



# Conditional GANs, Image-to-Image Translation

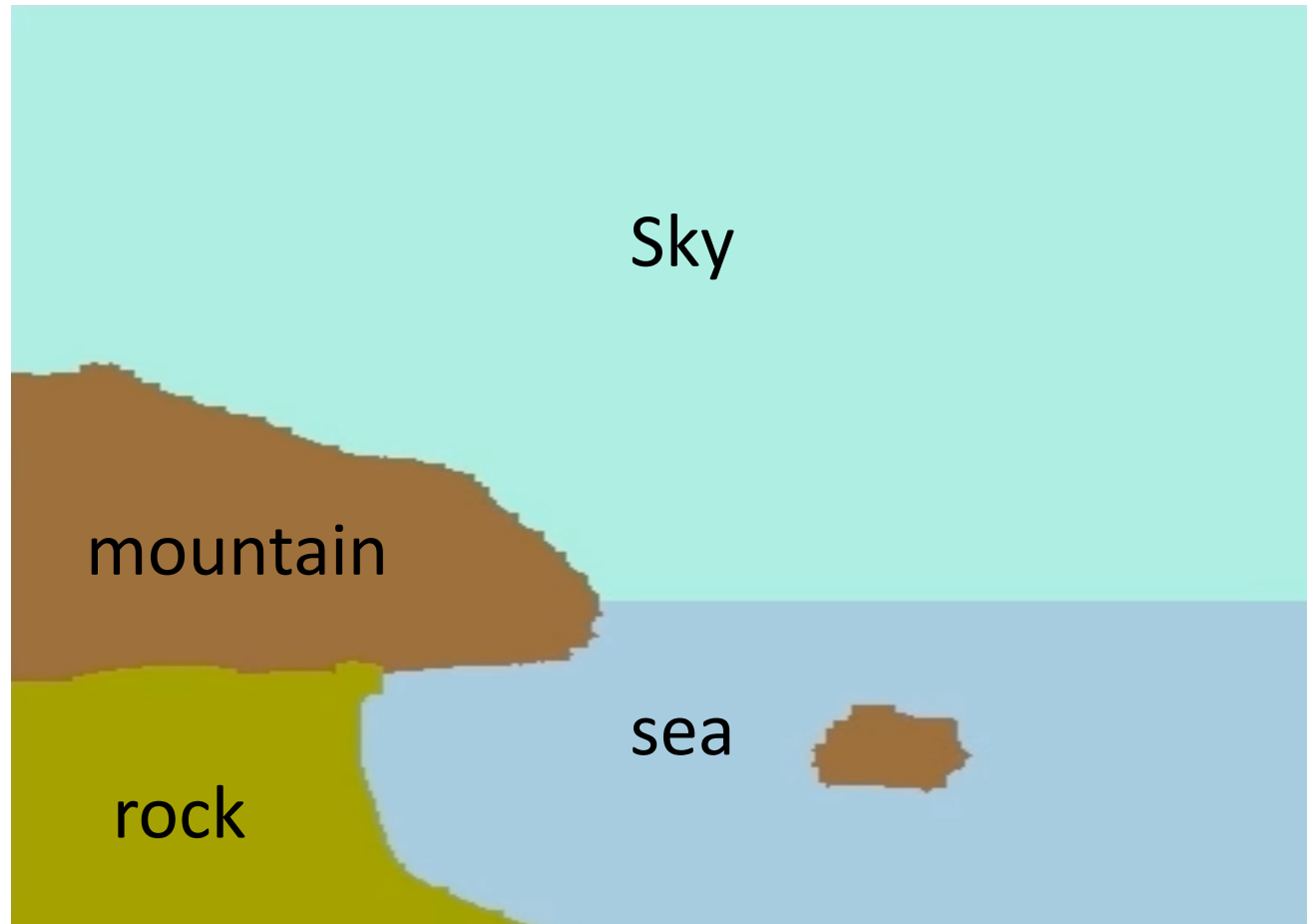
Jun-Yan Zhu

16-726, Spring 2022

# Logistics

- HW2 gather town party Mon 8-10 pm
- No class next week (due to Spring break)
- HW 1 Class Choice Award:
  - Vote by the end of Wed.
  - Winner will be announced on 03/14 (Mon)

# Problem Statement

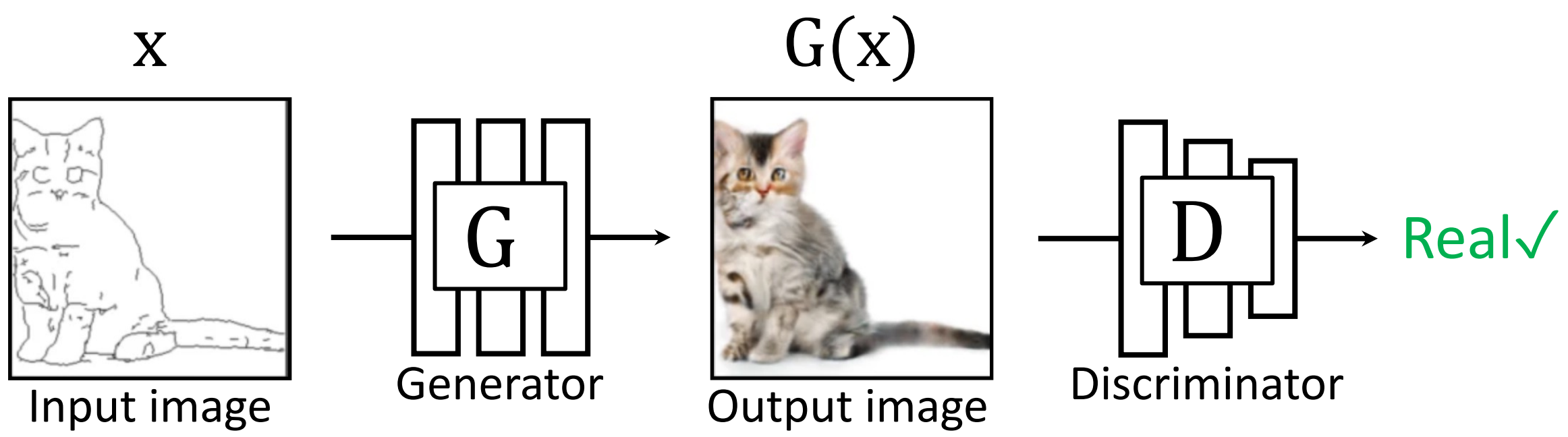


Input



Output

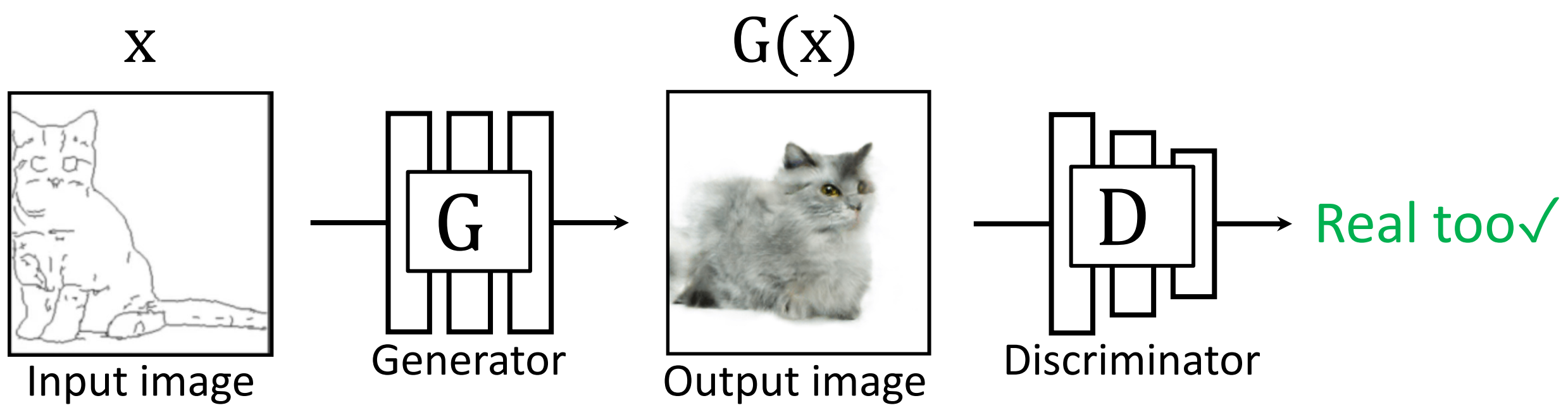
Goal: synthesize a realistic photograph given an input image



Learning objective

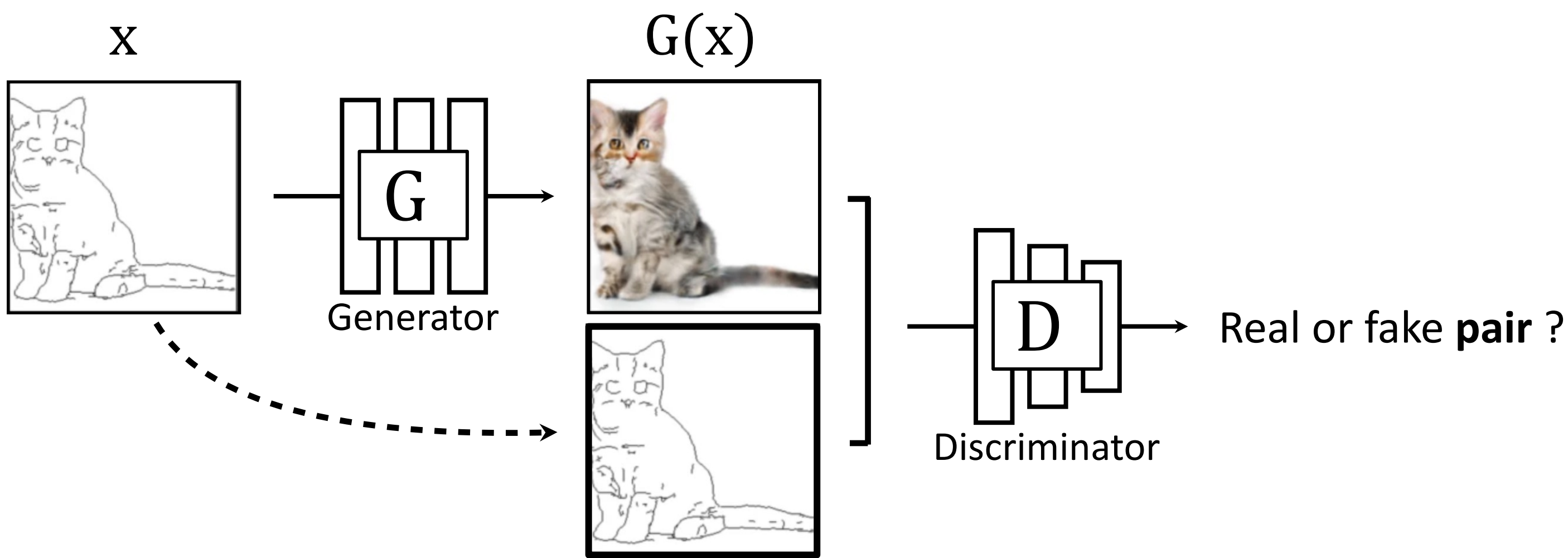
$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$





Learning objective

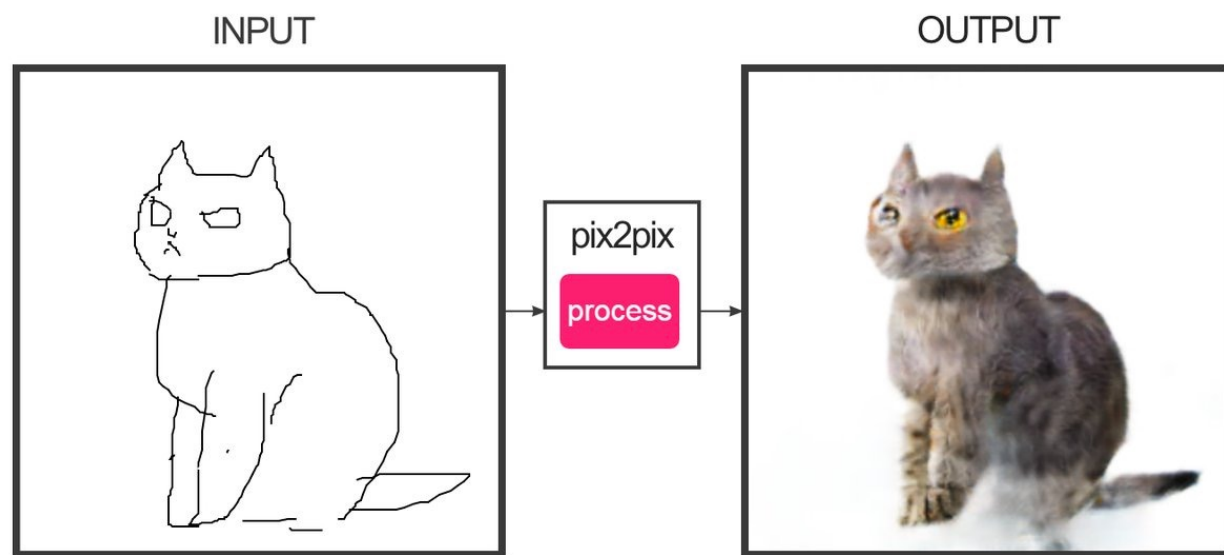
$$\min_G \max_D \mathbb{E}_x [\log(1 - D(G(x)))] + \mathbb{E}_y [\log D(y)]$$



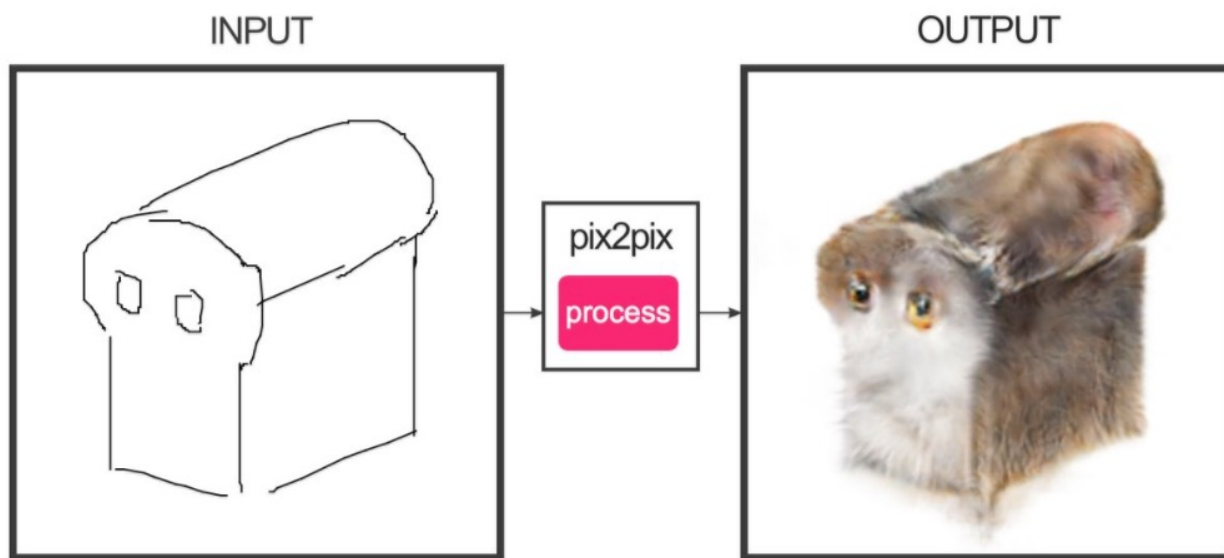
Learning objective

$$\min_G \max_D \mathbb{E}_x [\log(1 - D(x, G(x)))] + \mathbb{E}_{x,y} [\log D(x, y)]$$

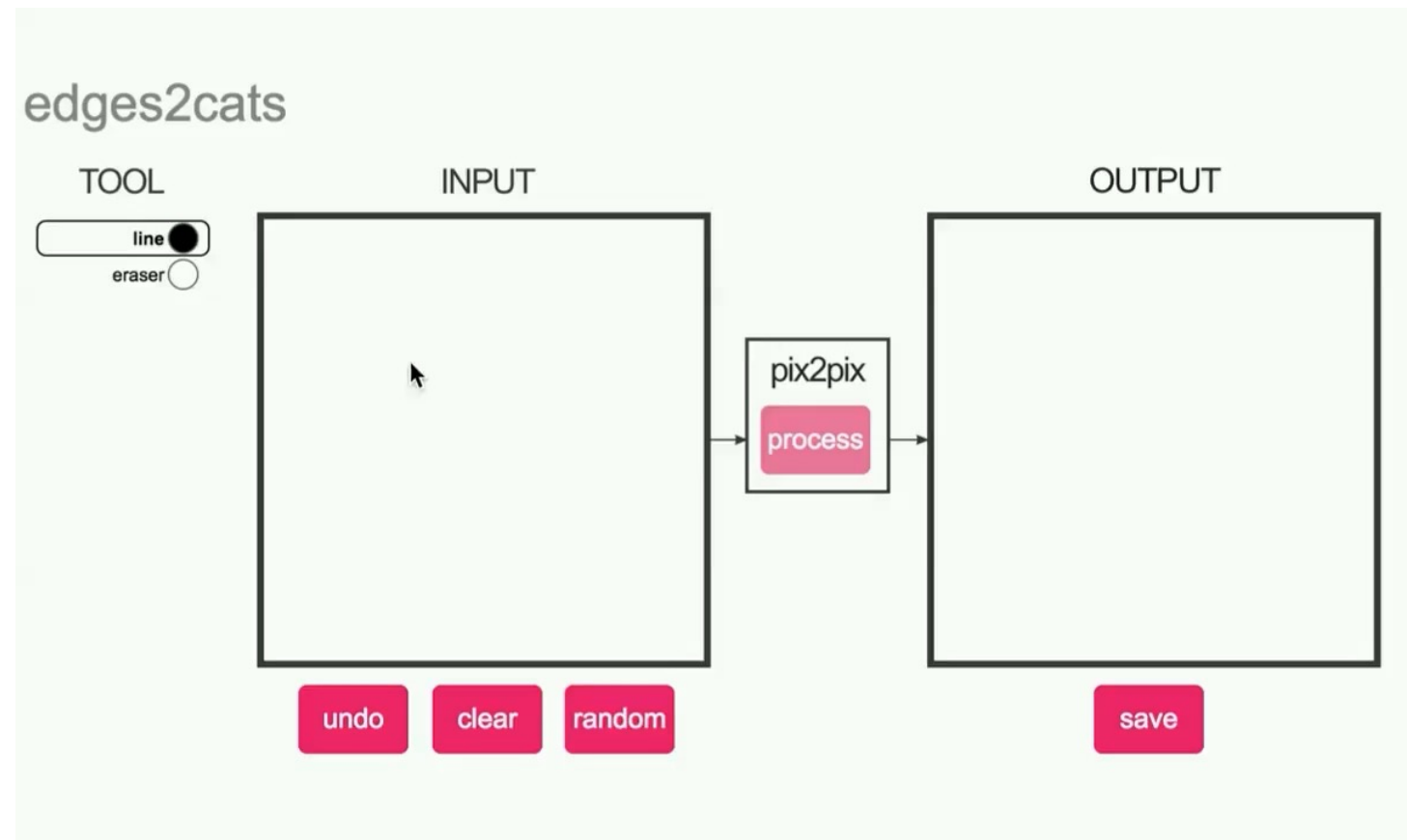
# #edges2cats [Christopher Hesse]



@gods\_tail



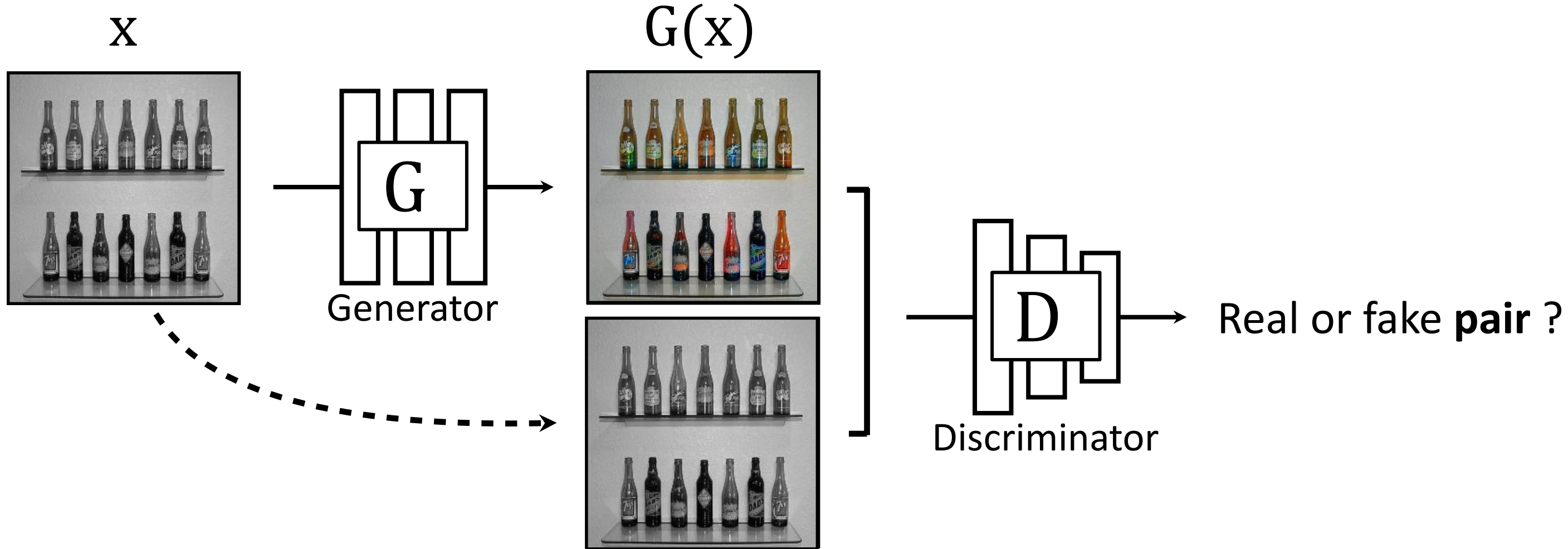
Ivy Tasi @ivymyt



@matthematician



Vitaly Vidmirov @vvid



Input: ~~Grayscale~~ **Grayscale**  $\rightarrow$  Output: ~~Color~~ **Color**



# Automatic Colorization with pix2pix

Input

Output

Input

Output

Input

Output

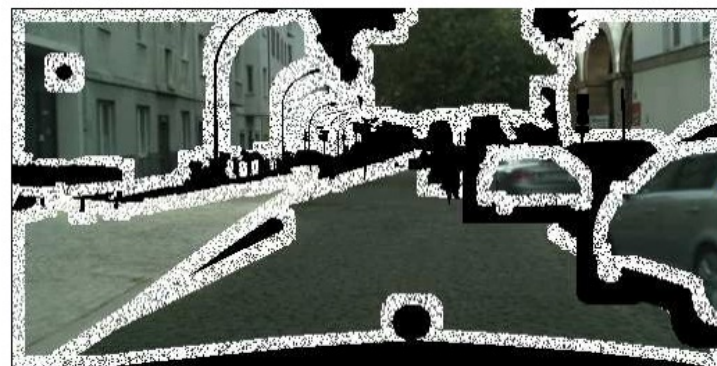


Learning vs. Exemplar-based

# Hybrid Method



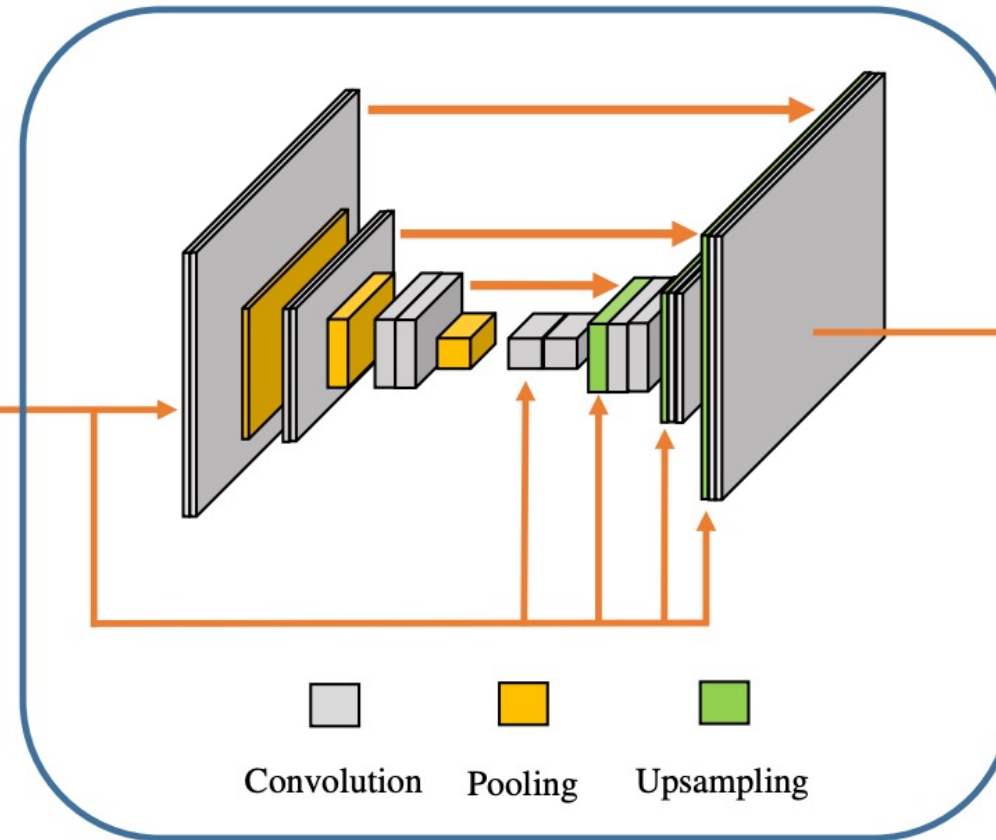
Semantic layout



Canvas

Output from

exemplar-based method



Convolution Pooling Upsampling

Synthesis network  $f$



Output

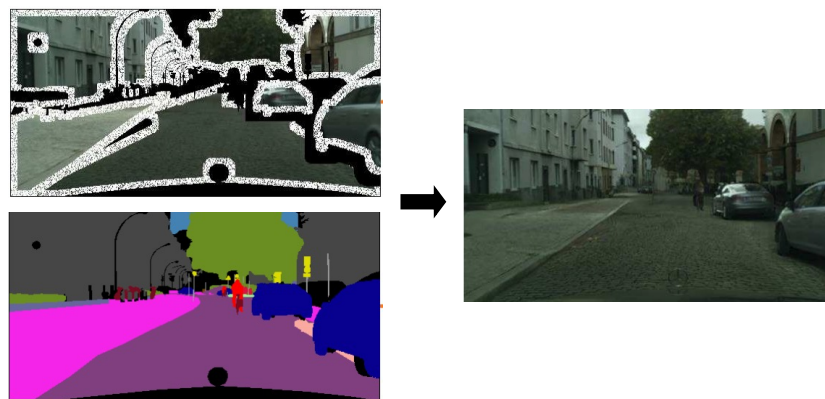


## Learning-based



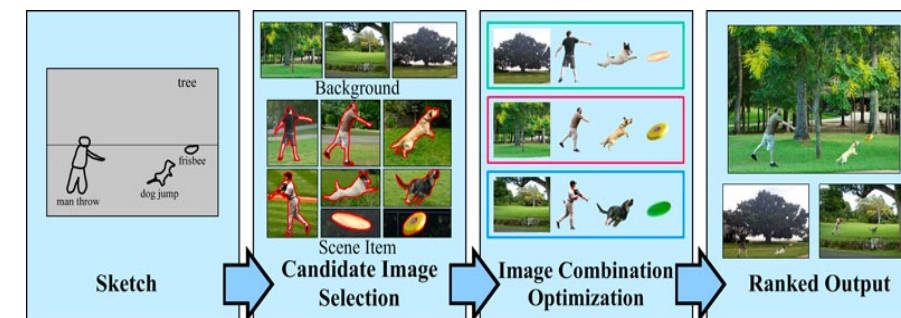
[Isola et al], [Wang et al]  
[Park et al], SEAN [Zhu et al]

## Hybrid method



SIMS [Qi et al]

## Exemplar-based



[Johnson et al], [Lalonde et al]  
[Tao et al], [Bansal et al]

Speed



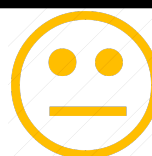
Local realism



Global realism



Match Input



# Discussion

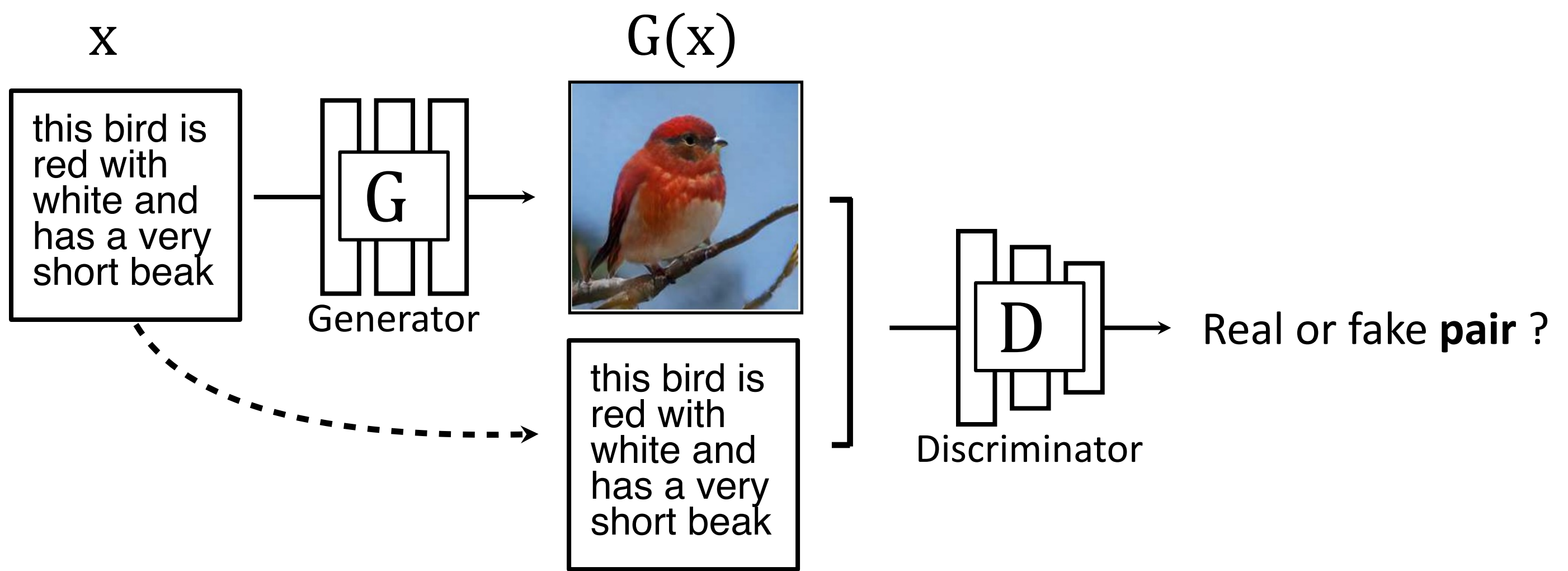
## Summary

- Intuitive user inputs.
- Realistic outputs.
- Used by visual artists.



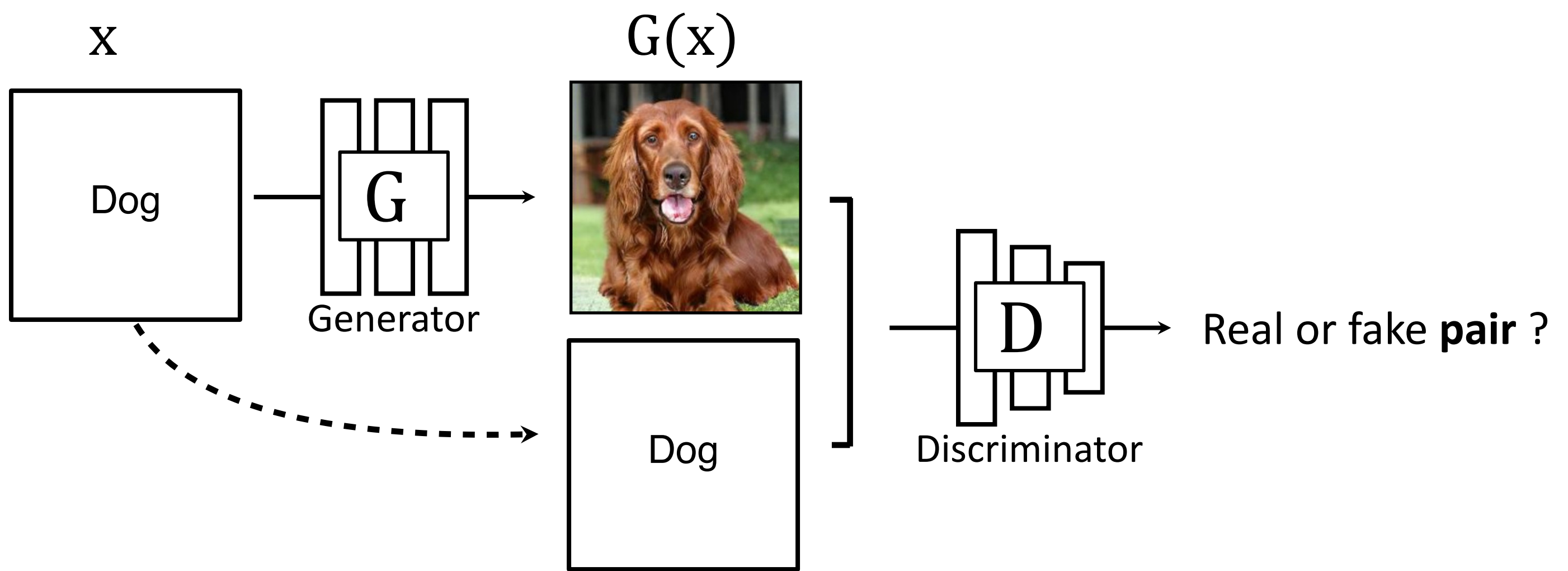
## Challenges

- Fine-grained controls (texture, 3D, and lighting).
- High-resolution output (4K).
- Model efficiency on mobile devices.
- Video Control.



Input: **Text** → Output: **Photo**

Text-to-Image Synthesis



Input: **Class** → Output: **Photo**

## Class-conditional GANs

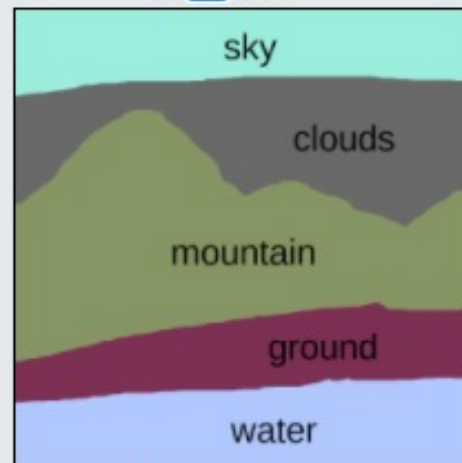
cGANs [Mirza and Osindero. 2014], SAGAN [Zhang et al., 2018], BigGAN [Brock et al., 2019]

StyleGAN-XL [Sauer et al., 2022]

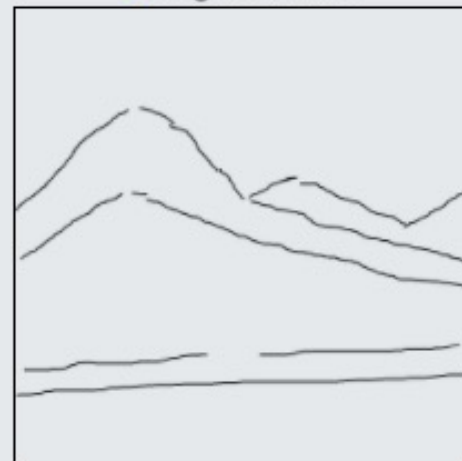


Snow mountains  
near a frozen lake  
with pink clouds in  
the sky.

Text



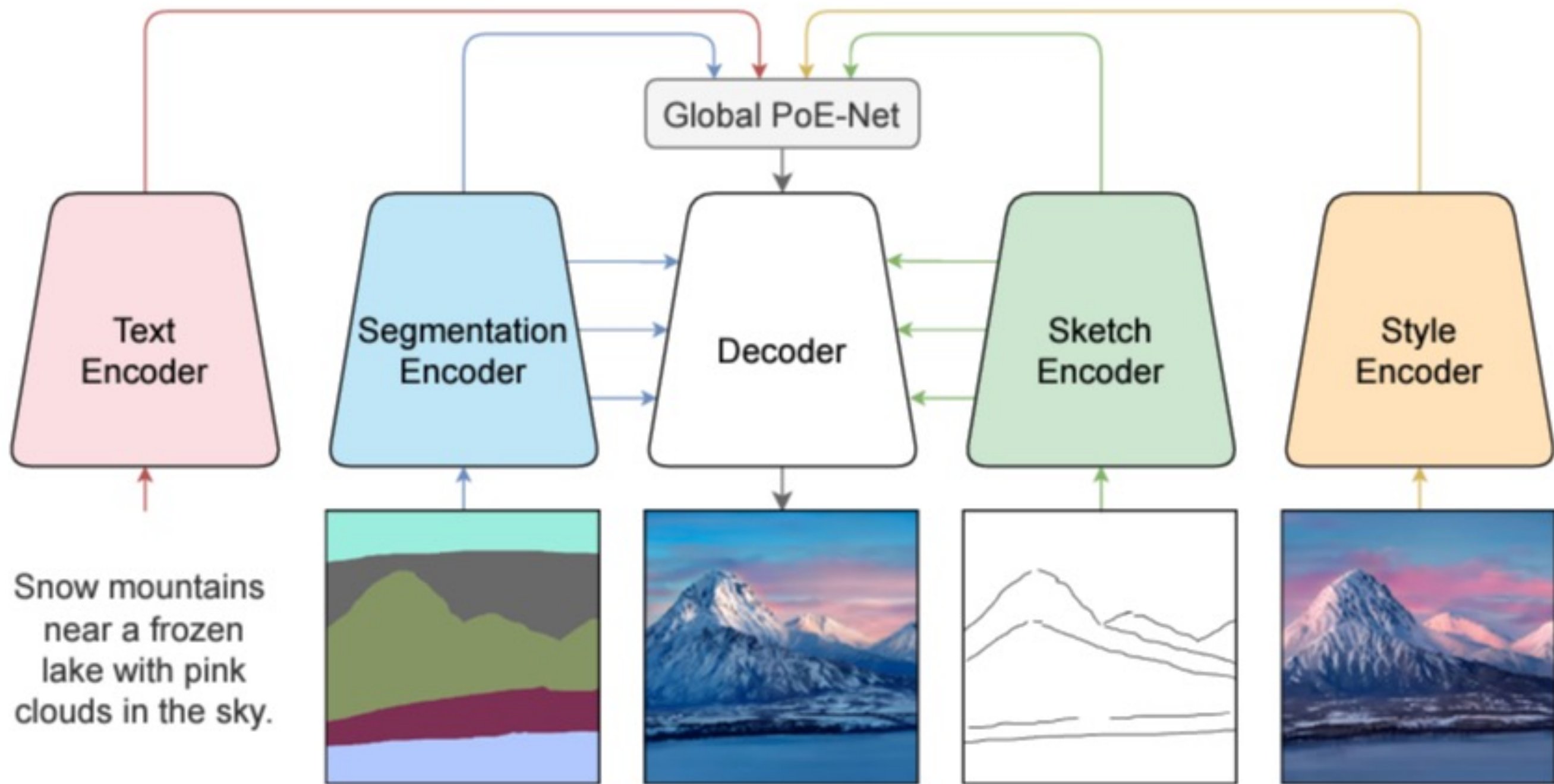
Segmentation



Sketch



Multimodal Conditional Image Synthesis with Product-of-Experts GANs [Huang et al., 2021]



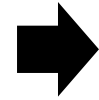
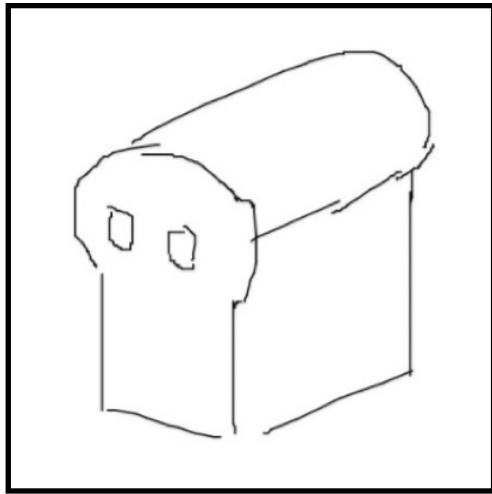
Multimodal Conditional Image Synthesis with Product-of-Experts GANs [Huang et al., 2021]

# GauGAN2 Demo

<http://gaugan.org/gaugan2/>



# Supervised Learning Approach



Edges2cats

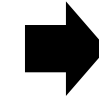
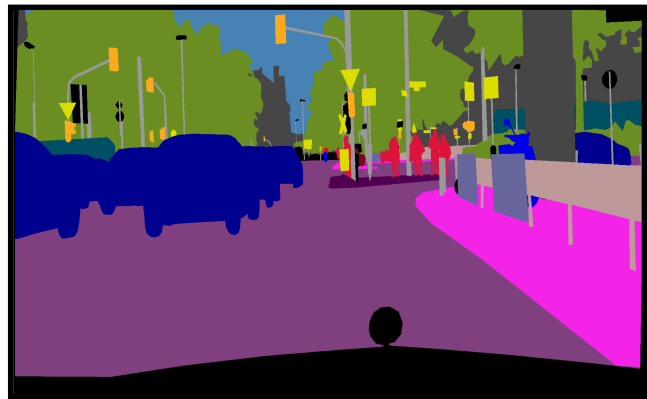
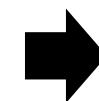


Image colorization

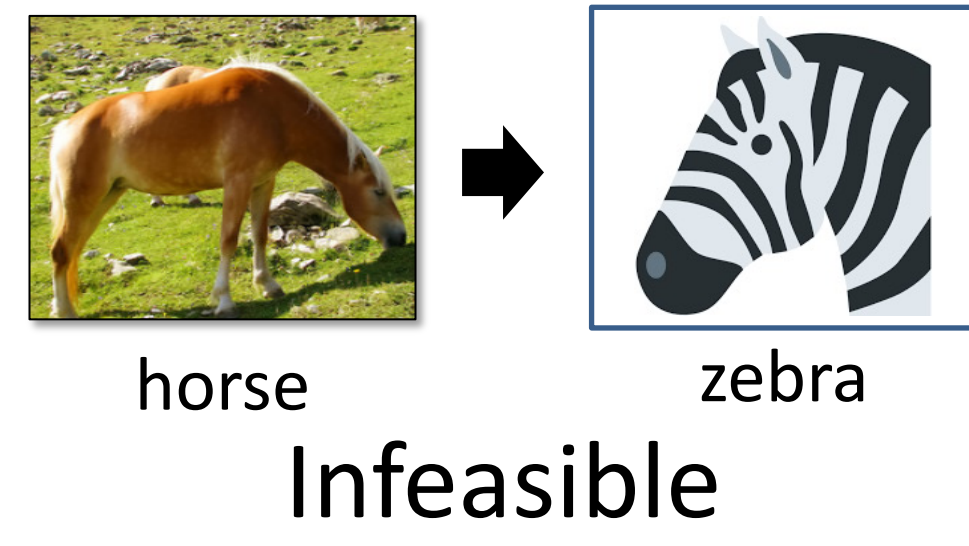
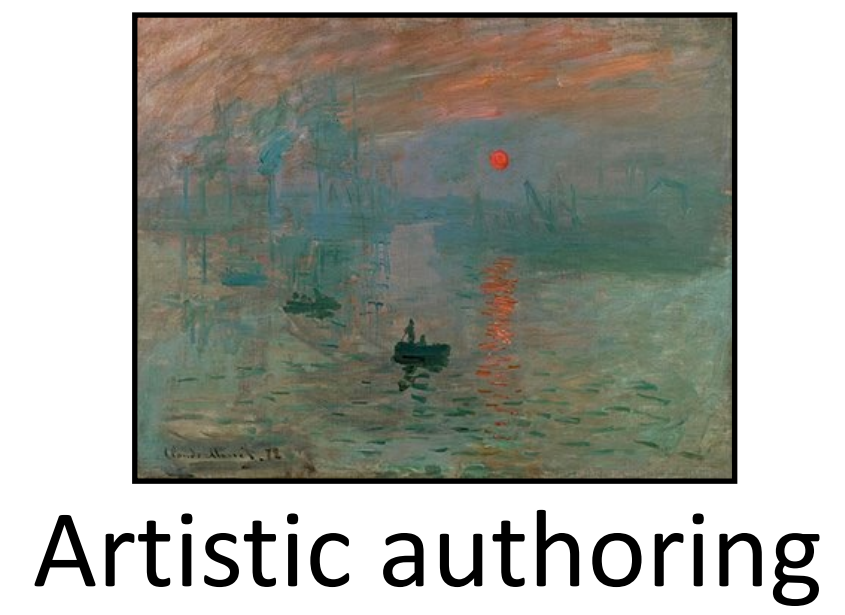
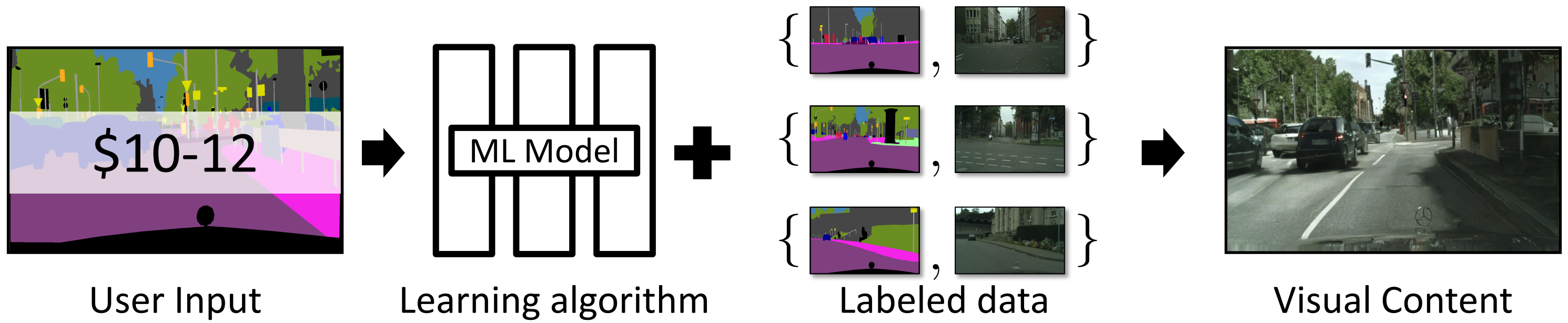


Street view images



Natural outdoor images

# Supervised Learning Approach

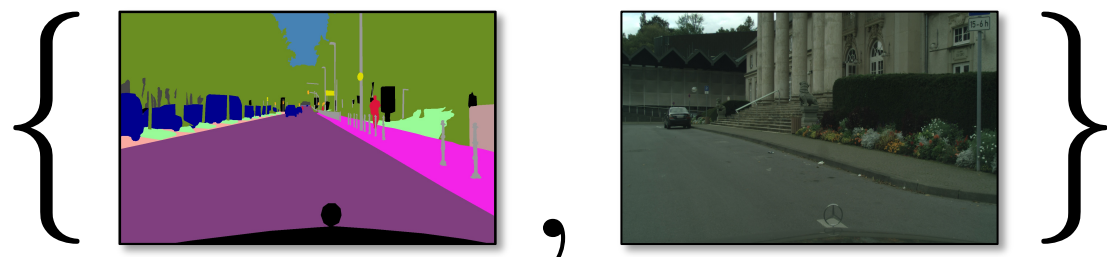
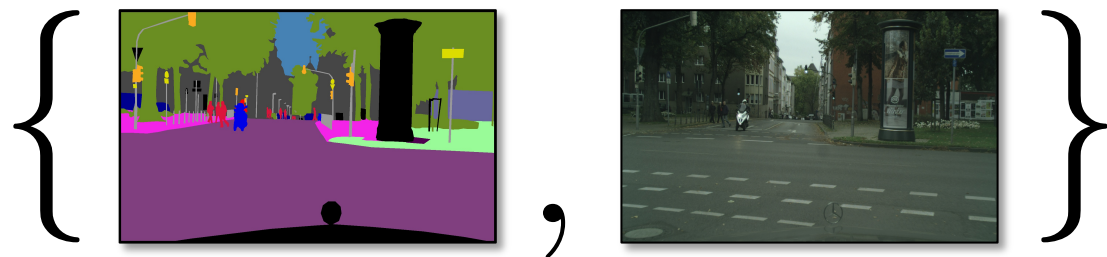
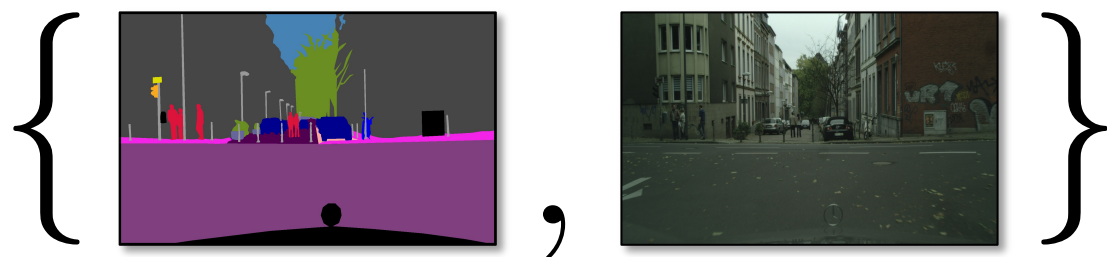




# Supervised

$x_i$

$y_i$



⋮

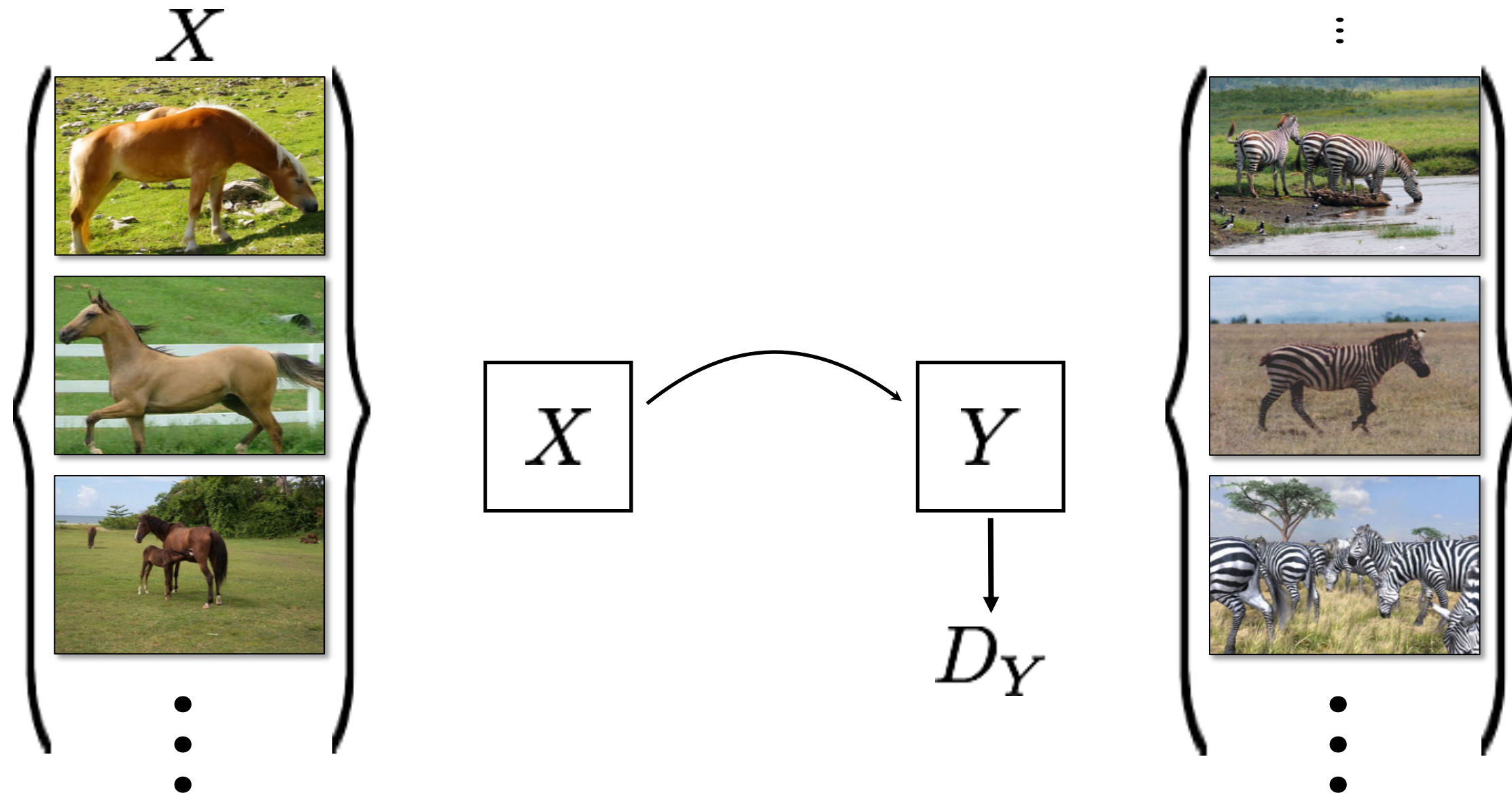
# Unsupervised

$X$

$Y$



# Unsupervised Learning of $p(y | x)$





# Unsupervised Learning of $p(y | x)$

$X$



⋮

$$\mathbb{E}_x \log(1 - D(\underline{G(x)})) + \mathbb{E}_y \log D(\underline{y})$$

$X$

$Y$

$D$

Discriminator

$Y$



⋮

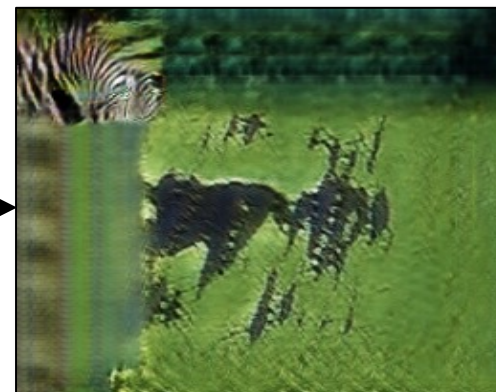
fake zebra



real zebra

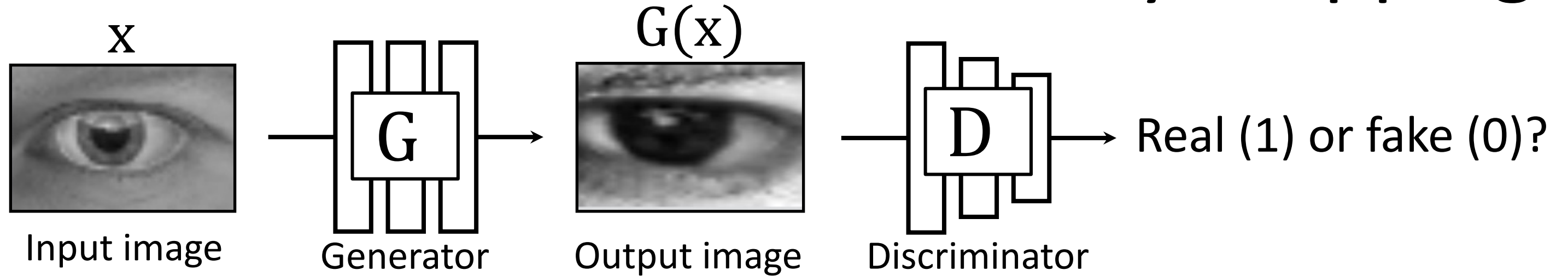


# Unsupervised Learning of $p(y | x)$



- artifacts
- ignore inputs

# Additional Constraint: Identity Mapping

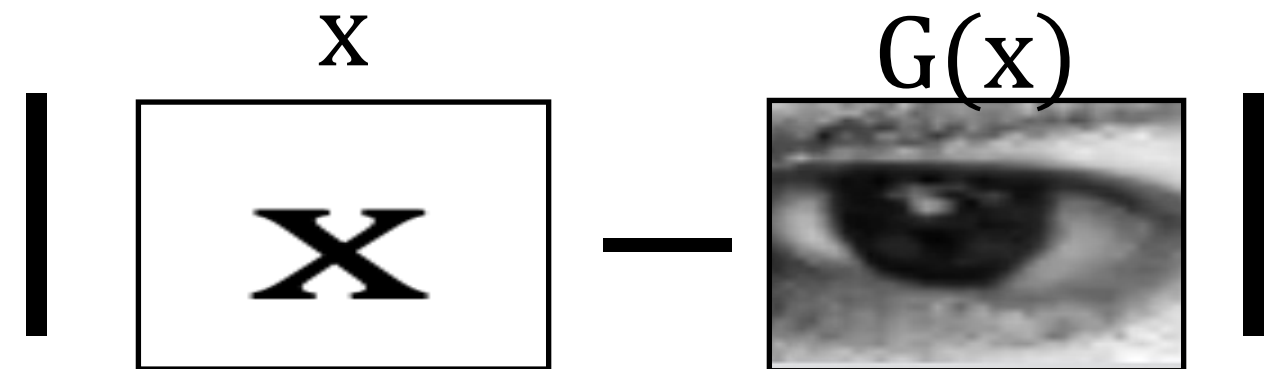


## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

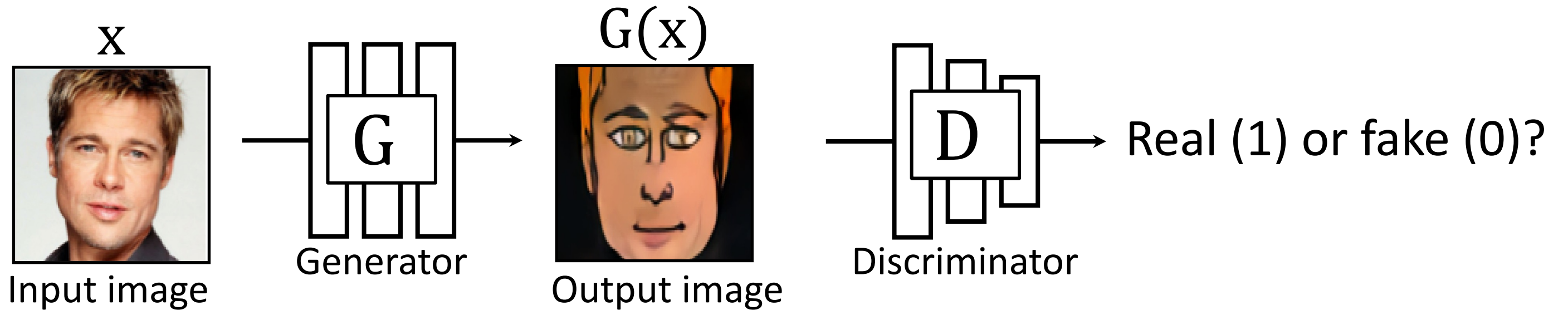
## Self-Regularization loss

$$\mathbb{E}_x \|G(x) - x\|_1$$





# Additional Constraint: Feature Loss



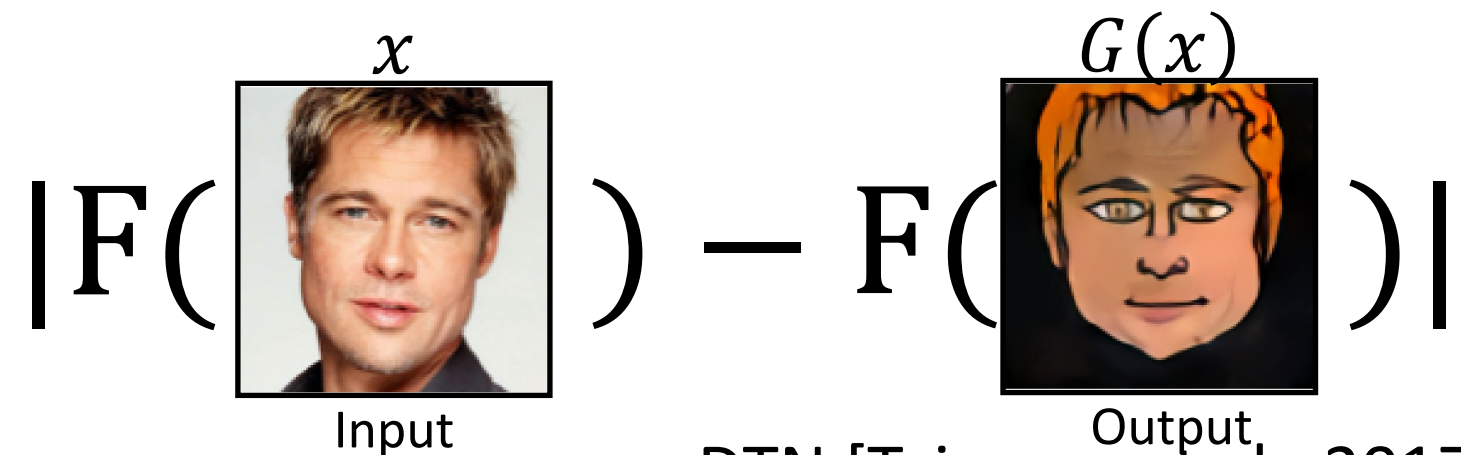
## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

## Feature loss

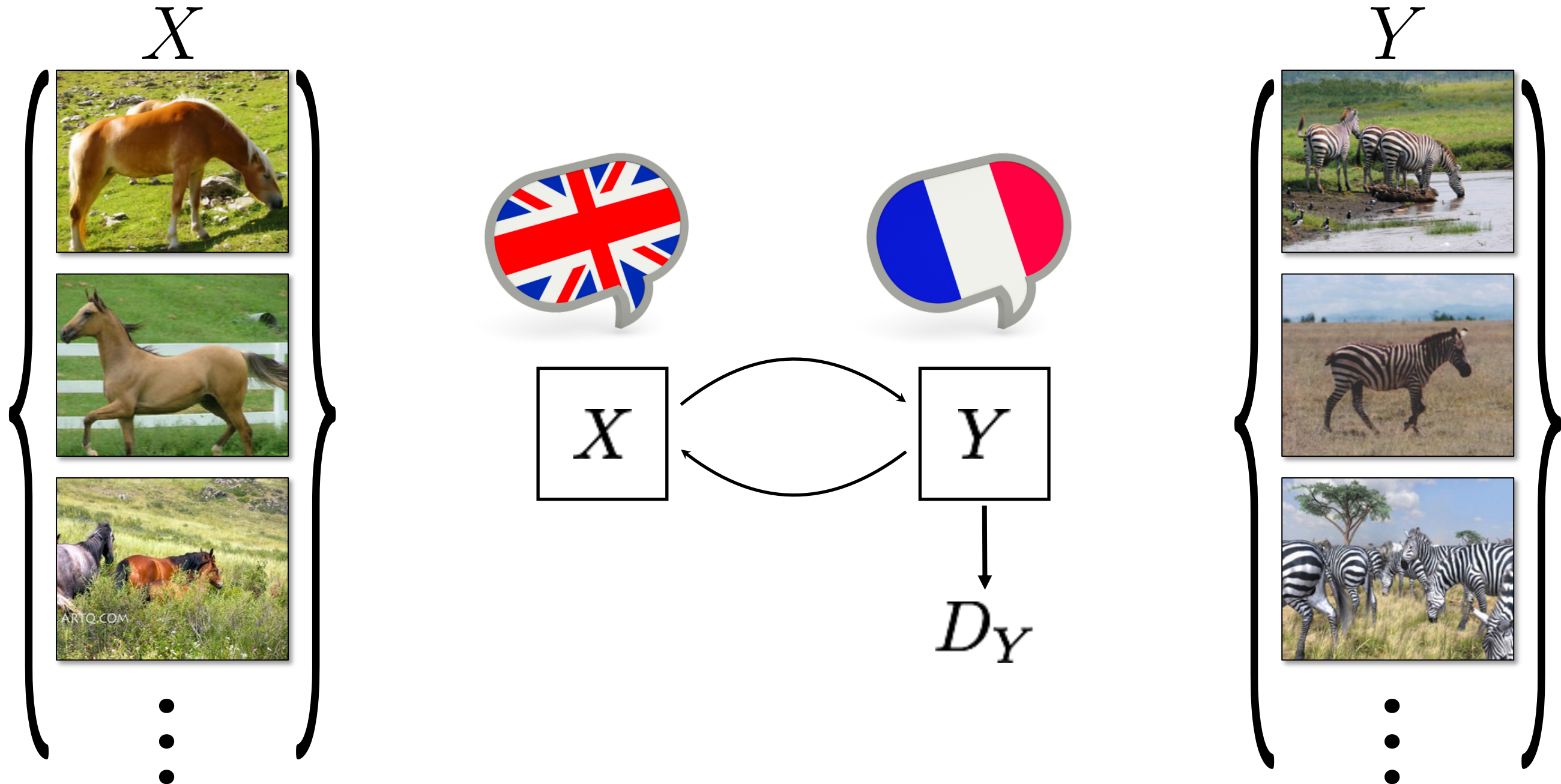
$$\mathbb{E}_x ||F(G(x)) - F(x)||$$

Requires  $F$  to work across two domains

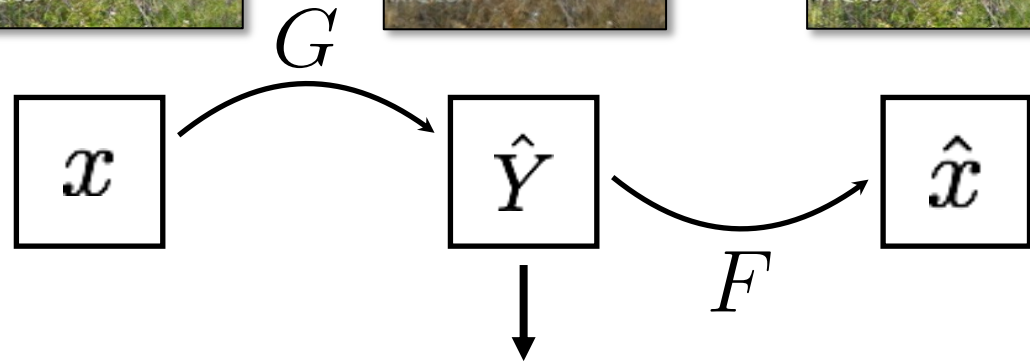
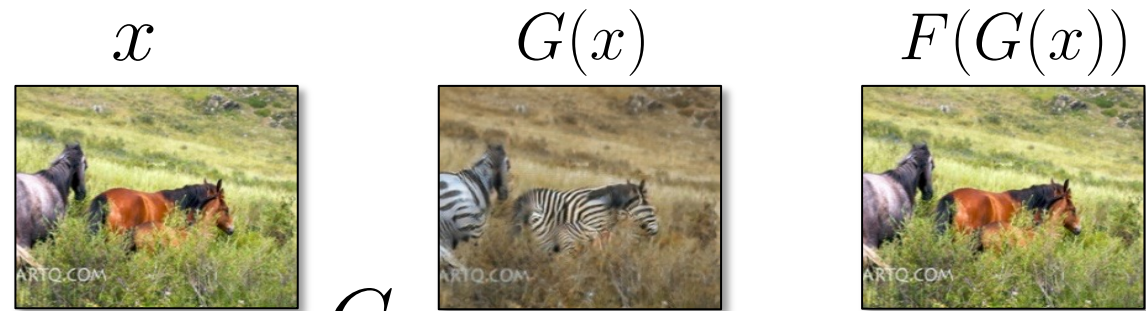


DTN [Taigman et al., 2017]

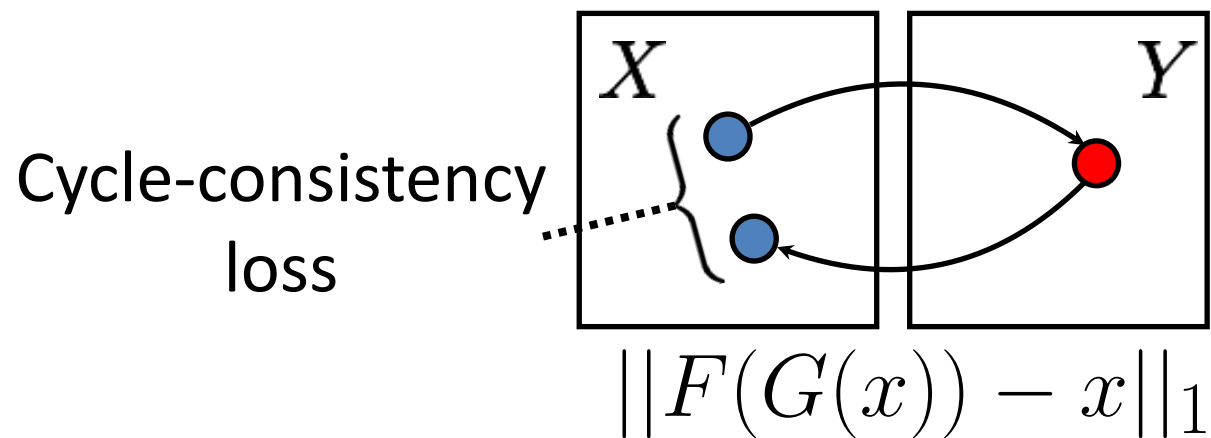
# Additional Constraint: Cycle-Consistency



# Cycle-Consistent Adversarial Networks



Adversarial loss  $D_Y(G(x))$



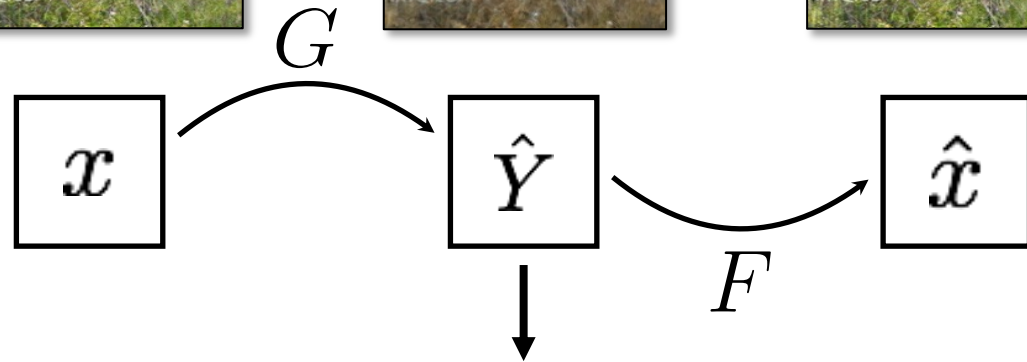
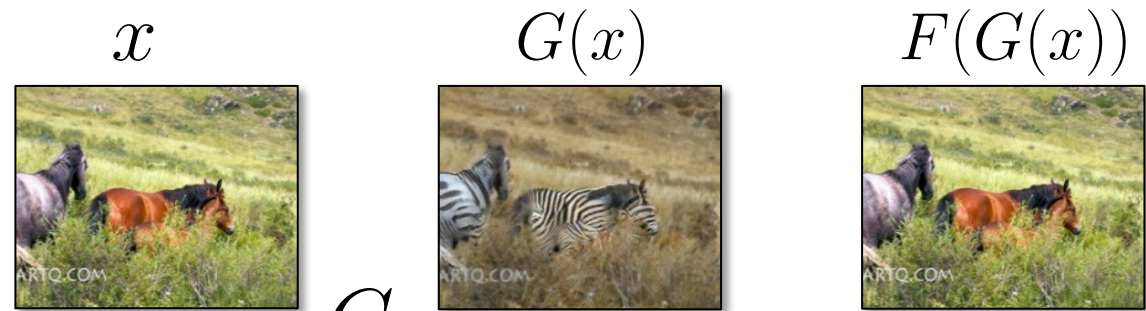
**Adversarial loss**

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

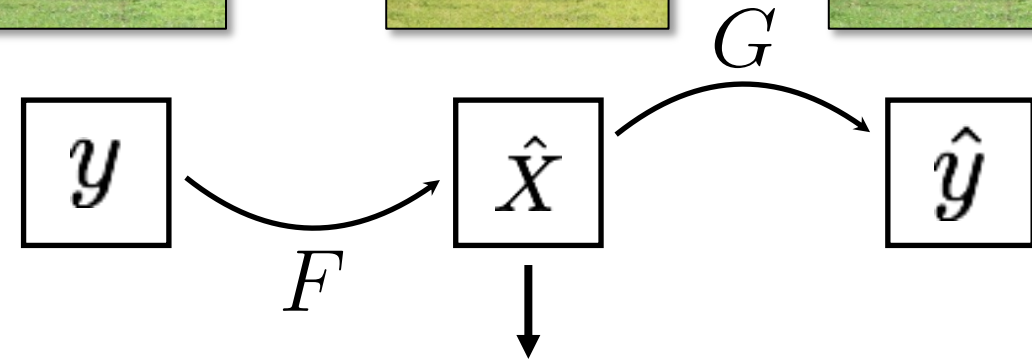
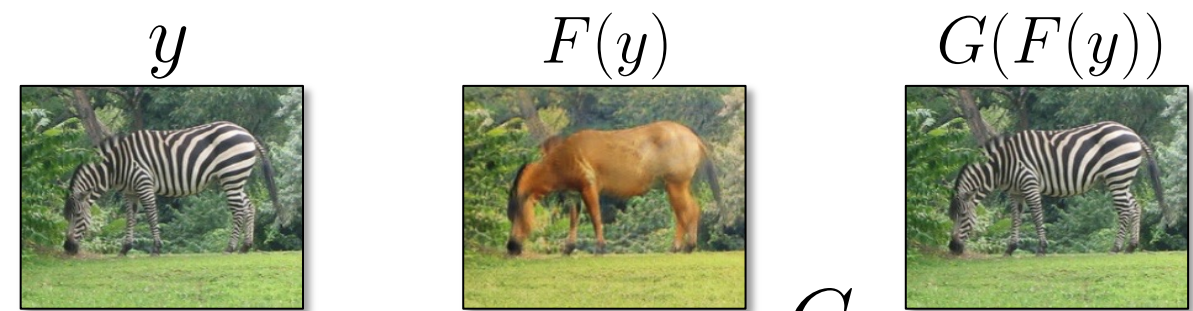
**Cycle-consistency loss**

$$\mathbb{E}_x \|F(G(x)) - x\|_1$$

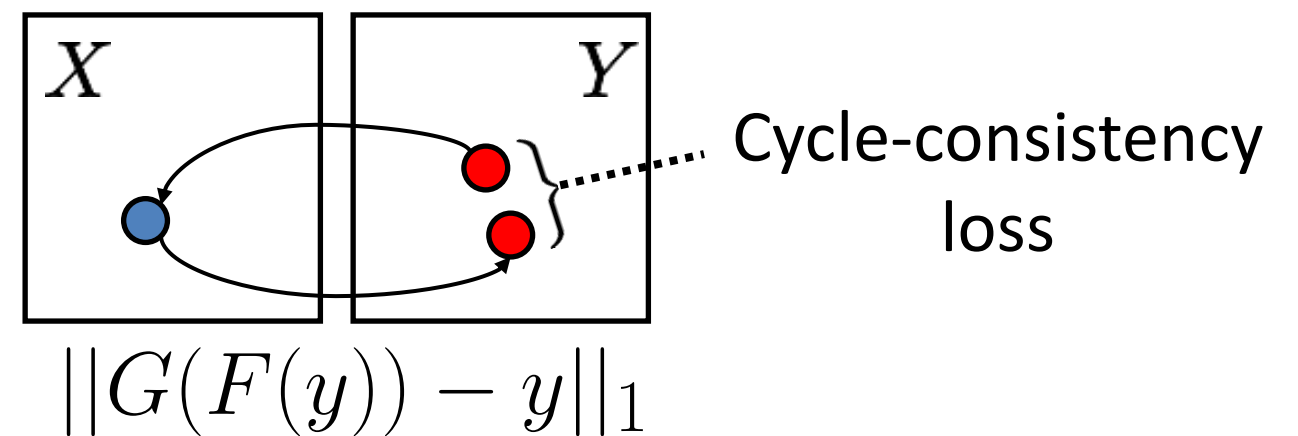
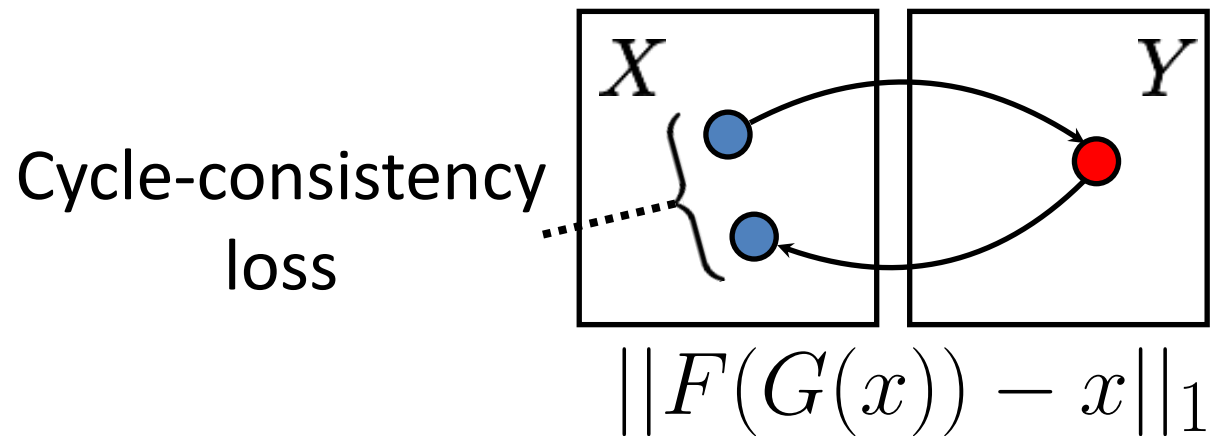
# Cycle-Consistent Adversarial Networks



Adversarial loss  $D_Y(G(x))$



$D_X(F(y))$  Adversarial loss



# Results



# Horse → Zebra





Orange → Apple





# Monet's paintings → photographic style





# Monet's paintings → photographic style





# Collection Style Transfer



Photograph ©Alexei Efros



Monet



Van Gogh



Cezanne



Ukiyo-e



# Improving the Realism of CG Rendering



CG Game: Grand Theft Auto



Street view images in German cities

Data from [Richter et al., 2016], [Cordts et al, 2016]

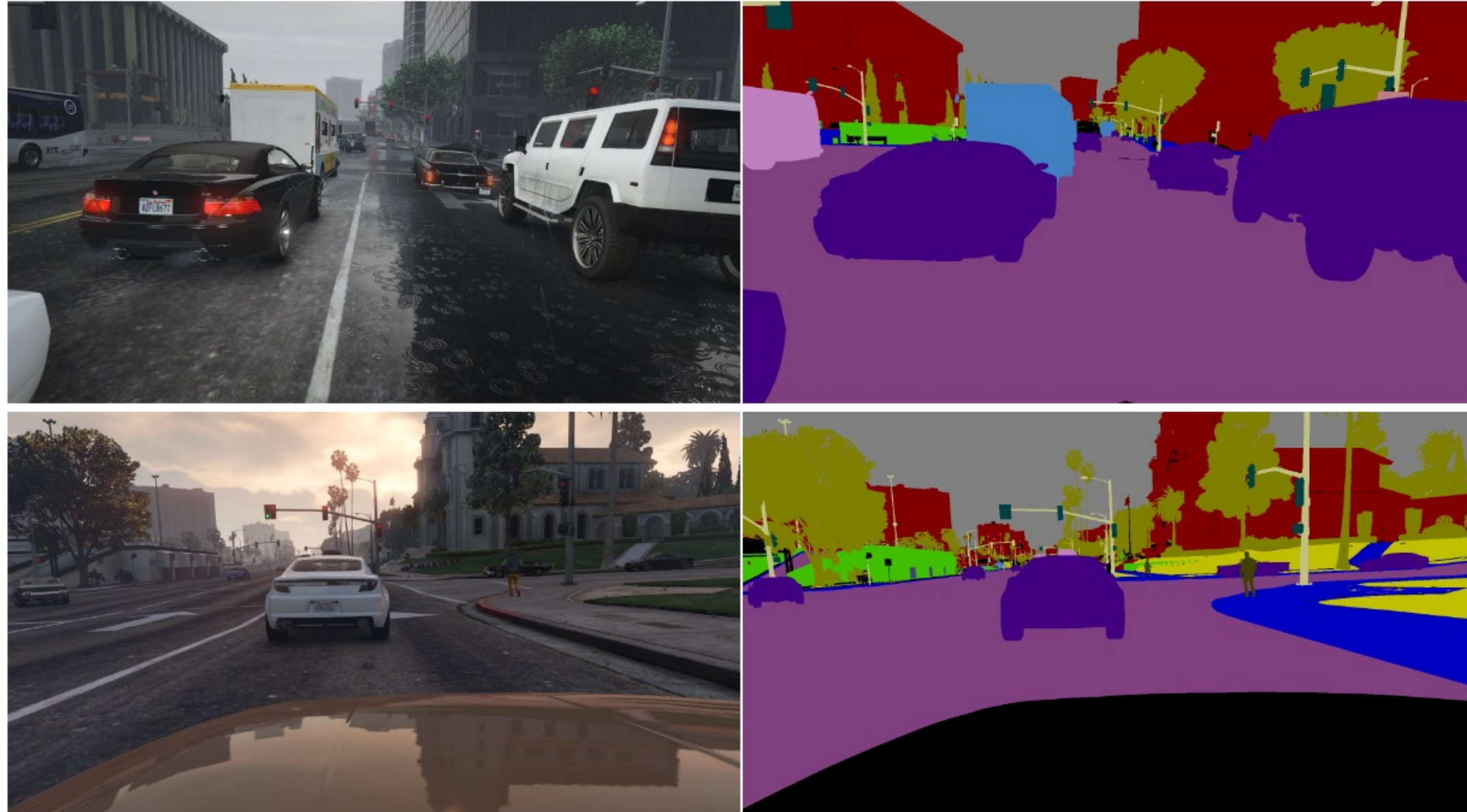
# Improving the Realism of CG Rendering



Output image with CG image street view style



# Domain Adaptation with CycleGAN

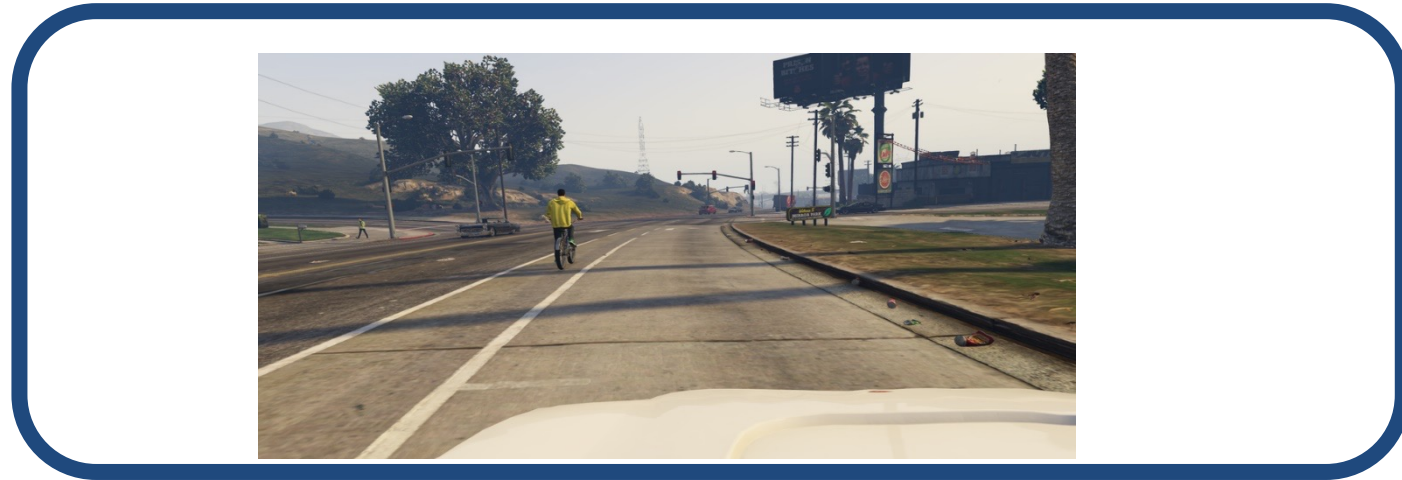


CG images

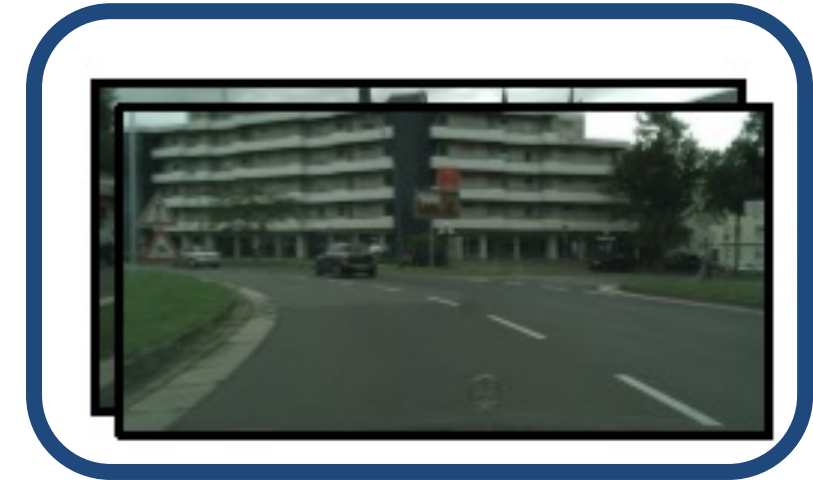
Free segmentation labels

Data and labels from [Richter et al. 2016]

# Domain Adaptation with CycleGAN



Train on CG data



Test on real images

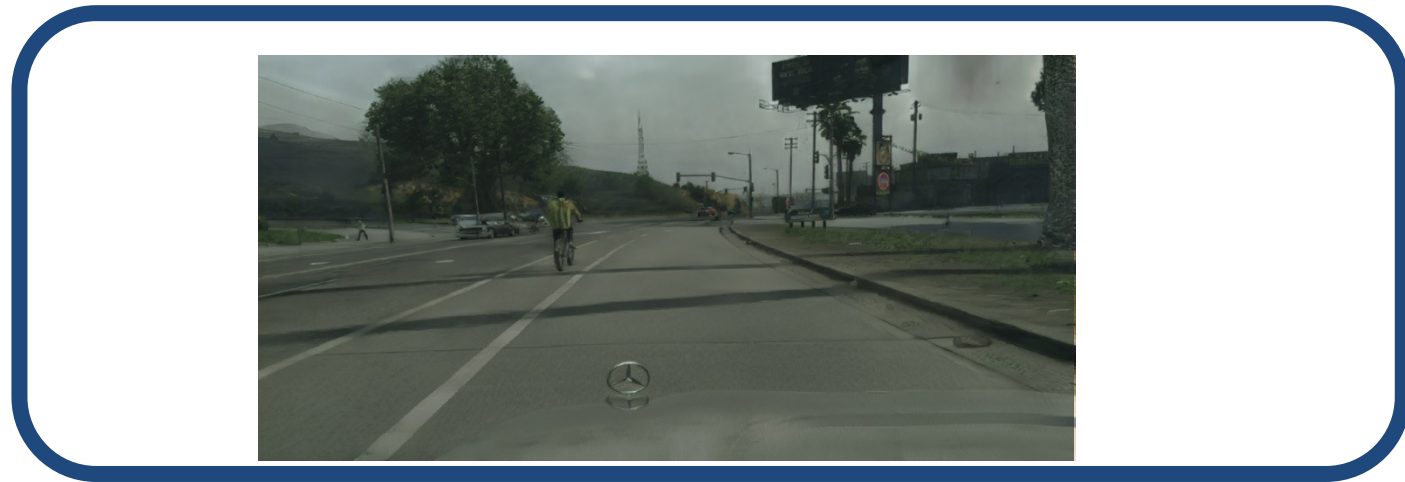
Class-weighted  
Accuracy

47.4

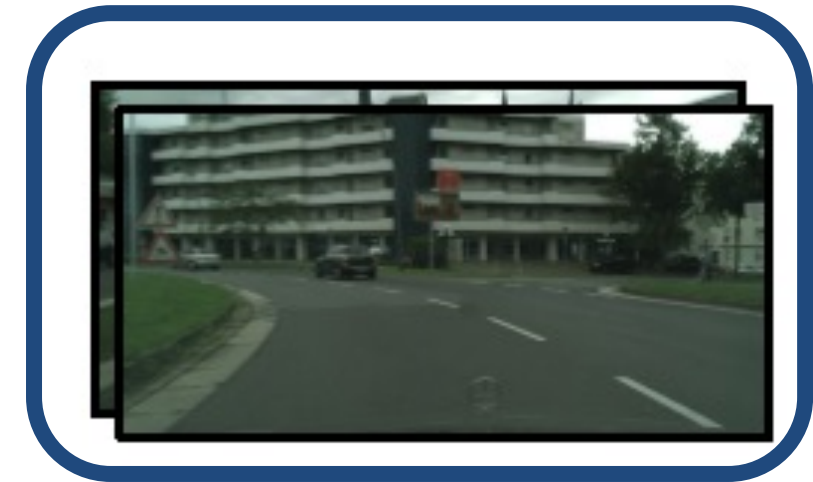
Train on CG

70  
60  
50  
40

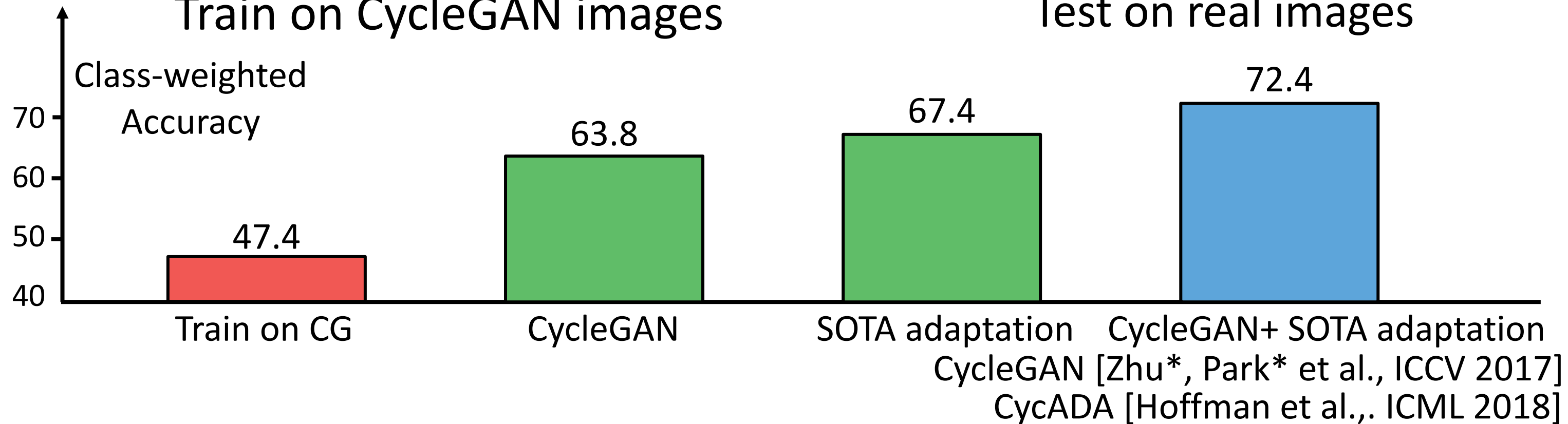
# Domain Adaptation with CycleGAN



Train on CycleGAN images



Test on real images

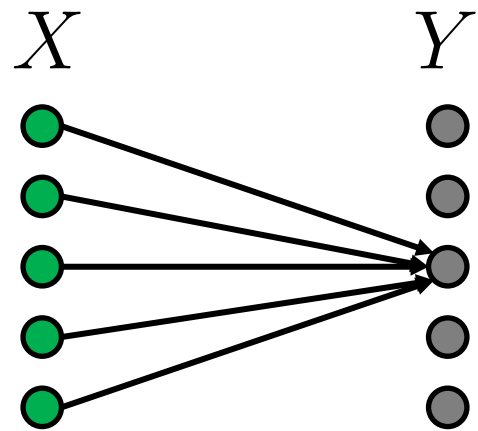
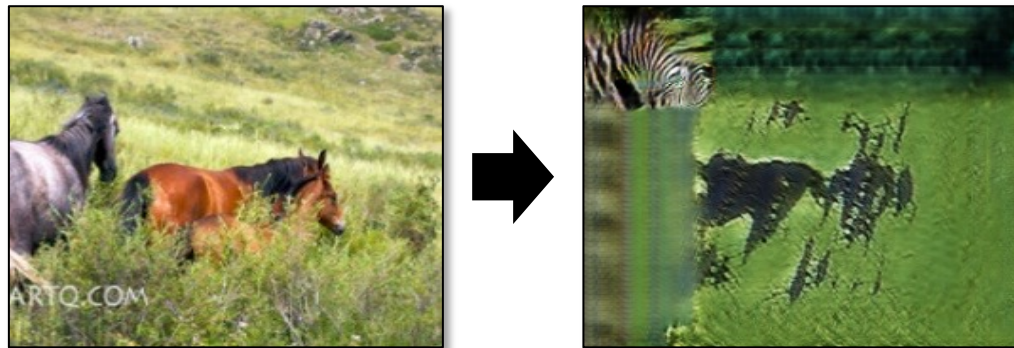


Why CycleGAN works

# Why CycleGAN works

## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

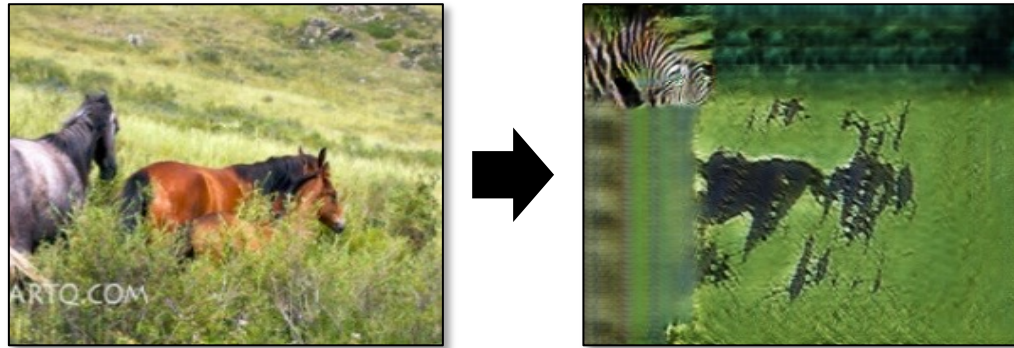




# Why CycleGAN works

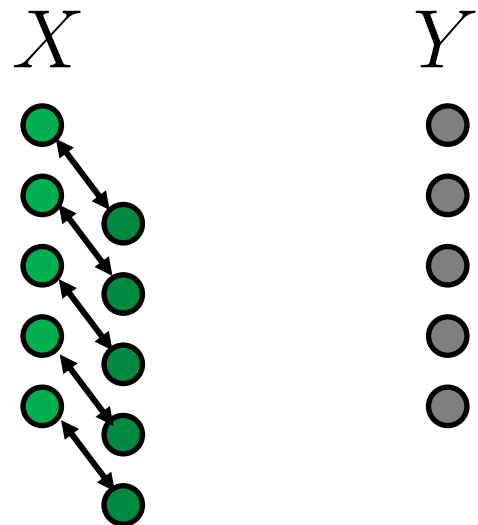
## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$



## Cycle-consistency loss

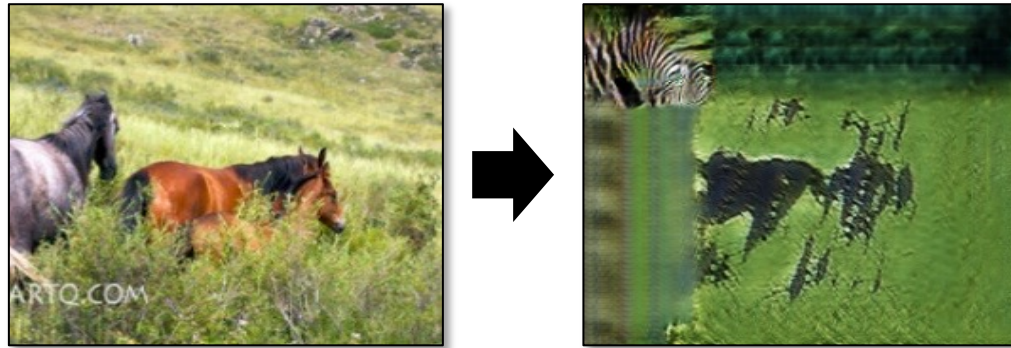
$$\mathbb{E}_x ||F(G(x)) - x||_1$$



# Why CycleGAN works

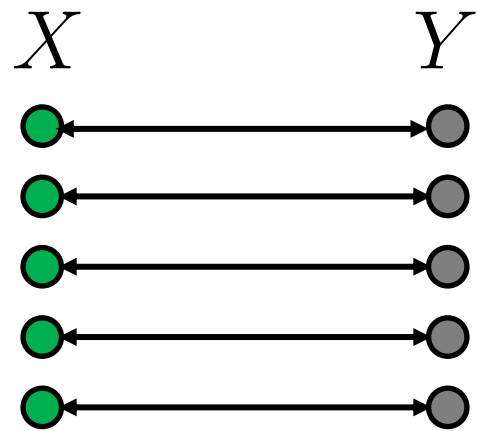
## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$



## Cycle-consistency loss

$$\mathbb{E}_x ||F(G(x)) - x||_1$$



## Full objective



# Why CycleGAN works

## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

## Cycle-consistency loss

$$\mathbb{E}_x ||F(G(x)) - x||_1$$

$x$



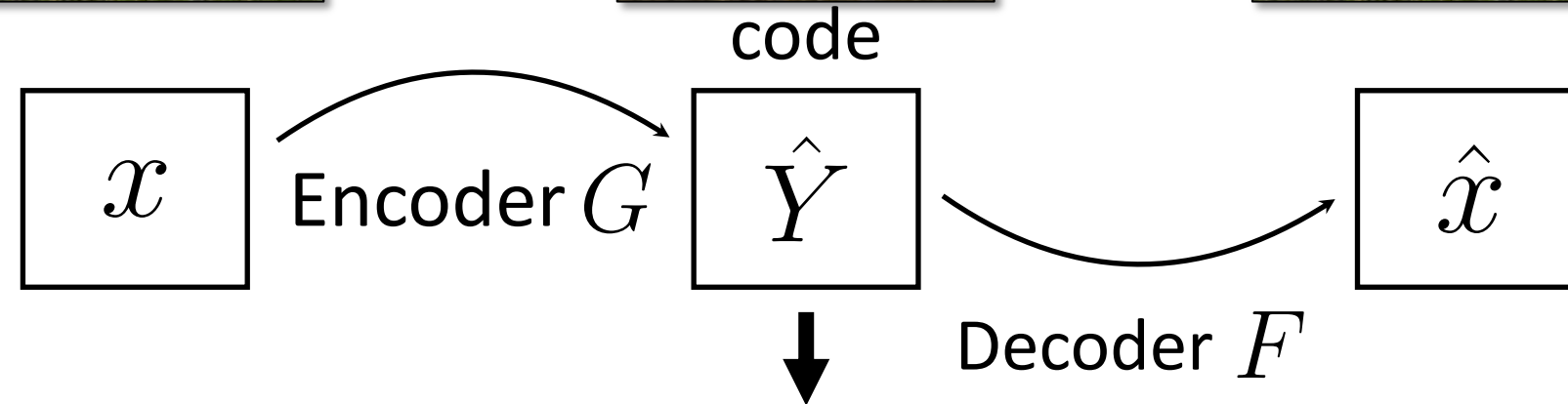
$G(x)$



$F(G(x))$



Auto-encoder  
w/ domain prior



Constraint:  $\mathbb{E}_x ||G(x) - p_{data}(Y)||_1$

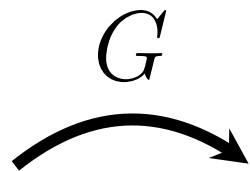
# Why CycleGAN works

## Adversarial loss

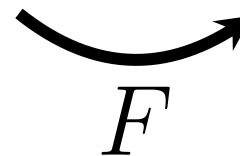
$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

Under-constrained problem

$x$



$\hat{Y}$



Prior of  $G$

$\hat{x}$



## Cycle-consistency loss

$$\mathbb{E}_x ||F(G(x)) - x||_1$$

A strong regularizer

**Assumption:** simple invertible function

**Probabilistic Interpretation :** Upper bound of conditional entropy  $H(y|x)$

[Li et al. 2017]



# Why CycleGAN works

## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

## Cycle-consistency loss

$$\mathbb{E}_x ||F(G(x)) - x||_1$$

flip the image



Invertible Perturbation

flip the image again

**Adversarial loss:** images are horizontally symmetric

**Cycle-consistency loss :**  $||F \circ P^{-1}(P \circ G(x)) - x||$

# Applications of CycleGAN

A blue oval with a black border, containing four lines of text. The text is centered within the oval and reads: "Natural images", "Computer Graphics", "Computational Photography", and "Computer Vision".

**Natural images**

Computer Graphics

Computational Photography

Computer Vision



# Photo Enhancement





# Image Dehazing



Foggy image



Clear image

[Engin et al. CVPRW 2018]

# Other Image data

## Natural images

Computer Graphics



Computer Vision

Biology

Medical Imaging

Robotics



[Bartha et al. 2018]

Remote Sensing

**Non-image data**

**Other Image data**

**Natural images**

Computer Graphics



Computer Vision

Biology  
Medical Imaging  
Robotics

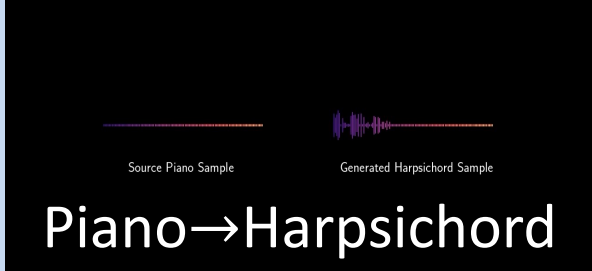


[Bartha et al. 2018]

Remote Sensing

Natural language (NLP)

Computer music



[Huang et al. 2019]

Audio processing

Cryptography

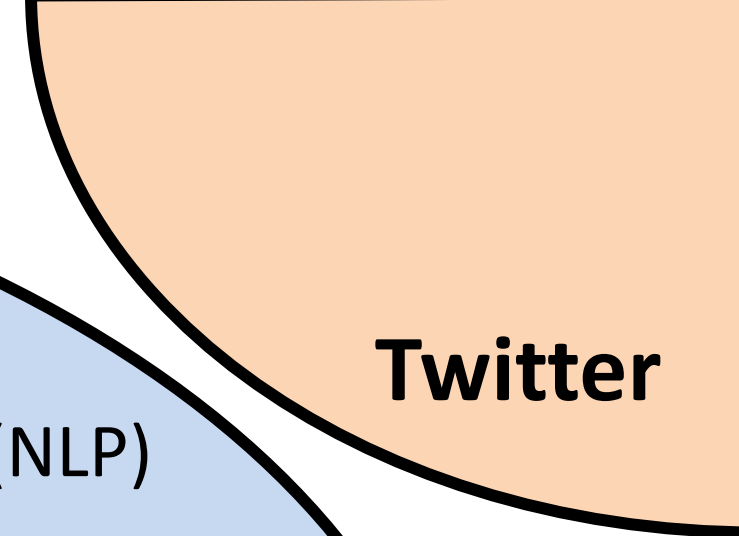
**Art**



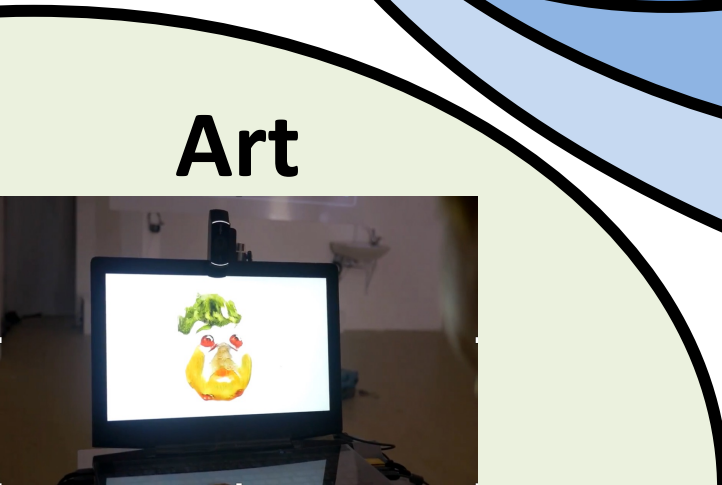
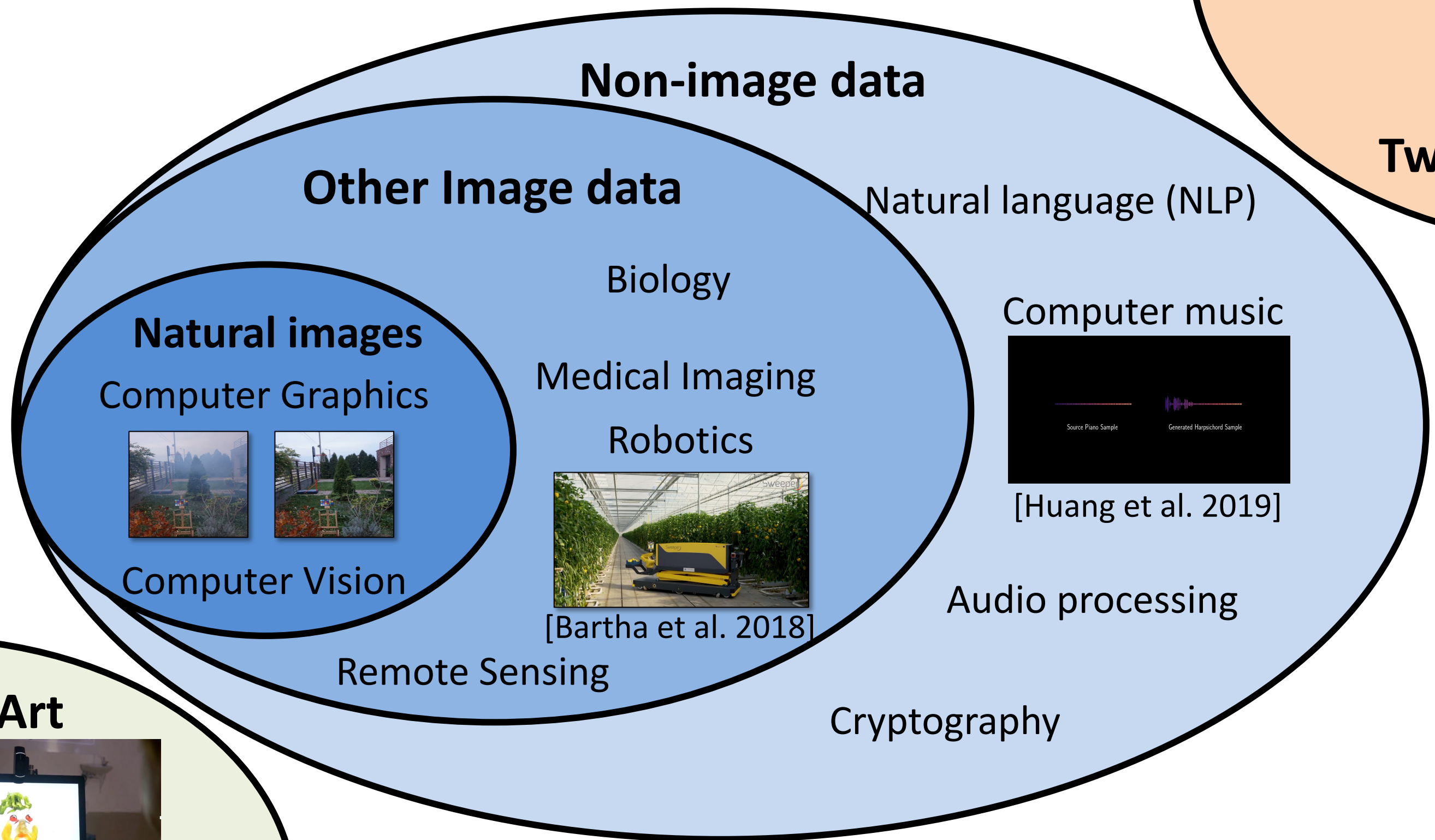
# Artistic Applications

## **The Electronic Curator**





Twitter



Art

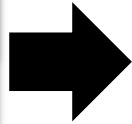


# Latest from #CycleGAN

Input dog



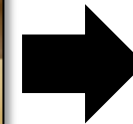
Output cat



Input cat



Output dog



# Style and Content Disentanglement

# Style and Content Separation

**A**

Classification

<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
A	B	C	D	E
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
B	C	A	E	D

Domain Adaptation

**B**

Extrapolation

<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
A	B	C	D	E
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>
<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>
?	?	C	D	E

Paired Image-to-Image Translation

**C**

Translation

<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	?	?	?
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>			
A	B	C	D	E			
<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>			
<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	?	?	?
?	—	—	—	?	F	G	H

Unpaired Image-to-Image Translation

Training
Generalization

Separating Style and Content  
[Tenenbaum and Freeman 1996]

$$y_k^{sc} = \sum_{i=1}^I \sum_{j=1}^J w_{ijk} a_i^s b_j^c.$$



# Style and Content

## Adversarial loss

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$



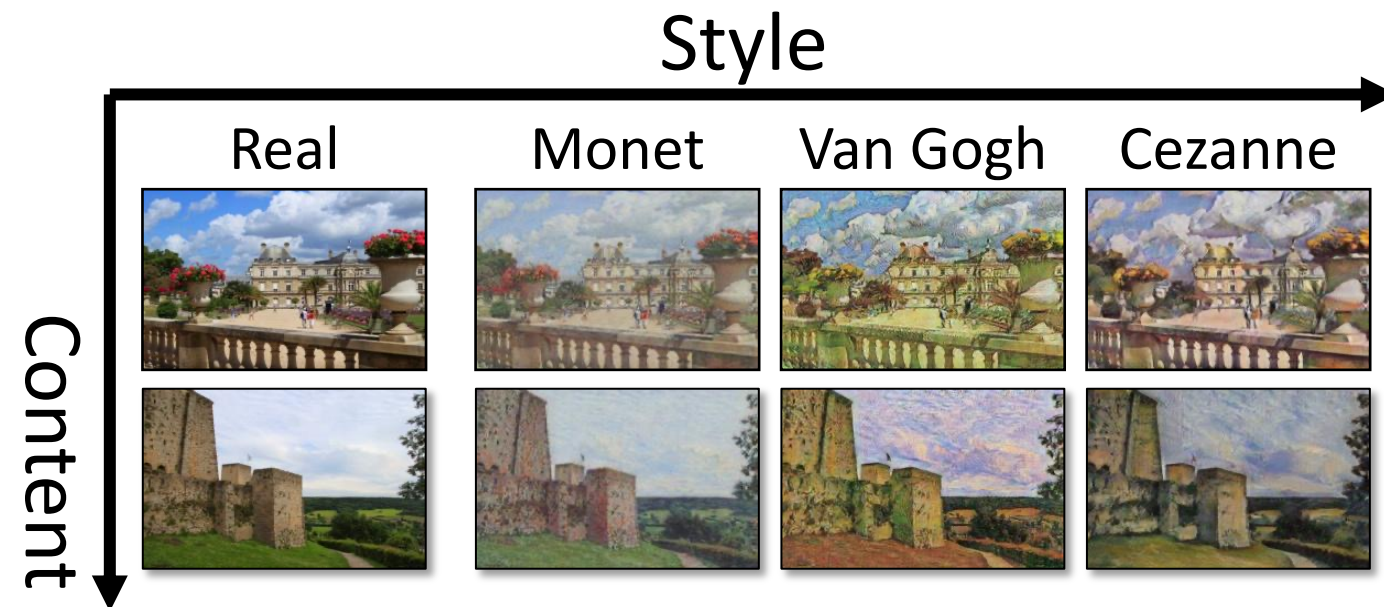
$p(x) \rightarrow p(y)$  change **style**

## Cycle-consistency loss

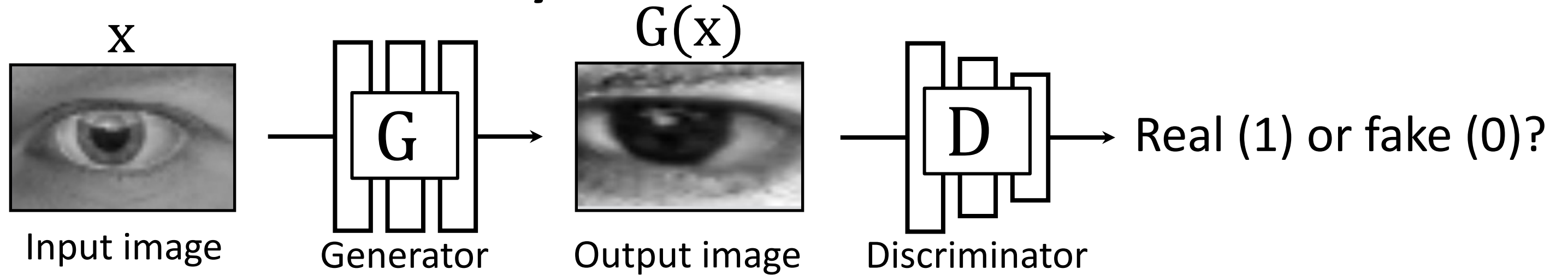
$$\mathbb{E}_x ||F(G(x)) - x||_1$$



Bidirectional: preserve **content**



# Style and Content

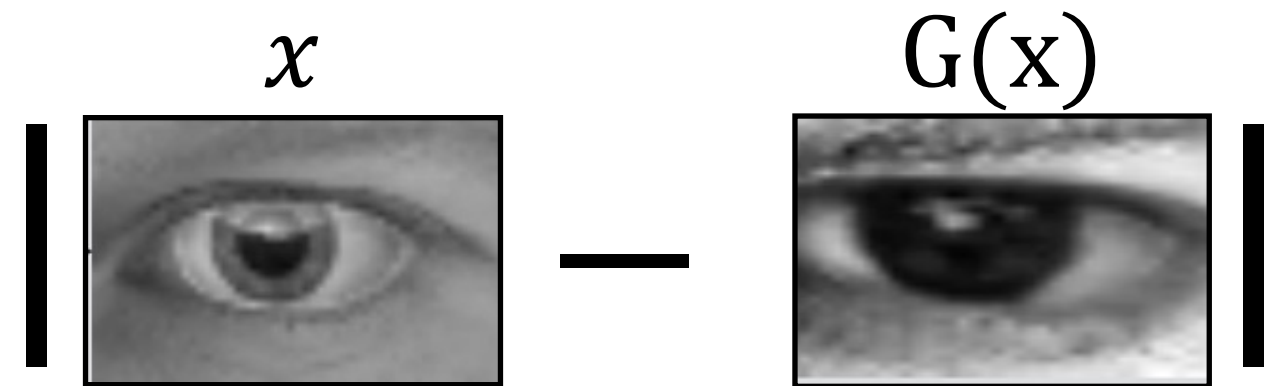


## Adversarial loss (change style)

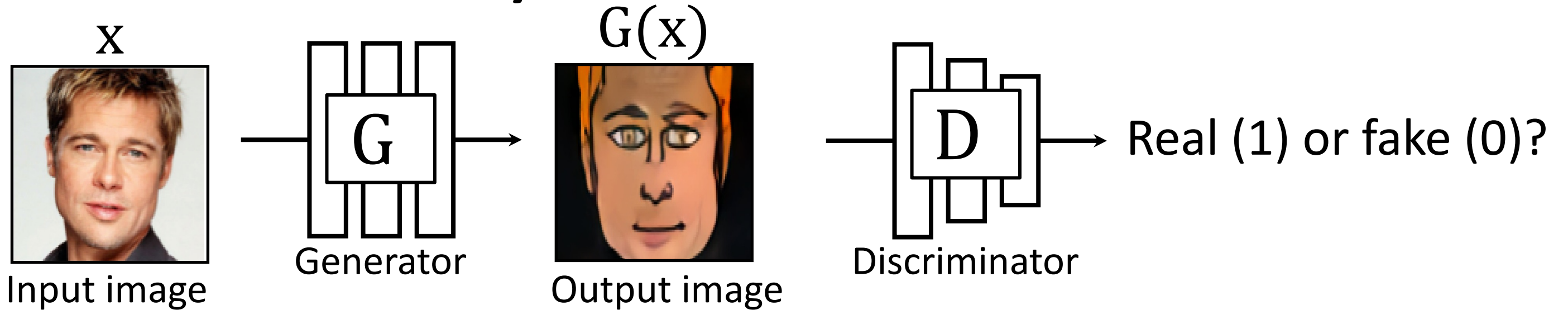
$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

## L1 loss (preserve content in pixel space)

$$\mathbb{E}_x ||G(x) - x||_1$$



# Style and Content

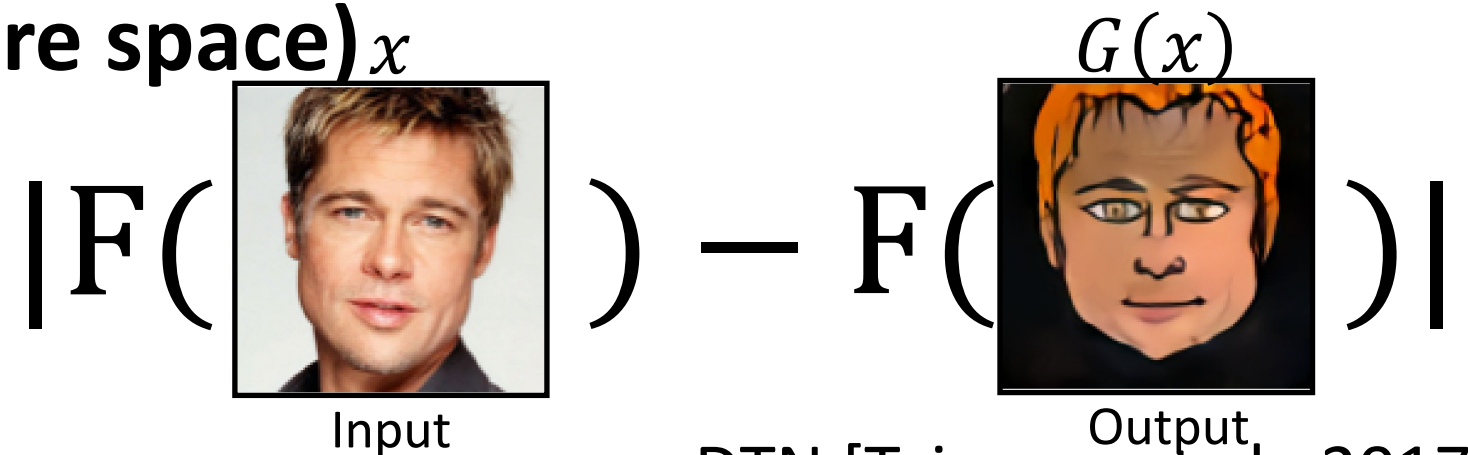


## Adversarial loss (change style)

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

## Feature loss (Preserve content in feature space)<sub>x</sub>

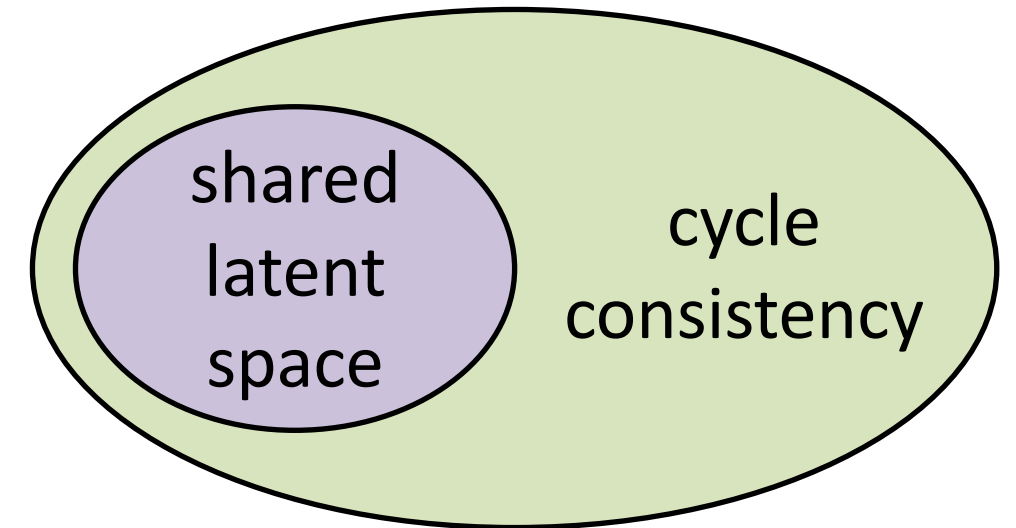
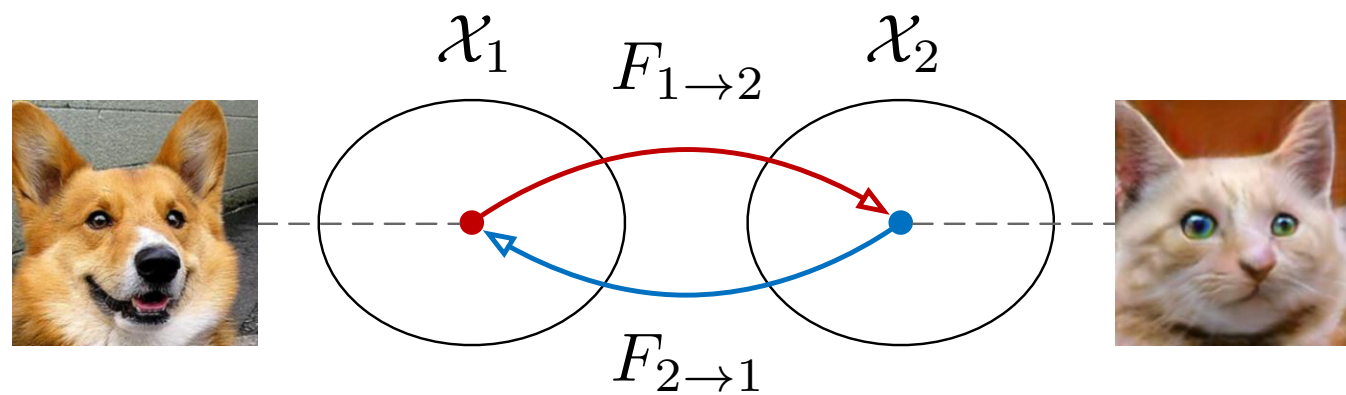
$$\mathbb{E}_x ||F(G(x)) - F(x)||$$



DTN [Taigman et al., 2017]

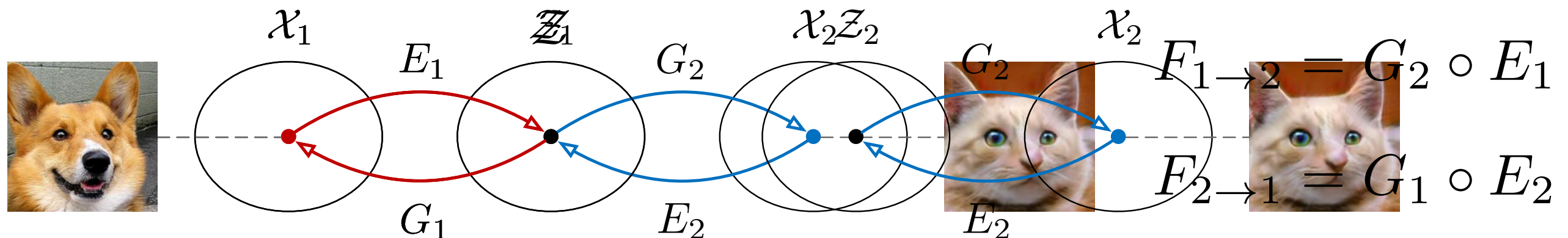
# CycleGAN and UNIT

- CycleGAN (cycle consistency)



- UNIT (shared latent space) [Liu et al. 2017]

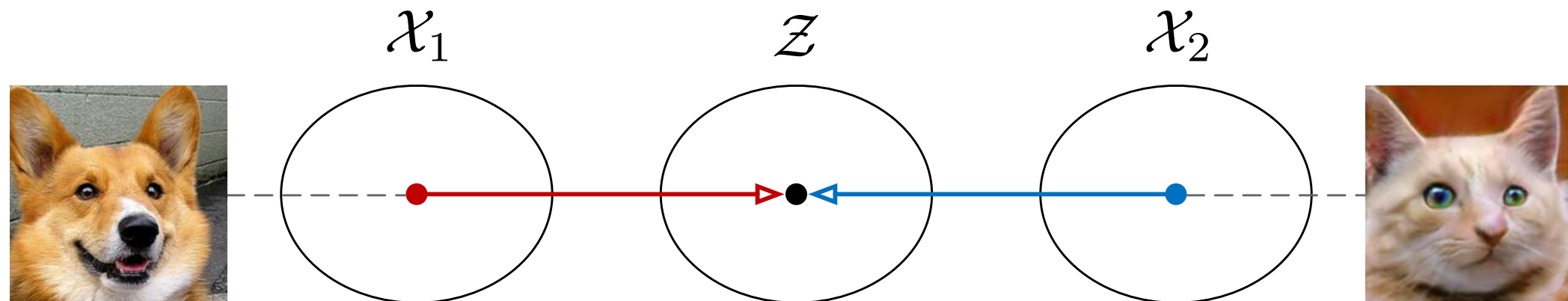
shared latent space  $\Rightarrow$  cycle consistency





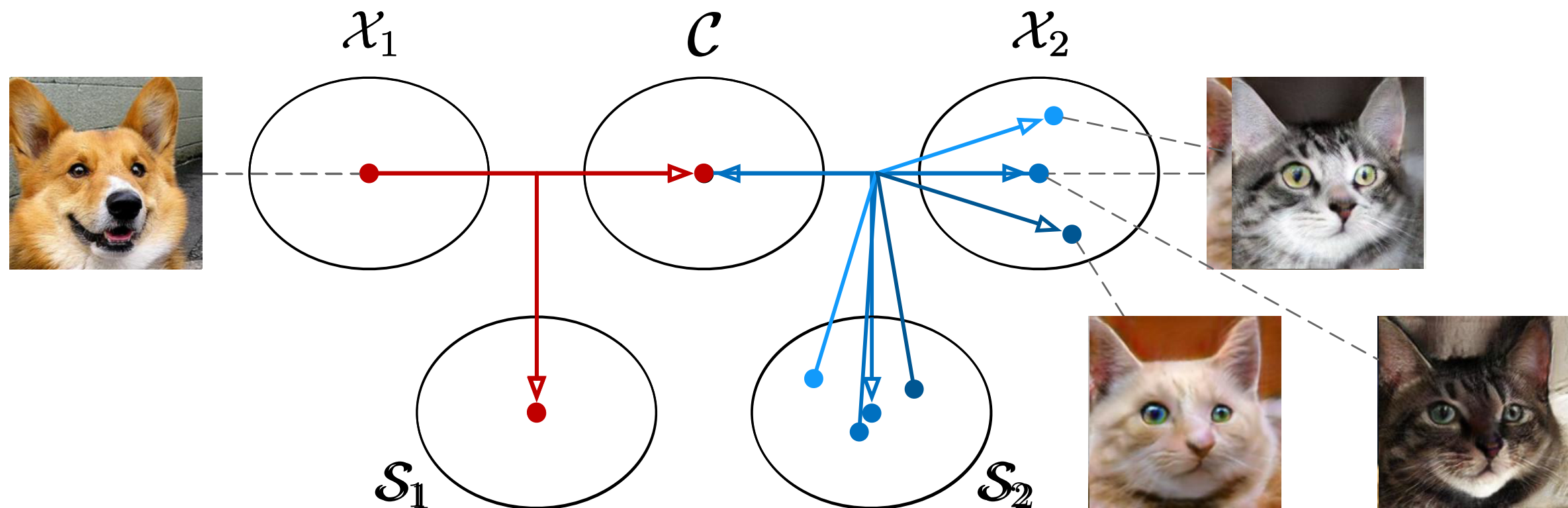
# Disentangling the Latent Space

- UNIT
  - A single **shared, domain-invariant** latent space  $\mathcal{Z}$

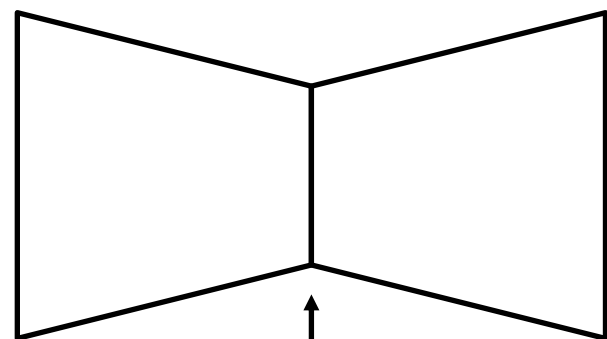
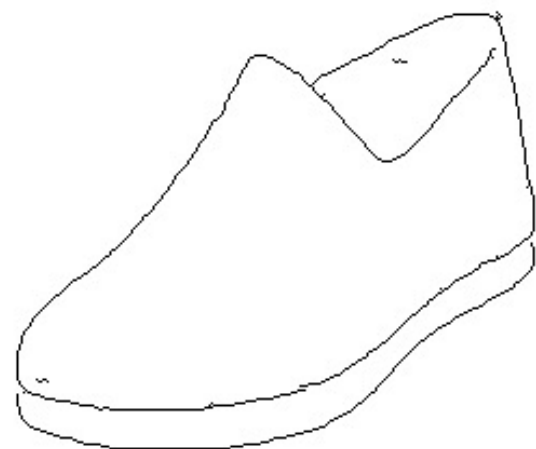


# Disentangling the Latent Space

- Multimodal UNIT (MUNIT)
  - A **content** space  $\mathcal{C}$  that is **shared, domain-invariant**
  - Two **style** spaces  $\mathcal{S}_1, \mathcal{S}_2$  that are **unshared, domain-specific**

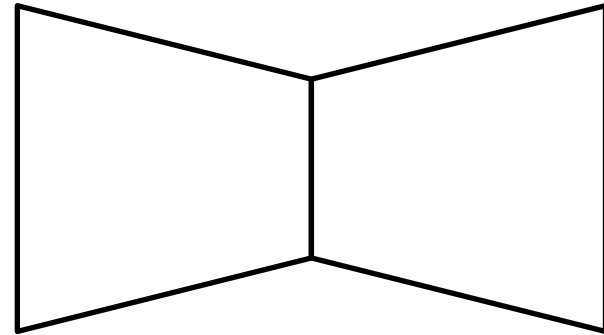
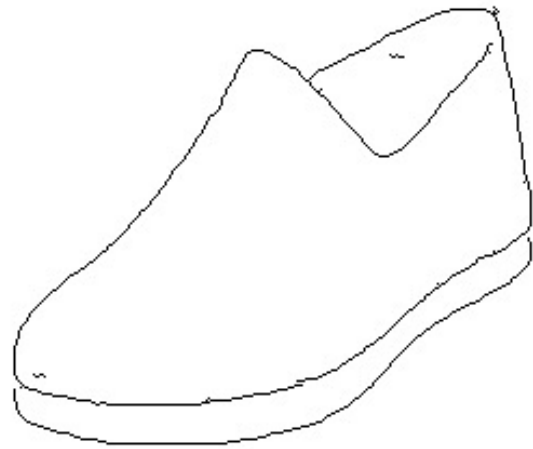


# Unimodality





# Towards Multimodality



...

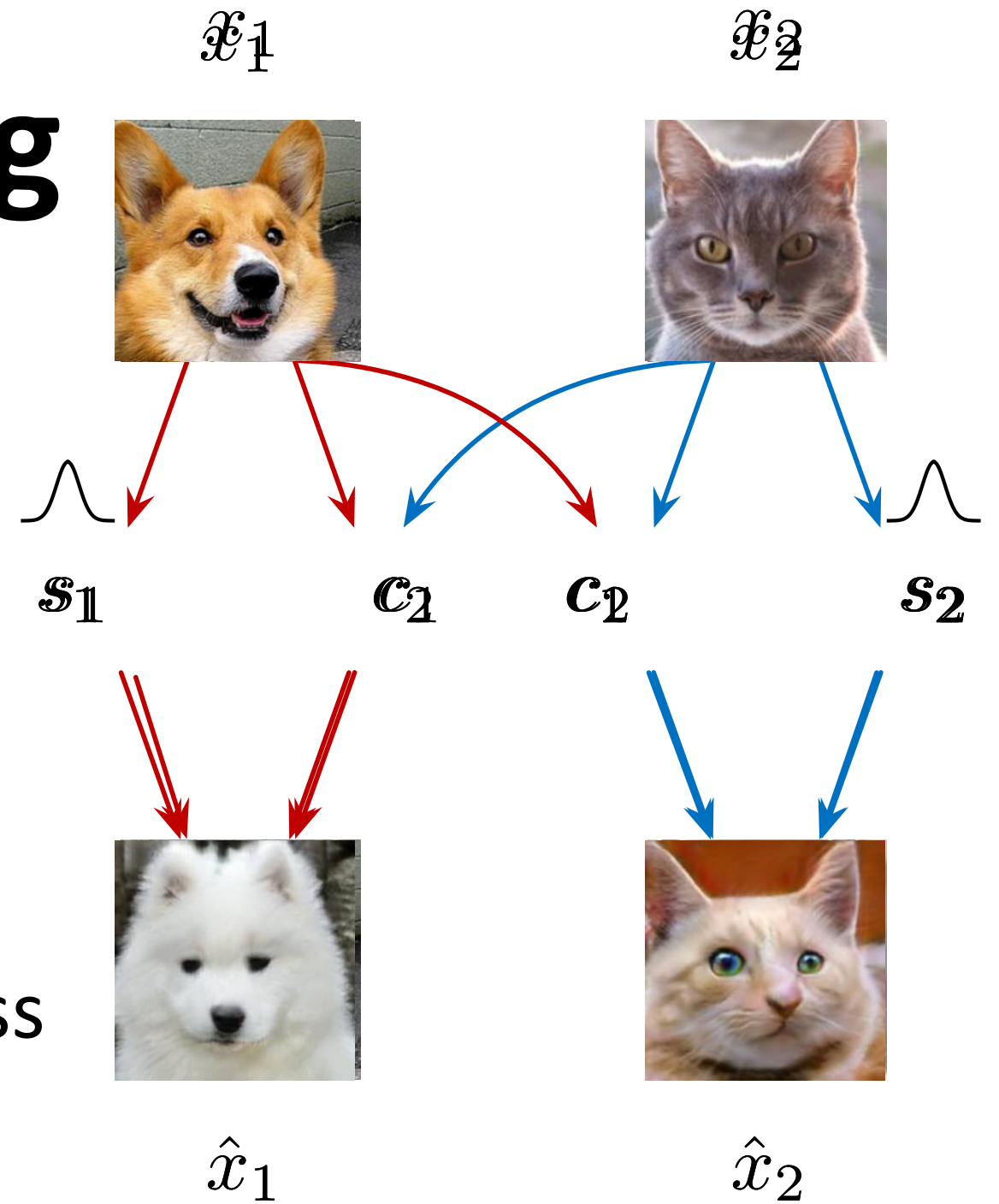
# Training

- Notations:

- $x$ : images
- $c$ : content
- $s$ : style

- Loss:

- Bidirectional reconstruction loss
  - Image reconstruction loss
  - Latent reconstruction loss
- GAN loss

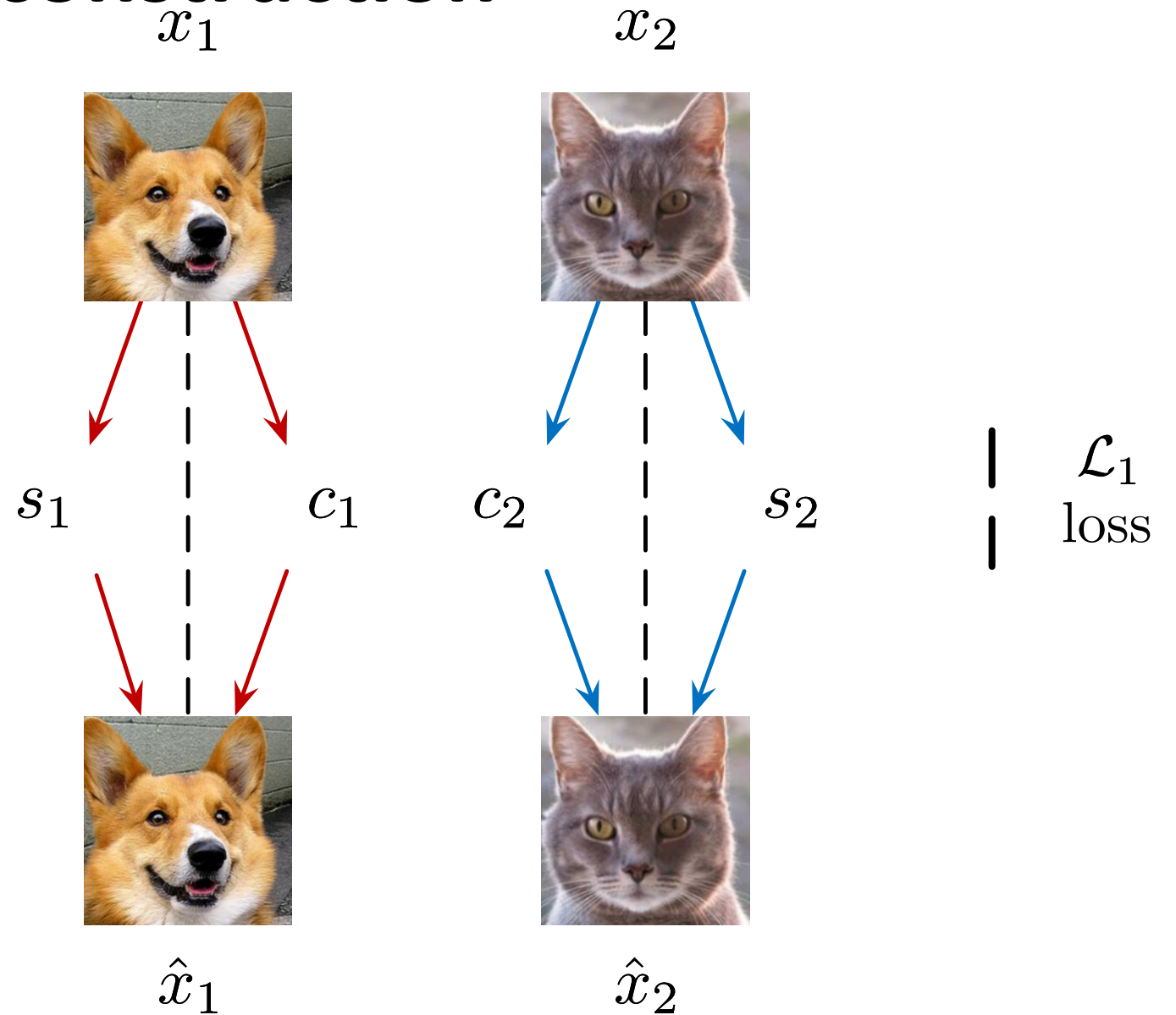


With cross-domain reconstruction

# Bidirectional Reconstruction Loss: Image Reconstruction

## Notations:

- $x$ : images
- $c$ : content
- $s$ : style

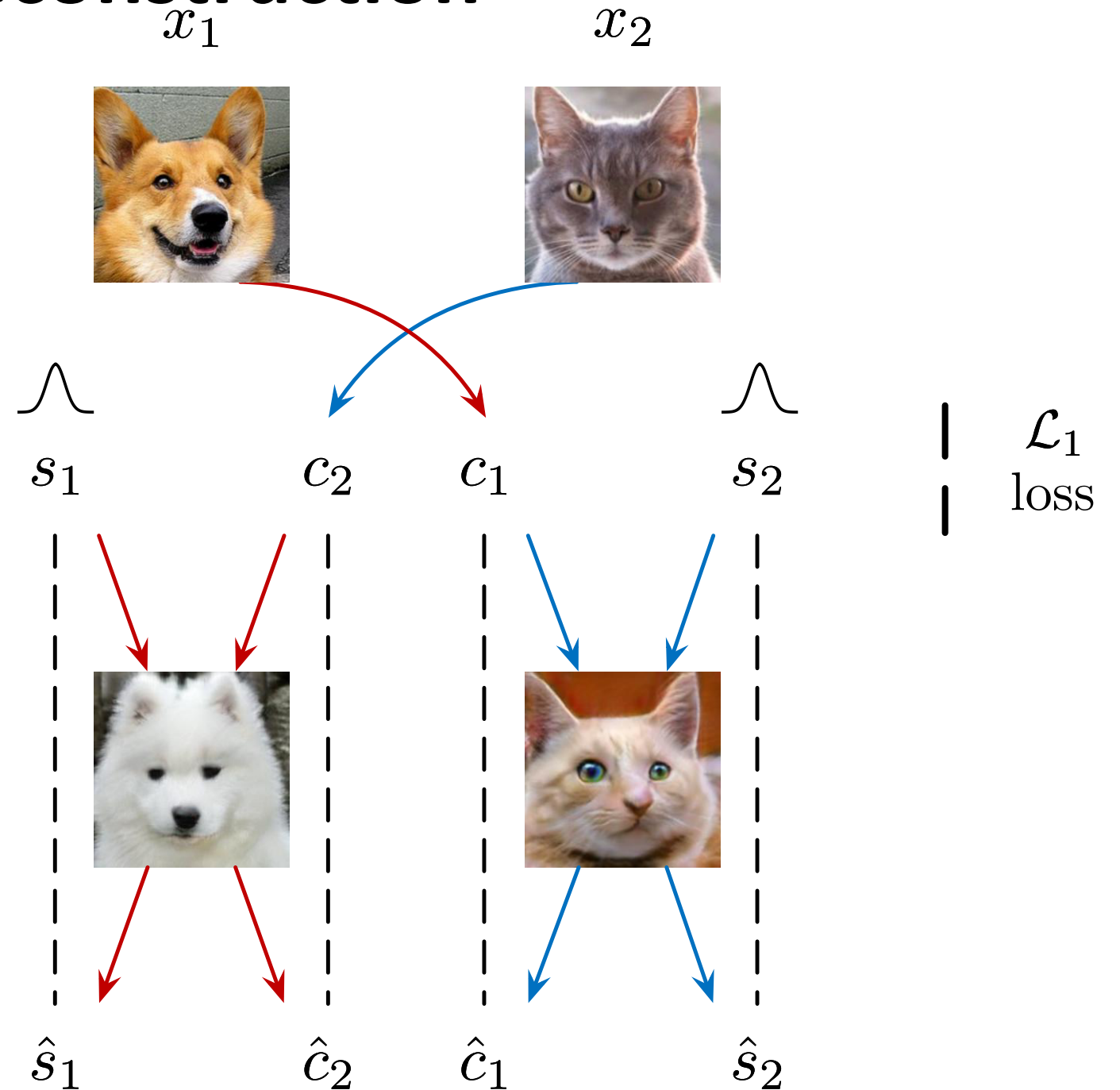




# Bidirectional Reconstruction Loss: Image Reconstruction

## Notations:

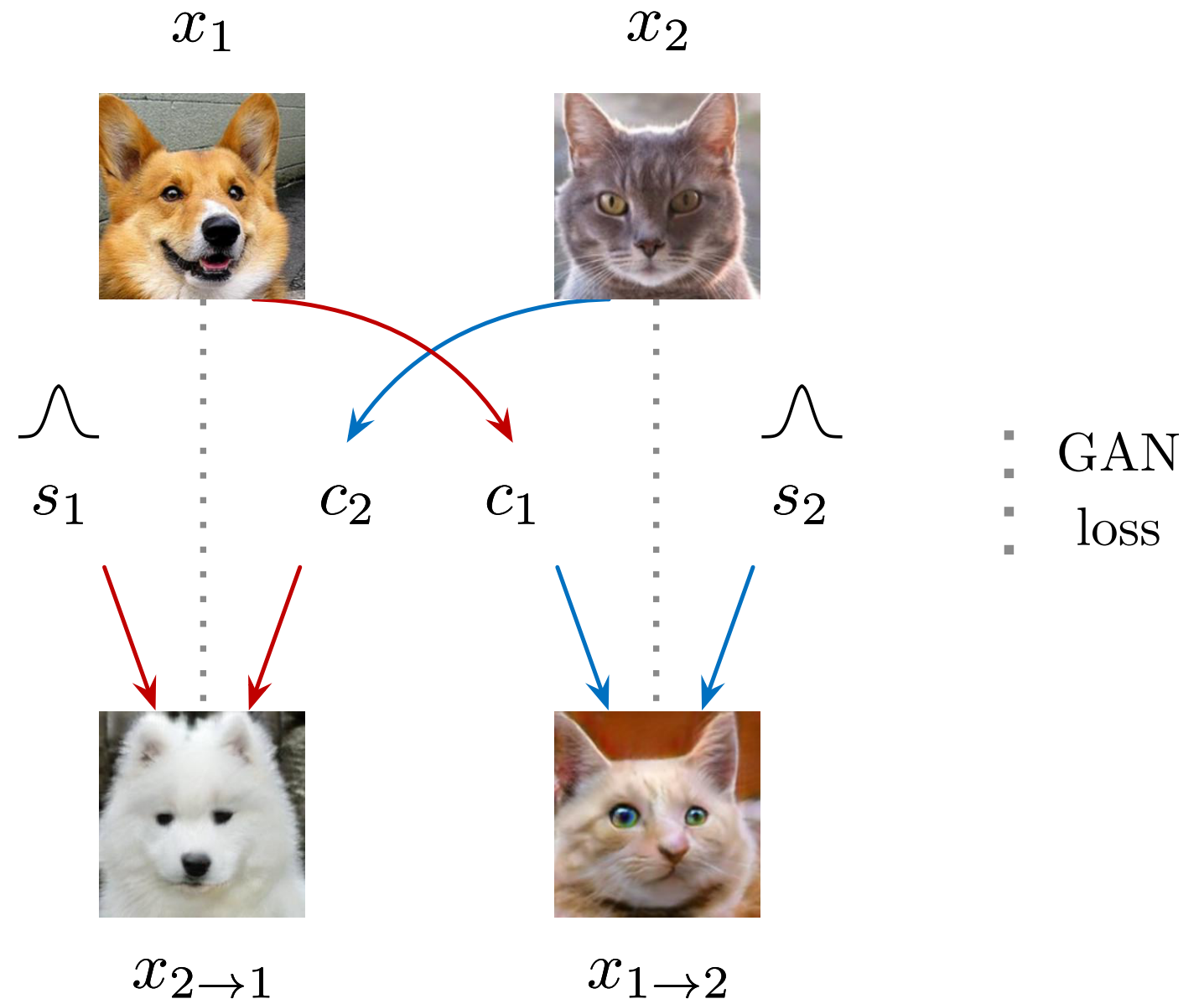
- $x$ : images
- $c$ : content
- $s$ : style



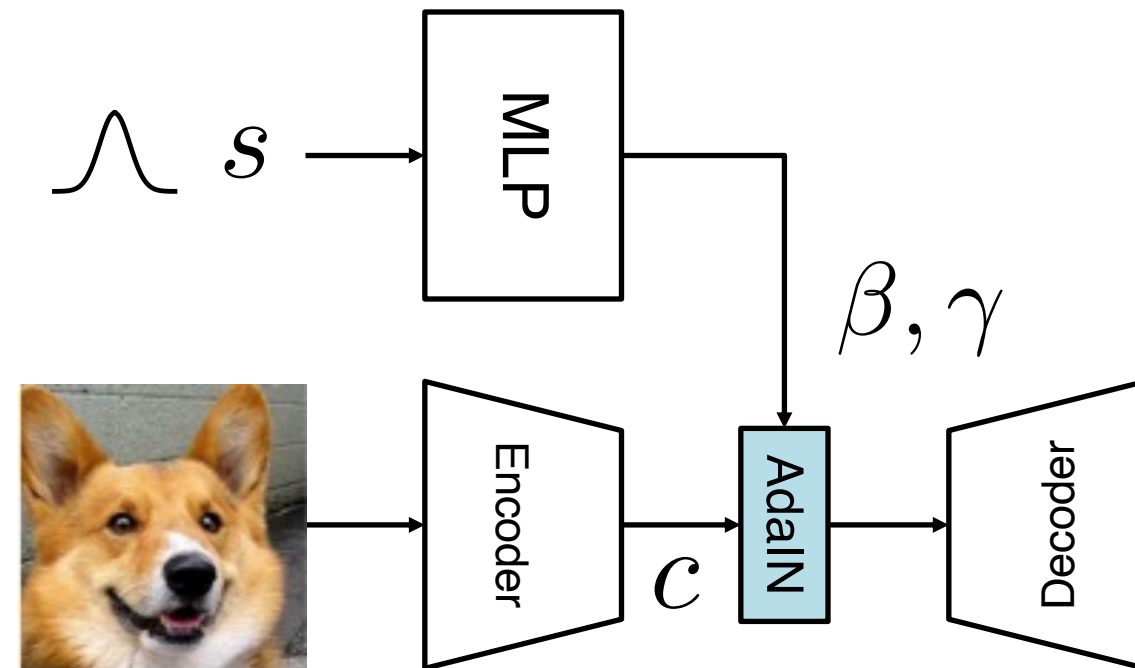
# GAN Loss

## Notations:

- $x$ : images
- $c$ : content
- $s$ : style



# AdaIN in a Generative Network

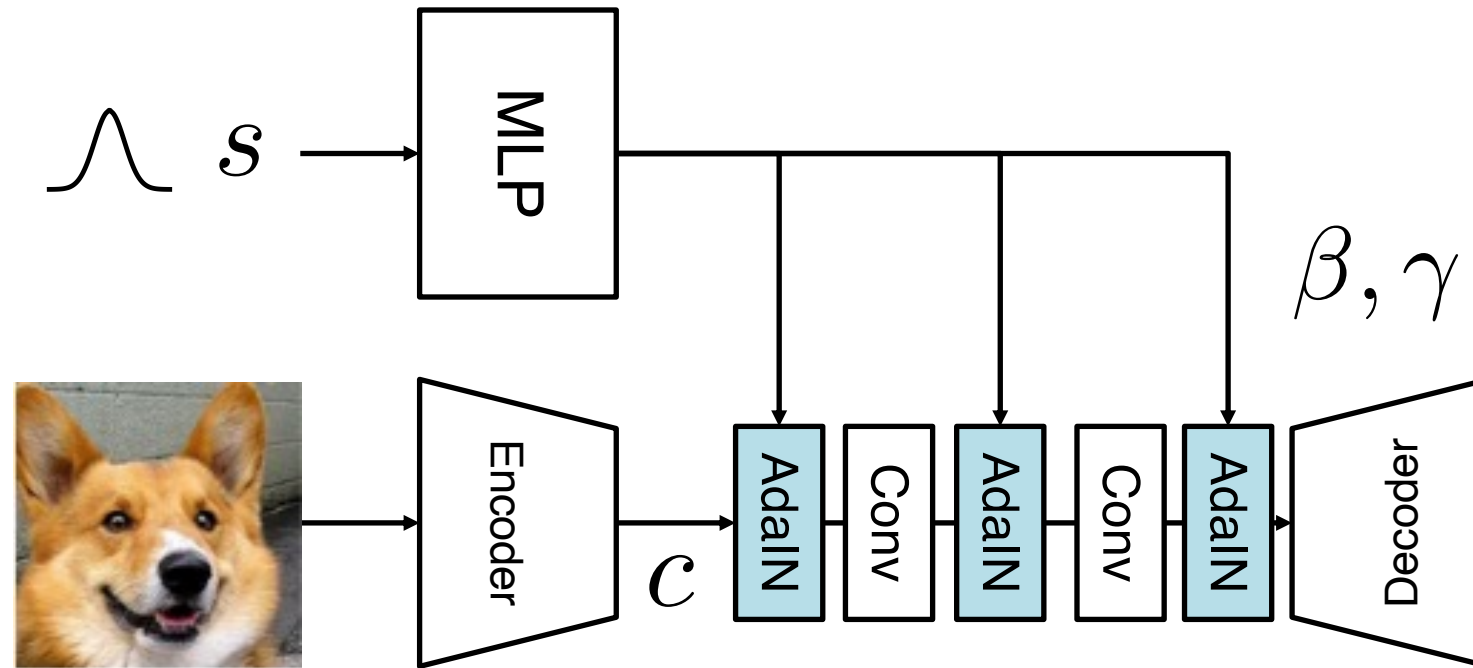


$$\text{AdaIN}(c, s) = \gamma \left( \frac{c - \mu(c)}{\sigma(c)} \right) + \beta$$

AdaIN in a generative network



# AdaIN in a Generative Network

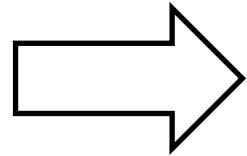


$$\text{AdaIN}(c, s) = \gamma \left( \frac{c - \mu(c)}{\sigma(c)} \right) + \beta$$

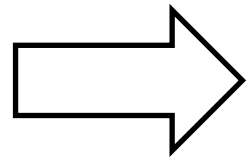
AdaIN in a generative network

# Sketches <-> Photo

Input

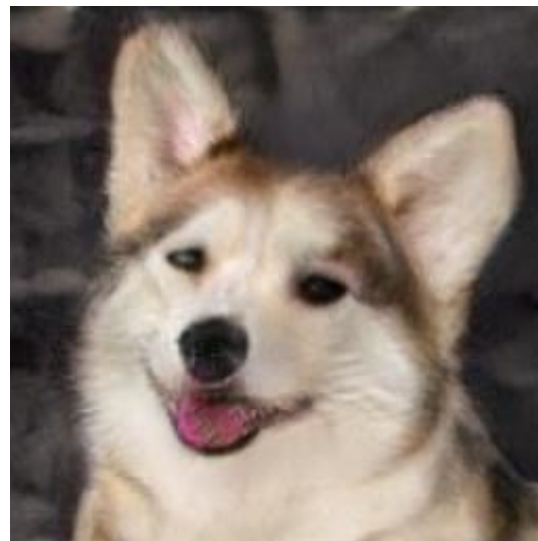
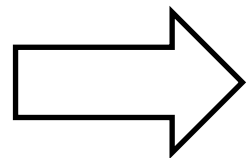


Outputs

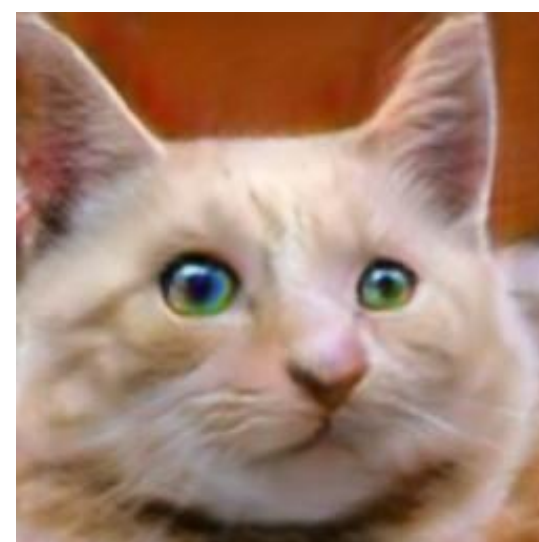
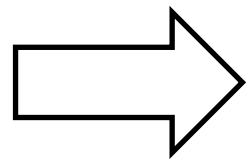


# Cats ↔ Dogs

Input



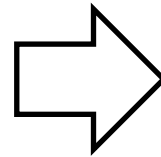
Outputs



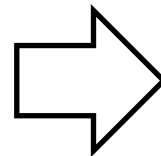


# Synthetic $\leftrightarrow$ Real

Input

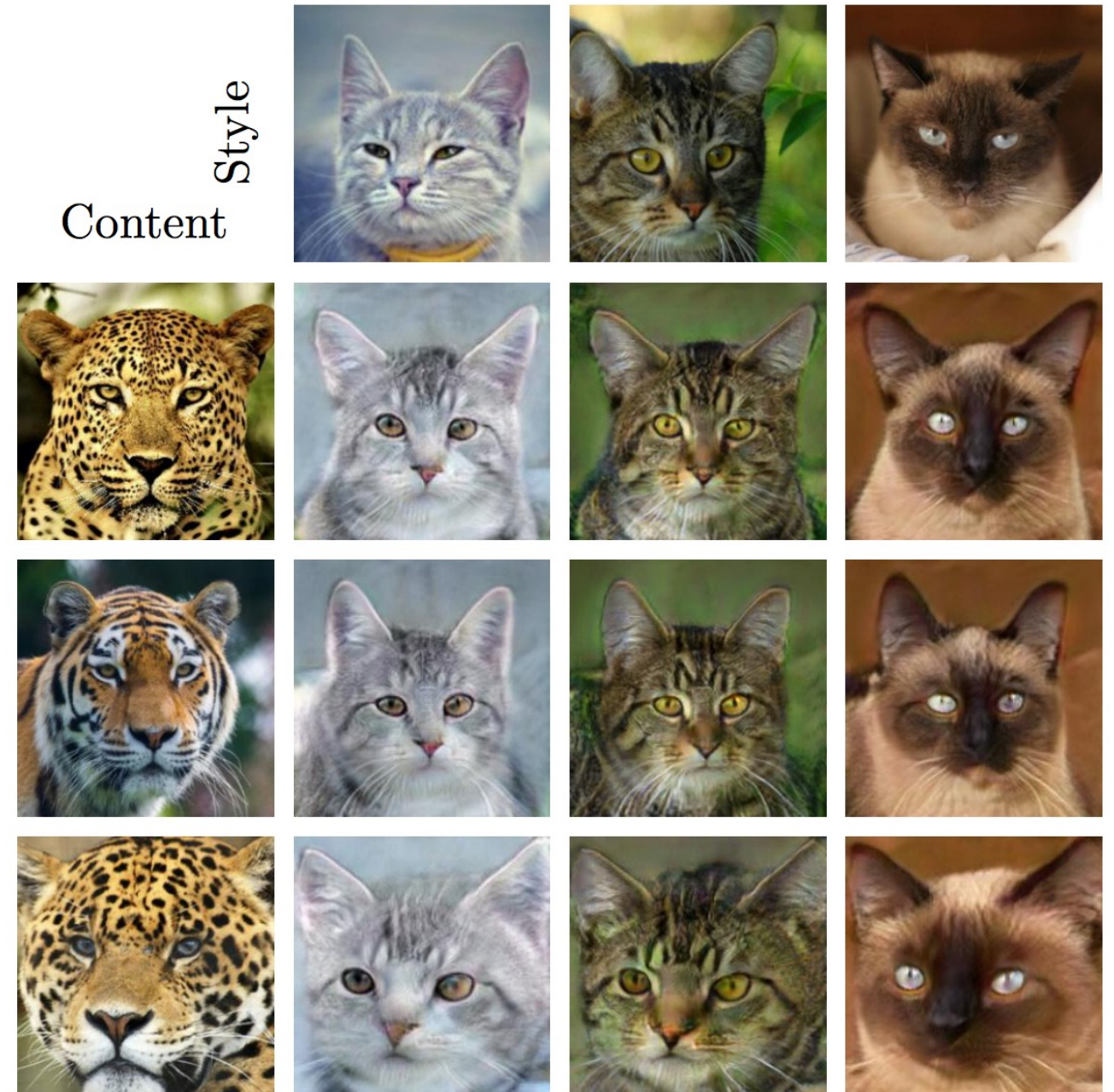


Outputs





# Example-guided Translation





# Example-guided Translation

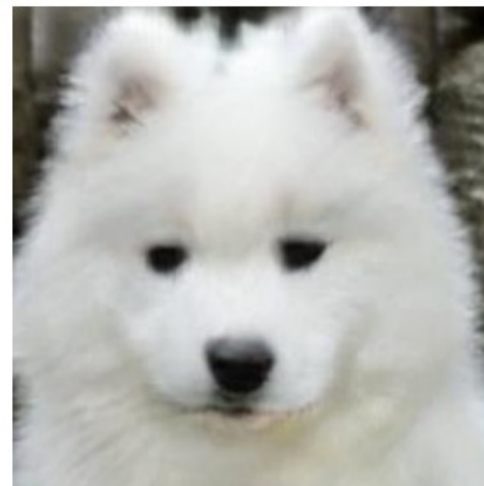
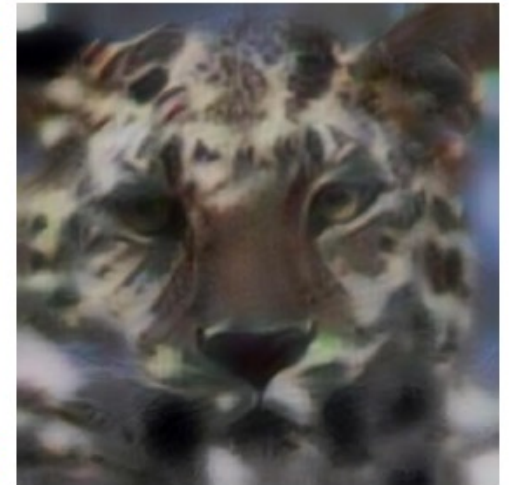
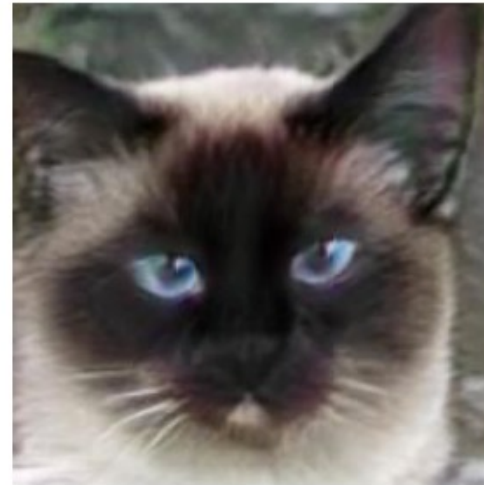
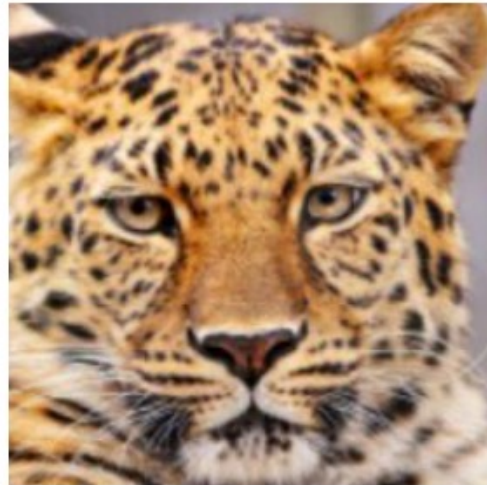
Content

Style

Ours

Gatys *et al.*

AdaIN





# Thank You!



16-726, Spring 2022

<https://learning-image-synthesis.github.io/sp22/>