

# Deep Generative Models II

## Generative Adversarial Networks

Machine Perception

Otmar Hilliges

30 April 2020

Last time

VAEs

Generative Adversarial Networks (Training objective)

# This Week

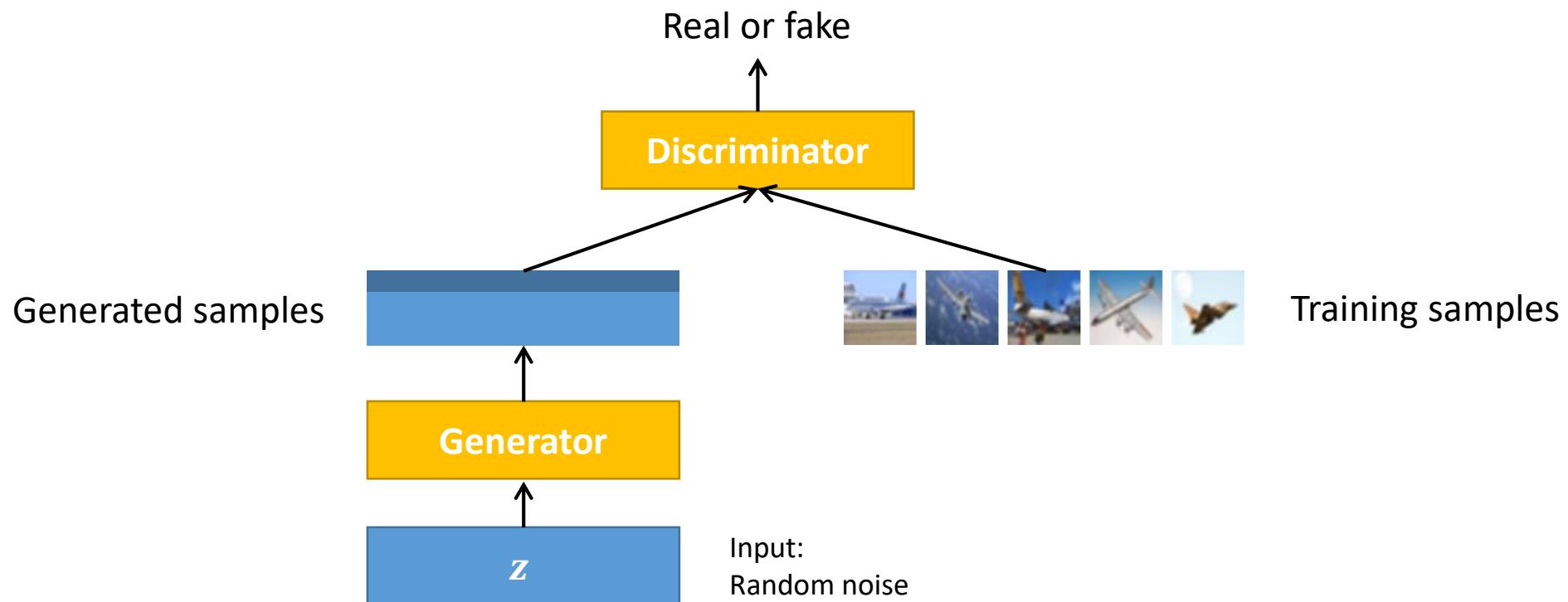
## Generative Adversarial Networks (GANs)

# Training GANs: Intuition

General idea is to have two networks play a two-player game

Generator: try to fool the discriminator by generating real-looking images

Discriminator: try to distinguish between real and fake images



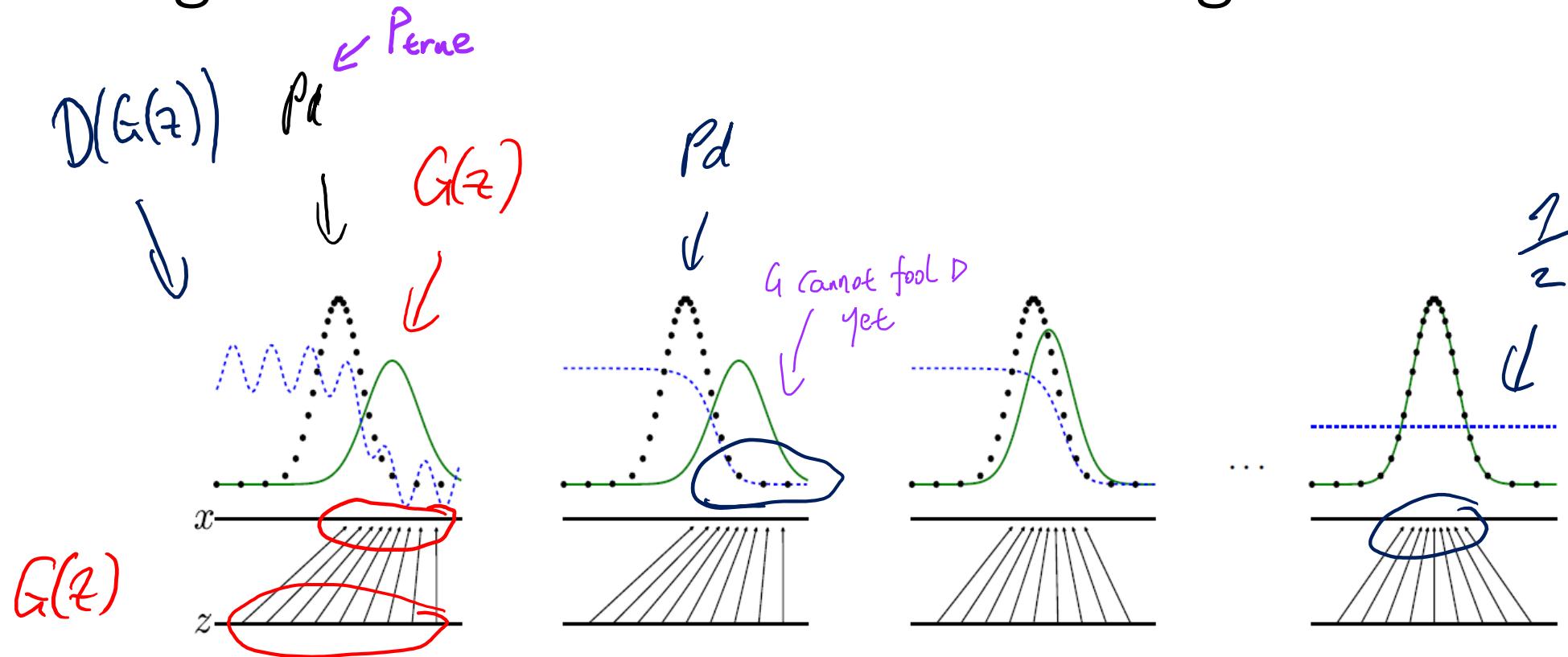
# Training GANs: Objective – Part II

$$\arg \min_G \max_D V(G, D) \text{ , where}$$

$$V(G, D) = \mathbb{E}_{x \sim p_d(x)}[\log(D(x))] + \mathbb{E}_{x \sim p_g(x)}[\log(1 - D(x))]$$

$G(z)$

# Training – Intuition and Idealized Progression



# Training in practice - issue

Objective:

$$\min_{\Theta_g} \max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D_{\Theta_d}(G_{\Theta_g}(z)))]$$

Alternate between:

1. Gradient *ascent* on  $D$ :

$$\max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D_{\Theta_d}(G_{\Theta_g}(z)))]$$

2. Gradient *descent* on  $G$ :

真假图片的数量要相等

$$\min_{\Theta_g} \mathbb{E}_{z \sim p_z(z)} [\log (1 - D_{\Theta_d}(G_{\Theta_g}(z)))]$$



# Training algorithm (pseudo-code)

---

**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:

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

**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**

[Ian Goodfellow et al., “Generative Adversarial Nets”, NIPS 2014]

# Training in practice - issue

## Objective:

$$\min_{\Theta_g} \max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} \left[ \log \left( 1 - D_{\Theta_d} (G_{\Theta_g}(z)) \right) \right]$$

in this region, generator produces bad images;  
and the same time, the gradient is very small, we might have a bit vanishing problem, the generator learns too slow.

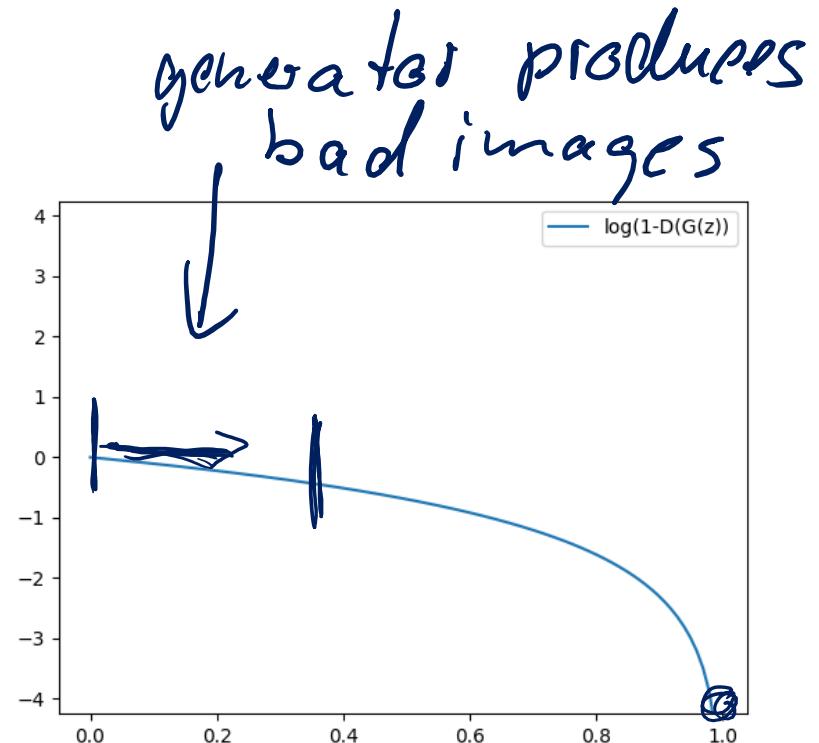
Alternate between:

1. Gradient *ascent* on  $D$ :

$$\max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} \left[ \log \left( 1 - D_{\Theta_d} (G_{\Theta_g}(z)) \right) \right]$$

2. Gradient *descent* on  $G$ :

$$\min_{\Theta_g} \mathbb{E}_{z \sim p_z(z)} \left[ \log \left( 1 - D_{\Theta_d} (G_{\Theta_g}(z)) \right) \right]$$



# Training in practice - fix

Objective:

$$\min_{\Theta_g} \max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} \left[ \log \left( 1 - D_{\Theta_d}(G_{\Theta_g}(z)) \right) \right]$$

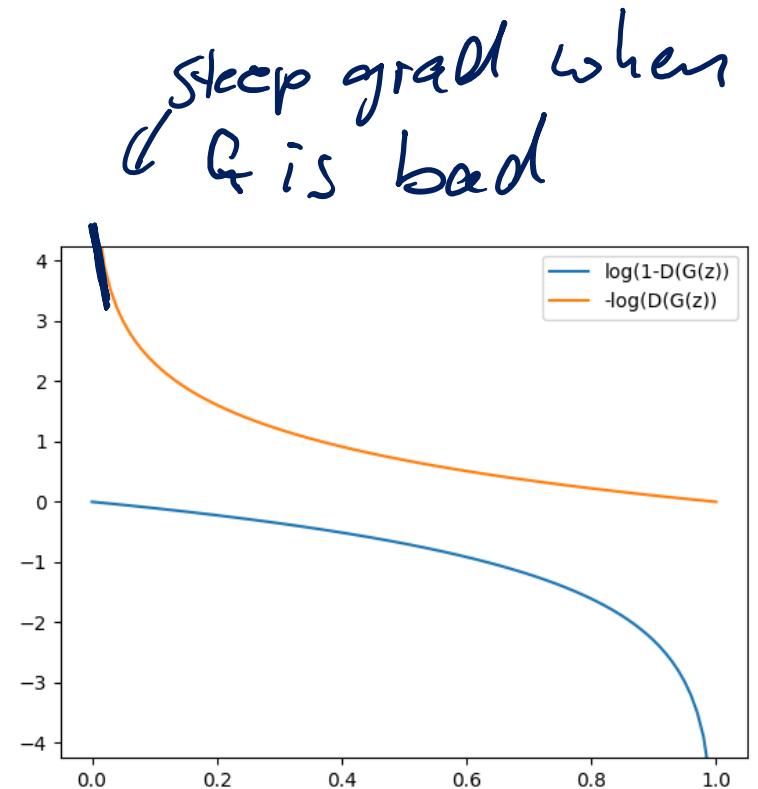
Alternate between:

1. Gradient *ascent* on D

$$\max_{\Theta_d} \mathbb{E}_{x \sim p_d(x)} [\log D_{\Theta_d}(x)] + \mathbb{E}_{z \sim p_z(z)} \left[ \log \left( 1 - D_{\Theta_d}(G_{\Theta_g}(z)) \right) \right]$$

2. Instead: gradient *ascent* on G:

$$\max_{\Theta_g} \mathbb{E}_{z \sim p_z(z)} \left[ \log \left( D_{\Theta_d}(G_{\Theta_g}(z)) \right) \right]$$



# Training algorithm (pseudo-code)

---

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

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

**end for**

- Sample minibatch of  $m$  noise samples  $\{z^{(1)}, \dots, z^{(m)}\}$  from noise prior  $p_g(z)$ .

- Update the generator by ascending its stochastic gradient (improved objective):

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log(D_{\theta_d}(G_{\theta_g}(z^{(i)})))$$

before: descending

**end for**

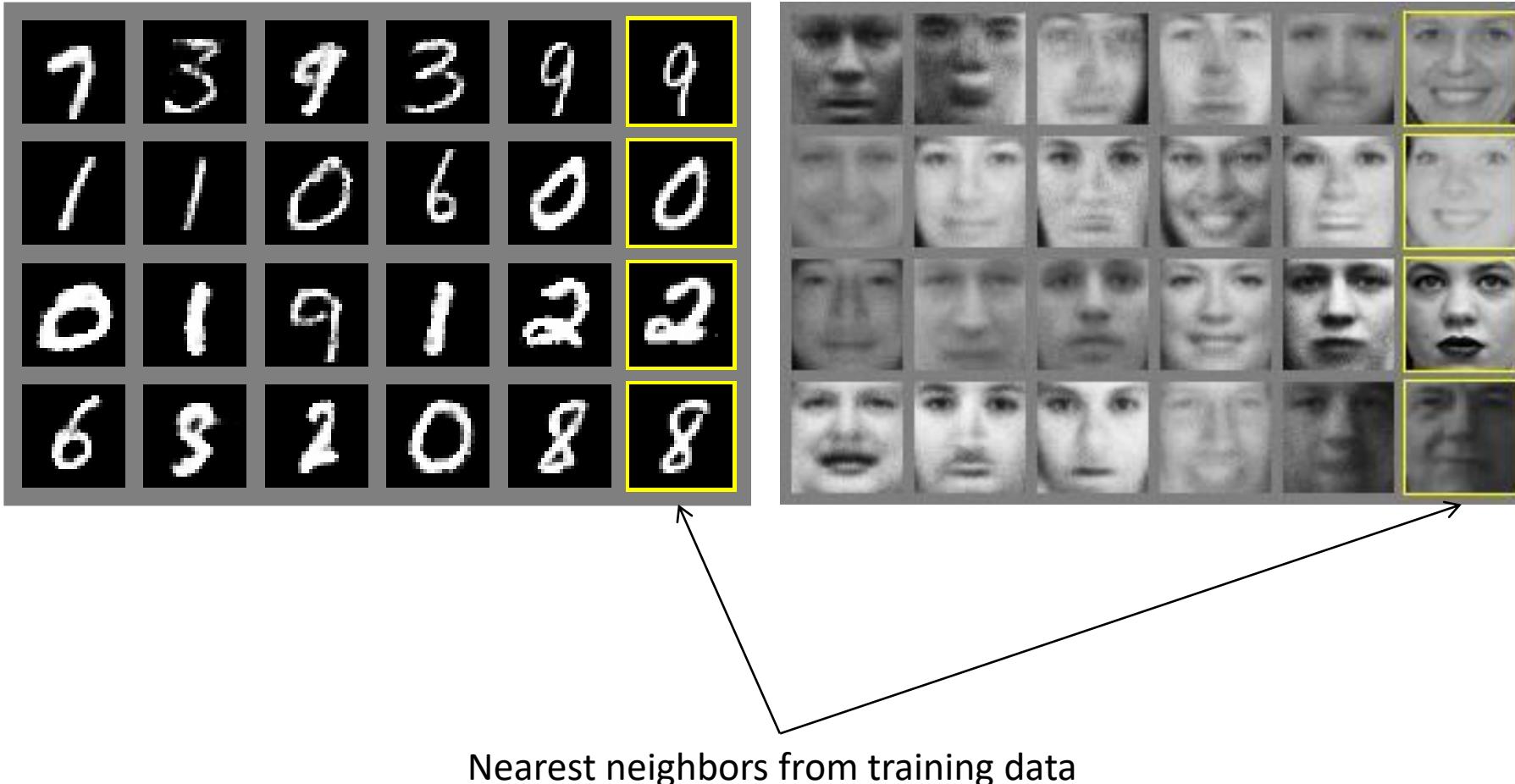
Update the generator by descending its stochastic gradient:

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

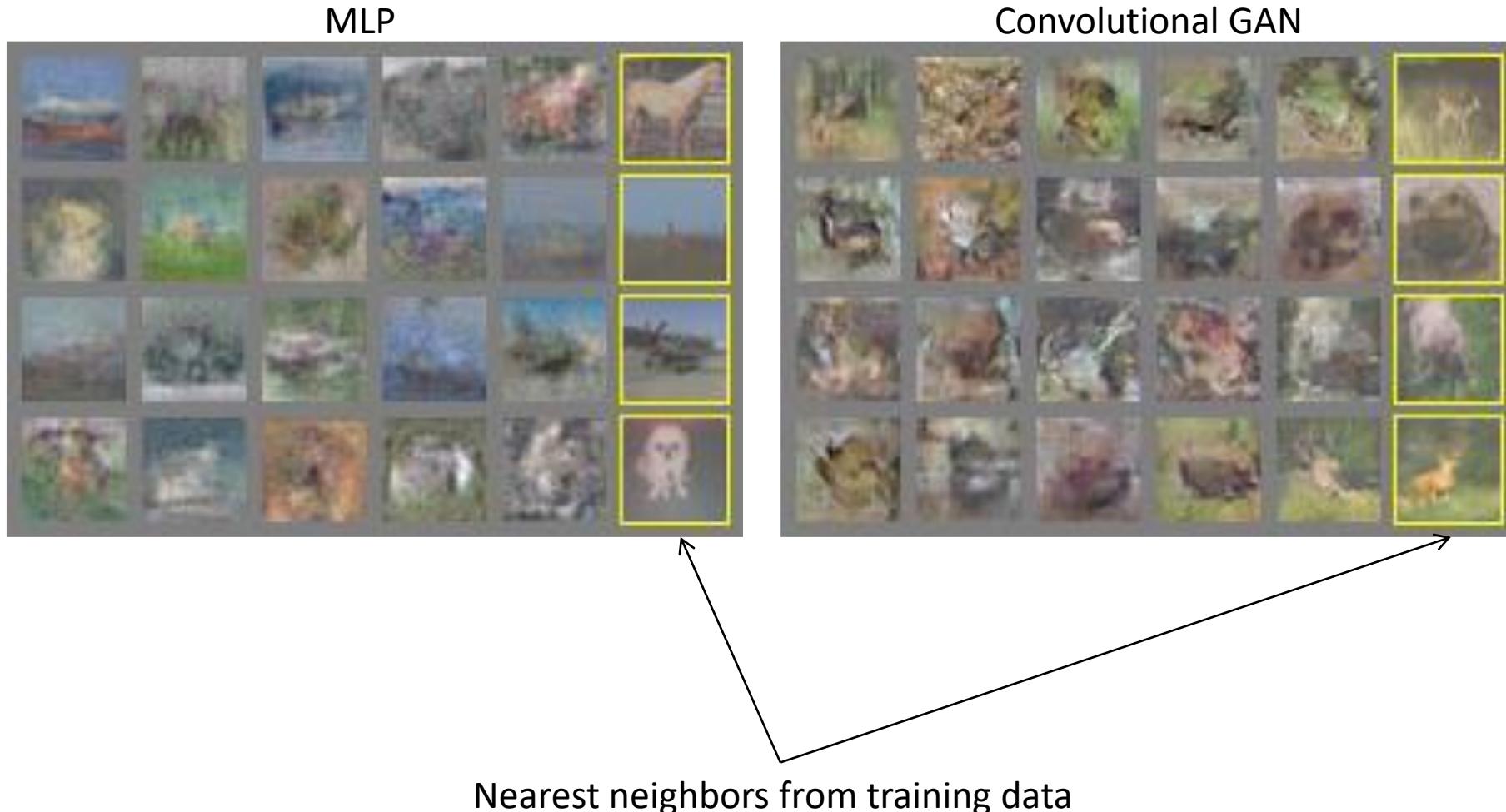


[Ian Goodfellow et al., “Generative Adversarial Nets”, NIPS 2014]

# GANs: Results MNIST & TFD



# GANs: Results CIFAR-10



# Issues

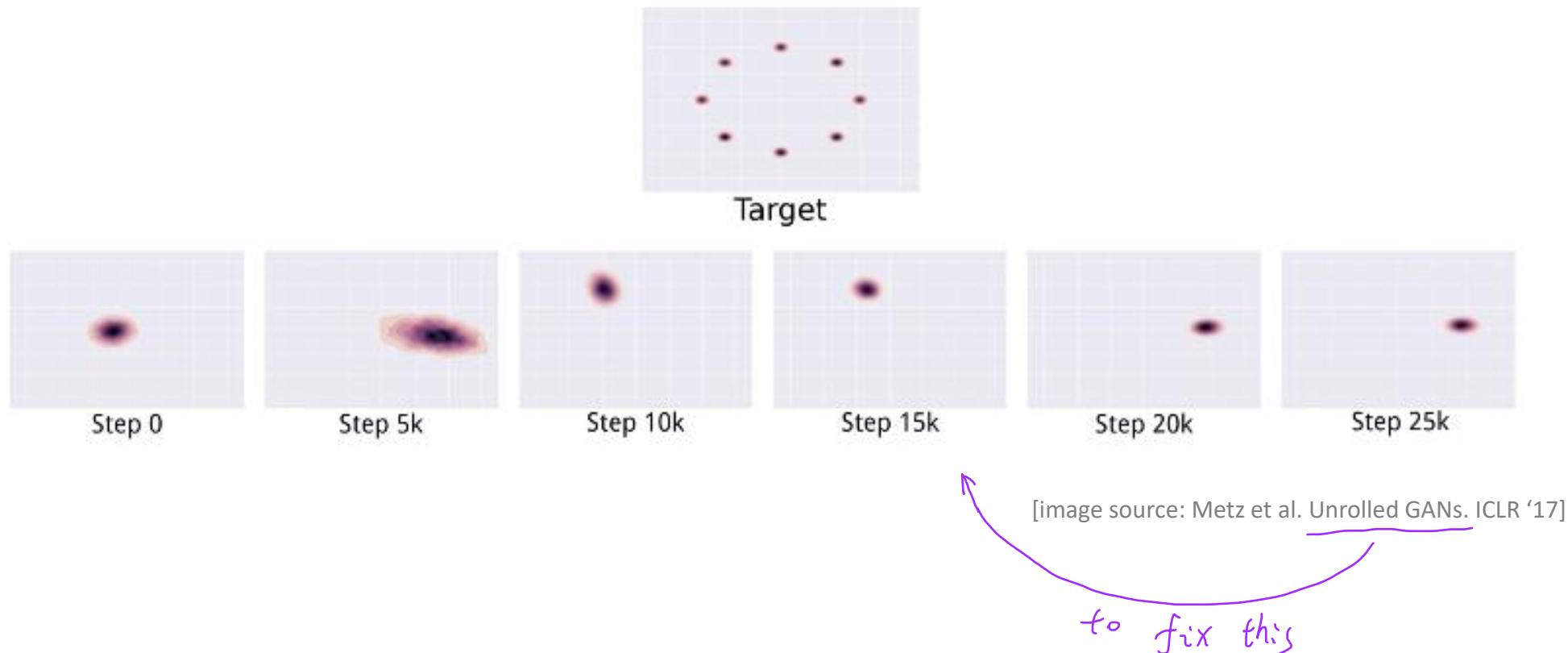
- **Difficult to train:**

Need to find Nash-Equilibrium in two-player game. Making downhill progress for one player may move the other player uphill.

Generator can trick discriminator due to underlying image statistics

# Issues

- Difficult to train:  
Mode collapse (no sample diversity): saddle point in dual energy landscape



# Comparison with VAE

- Implicit pdf in VAE with intractable data likelihood
- No variational bound is needed (*In GAN*)
- GANs only care about samples
- Sharper images
- (Specific model families usable within GANs are already known to be universal approximators, so GANs are already known to be asymptotically consistent. Some VAEs are conjectured to be asymptotically consistent, but this is not yet proven)

# Issues with Jensen–Shannon divergence => Wasserstein GAN

JSD correlates badly with sample quality → do not know when to stop training

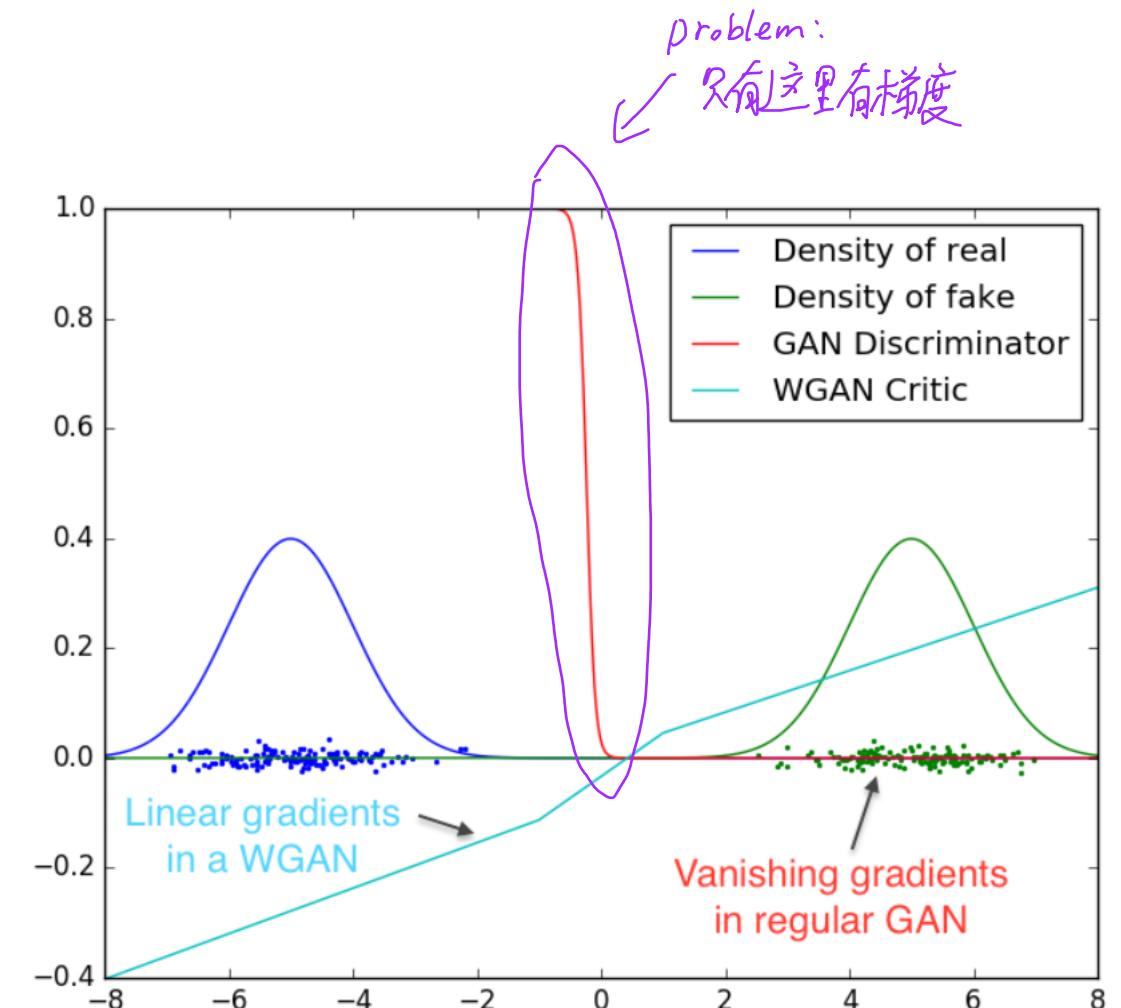
$D$  tends to be optimal, JSD saturates → bad gradients

One solution: Use Wasserstein distance

earth mover's distance

GAN objective can be generalized to entire family of divergences:

[Nowozin et al. NeurIPS '16. f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization]



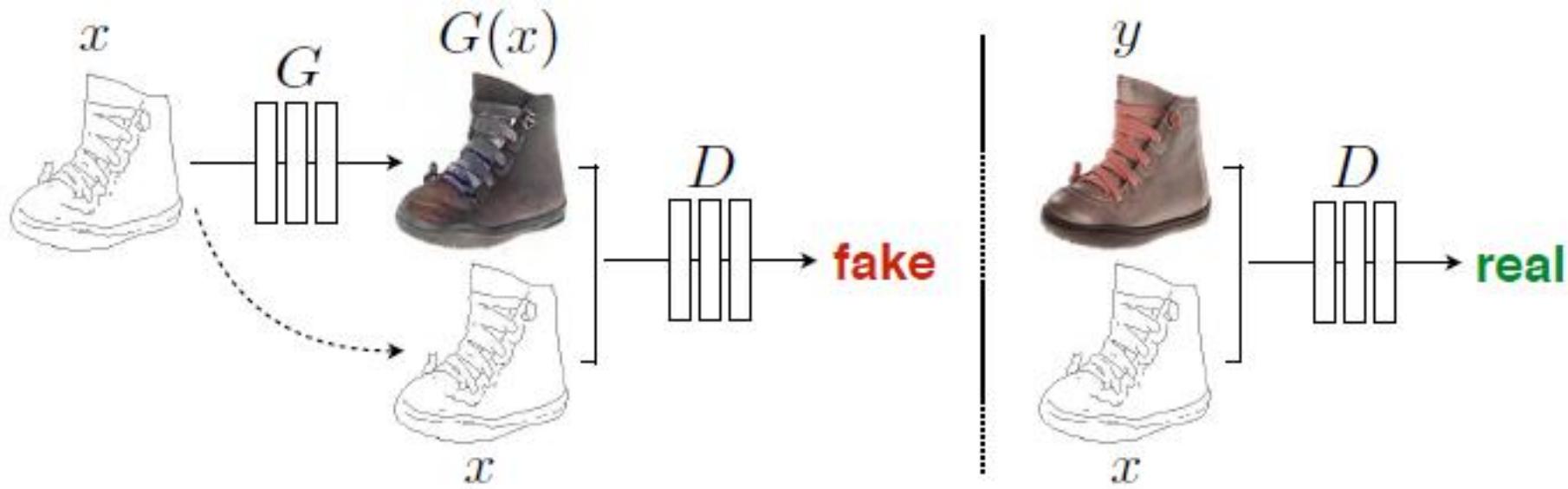
[image source: Arjovsky et al. Wasserstein GANs. '17]

# GAN Extensions & Applications

# Image-to-Image Translation with Conditional Adversarial Networks

[Isola et al. CVPR '17 – Pix2Pix]

pass sketch of img, translate it to real



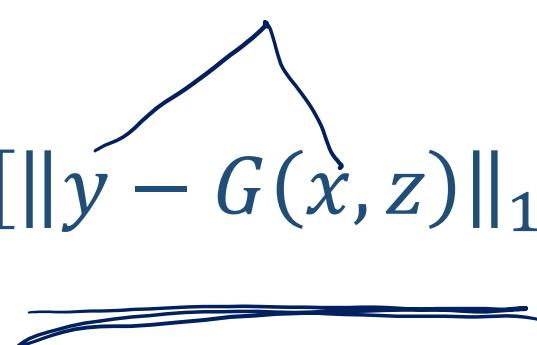
# Pix2Pix - Formulation

The Objective is:

$$\mathcal{L}(G, D) = \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G)$$

almost regular gan loss + reconstruction loss

where  $\mathcal{L}_{cGAN}$  is a conditional GAN objective:

$$\mathcal{L}_{cGAN}(G, D) = \mathbb{E}_{x,y}[\log D(x, y)] + \mathbb{E}_{x,z}[1 - \log D(x, G(x, z))]$$


and  $\mathcal{L}_{L1}$ :

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

# Pix2Pix - results

[Isola et al. CVPR '17 – Pix2Pix]

## Labels to Street Scene

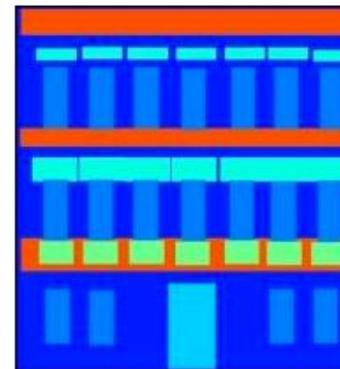


## input



## output

## Labels to Facade



input



outpu

## BW to Color



## input



## output

## Aerial to Map



## input

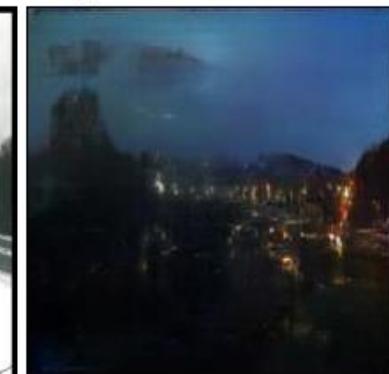


## output

## Day to Night



## input



## output

## Edges to Photo



## input



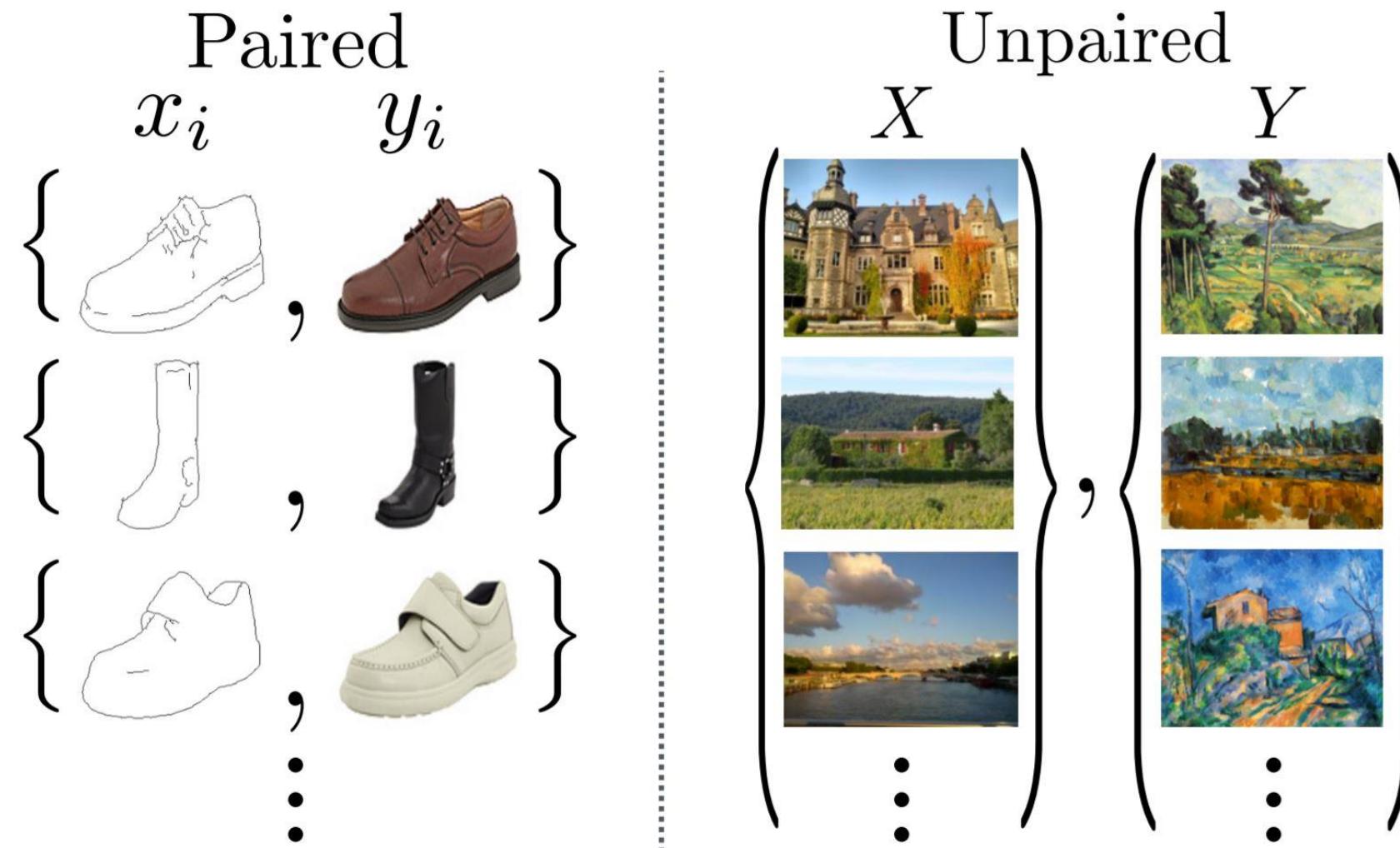
## output

# Pix2Pix - Limitations

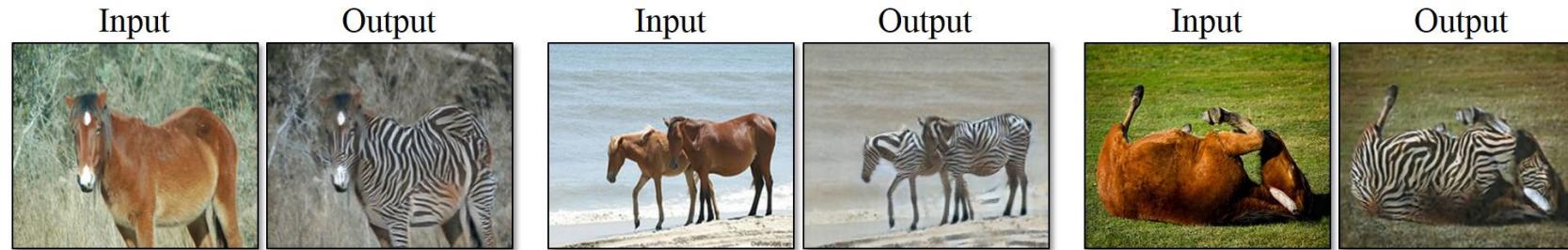
One obvious limitation is that this approach requires paired images as training data.

In many cases, one has at its disposal examples from two densities and wants to translate a sample from the first (“images of apples”) into a sample likely under the second (“images of oranges”).

# CycleGAN – unpaired image translation



# CycleGAN – unpaired image translation



horse → zebra



zebra → horse

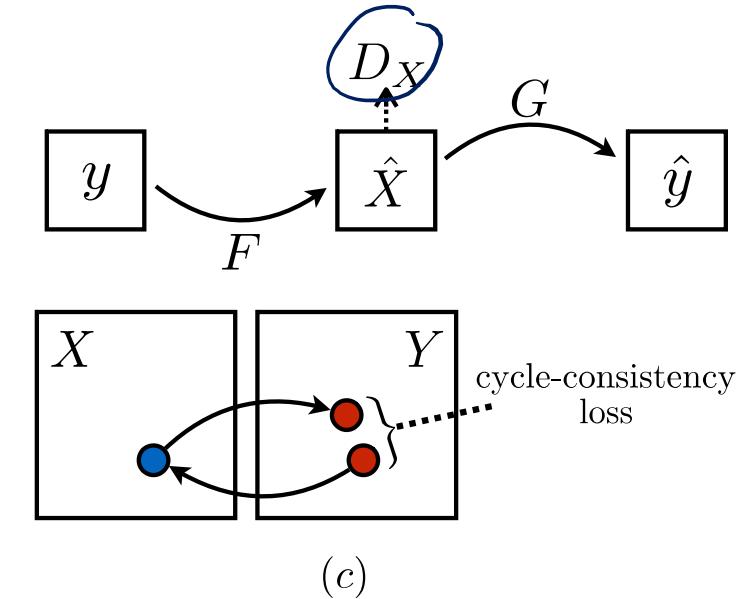
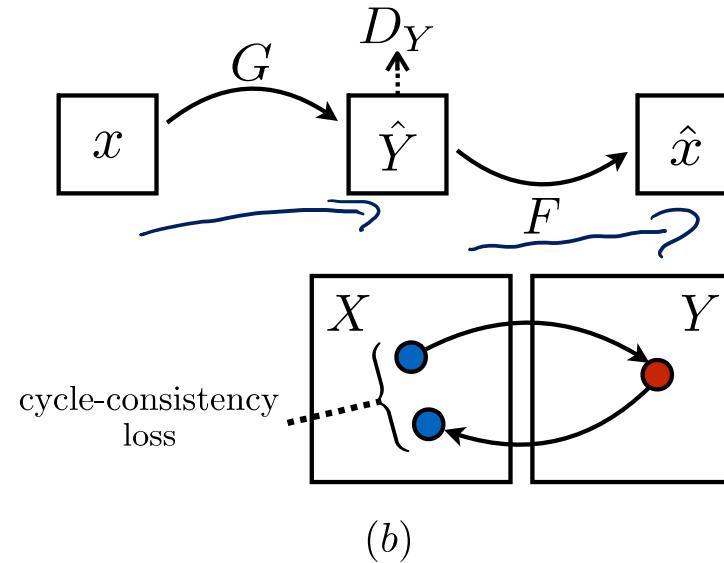
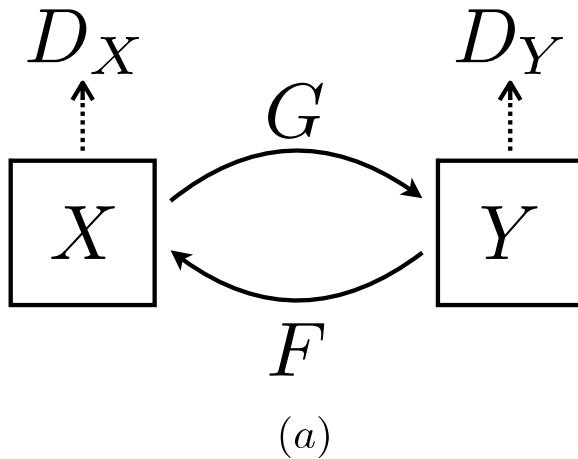


apple → orange



orange → apple

# CycleGAN - Overview



# CycleGAN - Formulation

The Objective is:

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

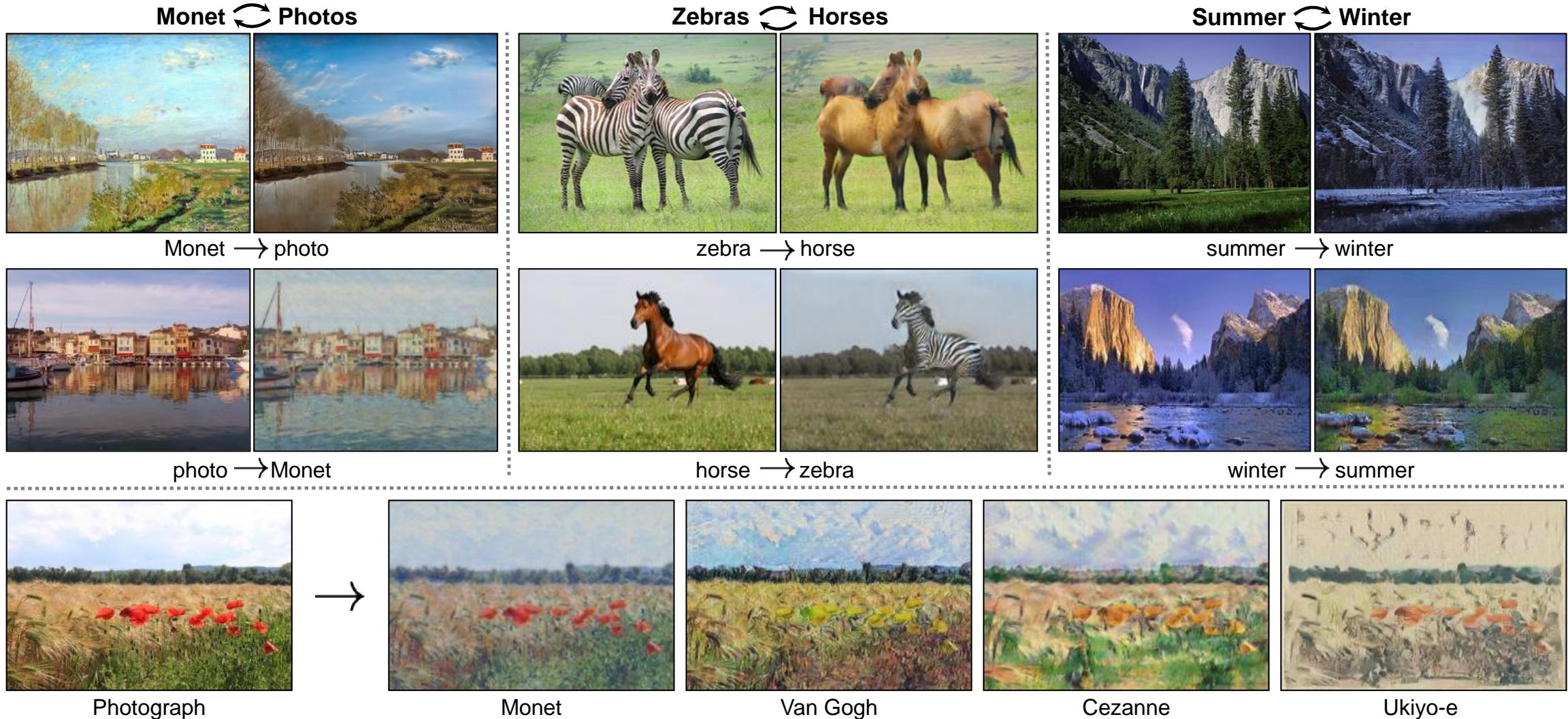
where  $\mathcal{L}_{GAN}$  is the standard GAN loss as previously and  $\mathcal{L}_{cyc}$  is:

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

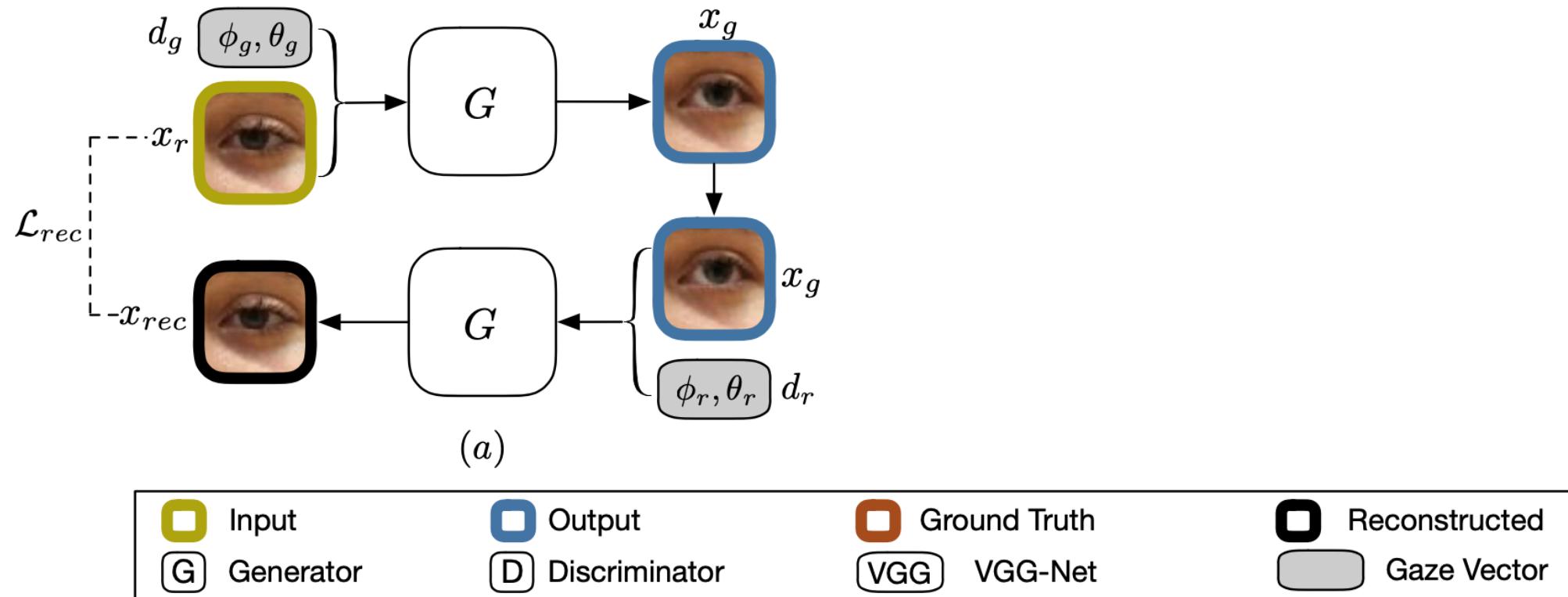
# CycleGAN – Effect of cycle consistency



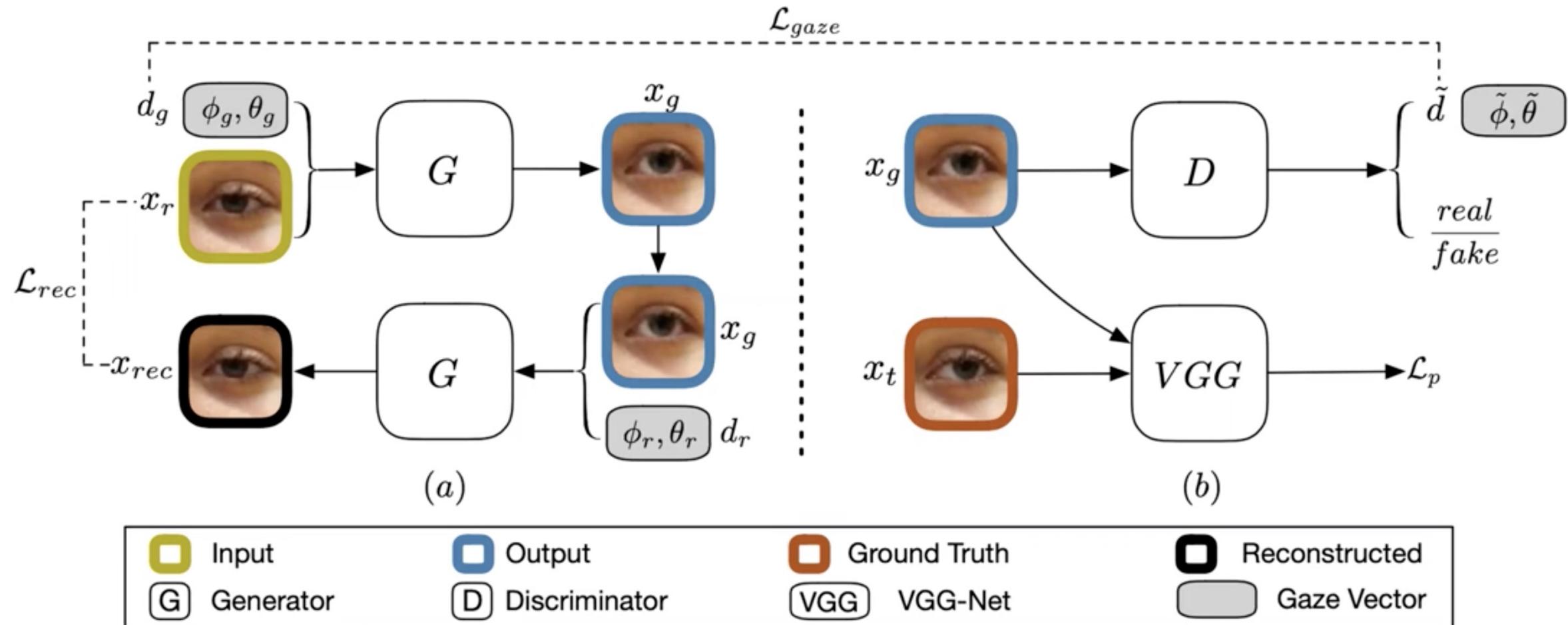
# CycleGAN – More results



# Gaze redirection



# Gaze redirection

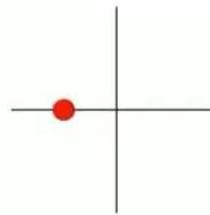


# Demo

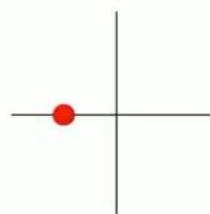
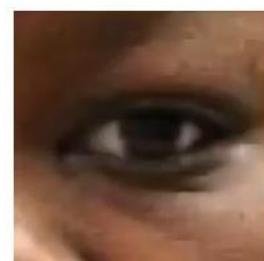
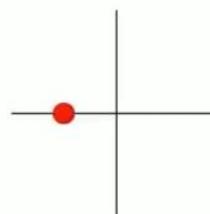
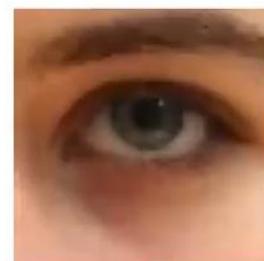
Input



Target gaze movement

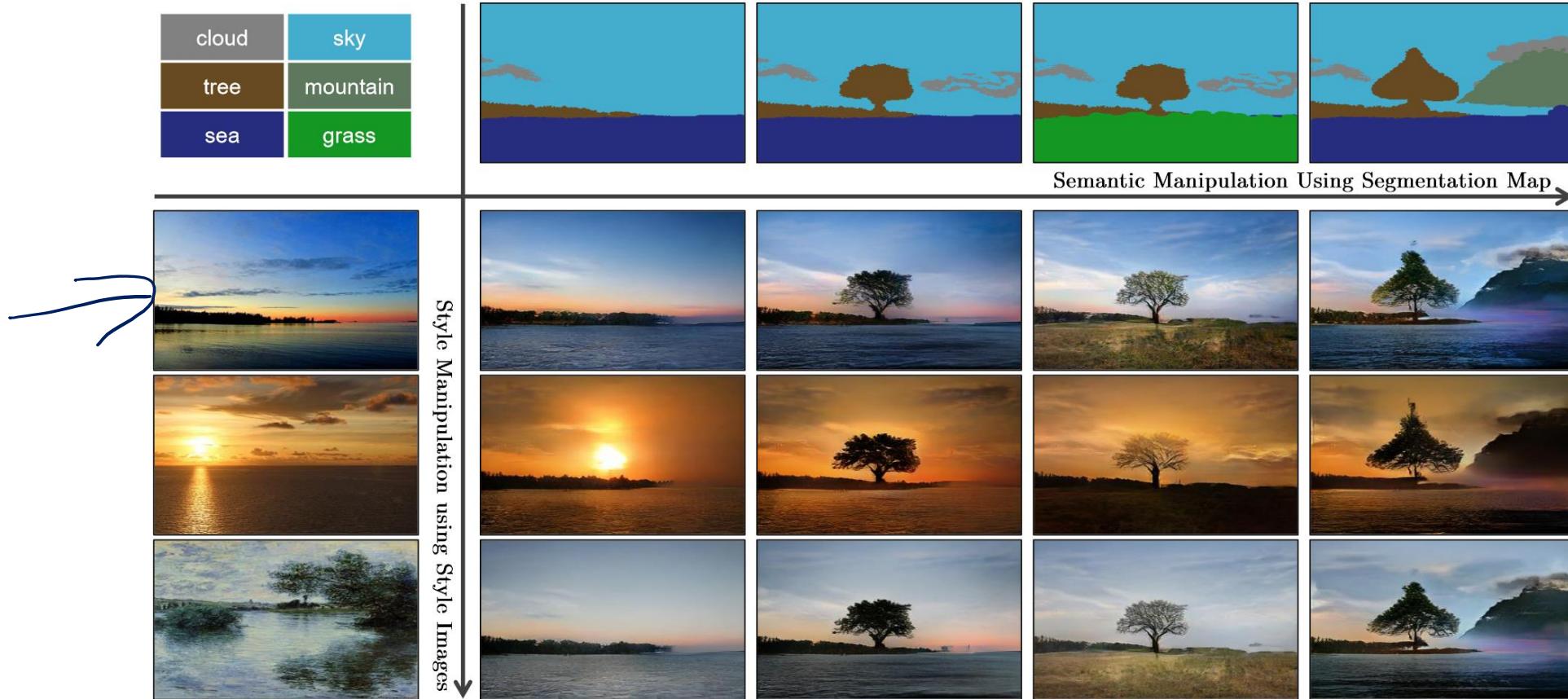


Output



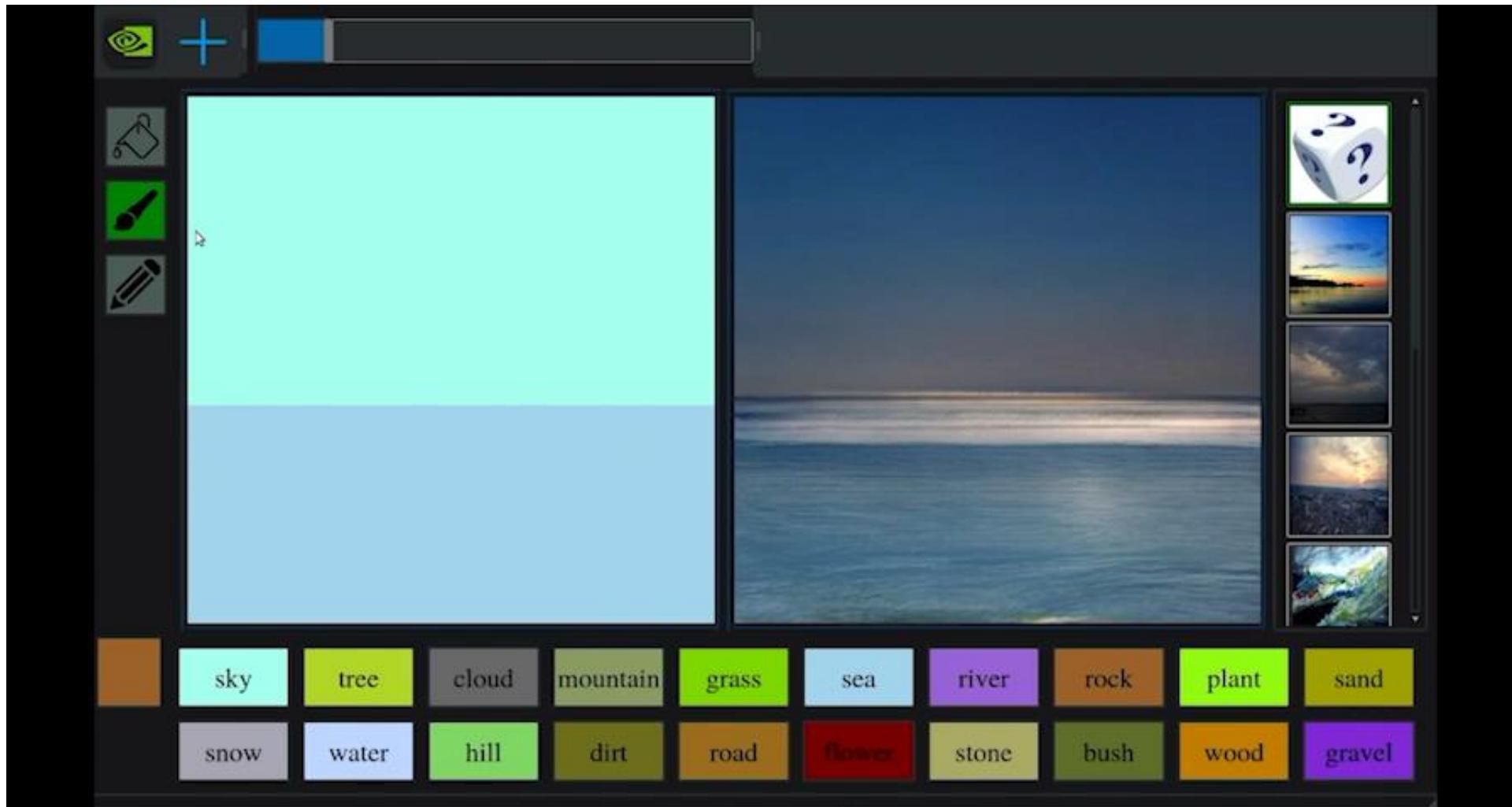
# GauGAN - Semantic Image Synthesis with Spatially-Adaptive Normalization

Task: Given semantic segmentation and reference style, synthesize image



[Park et al. CVPR '19 - Semantic image synthesis with spatially-adaptive normalization]

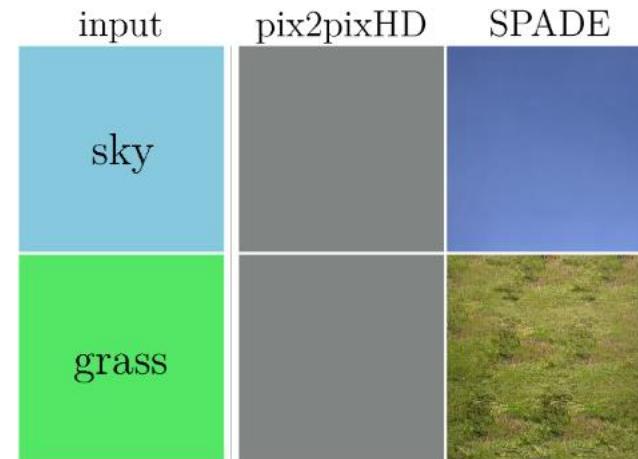
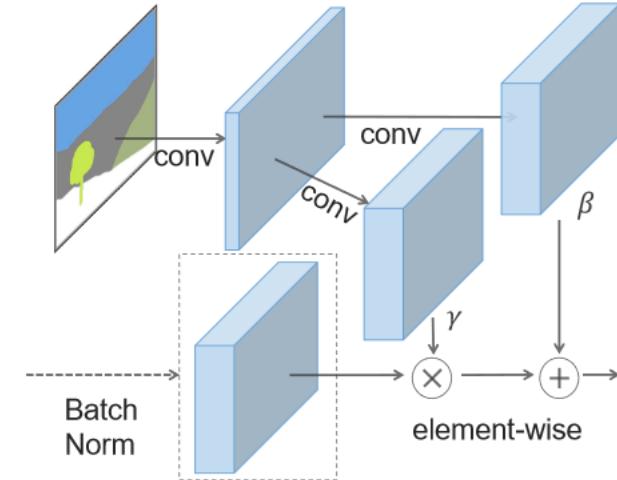
# GauGAN - Semantic Image Synthesis with Spatially-Adaptive Normalization



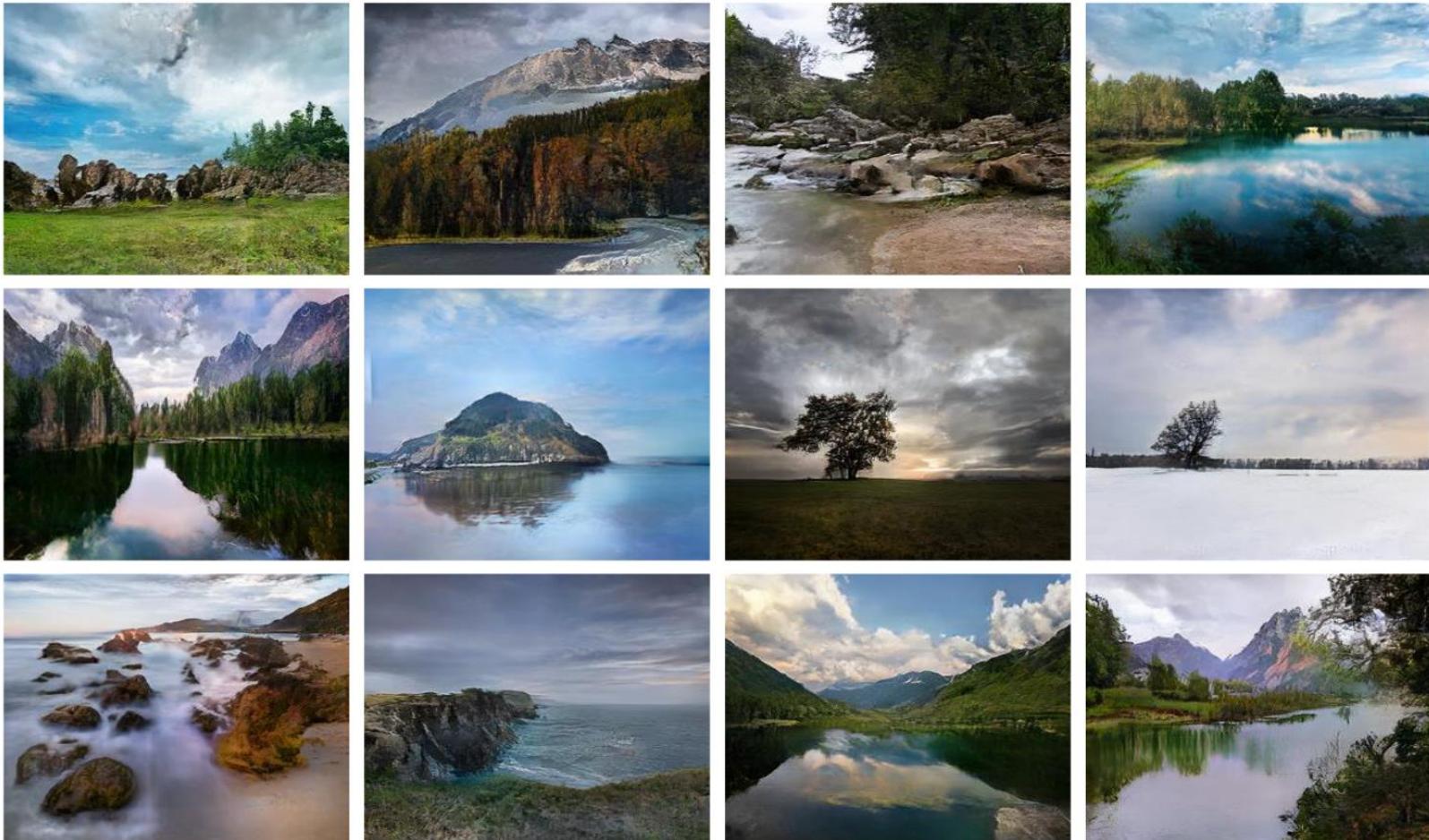
[Park et al. CVPR '19 - Semantic image synthesis with spatially-adaptive normalization]

# Problem with pix2pix:

- Unconditional normalization layer in generator “washes away” information of semantic labels.
- Proposal: Spatially-varying conditional normalization

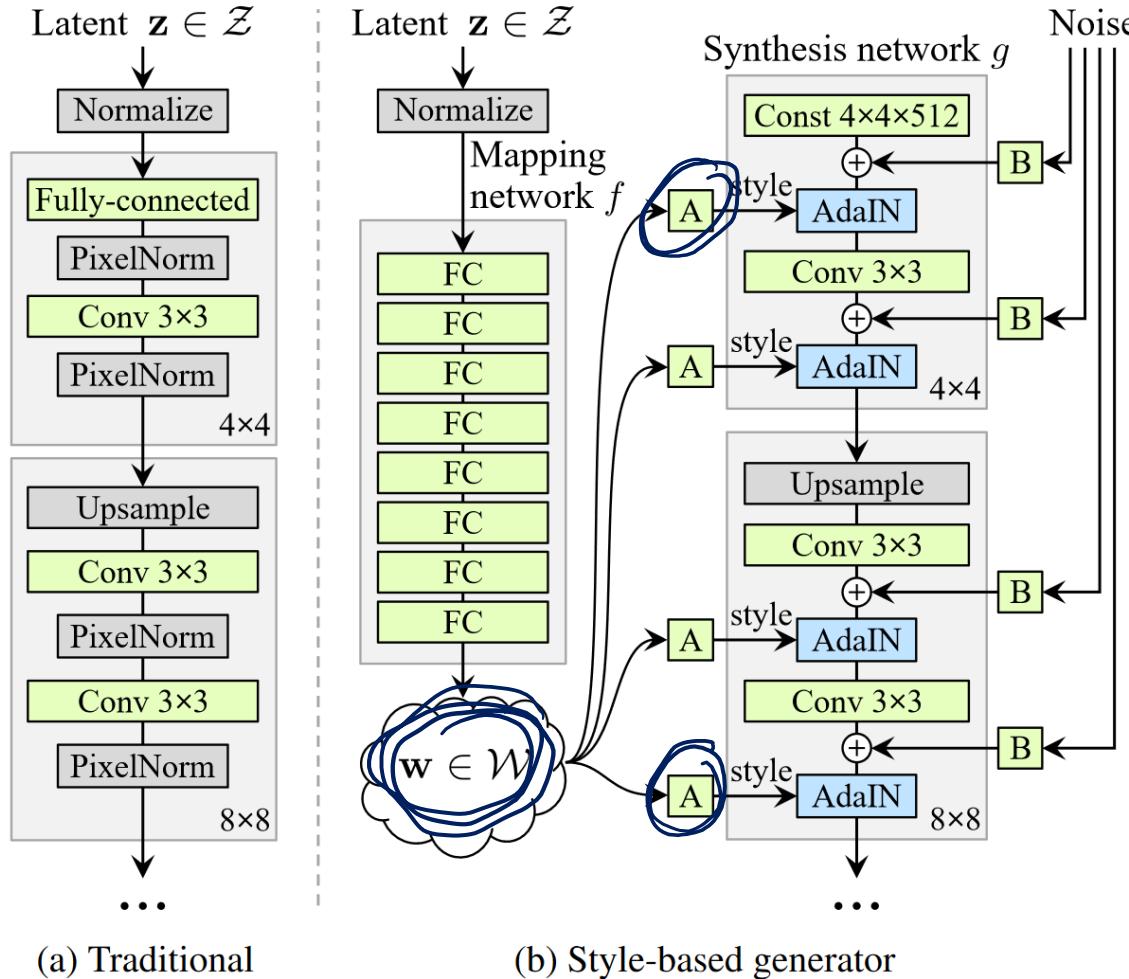


# Results on Flickr Landscapes dataset



Try yourself: <http://nvidia-research-mingyuliu.com/gaugan>

# A Style-Based Generator Architecture for Generative Adversarial Networks



# High-quality image generation



# Granular style control



# StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks

This bird is white with some black on its head and wings, and has a long orange beak



This bird has a yellow belly and tarsus, grey back, wings, and brown throat, nape with a black face

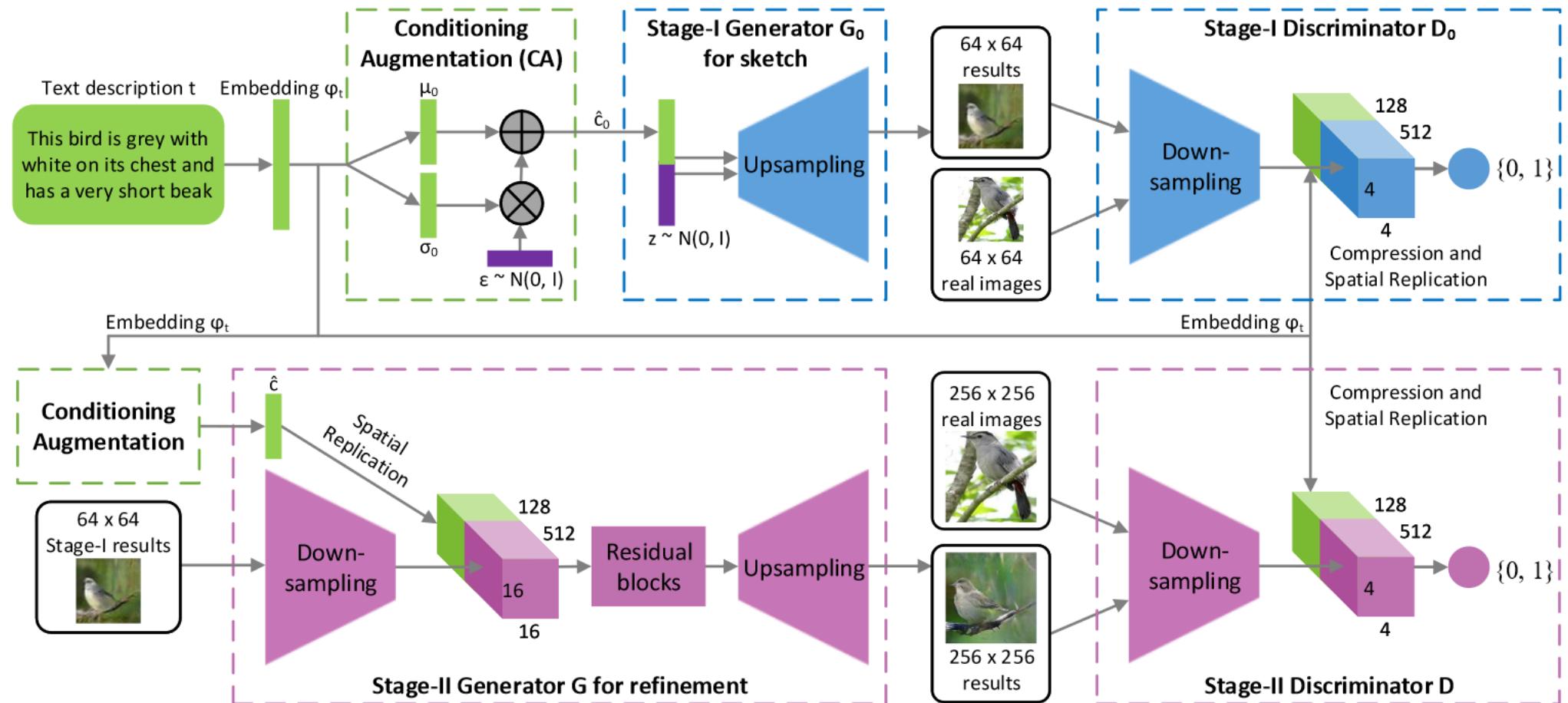


This flower has overlapping pink pointed petals surrounding a ring of short yellow filaments



[Zhang et al. ICCV '17 - StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks]

# Architecture



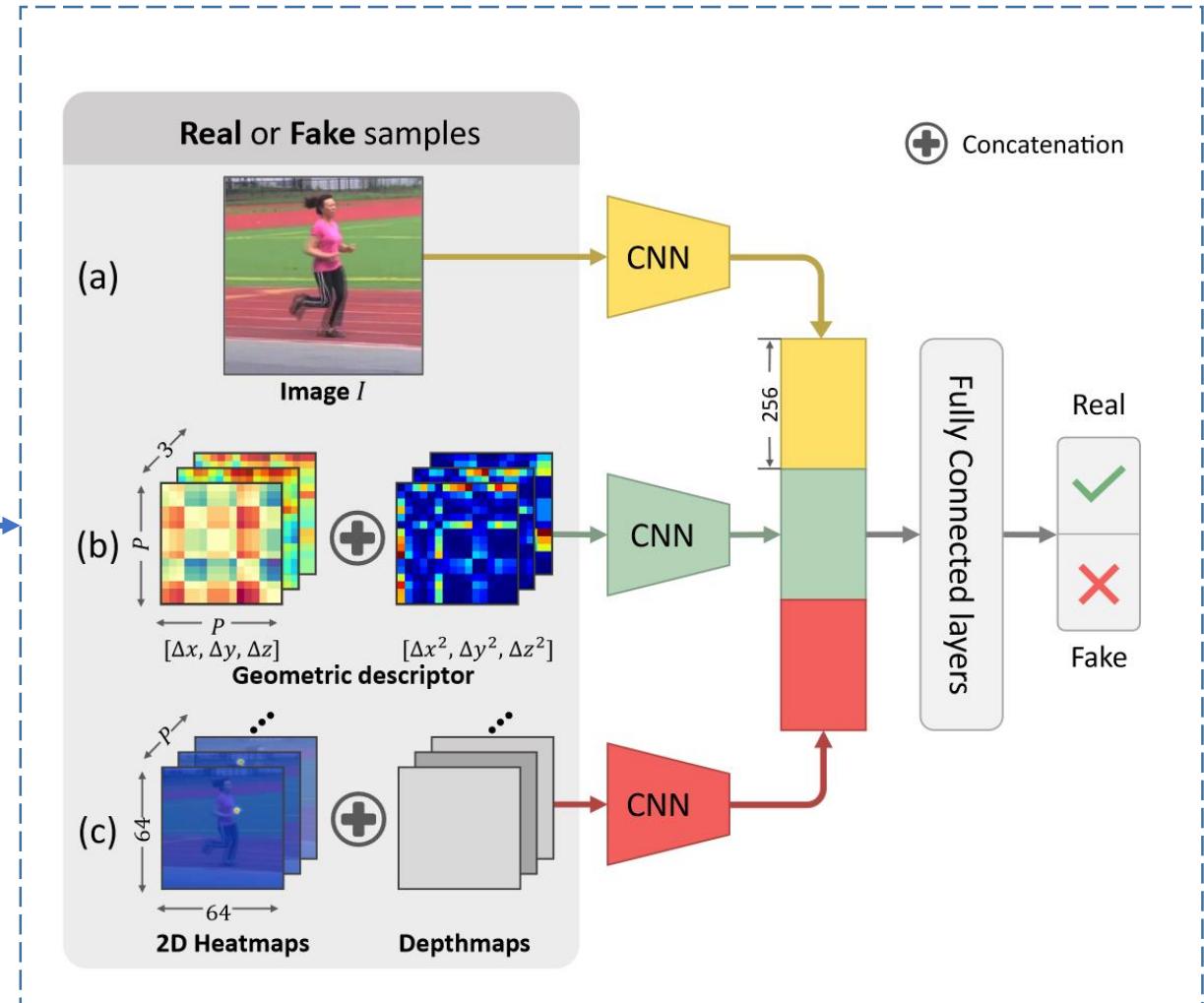
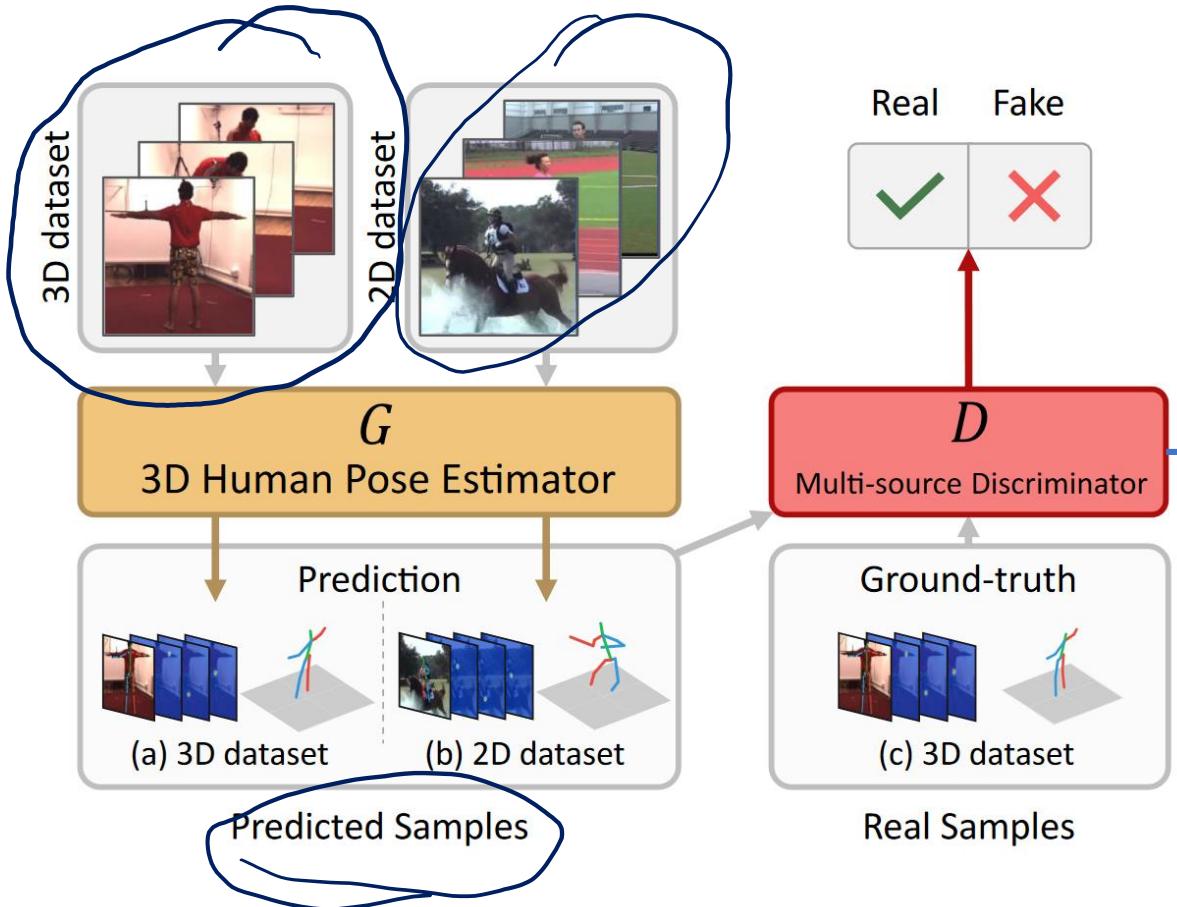
[Zhang et al. ICCV '17 - StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks]

# More qualitative examples

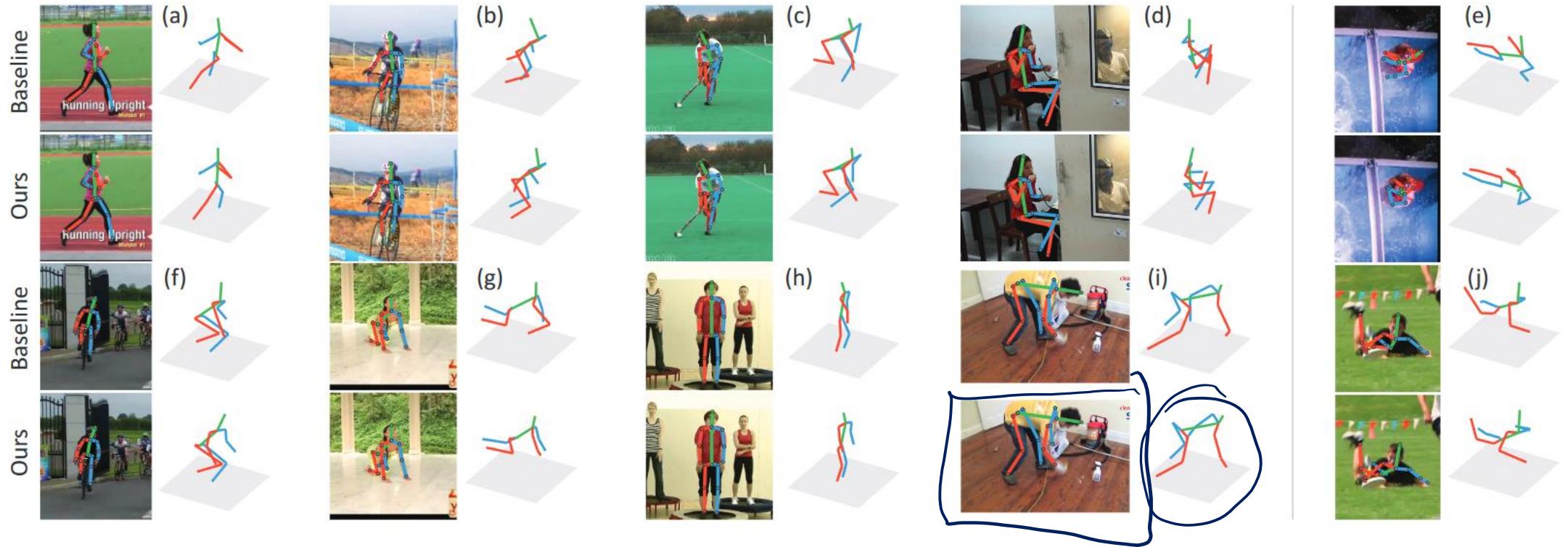
Text description	This bird is blue with white and has a very short beak	This bird has wings that are brown and has a yellow belly	A white bird with a black crown and yellow beak	This bird is white, black, and brown in color, with a brown beak	The bird has small beak, with reddish brown crown and gray belly	This is a small, black bird with a white breast and white on the wingbars.	This bird is white black and yellow in color, with a short black beak	
Stage-I images								
Stage-II images								

[Zhang et al. ICCV '17 - StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks]

# 3D Human Pose Estimation in the Wild by Adversarial Learning



# Result of adversarial training



# Summary

GANs have [theoretical] ability to model distributions fully

Often generating consistent samples is enough

State-of-the-art results on many tasks such as image-to-image translation

Higher image quality compared to VAEs

Key aspect is the “perceptual” loss emulated by the discriminator

# Next Week

GAN ++

Autoregressive Models