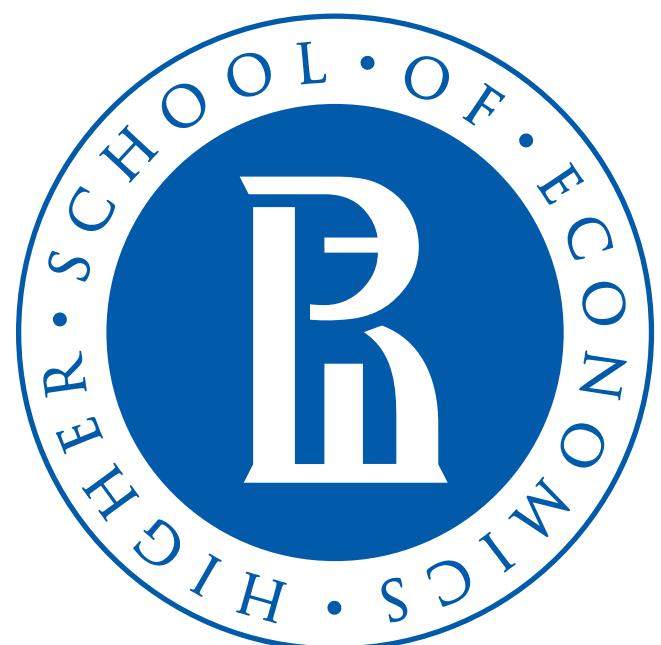


Generative Models, part II

GANs

Artem Maevskiy



Fifth Machine Learning in High Energy
Physics Summer School

DESY, Hamburg
1-10 July 2019

Outline

- Revising GANs
 - Basic idea, problems with GANs
- Wasserstein GAN
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- Side notes
 - GANs on weighted samples, conditionals
- Notable GAN architectures
 - pix2pix, CycleGAN, Progressive growing
- Evaluating GAN performance
- Applications in HEP

Outline

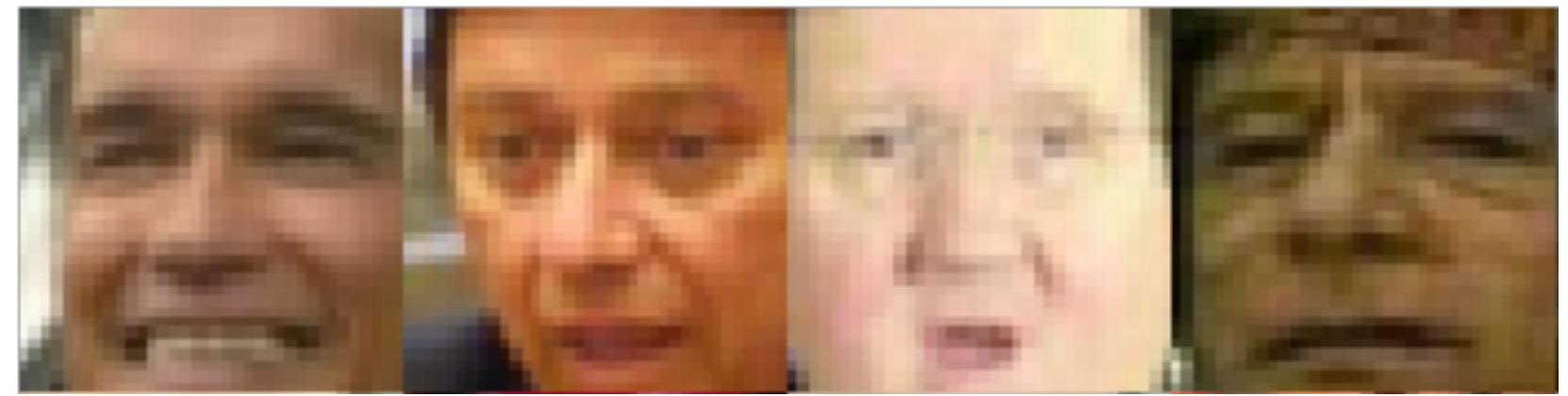
- **Revising GANs**
 - Basic idea, problems with GANs
- Wasserstein GAN
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- Side notes
 - GANs on weighted samples, conditionals
- Notable GAN architectures
 - pix2pix, CycleGAN, Progressive growing
- Evaluating GAN performance
- Applications in HEP

Revising Generative Adversarial Networks

Problem statement:

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

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

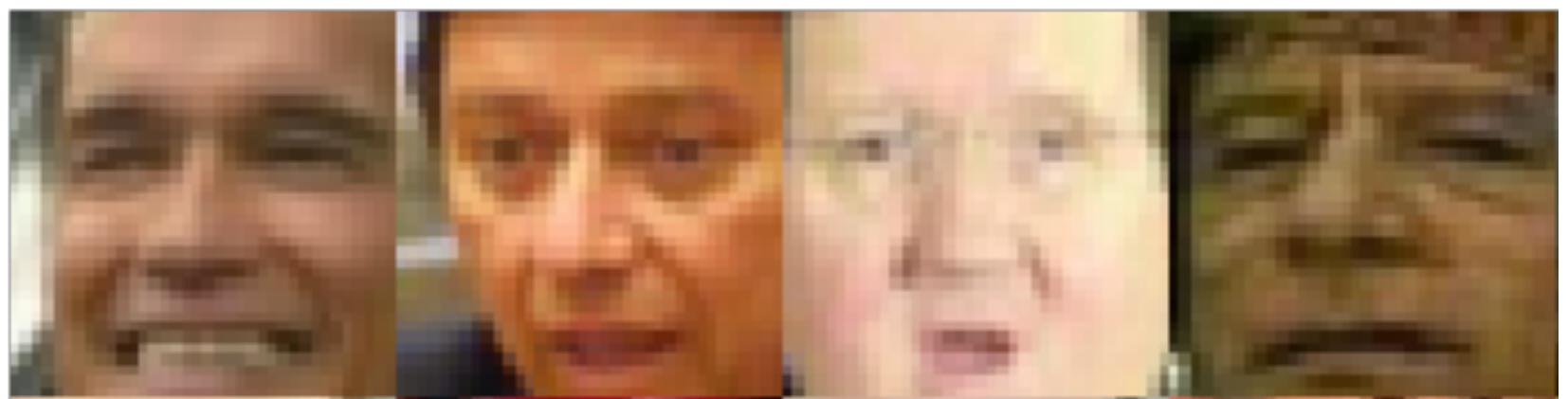


Revising Generative Adversarial Networks

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

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



Revising Generative Adversarial Networks

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

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



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

Revising Generative Adversarial Networks

Problem statement:

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

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

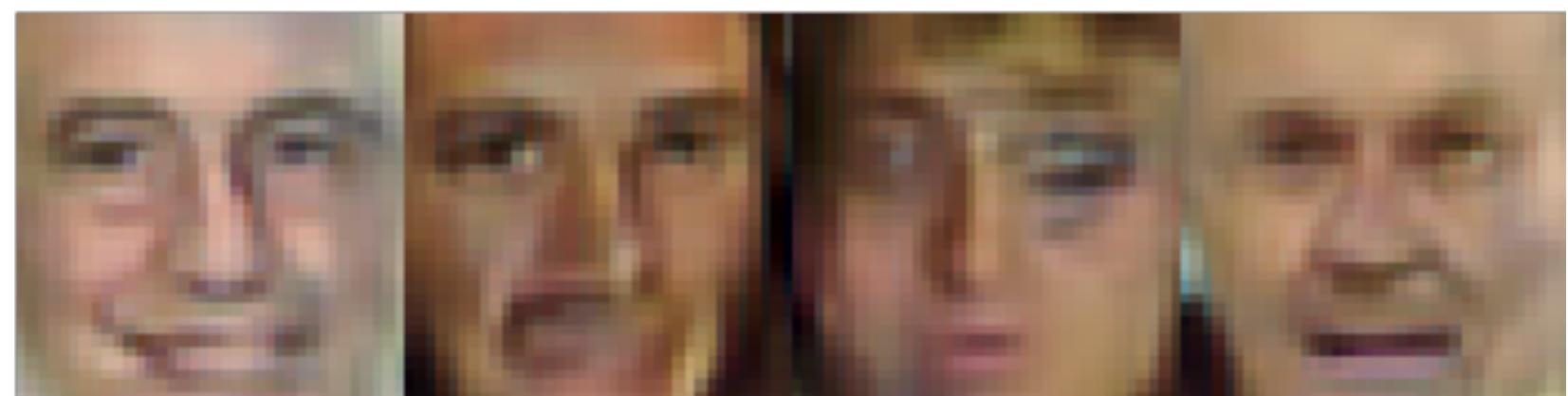
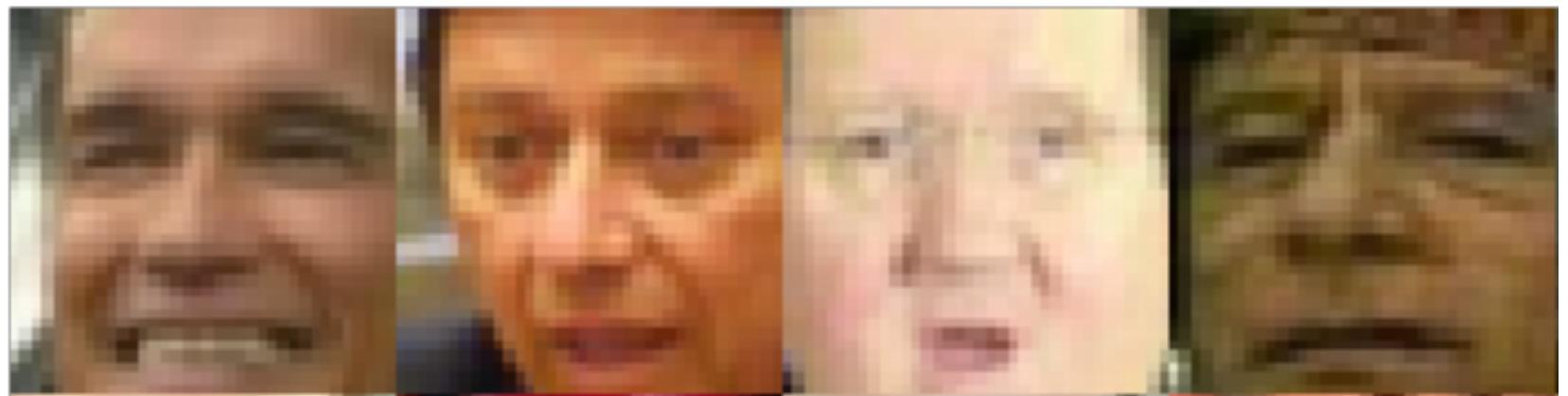
We want to learn to draw more samples from that distribution

a ‘similar’ distribution

$$\{\cancel{x}_j\} \sim \cancel{p}_{\text{data}}$$

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

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



Revising Generative Adversarial Networks

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

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

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

G_θ — generator neural network

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

Revising Generative Adversarial Networks

$\{\tilde{\mathbf{x}}_j\} \sim p_{gen, \theta}$ — this distribution is built 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.

Revising Generative Adversarial Networks

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

Probability the sample was drawn from data

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

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

G_θ — generator neural network

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.

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)

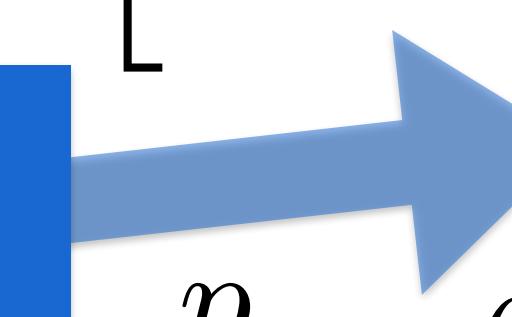
Let's have another look at the objective:

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

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

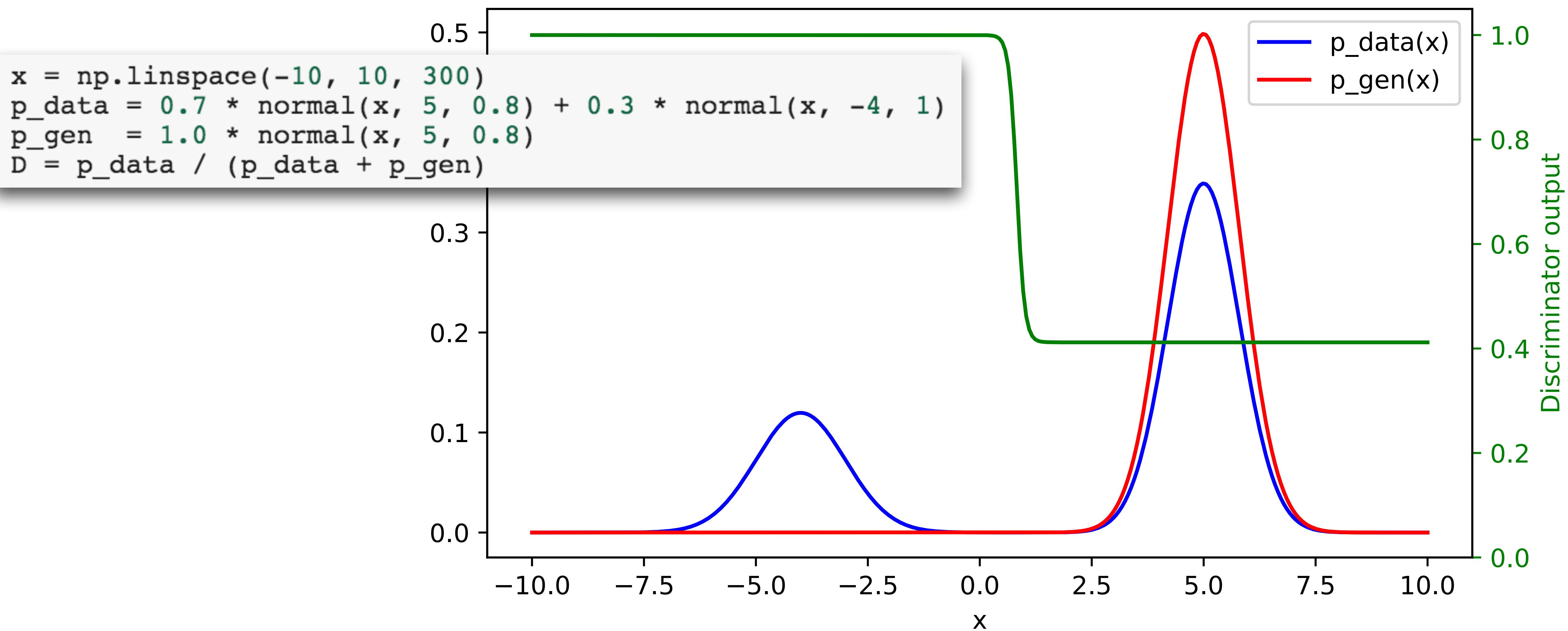
So the objective is: $V(\phi^*, \theta) = \int \left[p_{\text{data}}(\mathbf{x}) \log \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_{\text{gen}, \theta}(\mathbf{x})} + p_{\text{gen}, \theta}(\mathbf{x}) \log \frac{p_{\text{gen}, \theta}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) +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

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



Problems with GANs

(mode collapse)

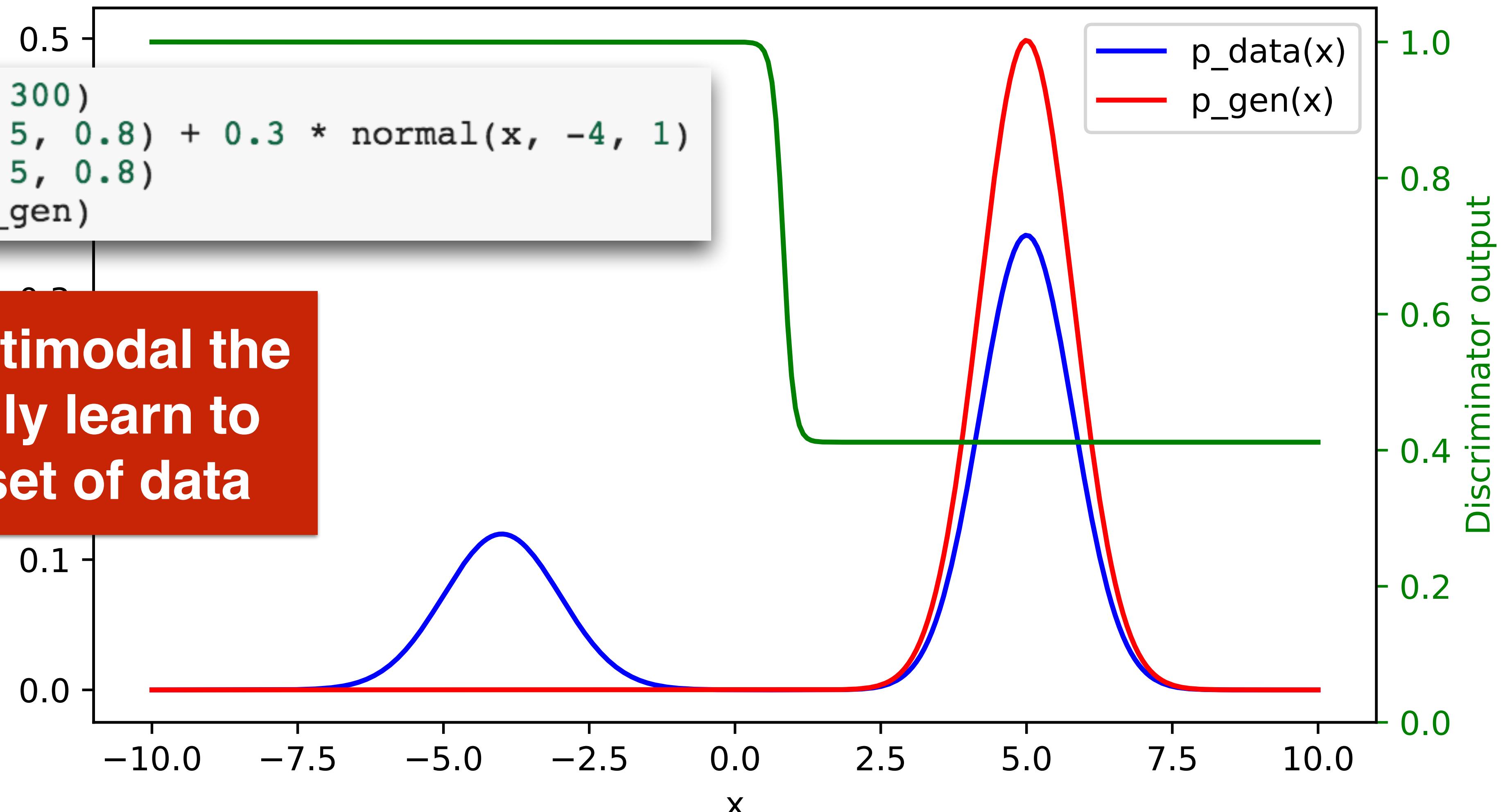


Problems with GANs

(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



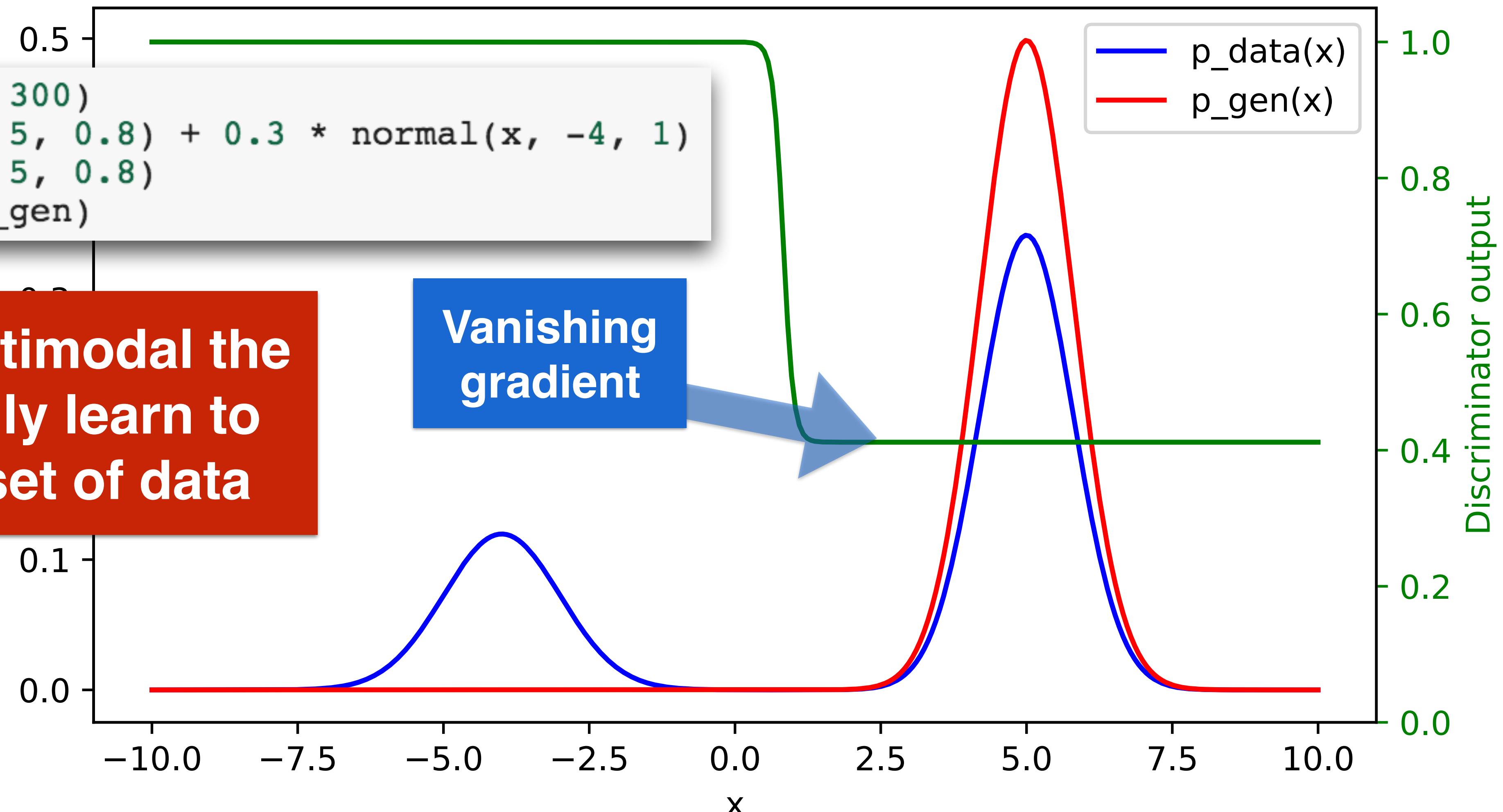
Problems with GANs

(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



Outline

- Revising GANs
 - Basic idea, problems with GANs
- **Wasserstein GAN**
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- Side notes
 - GANs on weighted samples, conditionals
- Notable GAN architectures
 - pix2pix, CycleGAN, Progressive growing
- Evaluating GAN performance
- Applications in HEP

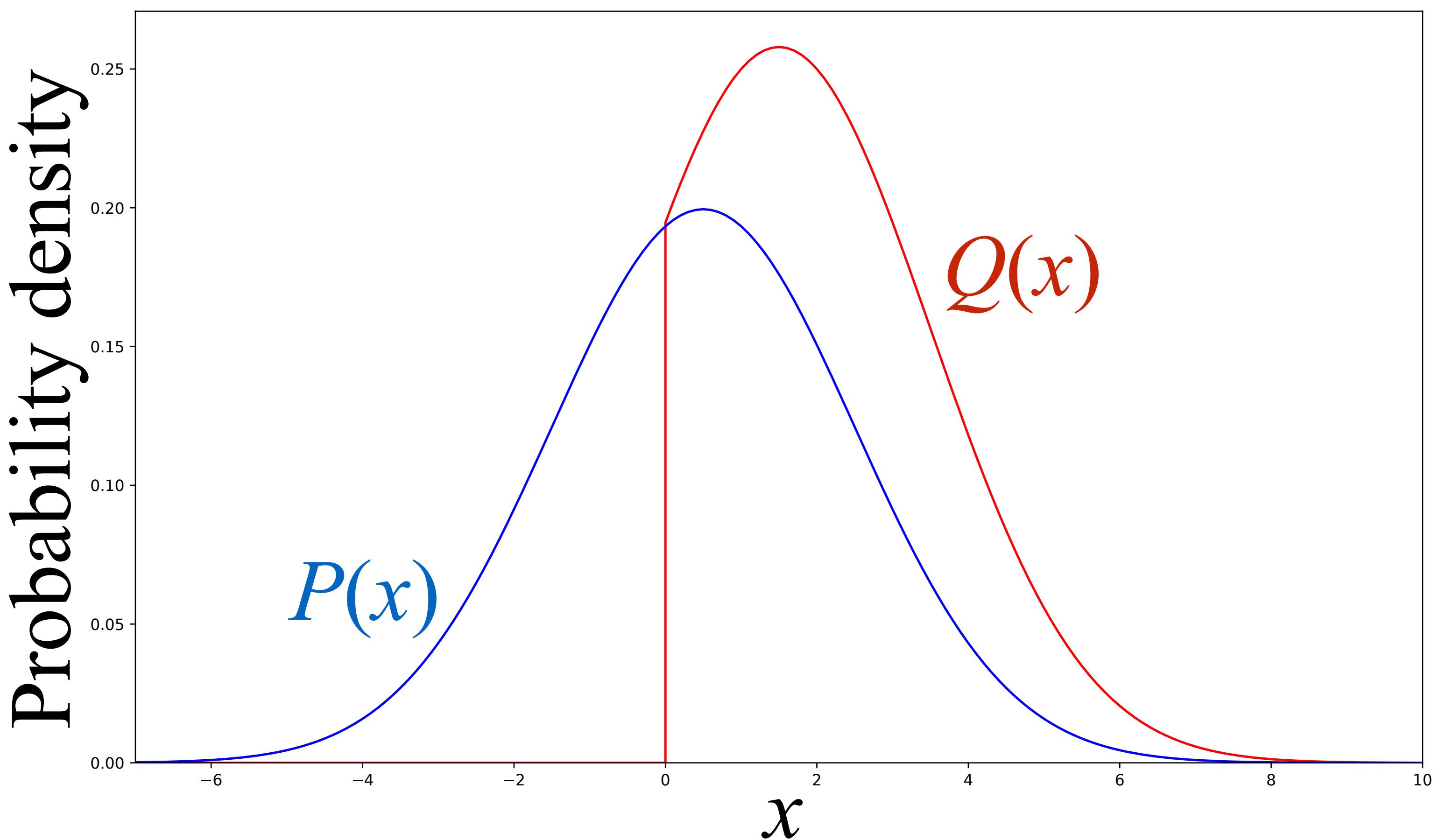
Wasserstein distance

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

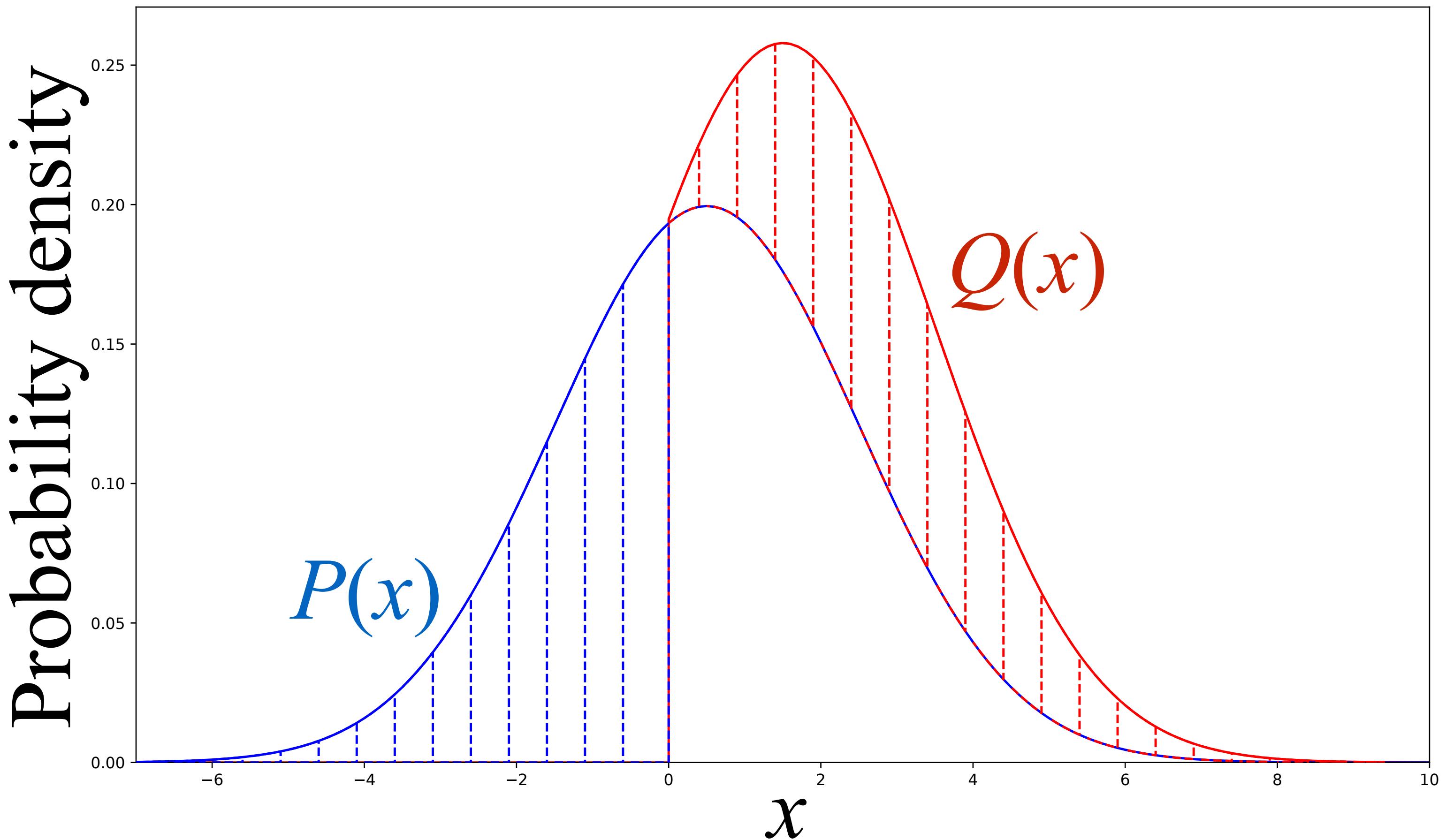


Wasserstein distance

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

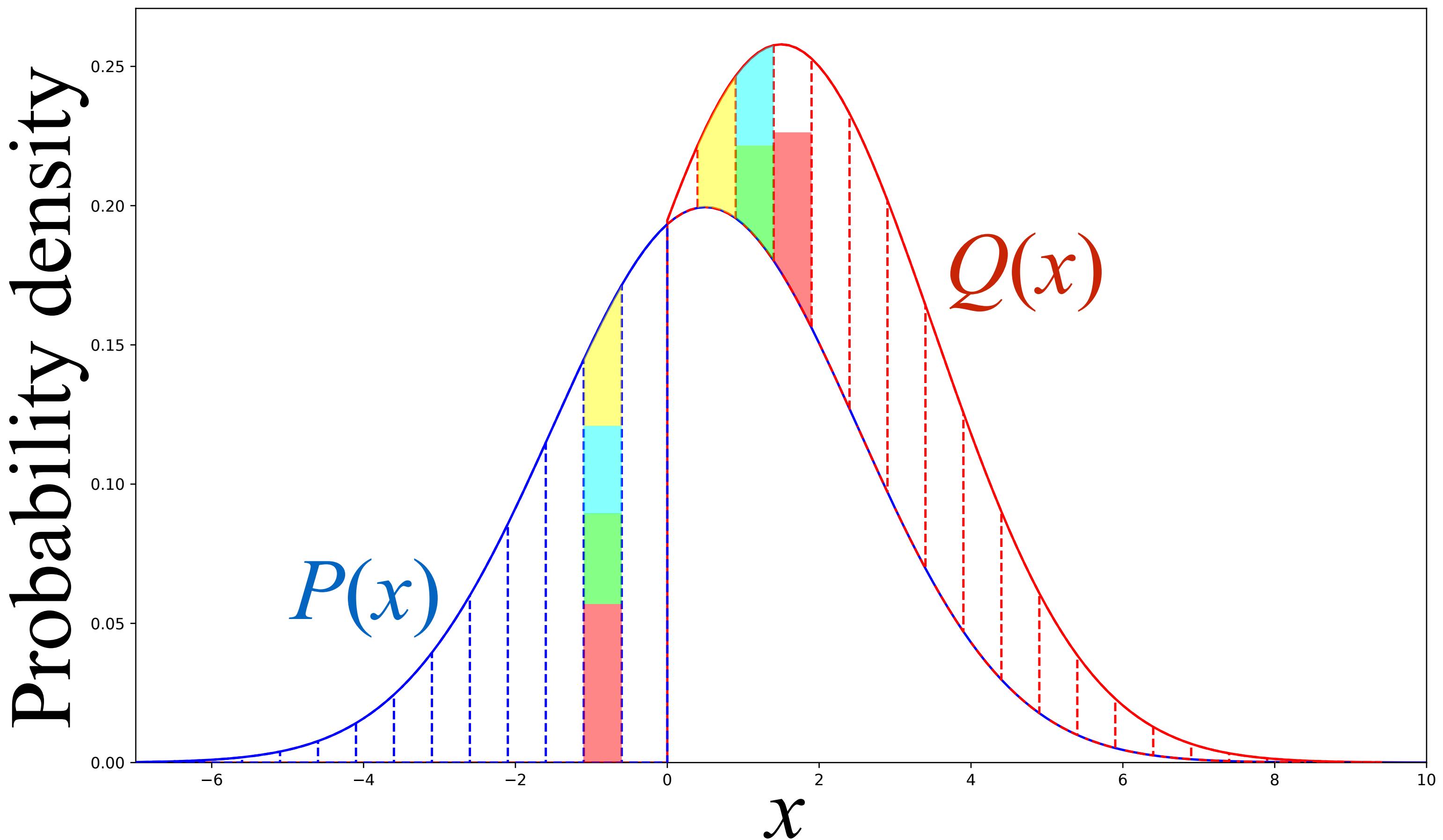


Wasserstein distance

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



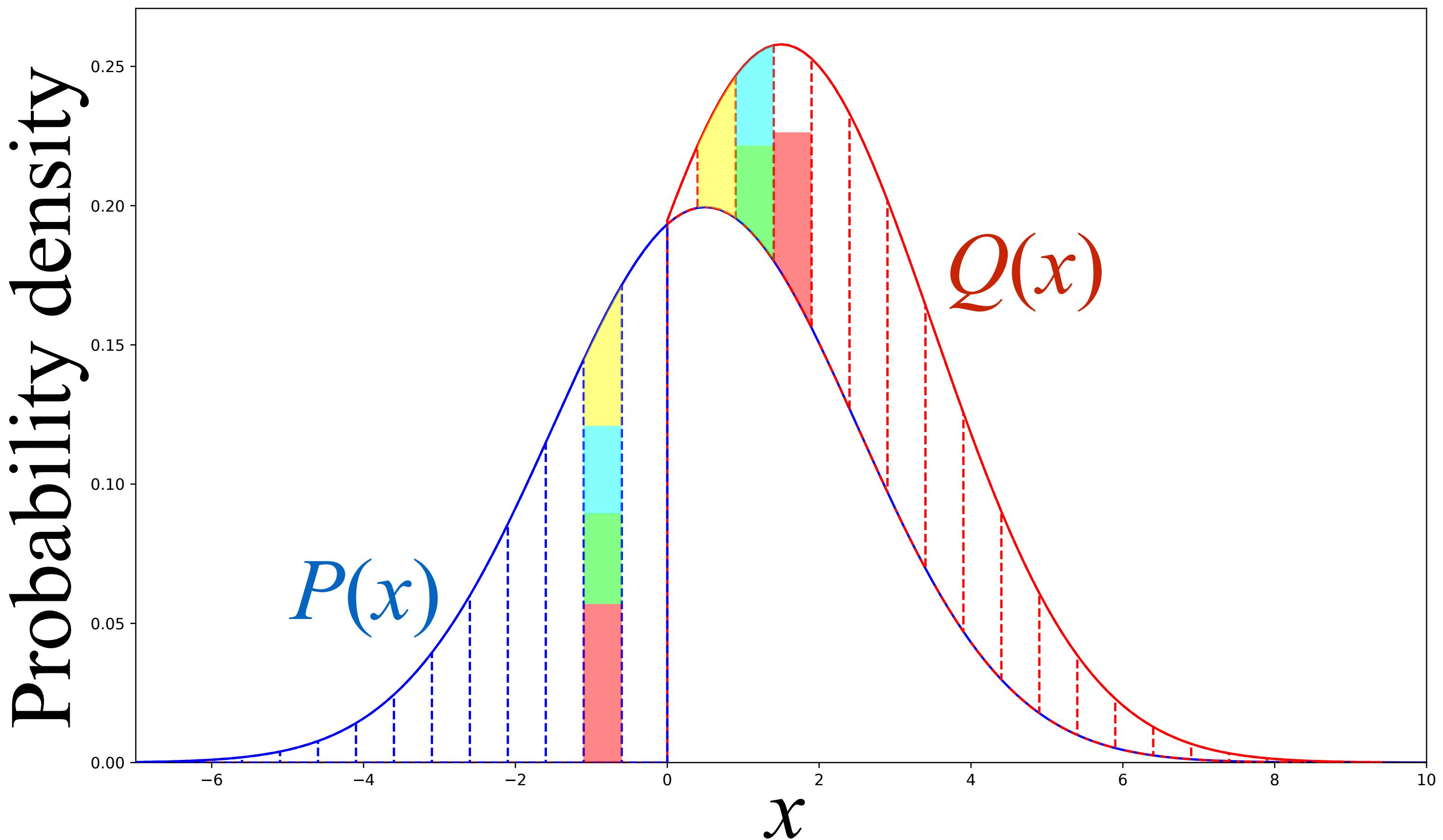
Wasserstein distance

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

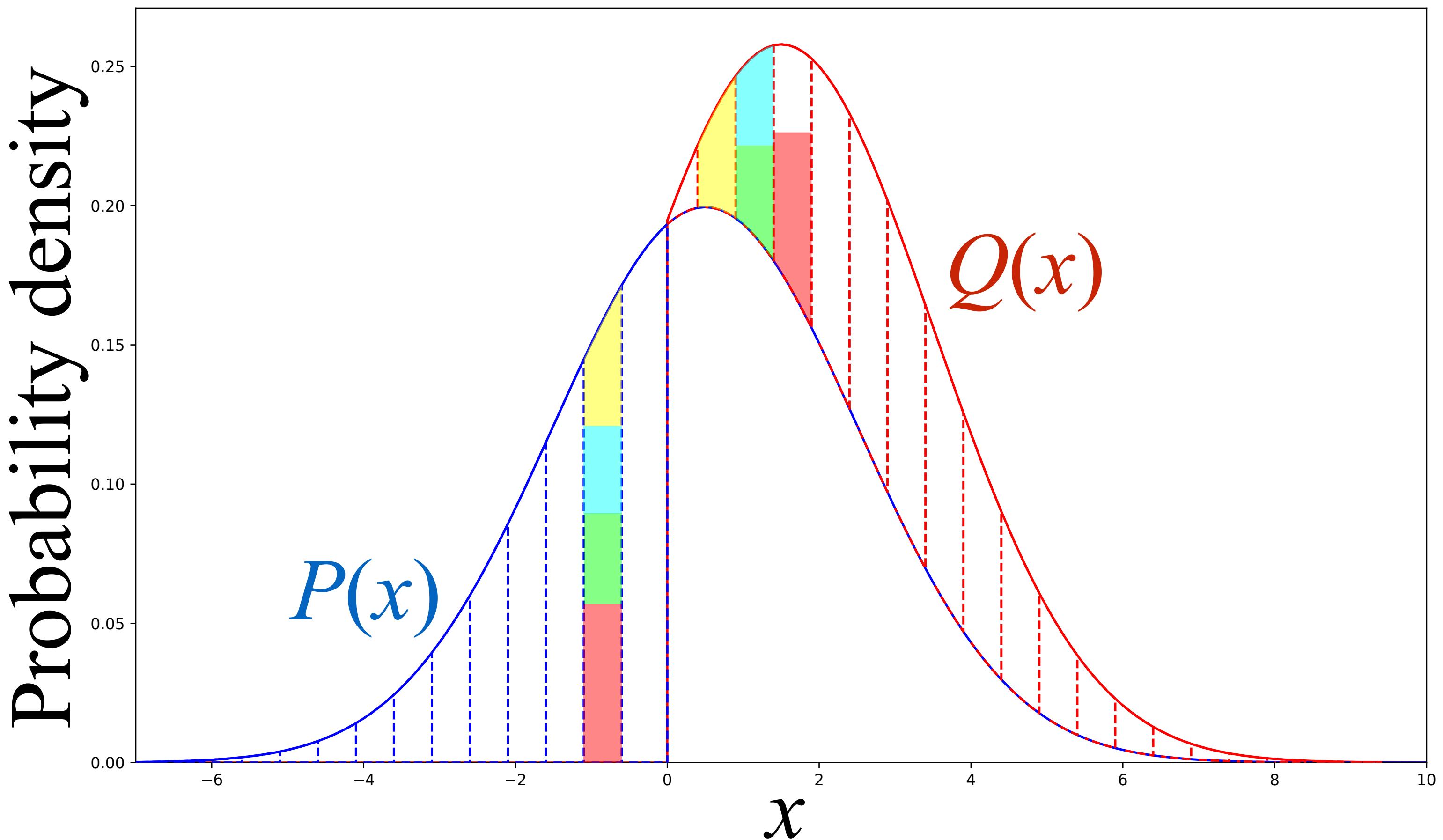
Also called «Earth mover's distance» (EMD)

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 $m \times \|x_2 - x_1\|$

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



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 x_1 to x_2 : $\gamma(x_1, x_2)$

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from x_1 to x_2 : $\gamma(x_1, x_2)$

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

$$\int \gamma(x_1, x_2) dx_1 = Q(x_2), \quad \int \gamma(x_1, x_2) dx_2 = P(x_1)$$

Wasserstein distance

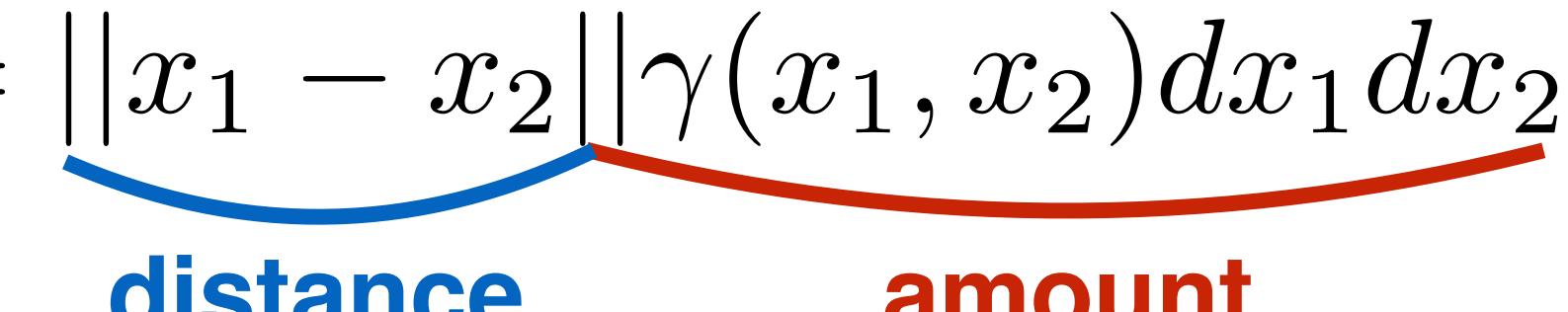
How to define it formally?

Say, we have a plan of how much earth to move from x_1 to x_2 : $\gamma(x_1, x_2)$

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

$$\int \gamma(x_1, x_2) dx_1 = Q(x_2), \quad \int \gamma(x_1, x_2) dx_2 = P(x_1)$$

The cost for some particular values x_1 and x_2 is: $dC = ||x_1 - x_2|| \gamma(x_1, x_2) dx_1 dx_2$



A diagram illustrating the cost term dC . It shows a blue curved arrow pointing from the expression $||x_1 - x_2||$ to the word "distance". A red curved arrow points from the expression $\gamma(x_1, x_2) dx_1 dx_2$ to the word "amount". The two arrows meet at a central point representing the product of the two terms.

Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from x_1 to x_2 : $\gamma(x_1, x_2)$

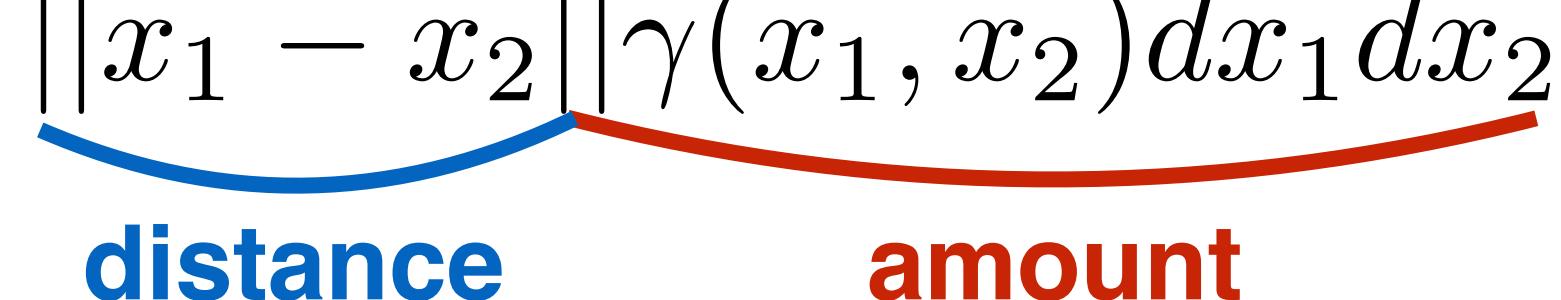
Since we want to get from P to Q , our plan has to satisfy these conditions:

$$\int \gamma(x_1, x_2) dx_1 = Q(x_2), \quad \int \gamma(x_1, x_2) dx_2 = P(x_1)$$

The cost for some particular values x_1 and x_2 is: $dC = ||x_1 - x_2|| \gamma(x_1, x_2) dx_1 dx_2$

Then, total cost can simply be written as:

$$C = \iint \gamma(x_1, x_2) ||x_1 - x_2|| dx_1 dx_2$$



Wasserstein distance

How to define it formally?

Say, we have a plan of how much earth to move from x_1 to x_2 : $\gamma(x_1, x_2)$

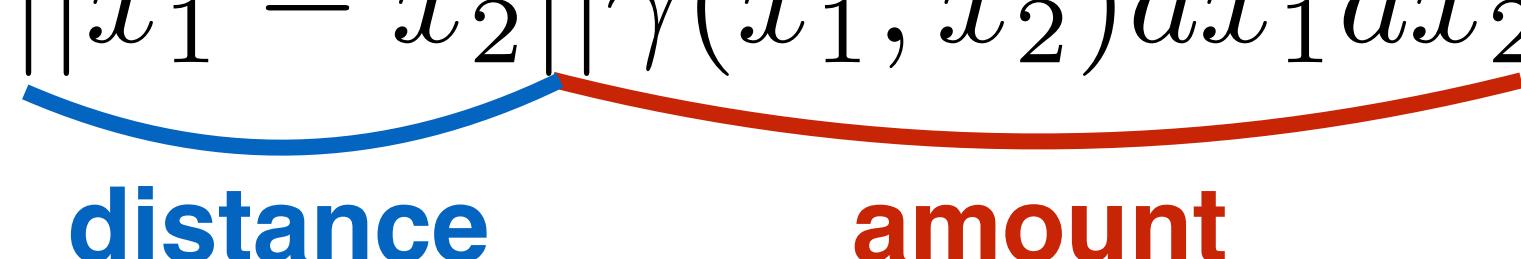
Since we want to get from P to Q , our plan has to satisfy these conditions:

$$\int \gamma(x_1, x_2) dx_1 = Q(x_2), \quad \int \gamma(x_1, x_2) dx_2 = P(x_1)$$

The cost for some particular values x_1 and x_2 is: $dC = ||x_1 - x_2|| \gamma(x_1, x_2) dx_1 dx_2$

Then, total cost can simply be written as:

$$C = \iint \gamma(x_1, x_2) ||x_1 - x_2|| dx_1 dx_2 = \mathbb{E}_{x_1, x_2 \sim \gamma(x_1, x_2)} ||x_1 - x_2||$$



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:

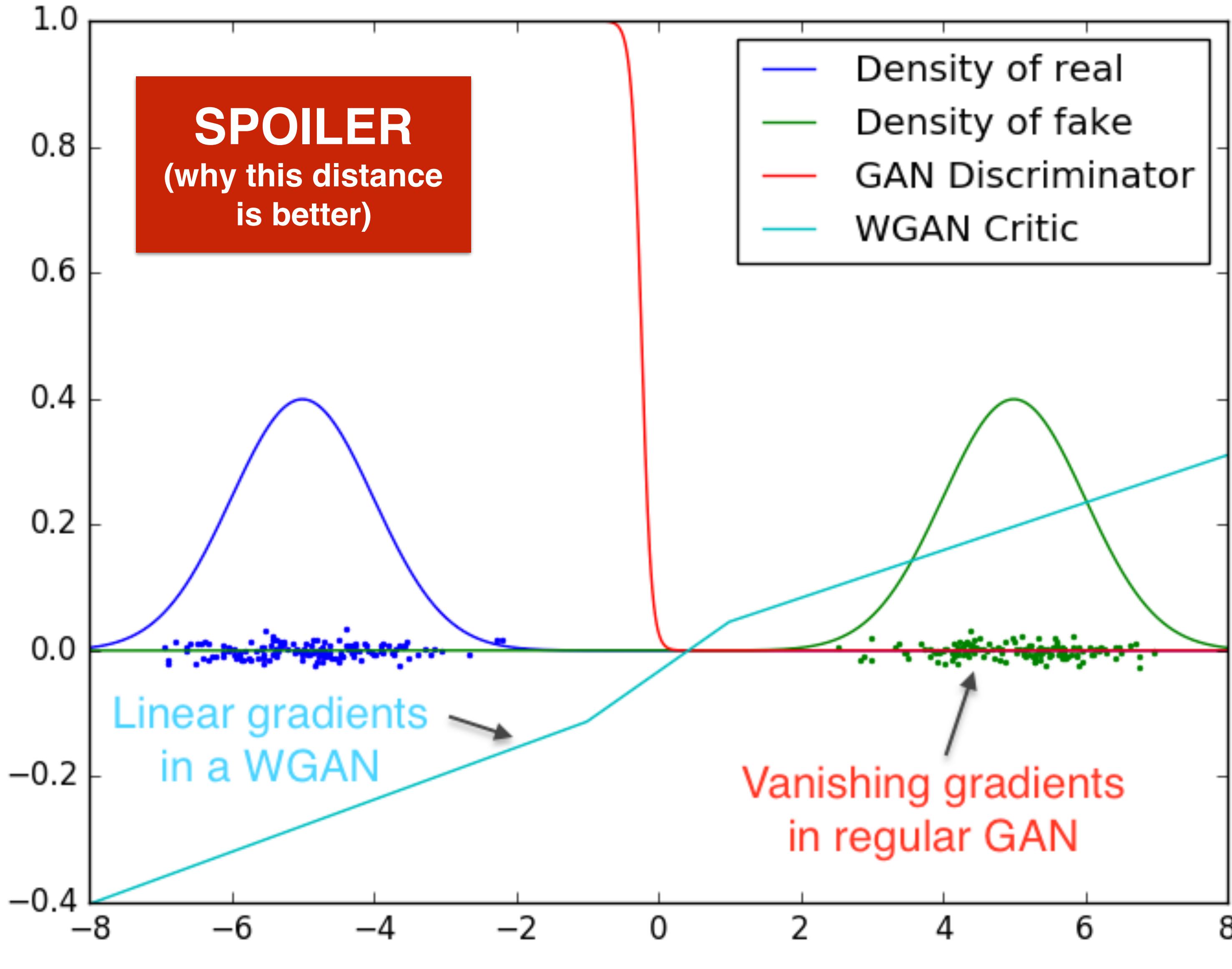
How would that be better than a regular GAN?

And more importantly...

How does one calculate this madness?!

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

WGAN



Looking ahead + answering
the first question

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

Kantorovich-Rubinstein duality

(or how to calculate EMD)

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

Kantorovich-Rubinstein duality (or how to calculate EMD)

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

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)

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

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)

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

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)

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

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)

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

Kantorovich-Rubinstein duality (or how to calculate EMD)

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

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

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

Kantorovich-Rubinstein duality (or how to calculate EMD)

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

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)

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

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)

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

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?



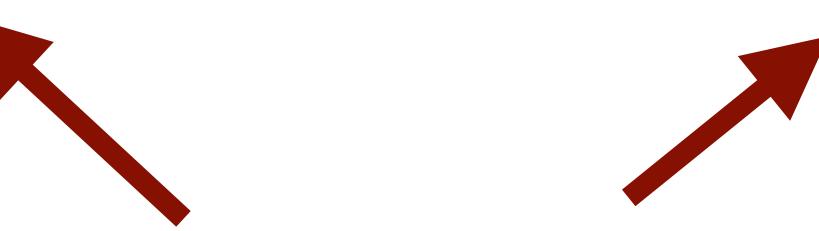
WGAN

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

<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)

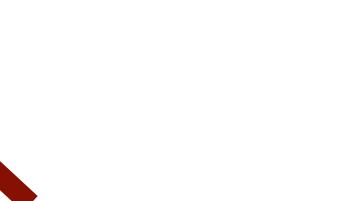
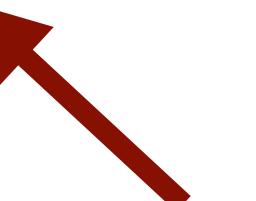
WGAN

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

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

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 $\text{EMD}(P, Q)$



Expectations can be estimated as sample mean (per batch)

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)

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)

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: $w \rightarrow \text{clip}(w, -c, c)$

Expectations can be estimated as sample mean (per batch)

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

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** θ 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: $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$
- 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: $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$
- 12: **end while**

<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

+

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

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

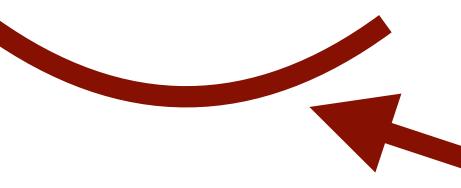
Optimal f should satisfy $\|\nabla f\| = 1$ almost
everywhere under P and Q

+

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

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

+

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

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

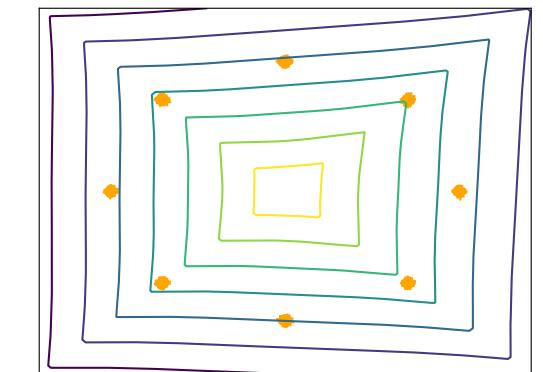
+

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

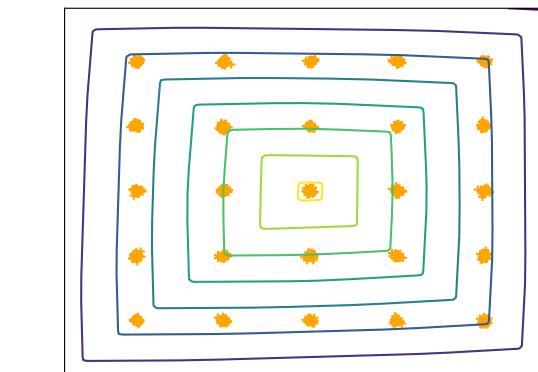
Weight clipping can be replaced by
penalizing the gradients to be of unit
magnitude

«Gradient penalty»

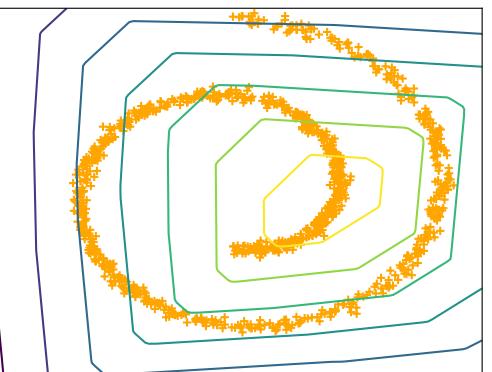
8 Gaussians



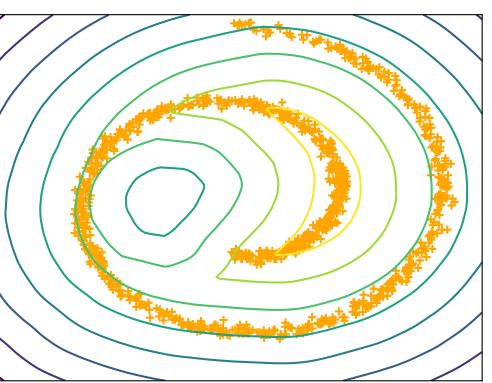
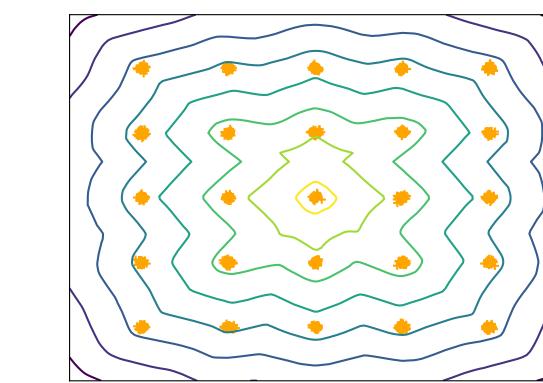
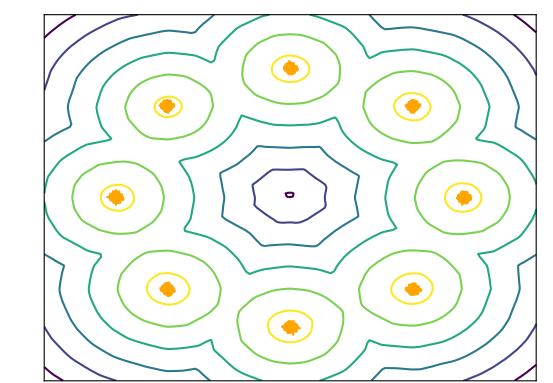
25 Gaussians



Swiss Roll



WGAN



WGAN-GP

Value surfaces of f

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

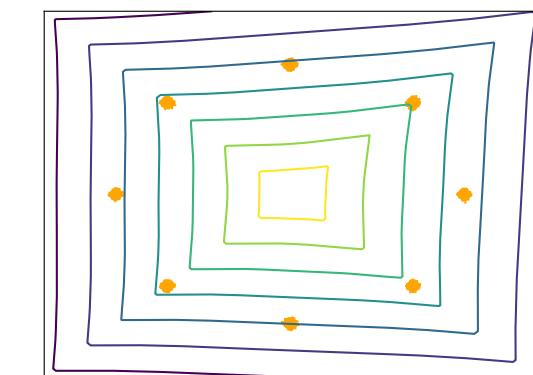
Optimal f should satisfy $\|\nabla f\| = 1$ almost
everywhere under P and Q

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

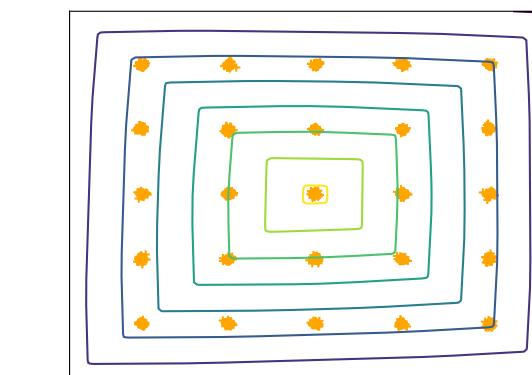
Weight clipping can be replaced by
penalizing the gradients to be of unit
magnitude

«Gradient penalty»

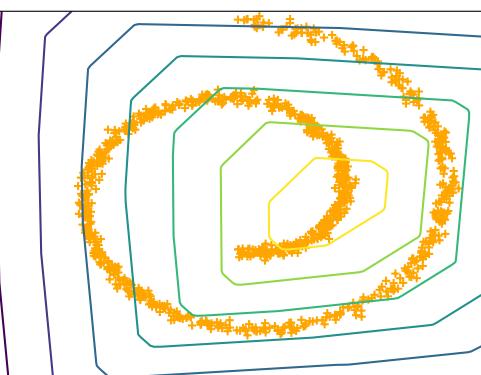
8 Gaussians



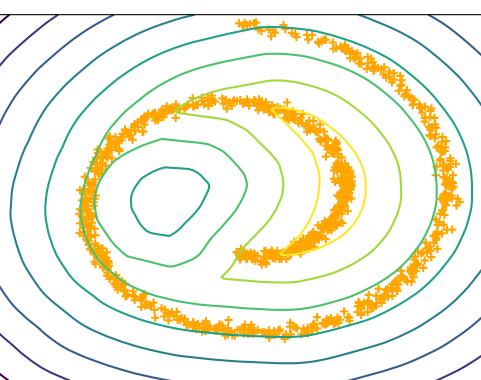
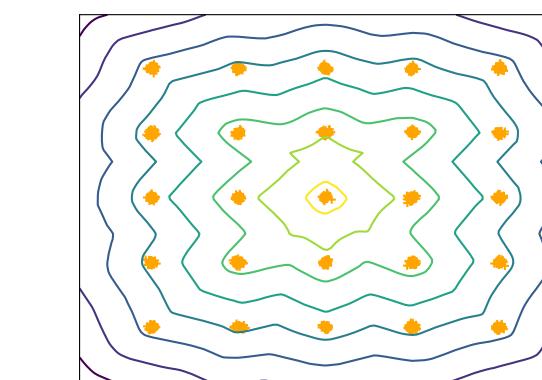
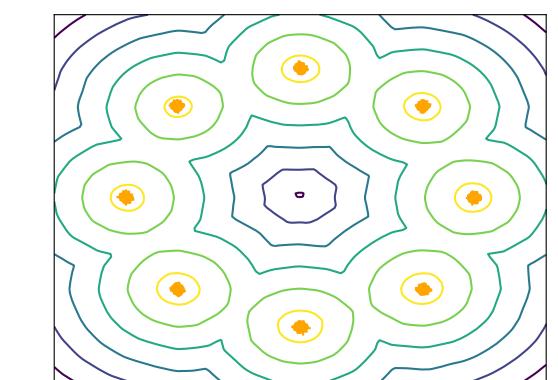
25 Gaussians



Swiss Roll



WGAN



WGAN-GP

Value surfaces of f

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

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

WGAN-GP

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

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

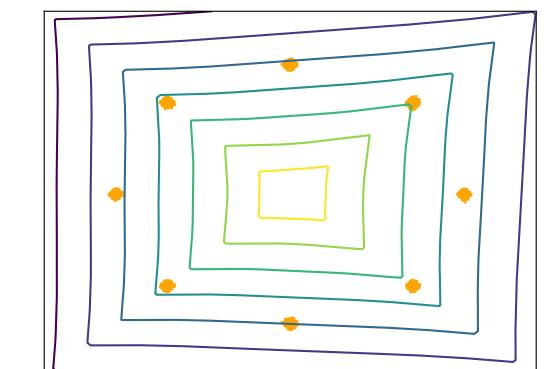
W
p
m

Alternative ‘one-sided’ penalty:

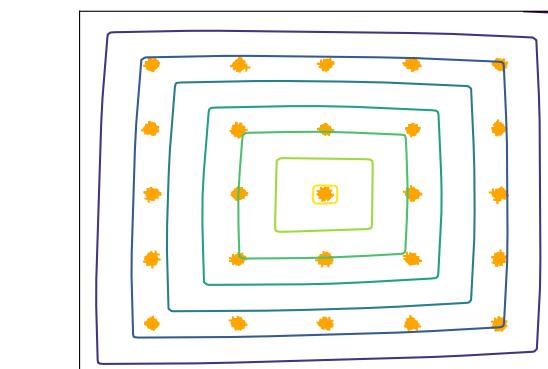
$$GP = \lambda \mathbb{E}_{\tilde{x} \sim \mathbb{P}_{\tilde{x}}} [\max(0, \|\nabla_{\tilde{x}} f(\tilde{x})\| - 1)^2]$$

«Gradient penalty»

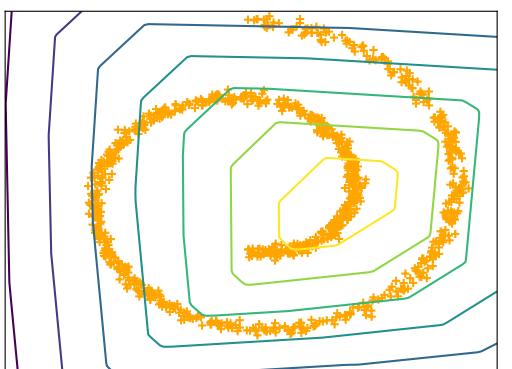
8 Gaussians



25 Gaussians



Swiss Roll



WGAN WGAN-GP

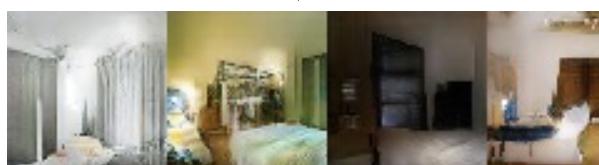
Value surfaces of f

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

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

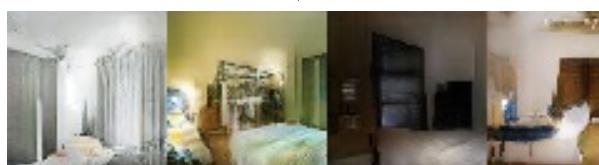
WGAN-GP

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			

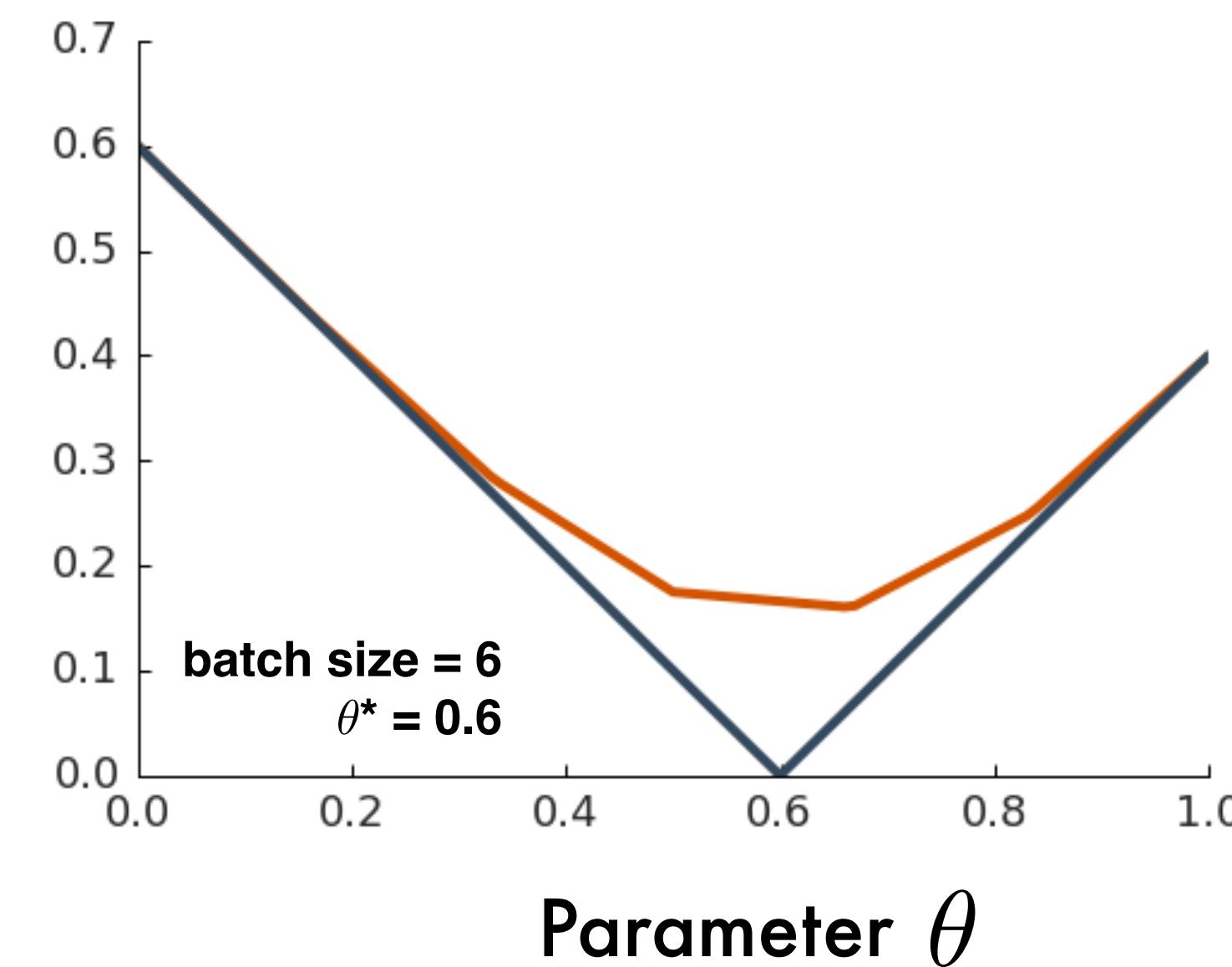
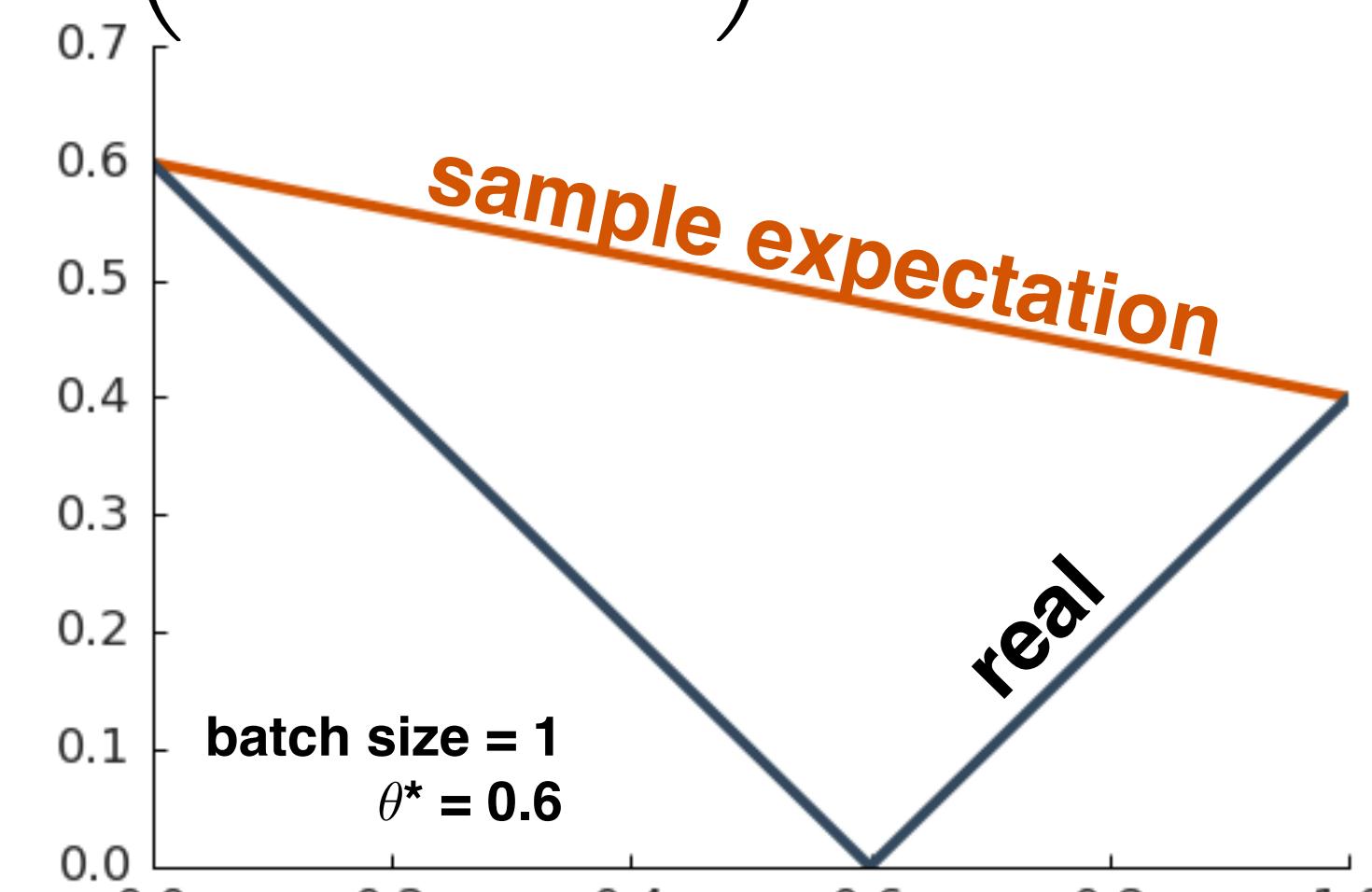
This technique allowed for very deep networks to be used for GANs

Biased Wasserstein gradients

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

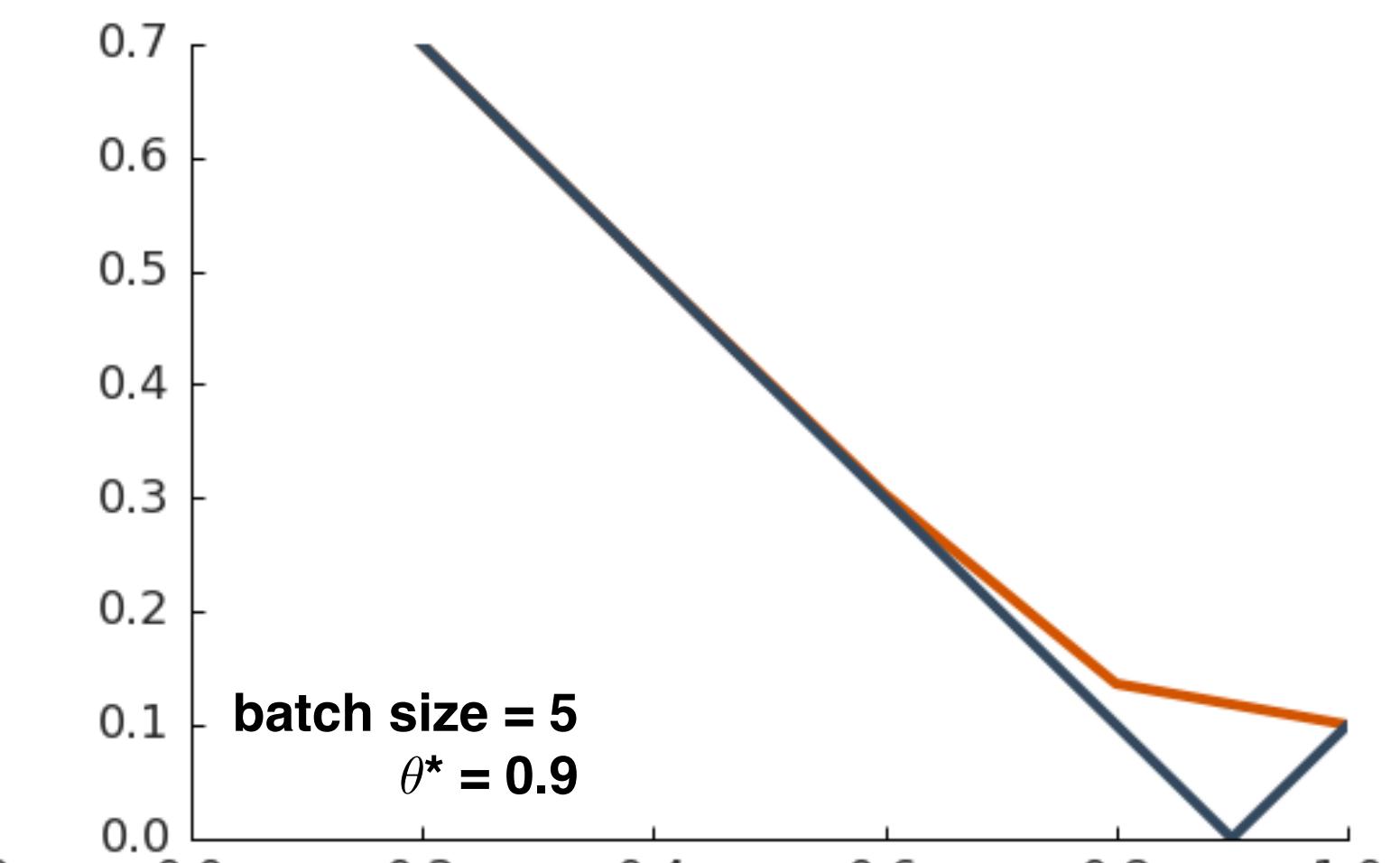
This may lead to a wrong minimum:

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



Bernoulli distribution

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



Cramer and energy distances

Cramer distance:

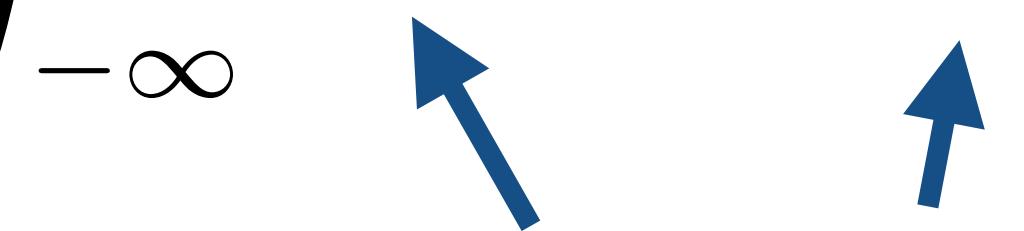
$$l_2^2(P, Q) := \int_{-\infty}^{\infty} (F_P(x) - F_Q(x))^2 dx$$


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>

Which Training Methods for GANs do actually Converge?

<https://arxiv.org/abs/1801.04406> — the authors analyzed GAN training behavior on a very simple system:

$$\begin{array}{ll} \text{target distribution:} & p_{\text{data}} = \delta(x) \\ \text{generator distribution:} & p_{\text{gen}} = \delta(x - \theta) \\ \text{discriminator:} & D_\psi(x) = \psi \cdot x \end{array}$$

Ideal solution: $\theta = \psi = 0$

Which Training Methods for GANs do actually Converge?

<https://arxiv.org/abs/1801.04406> — the authors analyzed GAN training behavior on a very simple system:

target distribution:	$p_{\text{data}} = \delta(x)$
generator distribution:	$p_{\text{gen}} = \delta(x - \theta)$
discriminator:	$D_\psi(x) = \psi \cdot x$

Single-parameter generator

Ideal solution: $\theta = \psi = 0$

Which Training Methods for GANs do actually Converge?

<https://arxiv.org/abs/1801.04406> — the authors analyzed GAN training behavior on a very simple system:

target distribution:

$$p_{\text{data}} = \delta(x)$$

generator distribution:

$$p_{\text{gen}} = \delta(x - \theta)$$

discriminator:

$$D_\psi(x) = \psi \cdot x$$

Ideal solution: $\theta = \psi = 0$

Single-parameter generator

Single-parameter discriminator

Which Training Methods for GANs do actually Converge?

<https://arxiv.org/abs/1801.04406> — the authors analyzed GAN training behavior on a very simple system:

target distribution:

$$p_{\text{data}} = \delta(x)$$

generator distribution:

$$p_{\text{gen}} = \delta(x - \theta)$$

discriminator:

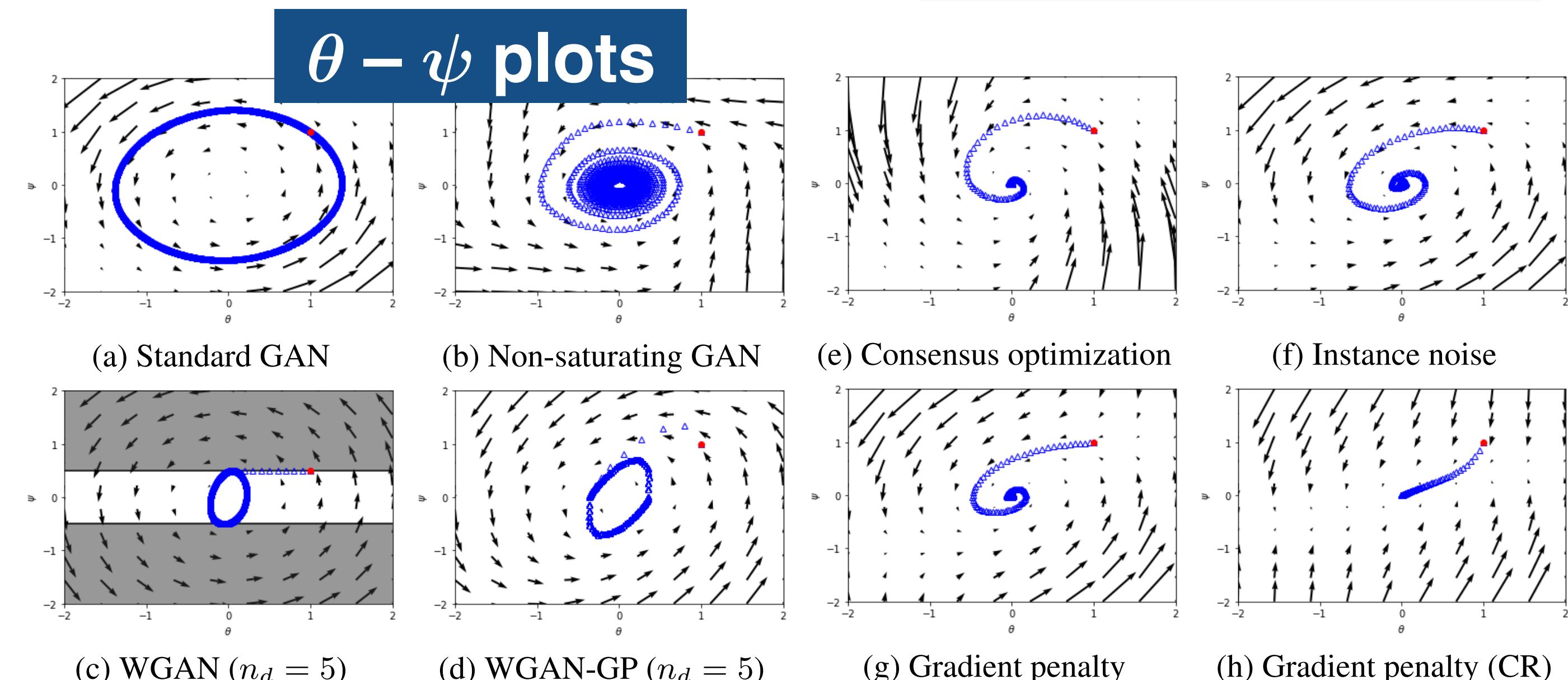
$$D_\psi(x) = \psi \cdot x$$

Single-parameter generator

Single-parameter discriminator

Ideal solution: $\theta = \psi = 0$

Examined various training schemes (with fixed learning rate):



Which Training Methods for GANs do actually Converge?

<https://arxiv.org/abs/1801.04406> — the authors analyzed GAN training behavior on a very simple system:

target distribution:

$$p_{\text{data}} = \delta(x)$$

generator distribution:

$$p_{\text{gen}} = \delta(x - \theta)$$

discriminator:

$$D_\psi(x) = \psi \cdot x$$

Single-parameter generator

Single-parameter discriminator

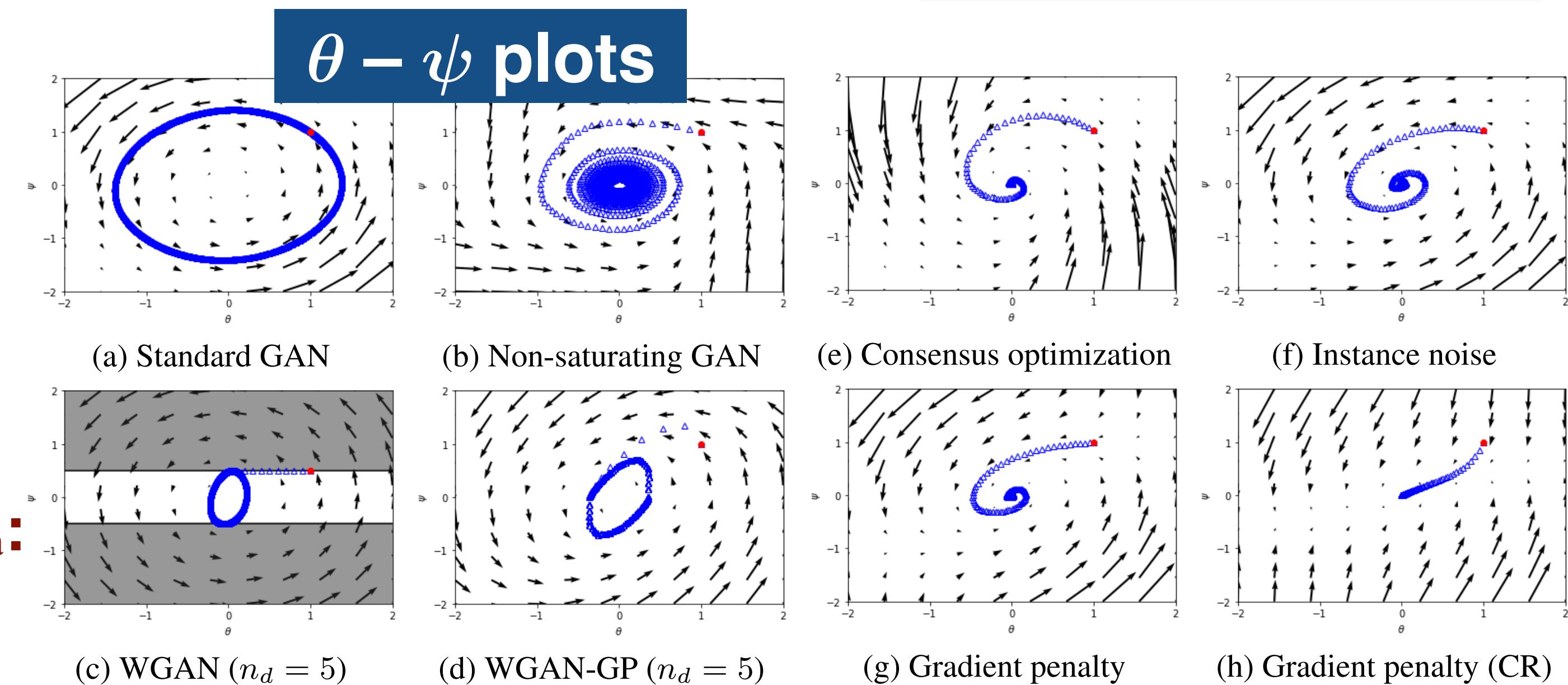
Ideal solution: $\theta = \psi = 0$

Examined various training schemes (with fixed learning rate):

Proposed theoretically motivated regularization schemes:

- **instance noise (adding noise to any discriminator input)**
- **zero-centered gradient penalty at $x \sim p_{\text{data}}$:**

$$\frac{\gamma}{2} \mathbb{E}_{x \sim p_{\text{data}}} \|\nabla D(x)\|^2$$



Which Training Methods for GANs do actually Converge?

<https://arxiv.org/abs/1801.04406> — the authors analyzed GAN training behavior on a very simple system:

target distribution:

$$p_{\text{data}} = \delta(x)$$

generator distribution:

$$p_{\text{gen}} = \delta(x - \theta)$$

discriminator:

$$D_\psi(x) = \psi \cdot x$$

Single-parameter generator

Single-parameter discriminator

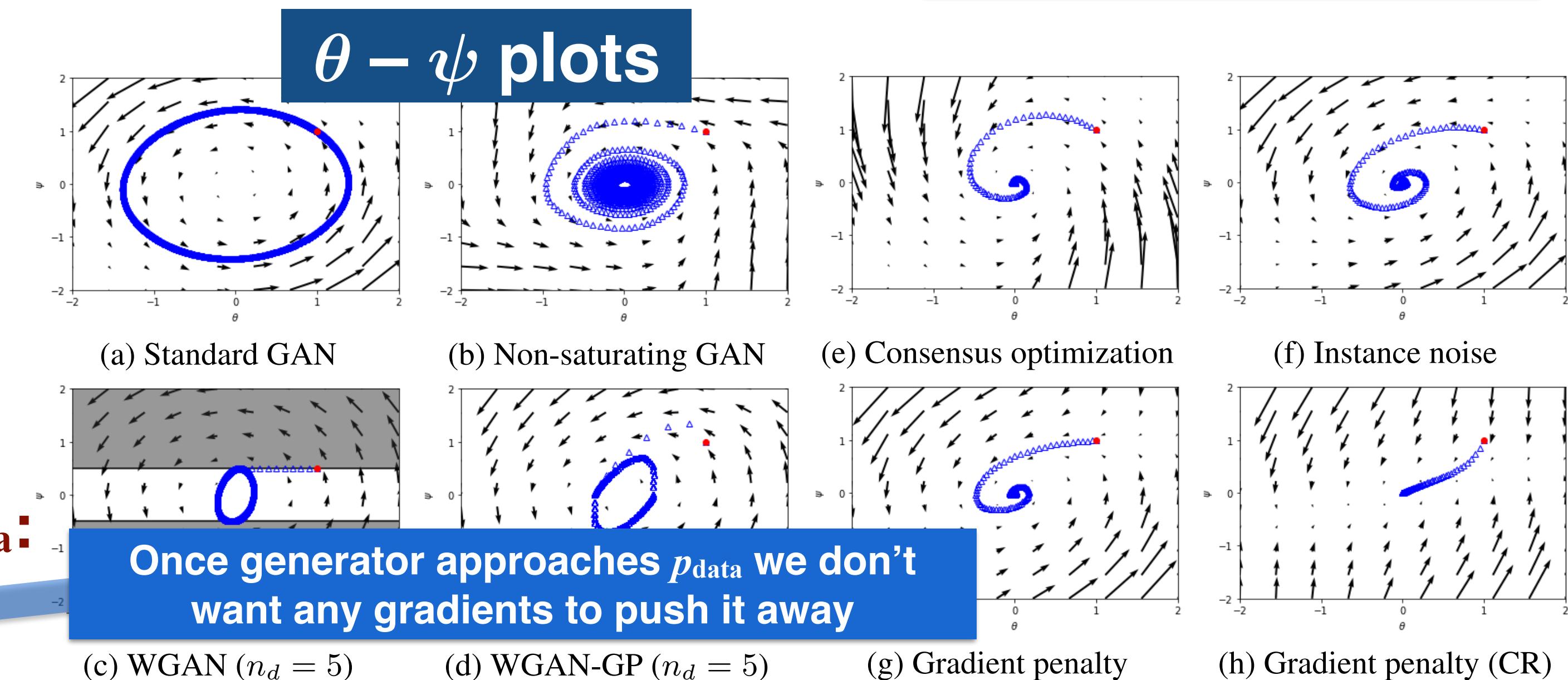
Ideal solution: $\theta = \psi = 0$

Examined various training schemes (with fixed learning rate):

Proposed theoretically motivated regularization schemes:

- **instance noise (adding noise to any discriminator input)**
- **zero-centered gradient penalty at $x \sim p_{\text{data}}$:**

$$\frac{\gamma}{2} \mathbb{E}_{x \sim p_{\text{data}}} \|\nabla D(x)\|^2$$



Boundary Equilibrium Generative Adversarial Networks (BEGAN)

<https://arxiv.org/abs/1703.10717>

Notice for Wasserstein between 1D distributions:

$$\inf \mathbb{E}[|x_1 - x_2|] \geq \inf |\mathbb{E}[x_1 - x_2]| = |m_1 - m_2|$$

Difference between
means

Boundary Equilibrium Generative Adversarial Networks (BEGAN)

<https://arxiv.org/abs/1703.10717>

Notice for Wasserstein between 1D distributions:

$$\inf \mathbb{E}[|x_1 - x_2|] \geq \inf |\mathbb{E}[x_1 - x_2]| = |m_1 - m_2|$$

The authors take auto-encoder loss:

$$\mathcal{L}(v) = |v - D(v)|^\eta \text{ where } \begin{cases} D : \mathbb{R}^{N_x} \mapsto \mathbb{R}^{N_x} \\ \eta \in \{1, 2\} \\ v \in \mathbb{R}^{N_x} \end{cases}$$

is the autoencoder function.
is the target norm.
is a sample of dimension N_x .

Difference between
means

Boundary Equilibrium Generative Adversarial Networks (BEGAN)

<https://arxiv.org/abs/1703.10717>

Notice for Wasserstein between 1D distributions:

$$\inf \mathbb{E}[|x_1 - x_2|] \geq \inf |\mathbb{E}[x_1 - x_2]| = |m_1 - m_2|$$

Difference between means

The authors take auto-encoder loss:

$$\mathcal{L}(v) = |v - D(v)|^\eta \text{ where } \begin{cases} D : \mathbb{R}^{N_x} \mapsto \mathbb{R}^{N_x} & \text{is the autoencoder function.} \\ \eta \in \{1, 2\} & \text{is the target norm.} \\ v \in \mathbb{R}^{N_x} & \text{is a sample of dimension } N_x. \end{cases}$$

Auto-encoder loss
for real data

Auto-encoder loss
for generated data

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x; \theta_D) - \mathcal{L}(G(z_D; \theta_G); \theta_D) & \text{for } \theta_D \\ \mathcal{L}_G = -\mathcal{L}_D & \text{for } \theta_G \end{cases}$$

And treat it as 1D distribution for the Wasserstein loss

- D is both auto-encoder and discriminator

Boundary Equilibrium Generative Adversarial Networks **(BEGAN)**

<https://arxiv.org/abs/1703.10717>

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x; \theta_D) - \mathcal{L}(G(z_D; \theta_G); \theta_D) & \text{for } \theta_D \\ \mathcal{L}_G = -\mathcal{L}_D & \text{for } \theta_G \end{cases}$$



Equilibrium at:
 $\mathbb{E} [\mathcal{L}(x)] = \mathbb{E} [\mathcal{L}(G(z))]$

Boundary Equilibrium Generative Adversarial Networks (BEGAN)

<https://arxiv.org/abs/1703.10717>

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x; \theta_D) - \mathcal{L}(G(z_D; \theta_G); \theta_D) & \text{for } \theta_D \\ \mathcal{L}_G = -\mathcal{L}_D & \text{for } \theta_G \end{cases}$$



Equilibrium at:
 $\mathbb{E} [\mathcal{L}(x)] = \mathbb{E} [\mathcal{L}(G(z))]$

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x) - k_t \cdot \mathcal{L}(G(z_D)) & \text{for } \theta_D \\ \mathcal{L}_G = \mathcal{L}(G(z_G)) & \text{for } \theta_G \\ k_{t+1} = k_t + \lambda_k (\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))) & \text{for each training step } t \end{cases}$$



Equilibrium at:
 $\gamma = \frac{\mathbb{E} [\mathcal{L}(G(z))]}{\mathbb{E} [\mathcal{L}(x)]}$

$\gamma \in [0, 1]$ — parameter to control the trade-off
between quality (low γ) and diversity (high γ)

Boundary Equilibrium Generative Adversarial Networks (BEGAN)

<https://arxiv.org/abs/1703.10717>



Figure 3: Random 64x64 samples at varying $\gamma \in \{0.3, 0.5, 0.7\}$

Outline

- Revising GANs
 - Basic idea, problems with GANs
- Wasserstein GAN
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- **Side notes**
 - GANs on weighted samples, conditionals
- Notable GAN architectures
 - pix2pix, CycleGAN, Progressive growing
- Evaluating GAN performance
- Applications in HEP

Weighted samples

Often in physics distributions come with sample weights

Weighted samples

Often in physics distributions come with sample weights

As our loss functions are defined in terms of expectations, the weights can be accommodated easily:

$$\mathcal{L} = \mathbb{E}_{x \sim P} f(x) + \dots = \frac{\sum_i w_i f(x_i)}{\sum_j w_j} + \dots$$

Weighted samples

Often in physics distributions come with sample weights

As our loss functions are defined in terms of expectations, the weights can be accommodated easily:

$$\mathcal{L} = \mathbb{E}_{x \sim P} f(x) + \dots = \frac{\sum_i w_i f(x_i)}{\sum_j w_j} + \dots$$

This can lead to less stable training if weights are allowed to be negative, e. g. from sPlot

Weighted samples

Often in physics distributions come with sample weights

As our loss functions are defined in terms of expectations, the weights can be accommodated easily:

$$\mathcal{L} = \mathbb{E}_{x \sim P} f(x) + \dots = \frac{\sum_i w_i f(x_i)}{\sum_j w_j} + \dots$$

This can lead to less stable training if weights are allowed to be negative, e. g. from sPlot

See this work for an alternative approach of treating negative sample weights: <https://arxiv.org/abs/1905.11719>

Conditional distributions

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

Conditional distributions

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

Most of the times

Conditional distributions

~~Somotimes~~ 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

~~Somotimes~~ 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

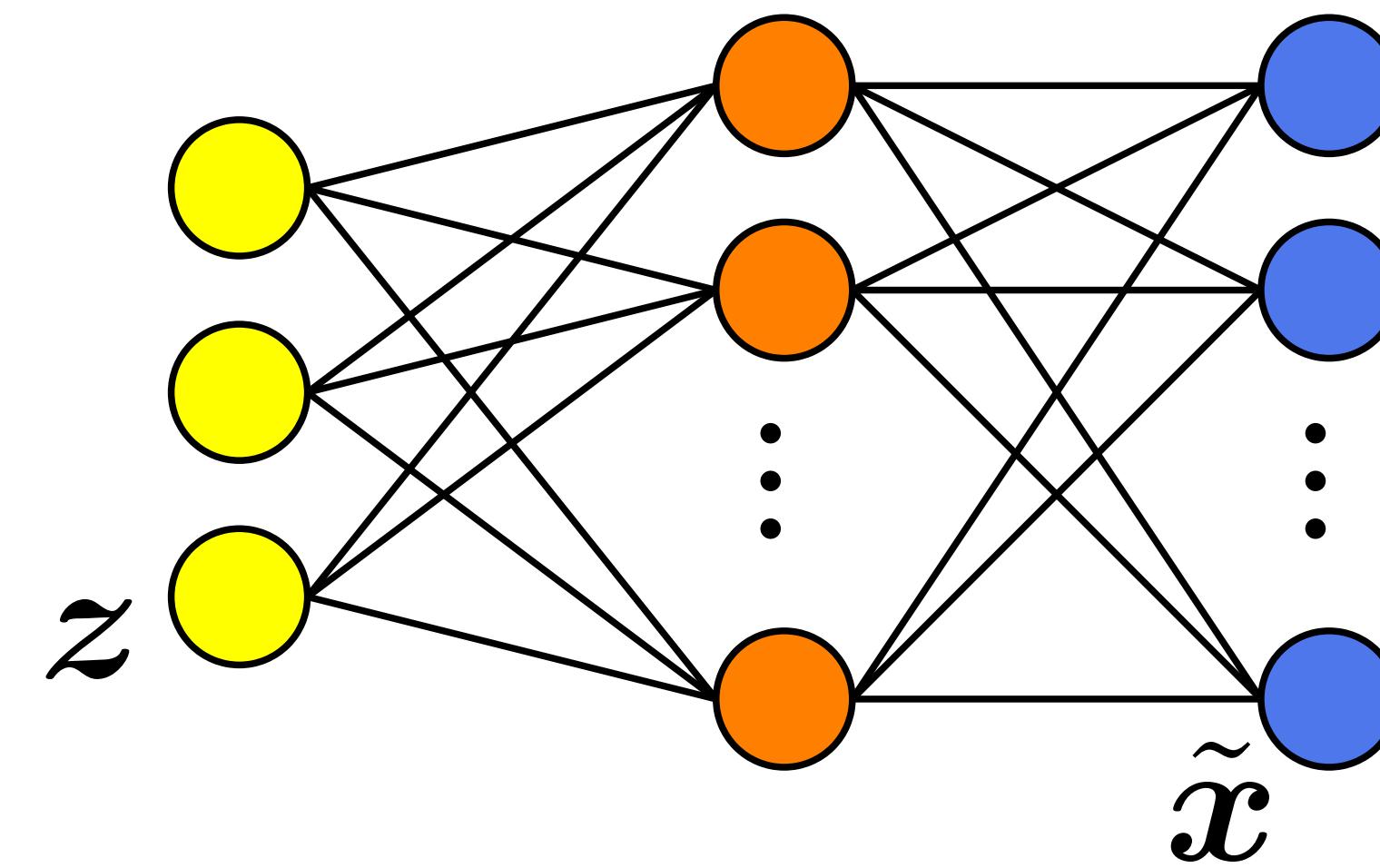
This can be achieved with any GAN architecture:

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

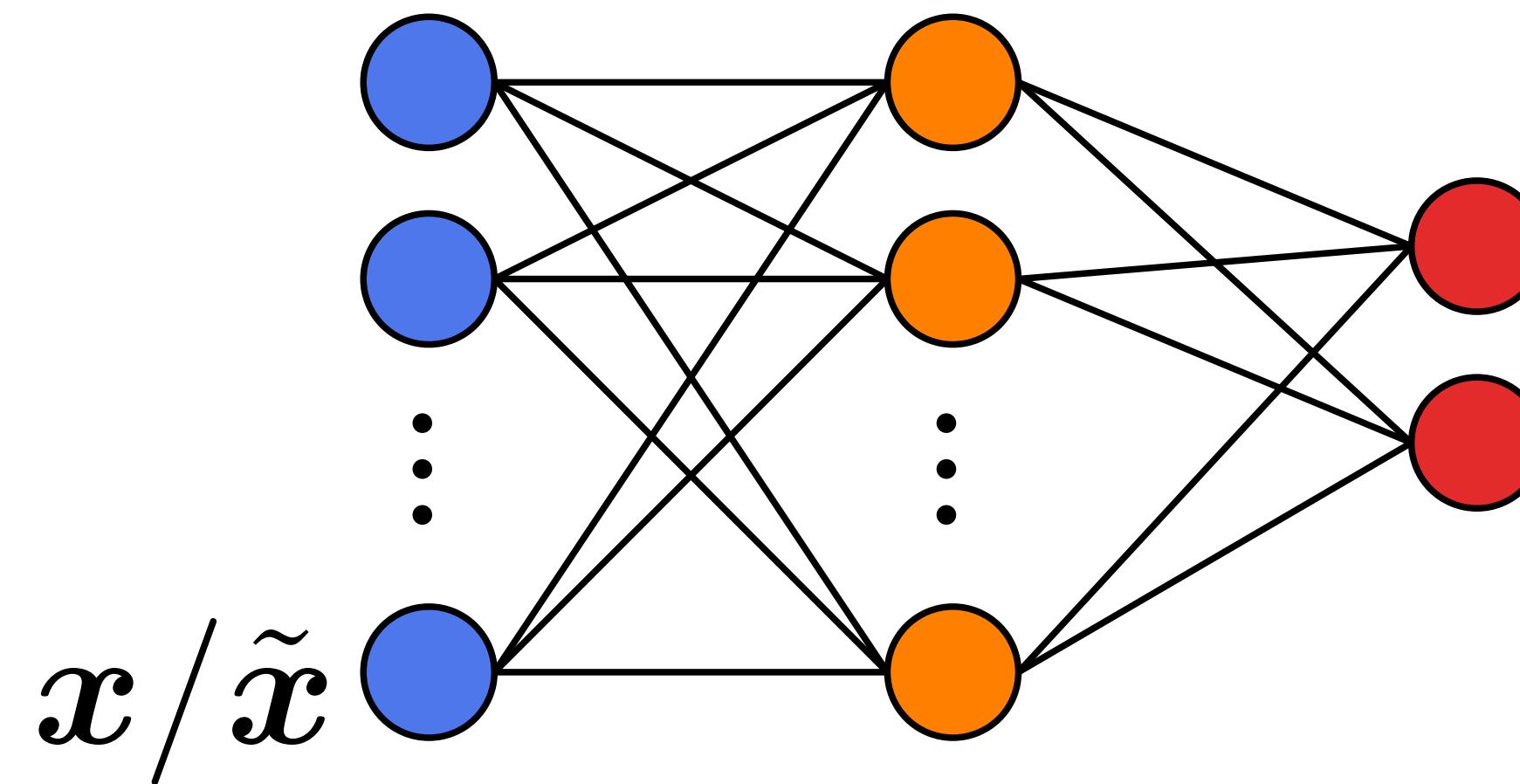
Conditional distributions

(simple case — fully-connected layers)

Generator

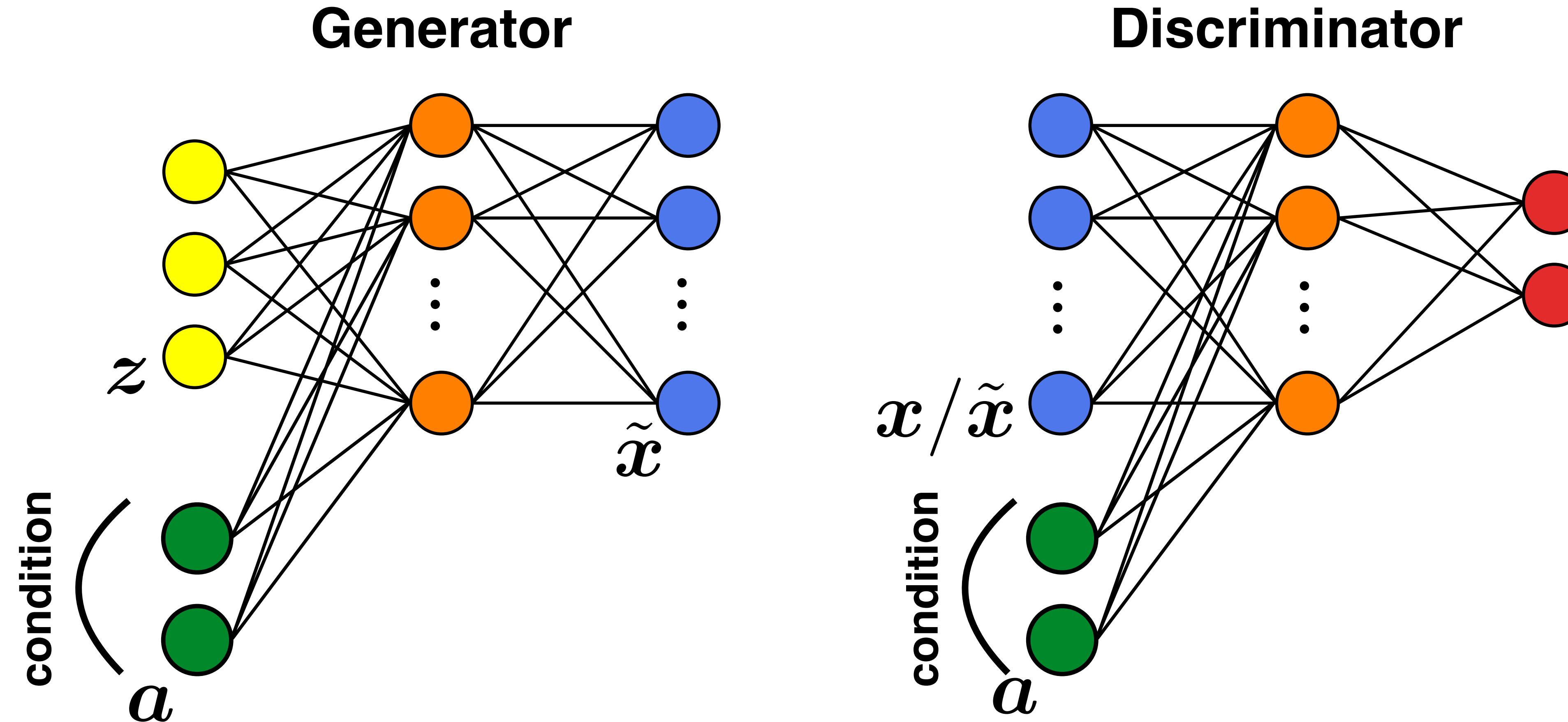


Discriminator



Conditional distributions

(simple case — fully-connected layers)

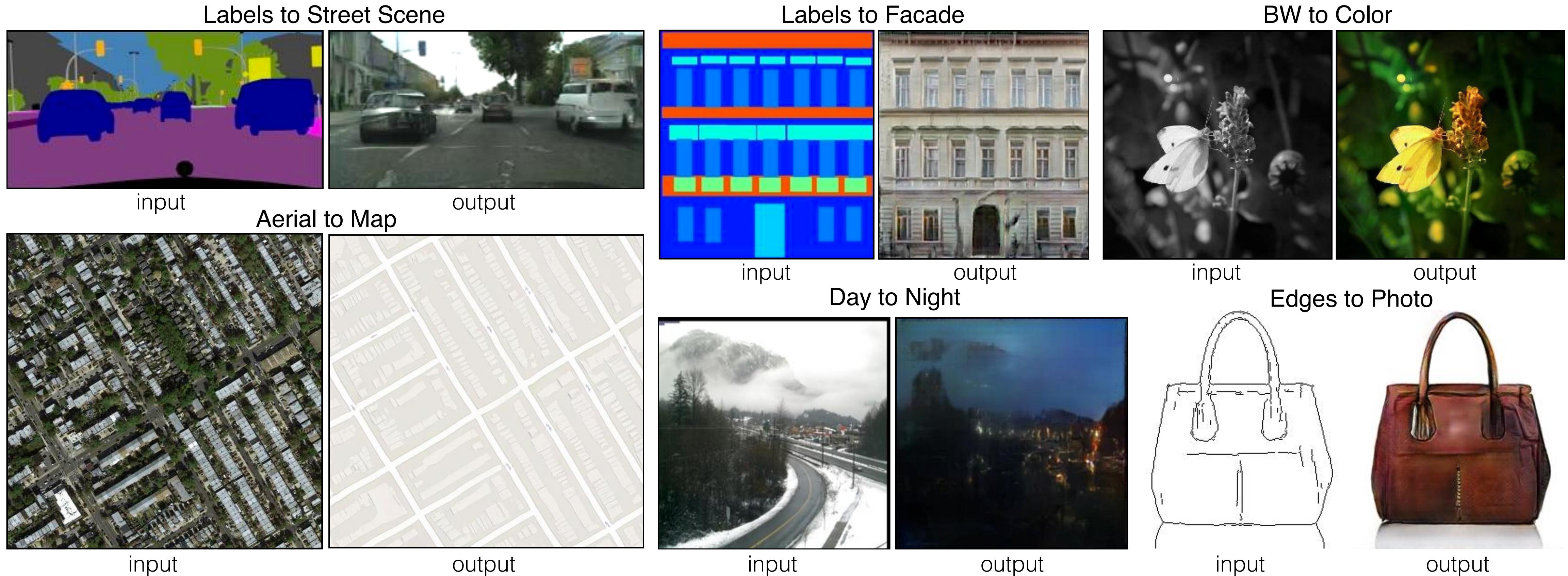


Outline

- Revising GANs
 - Basic idea, problems with GANs
- Wasserstein GAN
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- Side notes
 - GANs on weighted samples, conditionals
- **Notable GAN architectures**
 - pix2pix, CycleGAN, Progressive growing
- Evaluating GAN performance
- Applications in HEP

pix2pix

<https://arxiv.org/abs/1611.07004>



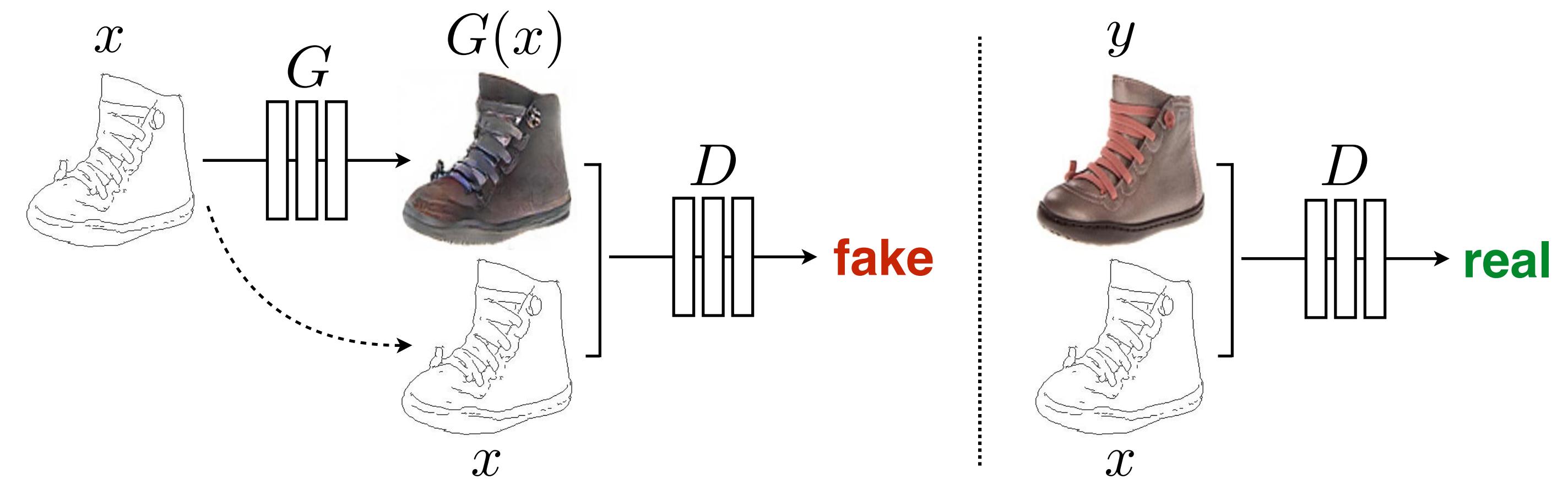
pix2pix

<https://arxiv.org/abs/1611.07004>

$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1].$$

Our final objective is

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$



pix2pix

<https://arxiv.org/abs/1611.07004>

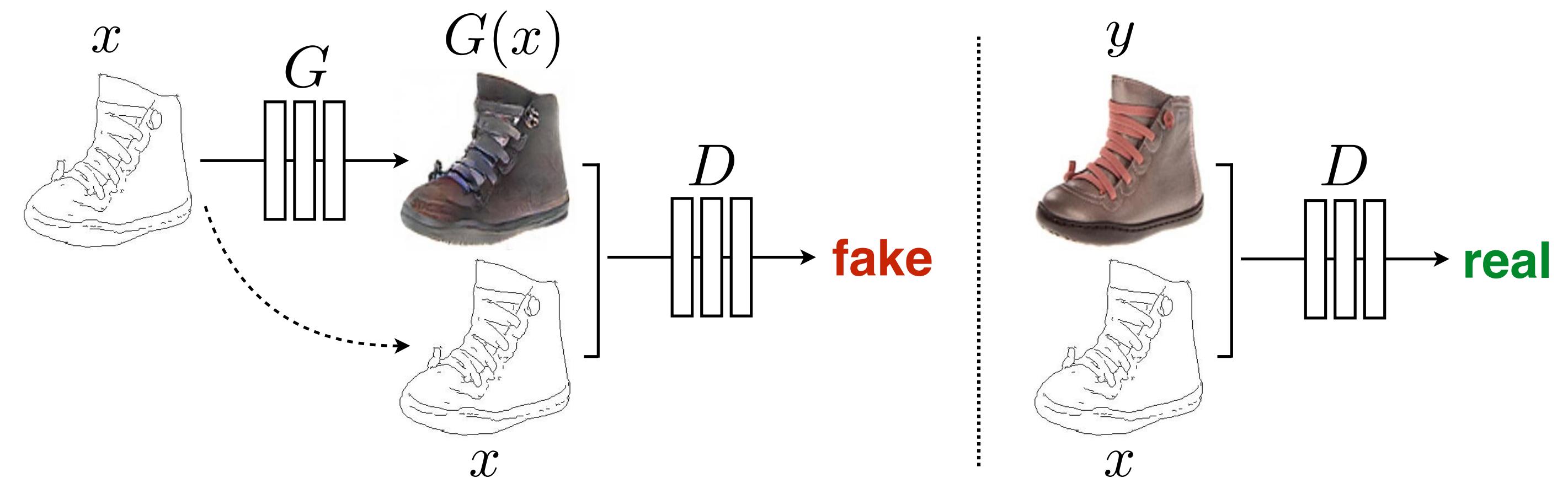
$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1].$$

Our final objective is

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$

L1 loss term to
capture low-frequency
information

- Discriminator doesn't need to classify the entire image
- Instead it 'convolutionally' scans smaller patches of the image and averages the result



pix2pix

<https://arxiv.org/abs/1611.07004>

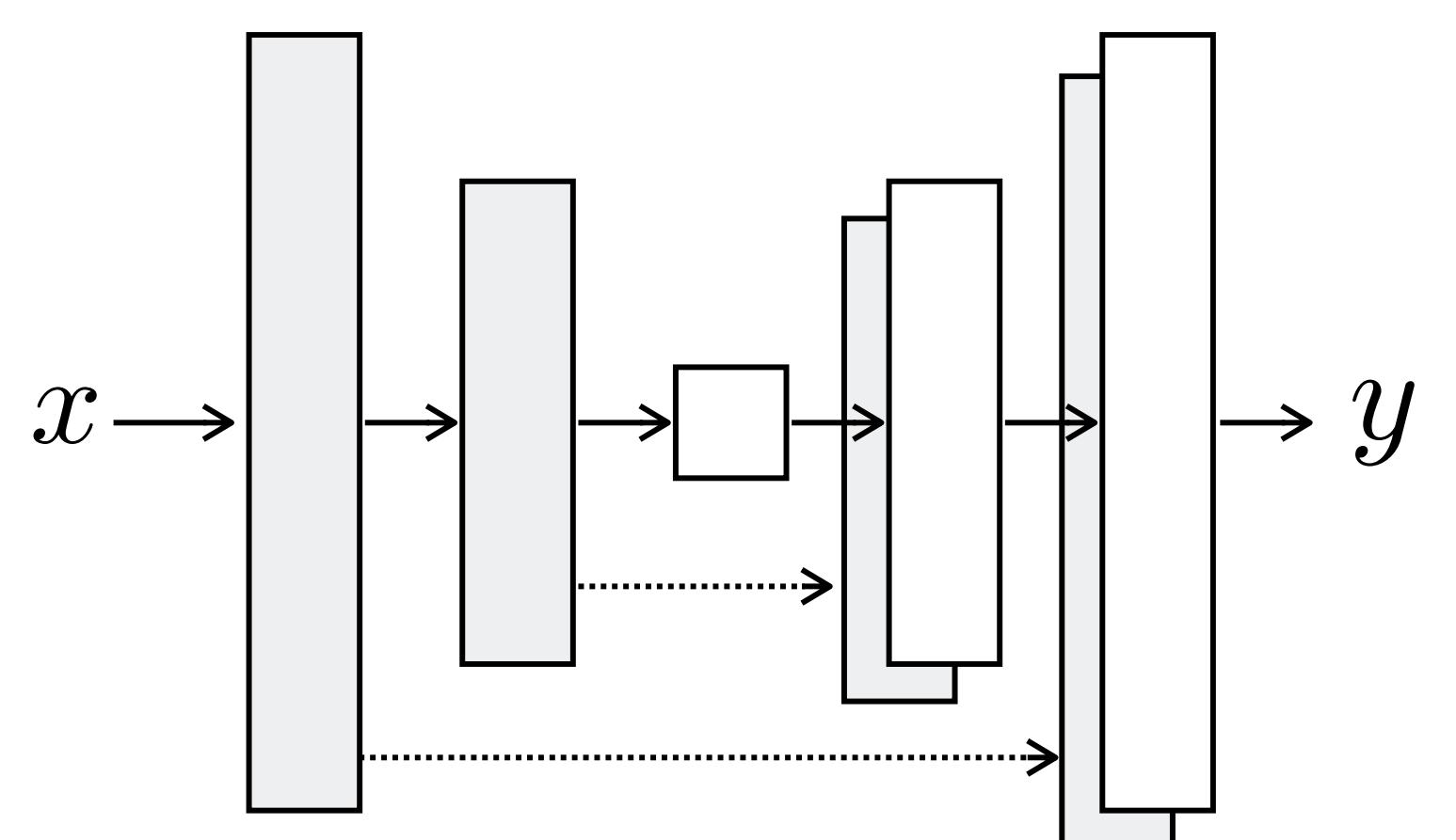
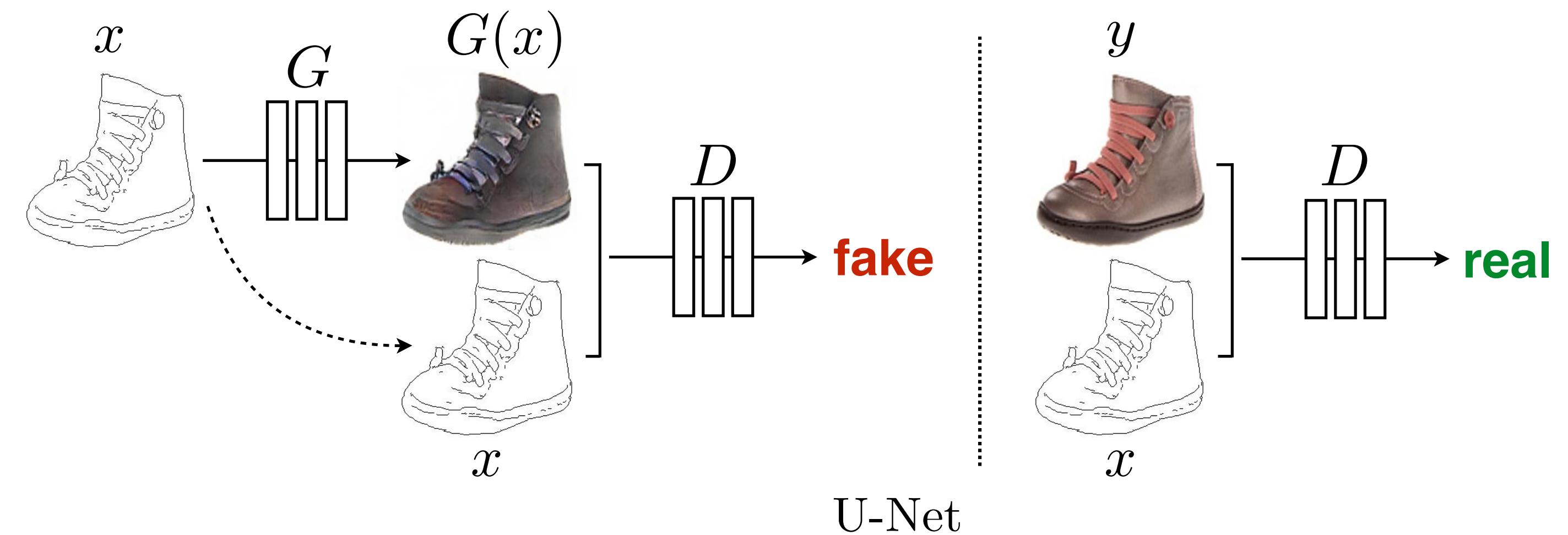
$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1].$$

Our final objective is

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$

L1 loss term to
capture low-frequency
information

- Discriminator doesn't need to classify the entire image
- Instead it 'convolutionally' scans smaller patches of the image and averages the result



'U-Net' generator architecture

pix2pix

<https://arxiv.org/abs/1611.07004>

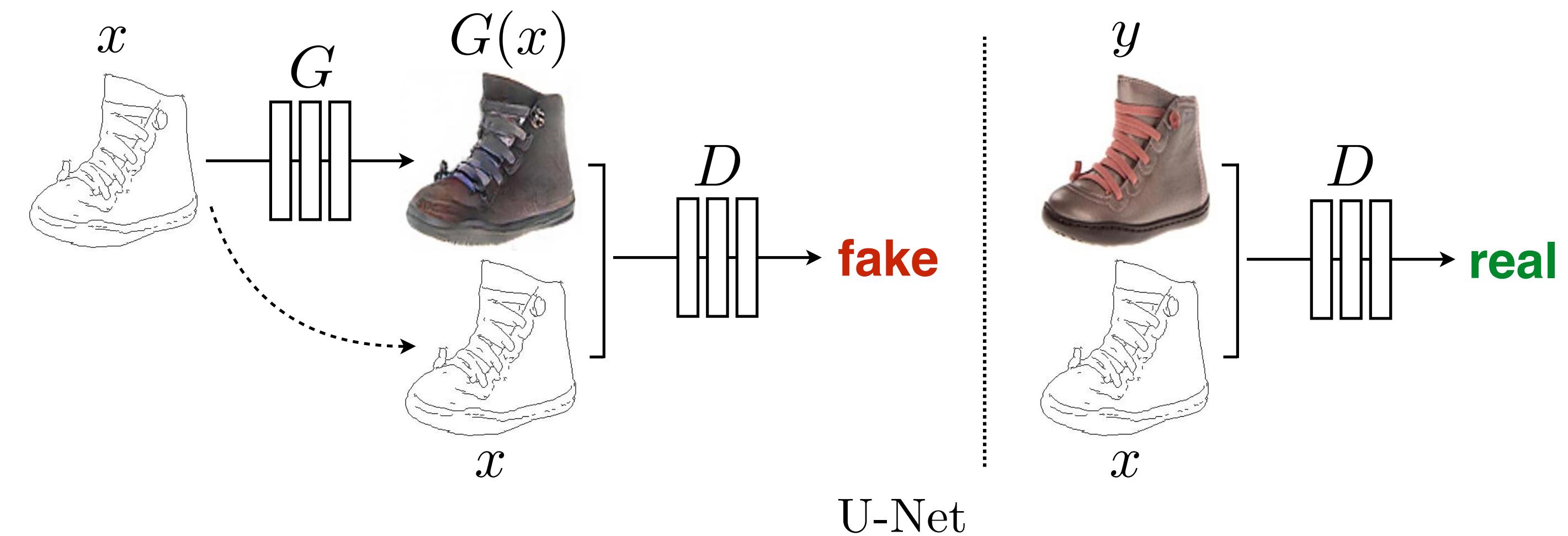
$$\mathcal{L}_{L1}(G) = \mathbb{E}_{x,y,z}[\|y - G(x, z)\|_1].$$

Our final objective is

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$

L1 loss term to
capture low-frequency
information

- Discriminator doesn't need to classify the entire image
- Instead it 'convolutionally' scans smaller patches of the image and averages the result



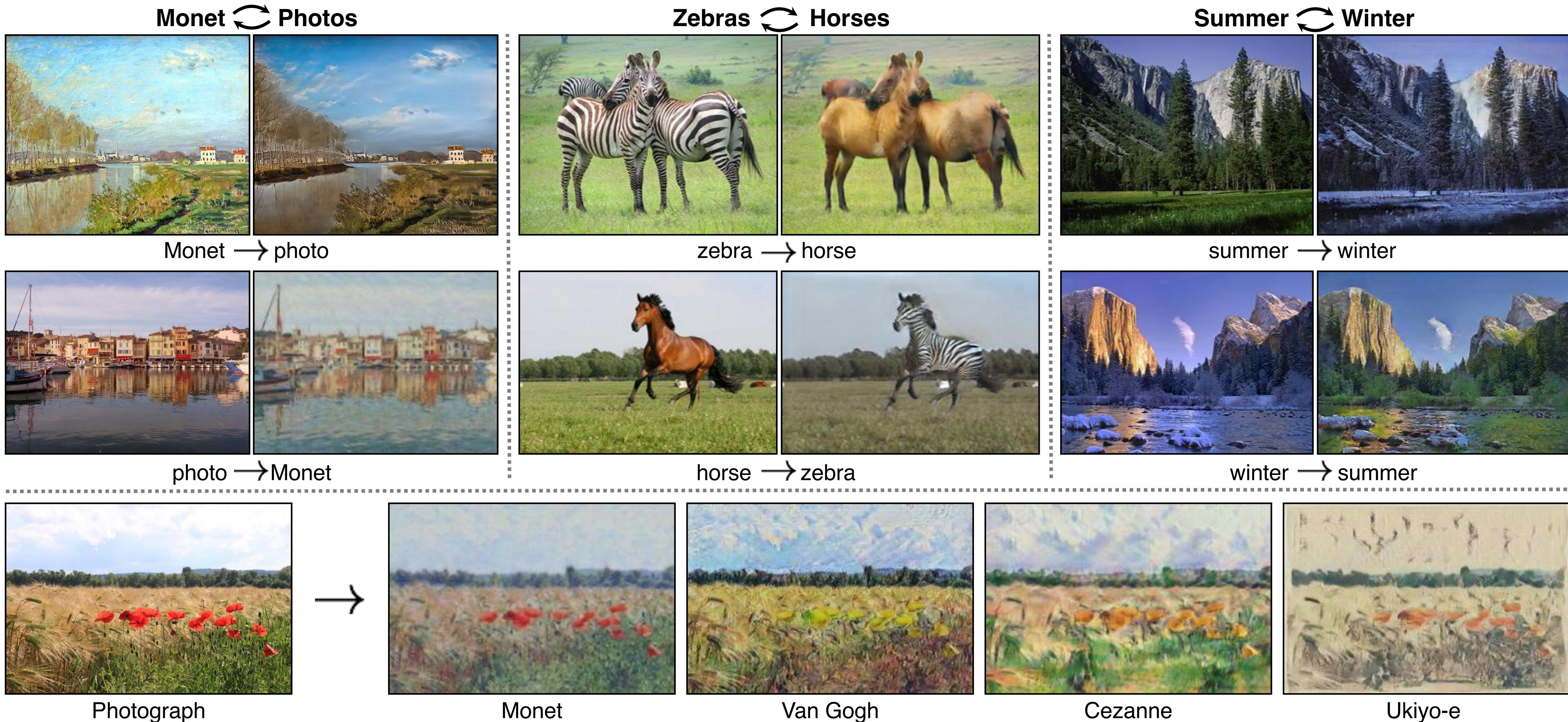
No noise input to the generator:
«The generator simply **learned to ignore the noise** ...
Instead we provide noise only in the form of **dropout**,
applied at both training and test time»



'U-Net' generator
architecture

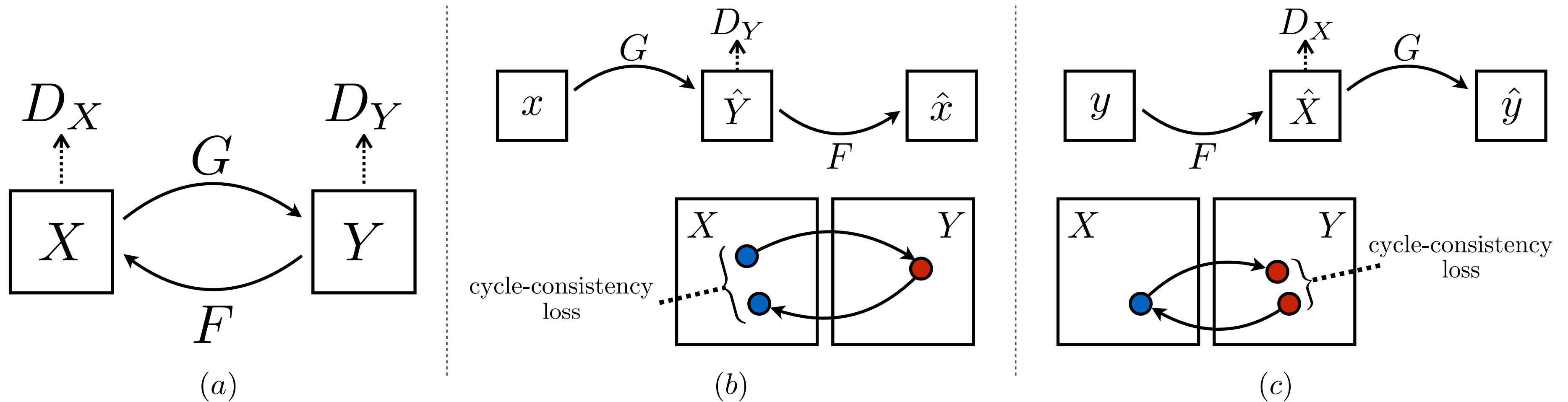
CycleGAN

<https://arxiv.org/abs/1703.10593>



CycleGAN

<https://arxiv.org/abs/1703.10593>



$$\begin{aligned} \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) &= \mathbb{E}_{y \sim p_{\text{data}}(y)} [\log D_Y(y)] \\ &+ \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log(1 - D_Y(G(x)))] \end{aligned} \quad (1)$$

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) &= \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ &+ \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1]. \end{aligned} \quad (2)$$

$$\begin{aligned} \mathcal{L}(G, F, D_X, D_Y) &= \mathcal{L}_{\text{GAN}}(G, D_Y, X, Y) \\ &+ \mathcal{L}_{\text{GAN}}(F, D_X, Y, X) \\ &+ \lambda \mathcal{L}_{\text{cyc}}(G, F), \end{aligned} \quad (3)$$

Progressive Growing of GANs

<https://arxiv.org/abs/1710.10196>

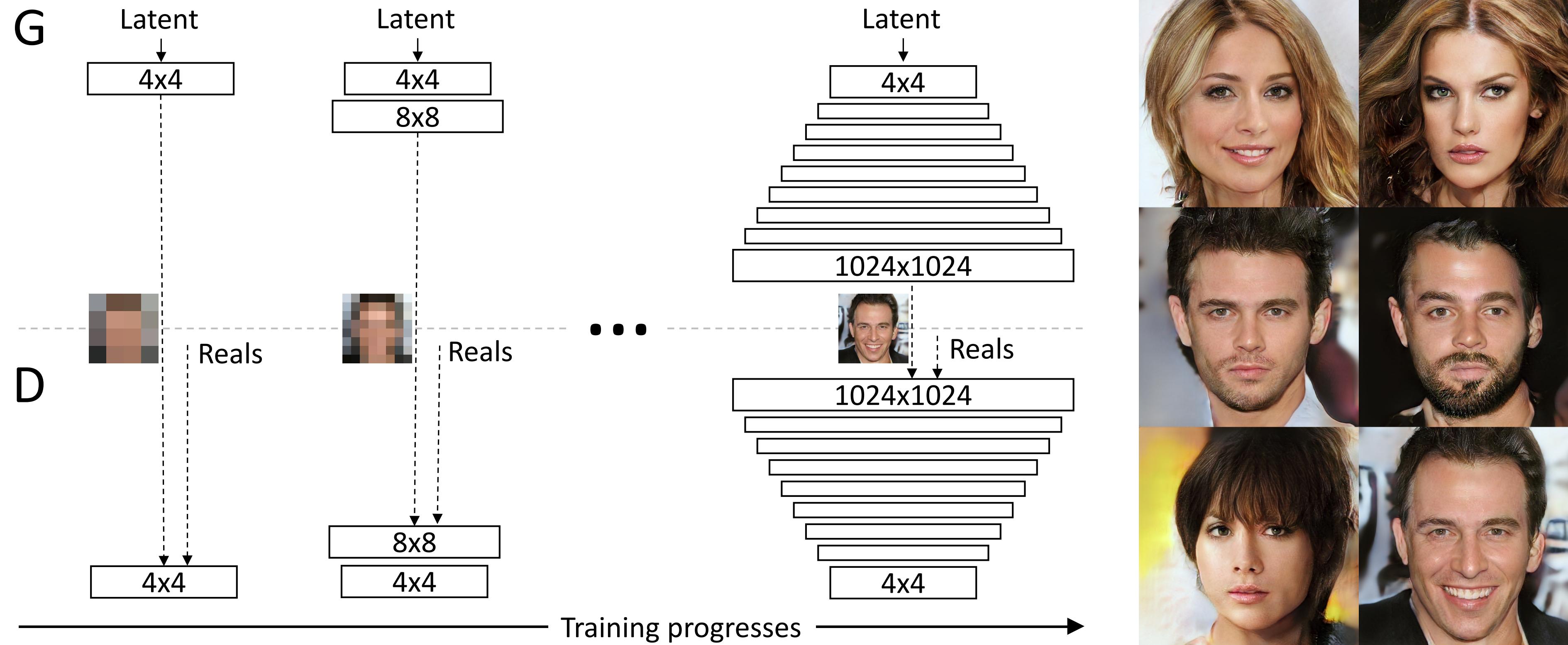


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here $N \times N$ refers to convolutional layers operating on $N \times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. On the right we show six example images generated using progressive growing at 1024×1024 .

Progressive Growing of GANs

<https://arxiv.org/abs/1710.10196>

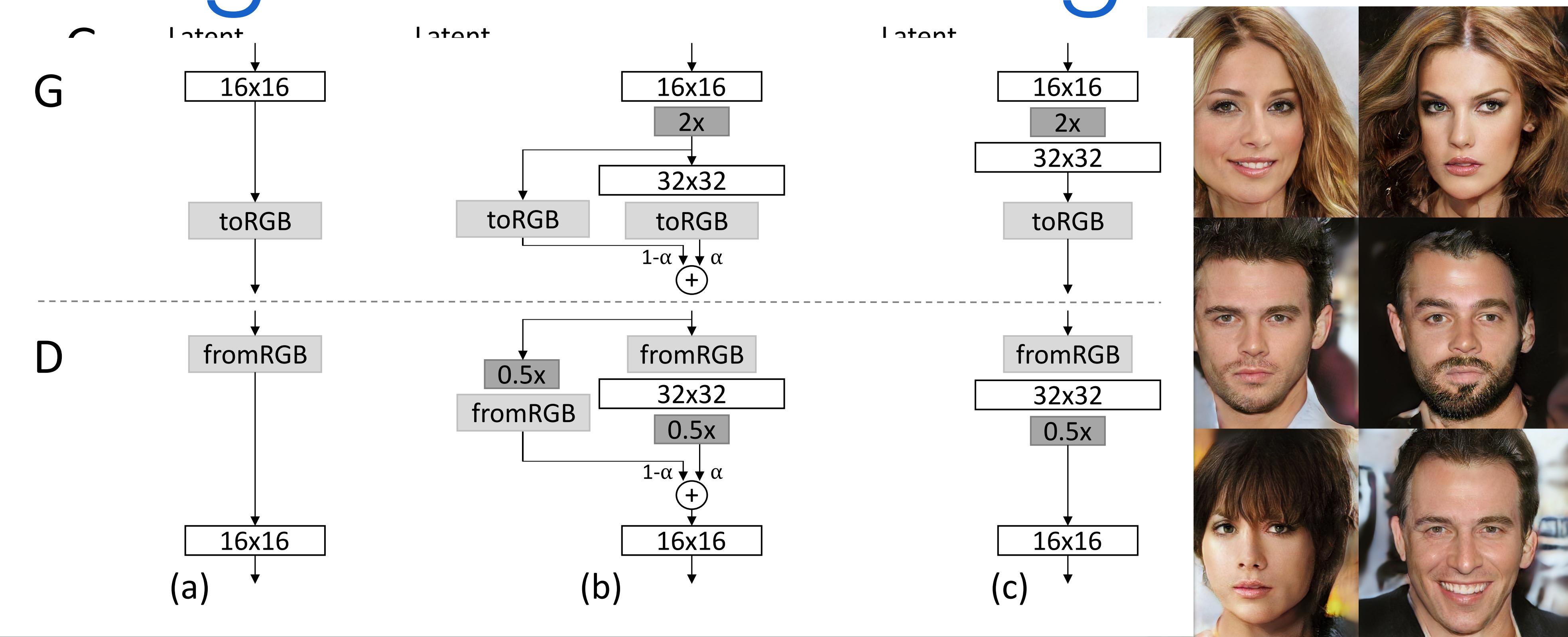


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of 4×4 pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here $N \times N$ refers to convolutional layers operating on $N \times N$ spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. On the right we show six example images generated using progressive growing at 1024×1024 .

Outline

- Revising GANs
 - Basic idea, problems with GANs
- Wasserstein GAN
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- Side notes
 - GANs on weighted samples, conditionals
- Notable GAN architectures
 - pix2pix, CycleGAN, Progressive growing
- **Evaluating GAN performance**
- Applications in HEP

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
 - Efficiencies
 - 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
 - Efficiencies
 - 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 $p(y) = \int p(y|x = G(z))dz$
 - 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 $p(y) = \int p(y|x = G(z))dz$
 - 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 $(0, 0)$.

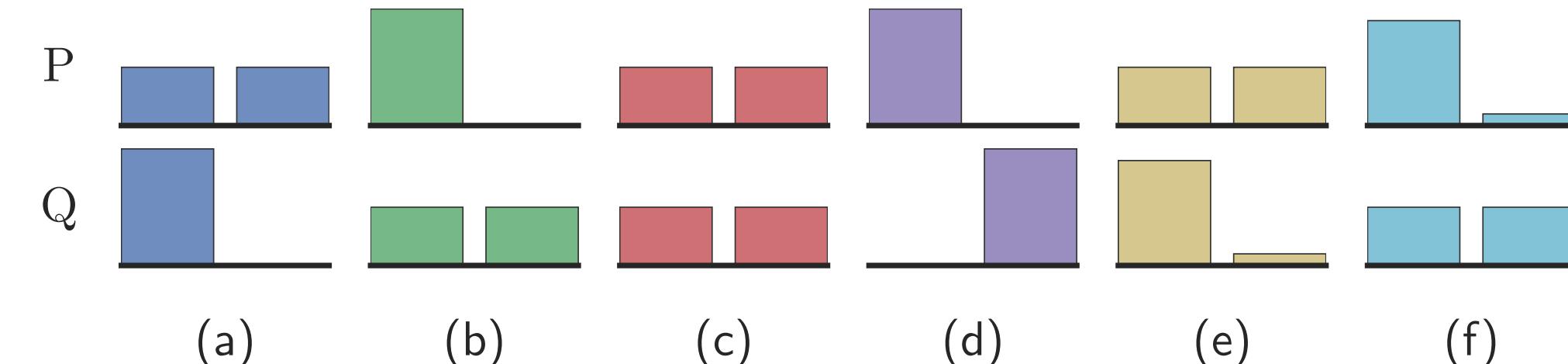


Figure 2: Intuitive examples of P and Q .

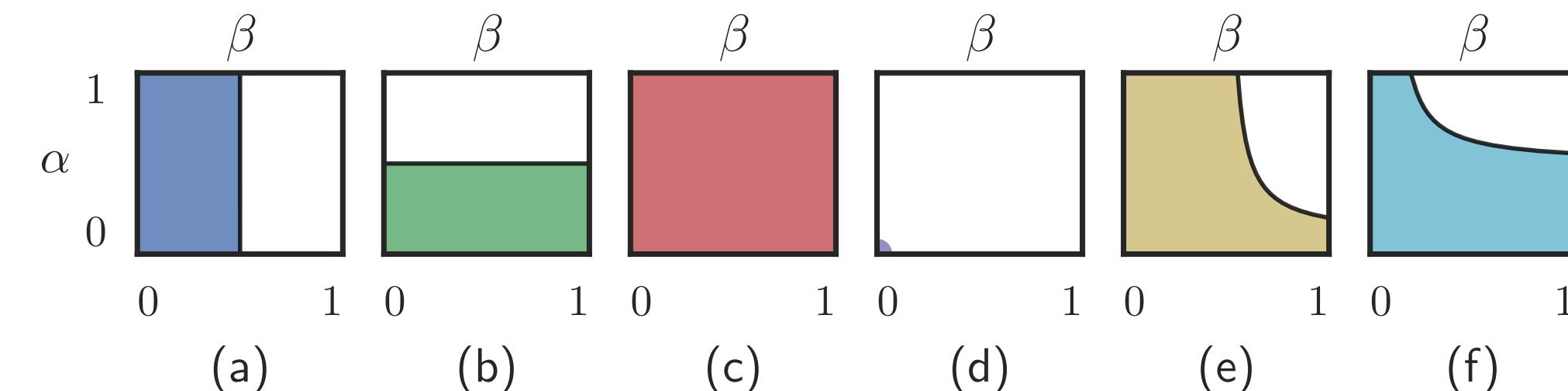


Figure 3: $\text{PRD}(Q, P)$ for the examples above.

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

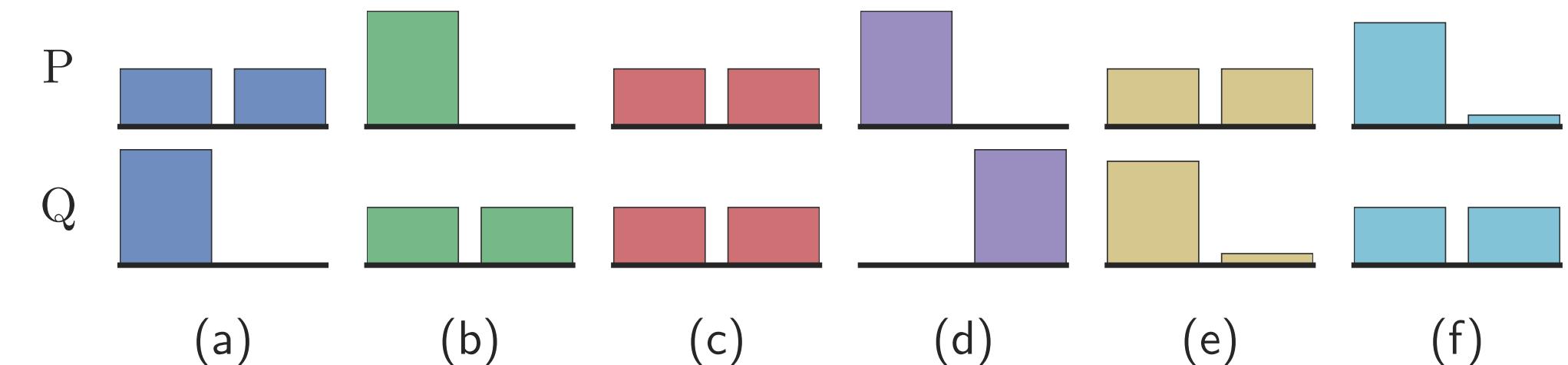


Figure 2: Intuitive examples of P and Q .

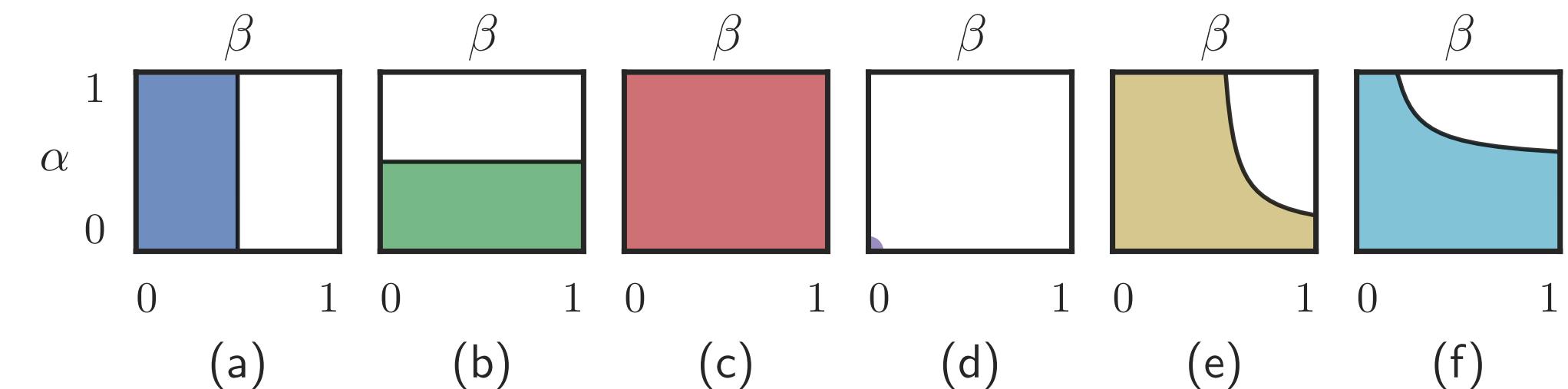


Figure 3: $\text{PRD}(Q, P)$ for the examples above.

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

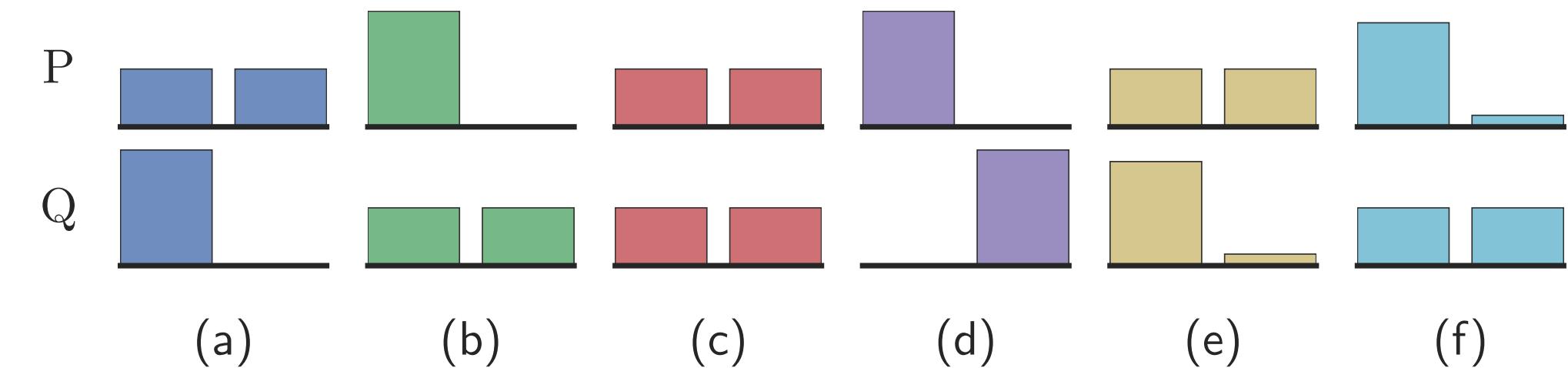


Figure 2: Intuitive examples of P and Q .

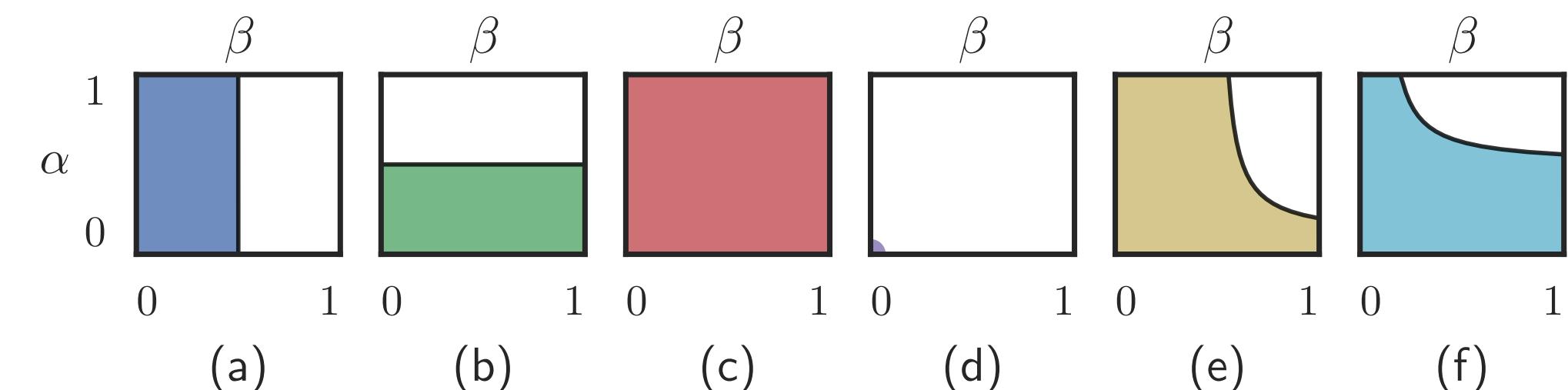


Figure 3: $\text{PRD}(Q, P)$ for the examples above.

Evaluating generative models

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

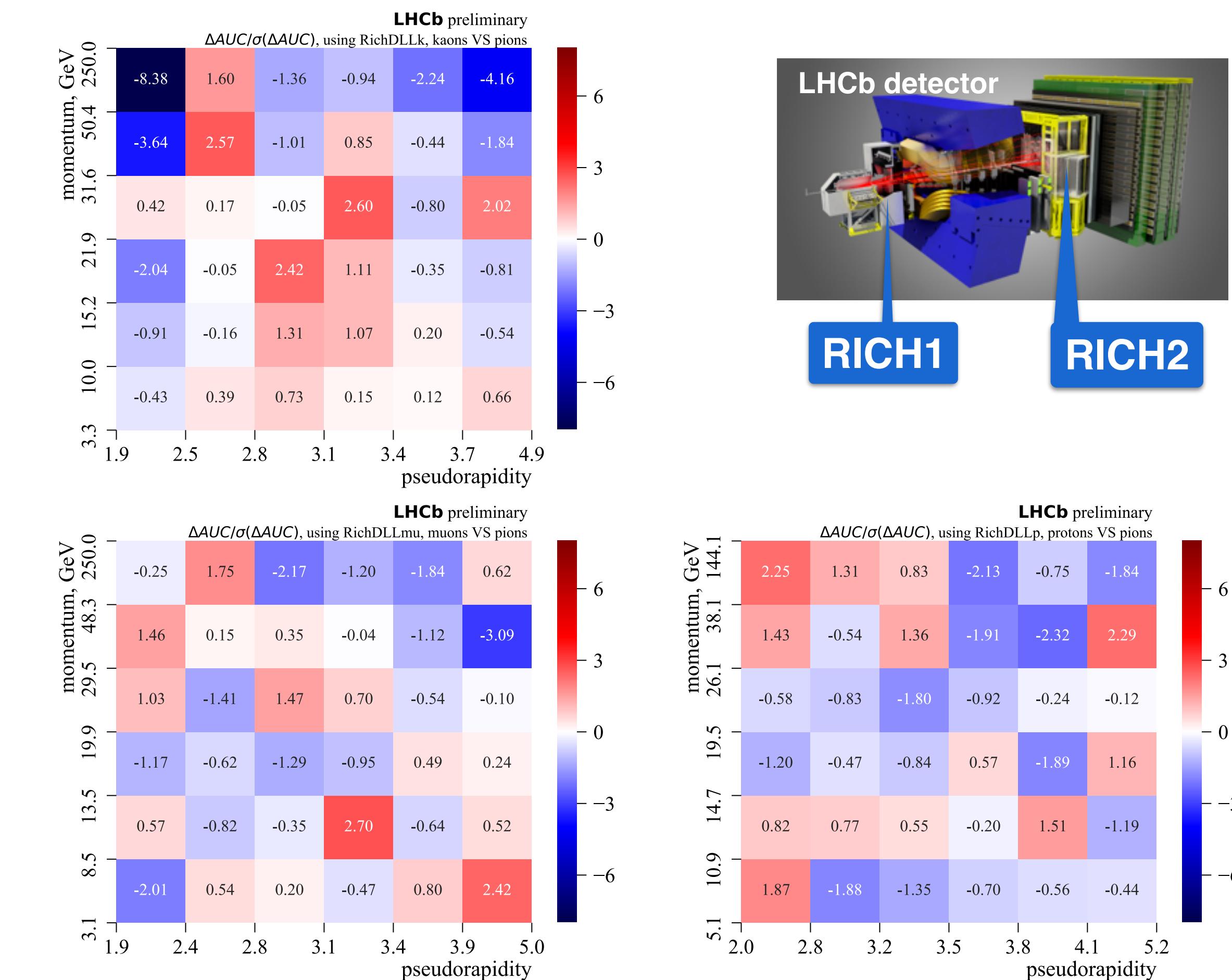
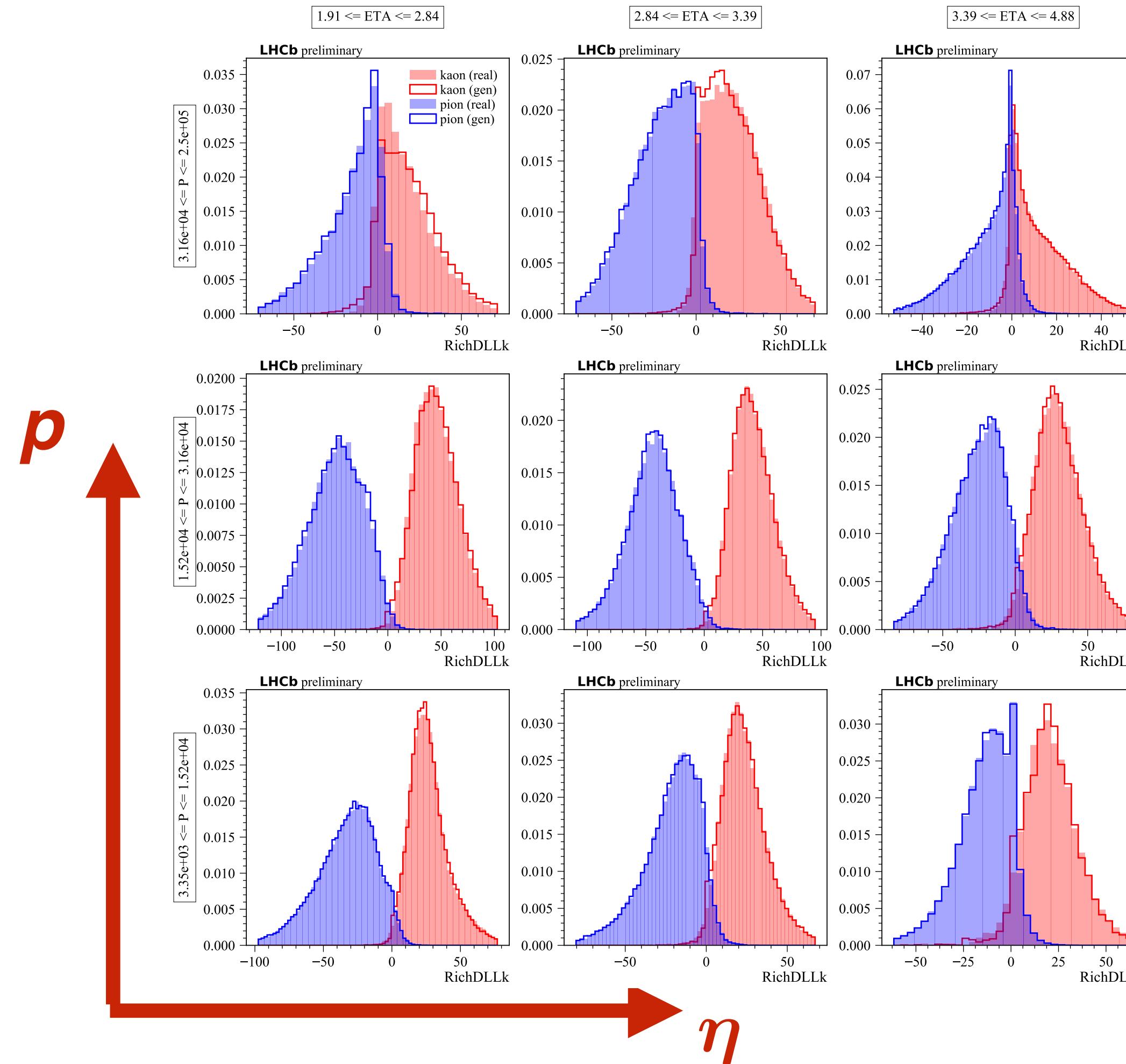
<https://arxiv.org/abs/1802.03446>

	Measure	Description
Quantitative	1. Average Log-likelihood [18, 22]	<ul style="list-style-type: none"> Log likelihood of explaining realworld held out/test data using a density estimated from the generated data (<i>e.g.</i> using KDE or Parzen window estimation). $L = \frac{1}{N} \sum_i \log P_{model}(\mathbf{x}_i)$
	2. Coverage Metric [33]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> KLD between conditional and marginal label distributions over generated data. $\exp(\mathbb{E}_{\mathbf{x}} [\text{KL}(p(y \mathbf{x}) \ p(y))])$
	4. Modified Inception Score (m-IS) [34]	<ul style="list-style-type: none"> Encourages diversity within images sampled from a particular category. $\exp(\mathbb{E}_{\mathbf{x}_i} [\mathbb{E}_{\mathbf{x}_j} ([\text{KL}(P(y \mathbf{x}_i) \ P(y \mathbf{x}_j))])])$
	5. Mode Score (MS) [35]	<ul style="list-style-type: none"> Similar to IS but also takes into account the prior distribution of the labels over real data. $\exp(\mathbb{E}_{\mathbf{x}} [\text{KL}(p(y \mathbf{x}) \ p(y^{train}))] - \text{KL}(p(y) \ p(y^{train})))$
	6. AM Score [36]	<ul style="list-style-type: none"> Takes into account the KLD between distributions of training labels vs. predicted labels, as well as the entropy of predictions. $\text{KL}(p(y \mathbf{x}) \ p(y)) + \mathbb{E}_{\mathbf{x}} [H(y \mathbf{x})]$
	7. Fréchet Inception Distance (FID) [37]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> The critic (<i>e.g.</i> 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]	<ul style="list-style-type: none"> Measures the support size of a discrete (continuous) distribution by counting the duplicates (near duplicates)
	11. Classifier Two Sample Test (C2ST) [40]	<ul style="list-style-type: none"> Answers whether two samples are drawn from the same distribution (<i>e.g.</i> by training a binary classifier)
	12. Classification Performance [1, 15]	<ul style="list-style-type: none"> An indirect technique for evaluating the quality of unsupervised representations (<i>e.g.</i> feature extraction; FCN score). See also the GAN Quality Index (GQI) [41].
	13. Boundary Distortion [42]	<ul style="list-style-type: none"> Measures diversity of generated samples and covariate shift using classification methods.
	14. Number of Statistically-Different Bins (NDB) [43]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> Measures the distributions of distances to the nearest neighbors of some query images (<i>i.e.</i> diversity)
	16. Generative Adversarial Metric (GAM) [31]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> 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 $\text{KL}(P_g(y \mathbf{x}), P_r(y \mathbf{x}))$
	20. Geometry Score [47]	<ul style="list-style-type: none"> Compares geometrical properties of the underlying data manifold between real and generated data.
	21. Reconstruction Error [48]	<ul style="list-style-type: none"> Measures the reconstruction error (<i>e.g.</i> L_2 norm) between a test image and its closest generated image by optimizing for z (<i>i.e.</i> $\min_{\mathbf{z}} \ G(\mathbf{z}) - \mathbf{x}^{(test)}\ ^2$)
	22. Image Quality Measures [49, 50, 51]	<ul style="list-style-type: none"> Evaluates the quality of generated images using measures such as SSIM, PSNR, and sharpness difference
	23. Low-level Image Statistics [52, 53]	<ul style="list-style-type: none"> 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]	<ul style="list-style-type: none"> These measures are used to quantify the degree of overfitting in GANs, often over toy datasets.
Qualitative	1. Nearest Neighbors	<ul style="list-style-type: none"> To detect overfitting, generated samples are shown next to their nearest neighbors in the training set
	2. Rapid Scene Categorization [18]	<ul style="list-style-type: none"> In these experiments, participants are asked to distinguish generated samples from real images in a short presentation time (<i>e.g.</i> 100 ms); <i>i.e.</i> real v.s fake
	3. Preference Judgment [54, 55, 56, 57]	<ul style="list-style-type: none"> Participants are asked to rank models in terms of the fidelity of their generated images (<i>e.g.</i> pairs, triples)
	4. Mode Drop and Collapse [58, 59]	<ul style="list-style-type: none"> Over datasets with known modes (<i>e.g.</i> 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]	<ul style="list-style-type: none"> Regards exploring and illustrating the internal representation and dynamics of models (<i>e.g.</i> space continuity) as well as visualizing learned features

Outline

- Revising GANs
 - Basic idea, problems with GANs
- Wasserstein GAN
 - Wasserstein distance and how to compute it
 - WGAN
 - Some issues and alternatives
- Side notes
 - GANs on weighted samples, conditionals
- Notable GAN architectures
 - pix2pix, CycleGAN, Progressive growing
- Evaluating GAN performance
- **Applications in HEP**

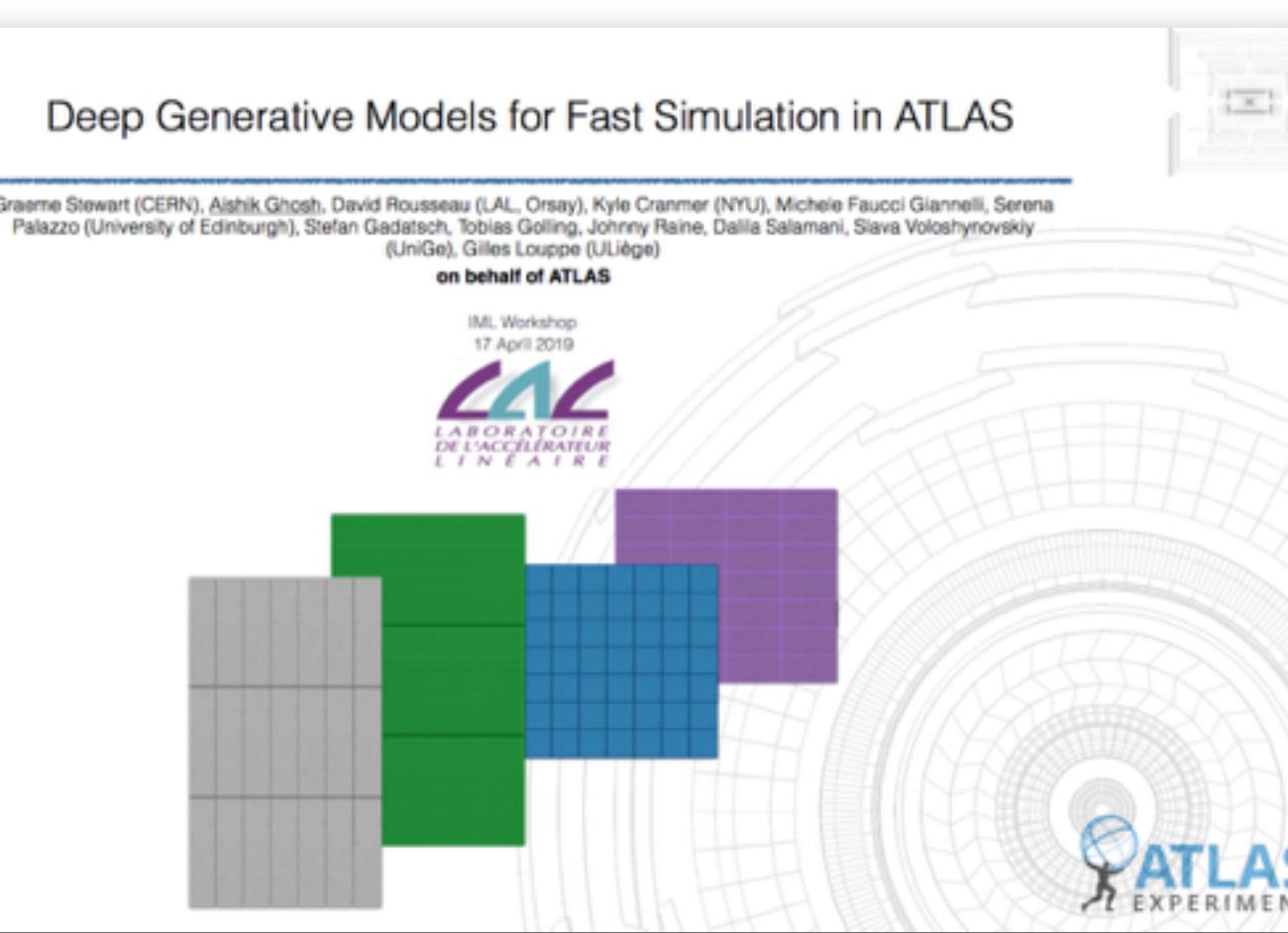
Fast simulation of Cherenkov detectors at LHCb



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

<https://arxiv.org/abs/1905.11825>

EM calorimeter shower fast simulation in ATLAS

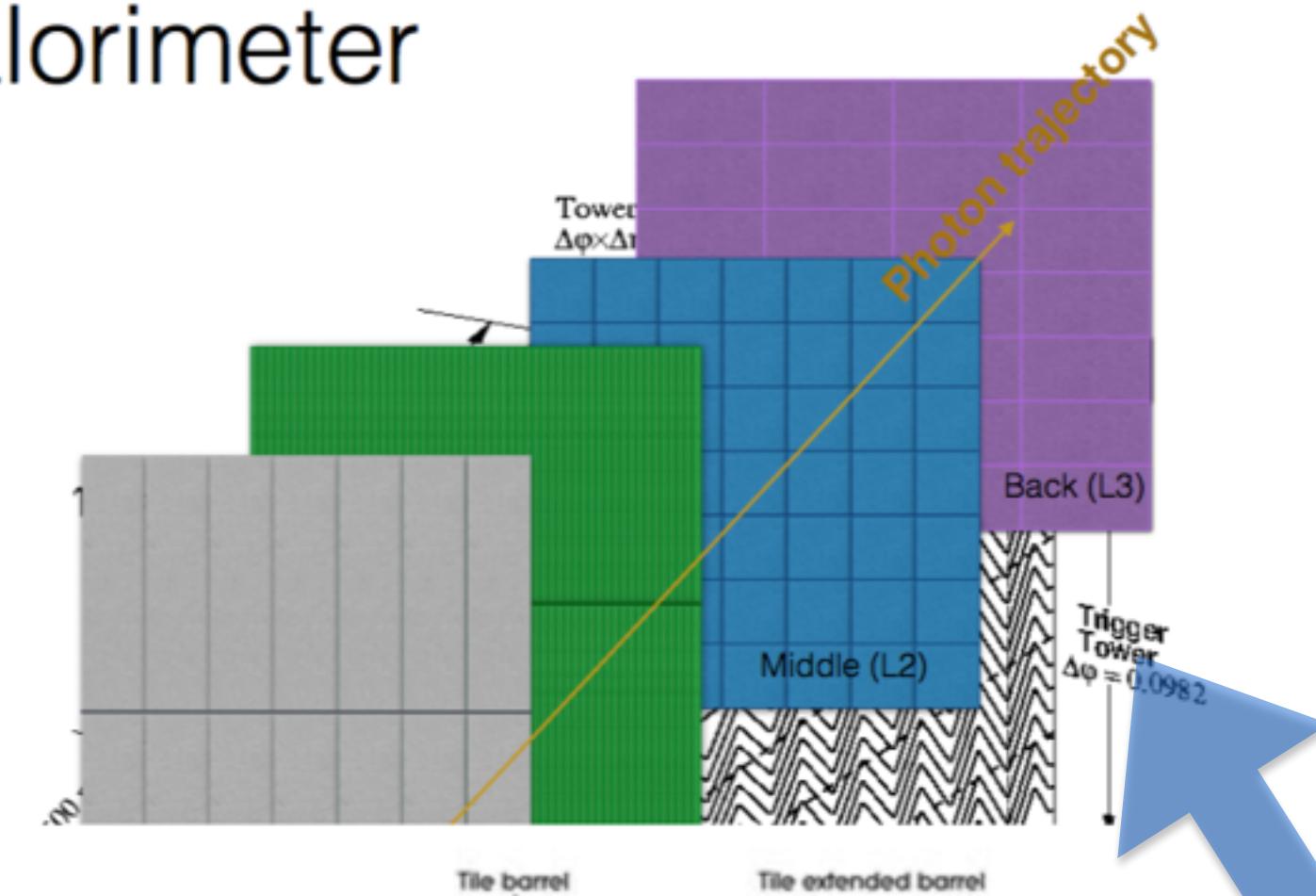


The Calorimeter

2-D Axis: ϕ vs η

Particle goes through 4 layers in this order:

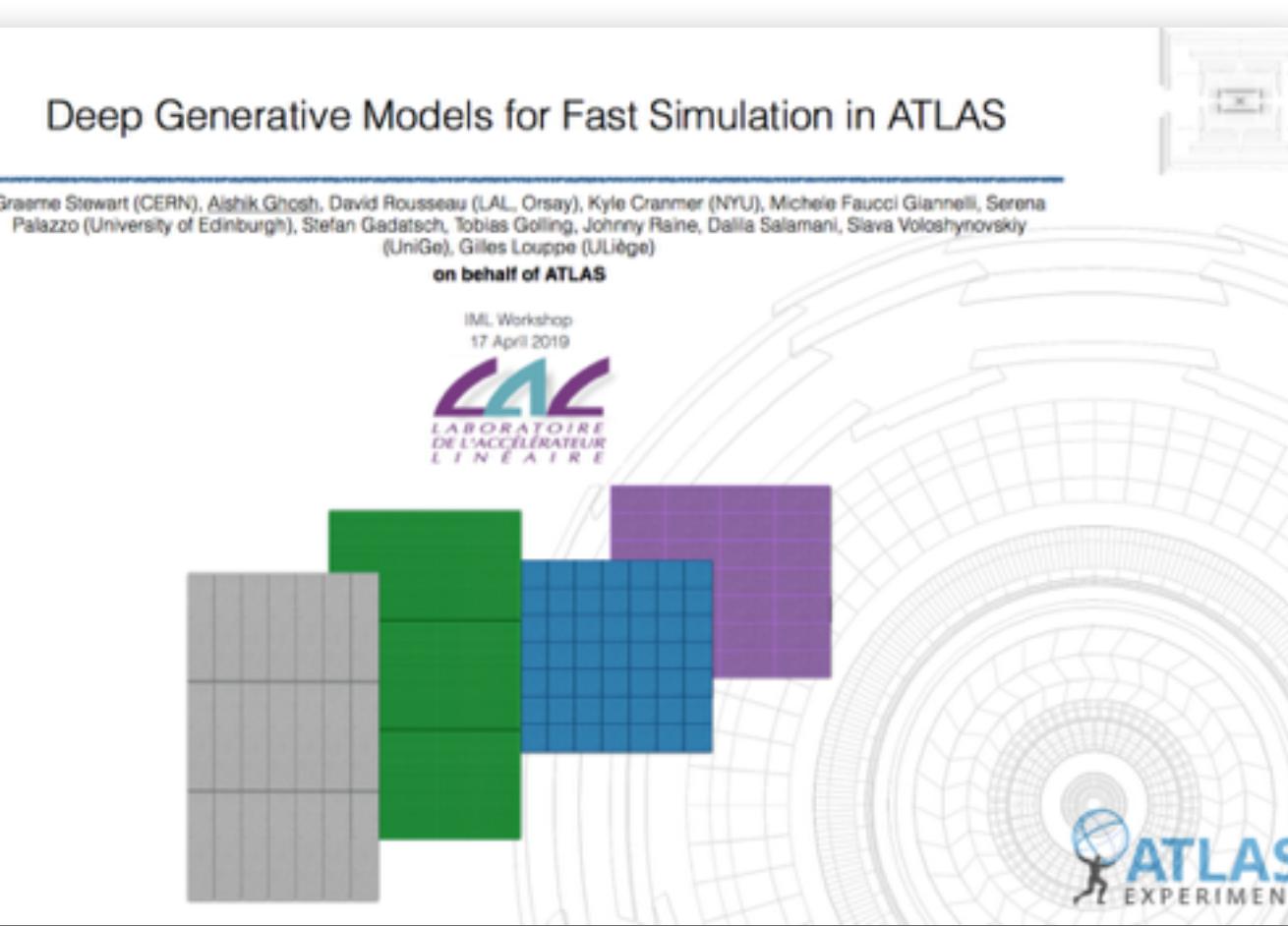
0. **Pre-Sampler**: Some energy deposit
1. **Strips**: Very granular in η ; more energy deposit
2. **Middle**: Thickest layer, maximum energy deposit
3. **Back**: Little Energy deposits



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



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

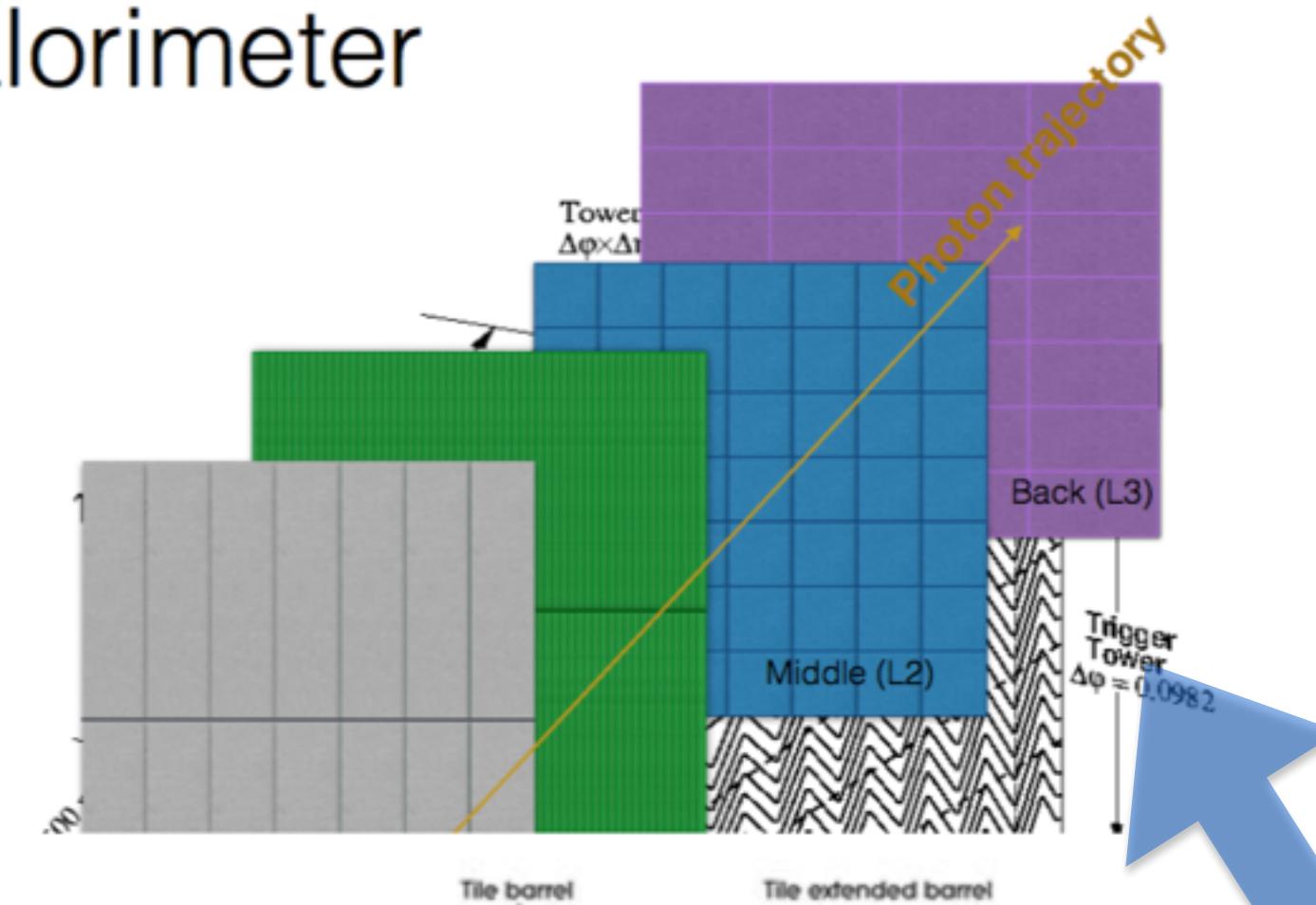
[https://indico.cern.ch/
event/766872/
contributions/3357991/](https://indico.cern.ch/event/766872/contributions/3357991/)

The Calorimeter

2-D Axis: ϕ vs η

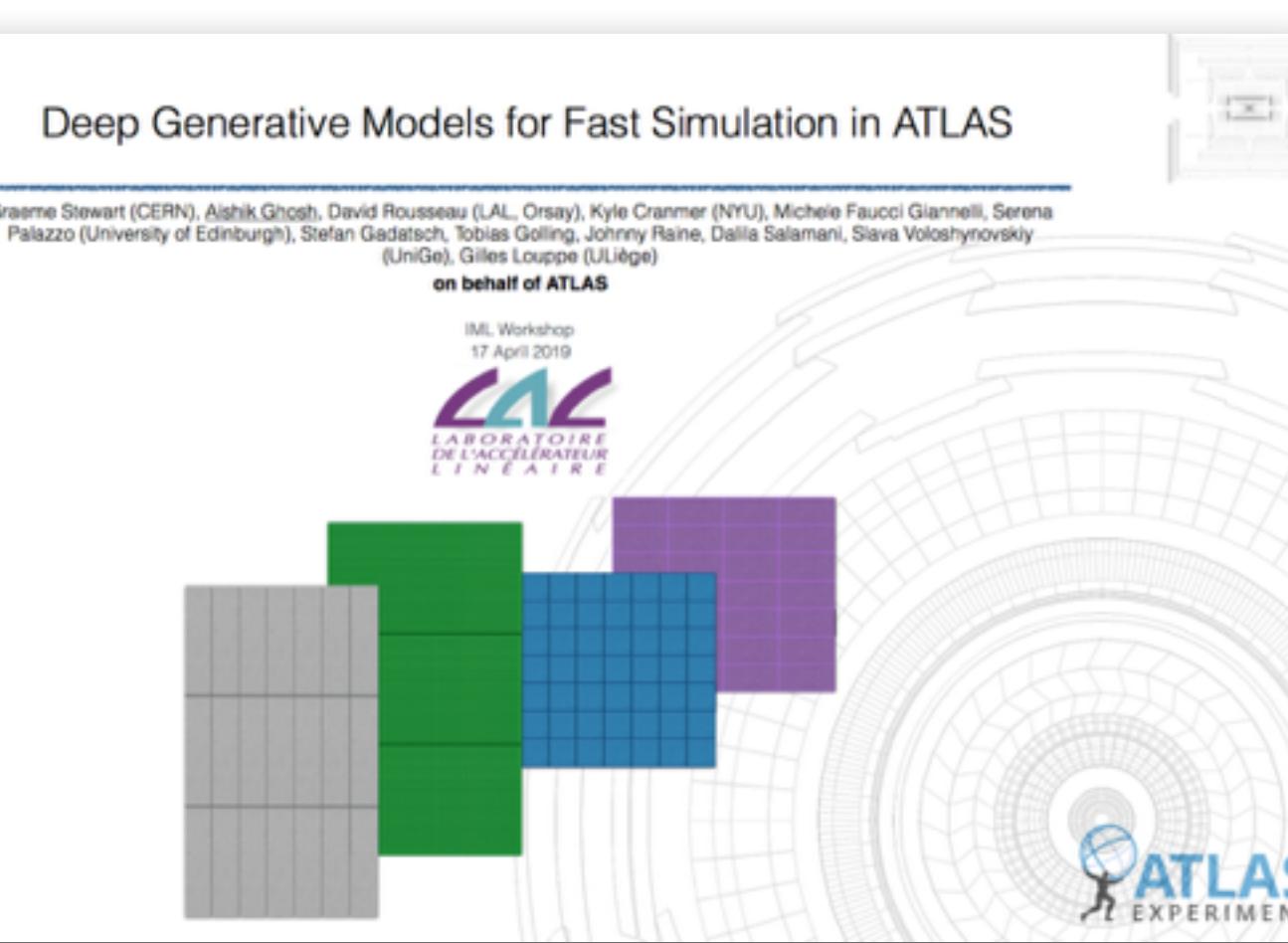
Particle goes through 4 layers in this order:

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)

EM calorimeter shower fast simulation in ATLAS



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

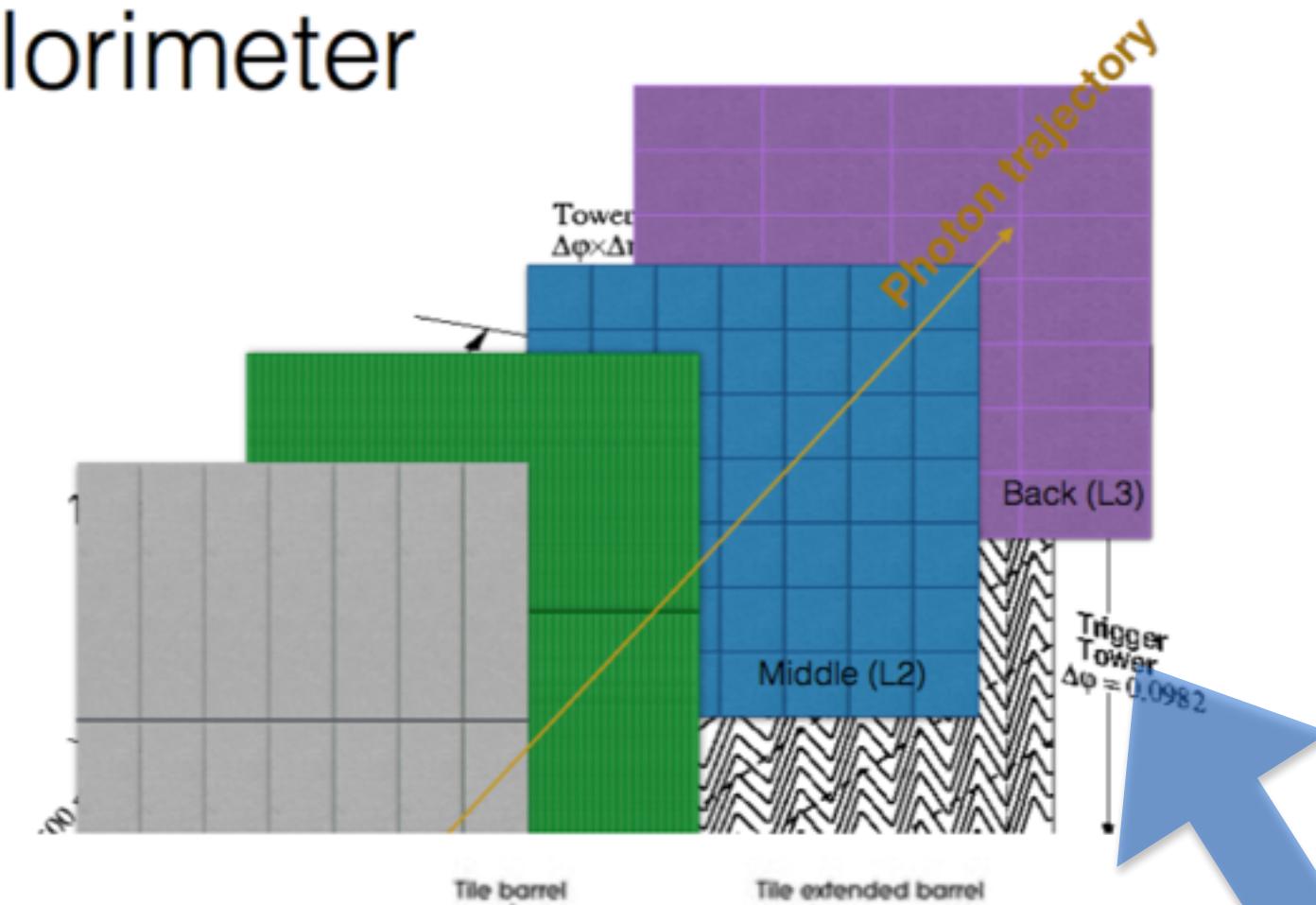
[https://indico.cern.ch/
event/766872/
contributions/3357991/](https://indico.cern.ch/event/766872/contributions/3357991/)

The Calorimeter

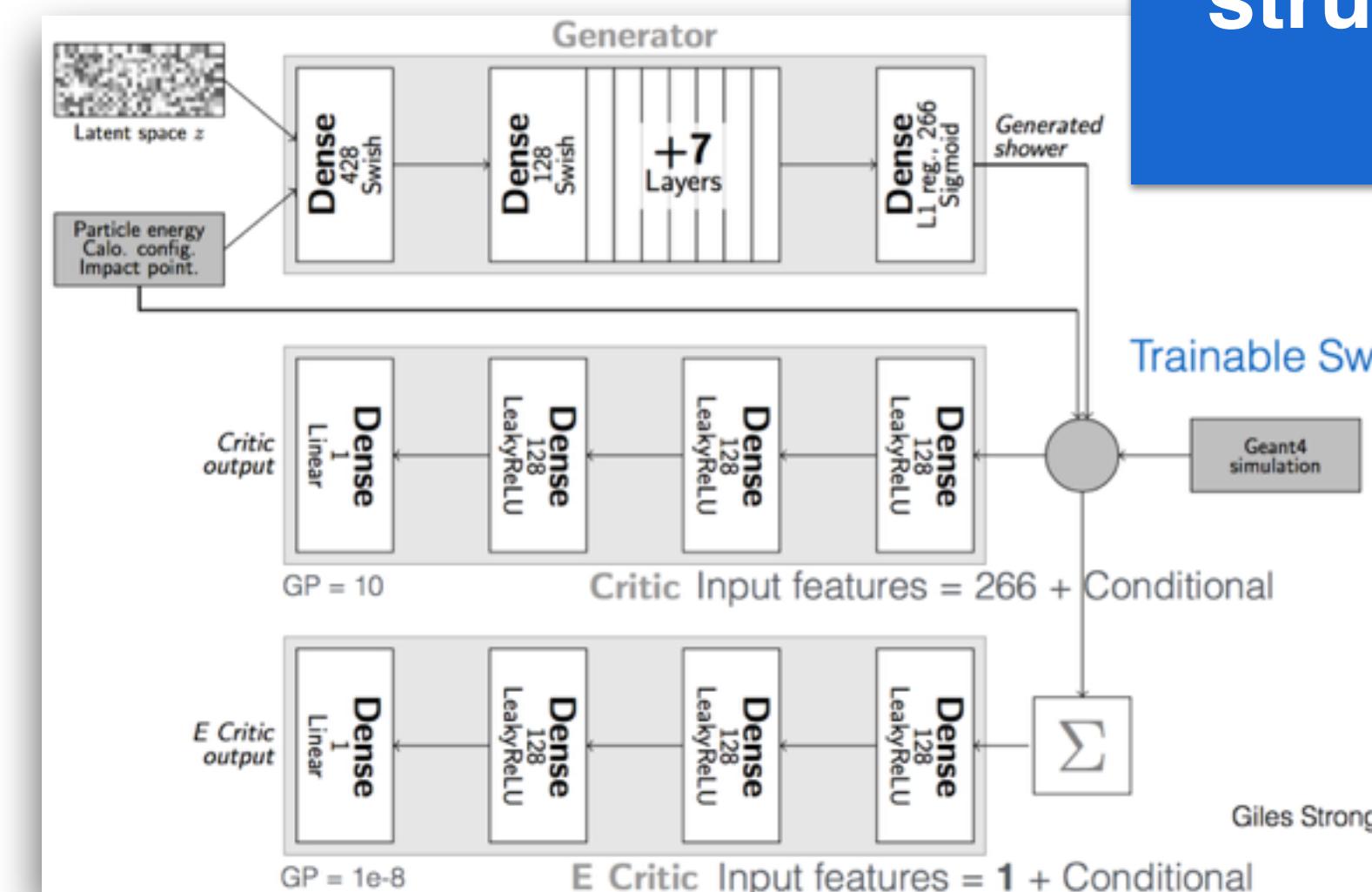
2-D Axis: ϕ vs η

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

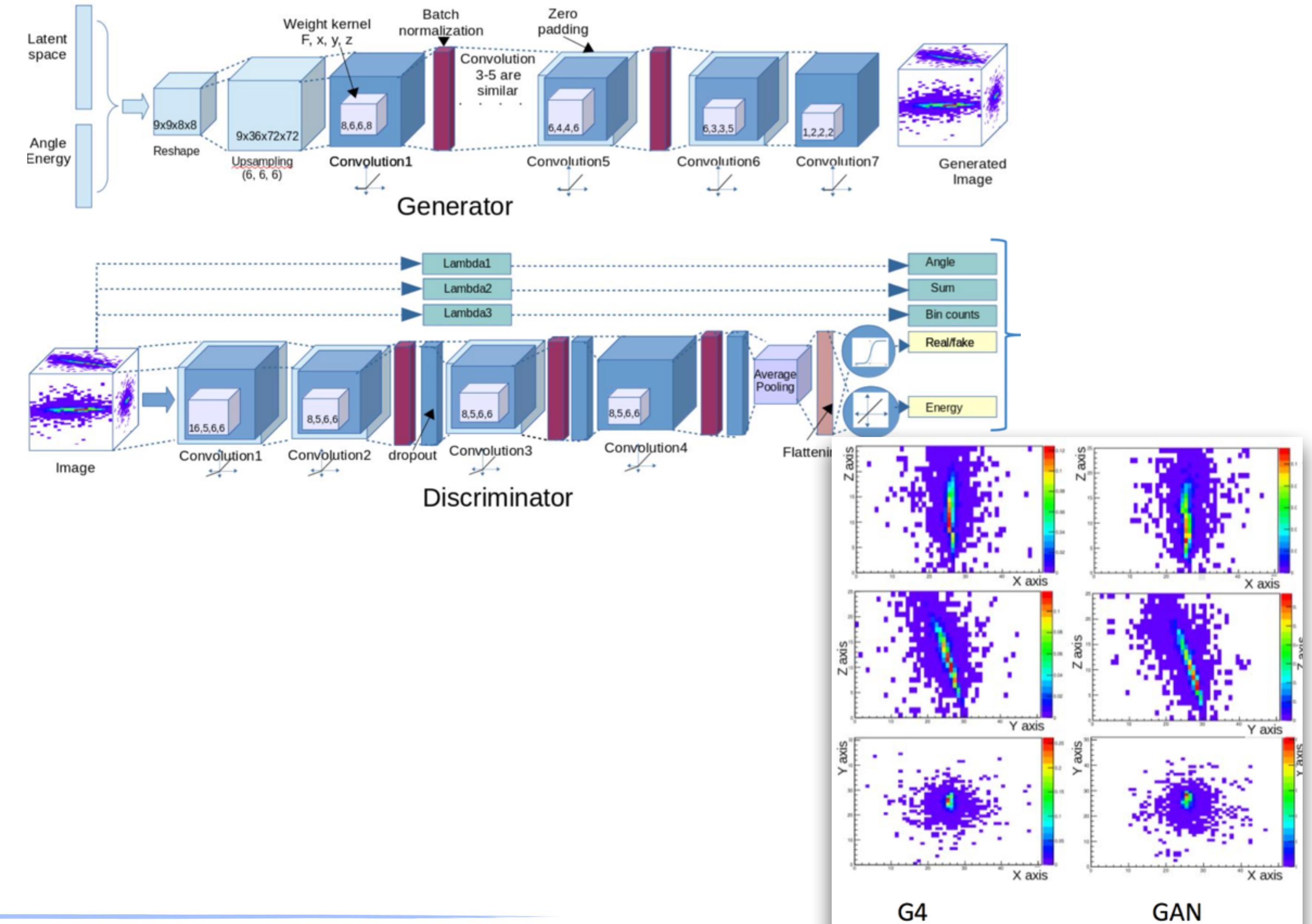
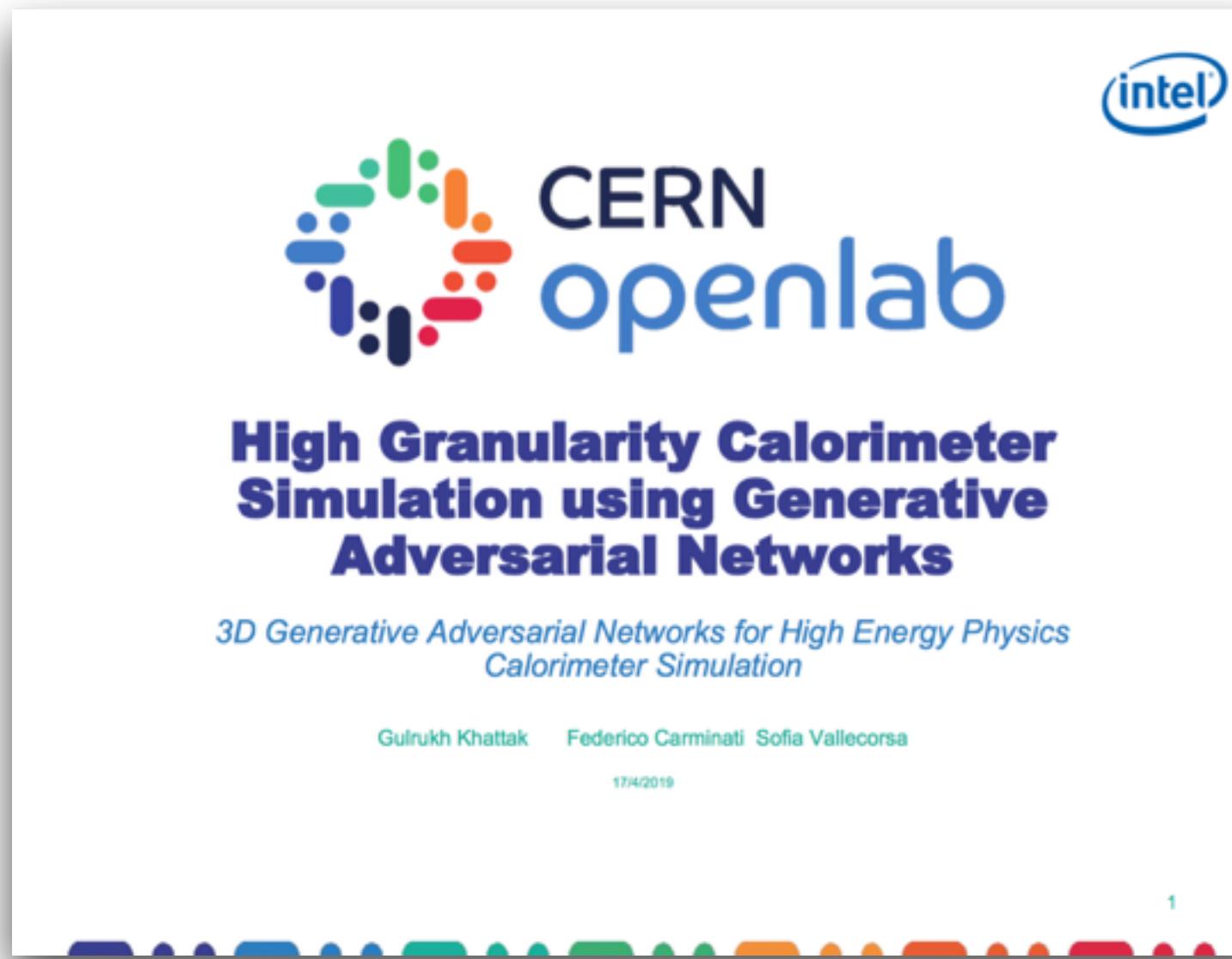


Generating output for 3D calorimeter structure (4 calorimeter layers, 266 output channels)



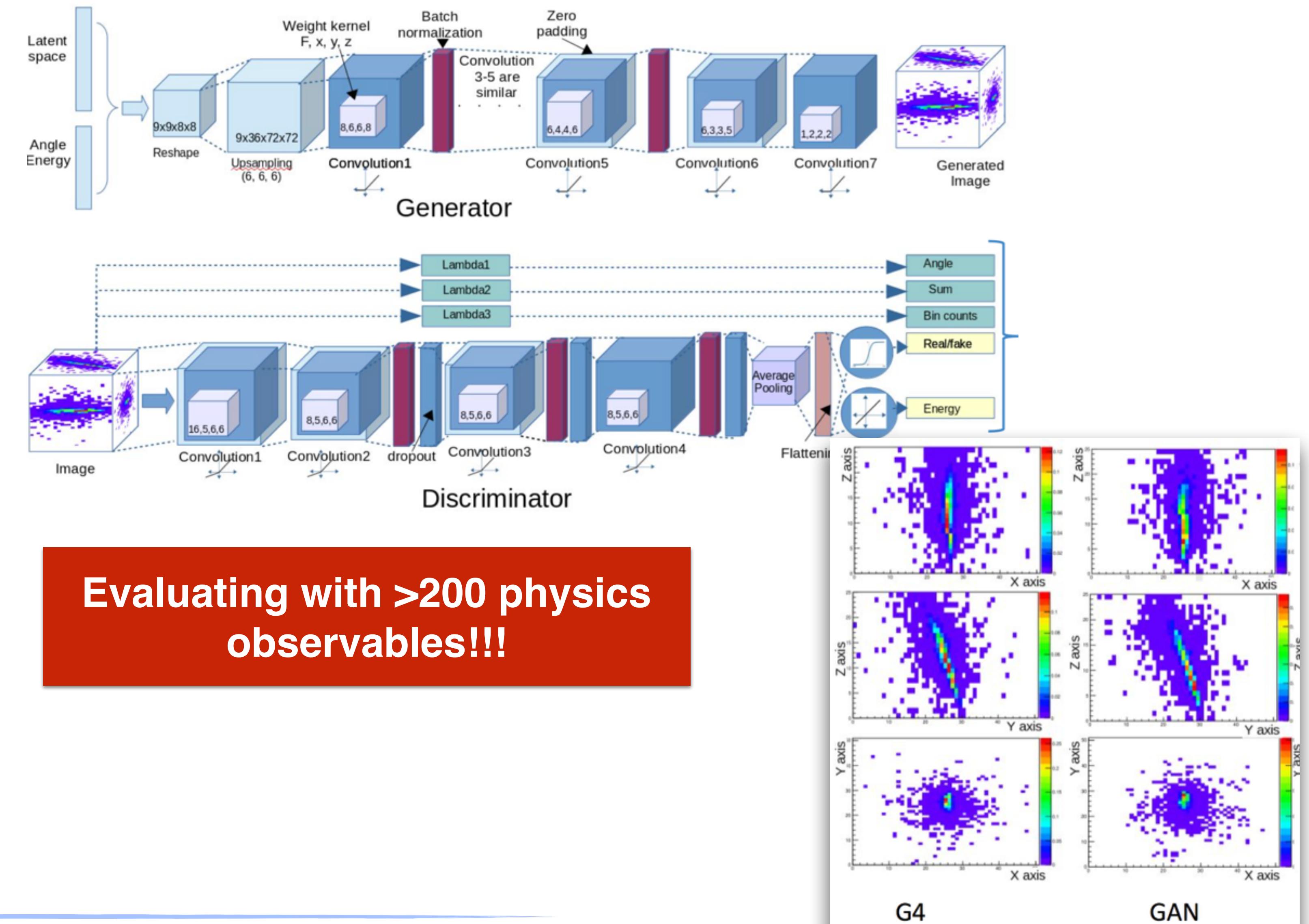
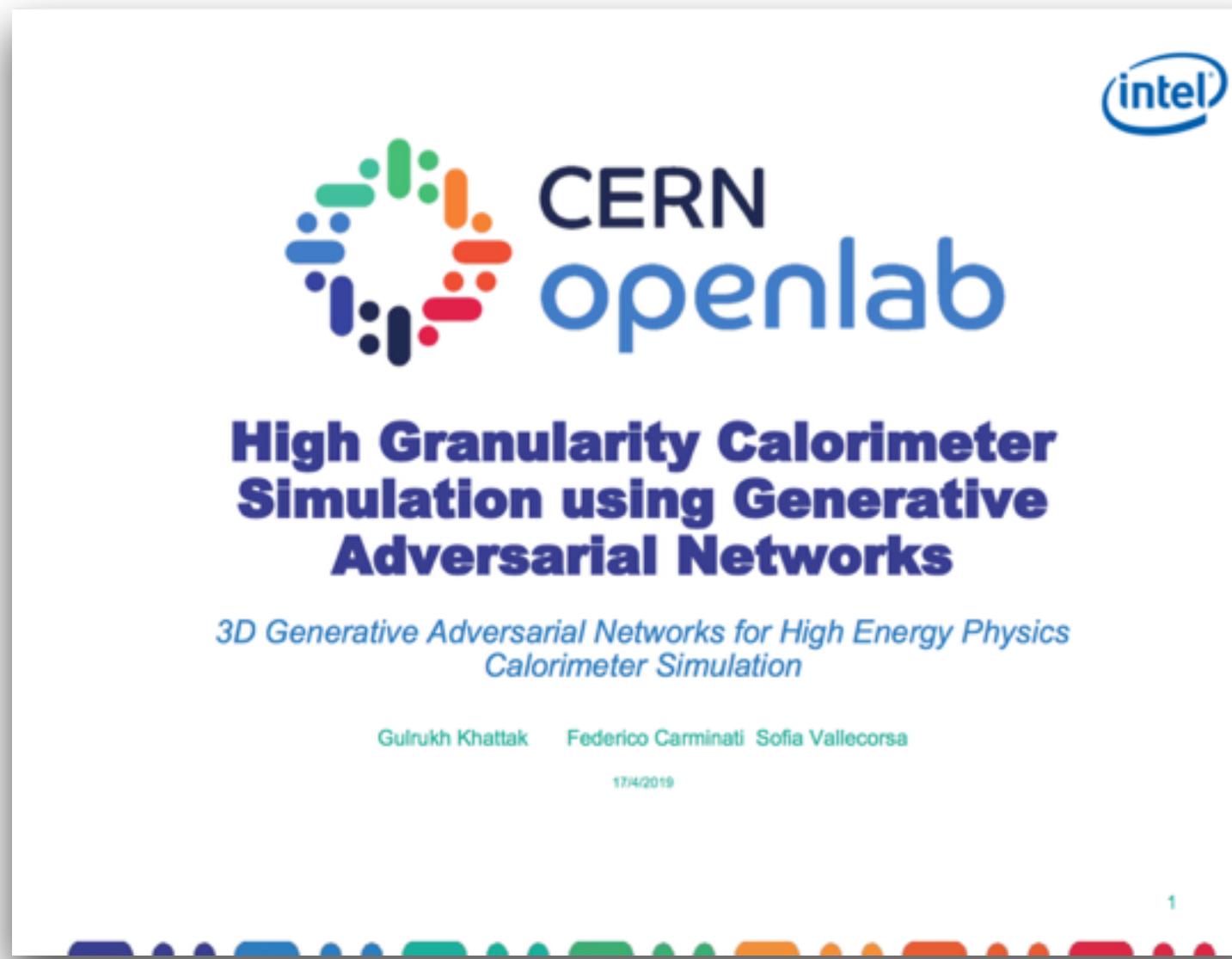
Two critics trained simultaneously (additional critic to look at the total energy distribution)

High granularity calorimeter fast simulation for CLIC



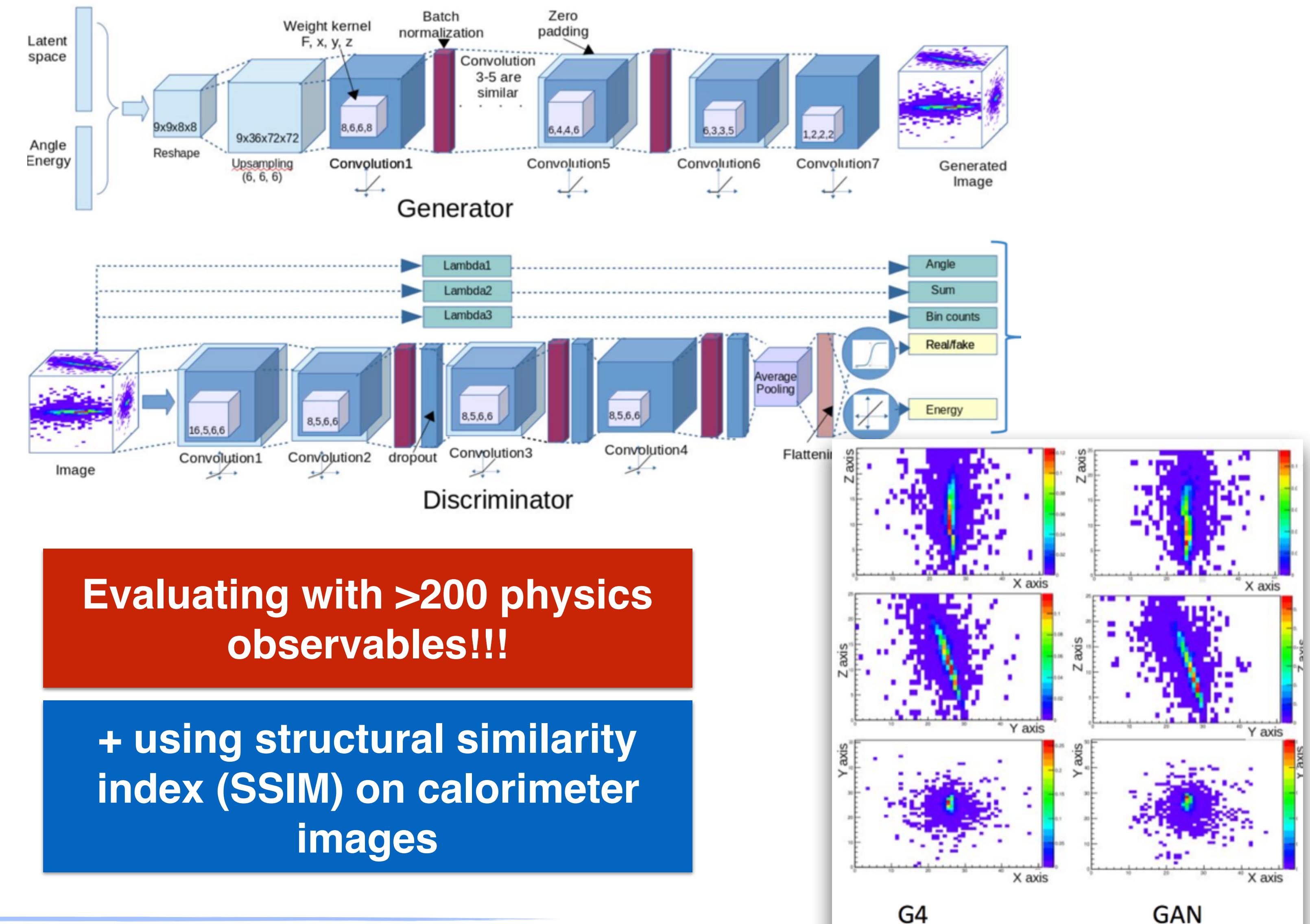
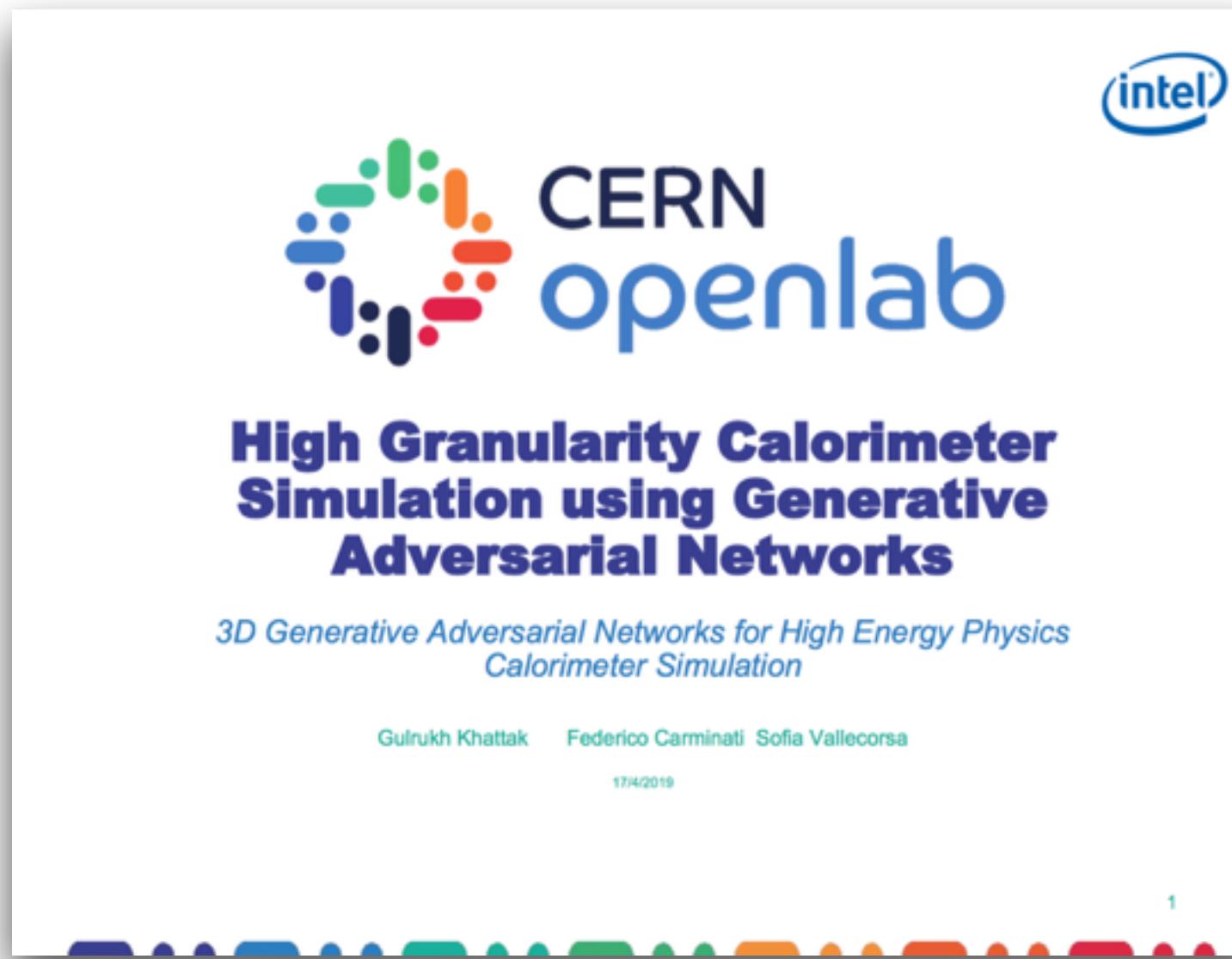
[https://indico.cern.ch/
event/766872/
contributions/3357987/](https://indico.cern.ch/event/766872/contributions/3357987/)

High granularity calorimeter fast simulation for CLIC



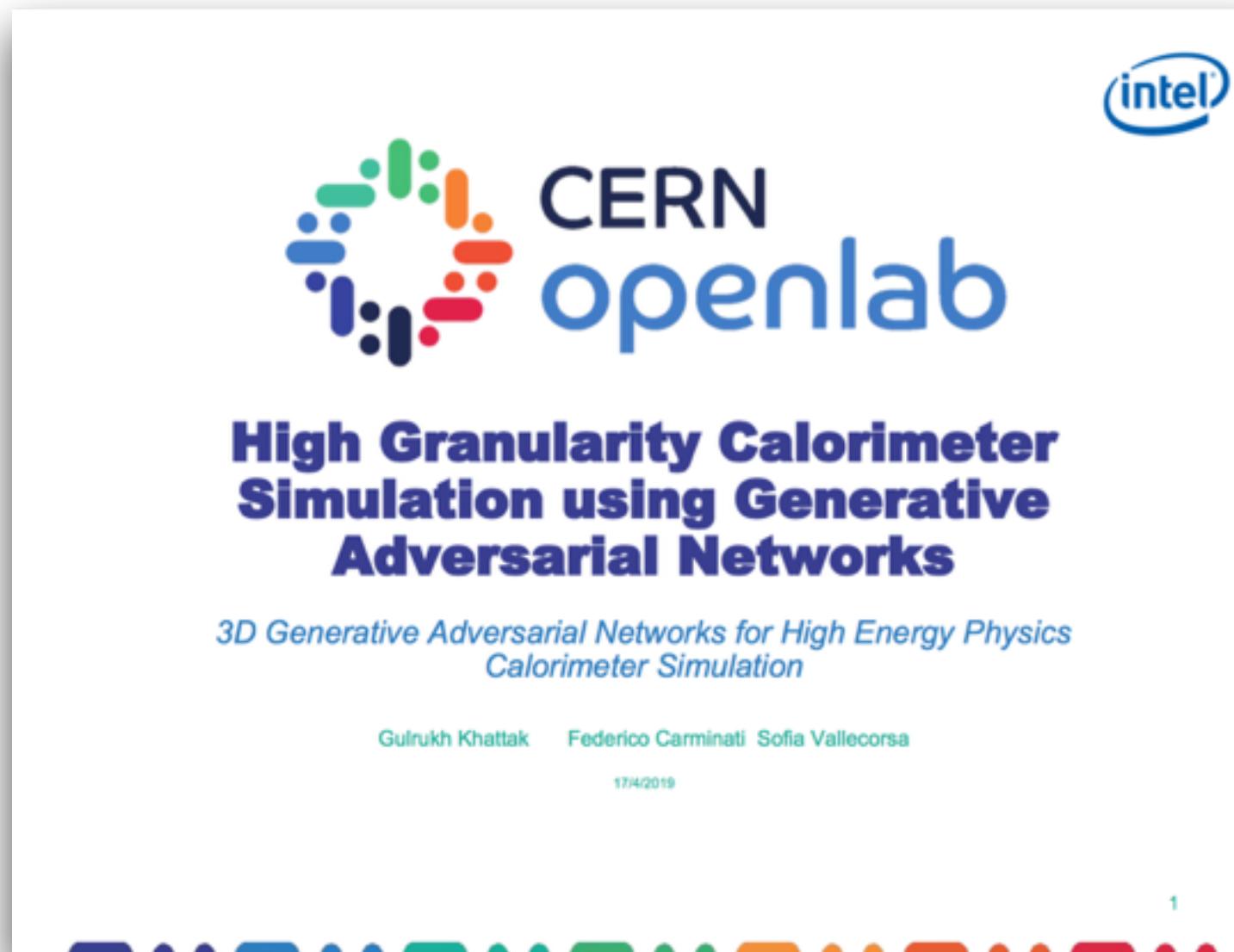
[https://indico.cern.ch/
event/766872/
contributions/3357987/](https://indico.cern.ch/event/766872/contributions/3357987/)

High granularity calorimeter fast simulation for CLIC



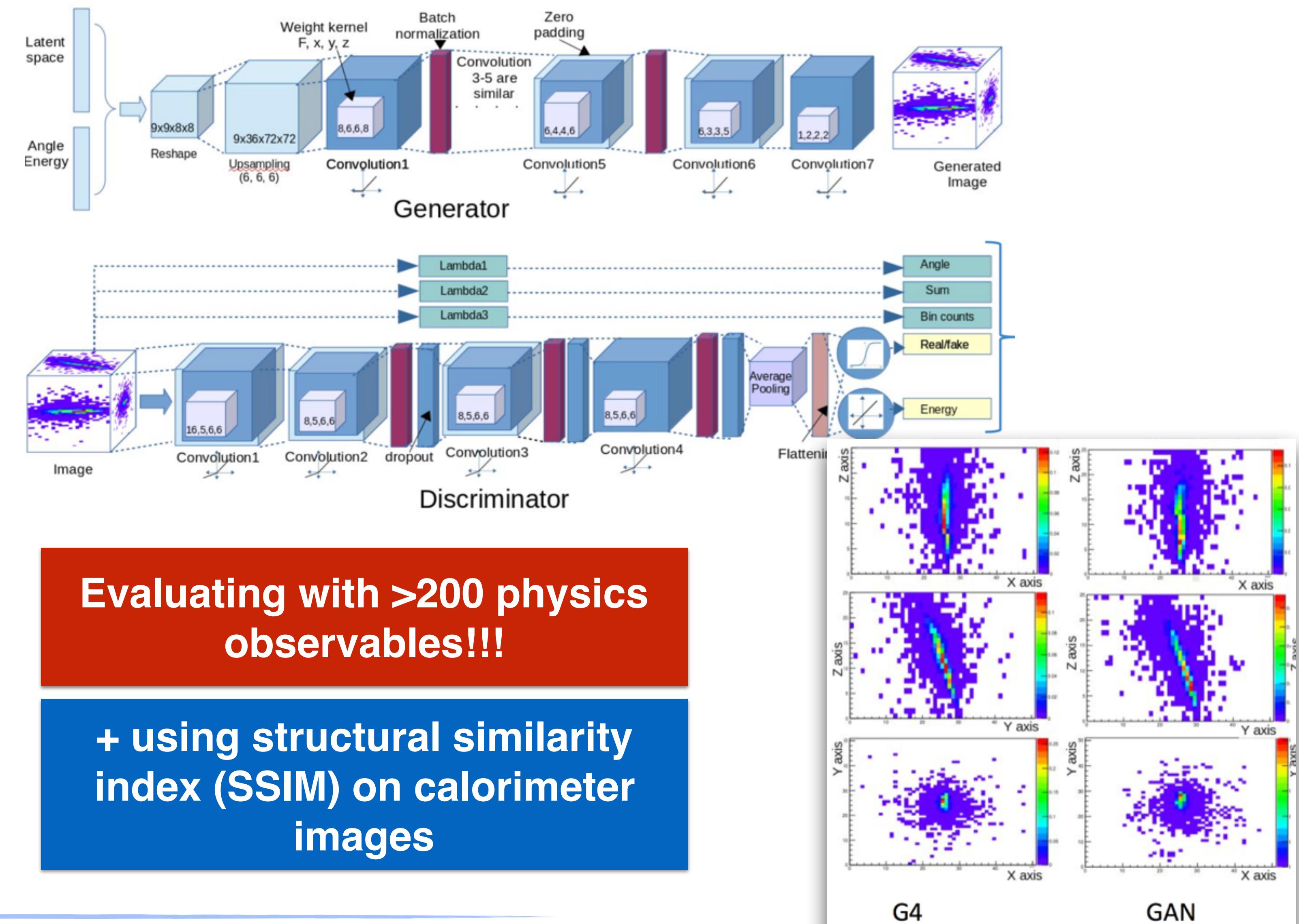
[https://indico.cern.ch/
event/766872/
contributions/3357987/](https://indico.cern.ch/event/766872/contributions/3357987/)

High granularity calorimeter fast simulation for CLIC

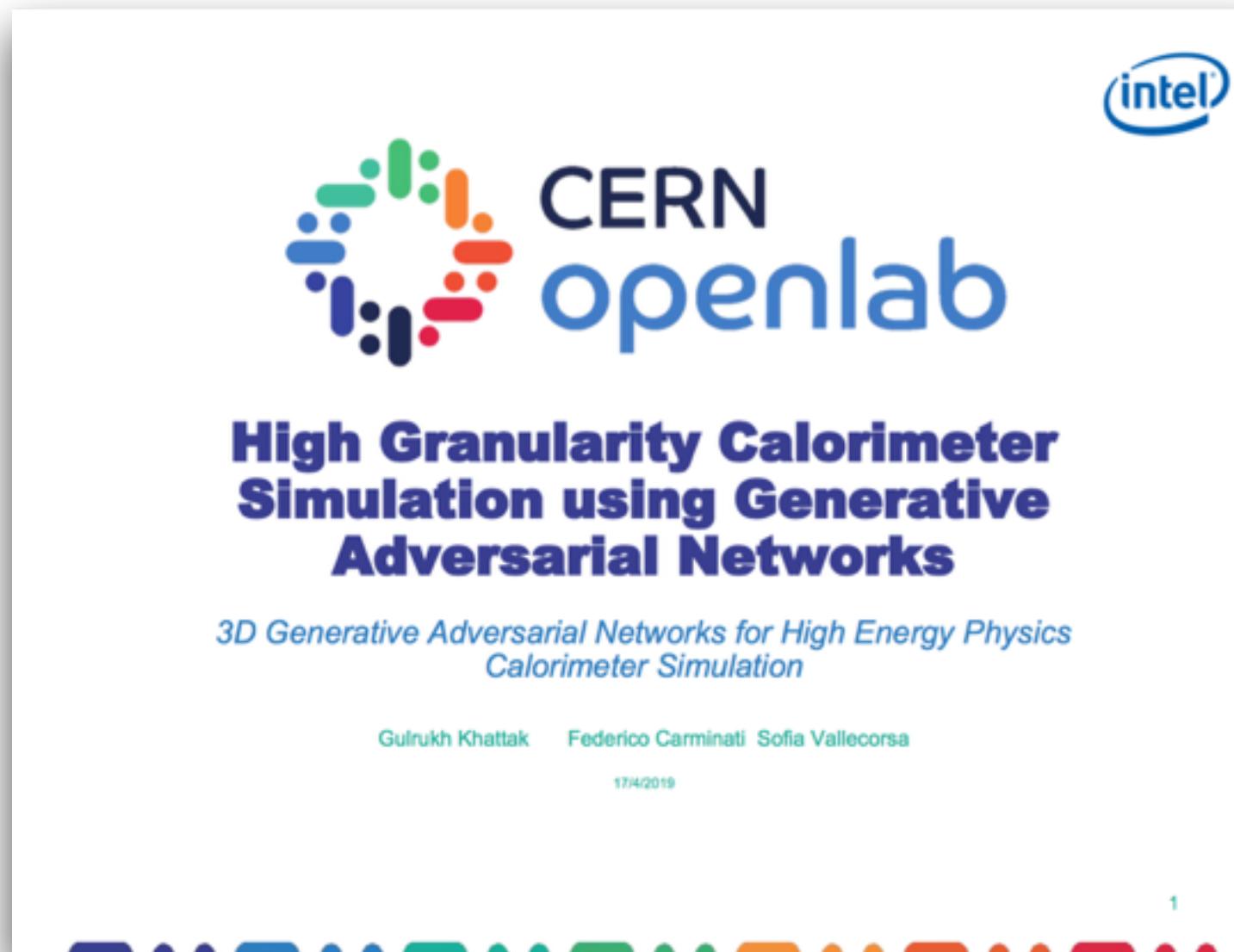


Sparse images with
~65000 pixels!!!

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

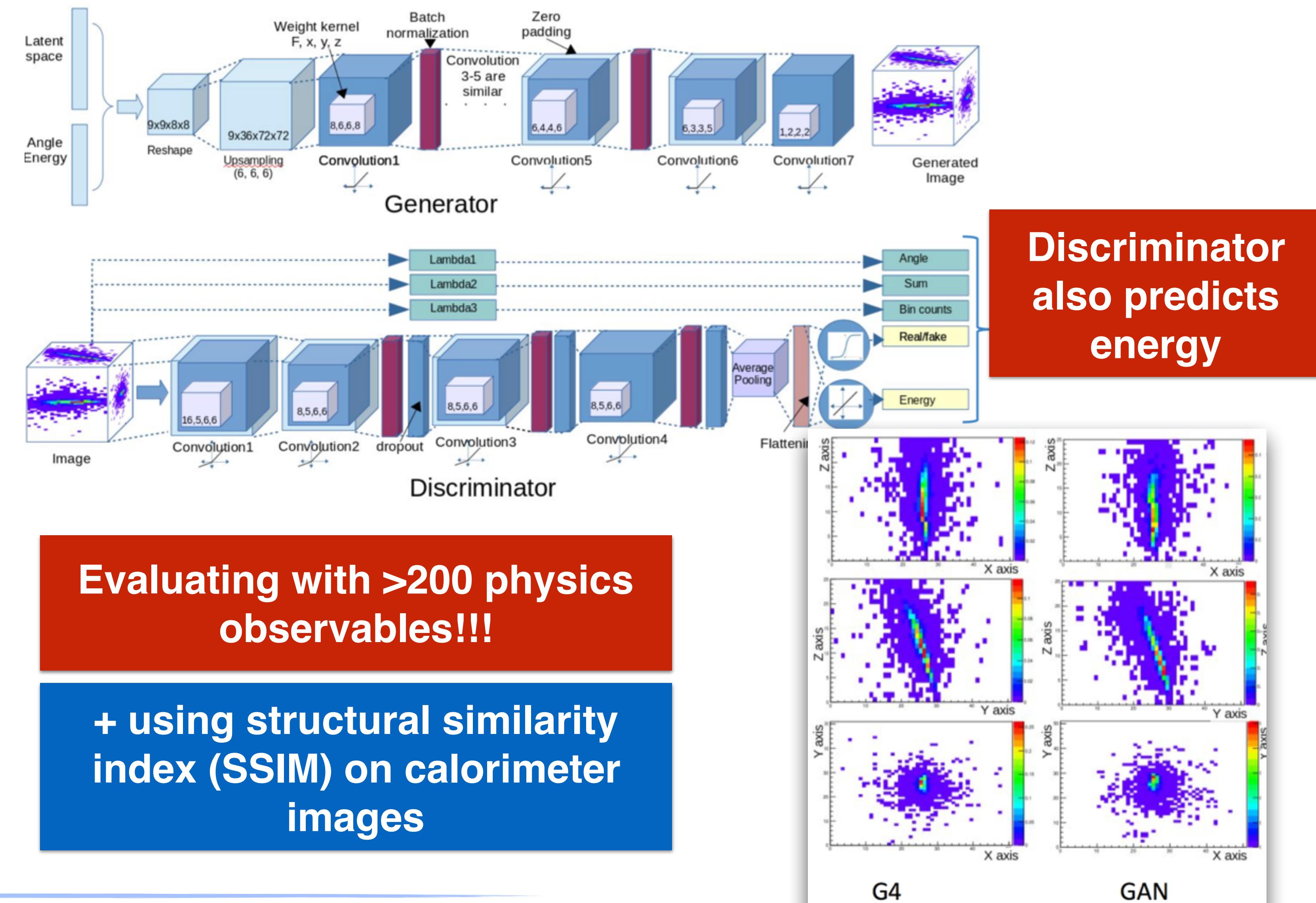


High granularity calorimeter fast simulation for CLIC

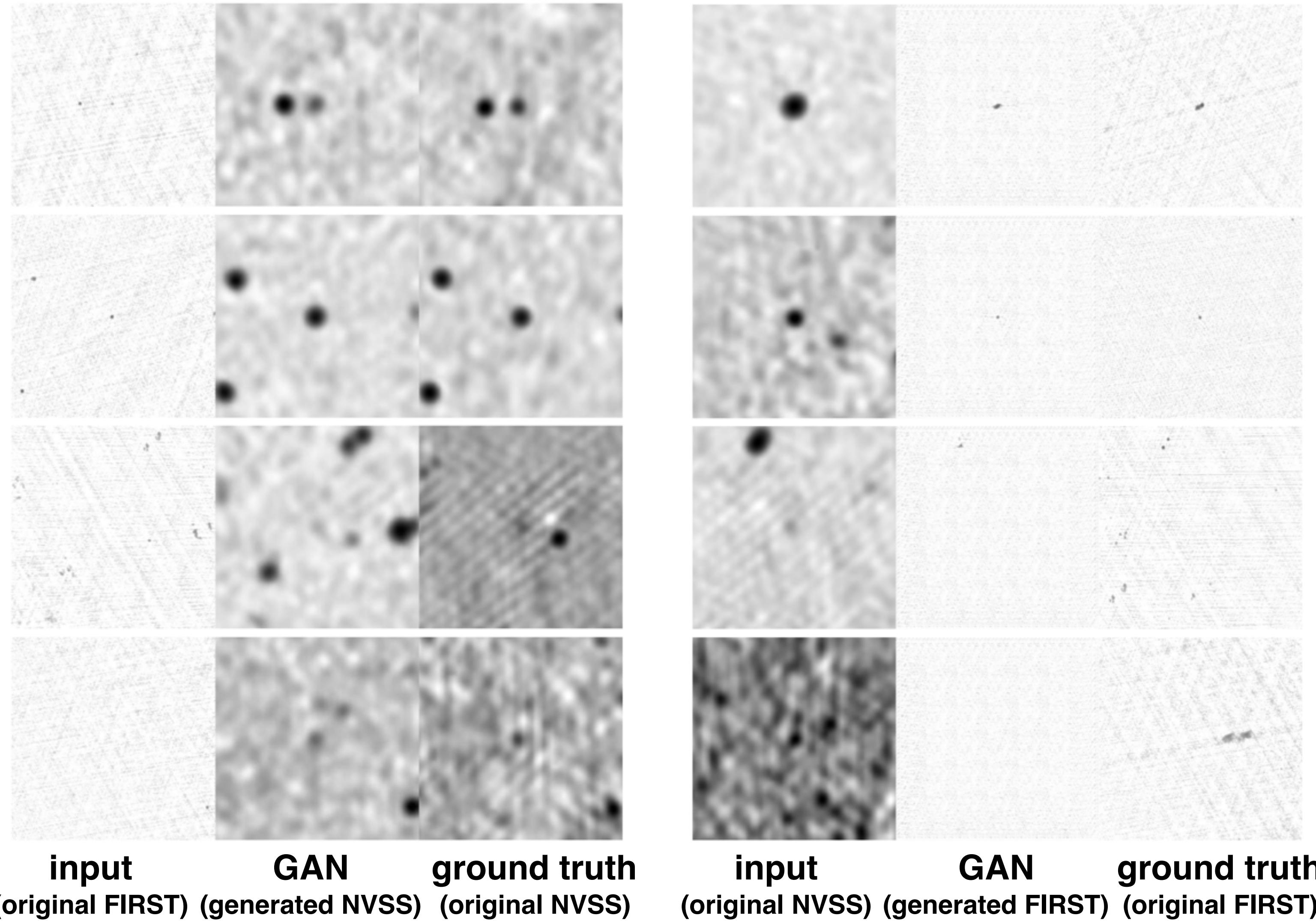


Sparse images with ~65000 pixels!!!

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



RadioGAN (Astrophysics)



GAN to compensate
for the trade-off
between brightness
sensitivity and angular
resolution

[https://arxiv.org/abs/
1906.03874](https://arxiv.org/abs/1906.03874)

Summary

- GANs – a rather broad field of ML
 - A lot of different architectures, this lecture doesn't pretend to be comprehensive
 - E.g. here's a single paper introducing 5 new GANs: <https://arxiv.org/abs/1606.00709>
- Wasserstein GAN is a useful technique that allows really deep networks to be used for data generation
- Finding a universal quality evaluation method is rather an open question