# Conditional GANs, Image-to-Image Translation

## Jun-Yan Zhu
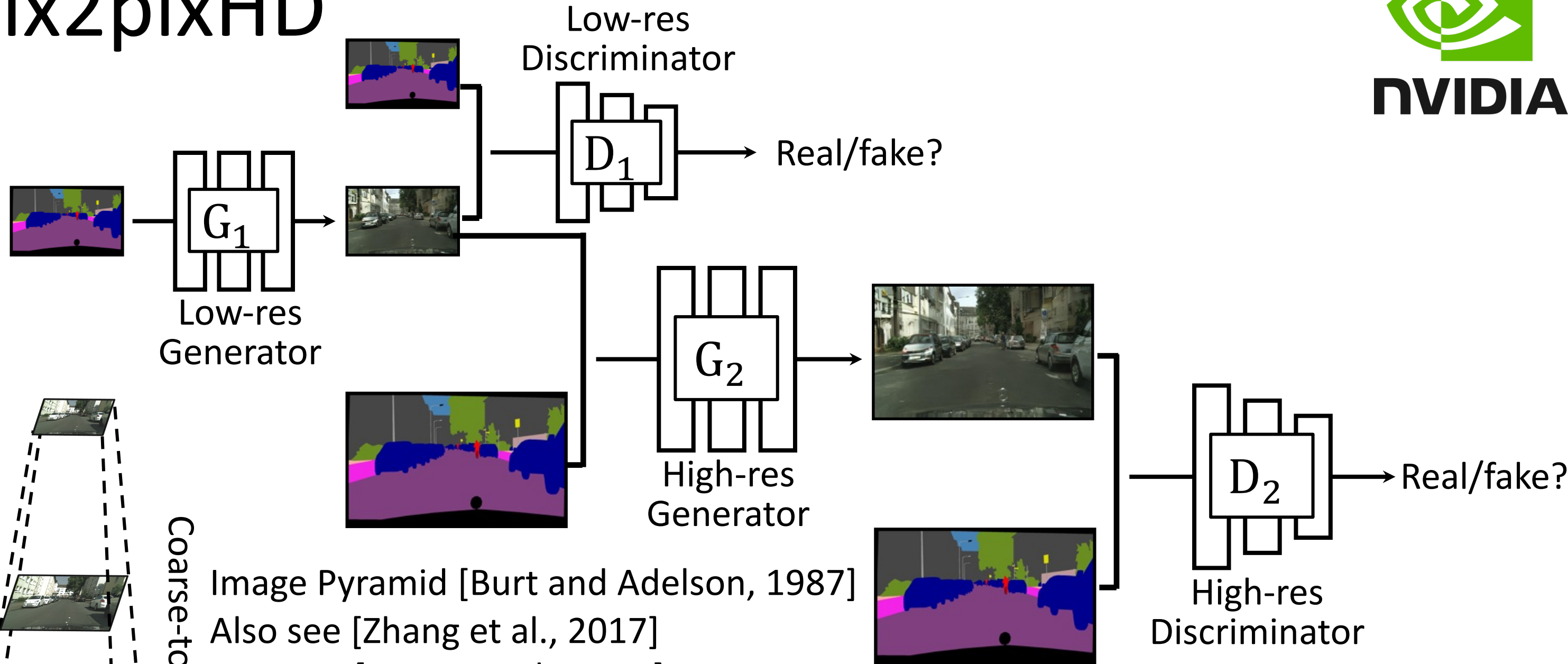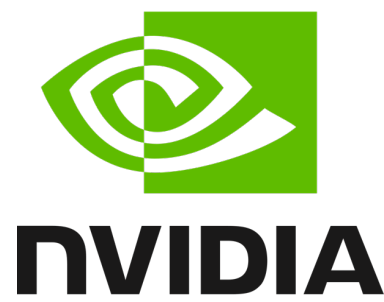
16-726, Spring 2023

# Improving Conditional GANs

- Multimodal synthesis.

- **High-resolution synthesis.**

- Model training without pairs

# The Curse of Dimensionality



Pix2pix output

# pix2pixHD



Low-res Discriminator

Low-res Generator

High-res Generator

High-res Discriminator

$D_1$ → Real/fake?

$D_2$ → Real/fake?

Coarse-to-fine

Image Pyramid [Burt and Adelson, 1987]
Also see [Zhang et al., 2017]
[Karras et al., 2018]

Objective: Multi-scale GANs loss + Perceptual Loss
+ Feature Matching Loss (with Discriminator's features)

pix2pixHD [Wang et al., 2018]
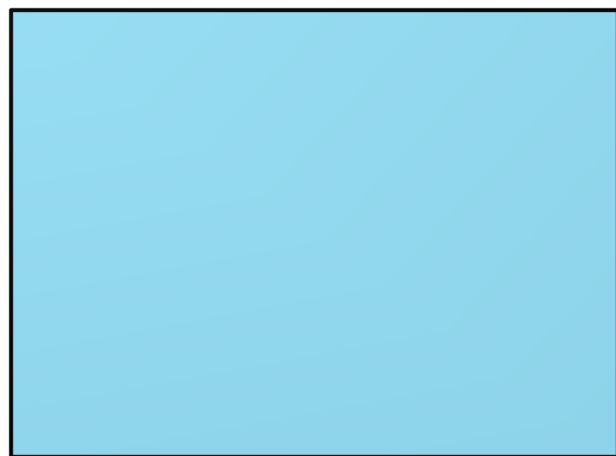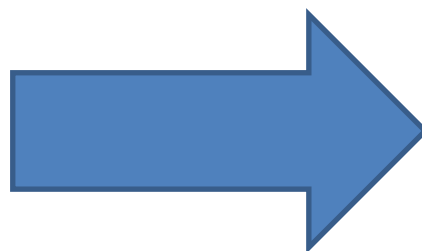
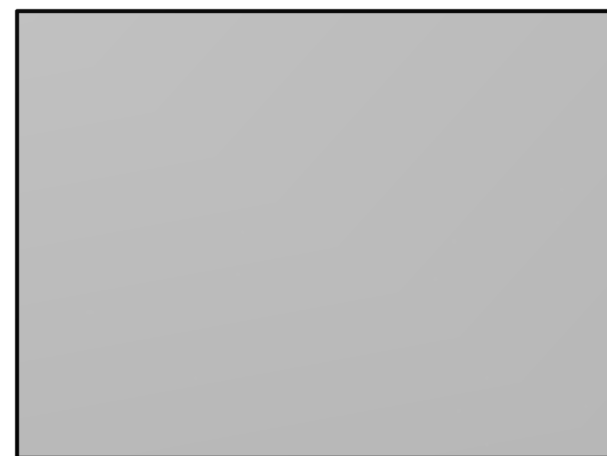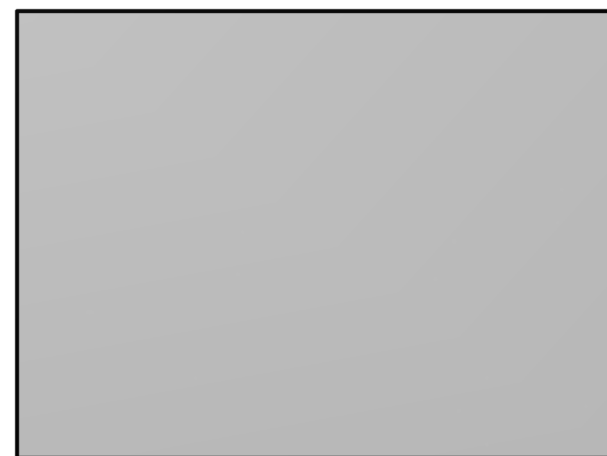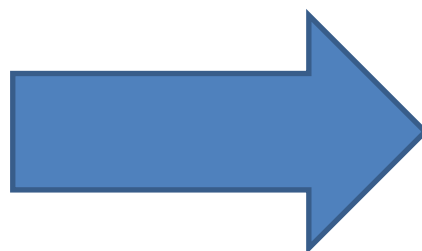# Conditional Image Synthesis in the Wild

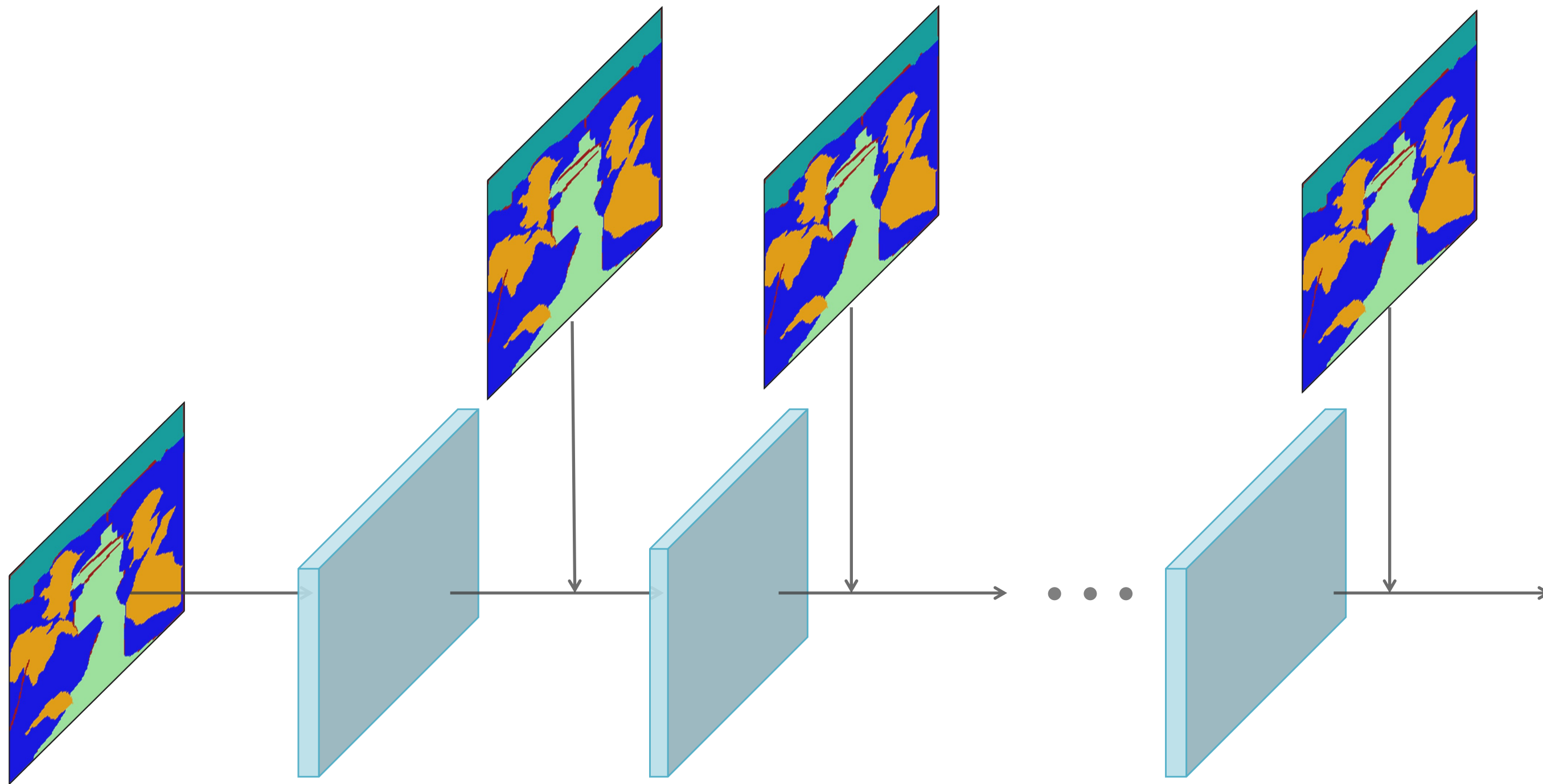

pix2pixHD [Wang et al., 2018]

input

output

sky

grass

pix2pixHD [Wang et al., 2018]

# Problem with standard networks



rock, water, moss, …
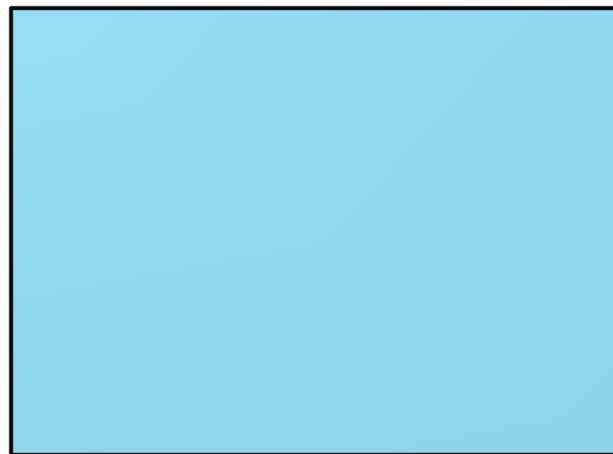
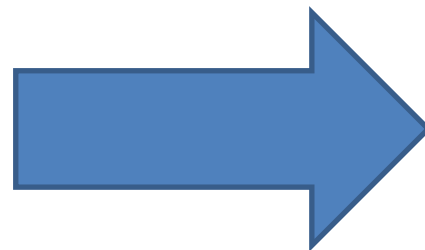conv

normalization

# Problem with standard networks
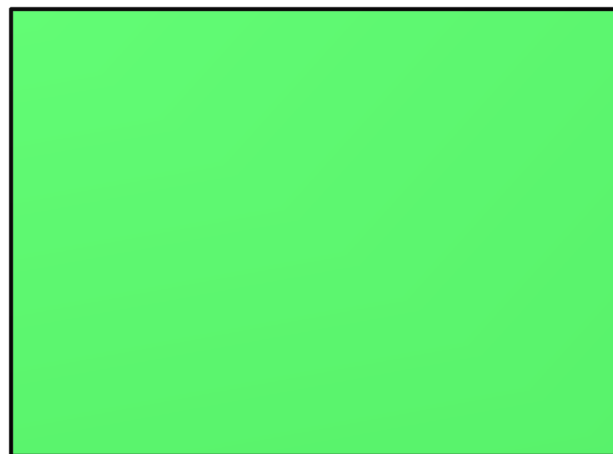
grass

conv
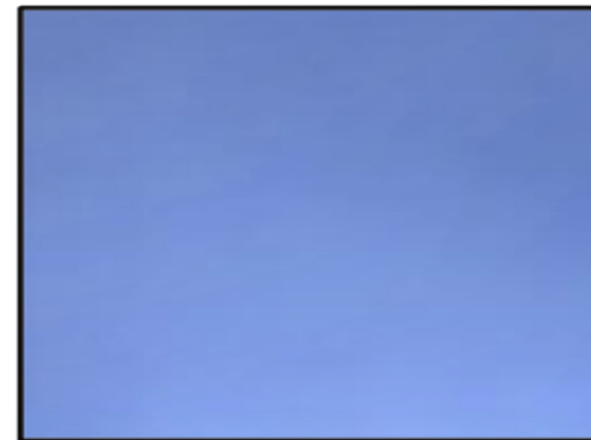
normalization

normalization normalization normalization

input

output

sky

grass

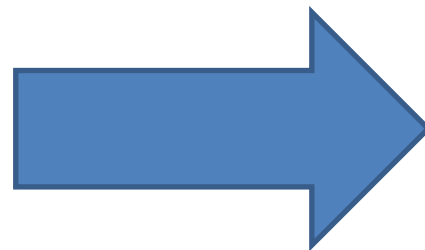SPADE (ours)

# SPADE(SPAtially ADaptive DEnormalization)

# SPADE(SPAtially ADaptive DEnormalization)

Batch Norm (Ioffe et al. 2015)

$$y = \underbrace{\frac{x - \mu}{\sigma}}_{\text{normalization}} \cdot \underbrace{\gamma + \beta}_{\text{affine transform}}$$

See other adaptive/conditional normalization: conditional BN (Dumoulin et al.),
AdaIN (Huang and Belongie), SFT (Wang et al.)

# Generator

# Semantic Control

# Semantic Control

# Semantic Control

# Style Control

# Style Control



Style Manipulation

# Style Control



Style Manipulation

By Darek Zabrocki, Concept Designer and Illustrator

# Learning vs. Exemplar-based

Learning-based

Exemplar-based
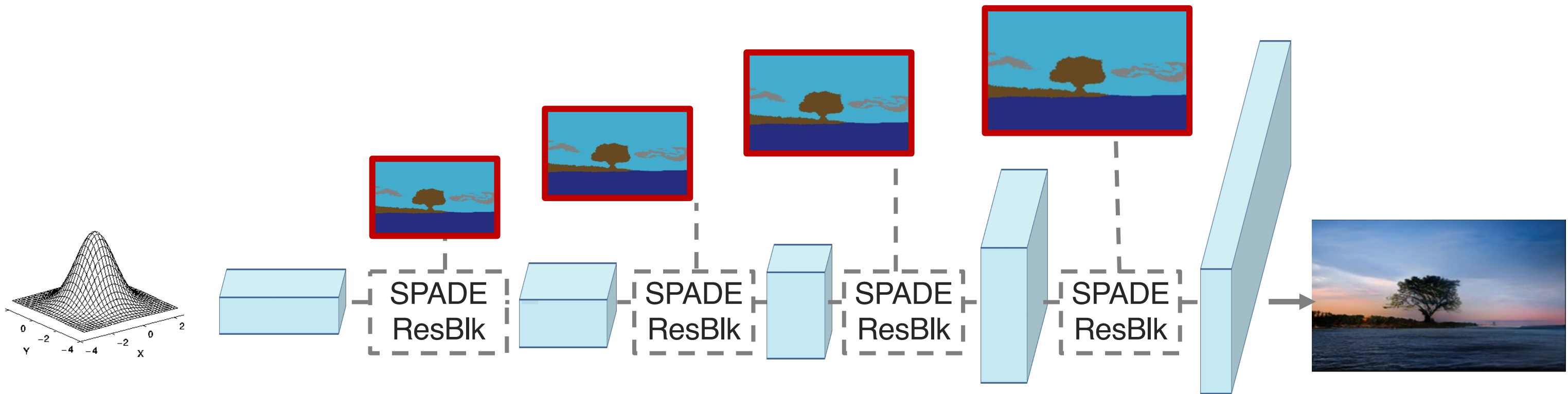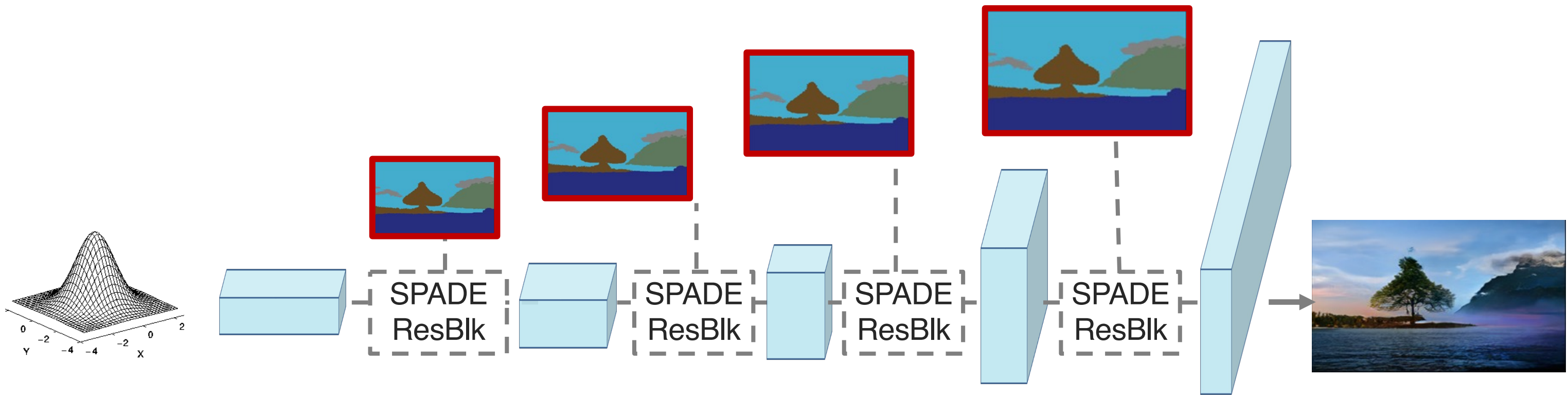
[Isola et al], [Wang et al]
[Park et al], SEAN [Zhu et al]

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

Speed

Local realism

Global realism

Match Input

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



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 $p(y \mid x)$



$X$

$X \to Y$

$D_Y$

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

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

$X$

$Y$

fake zebra

real zebra

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

$X$ → $Y$

$Y$ → $D$

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



Input image

Generator

G(x)



Output image

D

Discriminator

Real (1) or fake (0)?

**Adversarial loss**

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

**Feature loss**

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

$$|\mathrm{F}(\quad) - \mathrm{F}(\quad)|$$

x

Input

G(x)

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**

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

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

Adversarial loss $D_Y(G(x))$

Cycle-consistency loss

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

**Cycle-consistency loss**

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

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 $D_Y(G(x))$

$D_X(F(y))$ 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

# Domain Adaptation with CycleGAN



CG images

Free segmentation labels

Data and labels from [Richter et al. 2016]

# Domain Adaptation with CycleGAN



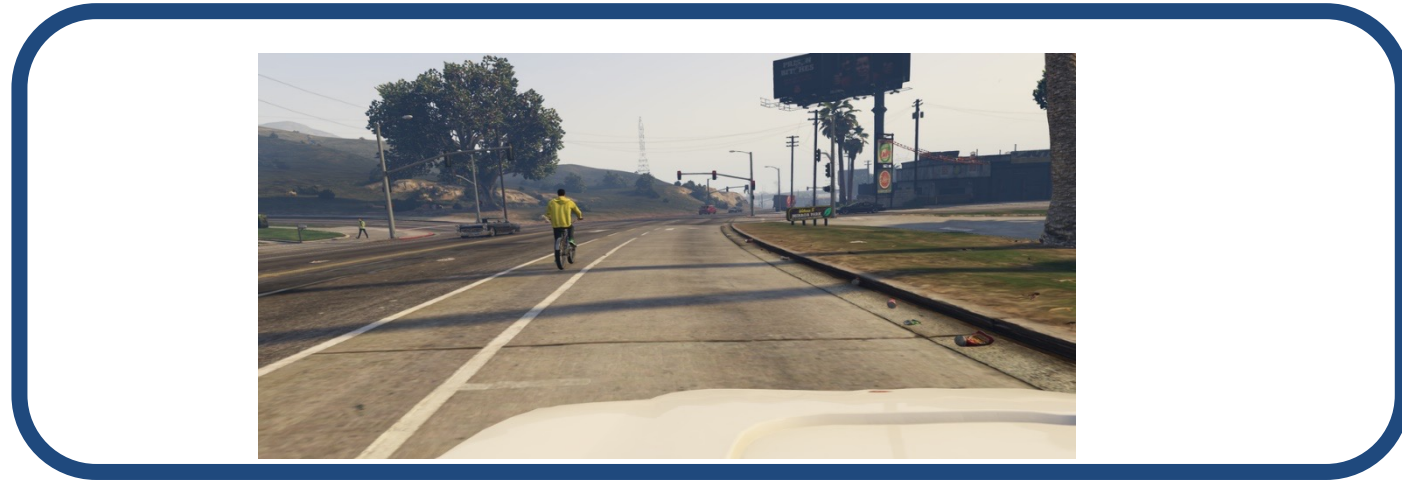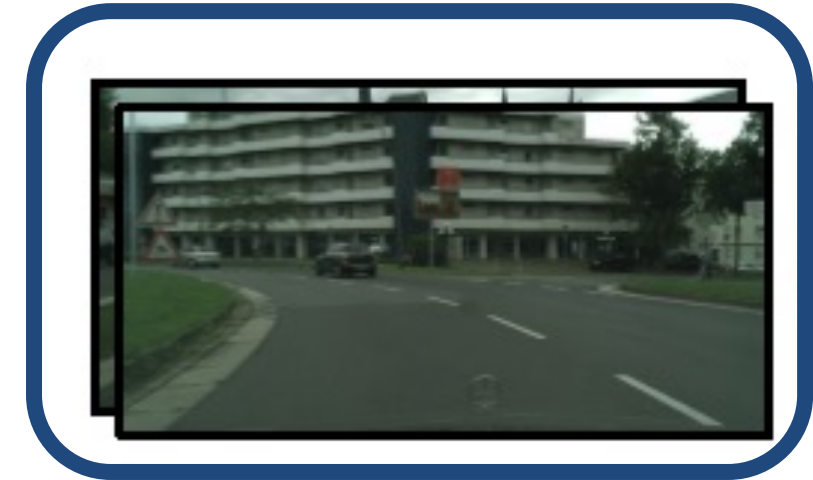Train on CG data

Test on real images

Class-weighted Accuracy

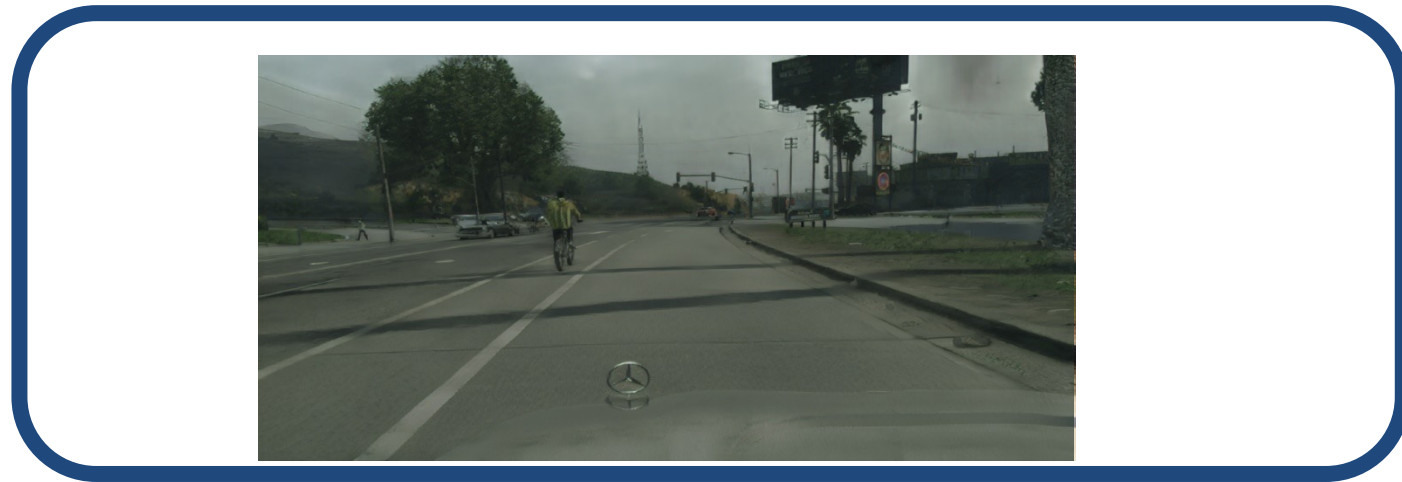| | |
|---|---|
| 70 | |
| 60 | |
| 50 | 47.4 |
| 40 | |

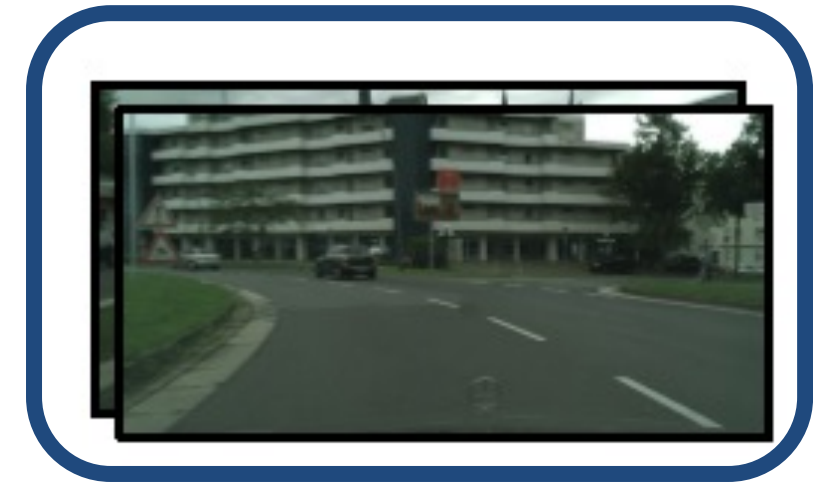Train on CG

# Domain Adaptation with CycleGAN



Train on CycleGAN images

Test on real images

Class-weighted Accuracy

| | |
|---|---|
| Train on CG | 47.4 |
| CycleGAN | 63.8 |
| SOTA adaptation | 67.4 |
| CycleGAN+ SOTA adaptation | 72.4 |

CycleGAN [Zhu*, Park* et al., ICCV 2017]
CycADA [Hoffman et al.,. ICML 2018]

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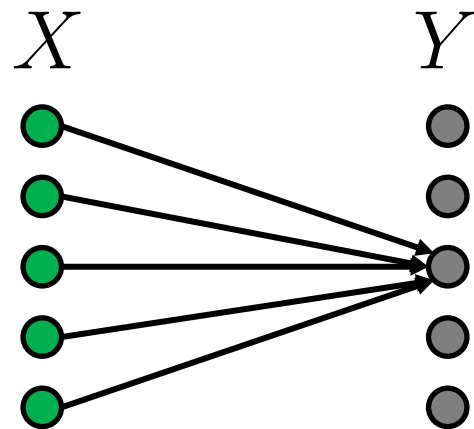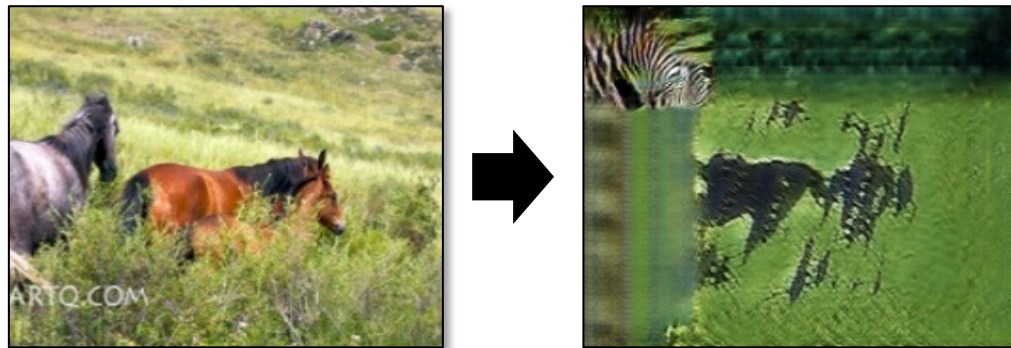
$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$

**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**

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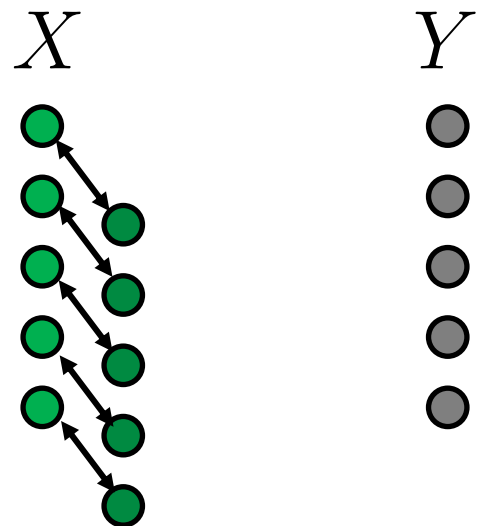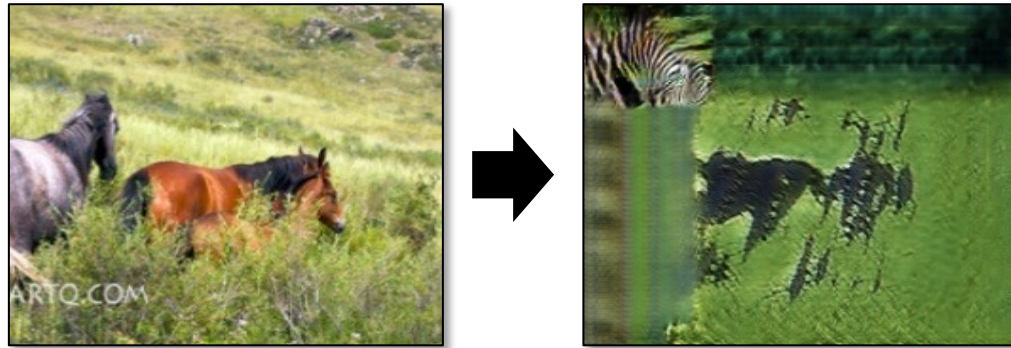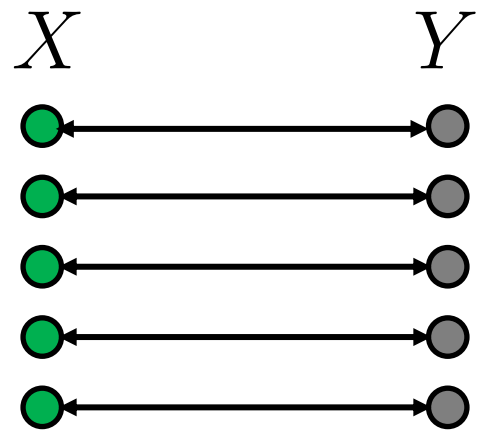


**Full objective**

# Why CycleGAN works

**Adversarial loss**

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

**Cycle-consistency loss**

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

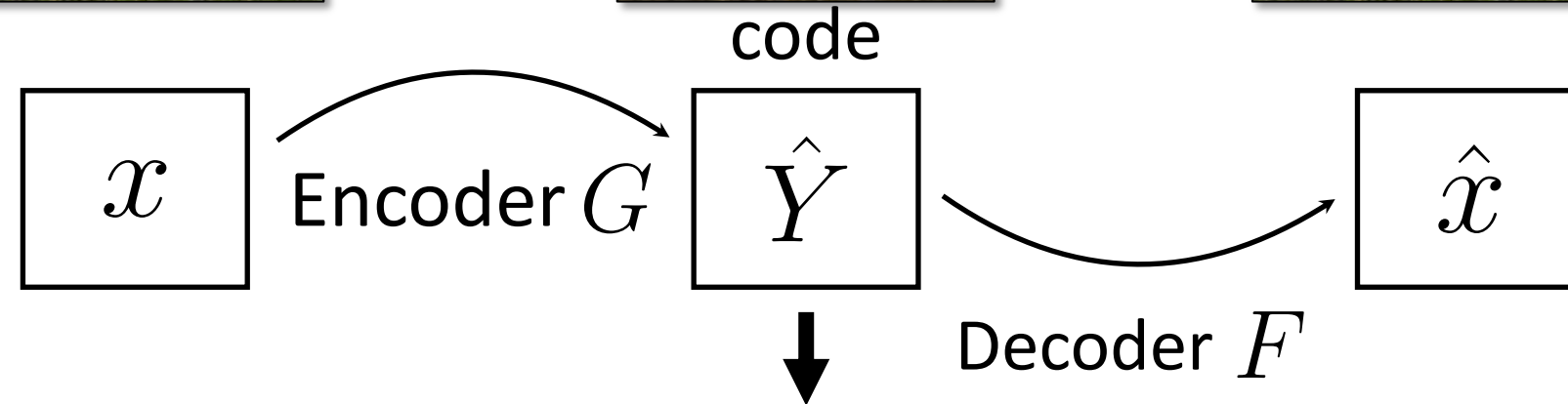$x$ ⟶ $G(x)$ ⟶ $F(G(x))$



**Auto-encoder w/ domain prior**

$x$ —Encoder $G$→ $\hat{Y}$ (code) —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 $G$

$x$ $\hat{Y}$ $\hat{x}$



$G$ $F$

<u>A strong regularizer</u>

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

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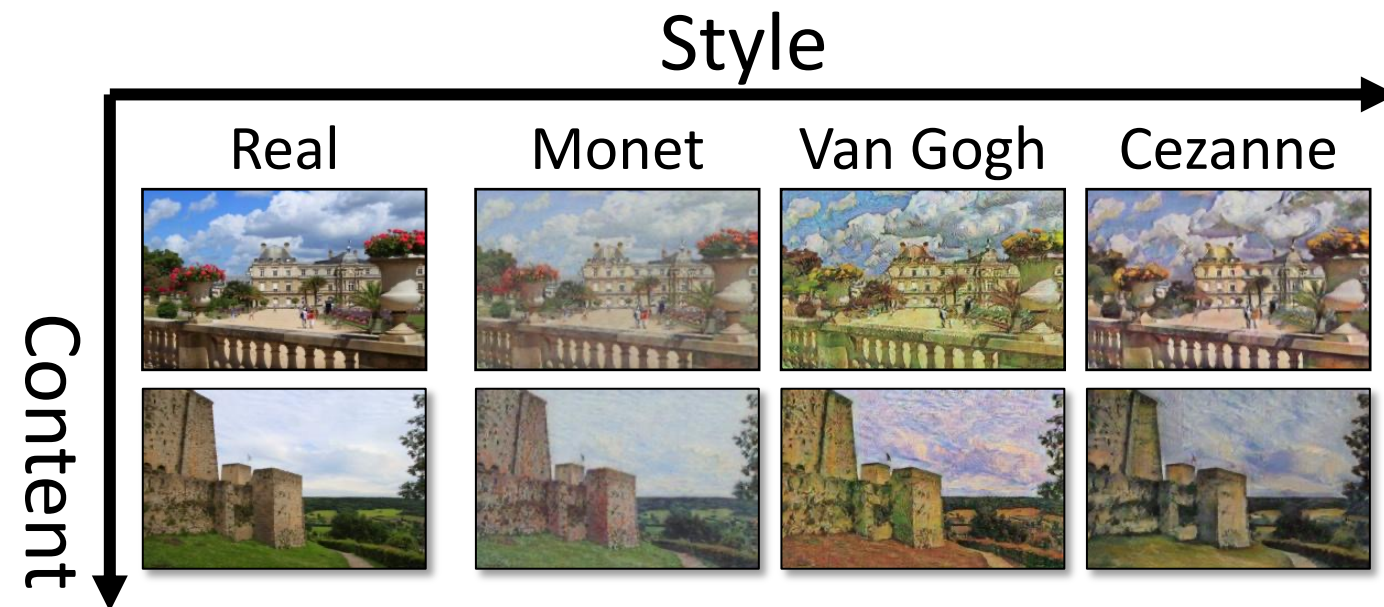
$$\mathbb{E}_x \log(1 - D_Y(G(x))) + \mathbb{E}_y \log D_Y(y)$$



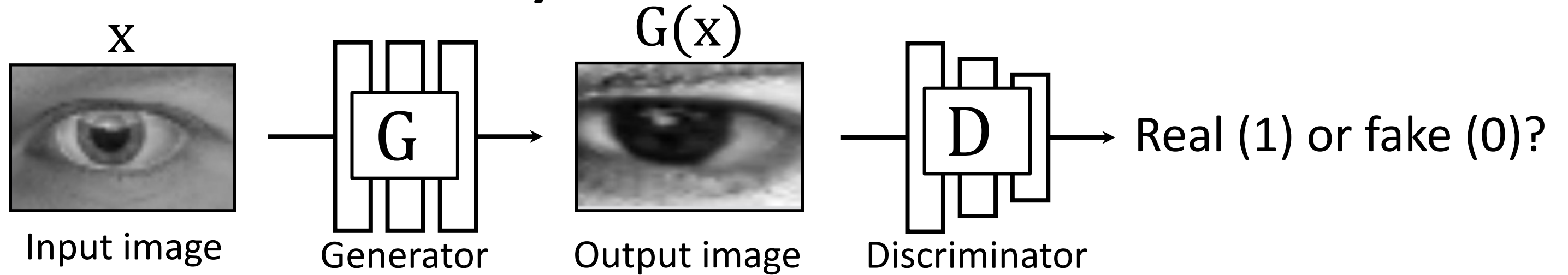$p(x) \rightarrow p(y)$ change **style**

**Cycle-consistency loss**

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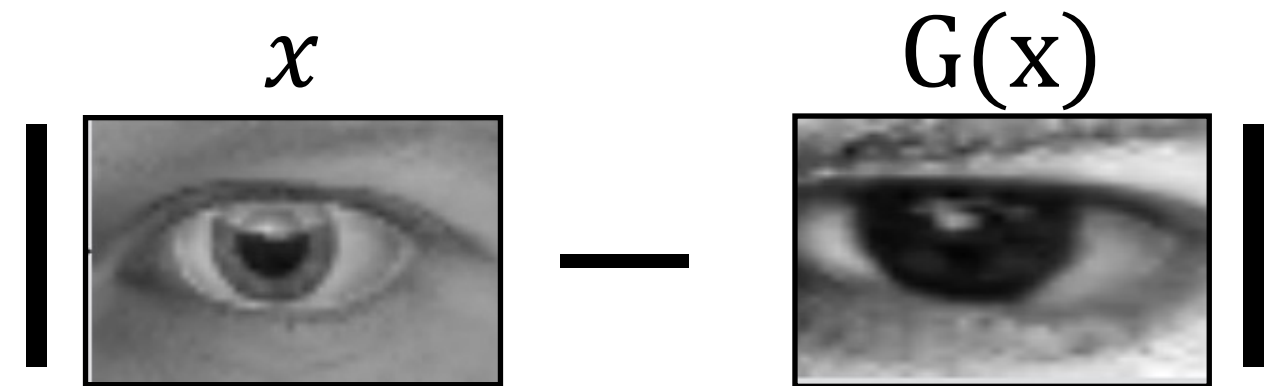


Bidirectional: preserve **content**



Style → : Real, Monet, Van Gogh, Cezanne

Content ↓

# Style and Content

x



Input image

G

Generator

G(x)



Output image

D

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)**

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

$x$



$G(x)$


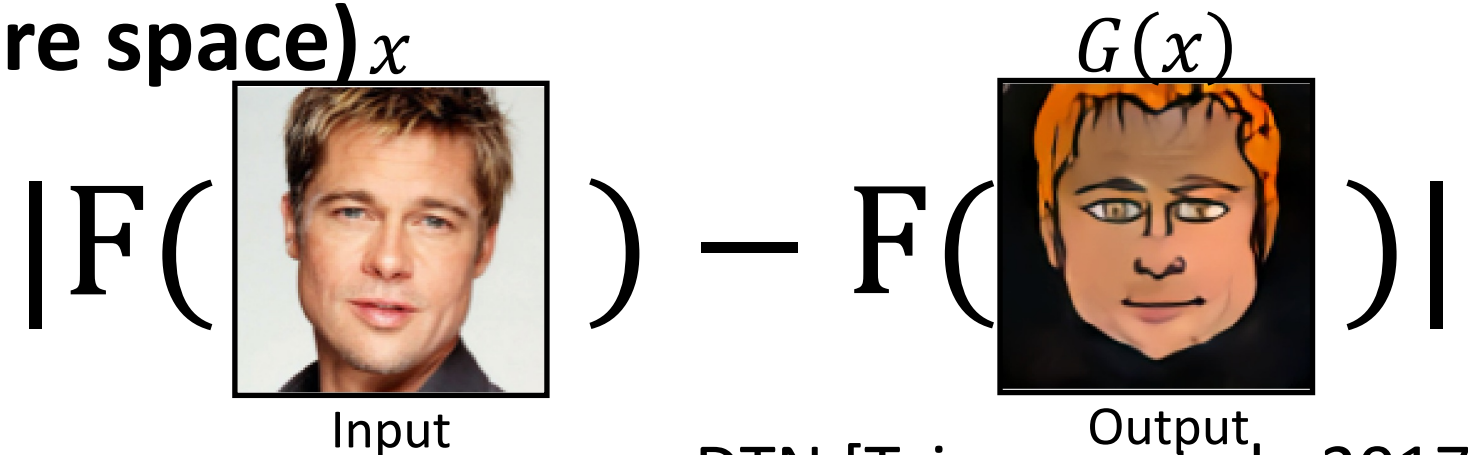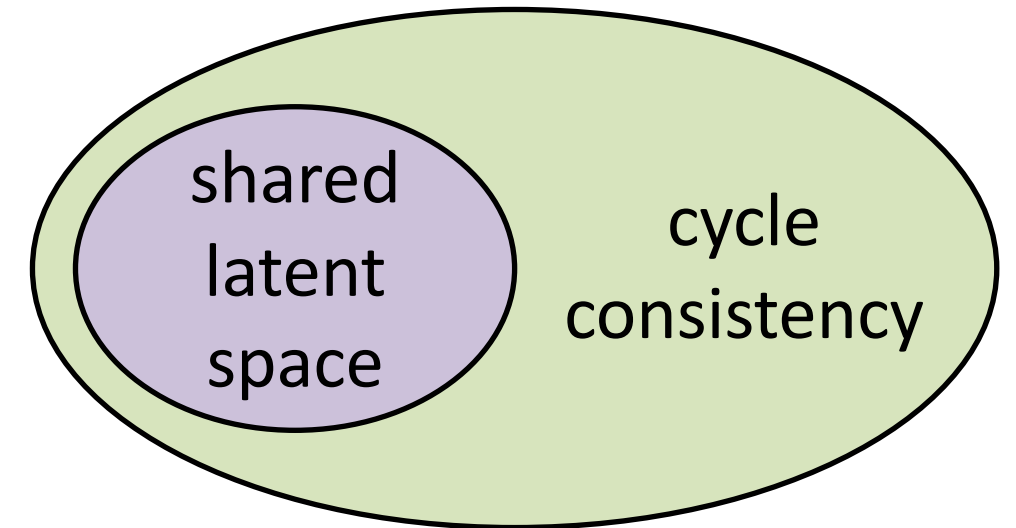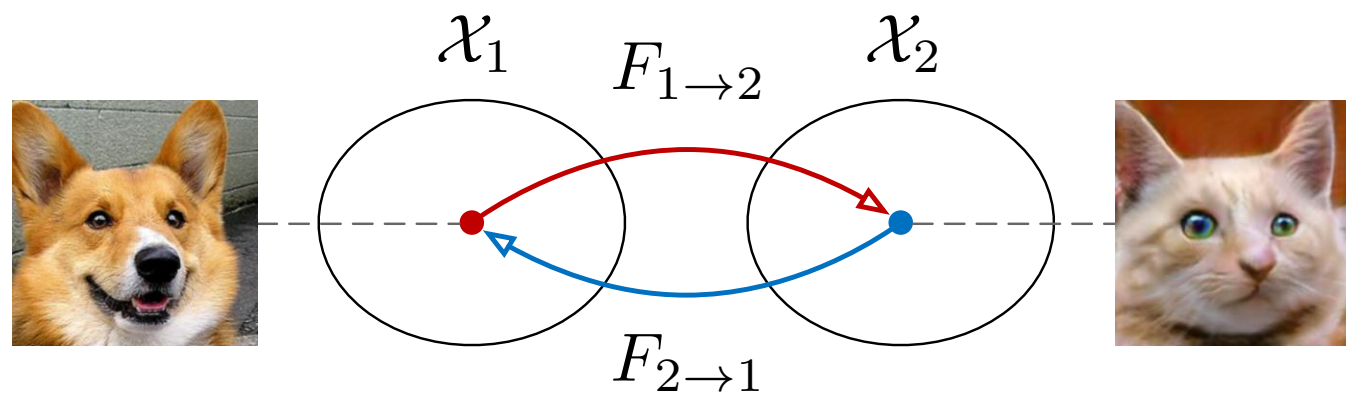
SimGAN [Shrivastava et al., 2017]

# Style and Content



x

G(x)

Real (1) or fake (0)?

Input image

Generator

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)** $x$ $G(x)$

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

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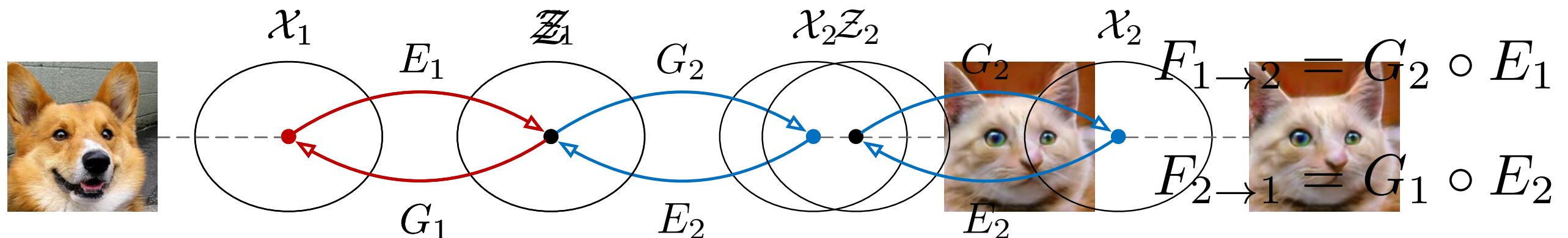
Input

Output

DTN [Taigman et al., 2017]

# CycleGAN and UNIT

- CycleGAN (**cycle consistency**)



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

shared latent space $\implies$ cycle consistency



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

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

# Disentangling the Latent Space

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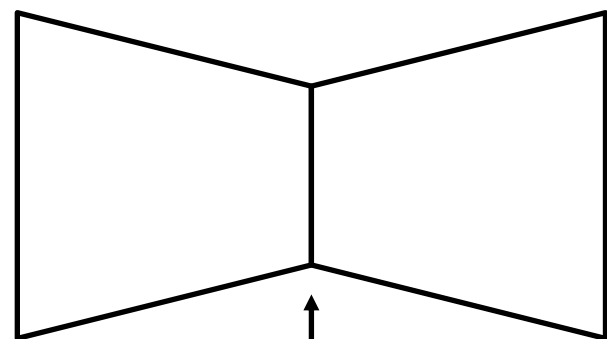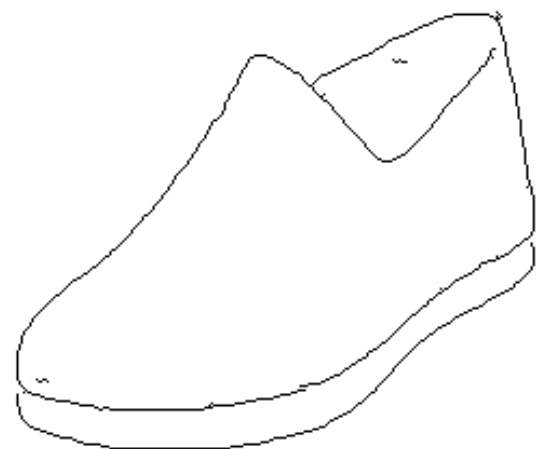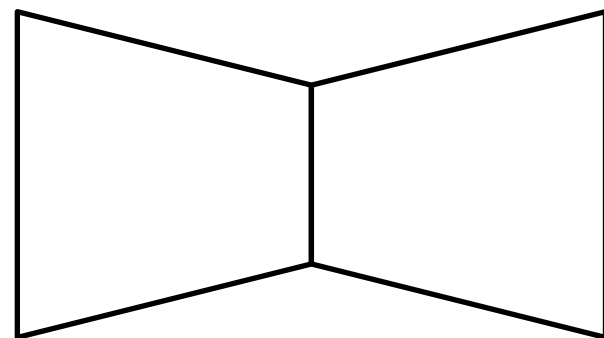  - A single **shared**, **domain-invariant** latent space $\mathcal{Z}$

# 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**
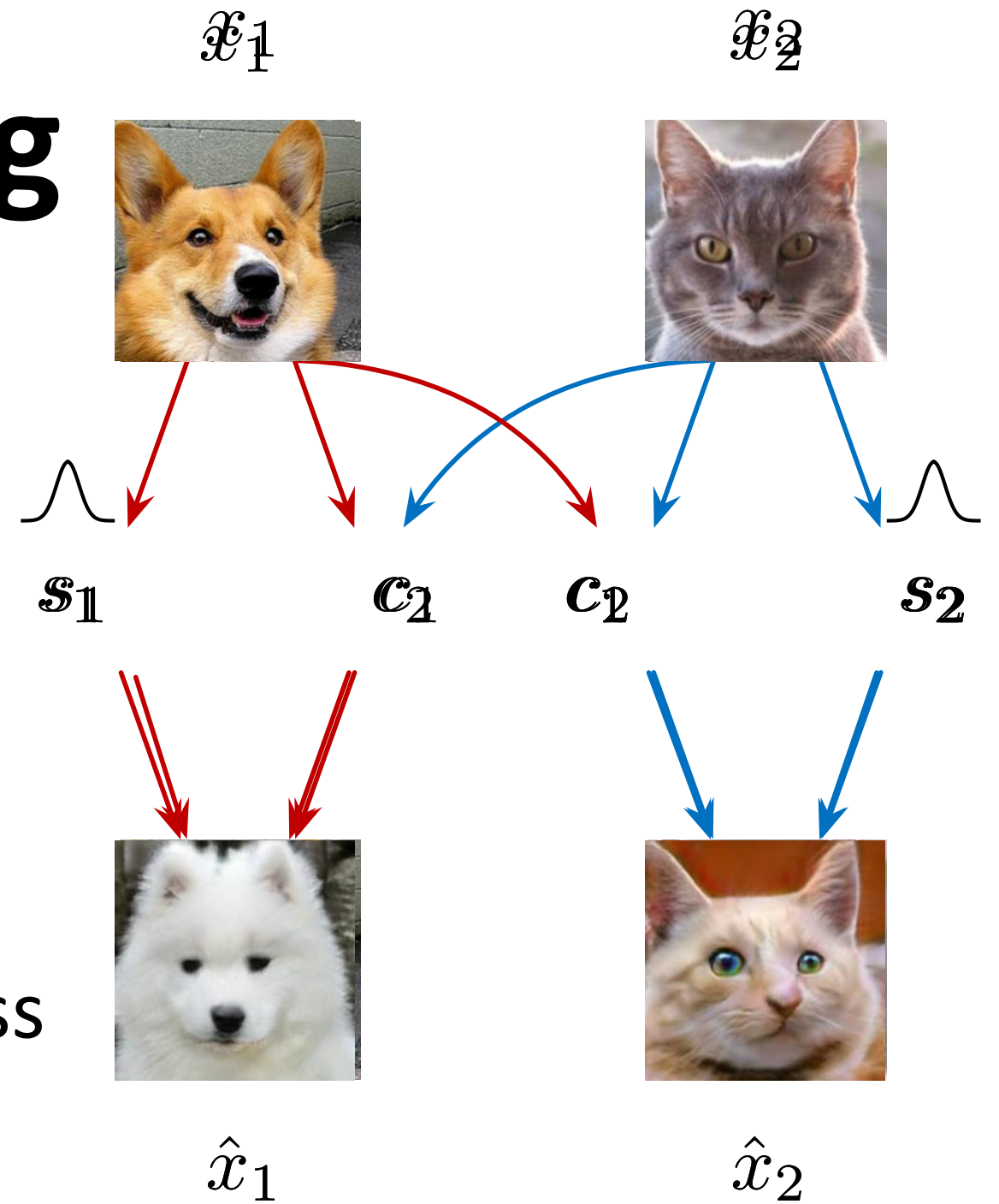
# Unimodality

# 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**

$x_1$

$x_2$

Notations:

− $x$: images

− $c$: content

− $s$: style



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

$\mathcal{L}_1$
loss

$\hat{x}_1$         $\hat{x}_2$

# Bidirectional Reconstruction Loss: Image Reconstruction

$x_1$       $x_2$



Notations:

– $x$: images

– $c$: content

– $s$: style

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

$\mathcal{L}_1$ loss

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

# GAN Loss

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



$x_1$　$x_2$

$s_1$　$c_2$　$c_1$　$s_2$

$x_{2 \to 1}$　$x_{1 \to 2}$

GAN loss

# 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
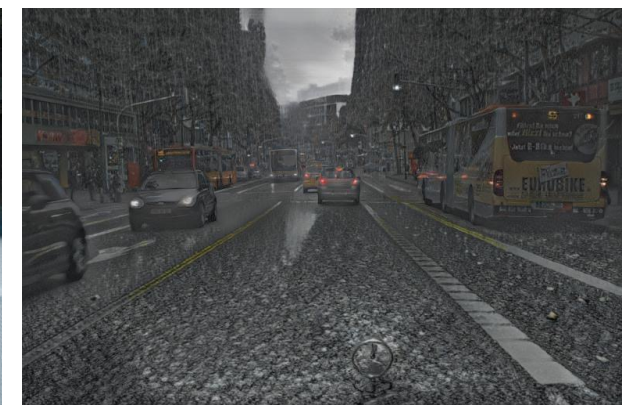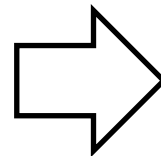
# Sketches <-> Photo

Input

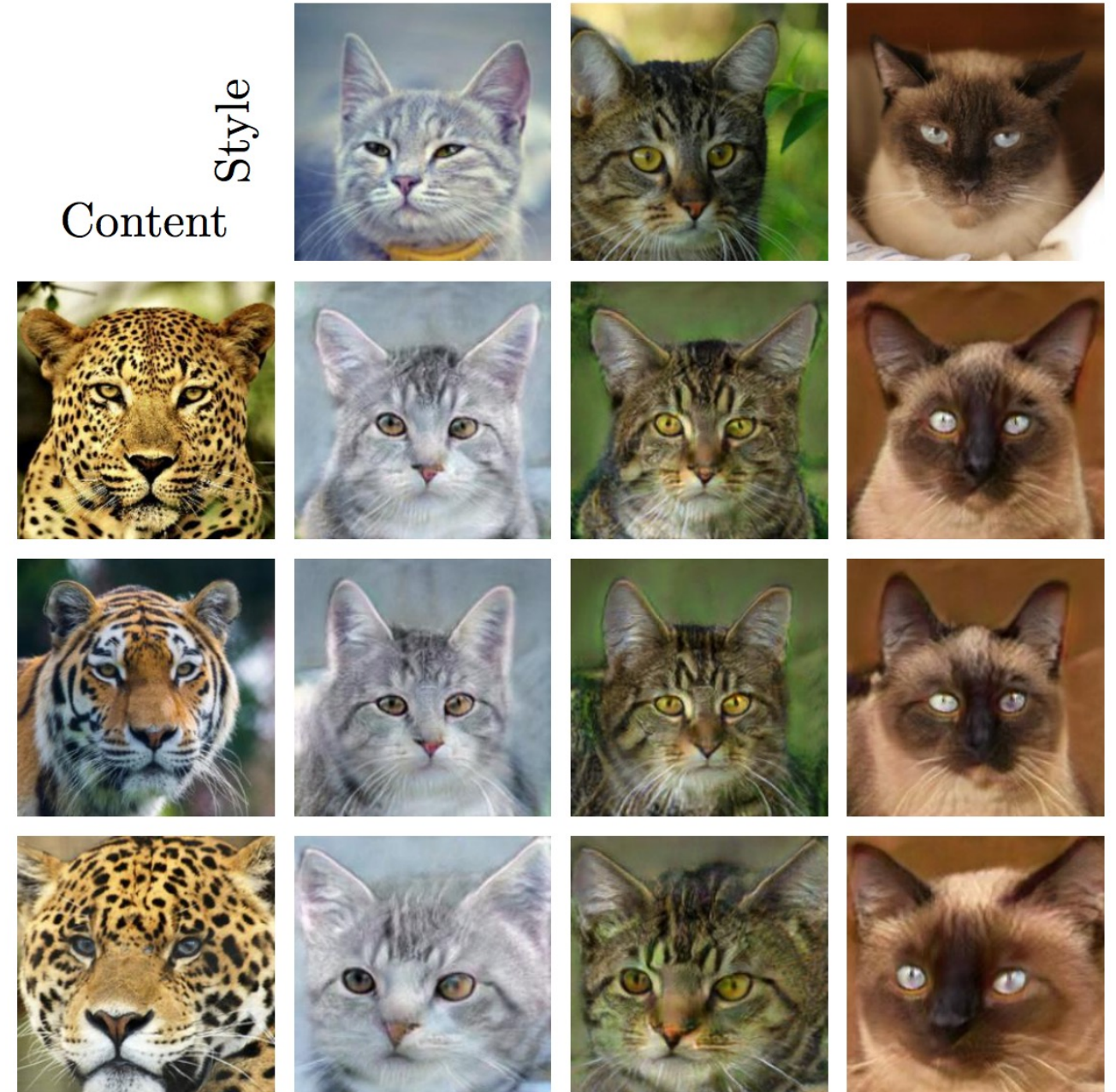Outputs

# Cats ↔ Dogs

Input

Outputs

# Synthetic ↔ Real

Input

Outputs
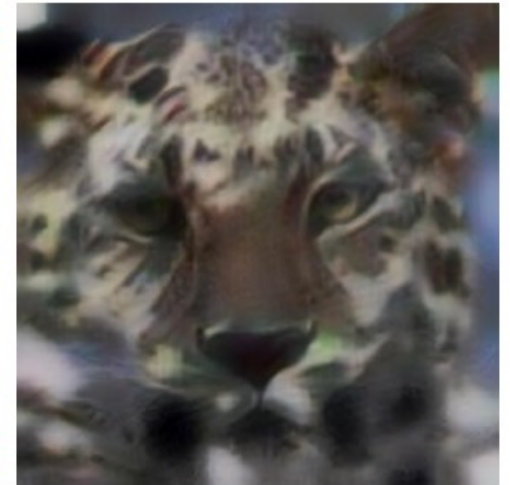
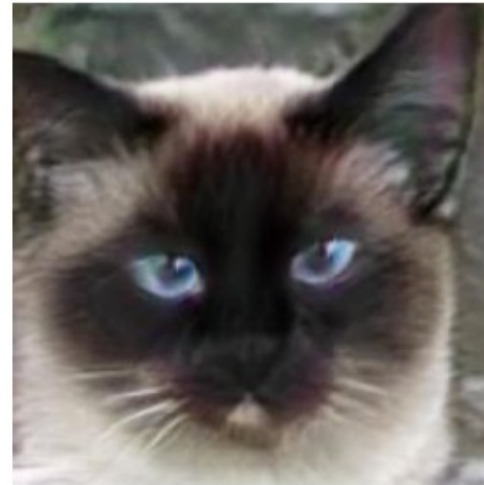# Example-guided Translation
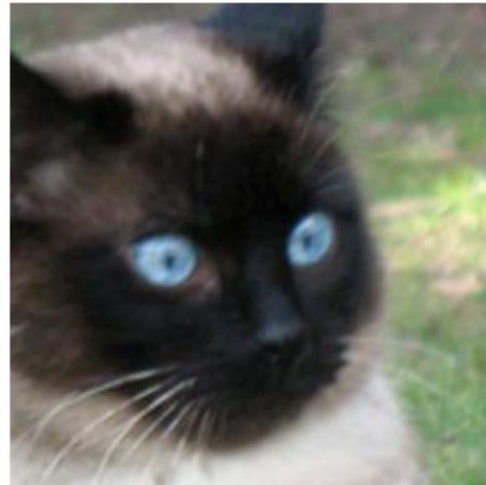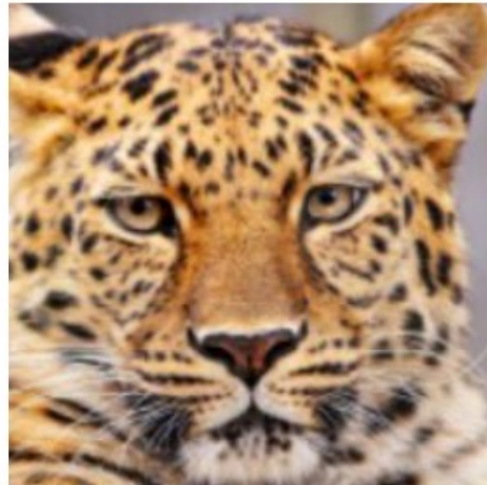
# Example-guided Translation



| Content | Style | Ours | Gatys *et al.* | AdaIN |

# Thank You!



16-726, Spring 2023