

# CS7150 Deep Learning

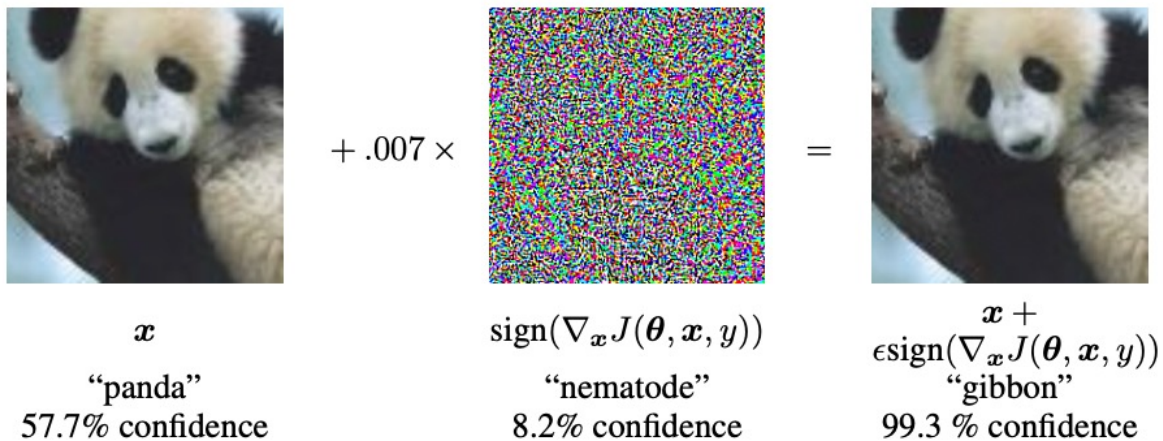
Jiaji Huang

<https://jiaji-huang.github.io>

03/23/2024

# Recap of Last Lecture

- Compare Networks by Feature Similarity
- Attribute a network's decision to
  - Input features
  - Training samples
- Construct adversarial samples

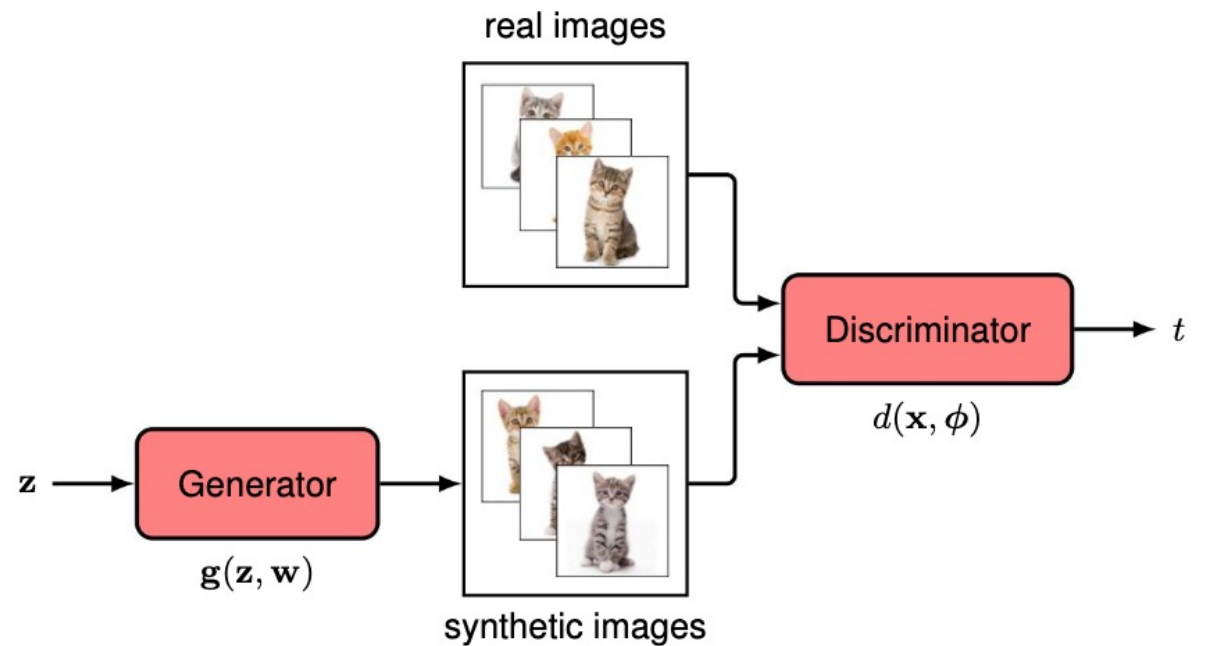


# Agenda

- Generative Adversarial Network (GAN)
- Earth Mover Distance and Wasserstein GAN

# Adversarial Training

- Generate datum “indistinguishable” from real ones
- Maximize a discriminator’s loss
- A competing game between **generator** and **discriminator**
- $z$ : latent variable



# Formalize

- Discriminator: a binary classifier,  $p(\mathbf{x} \text{ is real image}) = d(\mathbf{x}, \phi)$

$$\max_{\phi} \frac{1}{N_{real}} \sum_{\mathbf{x}:real} \log d(\mathbf{x}, \phi) + \frac{1}{N_{gen}} \sum_{\mathbf{x}:gen} \log(1 - d(\mathbf{x}, \phi))$$

- Generator: minimize the reward

$$\min_w \max_{\phi} \frac{1}{N_{real}} \sum_{\mathbf{x}:real} \log d(\mathbf{x}, \phi) + \frac{1}{N_{gen}} \sum_{\mathbf{z}} \log(1 - d(g(\mathbf{z}, w), \phi))$$

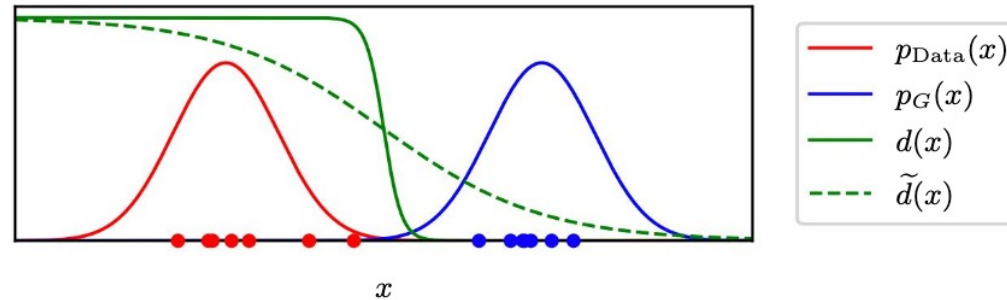
- Grad ascent on  $\phi$ . Grad descent on  $w$ :

$$\begin{aligned} \Delta \phi &= \eta \nabla_{\phi} L(w, \phi) \\ \Delta w &= -\eta \nabla_w L(w, \phi) \end{aligned}$$

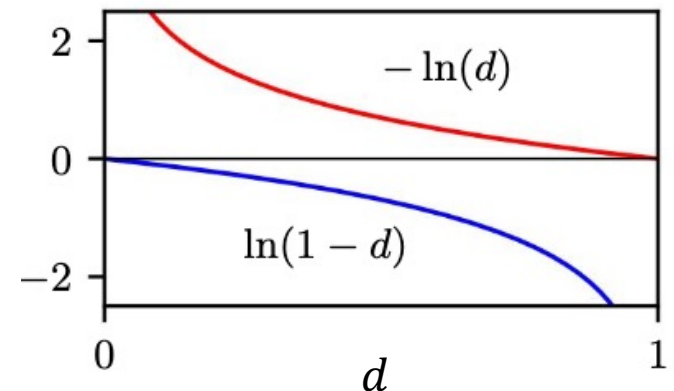
$g(\mathbf{z}, w)$

# Potential Issues

- Potential issue: “mode collapse”



- Too easy to discriminate,  $\nabla_{\phi} E = 0$ , stop learning
- Smooth the  $d(\mathbf{x})$  to  $\tilde{d}(\mathbf{x})$ , enables gradients
- Consider instead  $\log d(\mathbf{x}_{gen}, \phi)$



# How it works



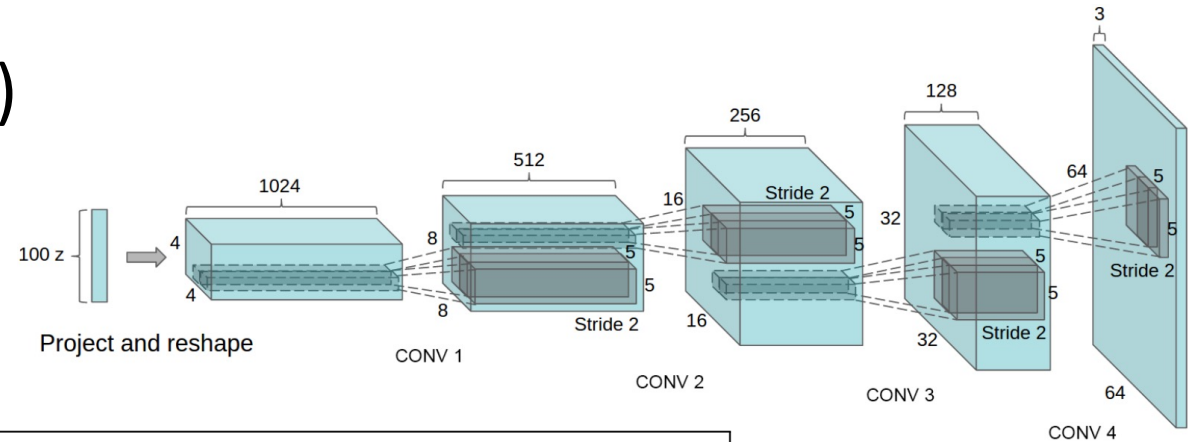
Generated Images and their nearest neighbors (in yellow box). Left: mnist; Middle: Toronto Face Dataset; Right: cifar-10



Generated images along a path of Interpolated z

# Generator Architectures

- Deep Convolutional GAN (DCGAN)



## Architecture guidelines for stable Deep Convolutional GANs

- Replace any pooling layers with strided convolutions (discriminator) and fractional-strided convolutions (generator).
- Use batchnorm in both the generator and the discriminator.
- Remove fully connected hidden layers for deeper architectures.
- Use ReLU activation in generator for all layers except for the output, which uses Tanh.
- Use LeakyReLU activation in the discriminator for all layers.

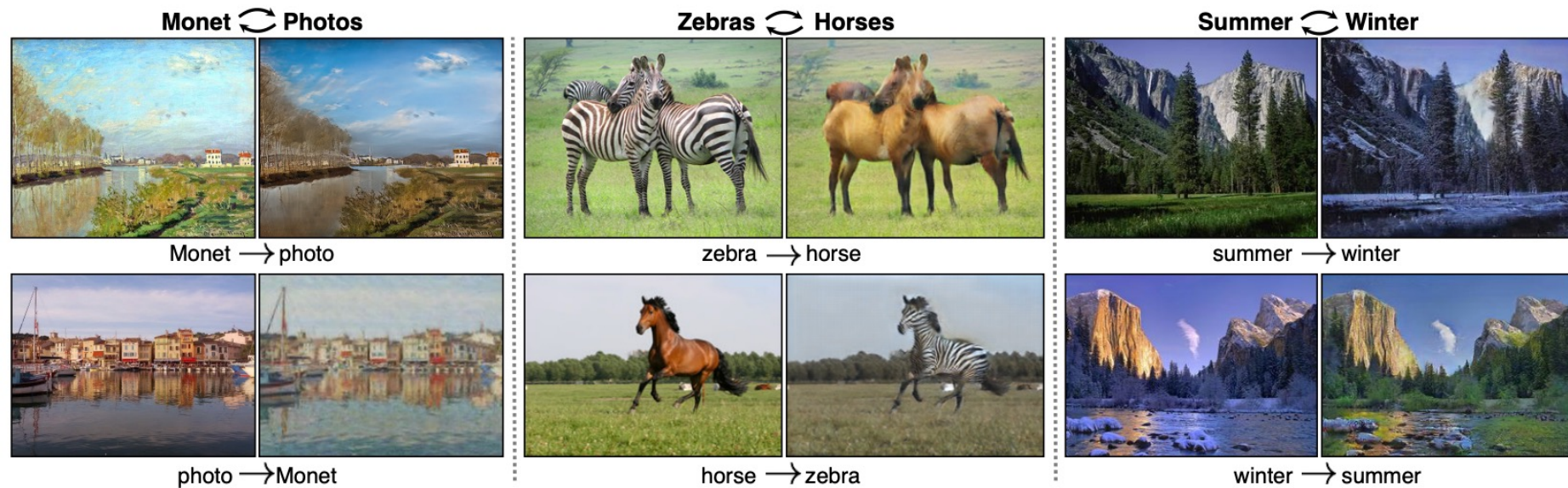


# DCGAN

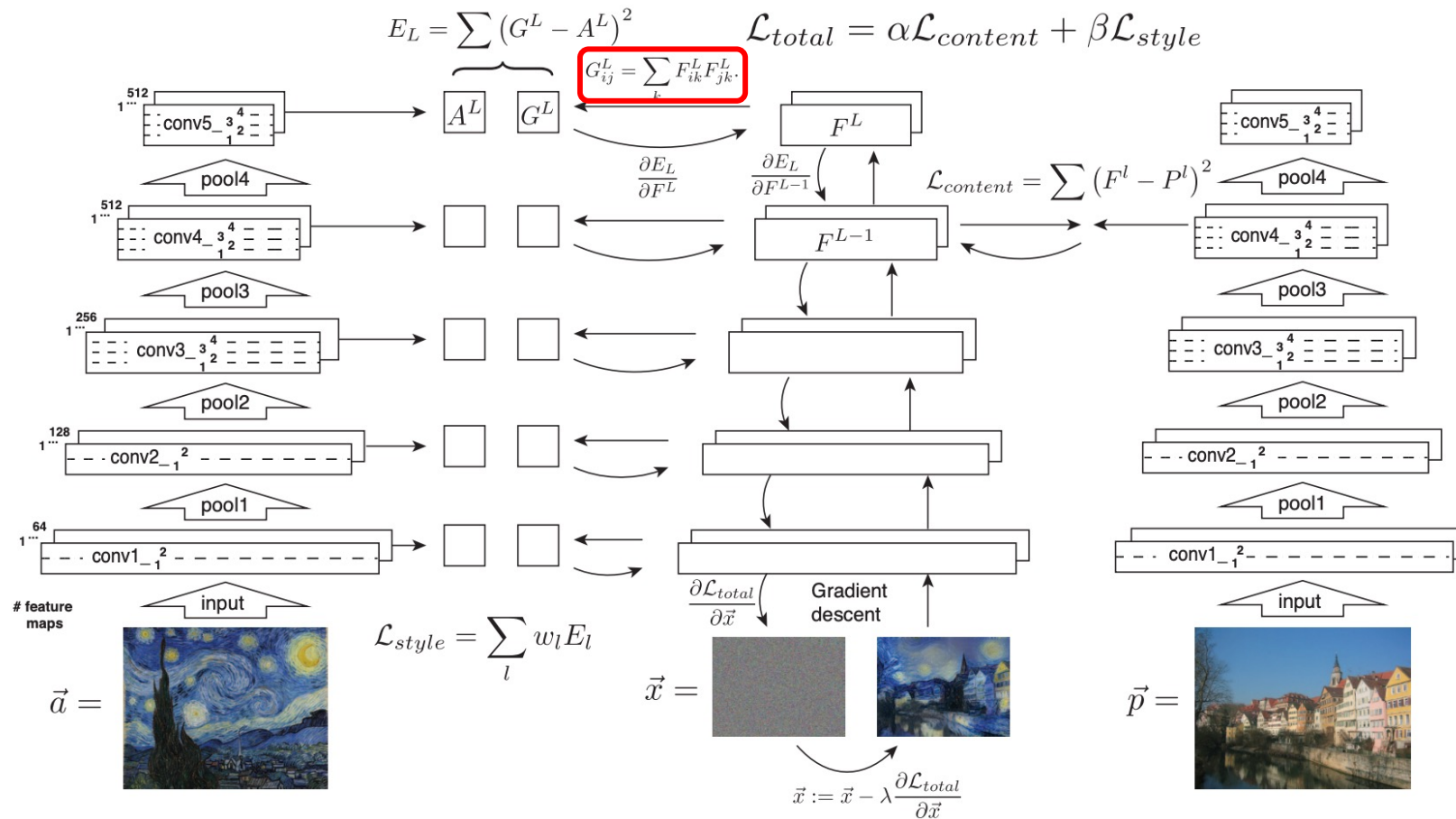


# Cycle GAN

- Learn mapping between two collections

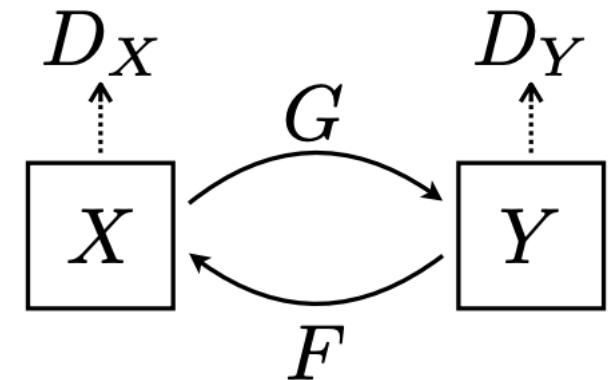


# Revisit Image Style Transfer (Lecture 3)



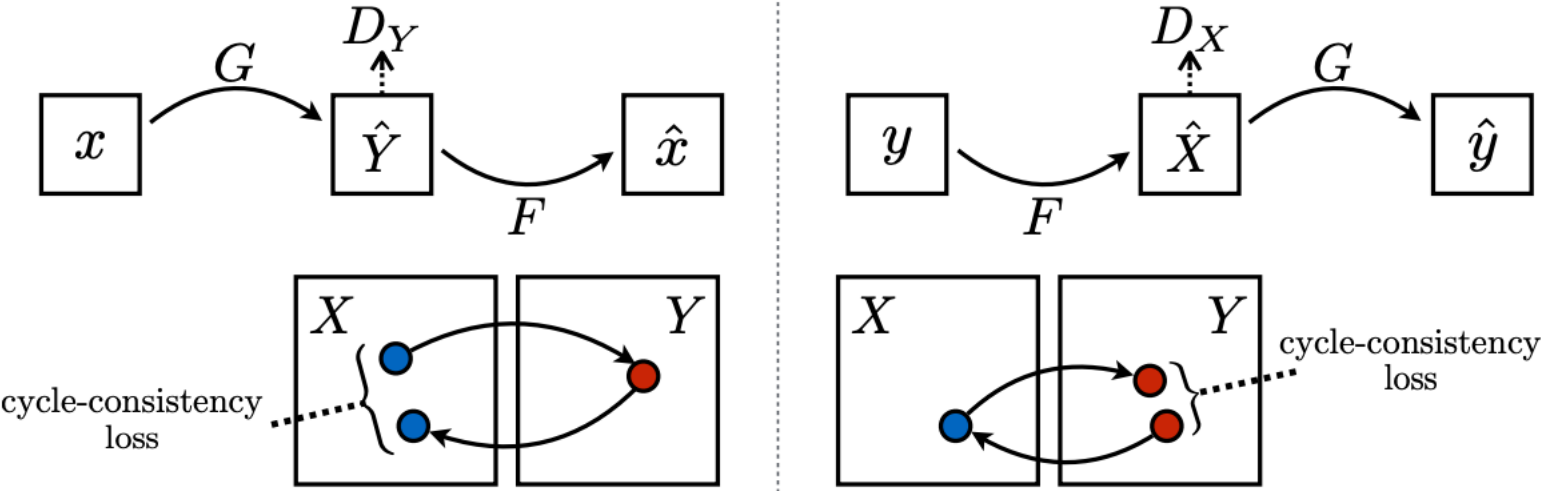
# Formalize: Cycle GAN

- Two collections of images:  $X$  and  $Y$
- Learn mapping  $G: X \mapsto Y$  to confuse the discriminator  $D_Y$ 
$$\min_G \max_{D_Y} \mathcal{L}(G, D_Y, X, Y)$$
$$= \mathbb{E}_{y \sim p(y)} [\log D_Y(y)] + \mathbb{E}_{x \sim p(x)} [\log(1 - D_Y(G(x)))]$$
- Likewise, Learn mapping  $F: Y \mapsto X$  via  $\min_F \max_{D_X} \mathcal{L}(F, D_X, X, Y)$
- Question: Key difference against style transfer?
  - Cycle GAN learns mapping between two styles
  - Style transfer only works for a specific image



# Regularize

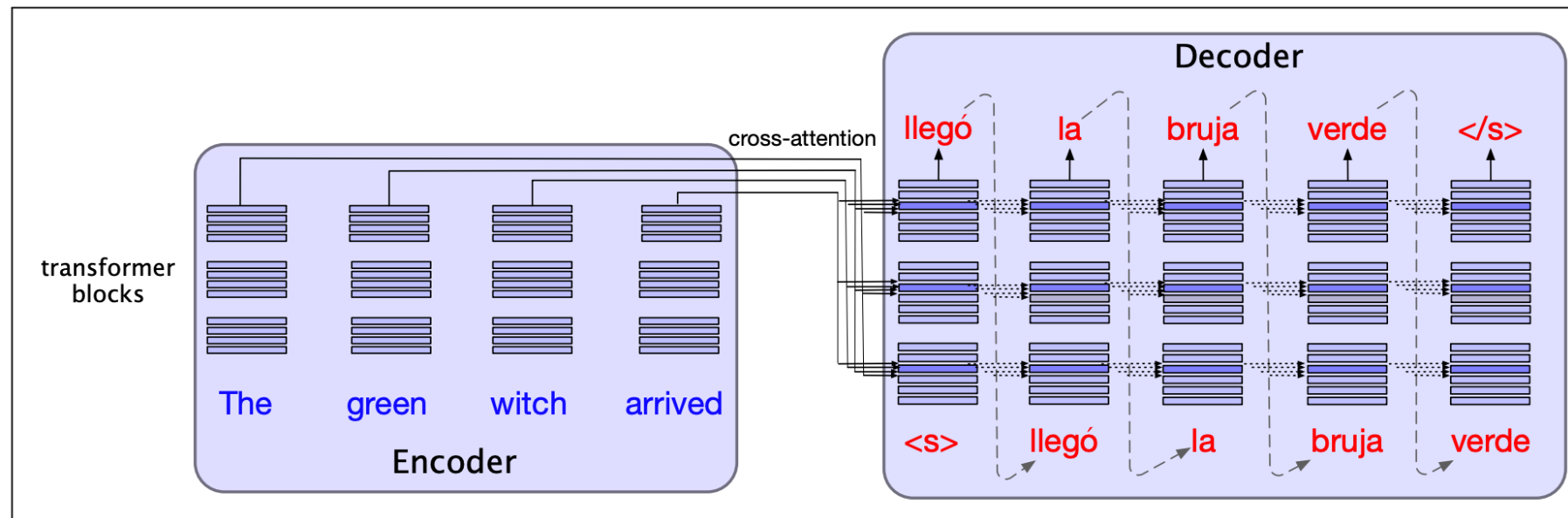
$$\mathcal{L}_{cycle} = \mathbb{E}_{x \sim p(x)} \|F(G(x)) - x\|_1 + \mathbb{E}_{y \sim p(y)} \|G(F(y)) - y\|_1$$



$$\min_{G,F} \left\{ \max_{D_Y} \mathcal{L}(G, D_Y, X, Y) + \max_{D_X} \mathcal{L}(F, D_X, X, Y) + \mathcal{L}_{cycle} \right\}$$

# GAN for Unsupervised Machine Translation

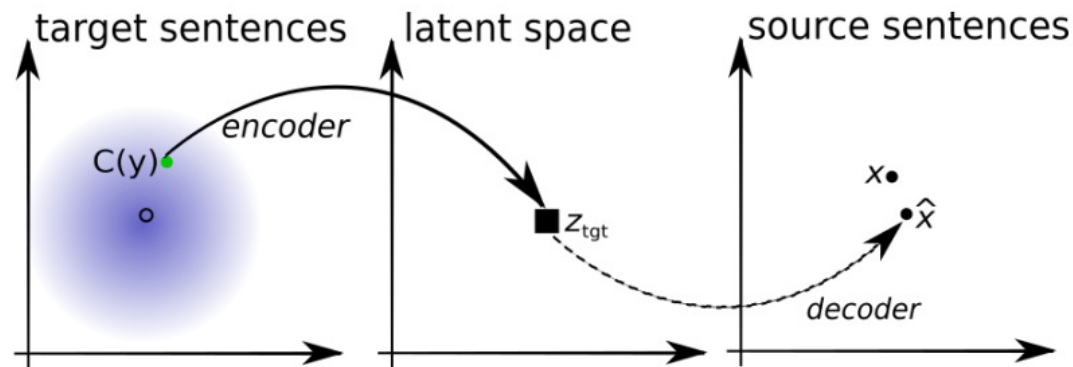
- Recap enc-dec models for machine translation



**Figure 10.5** The encoder-decoder transformer architecture for machine translation. The encoder uses the transformer blocks we saw in Chapter 9, while the decoder uses a more powerful block with an extra **cross-attention** layer that can attend to all the encoder words. We'll see this in more detail in the next section.

# GAN for Unsupervised Machine Translation

- Setup: monolingual corpora, no translation pair
- Encoder and decoder are language agnostic
- Back-translation:
  - Sample source sentence  $x$ , translate to target  $y$  using current model
  - Train as if supervised case, on  $(x, C(y))$ , where  $C(\cdot)$  adds noise



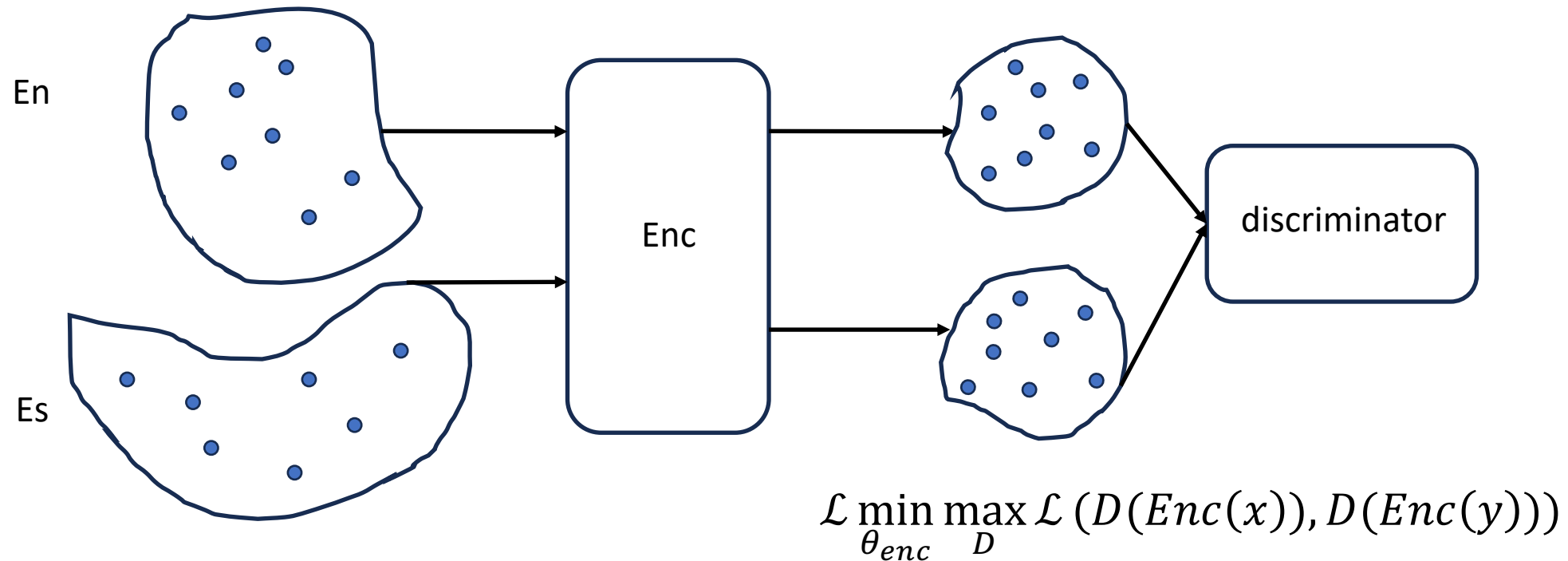
$$\mathcal{L}(\theta_{enc}, \theta_{dec}, C(y), x)$$

And likewise

$$\mathcal{L}(\theta_{enc}, \theta_{dec}, C(x), y)$$

# GAN for Unsupervised Machine Translation

- Languages should share some traits
- $Enc(x)$ 's and  $Enc(y)$ 's should be indistinguishable



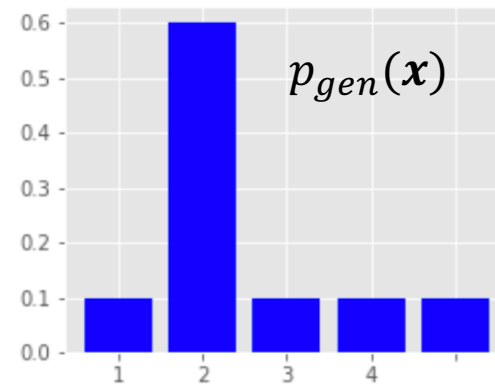
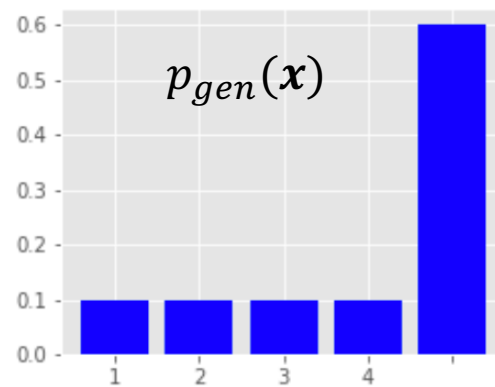
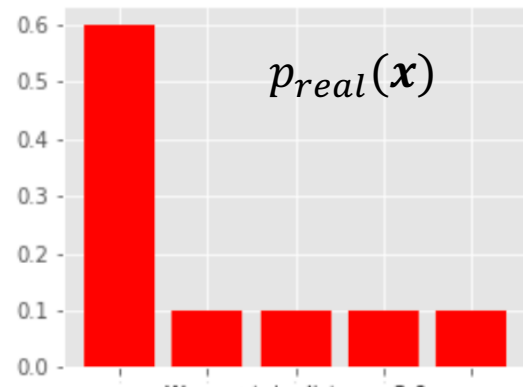


# Agenda

- Generative Adversarial Network (GAN)
- Wasserstein GAN

# Improve Objective

- Don't completely rely on discriminator
- Can we penalize the distance between  $p_{real}(\mathbf{x})$  and  $p_{gen}(\mathbf{x})$ ?
- e.g., KL divergence:  $\sum_{\mathbf{x} \sim p_{real}} \log \frac{p_{real}(\mathbf{x})}{p_{gen}(\mathbf{x})}$

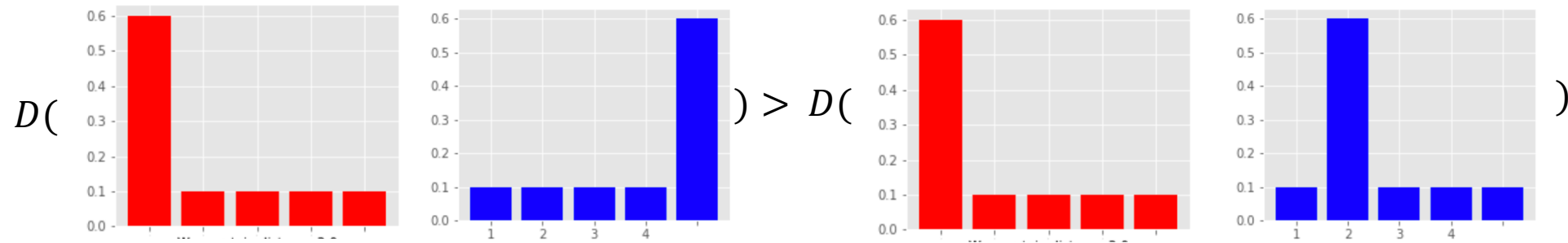


- Same KL-div 0.8959

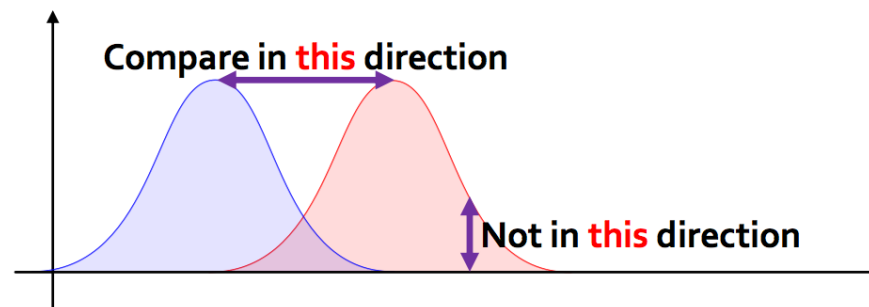
Example from [stackexchange](#)

# Issue with KL-div

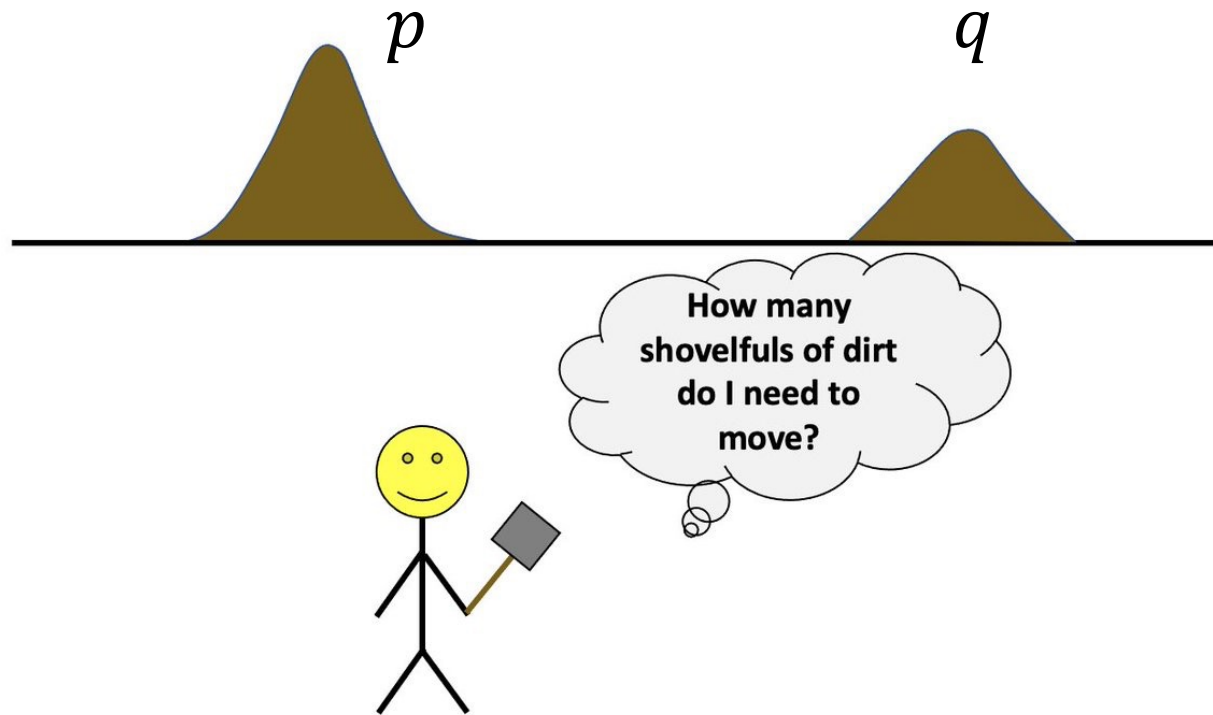
- But Intuition is



- We should consider the difference along the  $x$  axis also!



# Earth Mover Distance (EMD)



Earth mover's distance is an intuitive metric that denotes the minimum cost flow between two distributions (in other words, minimum work that needs to be done to turn one into the other)

Cartoon from [post on X](#)

- Discretized case

$$\min_M \sum_{i,j} \gamma_{i,j} D_{i,j}$$

$$\text{s.t. } \sum_i \gamma_{i,j} = q_j \text{ and} \\ \sum_j \gamma_{i,j} = p_i$$

- $\gamma$  is joint distribution with marginals  $p$  and  $q$

# Generalize: Wasserstein Distance (WD)

- Generalize the form of  $d_{i,j}$ , and continuous support

$$WD(p, q) = \min_{\gamma \in \Gamma(p, q)} \mathbb{E}_{(x, y) \sim \gamma} [d(x, y)^p]^{1/p}$$

- So for generative model, we seek to minimize

$$WD(p_{real}, p_{gen})$$

- Seems quite hard, but there is a duality

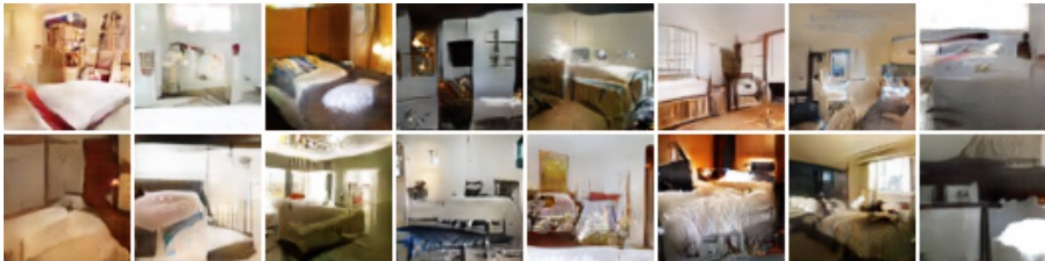
$$WD(p, q) = \max_f \mathbb{E}_{x \sim p} [f(x)] - \mathbb{E}_{x \sim q} [f(x)]$$

- Parameterize  $f(\cdot)$  as a deep network with parameter  $\phi$  (Wasserstein-GAN)

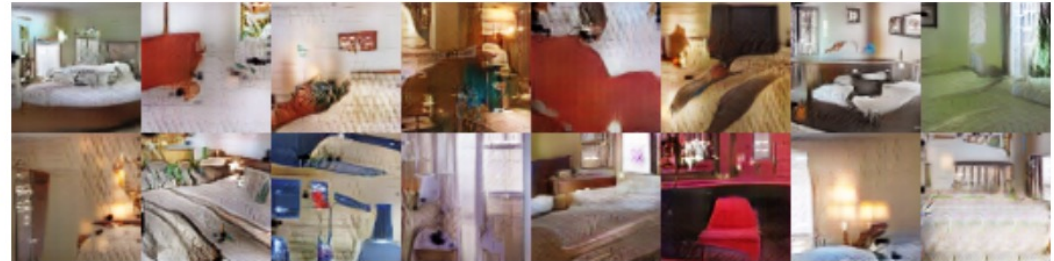
$$\min_w \max_{\phi} \frac{1}{N_{real}} \sum_{x:real} f(x; \phi) - \frac{1}{N_{gen}} \sum_z f(g(z, w))$$

# Compare to Traditional GAN

W + DC



DC



W + DC (tweaks removed)



DC (tweaks removed)

