# Lecture 8
# Generative Adversarial Networks (GANs)

## M2 Data Science and AI

## Yaohui Wang

http://www-sop.inria.fr/members/Yaohui.Wang/

- Generative Adversarial Networks: *Yaohui Wang*

- DeepFake Detection: *Dr. Antitza Dantcheva*

- Labs (TP): *David Anghelone*

# Question: VAE ?

# Ian Goodfellow



**Generative Adversarial Networks [NIPS 2014]**

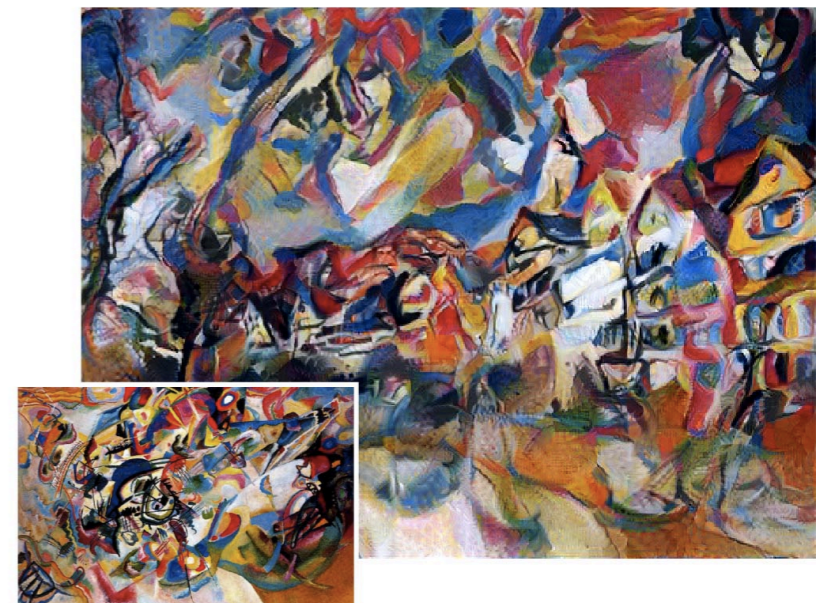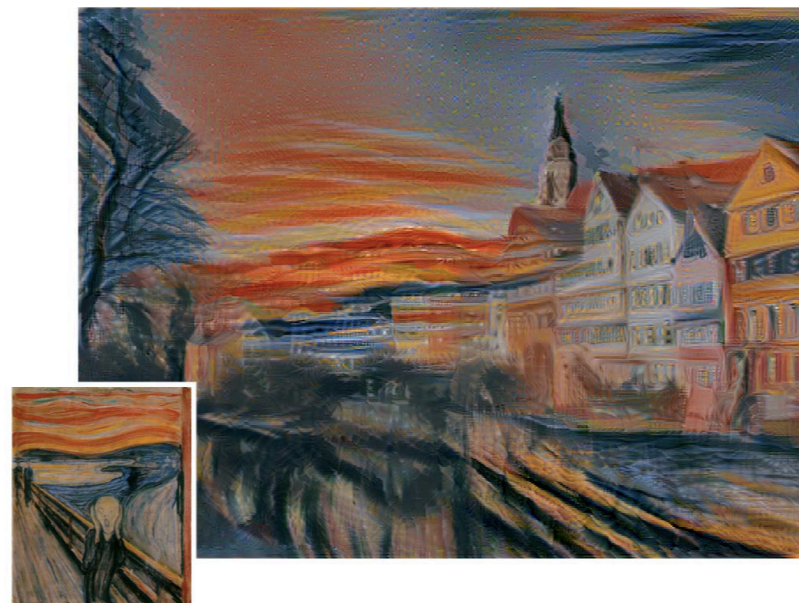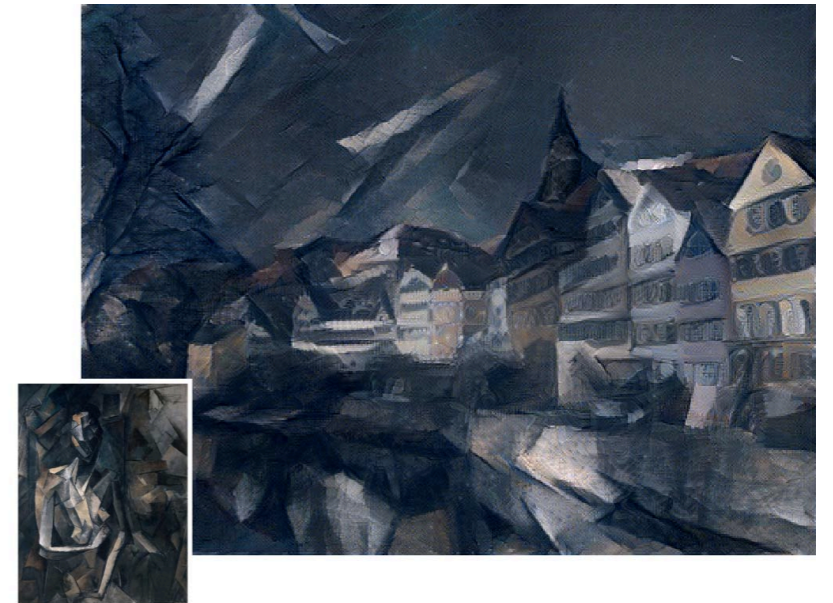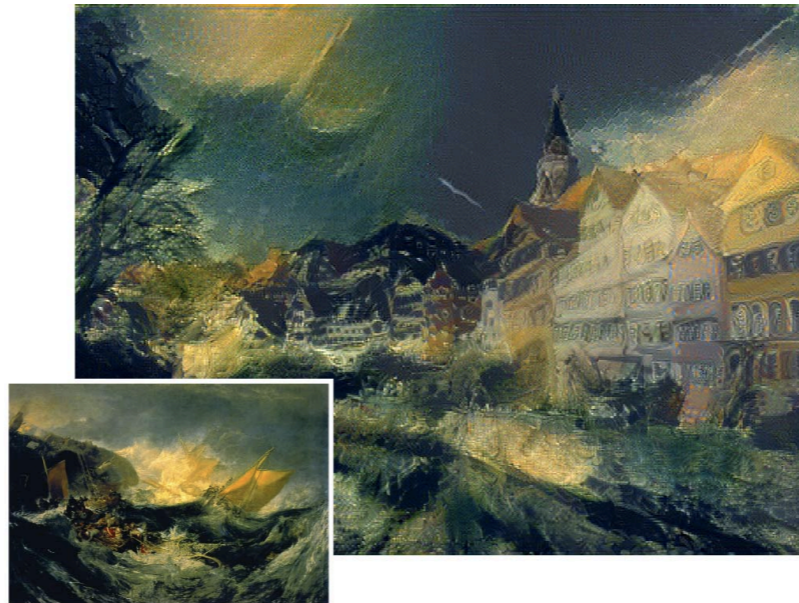"GANs are the most interesting idea in the last 10 years in ML"

- Yann LeCun

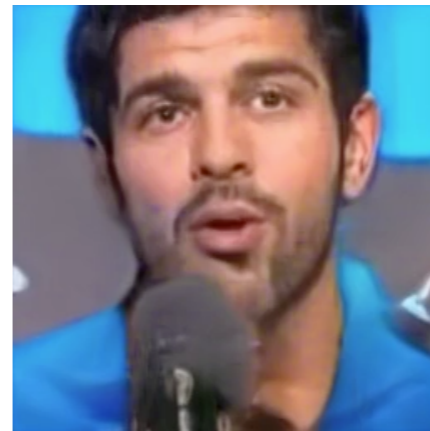# Image Generation

# Style Transfer

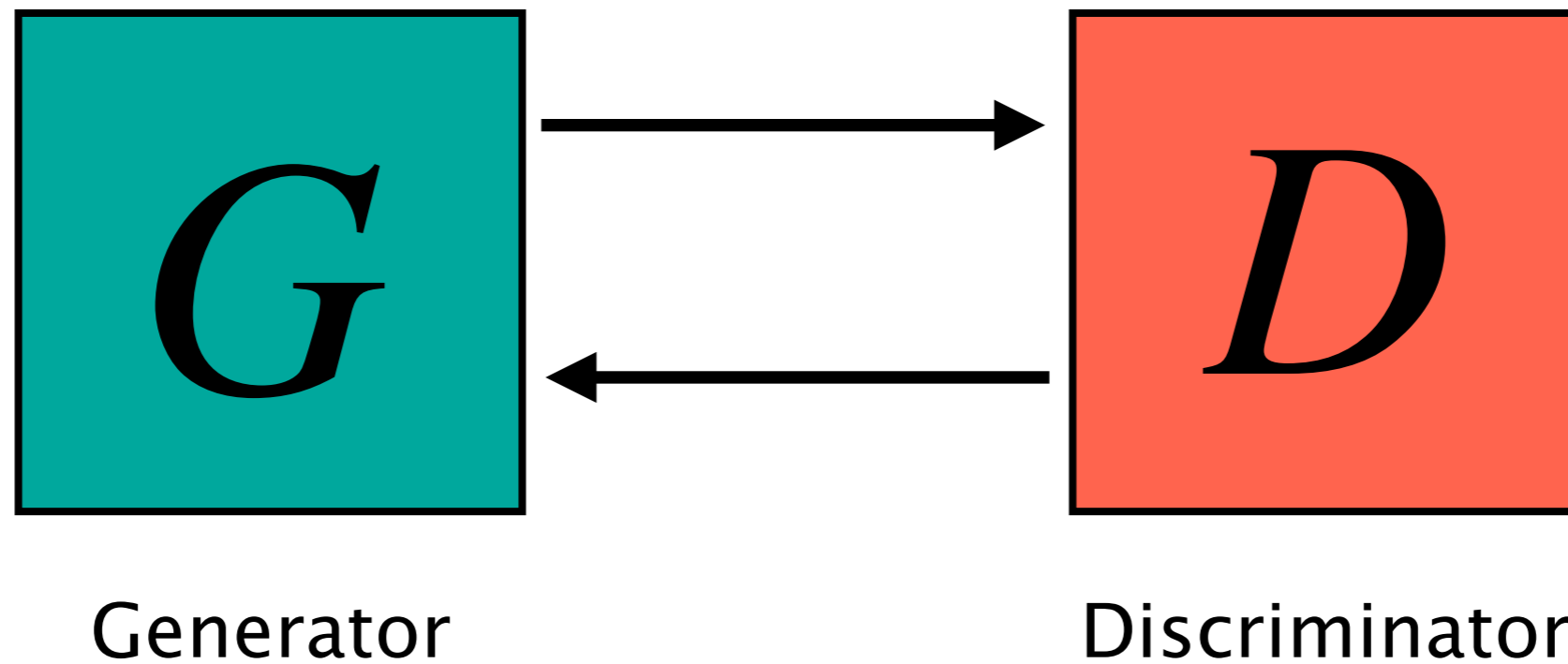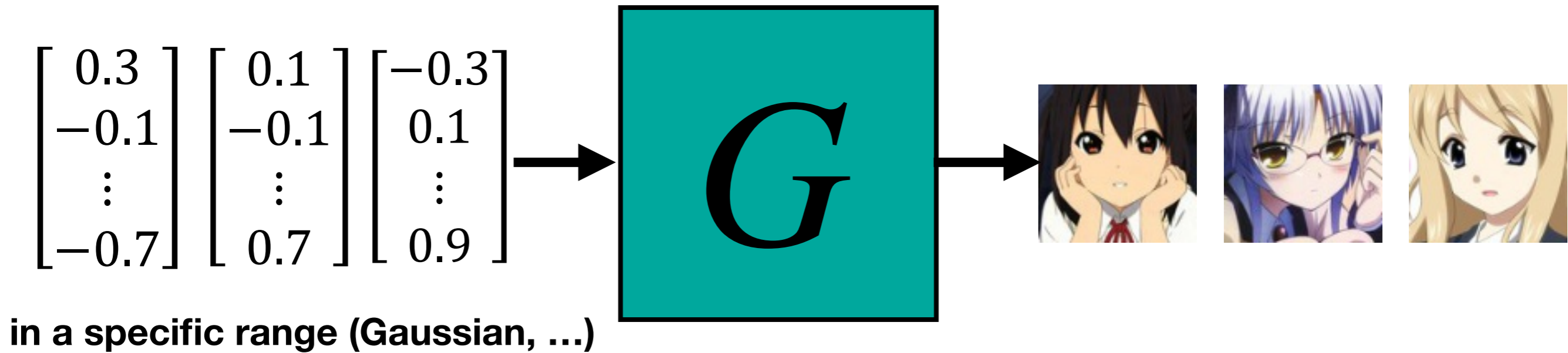# Video Generation

# Video Generation

# Outline

- Basic Idea of GAN

- Image Generation

  - Conditional GAN (CGAN, ACGAN)

  - Modern GANs (StyleGAN, BigGAN)

  - Image-to-image translation (Pix2Pix, CycleGAN)

- Video-to-video translation

- GANs Evaluation

- Video Generation

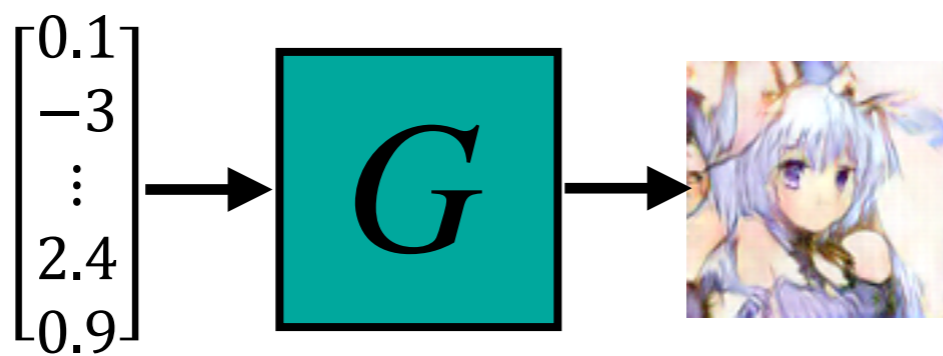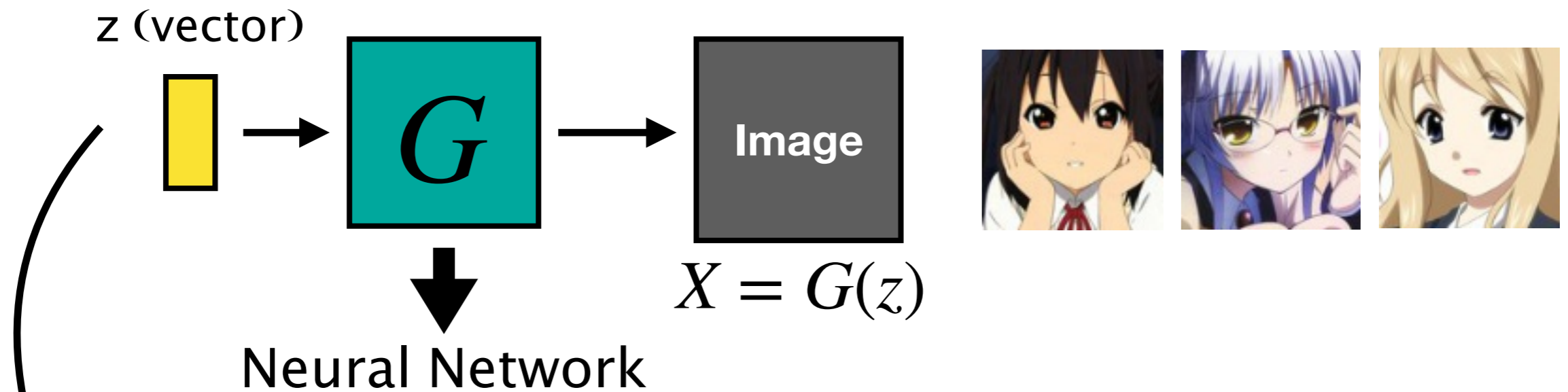- Lab (DCGAN for manga face generation)

# Basic idea of GAN

# Basic idea of GAN



G

Generator

D

Discriminator

# Basic idea of GAN

$$\begin{bmatrix} 0.3 \\ -0.1 \\ \vdots \\ -0.7 \end{bmatrix} \begin{bmatrix} 0.1 \\ -0.1 \\ \vdots \\ 0.7 \end{bmatrix} \begin{bmatrix} -0.3 \\ 0.1 \\ \vdots \\ 0.9 \end{bmatrix} \longrightarrow \boxed{G} \longrightarrow$$



**in a specific range (Gaussian, …)**

# Basic idea of GAN

z (vector)



$X = G(z)$

Neural Network

$$\begin{bmatrix} 0.1 \\ -3 \\ \vdots \\ 2.4 \\ 0.9 \end{bmatrix}$$

**Each dimension of input vector represents some characters**

$$\begin{bmatrix} 3 \\ -3 \\ \vdots \\ 2.4 \\ 0.9 \end{bmatrix}$$

**Longer hair**

$$\begin{bmatrix} 0.1 \\ 2.1 \\ \vdots \\ 5.4 \\ 0.9 \end{bmatrix}$$

**blue hair**

$$\begin{bmatrix} 0.1 \\ -3 \\ \vdots \\ 2.4 \\ -3.5 \end{bmatrix}$$

**open mouth**

14

# Basic idea of GAN



**Image** → **D** → scalar

$y = D(X)$

↓ Neural Network

higher value: more realistic
lower value: less realistic

 → **D** → **1.0**

 → **D** → **1.0**

 → **D** → **0.1**

 → **D** → **0.1**

# Basic idea of GAN

Adversarial Training (Generative <span style="color:red">Adversarial</span> Networks)



epoch 0

$z \to G \to$ (images) $\to D$

update $\theta_G$    update $\theta_D$

$z \to G \to$ (images) $\to D \leftarrow$ (images)

update $\theta_G$    update $\theta_D$

$z \to G \to$ (images) $\to D$

update $\theta_G$    update $\theta_D$

# Basic idea of GAN

## Adversarial Training (Generative Adversarial Networks)

**Algorithm**    Initialize $\theta_d$ for D and $\theta_g$ for G
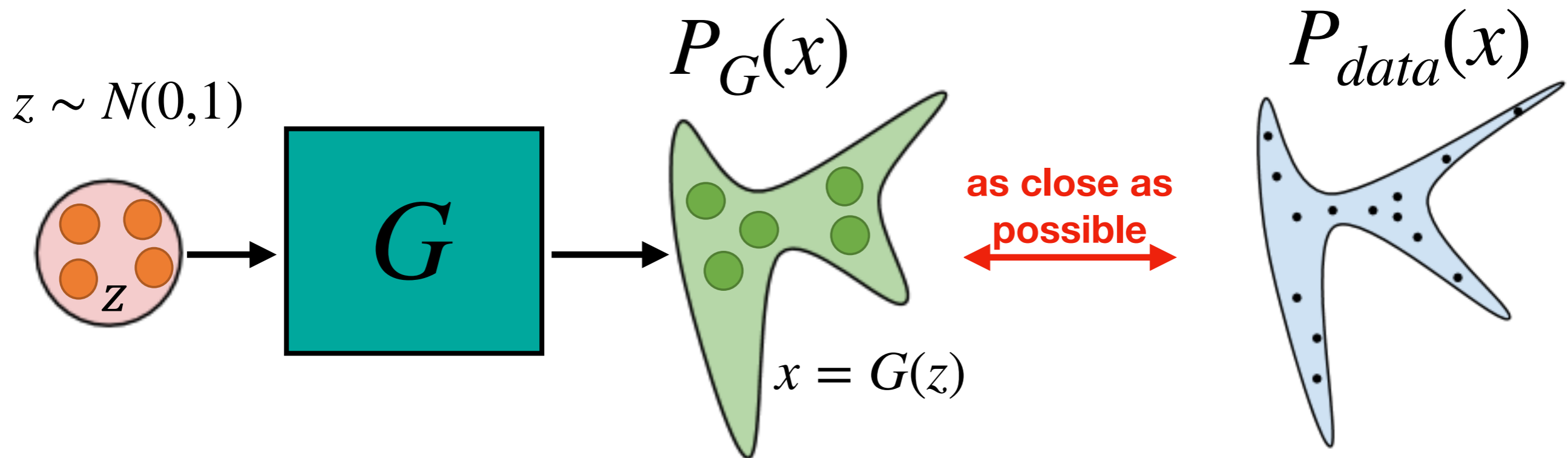
- In each training iteration:

**Learning D**
  - Sample m examples $\{x^1, x^2, \ldots, x^m\}$ from database
  - Sample m noise samples $\{z^1, z^2, \ldots, z^m\}$ from a distribution
  - Obtaining generated data $\{\tilde{x}^1, \tilde{x}^2, \ldots, \tilde{x}^m\}$, $\tilde{x}^i = G(z^i)$
  - Update discriminator parameters $\theta_d$ to maximize
    - $\tilde{V} = \frac{1}{m}\sum_{i=1}^m logD(x^i) + \frac{1}{m}\sum_{i=1}^m log\left(1 - D(\tilde{x}^i)\right)$
    - $\theta_d \leftarrow \theta_d + \eta\nabla\tilde{V}(\theta_d)$

**Learning G**
  - Sample m noise samples $\{z^1, z^2, \ldots, z^m\}$ from a distribution
  - Update generator parameters $\theta_g$ to maximize
    - $\tilde{V} = \frac{1}{m}\sum_{i=1}^m log\left(D\left(G(z^i)\right)\right)$
    - $\theta_g \leftarrow \theta_g - \eta\nabla\tilde{V}(\theta_g)$

# Basic idea of GAN

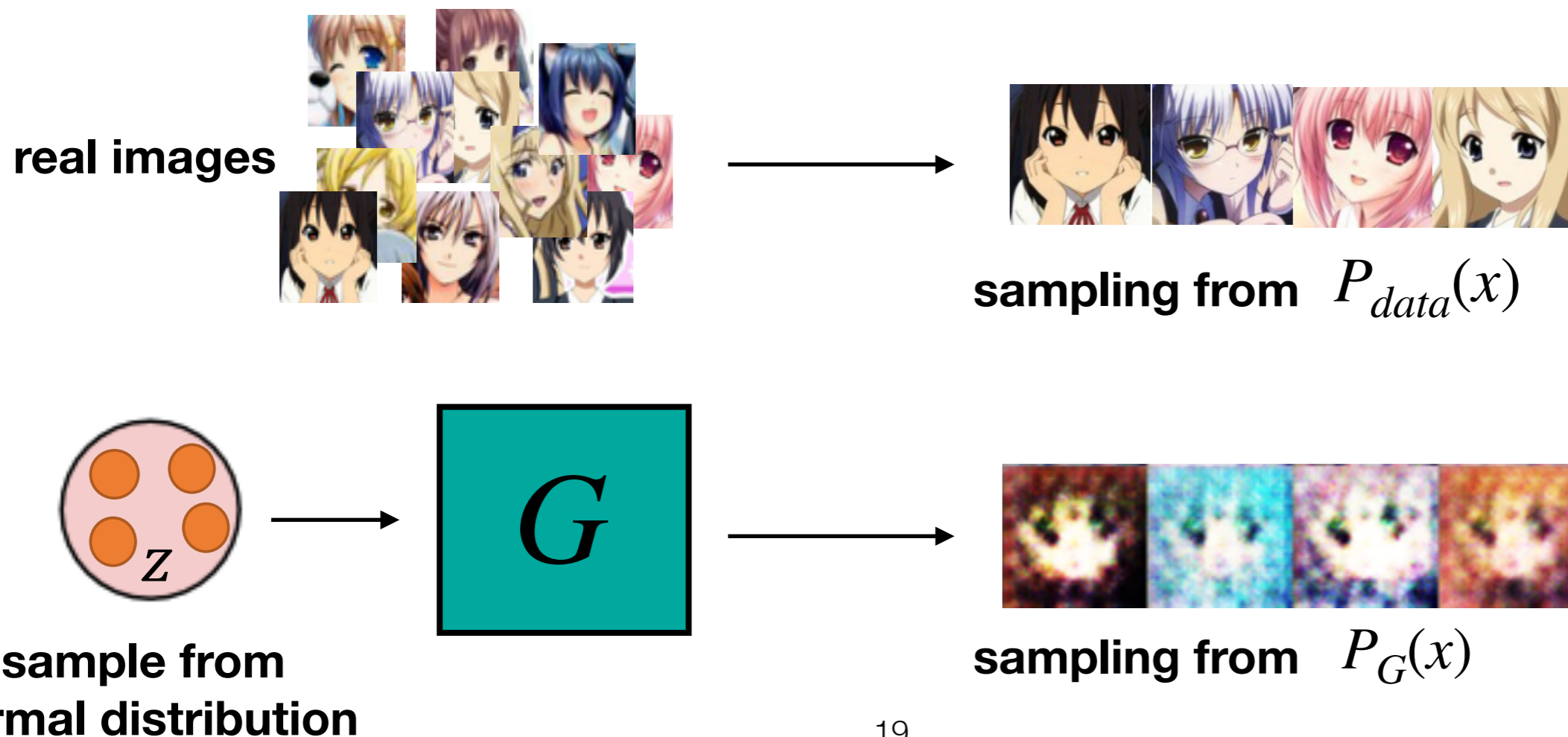**Generator**: G is a network. It defines a probability distribution $P_G$



$z \sim N(0,1)$

$G$

$P_G(x)$

$x = G(z)$

**as close as possible**

$P_{data}(x)$

$$G^* = \underset{G}{argmin}\ \boxed{Div(P_G, P_{data})}$$

**how to compute the divergence between two distributions ?**

# Basic idea of GAN

**Discriminator**
$$G^* = \underset{G}{argmin} \; Div(P_G, P_{data})$$

Although we do not know the distributions of $P_G(x)$ and $P_{data}(x)$ , we can still sample from them

**real images**



**sampling from** $P_{data}(x)$

**sample from normal distribution**



**sampling from** $P_G(x)$

# Basic idea of GAN

**Discriminator**

$$G* = \underset{G}{argmin} \, Div(P_G, P_{data})$$

**Objective function for D**

$$V(G, D) = E_{x \sim P_{data}}[logD(x)] + E_{x \sim P_G}[log(1 - D(x))]$$

**(G is fixed)**

$$D* = \underset{D}{argmax} \, V(G, D)$$

JS Divergence

**= binary classification**

# Basic idea of GAN

**Discriminator**

$$G* = \underset{G}{argmin} \ Div(P_G, P_{data})$$

**Objective function for G**

$$G* = \underset{G}{argmin}(E_{x \sim P_{data}}[logD(x)] + E_{x \sim P_G}[log(1 - D(G(z)))])$$

**(D is fixed)**

$$E_{x \sim P_G}[-log(D(G(z)))])$$

21

# Basic idea of GAN

$$E_{x \sim P_G}[log(1 - D(G(z)))])$$

**slow at the beginning**

$$E_{x \sim P_G}[-log(D(G(z)))])$$

**real implementation**

$-log(D(x))$

$0.5$    $D(x)$    $1$

$log(1 - D(x))$

# Basic idea of GAN

Different GANs

- Wasserstein GAN

- Wasserstein GAN-GP (gradient penalty)

- LSGAN

- ...

# Basic idea of GAN

$$V(G, D) = E_{x \sim P_{data}}[logD(x)] + E_{z \sim P_z}[log(1 - D(G(z)))]$$

$$G^* = \underset{G}{argmin} \ \underset{D}{max} V(G, D)$$

Training Steps:

- Initialize Generator and Discriminator

- In each training iteration:

    Step 1: Fix Generator G, and update Discriminator D

    Step 2: Fix Discriminator D, and update Generator G

# Vanilla GAN  (unconditional)

# Vanilla GAN [Ian Goodfellow, et al, NIPS 2014]

# Conditional GAN

# Conditional GAN

**CGAN**

D(X)

**D**

X

**real image**

X' = G(z)    **fake image**

**G**

C    z

# Conditional GAN

**ACGAN**

[August Odena, et al, ICML 2016]



D(X)    c

D

X
real image

X' = G(z)    fake image

G

c    z

# Conditional GAN

**male, with glasses**

**female, with glasses**

**male, without glasses**

**female, without glasses**

# Conditional GAN



without glasses, female, no black hair, no smiling, young

without glasses, male, no black hair, smiling, young

without glasses, female, black hair, smiling, young

with glasses, male, black hair, no smiling, young

with glasses, female, black hair, no smiling, old
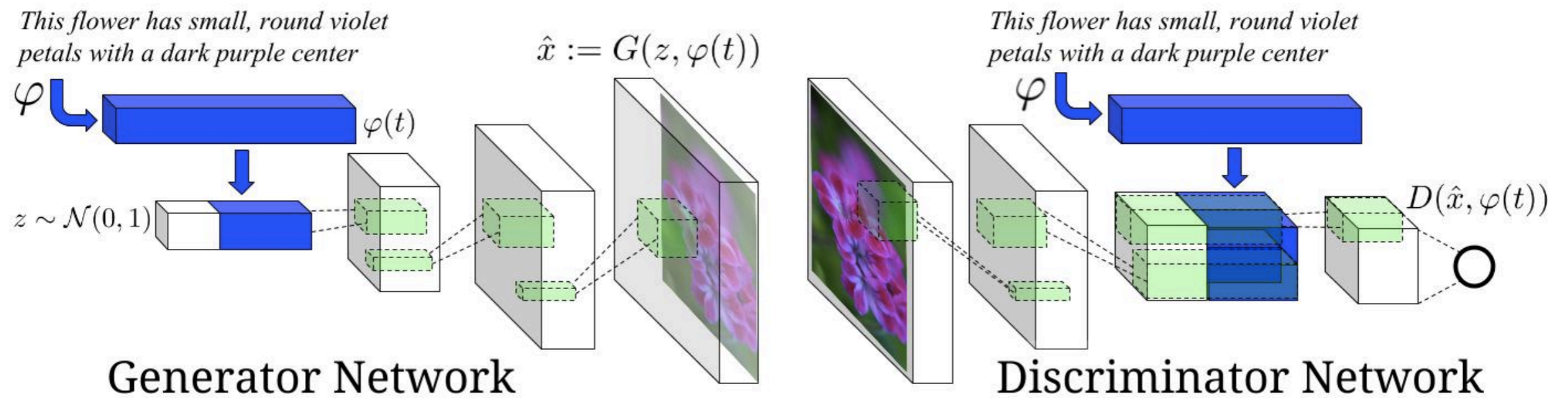
with glasses, male, black hair, smiling, old

with glasses, female, no black hair, smiling, old

without glasses, male, no black hair, no smiling, old
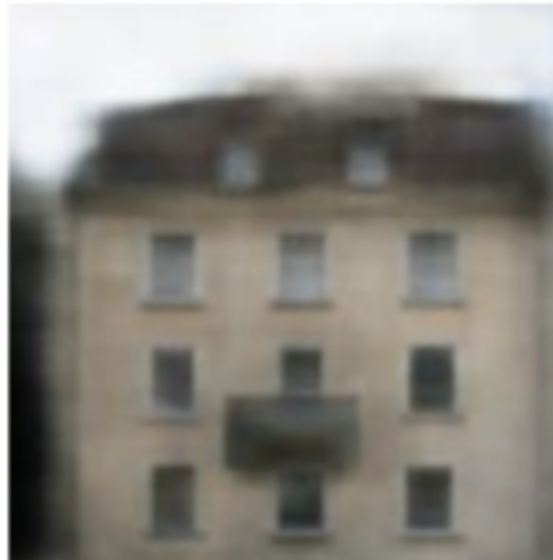
# Conditional GAN

## Text-to-image Generation



This flower has small, round violet petals with a dark purple center

$\hat{x} := G(z, \varphi(t))$

$\varphi(t)$

$z \sim \mathcal{N}(0, 1)$

**Generator Network**

This flower has small, round violet petals with a dark purple center

$D(\hat{x}, \varphi(t))$

**Discriminator Network**

# Image-to-image translation

# Image-to-image translation

- Traditional method



**NN**

**Image**

as close as possible

L1 / L2 loss

Testing:



**It is blurry,
what is the problem here ?**

# Image-to-image translation
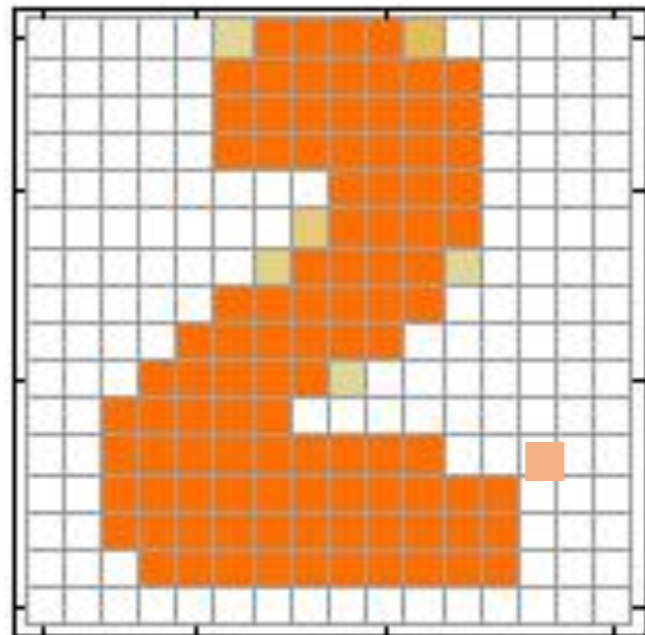
**generated image**                    **target**



as close as
possible

# Image-to-image translation



**target**

**1 pixel error**
**not realistic**

**1 pixel error**
**not realistic**

**6 pixel error**
**realistic**

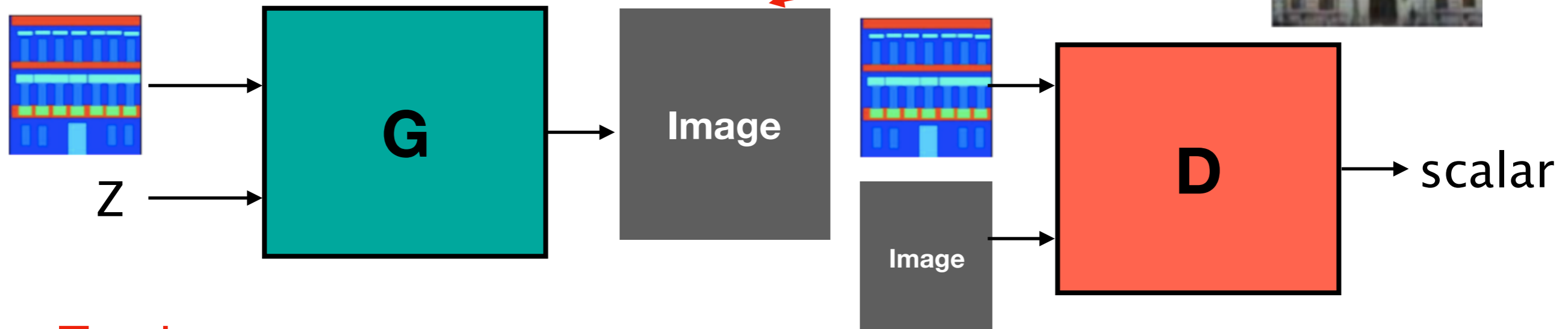**6 pixel error**
**realistic**

# Image-to-image translation

**Reconstruction loss can not provide a sharp generation, what should be the solution ?**

**Since we can not find a good metric, we can use GAN to learn the metric !**
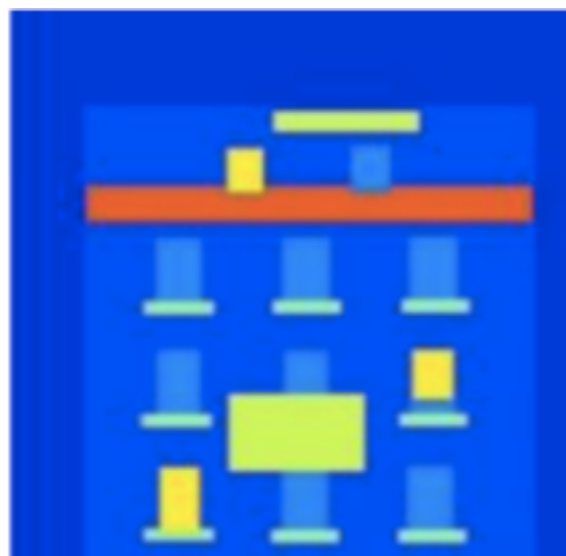
# Image-to-image translation

- GAN method (Pix2Pix)

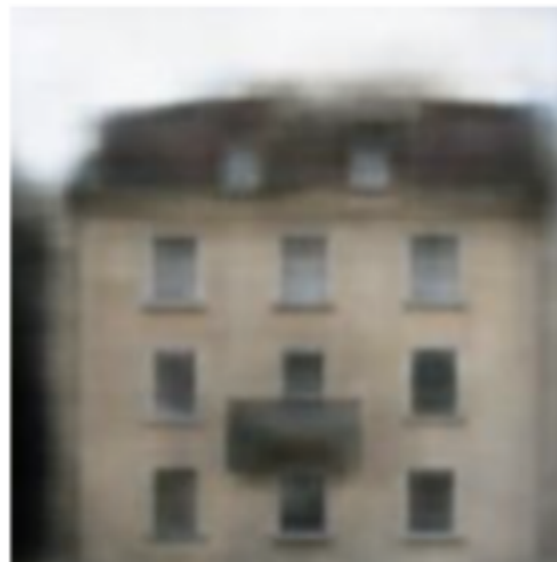**as close as possible**



G → Image

Z

D → scalar

Image

Testing:

**Input**          **Reconstruct**          **GAN**          **GAN + Reconstruct**
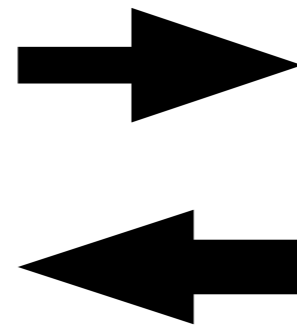
# Image-to-image translation

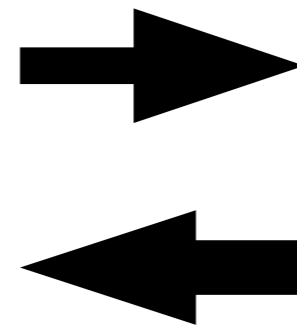- What about unpaired data (no ground truth of target image)?



**X: zebra**

**Y: horse**

**X: summer**
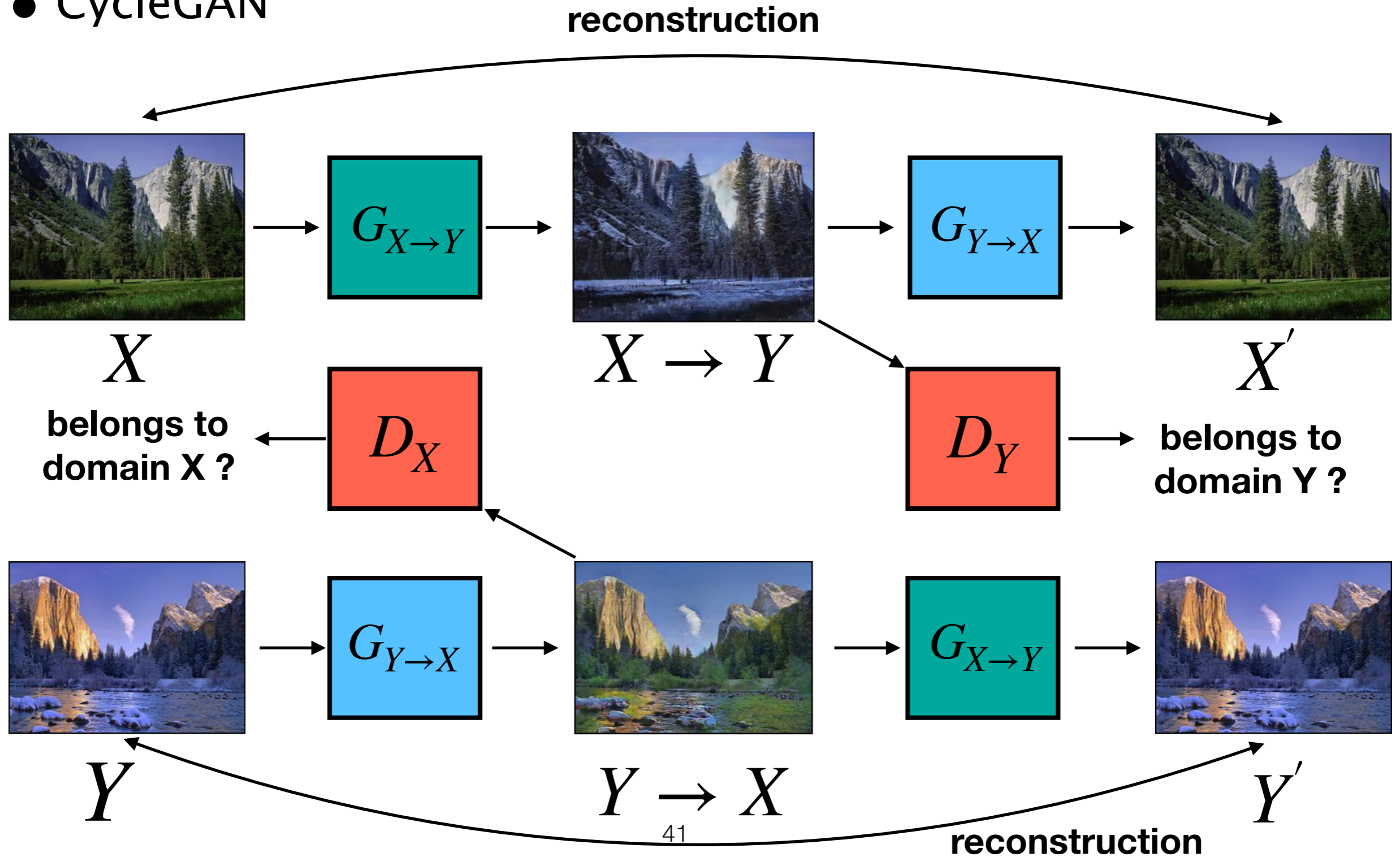
**Y: winter**

# Image-to-image translation

- CycleGAN

# Image-to-image translation

- CycleGAN

**reconstruction**



$X$

$X \rightarrow Y$

$X^{'}$

**belongs to domain X ?**

$G_{X \rightarrow Y}$

$G_{Y \rightarrow X}$

$D_X$

$D_Y$

**belongs to domain Y ?**

$Y$

$Y \rightarrow X$

$Y^{'}$

$G_{Y \rightarrow X}$

$G_{X \rightarrow Y}$

**reconstruction**

41

**Monet ⟳ Photos**

Monet → photo

photo → Monet

**Zebras ⟳ Horses**

zebra → horse

horse → zebra

**Summer ⟳ Winter**

summer → winter

winter → summer

Photograph → Monet   Van Gogh   Cezanne   Ukiyo-e

# Image-to-image translation

- UNIT

- MUNIT

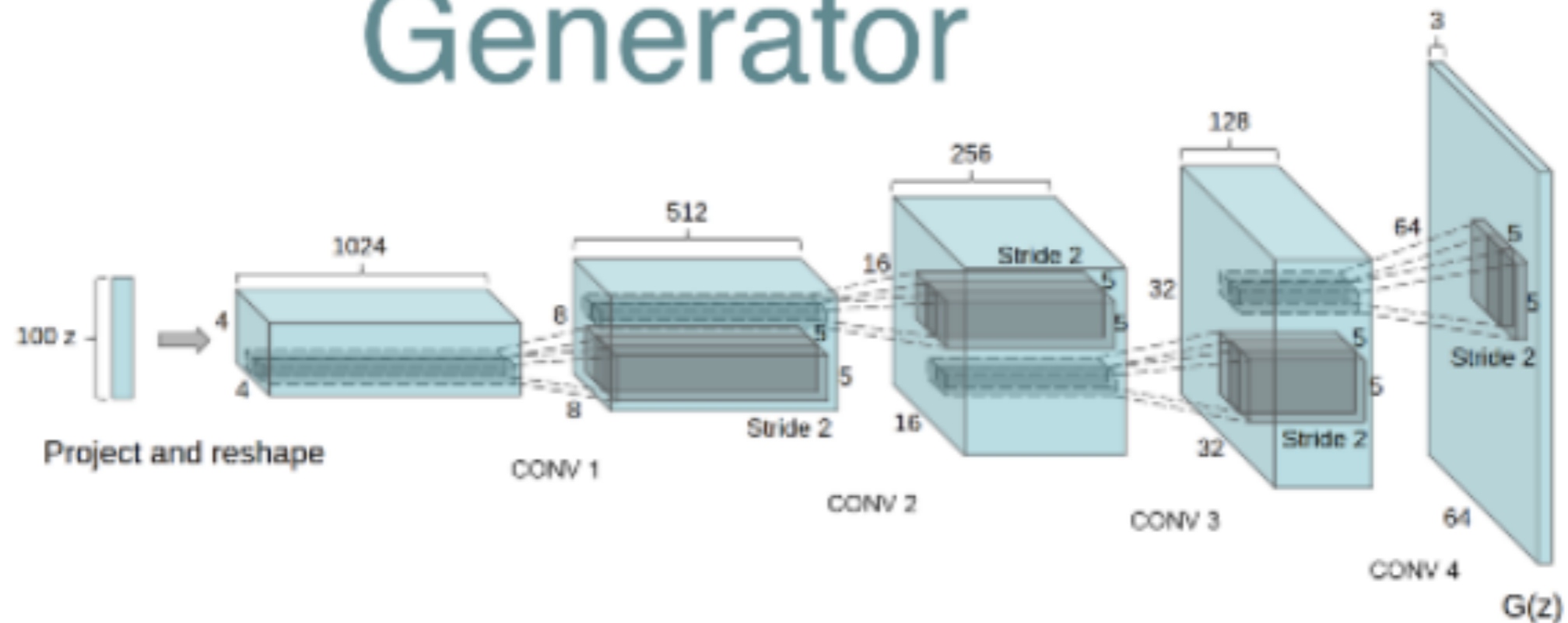- …

# Modern GAN Architectures
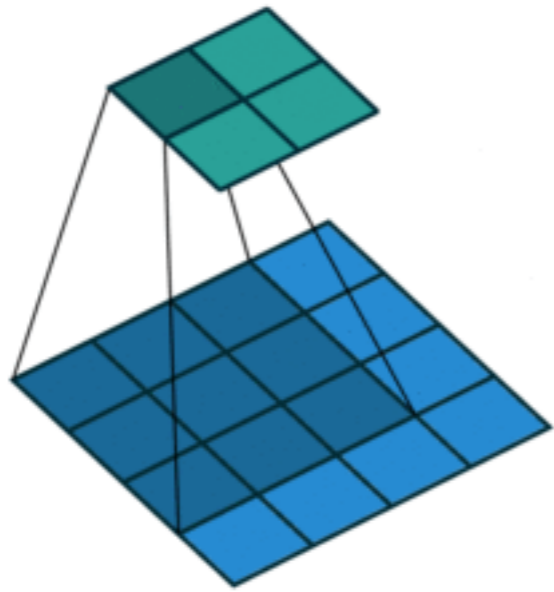
# Modern Architectures

## DCGAN

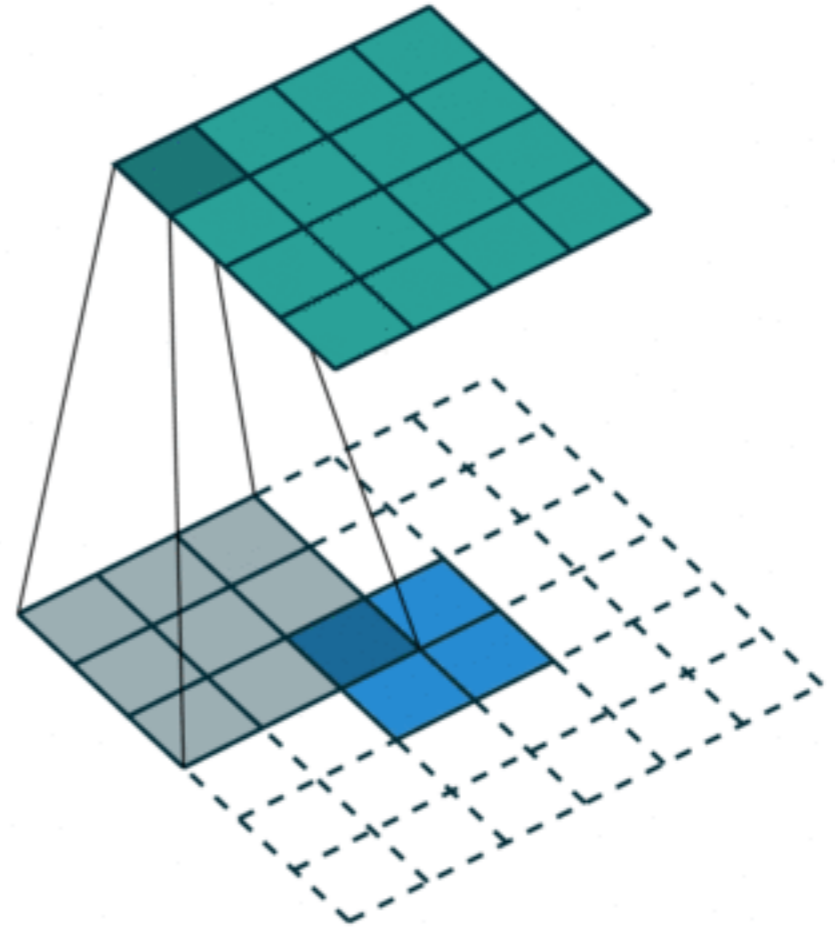[A Radford, et al, arXiv 2015]

**convolution**

**transposed convolution**

# Modern Architectures



Discriminator

[A Radford, et al, arXiv 2015]

# Results - MNIST

# Results - CelebA (faces)

# Results - LSUN (bedrooms)

# Modern Architectures

## StyleGAN (NVIDIA)

https://github.com/NVlabs/stylegan



[T Karras, et al, CVPR 2019]

# Modern Architectures

## StyleGAN



https://www.youtube.com/watch?v=kSLJriaOumA

Karras et al, A Style-Based Generator Architecture for Generative Adversarial Networks, CVPR 2019

# Modern Architectures

## StyleGAN

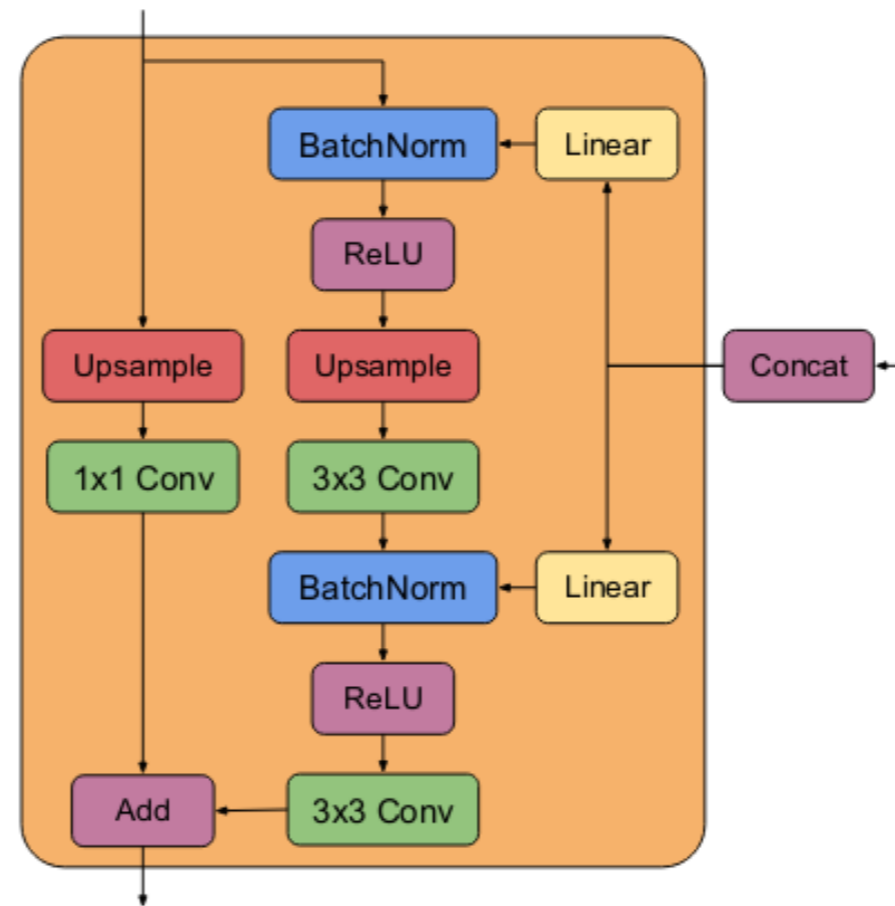| GPUs | 1024×1024 | 512×512 | 256×256 |
|---|---|---|---|
| 1 | 41 days 4 hours | 24 days 21 hours | 14 days 22 hours |
| 2 | 21 days 22 hours | 13 days 7 hours | 9 days 5 hours |
| 4 | 11 days 8 hours | 7 days 0 hours | 4 days 21 hours |
| 8 | 6 days 14 hours | 4 days 10 hours | 3 days 8 hours |

# Modern Architectures
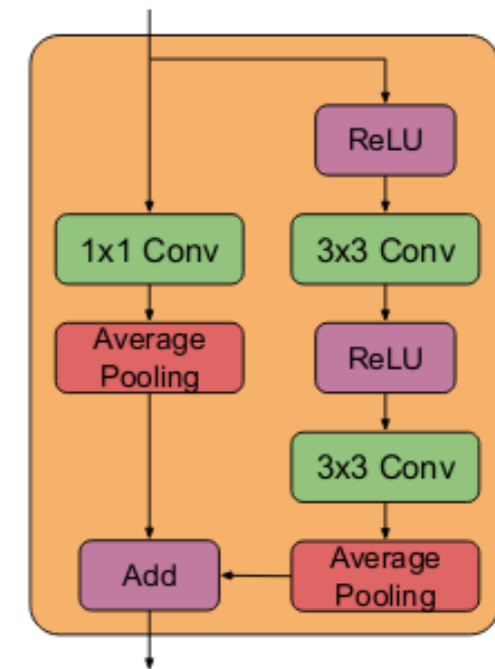
## BigGAN (DeepMind)

https://github.com/ajbrock/BigGAN-PyTorch



[A Brock, et al, ICLR 2019]

# Modern Architectures

## BigGAN

On 8xV100 with full-precision training (no Tensor cores), this script takes 15 days to train to 150k iterations.
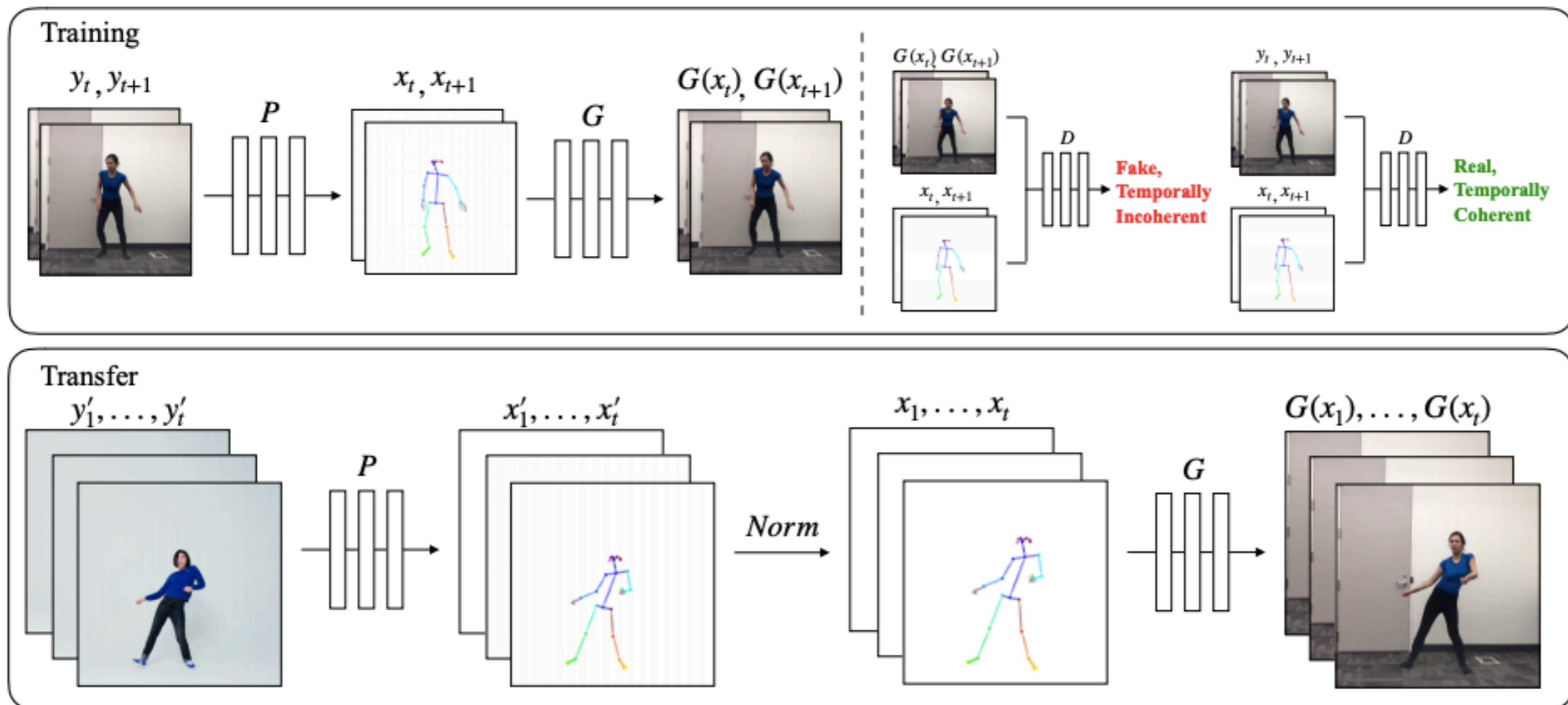
# Vid-to-vid translation

# Vid-to-vid translation
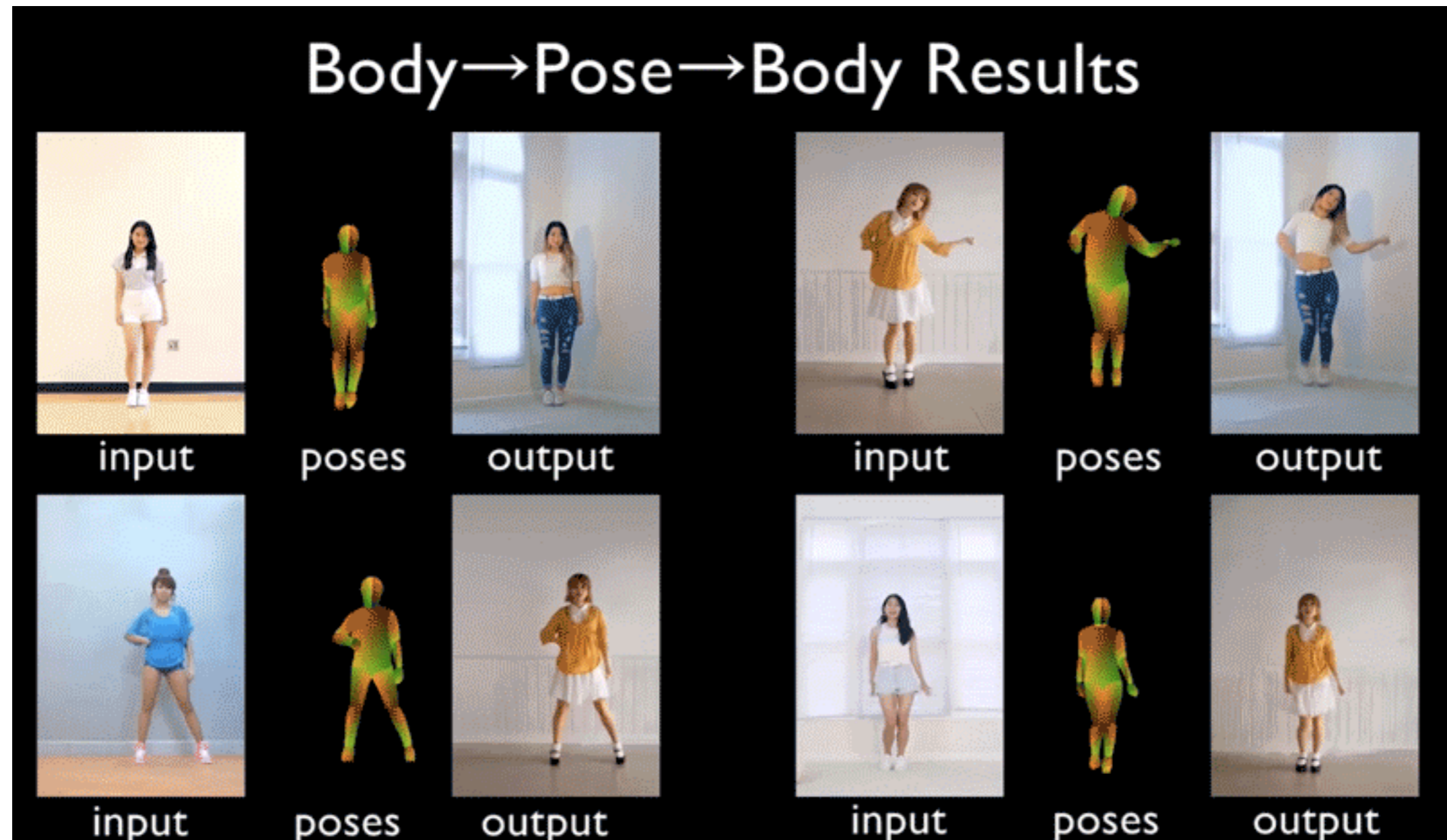
[Carolin Chan, et al, ICCV 2019]

- Everybody dance now



https://www.youtube.com/watch?v=PCBTZh41Ris

# Vid-to-vid translation

- Video–to–video synthesis

[Ting-chun Wang, et al, NIPS 2018]

https://github.com/NVIDIA/vid2vid

# Outline

- Basic Idea of GAN

- Image Generation

  - Conditional GAN (CGAN, ACGAN)

  - Modern GANs (StyleGAN, BigGAN)

  - Image-to-image translation (Pix2Pix, CycleGAN)

- Video-to-video translation

- GANs Evaluation

- Video Generation

- Lab (DCGAN for manga face generation)

# GANs Evaluation

# GANs Evaluation

Two Metrics:

- Inception Score (IS) ⬆

- Fréchet Inception Distance (FID) ⬇

# GANs Evaluation-IS

Requirements:

- High-quality (clear contents, sharp images)

- Diversity (different contents)

# GANs Evaluation - IS

## Conditional generation

**Definition:** $IS(G) = exp(\mathbb{E}_{x \sim p_g} D_{KL}(p(y \,|\, x) \,||\, p(y)))$

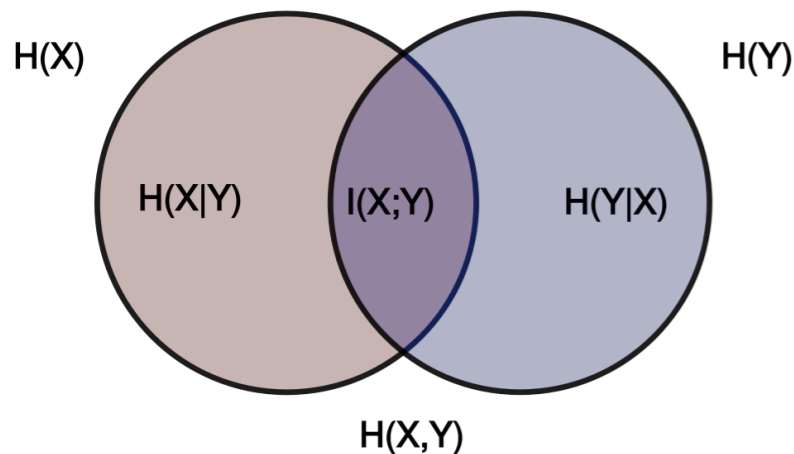$x$ **generated sample**

$p(y \,|\, x)$ **Conditional distribution**

$p(y) = \int_x p(y \,|\, x) p_g(x)$ **Marginal distribution**

# GANs Evaluation - IS

**Definition**: $IS(G) = exp(\mathbb{E}_{x \sim p_g} D_{KL}(p(y|x) \mid\mid p(y)))$

$$ln(IS(G)) = E_{x \sim p_g} D_{KL}(p(y|x)||p(y))$$

$$= \sum_x p(x) D_{KL}(p(y|x)||p(y))$$

$$= \sum_x p(x) \sum_i p(y=i|x) ln(\frac{p(y=i|x)}{p(y=i)})$$

$$= \sum_x \sum_i p(x, y=i) ln \frac{p(x, y=i)}{p(x)p(y=i)}$$

$$= I(y; x) \quad \text{Mutual Information}$$
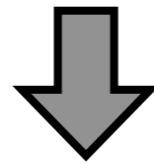
$$= H(y) - H(y|x) \quad \text{Entropy difference}$$

H(X)

H(Y)

H(X|Y)  I(X;Y)  H(Y|X)

H(X,Y)

# GANs Evaluation - IS

- The generative algorithm should output a high diversity of images from all the different classes in ImageNet —> H(y) should be high

- The images generated should contain clear objects (i.e. the images are sharp rather than blurry) —> H(y|x) should be low
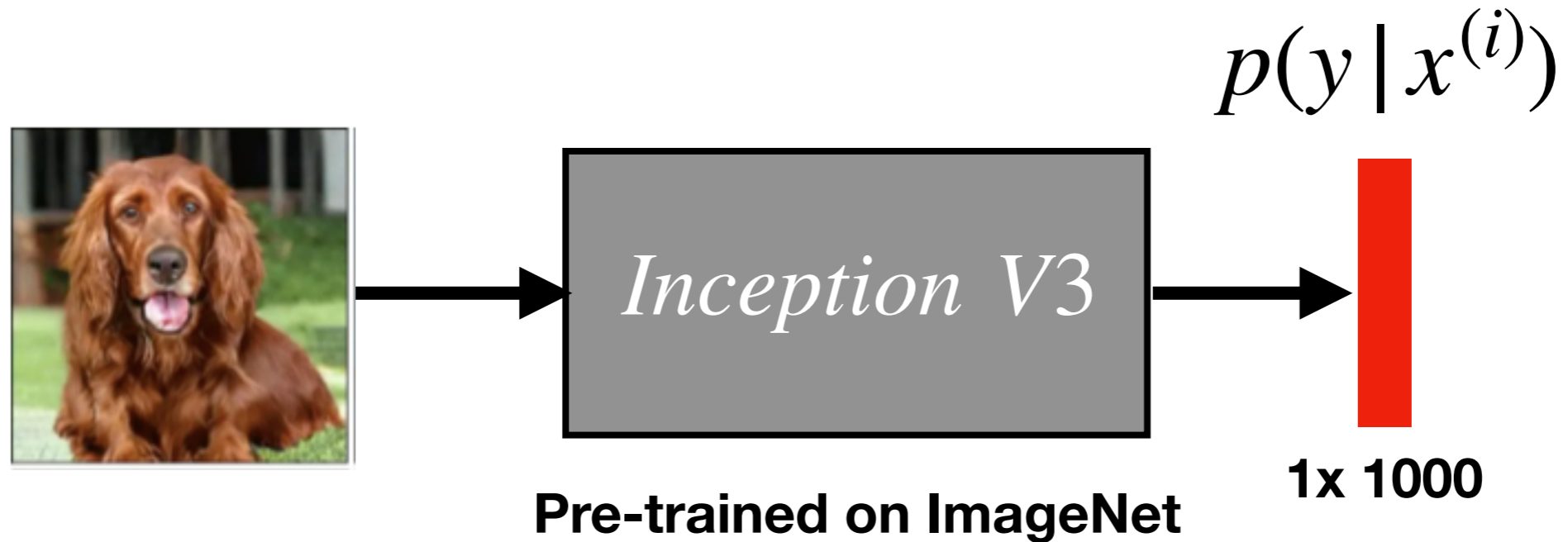
# GANs Evaluation - IS

$$IS(G) = exp(\mathbb{E}_{x \sim p_g} D_{KL}(p(y|x) \,||\, p(y)))$$

$$IS(G) \approx exp(\mathbb{E}_{x \sim p_g} D_{KL}(p(y|x^{(i)}) \,||\, \hat{p}(y)))$$

$$\hat{p}(y) = \frac{1}{N} \sum_{i=1}^{N} p(y|x^{(i)})$$

# GANs Evaluation - IS

$$p(y|x^{(i)})$$



**Pre-trained on ImageNet**

**1x 1000**

**Sampled 5000 images from GANs**

**Attention: training and evaluation must use the same dataset**

# GANs Evaluation - IS

Problem of IS ?

Not considering the distribution of training dataset

# GANs Evaluation - FID

$$FID = |\mu - \mu_w|^2 + tr(\Sigma + \Sigma_w - 2(\Sigma\Sigma_w)^{1/2})$$

$N(\mu, \Sigma)$   distribution of generated set

$N(\mu_w, \Sigma_w)$   distribution of training set

Computer Fréchet distance
between two distributions

# GANs Evaluation - FID

$$FID = |\mu - \mu_w|^2 + tr(\Sigma + \Sigma_w - 2(\Sigma\Sigma_w)^{1/2})$$



**Inception V3**

**Pre-trained on ImageNet**

$p(y|x^{(i)})$

**feature**

# GANs Evaluation - FID



*Inception V3*

**Pre-trained on ImageNet**

$p(y|x^{(i)})$

**feature**

**Sampled approximately same number images from GANs as original training set, ideally 10000 images for ImageNet**

# GANs Evaluation

- IS and FID are two metrics for GANs evaluation. FID is more widely used than IS. However, both methods require large-scale generated samples.

- New metric needs to be proposed to guide training process.

# Video Generation

# Recall: Spatio-temporal Modeling

- *Spatiotemporal CNN (3D CNN)*

- *LSTM and GRU*

# Spatiotemporal CNN (3D CNN)

# LSTM and GRU



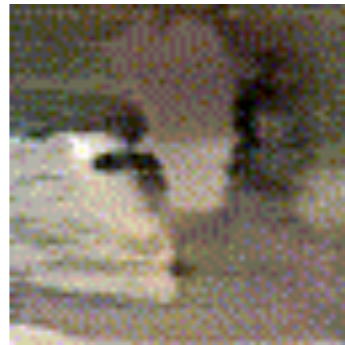**LSTM**                    **GRU**

# Spatiotemporal CNN (3D CNN)



**2D Generator**

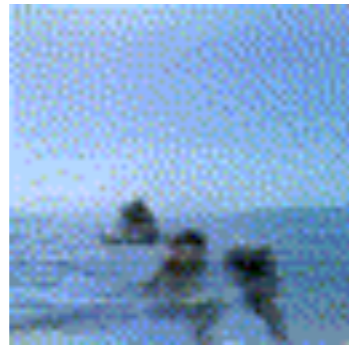image



**3D Generator**

video
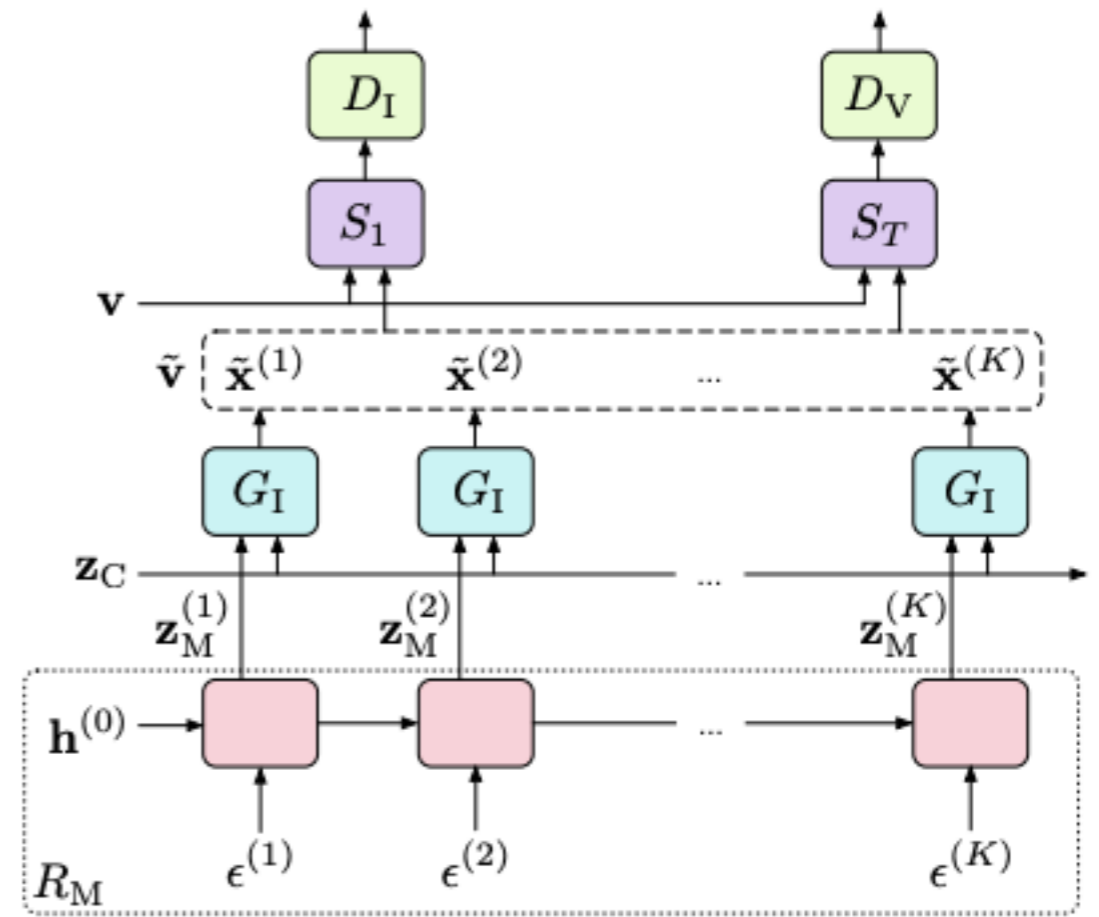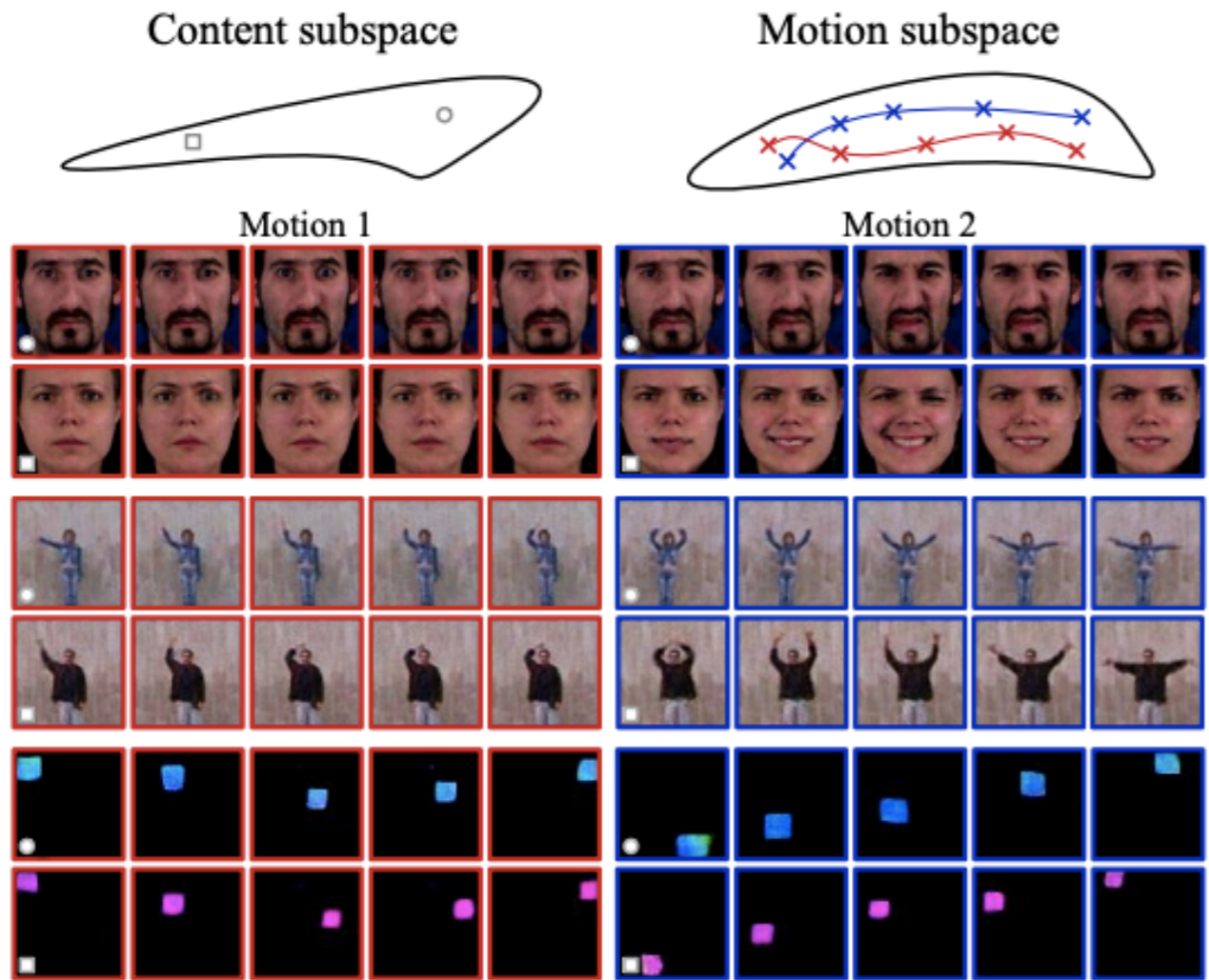
# LSTM and GRU

# VGAN



VGAN [NeurIPS'16]

# VGAN



**http://www.cs.columbia.edu/~vondrick/tinyvideo/**

# MoCoGAN



MoCoGAN [CVPR'19]

# MoCoGAN

# MoCoGAN

# G³AN



G3AN [CVPR'20]

# G³AN

**What I can not create, I do not understand**

*- R. Feynman*

# Thank You !