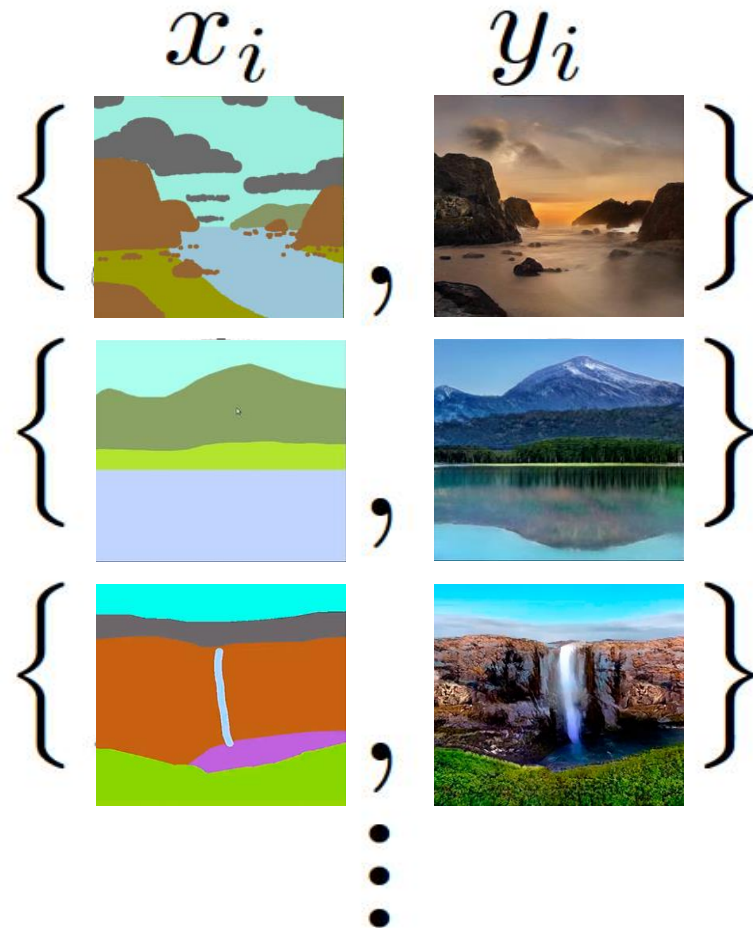# Multimodal Unsupervised Image-to-Image Translation

Ming-Yu Liu
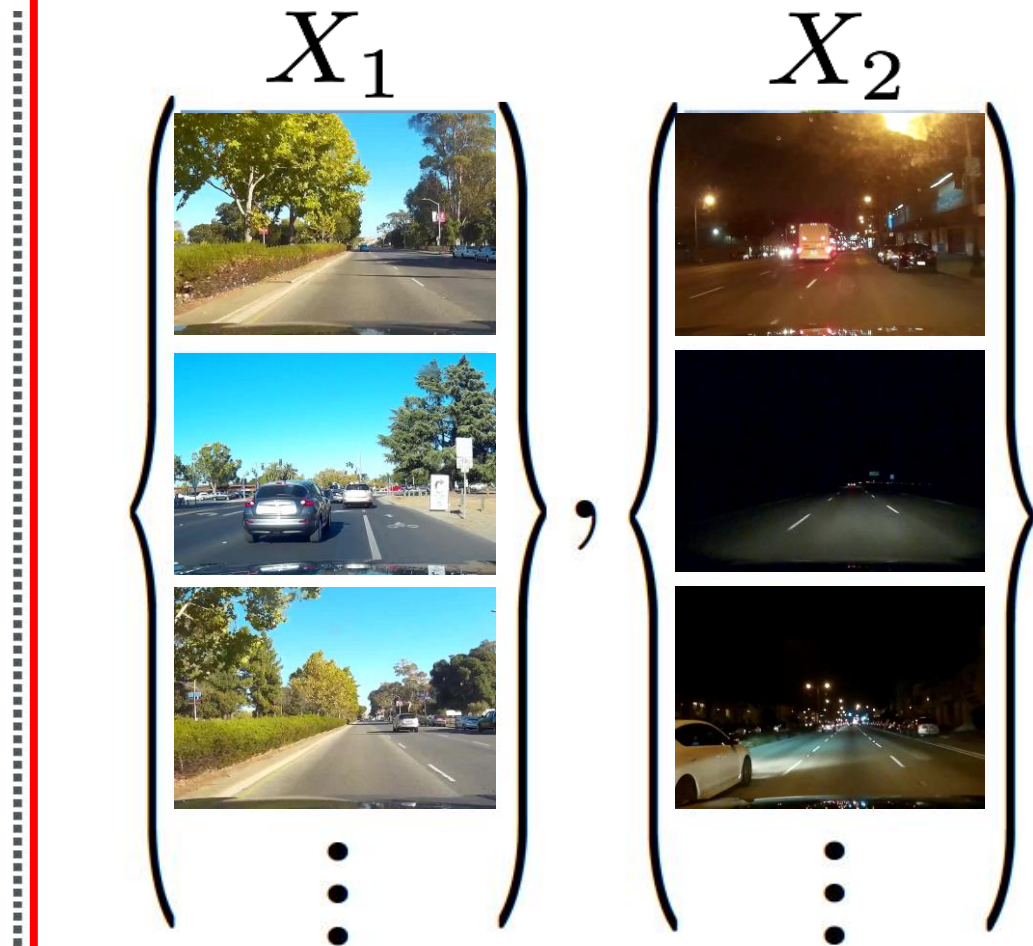
NVIDIA

# Supervised vs Unsupervised

Supervised/Paired/Aligned/Registered

$x_i$  $y_i$



Unsupervised/Unpaired/Unaligned/Unregistered

$X_1$  $X_2$

# Image Domain Transfer

Given an input image
in one domain

Output a corresponding image
in differerent domain



Image
Translator

$F$

Summer image domain

Winter image domain

# Example Applications
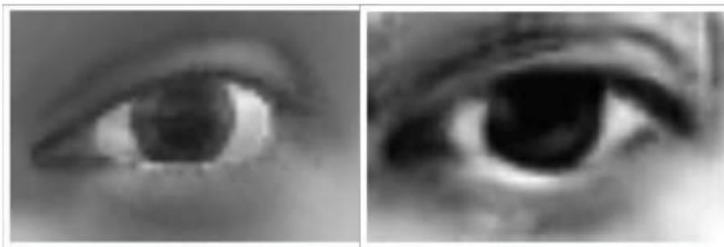


Low-res to high-res

Blurry to sharp

Image to painting
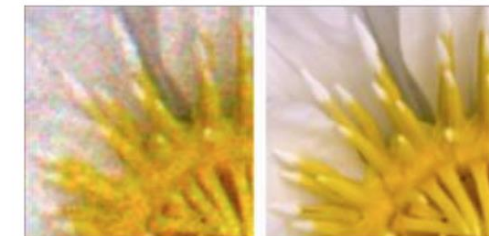
LDR to HDR

Synthetic to real
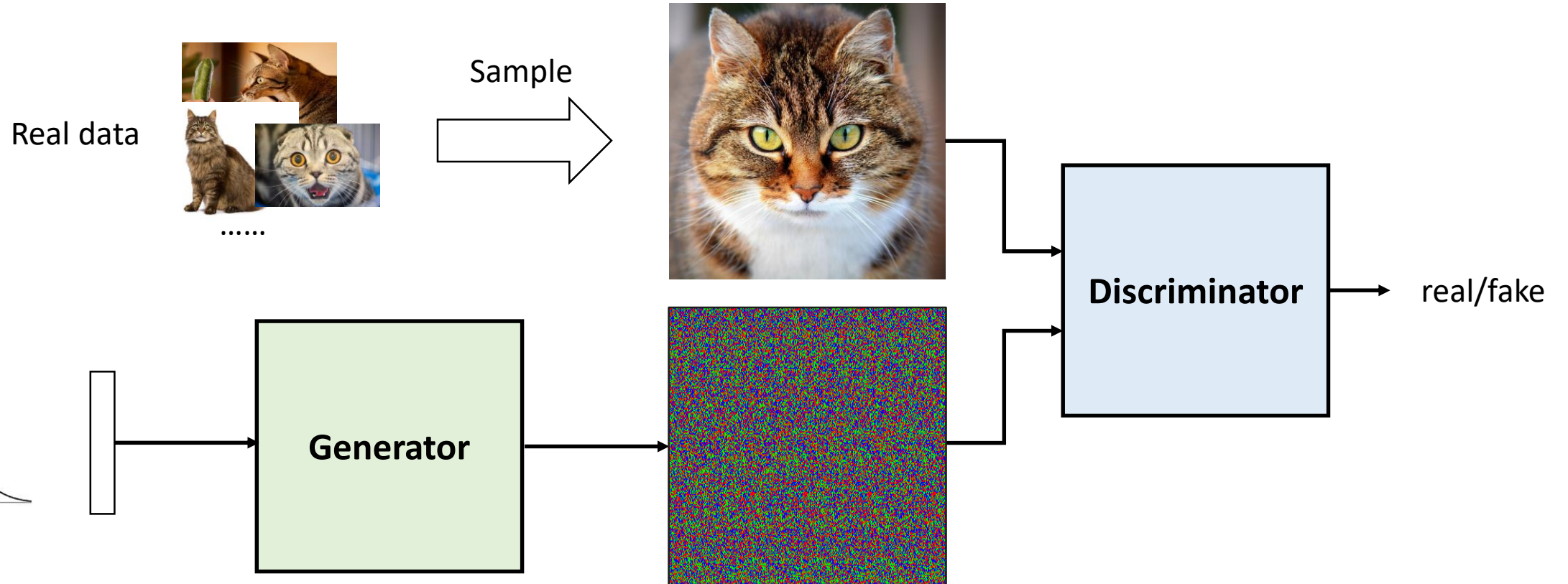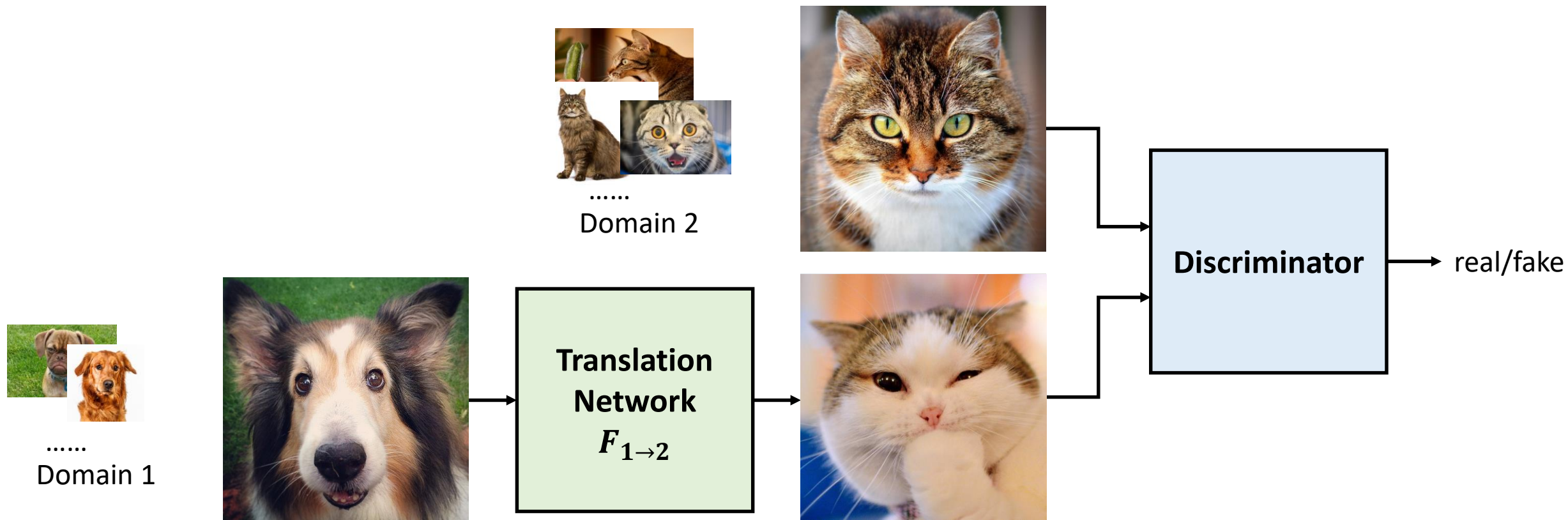
Thermal to color

Day to night

Summer to winter

Noisy to clean

# Generative Adversarial Networks (GANs)
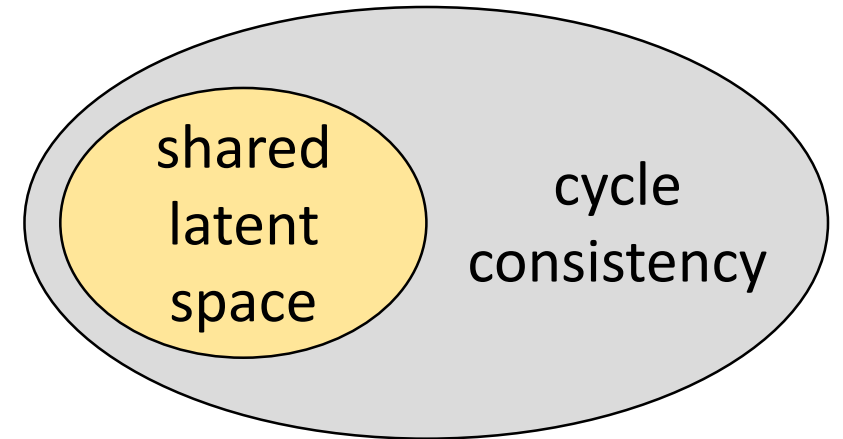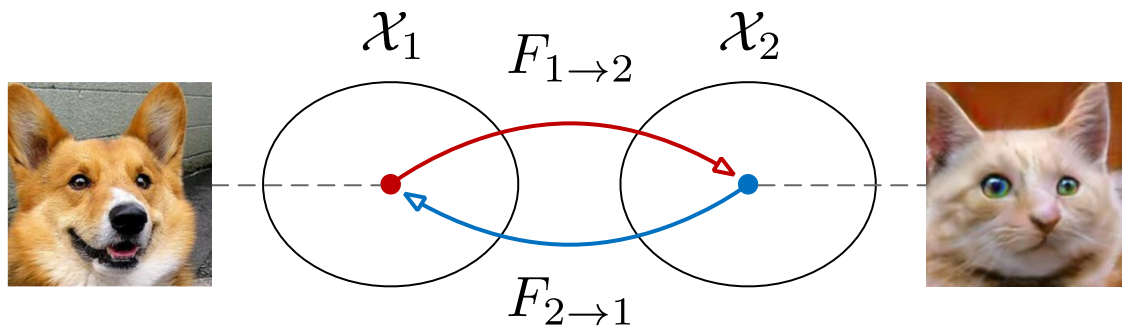


Goodfellow et al. 2014
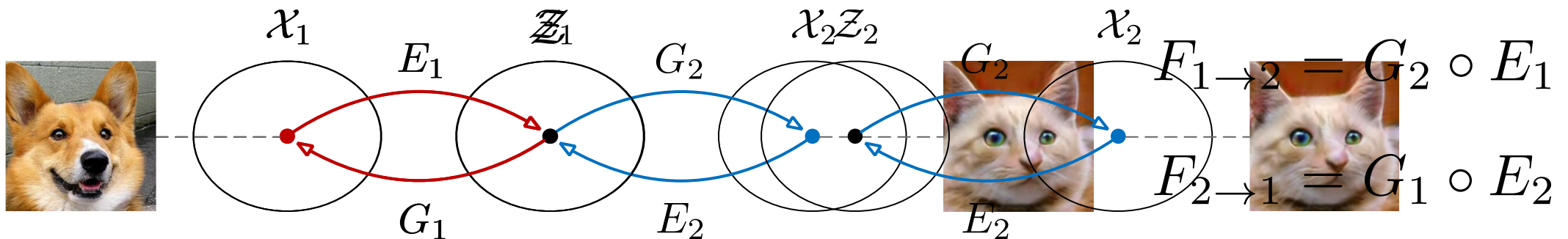
# Plain GAN for Unsupervised Image-to-Image Translation

# CycleGAN and UNIT
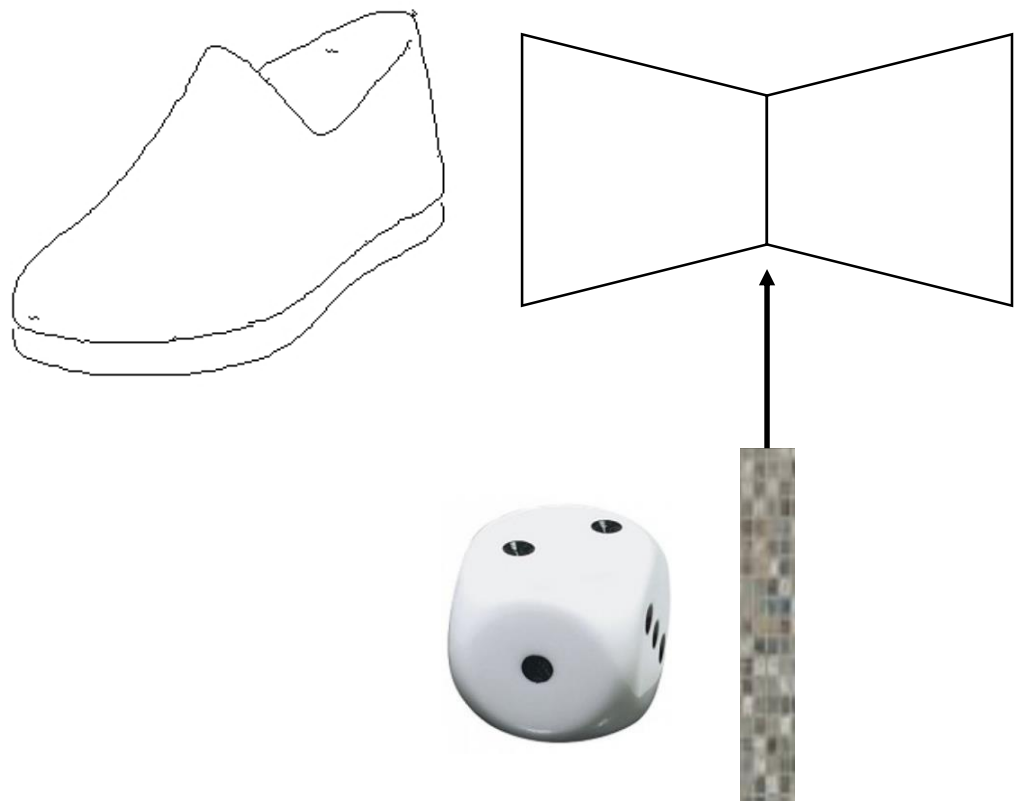
- CycleGAN (**cycle consistency**) [Zhu et al. 2017]



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

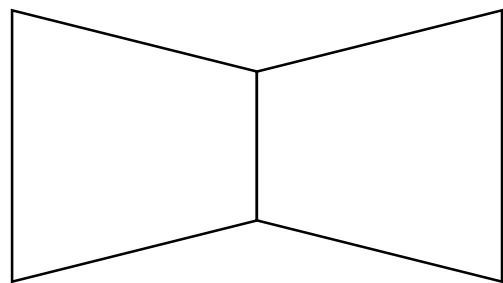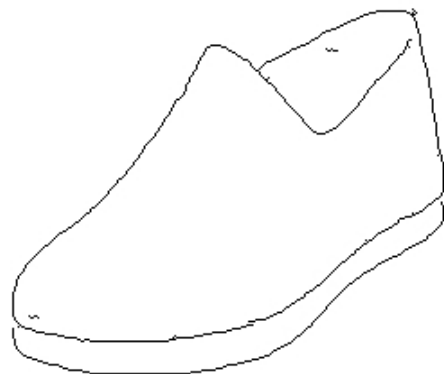shared latent space $\Longrightarrow$ cycle consistency



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

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

# Unimodality
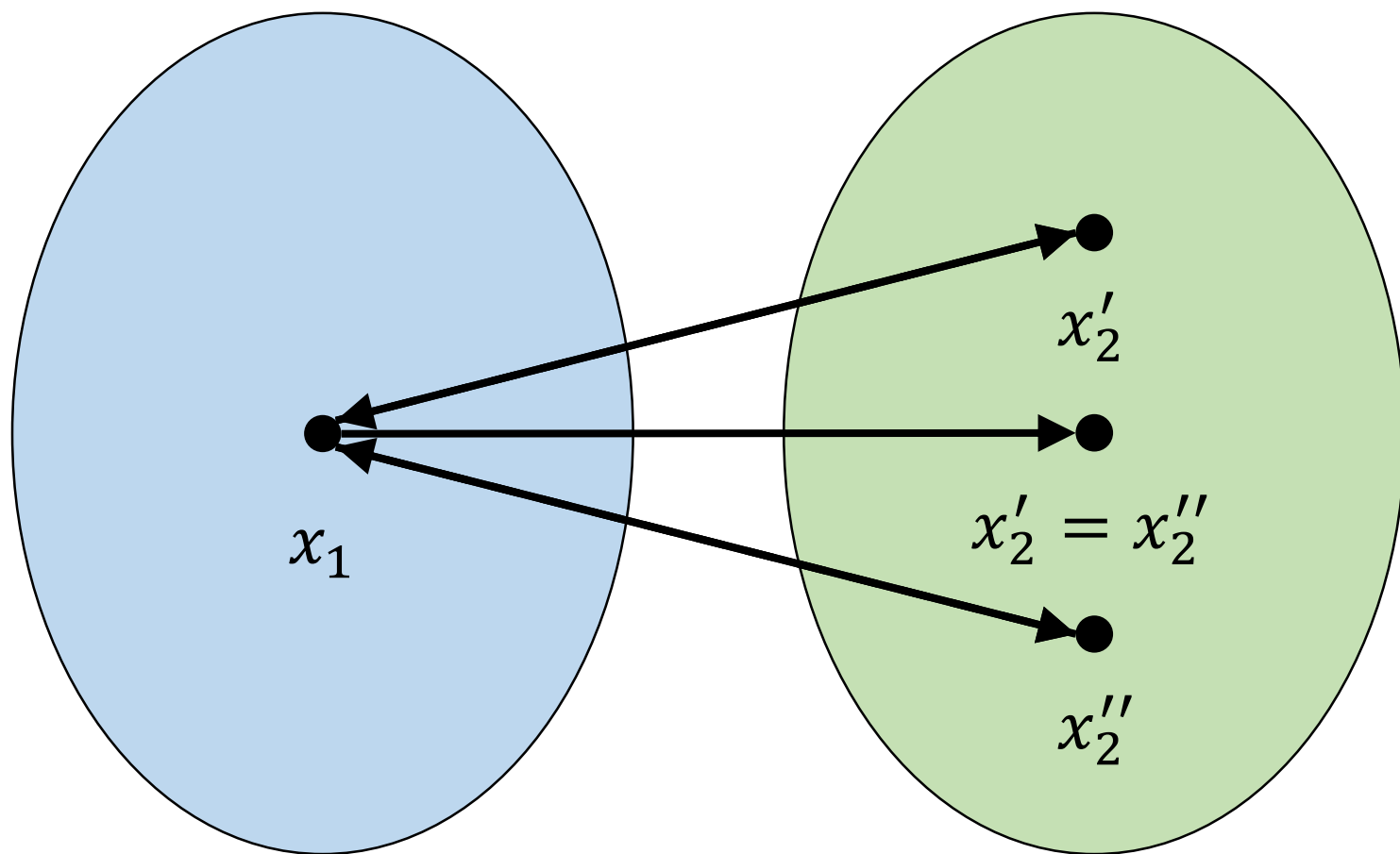
# Towards Multimodality



...

# Shared latent space does not allow multimodality



Domain $\mathcal{X}_1$

Domain $\mathcal{X}_2$
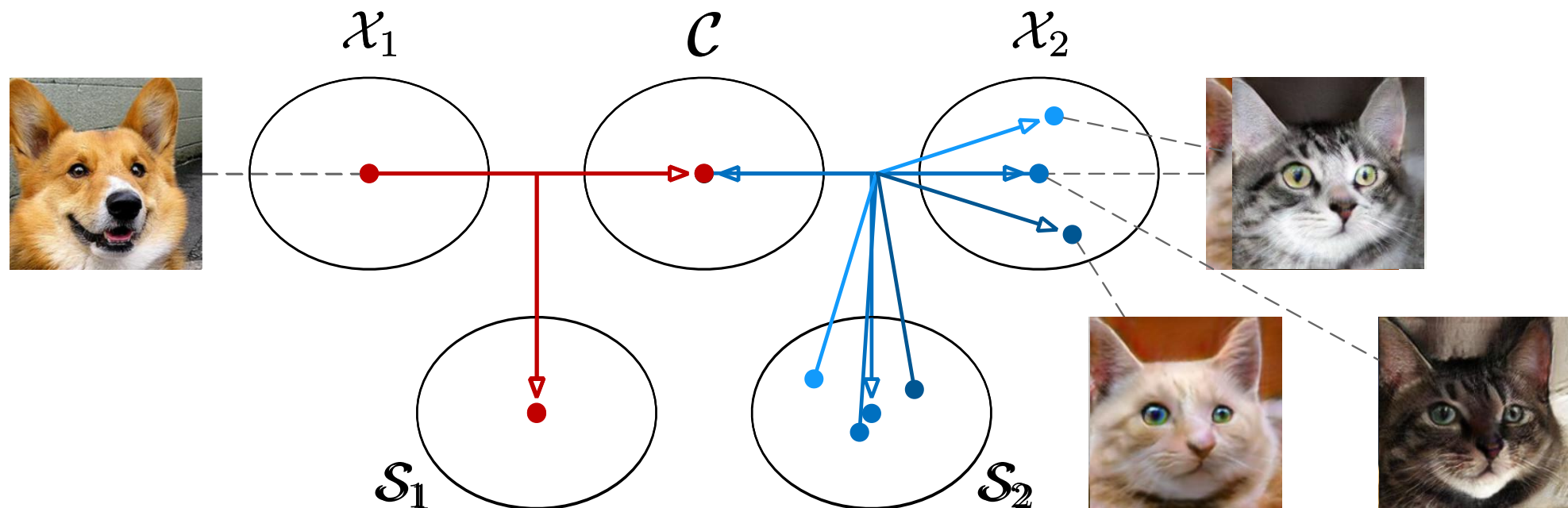
# Disentangling the Latent Space

- UNIT
  - A single **shared**, **domain-invariant** latent space $\mathcal{Z}$
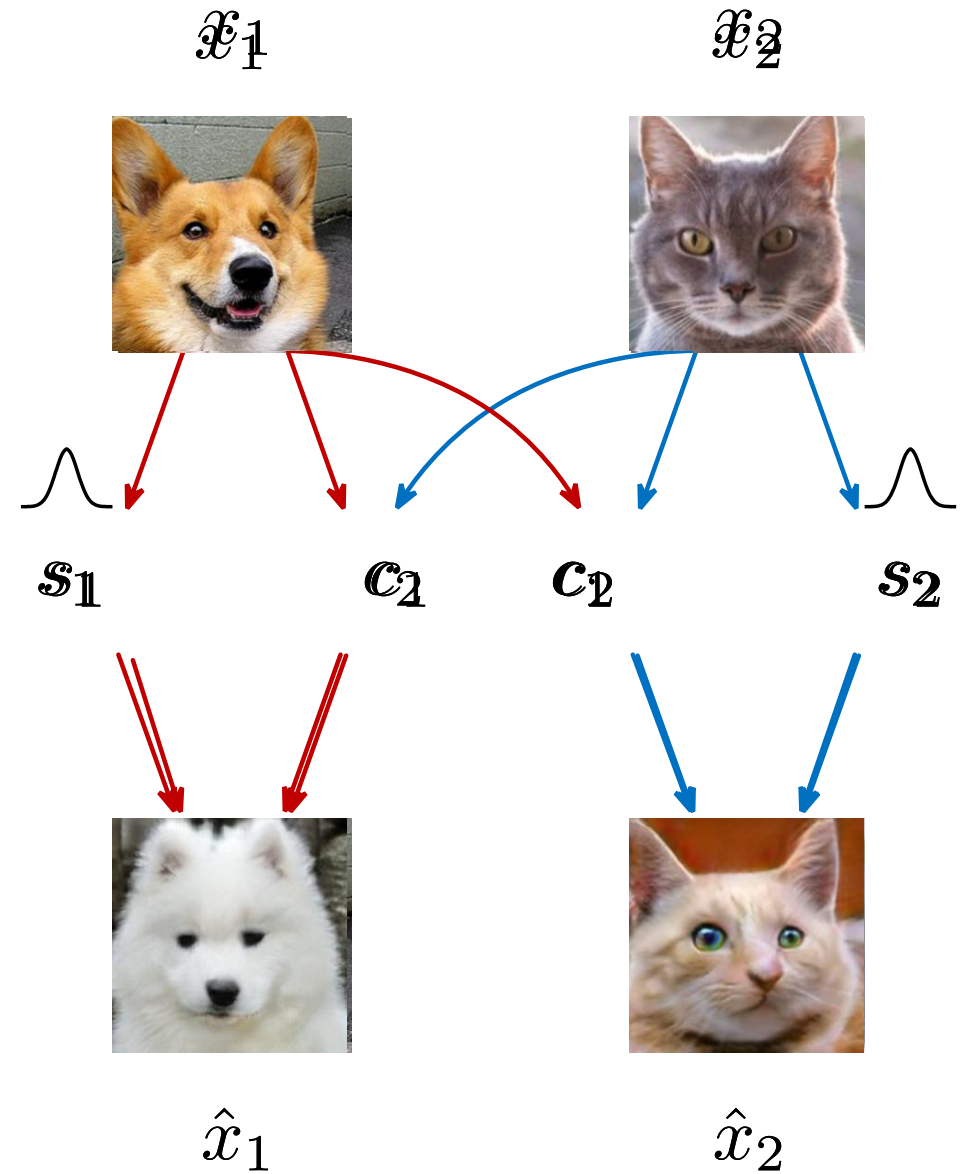
# Disentangling the Latent Space

- 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**

# Training

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

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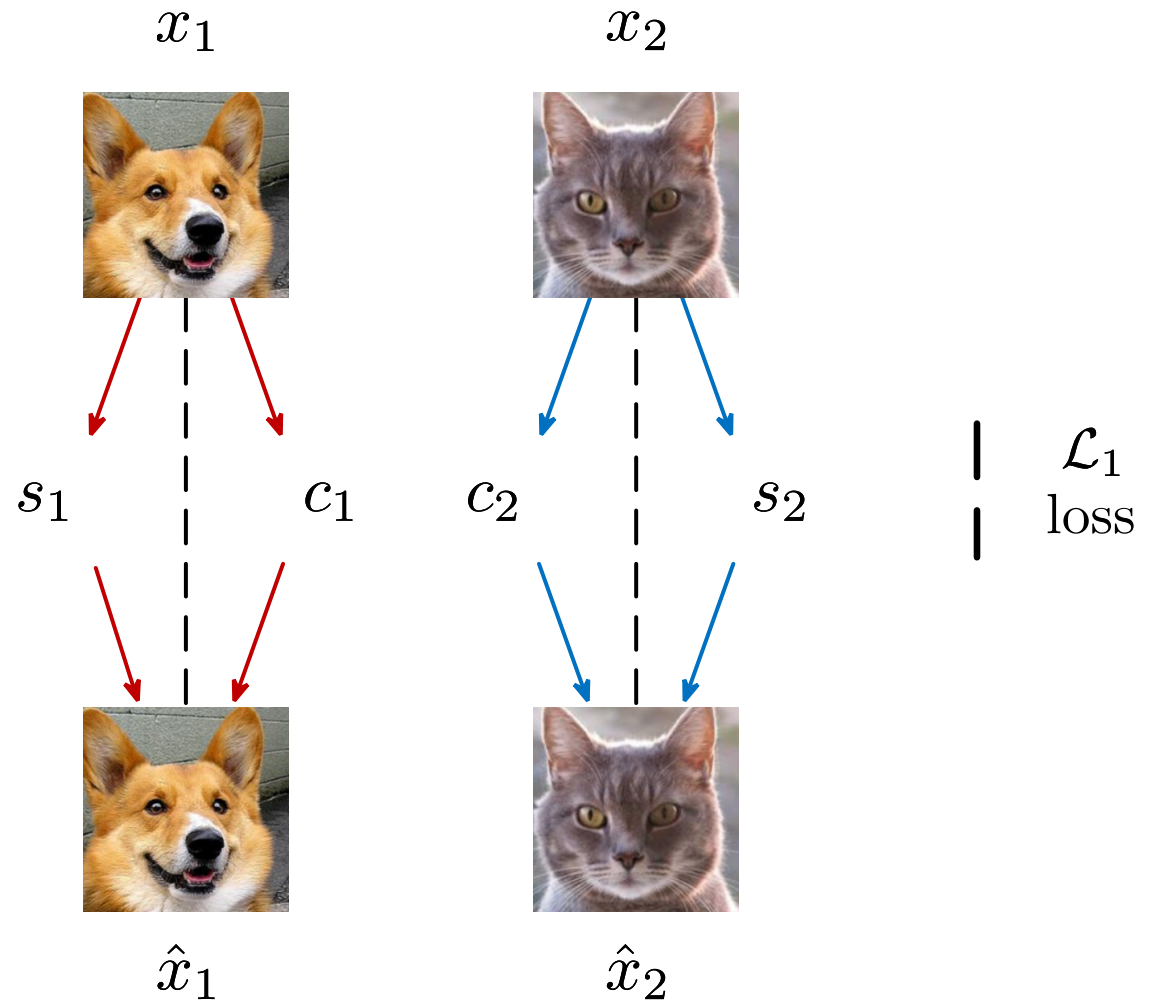    - Latent reconstruction loss
  - GAN loss



Within-domain reconstruction
Cross-domain translation

# Bidirectional Reconstruction Loss:
# Image Reconstruction

Notations:

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

$x_1$ $x_2$



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

$| \quad \mathcal{L}_1$
$| \quad$ loss

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

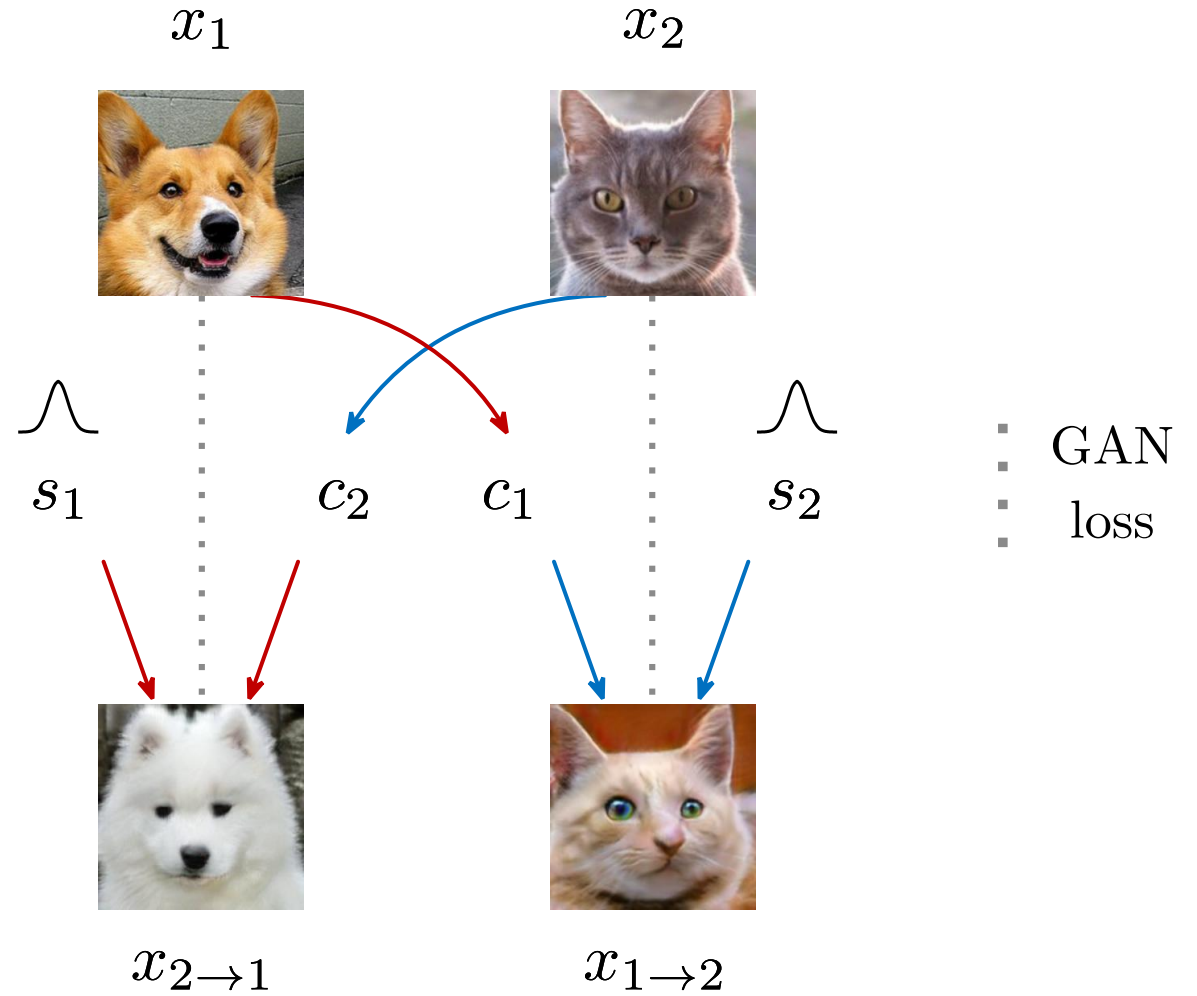# Bidirectional Reconstruction Loss:
# Latent Reconstruction

Notations:

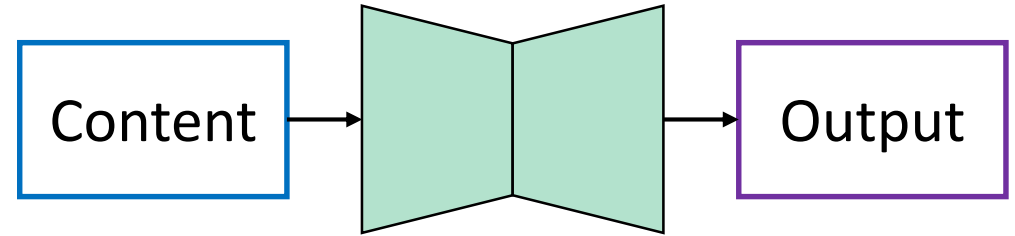- $x$: images
- $c$: content
- $s$: style

# GAN Loss

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

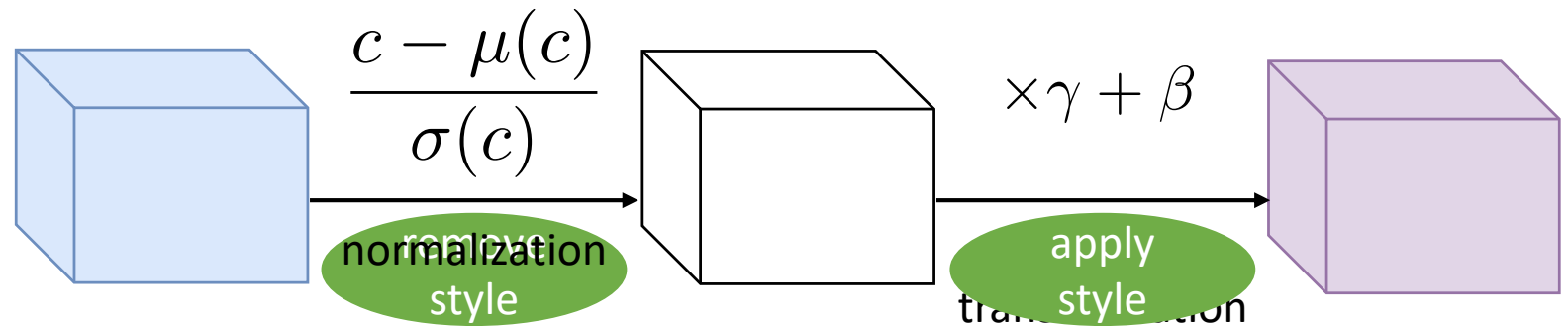# Background: Instance Normalization (IN)
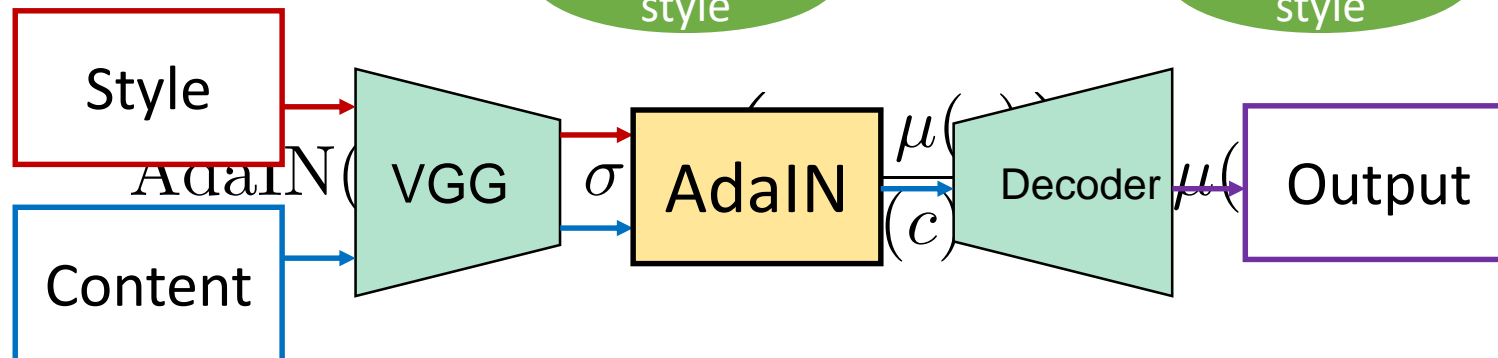
Feedforward transfer of a single style

Content → Output

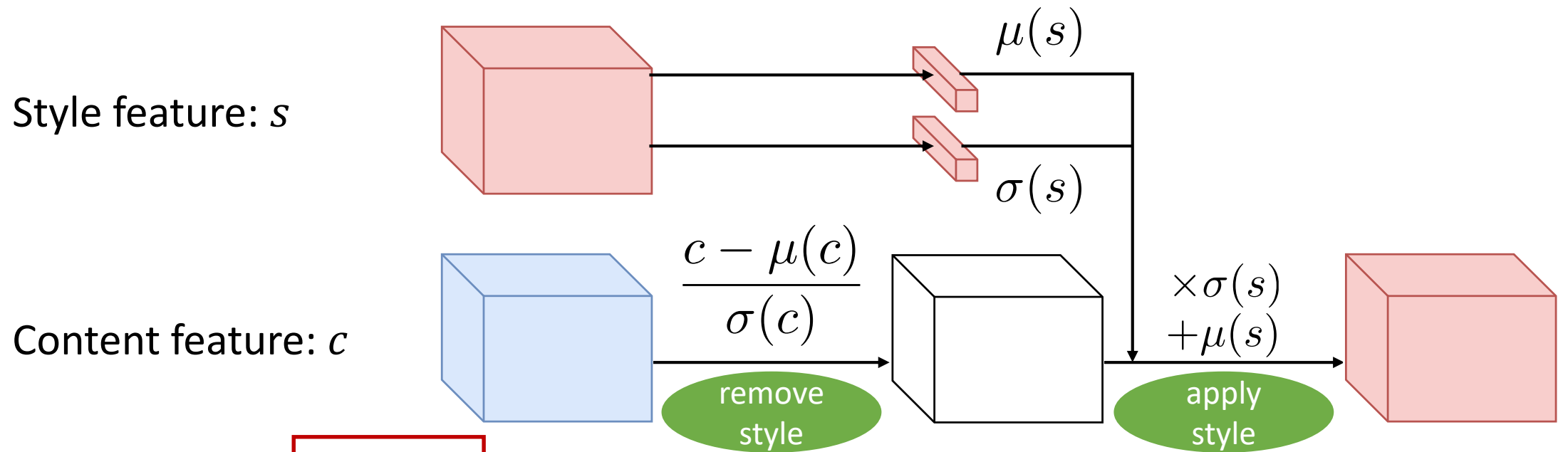Content feature: $c$

$$\frac{c - \mu(c)}{\sigma(c)}$$

$$\times \gamma + \beta$$

normalization style

apply style

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

"Instance Normalization: The Missing Ingredient for Fast Stylization", Ulyanov et al. 2017

# Adaptive Instance Normalization (AdaIN)

Feedforward transfer of **arbitrary** styles

# AdaIN in a Generative Network



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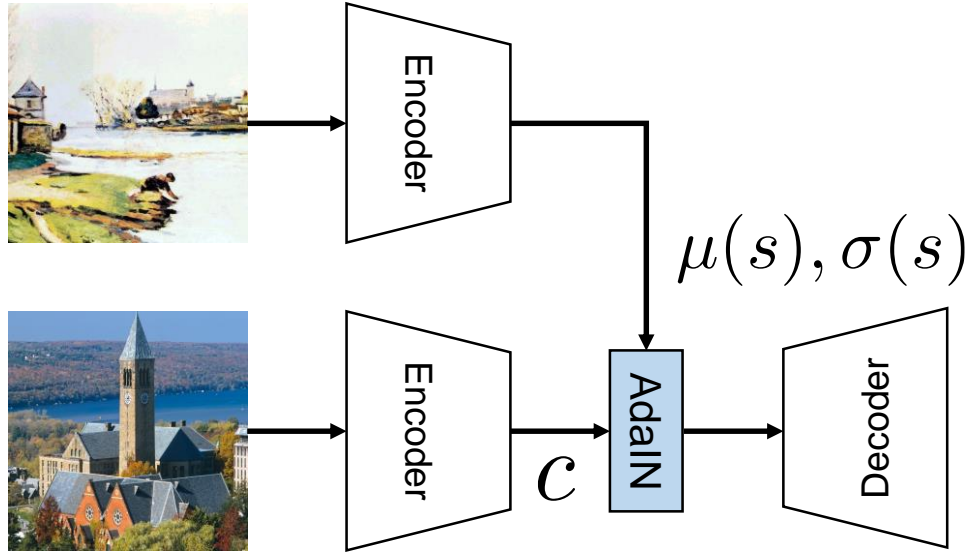
$$\text{AdaIN}(c, s) = \gamma\left(\frac{c - \mu(c)}{\sigma(c)}\right) + \beta$$
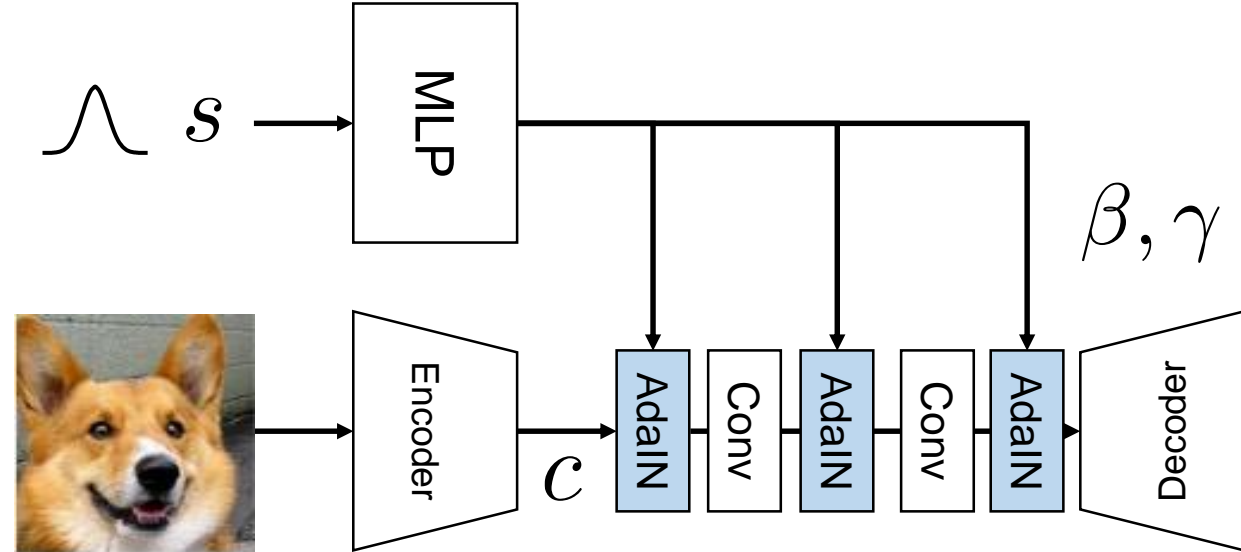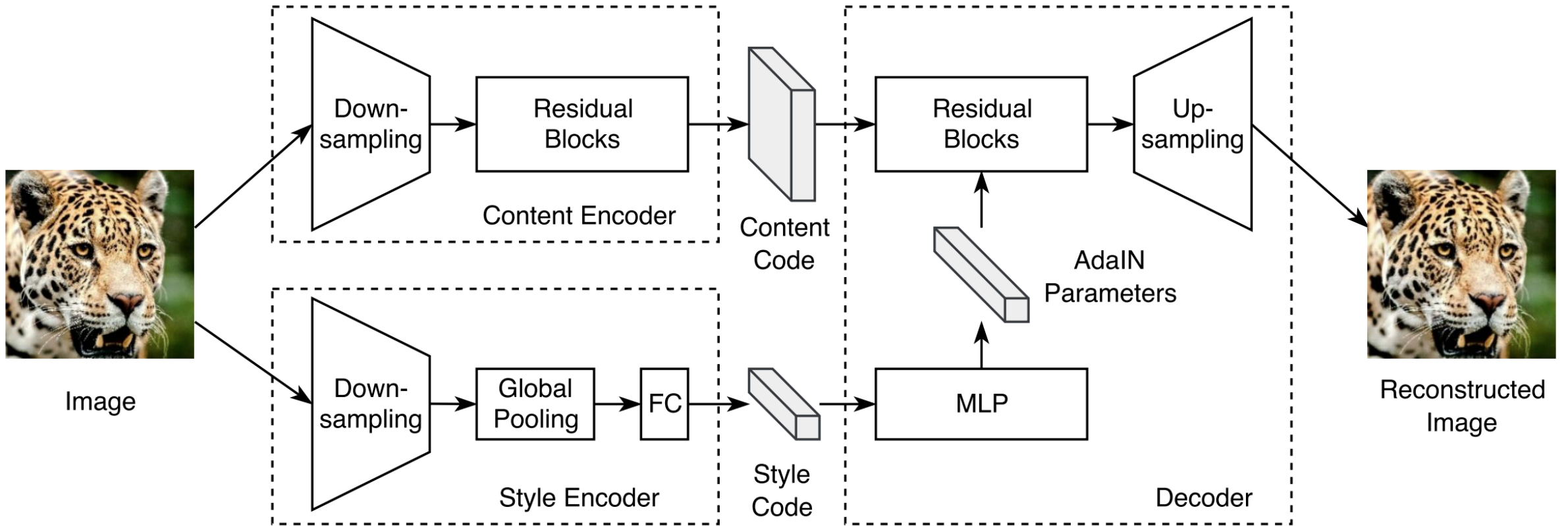
AdaIN in style transfer

AdaIN in a generative network

# AdaIN in a Generative Network



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

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

AdaIN in style transfer

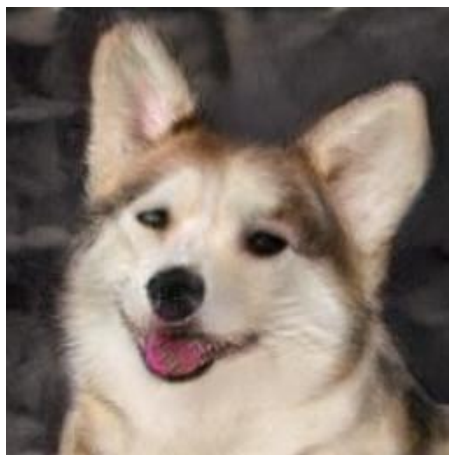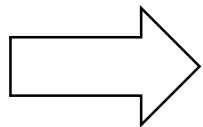AdaIN in a generative network

# Architectural Implementation

# Sketches <-> Photo

Input

Outputs

# Cats ↔ Dogs

Input

Outputs
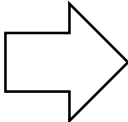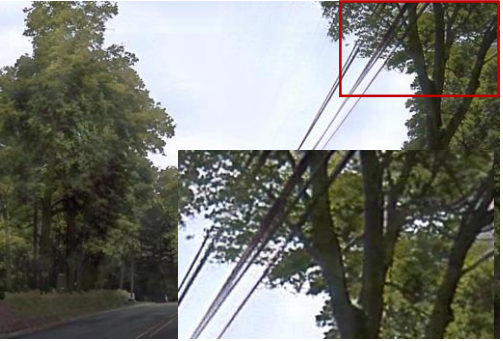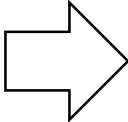
# Synthetic ↔ Real

Input

Outputs

# Summer ↔ Winter

Input

Outputs

# Example-guided Translation

# Example-guided Translation



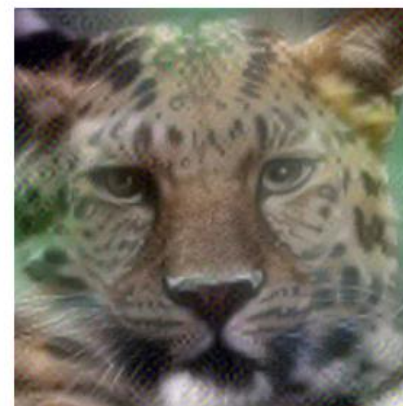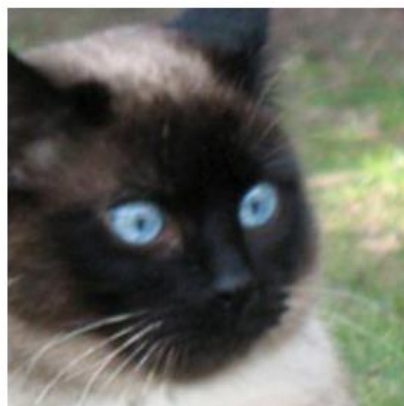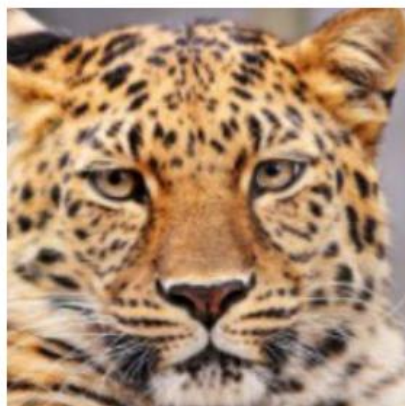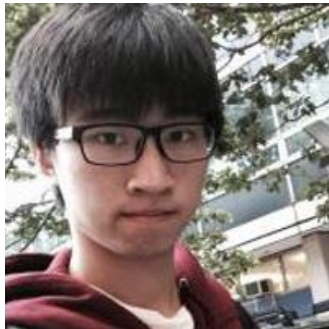| Content | Style | Ours | Gatys *et al.* | AdaIN |

# Conclusion

- Translate one input image to multiple corresponding images in the target domain.

- Content and style decomposition via the AdaIN design

- ECCV 2018

- MUNIT code: https://github.com/nvlabs/munit/

- Paper: https://arxiv.org/abs/1804.04732

Xun Huang
NVIDIA, Cornell

Ming-Yu Liu
NVIDIA

Serge Belongie
Cornell

Jan Kautz
NVIDIA