

# Introduction of Generative Adversarial Network (GAN)

李宏毅

Hung-yi Lee

# Generative Adversarial Network (GAN)

- How to pronounce “GAN”?



Google 小姐

# Yann LeCun's comment

## What are some recent and potentially upcoming breakthroughs in unsupervised learning?



**Yann LeCun**, Director of AI Research at Facebook and Professor at NYU



Written Jul 29 · Upvoted by Joaquin Quiñonero Candela, Director Applied Machine Learning at Facebook and Huang Xiao

Adversarial training is the coolest thing since sliced bread.

I've listed a bunch of relevant papers in a previous answer.

Expect more impressive results with this technique in the coming years.

What's missing at the moment is a good understanding of it so we can make it work reliably. It's very finicky. Sort of like ConvNet were in the 1990s, when I had the reputation of being the only person who could make them work (which wasn't true).

# Yann LeCun's comment

## What are some recent and potentially upcoming breakthroughs in deep learning?



**Yann LeCun**, Director of AI Research at Facebook and Professor at NYU

Written Jul 29 · Upvoted by Joaquin Quiñonero Candela, Director Applied Machine Learning at Facebook and Nikhil Garg, I lead a team of Quora engineers working on ML/NLP problems



.....

The most important one, in my opinion, is adversarial training (also called GAN for Generative Adversarial Networks). This is an idea that was originally proposed by Ian Goodfellow when he was a student with Yoshua Bengio at the University of Montreal (he since moved to Google Brain and recently to OpenAI).

This, and the variations that are now being proposed is the most interesting idea in the last 10 years in ML, in my opinion.

<https://www.quora.com/What-are-some-recent-and-potentially-upcoming-breakthroughs-in-deep-learning>

# All Kinds of GAN ...

<https://github.com/hindupuravinash/the-gan-zoo>

GAN

ACGAN

BGAN

CGAN

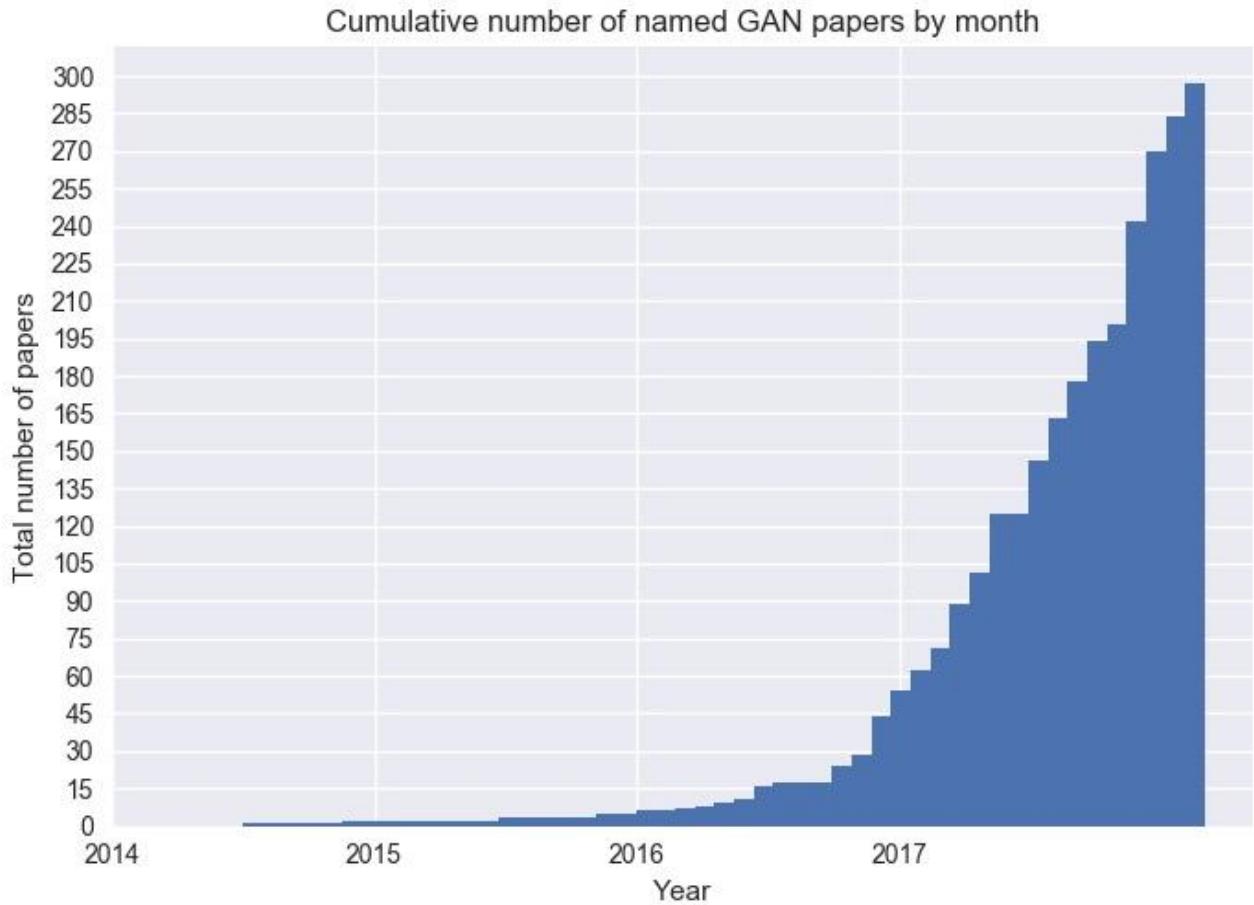
DCGAN

EBGAN

fGAN

GoGAN

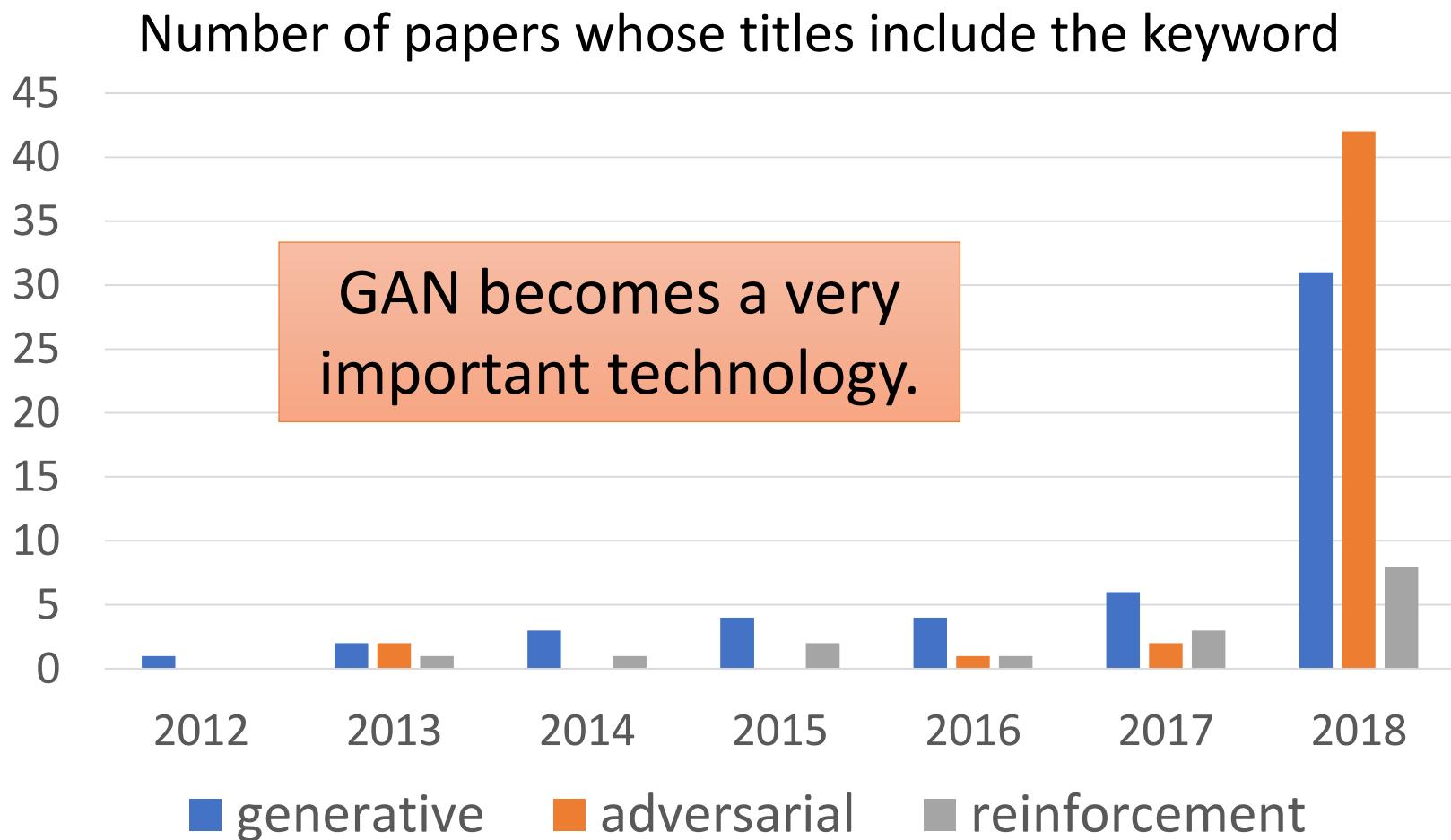
⋮



Mihaela Rosca, Balaji Lakshminarayanan, David Warde-Farley, Shakir Mohamed, "Variational Approaches for Auto-Encoding Generative Adversarial Networks", arXiv, 2017

<sup>2</sup>We use the Greek  $\alpha$  prefix for  $\alpha$ -GAN, as AEGAN and most other Latin prefixes seem to have been taken  
<https://deephunt.in/the-gan-zoo-79597dc8c347>.

Keyword search on session index page,  
so session names are included.



# Outline

Basic Idea of GAN

GAN as structured learning

Can Generator learn by itself?

Can Discriminator generate?

A little bit theory

# Generation

We will control what to generate latter. → Conditional Generation

## *Image Generation*

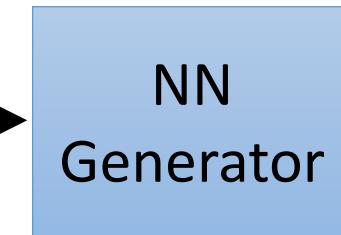
$$\begin{bmatrix} 0.3 \\ -0.1 \\ \vdots \\ -0.7 \end{bmatrix} \begin{bmatrix} 0.1 \\ -0.1 \\ \vdots \\ 0.7 \end{bmatrix} \begin{bmatrix} -0.3 \\ 0.1 \\ \vdots \\ 0.9 \end{bmatrix}$$

In a specific range



## *Sentence Generation*

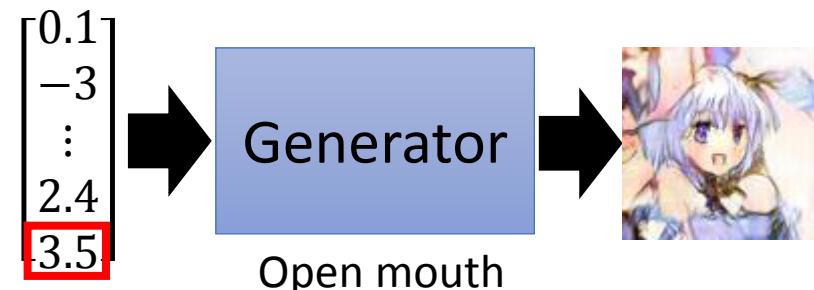
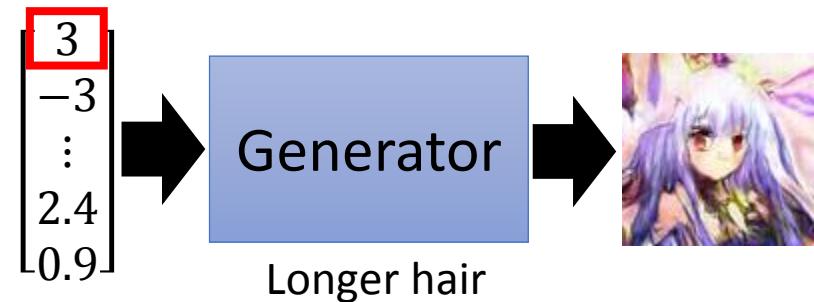
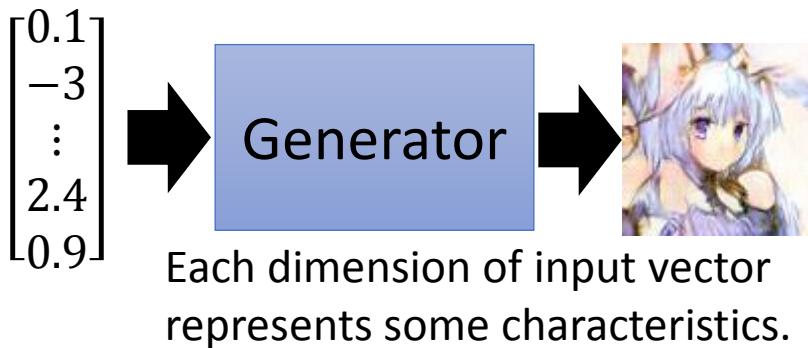
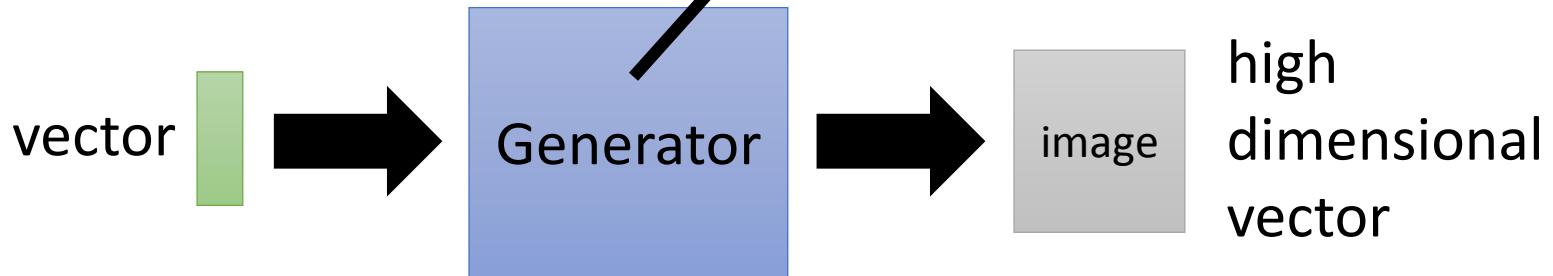
$$\begin{bmatrix} 0.3 \\ -0.1 \\ \vdots \\ -0.7 \end{bmatrix} \begin{bmatrix} 0.1 \\ -0.1 \\ \vdots \\ 0.2 \end{bmatrix} \begin{bmatrix} -0.3 \\ 0.1 \\ \vdots \\ 0.5 \end{bmatrix}$$



How are you?  
Good morning.  
Good afternoon.

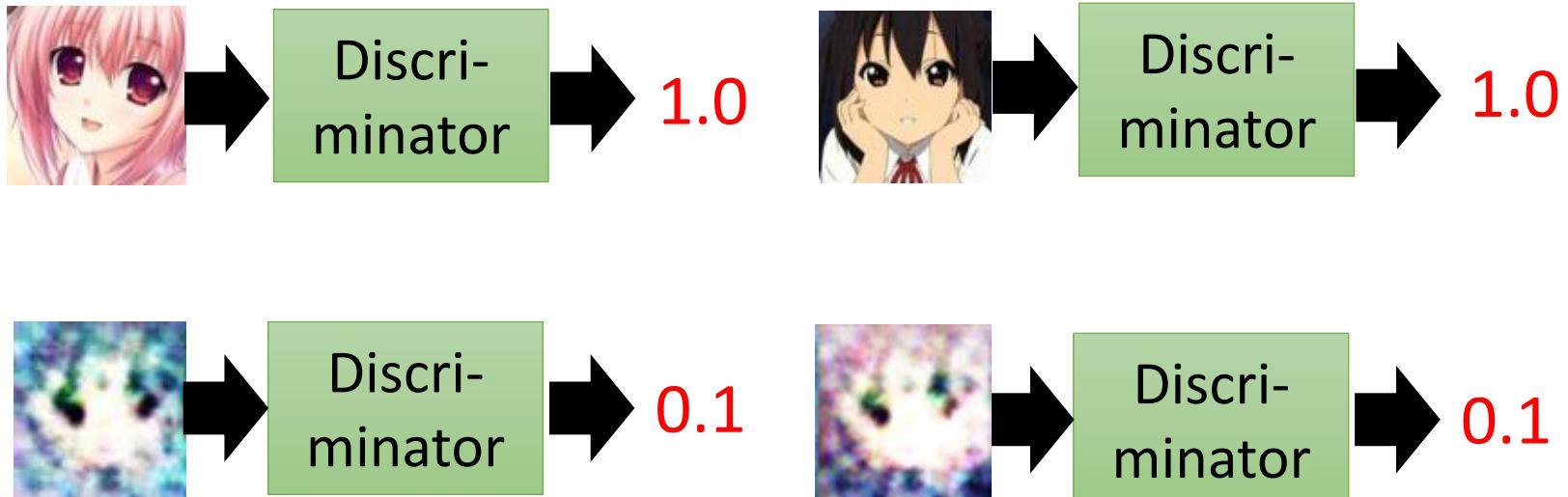
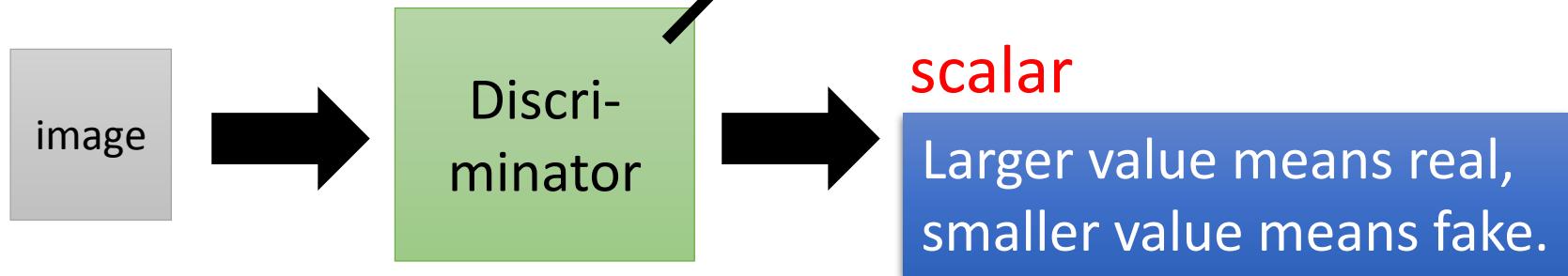
# Basic Idea of GAN

It is a neural network (NN), or a function.

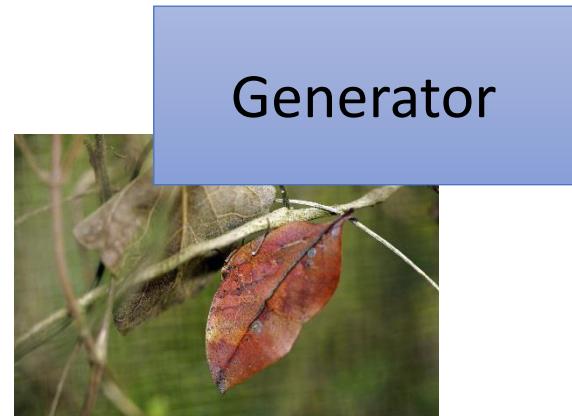
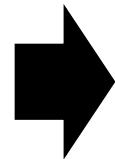
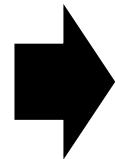


# Basic Idea of GAN

It is a neural network (NN), or a function.



# Basic Idea of GAN



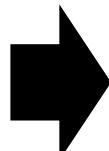
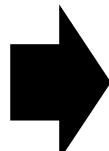
Brown

veins

Butterflies are  
not brown

Butterflies do  
not have veins

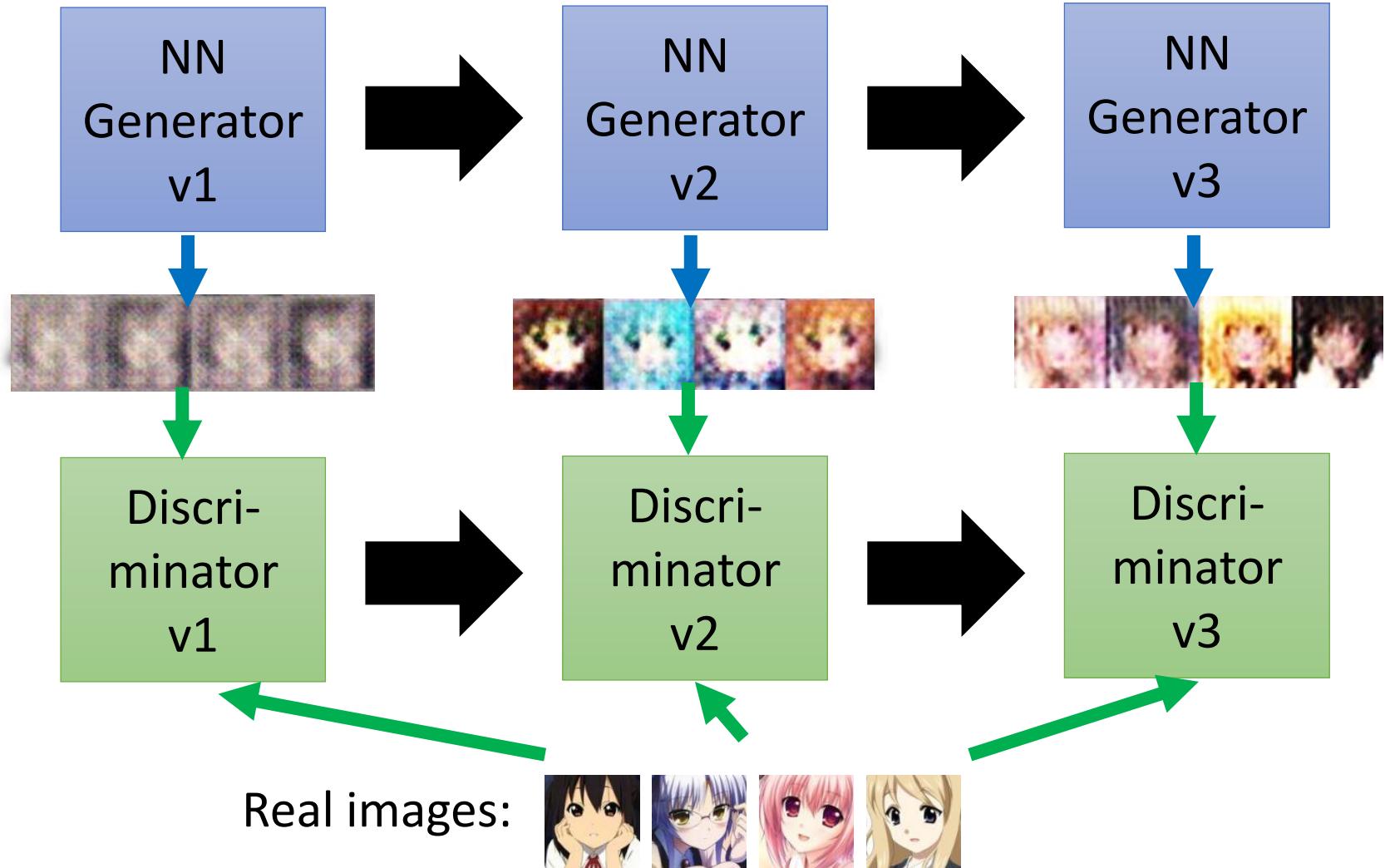
.....



Discriminator

# Basic Idea of GAN

This is where the term  
***“adversarial”*** comes from.  
You can explain the process  
in different ways.....



# Basic Idea of GAN (和平的比喻)

Generator  
(student)

Discriminator  
(teacher)



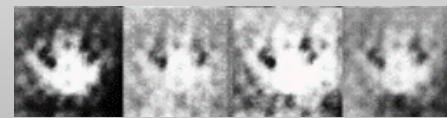
Generator  
v1



Discriminator  
v1

沒有兩個圈

Generator  
v2



Discriminator  
v2

沒有彩色

Generator  
v3



為什麼不自己學？

為什麼不自己做？

# Generator v.s. Discriminator

- 寫作敵人，唸做朋友

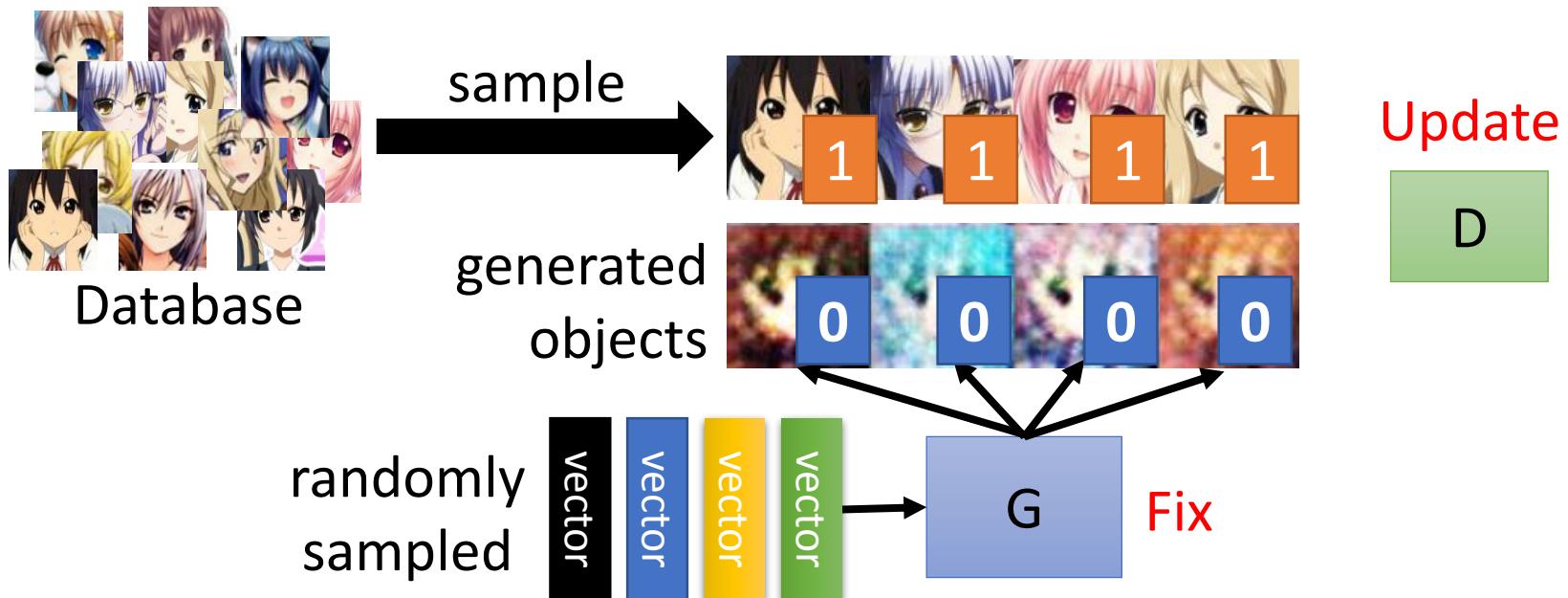


# Algorithm

- Initialize generator and discriminator
- In each training iteration:



**Step 1:** Fix generator G, and update discriminator D



Discriminator learns to assign high scores to real objects and low scores to generated objects.

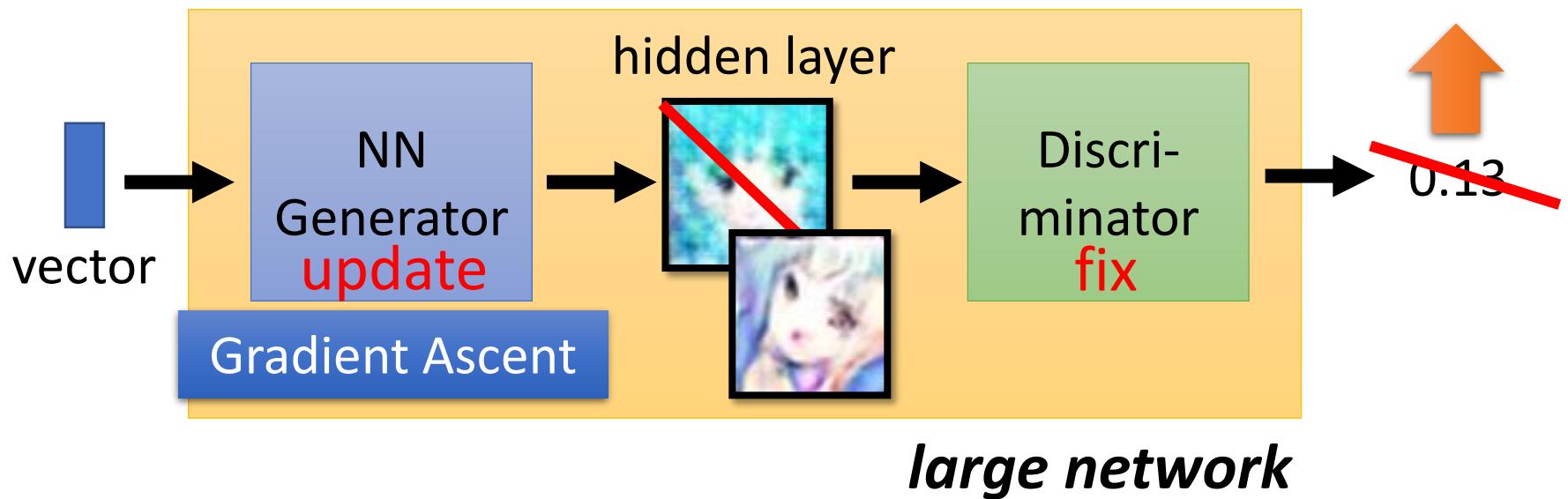
# Algorithm

- Initialize generator and discriminator
- In each training iteration:



**Step 2:** Fix discriminator D, and update generator G

Generator learns to “fool” the discriminator



# Algorithm

Initialize  $\theta_d$  for D and  $\theta_g$  for G

- In each training iteration:

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

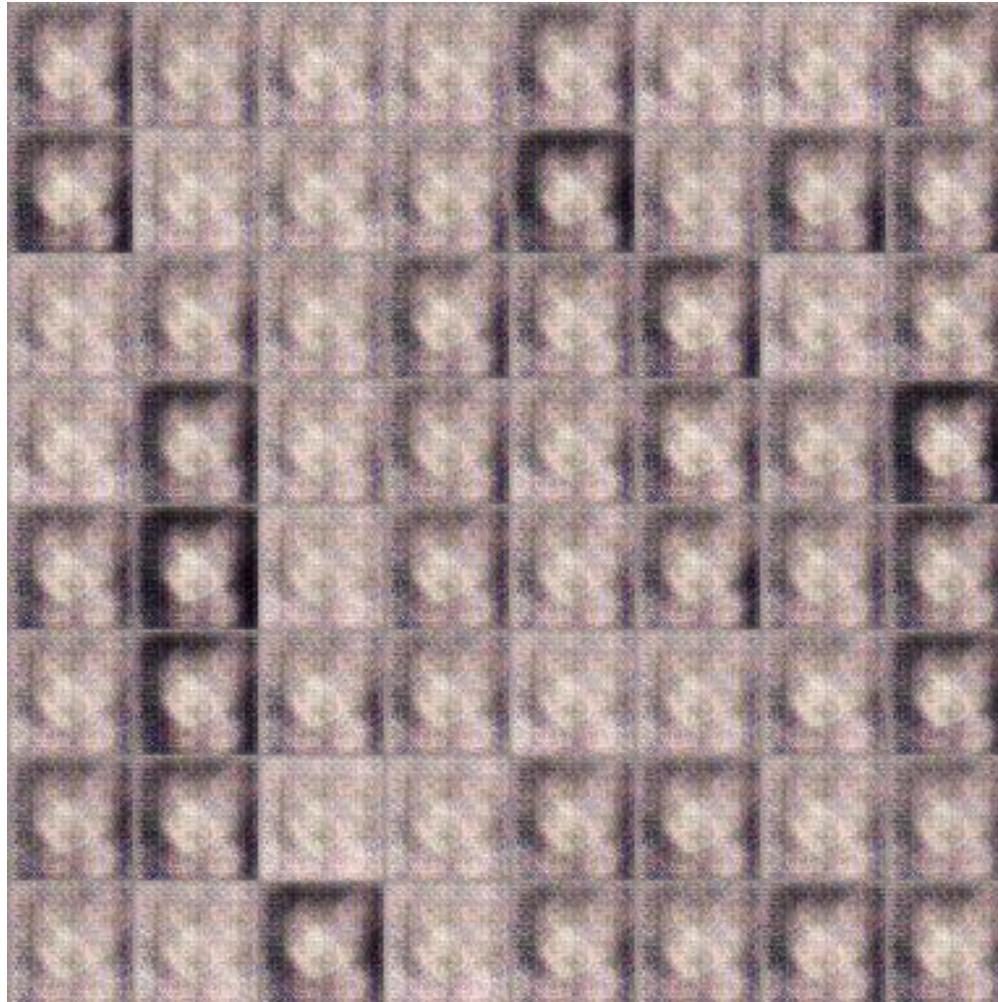
- Sample m noise samples  $\{z^1, z^2, \dots, z^m\}$  from a distribution
- Update generator parameters  $\theta_g$  to maximize
  - $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log (D(G(z^i)))$
  - $\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$

Learning  
D

Learning  
G

# Anime Face Generation

100 updates



Source of training data: <https://zhuanlan.zhihu.com/p/24767059>

# Anime Face Generation



1000 updates

# Anime Face Generation

2000 updates



# Anime Face Generation

5000 updates



# Anime Face Generation

10,000 updates



# Anime Face Generation

20,000 updates



# Anime Face Generation

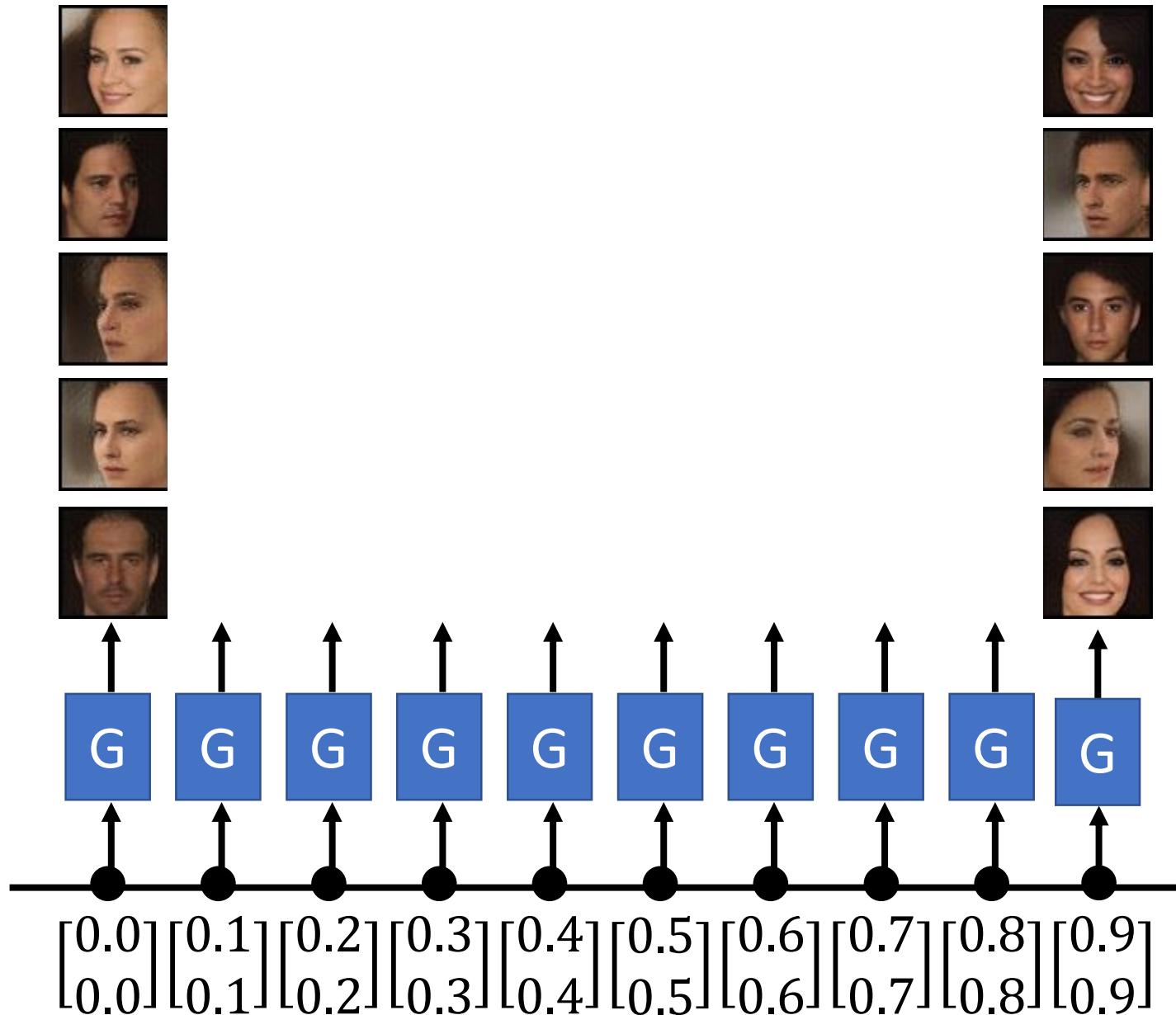
50,000 updates





The faces  
generated by  
machine.

圖片生成：  
吳宗翰、謝濬丞、  
陳延昊、錢柏均



感謝陳柏文同學提供實驗結果

# Outline

Basic Idea of GAN

GAN as structured learning

Can Generator learn by itself?

Can Discriminator generate?

A little bit theory

# Structured Learning

Machine learning is to find a function  $f$

$$f : X \rightarrow Y$$

**Regression:** output a scalar

**Classification:** output a “class” (one-hot vector)

1	0	0
---	---	---

Class 1

0	1	0
---	---	---

Class 2

0	0	1
---	---	---

Class 3

**Structured Learning/Prediction:** output a sequence, a matrix, a graph, a tree .....

Output is composed of components with dependency

# Output Sequence

$$f : X \rightarrow Y$$

## Machine Translation

$X$ ：“機器學習及其深層與  
結構化”  
(sentence of language 1)

$Y$ ：“Machine learning and  
having it deep and structured”  
(sentence of language 2)

## Speech Recognition

$X$ ：  
(speech)

$Y$ ：感謝大家來上課”  
(transcription)

## Chat-bot

$X$ ：“How are you?”  
(what a user says)

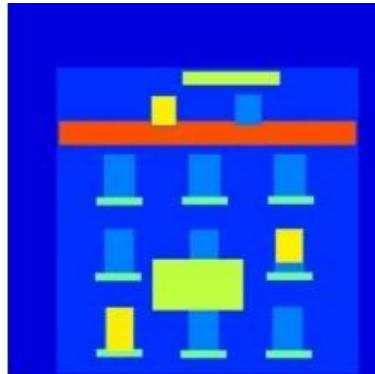
$Y$ ：“I'm fine.”  
(response of machine)

# Output Matrix

$$f : X \rightarrow Y$$

Image to Image

$X :$



$Y :$



Colorization:



Ref: <https://arxiv.org/pdf/1611.07004v1.pdf>

Text to Image

$X :$  “this white and yellow flower  
have thin white petals and a  
round yellow stamen”

$Y :$



ref: <https://arxiv.org/pdf/1605.05396.pdf>

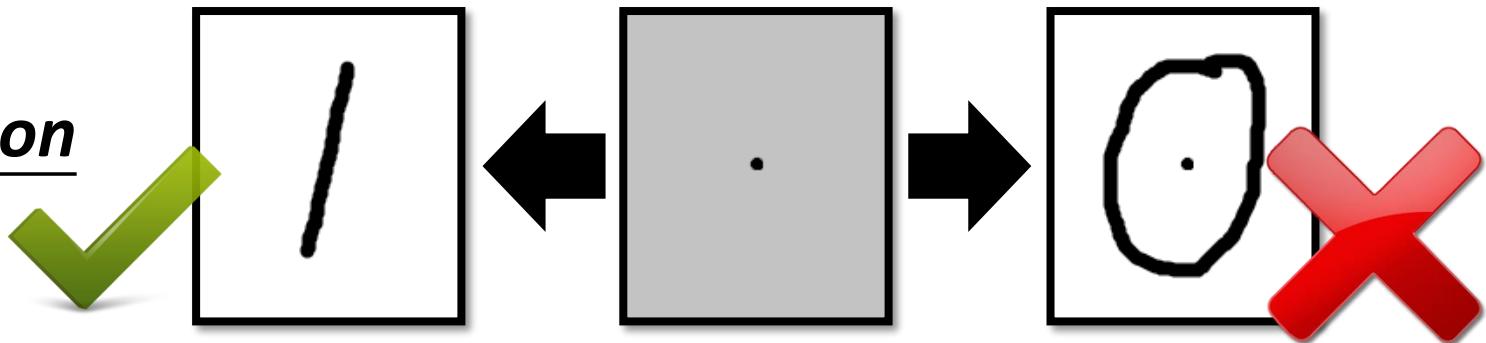
# Why Structured Learning Challenging?

- **One-shot/Zero-shot Learning:**
  - In classification, each class has some examples.
  - In structured learning,
    - If you consider each possible output as a “class” .....
    - Since the output space is huge, most “classes” do not have any training data.
    - Machine has to create new stuff during testing.
    - Need more intelligence

# Why Structured Learning Challenging?

- Machine has to learn to do *planning*
  - Machine generates objects component-by-component, but it should have a big picture in its mind.
  - Because the output components have dependency, they should be considered globally.

Image Generation



Sentence Generation

這個婆娘不是人

九天玄女下凡塵



# Structured Learning Approach

## **Generator**

Learn to generate  
the object at the  
component level



## **Discriminator**

Evaluating the  
whole object, and  
find the best one



# Outline

Basic Idea of GAN

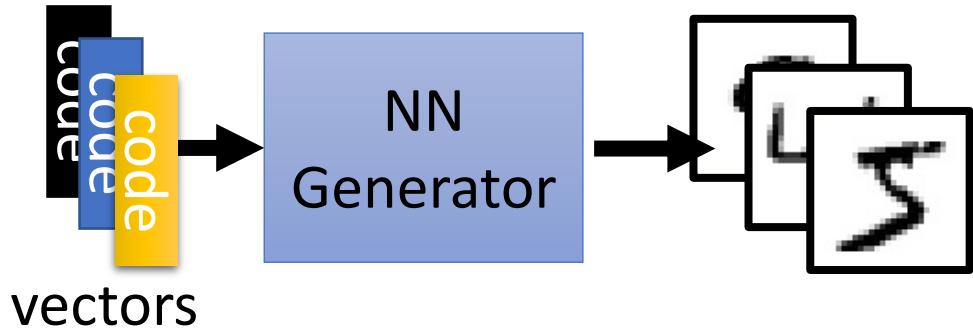
GAN as structured learning

Can Generator learn by itself?

Can Discriminator generate?

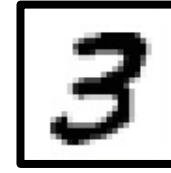
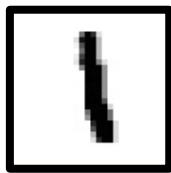
A little bit theory

# Generator



code:  
(where does they  
come from?)

Image:

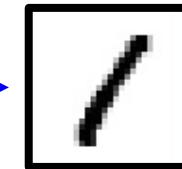


$$\begin{bmatrix} 0.1 \\ 0.9 \end{bmatrix}$$

NN  
Generator

As close as possible

image



c.f.



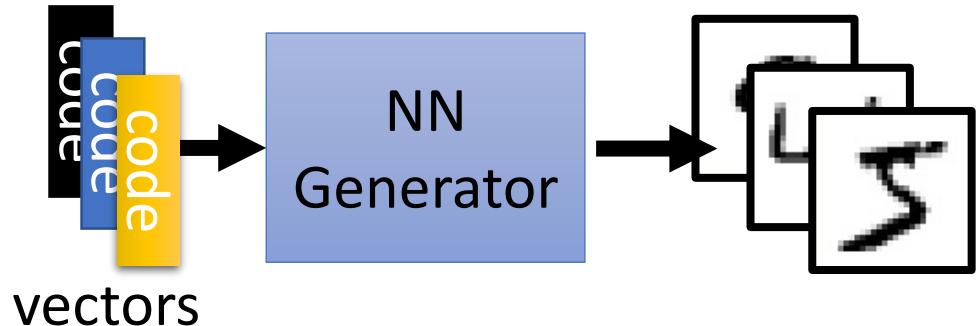
NN  
Classifier

As close as possible

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 0 \\ \vdots \end{bmatrix}$$

# Generator



code:  
[ 0.1  
-0.5 ]

Image:

[ 0.1  
0.9 ]



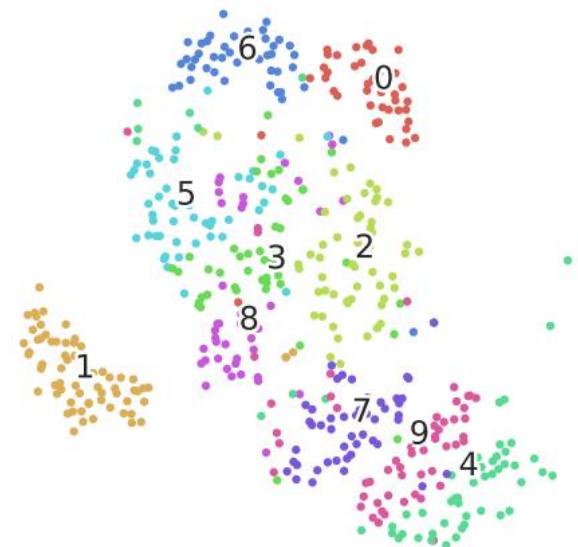
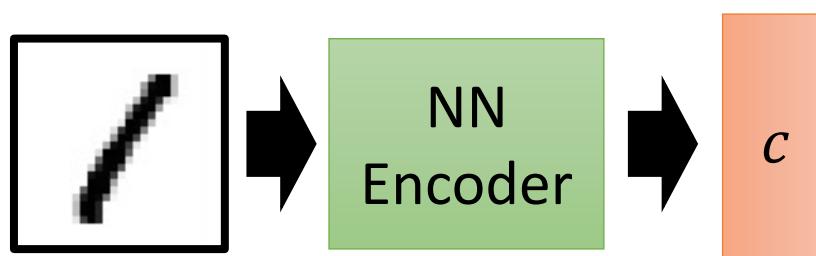
[ 0.2  
-0.1 ]



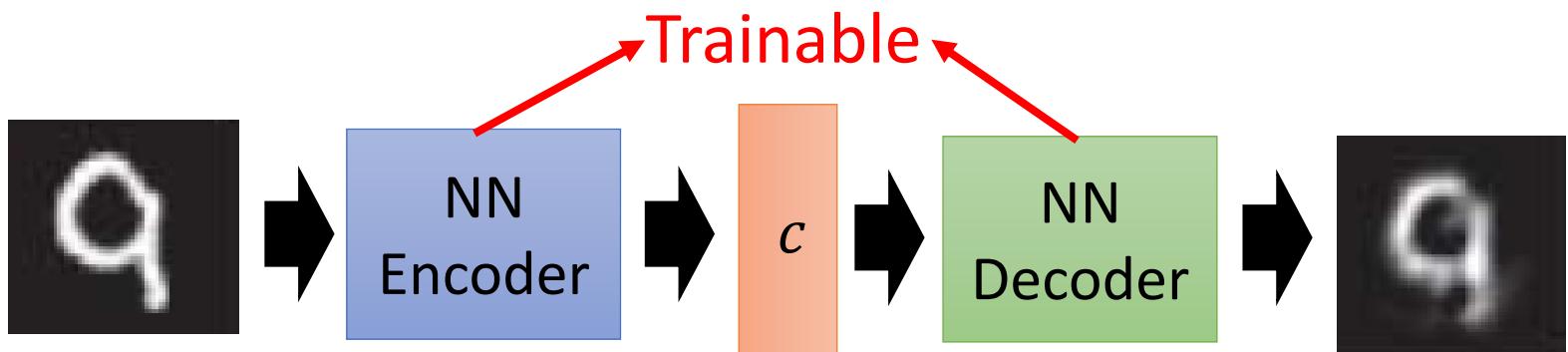
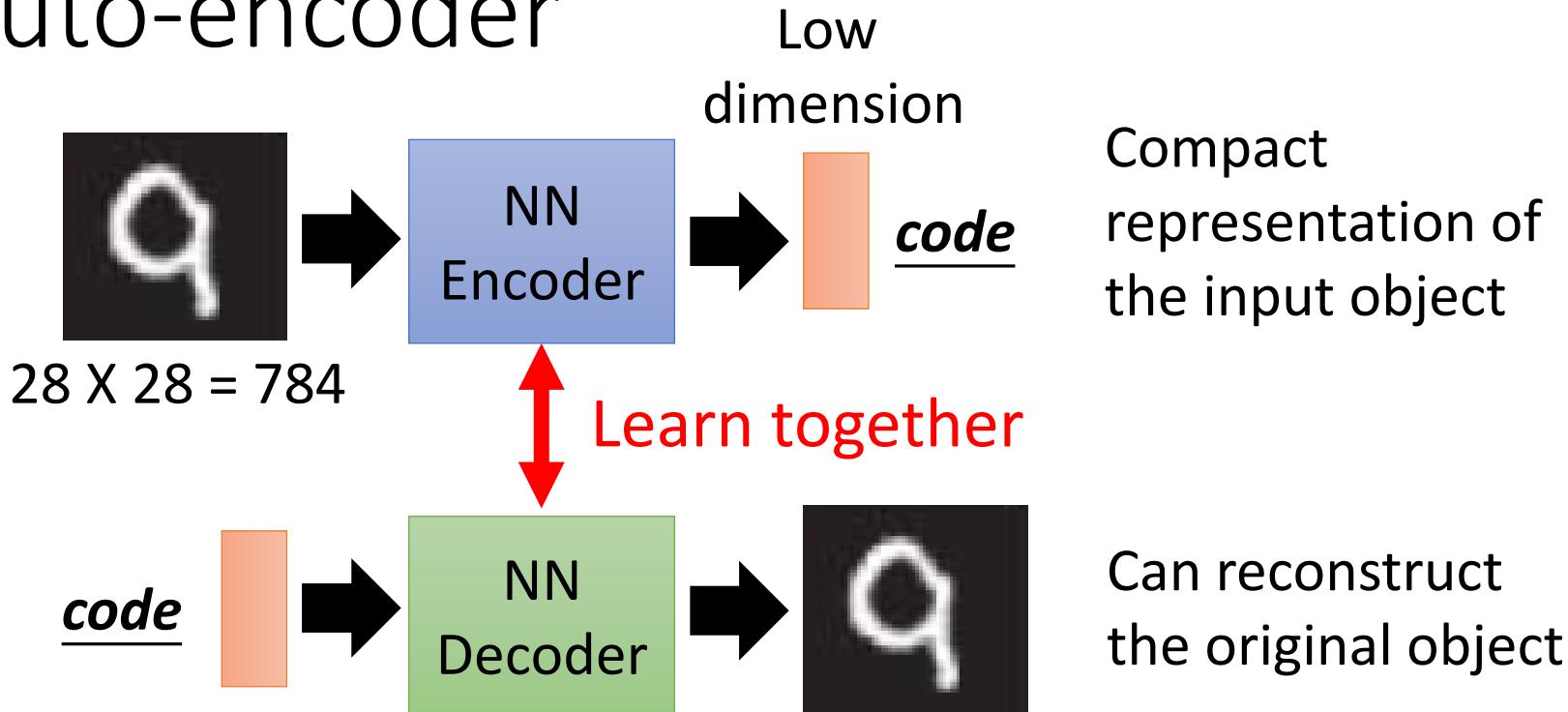
[ 0.3  
0.2 ]



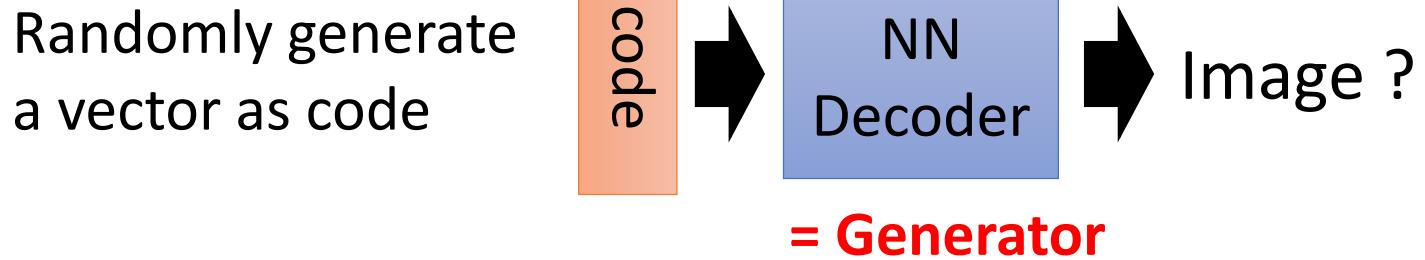
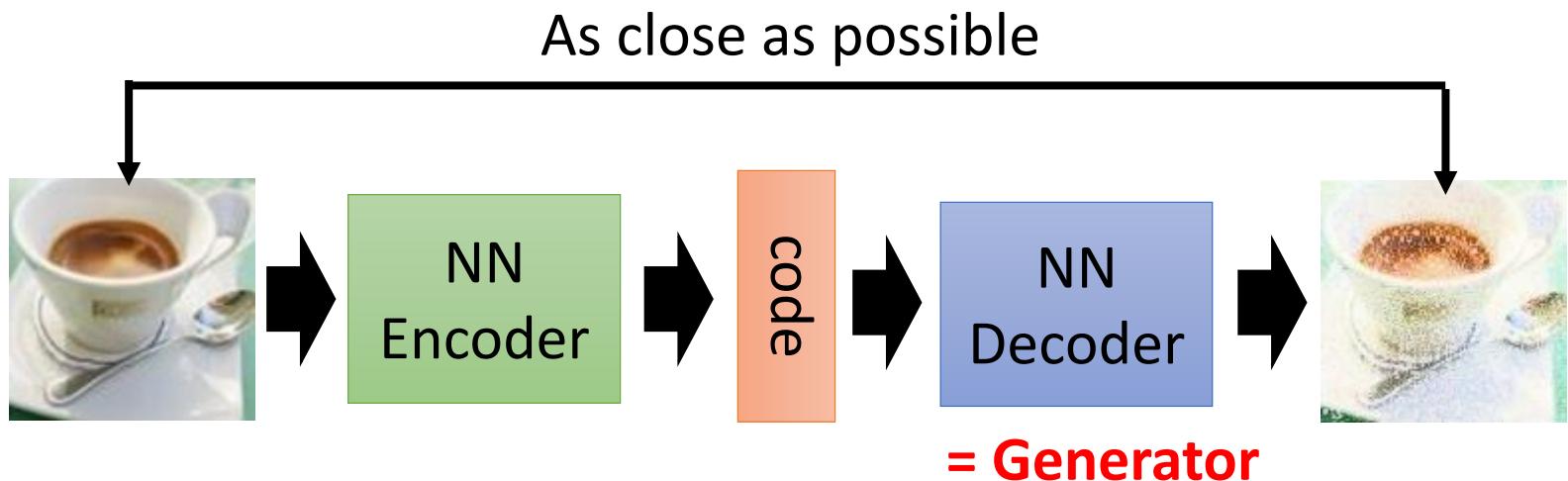
Encoder in auto-encoder  
provides the code 😊



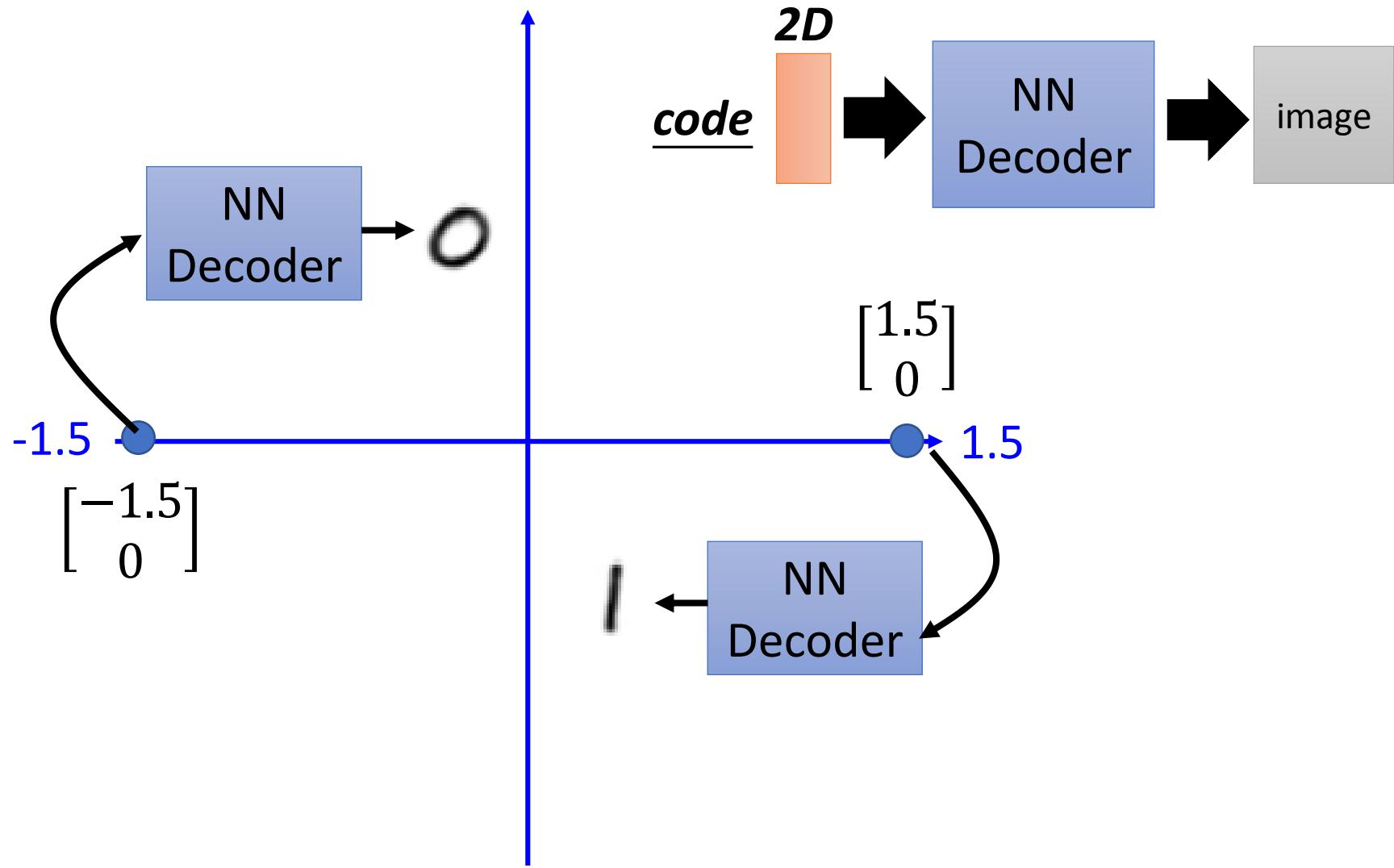
# Auto-encoder



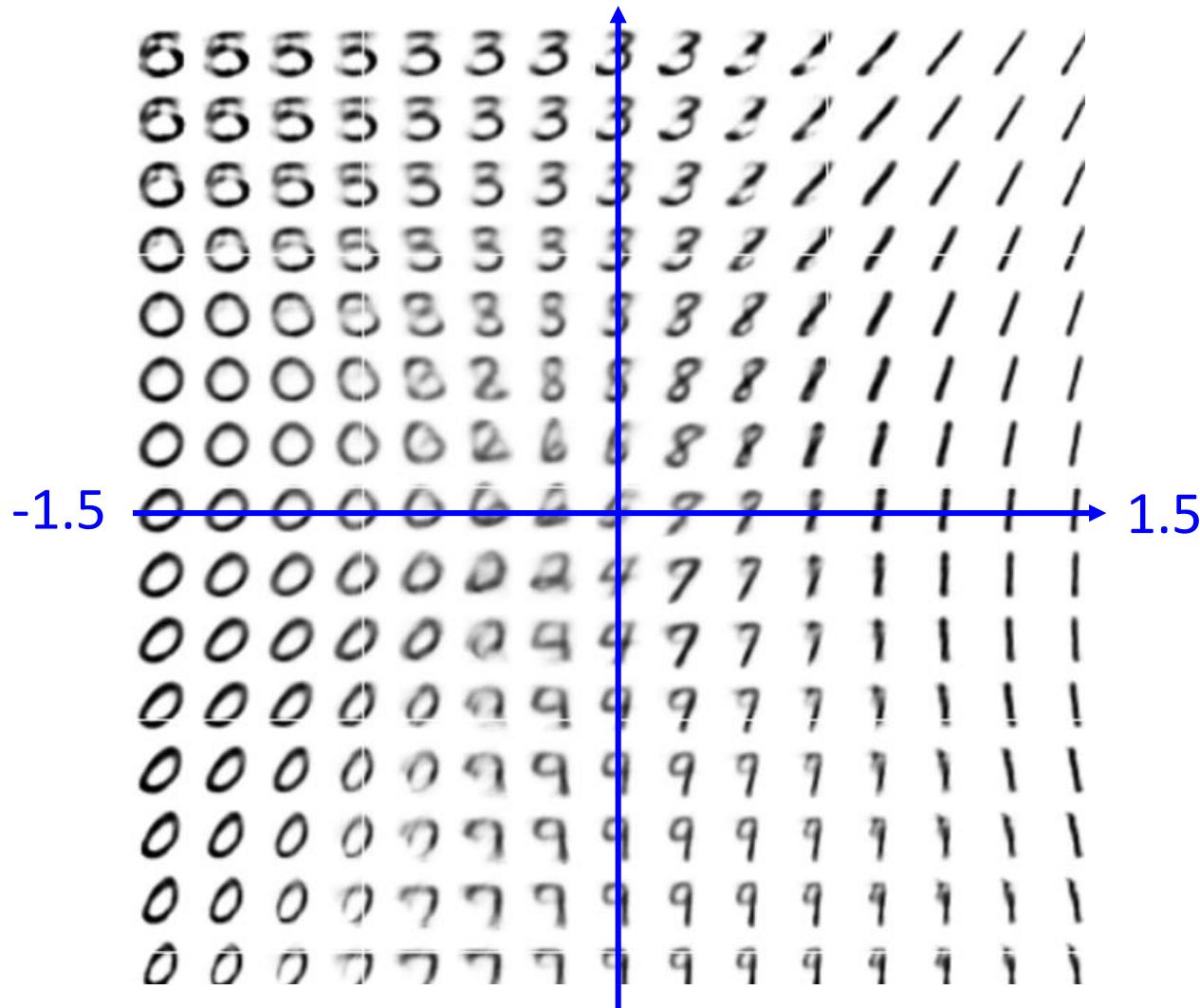
# Auto-encoder



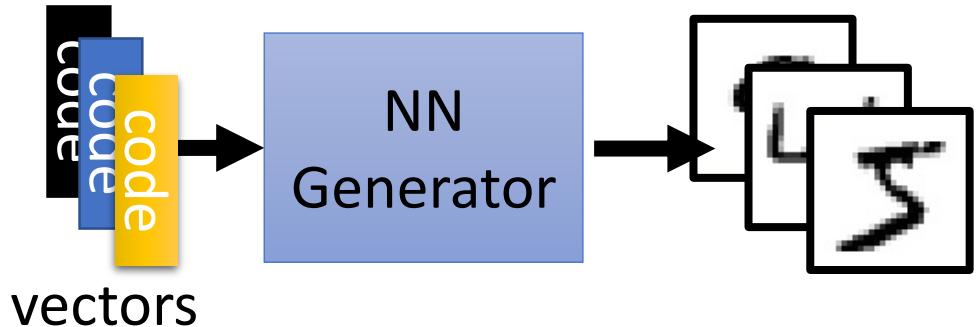
# Auto-encoder



# Auto-encoder



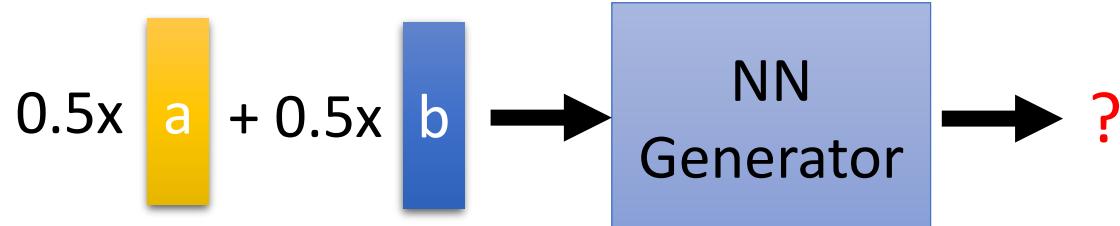
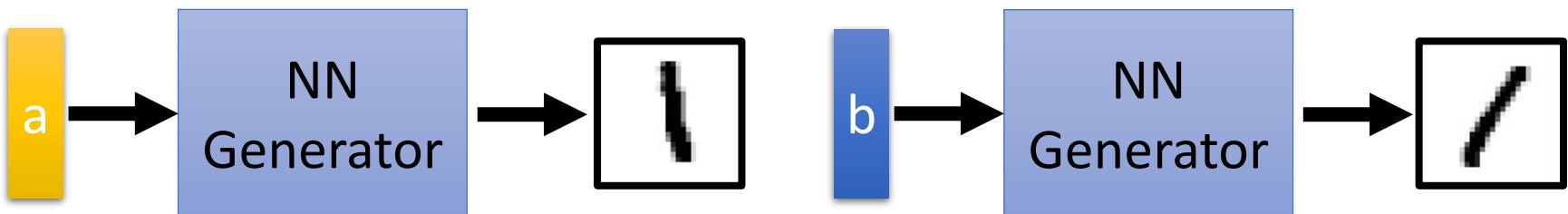
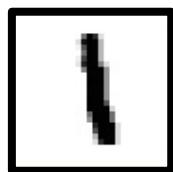
# Auto-encoder



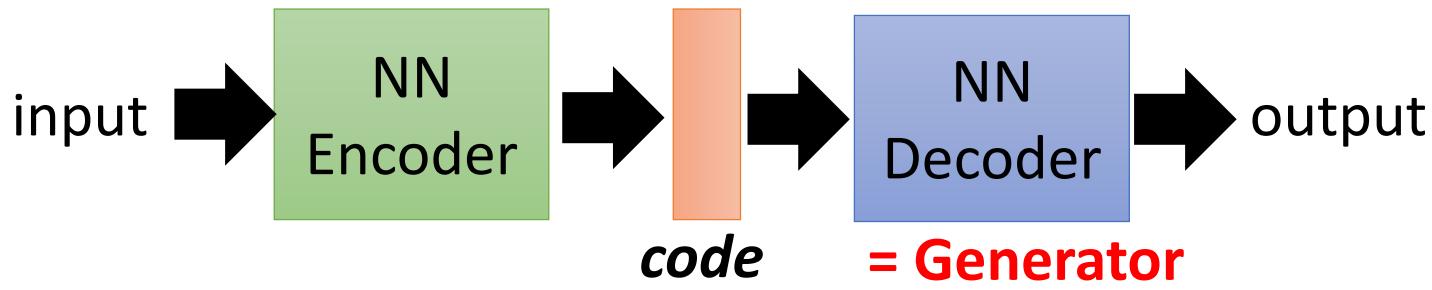
code:  
[ 0.1 ]  
[ -0.5 ]

(where does them  
come from?)

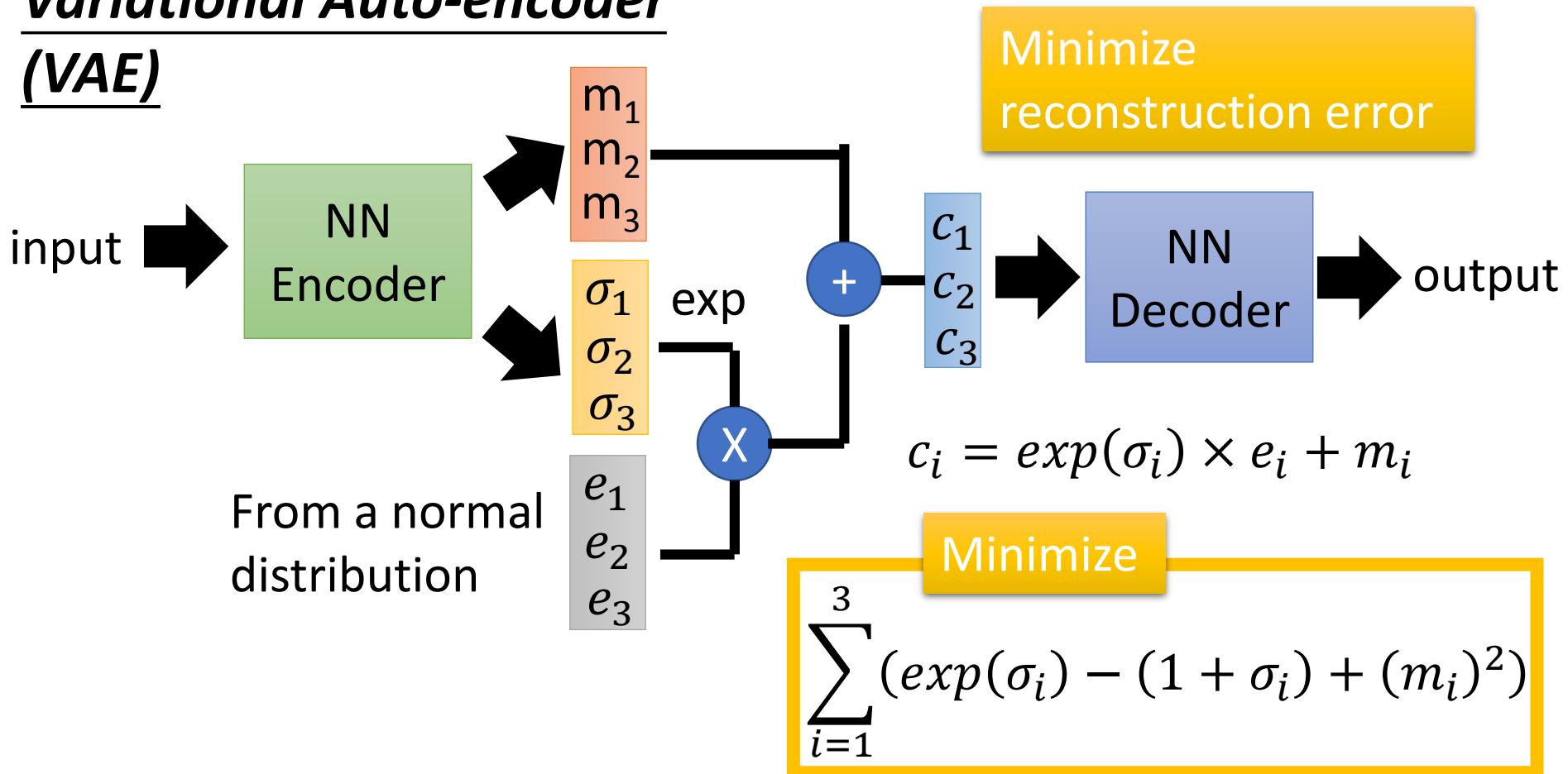
Image:



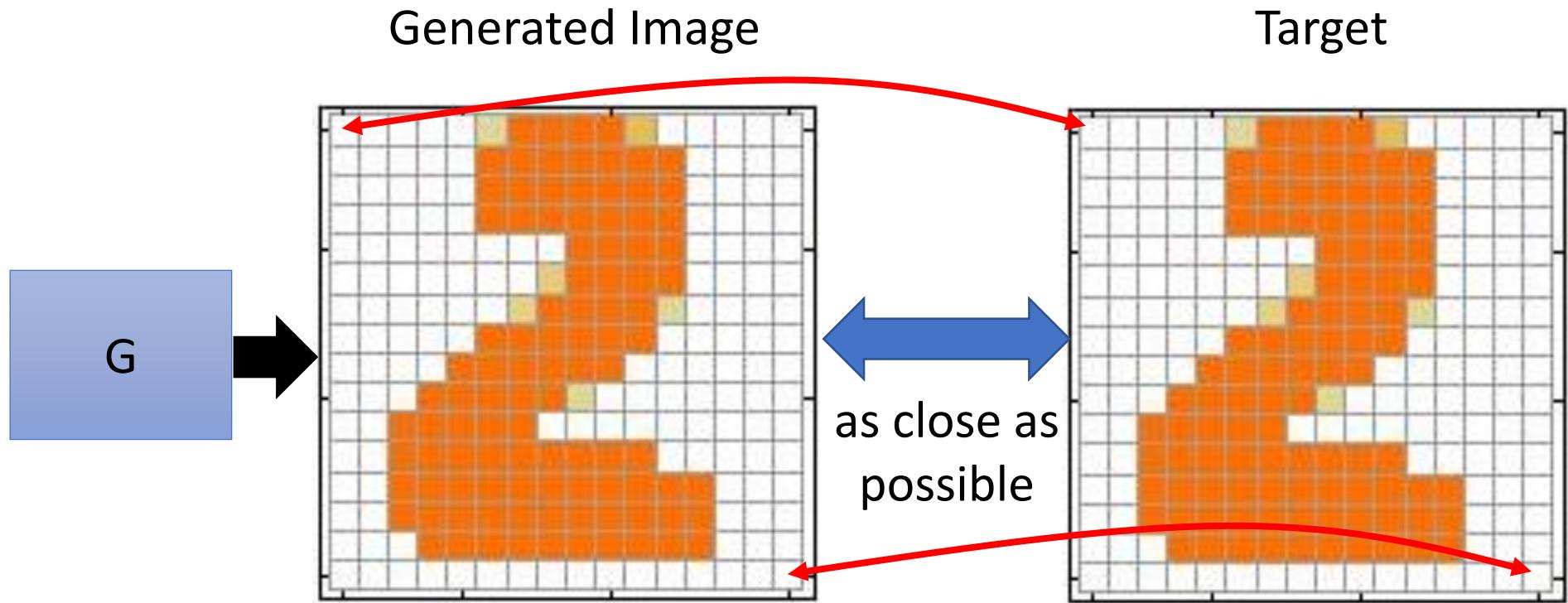
# Auto-encoder



# Variational Auto-encoder (VAE)



# What do we miss?



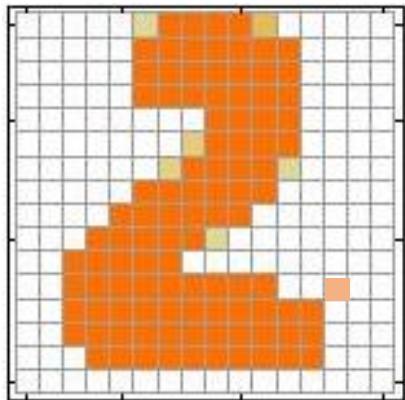
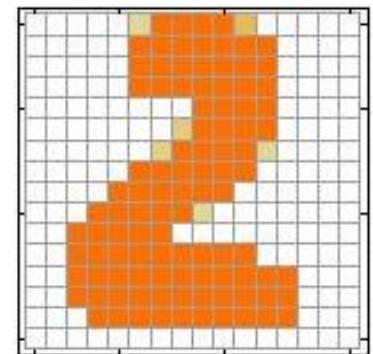
It will be fine if the generator can truly copy the target image.

What if the generator makes some mistakes .....

Some mistakes are serious, while some are fine.

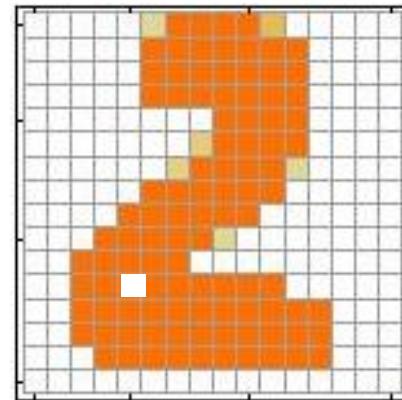
# What do we miss?

Target



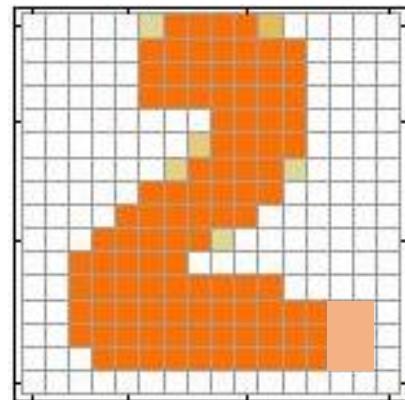
1 pixel error

我覺得不行



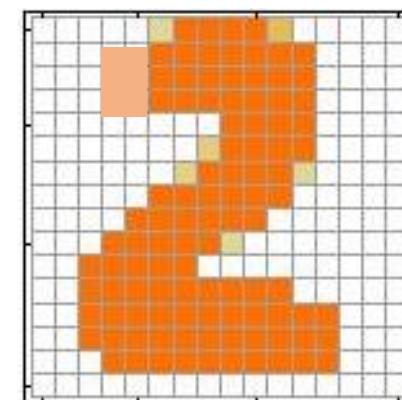
1 pixel error

我覺得不行



6 pixel errors

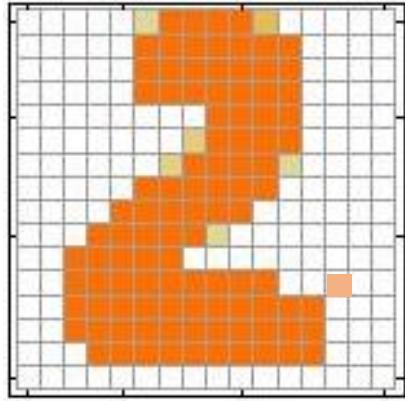
我覺得其實  
可以



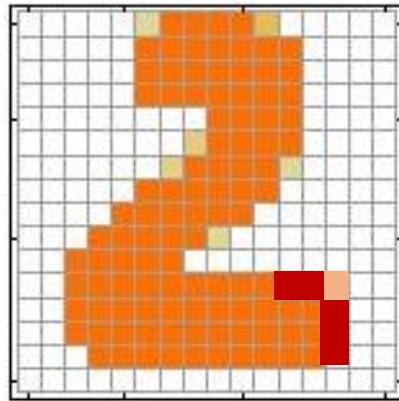
6 pixel errors

我覺得其實  
可以

# What do we miss?

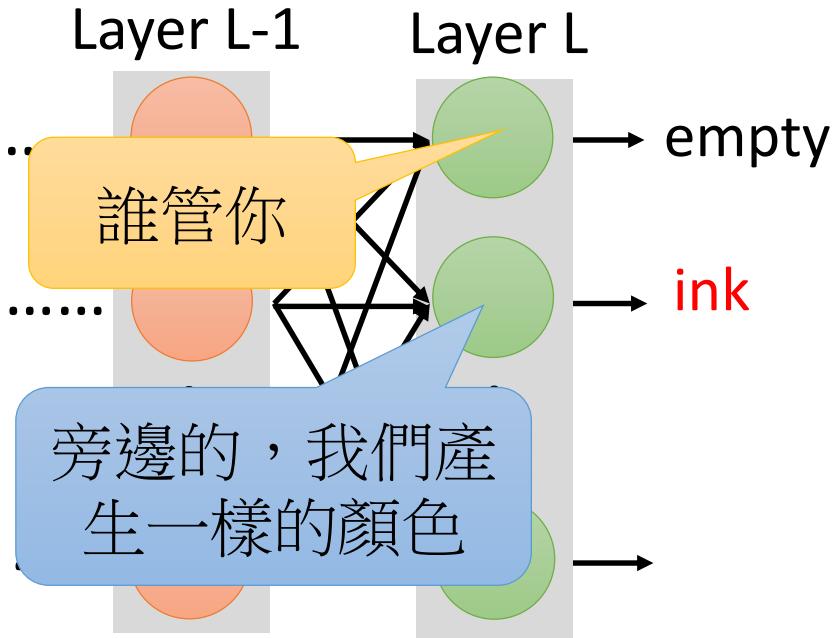


我覺得不行



我覺得其實可以

Each neural in output layer corresponds to a pixel.

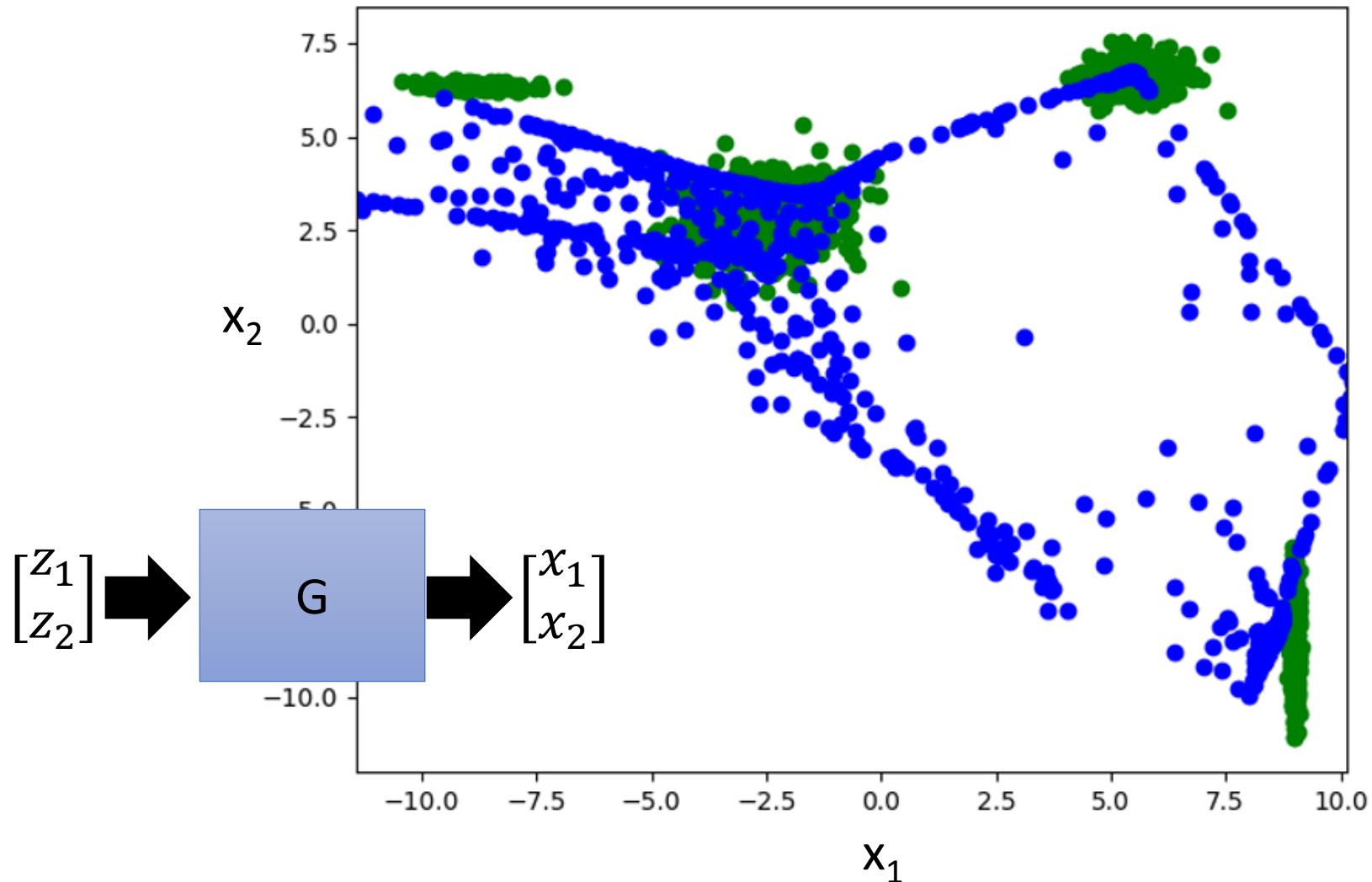


The relation between the components are critical.

Although highly correlated, they cannot influence each other.

Need deep structure to catch the relation between components.

# (Variational) Auto-encoder



# Outline

Basic Idea of GAN

GAN as structured learning

Can Generator learn by itself?

Can Discriminator generate?

A little bit theory

# Discriminator

Evaluation function, Potential Function, Energy Function ...

- Discriminator is a function  $D$  (network, can deep)

$$D: X \rightarrow \mathbb{R}$$

- Input  $x$ : an object  $x$  (e.g. an image)
- Output  $D(x)$ : scalar which represents how “good” an object  $x$  is

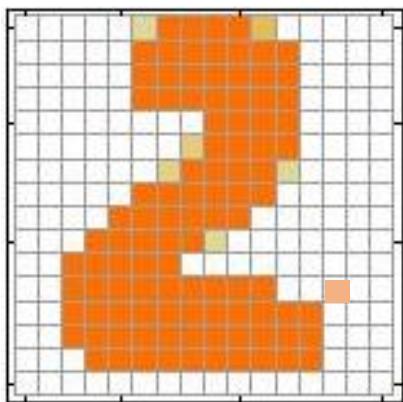


Can we use the discriminator to generate objects?

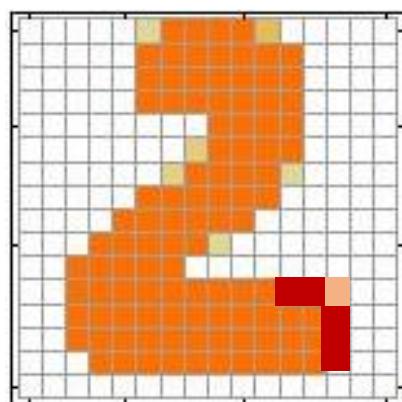
Yes.

# Discriminator

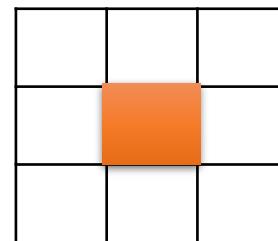
- It is easier to catch the relation between the components by top-down evaluation.



我覺得不行



我覺得其實 OK



This CNN filter is  
good enough.

# Discriminator

- Suppose we already have a good discriminator  
 $D(x)$  ...

Inference

- Generate object  $\tilde{x}$  that

$$\tilde{x} = \arg \max_{x \in X} D(x)$$

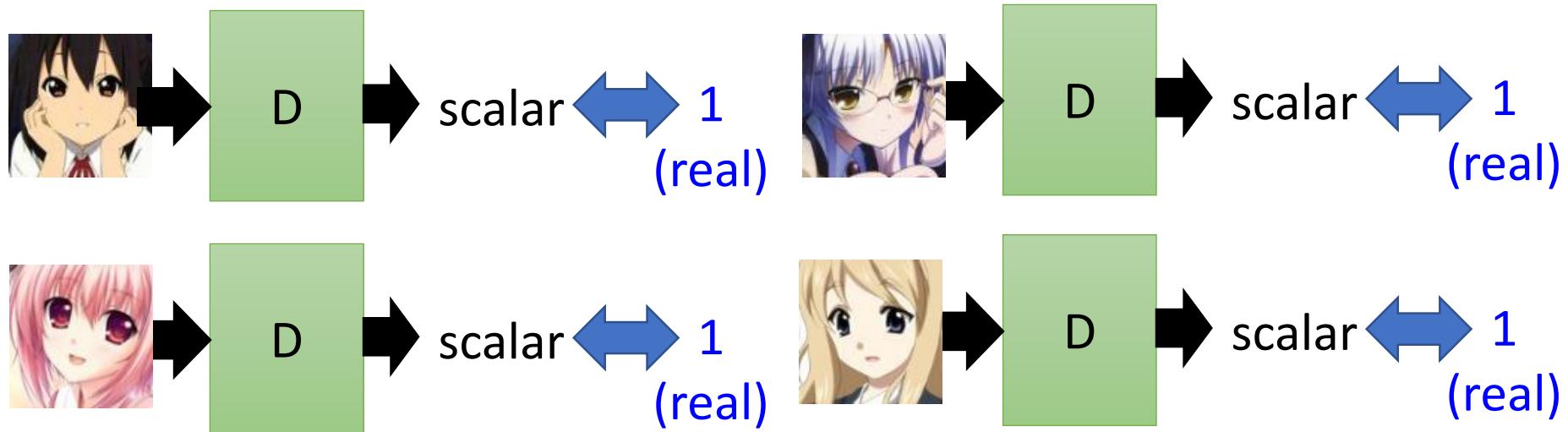
Enumerate all possible  $x$  !!!

It is feasible ???

How to learn the discriminator?

# Discriminator - Training

- I have some real images

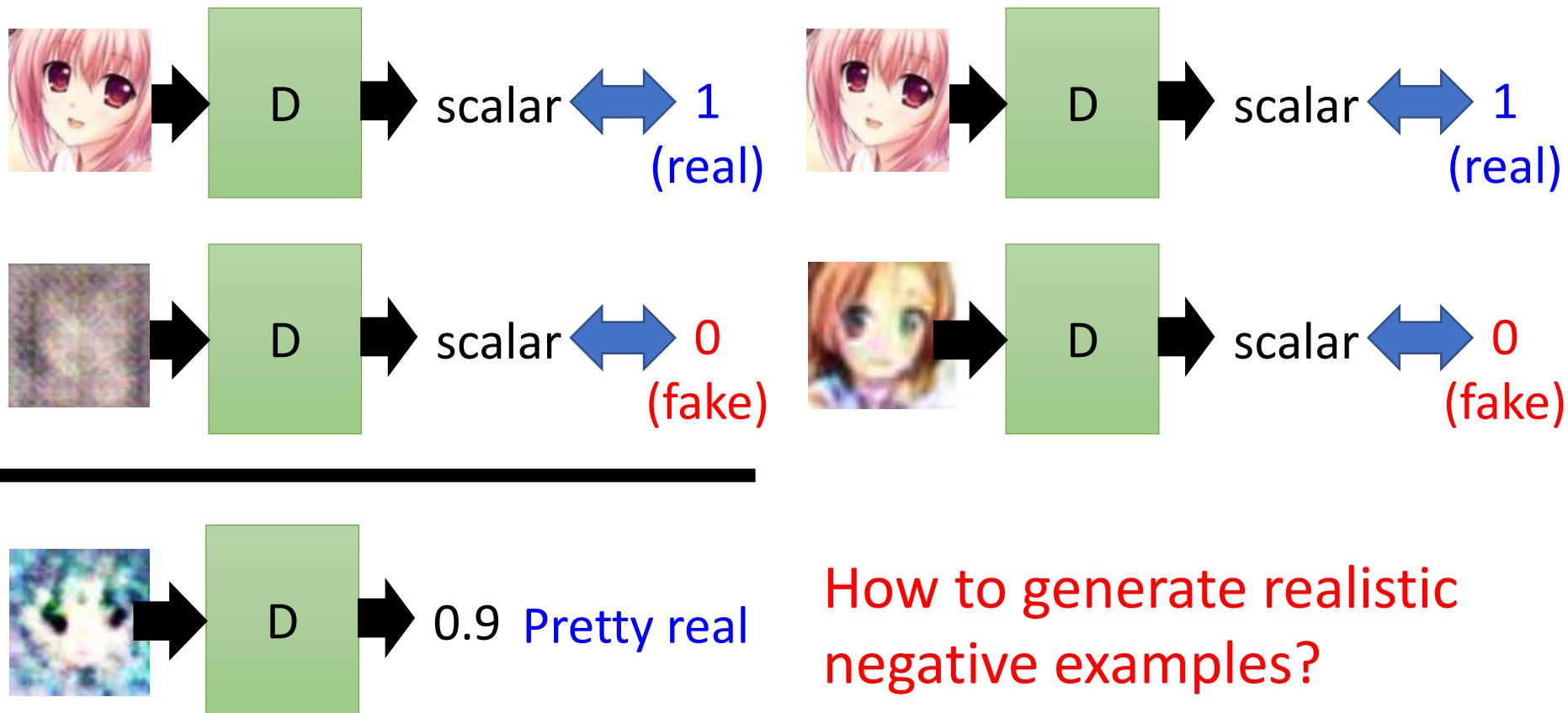


Discriminator only learns to output “1” (real).

Discriminator training needs some negative examples.

# Discriminator - Training

- Negative examples are critical.



# Discriminator - Training

- General Algorithm

- Given a set of **positive examples**, randomly generate a set of **negative examples**.

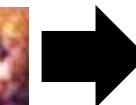


- In each iteration

- Learn a discriminator D that can discriminate positive and negative examples.



v.s.



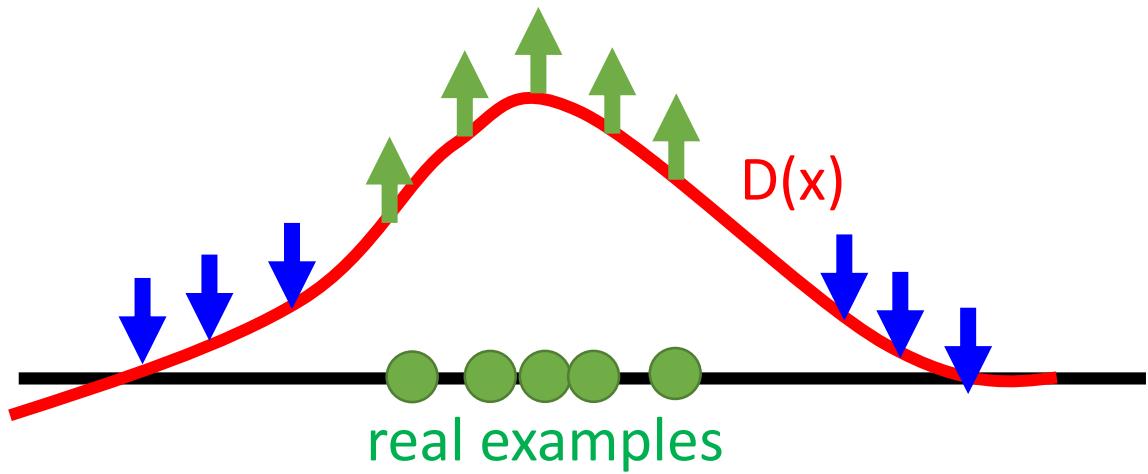
D

- Generate negative examples by discriminator D



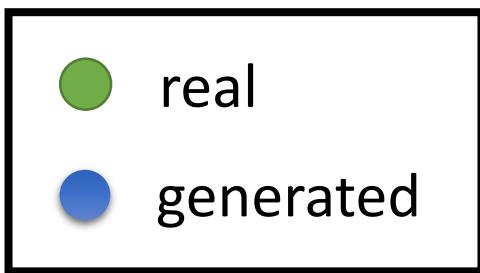
$$\tilde{x} = \arg \max_{x \in X} D(x)$$

# Discriminator - Training

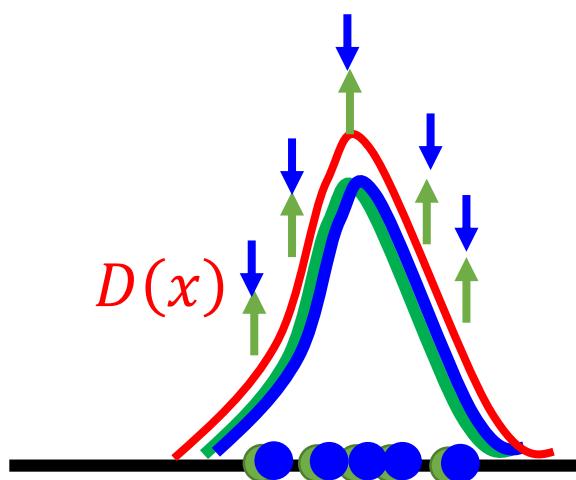
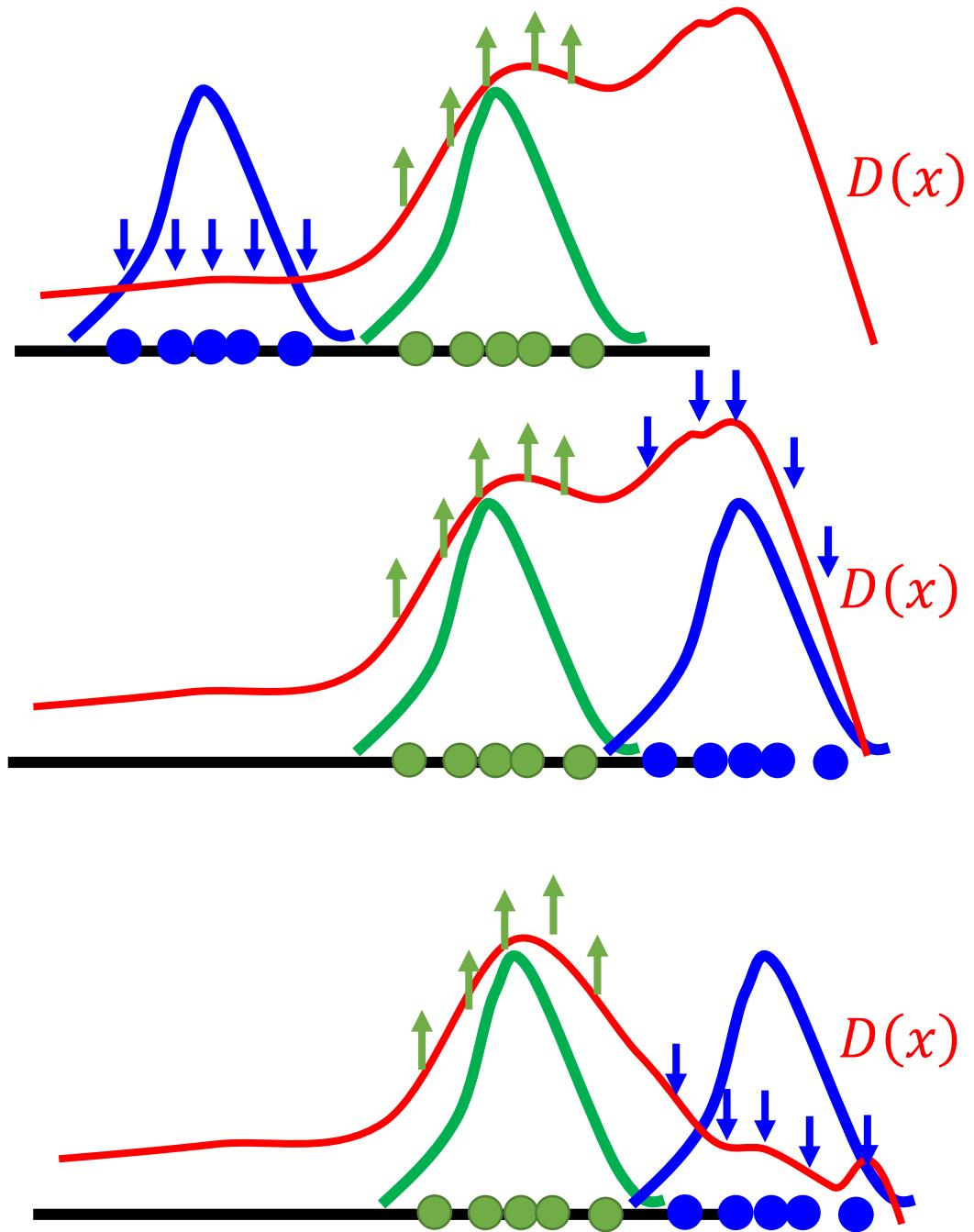


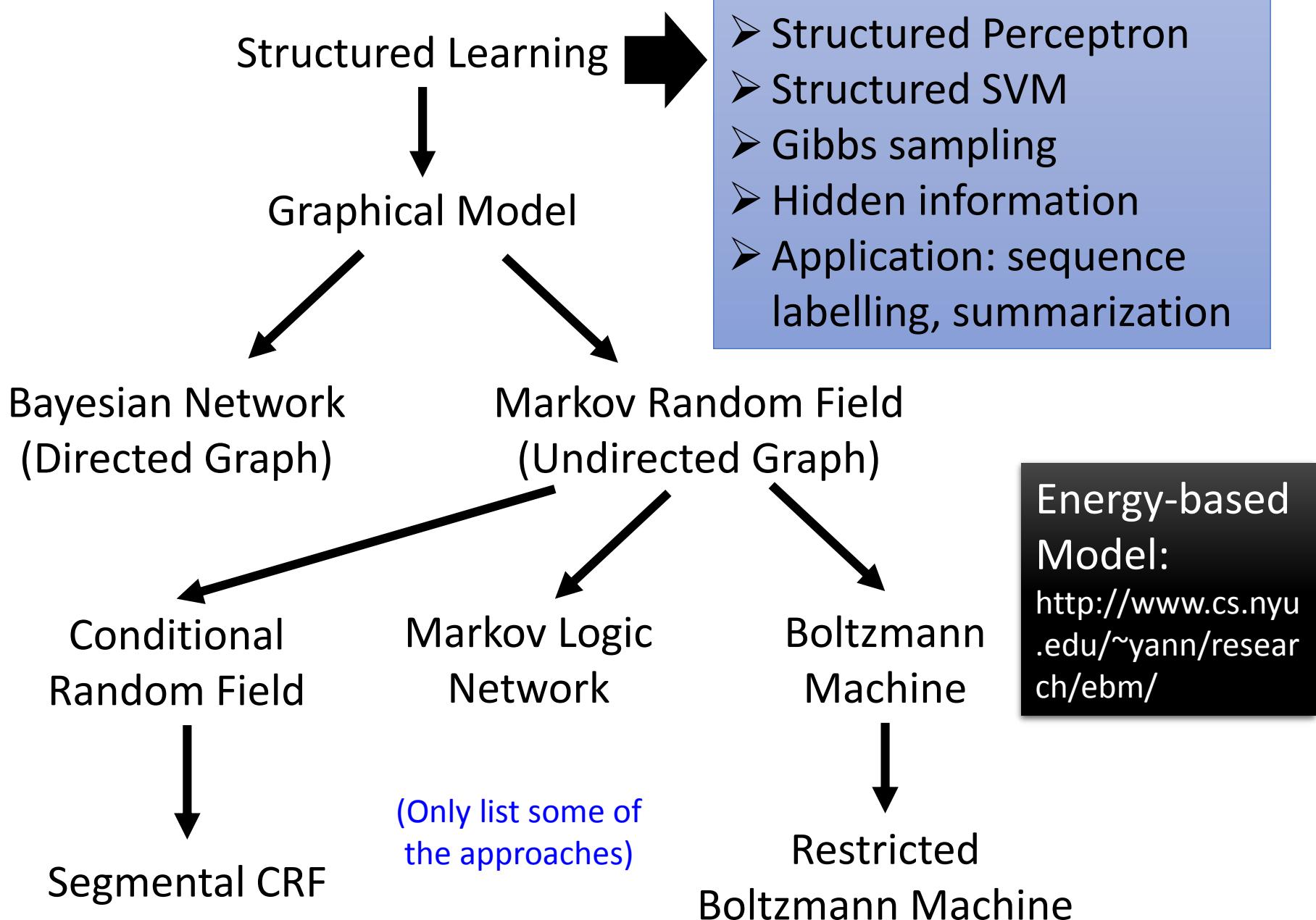
In practice, you cannot decrease all the  $x$  other than real examples.

# Discriminator - Training



In the end .....





# Generator v.s. Discriminator

- **Generator**

- Pros:
  - Easy to generate even with deep model
- Cons:
  - Imitate the appearance
  - Hard to learn the correlation between components

- **Discriminator**

- Pros:
  - Considering the big picture
- Cons:
  - Generation is not always feasible
    - Especially when your model is deep
  - How to do negative sampling?

# Generator + Discriminator

- General Algorithm

- Given a set of **positive examples**, randomly generate a set of **negative examples**.



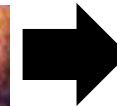
- In each iteration



- Learn a discriminator D that can discriminate positive and negative examples.



v.s.



D

- Generate negative examples by discriminator D

$$\boxed{G \rightarrow \tilde{x}}$$

$$= \boxed{\tilde{x} = \arg \max_{x \in X} D(x)}$$

# Benefit of GAN

- From Discriminator's point of view
  - Using generator to generate negative samples

$$\boxed{\begin{array}{c} \text{G} \\ \longrightarrow \tilde{x} \end{array}} = \boxed{\tilde{x} = \arg \max_{x \in X} D(x)}$$

efficient

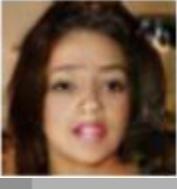
- From Generator's point of view
  - Still generate the object component-by-component
  - But it is learned from the discriminator with global view.

# GAN

VAE

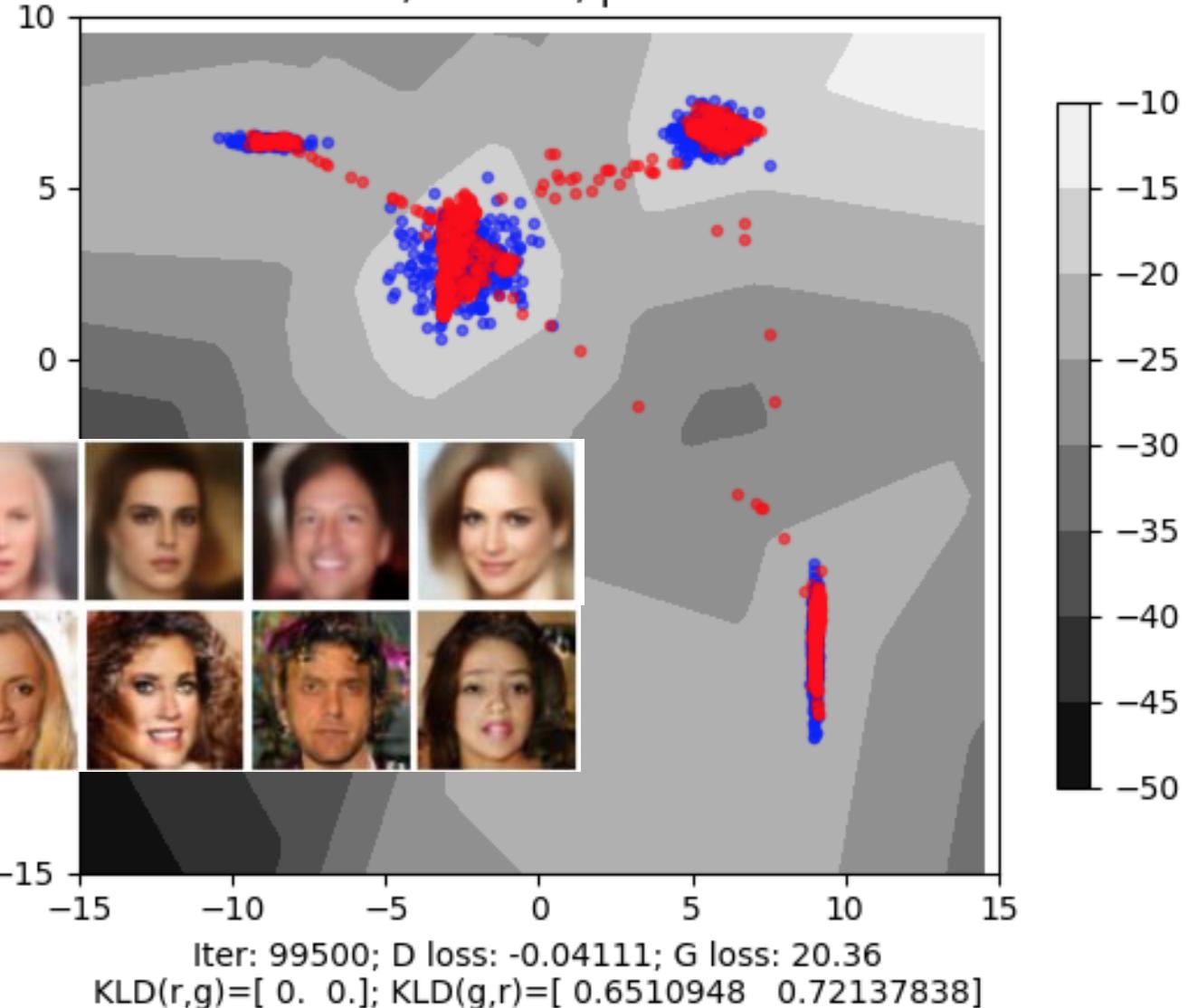


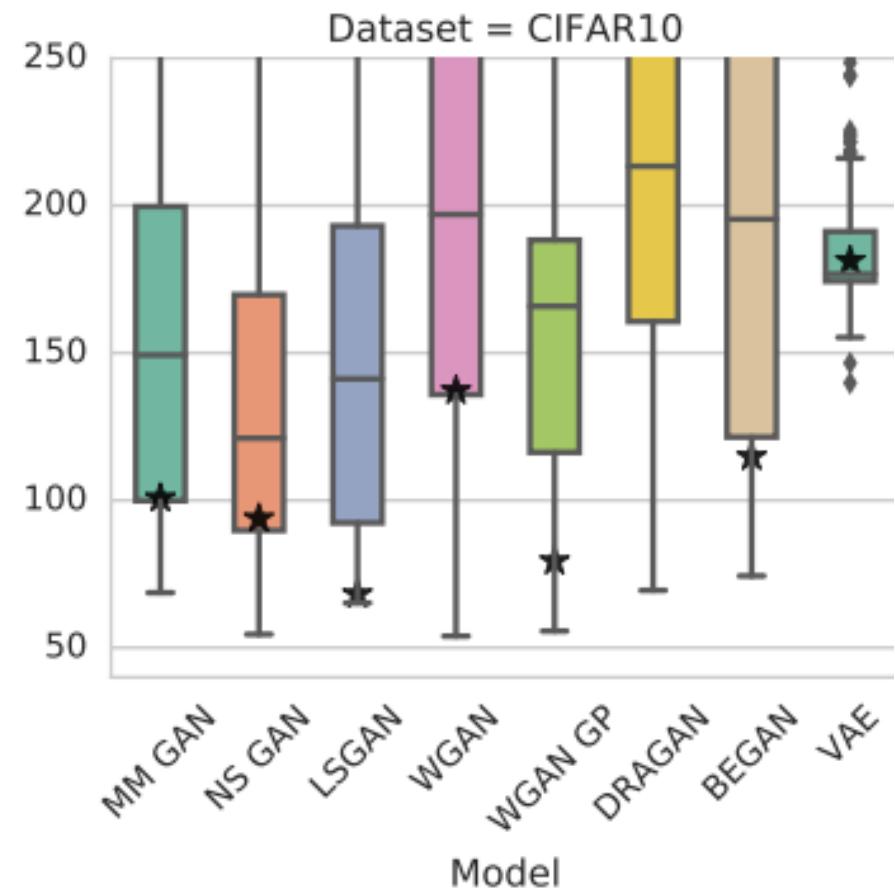
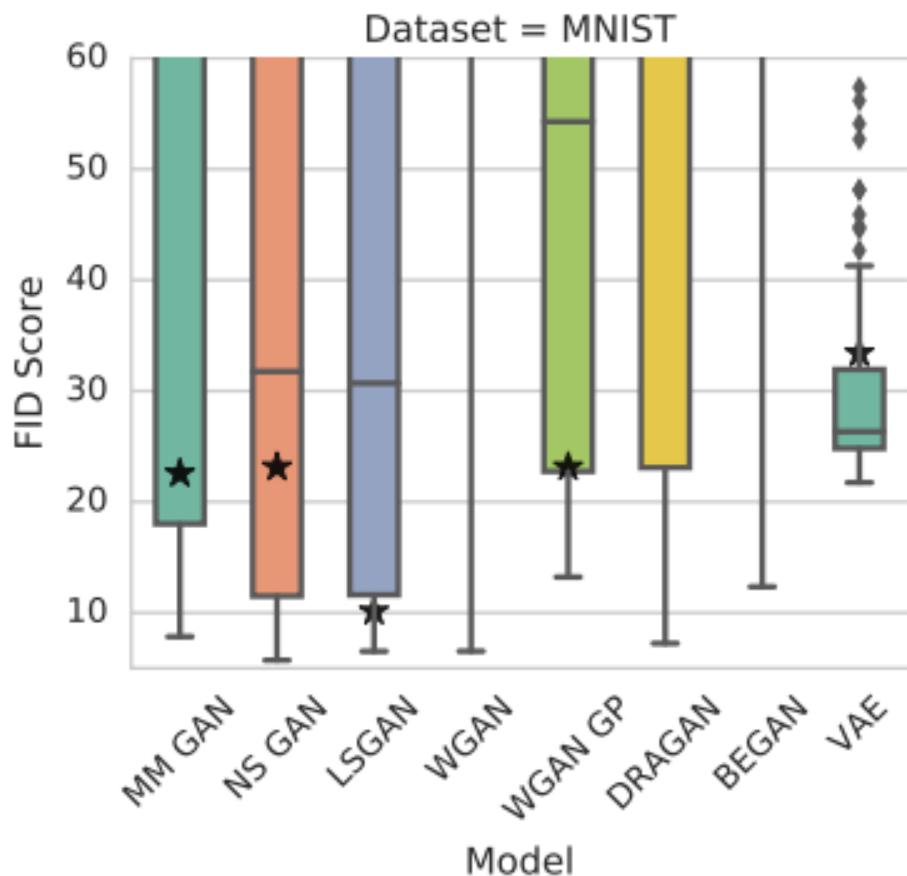
GAN



<https://arxiv.org/abs/1512.09300>

wgan-gp-sub1000-gauss4  
Samples and Decision Boundary  
G: 2\*20; D: 4\*10; prior dim: 2





FID [[Martin Heusel, et al., NIPS, 2017](#)]: Smaller is better

# Next Time

- Preview
  - <https://youtu.be/0CKeqXI5IY0>
  - <https://youtu.be/KSN4QYgAtao>

# Outline

Basic Idea of GAN

GAN as structured learning

Can Generator learn by itself?

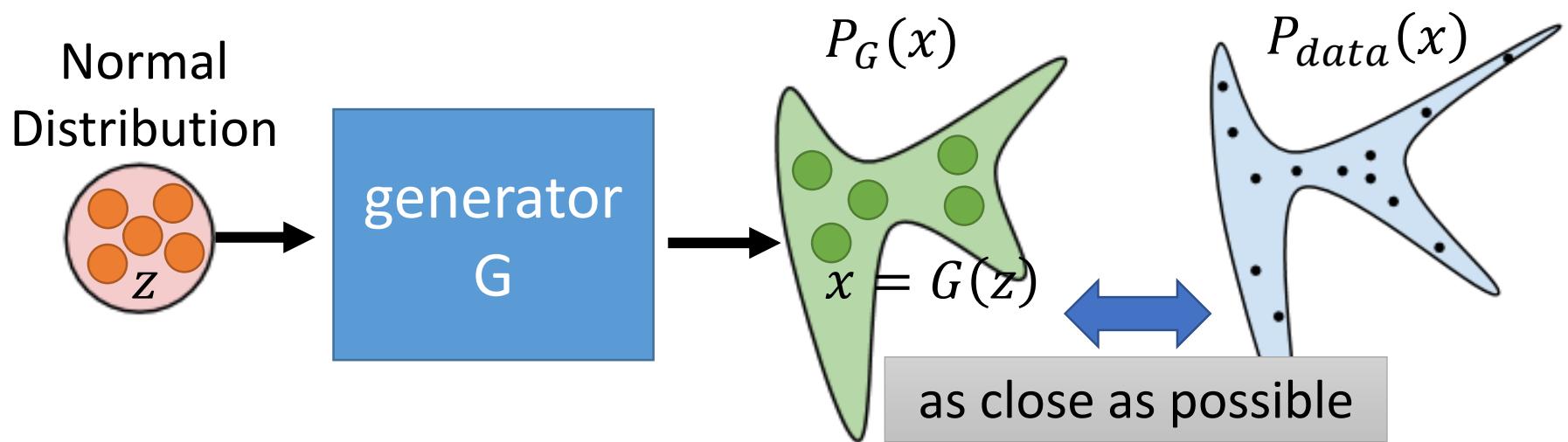
Can Discriminator generate?

A little bit theory

# Generator

$x$ : an image (a high-dimensional vector)

- A generator  $G$  is a network. The network defines a probability distribution  $P_G$



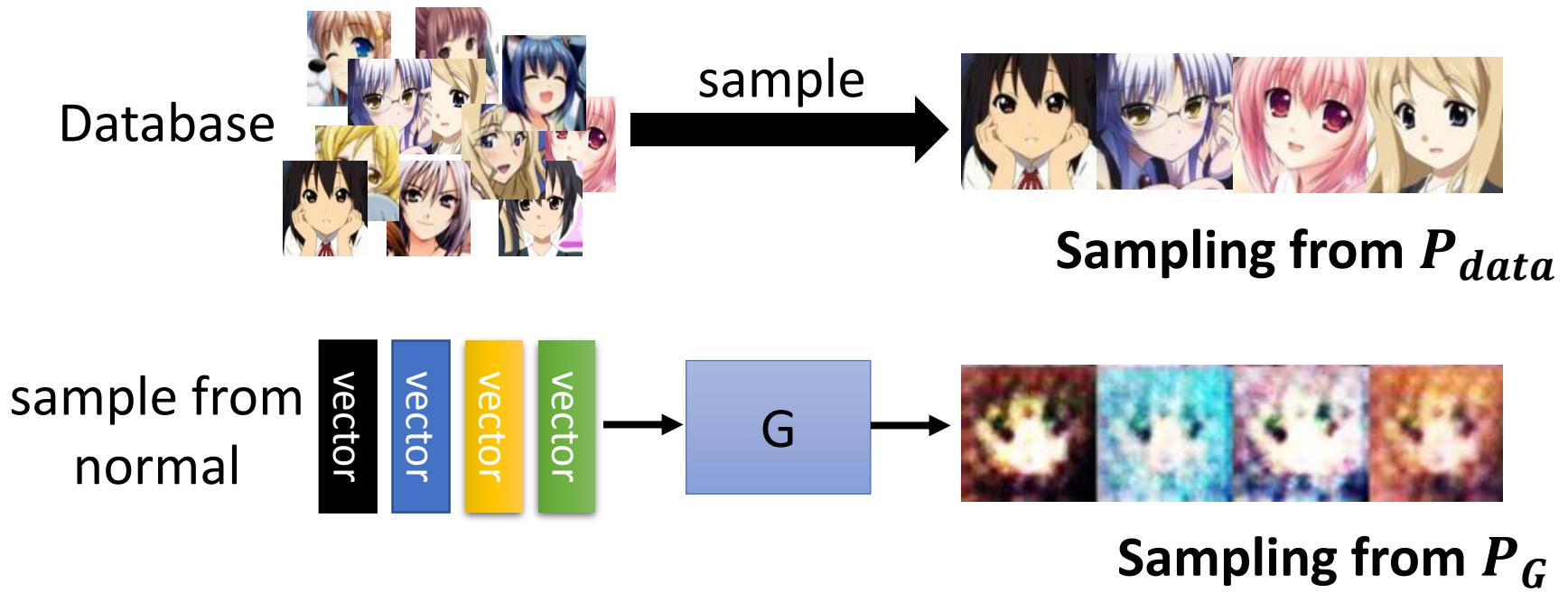
$$G^* = \arg \min_G \underline{Div(P_G, P_{data})}$$

Divergence between distributions  $P_G$  and  $P_{data}$   
How to compute the divergence?

# Discriminator

$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

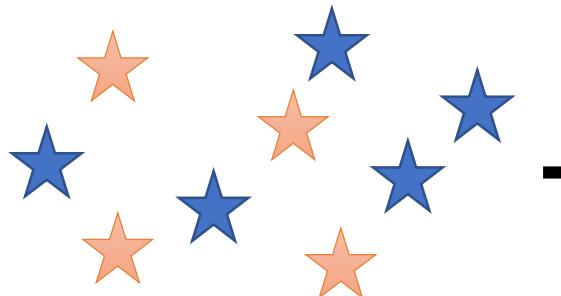
Although we do not know the distributions of  $P_G$  and  $P_{data}$ , we can sample from them.



# Discriminator

$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

- ★ : data sampled from  $P_{data}$
- ☆ : data sampled from  $P_G$



Using the example objective function is exactly the same as training a binary classifier.

## Example Objective Function for D

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

↑  
(G is fixed)

**Training:**  $D^* = \arg \max_D V(D, G)$

The maximum objective value is related to JS divergence.

# Discriminator

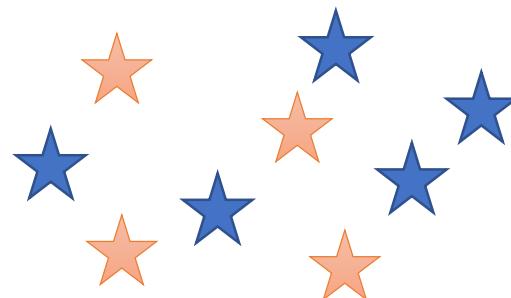
$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

★ : data sampled from  $P_{data}$

☆ : data sampled from  $P_G$

**Training:**

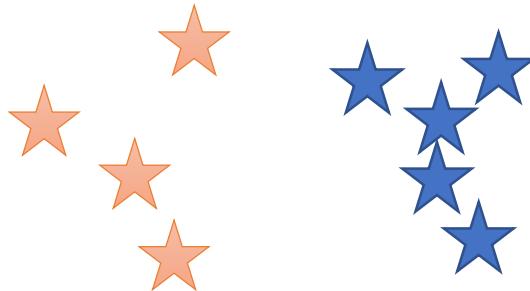
$$D^* = \arg \max_D V(D, G)$$



small divergence

train

Discriminator



large divergence

train

hard to discriminate  
(cannot make objective large)

Discriminator

easy to discriminate

$$G^* = \arg \min_G \max_D V(G, D)$$

$$D^* = \arg \max_D V(D, G)$$

The maximum objective value is related to JS divergence.

- 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  $G$

# Can we use other divergence?

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

Name	Conjugate $f^*(t)$
Total variation	$t$
Kullback-Leibler (KL)	$\exp(t - 1)$
Reverse KL	$-1 - \log(-t)$
Pearson $\chi^2$	$\frac{1}{4}t^2 + t$
Neyman $\chi^2$	$2 - 2\sqrt{1-t}$
Squared Hellinger	$\frac{t}{1-t}$
Jeffrey	$W(e^{1-t}) + \frac{1}{W(e^{1-t})} + t - 2$
Jensen-Shannon	$-\log(2 - \exp(t))$
Jensen-Shannon-weighted	$(1 - \pi) \log \frac{1-\pi}{1-\pi e^{t/\pi}}$
GAN	$-\log(1 - \exp(t))$

Using the divergence  
you like ☺

# Conditional Generation by GAN

李宏毅

Hung-yi Lee

# Text-to-Image

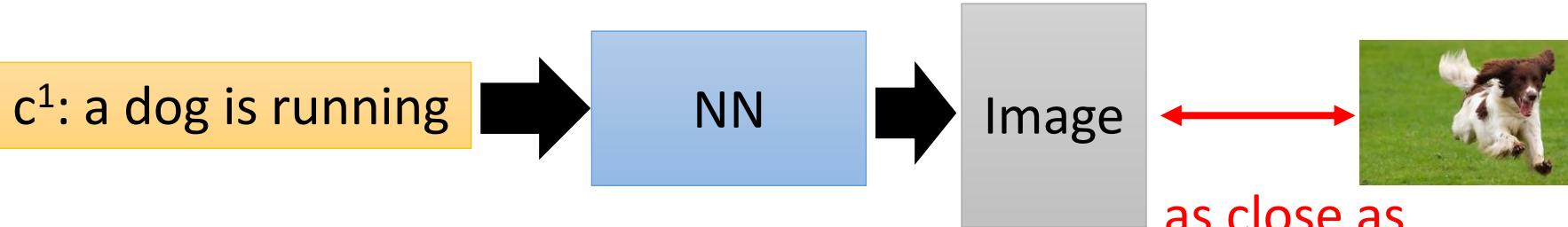
a dog is running



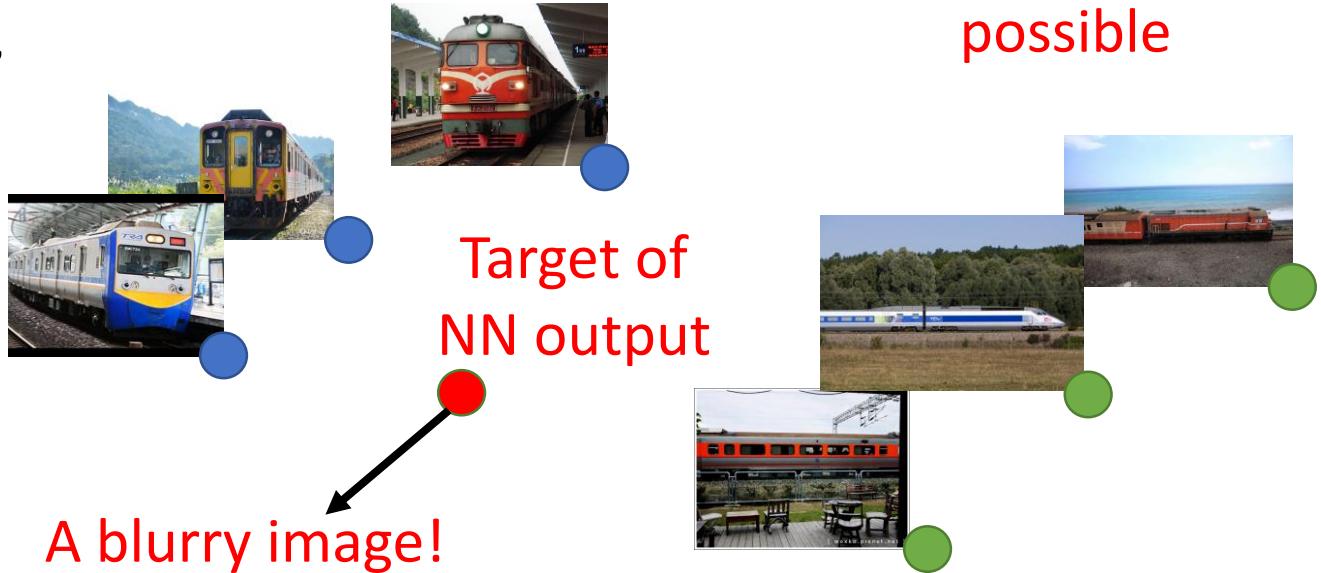
a bird is flying



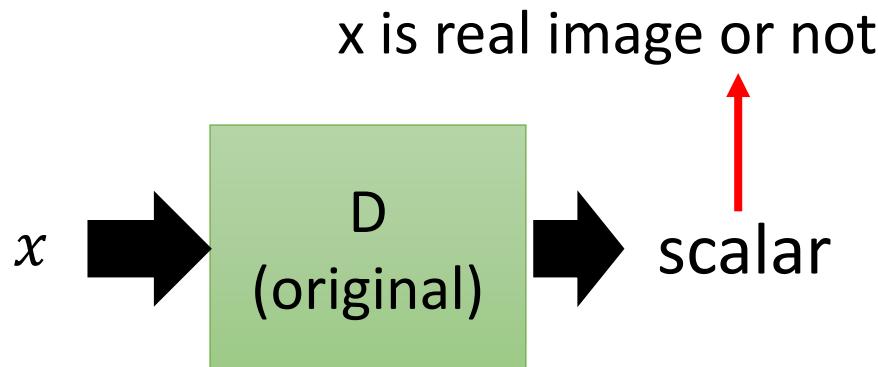
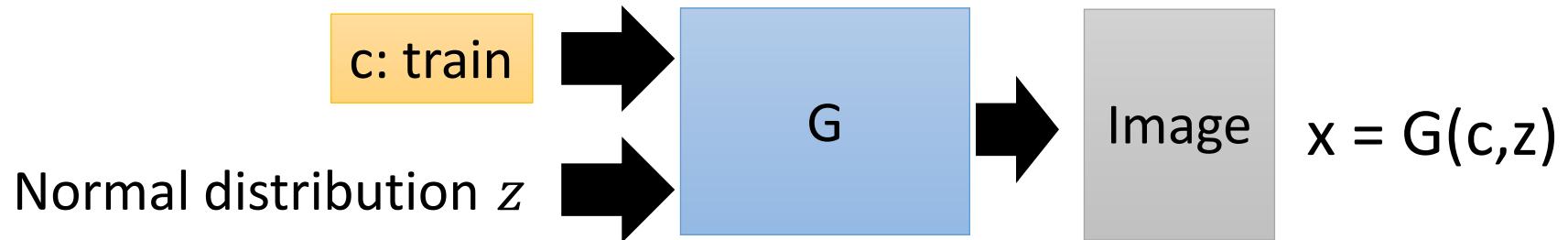
- Traditional supervised approach



Text: “train”

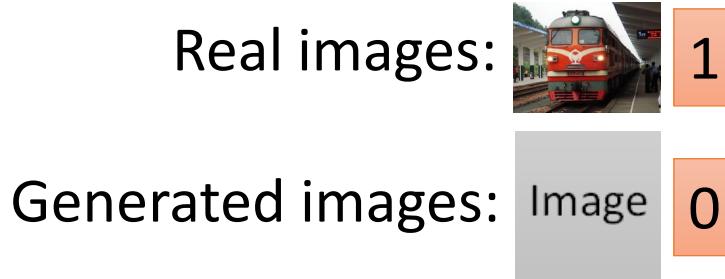


# Conditional GAN

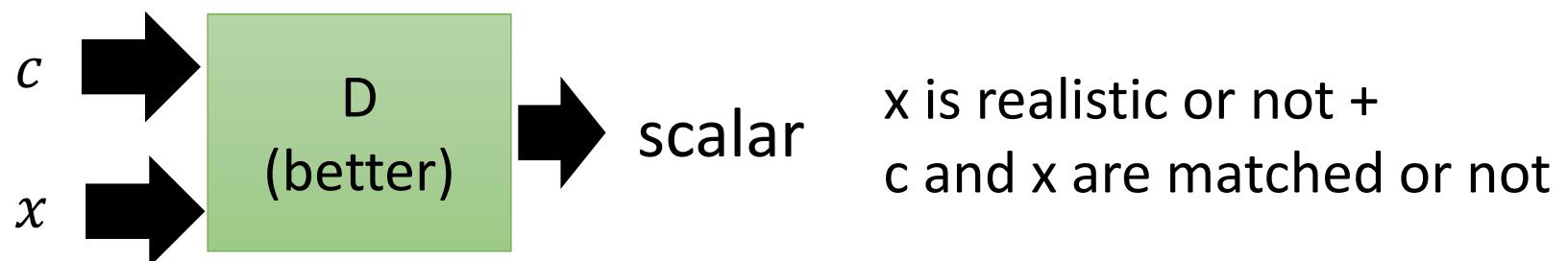
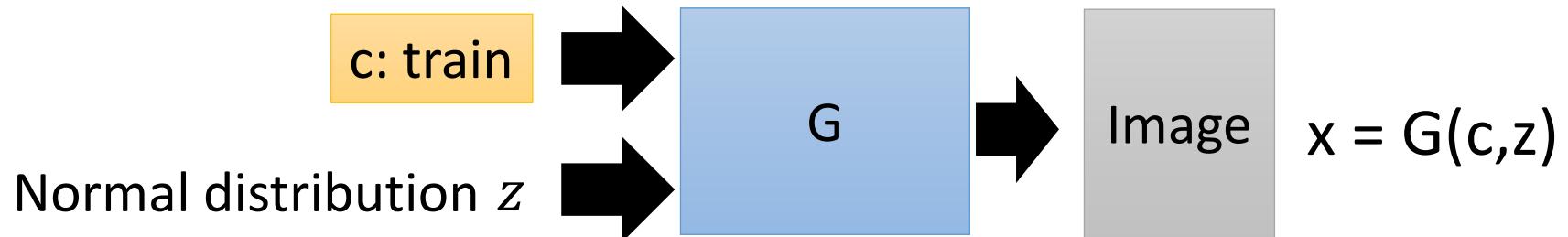


Generator will learn to  
generate realistic images ...

But completely ignore the  
input conditions.



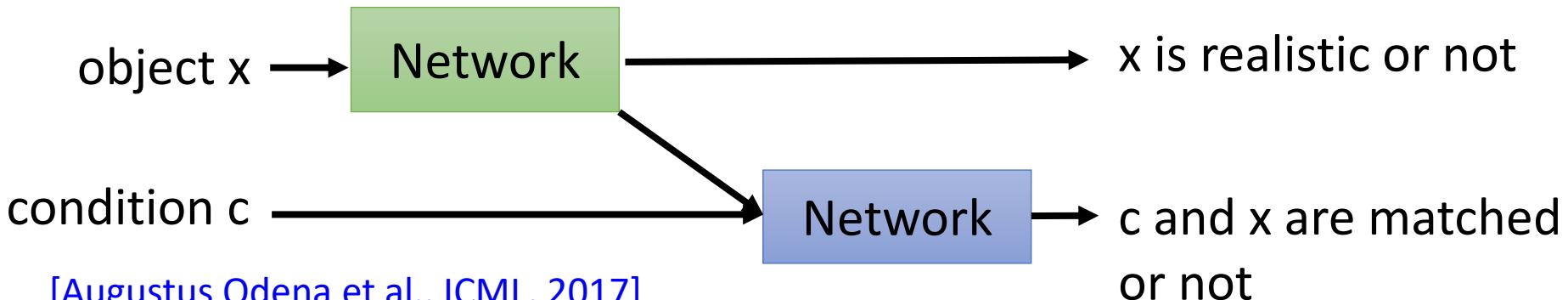
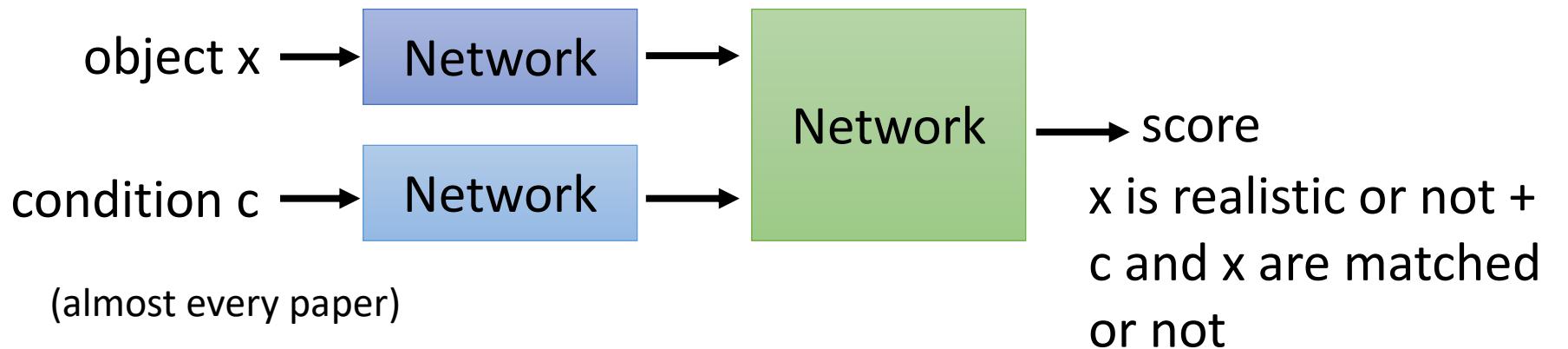
# Conditional GAN



True text-image pairs: (train ,  ) 1

(cat ,  ) 0      (train ,  ) 0

# Conditional GAN - Discriminator



[Augustus Odena et al., ICML, 2017]

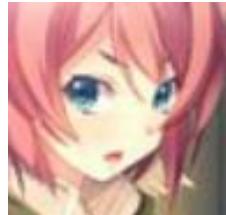
[Takeru Miyato, et al., ICLR, 2018]

[Han Zhang, et al., arXiv, 2017]

# Conditional GAN

The images are generated by  
Yen-Hao Chen, Po-Chun Chien,  
Jun-Chen Xie, Tsung-Han Wu.

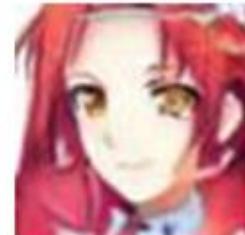
## *paired data*



blue eyes  
red hair  
short hair

Collecting anime faces  
and the description of its  
characteristics

red hair,  
green eyes

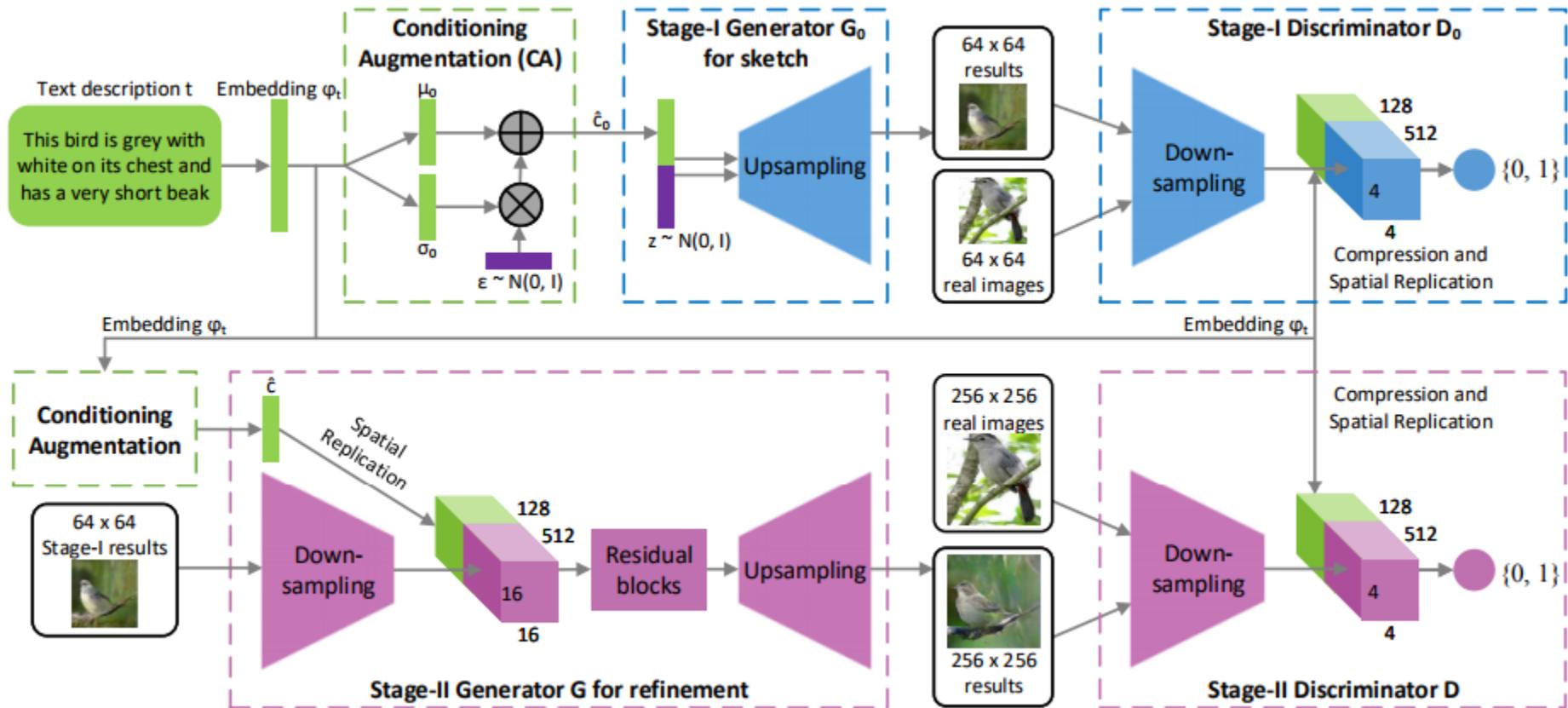


blue hair,  
red eyes

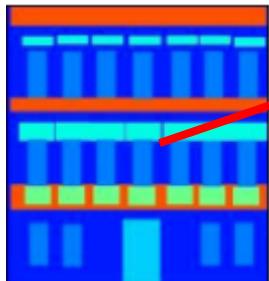


# Stack GAN

Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, Dimitris Metaxas, "StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks", ICCV, 2017



# Image-to-image



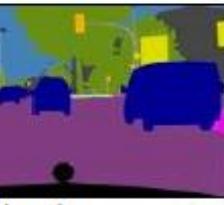
$c$   
 $z$



$$x = G(c, z)$$



Labels to Street Scene

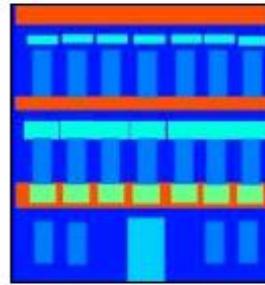


input



output

Labels to Facade



input



output

BW to Color



input



output

Aerial to Map



input

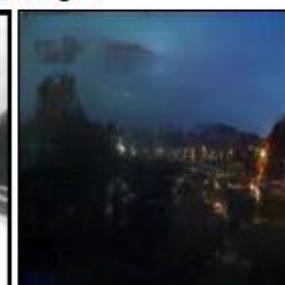


output

Day to Night



input



output

Edges to Photo

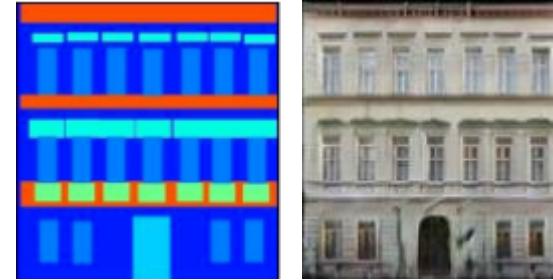


input

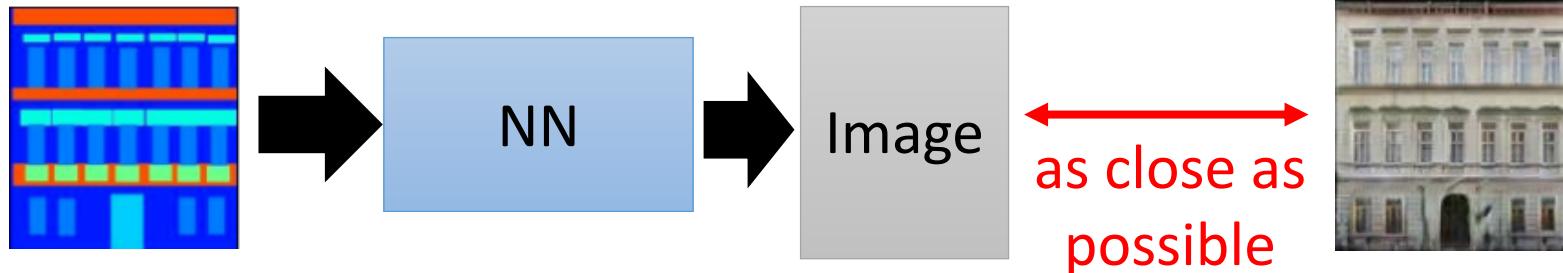


output

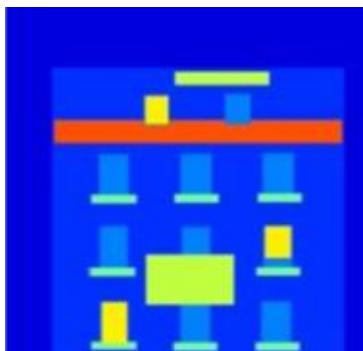
# Image-to-image



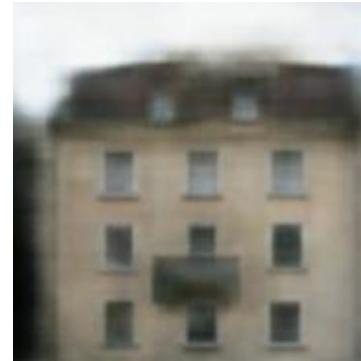
- Traditional supervised approach



Testing:



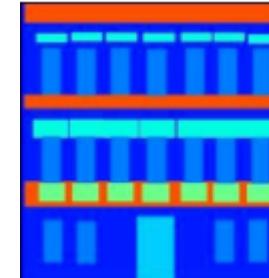
input



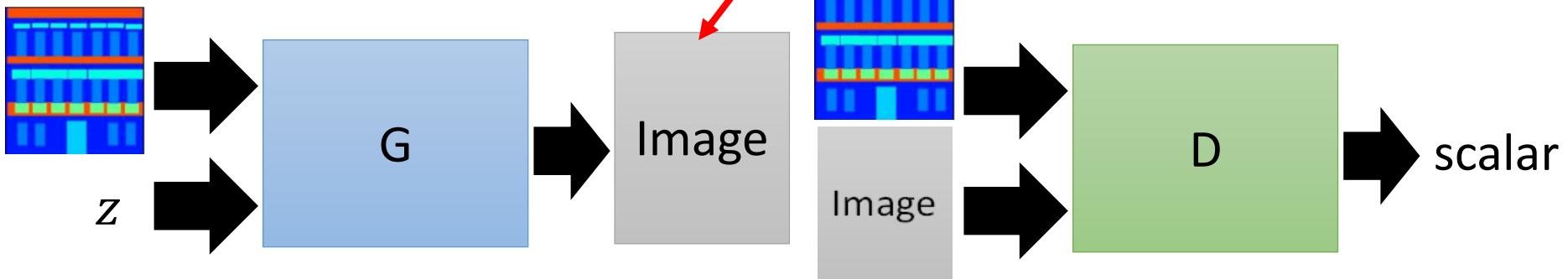
close

It is blurry because it is the average of several images.

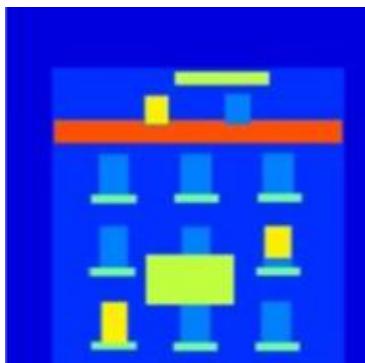
# Image-to-image



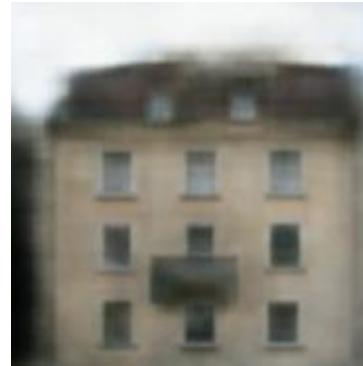
- Experimental results



Testing:



input



close



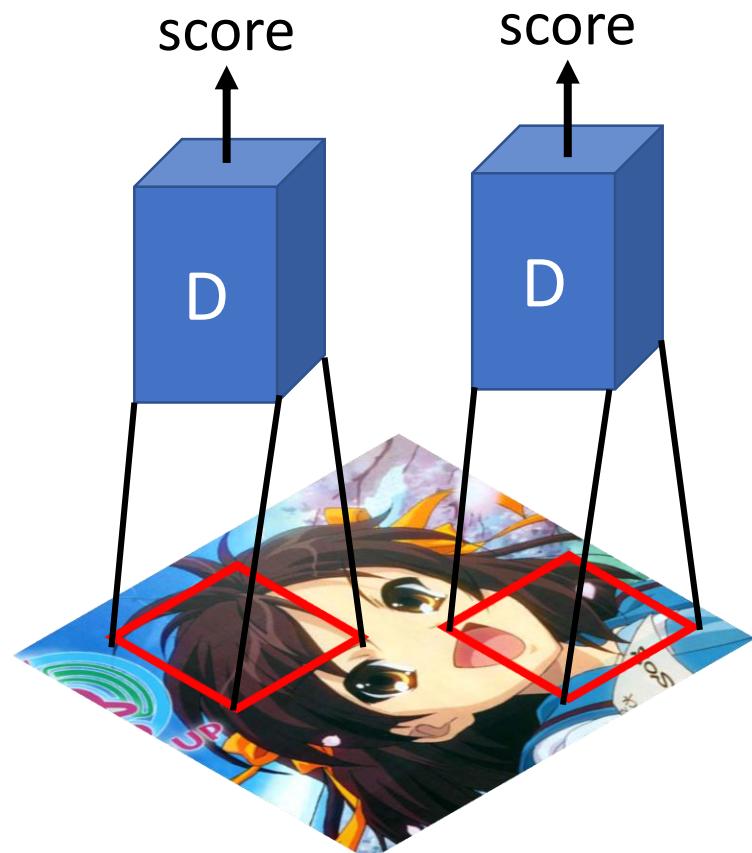
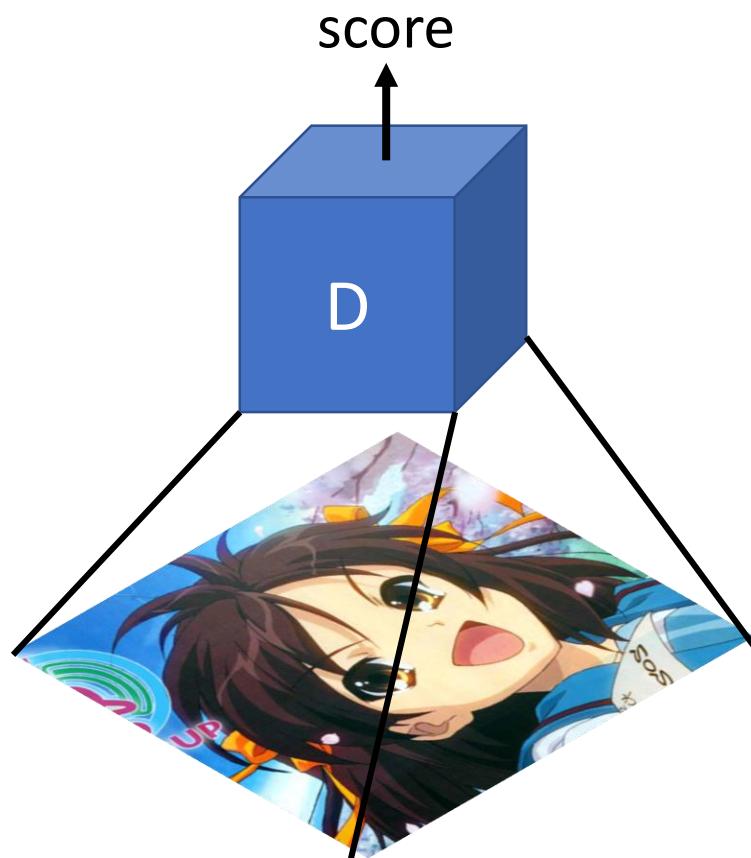
GAN



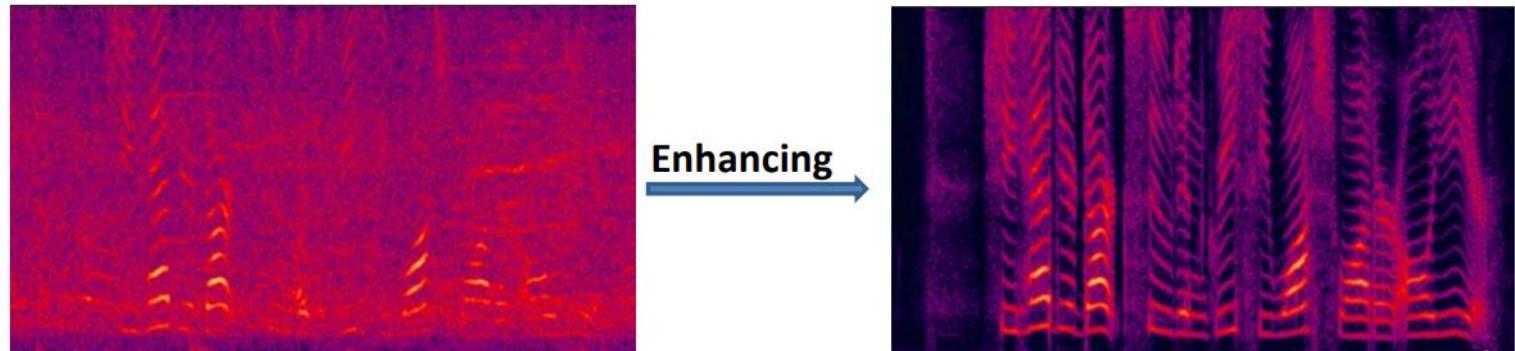
GAN + close

# Patch GAN

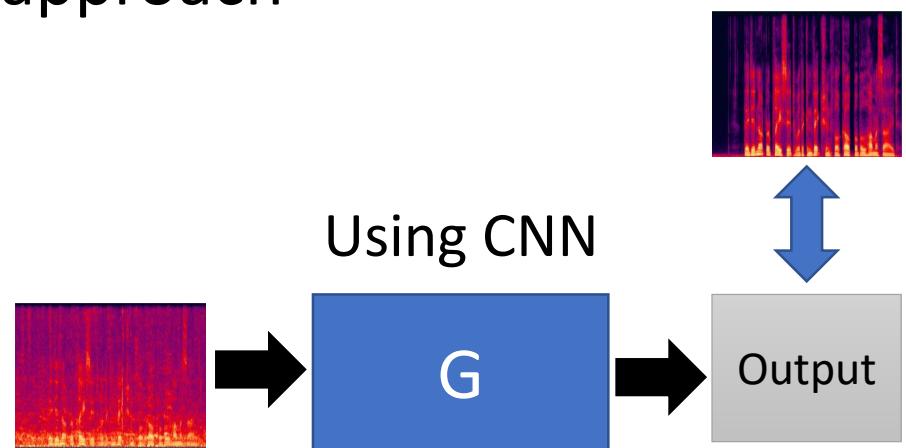
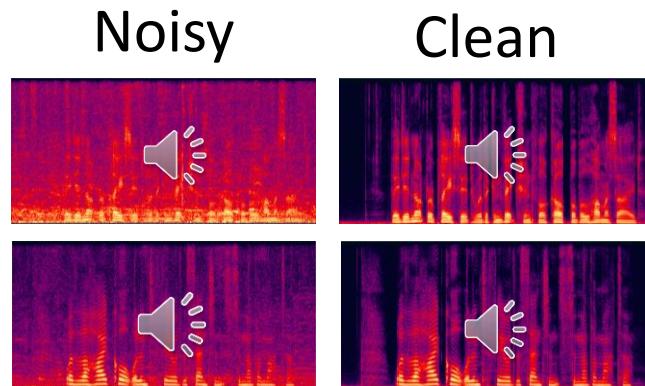
<https://arxiv.org/pdf/1611.07004.pdf>



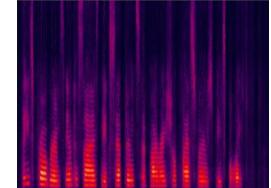
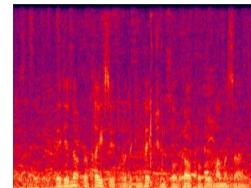
# Speech Enhancement



- Typical deep learning approach

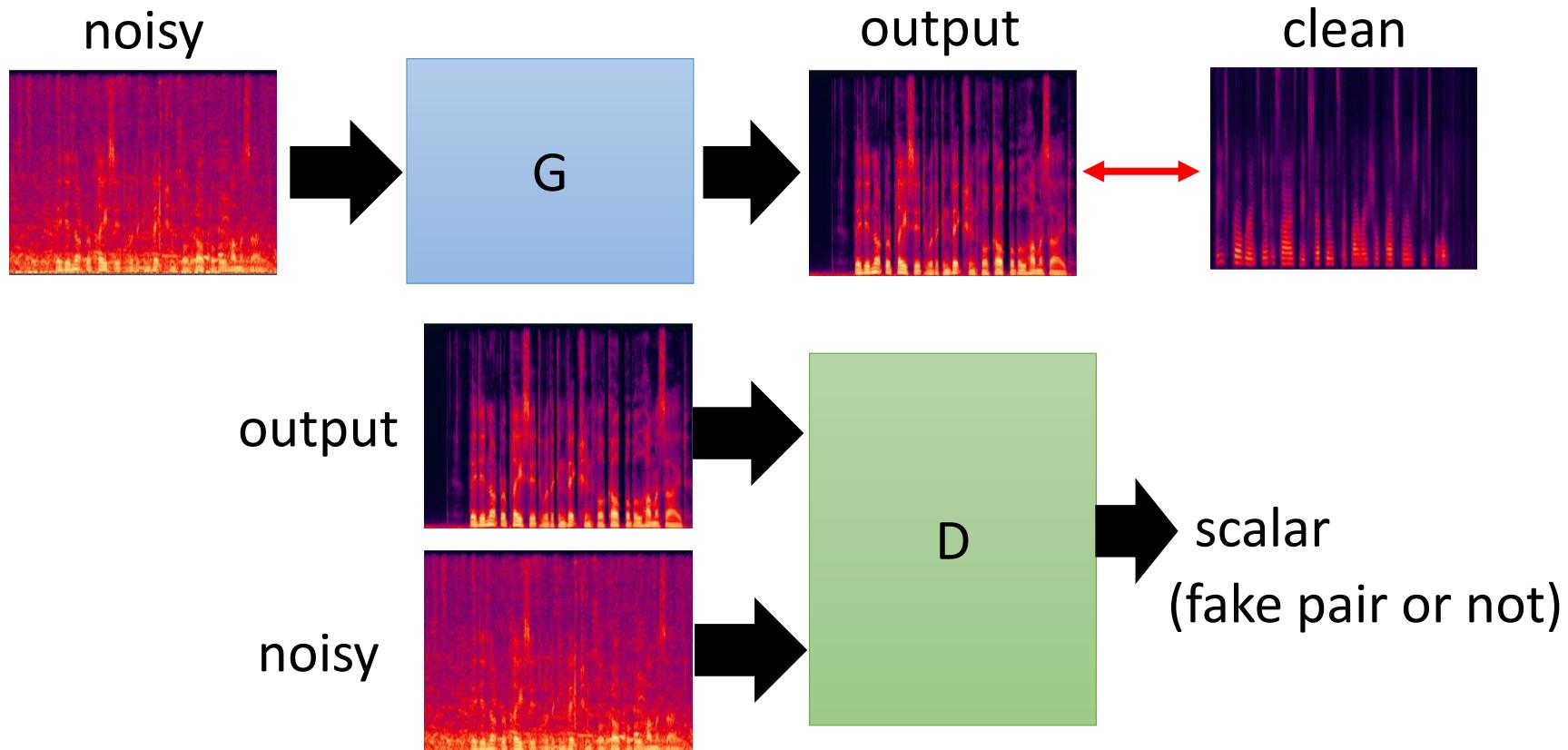


training data

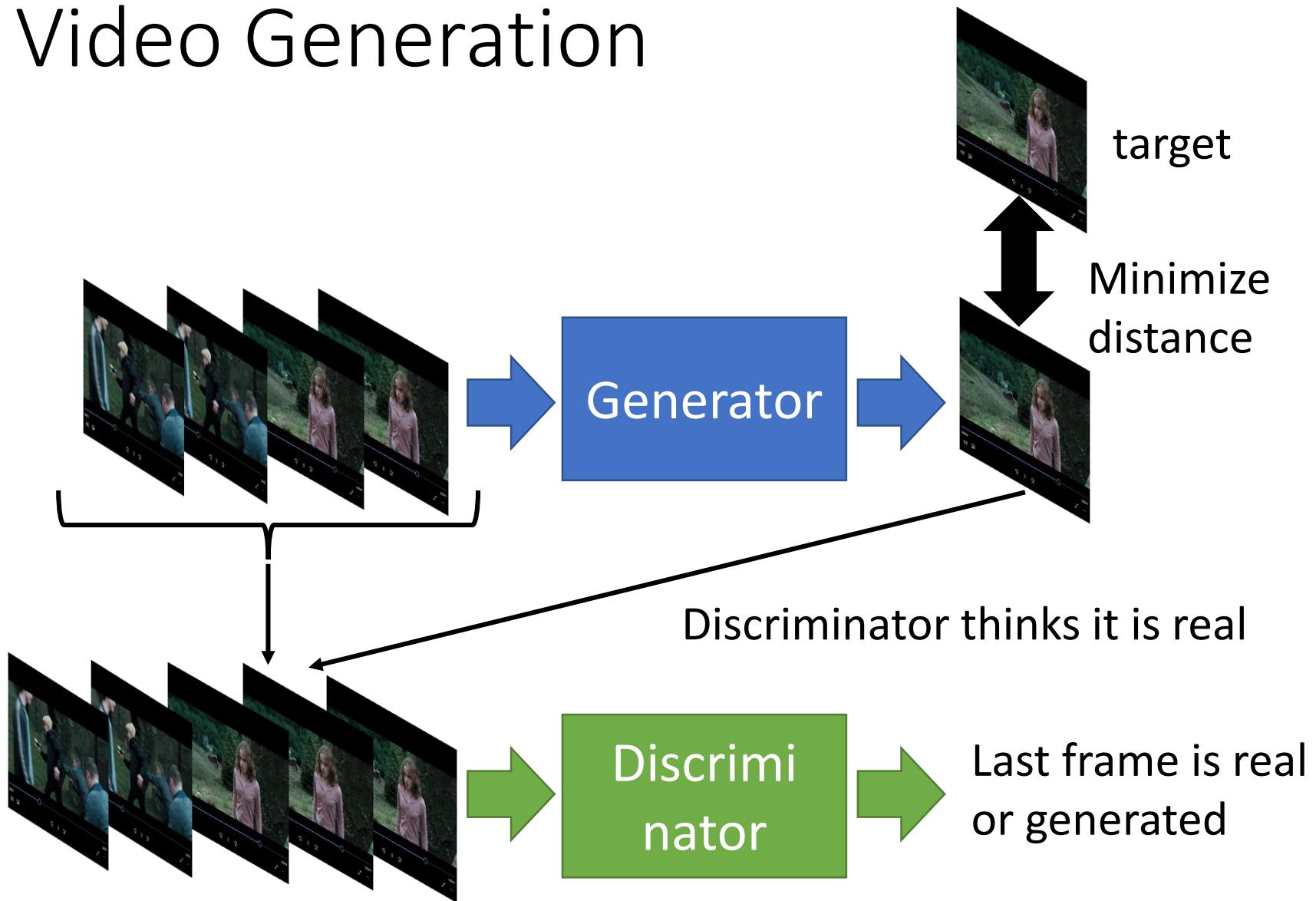


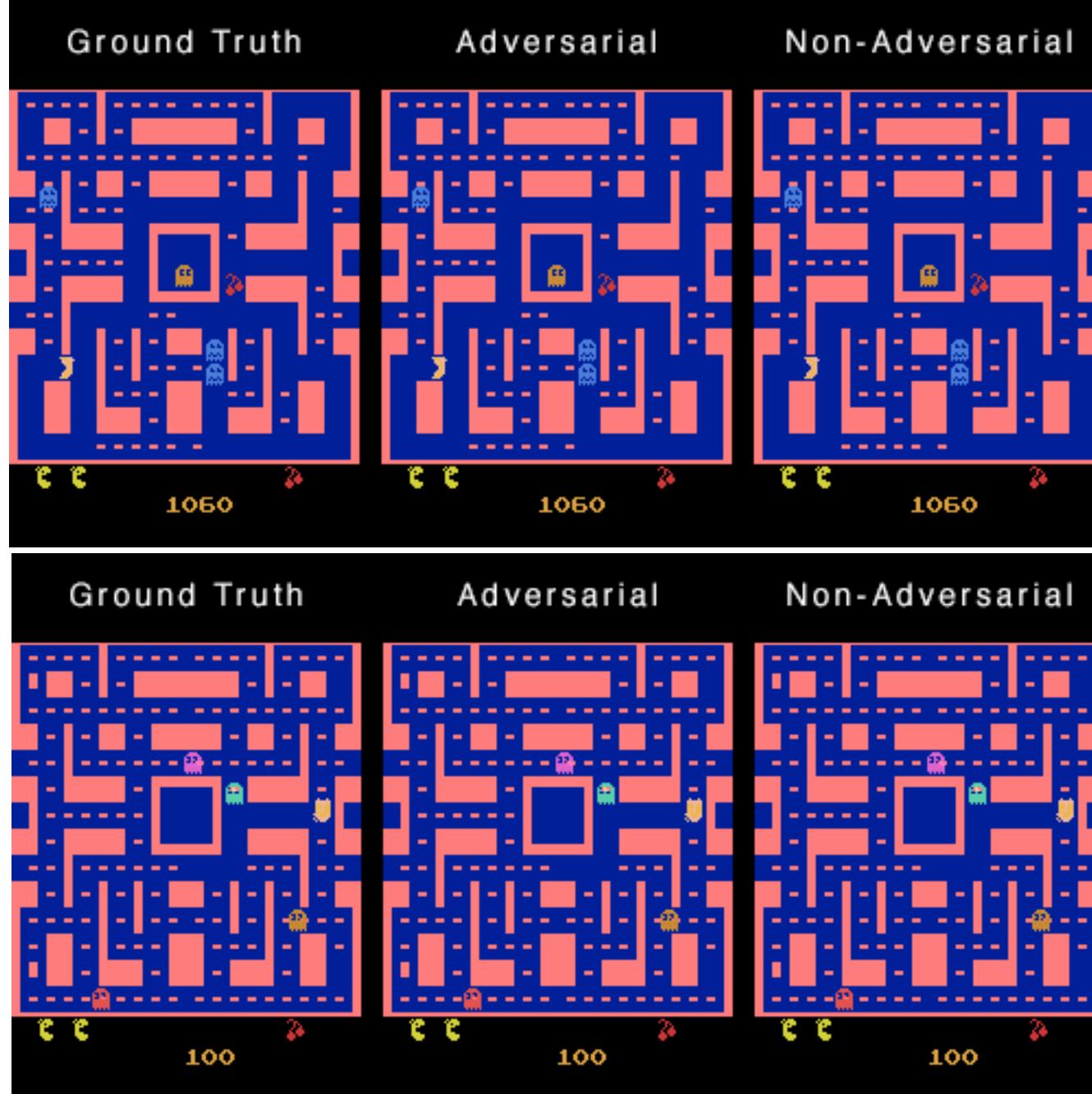
# Speech Enhancement

- Conditional GAN



# Video Generation



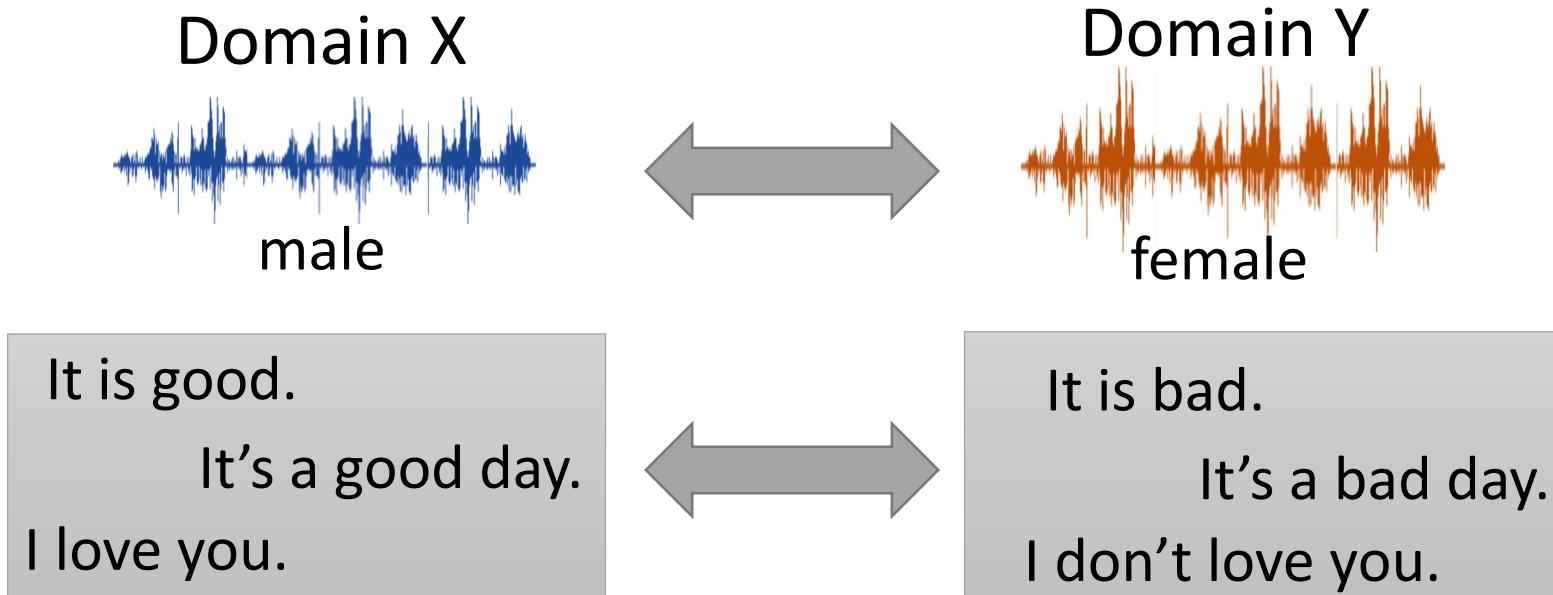


# Unsupervised Conditional Generation

# ***Unsupervised Conditional Generation***

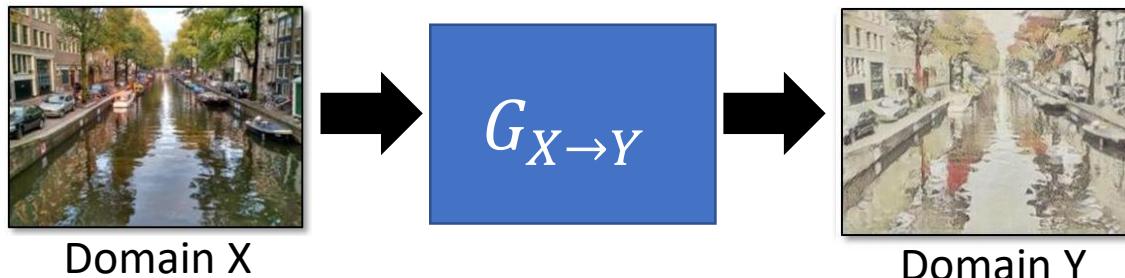


Transform an object from one domain to another  
***without paired data*** (e.g. style transfer)

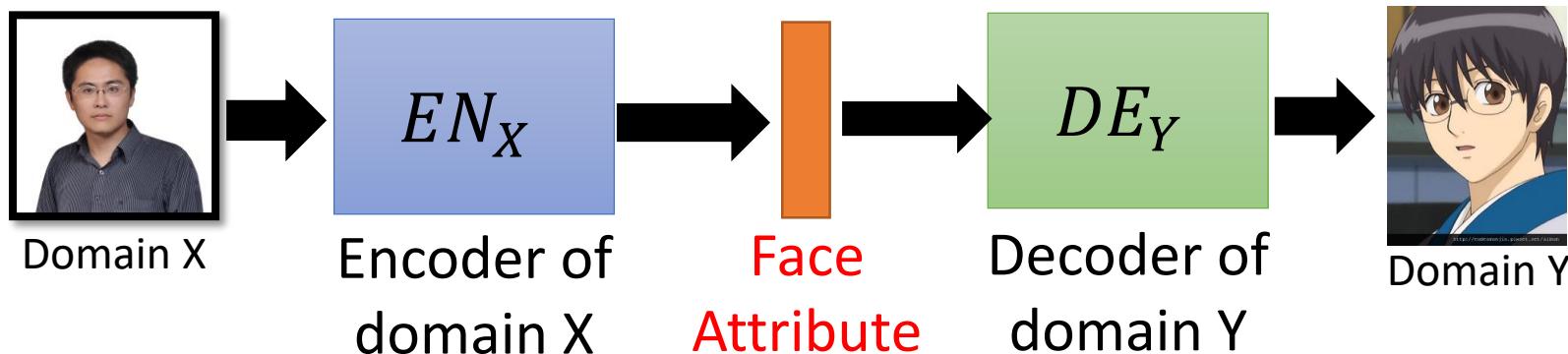


# Unsupervised Conditional Generation

- Approach 1: Direct Transformation

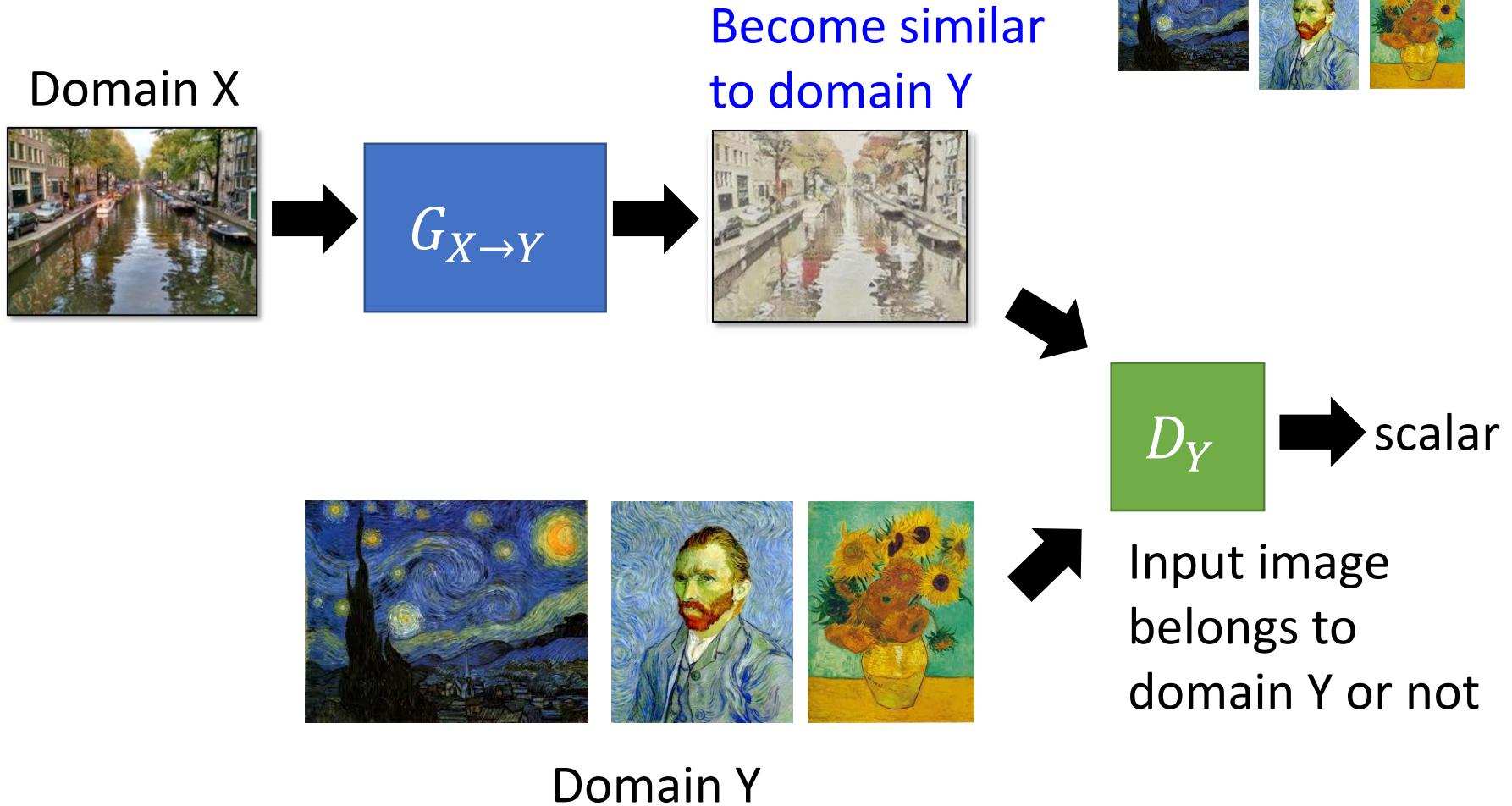


- Approach 2: Projection to Common Space

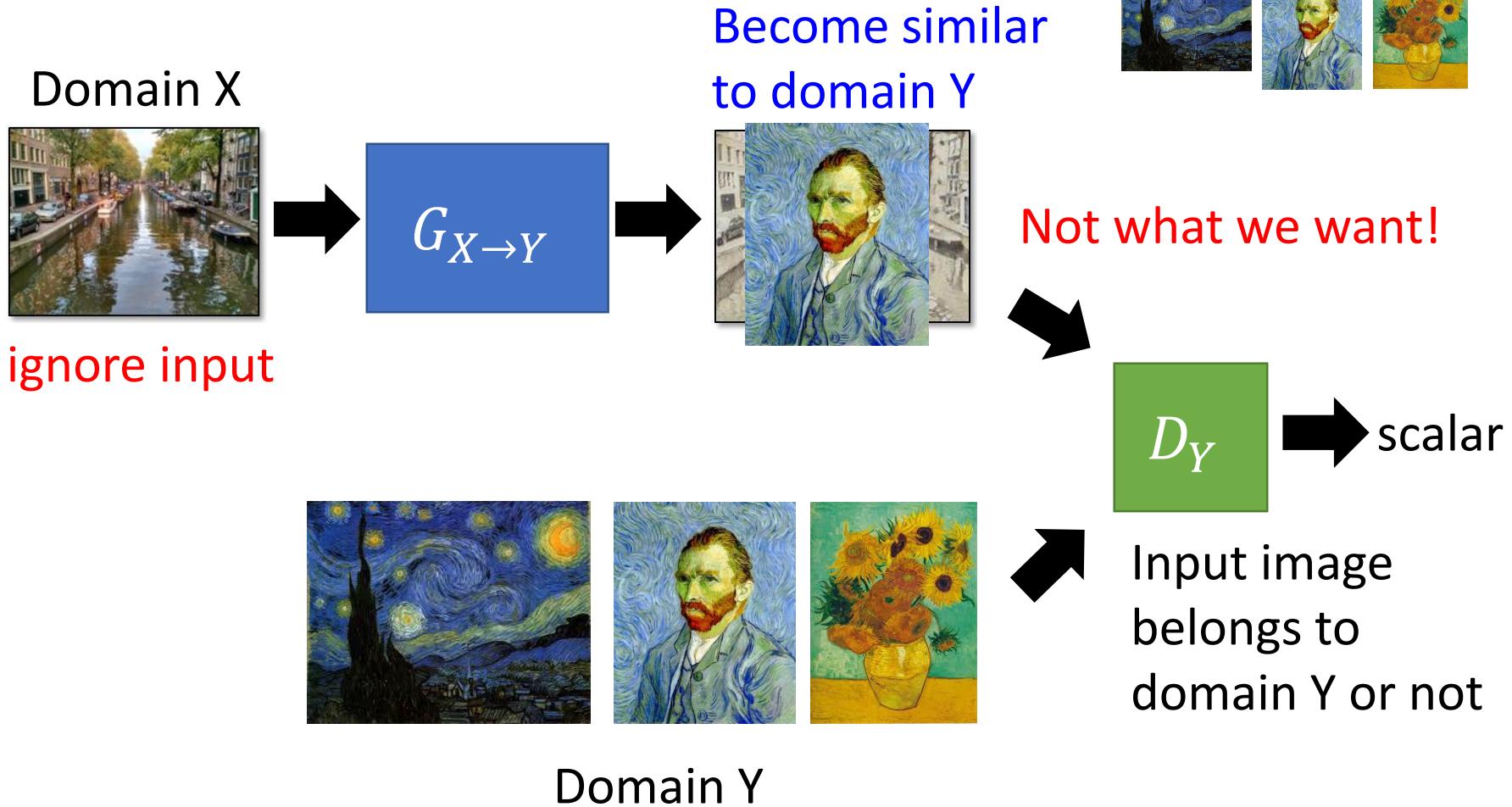


Larger change, only keep the semantics

# Direct Transformation



# Direct Transformation

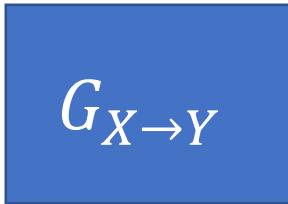


# Direct Transformation

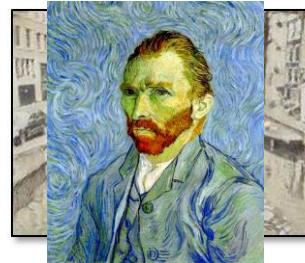
Domain X



ignore input



Become similar  
to domain Y



Not what we want!



scalar

The issue can be avoided by network design.

Simpler generator makes the input and output more closely related.



Domain Y



Input image  
belongs to  
domain Y or not

[Tomer Galanti, et al. ICLR, 2018]

# Direct Transformation

Domain X



Domain Y



Domain X



$$G_{X \rightarrow Y}$$

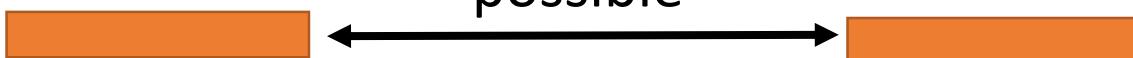
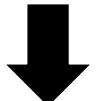
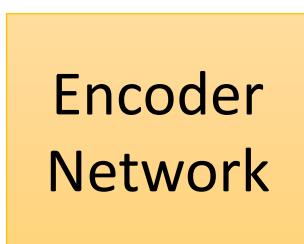


Become similar  
to domain Y



← pre-trained →

as close as  
possible

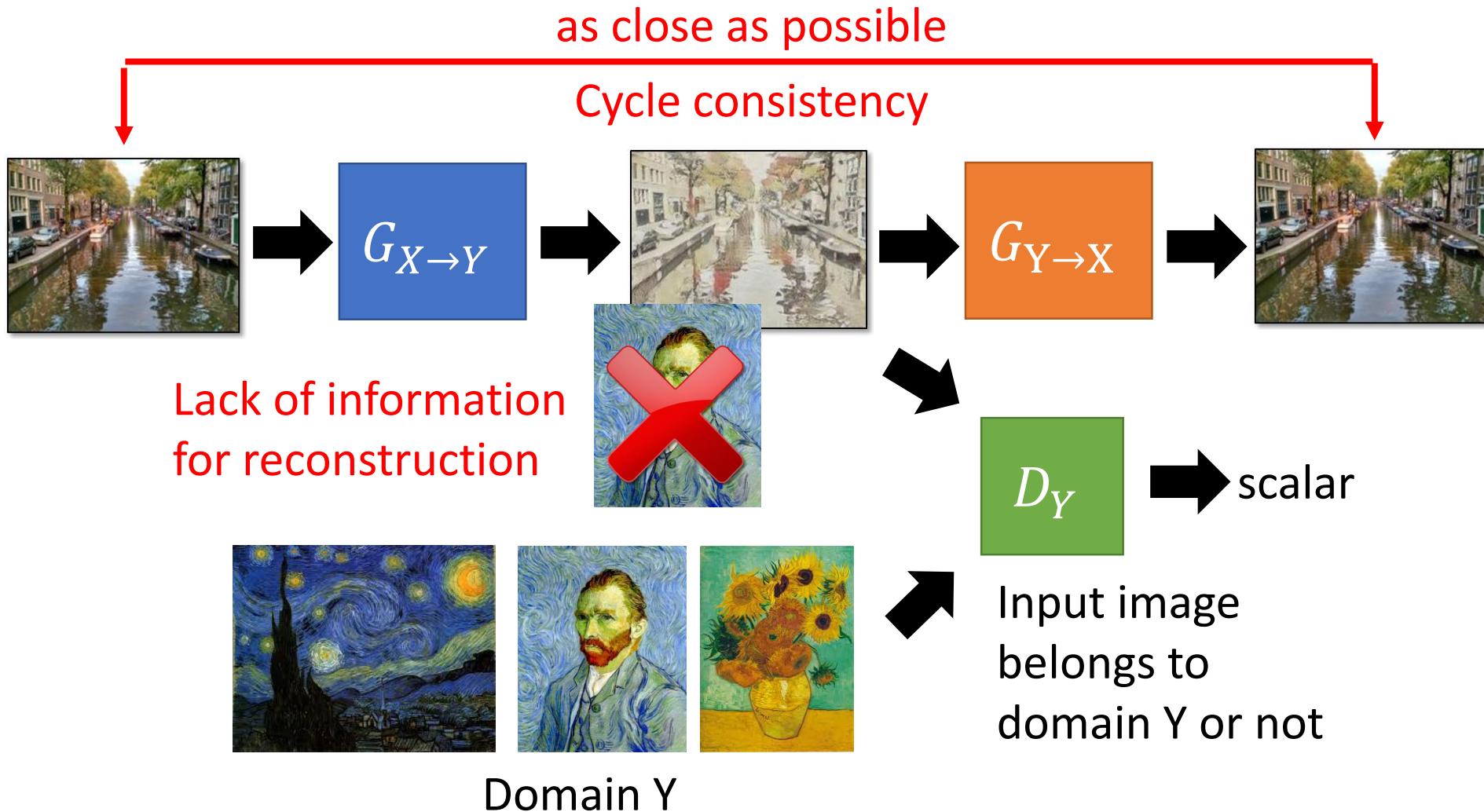


$$D_Y \rightarrow \text{scalar}$$

Input image  
belongs to  
domain Y or not

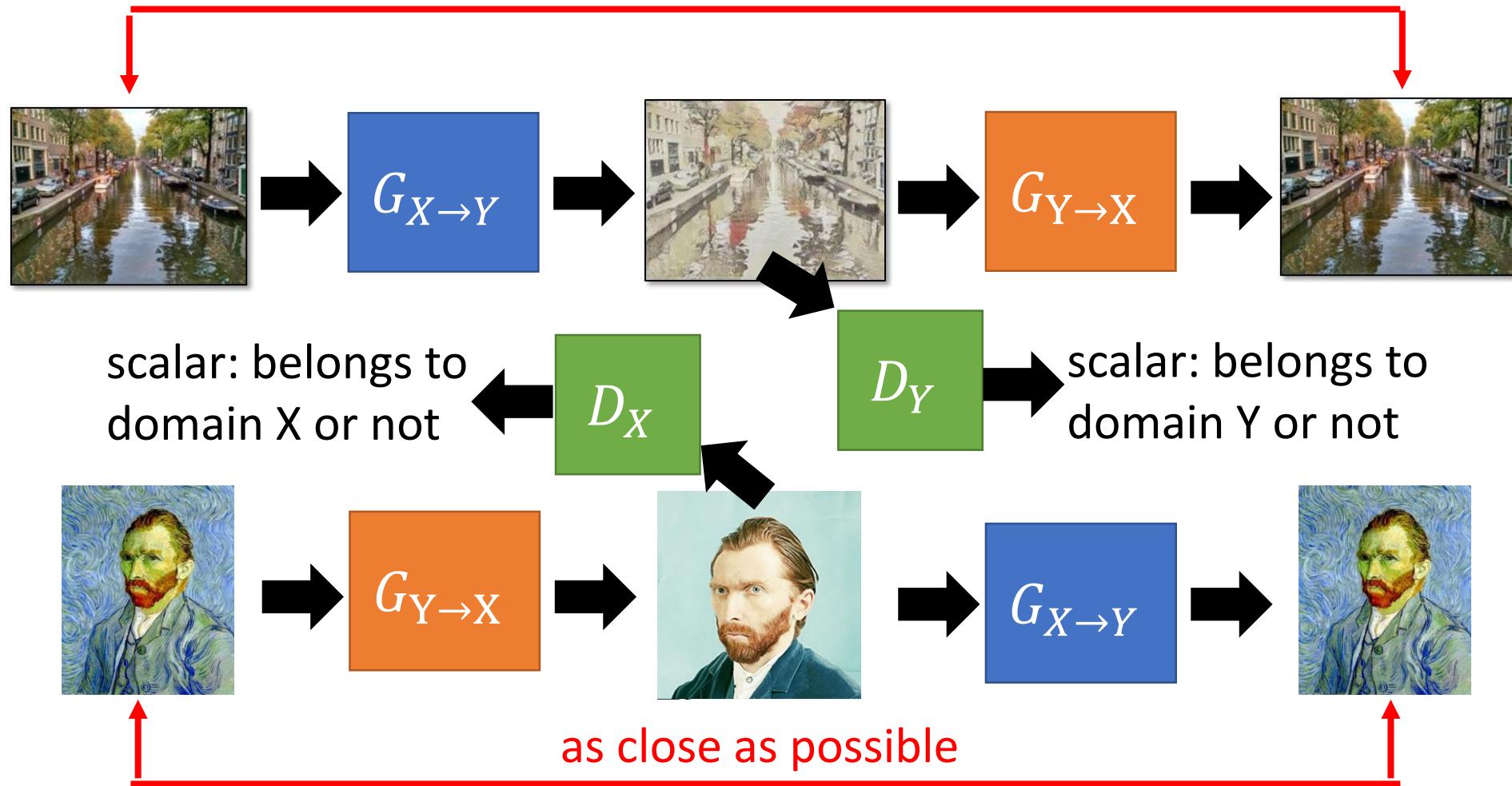
Baseline of DTN [Yaniv Taigman, et al., ICLR, 2017]

# Direct Transformation



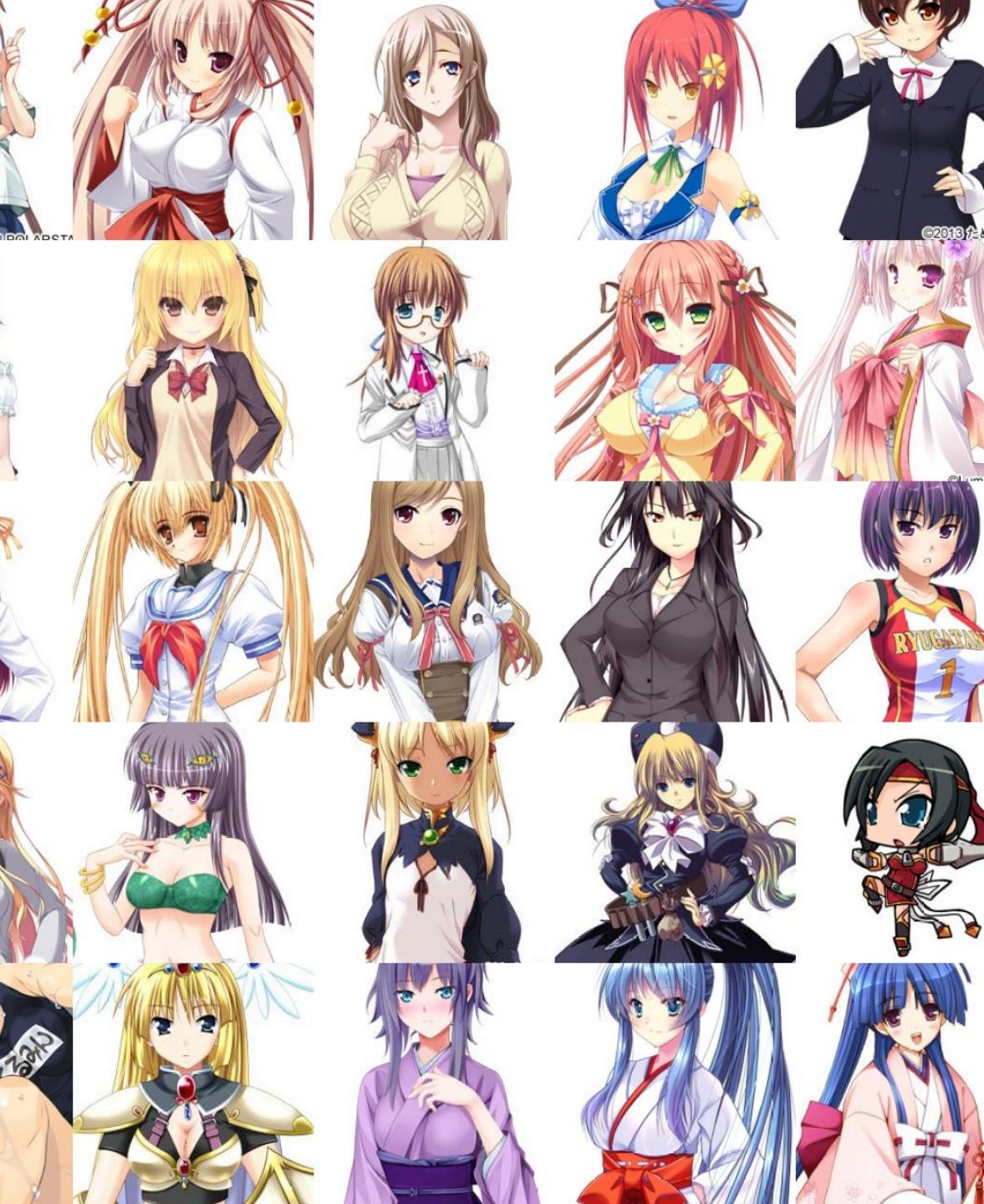
# Direct Transformation

as close as possible



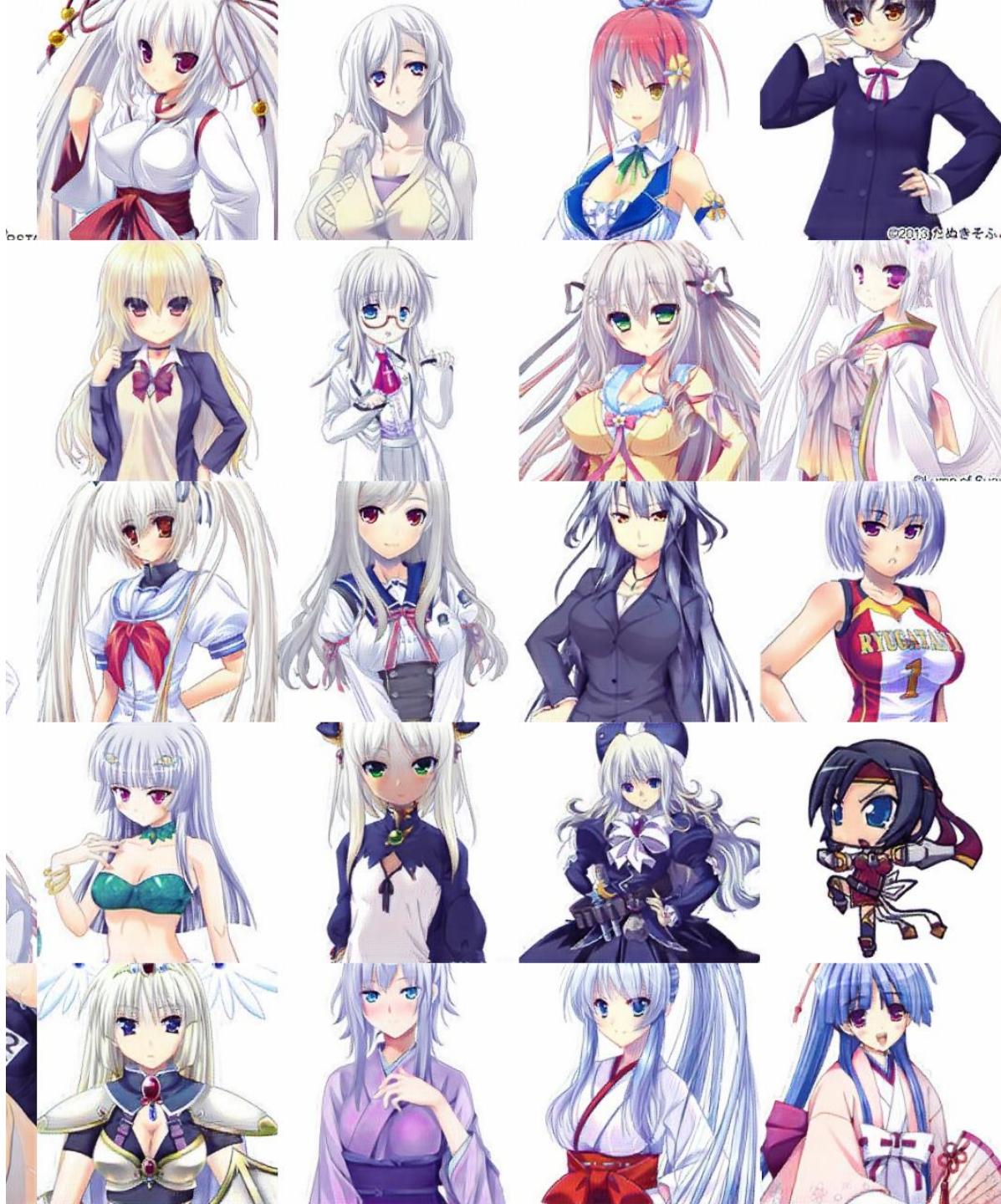
# Cycle GAN – Silver Hair

- <https://github.com/Aixile/chainer-cyclegan>



# Cycle GAN – Silver Hair

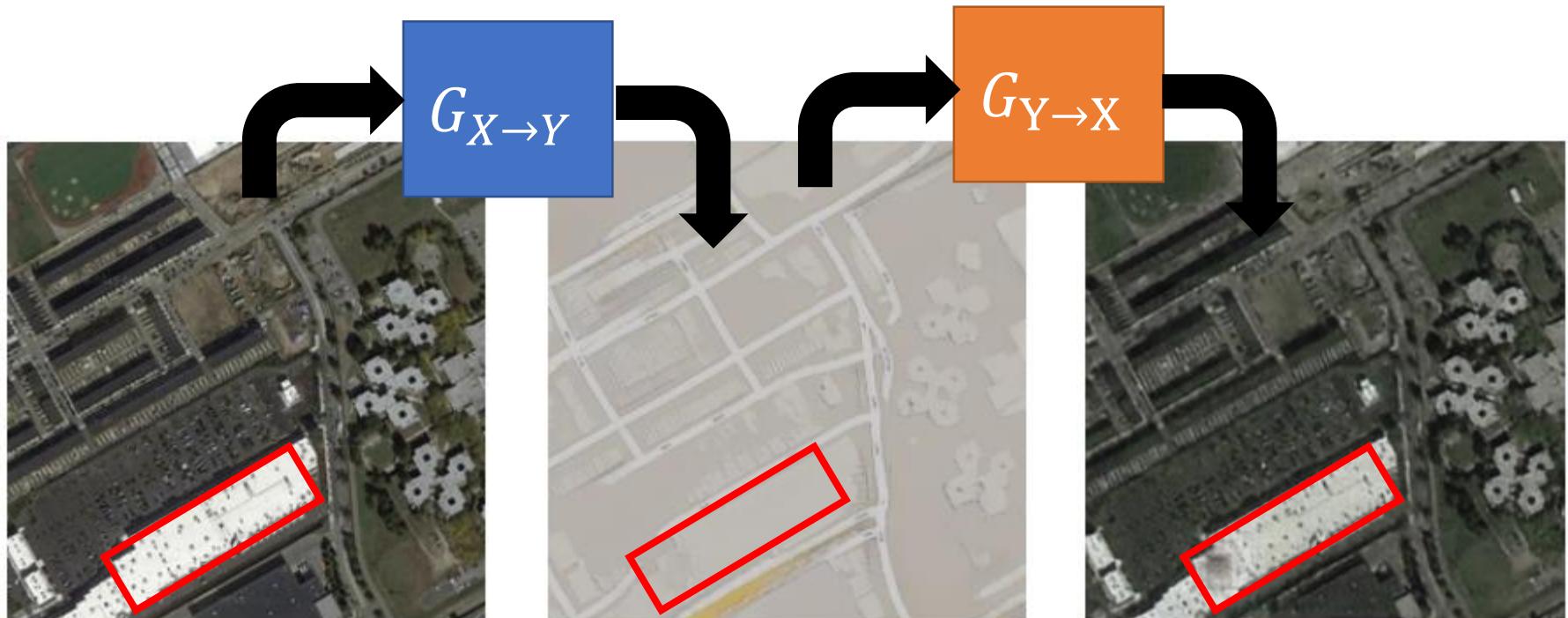
- <https://github.com/Aixile/chainer-cyclegan>



# Issue of Cycle Consistency

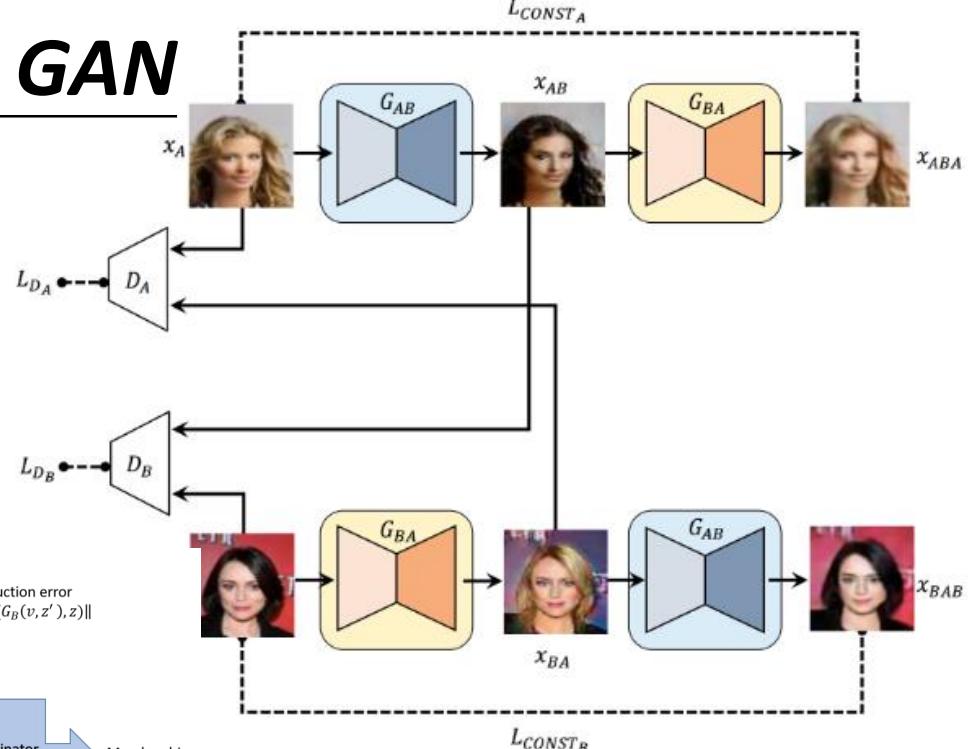
- **CycleGAN: a Master of Steganography (隱寫術)**

[Casey Chu, et al., NIPS workshop, 2017]



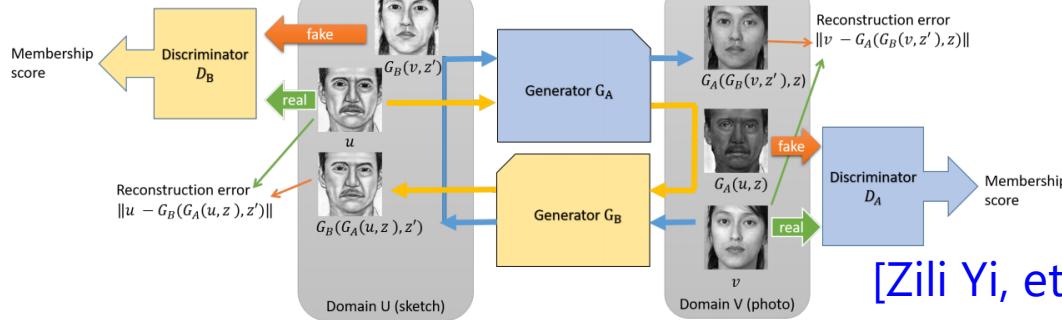
The information is hidden.

# Disco GAN



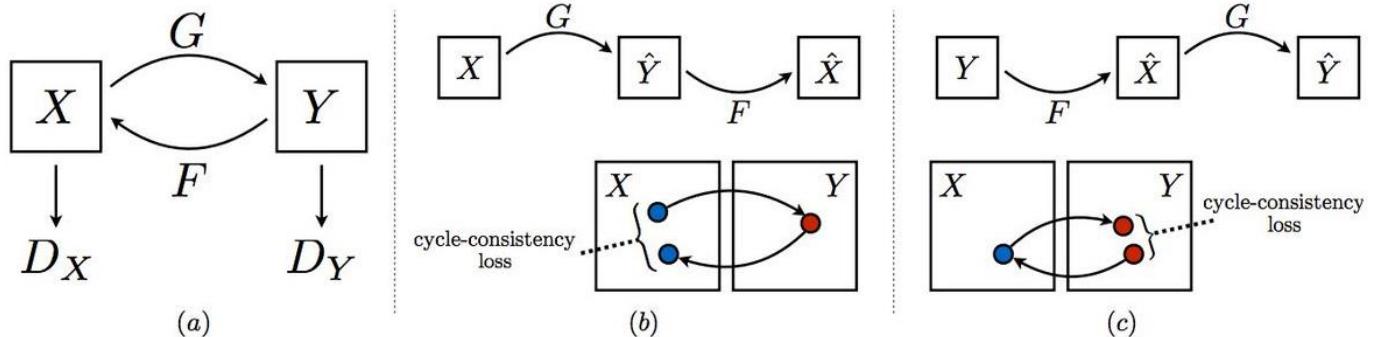
[Taeksoo Kim, et al., ICML, 2017]

# Dual GAN



[Zili Yi, et al., ICCV, 2017]

# Cycle GAN



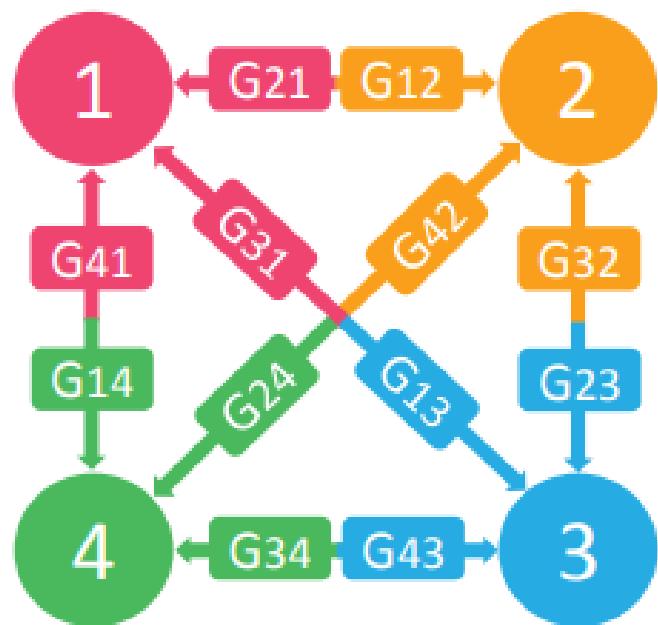
[Jun-Yan Zhu, et al., ICCV, 2017]

# StarGAN

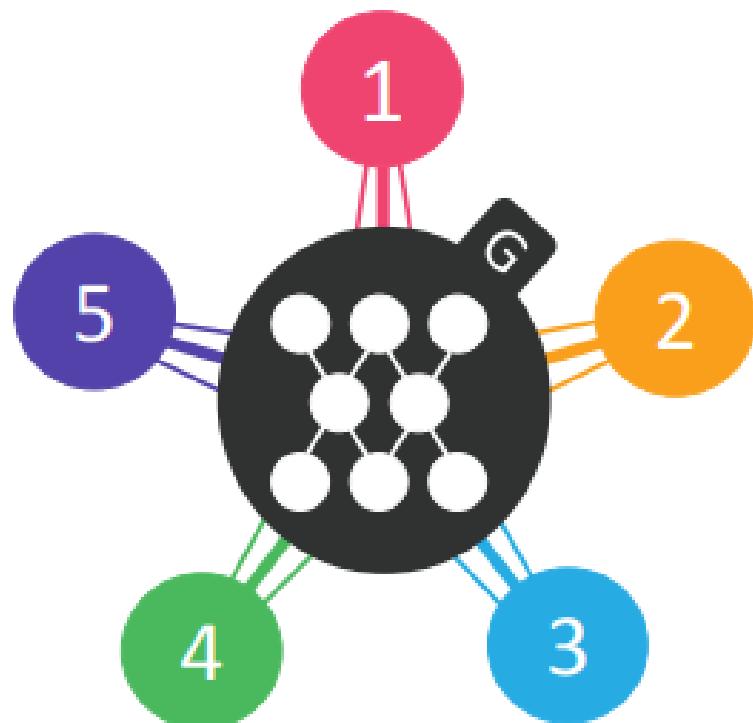
For multiple domains,  
considering starGAN

[Yunjey Choi, arXiv, 2017]

(a) Cross-domain models

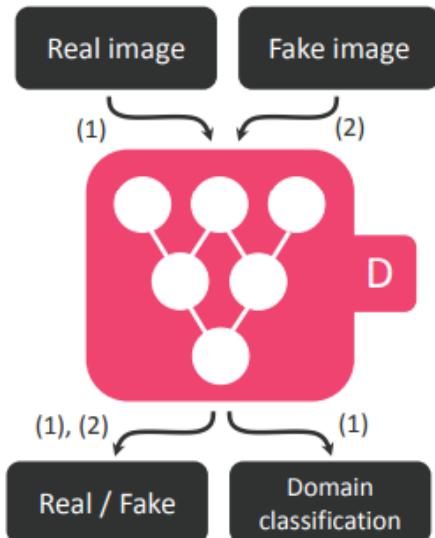


(b) StarGAN

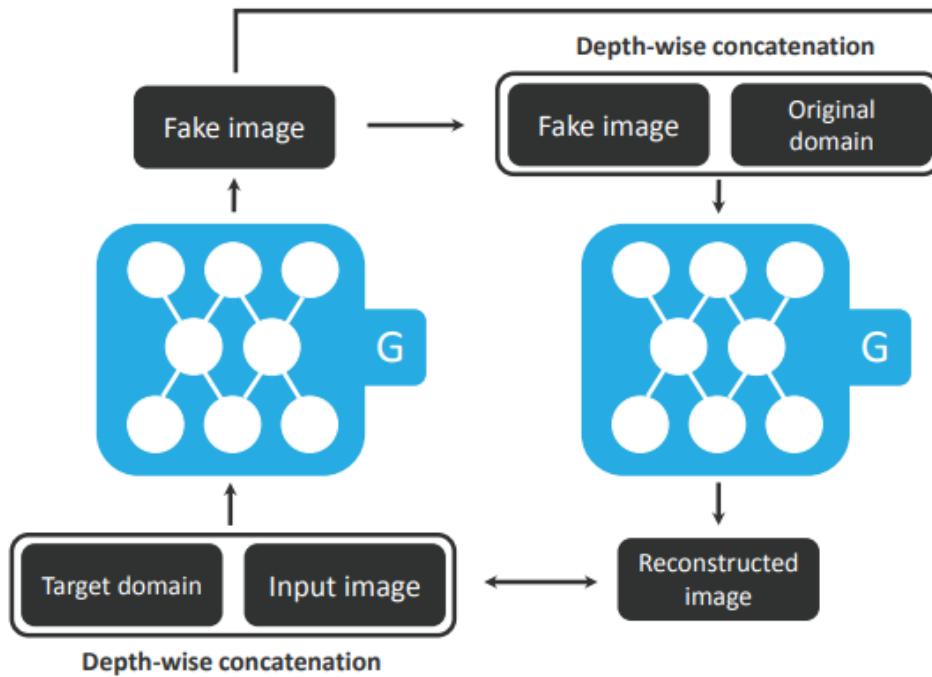


# StarGAN

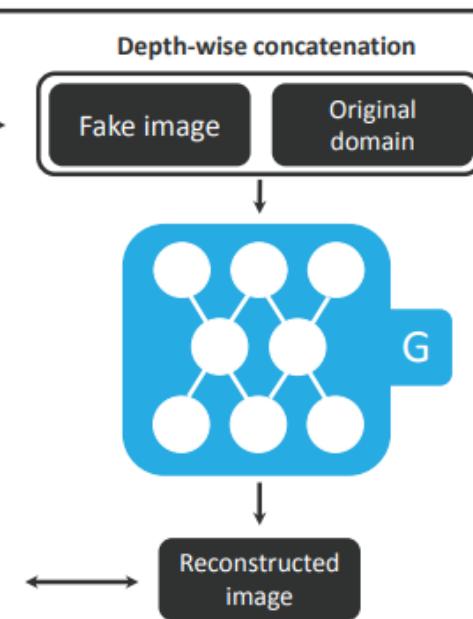
(a) Training the discriminator



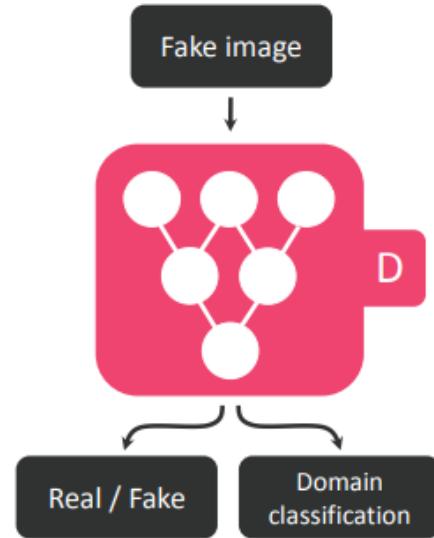
(b) Original-to-target domain



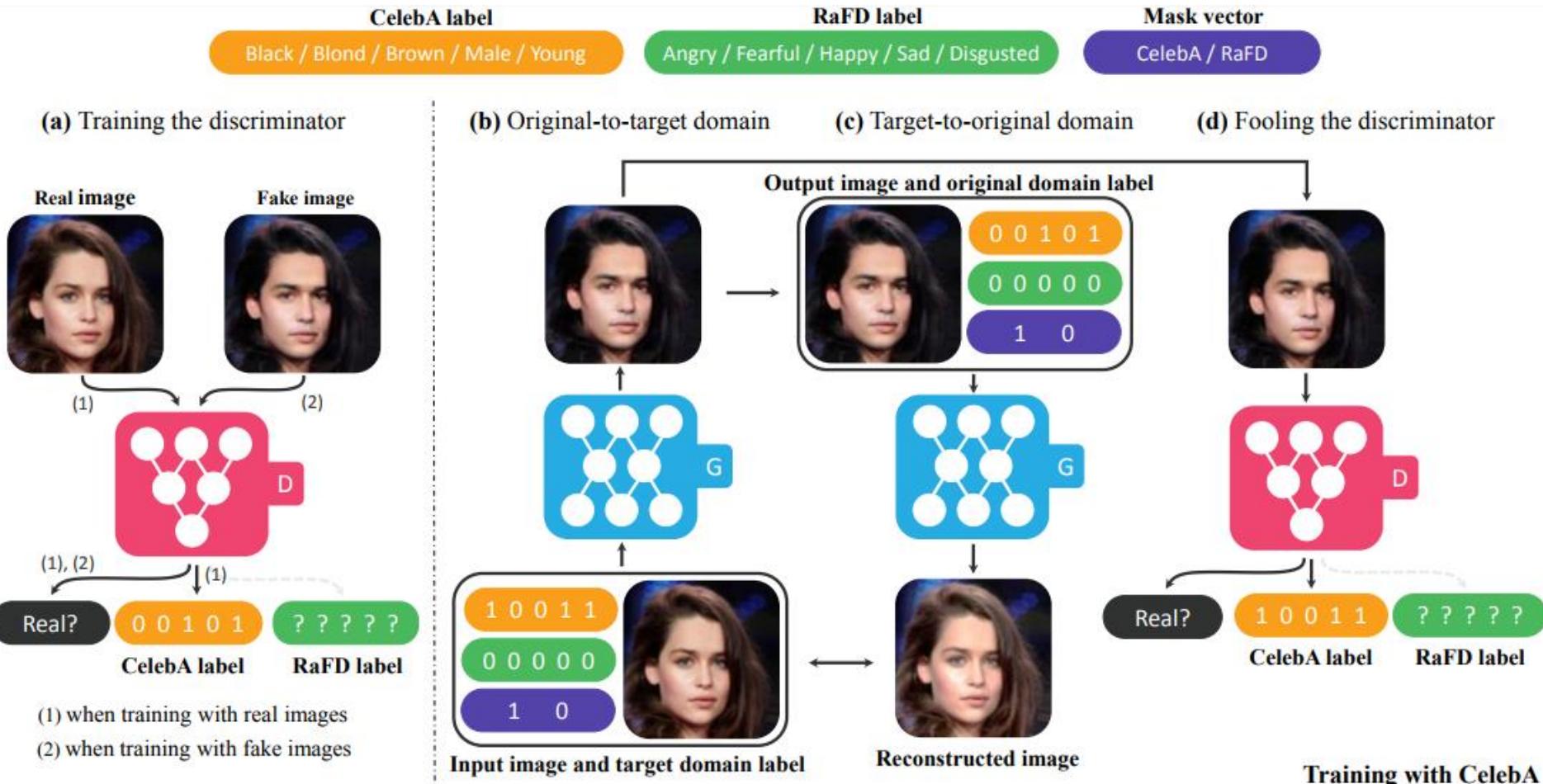
(c) Target-to-original domain



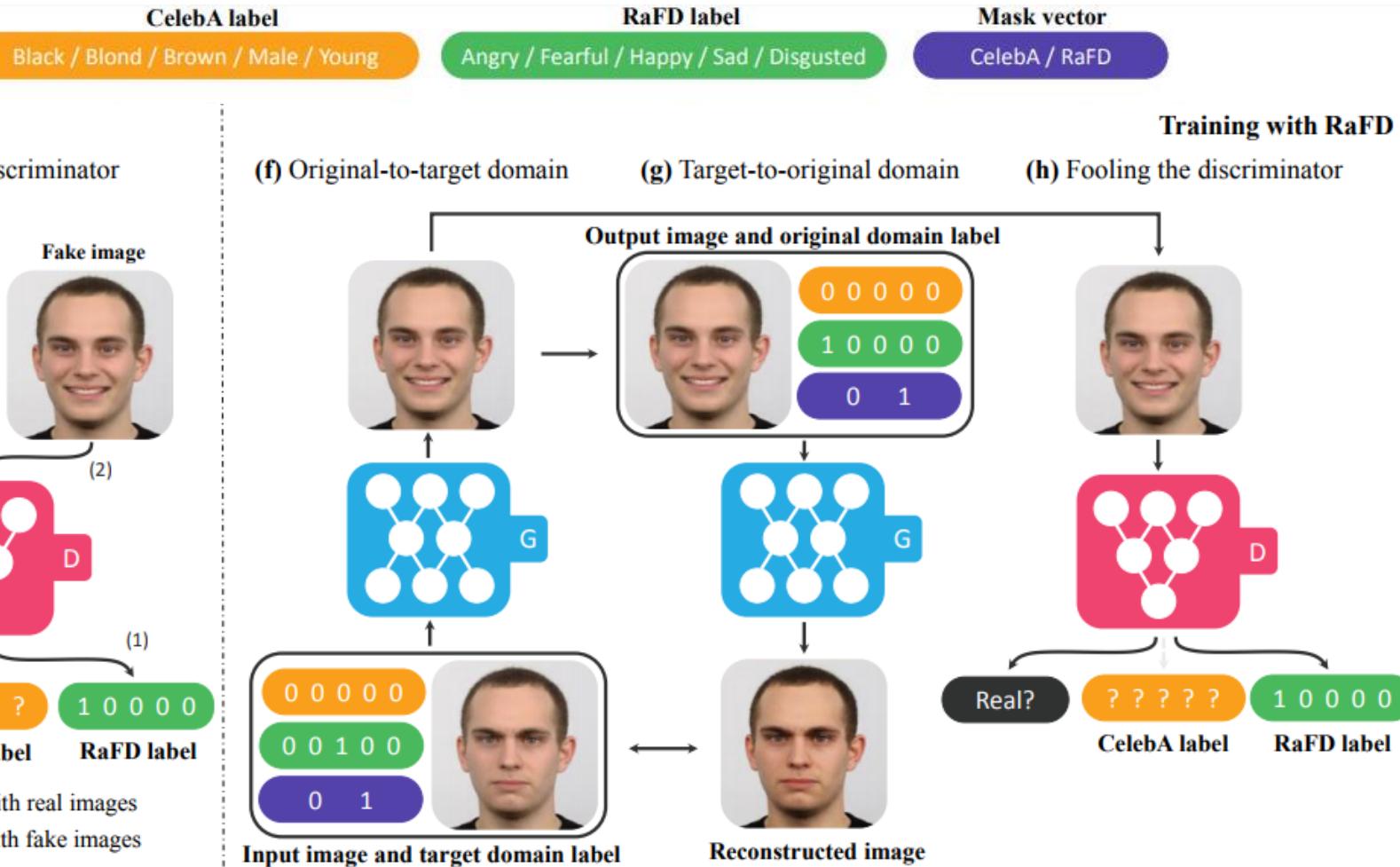
(d) Fooling the discriminator



# StarGAN

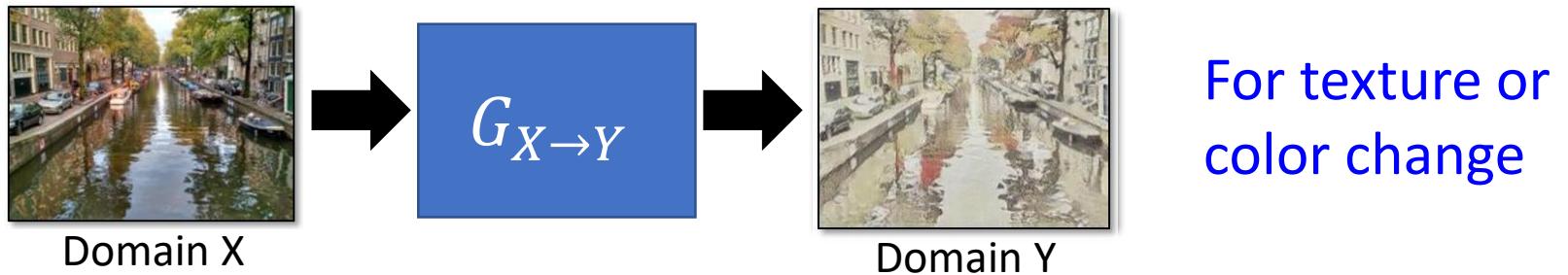


# StarGAN

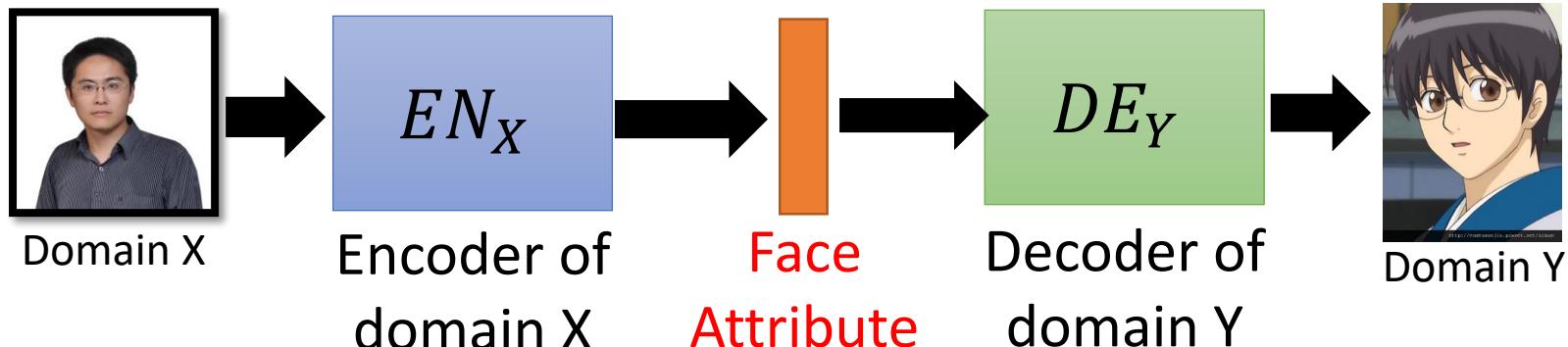


# Unsupervised Conditional Generation

- Approach 1: Direct Transformation



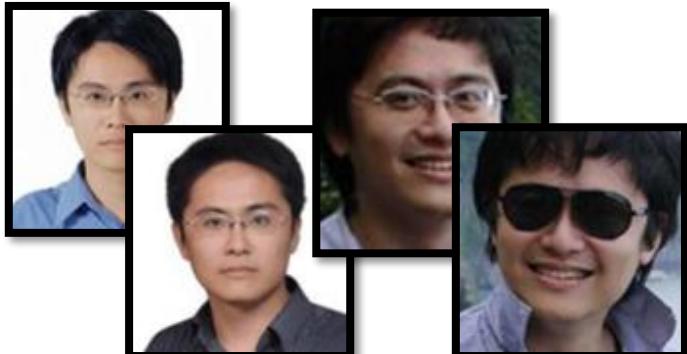
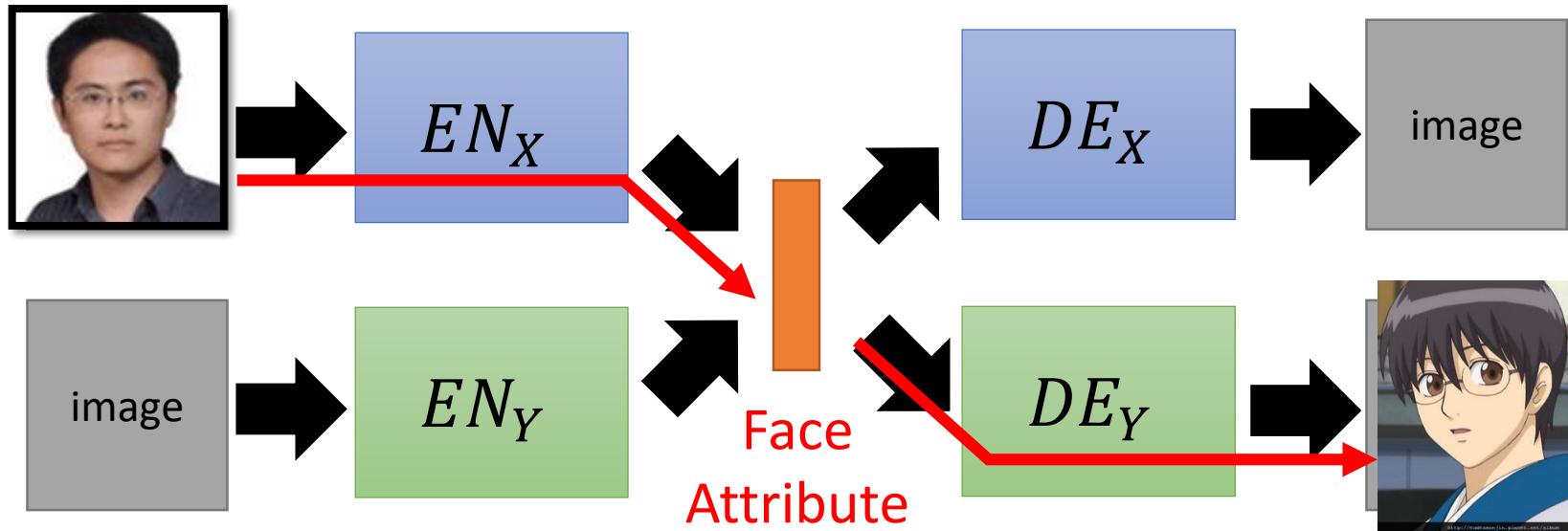
- Approach 2: Projection to Common Space



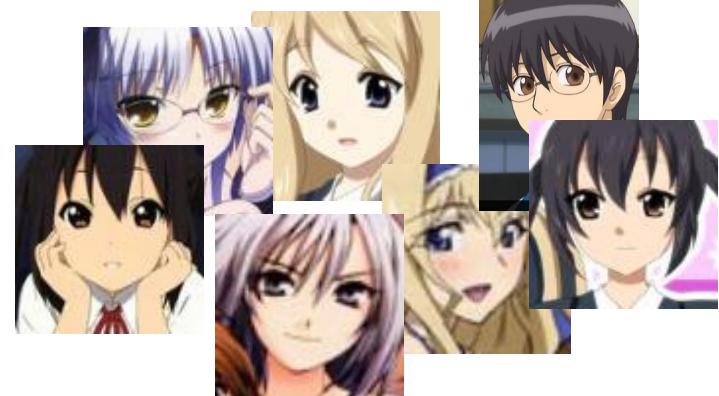
Larger change, only keep the semantics

# Projection to Common Space

Target



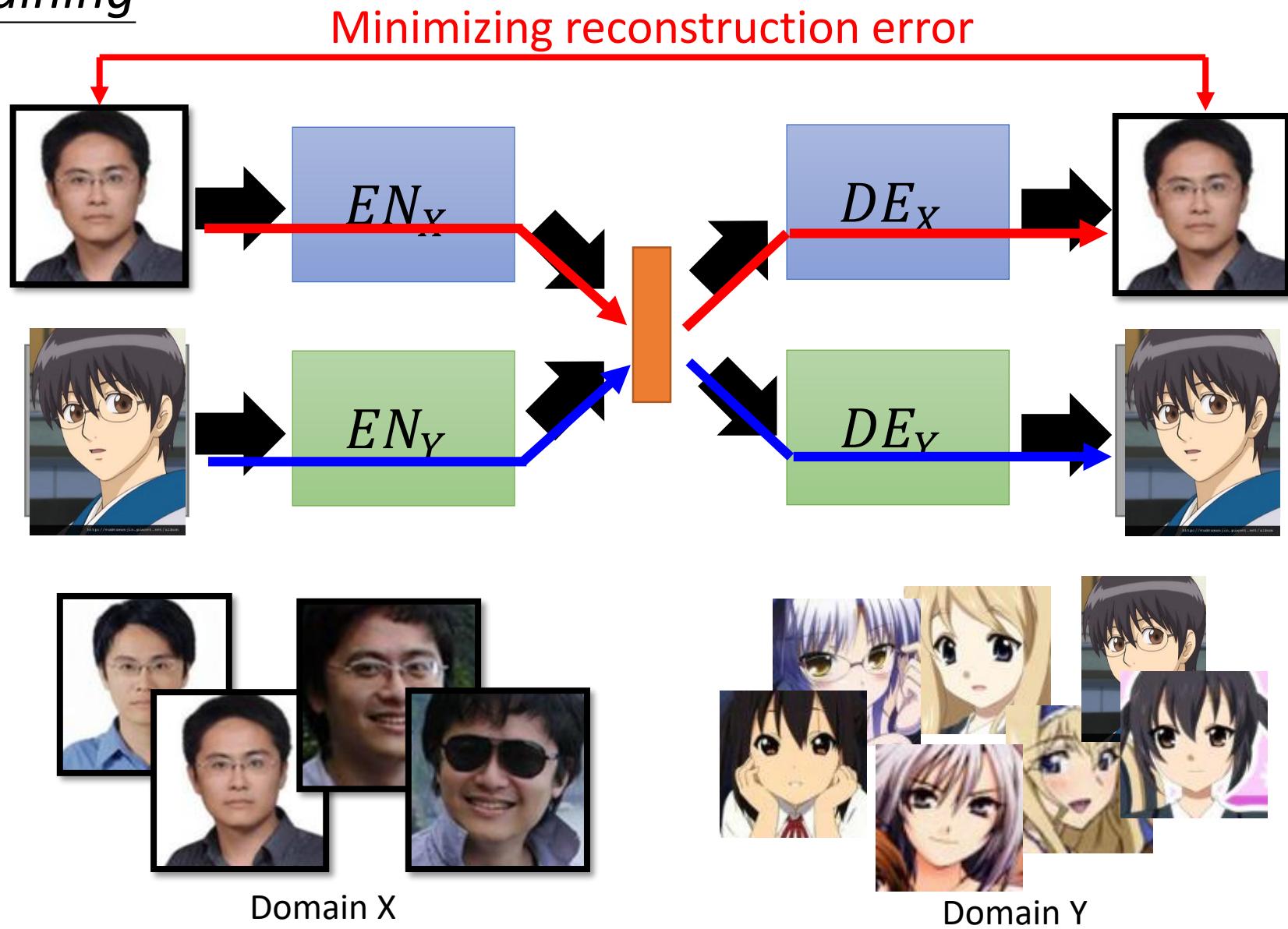
Domain X



Domain Y

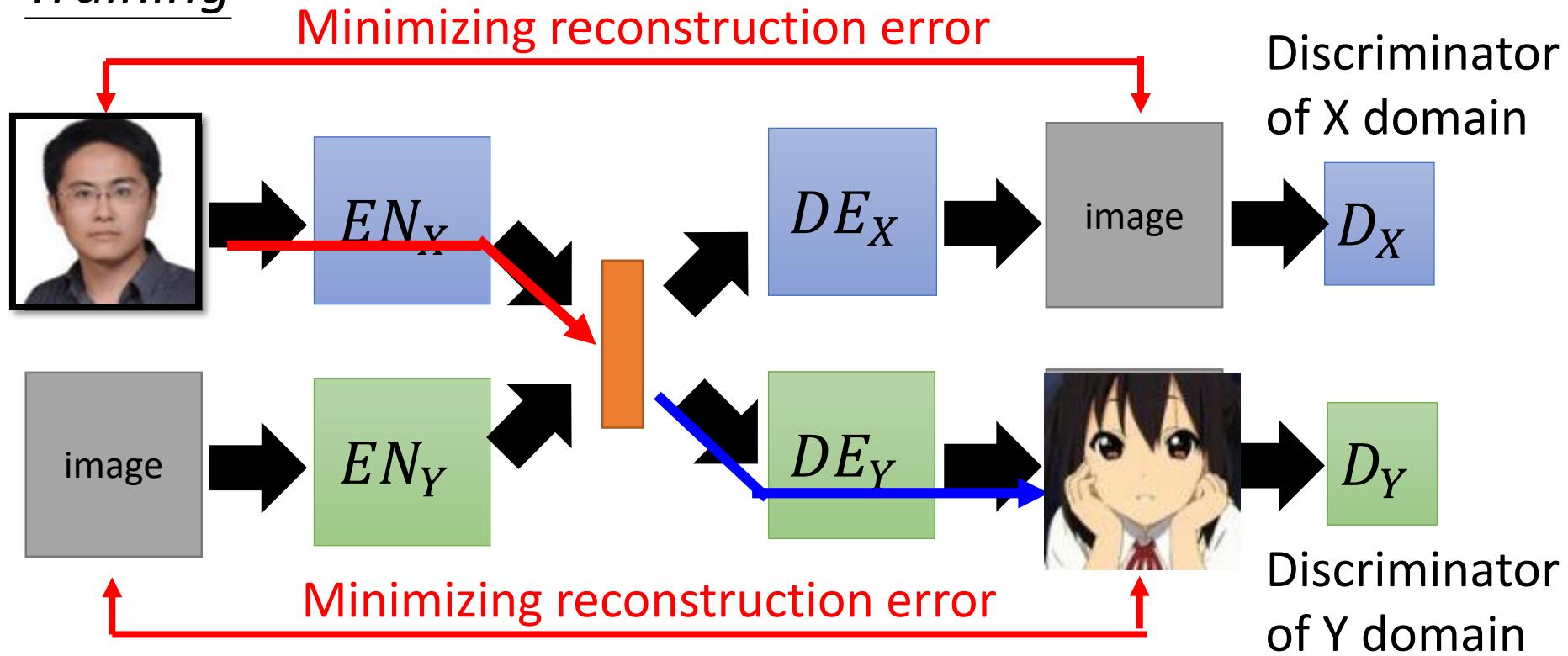
# *Projection to Common Space*

## Training



# Projection to Common Space

## Training

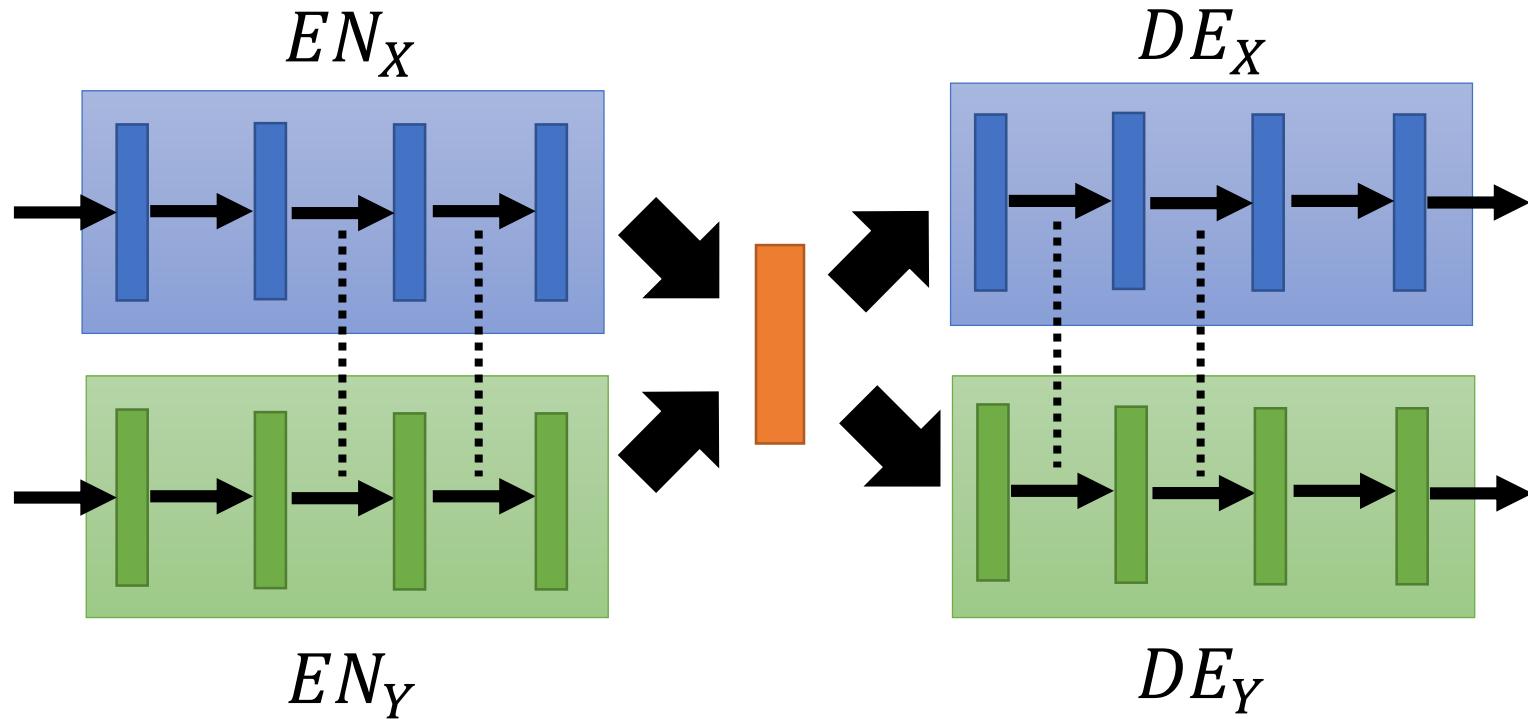


Because we train two auto-encoders separately ...

The images with the same attribute may not project to the same position in the latent space.

# *Projection to Common Space*

## Training

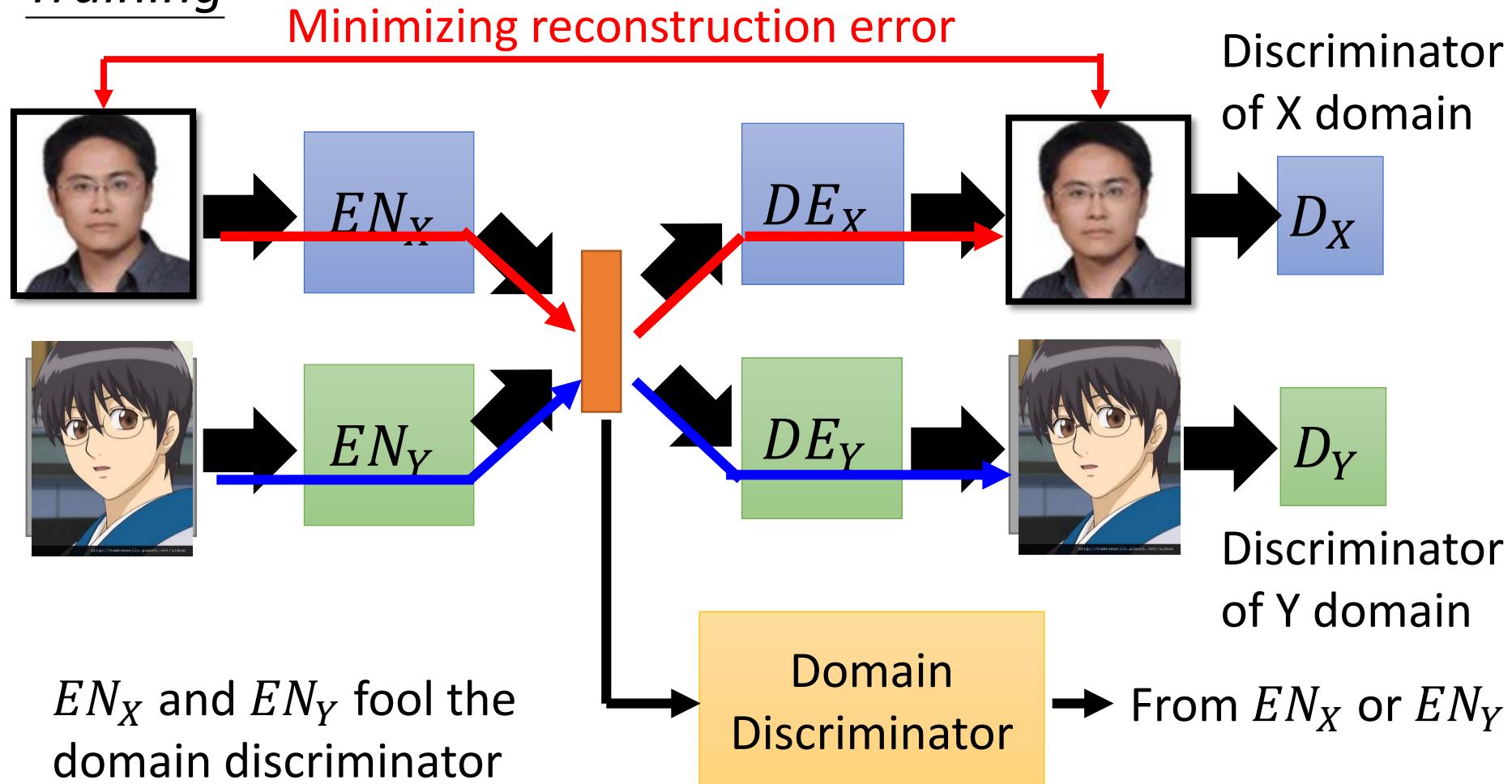


Sharing the parameters of encoders and decoders

Couple GAN [Ming-Yu Liu, et al., NIPS, 2016]  
UNIT [Ming-Yu Liu, et al., NIPS, 2017]

# Projection to Common Space

## Training

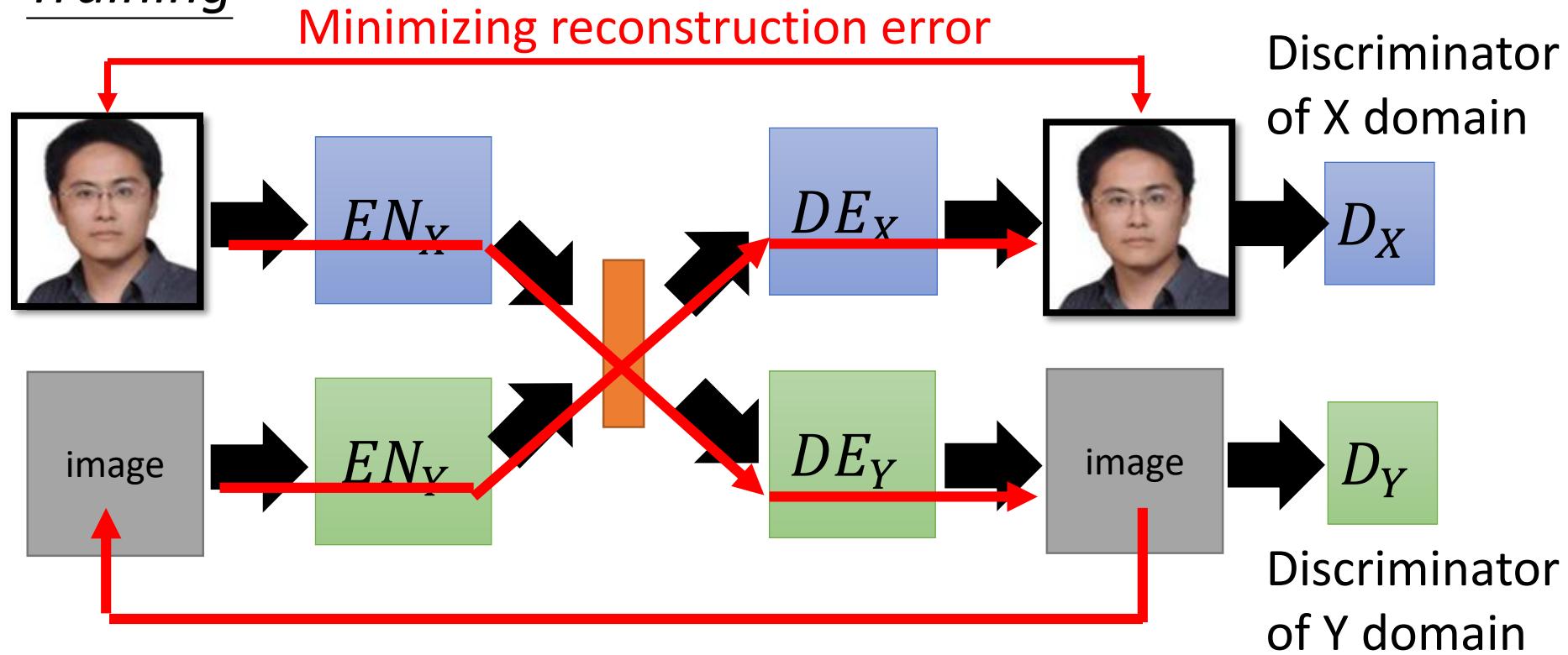


The domain discriminator forces the output of  $EN_X$  and  $EN_Y$  have the same distribution.

[Guillaume Lample, et al., NIPS, 2017]

# Projection to Common Space

## Training

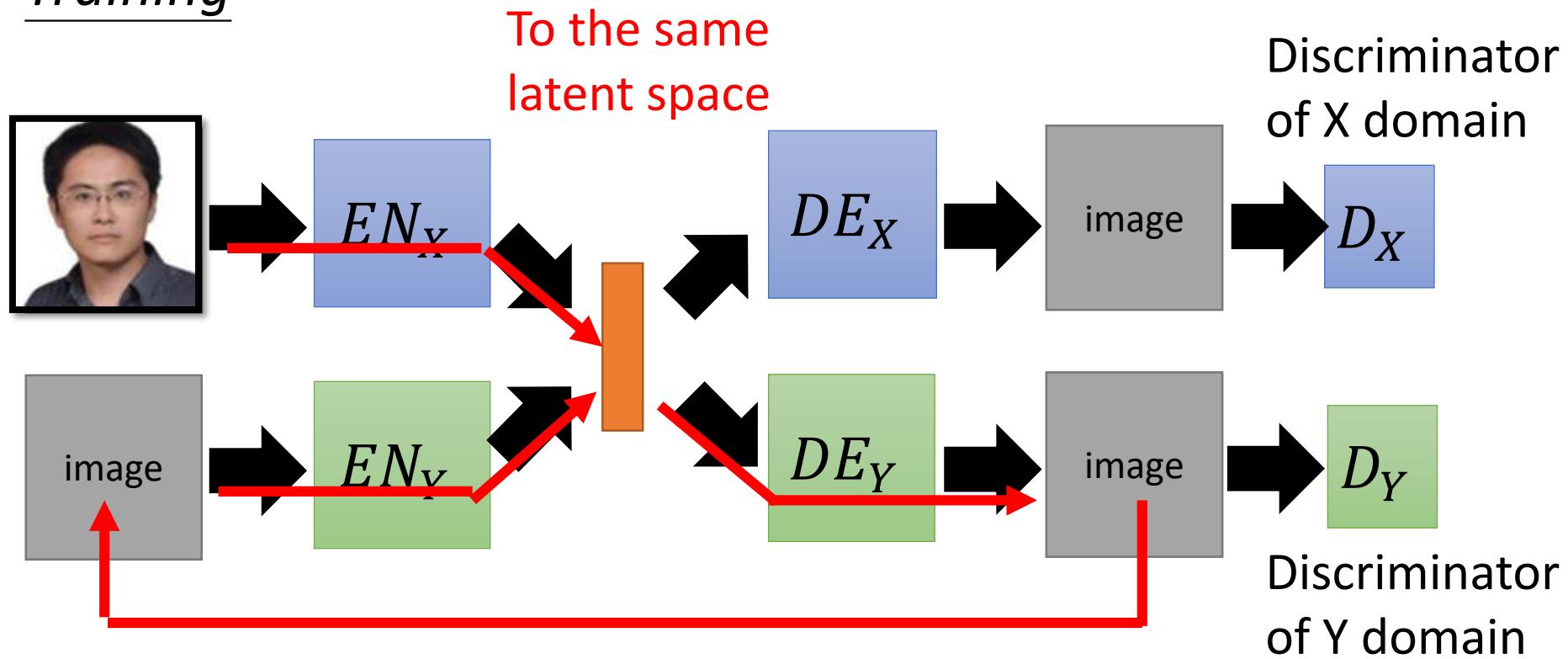


Cycle Consistency:

Used in ComboGAN [Asha Anoosheh, et al., arXiv, 017]

# Projection to Common Space

## Training



Semantic Consistency:

Used in DTN [Yaniv Taigman, et al., ICLR, 2017] and  
XGAN [Amélie Royer, et al., arXiv, 2017]

# 世界二次元化

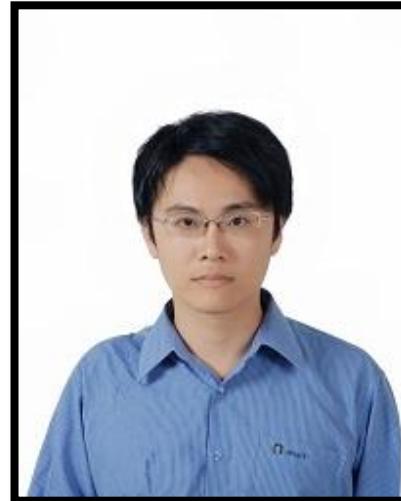
- Using the code:  
[https://github.com/Hiking/kawaii\\_creator](https://github.com/Hiking/kawaii_creator)
- It is not cycle GAN,  
Disco GAN



input



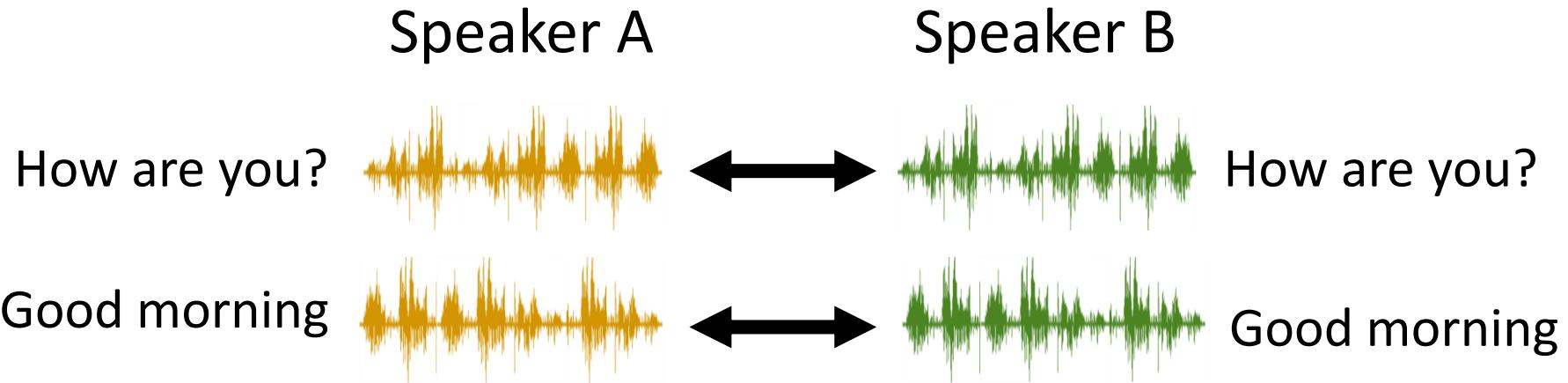
output domain



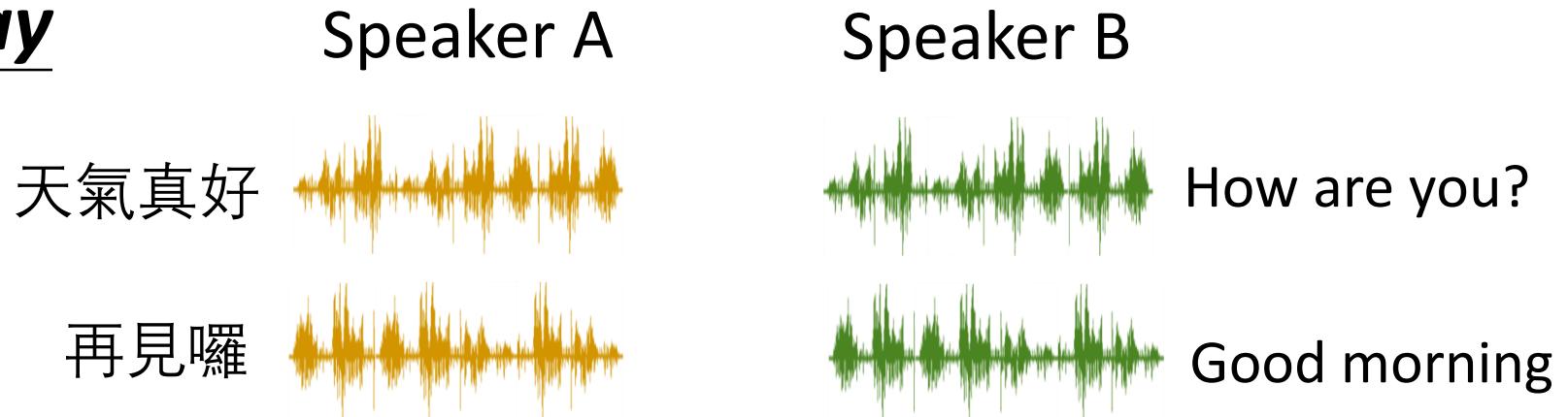
# Voice Conversion



## In the past



## Today



Speakers A and B are talking about completely different things.

Speaker A

我



Speaker B



感謝周儒杰同學提供實驗結果

# Reference

- Jun-Yan Zhu, Taesung Park, Phillip Isola, Alexei A. Efros, Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks, ICCV, 2017
- Zili Yi, Hao Zhang, Ping Tan, Minglun Gong, DualGAN: Unsupervised Dual Learning for Image-to-Image Translation, ICCV, 2017
- Tomer Galanti, Lior Wolf, Sagie Benaim, The Role of Minimal Complexity Functions in Unsupervised Learning of Semantic Mappings, ICLR, 2018
- Yaniv Taigman, Adam Polyak, Lior Wolf, Unsupervised Cross-Domain Image Generation, ICLR, 2017
- Asha Anoosheh, Eirikur Agustsson, Radu Timofte, Luc Van Gool, ComboGAN: Unrestrained Scalability for Image Domain Translation, arXiv, 2017
- Amélie Royer, Konstantinos Bousmalis, Stephan Gouws, Fred Bertsch, Inbar Mosseri, Forrester Cole, Kevin Murphy, XGAN: Unsupervised Image-to-Image Translation for Many-to-Many Mappings, arXiv, 2017

# Reference

- Guillaume Lample, Neil Zeghidour, Nicolas Usunier, Antoine Bordes, Ludovic Denoyer, Marc'Aurelio Ranzato, Fader Networks: Manipulating Images by Sliding Attributes, NIPS, 2017
- Taeksoo Kim, Moonsu Cha, Hyunsoo Kim, Jung Kwon Lee, Jiwon Kim, Learning to Discover Cross-Domain Relations with Generative Adversarial Networks, ICML, 2017
- Ming-Yu Liu, Oncel Tuzel, “Coupled Generative Adversarial Networks”, NIPS, 2016
- Ming-Yu Liu, Thomas Breuel, Jan Kautz, Unsupervised Image-to-Image Translation Networks, NIPS, 2017
- Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, Jaegul Choo, StarGAN: Unified Generative Adversarial Networks for Multi-Domain Image-to-Image Translation, arXiv, 2017

# Theory behind GAN

# Generation

Using Generative  
Adversarial  
Network (GAN)

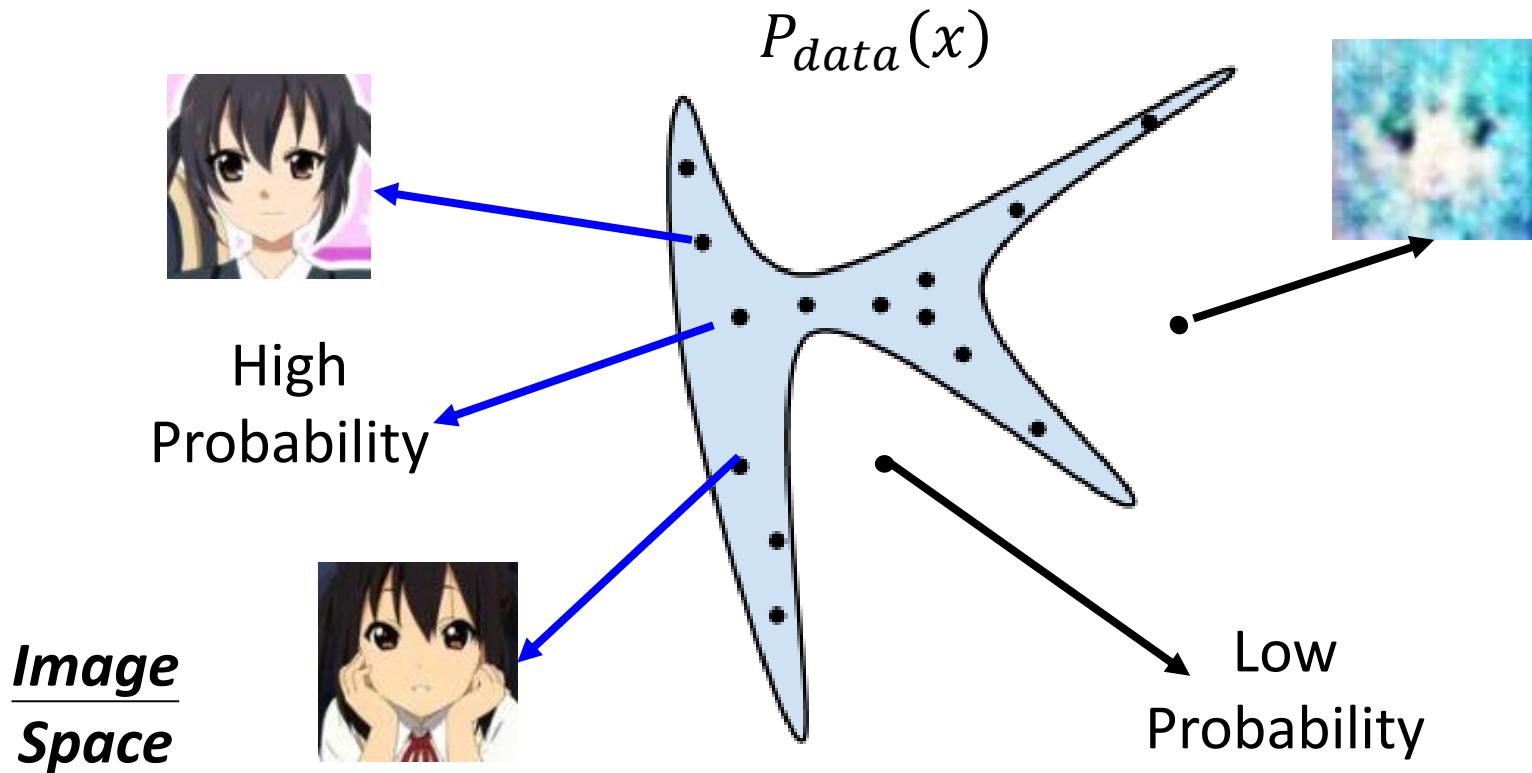


Drawing?

# Generation

$x$ : an image (a high-dimensional vector)

- We want to find data distribution  $P_{data}(x)$



# Maximum Likelihood Estimation

- Given a data distribution  $P_{data}(x)$  (We can sample from it.)
- We have a distribution  $P_G(x; \theta)$  parameterized by  $\theta$ 
  - We want to find  $\theta$  such that  $P_G(x; \theta)$  close to  $P_{data}(x)$
  - E.g.  $P_G(x; \theta)$  is a Gaussian Mixture Model,  $\theta$  are means and variances of the Gaussians

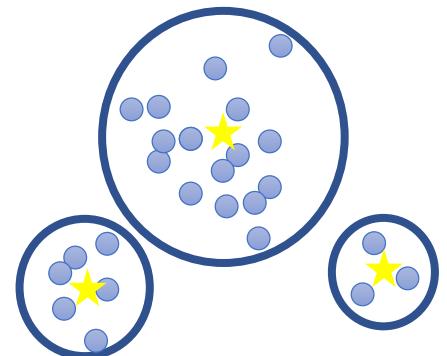
Sample  $\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$

We can compute  $P_G(x^i; \theta)$

Likelihood of generating the samples

$$L = \prod_{i=1}^m P_G(x^i; \theta)$$

Find  $\theta^*$  maximizing the likelihood



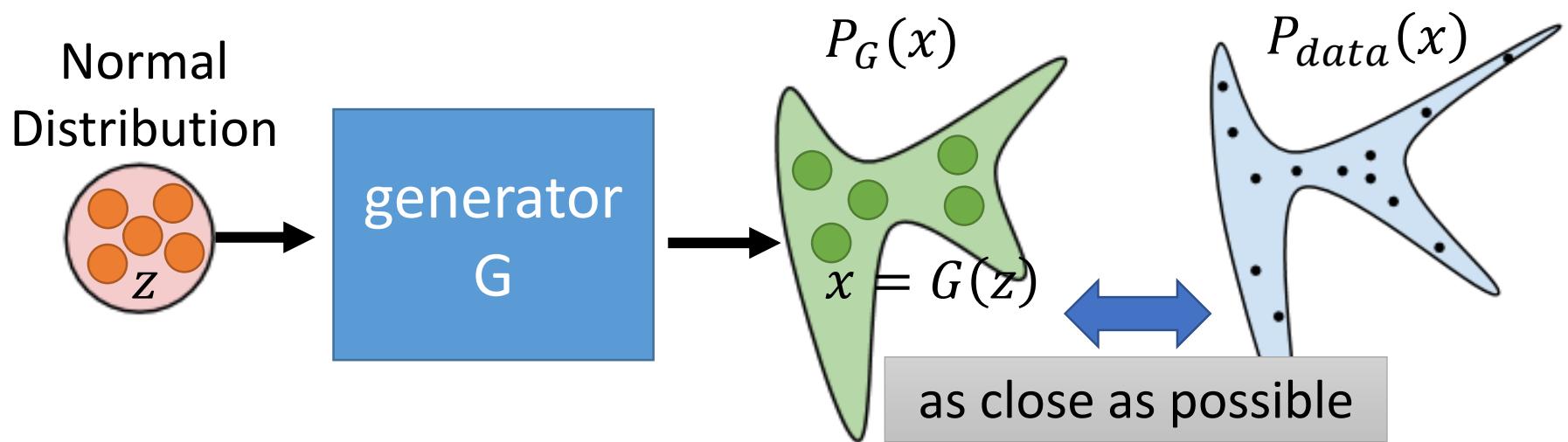
# Maximum Likelihood Estimation = Minimize KL Divergence

$$\begin{aligned}\theta^* &= \arg \max_{\theta} \prod_{i=1}^m P_G(x^i; \theta) = \arg \max_{\theta} \log \prod_{i=1}^m P_G(x^i; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^m \log P_G(x^i; \theta) \quad \{x^1, x^2, \dots, x^m\} \text{ from } P_{data}(x) \\ &\approx \arg \max_{\theta} E_{x \sim P_{data}} [\log P_G(x; \theta)] \\ &= \arg \max_{\theta} \int_x P_{data}(x) \log P_G(x; \theta) dx - \int_x P_{data}(x) \log P_{data}(x) dx \\ &= \arg \min_{\theta} KL(P_{data} || P_G) \quad \text{How to define a general } P_G?\end{aligned}$$

# Generator

$x$ : an image (a high-dimensional vector)

- A generator  $G$  is a network. The network defines a probability distribution  $P_G$



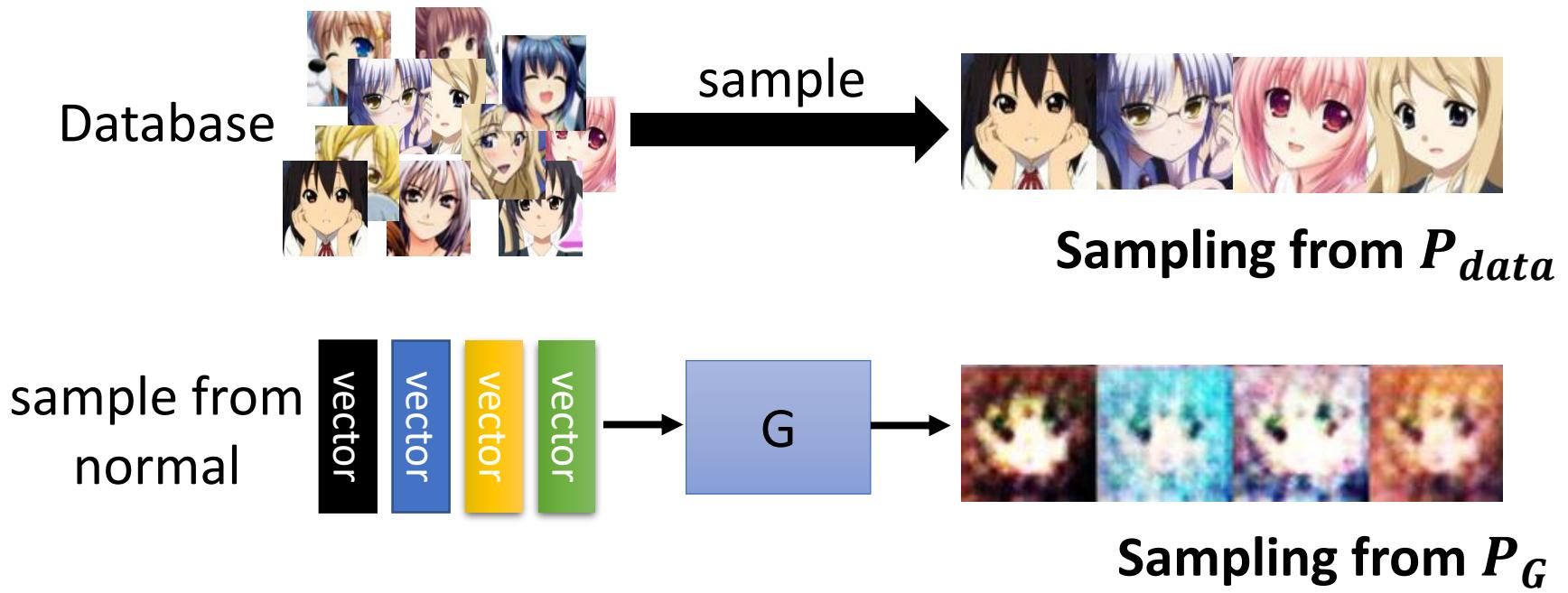
$$G^* = \arg \min_G \underline{Div(P_G, P_{data})}$$

Divergence between distributions  $P_G$  and  $P_{data}$   
How to compute the divergence?

# Discriminator

$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

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



# Discriminator

$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

- ★ : data sampled from  $P_{data}$
- ☆ : data sampled from  $P_G$



Using the example objective function is exactly the same as training a binary classifier.



## Example Objective Function for D

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

↑  
(G is fixed)

**Training:**  $D^* = \arg \max_D V(D, G)$

The maximum objective value is related to JS divergence.

# Discriminator

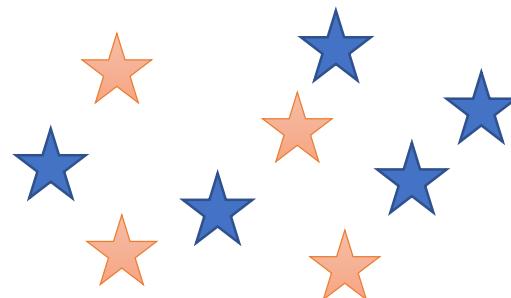
$$G^* = \arg \min_G \text{Div}(P_G, P_{data})$$

★ : data sampled from  $P_{data}$

☆ : data sampled from  $P_G$

**Training:**

$$D^* = \arg \max_D V(D, G)$$

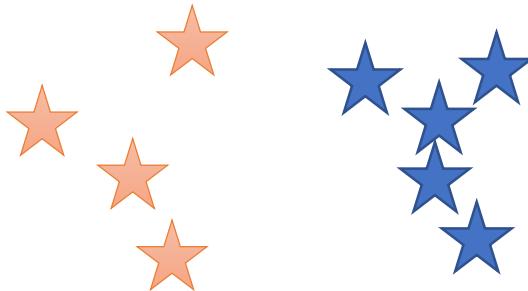


small divergence

train

Discriminator

hard to discriminate  
(cannot make objective large)



large divergence

train

Discriminator

easy to discriminate

$$\max_D V(G, D)$$

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

- Given G, what is the optimal D\* maximizing

$$\begin{aligned} V &= E_{x \sim P_{data}}[\log D(x)] + E_{x \sim P_G}[\log(1 - D(x))] \\ &= \int_x P_{data}(x) \log D(x) dx + \int_x P_G(x) \log(1 - D(x)) dx \\ &= \int_x [P_{data}(x) \log D(x) + P_G(x) \log(1 - D(x))] dx \end{aligned}$$

Assume that D(x) can be any function

- Given x, the optimal D\* maximizing

$$P_{data}(x) \log D(x) + P_G(x) \log(1 - D(x))$$

$$\max_D V(G, D)$$

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

- Given  $x$ , the optimal  $D^*$  maximizing

$$P_{data}(x) \log D(x) + P_G(x) \log(1 - D(x))$$

**a**                    **D**                    **b**                    **D**

- Find  $D^*$  maximizing:  $f(D) = a \log(D) + b \log(1 - D)$

$$\frac{df(D)}{dD} = a \times \frac{1}{D} + b \times \frac{1}{1-D} \times (-1) = 0$$

$$a \times \frac{1}{D^*} = b \times \frac{1}{1 - D^*} \quad a \times (1 - D^*) = b \times D^*$$

$$a - aD^* = bD^* \quad a = (a + b)D^*$$

$$D^* = \frac{a}{a + b} \quad \rightarrow \quad 0 <$$

$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} < 1$$

$$\max_D V(G, D)$$

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

$$\max_D V(G, D) = V(G, D^*)$$

$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)}$$

$$= E_{x \sim P_{data}} \left[ \log \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} \right]$$

$$+ E_{x \sim P_G} \left[ \log \frac{P_G(x)}{P_{data}(x) + P_G(x)} \right]$$

$$= \int_x P_{data}(x) \log \frac{\frac{1}{2} P_{data}(x)}{\underline{\frac{P_{data}(x) + P_G(x)}{2}}} dx$$

$$+ 2 \log \frac{1}{2} - 2 \log 2$$

$$+ \int_x P_G(x) \log \frac{\frac{1}{2} P_G(x)}{\underline{\frac{P_{data}(x) + P_G(x)}{2}}} dx$$

$$\max_D V(G, D)$$

$$\begin{aligned} \text{JSD}(P \parallel Q) &= \frac{1}{2}D(P \parallel M) + \frac{1}{2}D(Q \parallel M) \\ M &= \frac{1}{2}(P + Q) \end{aligned}$$

$$\max_D V(G, D) = V(G, D^*)$$

$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)}$$

$$\begin{aligned} &= -2\log 2 + \int_x P_{data}(x) \log \frac{P_{data}(x)}{(P_{data}(x) + P_G(x))/2} dx \\ &\quad + \int_x P_G(x) \log \frac{P_G(x)}{(P_{data}(x) + P_G(x))/2} dx \end{aligned}$$

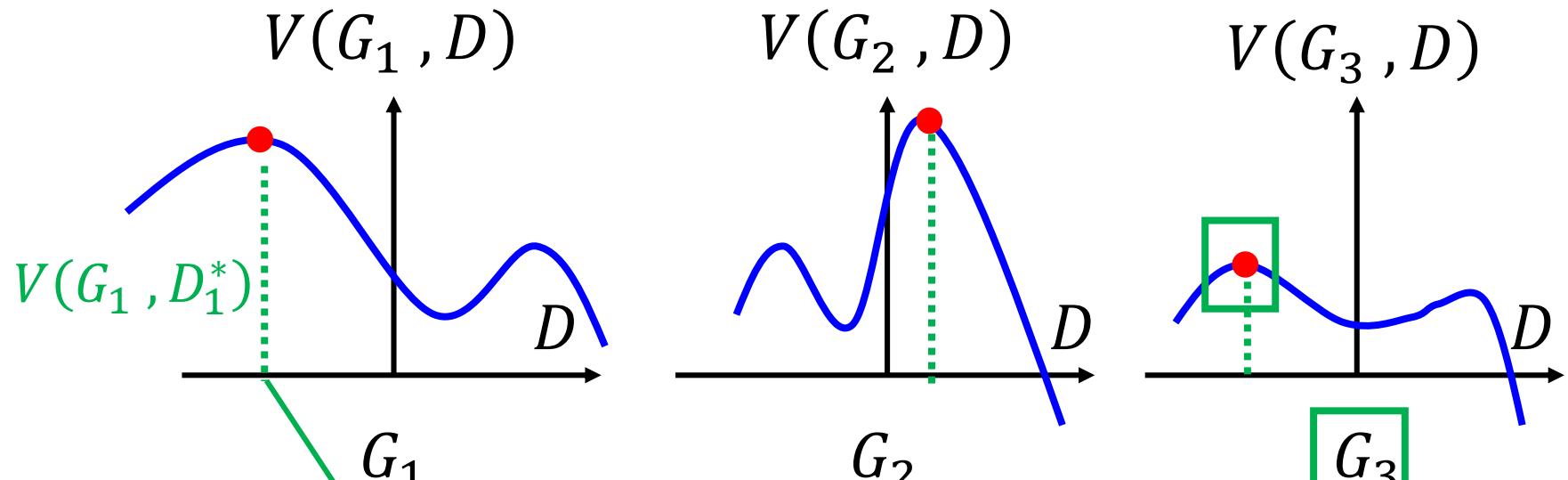
$$= -2\log 2 + \text{KL}\left(P_{data} \parallel \frac{P_{data} + P_G}{2}\right) + \text{KL}\left(P_G \parallel \frac{P_{data} + P_G}{2}\right)$$

$$= -2\log 2 + 2\text{JSD}(P_{data} \parallel P_G) \quad \text{Jensen-Shannon divergence}$$

$$G^* = \arg \min_G \max_D V(G, D)$$

$$D^* = \arg \max_D V(D, G)$$

The maximum objective value is related to JS divergence.



Divergence between  $P_{G_1}$  and  $P_{data}$

$$G^* = \arg \min_G \max_D V(G, D)$$

$$D^* = \arg \max_D V(D, G)$$

The maximum objective value is related to JS divergence.

- 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  $G$

# Algorithm

$$G^* = \arg \min_G \max_D V(G, D)$$

$L(G)$

- To find the best G minimizing the loss function  $L(G)$ ,

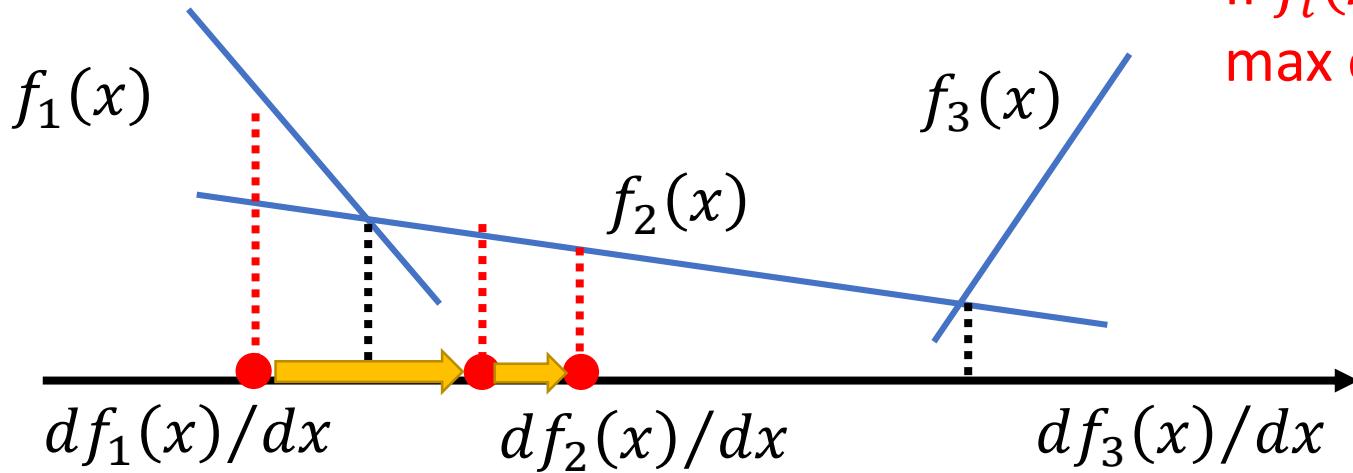
$$\theta_G \leftarrow \theta_G - \eta \partial L(G) / \partial \theta_G$$

$\theta_G$  defines G

$$f(x) = \max\{f_1(x), f_2(x), f_3(x)\}$$

$$\frac{df(x)}{dx} = ? \quad df_i(x)/dx$$

If  $f_i(x)$  is the  
max one



# Algorithm

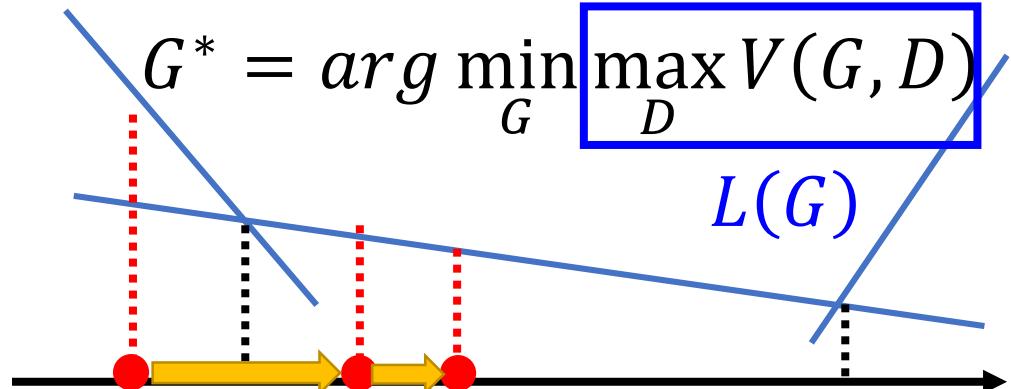
- Given  $G_0$
- Find  $D_0^*$  maximizing  $V(G_0, D)$  **Using Gradient Ascent**

$V(G_0, D_0^*)$  is the JS divergence between  $P_{data}(x)$  and  $P_{G_0}(x)$

- $\theta_G \leftarrow \theta_G - \eta \partial V(G, D_0^*) / \partial \theta_G \rightarrow$  Obtain  $G_1$  Decrease JS divergence(?)
- Find  $D_1^*$  maximizing  $V(G_1, D)$

$V(G_1, D_1^*)$  is the JS divergence between  $P_{data}(x)$  and  $P_{G_1}(x)$

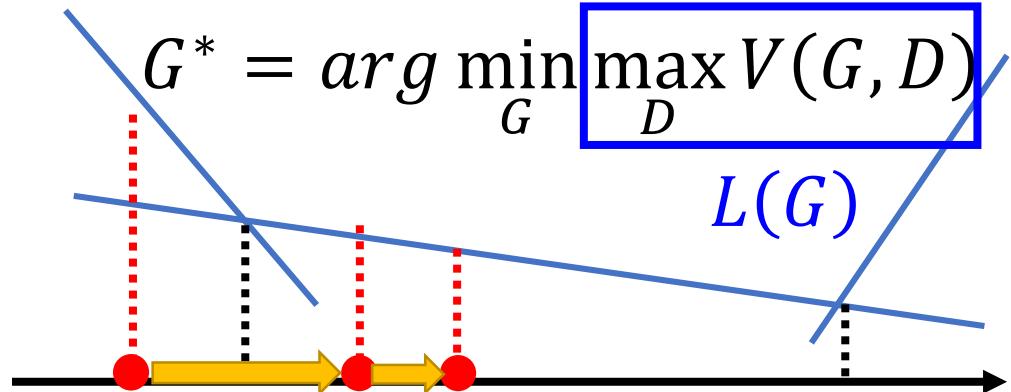
- $\theta_G \leftarrow \theta_G - \eta \partial V(G, D_1^*) / \partial \theta_G \rightarrow$  Obtain  $G_2$  Decrease JS divergence(?)
- .....



# Algorithm

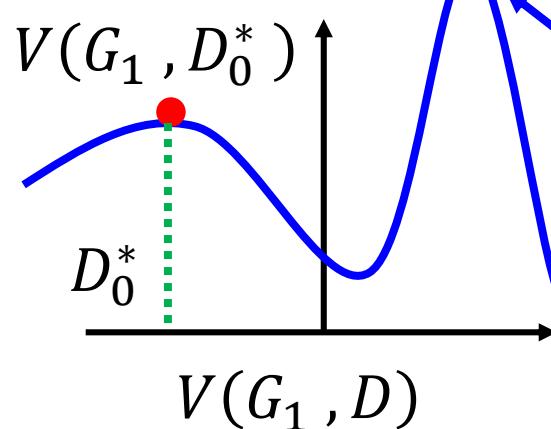
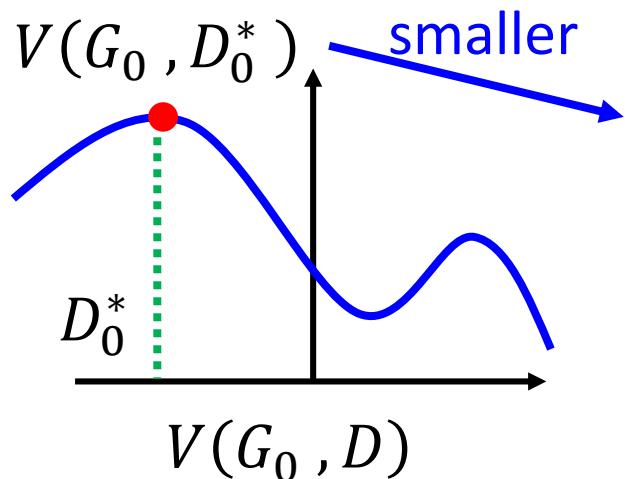
- Given  $G_0$
- Find  $D_0^*$  maximizing  $V(G_0, D)$

$V(G_0, D_0^*)$  is the JS divergence between  $P_{data}(x)$  and  $P_{G_0}(x)$



- $\theta_G \leftarrow \theta_G - \eta \partial V(G, D_0^*) / \partial \theta_G \rightarrow$  Obtain  $G_1$

Decrease JS divergence(?)



$V(G_1, D_1^*) \dots$

Assume  $D_0^* \approx D_1^*$

**Don't update G too much**

# In practice ...

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

- Given  $G$ , how to compute  $\max_D V(G, D)$ 
  - Sample  $\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$ , sample  $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$  from generator  $P_G(x)$

Maximize  $\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log(1 - D(\tilde{x}^i))$

||

## Binary Classifier

$D$  is a binary classifier with sigmoid output (can be deep)

$\{x^1, x^2, \dots, x^m\}$  from  $P_{data}(x)$   $\rightarrow$  Positive examples

$\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$  from  $P_G(x)$   $\rightarrow$  Negative examples

Minimize Cross-entropy

# Algorithm

Initialize  $\theta_d$  for D and  $\theta_g$  for G

- In each training iteration:

Can only find  
lower bound of

$$\max_D V(G, D)$$

- Sample m examples  $\{x^1, x^2, \dots, x^m\}$  from data distribution  $P_{data}(x)$
- Sample m noise samples  $\{z^1, z^2, \dots, z^m\}$  from the prior  $P_{prior}(z)$
- Obtaining generated data  $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$ ,  $\tilde{x}^i = G(z^i)$
- Update discriminator parameters  $\theta_d$  to maximize

- $$\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(\tilde{x}^i))$$
- $$\theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$$

- Sample another m noise samples  $\{z^1, z^2, \dots, z^m\}$  from the prior  $P_{prior}(z)$

- Update generator parameters  $\theta_g$  to minimize

- $$\tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) + \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^i)))$$
- $$\theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$$

Learning  
D

Repeat  
k times

Learning  
G

Only  
Once

# Objective Function for Generator in Real Implementation

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

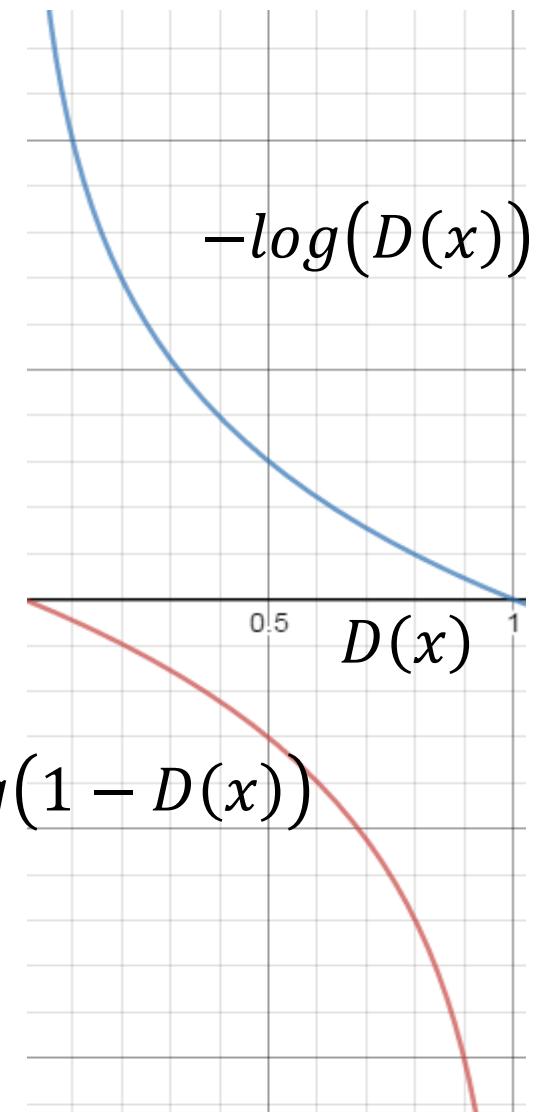
Slow at the beginning

Minimax GAN (MMGAN)

$$V = E_{x \sim P_G} [-\log(D(x))]$$

Real implementation:  
label  $x$  from  $P_G$  as positive

Non-saturating GAN (NSGAN)

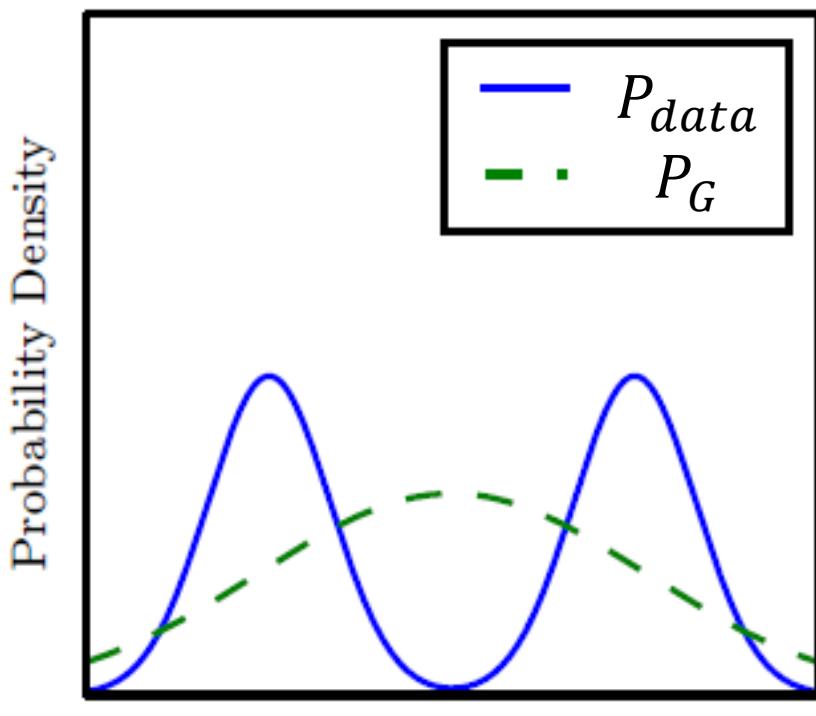


# fGAN: General Framework of GAN

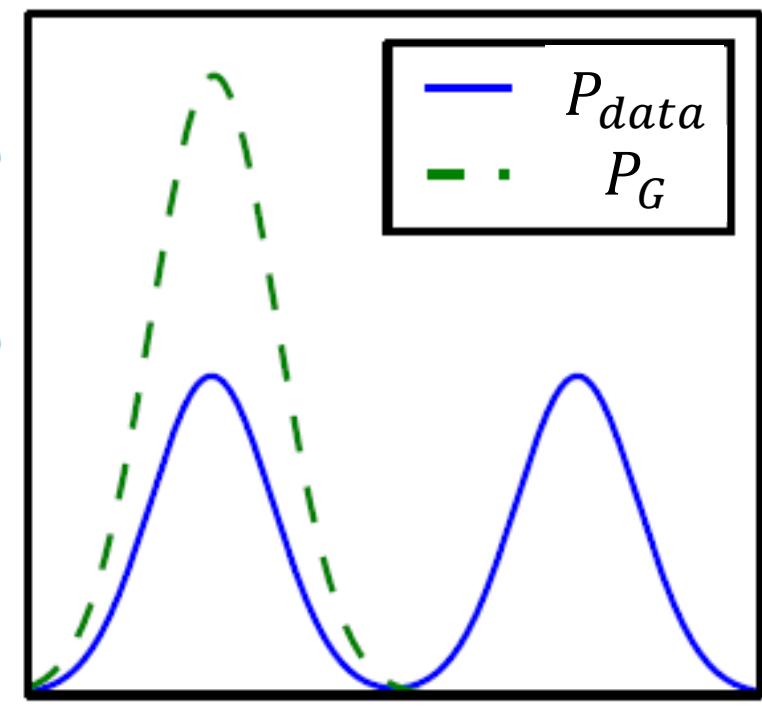
# Flaw in Optimization?

$$KL = \int P_{data} \log \frac{P_{data}}{P_G} dx$$

$$\text{Reverse } KL = \int P_G \log \frac{P_G}{P_{data}} dx$$



$x$   
Maximum likelihood  
(minimize  $KL(P_{data} || P_G)$ )



$x$   
Minimize  $KL(P_G || P_{data})$   
(reverse KL)

## f-divergence

$P$  and  $Q$  are two distributions.  $p(x)$  and  $q(x)$  are the probability of sampling  $x$ .

$$D_f(P||Q) = \int_x q(x)f\left(\frac{p(x)}{q(x)}\right)dx \quad \begin{array}{l} f \text{ is convex} \\ f(1) = 0 \end{array}$$

$D_f(P||Q)$  evaluates the difference of P and Q

If  $p(x) = q(x)$  for all  $x$

smallest

$$D_f(P||Q) = \int_x q(x)f\left(\frac{p(x)}{q(x)}\right)dx = 0$$

$$D_f(P||Q) = \int_x q(x)f\left(\frac{p(x)}{q(x)}\right)dx$$

Because  $f$  is convex

$$\geq f\left(\int_x q(x) \frac{p(x)}{q(x)} dx\right)$$

$$= f(1) = 0$$

If P and Q are the same distributions,  
 $D_f(P||Q)$  has the smallest value, which is 0

## f-divergence

$$D_f(P||Q) = \int_x q(x)f\left(\frac{p(x)}{q(x)}\right)dx \quad \begin{array}{l} f \text{ is convex} \\ f(1) = 0 \end{array}$$

---

$$f(x) = x \log x$$

$$D_f(P||Q) = \int_x q(x) \frac{p(x)}{q(x)} \log \left( \frac{p(x)}{q(x)} \right) dx = \int_x p(x) \log \left( \frac{p(x)}{q(x)} \right) dx \quad \text{KL}$$

$$f(x) = -\log x$$

$$D_f(P||Q) = \int_x q(x) \left( -\log \left( \frac{p(x)}{q(x)} \right) \right) dx = \int_x q(x) \log \left( \frac{q(x)}{p(x)} \right) dx \quad \text{Reverse KL}$$

$$f(x) = (x - 1)^2$$

$$D_f(P||Q) = \int_x q(x) \left( \frac{p(x)}{q(x)} - 1 \right)^2 dx = \int_x \frac{(p(x) - q(x))^2}{q(x)} dx \quad \text{Chi Square}$$

# Fenchel Conjugate

$$D_f(P||Q) = \int_x q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

f is convex,  $f(1) = 0$

- Every convex function f has a conjugate function  $f^*$

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\}$$

$$f^*(\mathbf{t}_1) = \max_{x \in \text{dom}(f)} \{x\mathbf{t}_1 - f(x)\}$$

$$x_1 \mathbf{t}_1 - f(x_1)$$



$$f^*(\mathbf{t}_2) = \max_{x \in \text{dom}(f)} \{x\mathbf{t}_2 - f(x)\}$$

$$x_2 \mathbf{t}_1 - f(x_2)$$



$$x_3 \mathbf{t}_2 - f(x_3)$$



$$x_3 \mathbf{t}_1 - f(x_3)$$



$$x_2 \mathbf{t}_2 - f(x_2)$$



$$x_1 \mathbf{t}_2 - f(x_1)$$



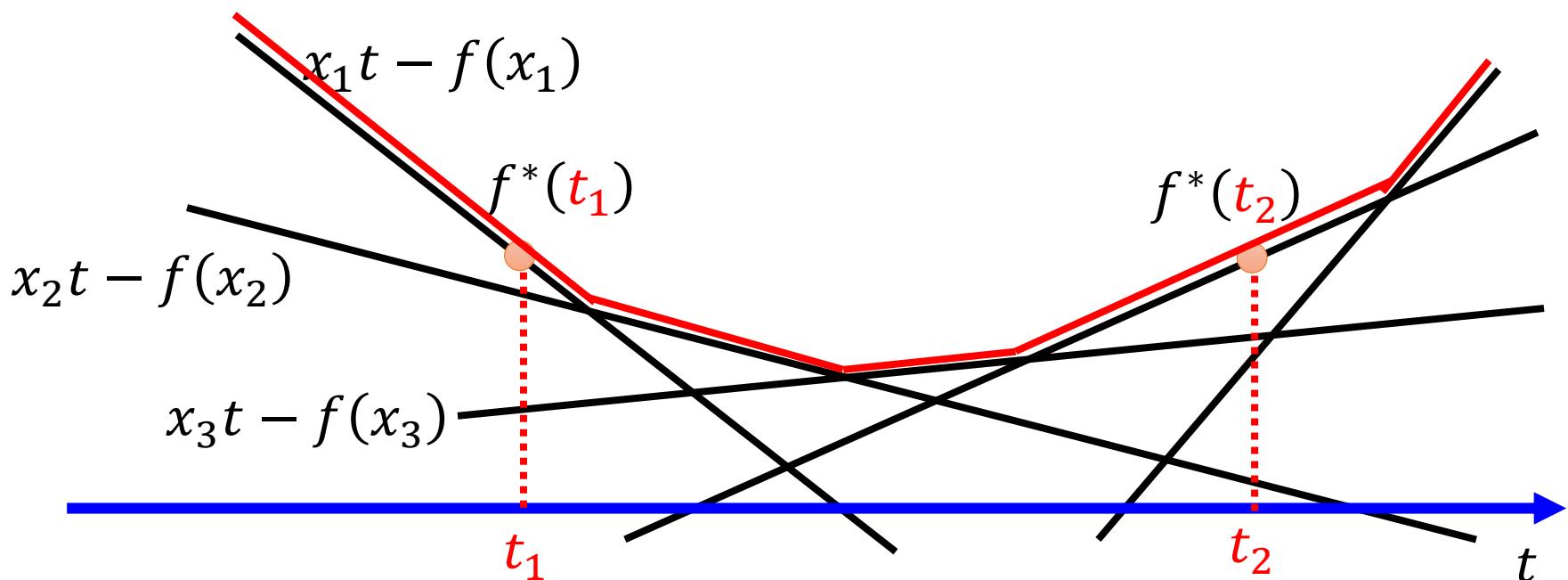
# Fenchel Conjugate

$$D_f(P||Q) = \int_x q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

f is convex,  $f(1) = 0$

- Every convex function f has a conjugate function  $f^*$

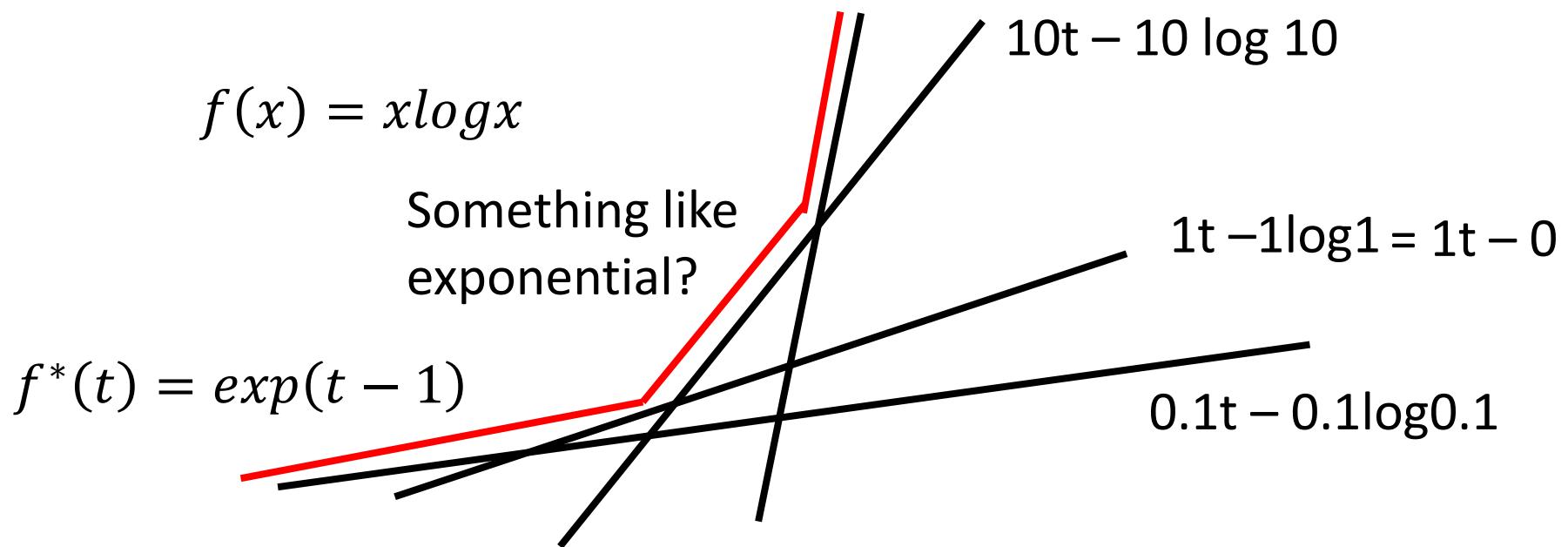
$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\}$$



# Fenchel Conjugate

- Every convex function  $f$  has a conjugate function  $f^*$

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\}$$



# Fenchel Conjugate

- Every convex function  $f$  has a conjugate function  $f^*$
- $(f^*)^* = f$

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\}$$

$$f(x) = x \log x \quad \longleftrightarrow \quad f^*(t) = \exp(t - 1)$$

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - x \log x\}$$

$$g(x) = xt - x \log x \quad \text{Given } t, \text{ find } x \text{ maximizing } g(x)$$

$$t - \log x - 1 = 0 \quad x = \exp(t - 1)$$

$$f^*(t) = \exp(t - 1) \times t - \exp(t - 1) \times (t - 1) = \exp(t - 1)$$

# Connection with GAN

$$f^*(t) = \max_{x \in \text{dom}(f)} \{xt - f(x)\} \leftrightarrow f(\underline{x}) = \max_{t \in \text{dom}(f^*)} \{\underline{x}t - f^*(t)\}$$

$$D_f(P||Q) = \int_x q(x) f\left(\frac{p(x)}{q(x)}\right) dx \quad \boxed{\frac{p(x)}{q(x)}} \quad \boxed{\frac{p(x)}{q(x)}}$$

$$= \int_x q(x) \left( \max_{t \in \text{dom}(f^*)} \left\{ \frac{p(x)}{q(x)} \underline{t} - f^*(\underline{t}) \right\} \right) dx$$

$$\approx \max_D \int_x p(x) D(x) dx - \int_x q(x) f^*(D(x)) dx$$

D is a function  
whose input is x,  
and output is t

$$D_f(P||Q) \geq \int_x q(x) \left( \frac{p(x)}{q(x)} \underline{D(x)} - f^*(\underline{D(x)}) \right) dx$$

$$= \int_x p(x) D(x) dx - \int_x q(x) f^*(D(x)) dx$$

# Connection with GAN

$$\begin{aligned} D_f(P||Q) &\approx \max_{\text{D}} \int_x p(x)D(x)dx - \int_x q(x)f^*(D(x))dx \\ &= \max_{\text{D}} \left\{ E_{x \sim P}[D(x)] - E_{x \sim Q}[f^*(D(x))] \right\} \end{aligned}$$

Samples from P              Samples from Q

$$D_f(P_{data}||P_G) = \max_{\text{D}} \left\{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[f^*(D(x))] \right\}$$

$$\begin{aligned} G^* &= \arg \min_G D_f(P_{data}||P_G) && \text{Original GAN has} \\ &= \arg \min_G \max_D \left\{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[f^*(D(x))] \right\} && \text{different V(G,D)} \\ &= \arg \min_G \max_D V(G, D) \quad \text{familiar? } \smiley \end{aligned}$$

$$D_f(P_{data} || P_G) = \max_D \{ E_{x \sim P_{data}} [D(x)] - E_{x \sim P_G} [f^*(D(x))] \}$$

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

Using the f-divergence  
you like ☺

<https://arxiv.org/pdf/1606.00709.pdf>

Name	Conjugate $f^*(t)$
Total variation	$t$
Kullback-Leibler (KL)	$\exp(t - 1)$
Reverse KL	$-1 - \log(-t)$
Pearson $\chi^2$	$\frac{1}{4}t^2 + t$
Neyman $\chi^2$	$2 - 2\sqrt{1-t}$
Squared Hellinger	$\frac{t}{1-t}$
Jeffrey	$W(e^{1-t}) + \frac{1}{W(e^{1-t})} + t - 2$
Jensen-Shannon	$-\log(2 - \exp(t))$
Jensen-Shannon-weighted	$(1 - \pi) \log \frac{1 - \pi}{1 - \pi e^{t/\pi}}$
GAN	$-\log(1 - \exp(t))$

# Tips for Improving GAN

Martin Arjovsky, Soumith Chintala, Léon Bottou, Wasserstein GAN, arXiv preprint, 2017

Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, Aaron Courville,  
“Improved Training of Wasserstein GANs”, arXiv preprint, 2017

# JS divergence is not suitable

- In most cases,  $P_G$  and  $P_{data}$  are not overlapped.
- 1. The nature of data

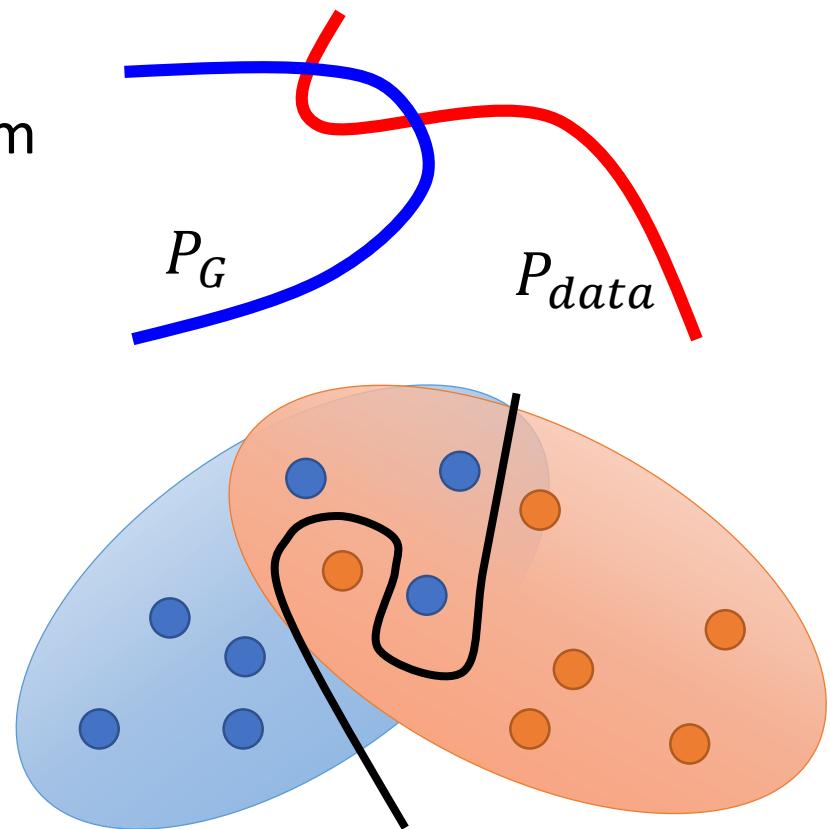
Both  $P_{data}$  and  $P_G$  are low-dim manifold in high-dim space.

The overlap can be ignored.

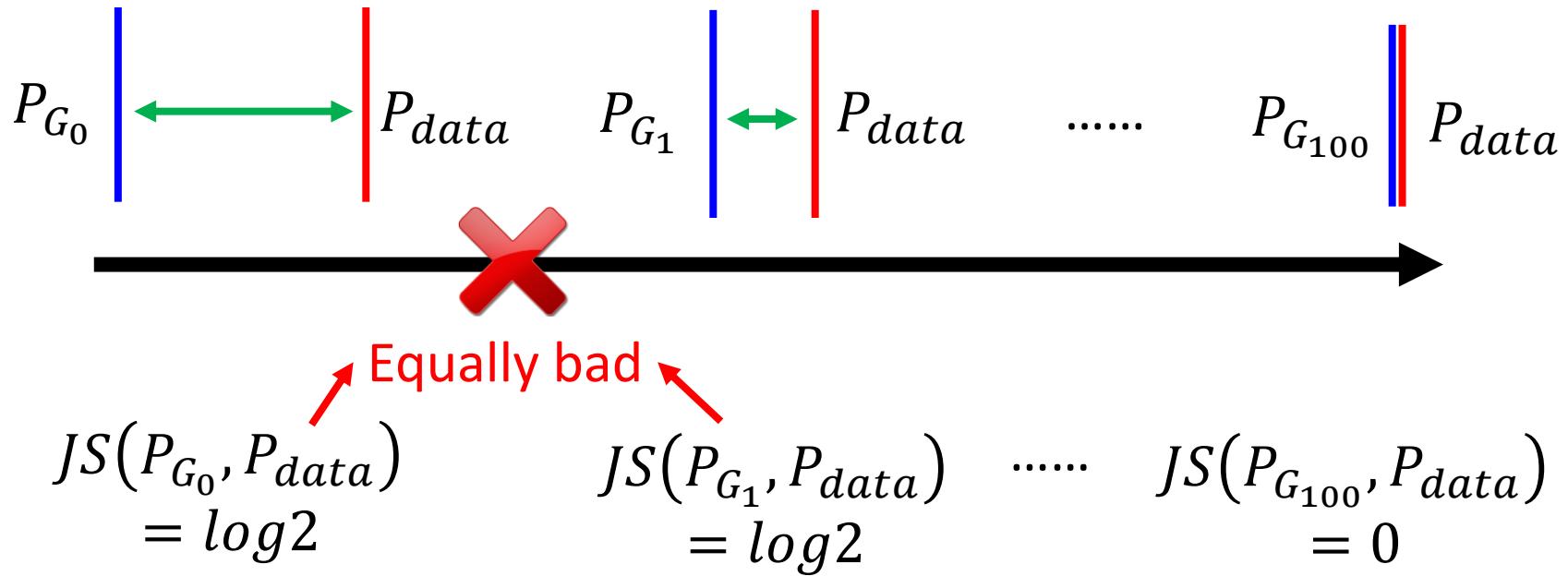
- 2. Sampling

Even though  $P_{data}$  and  $P_G$  have overlap.

If you do not have enough sampling .....



# What is the problem of JS divergence?

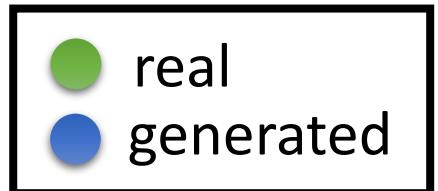


JS divergence is  $\log 2$  if two distributions do not overlap.

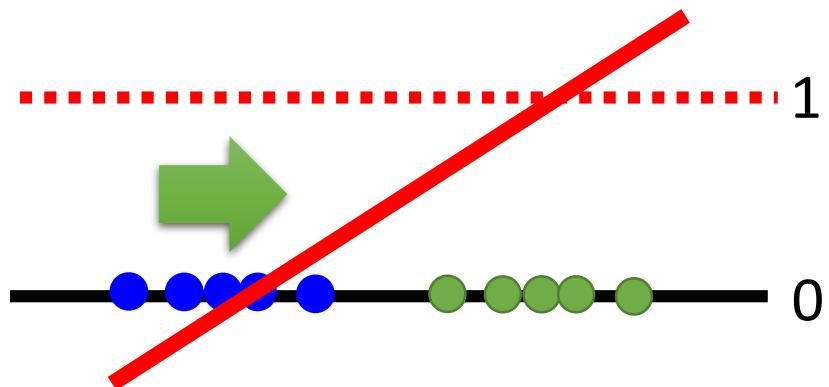
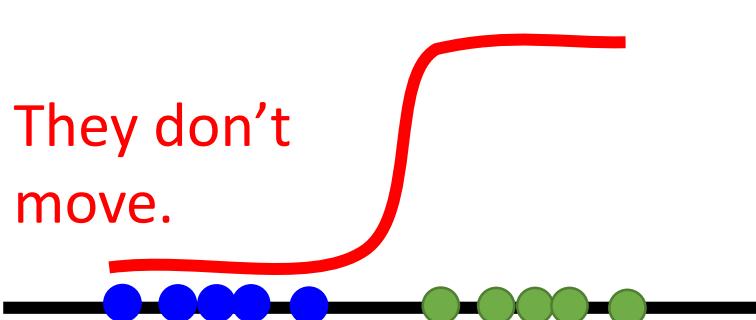
Intuition: If two distributions do not overlap, binary classifier achieves 100% accuracy

→ Same objective value is obtained. → Same divergence

# Least Square GAN (LSGAN)

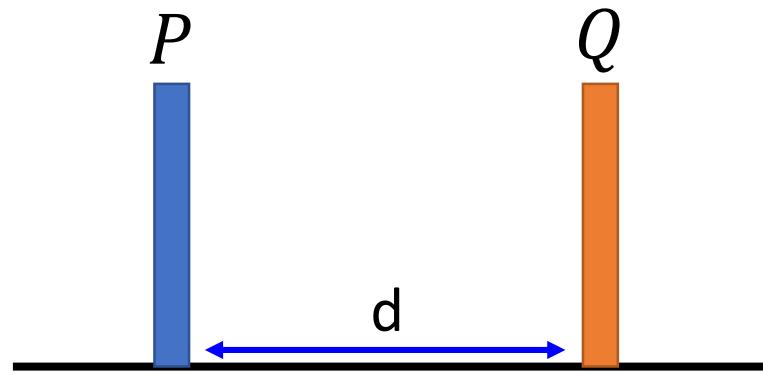


- Replace sigmoid with linear (replace classification with regression)



# Wasserstein GAN (WGAN): Earth Mover's Distance

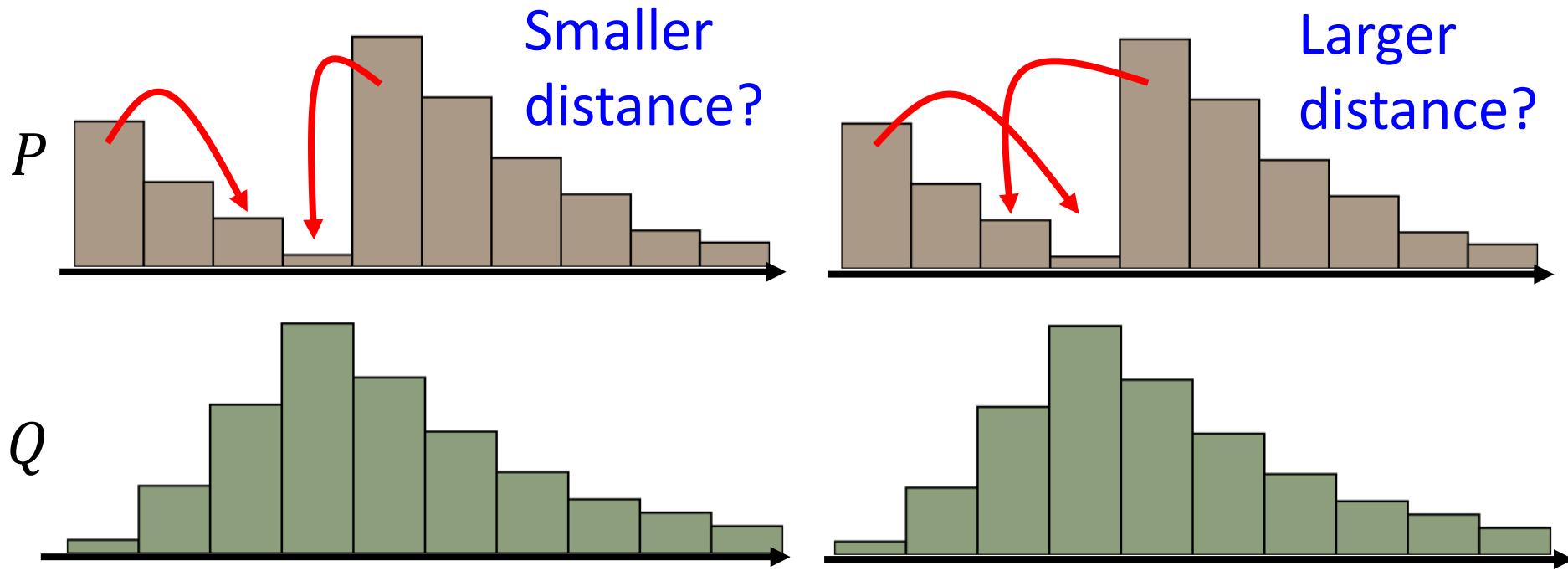
- Considering one distribution P as a pile of earth, and another distribution Q as the target
- The average distance the earth mover has to move the earth.



$$W(P, Q) = d$$



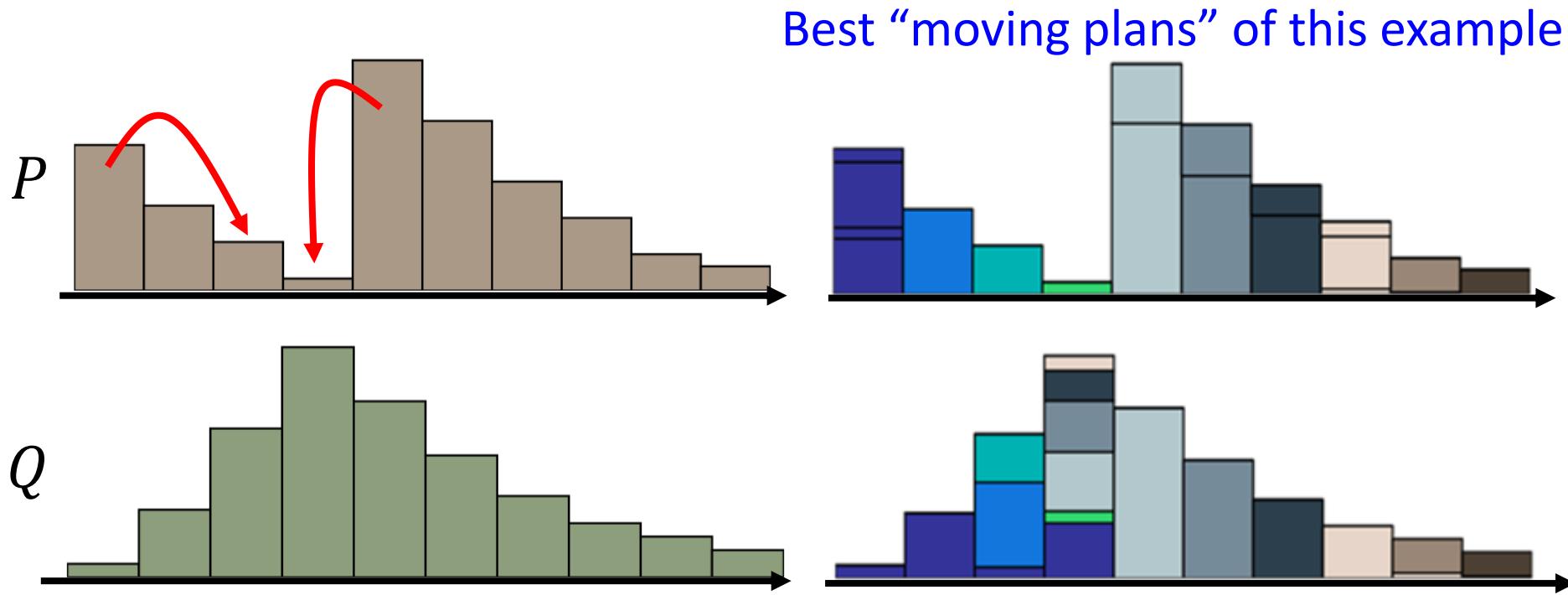
# WGAN: Earth Mover's Distance



There many possible “moving plans”.

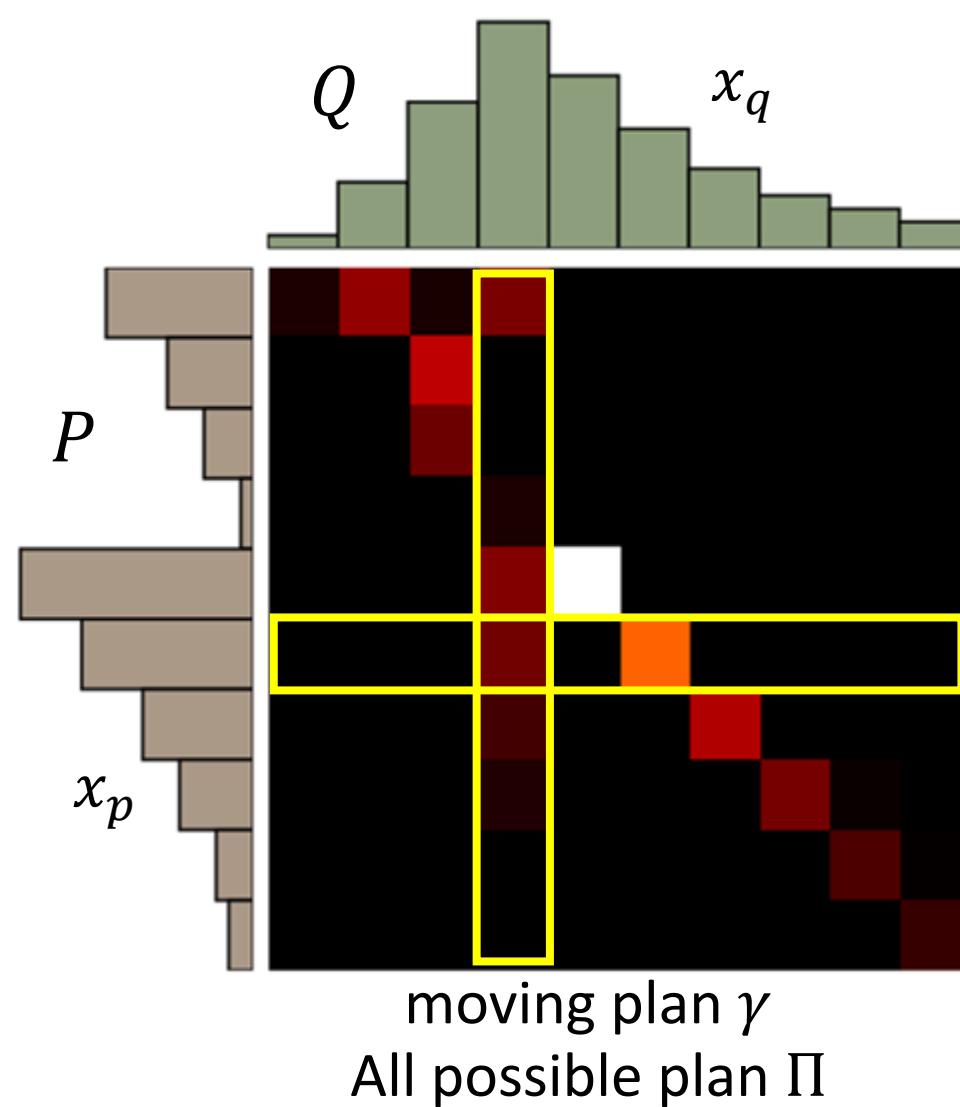
Using the “moving plan” with the smallest average distance to define the earth mover’s distance.

# WGAN: Earth Mover's Distance



There many possible “moving plans”.

Using the “moving plan” with the smallest average distance to define the earth mover’s distance.



A “moving plan” is a matrix  
The value of the element is the  
amount of earth from one  
position to another.

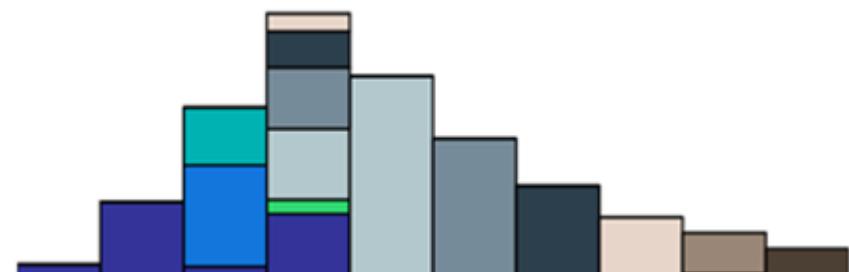
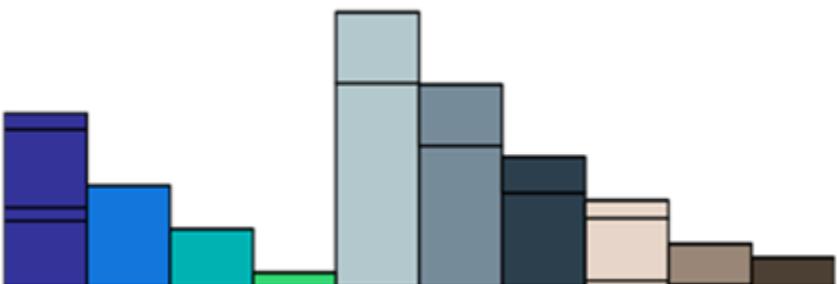
Average distance of a plan  $\gamma$ :

$$B(\gamma) = \sum_{x_p, x_q} \gamma(x_p, x_q) \|x_p - x_q\|$$

Earth Mover’s Distance:

$$W(P, Q) = \min_{\gamma \in \Pi} B(\gamma)$$

The best plan

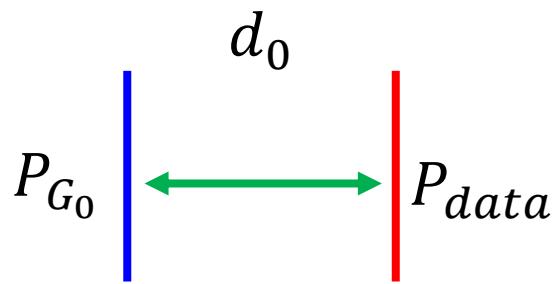
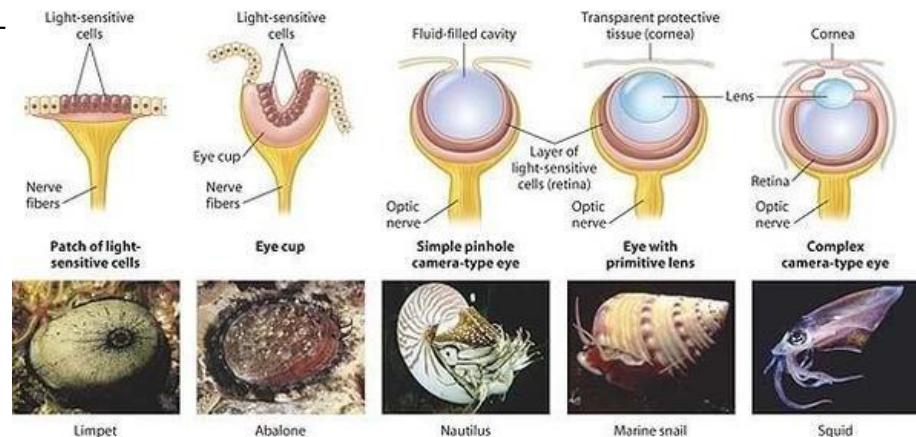


# Why Earth Mover's Distance?

$$D_f(P_{data} || P_G)$$

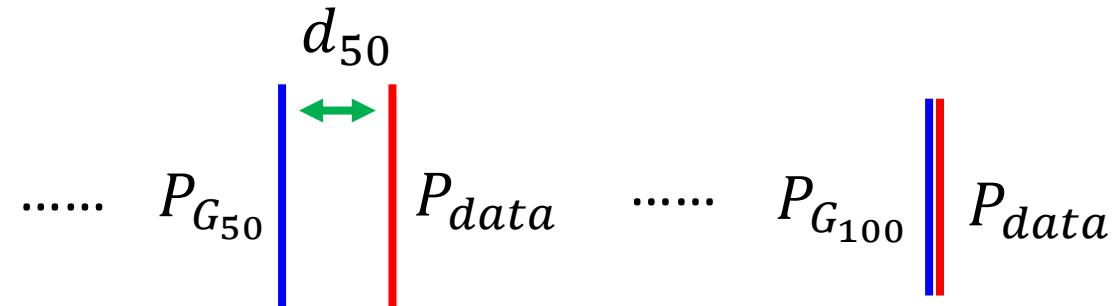


$$W(P_{data}, P_G)$$



$$JS(P_{G_0}, P_{data}) = \log 2$$

$$W(P_{G_0}, P_{data}) = d_0$$



$$JS(P_{G_{50}}, P_{data}) = \log 2$$

$$W(P_{G_{50}}, P_{data}) = d_{50}$$

$$JS(P_{G_{100}}, P_{data}) = 0$$

$$W(P_{G_{100}}, P_{data}) = 0$$

# WGAN

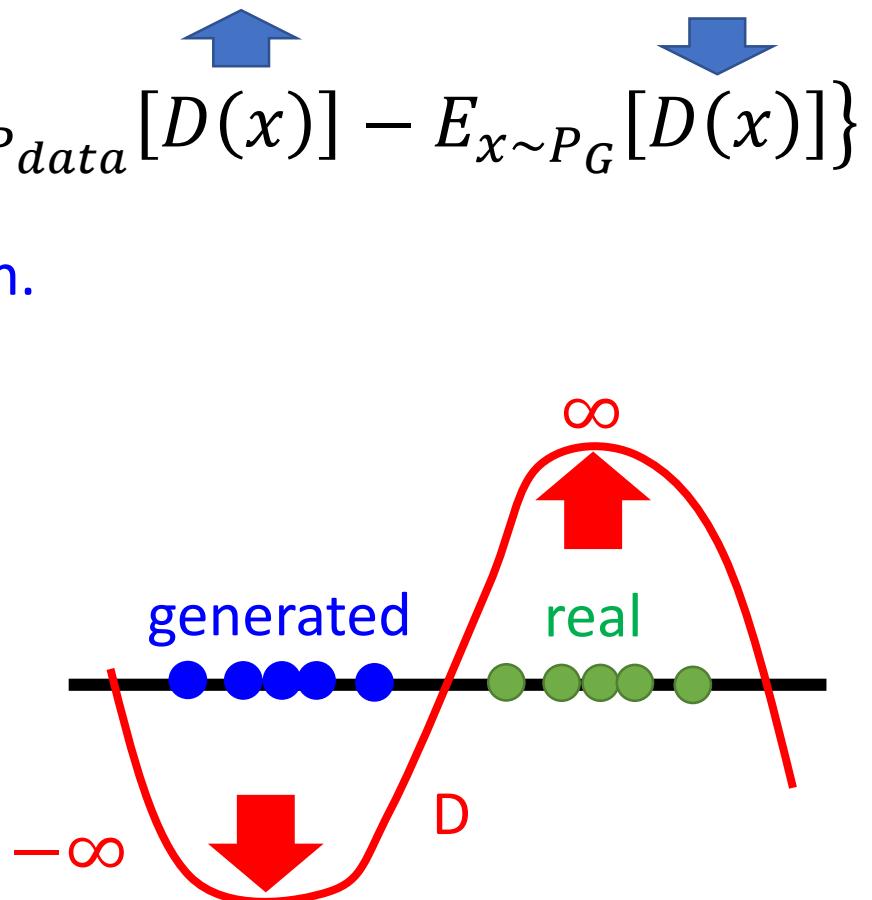
Evaluate wasserstein distance between  $P_{data}$  and  $P_G$

$$V(G, D) = \max_{D \in \text{1-Lipschitz}} \left\{ E_{x \sim P_{data}} [D(x)] - E_{x \sim P_G} [D(x)] \right\}$$

D has to be smooth enough.

Without the constraint, the training of D will not converge.

Keeping the D smooth forces  $D(x)$  become  $\infty$  and  $-\infty$



## Weight Clipping [Martin Arjovsky, et al., arXiv, 2017]

WGAN

Force the parameters w between c and -c  
After parameter update, if  $w > c$ ,  $w = c$ ;  
if  $w < -c$ ,  $w = -c$

Evaluate wasserstein distance between  $P_{data}$  and  $P_G$

$$V(G, D) = \max_{D \in \text{1-Lipschitz}} \left\{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)] \right\}$$

D has to be smooth enough. How to fulfill this constraint?

### Lipschitz Function

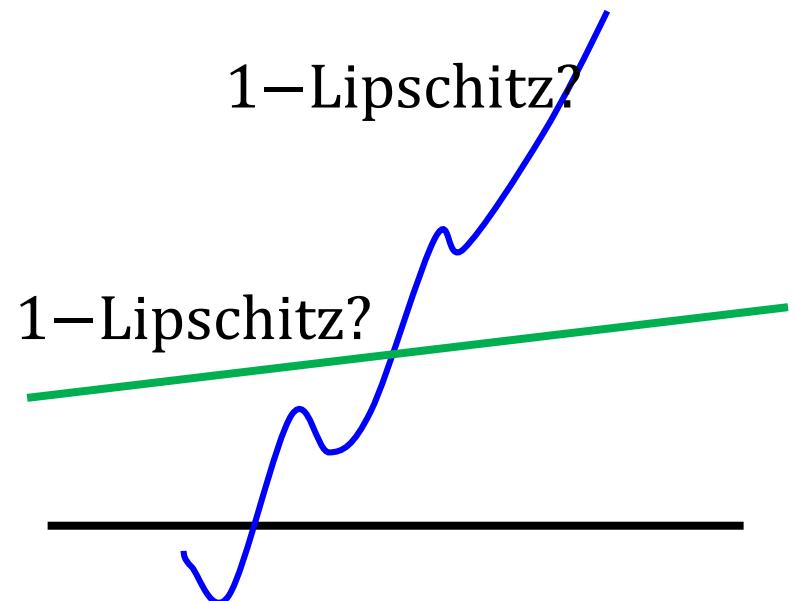
$$\|f(x_1) - f(x_2)\| \leq K \|x_1 - x_2\|$$

Output  
change

Input  
change

K=1 for "1 – Lipschitz"

Do not change fast



## Improved WGAN (WGAN-GP)

$$V(G, D) = \max_{D \in 1-\text{Lipschitz}} \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$

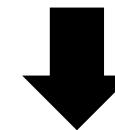
A differentiable function is 1-Lipschitz if and only if it has gradients with norm less than or equal to 1 everywhere.

$$D \in 1 - \text{Lipschitz} \iff \|\nabla_x D(x)\| \leq 1 \text{ for all } x$$

$$V(G, D) \approx \max_D \{E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)]\}$$

$$- \lambda \int_x \max(0, \|\nabla_x D(x)\| - 1) dx\}$$

Prefer  $\|\nabla_x D(x)\| \leq 1$  for all x

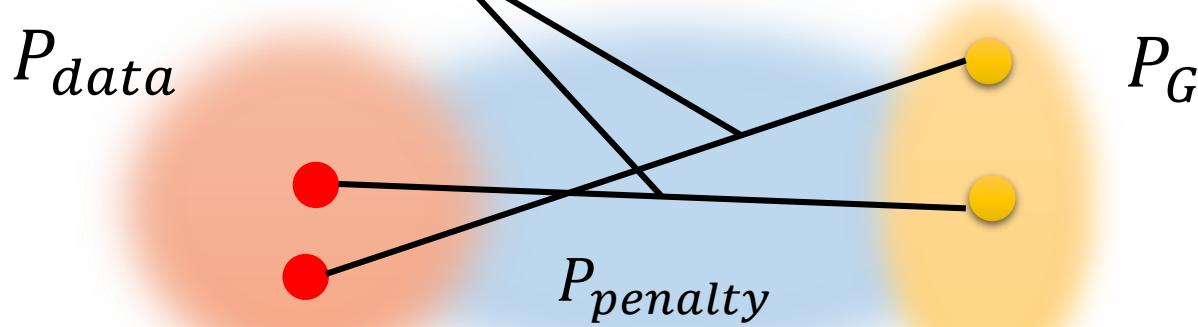


$$- \lambda E_{x \sim P_{penalty}} [\max(0, \|\nabla_x D(x)\| - 1)]\}$$

Prefer  $\|\nabla_x D(x)\| \leq 1$  for x sampling from  $x \sim P_{penalty}$

# Improved WGAN (WGAN-GP)

$$V(G, D) \approx \max_D \{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)] \\ - \lambda E_{x \sim P_{penalty}}[\max(0, \|\nabla_x D(x)\| - 1)] \}$$

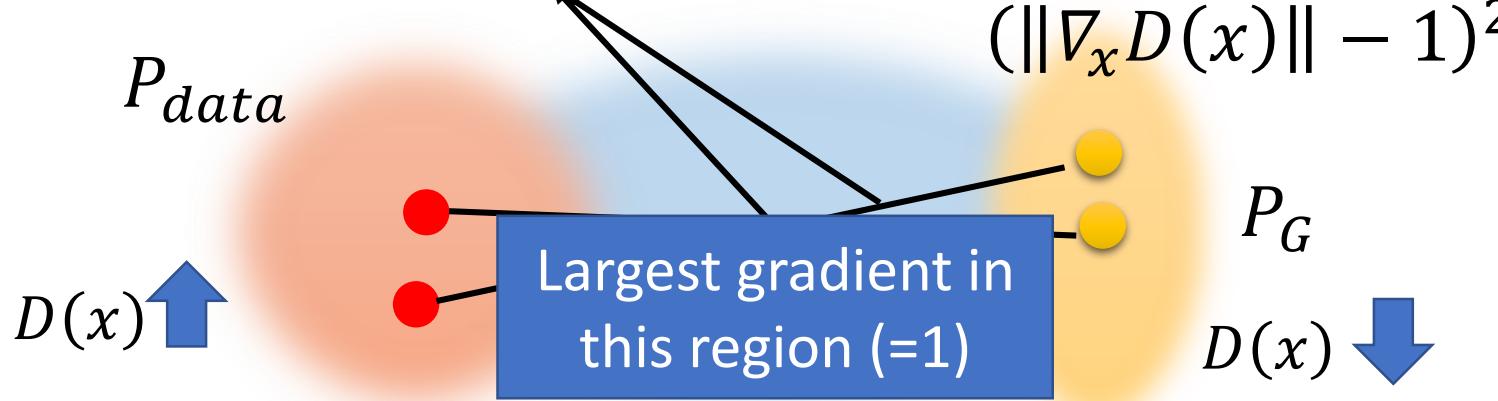


“Given that enforcing the Lipschitz constraint everywhere is intractable, enforcing it ***only along these straight lines*** seems sufficient and experimentally results in good performance.”

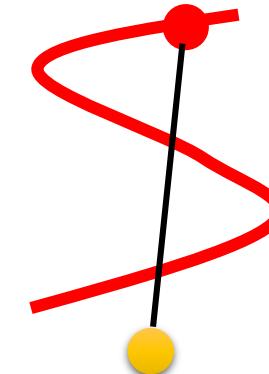
Only give gradient constraint to the region between  $P_{data}$  and  $P_G$  because they influence how  $P_G$  moves to  $P_{data}$

## Improved WGAN (WGAN-GP)

$$V(G, D) \approx \max_D \{ E_{x \sim P_{data}}[D(x)] - E_{x \sim P_G}[D(x)] \\ - \lambda E_{x \sim P_{penalty}} [\max(0, \|\nabla_x D(x)\| - 1)] \}$$



“Simply penalizing overly large gradients also works in theory, but experimentally we found that this approach converged faster and to better optima.”



# Spectrum Norm

Spectral Normalization → Keep gradient norm smaller than 1 everywhere [Miyato, et al., ICLR, 2018]



# Algorithm of

# WGAN

- In each training iteration:

No sigmoid for the output of D

- Sample m examples  $\{x^1, x^2, \dots, x^m\}$  from data distribution  $P_{data}(x)$
- Sample m noise samples  $\{z^1, z^2, \dots, z^m\}$  from the prior  $P_{prior}(z)$
- Obtaining generated data  $\{\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^m\}$ ,  $\tilde{x}^i = G(z^i)$
- Update discriminator parameters  $\theta_d$  to maximize

$$\cdot \tilde{V} = \frac{1}{m} \sum_{i=1}^m D(x^i) - \frac{1}{m} \sum_{i=1}^m D(\tilde{x}^i)$$

$$\cdot \theta_d \leftarrow \theta_d + \eta \nabla \tilde{V}(\theta_d)$$

- Sample another m noise samples  $\{z^1, z^2, \dots, z^m\}$  from the prior  $P_{prior}(z)$

Weight clipping /  
Gradient Penalty ...

- Update generator parameters  $\theta_g$  to minimize

$$\cdot \tilde{V} = \frac{1}{m} \sum_{i=1}^m \log D(x^i) - \frac{1}{m} \sum_{i=1}^m \log D(G(z^i))$$

$$\cdot \theta_g \leftarrow \theta_g - \eta \nabla \tilde{V}(\theta_g)$$

Learning  
D

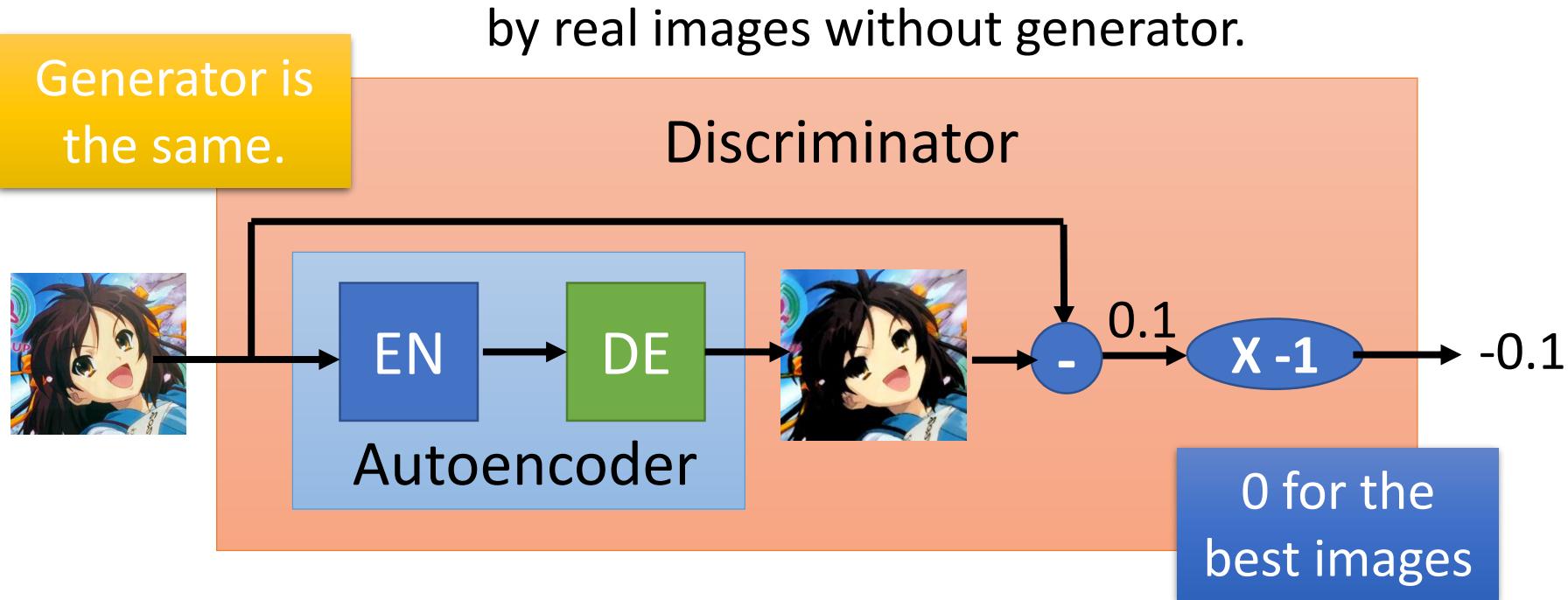
Repeat  
k times

Learning  
G

Only  
Once

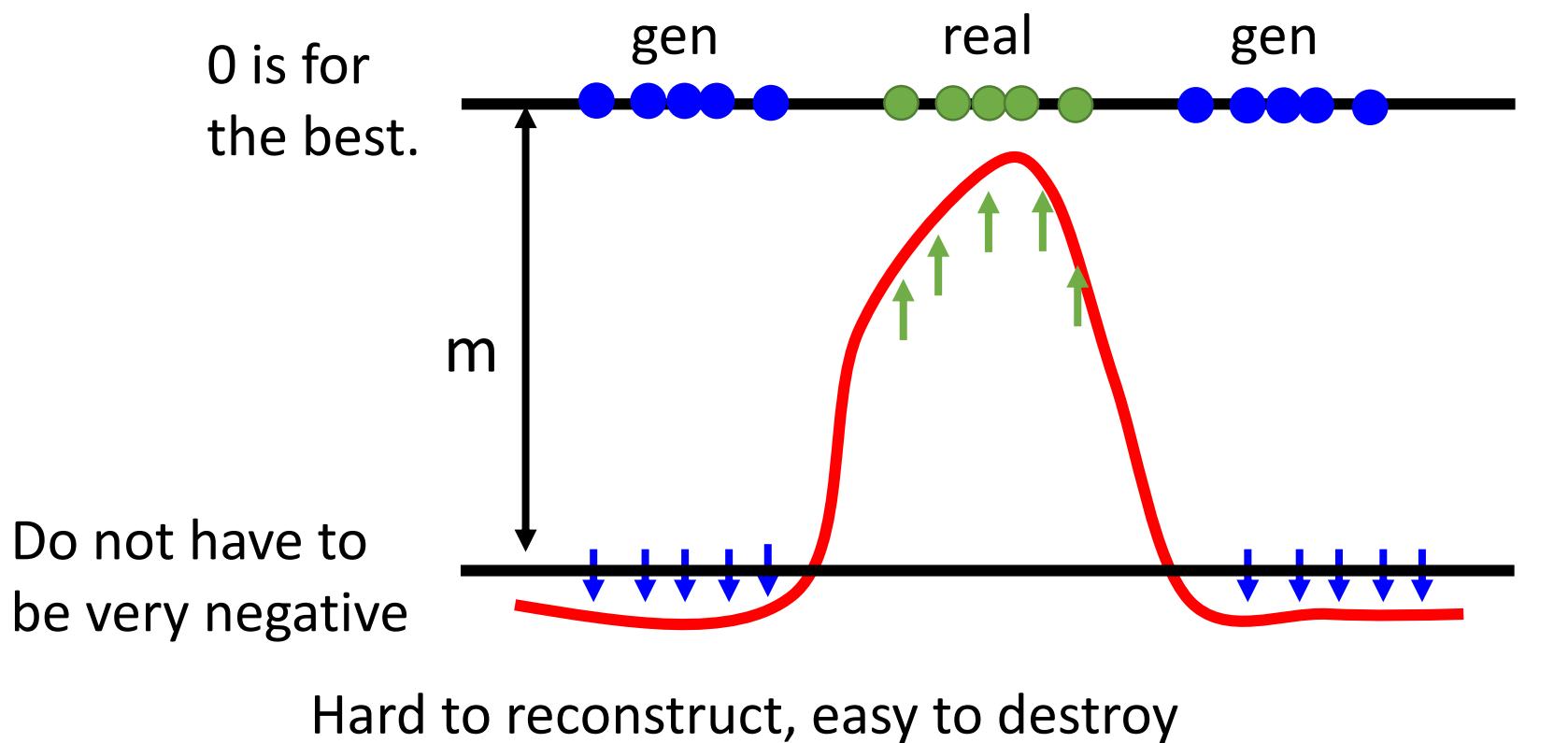
# Energy-based GAN (EBGAN)

- Using an autoencoder as discriminator D
  - Using the negative reconstruction error of auto-encoder to determine the goodness
  - **Benefit:** The auto-encoder can be pre-train by real images without generator.



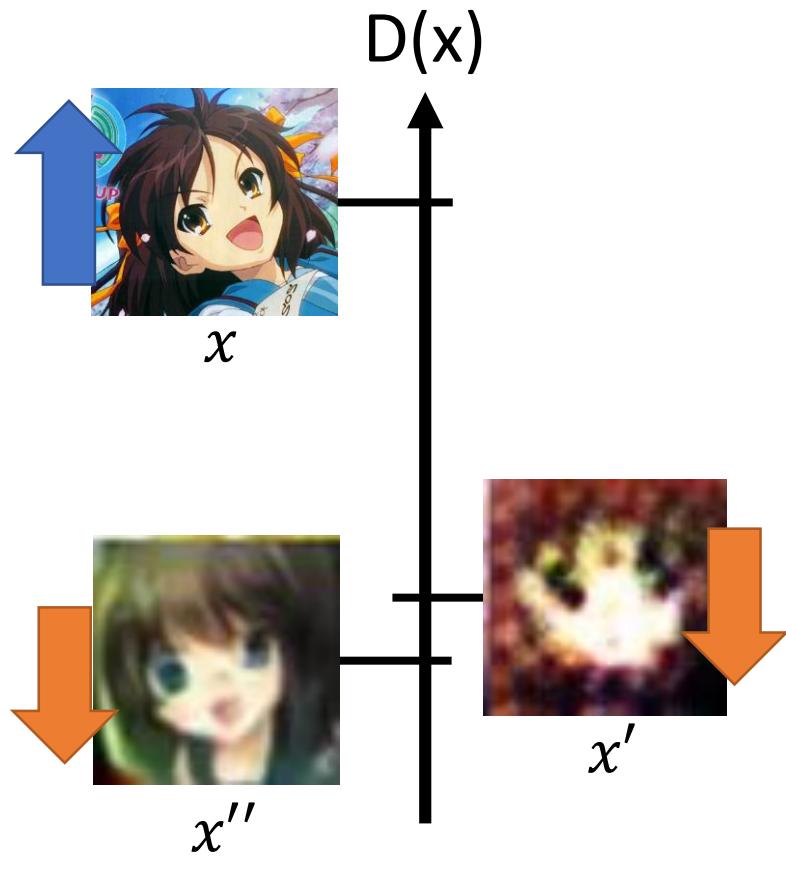
# EBGAN

Auto-encoder based discriminator  
only gives limited region large value.

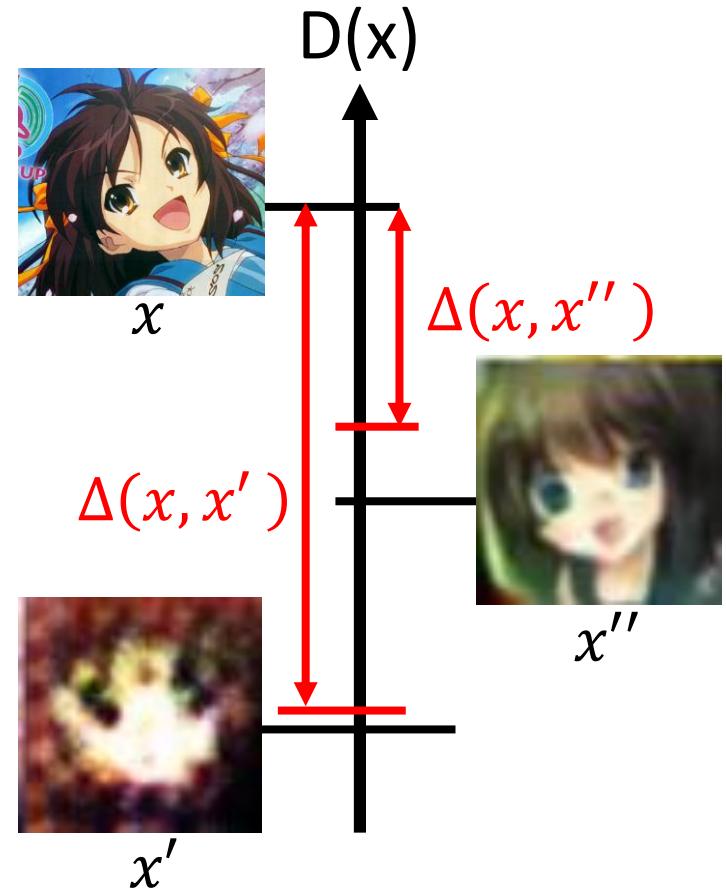


# Outlook: Loss-sensitive GAN (LSGAN)

**WGAN**



**LSGAN**



# Reference

- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, Yoshua Bengio, Generative Adversarial Networks, NIPS, 2014
- Sebastian Nowozin, Botond Cseke, Ryota Tomioka, “f-GAN: Training Generative Neural Samplers using Variational Divergence Minimization”, NIPS, 2016
- Martin Arjovsky, Soumith Chintala, Léon Bottou, Wasserstein GAN, arXiv, 2017
- Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, Aaron Courville, Improved Training of Wasserstein GANs, NIPS, 2017
- Junbo Zhao, Michael Mathieu, Yann LeCun, Energy-based Generative Adversarial Network, arXiv, 2016
- Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, Olivier Bousquet, “Are GANs Created Equal? A Large-Scale Study”, arXiv, 2017
- Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, Xi Chen Improved Techniques for Training GANs, NIPS, 2016

# Reference

- Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Sepp Hochreiter, GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium, NIPS, 2017
- Naveen Kodali, Jacob Abernethy, James Hays, Zsolt Kira, “On Convergence and Stability of GANs”, arXiv, 2017
- Xiang Wei, Boqing Gong, Zixia Liu, Wei Lu, Liqiang Wang, Improving the Improved Training of Wasserstein GANs: A Consistency Term and Its Dual Effect, ICLR, 2018
- Takeru Miyato, Toshiki Kataoka, Masanori Koyama, Yuichi Yoshida, Spectral Normalization for Generative Adversarial Networks, ICLR, 2018

# Feature Extraction

# InfoGAN

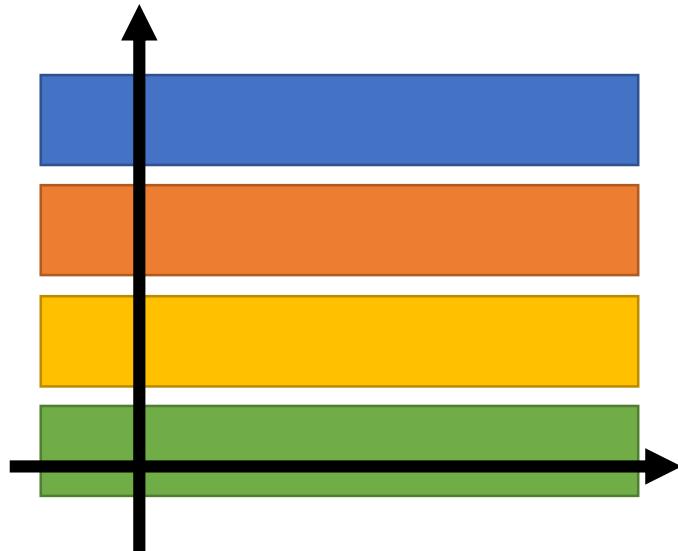
(The colors represents the characteristics.)

Regular  
GAN

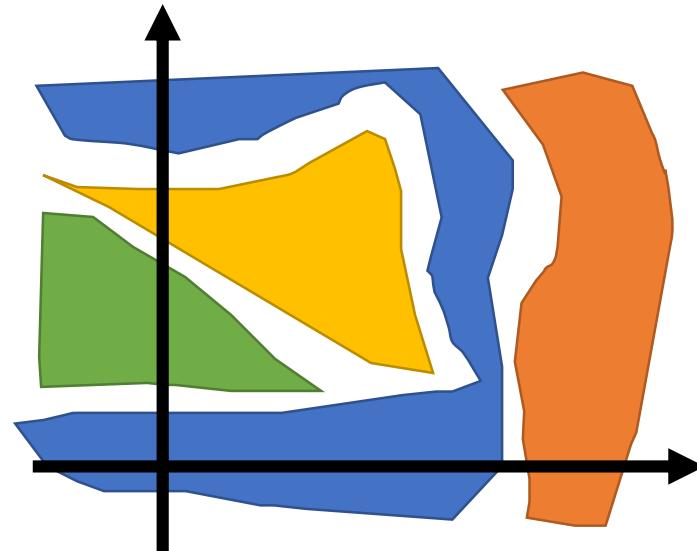


Modifying a specific dimension,  
no clear meaning

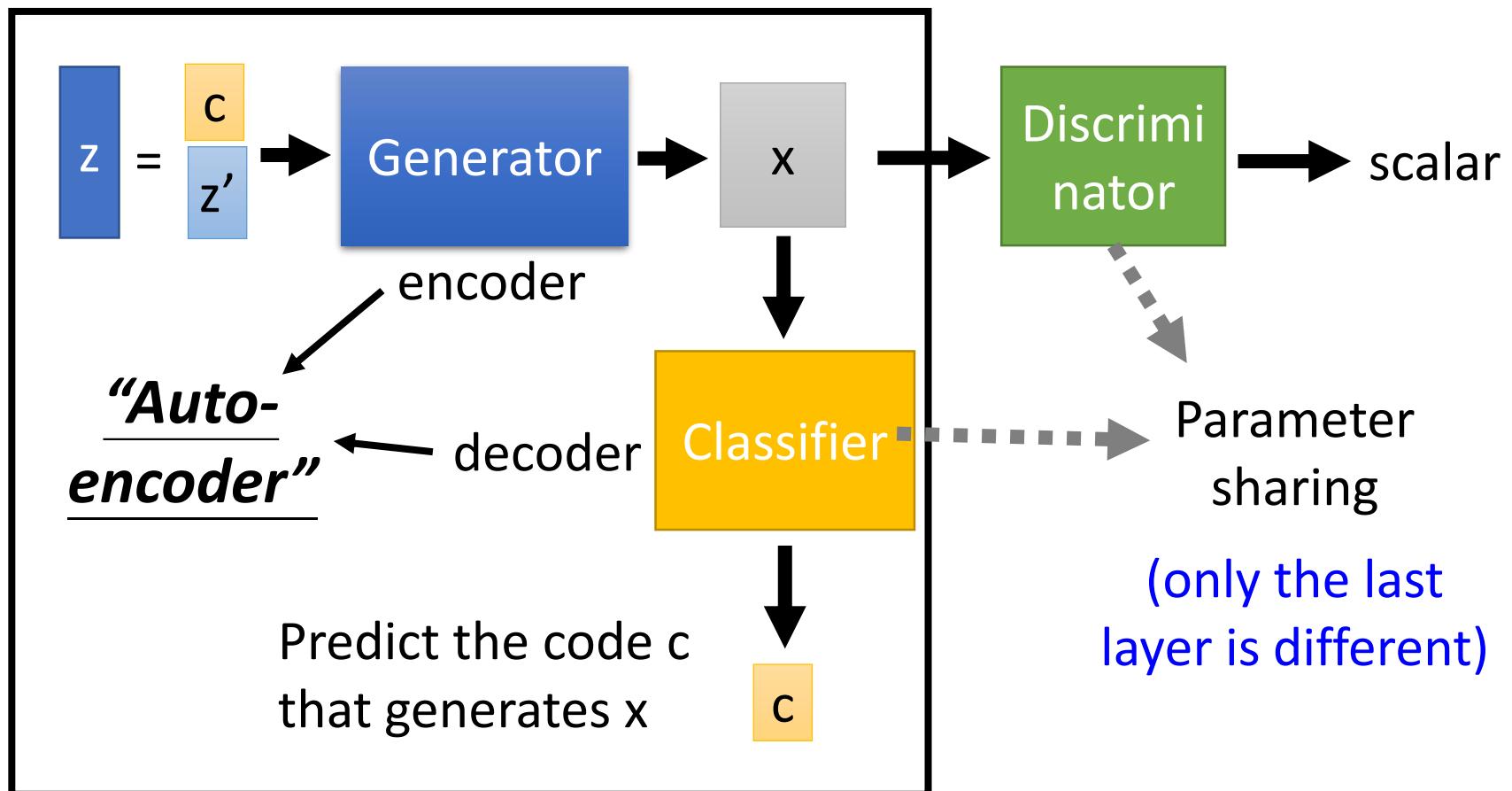
What we expect



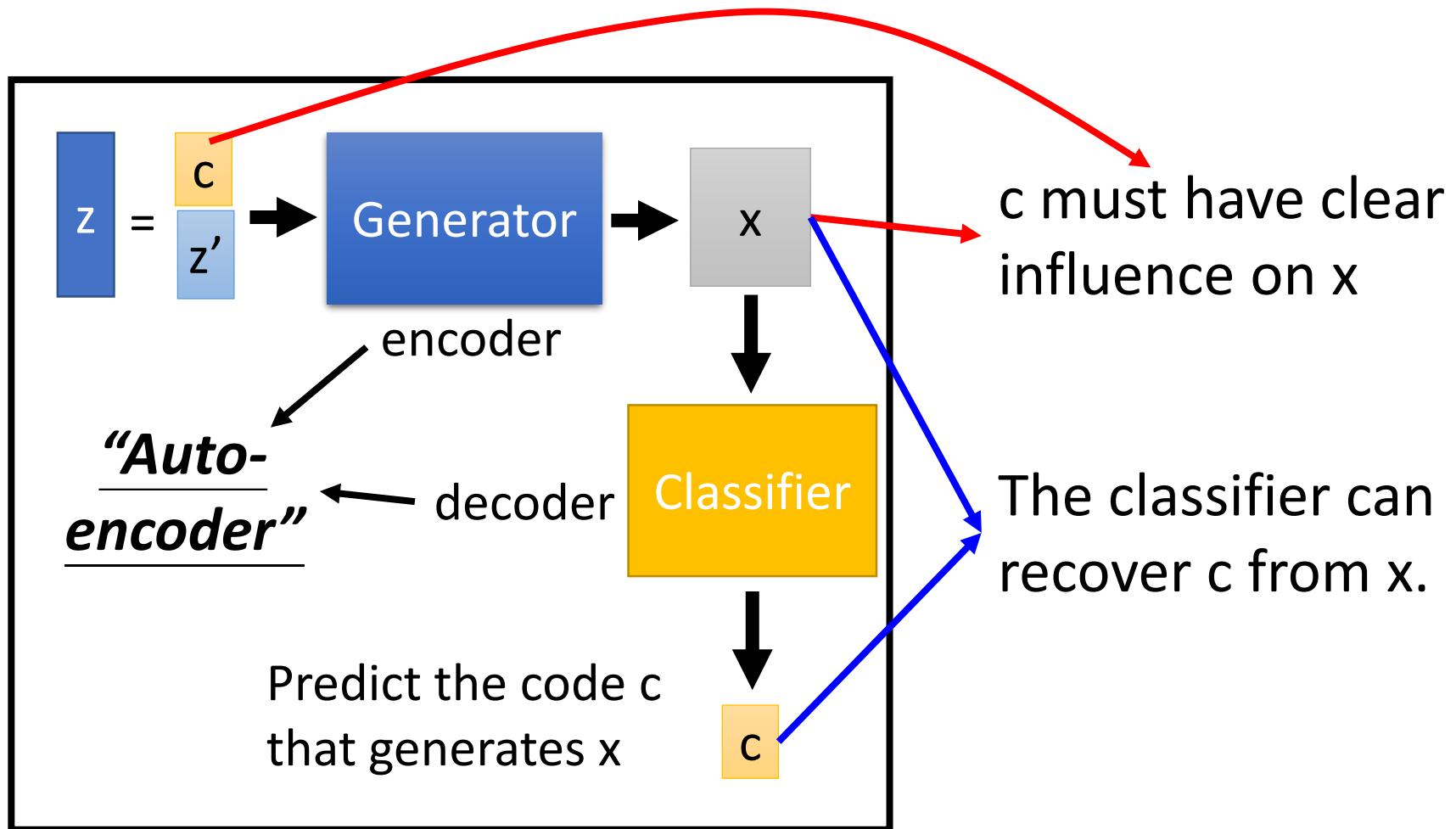
Actually ...

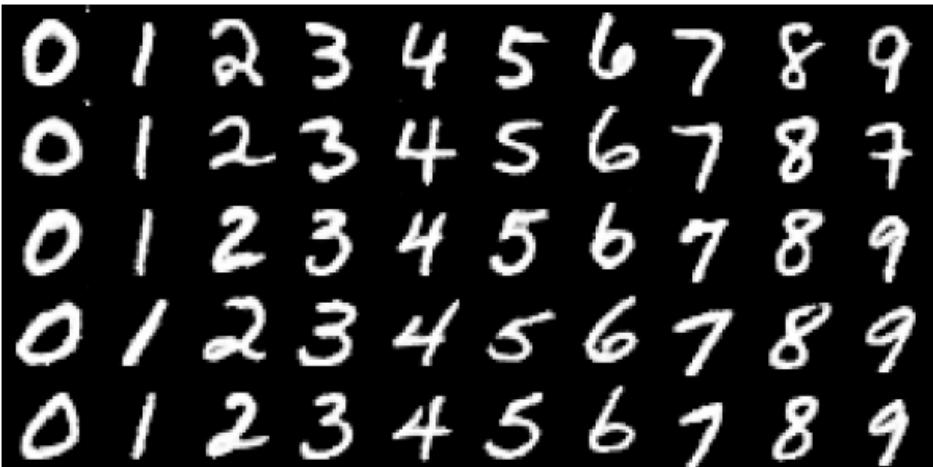


# What is InfoGAN?

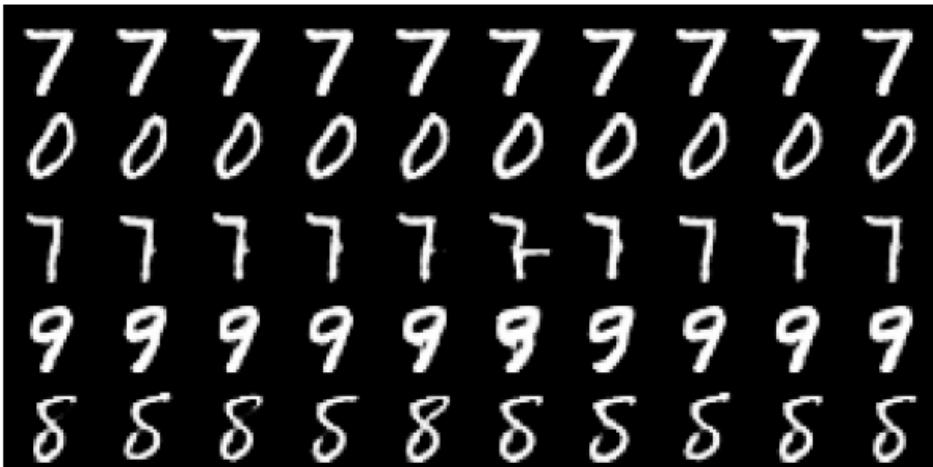


# What is InfoGAN?

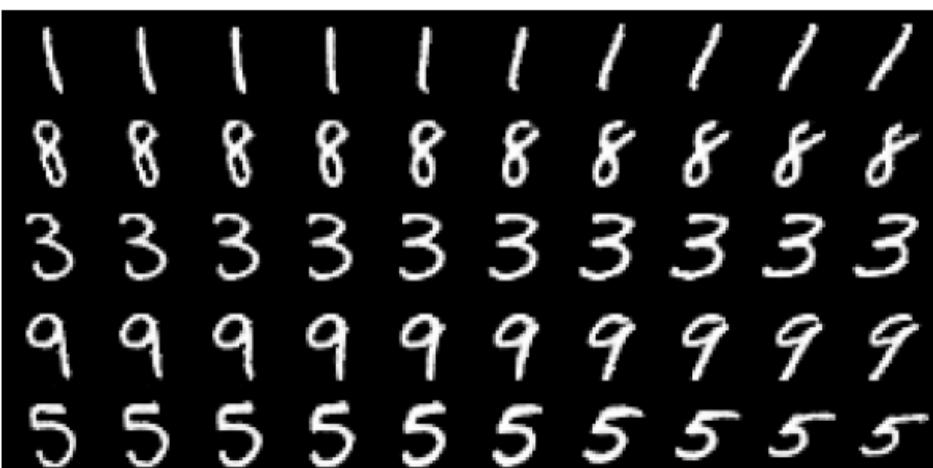




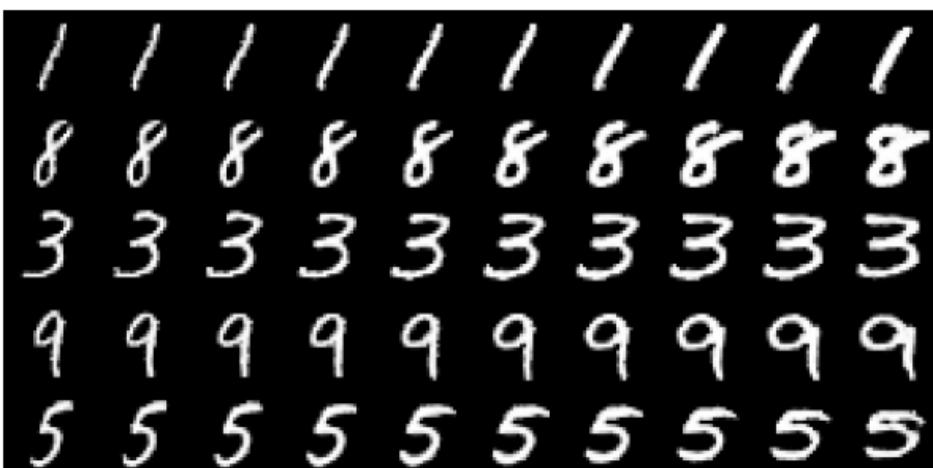
(a) Varying  $c_1$  on InfoGAN (Digit type)



(b) Varying  $c_1$  on regular GAN (No clear meaning)



(c) Varying  $c_2$  from  $-2$  to  $2$  on InfoGAN (Rotation)

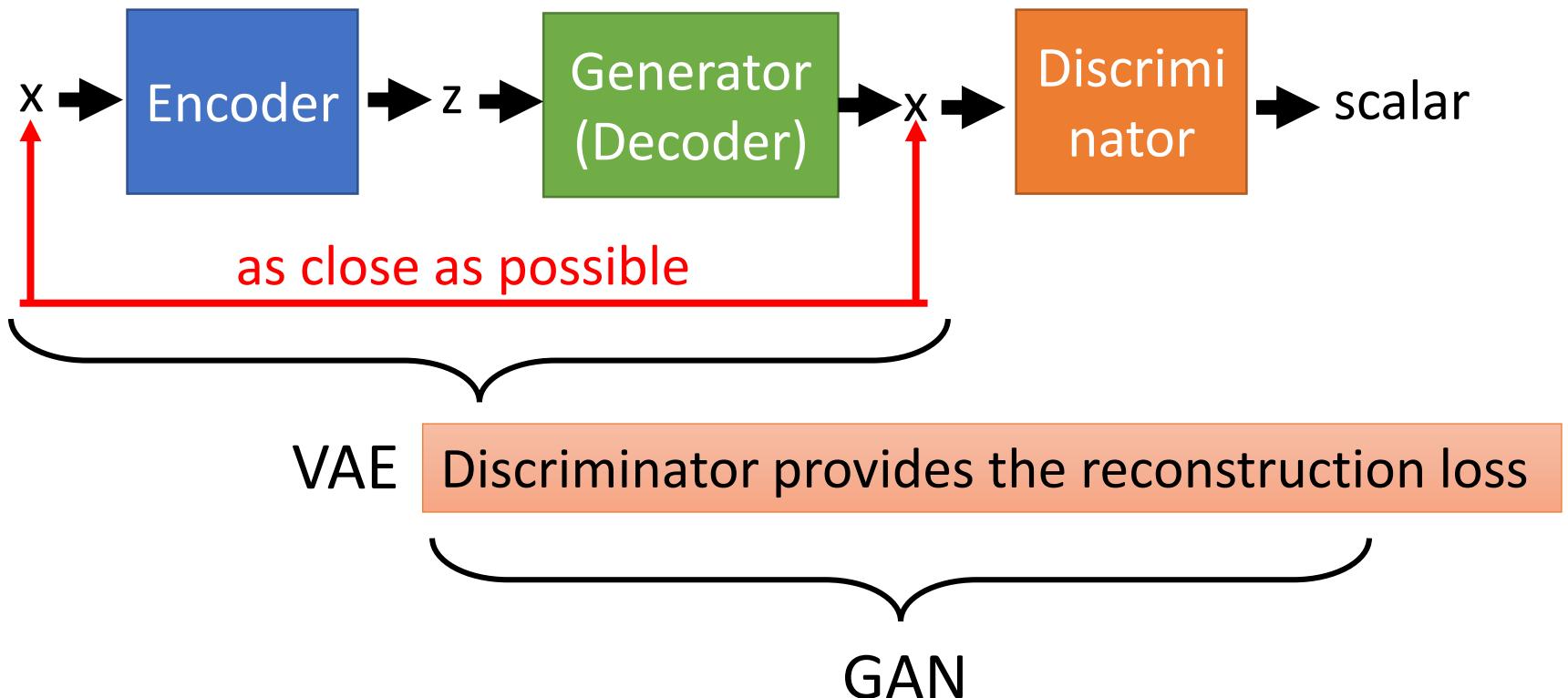


(d) Varying  $c_3$  from  $-2$  to  $2$  on InfoGAN (Width)

# VAE-GAN

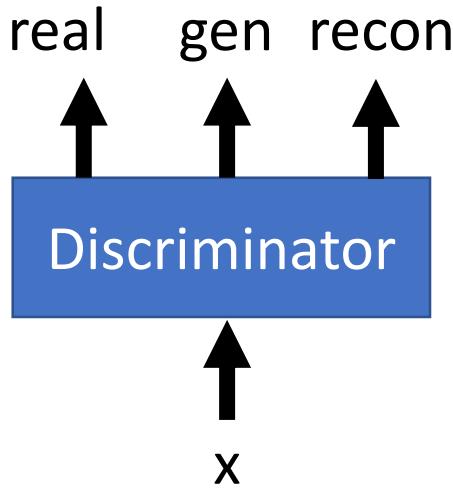
Anders Boesen, Lindbo Larsen, Søren Kaae  
Sønderby, Hugo Larochelle, Ole Winther, "Autoencoding  
beyond pixels using a learned similarity metric", ICML. 2016

- Minimize reconstruction error
- $z$  close to normal
- Minimize reconstruction error
- Cheat discriminator
- Discriminate real, generated and reconstructed images



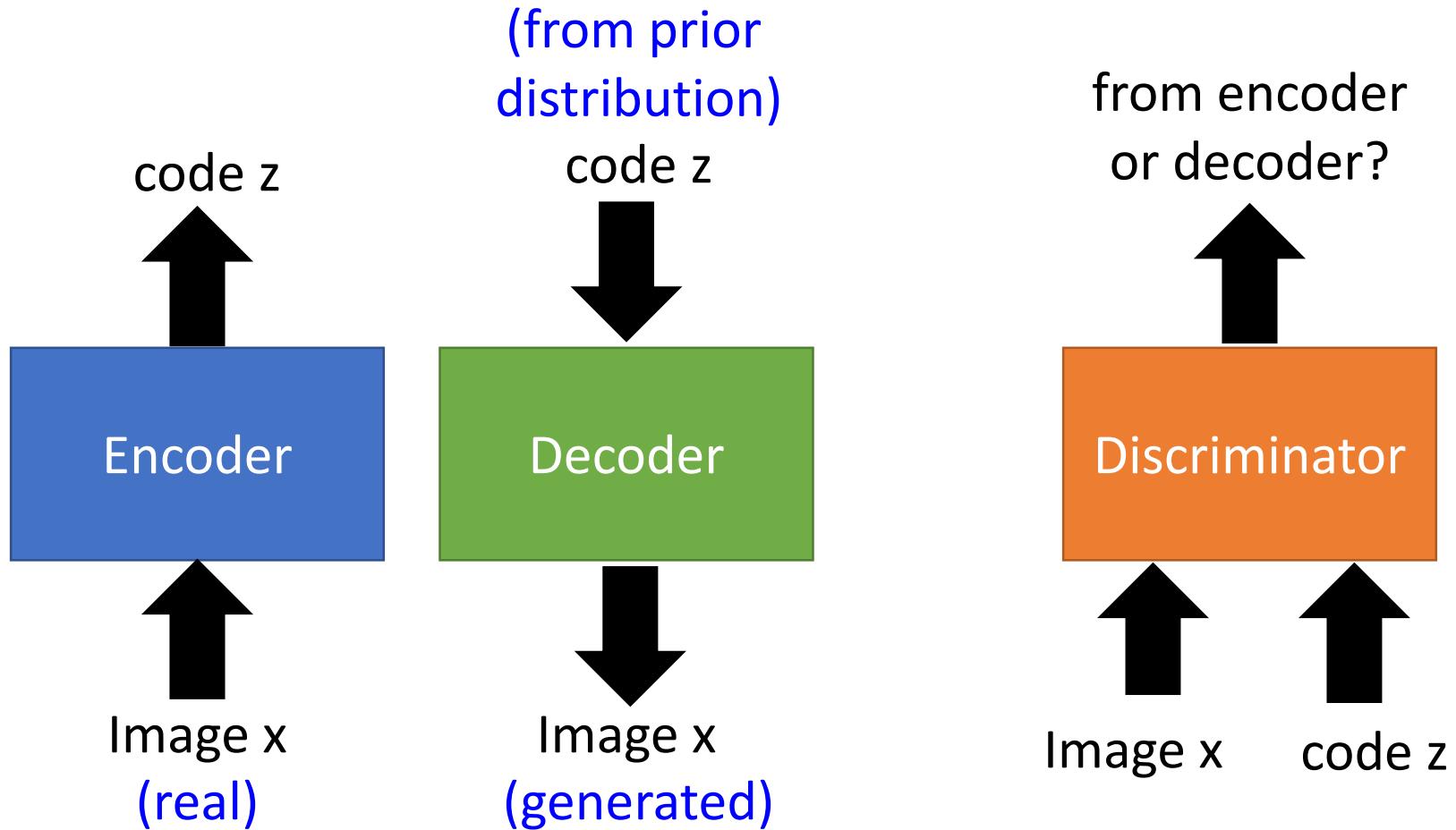
# Algorithm

Another kind of discriminator:



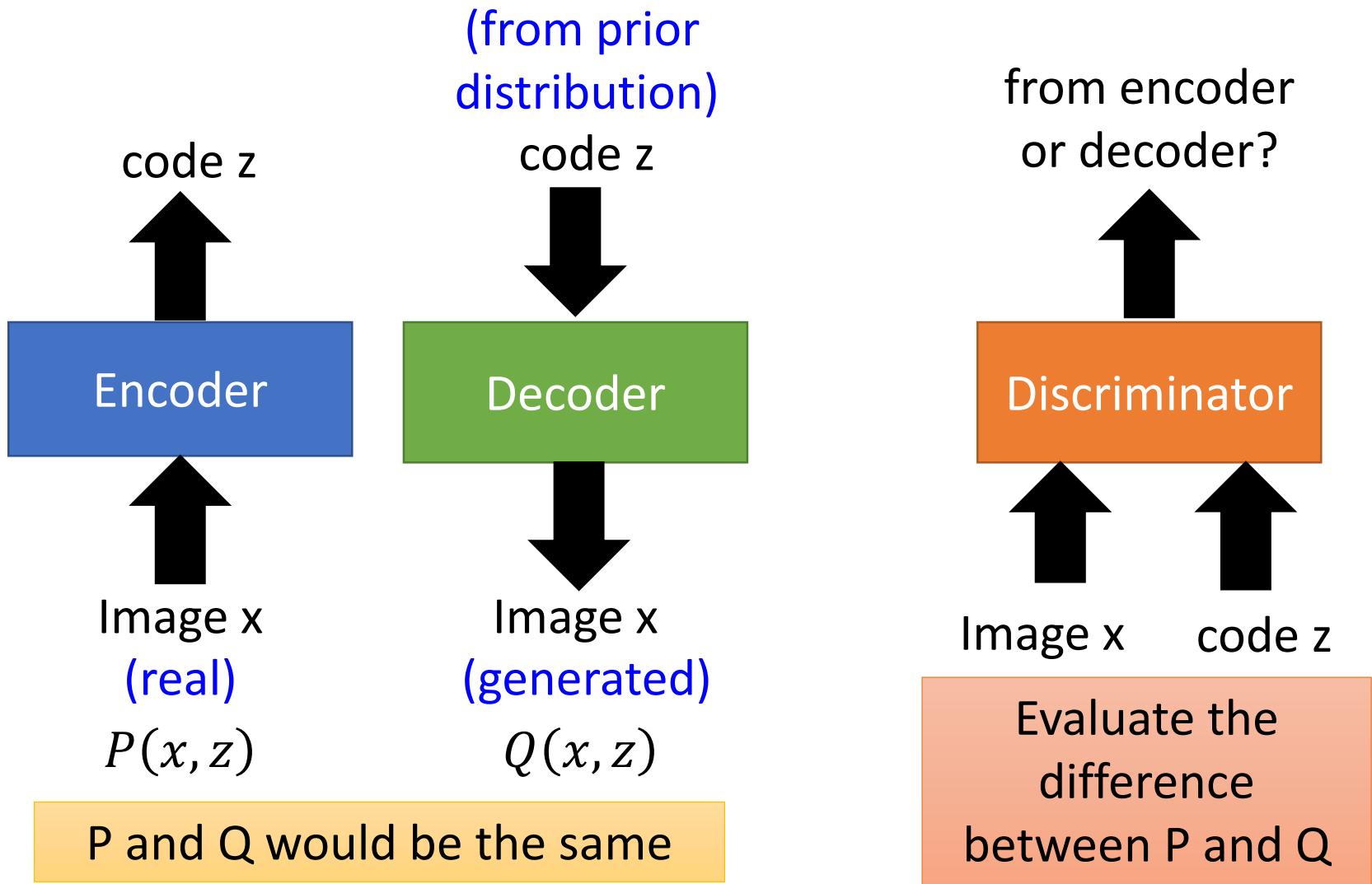
- Initialize En, De, Dis
- In each iteration:
  - Sample M images  $x^1, x^2, \dots, x^M$  from database
  - Generate M codes  $\tilde{z}^1, \tilde{z}^2, \dots, \tilde{z}^M$  from encoder
    - $\tilde{z}^i = En(x^i)$
  - Generate M images  $\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^M$  from decoder
    - $\tilde{x}^i = De(\tilde{z}^i)$
  - Sample M codes  $z^1, z^2, \dots, z^M$  from prior P(z)
  - Generate M images  $\hat{x}^1, \hat{x}^2, \dots, \hat{x}^M$  from decoder
    - $\hat{x}^i = De(z^i)$
  - Update En to decrease  $\|\tilde{x}^i - x^i\|$ , decrease  $KL(P(\tilde{z}^i | x^i) || P(z))$
  - Update De to decrease  $\|\tilde{x}^i - x^i\|$ , increase  $Dis(\tilde{x}^i)$  and  $Dis(\hat{x}^i)$
  - Update Dis to increase  $Dis(x^i)$ , decrease  $Dis(\tilde{x}^i)$  and  $Dis(\hat{x}^i)$

# BiGAN



# Algorithm

- Initialize encoder En, decoder De, discriminator Dis
- In each iteration:
  - Sample M images  $x^1, x^2, \dots, x^M$  from database
  - Generate M codes  $\tilde{z}^1, \tilde{z}^2, \dots, \tilde{z}^M$  from encoder
    - $\tilde{z}^i = En(x^i)$
  - Sample M codes  $z^1, z^2, \dots, z^M$  from prior  $P(z)$
  - Generate M codes  $\tilde{x}^1, \tilde{x}^2, \dots, \tilde{x}^M$  from decoder
    - $\tilde{x}^i = De(z^i)$
  - Update Dis to increase  $Dis(x^i, \tilde{z}^i)$ , decrease  $Dis(\tilde{x}^i, z^i)$
  - Update En and De to decrease  $Dis(x^i, \tilde{z}^i)$ , increase  $Dis(\tilde{x}^i, z^i)$



Optimal encoder  
and decoder:

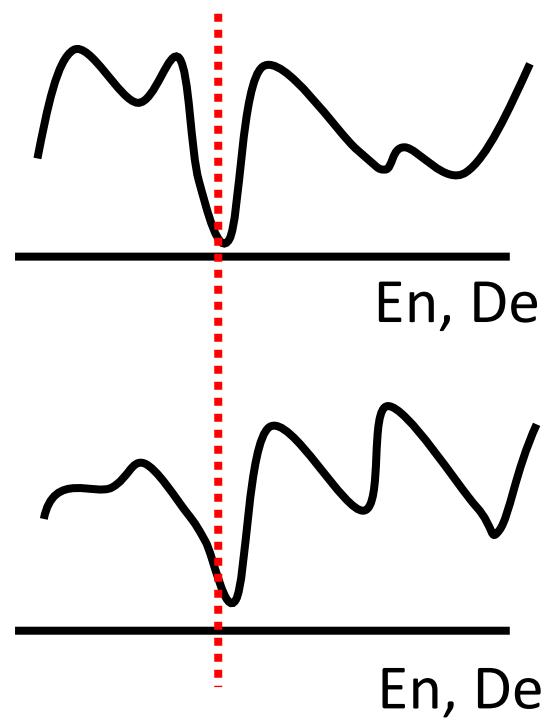
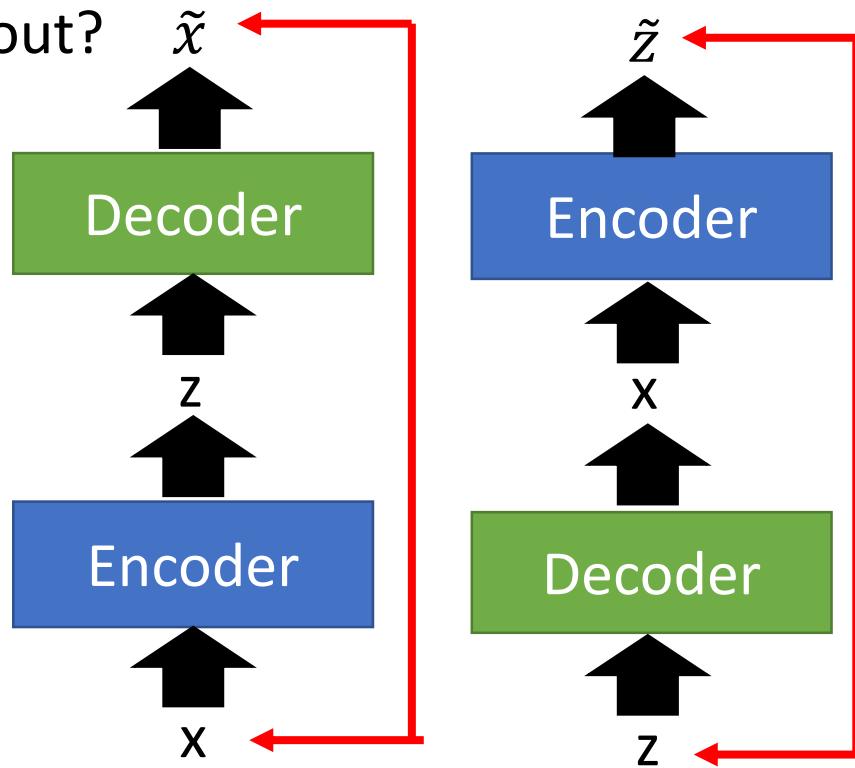
$$\begin{aligned}
 \text{En}(x') &= z' & \xrightarrow{\quad} \text{De}(z') &= x' & \text{For all } x' \\
 \text{De}(z'') &= x'' & \xrightarrow{\quad} \text{En}(x'') &= z'' & \text{For all } z''
 \end{aligned}$$

# BiGAN

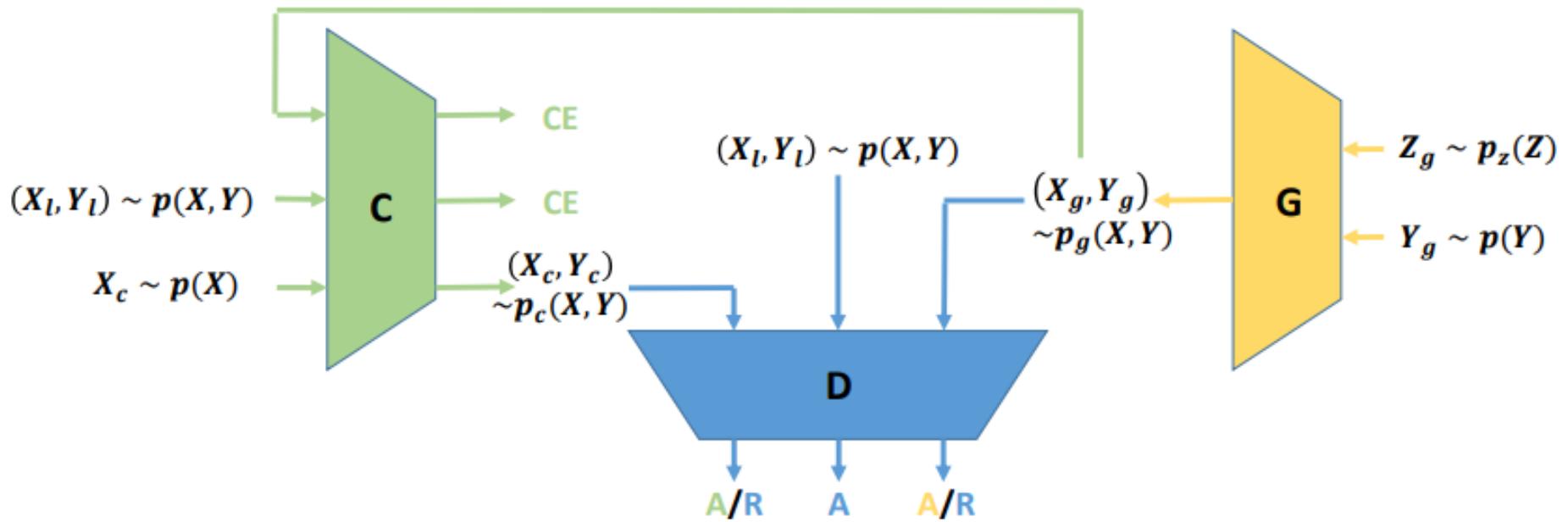
Optimal encoder  
and decoder:

$$\text{En}(x') = z' \rightarrow \text{De}(z') = x' \quad \text{For all } x'$$
$$\text{De}(z'') = x'' \rightarrow \text{En}(x'') = z'' \quad \text{For all } z''$$

How about?



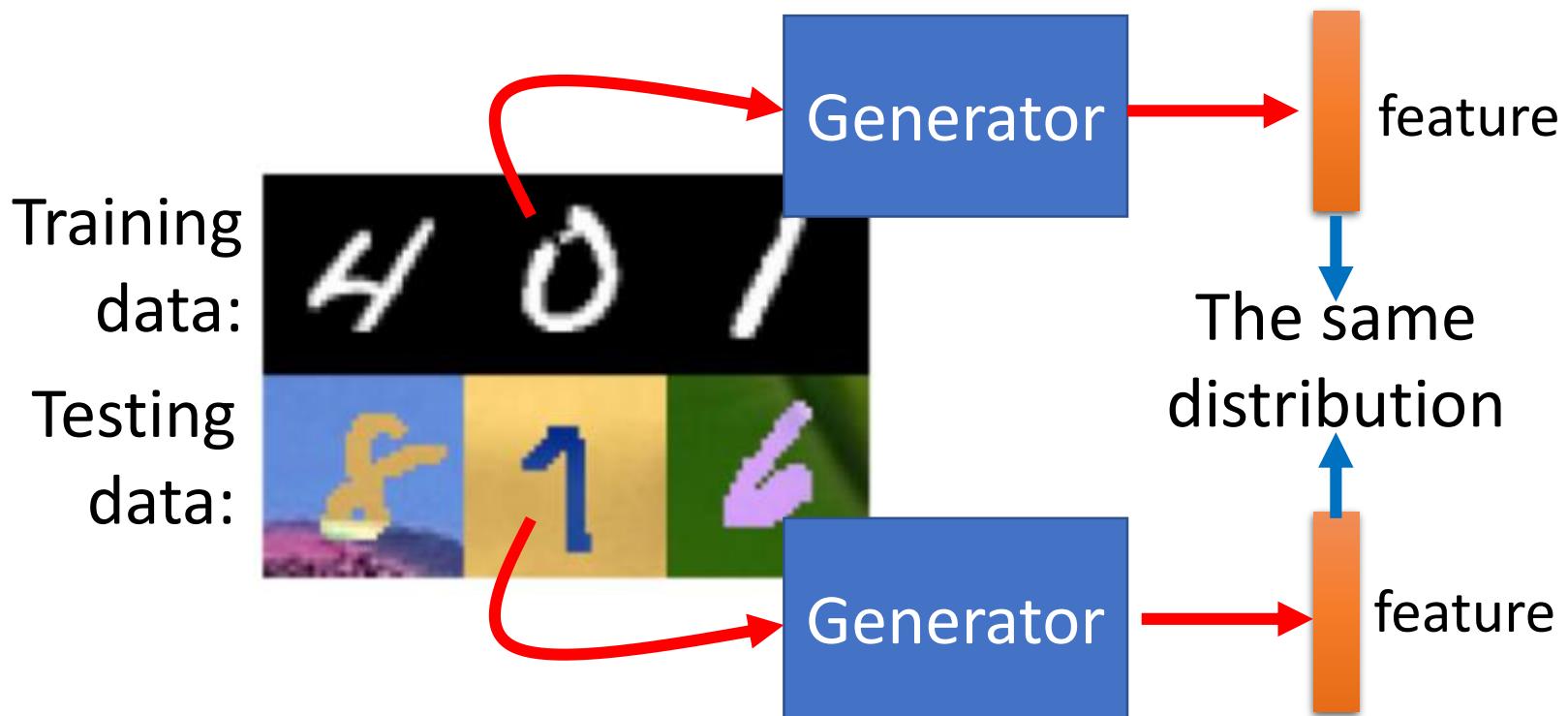
# Triple GAN



Chongxuan Li, Kun Xu, Jun Zhu, Bo Zhang, "Triple Generative Adversarial Nets", arXiv 2017

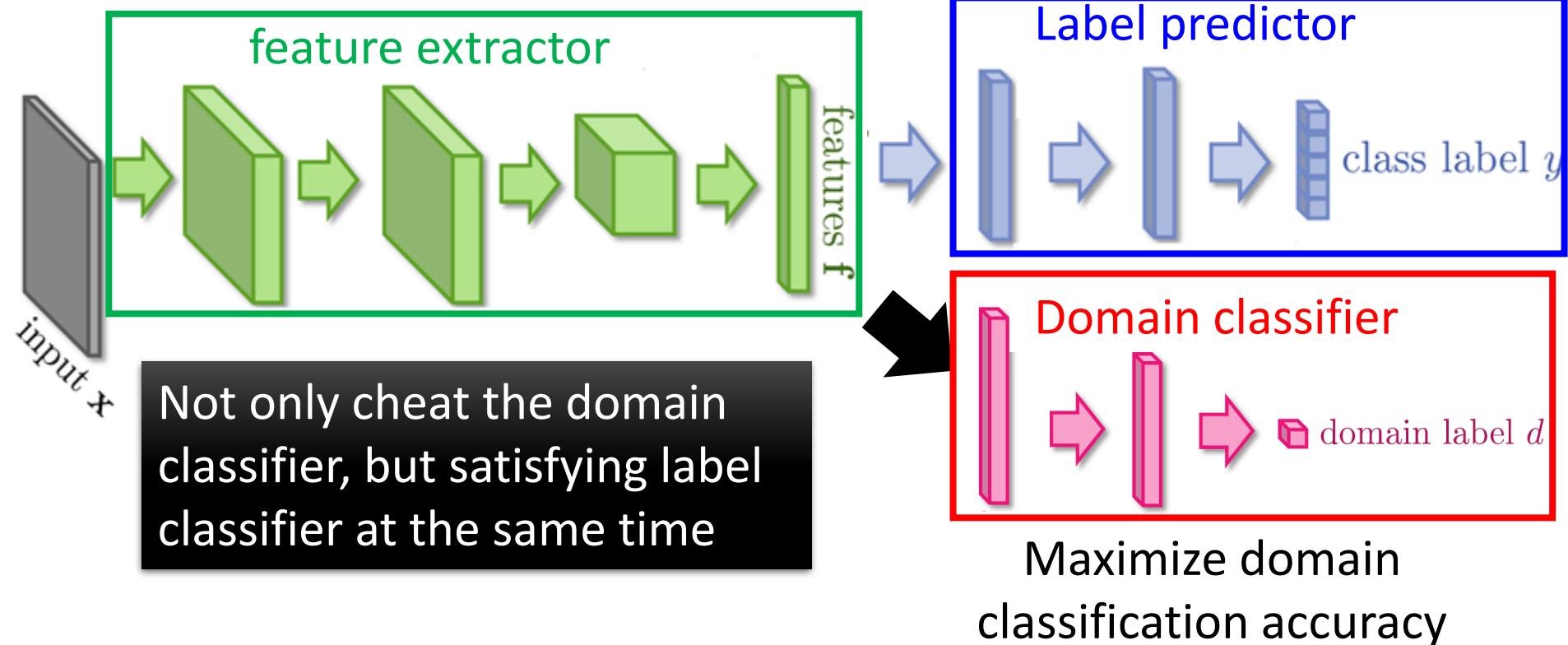
# Domain-adversarial training

- Training and testing data are in different domains



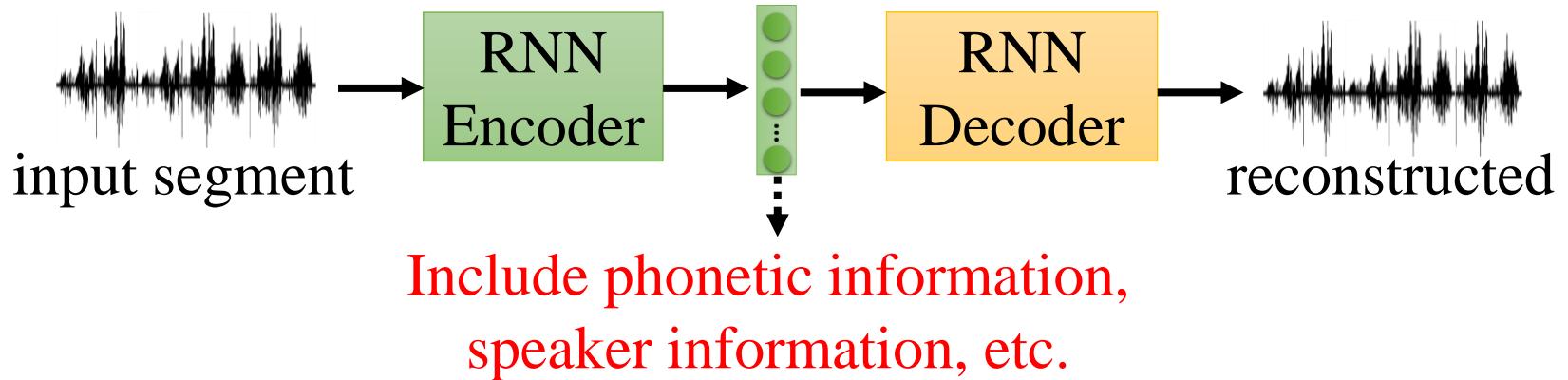
# Domain-adversarial training

Maximize label classification accuracy +  
minimize domain classification accuracy

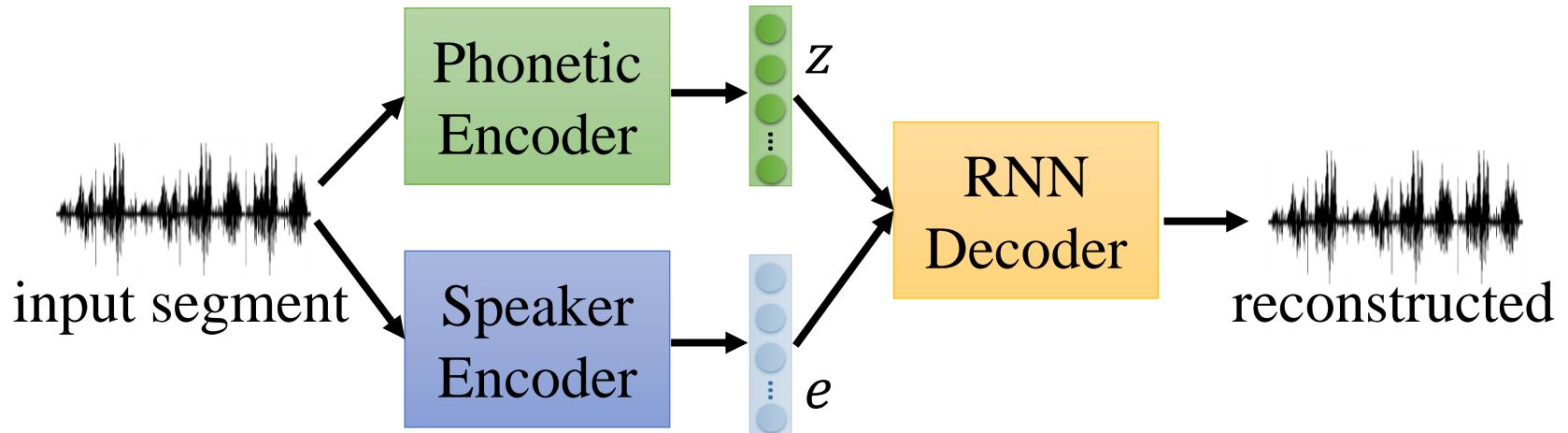


This is a big network, but different parts have different goals.

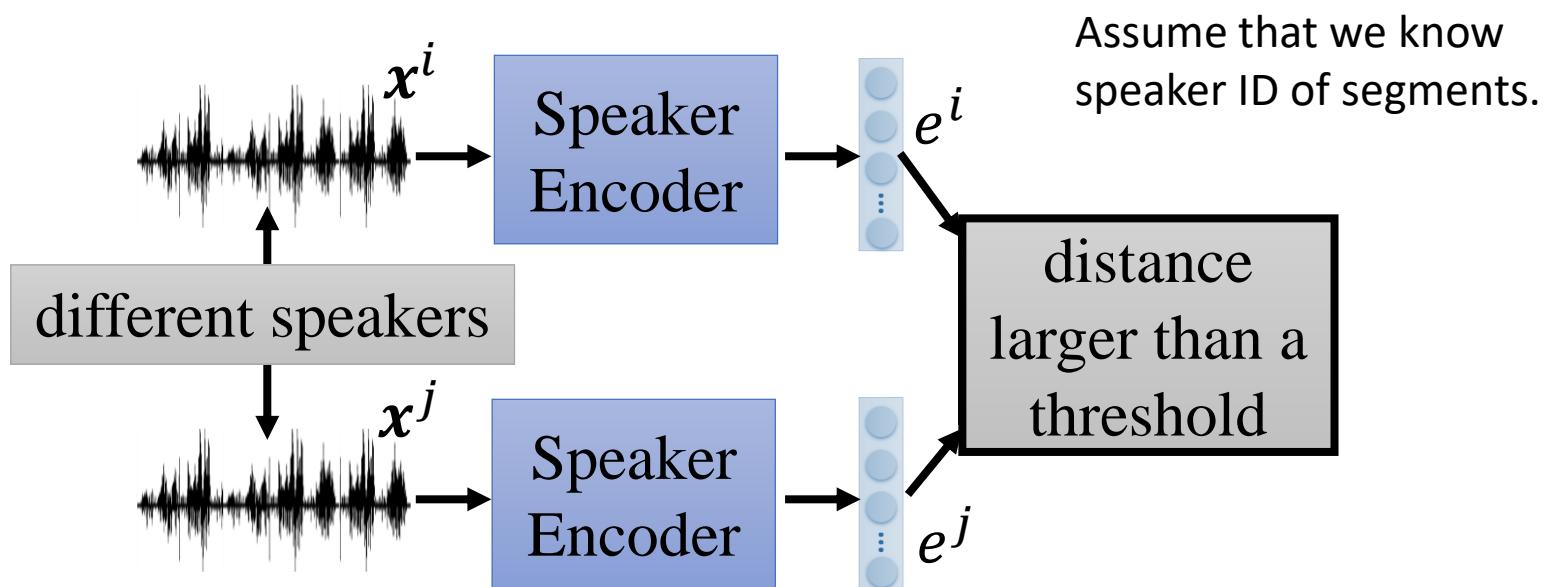
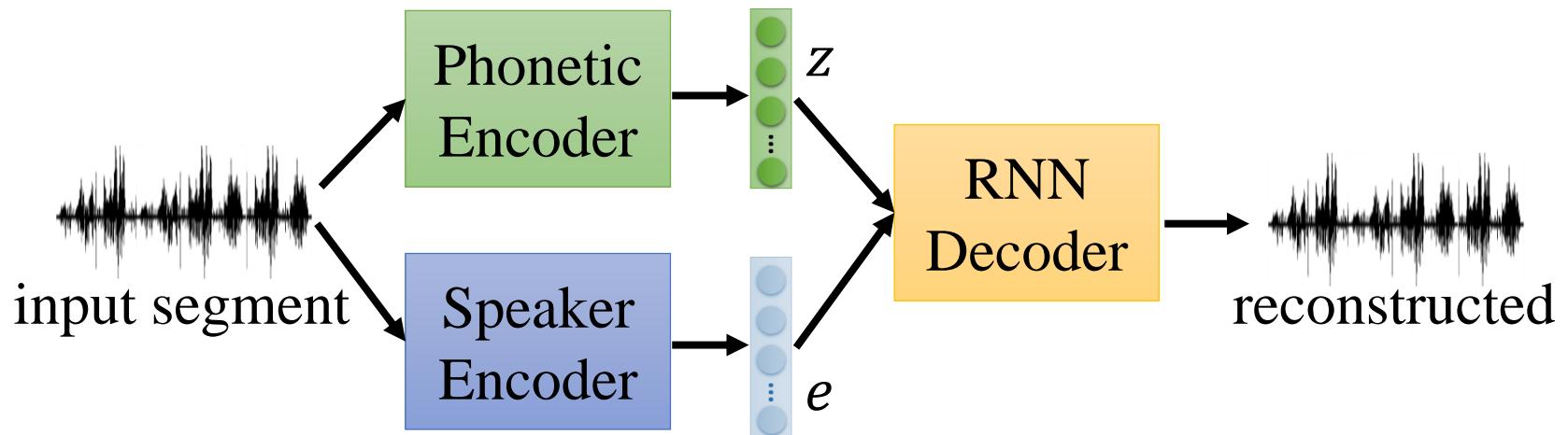
## Original Seq2seq Auto-encoder



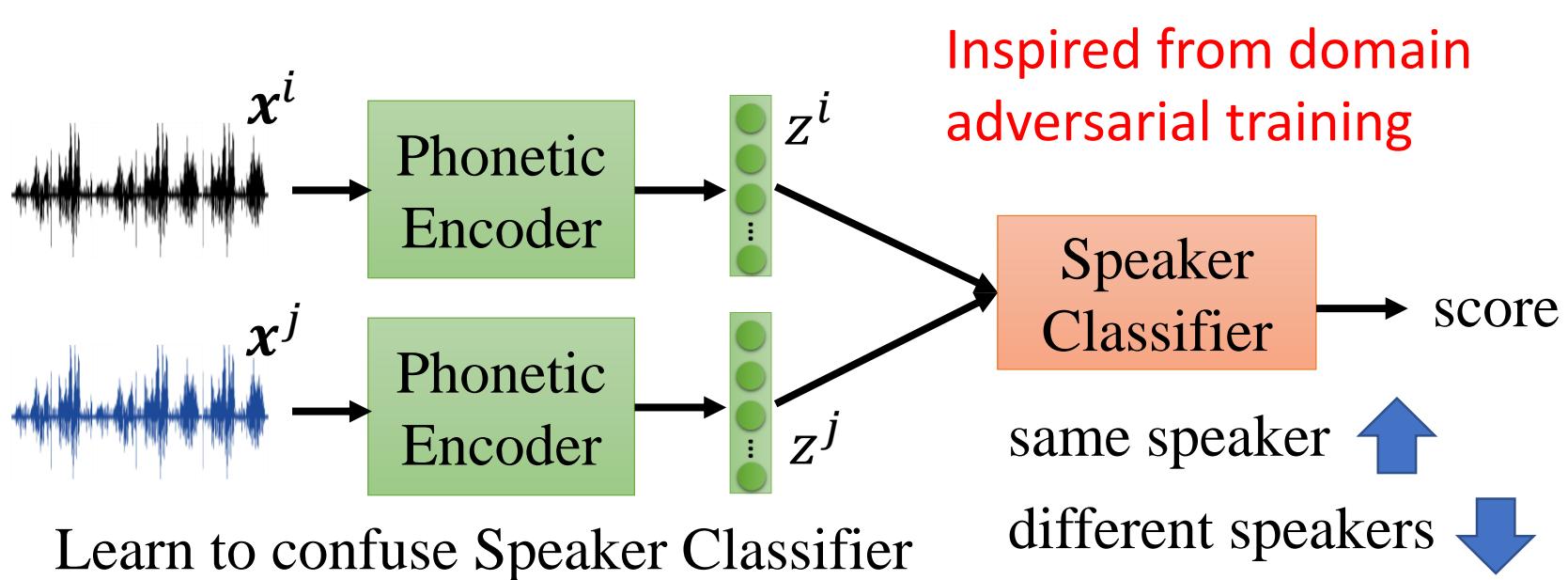
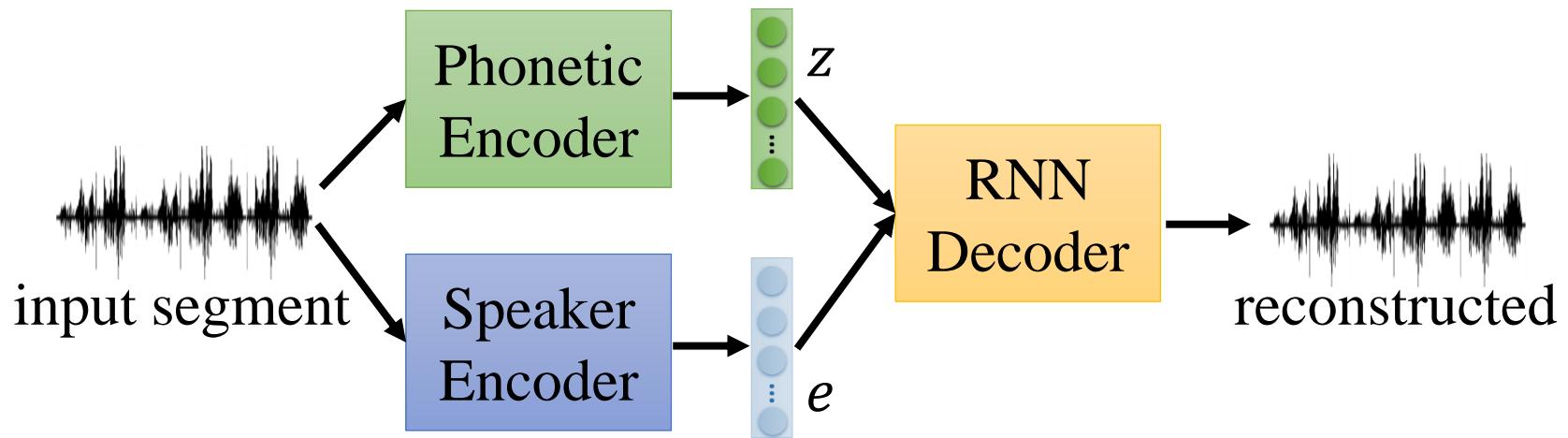
## Feature Disentangle

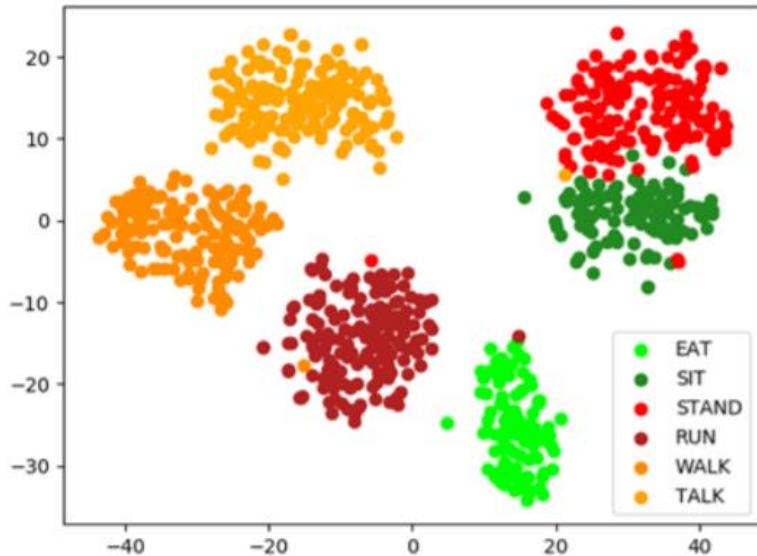


# Feature Disentangle

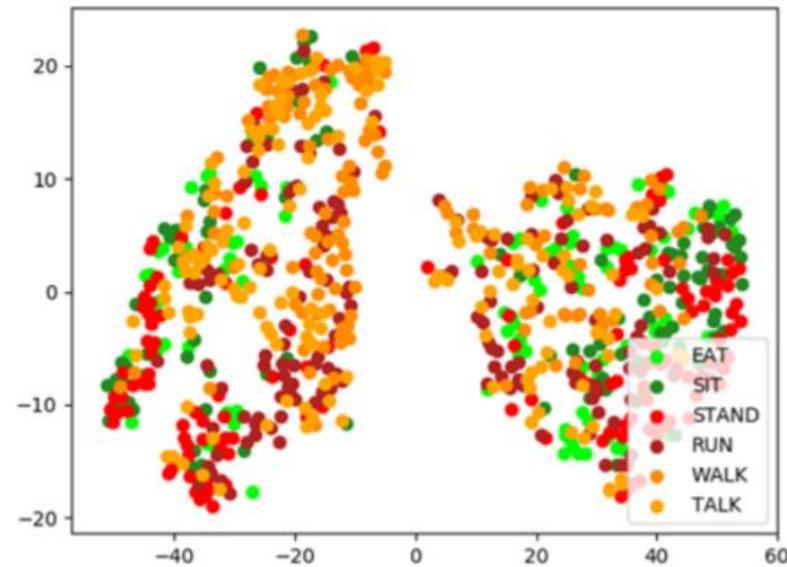


# Feature Disentangle

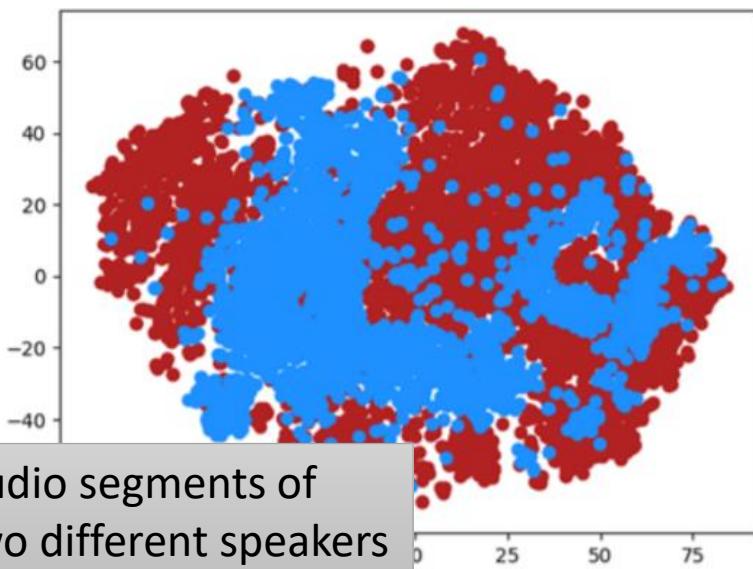




phonetic embedding  $z$

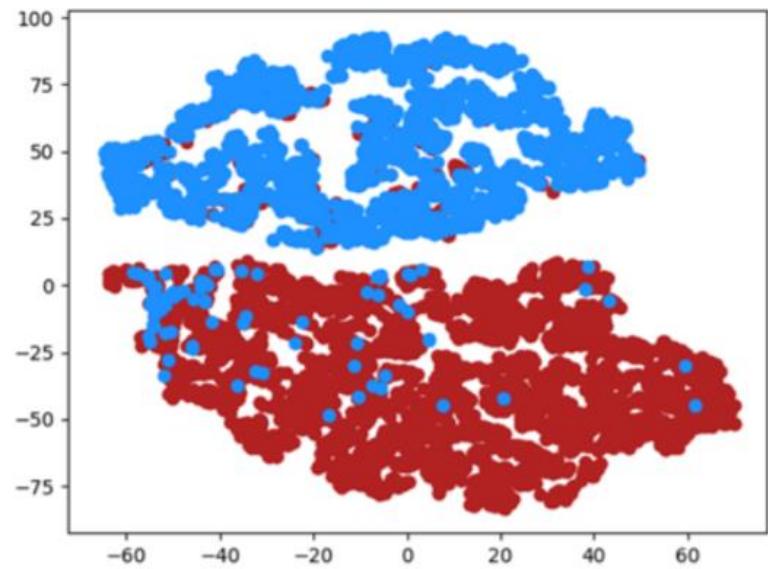


speaker embedding  $e$



Audio segments of  
two different speakers

phonetic embedding  $z$



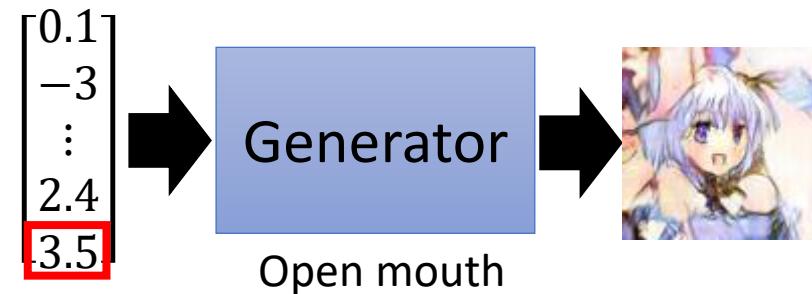
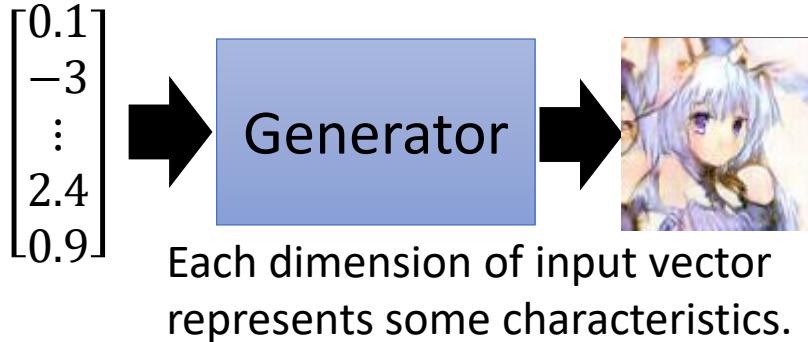
speaker embedding  $e$

# Acknowledgement

- 感謝 許傑盛 同學發現投影片上的打字錯誤

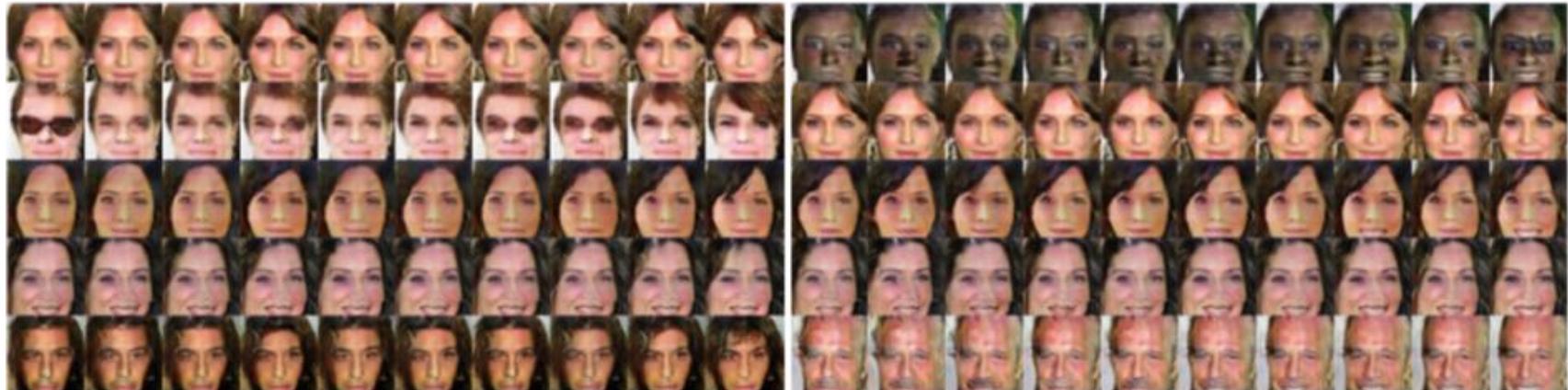
# Intelligent Photo Editing

# Modifying Input Code



- The input code determines the generator output.
- Understand the meaning of each dimension to control the output.

# Connecting Code and Attribute



(c) Hair style

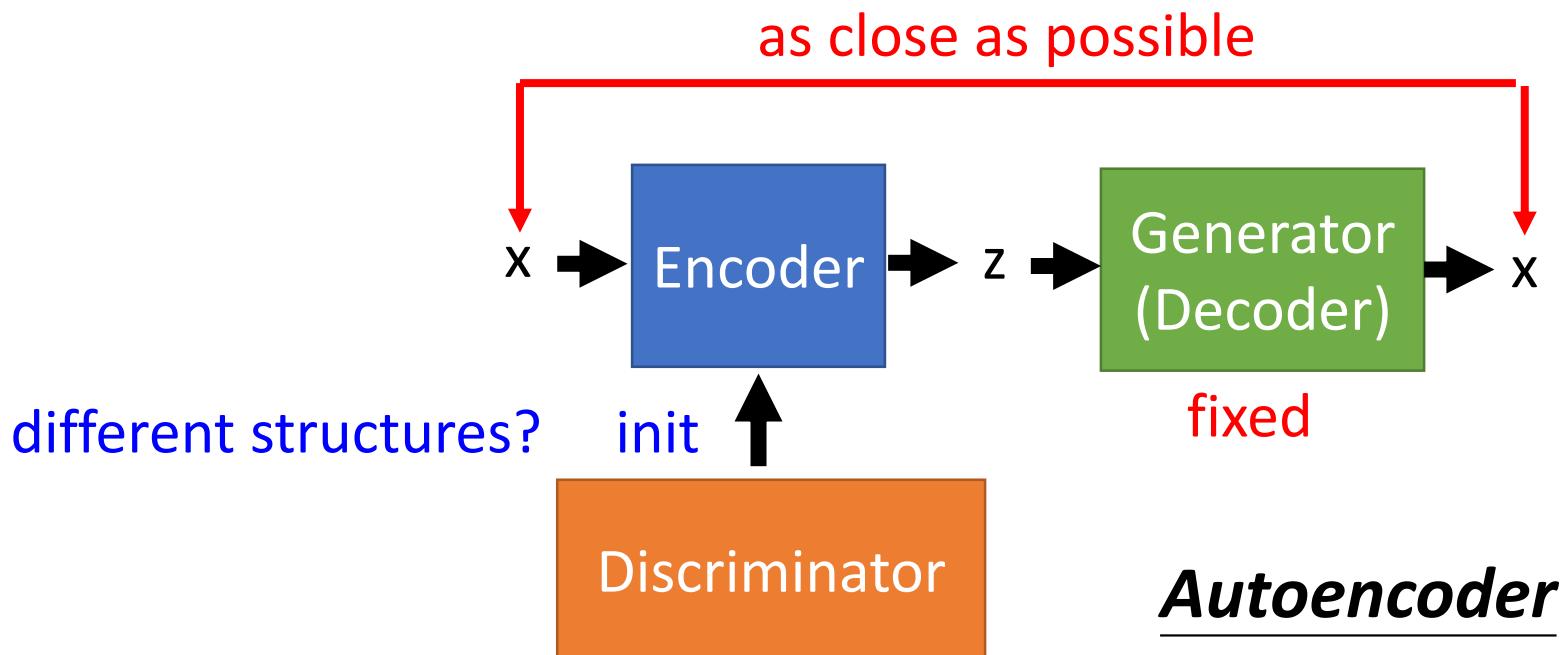
(d) Emotion

CelebA

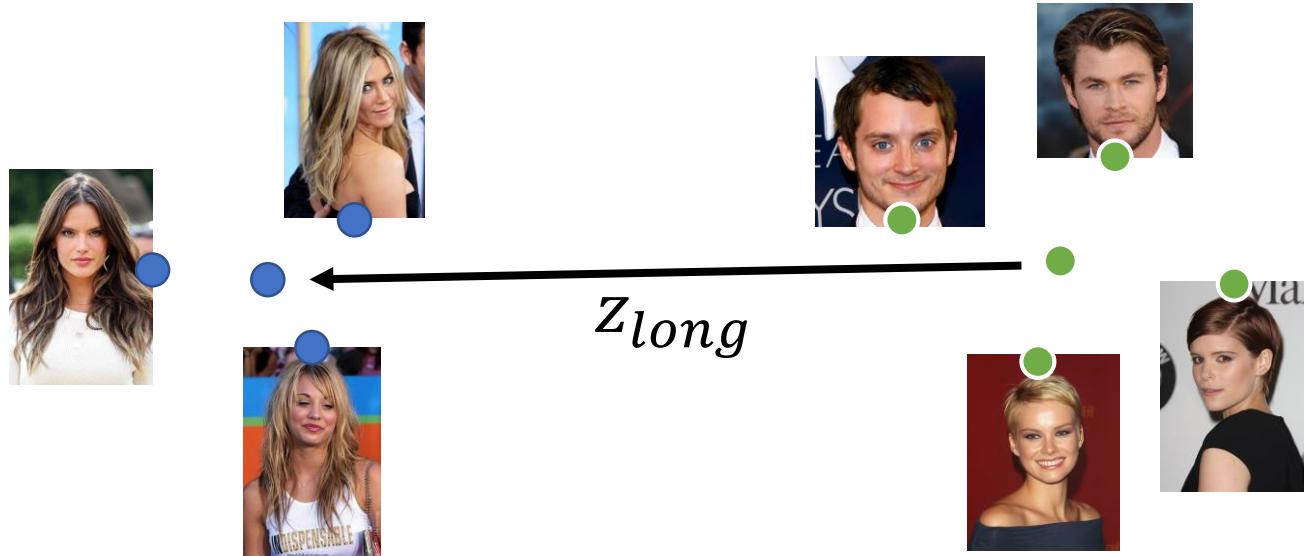
Image	Attributes
	Arched eyebrows, attractive, brown hair, heavy makeup, high cheekbones, mouth slightly open, no beard, pointy nose, smiling, straight hair, wearing earrings, wearing lipstick, young.
	5 o'clock shadows, attractive, bags under eyes, big lips, big nose, black hair, bushy eyebrows, male, no beard, pointy nose, straight hair, young.

# GAN+Autoencoder

- We have a generator (input  $z$ , output  $x$ )
- However, given  $x$ , how can we find  $z$ ?
  - Learn an encoder (input  $x$ , output  $z$ )



# Attribute Representation



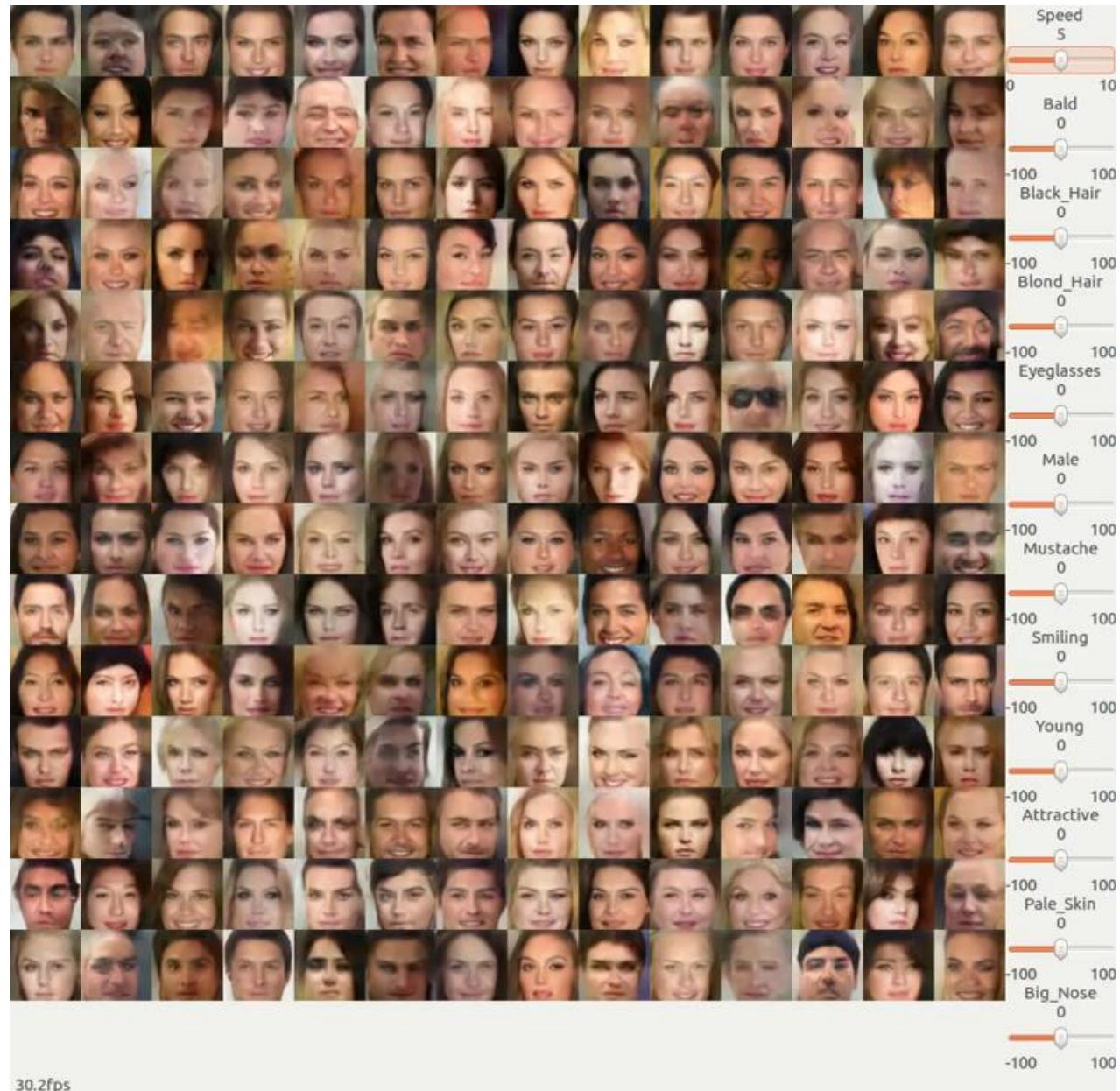
$$z_{long} = \frac{1}{N_1} \sum_{x \in long} En(x) - \frac{1}{N_2} \sum_{x' \notin long} En(x')$$

Short  
Hair

$$x \rightarrow En(x) + z_{long} = z' \rightarrow Gen(z')$$

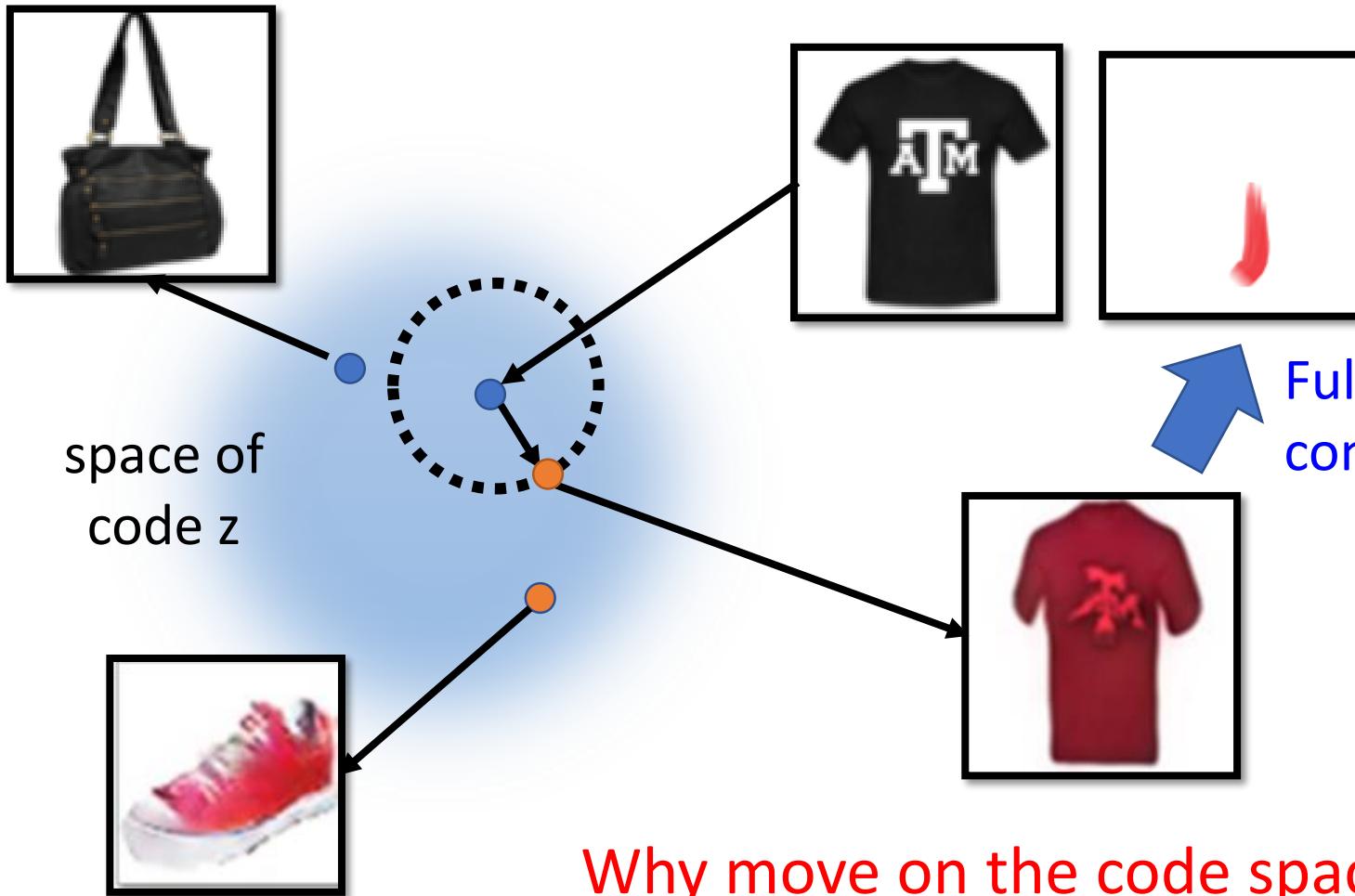
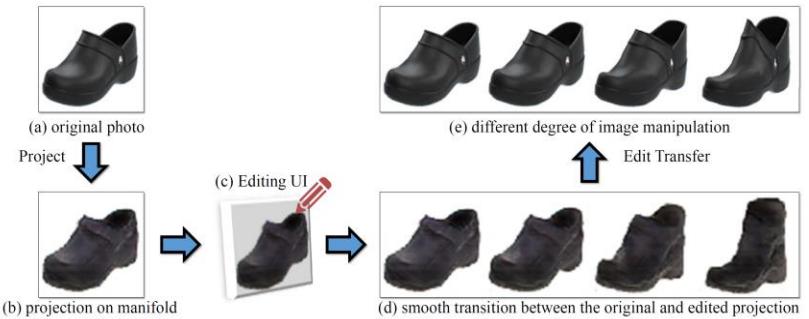
Long  
Hair

# Photo Editing

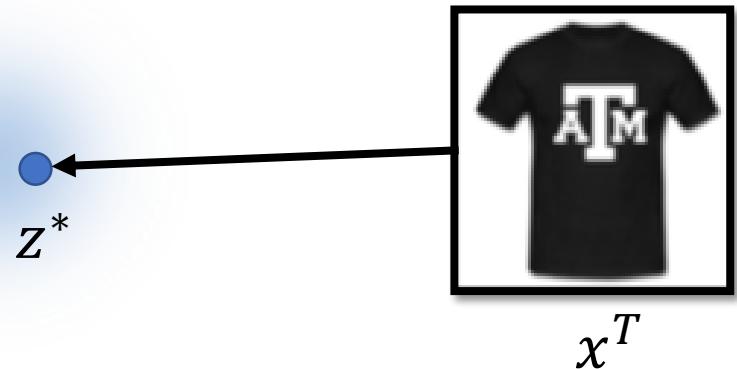


<https://www.youtube.com/watch?v=kPEIJJsQr7U>

# Basic Idea



# Back to z



- **Method 1**

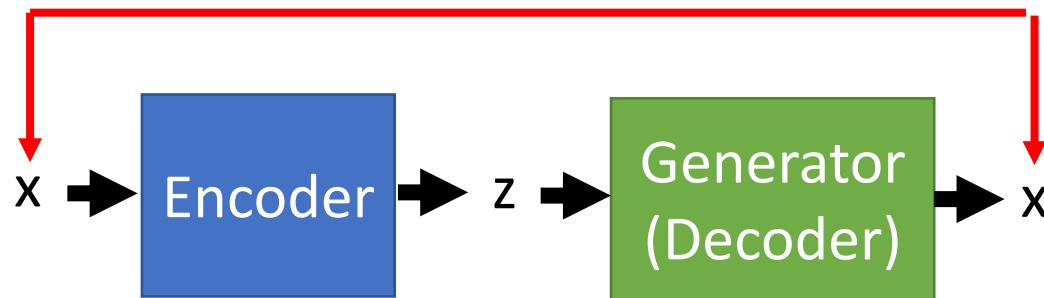
$$z^* = \arg \min_z L(G(z), x^T)$$

Gradient Descent

→ Difference between  $G(z)$  and  $x^T$   
➤ Pixel-wise  
➤ By another network

- **Method 2**

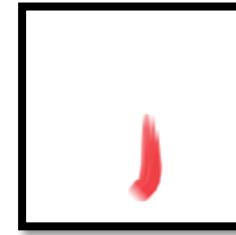
as close as possible



- **Method 3**

Using the results from method 2 as the initialization of method 1

# Editing Photos



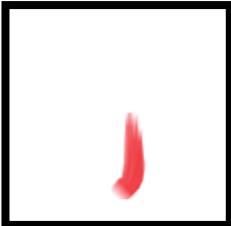
- $z_0$  is the code of the input image



Using discriminator  
to check the image is  
realistic or not

$$z^* = \arg \min_z U(G(z)) + \lambda_1 \|z - z_0\|^2 - \lambda_2 D(G(z))$$

Not too far away from  
the original image



Does it fulfill the constraint of editing?

# **Generative Visual Manipulation on the Natural Image Manifold**

**Jun-Yan Zhu**

**Philipp Krähenbühl**

**Eli Shechtman**

**Alexei A. Efros**



<https://www.youtube.com/watch?v=9c4z6YsBGQ0>

Jun-Yan Zhu, Philipp Krähenbühl, Eli Shechtman and Alexei A. Efros. "Generative Visual Manipulation on the Natural Image Manifold", ECCV, 2016.



## Neural Photo Editing

Andrew Brock



Andrew Brock, Theodore Lim, J.M. Ritchie, Nick Weston, **Neural Photo Editing with Introspective Adversarial Networks**, arXiv preprint, 2017

# Image super resolution

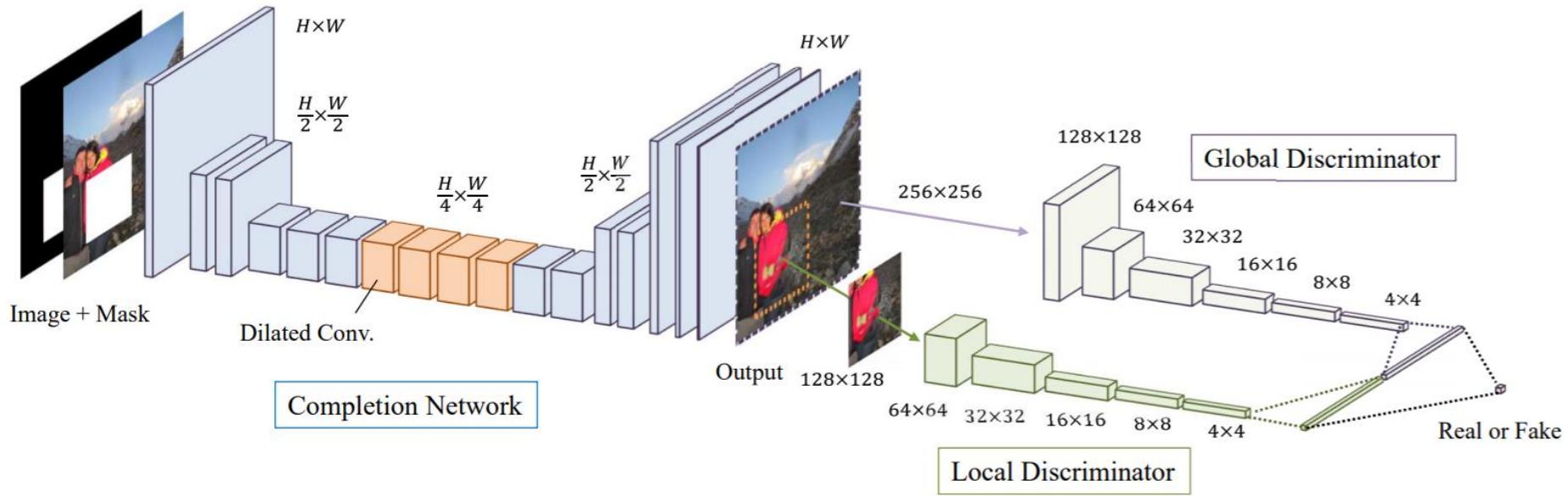
- Christian Ledig, Lucas Theis, Ferenc Huszar, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, Wenzhe Shi, “Photo-Realistic Single Image Super-Resolution Using a Generative Adversarial Network”, CVPR, 2016



Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [4× upscaling]

# Image Completion

<http://hi.cs.waseda.ac.jp/~iizuka/projects/completion/en/>



# Demo

Image completion is a very complicated task...



Previous approach



Previous approach

<https://www.youtube.com/watch?v=5Ua4NUKowPU>

# Improving Sequence Generation by GAN

李宏毅

Hung-yi Lee

# Outline

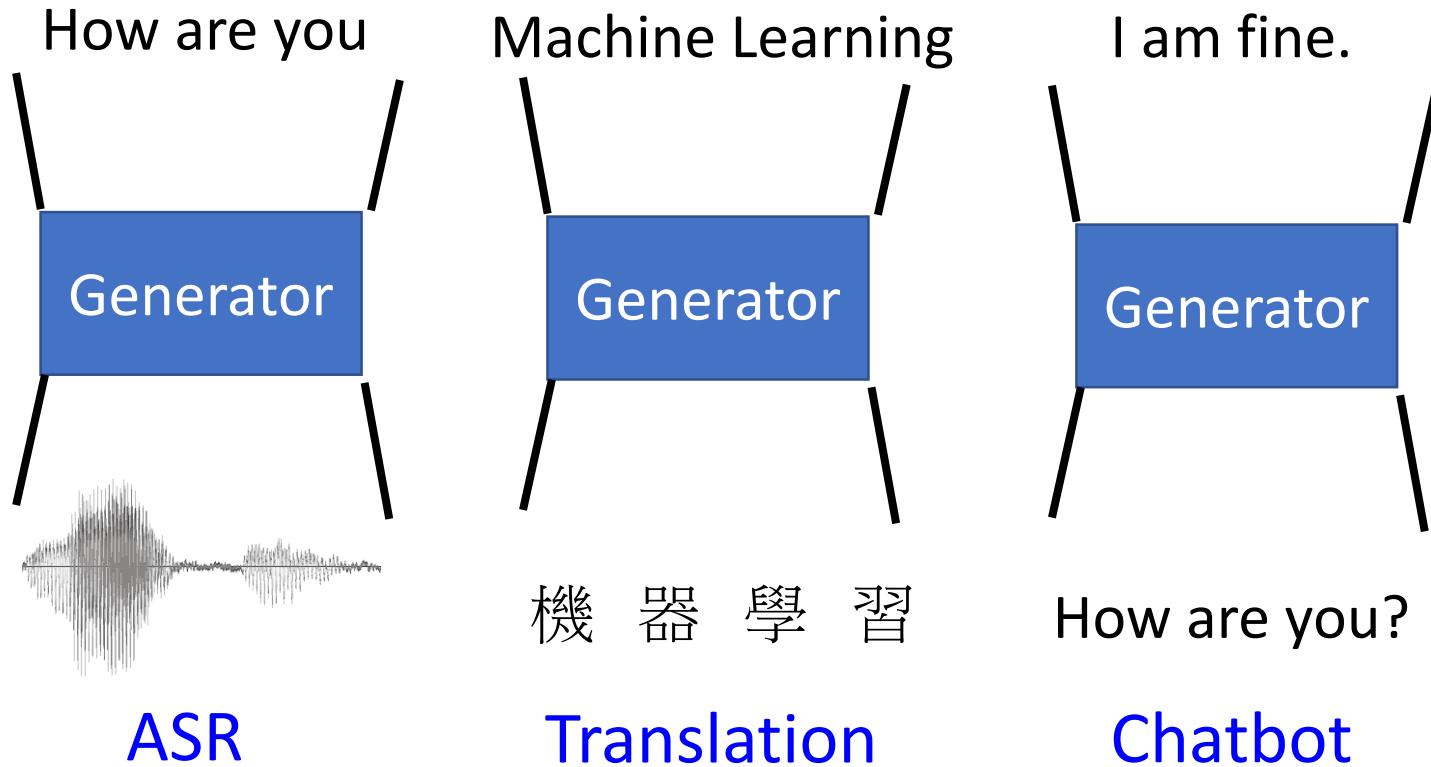
## Conditional Sequence Generation

- RL (human feedback)
- GAN (discriminator feedback)

## Unsupervised Conditional Sequence Generation

- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

# Conditional Sequence Generation



The generator is a typical seq2seq model.

With GAN, you can train seq2seq model in another way.

# Review: Sequence-to-sequence

- Chat-bot as example

Output:	Not bad	I'm John.
Human	better	
Training Criterion		better

Training  
data:  
⋮

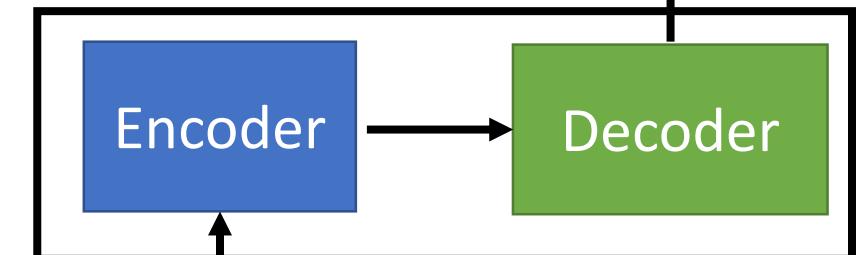
A: How are you ?

B: I'm good.  
⋮

Maximize  
likelihood

I'm good.

output  
sentence x



Input sentence c  
How are you ?

Generator

# Outline of Part III

## Improving Supervised Seq-to-seq Model

- RL (human feedback)
- GAN (discriminator feedback)

## Unsupervised Seq-to-seq Model

- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

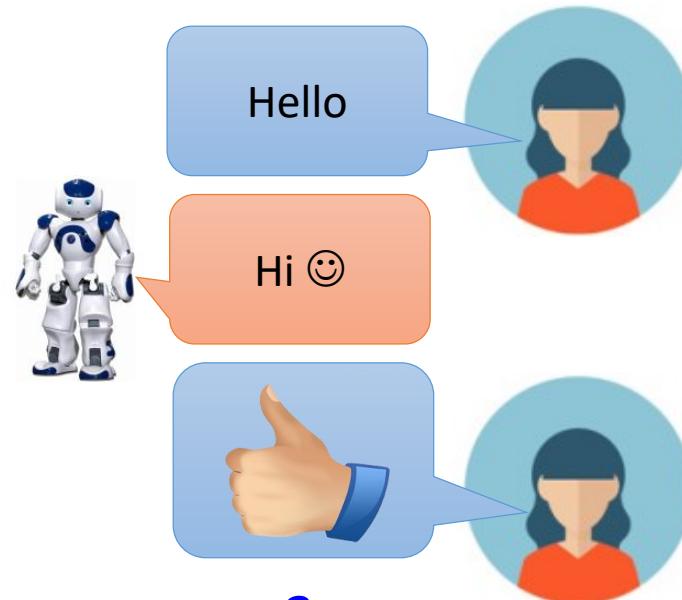
# Introduction

- Machine obtains feedback from user



-10

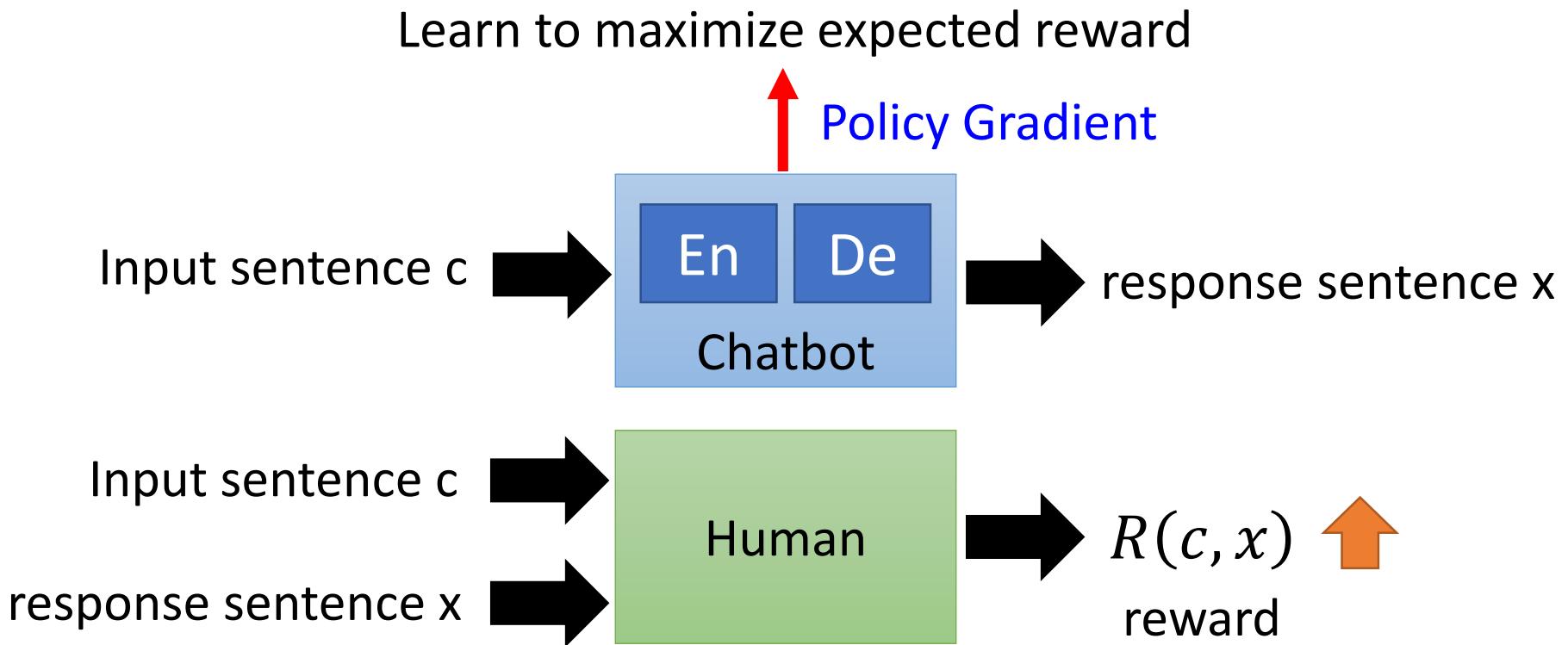
[https://image.freepik.com/free-vector/variety-of-human-avatars\\_23-2147506285.jpg](https://image.freepik.com/free-vector/variety-of-human-avatars_23-2147506285.jpg)  
[http://www.freepik.com/free-vector/variety-of-human-avatars\\_766615.htm](http://www.freepik.com/free-vector/variety-of-human-avatars_766615.htm)



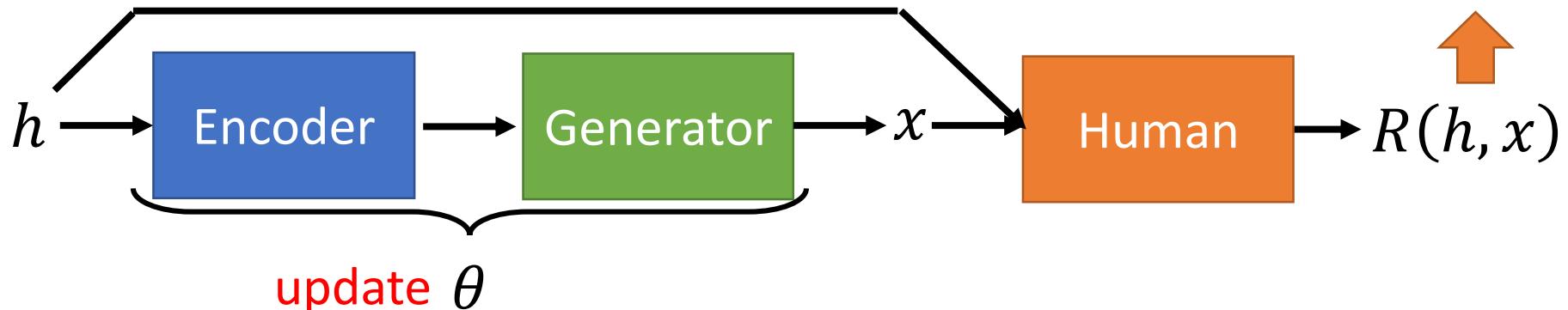
3

- Chat-bot learns to maximize the *expected reward*

# Maximizing Expected Reward



# Maximizing Expected Reward

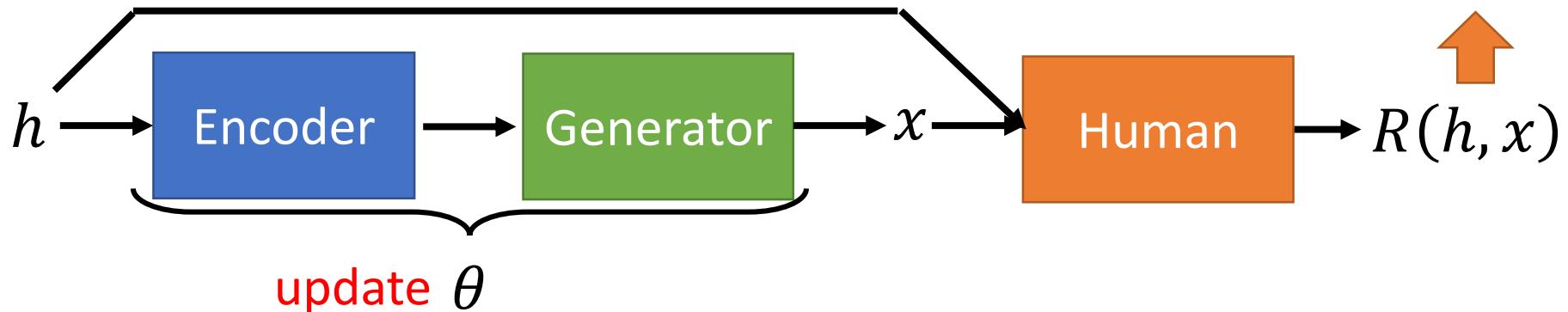


$$\theta^* = \arg \max_{\theta} \bar{R}_{\theta} \quad \leftarrow \text{Maximizing expected reward}$$

$$\bar{R}_{\theta} = \sum_h P(h) \sum_x R(h, x) \frac{P_{\theta}(x|h)}{\text{Randomness in generator}}$$

Probability that the input/history is  $h$

# Maximizing Expected Reward



$$\theta^* = \arg \max_{\theta} \bar{R}_{\theta} \quad \leftarrow \text{Maximizing expected reward}$$

$$\begin{aligned}\bar{R}_{\theta} &= \sum_h P(h) \sum_x R(h, x) P_{\theta}(x|h) = E_{h \sim P(h)} \left[ E_{x \sim P_{\theta}(x|h)} [R(h, x)] \right] \\ &= E_{h \sim P(h), x \sim P_{\theta}(x|h)} [R(h, x)] \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i)\end{aligned}$$

**Sample:**  $(h^1, x^1), (h^2, x^2), \dots, (h^N, x^N)$

Where  
is  $\theta$ ?

# Policy Gradient

$$\frac{d \log(f(x))}{dx} = \frac{1}{f(x)} \frac{df(x)}{dx}$$

$$\bar{R}_\theta = \sum_h P(h) \sum_x R(h, x) P_\theta(x|h) \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i)$$

$$\nabla \bar{R}_\theta = \sum_h P(h) \sum_x R(h, x) \nabla P_\theta(x|h) \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \nabla \log P_\theta(x|h)$$

$$= \sum_h P(h) \sum_x R(h, x) P_\theta(x|h) \frac{\nabla P_\theta(x|h)}{P_\theta(x|h)}$$

$$= \sum_h P(h) \sum_x R(h, x) P_\theta(x|h) \boxed{\nabla \log P_\theta(x|h)}$$

$$= E_{h \sim P(h), x \sim P_\theta(x|h)} [R(h, x) \nabla \log P_\theta(x|h)]$$

Sampling



# Policy Gradient

- Gradient Ascent

$$\theta^{new} \leftarrow \theta^{old} + \eta \nabla \bar{R}_{\theta^{old}}$$

$$\nabla \bar{R}_{\theta} \approx \frac{1}{N} \sum_{i=1}^N R(h^i, x^i) \nabla \log P_{\theta}(x^i | h^i)$$

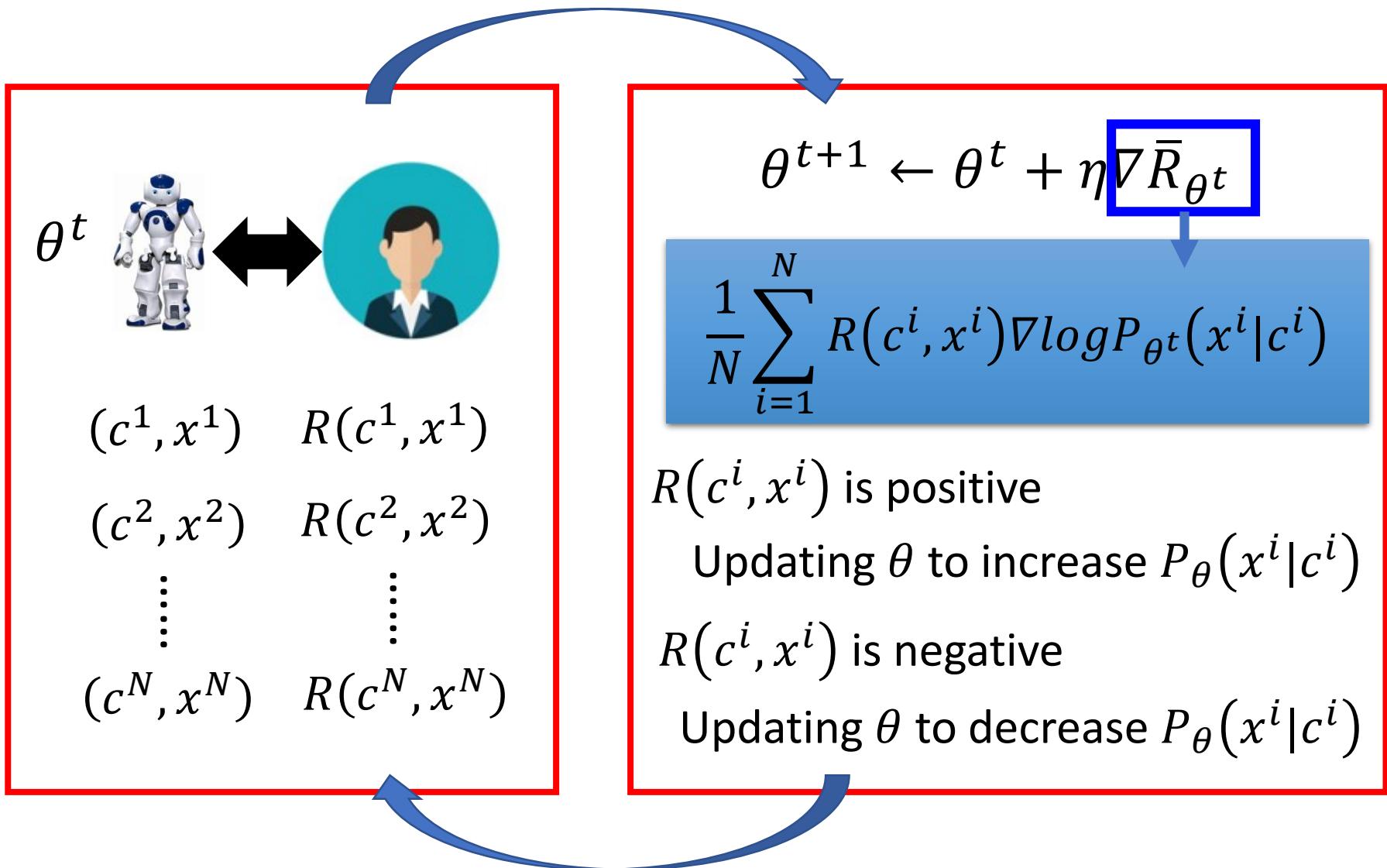
$R(h^i, x^i)$  is positive

→ After updating  $\theta$ ,  $P_{\theta}(x^i | h^i)$  will increase

$R(h^i, x^i)$  is negative

→ After updating  $\theta$ ,  $P_{\theta}(x^i | h^i)$  will decrease

# Policy Gradient - Implementation



# Comparison

	Maximum Likelihood	Reinforcement Learning
Objective Function	$\frac{1}{N} \sum_{i=1}^N \log P_\theta(\hat{x}^i   c^i)$	$\frac{1}{N} \sum_{i=1}^N R(c^i, x^i) \log P_\theta(x^i   c^i)$
Gradient	$\frac{1}{N} \sum_{i=1}^N \nabla \log P_\theta(\hat{x}^i   c^i)$	$\frac{1}{N} \sum_{i=1}^N R(c^i, x^i) \nabla \log P_\theta(x^i   c^i)$
Training Data	$\{(c^1, \hat{x}^1), \dots, (c^N, \hat{x}^N)\}$ $R(c^i, \hat{x}^i) = 1$	$\{(c^1, x^1), \dots, (c^N, x^N)\}$ obtained from interaction weighted by $R(c^i, x^i)$

# Alpha GO style training !



- Let two agents talk to each other



How old are you?



See you.



How old are you?



I am 16.



See you.



See you.



I thought you were 12.



What make you  
think so?

Using a pre-defined evaluation function to compute  $R(h,x)$

# Outline of Part III

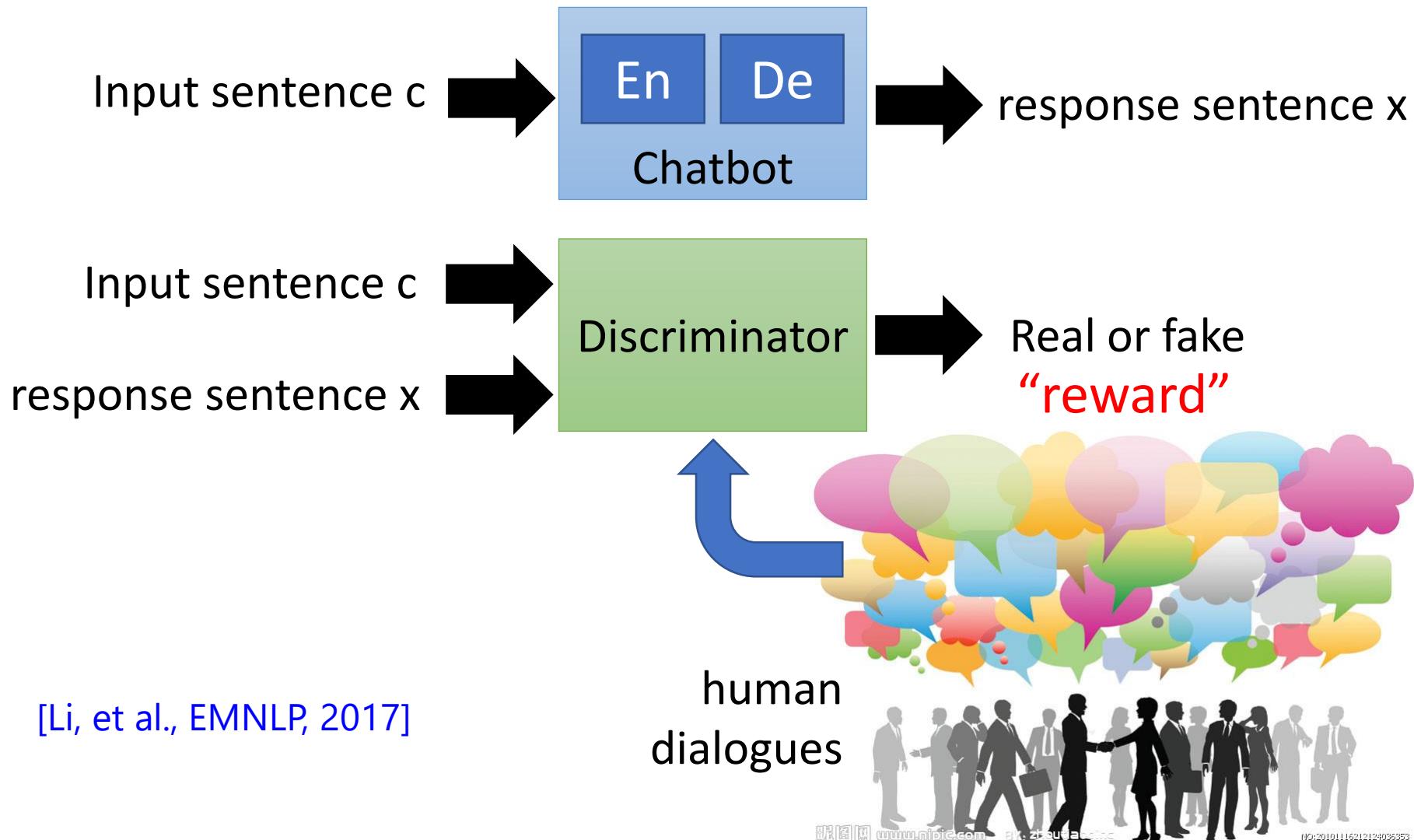
## Improving Supervised Seq-to-seq Model

- RL (human feedback)
- GAN (discriminator feedback)

## Unsupervised Seq-to-seq Model

- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

# Conditional GAN



# Algorithm

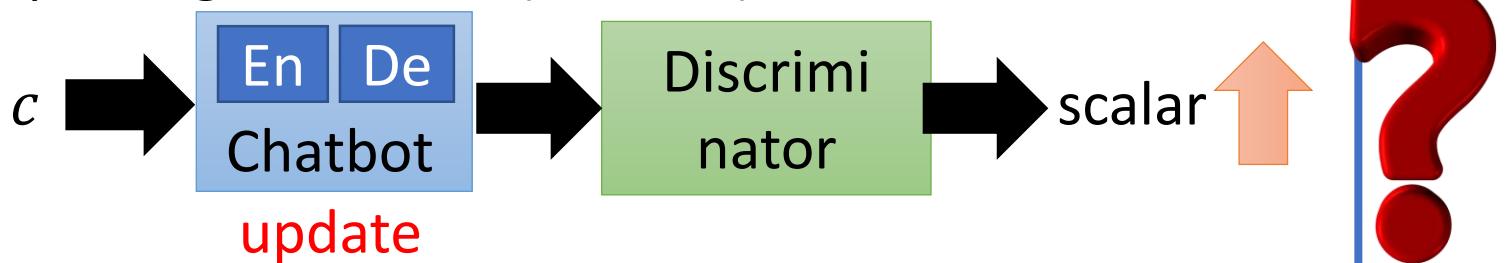
Training data:

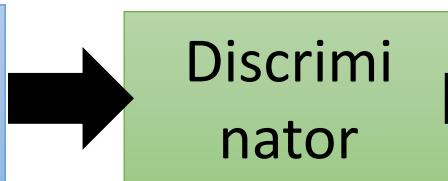
Pairs of conditional input  $c$  and response  $x$

- Initialize generator  $G$  (chatbot) and discriminator  $D$
- In each iteration:

- Sample input  $c$  and response  $x$  from training set
- Sample input  $c'$  from training set, and generate response  $\tilde{x}$  by  $G(c')$
- Update  $D$  to increase  $D(c, x)$  and decrease  $D(c', \tilde{x})$

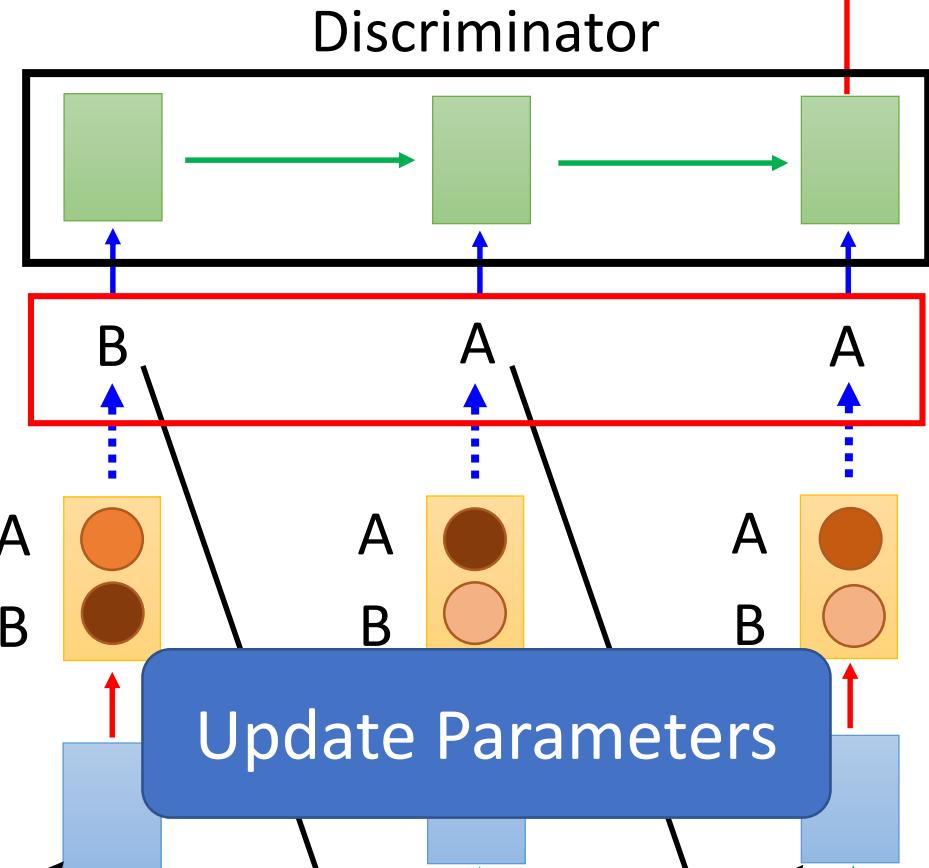
- Update generator  $G$  (chatbot) such that





Can we use  
gradient ascent?

# NO!



Due to the sampling process, “discriminator+ generator”  
is not differentiable



# Three Categories of Solutions

## Gumbel-softmax

- [Matt J. Kusner, et al, arXiv, 2016]

## Continuous Input for Discriminator

- [Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

## “Reinforcement Learning”

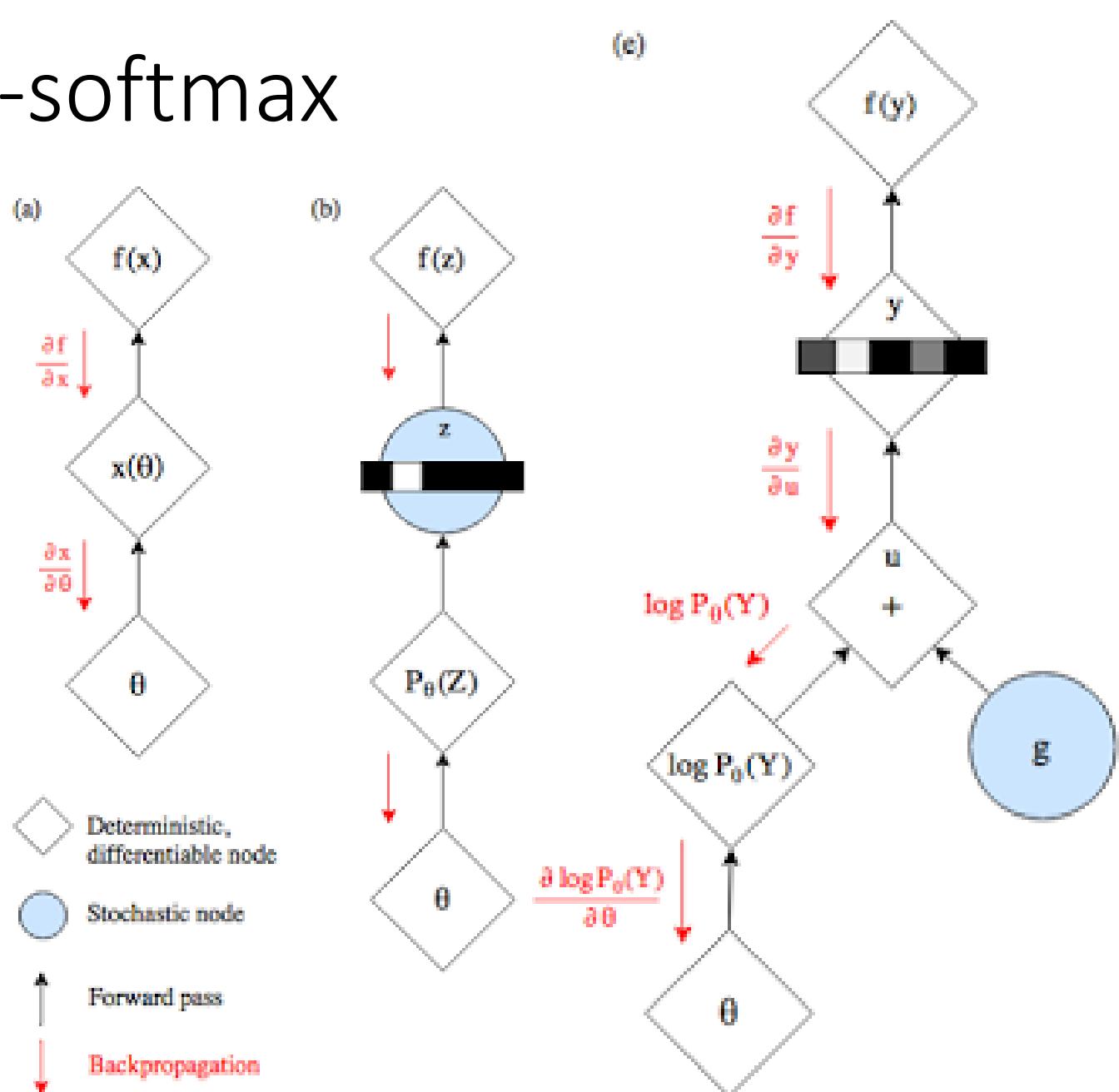
- [Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]

# Gumbel-softmax

<https://gabrielhuang.gitbooks.io/machine-learning/reparametrization-trick.html>

<https://casmls.github.io/general/2017/02/01/GumbelSoftmax.html>

<http://blog.evjang.com/2016/11/tutorial-categorical-variational.html>



# Three Categories of Solutions

## Gumbel-softmax

- [Matt J. Kusner, et al, arXiv, 2016]

## Continuous Input for Discriminator

- [Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

## “Reinforcement Learning”

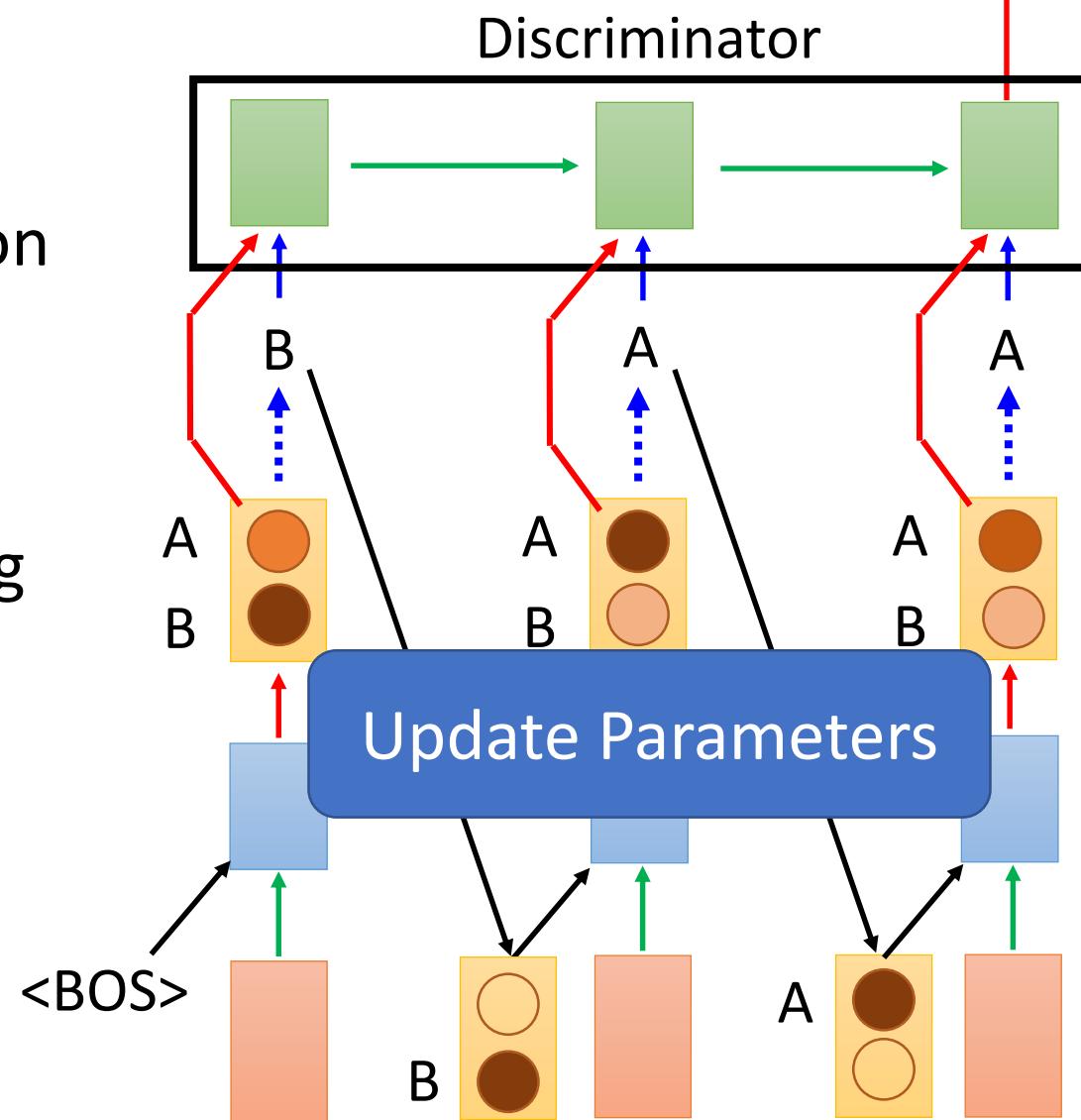
- [Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]



Use the distribution as the input of discriminator

Avoid the sampling process

We can do backpropagation now.



# What is the problem?

- Real sentence

1	0	0	0	0
0	1	0	0	0
0	0	1	0	0
0	0	0	1	0
0	0	0	0	1

Discriminator can immediately find the difference.

- Generated

Can never be 1-of-N

0.9	0.1	0.1	0	0
0.1	0.9	0.1	0	0
0	0	0.7	0.1	0
0	0	0.1	0.8	0.1
0	0	0	0.1	0.9

WGAN is helpful

# Three Categories of Solutions

## Gumbel-softmax

- [Matt J. Kusner, et al, arXiv, 2016]

## Continuous Input for Discriminator

- [Sai Rajeswar, et al., arXiv, 2017][Ofir Press, et al., ICML workshop, 2017][Zhen Xu, et al., EMNLP, 2017][Alex Lamb, et al., NIPS, 2016][Yizhe Zhang, et al., ICML, 2017]

## “Reinforcement Learning”

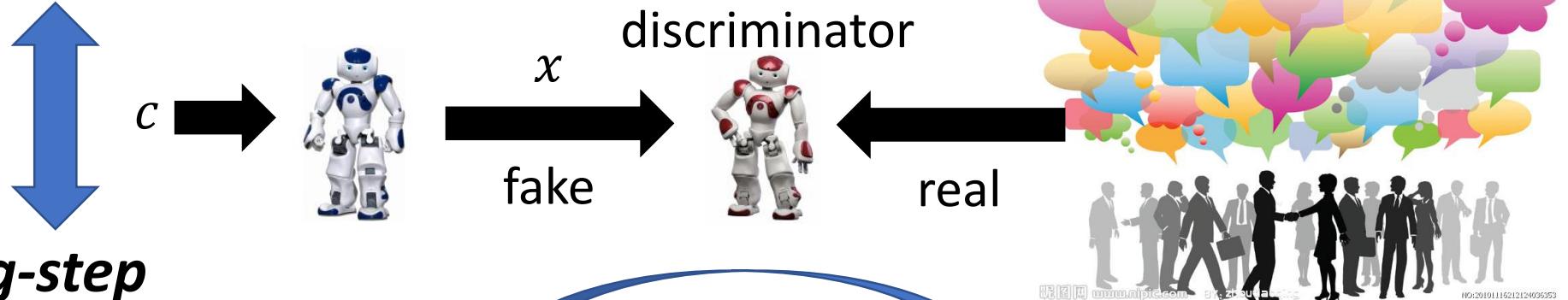
- [Yu, et al., AAAI, 2017][Li, et al., EMNLP, 2017][Tong Che, et al, arXiv, 2017][Jiaxian Guo, et al., AAAI, 2018][Kevin Lin, et al, NIPS, 2017][William Fedus, et al., ICLR, 2018]

# Reinforcement Learning?

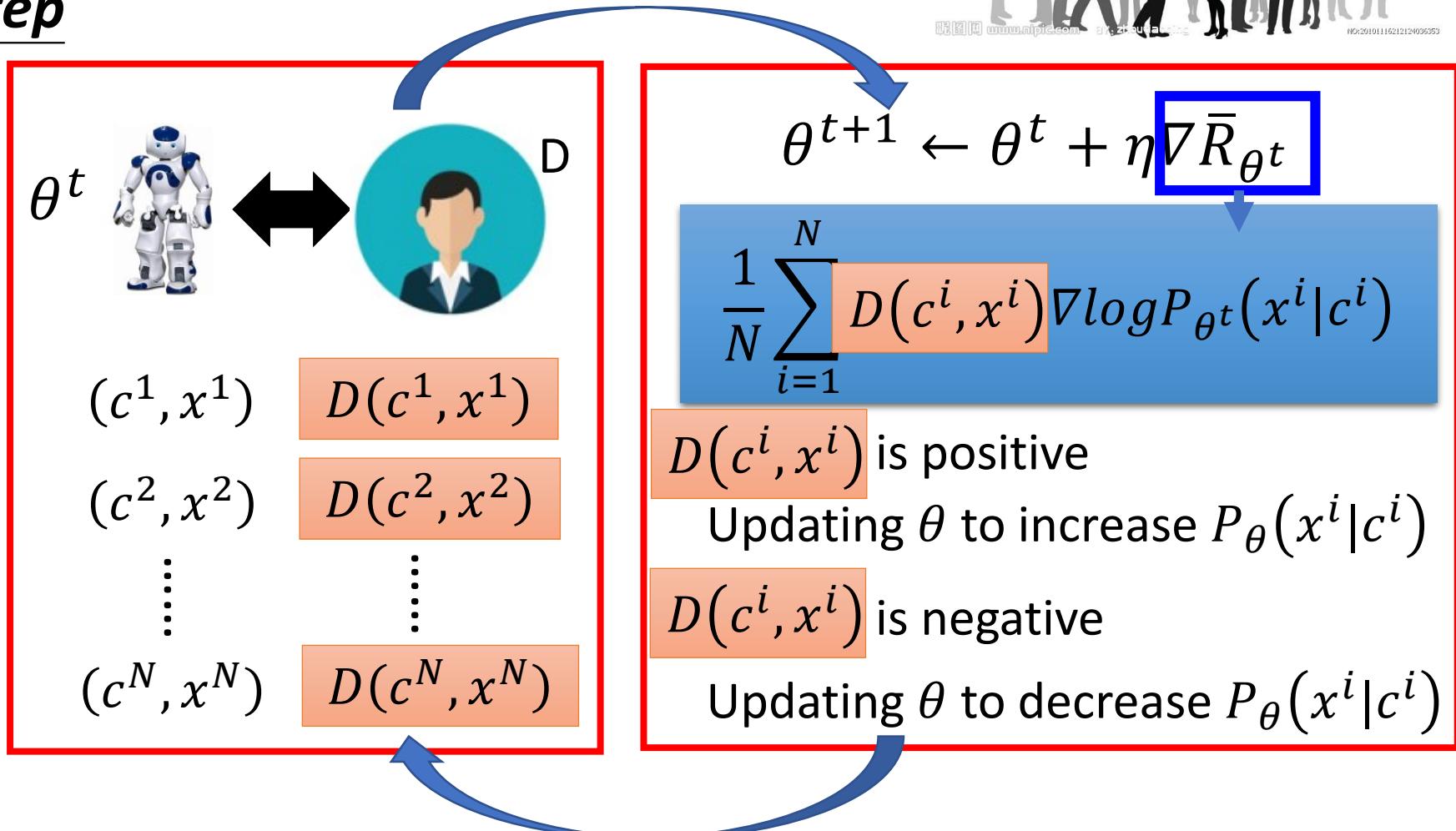


- Consider the output of discriminator as **reward**
  - Update generator to increase discriminator = to get maximum reward
  - Using the formulation of policy gradient, replace reward  $R(c, x)$  with discriminator output  $D(c, x)$
- Different from typical RL
  - The discriminator would update

## d-step



## g-step



# Reward for Every Generation Step

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{i=1}^N D(c^i, x^i) \nabla \log P_\theta(x^i | c^i)$$

$c^i$  = “What is your name?”       $D(c^i, x^i)$  is negative

$x^i$  = “I don’t know”      Update  $\theta$  to decrease  $\log P_\theta(x^i | c^i)$

$$\log P_\theta(x^i | c^i) = \log P(x_1^i | c^i) + \log P(x_2^i | c^i, x_1^i) + \log P(x_3^i | c^i, x_{1:2}^i)$$

$P("I" | c^i)$     

$c^i$  = “What is your name?”       $D(c^i, x^i)$  is positive

$x^i$  = “I am John”      Update  $\theta$  to increase  $\log P_\theta(x^i | c^i)$

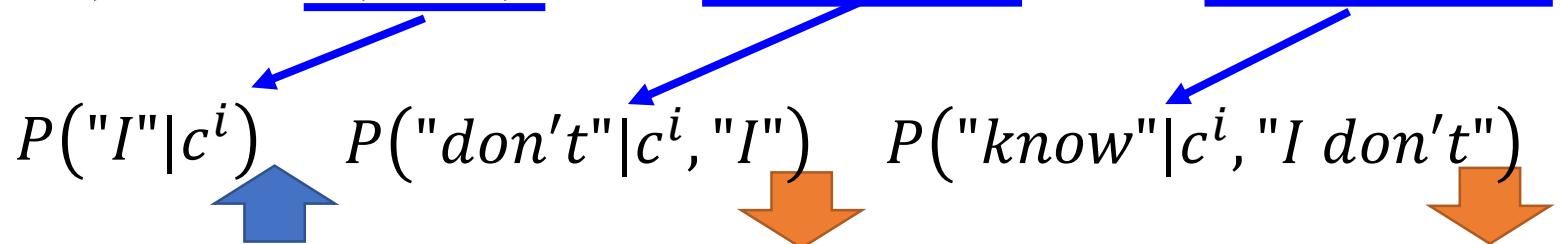
$$\log P_\theta(x^i | c^i) = \log P(x_1^i | c^i) + \log P(x_2^i | c^i, x_1^i) + \log P(x_3^i | c^i, x_{1:2}^i)$$

$P("I" | c^i)$    

# Reward for Every Generation Step

$h^i = \text{"What is your name?"}$      $x^i = \text{"I don't know"}$

$$\log P_{\theta}(x^i | h^i) = \underbrace{\log P(x_1^i | c^i)} + \underbrace{\log P(x_2^i | c^i, x_1^i)} + \underbrace{\log P(x_3^i | c^i, x_{1:2}^i)}$$



$$\nabla \bar{R}_{\theta} \approx \frac{1}{N} \sum_{i=1}^N \underbrace{D(c^i, x^i)}_{\text{green}} \nabla \underbrace{\log P_{\theta}(x^i | c^i)}_{\text{red}}$$

$$\nabla \bar{R}_{\theta} \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \underbrace{(Q(c^i, x_{1:t}^i) - b)}_{\text{green}} \nabla \underbrace{\log P_{\theta}(x_t^i | c^i, x_{1:t-1}^i)}_{\text{red}}$$

Method 1. Monte Carlo (MC) Search [Yu, et al., AAAI, 2017]

Method 2. Discriminator For Partially Decoded Sequences

[Li, et al., EMNLP, 2017]

# Tips: RankGAN

Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, Ming-Ting Sun,  
"Adversarial Ranking for Language Generation", NIPS 2017

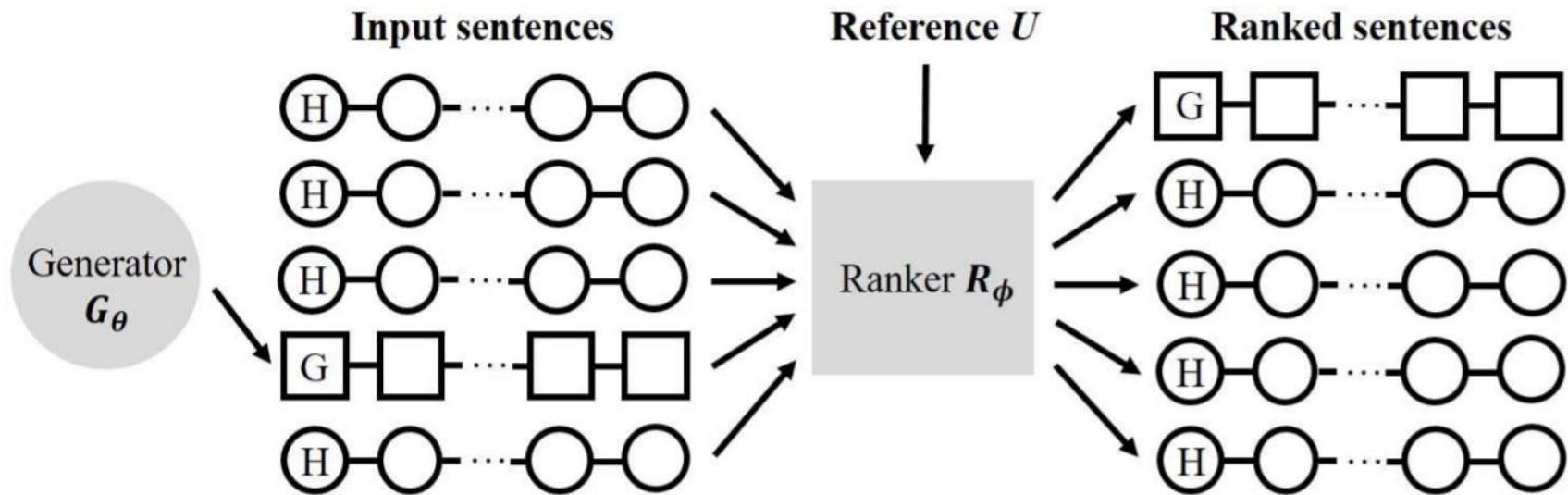


Image caption generation:

Method	BLEU-2	BLEU-3	BLEU-4
MLE	0.781	0.624	0.589
SeqGAN	0.815	0.636	0.587
RankGAN	<b>0.845</b>	<b>0.668</b>	<b>0.614</b>

Method	Human score
SeqGAN	3.44
RankGAN	4.61
Human-written	<b>6.42</b>

# *Experimental Results*

Input	We've got to look for another route.
MLE	I'm sorry.
GAN	You're not going to be here for a while.
Input	You can save him by talking.
MLE	I don't know.
GAN	You know what's going on in there, you know what I mean?

- MLE frequently generates “I’m sorry”, “I don’t know”, etc. (corresponding to fuzzy images?)
- GAN generates longer and more complex responses (however, no strong evidence shows that they are better)

---

Find more comparison in the survey papers.

[Lu, et al., arXiv, 2018][Zhu, et al., arXiv, 2018]

# More Applications

- Supervised machine translation [Wu, et al., arXiv 2017][Yang, et al., arXiv 2017]
- Supervised abstractive summarization [Liu, et al., AAAI 2018]
- Image/video caption generation [Rakshith Shetty, et al., ICCV 2017][Liang, et al., arXiv 2017]

If you are using seq2seq models,  
consider to improve them by GAN.

# Outline of Part III

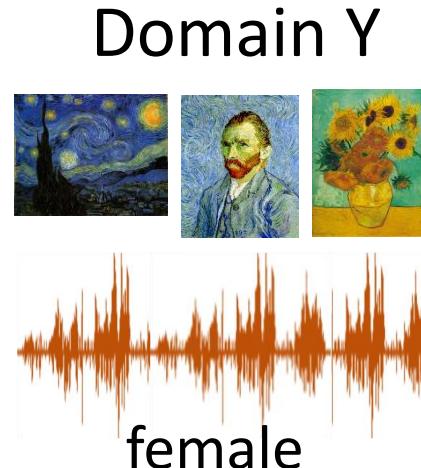
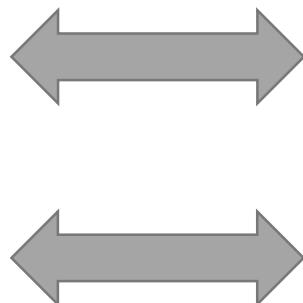
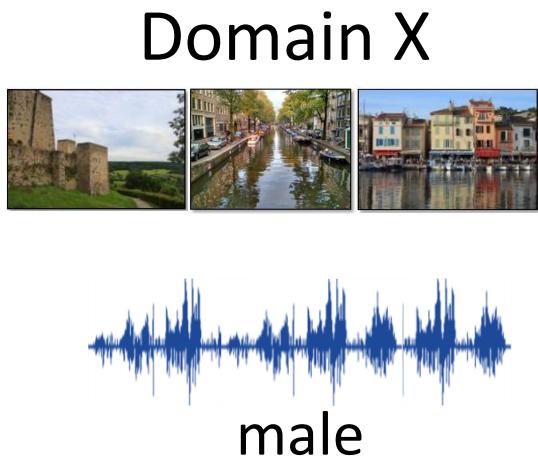
## Conditional Sequence Generation

- RL (human feedback)
- GAN (discriminator feedback)

## Unsupervised Conditional Sequence Generation

- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

# Text Style Transfer



positive sentences

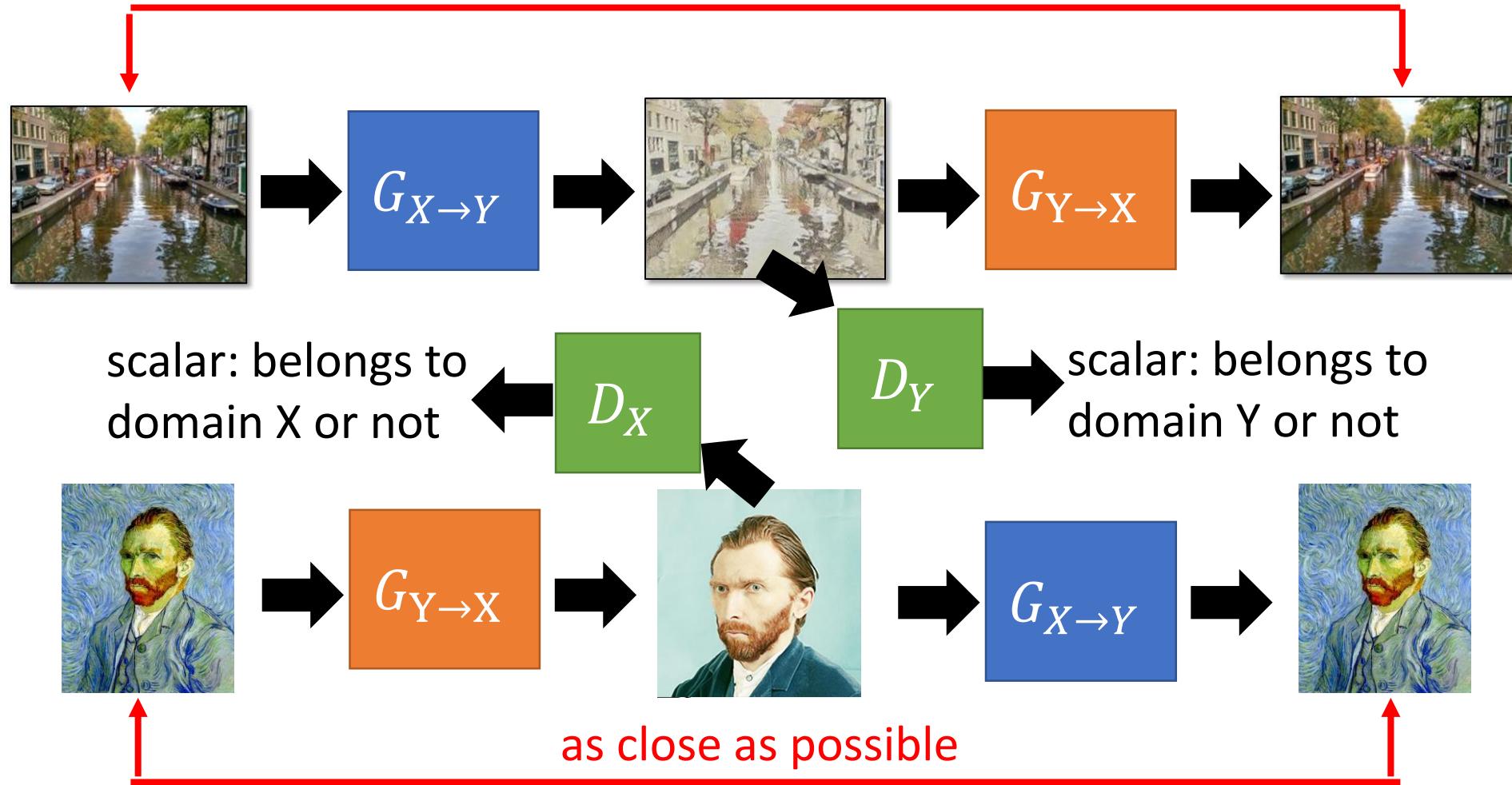
It is good.  
It's a good day.  
I love you.

negative sentences

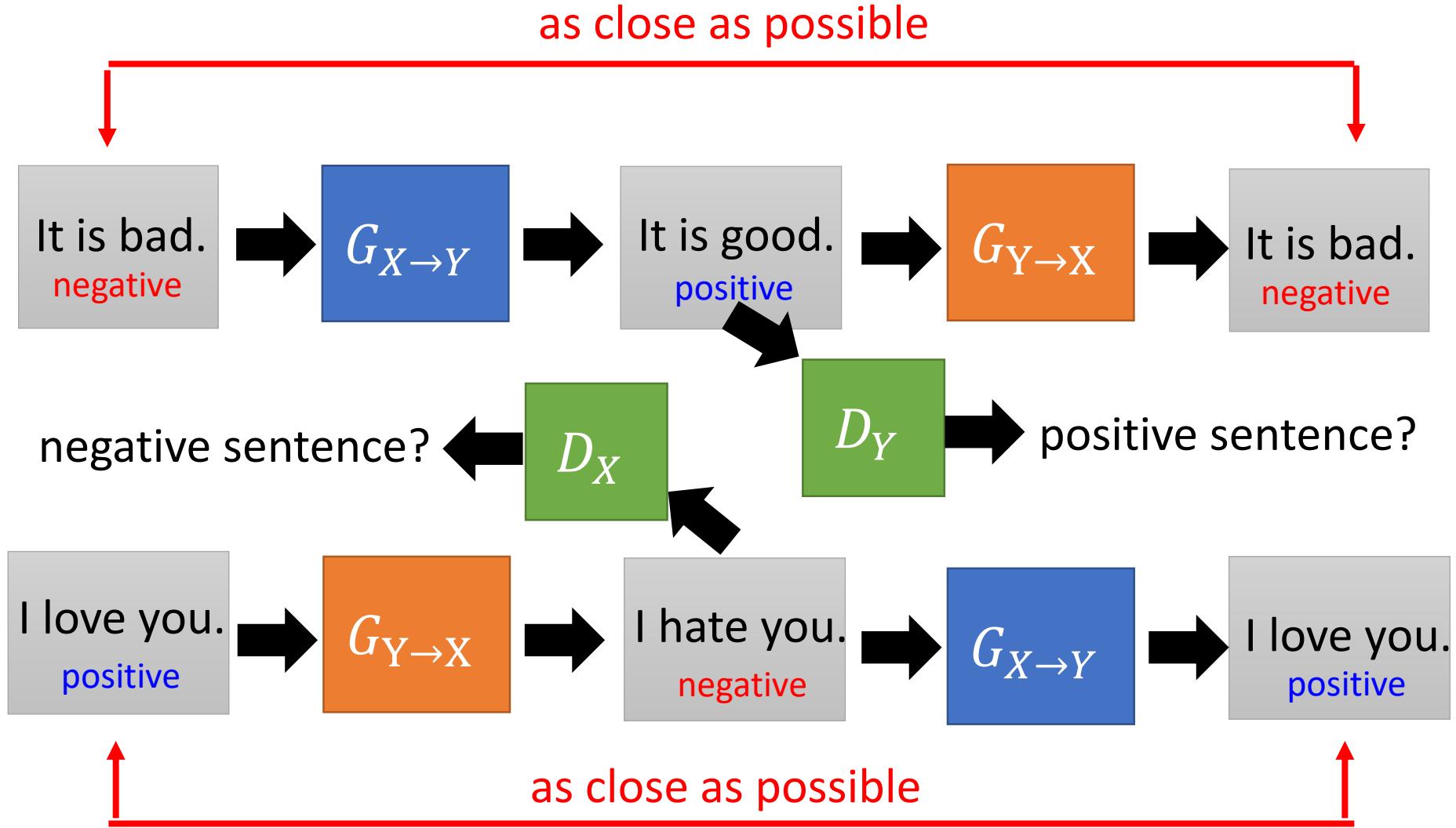
It is bad.  
It's a bad day.  
I don't love you.

# Direct Transformation

as close as possible



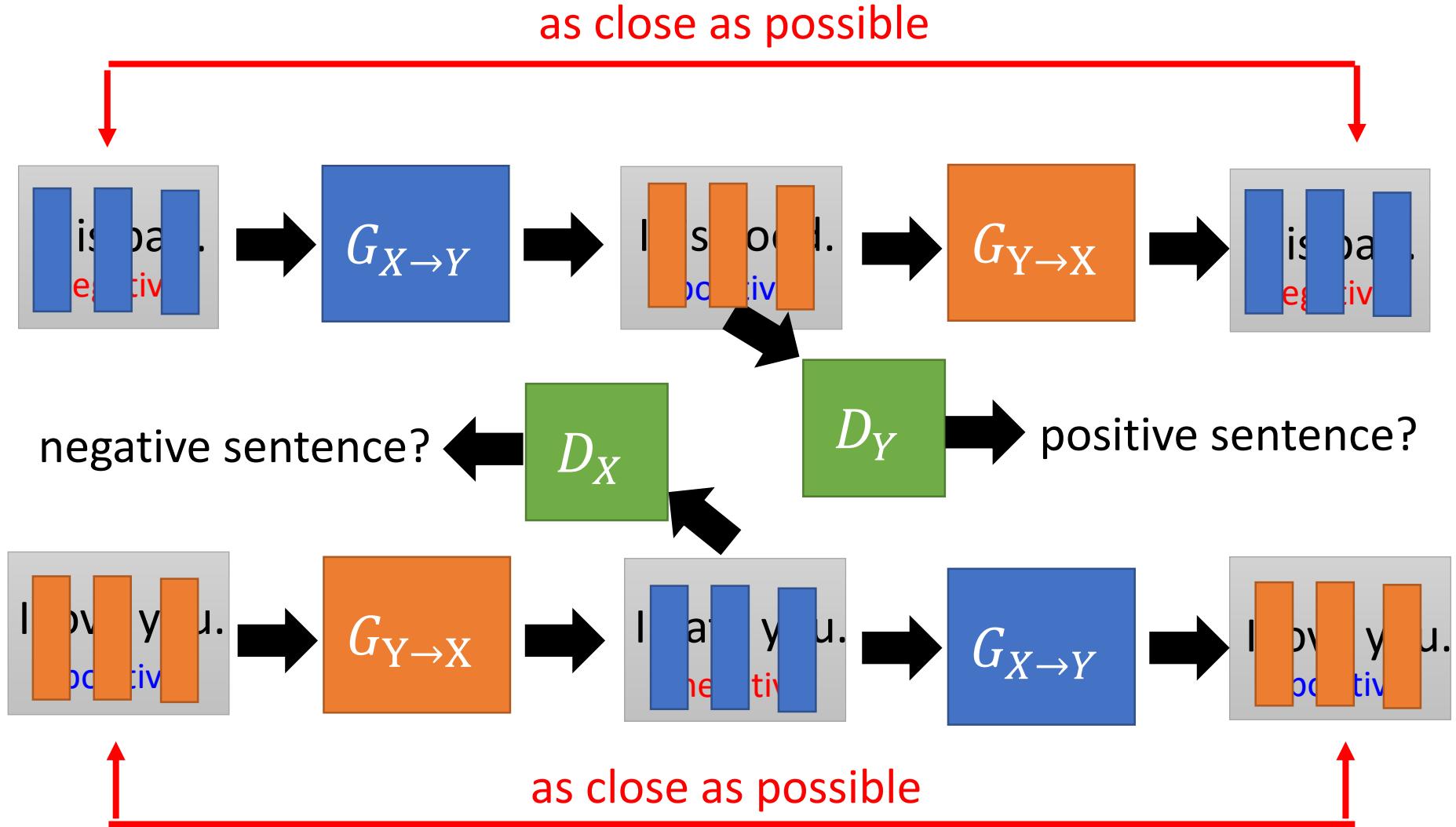
# Direct Transformation



# Direct Transformation

Discrete?

Word embedding  
[Lee, et al., ICASSP, 2018]



- **Negative sentence to **positive** sentence:**

it's a crappy day → it's a great day

i wish you could be here → you could be here

it's not a good idea → it's good idea

i miss you → i love you

i don't love you → i love you

i can't do that → i can do that

i feel so sad → i happy

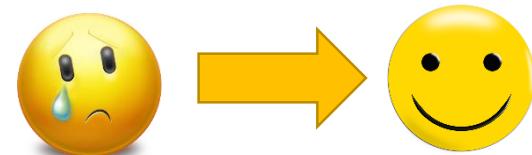
it's a bad day → it's a good day

it's a dummy day → it's a great day

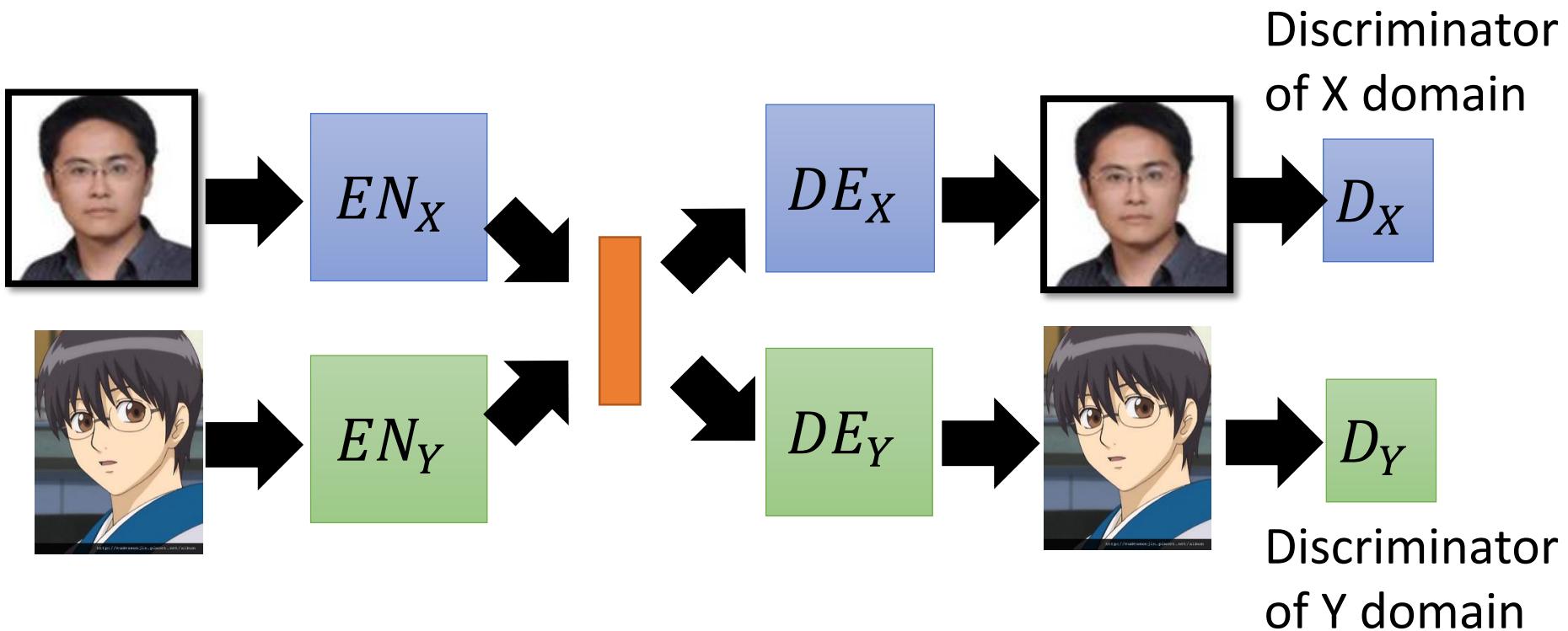
sorry for doing such a horrible thing → thanks for doing a great thing

my doggy is sick → my doggy is my doggy

my little doggy is sick → my little doggy is my little doggy



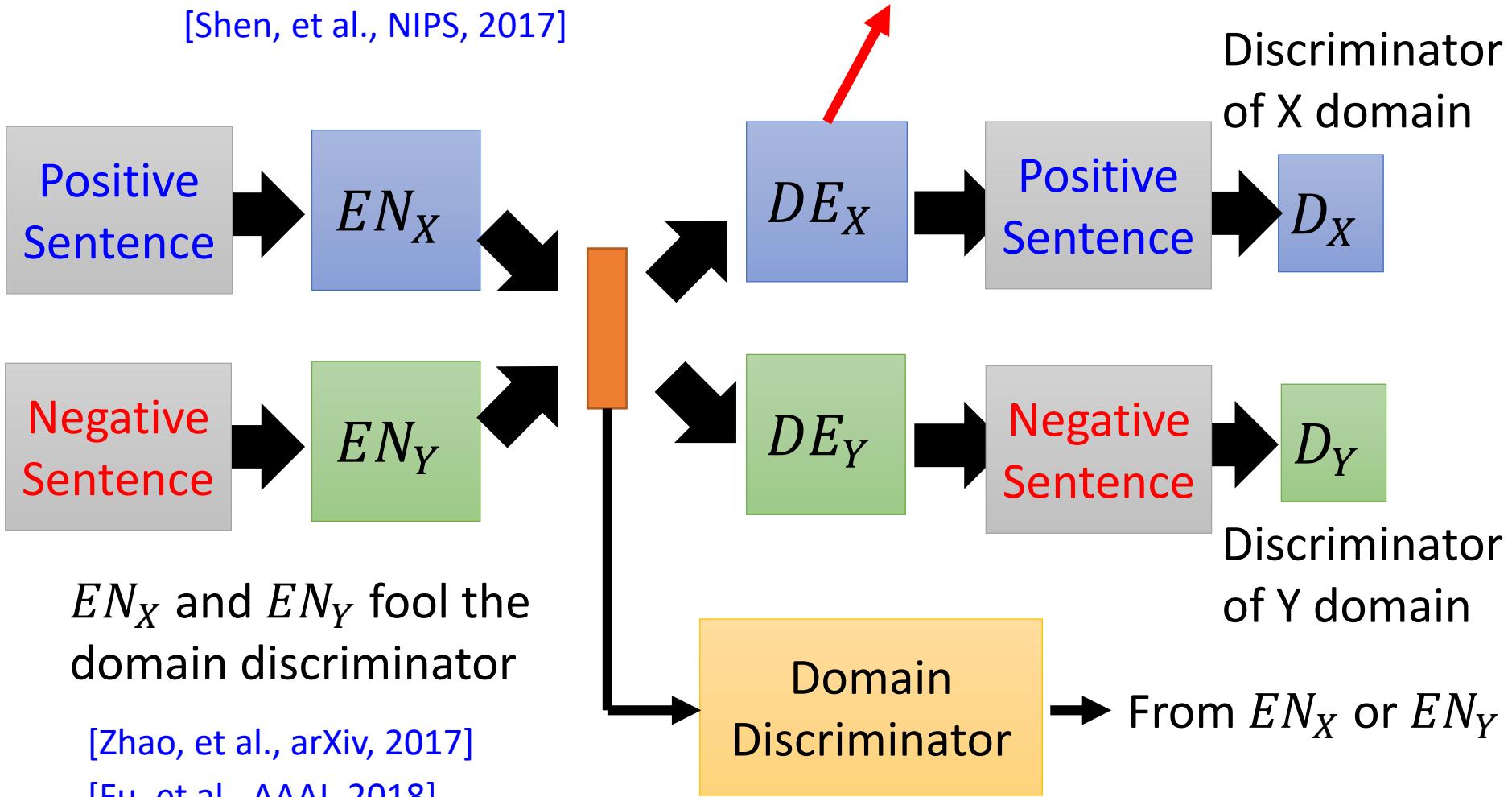
# Projection to Common Space



# Projection to Common Space

Decoder hidden layer as discriminator input

[Shen, et al., NIPS, 2017]



# Outline of Part III

## Improving Supervised Seq-to-seq Model

- RL (human feedback)
- GAN (discriminator feedback)

## Unsupervised Seq-to-seq Model

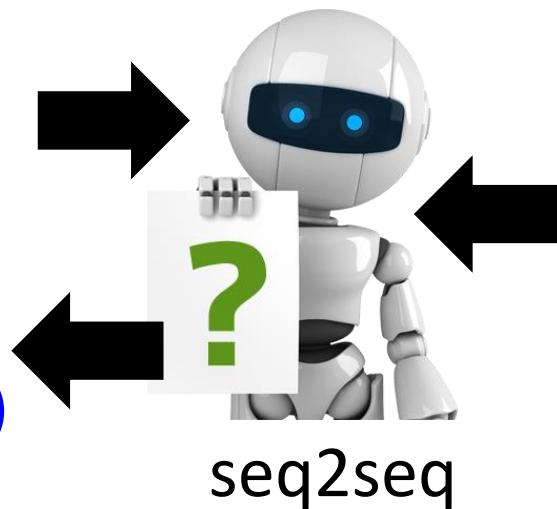
- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

# Abstractive Summarization

- Now machine can do **abstractive summary** by seq2seq (write summaries in its own words)



summary  
(in its own words)



summary 1



summary 2

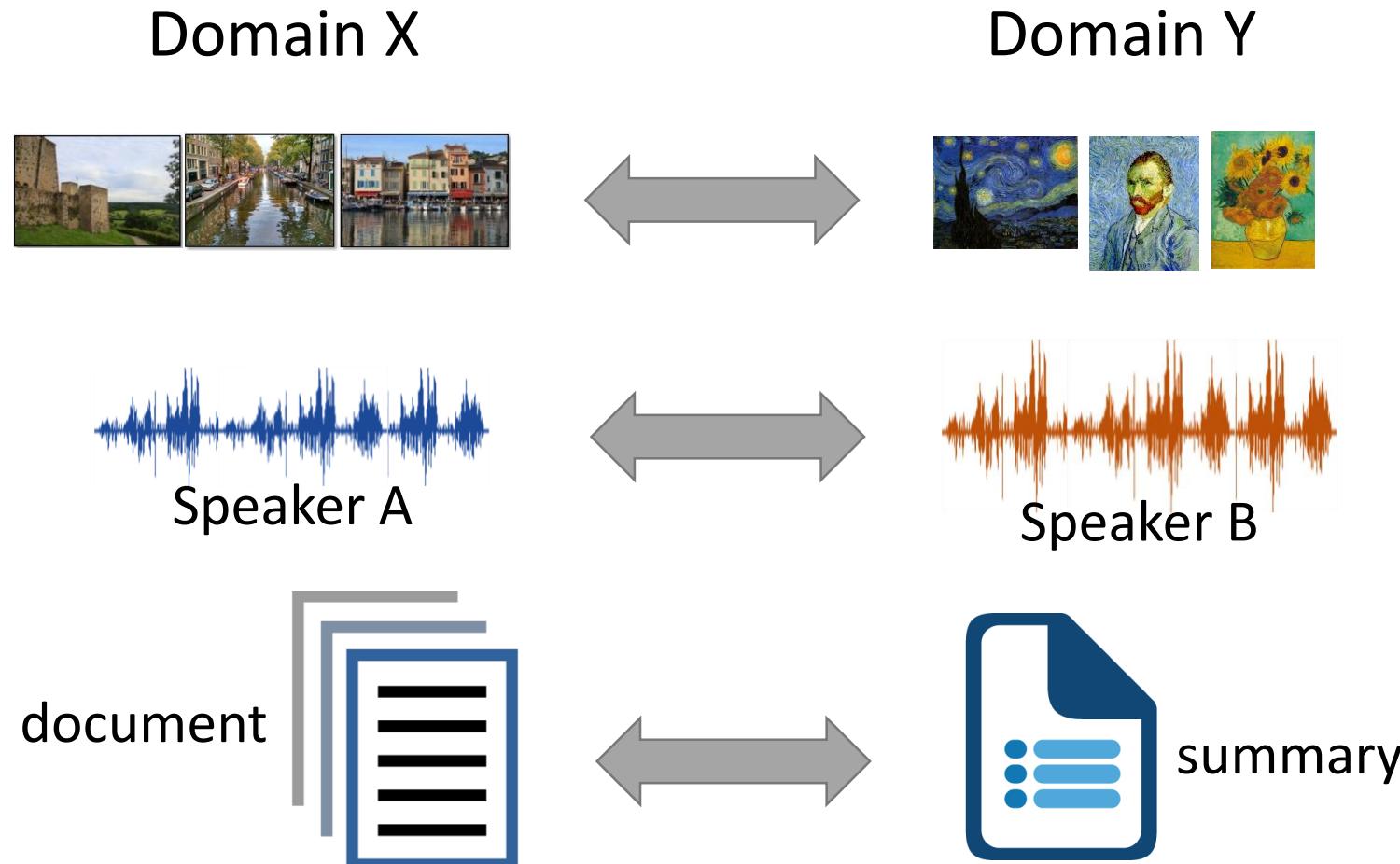


summary 3

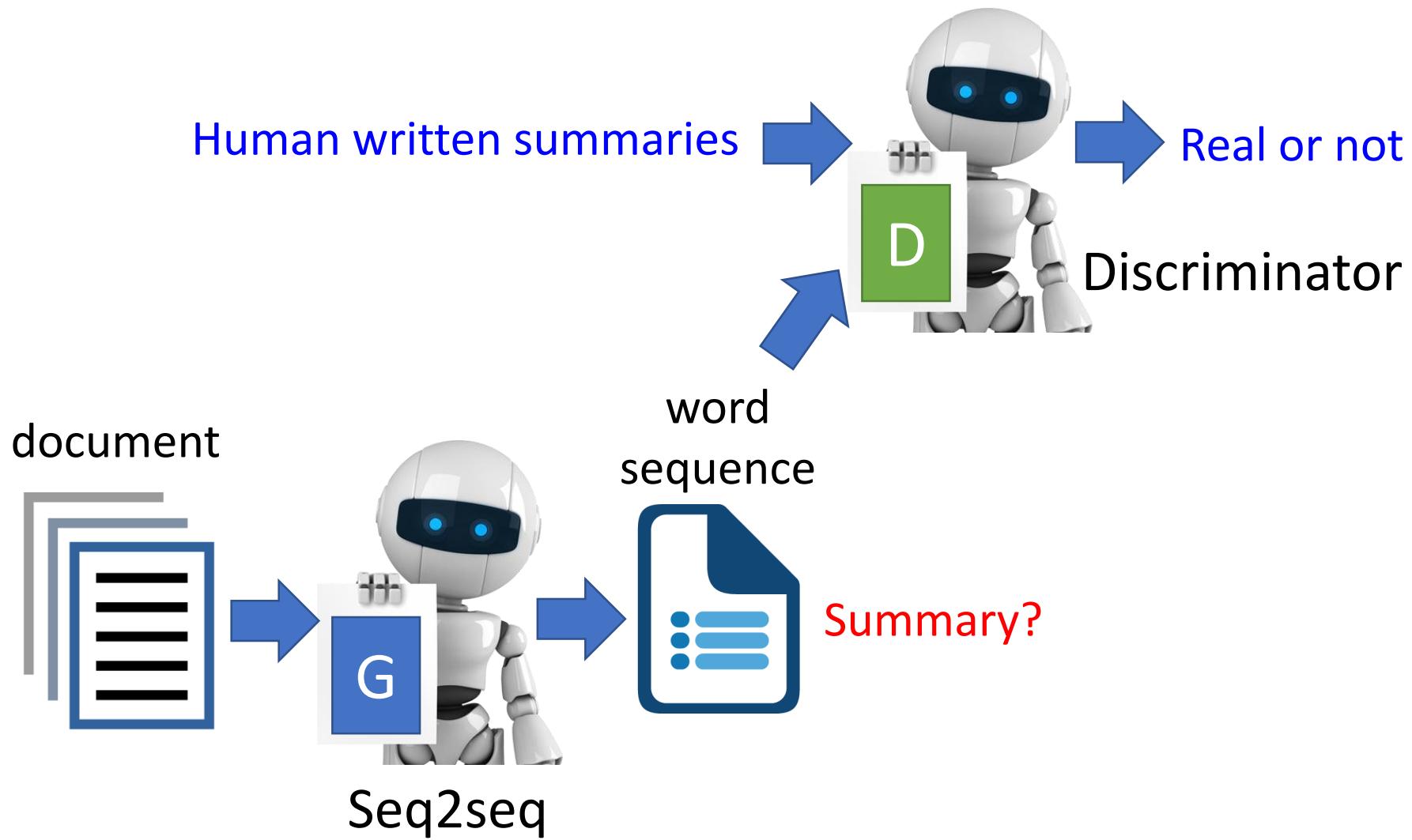
**Supervised: We need lots of  
labelled training data.**

Training Data

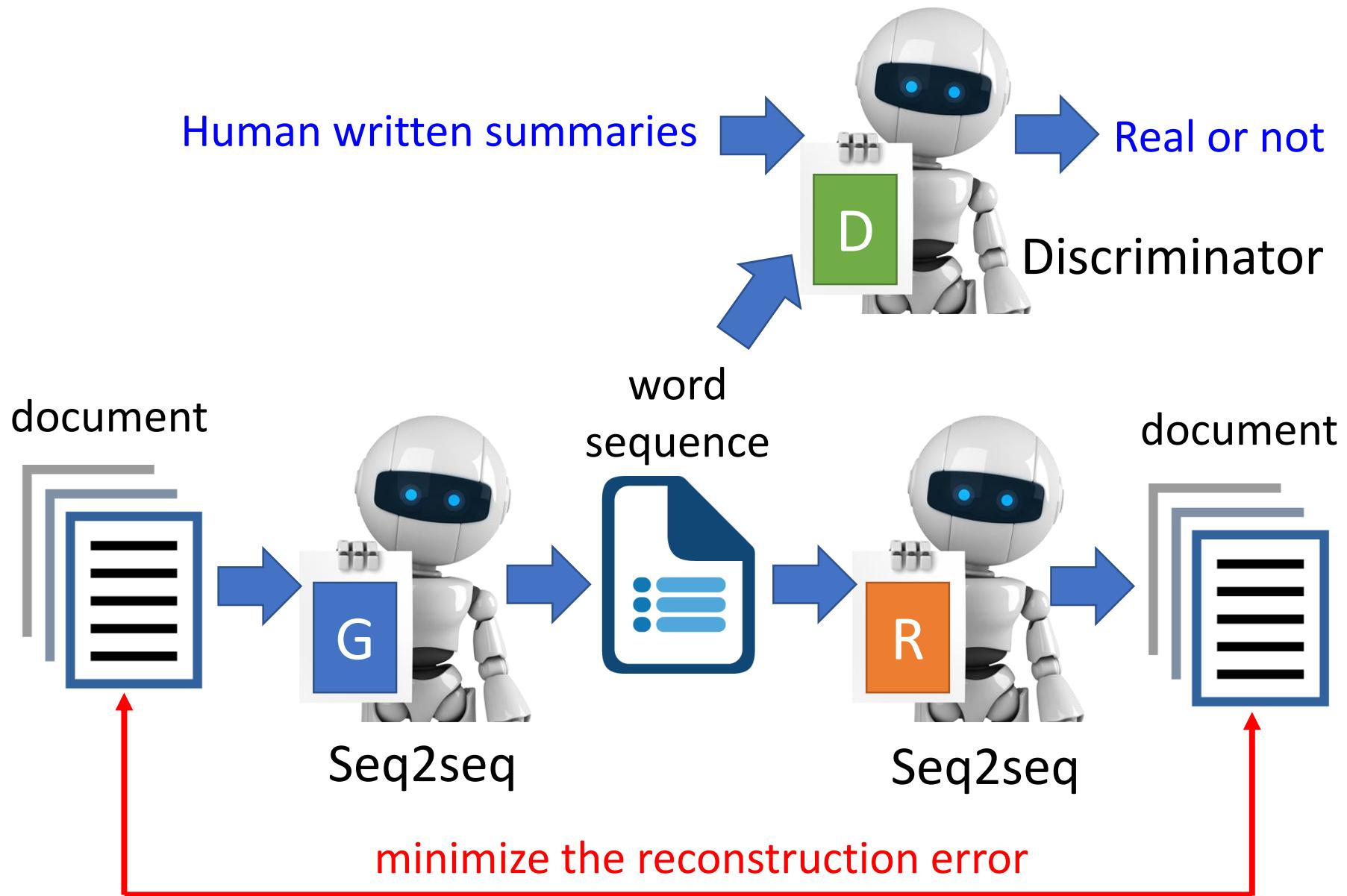
# Review: Unsupervised Conditional Generation



# Unsupervised Abstractive Summarization



# Unsupervised Abstractive Summarization



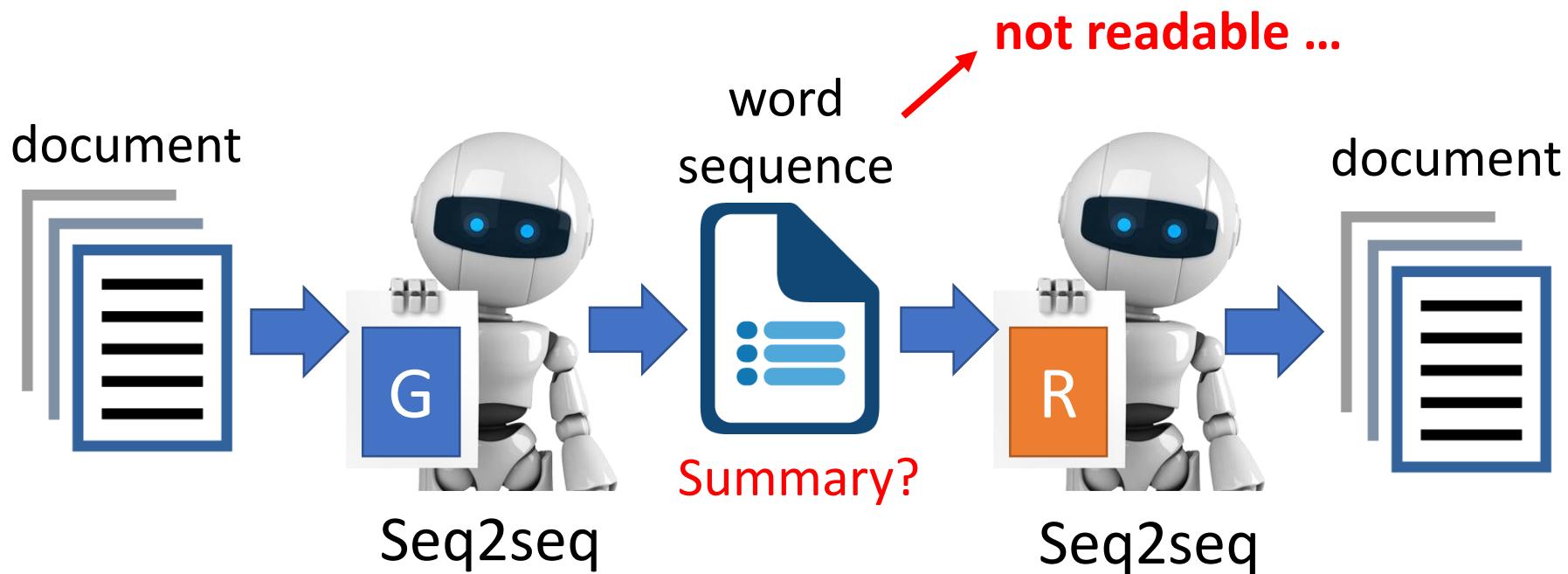
# Unsupervised Abstractive Summarization

Only need a lot  
of documents to  
train the model



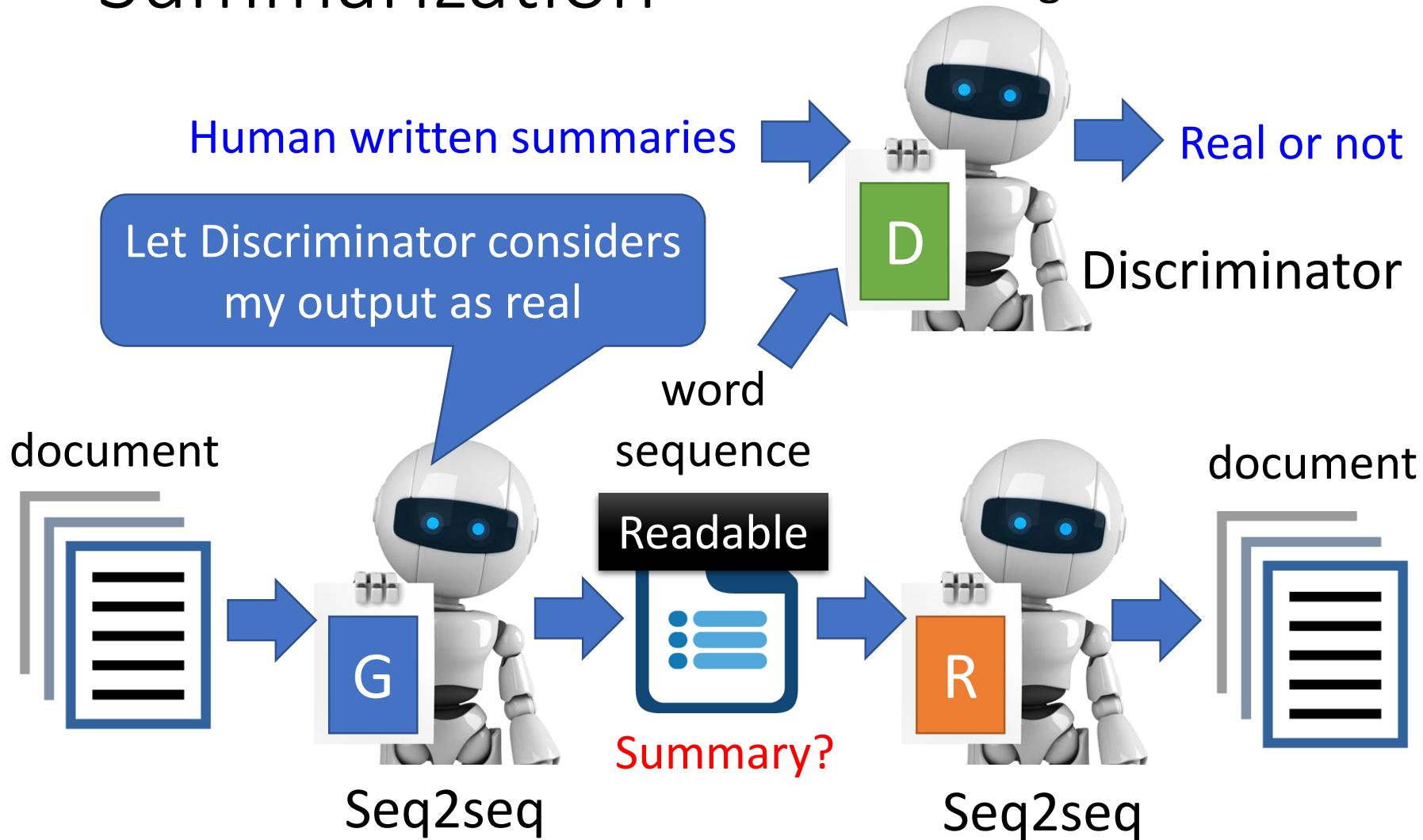
This is a ***seq2seq2seq auto-encoder***.

Using a sequence of words as latent representation.



# Unsupervised Abstractive Summarization

REINFORCE algorithm is used.



# Unsupervised Abstractive Summarization

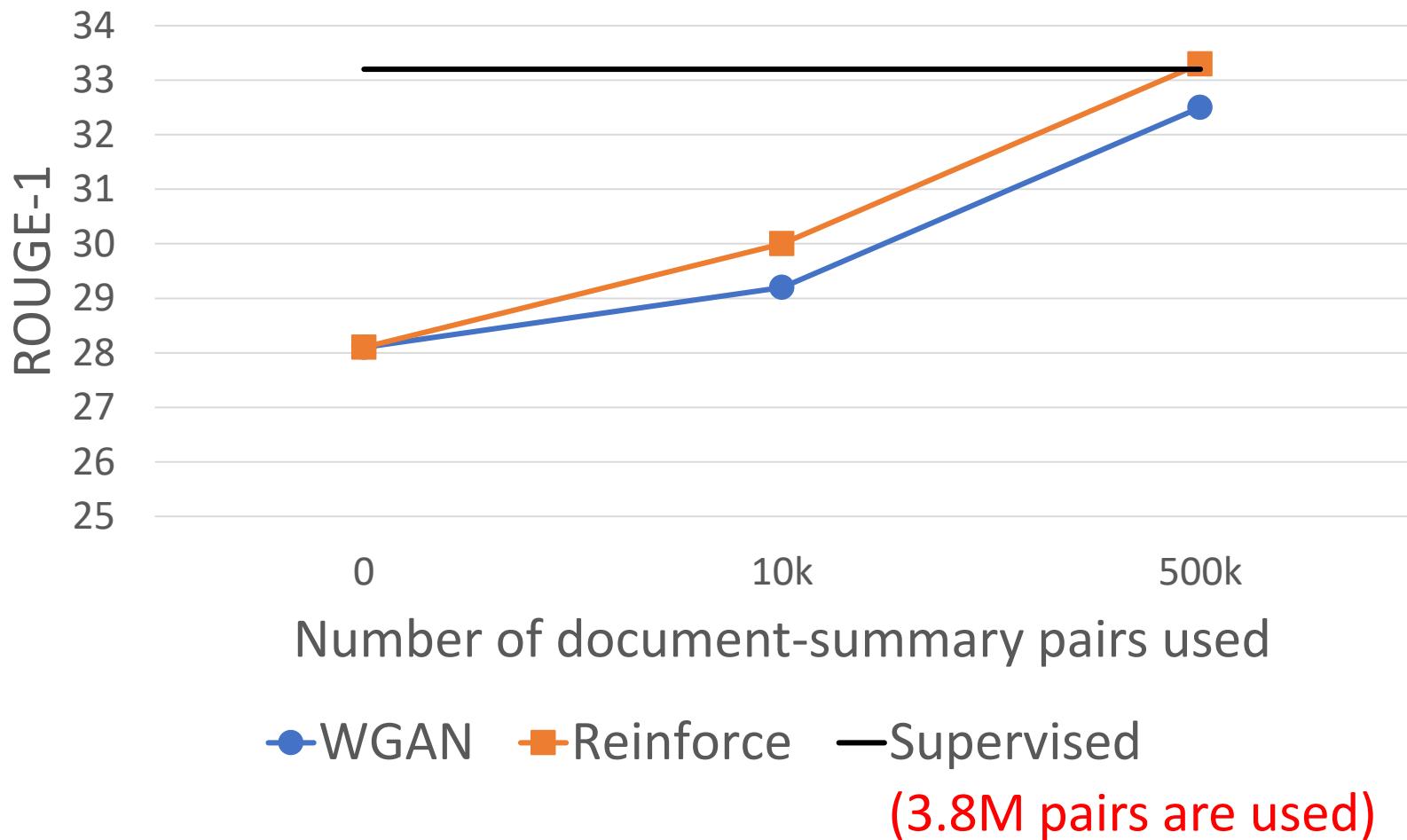
- **Document:** 澳大利亞今天與13個國家簽署了反興奮劑雙邊協議，旨在加強體育競賽之外的藥品檢查並共享研究成果 .....
- **Summary:**
  - **Human:** 澳大利亞與13國簽署反興奮劑協議
  - **Unsupervised:** 澳大利亞加強體育競賽之外的藥品檢查
- **Document:** 中華民國奧林匹克委員會今天接到一九九二年冬季奧運會邀請函，由於主席張豐緒目前正在中南美洲進行友好訪問，因此尚未決定是否派隊赴賽 .....
- **Summary:**
  - **Human:** 一九九二年冬季奧運會函邀我參加
  - **Unsupervised:** 奧委會接獲冬季奧運會邀請函

# Unsupervised Abstractive Summarization

- **Document:**據此間媒體27日報道,印度尼西亞蘇門答臘島的兩個省近日來連降暴雨,洪水泛濫導致塌方,到26日為止至少已有60人喪生,100多人失蹤 .....
- **Summary:**
  - **Human:**印尼水災造成60人死亡
  - **Unsupervised:**印尼門洪水泛濫導致塌雨
- **Document:**安徽省合肥市最近為領導幹部下基層做了新規定:一律輕車簡從,不準搞迎來送往、不準搞層層陪同 .....
- **Summary:**
  - **Human:**合肥規定領導幹部下基層活動從簡
  - **Unsupervised:**合肥領導幹部下基層做搞迎來送往規定:一律簡

# Semi-supervised Learning

Using  
matched data



# Outline of Part III

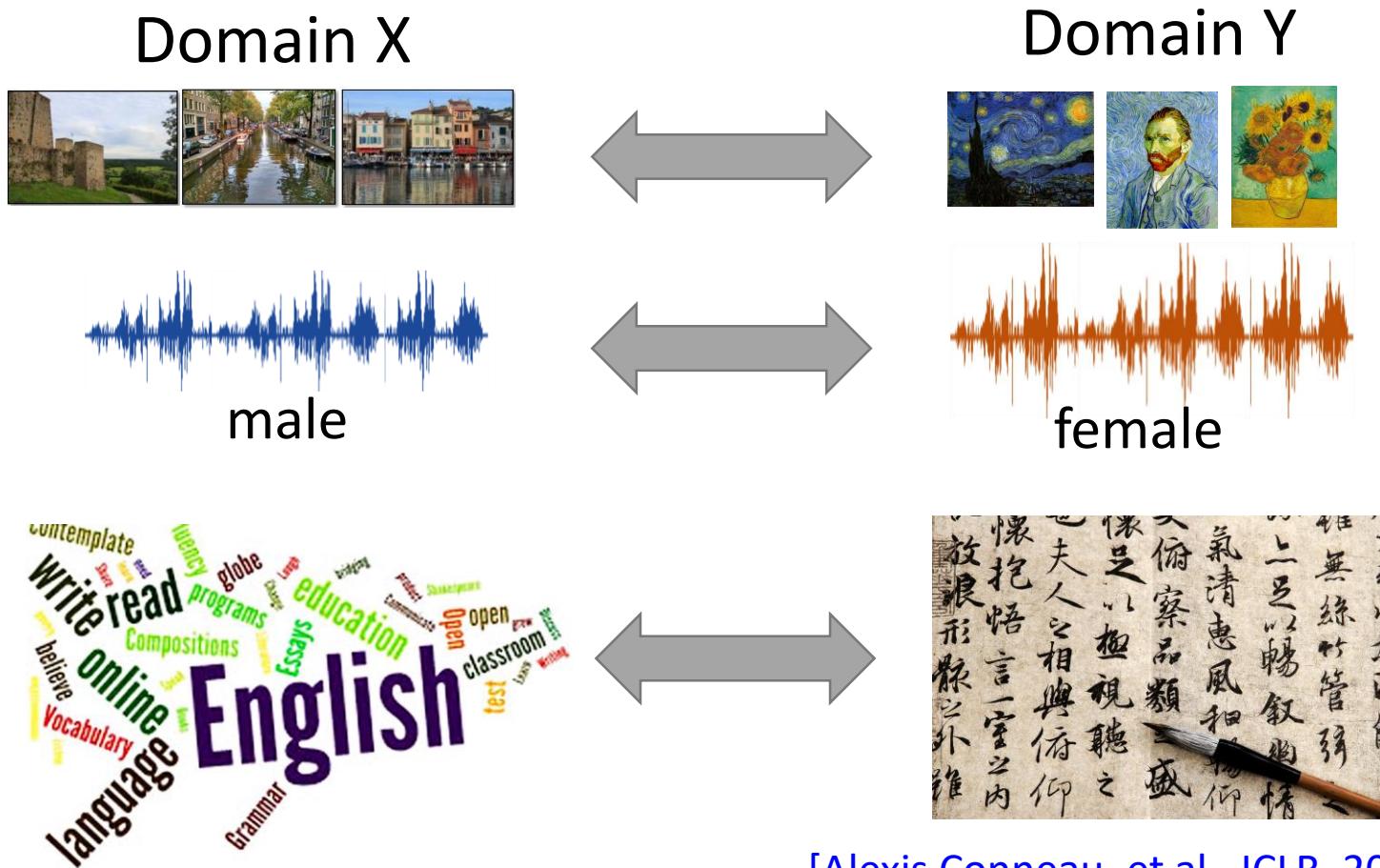
## Improving Supervised Seq-to-seq Model

- RL (human feedback)
- GAN (discriminator feedback)

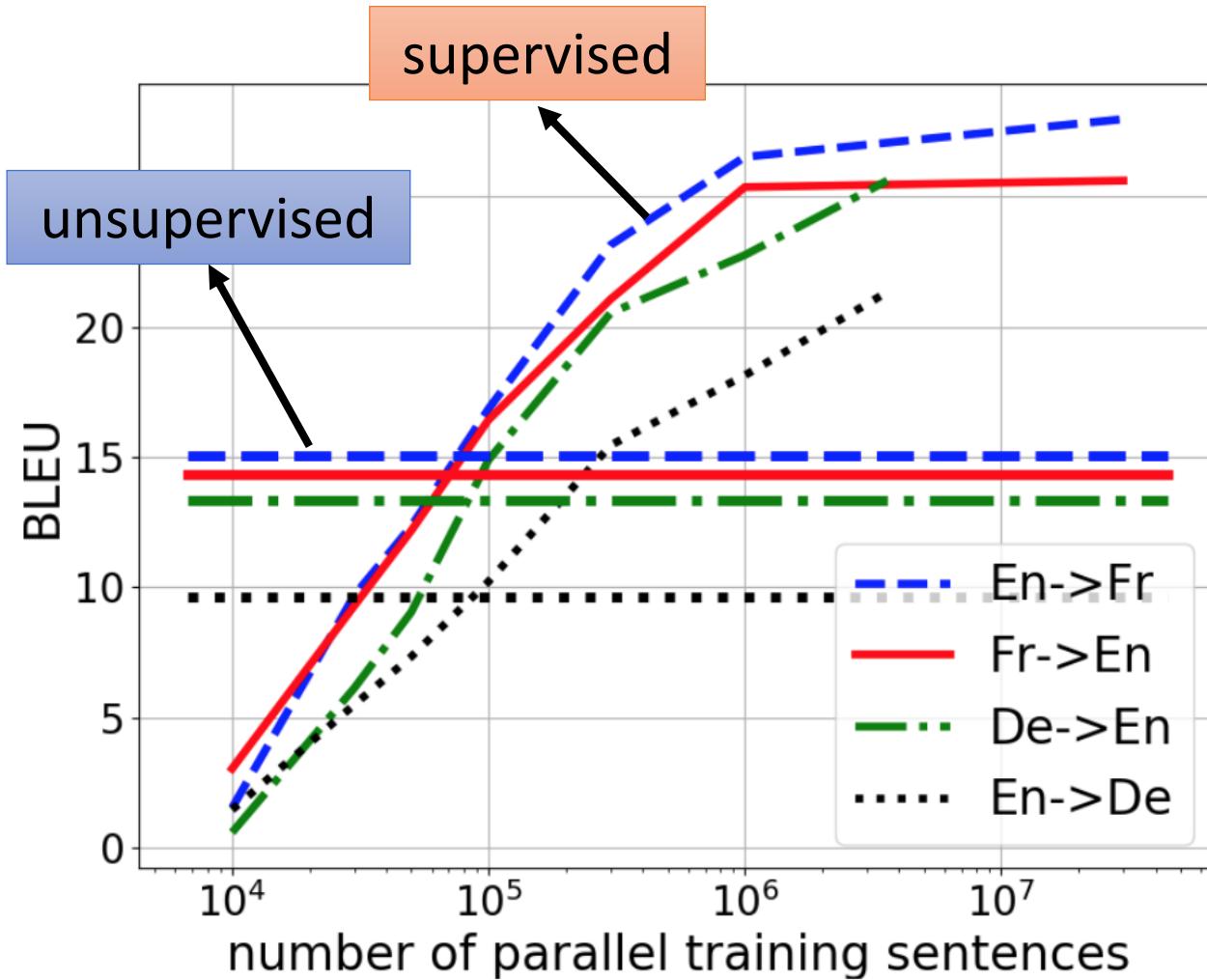
## Unsupervised Seq-to-seq Model

- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

# Unsupervised Machine Translation



[Alexis Conneau, et al., ICLR, 2018]  
[Guillaume Lample, et al., ICLR, 2018]

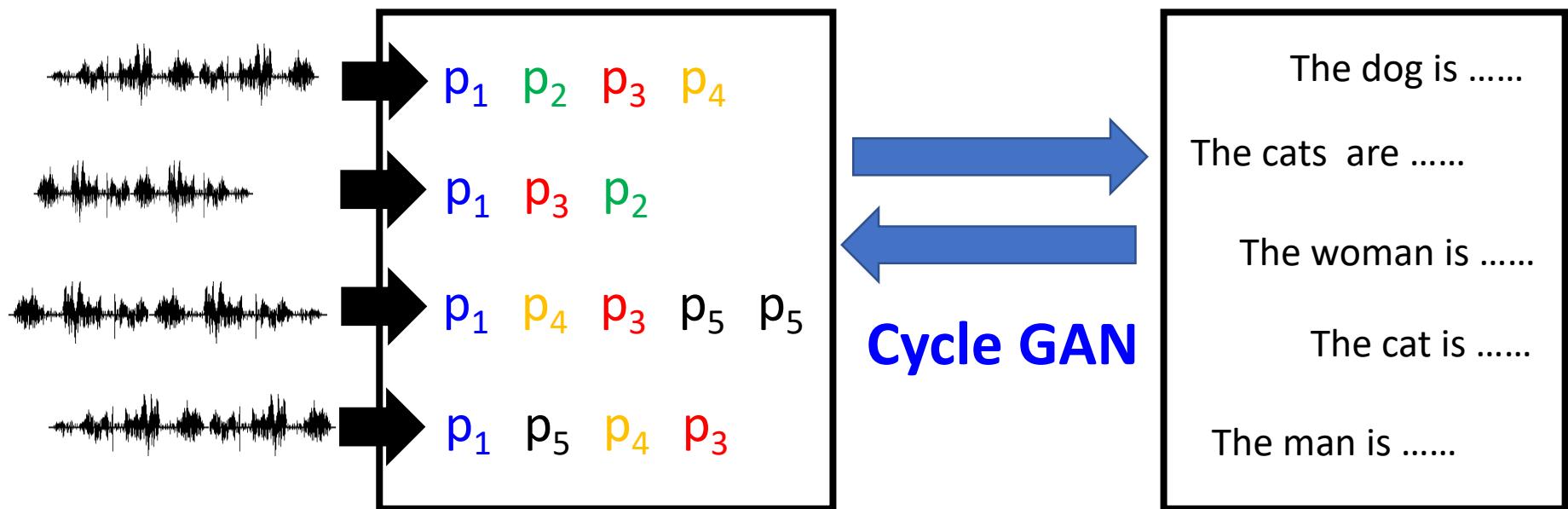


**Unsupervised learning  
with 10M sentences**

=

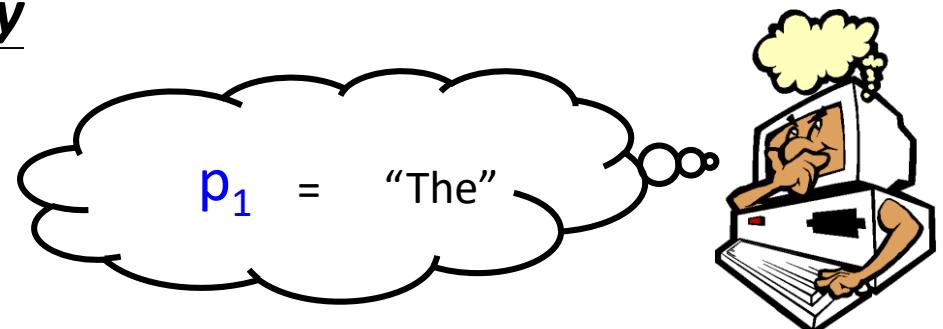
**Supervised learning with  
100K sentence pairs**

# Unsupervised Speech Recognition



## Acoustic Pattern Discovery

Can we achieve  
unsupervised speech  
recognition?

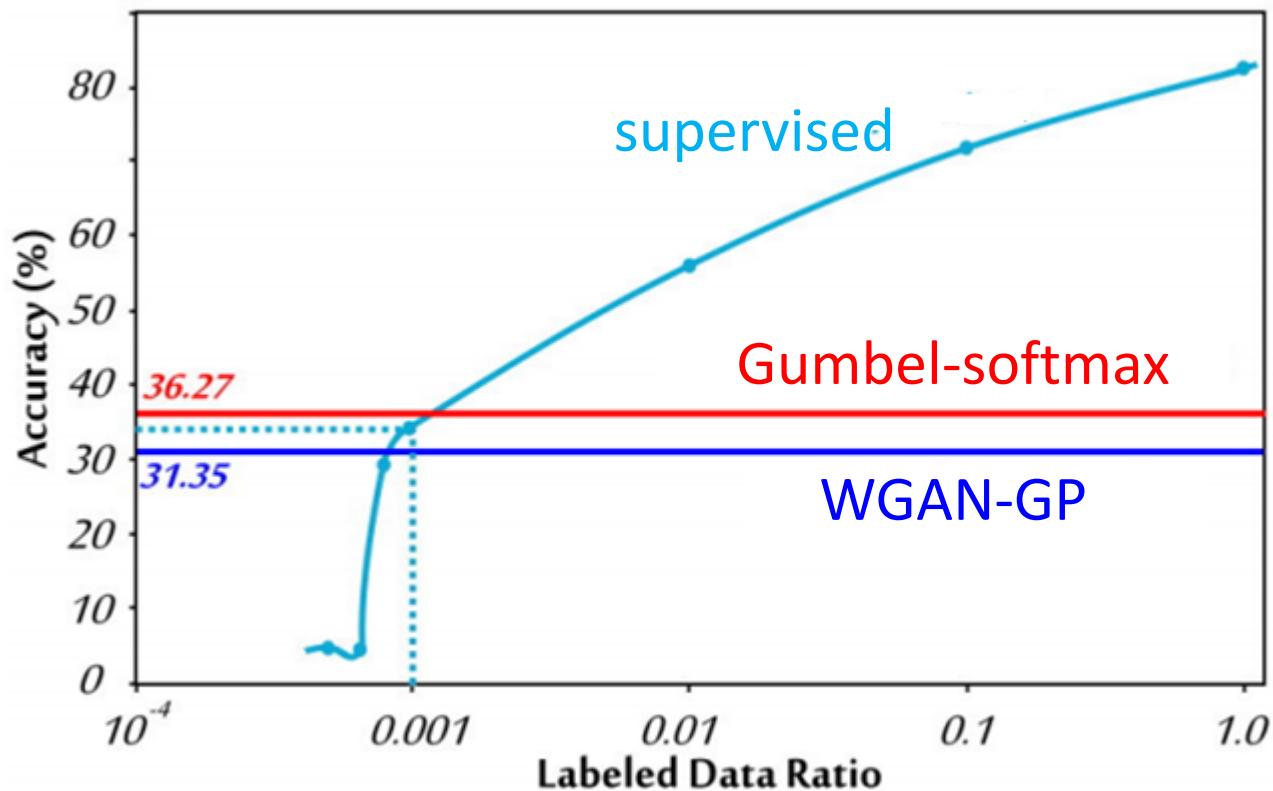


[Liu, et al., arXiv, 2018] [Chen, et al., arXiv, 2018]

# Unsupervised Speech Recognition

- Phoneme recognition

Audio: TIMIT  
Text: WMT



# Concluding Remarks

## Conditional Sequence Generation

- RL (human feedback)
- GAN (discriminator feedback)

## Unsupervised Conditional Sequence Generation

- Text Style Transfer
- Unsupervised Abstractive Summarization
- Unsupervised Translation

# Concluding Remarks

from A to Z

**A**

**B**

**C**

**D**

**E**

**F**

ACGAN

BiGAN

CycleGAN

DCGAN

EBGAN

fGAN

DuelGAN

**G**

**H**

**I**

**J**

**K**

**L**

GAN

?

InfoGAN

?

?

LSGAN

(only list those mentioned in class)

**M****N****O****P****Q****R**

MMGAN

NSGAN

?

Progressive  
GAN

?

Rank  
GAN**S****T****U****V****W****X**

StackGAN

Triple  
GANUnroll  
GAN

VAEGAN

WGAN

XGAN

StarGAN

SeqGAN

**Y**

?

**Z**

?

# Reference

- **Conditional Sequence Generation**
  - Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, Dan Jurafsky, Deep Reinforcement Learning for Dialogue Generation, EMNLP, 2016
  - Jiwei Li, Will Monroe, Tianlin Shi, Sébastien Jean, Alan Ritter, Dan Jurafsky, Adversarial Learning for Neural Dialogue Generation, EMNLP, 2017
  - Matt J. Kusner, José Miguel Hernández-Lobato, GANS for Sequences of Discrete Elements with the Gumbel-softmax Distribution, arXiv 2016
  - Tong Che, Yanran Li, Ruixiang Zhang, R Devon Hjelm, Wenjie Li, Yangqiu Song, Yoshua Bengio, Maximum-Likelihood Augmented Discrete Generative Adversarial Networks, arXiv 2017
  - Lantao Yu, Weinan Zhang, Jun Wang, Yong Yu, SeqGAN: Sequence Generative Adversarial Nets with Policy Gradient, AAAI 2017

# Reference

- **Conditional Sequence Generation**

- Sai Rajeswar, Sandeep Subramanian, Francis Dutil, Christopher Pal, Aaron Courville, Adversarial Generation of Natural Language, arXiv, 2017
- Ofir Press, Amir Bar, Ben Bogin, Jonathan Berant, Lior Wolf, Language Generation with Recurrent Generative Adversarial Networks without Pre-training, ICML workshop, 2017
- Zhen Xu, Bingquan Liu, Baoxun Wang, Chengjie Sun, Xiaolong Wang, Zhuoran Wang, Chao Qi , Neural Response Generation via GAN with an Approximate Embedding Layer, EMNLP, 2017
- Alex Lamb, Anirudh Goyal, Ying Zhang, Saizheng Zhang, Aaron Courville, Yoshua Bengio, Professor Forcing: A New Algorithm for Training Recurrent Networks, NIPS, 2016
- Yizhe Zhang, Zhe Gan, Kai Fan, Zhi Chen, Ricardo Henao, Dinghan Shen, Lawrence Carin, Adversarial Feature Matching for Text Generation, ICML, 2017
- Jiaxian Guo, Sidi Lu, Han Cai, Weinan Zhang, Yong Yu, Jun Wang, Long Text Generation via Adversarial Training with Leaked Information, AAAI, 2018
- Kevin Lin, Dianqi Li, Xiaodong He, Zhengyou Zhang, Ming-Ting Sun, Adversarial Ranking for Language Generation, NIPS, 2017
- William Fedus, Ian Goodfellow, Andrew M. Dai, MaskGAN: Better Text Generation via Filling in the \_\_\_\_\_, ICLR, 2018

# Reference

- **Conditional Sequence Generation**
  - Sidi Lu, Yaoming Zhu, Weinan Zhang, Jun Wang, Yong Yu, Neural Text Generation: Past, Present and Beyond, arXiv, 2018
  - Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, Yong Yu, Texygen: A Benchmarking Platform for Text Generation Models, arXiv, 2018
  - Zhen Yang, Wei Chen, Feng Wang, Bo Xu, Improving Neural Machine Translation with Conditional Sequence Generative Adversarial Nets, NAACL, 2018
  - Lijun Wu, Yingce Xia, Li Zhao, Fei Tian, Tao Qin, Jianhuang Lai, Tie-Yan Liu, Adversarial Neural Machine Translation, arXiv 2017
  - Linqing Liu, Yao Lu, Min Yang, Qiang Qu, Jia Zhu, Hongyan Li, Generative Adversarial Network for Abstractive Text Summarization, AAAI 2018
  - Rakshith Shetty, Marcus Rohrbach, Lisa Anne Hendricks, Mario Fritz, Bernt Schiele, Speaking the Same Language: Matching Machine to Human Captions by Adversarial Training, ICCV 2017
  - Xiaodan Liang, Zhiting Hu, Hao Zhang, Chuang Gan, Eric P. Xing, Recurrent Topic-Transition GAN for Visual Paragraph Generation, arXiv 2017

# Reference

- **Text Style Transfer**
  - Zhenxin Fu, Xiaoye Tan, Nanyun Peng, Dongyan Zhao, Rui Yan, Style Transfer in Text: Exploration and Evaluation, AAAI, 2018
  - Tianxiao Shen, Tao Lei, Regina Barzilay, and Tommi Jaakkola, Style Transfer from Non-Parallel Text by Cross-Alignment, NIPS 2017
  - Chih-Wei Lee, Yau-Shian Wang, Tsung-Yuan Hsu, Kuan-Yu Chen, Hung-Yi Lee, Lin-shan Lee, Scalable Sentiment for Sequence-to-sequence Chatbot Response with Performance Analysis, ICASSP, 2018
  - Junbo (Jake) Zhao, Yoon Kim, Kelly Zhang, Alexander M. Rush, Yann LeCun, Adversarially Regularized Autoencoders, arxiv, 2017

# Reference

- **Unsupervised Machine Translation**
  - Alexis Conneau, Guillaume Lample, Marc'Aurelio Ranzato, Ludovic Denoyer, Hervé Jégou, Word Translation Without Parallel Data, ICRL 2018
  - Guillaume Lample, Ludovic Denoyer, Marc'Aurelio Ranzato, Unsupervised Machine Translation Using Monolingual Corpora Only, ICRL 2018
- **Unsupervised Speech Recognition**
  - Da-Rong Liu, Kuan-Yu Chen, Hung-Yi Lee, Lin-shan Lee, Completely Unsupervised Phoneme Recognition by Adversarially Learning Mapping Relationships from Audio Embeddings, arXiv, 2018
  - Yi-Chen Chen, Chia-Hao Shen, Sung-Feng Huang, Hung-yi Lee, Towards Unsupervised Automatic Speech Recognition Trained by Unaligned Speech and Text only, arXiv, 2018

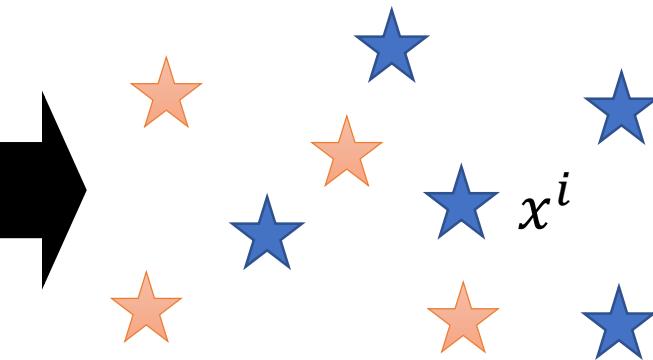
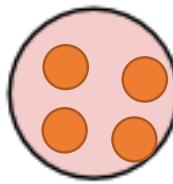
# Evaluation

Ref: Lucas Theis, Aäron van den Oord, Matthias Bethge, “A note on the evaluation of generative models”, arXiv preprint, 2015

# Likelihood

- ★ : real data (not observed during training)
- ★ : generated data

Prior  
Distribution



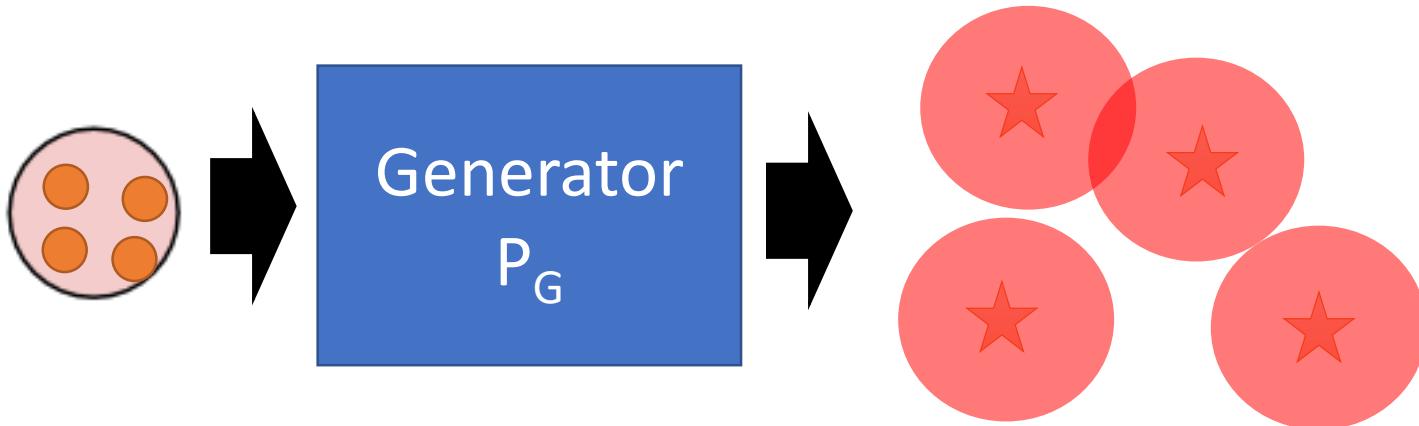
$$\text{Log Likelihood: } L = \frac{1}{N} \sum_i \log P_G(x^i)$$

We cannot compute  $P_G(x^i)$ . We can only sample from  $P_G$ .

# Likelihood

## - Kernel Density Estimation

- Estimate the distribution of  $P_G(x)$  from sampling



Each sample is the mean of a Gaussian with the same covariance.

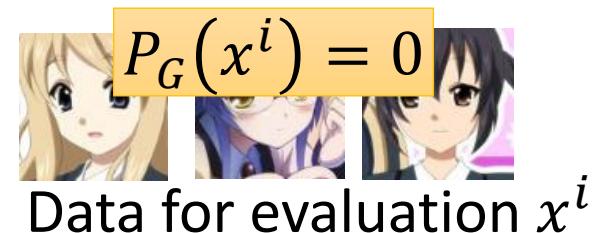
Now we have an approximation of  $P_G$ , so we can compute  $P_G(x^i)$  for each real data  $x^i$

Then we can compute the likelihood.

# Likelihood v.s. Quality

- Low likelihood, high quality?

Considering a model generating good images (small variance)



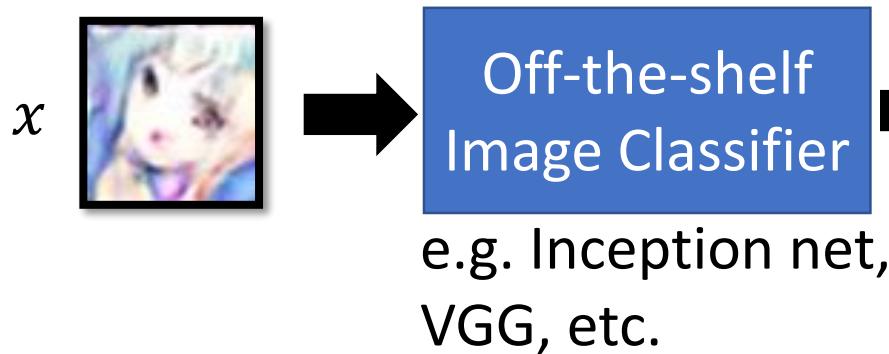
- High likelihood, low quality?



$$L = \frac{1}{N} \sum_i \log \frac{P_G(x^i)}{100} = -\log 100 + \frac{1}{N} \sum_i \log P_G(x^i)$$

4.6

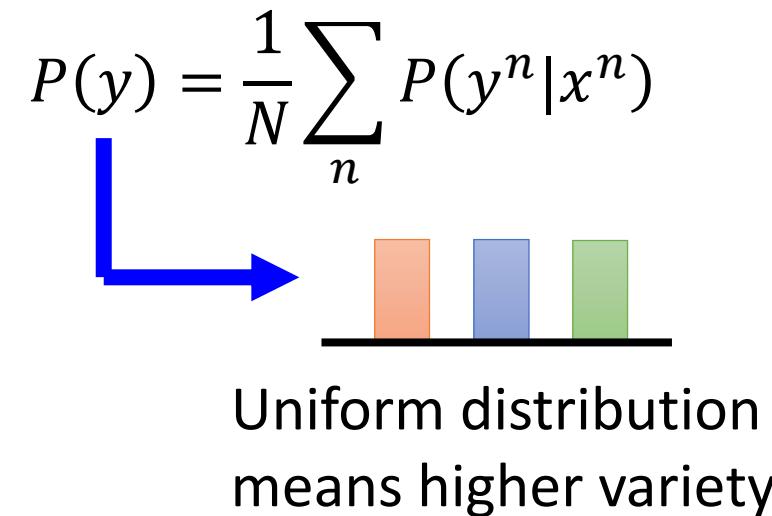
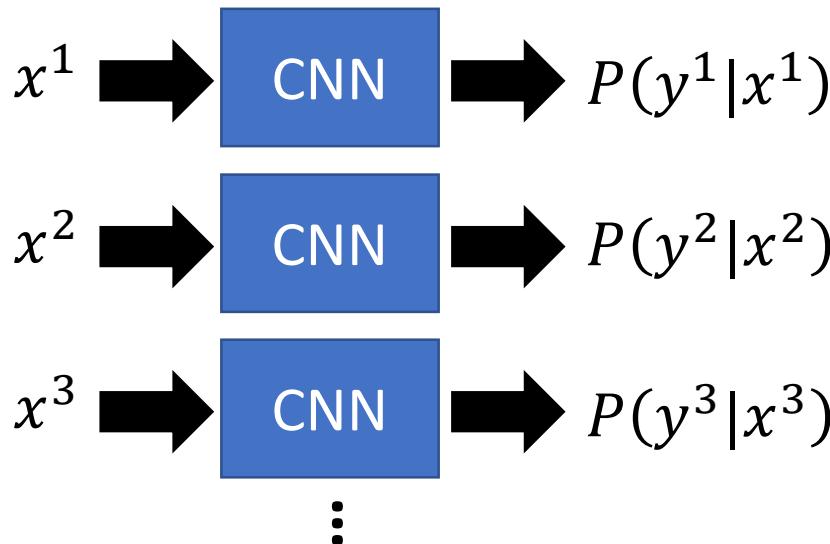
# Objective Evaluation



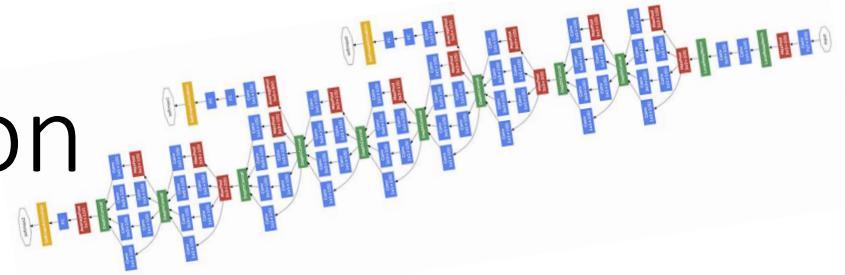
$x$ : image  
 $y$ : class (output of CNN)



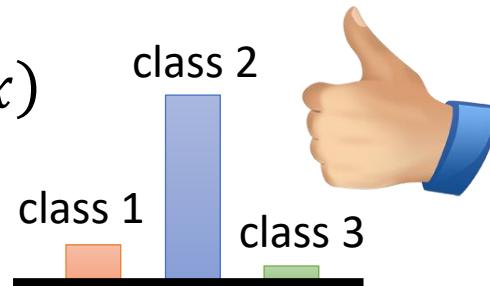
Concentrated distribution means higher visual quality



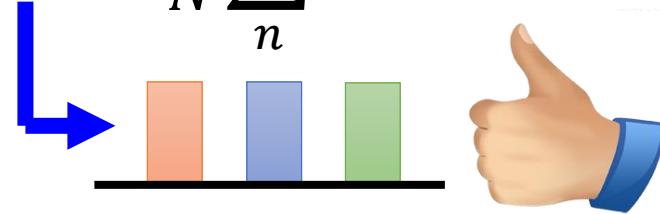
# Objective Evaluation



$$P(y|x)$$



$$P(y) = \frac{1}{N} \sum_n P(y^n|x^n)$$



## Inception Score

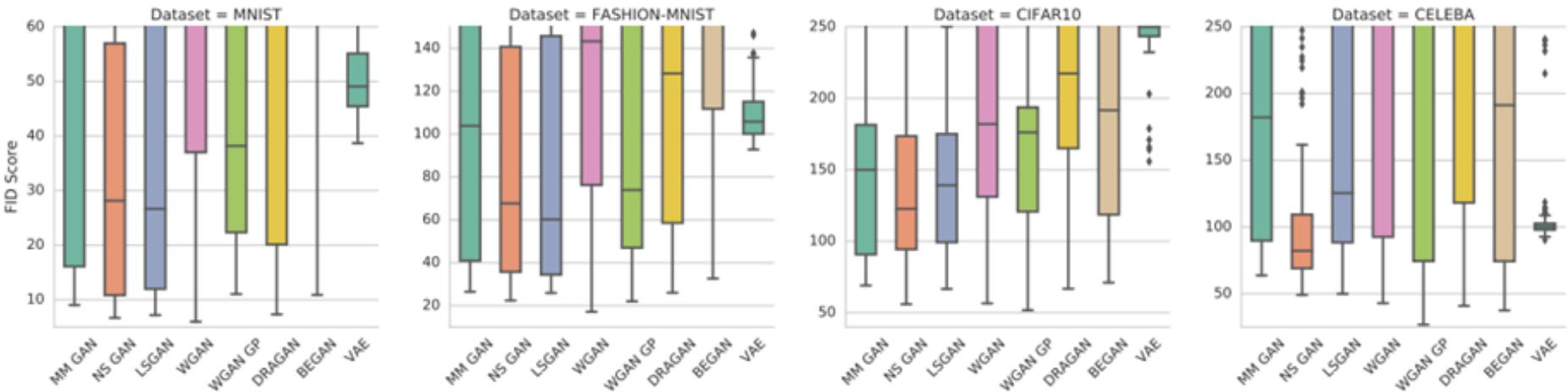
$$= \sum_x \sum_y P(y|x) \log P(y|x)$$

Negative entropy of  $P(y|x)$

$$- \sum_y P(y) \log P(y)$$

Entropy of  $P(y)$

GAN	DISCRIMINATOR LOSS	GENERATOR LOSS
MM GAN	$\mathcal{L}_D^{GAN} = -\mathbb{E}_{x \sim p_d} [\log(D(x))] + \mathbb{E}_{\hat{x} \sim p_g} [\log(1 - D(\hat{x}))]$	$\mathcal{L}_G^{GAN} = -\mathcal{L}_D^{GAN}$
NS GAN	$\mathcal{L}_D^{NSGAN} = \mathcal{L}_D^{GAN}$	$\mathcal{L}_G^{NSGAN} = \mathbb{E}_{\hat{x} \sim p_g} [\log(D(\hat{x}))]$
WGAN	$\mathcal{L}_D^{WGAN} = -\mathbb{E}_{x \sim p_d} [D(x)] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$	$\mathcal{L}_G^{WGAN} = -\mathcal{L}_D^{WGAN}$
WGAN GP	$\mathcal{L}_D^{WGAN} = \mathcal{L}_D^{WGAN} + \lambda \mathbb{E}_{\hat{x} \sim p_g} [(  \nabla D(\alpha x + (1 - \alpha)\hat{x})  _2 - 1)^2]$	$\mathcal{L}_G^{WGAN} = -\mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})]$
LS GAN	$\mathcal{L}_D^{LSGAN} = -\mathbb{E}_{x \sim p_d} [(D(x) - 1)^2] + \mathbb{E}_{\hat{x} \sim p_g} [D(\hat{x})^2]$	$\mathcal{L}_G^{LSGAN} = -\mathbb{E}_{\hat{x} \sim p_g} [(D(\hat{x}) - 1)^2]$
DRAGAN	$\mathcal{L}_D^{DRAGAN} = \mathcal{L}_D^{GAN} + \lambda \mathbb{E}_{\hat{x} \sim p_d + \mathcal{N}(0, c)} [(  \nabla D(\hat{x})  _2 - 1)^2]$	$\mathcal{L}_G^{DRAGAN} = -\mathcal{L}_D^{NSGAN}$
BEGAN	$\mathcal{L}_D^{BEGAN} = \mathbb{E}_{x \sim p_d} [      x - AE(x)      _1] - k_t \mathbb{E}_{\hat{x} \sim p_g} [      \hat{x} - AE(\hat{x})      _1]$	$\mathcal{L}_G^{BEGAN} = \mathbb{E}_{\hat{x} \sim p_g} [      \hat{x} - AE(\hat{x})      _1]$



Smaller is better

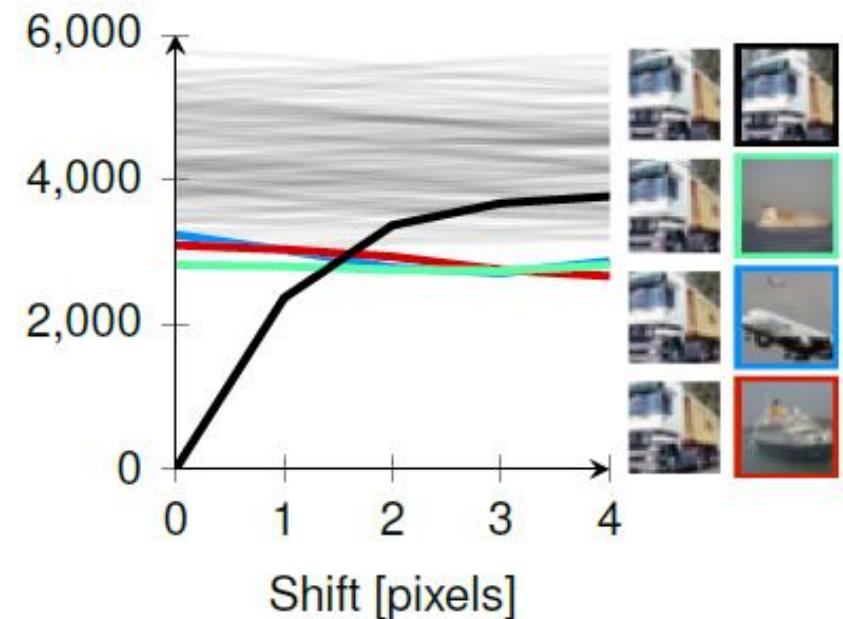
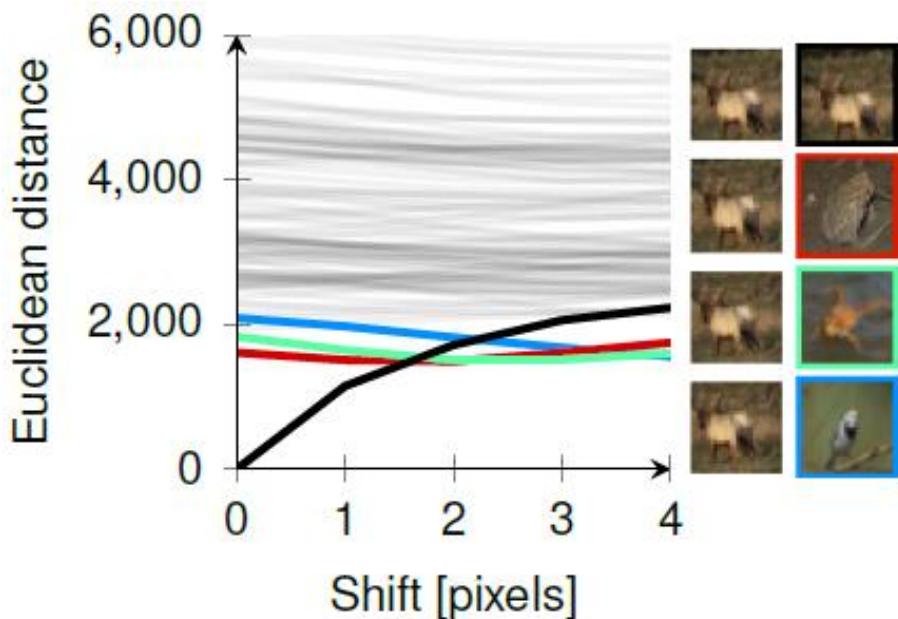
FIT:

<https://arxiv.org/pdf/1706.08500.pdf>

Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, Olivier Bousquet, "Are GANs Created Equal? A Large-Scale Study", arXiv, 2017

# We don't want memory GAN.

- Using k-nearest neighbor to check whether the generator generates new objects



# Missing Mode ?

Mode collapse is  
easy to detect.



?

