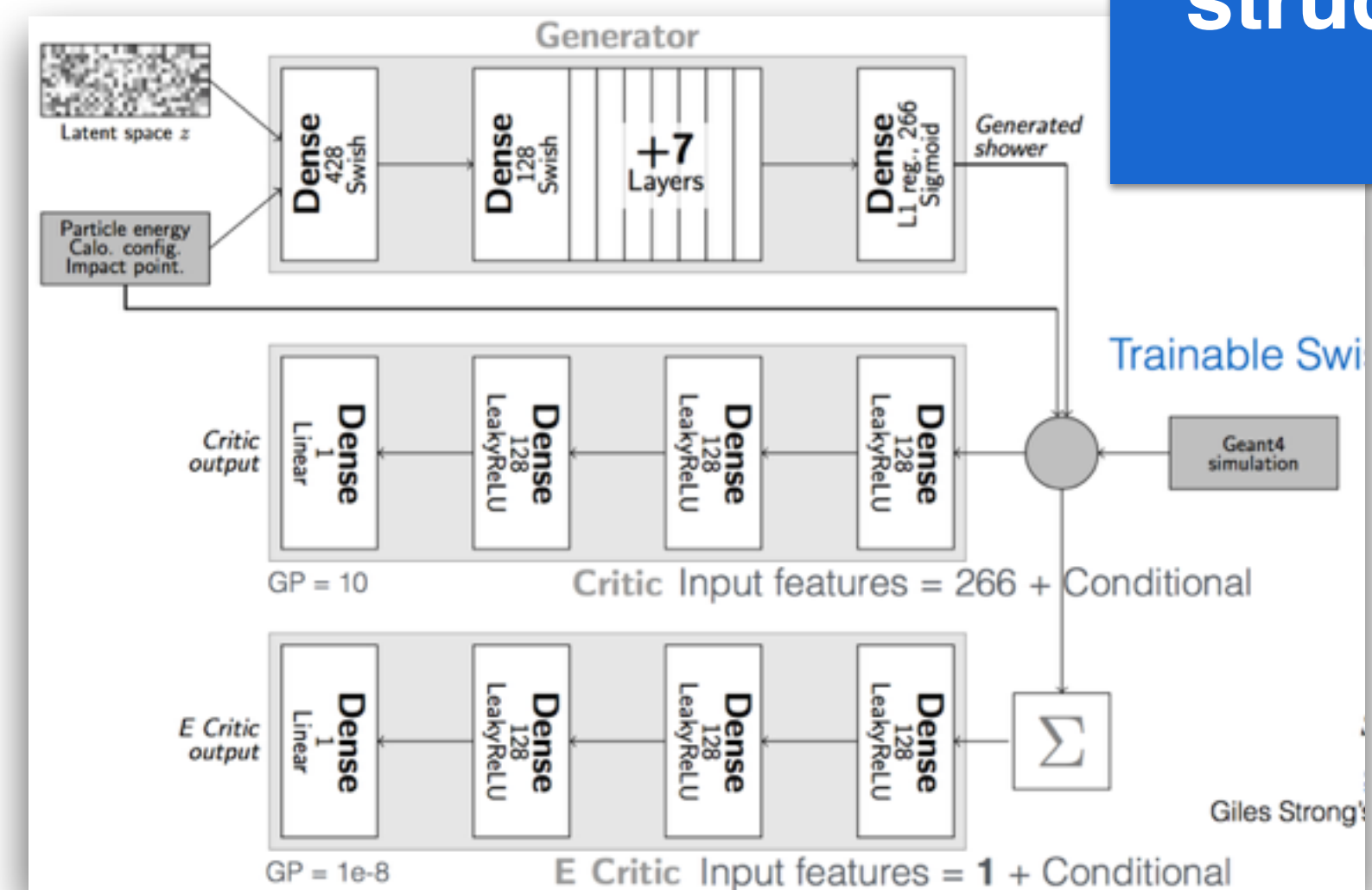
2-D Axis: ϕ vs η

Particle goes through 4 layers in this order:

0. **Pre-Sampler** : Some energy deposit
1. **Strips**: Very granular in η ; more energy deposit
2. **Middle**: Thickest layer, maximum energy deposit
3. **Back**: Little Energy deposits



Generating output for 3D calorimeter structure (4 calorimeter layers, 266 output channels)

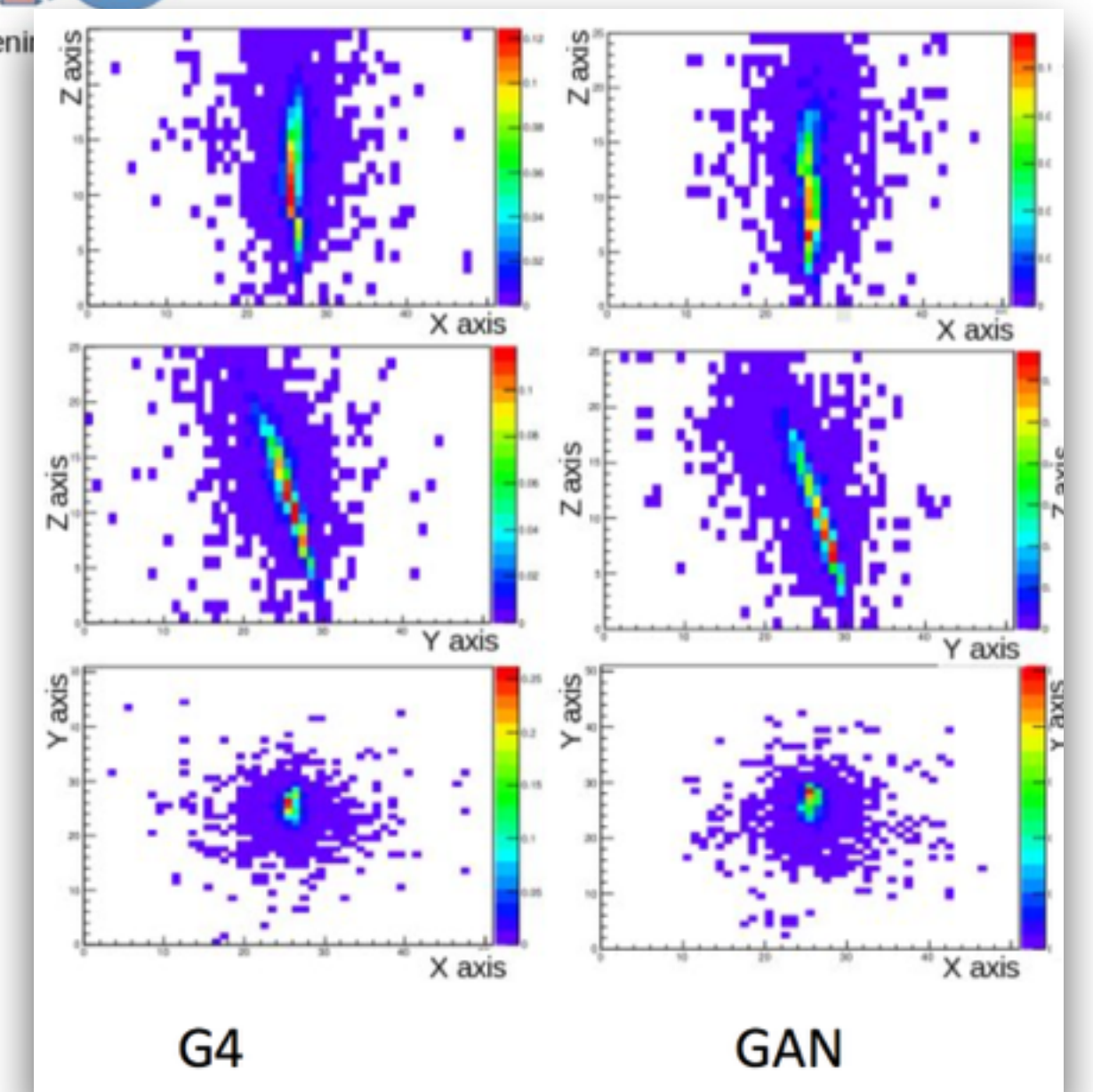
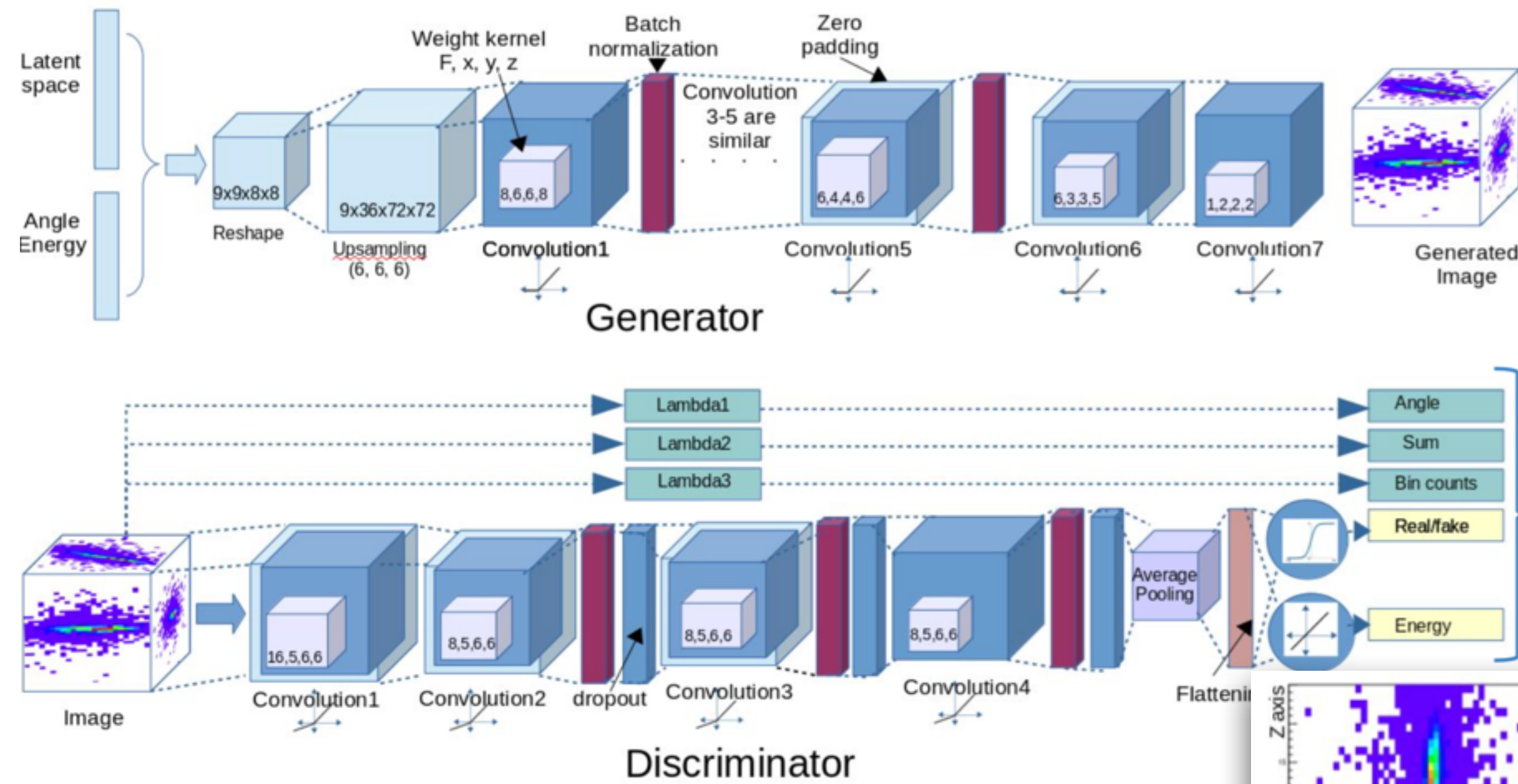


Two critics trained simultaneously (additional critic to look at the total energy distribution)

Evaluating GAN performance through physics observables ($\Delta\eta$, $\Delta\phi$, E_{sim}/E_{truth} , etc.)

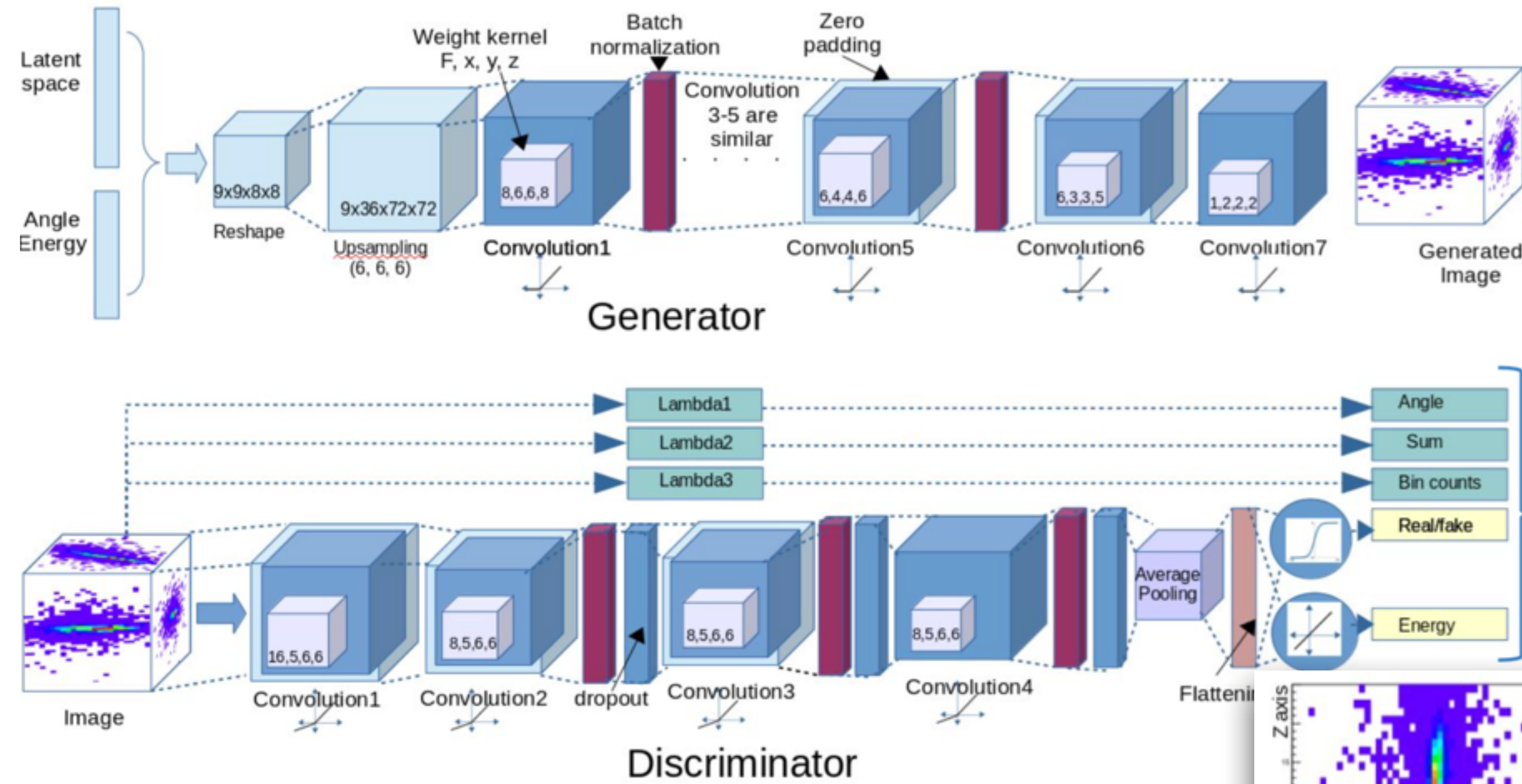
<https://indico.cern.ch/event/766872/contributions/3357991/>

High granularity calorimeter fast simulation for CLIC

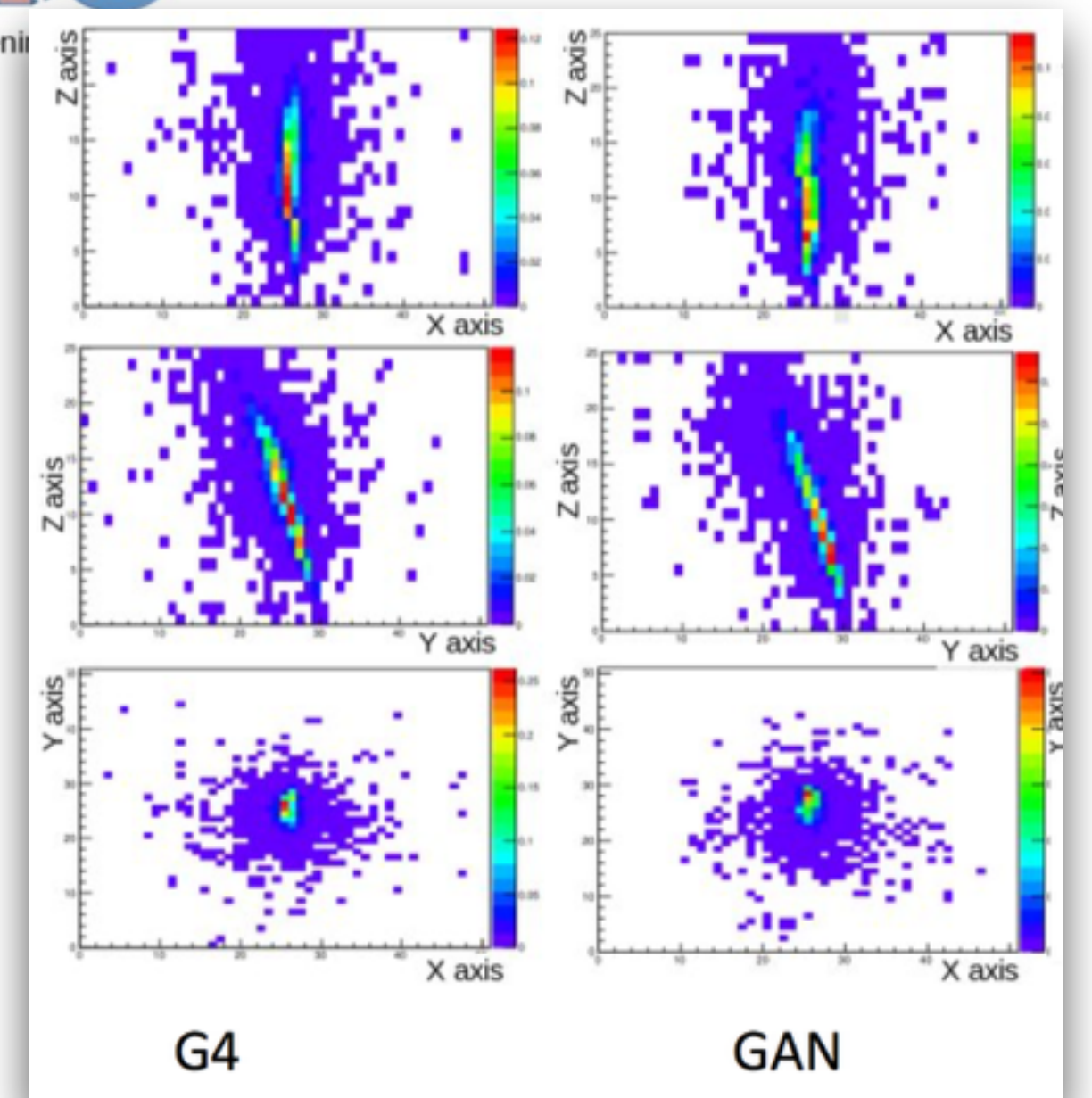


[https://indico.cern.ch/
event/766872/
contributions/3357987/](https://indico.cern.ch/event/766872/contributions/3357987/)

High granularity calorimeter fast simulation for CLIC

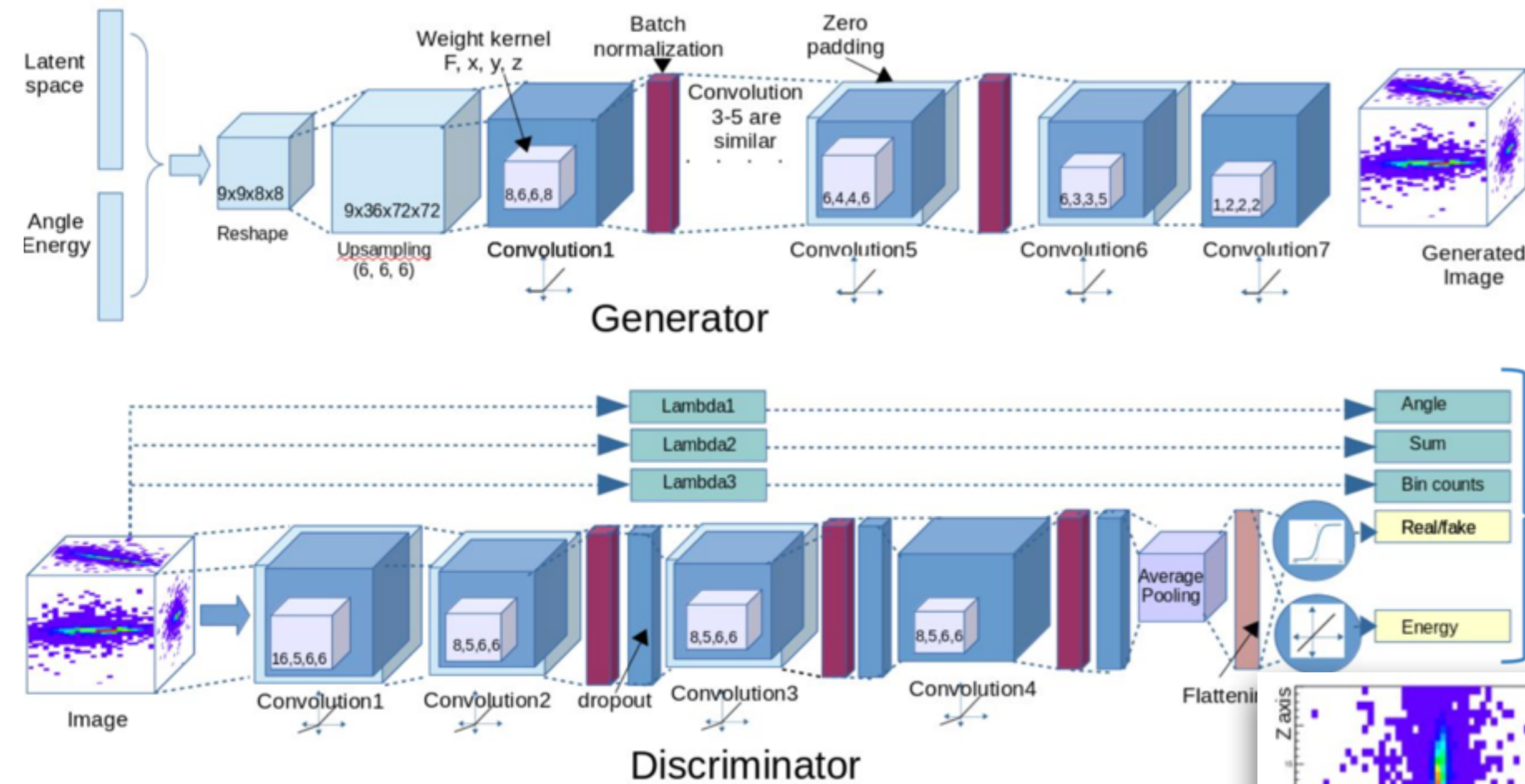
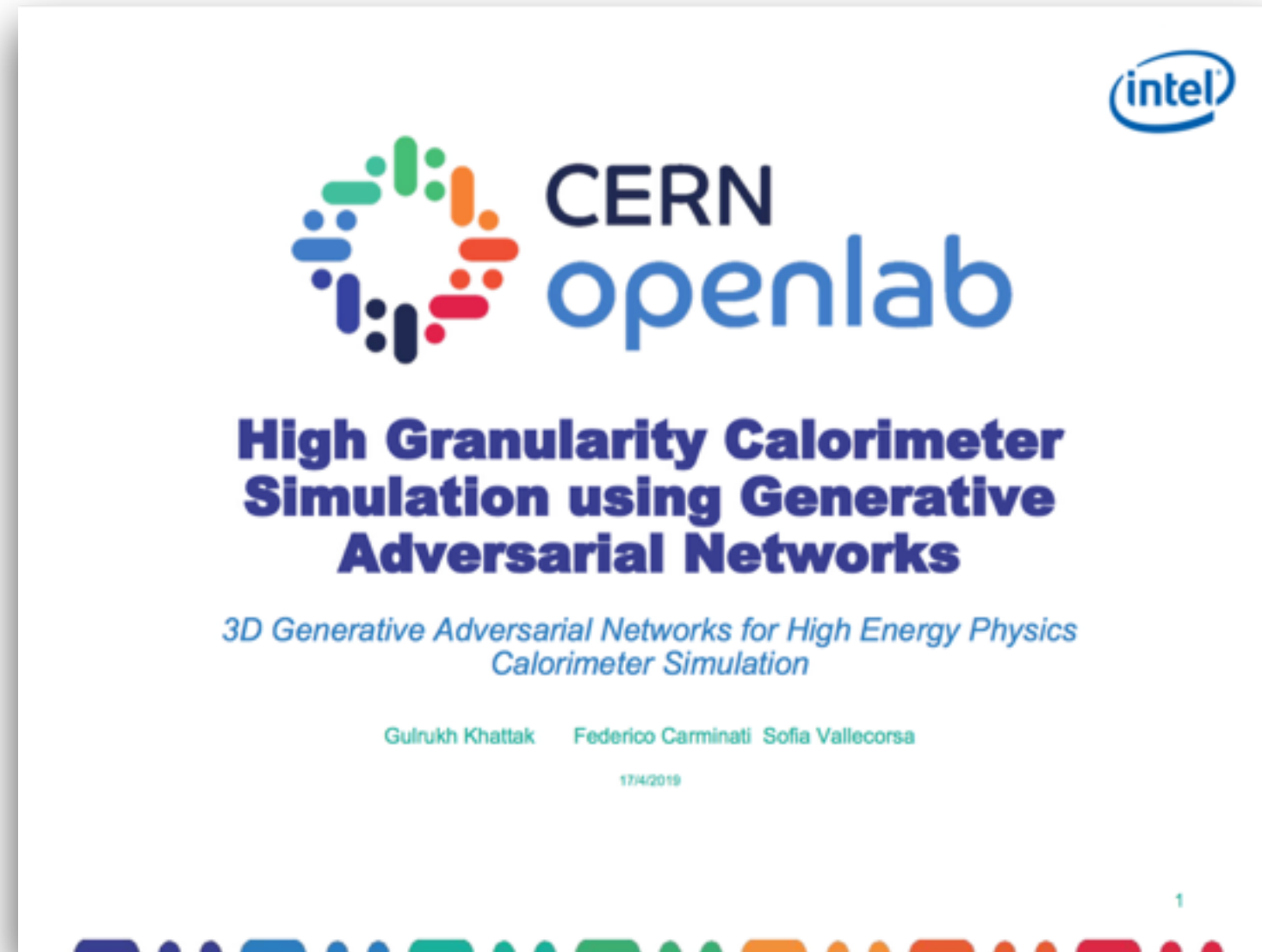


Evaluating with >200 physics observables!!!



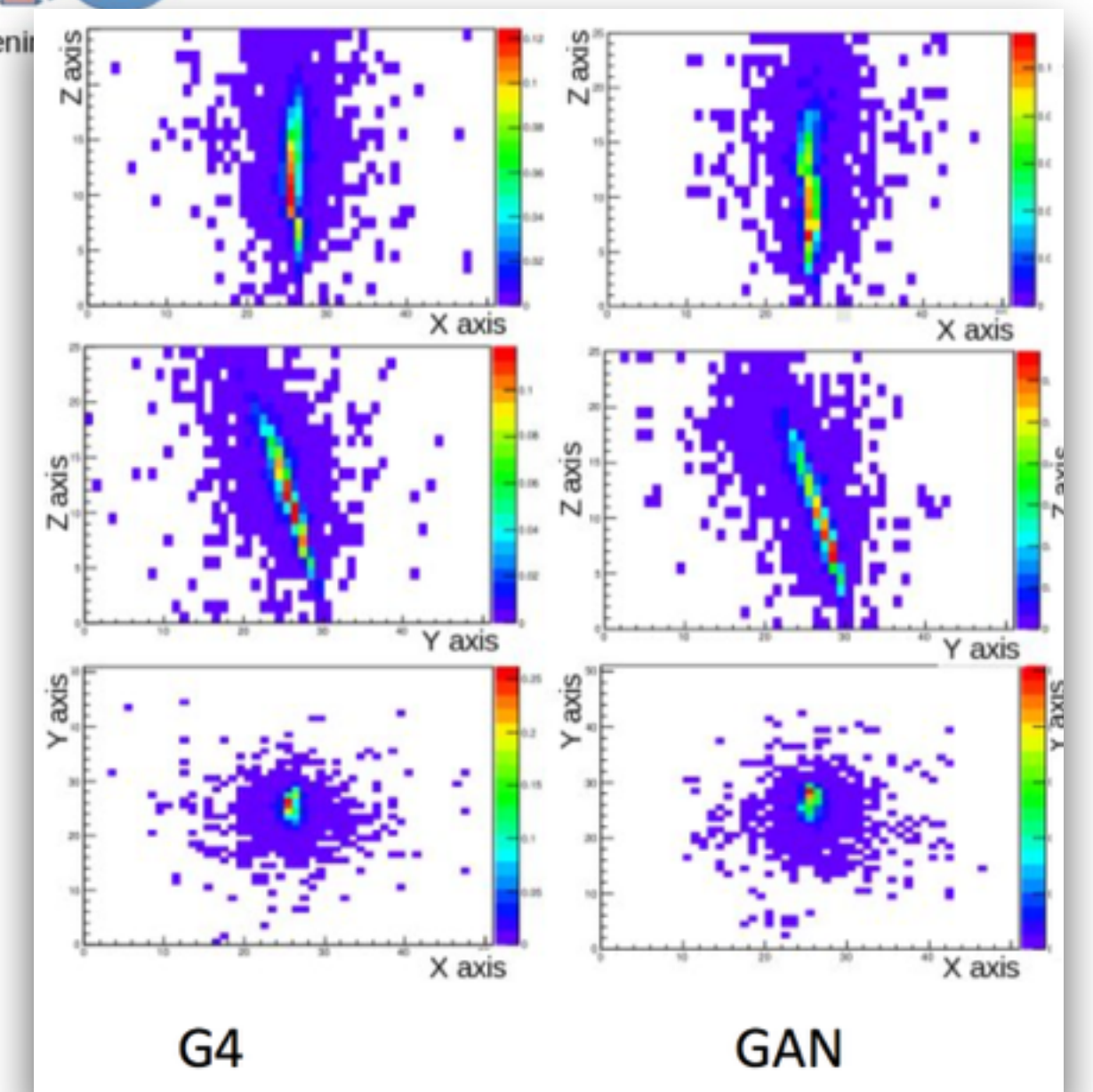
[https://indico.cern.ch/
event/766872/
contributions/3357987/](https://indico.cern.ch/event/766872/contributions/3357987/)

High granularity calorimeter fast simulation for CLIC



Evaluating with >200 physics observables!!!

+ using structural similarity index (SSIM) on calorimeter images

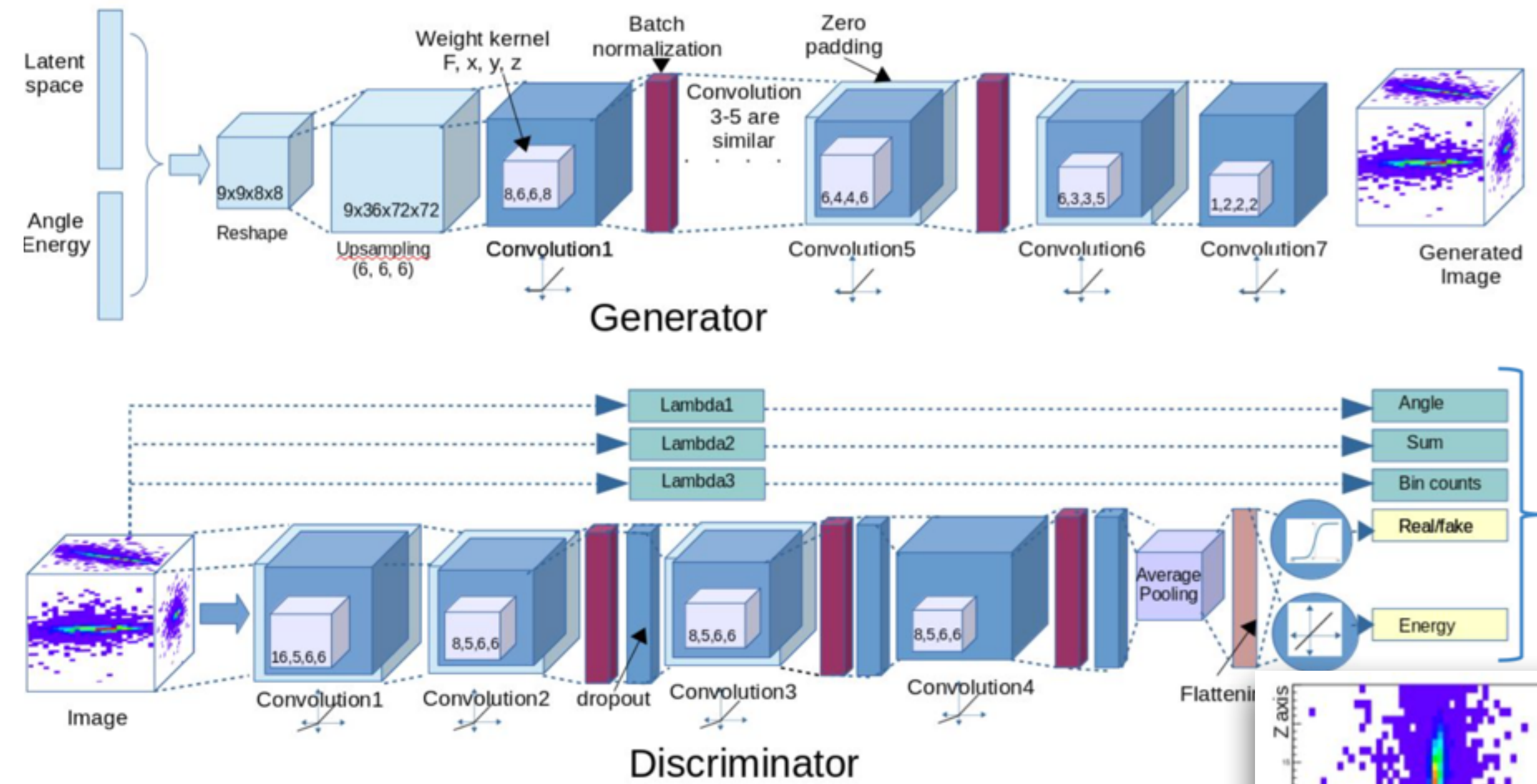


<https://indico.cern.ch/event/766872/contributions/3357987/>

High granularity calorimeter fast simulation for CLIC

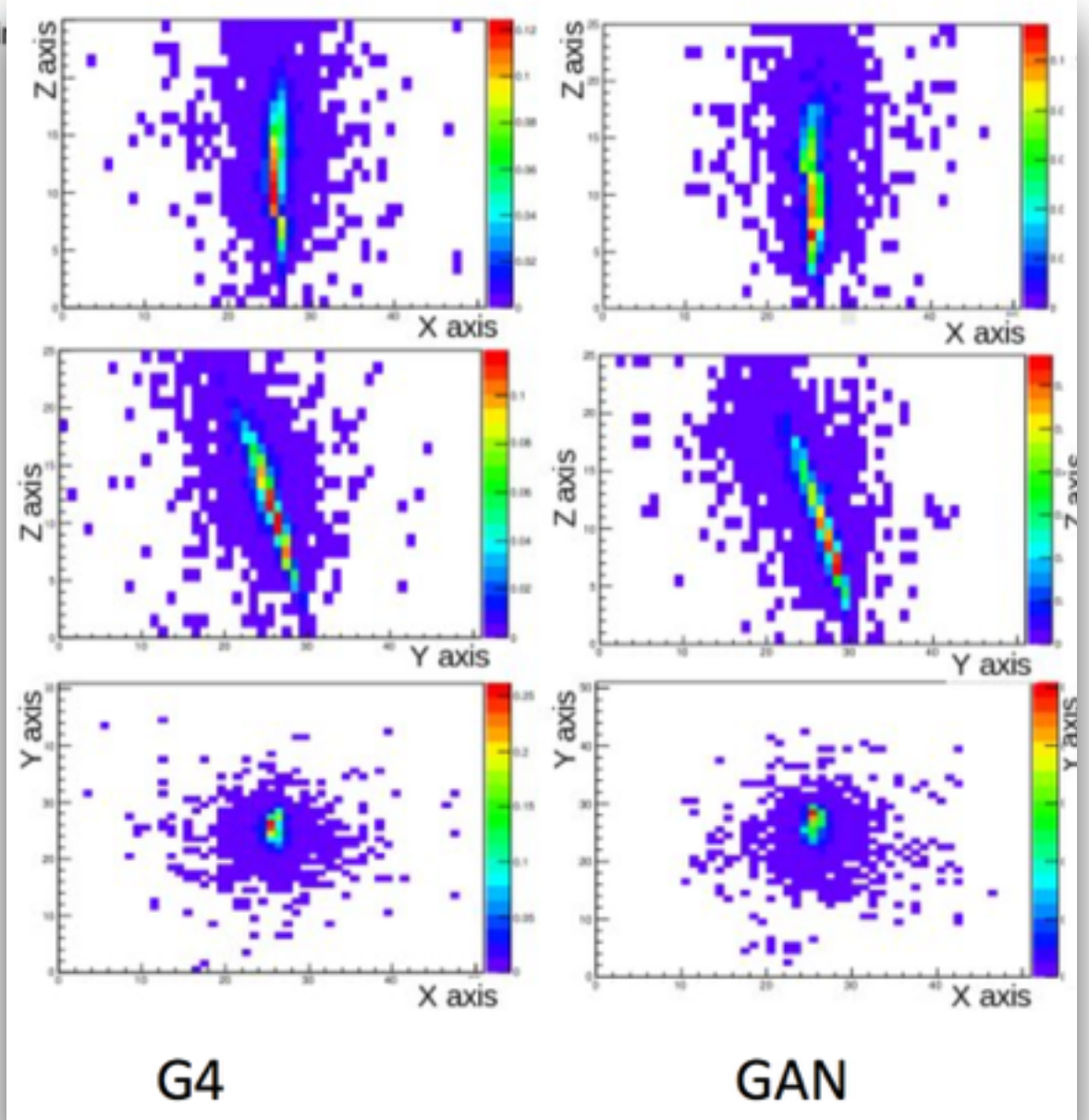


Sparse images with ~65000 pixels!!!



Evaluating with >200 physics observables!!!

+ using structural similarity index (SSIM) on calorimeter images

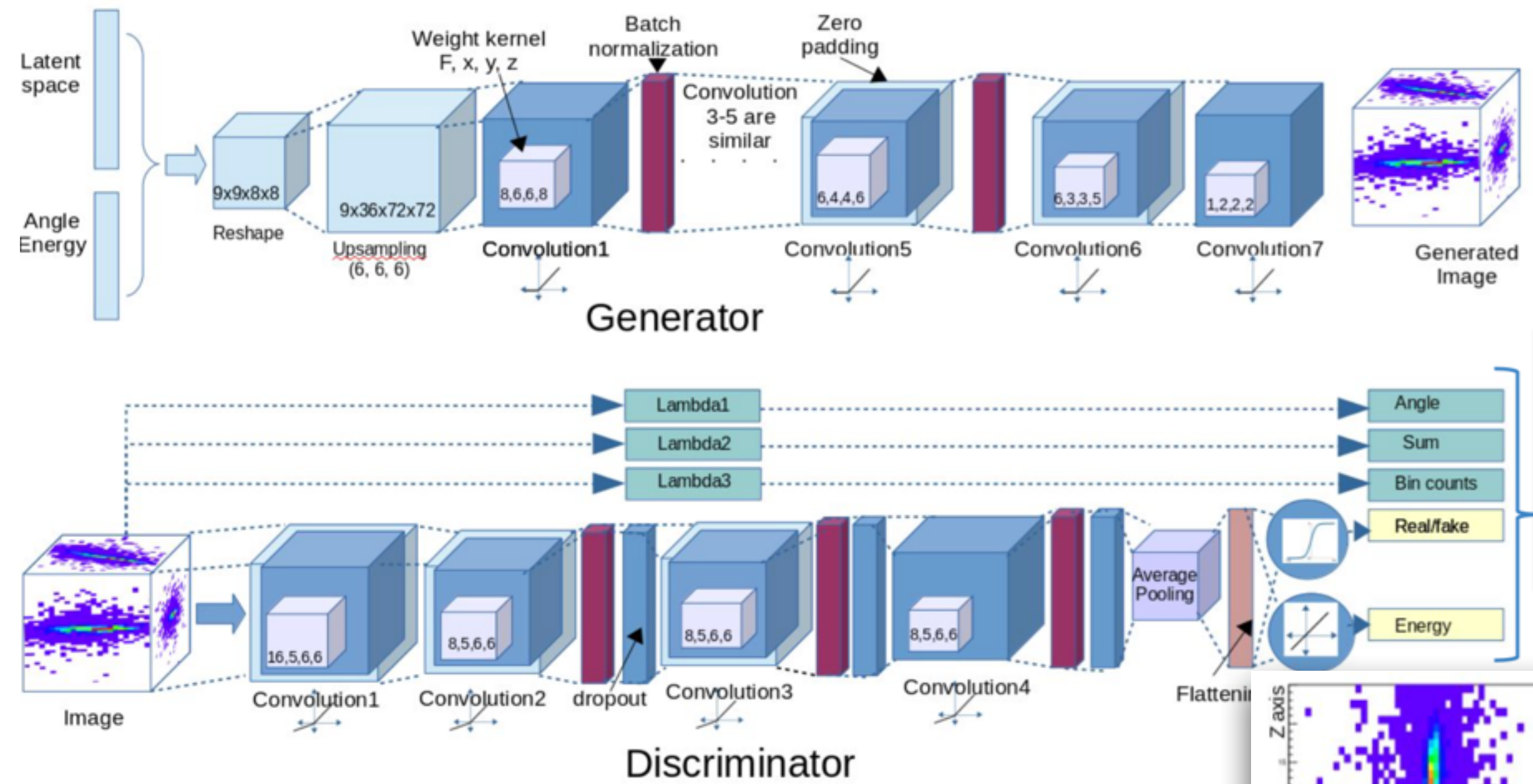


High granularity calorimeter fast simulation for CLIC



Sparse images with ~65000 pixels!!!

<https://indico.cern.ch/event/766872/contributions/3357987/>



Discriminator also predicts energy

Evaluating with >200 physics observables!!!

+ using structural similarity index (SSIM) on calorimeter images

