



AIR5101 / CIE6021

Generative AI

Lecture 2: Generative Adversarial Networks

Instructor: Zhen Liu

Spring 2025, CUHK-Shenzhen

Generative Adversarial Nets

Ian J. Goodfellow, Jean Pouget-Abadie*, Mehdi Mirza, Bing Xu, David Warde-Farley,
Sherjil Ozair†, Aaron Courville, Yoshua Bengio‡

Département d'informatique et de recherche opérationnelle
Université de Montréal
Montréal, QC H3C 3J7

Abstract

We propose a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models: a generative model G that captures the data distribution, and a discriminative model D that estimates the probability that a sample came from the training data rather than G . The training procedure for G is to maximize the probability of D making a mistake. This framework corresponds to a minimax two-player game. In the space of arbitrary functions G and D , a unique solution exists, with G recovering the training data distribution and D equal to $\frac{1}{2}$ everywhere. In the case where G and D are defined by multilayer perceptrons, the entire system can be trained with backpropagation. There is no need for any Markov chains or unrolled approximate inference networks during either training or generation of samples. Experiments demonstrate the potential of the framework through qualitative and quantitative evaluation of the generated samples.

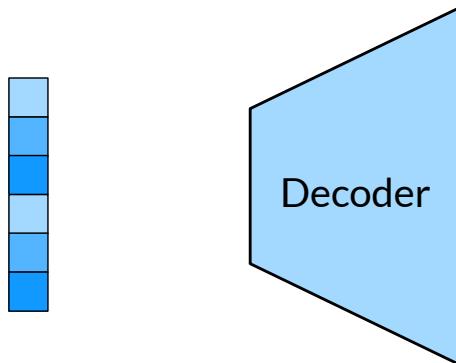


Agenda

- Vanilla GAN
- Variants (W-GAN, etc.)
- Adversary as a Loss

Vanilla GAN

Recap: Latent Variable Model



$z \sim p(z)$
Latent features
Shape, Pose, ...

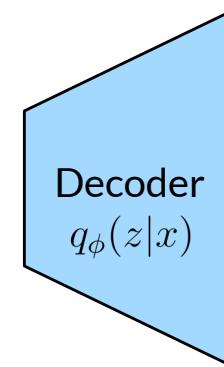
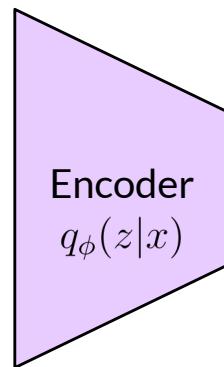


$p_\theta(x)$
Target Distribution

Recap: Variational Autoencoder (VAE)



$x \sim p_{\mathcal{D}}$

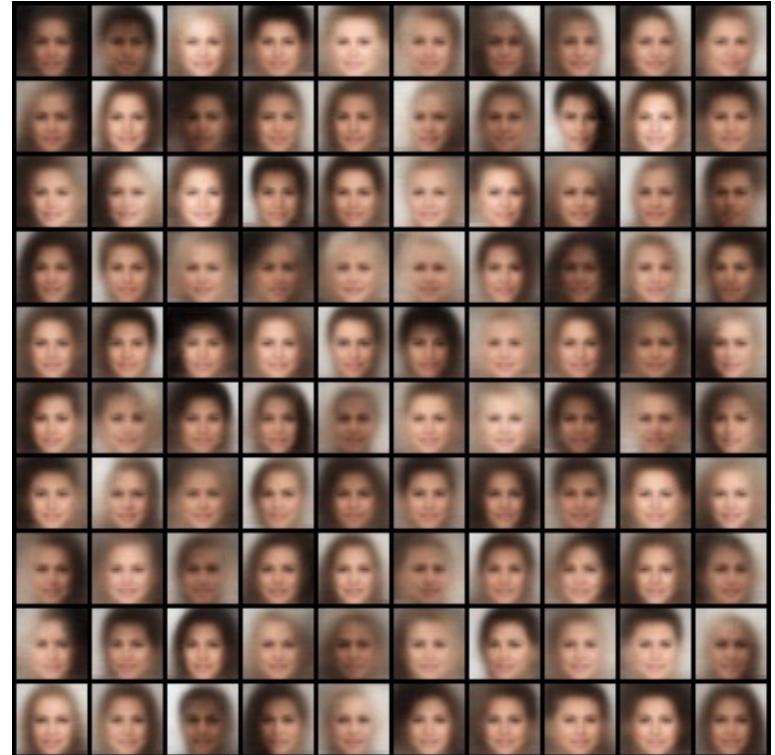


\hat{x}

Recall: Vanilla VAE on Images

Blurry reconstructed images even on simple datasets

- Canonical pose
- Canonical size
- Single-category objects
- Consistent lighting



Why VAEs fail in high-dimensional data?

Strong assumptions

- Gaussian encoder + Gaussian latent distribution

Why VAEs fail in high-dimensional data?

Strong assumptions

- Gaussian encoder + Gaussian latent distribution
- Gaussian decoder

Why VAEs fail in high-dimensional data?

Strong assumptions

- Gaussian encoder + Gaussian latent distribution
- Gaussian decoder

Any other metrics for measuring difference b/w distributions?

Why VAEs fail in high-dimensional data?

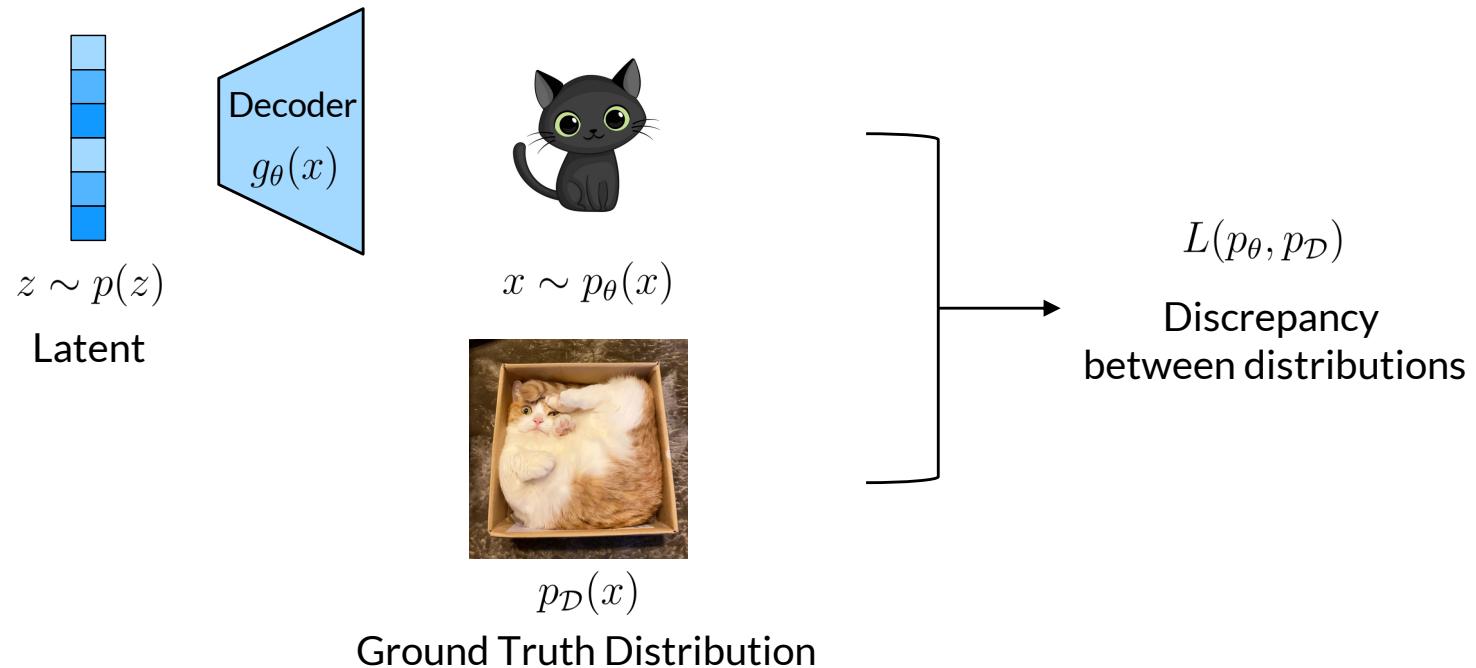
Strong assumptions

- Gaussian encoder + Gaussian latent distribution
- Gaussian decoder

Any other metrics for measuring difference b/w distributions?

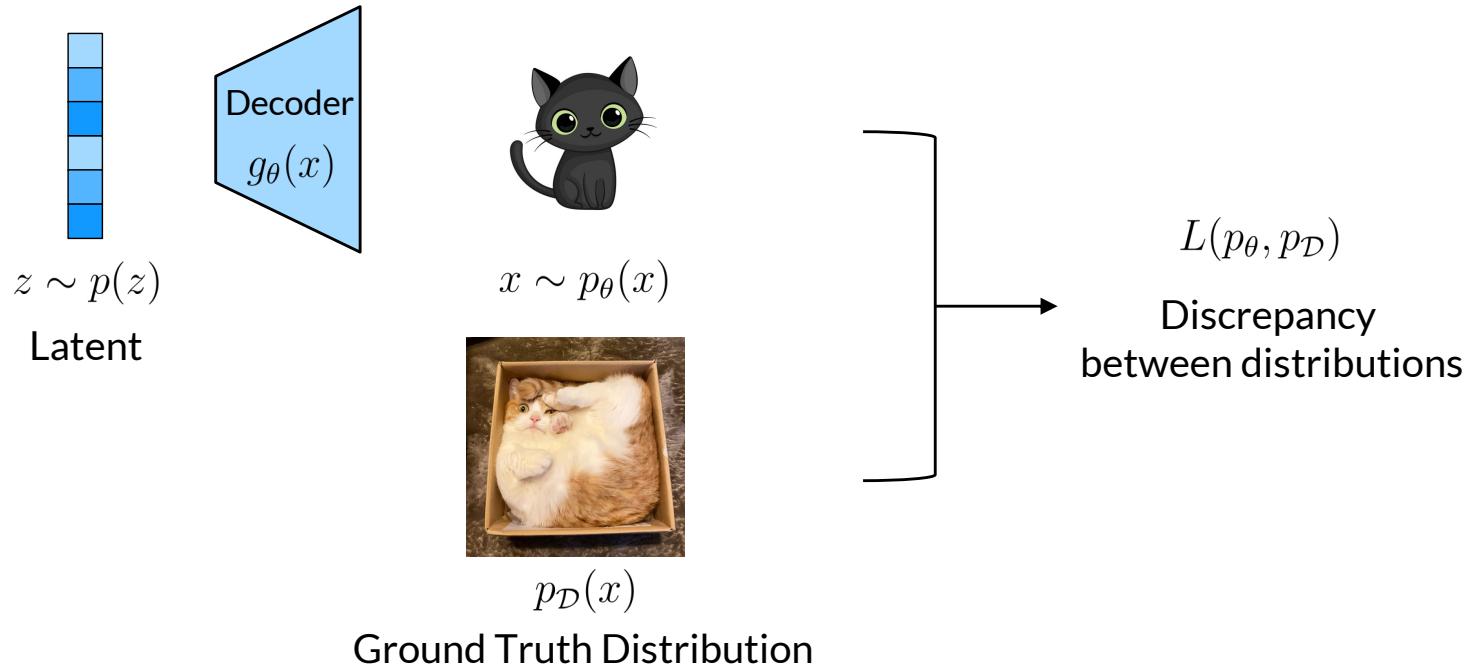
Generative Adversarial Networks

What should be a good discrepancy function?



Generative Adversarial Networks

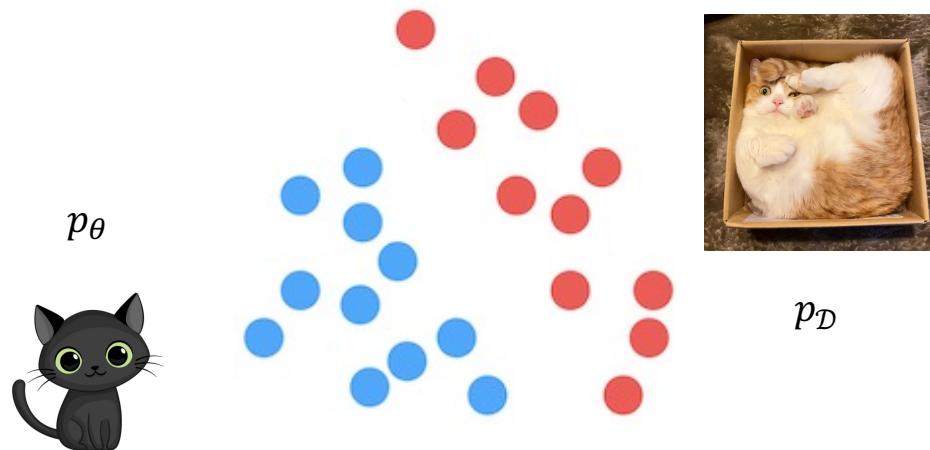
What should be a good discrepancy function?



Discriminator

Recall that $L(p_\theta, p_{\mathcal{D}}) = 0$ when $p_\theta = p_{\mathcal{D}}$.

What is the easy thing to distinguish two sets?

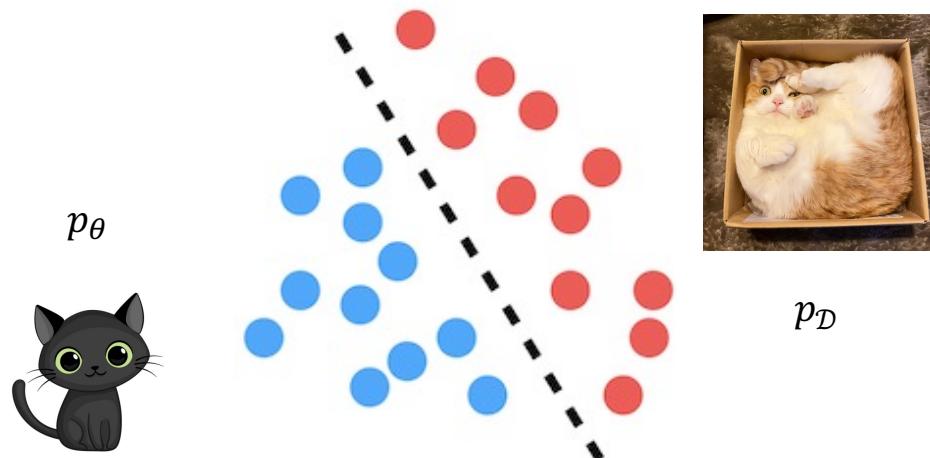


Discriminator

Recall that $L(p_\theta, p_{\mathcal{D}}) = 0$ when $p_\theta = p_{\mathcal{D}}$.

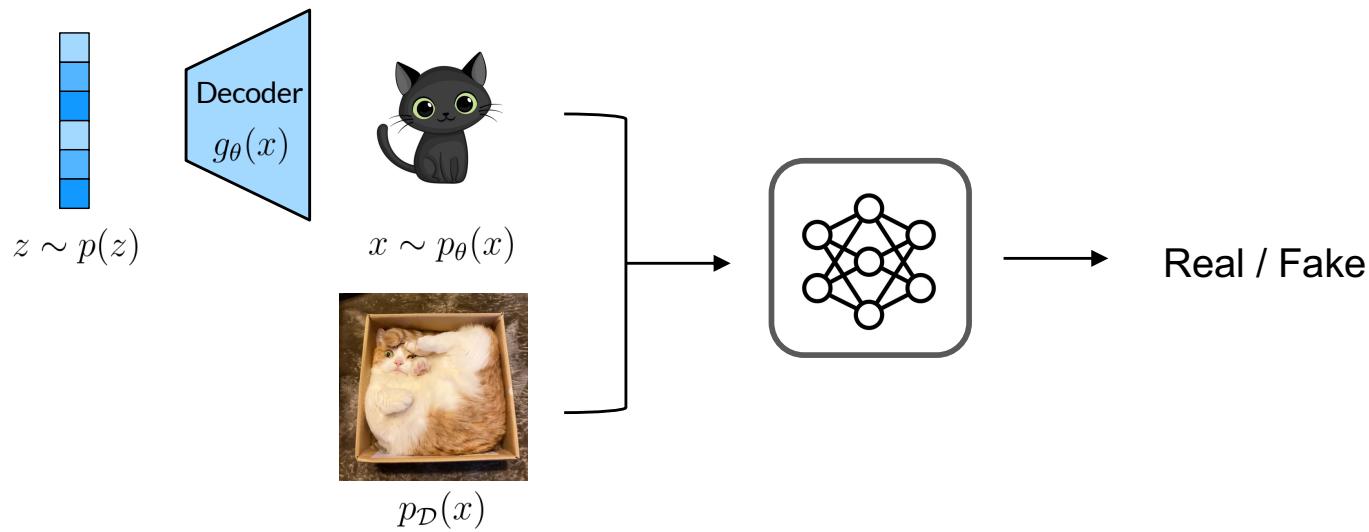
What is the easy thing to distinguish two sets?

A binary classifier!



Generative Adversarial Networks

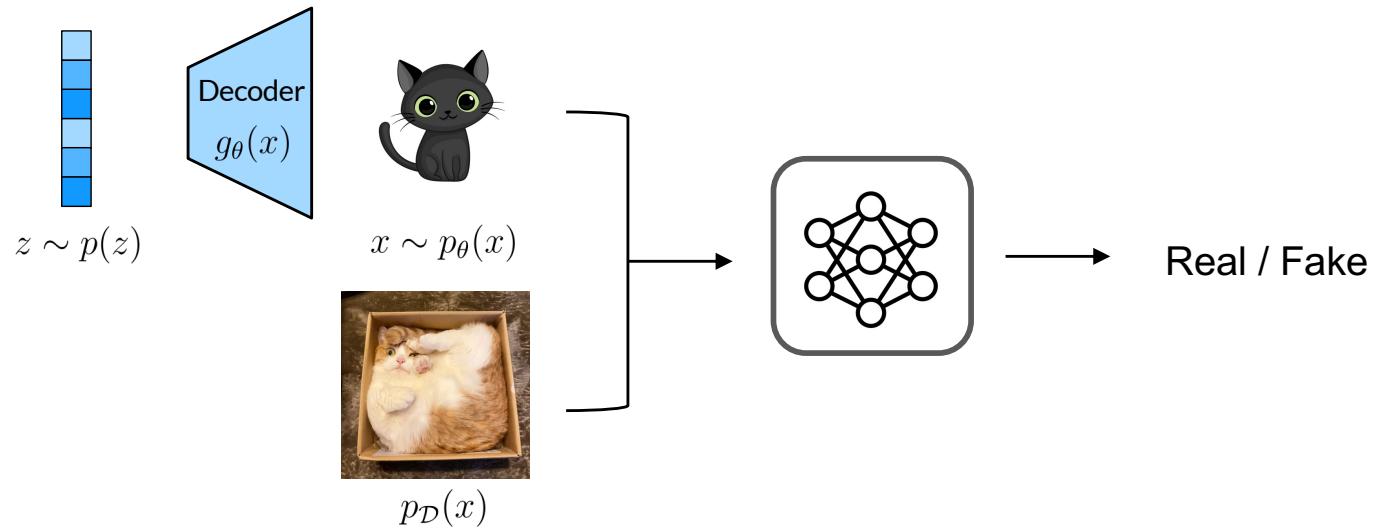
$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$



Generative Adversarial Networks

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Min-max optimization



Generative Adversarial Networks

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

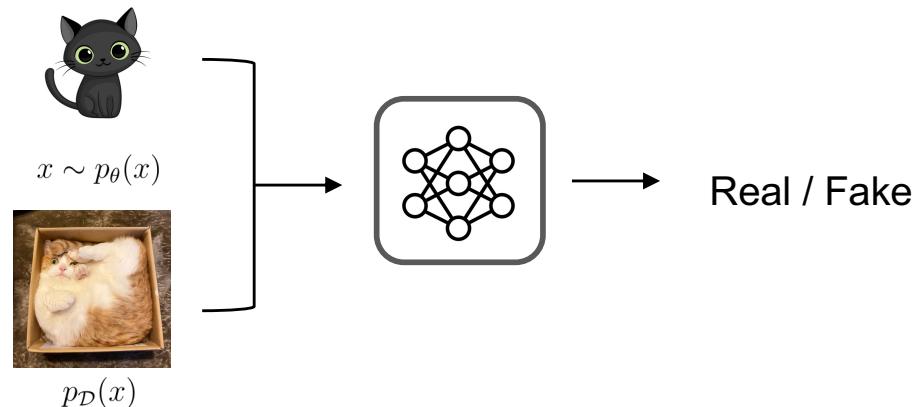
Optimizing D (G fixed):

Data: label 1

Generated: label 0

Logistic regression

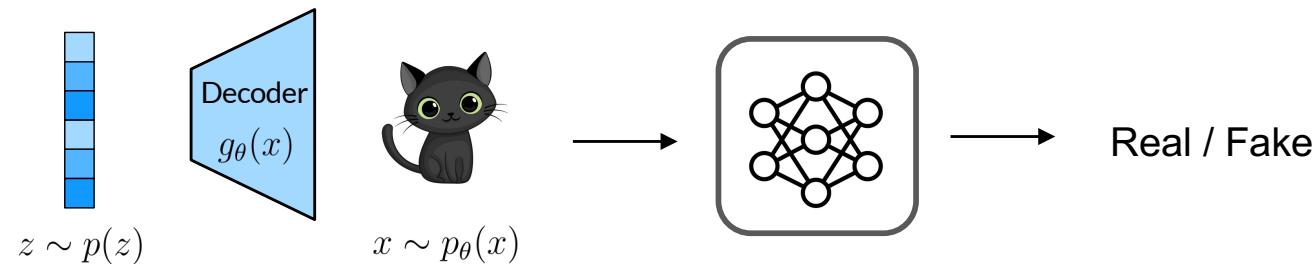
$$\mathcal{L} = -\mathbb{E}_{x \sim p_{\text{data}}} [\log f(x)] - \mathbb{E}_{x \sim p_g} [\log(1 - f(x))]$$



Generative Adversarial Networks

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Optimizing G (D fixed)
Generator as adversary

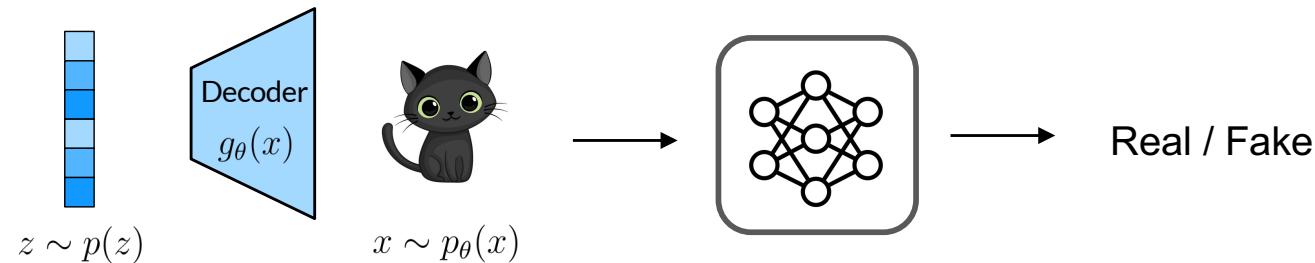


Generative Adversarial Networks

$$\min_G \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Deceive the discriminator
Push $D(G(\mathbf{z}))$ to 1

Optimizing G (D fixed)
Generator as adversary

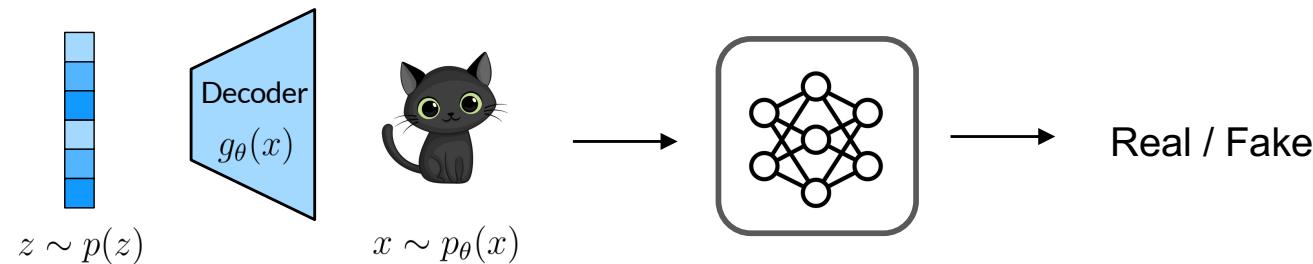


Generative Adversarial Networks

$$\min_G \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Optimizing G (D fixed)
Generator as adversary

Almost equivalent to $\max_G \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(D(G(\mathbf{z})))]$

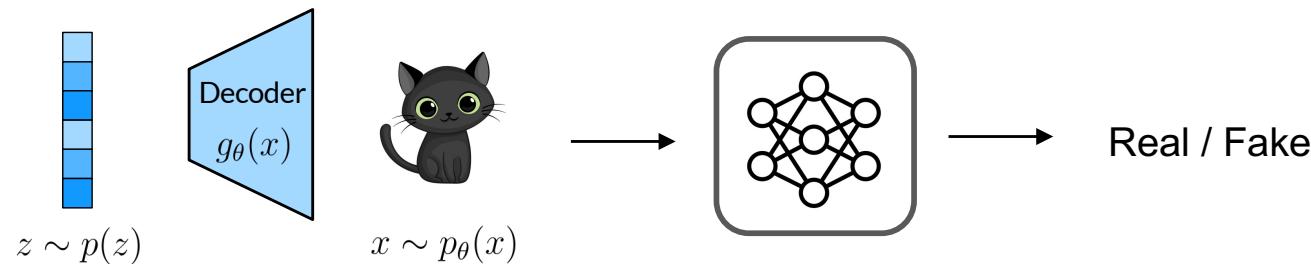


Generative Adversarial Networks

$$\min_G \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

Optimizing G (D fixed)
Generator as adversary

Almost equivalent to $\max_G \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(D(G(\mathbf{z})))]$



Algorithm

Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Sample minibatch of m examples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ from data generating distribution $p_{\text{data}}(\mathbf{x})$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

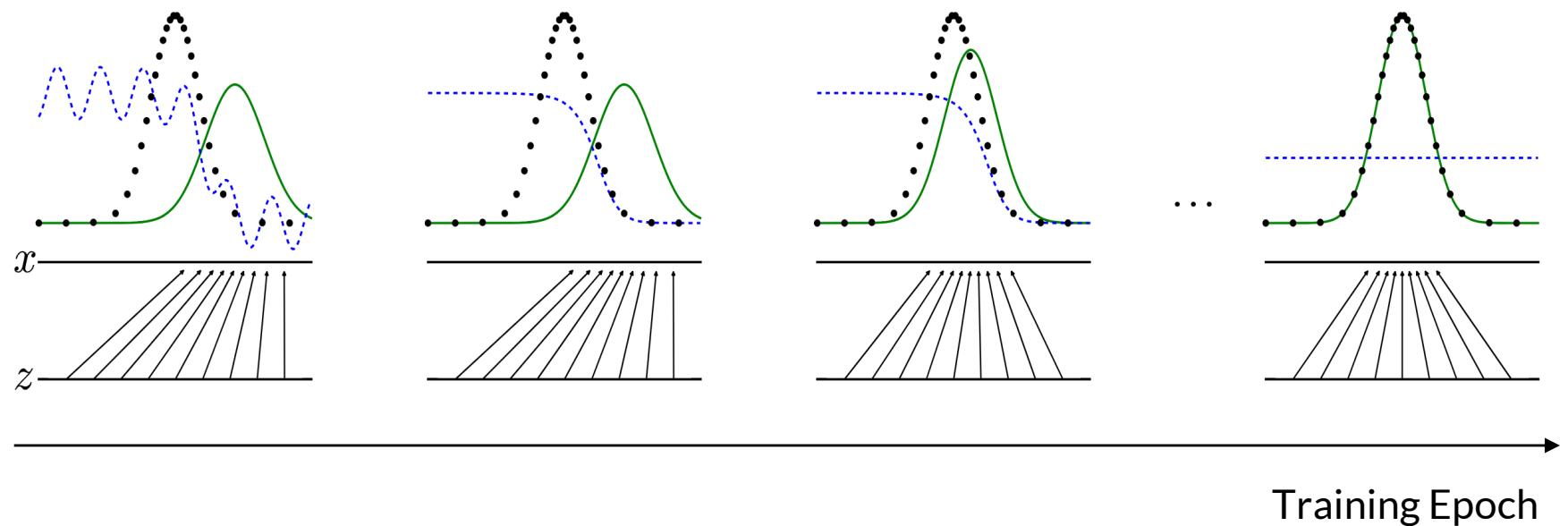
end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Critic Update

Generator
Update

Illustration



Black: data points

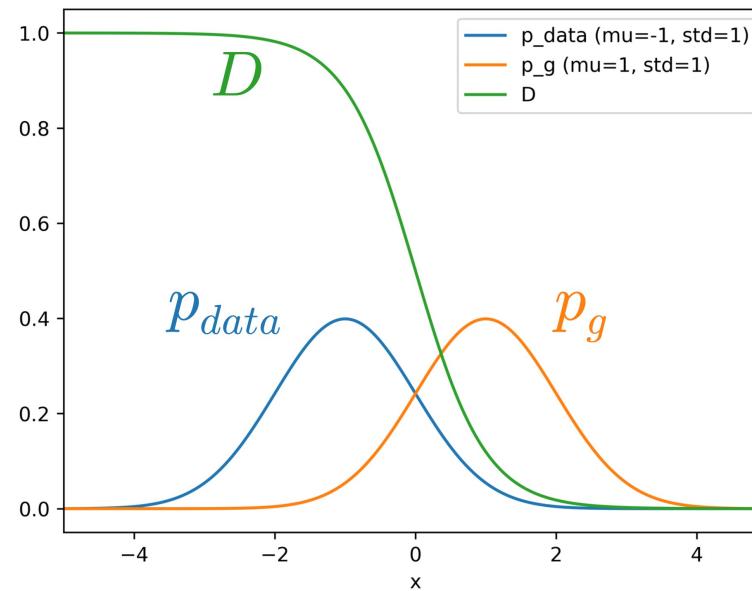
Green: modeled distribution

Blue: Discriminator

Theory

Proposition 1. *For G fixed, the optimal discriminator D is*

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}$$



Theory

Proposition 1. *For G fixed, the optimal discriminator D is*

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}$$

Indeed, $D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})}$ is a generative model itself.

Alternative interpretation: **D = generative model, G = sampler**

Will revisit briefly in the next lecture

Theory

What is the optimal discriminator given G?

$$C(G) = -\log(4) + 2 \cdot JSD(p_{\text{data}} \| p_g)$$

Jensen-Shannon Divergence

JS Divergence

$$D_{JS}(p\|q) \triangleq \frac{1}{2}D_{KL}(p\|\frac{p+q}{2}) + \frac{1}{2}D_{KL}(q\|\frac{p+q}{2})$$

Symmetric difference to the “middle point”

Theory

JS Divergence

$$D_{JS}(p\|q) \triangleq \frac{1}{2}D_{KL}(p\|\frac{p+q}{2}) + \frac{1}{2}D_{KL}(q\|\frac{p+q}{2})$$

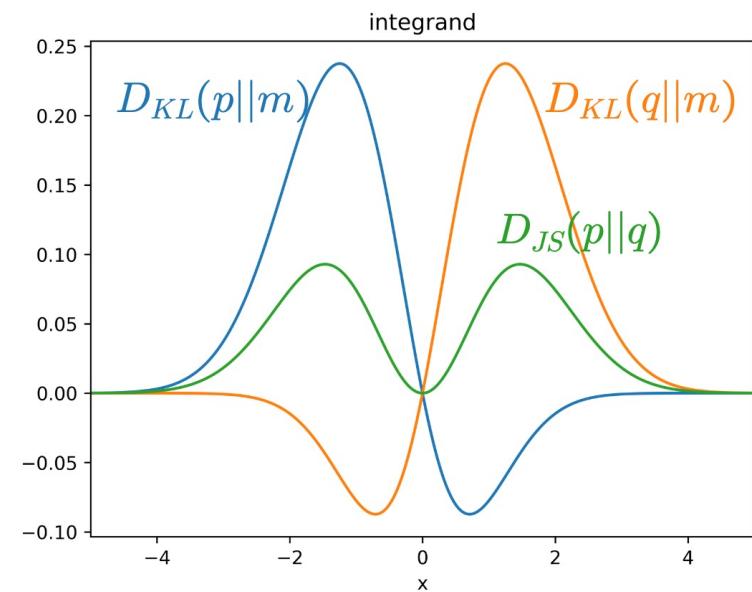
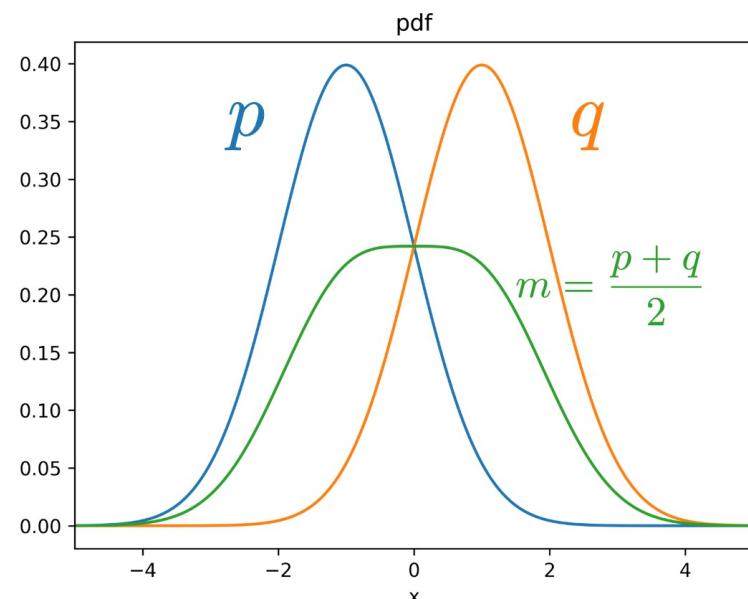
Properties:

- Symmetric (KL divergence is not)
- Bounded in $[0, \log 2]$
- Smoother gradients

Theory

JS Divergence

$$D_{JS}(p\|q) \triangleq \frac{1}{2}D_{KL}(p\|\frac{p+q}{2}) + \frac{1}{2}D_{KL}(q\|\frac{p+q}{2})$$



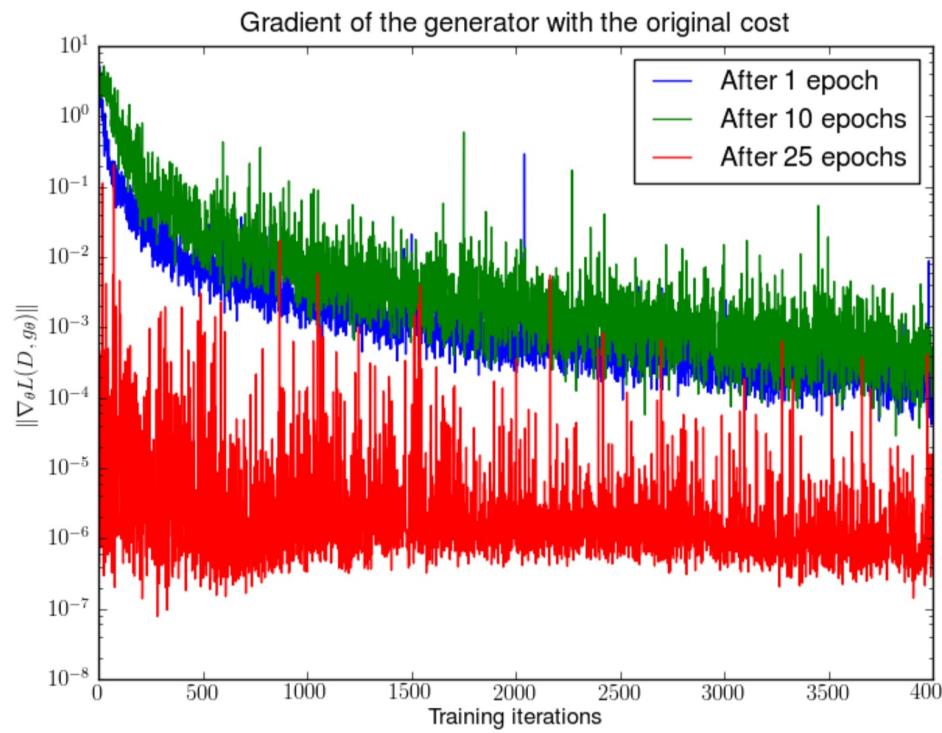
https://mit-6s978.github.io/assets/pdfs/lec4_gan.pdf

DCGAN

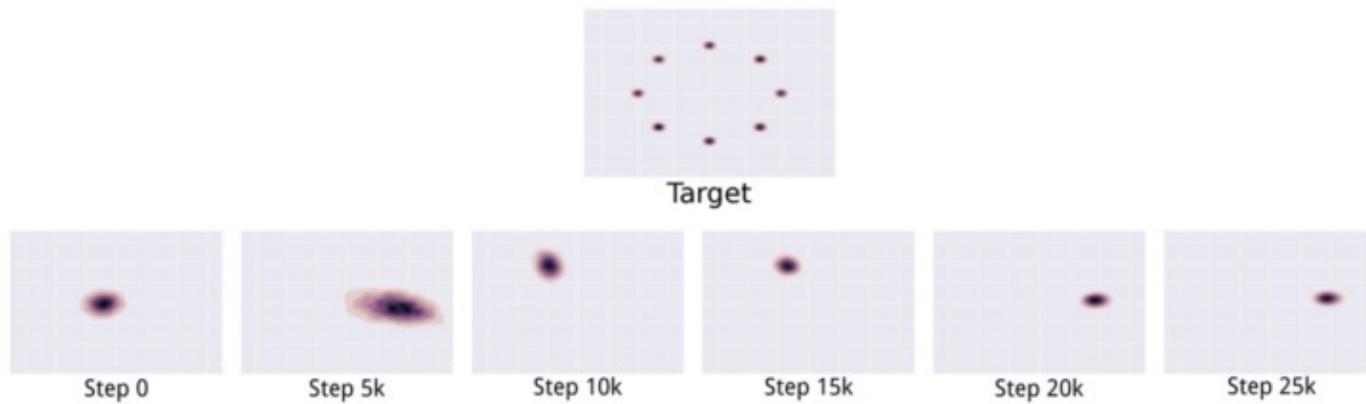
Results from the first GAN model with ConvNets



Vanishing Gradients



Mode Collapse



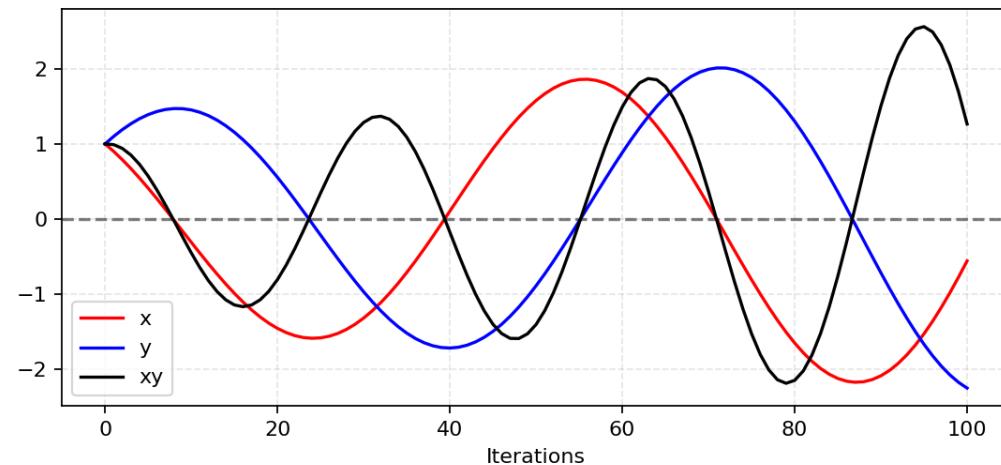
Many modes not captured

Training Instability

Min-max optimization = differential game

Example: finding a Nash equilibrium in a 1D differentiable game

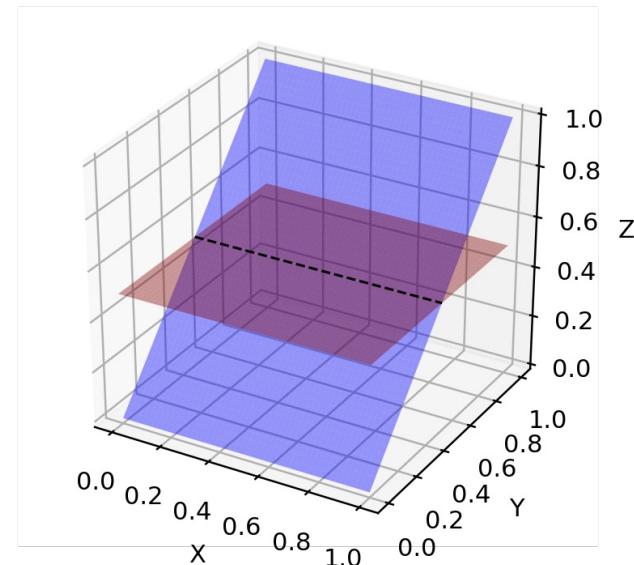
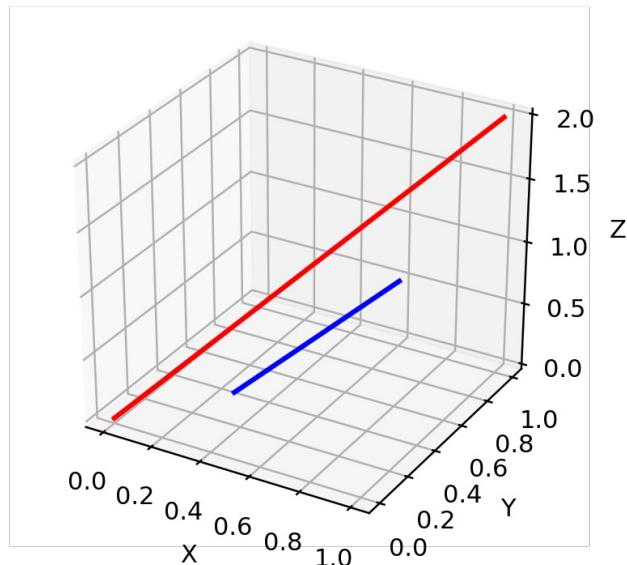
$$\min_x \max_y xy$$



Trivial Solutions

Low dimensional support of both data and modeled distributions

Need careful tuning to avoid trivial discriminator



Visualize GAN Training in Your Browser

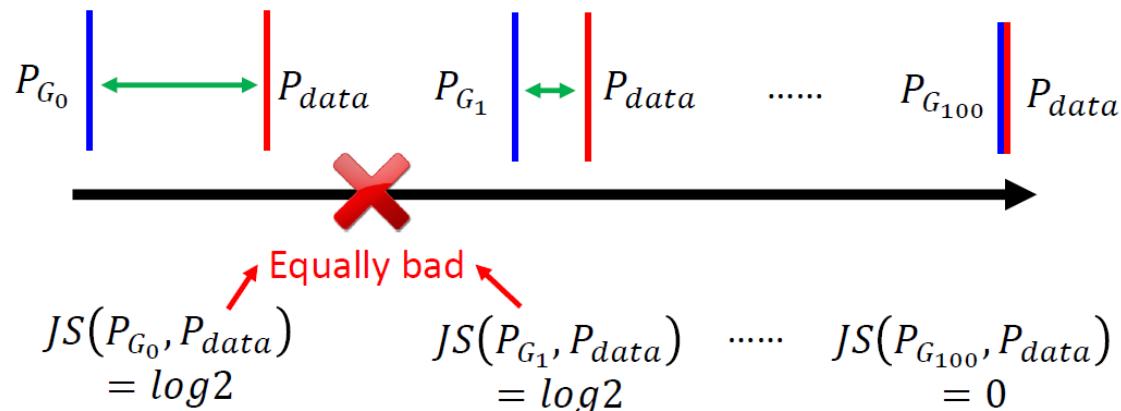
<https://poloclub.github.io/ganlab/>



GAN Variants

Issue with JS Divergence

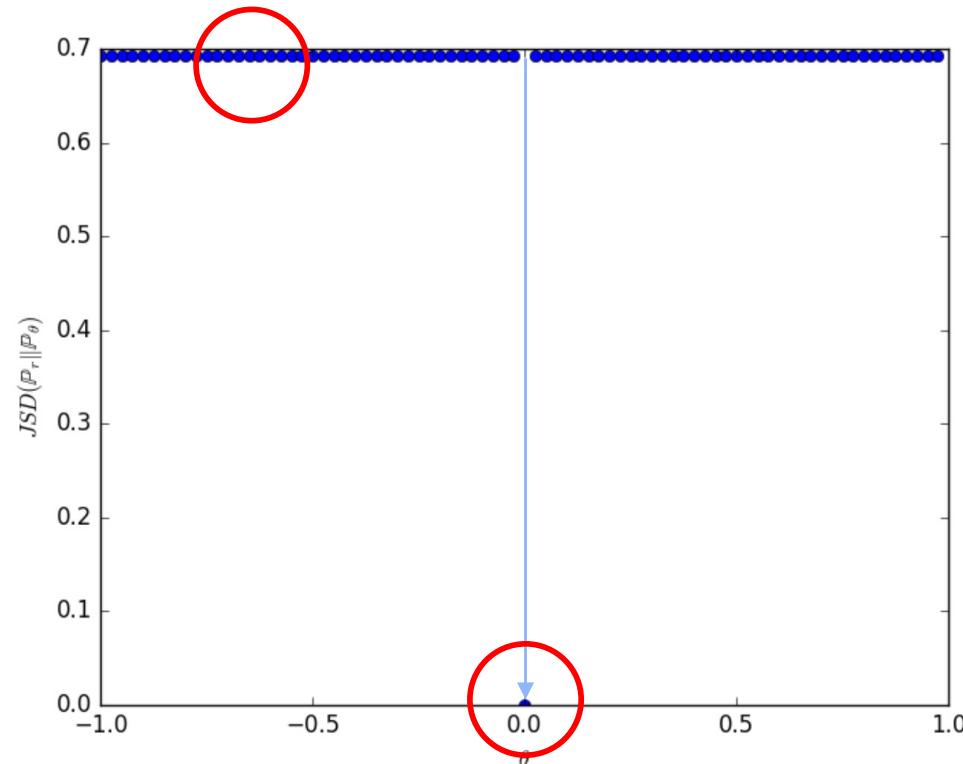
Suppose p_θ and $p_{\mathcal{D}}$ are parallel lines in 2D
We want to move p_θ parallelly to $p_{\mathcal{D}}$



JS divergence is $\log 2$ if two distributions do not overlap.

Issue with JS Divergence

Zero gradient elsewhere



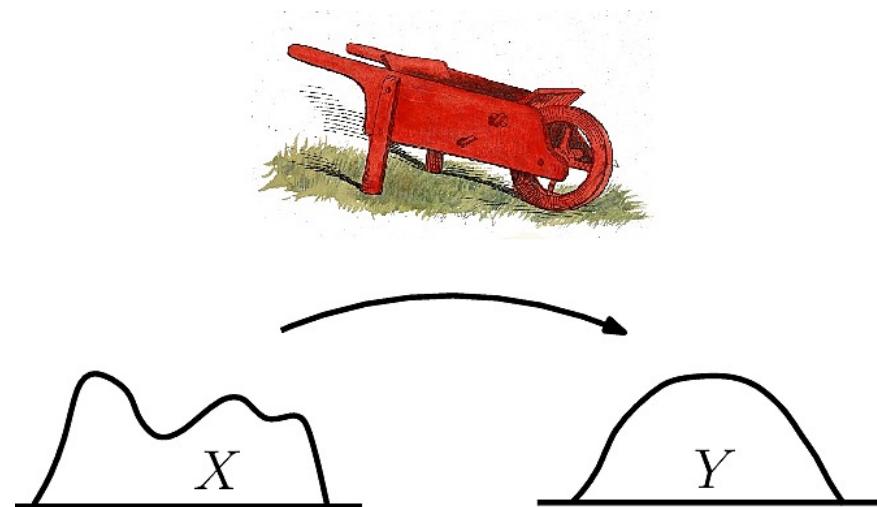
Any other distance function?

Infinite gradient when overlapped

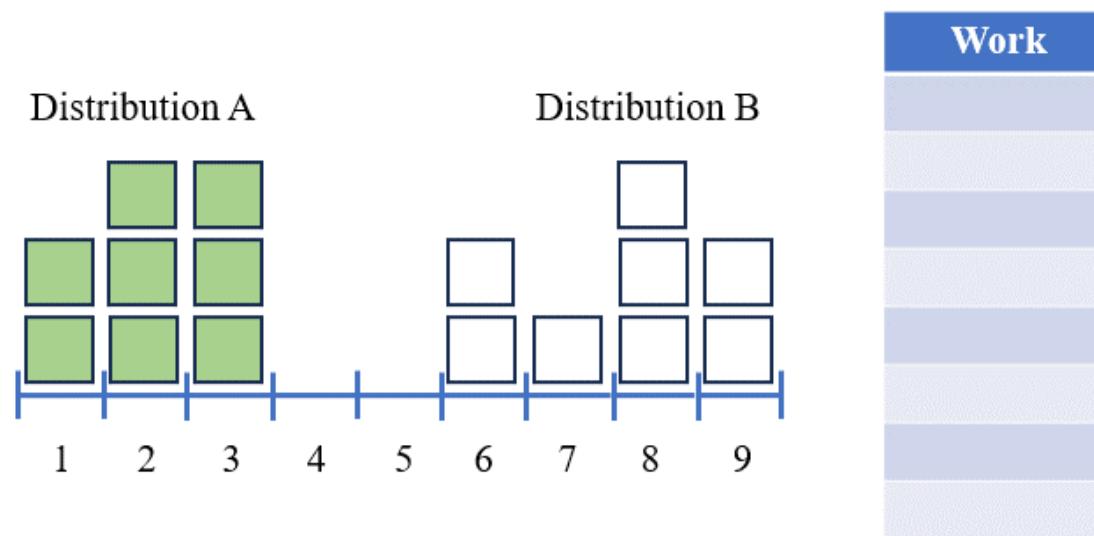
Wasserstein Distance

a.k.a. Earth Mover's Distance

How can you move one pile of earth to the other shape (at different location) with the least cost?

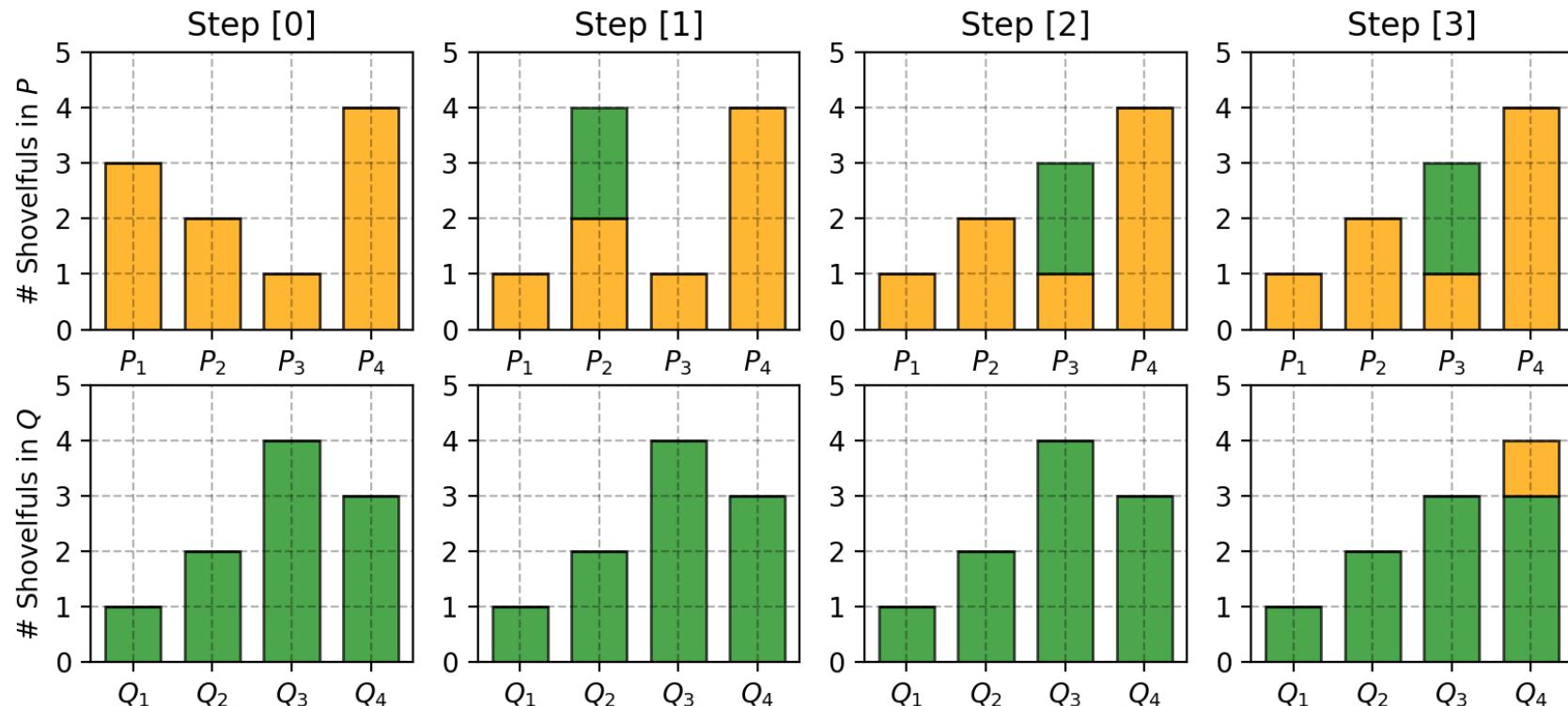


Wasserstein Distance



<https://towardsdatascience.com/comparison-of-distributions-with-earth-movers-distance-71f714440923?gi=3b6f5b39cceb>

Wasserstein Distance



Wasserstein Distance

Formal Definition

$$W_p(\mu, \nu) = \inf_{\gamma \in \Gamma(\mu, \nu)} \left(\mathbf{E}_{(x,y) \sim \gamma} d(x, y)^p \right)^{\frac{1}{p}},$$

$\Gamma(\mu, \nu)$ Set of all possible couplings
(continuous transportation plans)

$d(x, y)$ Distance function between two points

$\inf_{\gamma \in \Gamma(\mu, \nu)}$ Find the cheapest transportation plan

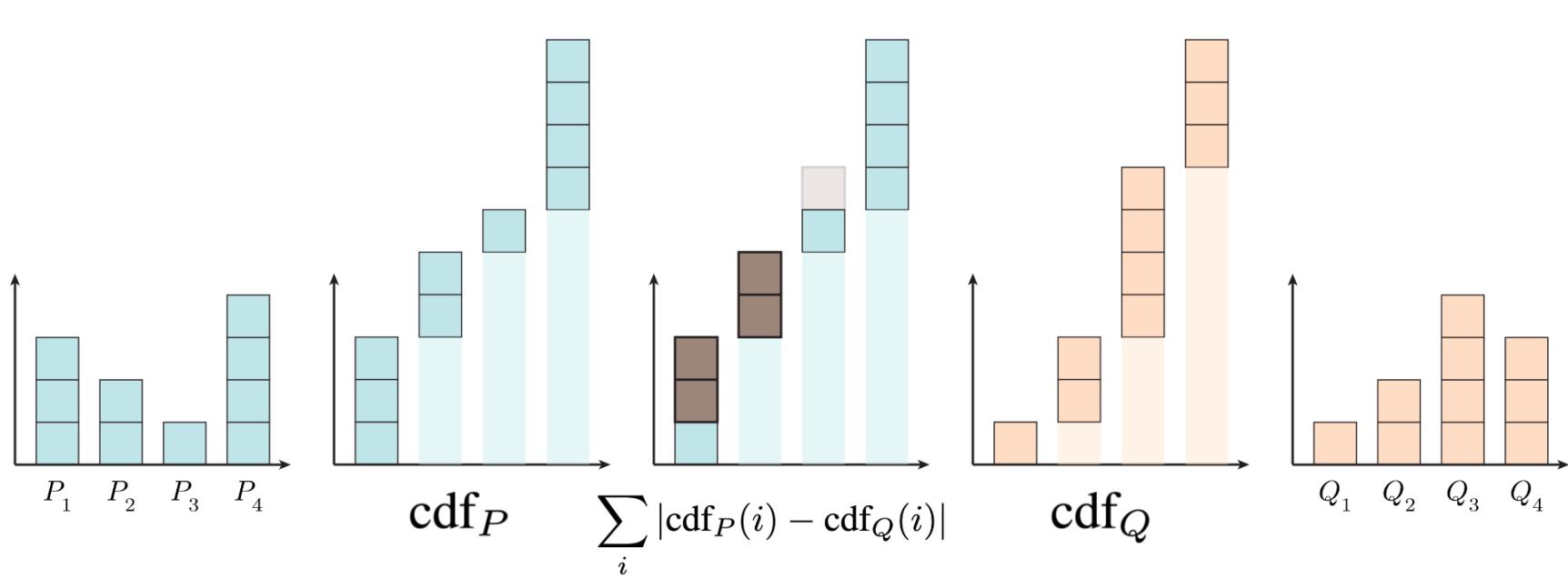
Wasserstein Distance in 1D

If we have the CDF of distributions for 1D distributions p and q

1-Wasserstein distance:

$$W_1(p, q) = \int_x |\text{cdf}_p(x) - \text{cdf}_q(x)| dx$$

Wasserstein Distance in 1D



$$W(P, Q) = 5 \times \blacksquare$$

Wasserstein GAN

In a nutshell, to replace JS divergence with **Wasserstein distance**

Do not need to exactly compute Wasserstein distance

Critic update:

$$\begin{array}{ll} & \text{Logits for binary classification} \\ \text{Vanilla} & \max_D \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))] \\ \\ \text{W-GAN} & \max_{w \in \mathcal{W}} \mathbb{E}_{x \sim \mathbb{P}_r} [f_w(x)] - \mathbb{E}_{z \sim p(z)} [f_w(g_\theta(z))] \\ & \quad \quad \quad \text{Any real value} \\ & \quad \quad \quad \text{Same formulation for discriminator} \end{array}$$

Key Difference: (asymmetric) logits to symmetric critic values

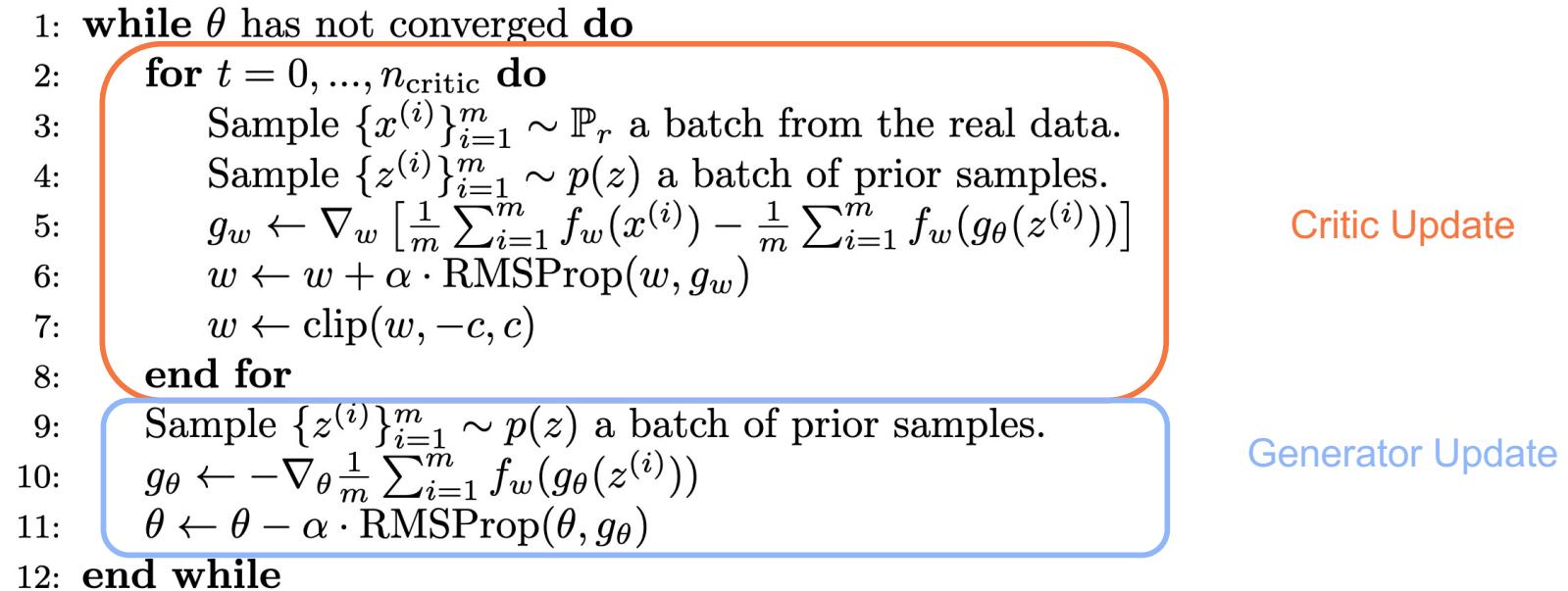
Wasserstein GAN

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.

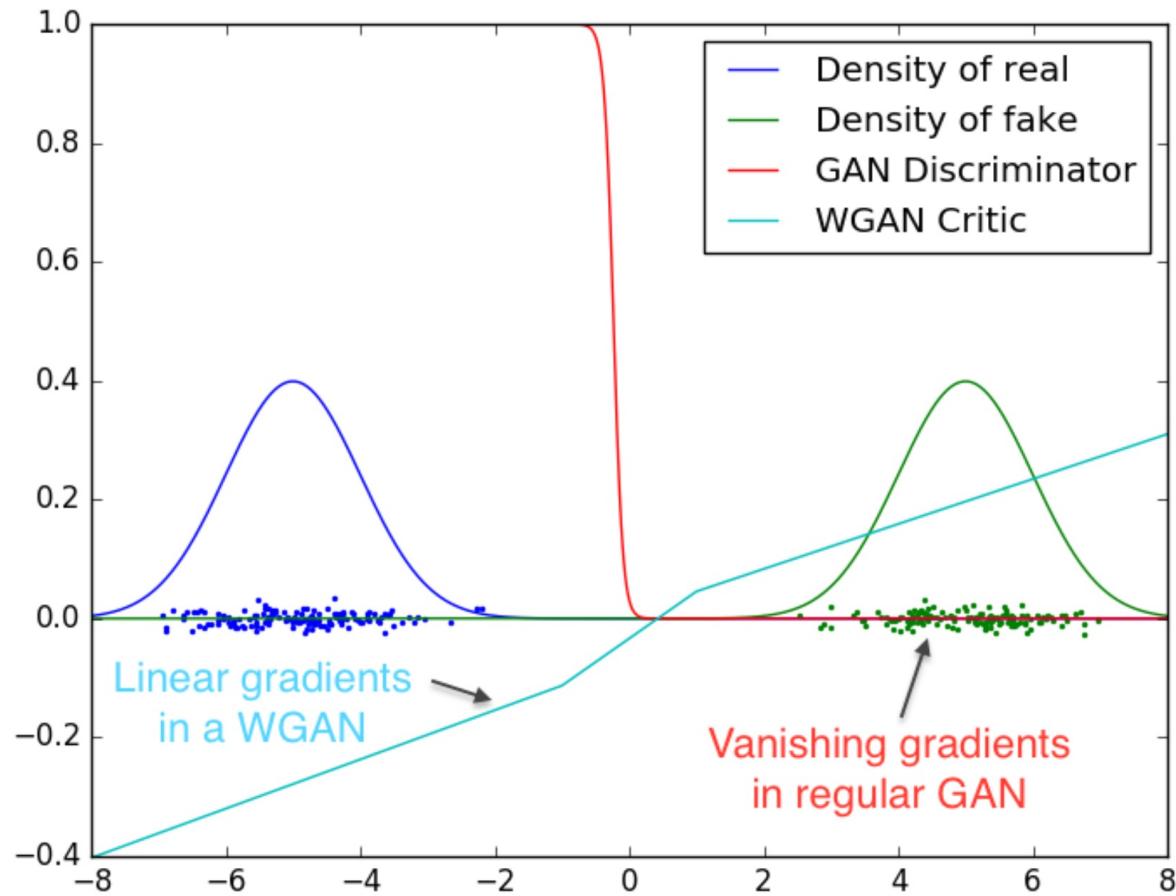
```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m}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



Critic Update

Generator Update

W-GAN vs. Vanilla GAN



WGAN-GP

Use gradient penalty instead of weight clipping to stabilize training

$$L = \underbrace{\mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}} [(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2]}_{\text{Our gradient penalty}}.$$

f-GAN

Instead of JS divergence, use the more general f-divergence

Name	$D_f(P\ Q)$	Generator $f(u)$
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} dx$	$u \log u$
Reverse KL	$\int q(x) \log \frac{q(x)}{p(x)} dx$	$-\log u$
Pearson χ^2	$\int \frac{(q(x)-p(x))^2}{p(x)} dx$	$(u - 1)^2$
Squared Hellinger	$\int \left(\sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx$	$(\sqrt{u} - 1)^2$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$	$-(u + 1) \log \frac{1+u}{2} + u \log u$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx - \log(4)$	$u \log u - (u + 1) \log(u + 1)$

f-GAN

Instead of JS divergence, use the more general f-divergence

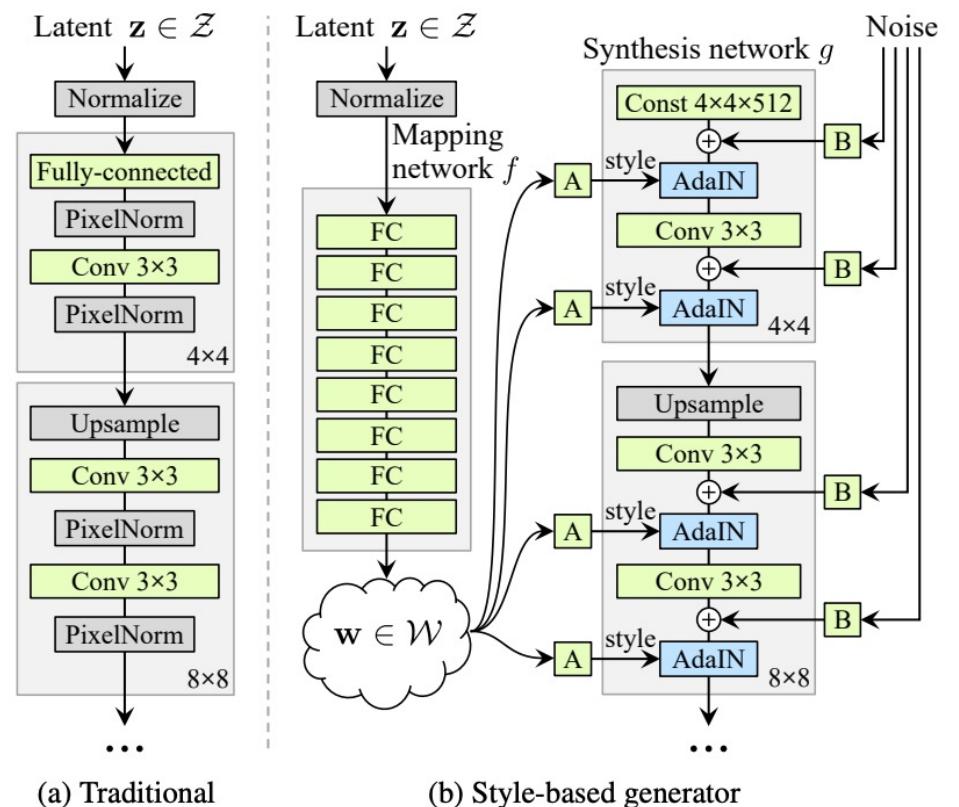
$$F(\theta, \omega) = \underbrace{\mathbb{E}_{x \sim P} [g_f(V_\omega(x))]}_{\text{Real}} + \underbrace{\mathbb{E}_{x \sim Q_\theta} [-f^*(g_f(V_\omega(x)))]}_{\text{Fake}}$$

Name	Output activation g_f	dom_{f^*}	Conjugate $f^*(t)$	$f'(1)$
Kullback-Leibler (KL)	v	\mathbb{R}	$\exp(t - 1)$	1
Reverse KL	$-\exp(-v)$	\mathbb{R}_-	$-1 - \log(-t)$	-1
Pearson χ^2	v	\mathbb{R}	$\frac{1}{4}t^2 + t$	0
Squared Hellinger	$1 - \exp(-v)$	$t < 1$	$\frac{t}{1-t}$	0
Jensen-Shannon	$\log(2) - \log(1 + \exp(-v))$	$t < \log(2)$	$-\log(2 - \exp(t))$	0
GAN	$-\log(1 + \exp(-v))$	\mathbb{R}_-	$-\log(1 - \exp(t))$	$-\log(2)$

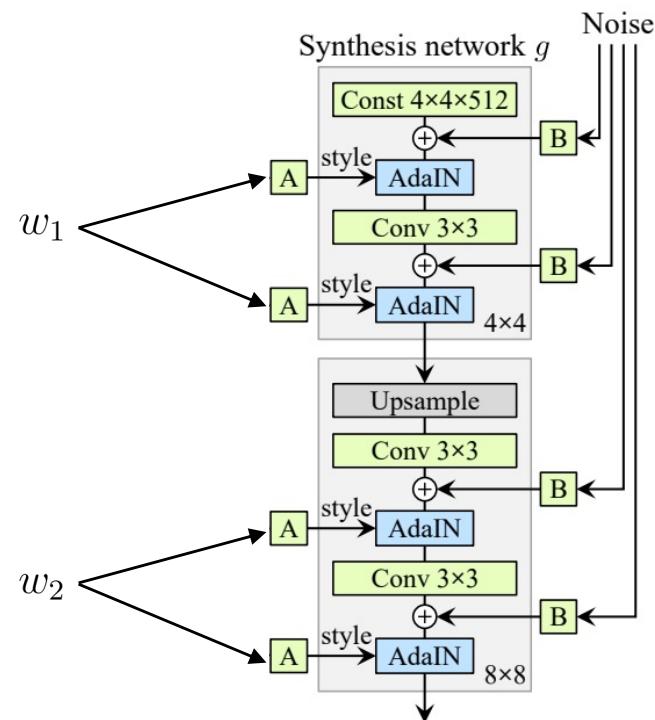
StyleGAN

- Latent as style condition
- Multi-resolution noise
- Blending with AdaIN
(module for style transfer)

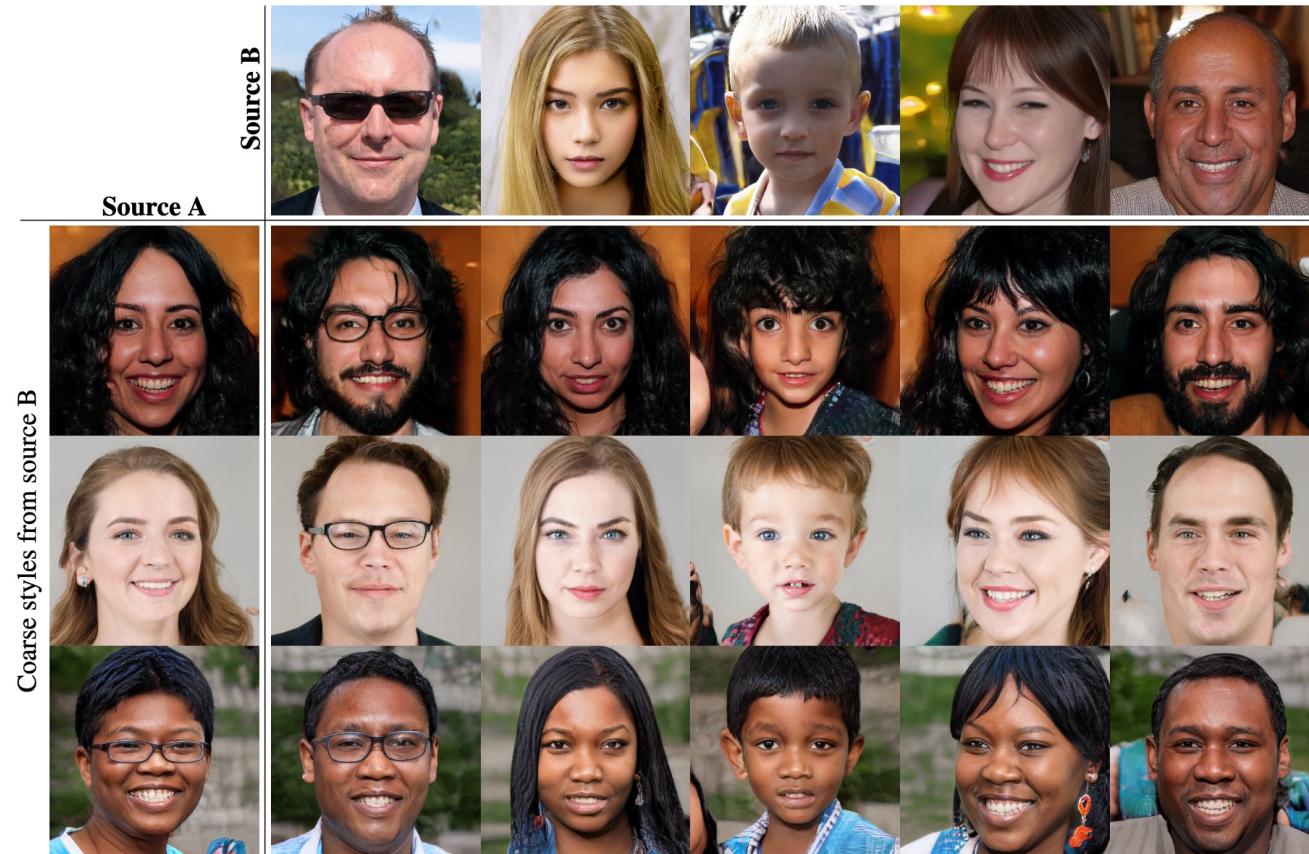
$$\text{AdaIN}(x, y) = \sigma(y) \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$



Style Mixing with StyleGAN

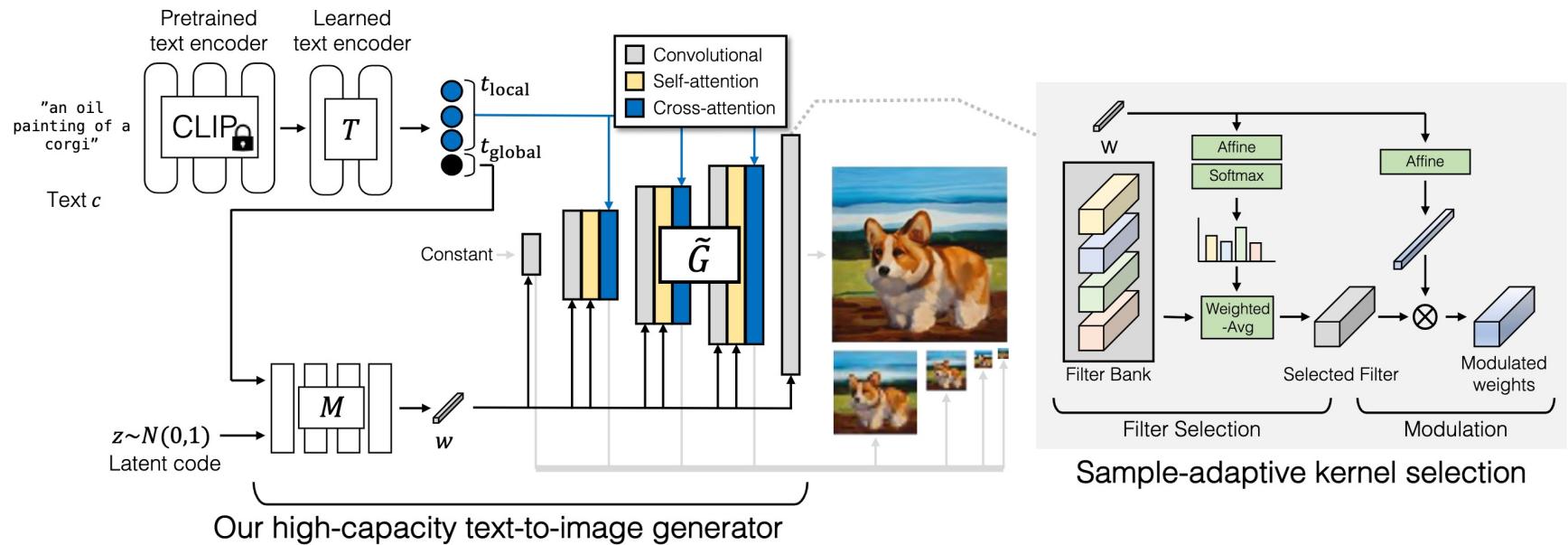


StyleGAN Families



GigaGAN

Heavily engineered architecture and losses
for large-scale text-conditioned models



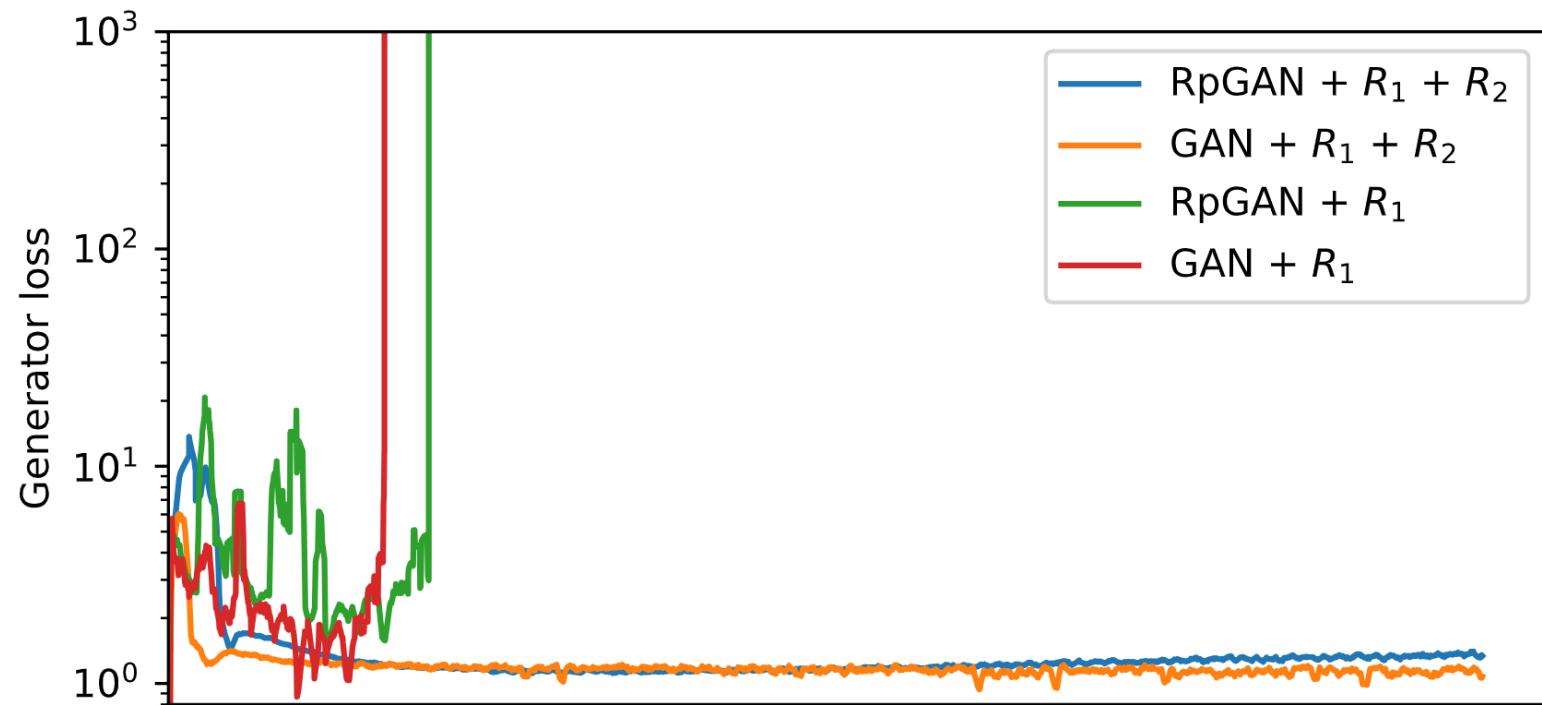
Latest: Relativistic GANs + Regularizations

$$\mathcal{L}(\theta, \psi) = \mathbb{E}_{\substack{z \sim p_z \\ x \sim p_{\mathcal{D}}}} [f(D_{\psi}(G_{\theta}(z)) - D_{\psi}(x))] \quad (\text{typically with a saturating f function})$$

$$R_1(\psi) = \frac{\gamma}{2} \mathbb{E}_{x \sim p_{\mathcal{D}}} \left[\|\nabla_x D_{\psi}\|^2 \right]$$
$$R_2(\theta, \psi) = \frac{\gamma}{2} \mathbb{E}_{x \sim p_{\theta}} \left[\|\nabla_x D_{\psi}\|^2 \right]$$

Zero-centered Gradient Penalty

Latest: Relativistic GANs + Regularizations



Latest: Relativistic GANs + Regularizations

Competitive performance

Remains to be verified on larger models

Model	NFE↓	FID↓
BigGAN-deep [3]	1	4.06
DDPM [20]	250	11.0
DDIM [76]	50	13.7
ADM [7]	§250	2.91
EDM [33]	79	2.23
CT [79]	2	11.1
CD [79]	3	4.32
iCT-deep [77]	2	2.77
DMD [96]	1	2.62
Ours—Config E	1	2.09
<i>With ImageNet feature leakage [41]:</i>		
StyleGAN-XL* [69]	1	1.52

Adversary as a Loss

Adversarial loss

Not necessary to have a generator $g(z)$

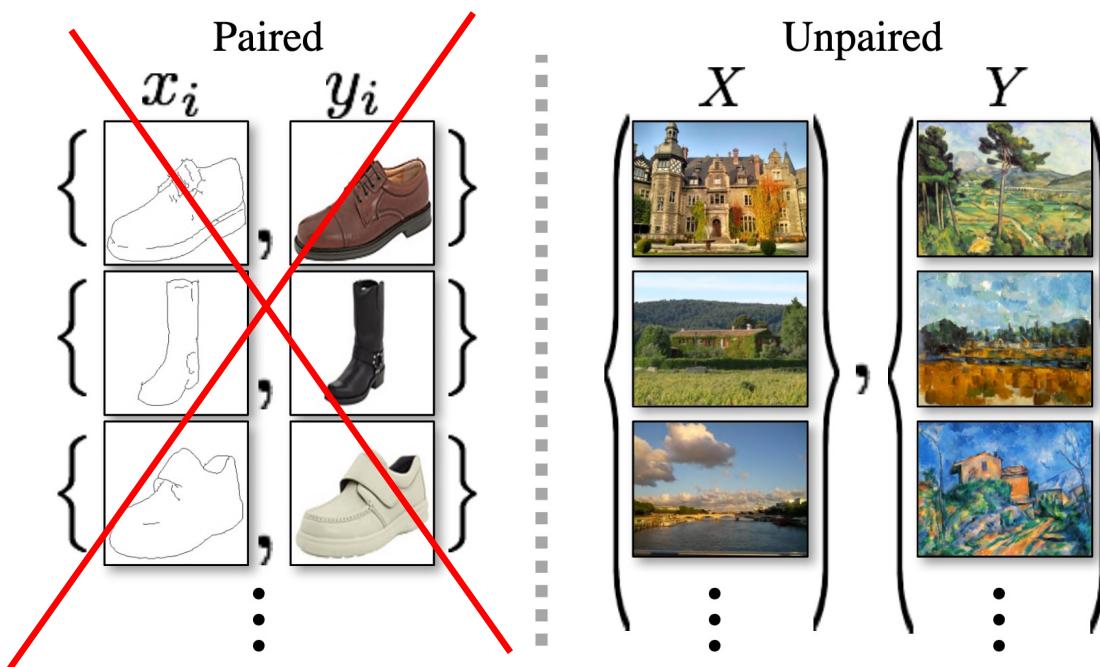
- Can be used with VAEs, diffusion models and other generative models

Can use discriminator to refine results

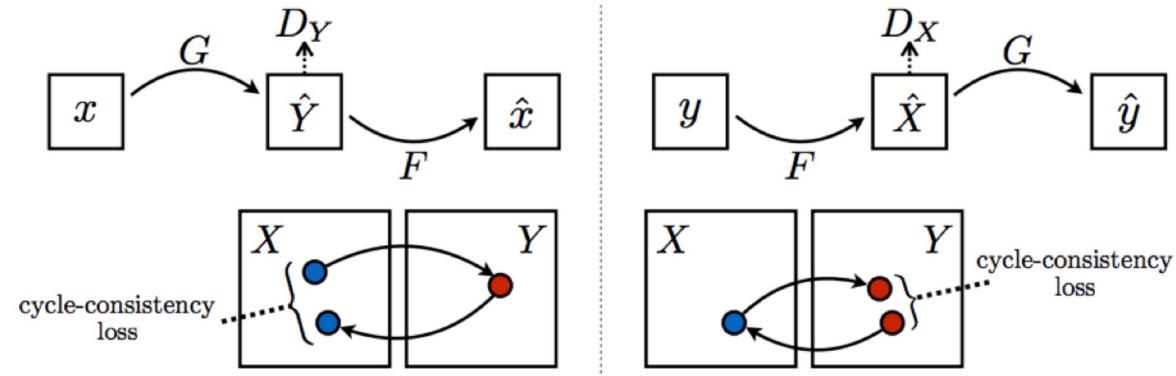
- Supplement to reconstruction losses for better visual quality
- Leverage implicit priors in neural nets (and adversarial training)

CycleGAN

How to learn style transfer from unpaired data?



CycleGAN



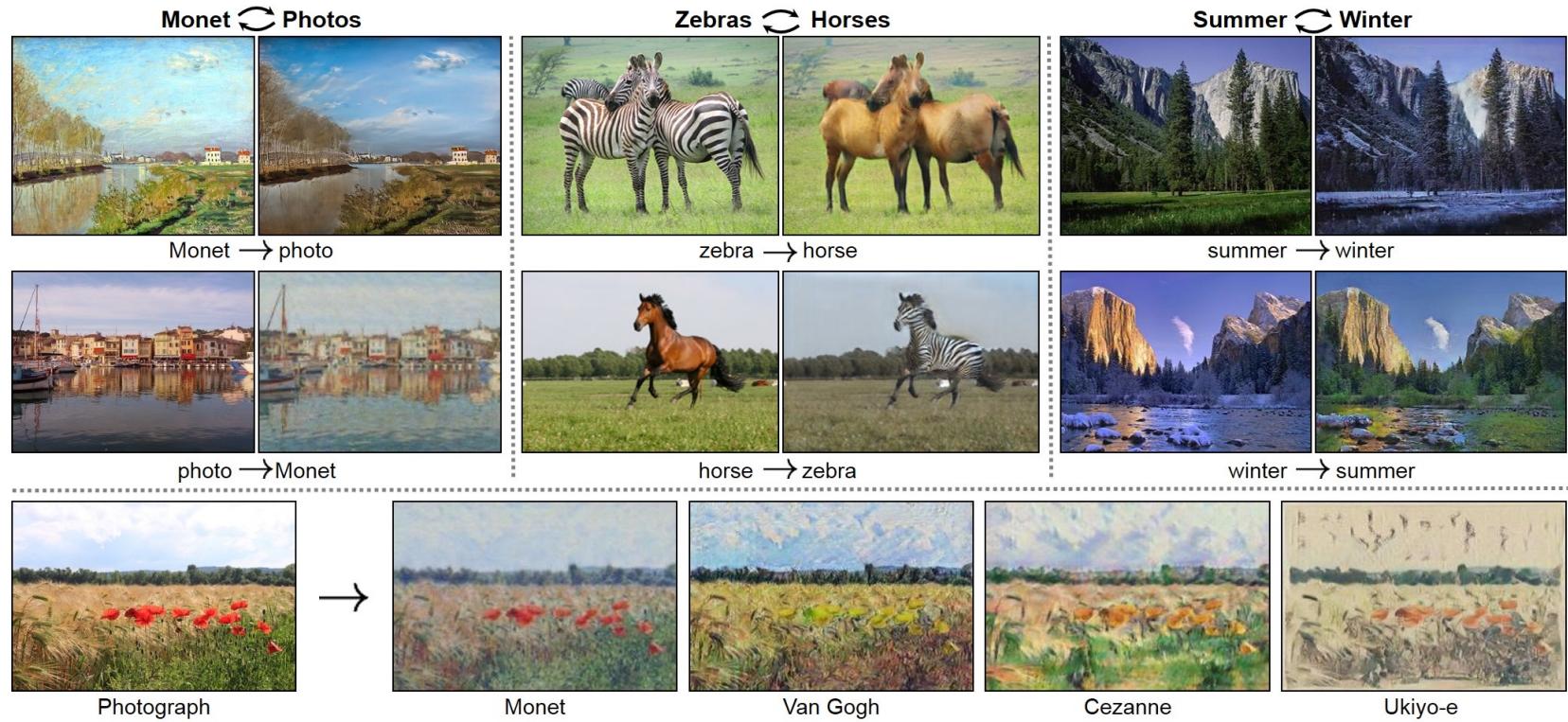
$$\begin{aligned}\mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] \\ & + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1].\end{aligned}$$

Cycle Consistency

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

Full loss = GAN losses + Cycle consistency

CycleGAN



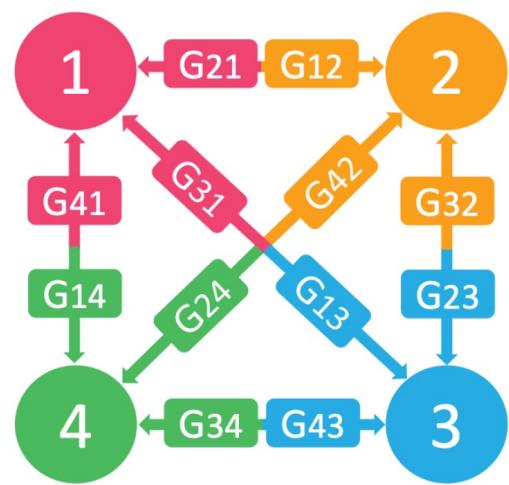
<https://junyanz.github.io/CycleGAN/>

63

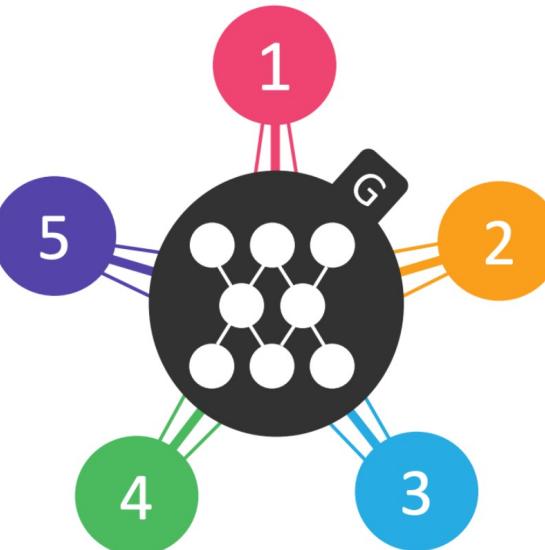
StarGAN

More than one domain?

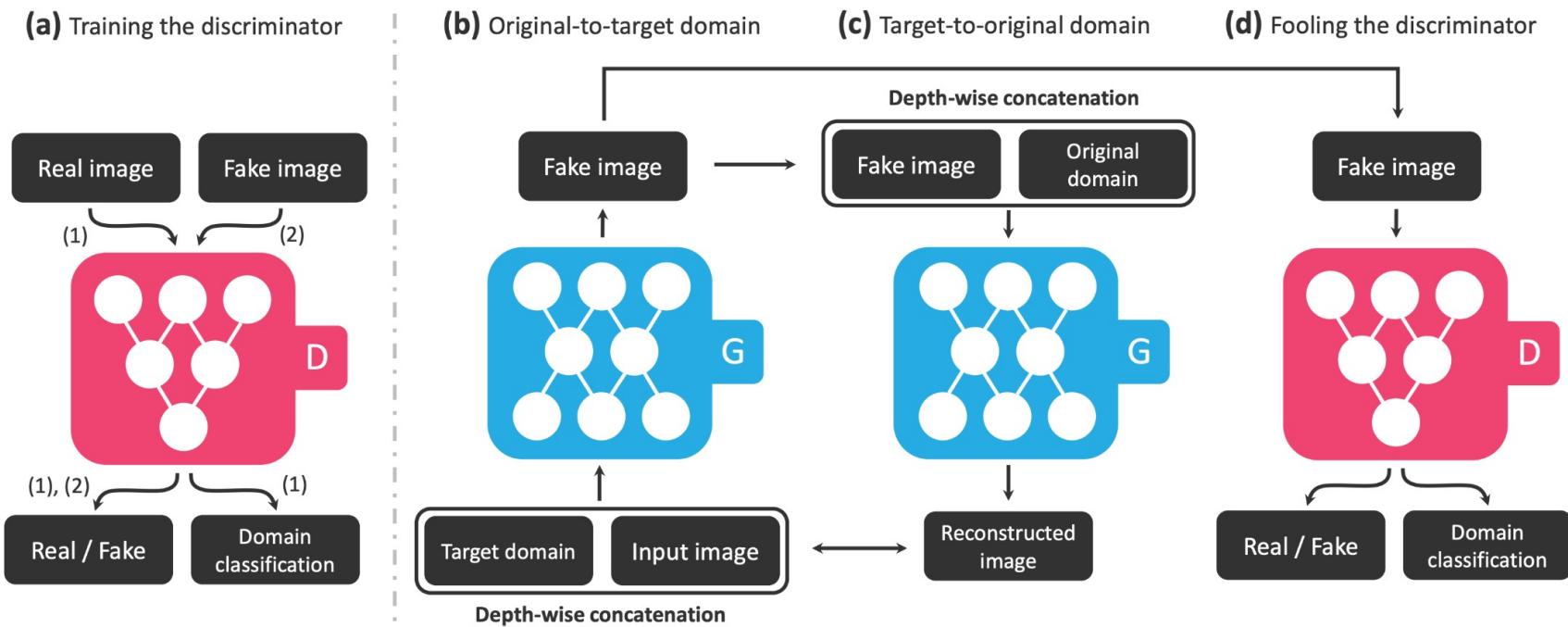
(a) Cross-domain models



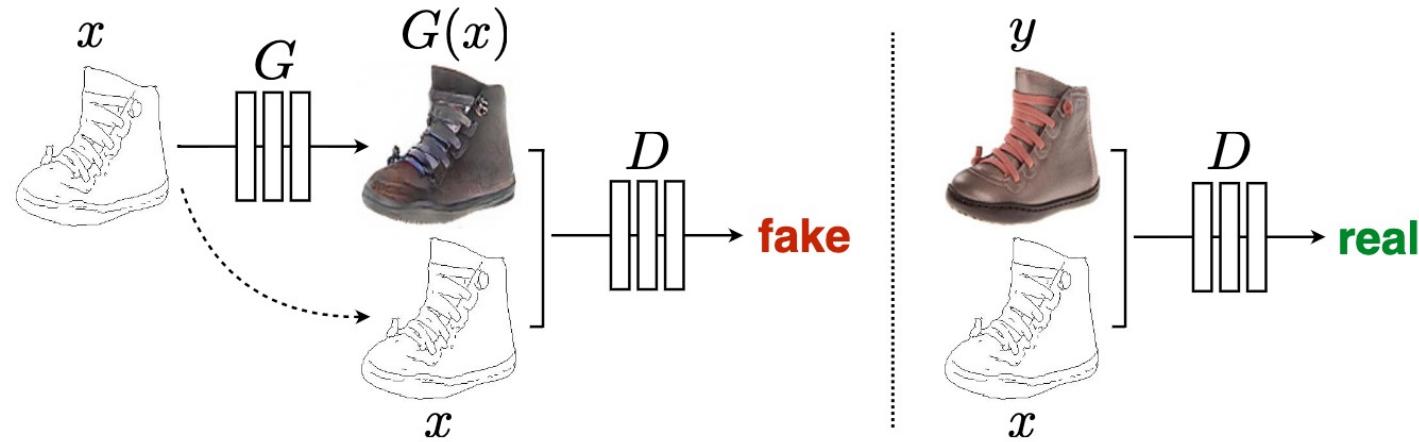
(b) StarGAN



StarGAN



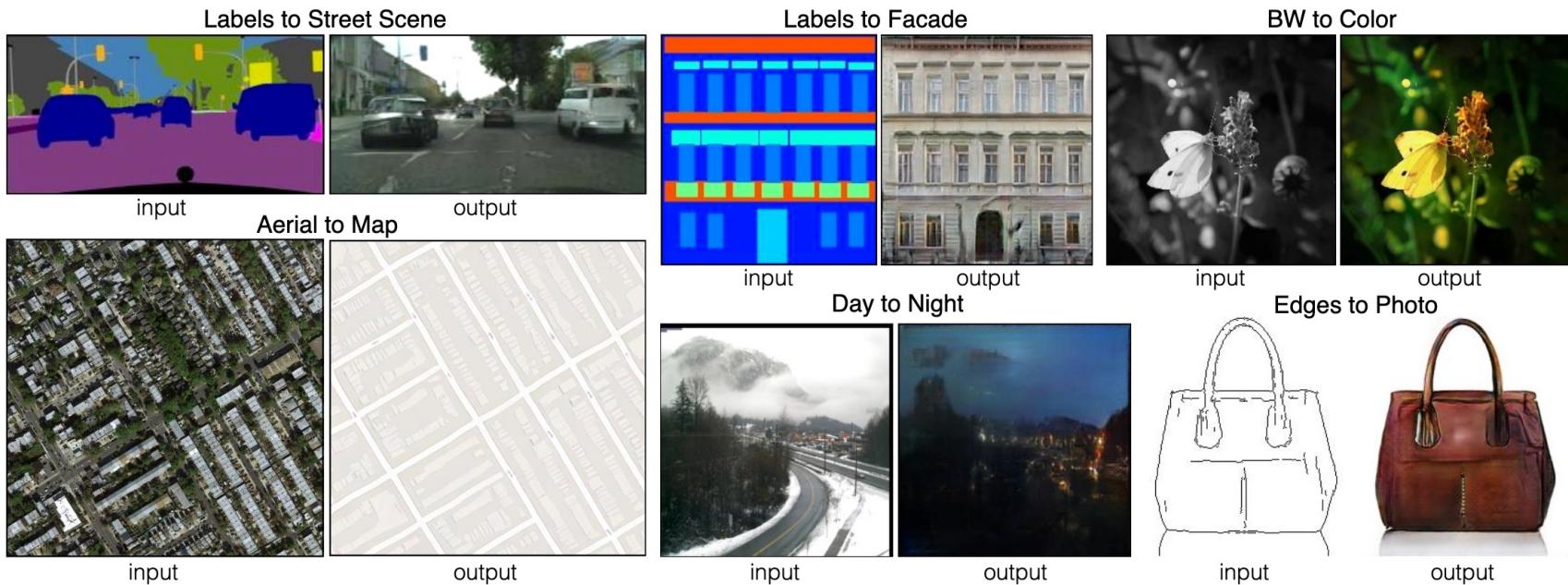
Pix2Pix



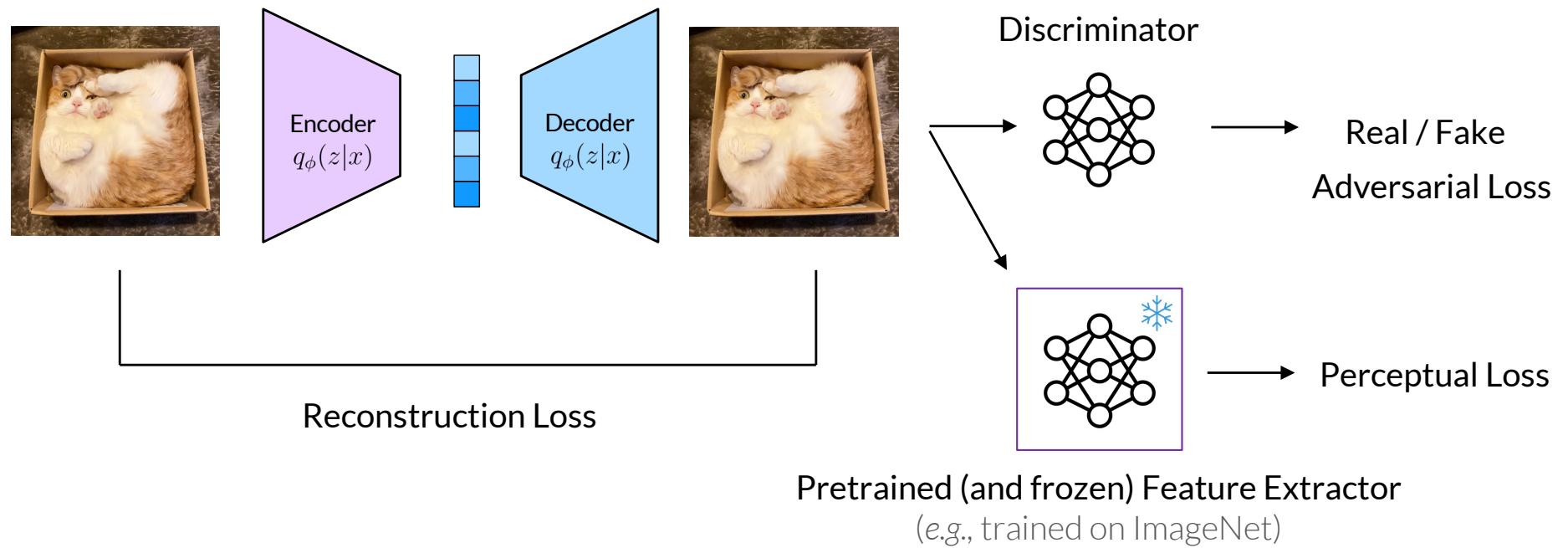
$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)$$

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

Pix2Pix



VAE + GAN + Perceptual Loss



Used for modern image encoders / tokenizers for images
(StableDiffusion, Sora, etc.)

Summary

- Vanilla GAN
- Variants (W-GAN, etc.)
- Adversary as a Loss