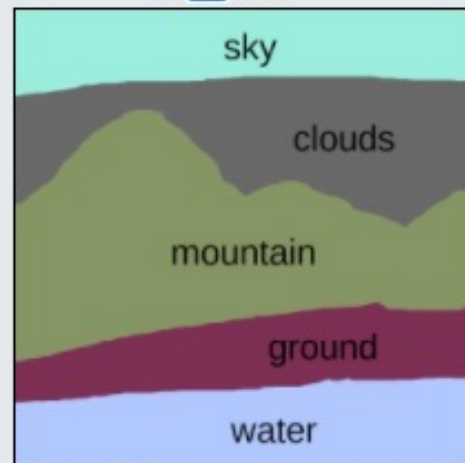# Image-to-Image Translation and Conditional Generative Models (part II)

Jun-Yan Zhu
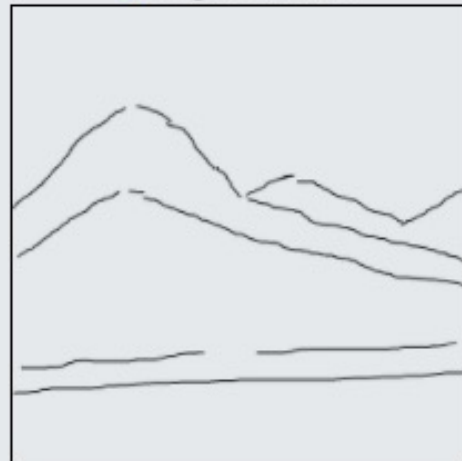
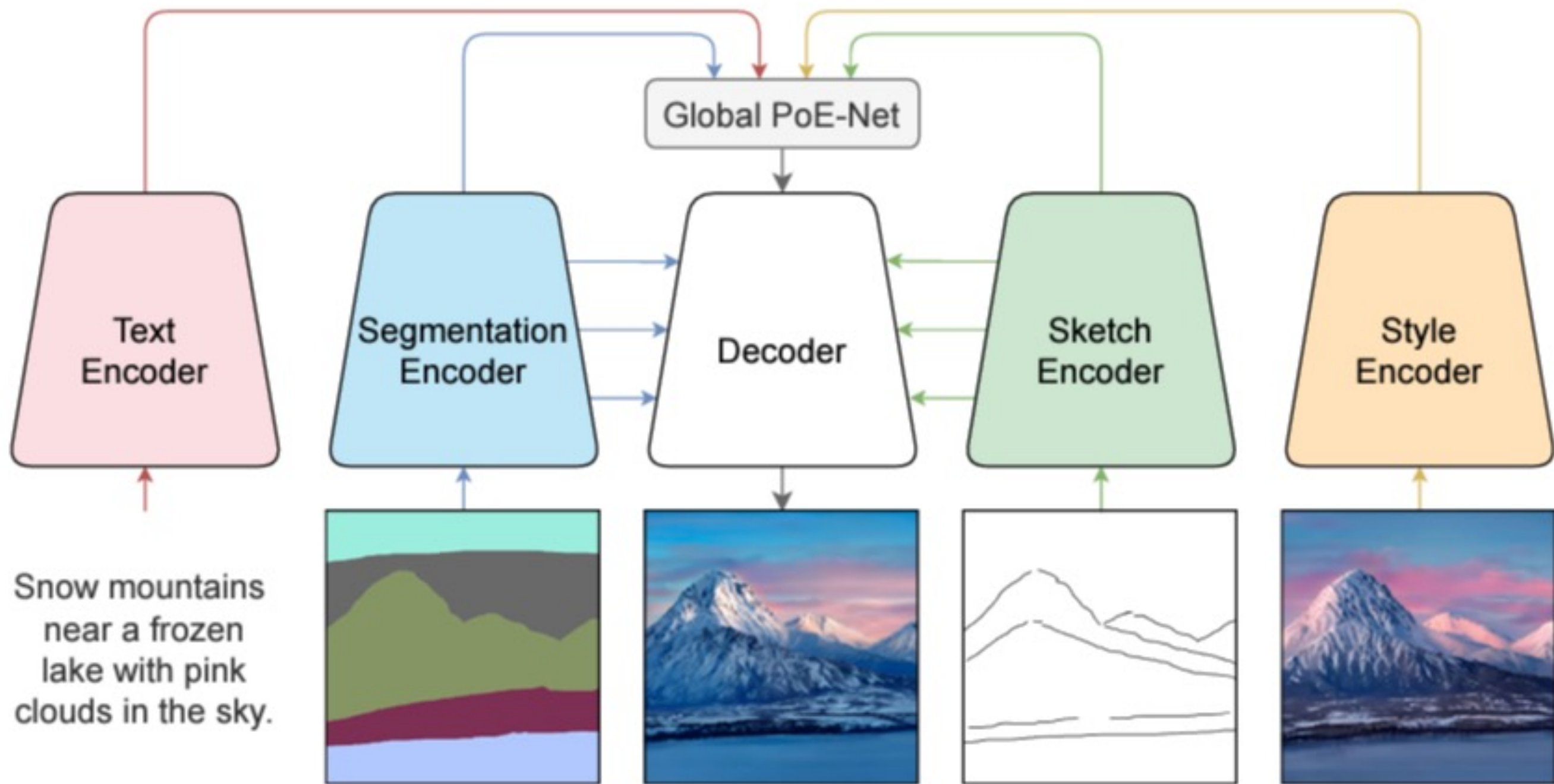16-726, Spring 2025

Many slides from Phillip Isola, Ming-Yu Liu, Xun Huang, etc.

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

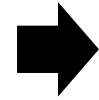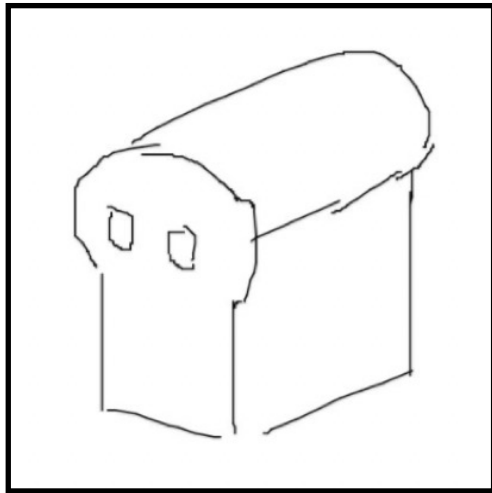Multimodal Conditional Image Synthesis with Product-of-Experts GANs [Huang et al., 2021]

# Supervised Learning Approach



Edges2cats

Image colorization

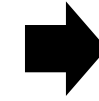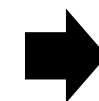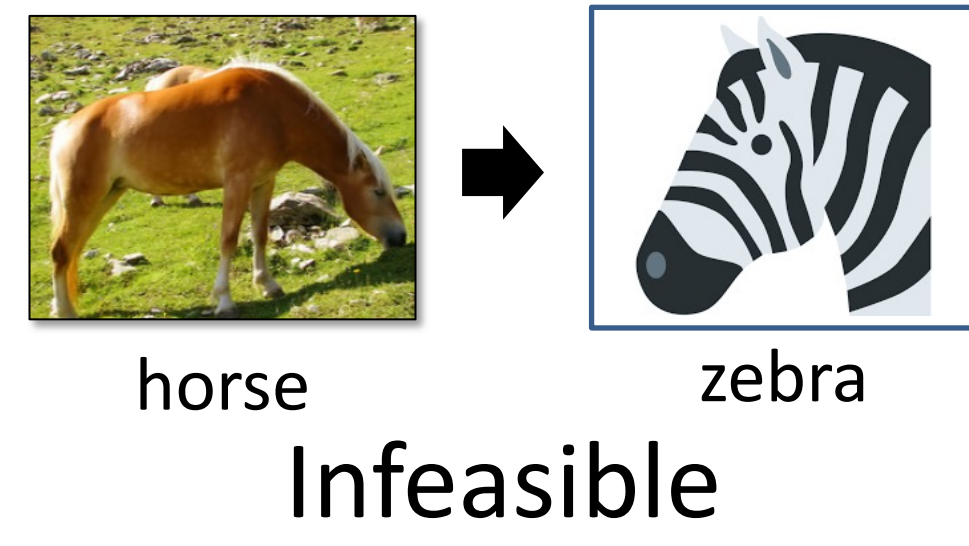Street view images

Natural outdoor images

# Supervised Learning Approach



User Input

Learning algorithm

Labeled data

Visual Content

Expensive labor

Artistic authoring

horse

zebra

Infeasible

# Supervised

$x_i$      $y_i$



# Unsupervised

$X$      $Y$

# Unsupervised Learning of $p(y \mid x)$



[Zhu*, Park*, Isola, and Efros, 2017]

# Unsupervised Learning of $p(y \mid x)$

$X$

$Y$

fake zebra

real zebra

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

$X$ → $Y$

$Y$ → $D$

Discriminator

# Unsupervised Learning of $p(y \mid x)$



- artifacts
- ignore inputs

[Goodfellow et al. 2014]

# Additional Constraint: Identity Mapping

x



Input image

Generator

G(x)

Output image

Discriminator

Real (1) or fake (0)?

**Adversarial loss**

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

**Self-Regularization loss**

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

x        G(x)

$$\Big| \quad - \quad \Big|$$

SimGAN [Shrivastava et al., 2017]

# Additional Constraint: Feature Loss



x
G(x)

G — Generator
D — Discriminator

Real (1) or fake (0)?

Input image
Output image

**Adversarial loss**

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

**Feature loss**

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

$$|F(\ \ ) - F(\ \ )|$$

x
G(x)

Input
Output

Requires F to work across two domains

DTN [Taigman et al., 2017]

# Additional Constraint: Cycle-Consistency



CycleGAN [Zhu*, Park* et al., ICCV 2017]

# Cycle-Consistent Adversarial Networks

$x$        $G(x)$        $F(G(x))$



**Adversarial loss**

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

Adversarial loss $D_Y(G(x))$

**Cycle-consistency loss**

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

Cycle-consistency loss

$$||F(G(x)) - x||_1$$

CycleGAN [Zhu*, Park* et al., ICCV 2017]

# Cycle-Consistent Adversarial Networks

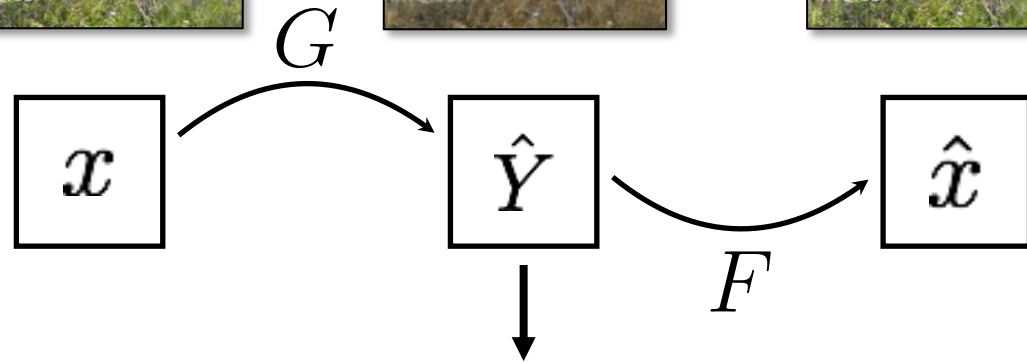$x$       $G(x)$       $F(G(x))$       $y$       $F(y)$       $G(F(y))$



$x \xrightarrow{G} \hat{Y} \xrightarrow{F} \hat{x}$

$y \xrightarrow{F} \hat{X} \xrightarrow{G} \hat{y}$

Adversarial loss $D_{\mathrm{Y}}(G(x))$

$D_X(F(y))$ Adversarial loss

Cycle-consistency loss

Cycle-consistency loss

$$||F(G(x)) - x||_1$$

$$||G(F(y)) - y||_1$$

CycleGAN [Zhu*, Park* et al., ICCV 2017]

# 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

# Why CycleGAN works

# Why CycleGAN works

**Adversarial loss**

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



$X$      $Y$

# Why CycleGAN works

**Adversarial loss**

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

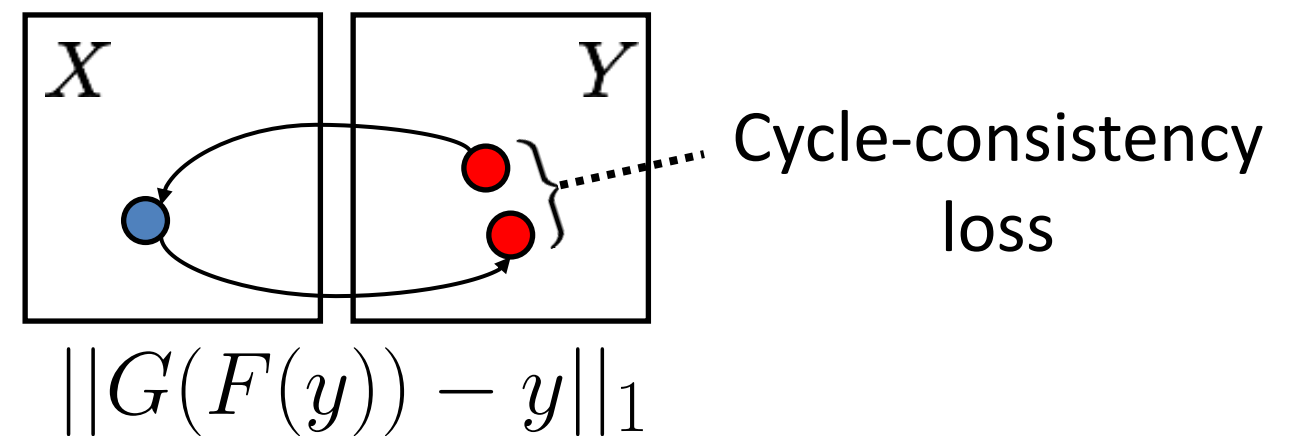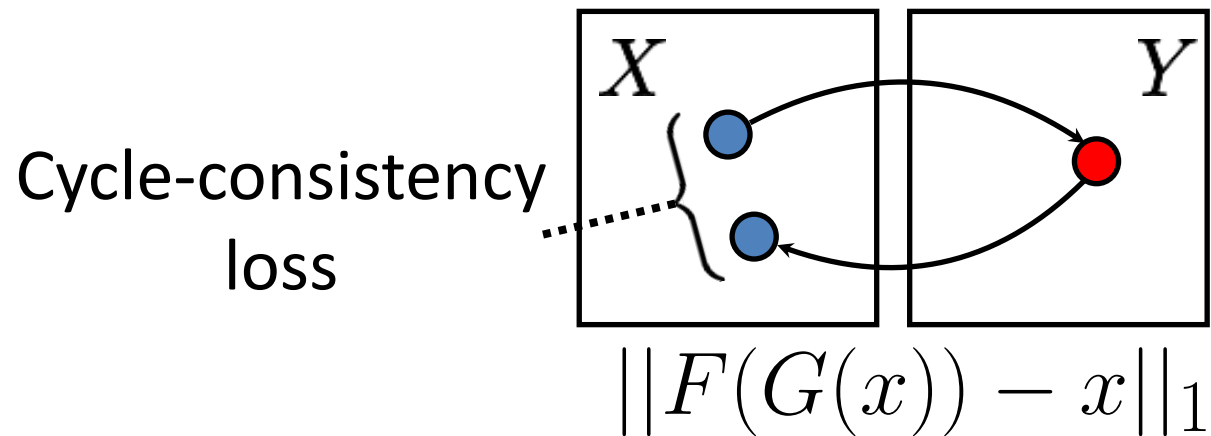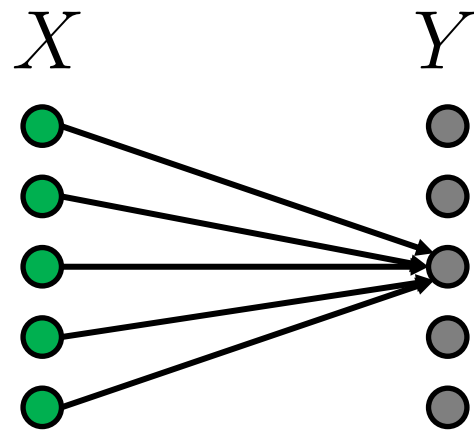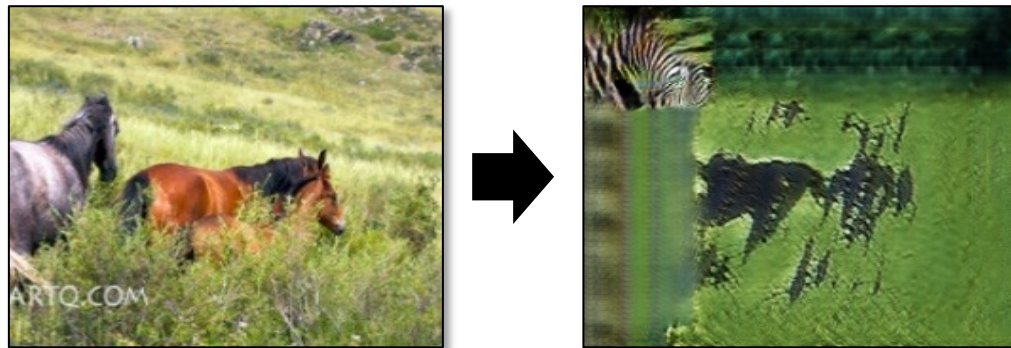**Cycle-consistency loss**

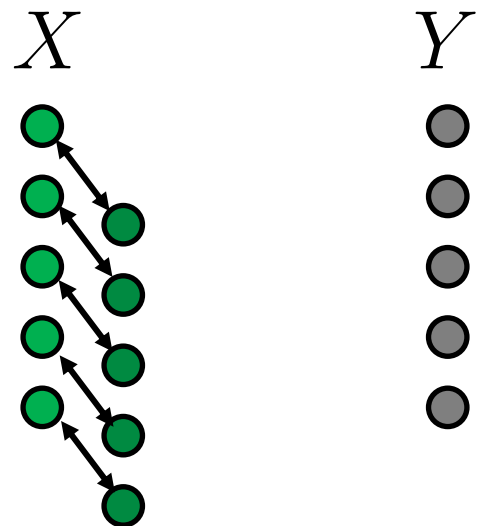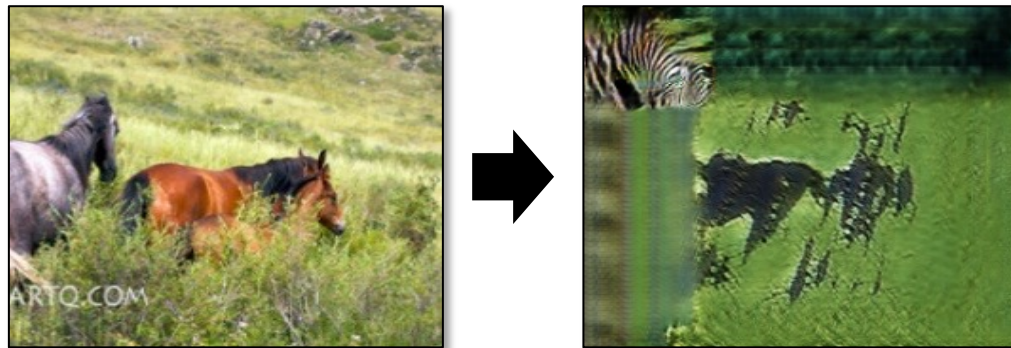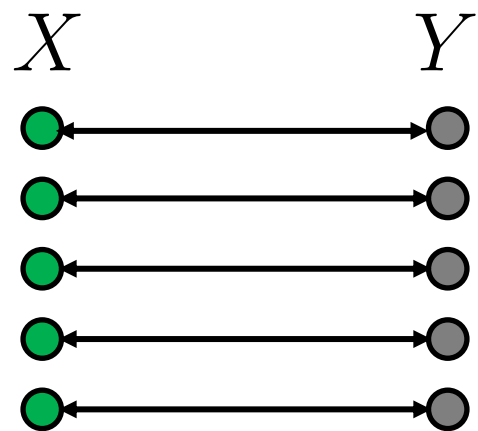$$\mathbb{E}_x ||F(G(x)) - x||_1$$



$X$    $Y$

# Why CycleGAN works

## Adversarial loss

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

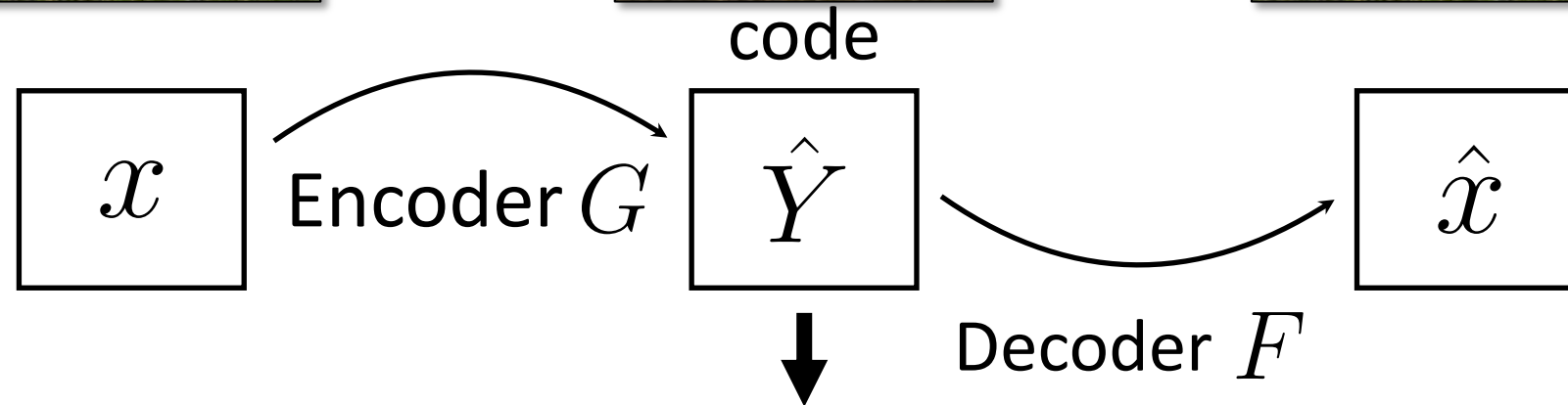## Cycle-consistency loss

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



$X \qquad Y$

## 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**

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

$x$     $G(x)$     $F(G(x))$



code

Auto-encoder
w/ domain prior

$x$   Encoder $G$   $\hat{Y}$   Decoder $F$   $\hat{x}$

Constraint: $G(x) \sim p_{data}(Y)$

[Hinton and Salakhutdinov. Science 2006]

# Why CycleGAN works

**Adversarial loss**

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

Under-constrained problem

**Cycle-consistency loss**

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

Prior of $G$



$x$      $G$      $\hat{Y}$      $F$      $\hat{x}$

## 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



$P \circ G$

$F \circ P^{-1}$

flip the image again

## Invertible Perturbation

**Adversarial loss**: images are horizontally symmetric

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

# Style and Content Disentanglement

# Style and Content Separation



**A** Classification — Domain Adaptation

**B** Extrapolation — Paired Image-to-Image Translation

**C** Translation — 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**

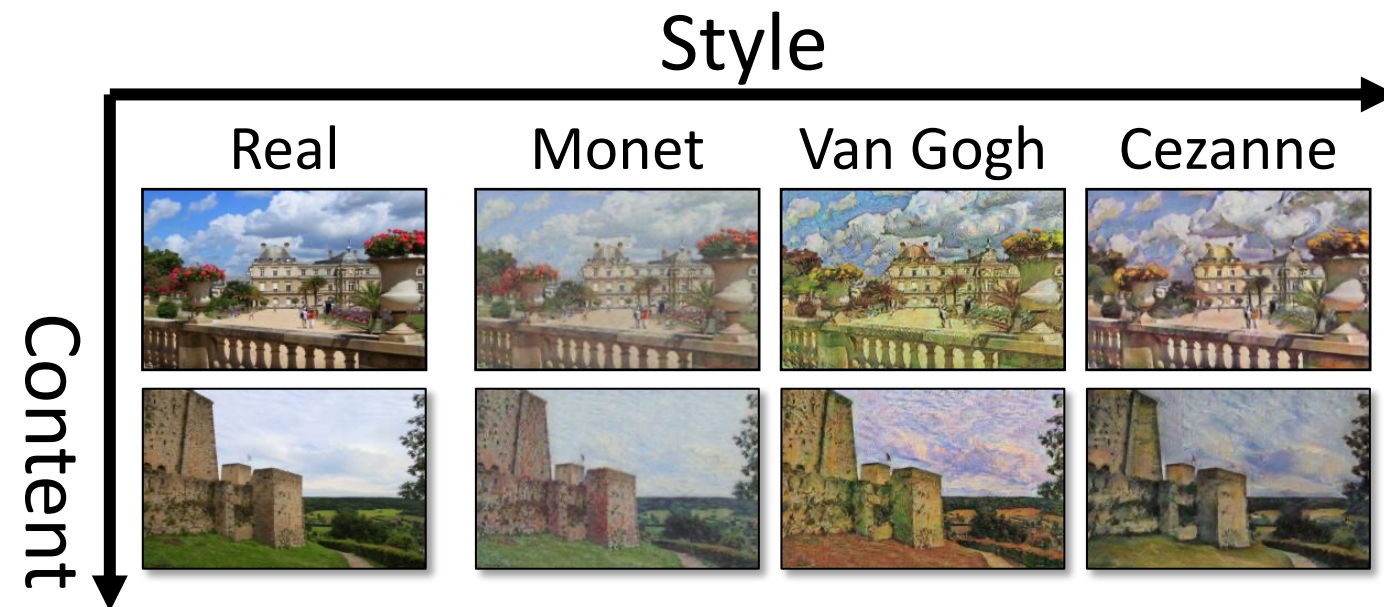$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$


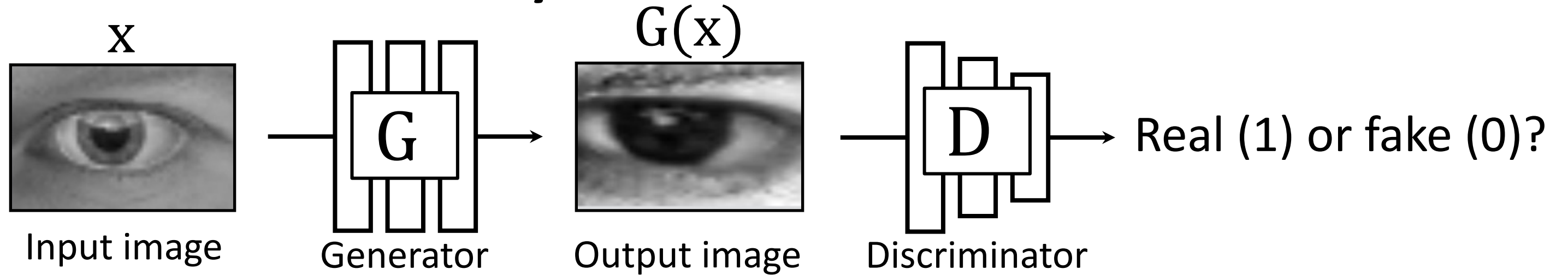
$p(x) \to p(y)$ change **style**

**Cycle-consistency loss**

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



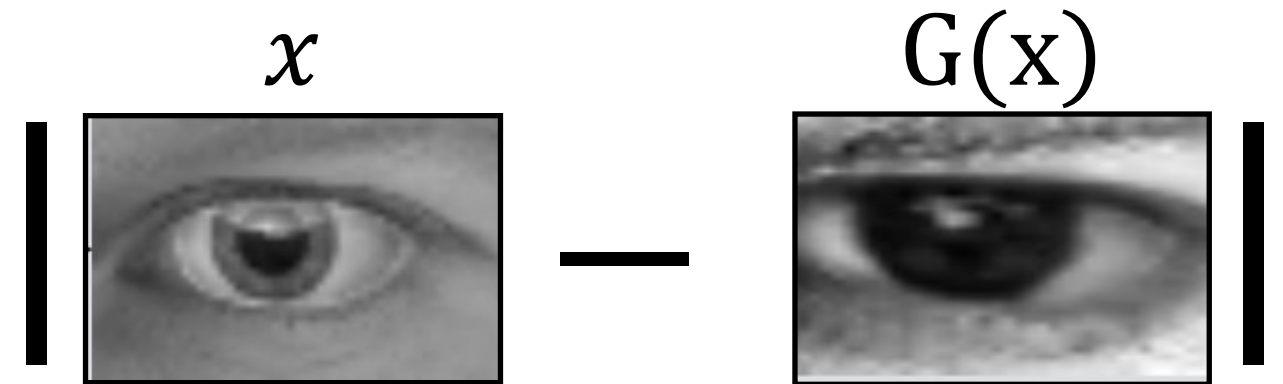Bidirectional:  preserve **content**

Style



Real    Monet    Van Gogh    Cezanne

Content

Separating Style and Content
[Tenenbaum and Freeman 1996]

# Style and Content

x



Input image

Generator

G(x)



Output image

Discriminator
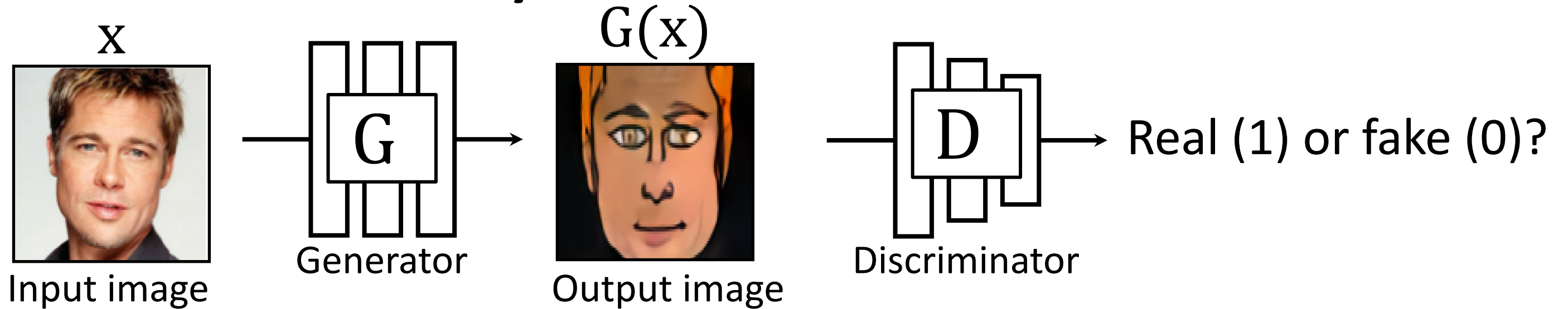
Real (1) or fake (0)?

**Adversarial loss (change style)**

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

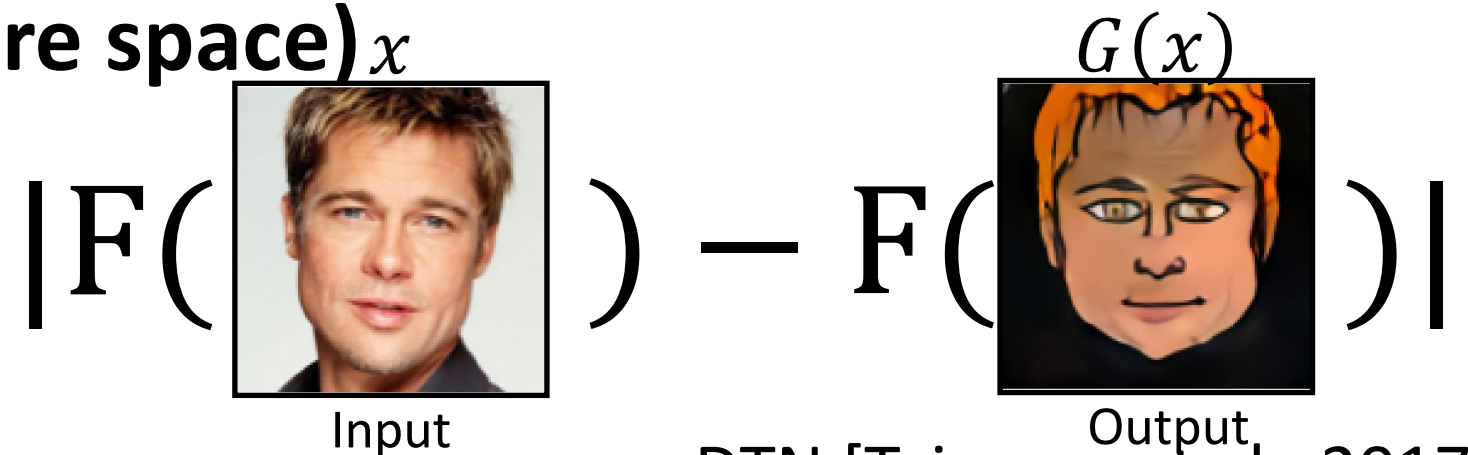**L1 loss (preserve content in pixel space)**

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

$x$



$G(x)$



SimGAN [Shrivastava et al., 2017]

# Style and Content



x

G(x)

Real (1) or fake (0)?

Generator

Input image

Output image

Discriminator

**Adversarial loss (change style)**

$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$
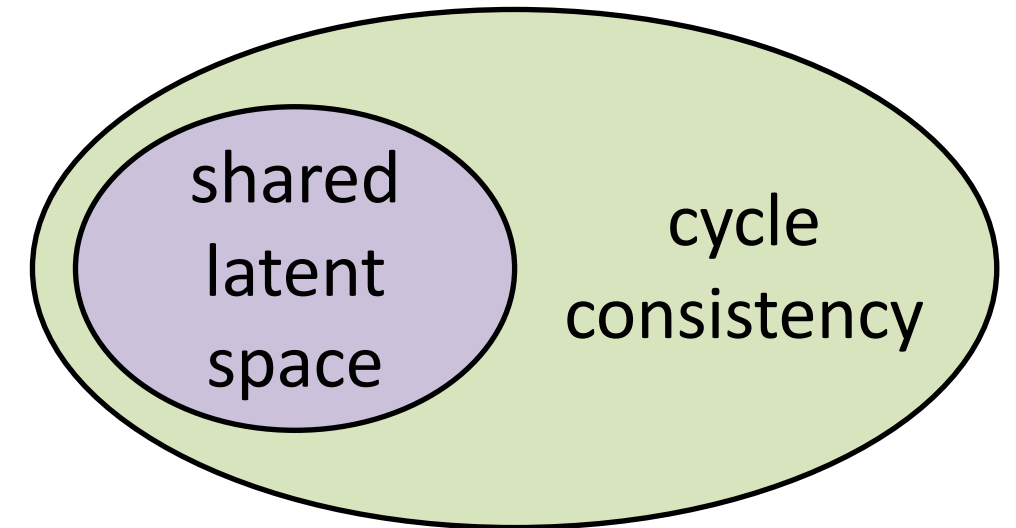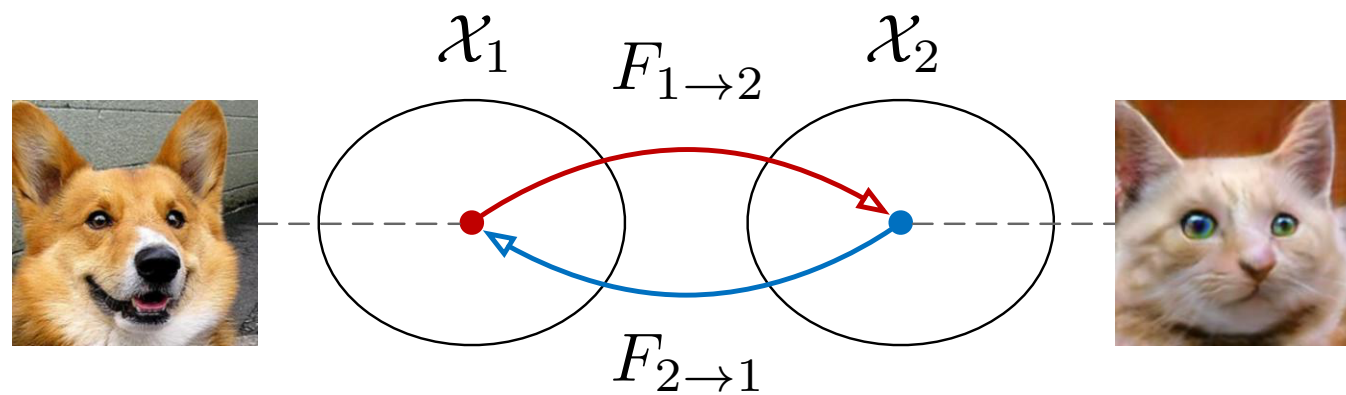
**Feature loss (Preserve content in feature space)**
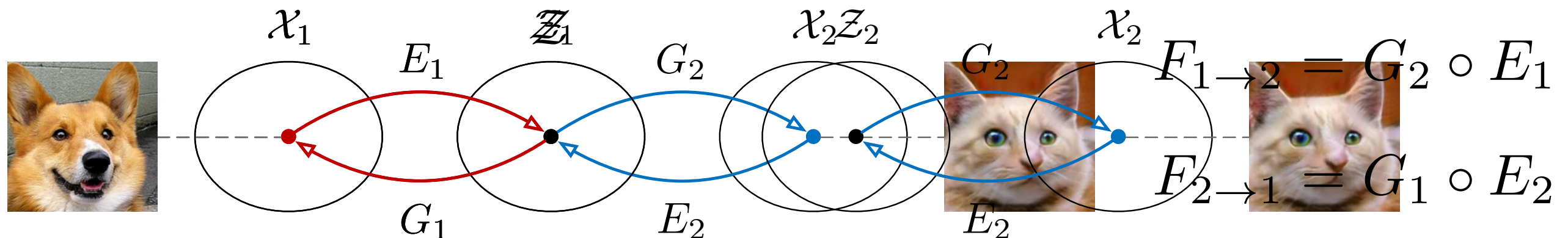
$$\mathbb{E}_x || F(G(x)) - F(x) ||$$

$$| F(\quad) - F(\quad) |$$

Input

Output

DTN [Taigman et al., 2017]

# CycleGAN and UNIT

- CycleGAN (**cycle consistency**)



$\mathcal{X}_1$ $\quad F_{1\to 2}$ $\quad \mathcal{X}_2$

$F_{2\to 1}$

shared latent space

cycle consistency

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

shared latent space $\implies$ cycle consistency



$\mathcal{X}_1$ $\quad E_1$ $\quad \mathbb{Z}_1$ $\quad G_2$ $\quad \mathcal{X}_2\mathcal{Z}_2$ $\quad G_2$ $\quad \mathcal{X}_2$

$G_1$ $\quad E_2$ $\quad E_2$

$F_{1\to 2} = G_2 \circ E_1$

$F_{2\to 1} = G_1 \circ E_2$

# Disentangling the Latent Space

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

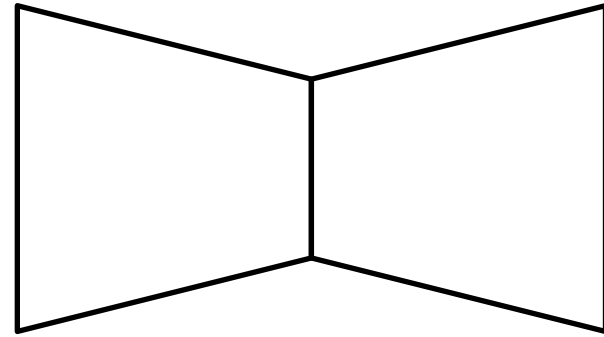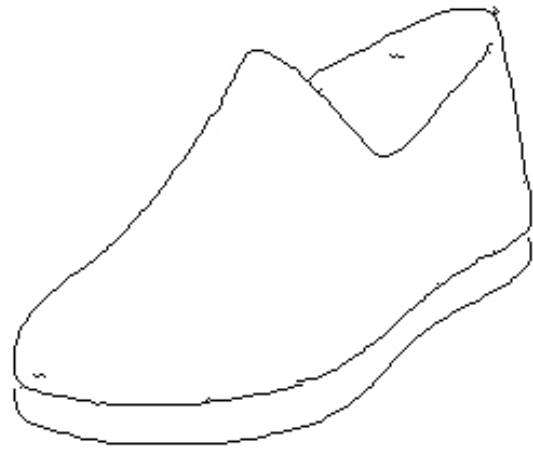# Disentangling the Latent Space

- 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**
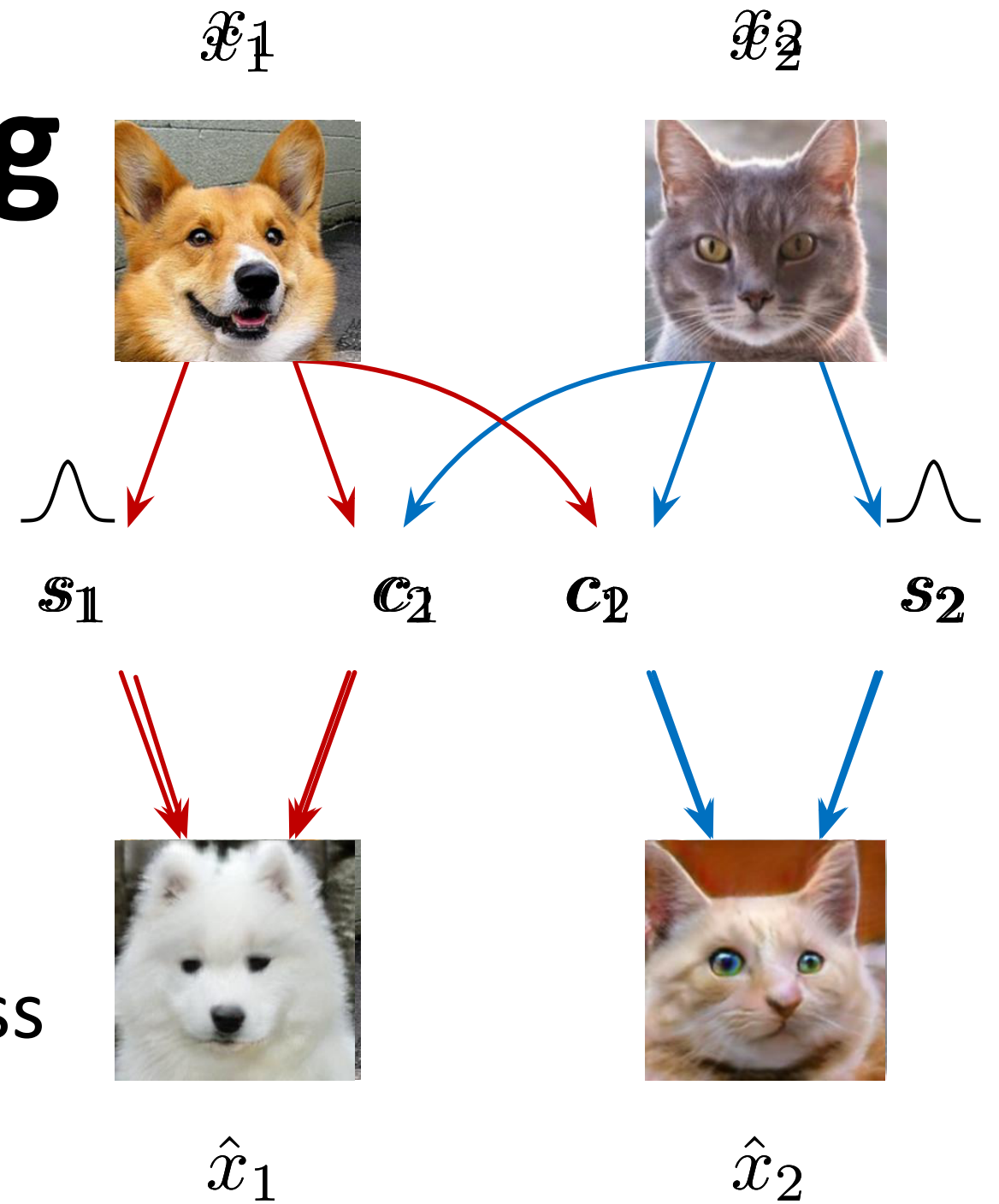
# Unimodality

# Towards Multimodality

# Training



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

- Loss:
  - Bidirectional reconstruction loss
    - Image reconstruction loss
    - Latent reconstruction loss
  - GAN loss

# Bidirectional Reconstruction Loss:
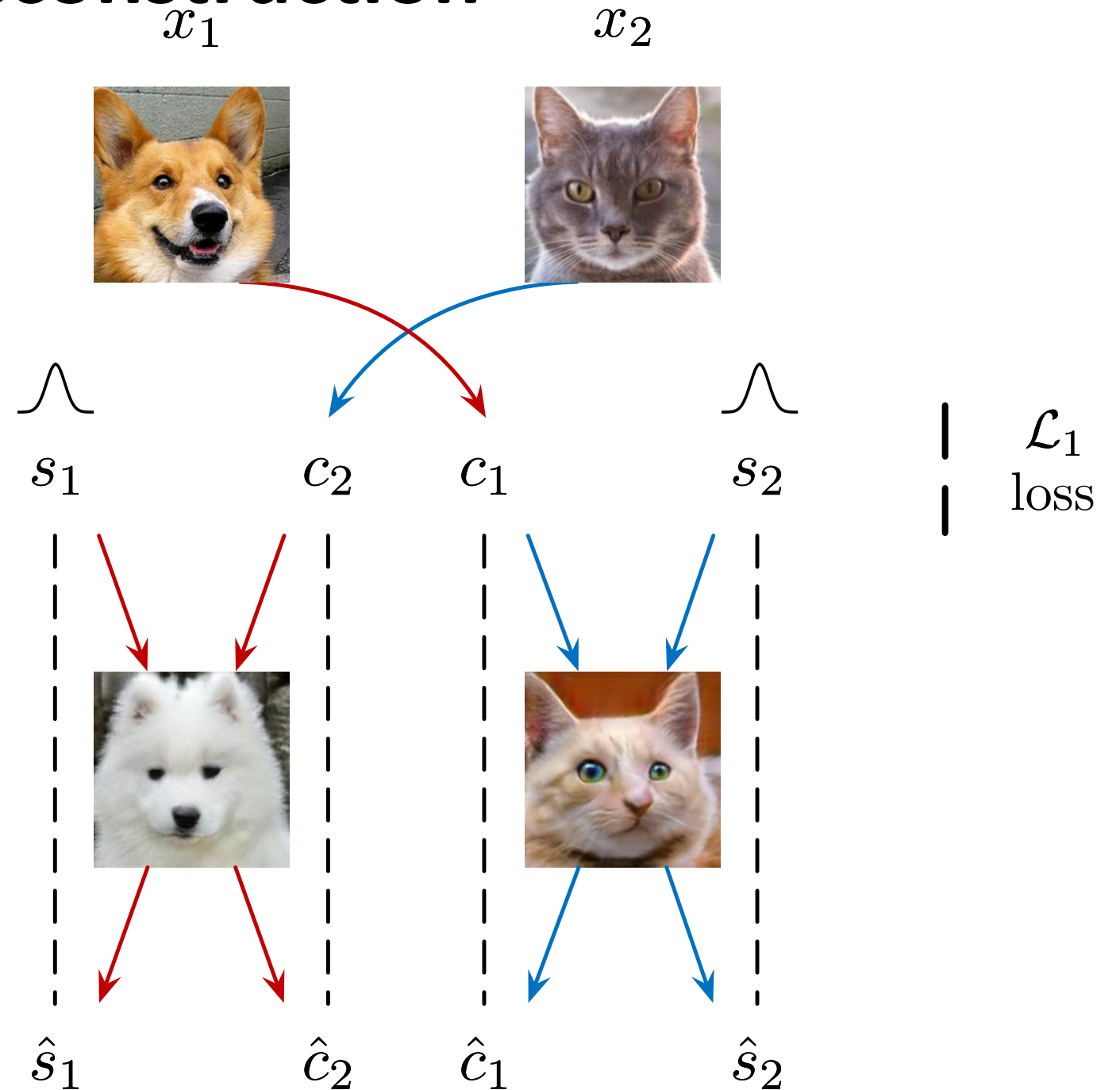# Image Reconstruction

Notations:

$-$ $x$: images

$-$ $c$: content

$-$ $s$: style

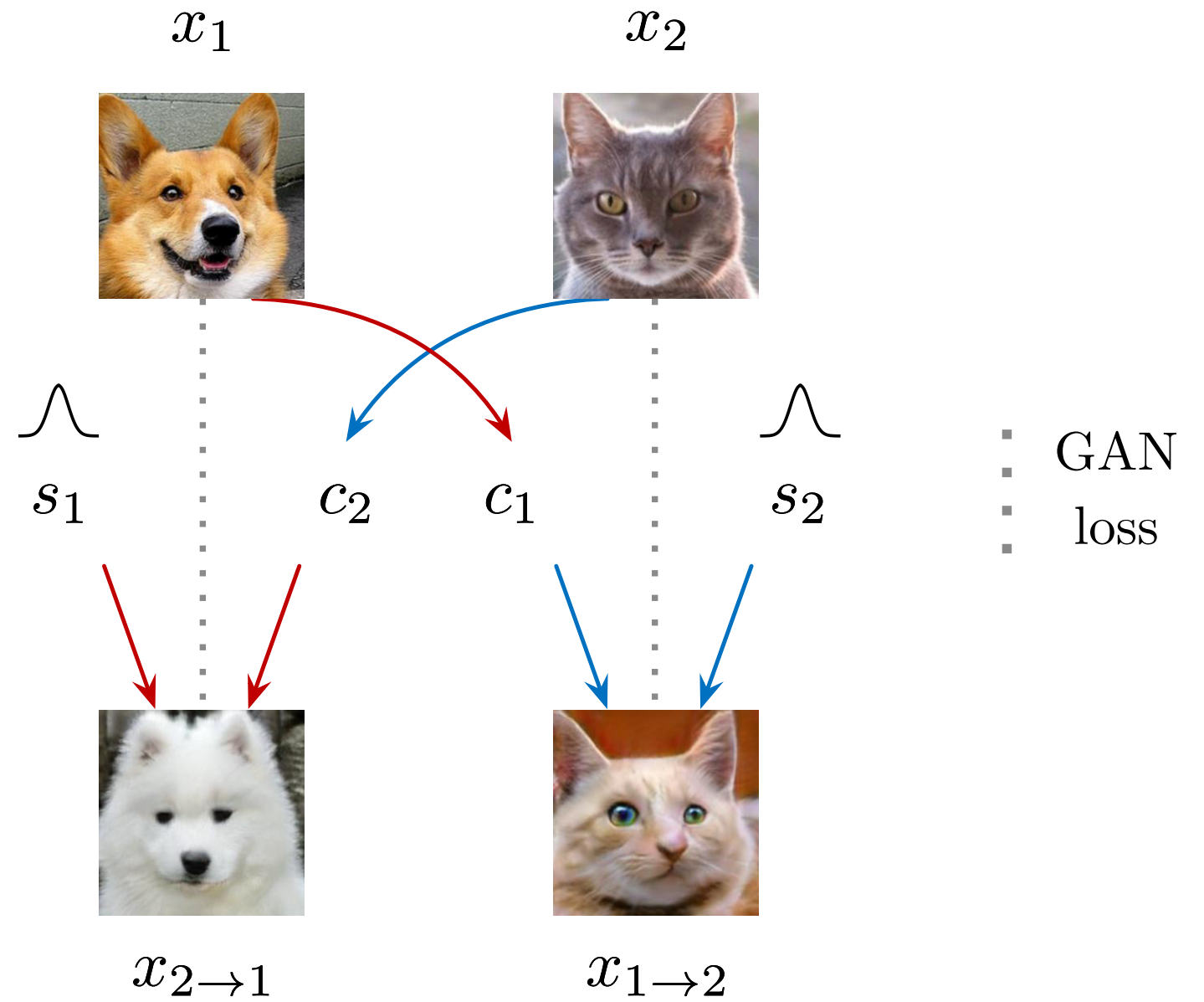# Bidirectional Reconstruction Loss:
# Image Reconstruction
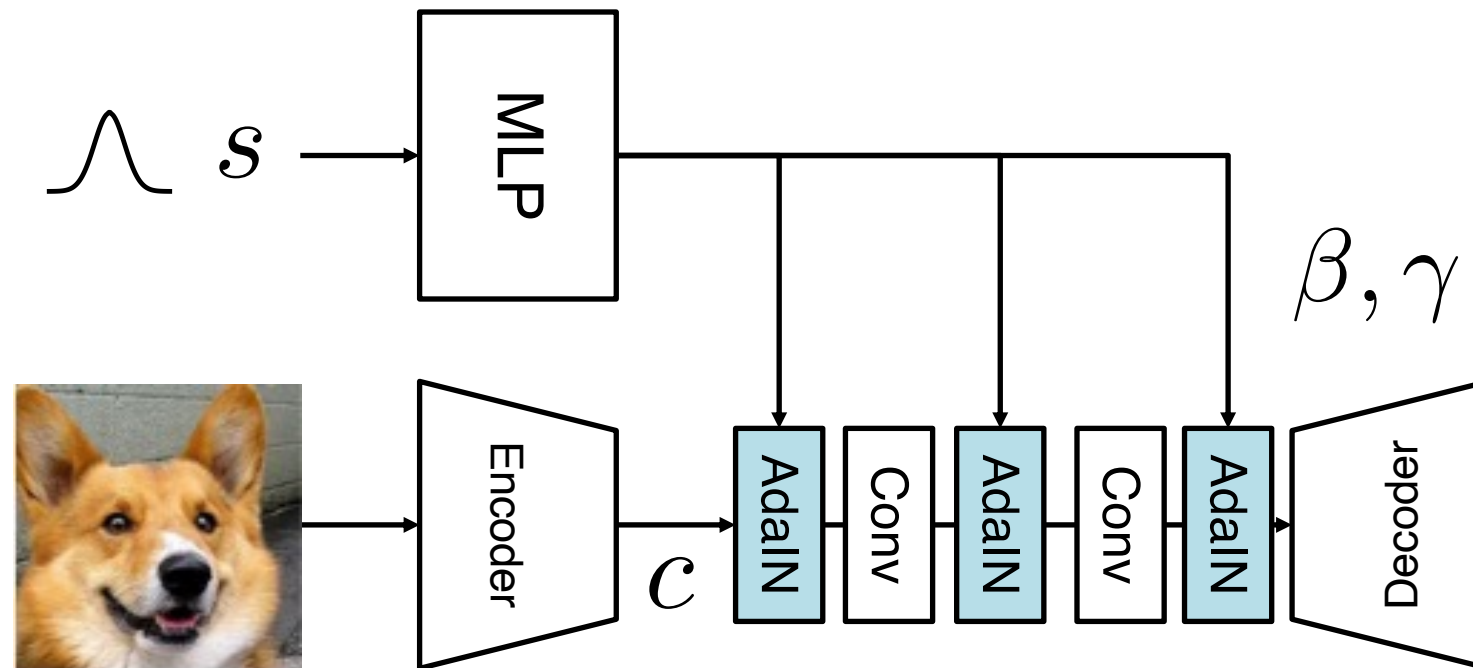
Notations:

— $x$: images

— $c$: content

— $s$: style



$x_1$

$x_2$

$s_1$     $c_2$     $c_1$     $s_2$

$\mathcal{L}_1$ loss

$\hat{s}_1$     $\hat{c}_2$     $\hat{c}_1$     $\hat{s}_2$

# GAN Loss

Notations:

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

# AdaIN in a Generative Network



$$\mathrm{AdaIN}(c, s) = \gamma \left( \frac{c - \mu(c)}{\sigma(c)} \right) + \beta$$
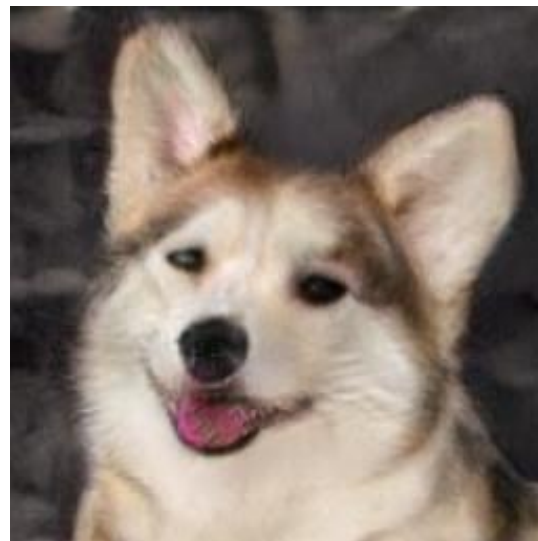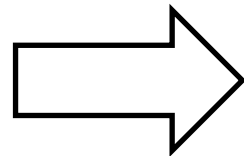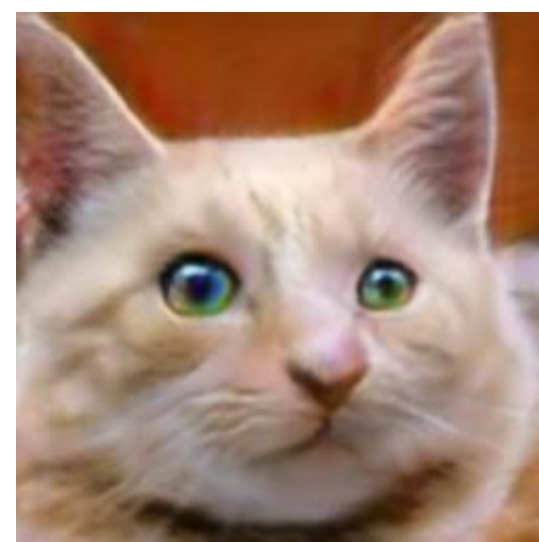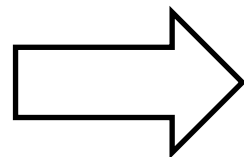
AdaIN in a generative network

# Sketches <-> Photo

Input

Outputs

# Cats ↔ Dogs

Input

Outputs
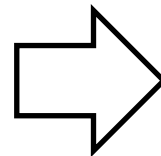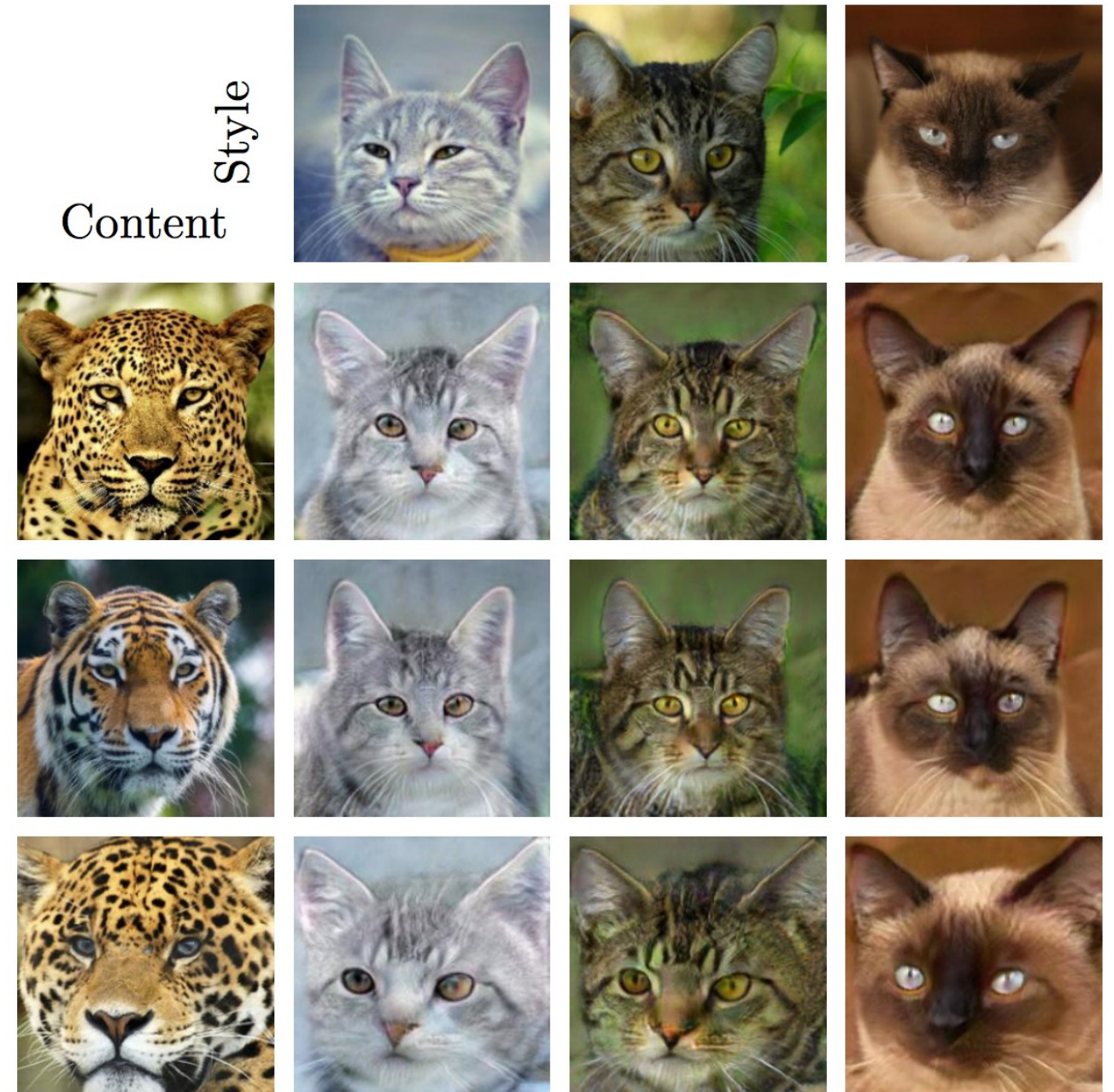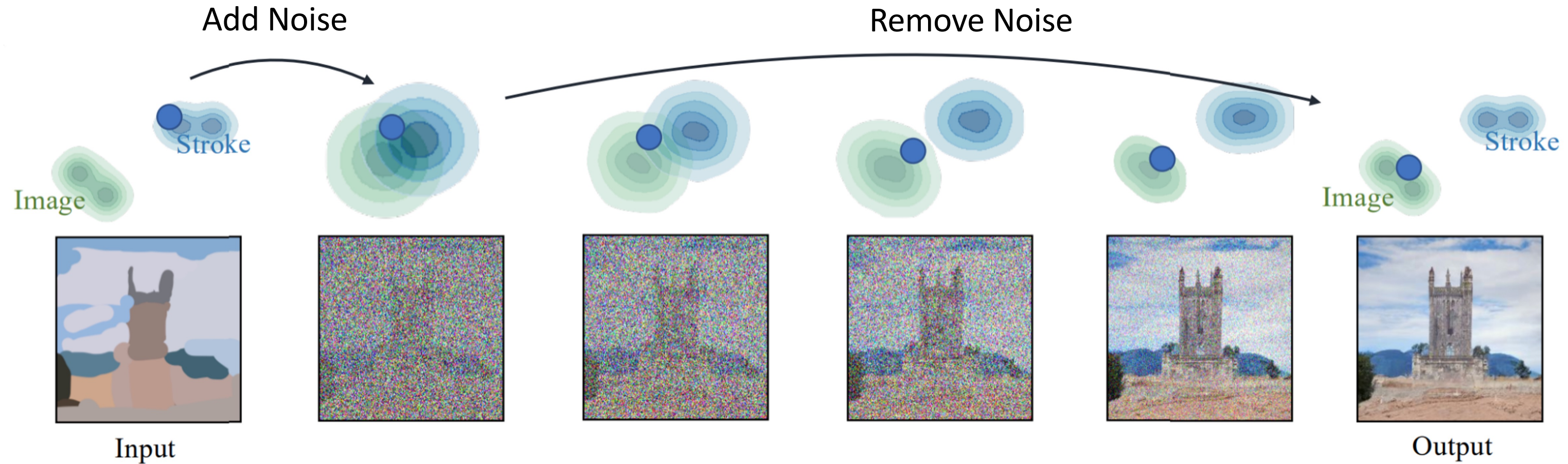
# Synthetic ↔ Real

Input

Outputs

# Example-guided Translation

# Image-to-Image Translation with Diffusion Models

# Guided Image Synthesis

SDEdit (https://arxiv.org/abs/2108.01073) recipe: diffuse → denoise
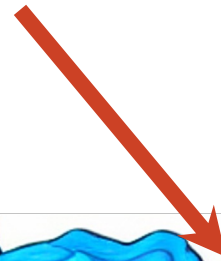
# Guided Image Synthesis



*input*

# "Upgrade" your child's artwork

# abstract art from photos



original post by u/Pereulkov
https://www.reddit.com/r/StableDiffusion/comments/xhhyad/i_made_abstract_art_from_my_photos/

# Thank You!



16-726, Spring 2025

https://learning-image-synthesis.github.io/