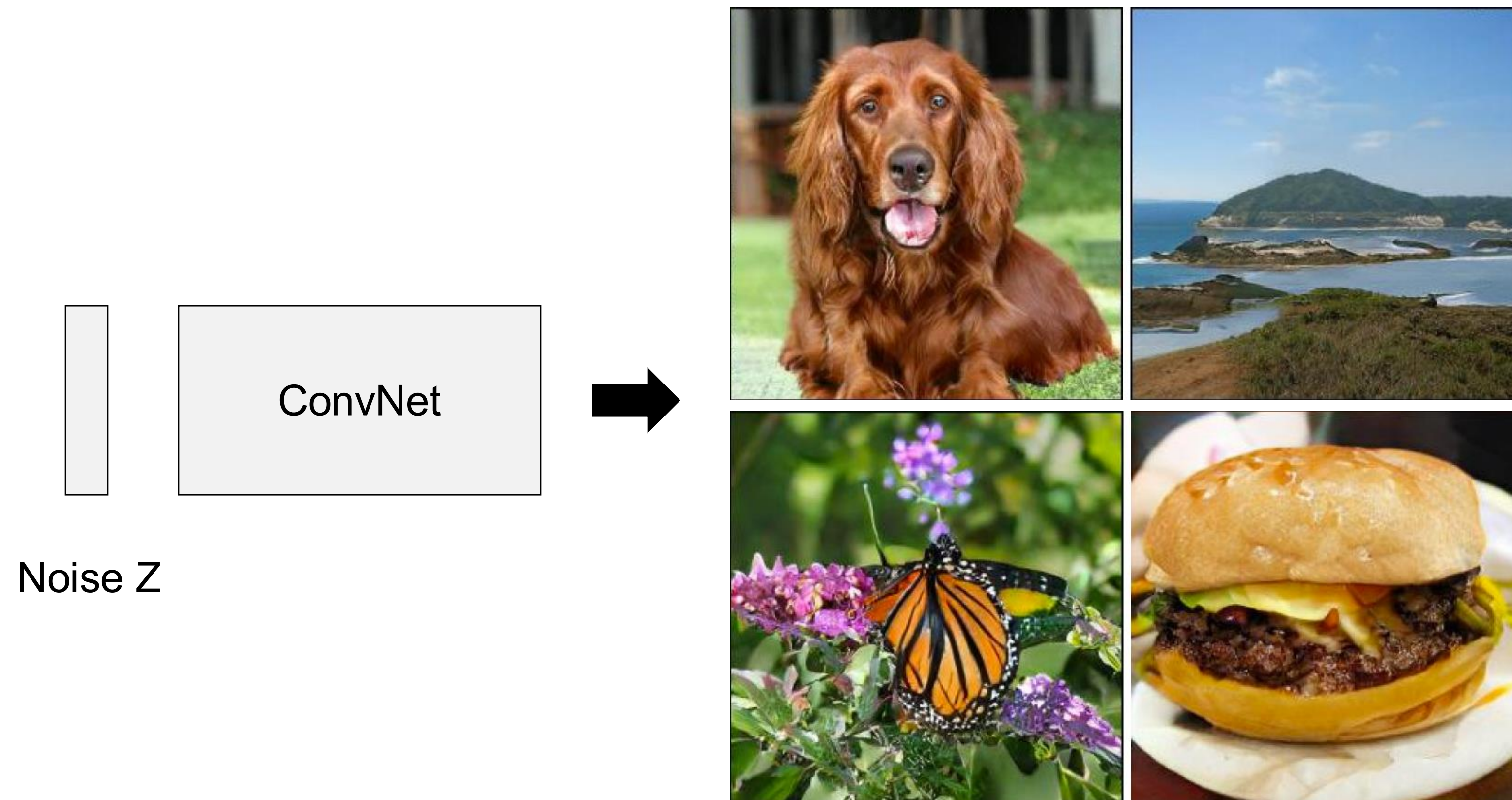


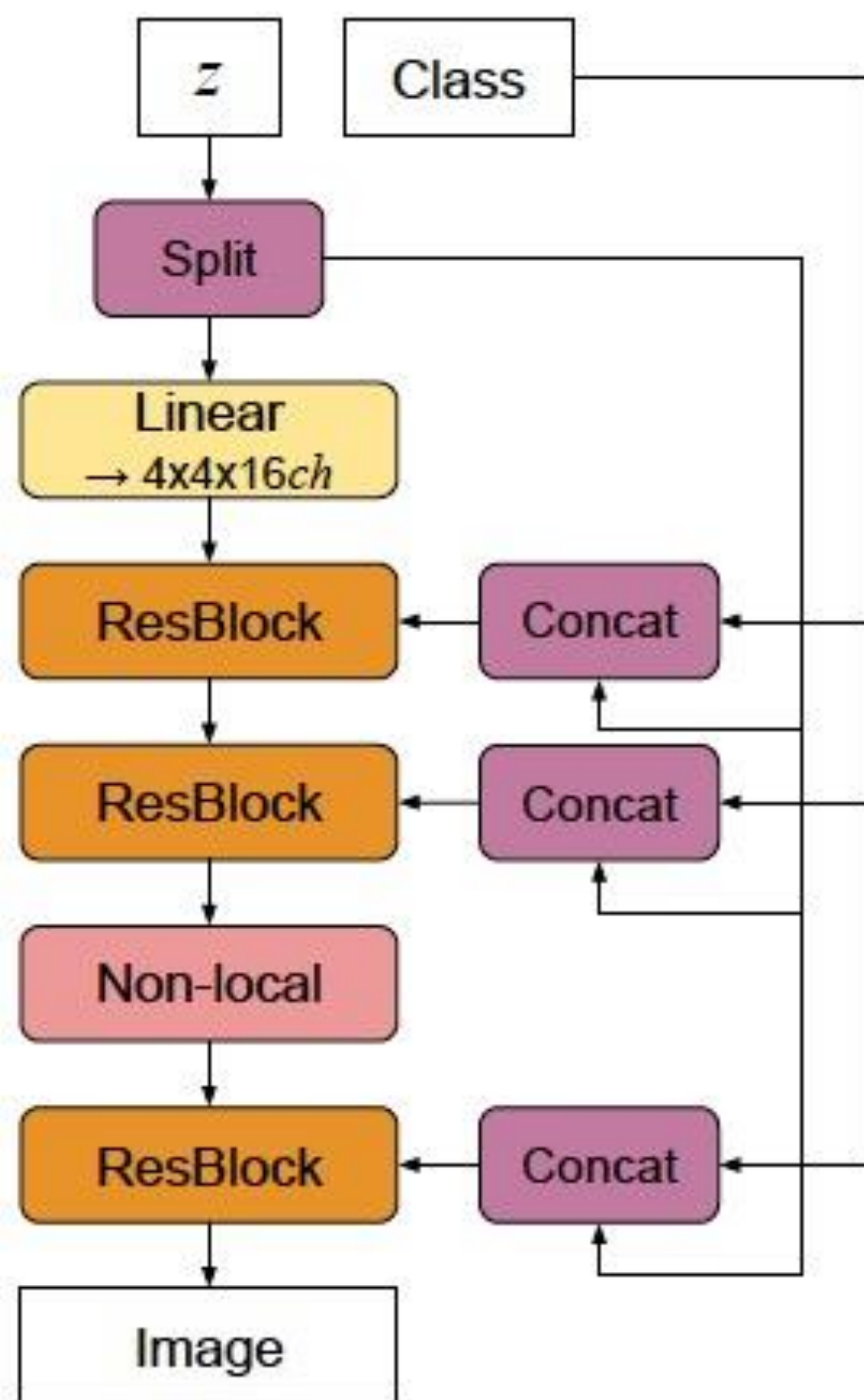
Conditional GAN and Variational Auto-Encoders

Xiaolong Wang

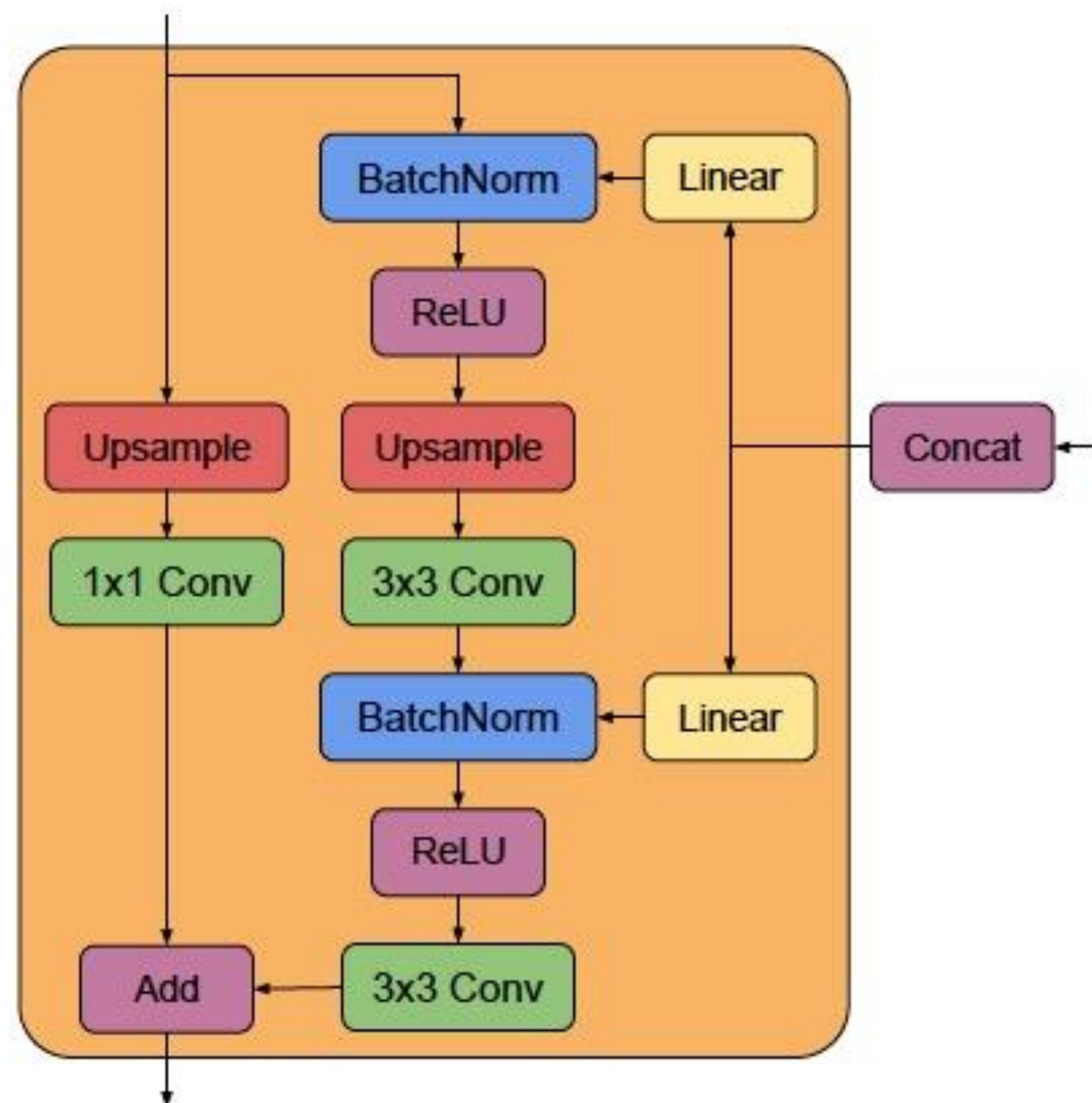
Last class



BigGAN: Class-Conditioned



(a)



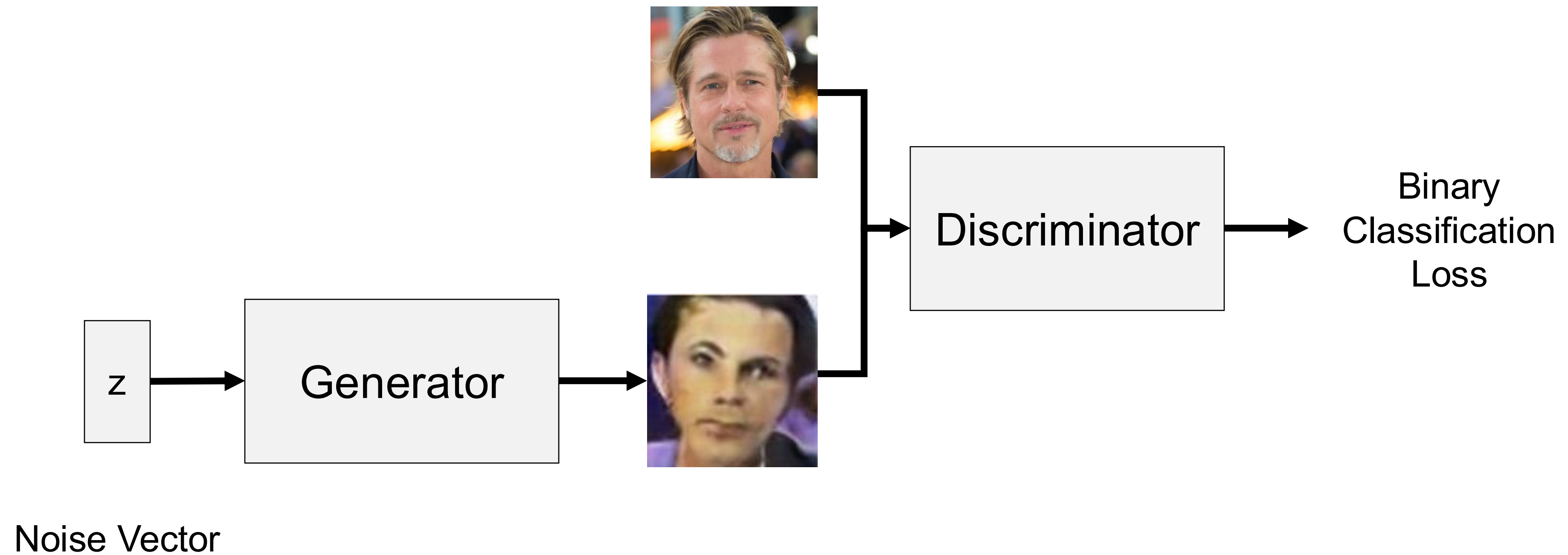
(b)

This Class

- Image-to-Image Translation: pix2pix
- Unpaired Image-to-Image Translation: CycleGAN
- Variational Autoencoder (VAE)

Image-to-Image Translation: pix2pix

GANs



Conditional GANs

Edges to Photo



input



output

BW to Color

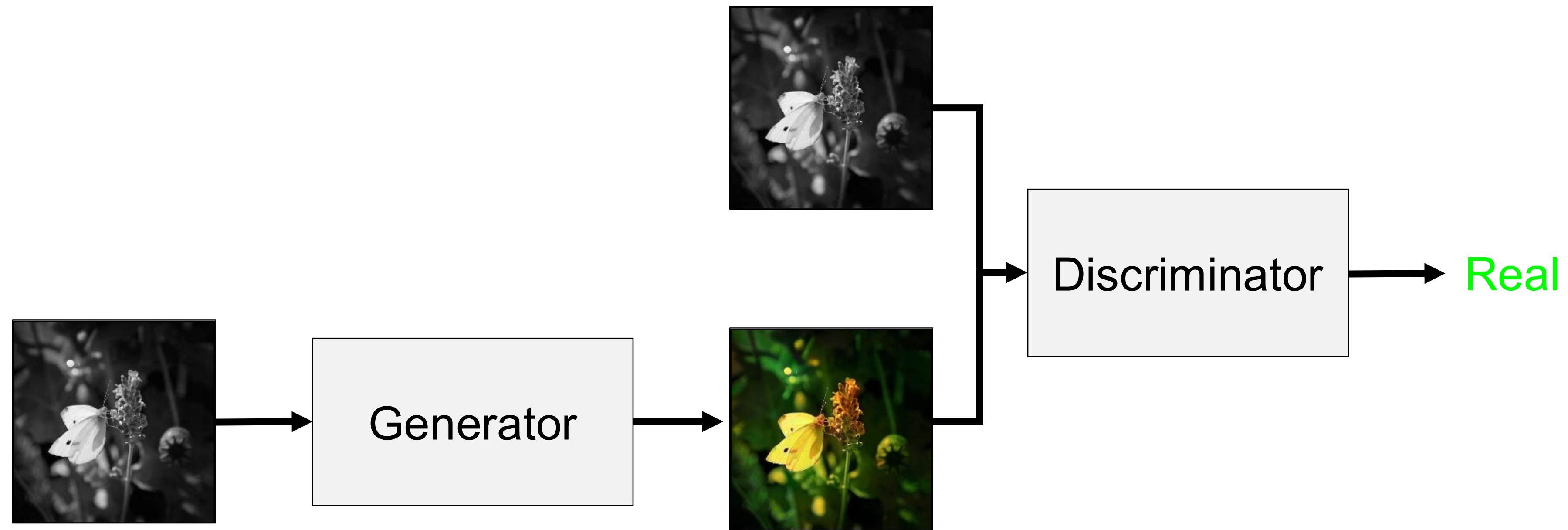


input



output

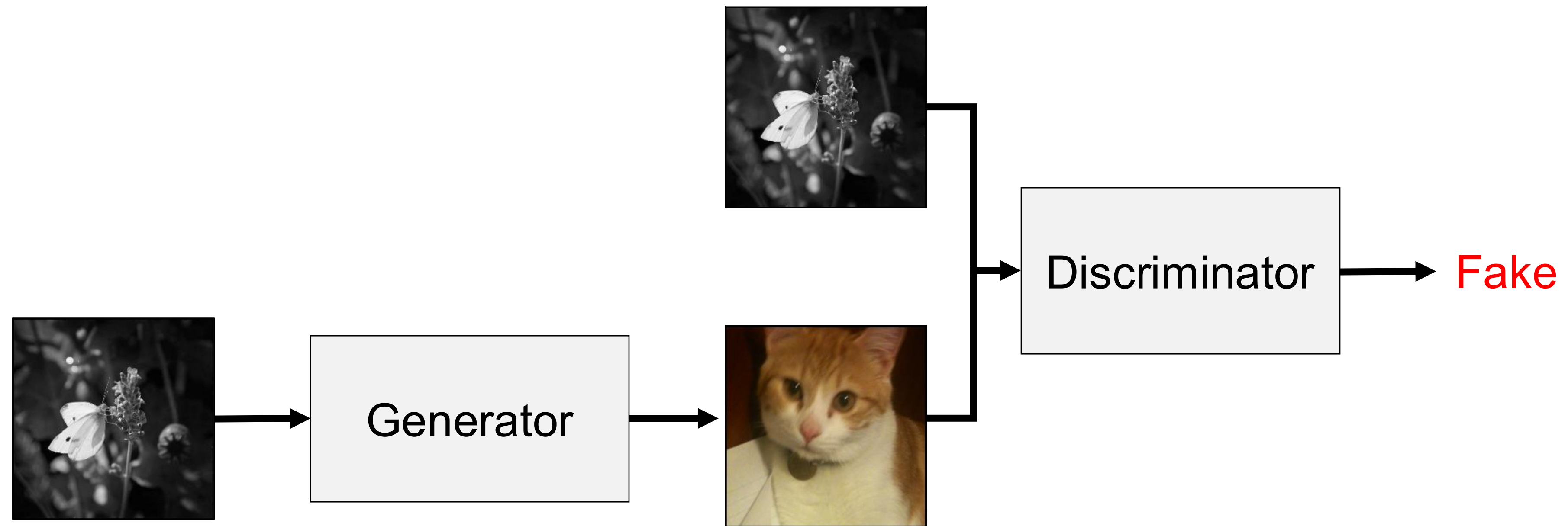
Conditional GANs



Generator takes an image as input, not noise.

Discriminator takes a pair of images as inputs, not just one image.

Conditional GANs



Generator takes an image as input, not noise.

Discriminator takes a pair of images as inputs, not just one image.

Pix2Pix

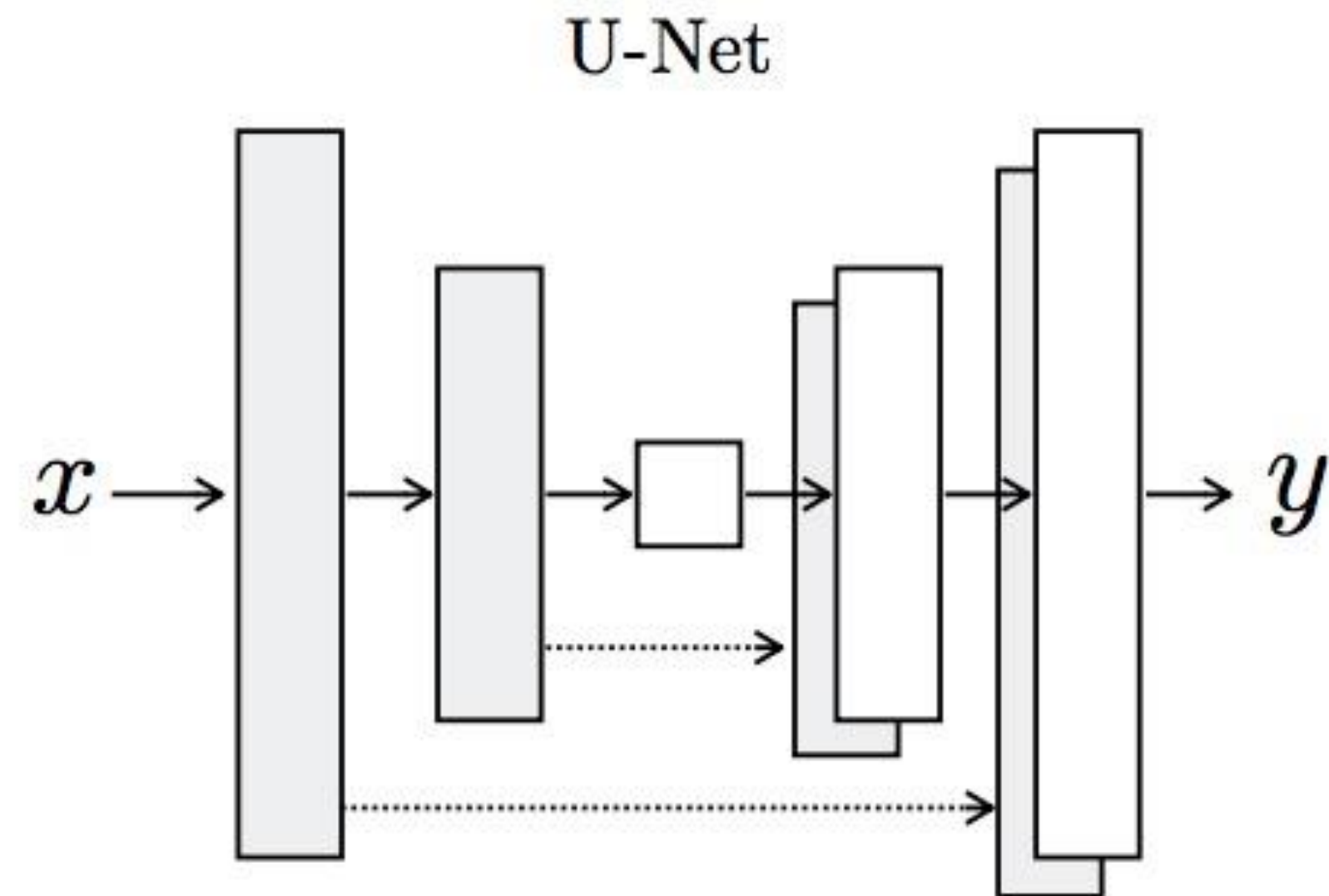
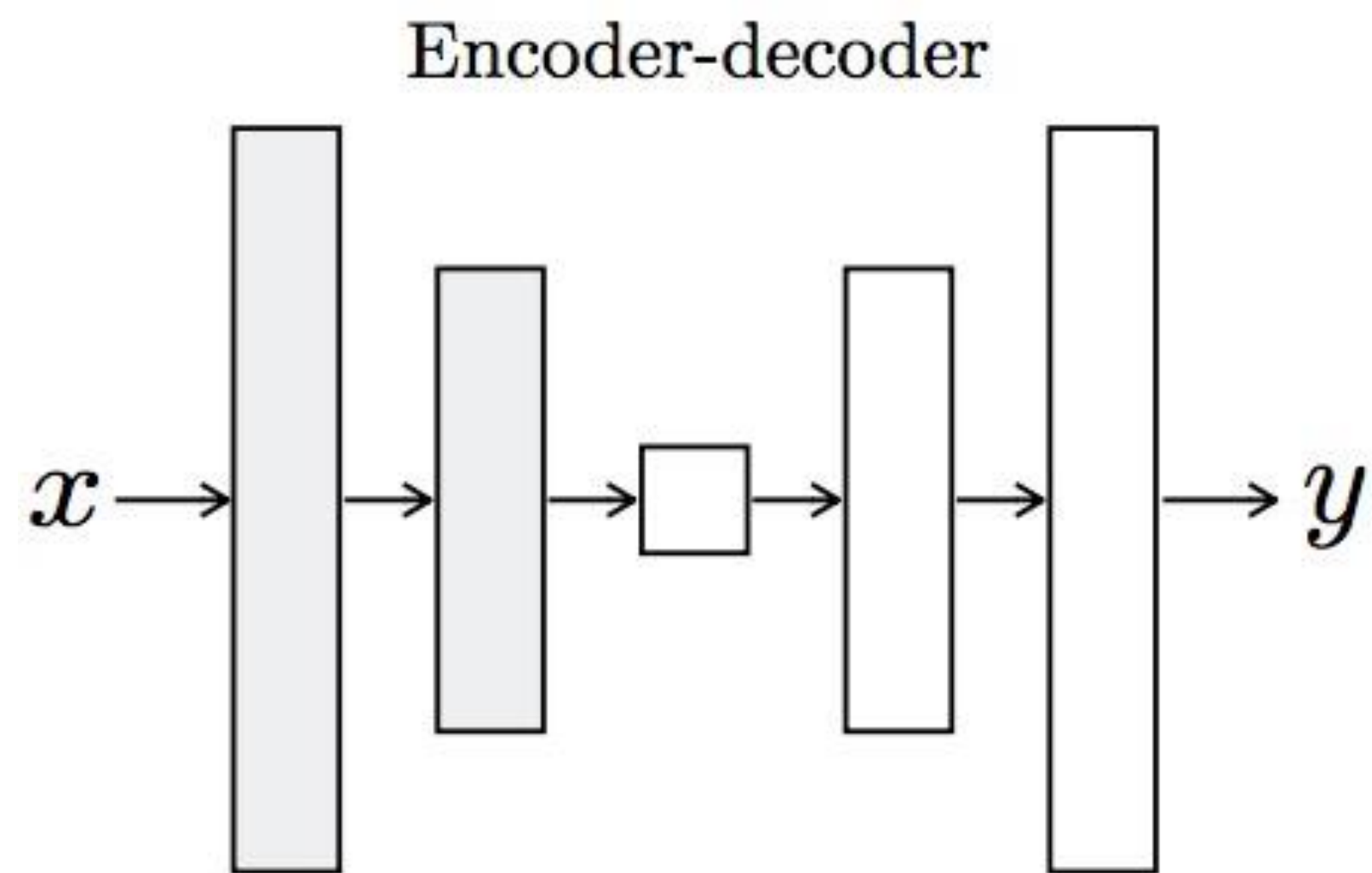
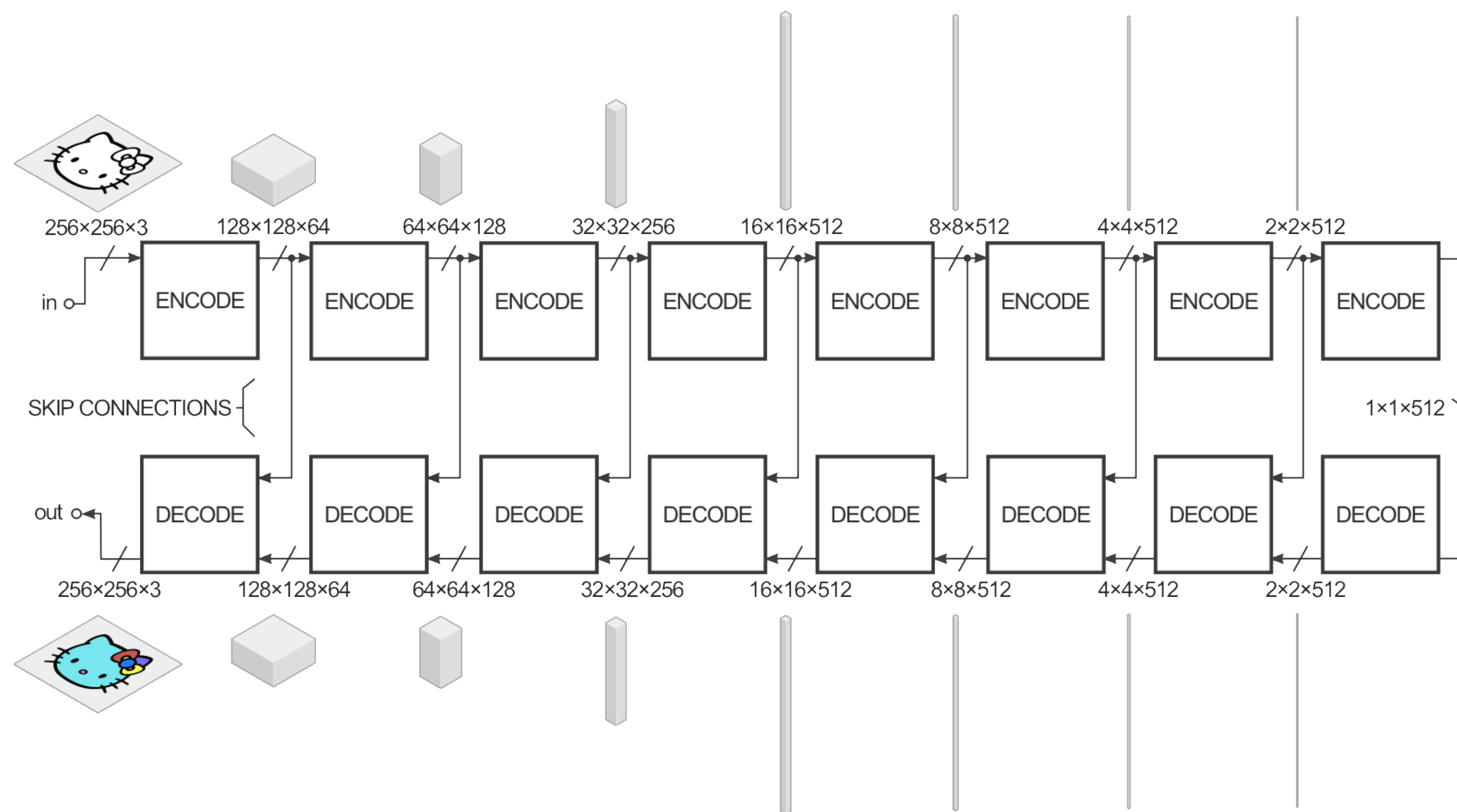


Image-to-image translation



Encode: convolution → BatchNorm → ReLU

Decode: transposed convolution → BatchNorm → ReLU

Image-to-image translation

Effect of adding skip connections to the generator



Image-to-image translation

- Generator loss: GAN loss plus L1 reconstruction penalty

$$G^* = \arg \min_G (\max_D \mathcal{L}_{GAN}(G, D) + \lambda \sum_i \|y_i - G(x_i)\|_1)$$

Generated output
 $G(x_i)$ should be close to
ground truth target y_i

Image-to-image translation

- Generator loss: GAN loss plus L1 reconstruction penalty

$$G^* = \arg \min_G (\max_D \mathcal{L}_{GAN}(G, D) + \lambda \sum \|y_i - G(x_i)\|_1)$$



Image-to-image translation: Results

- Day to night

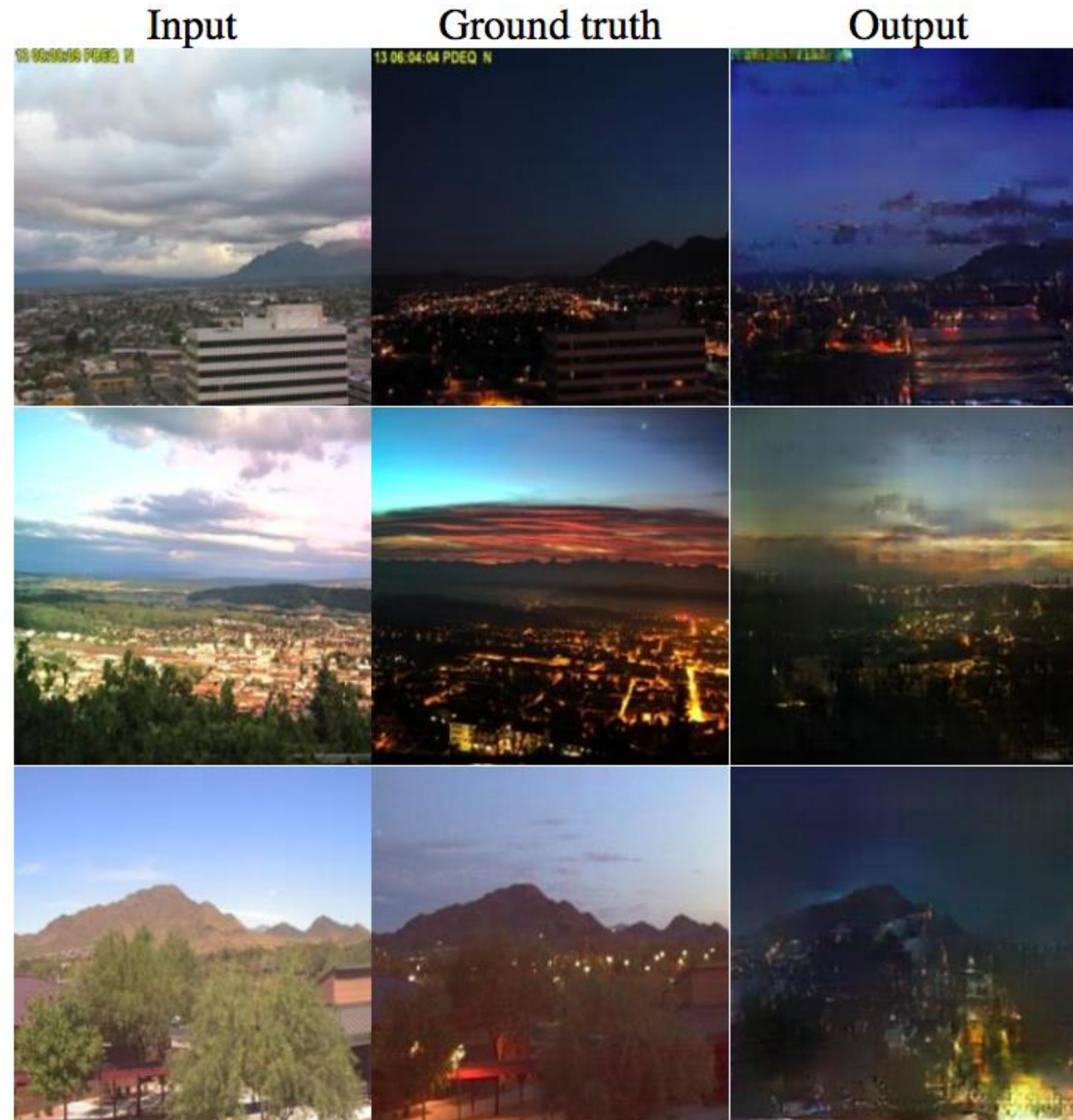


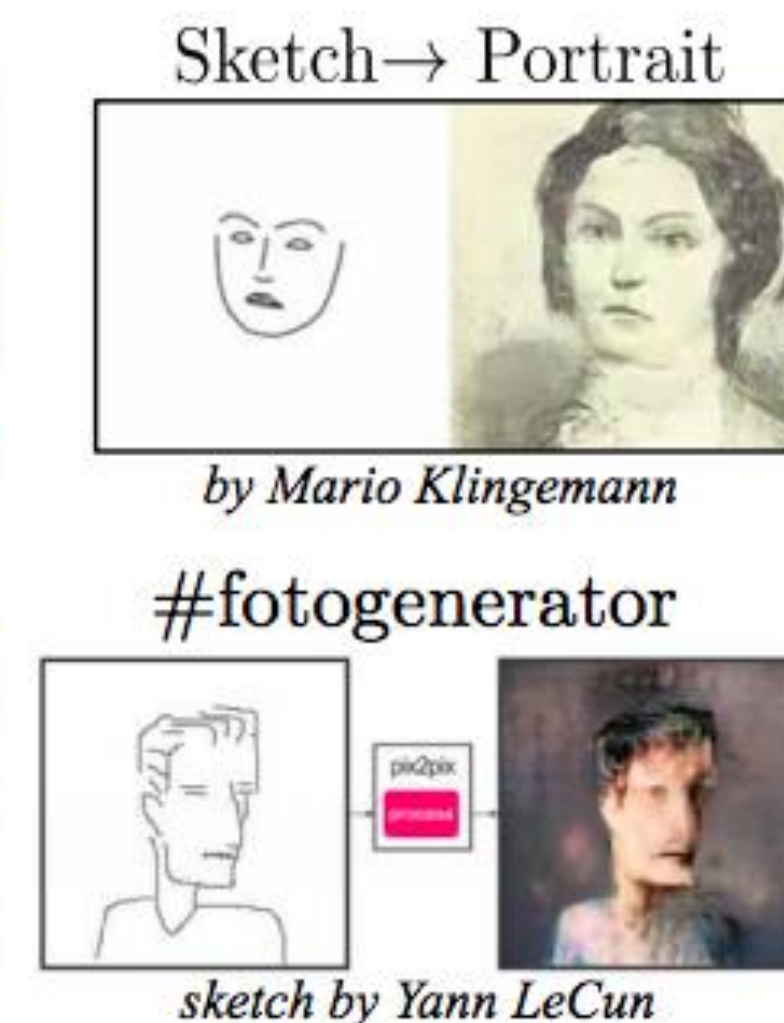
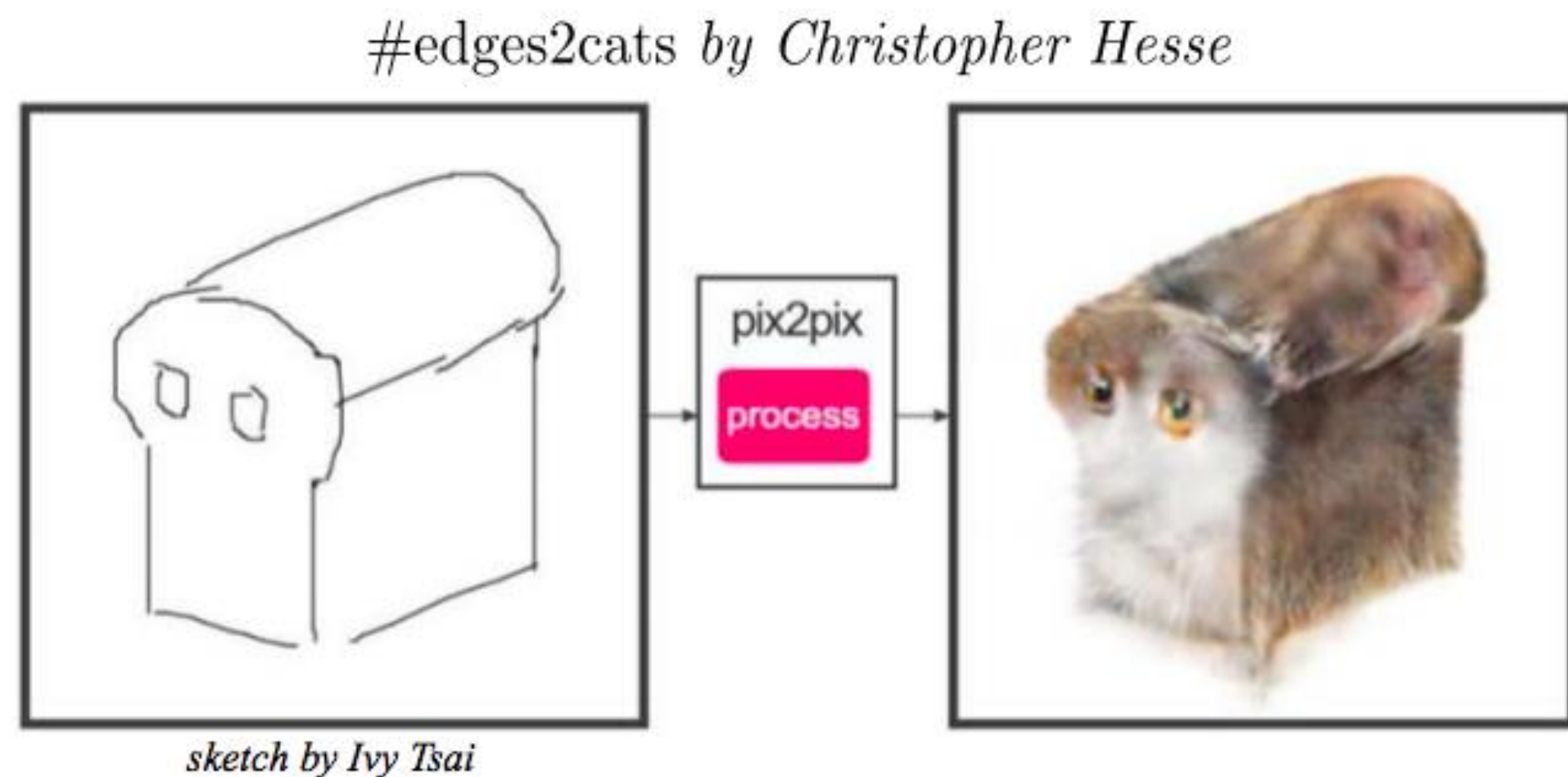
Image-to-image translation: Results

- Edges



Image-to-image translation: Results

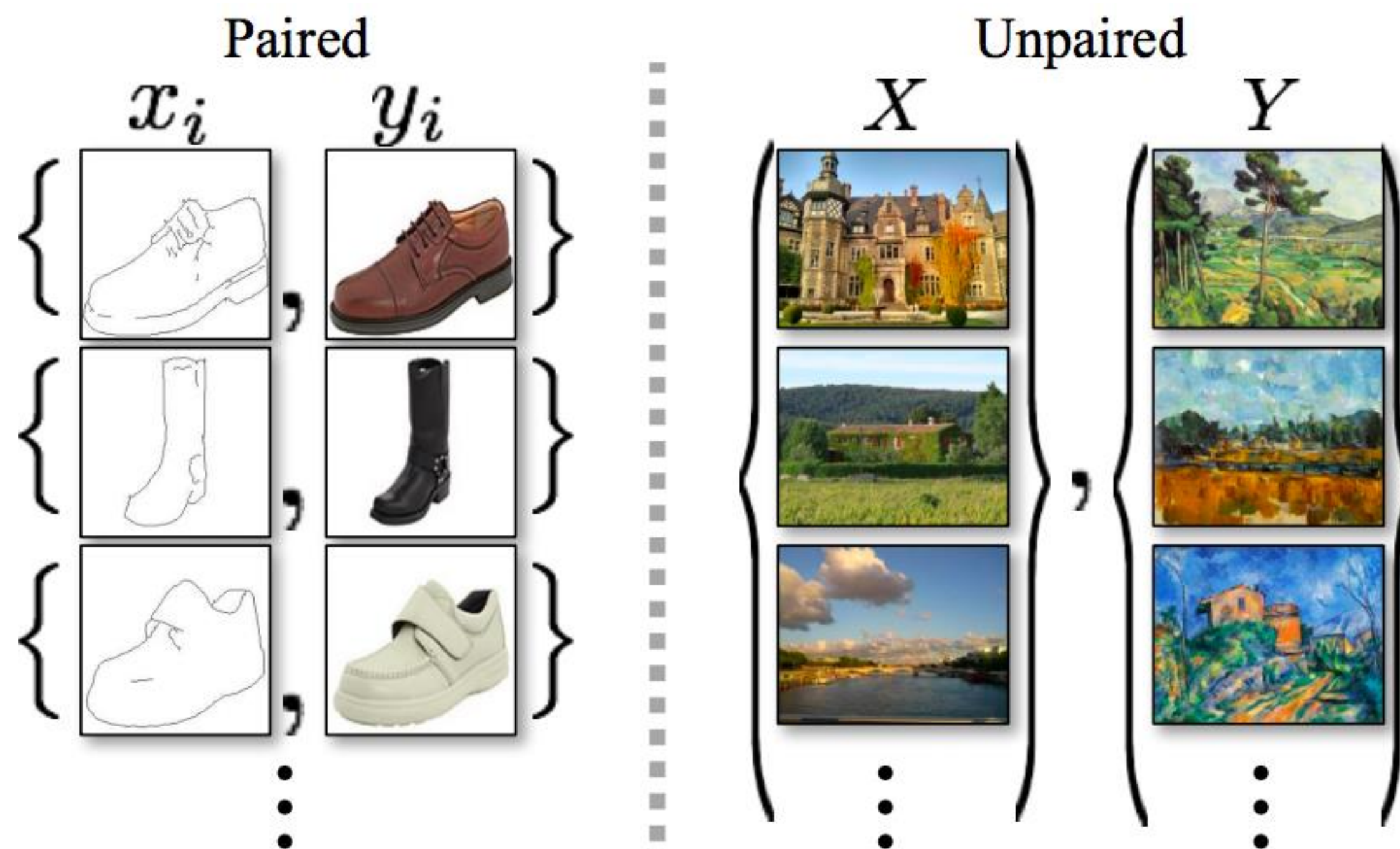
- [pix2pix demo](#)



Unpaired Image-to-Image Translation: CycleGAN

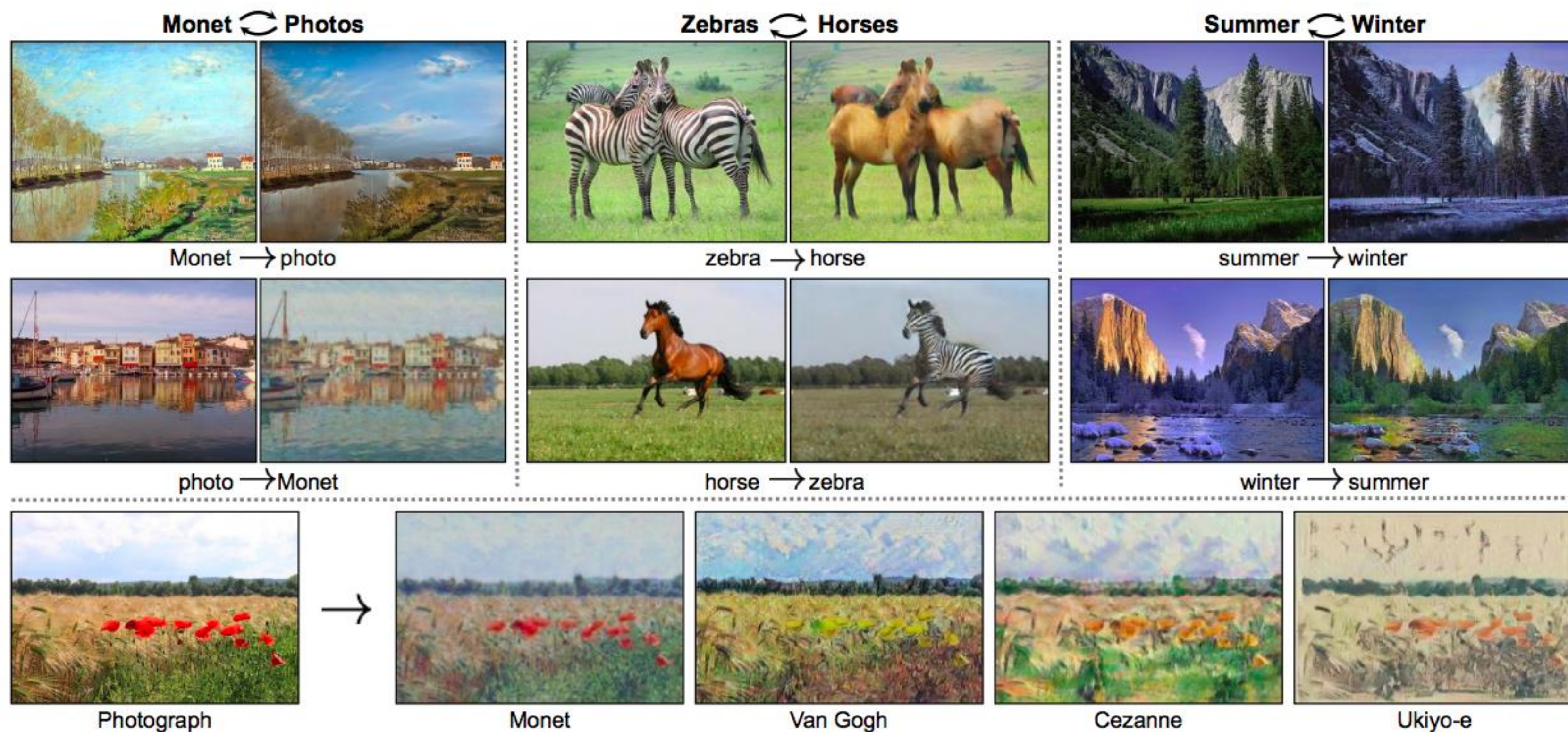
Unpaired image-to-image translation

- Given two unordered image collections X and Y , learn to “translate” an image from one into the other and vice versa

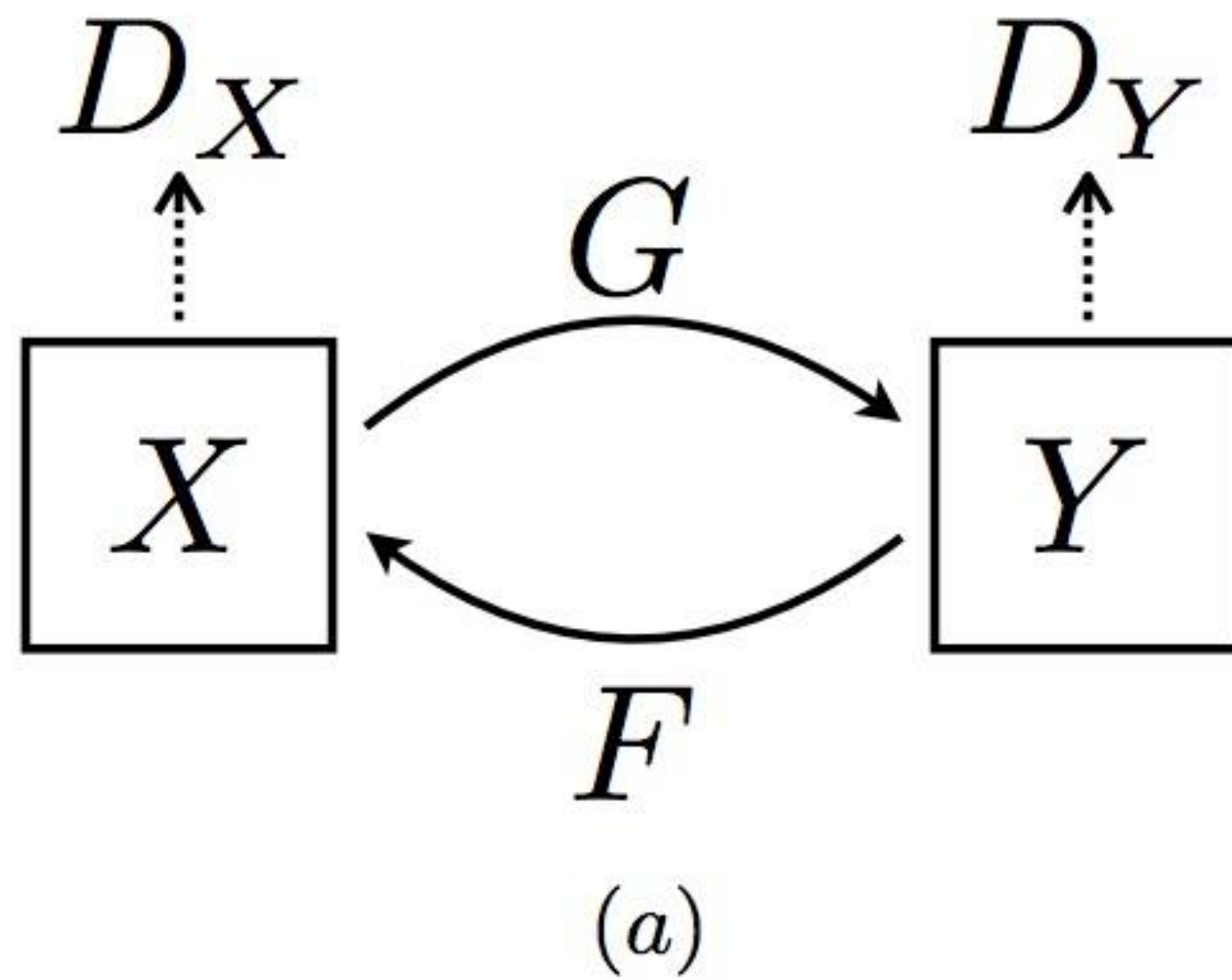


Unpaired image-to-image translation

- Given two unordered image collections X and Y , learn to “translate” an image from one into the other and vice versa



CycleGAN



CycleGAN: Loss

- Requirements:
 - G translates from X to Y , F translates from Y to X
 - D_X recognizes images from X , D_Y from Y
 - We want $F(G(x)) \approx x$ and $G(F(y)) \approx y$
- CycleGAN discriminator loss: LSGAN

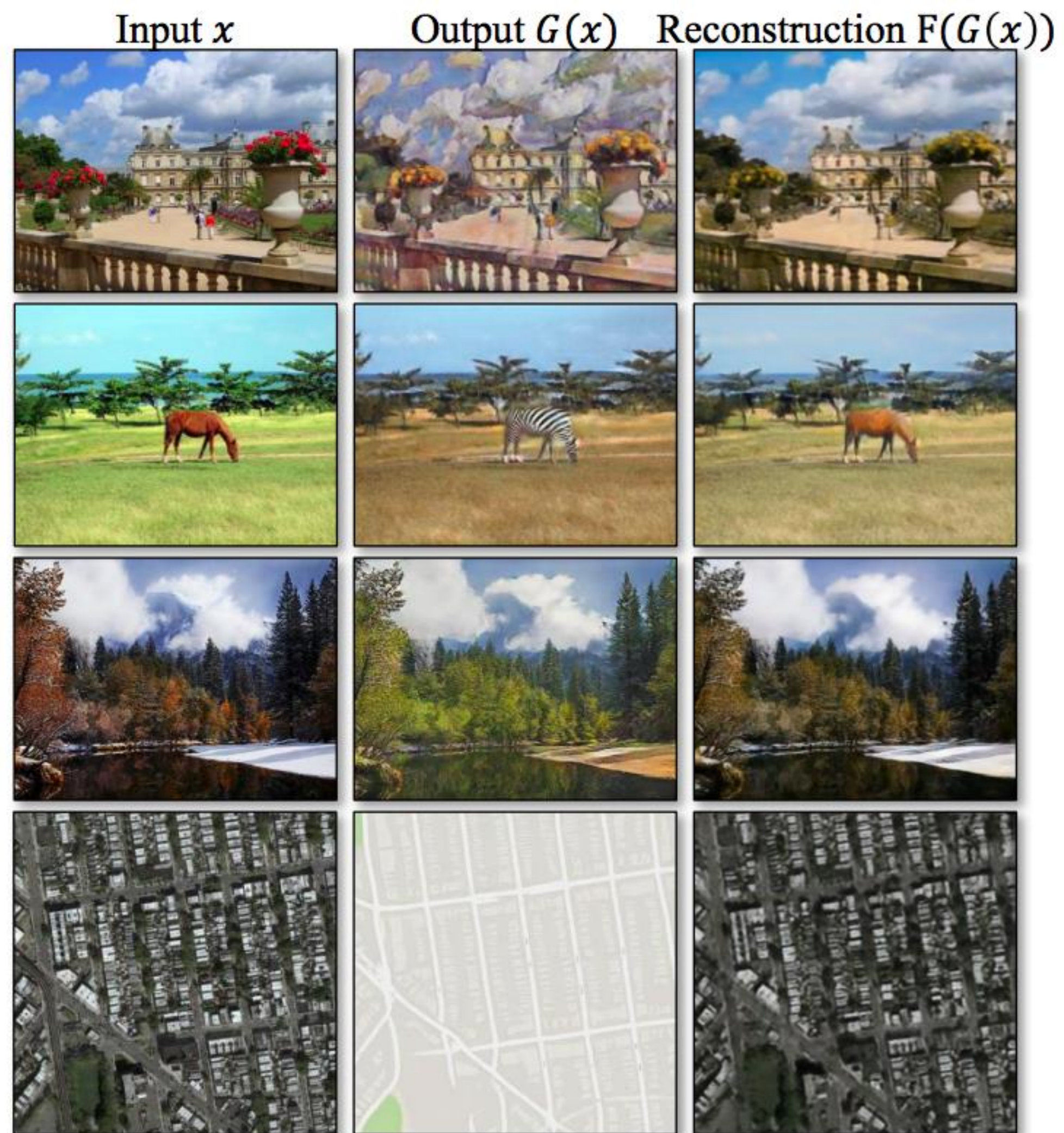
$$\mathcal{L}_{\text{GAN}}(D_Y) = \mathbb{E}_{y \sim p_{\text{data}}(y)} [(D_Y(y) - 1)^2] + \mathbb{E}_{x \sim p_{\text{data}}(x)} [D_Y(G(x))^2]$$

$$\mathcal{L}_{\text{GAN}}(D_X) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D_X(x) - 1)^2] + \mathbb{E}_{y \sim p_{\text{data}}(y)} [D_X(F(y))^2]$$

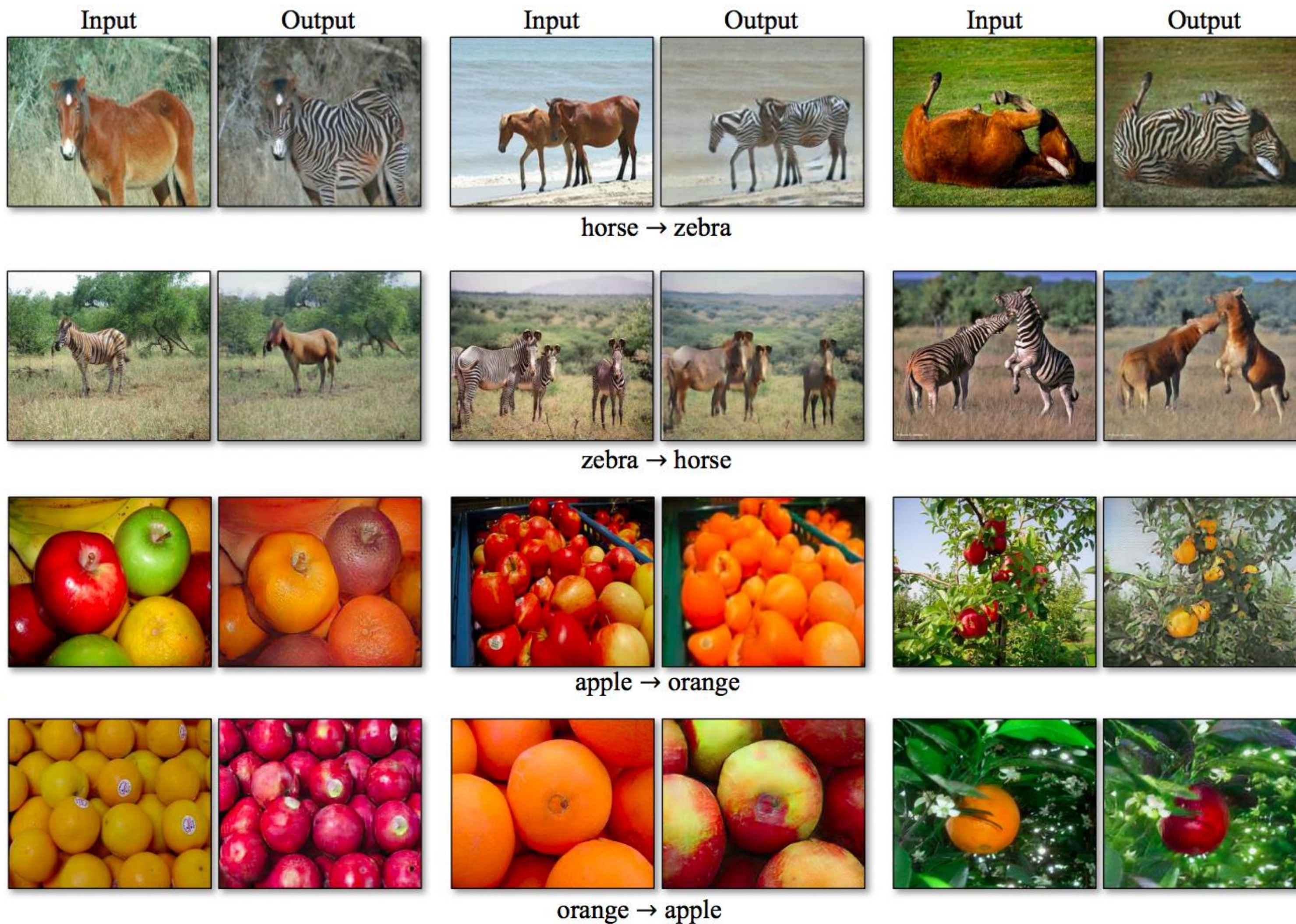
- CycleGAN generator loss:

$$\begin{aligned} \mathcal{L}_{\text{cyc}}(G, F) = & \mathbb{E}_{x \sim p_{\text{data}}(x)} [(D_Y(G(x)) - 1)^2] + \mathbb{E}_{y \sim p_{\text{data}}(y)} [(D_X(F(y)) - 1)^2] \\ & + \mathbb{E}_{x \sim p_{\text{data}}(x)} [\|F(G(x)) - x\|_1] + \mathbb{E}_{y \sim p_{\text{data}}(y)} [\|G(F(y)) - y\|_1] \end{aligned}$$

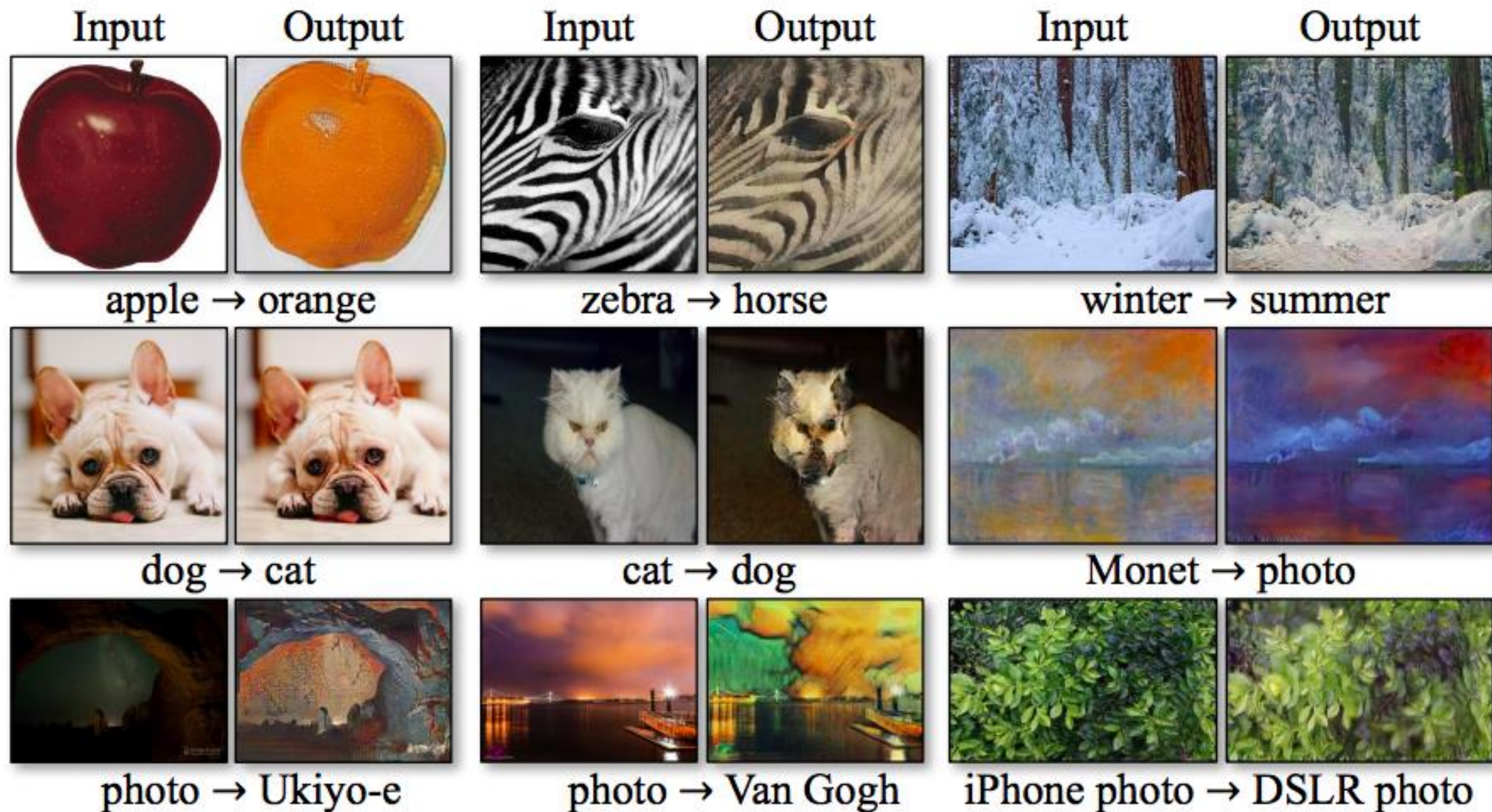
CycleGAN



CycleGAN: Results

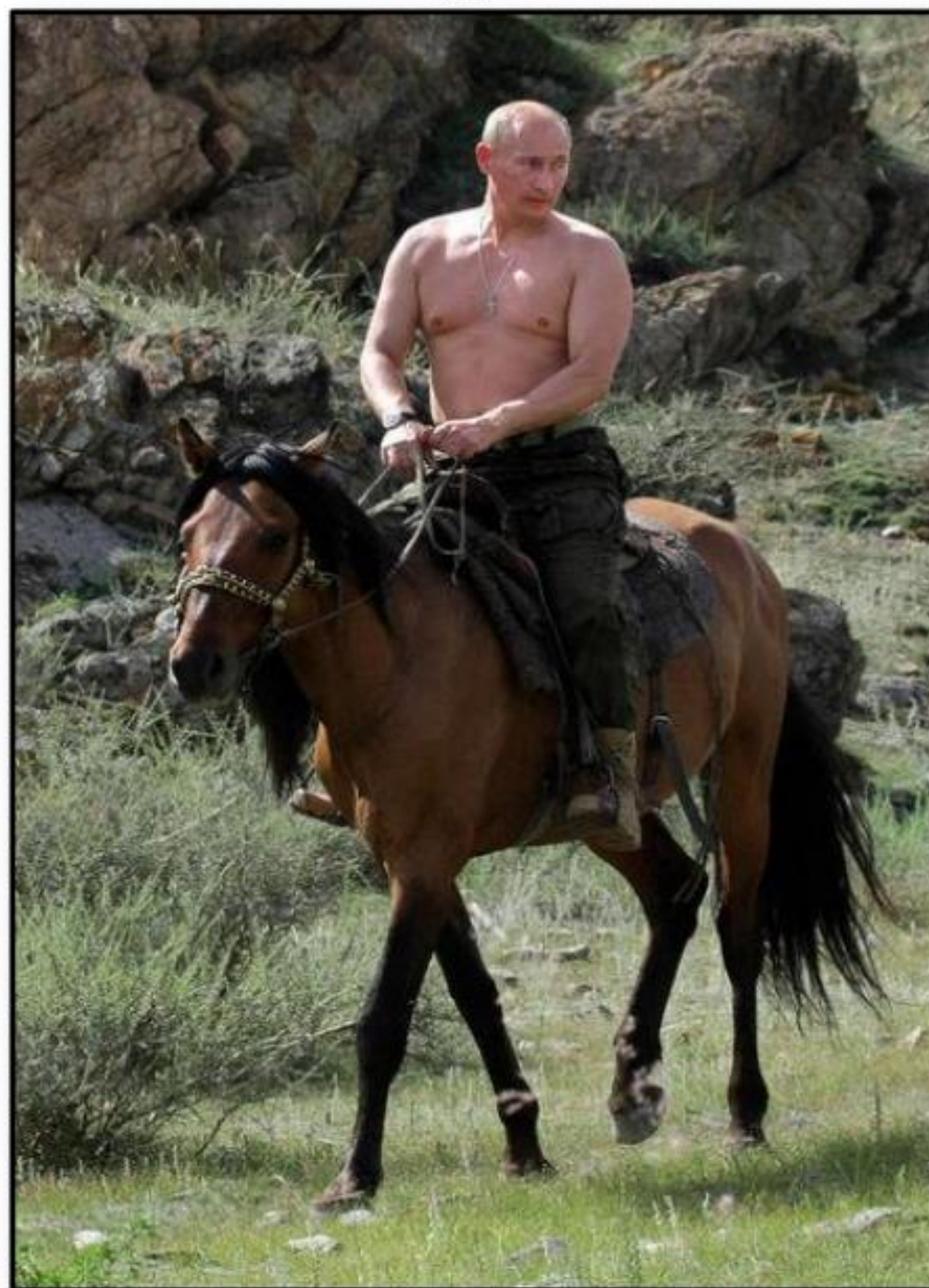


CycleGAN: Failure cases

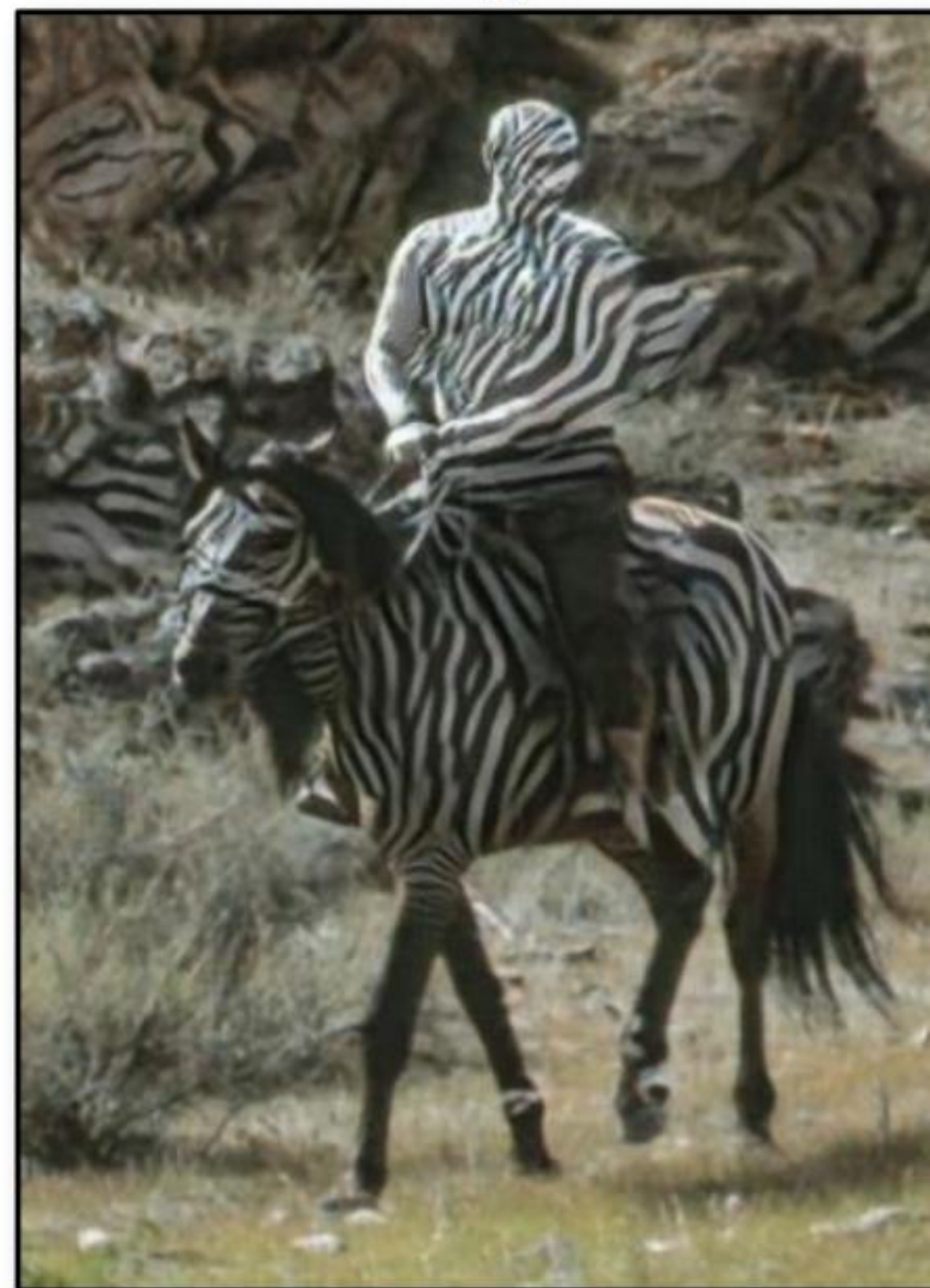


CycleGAN: Failure cases

Input



Output

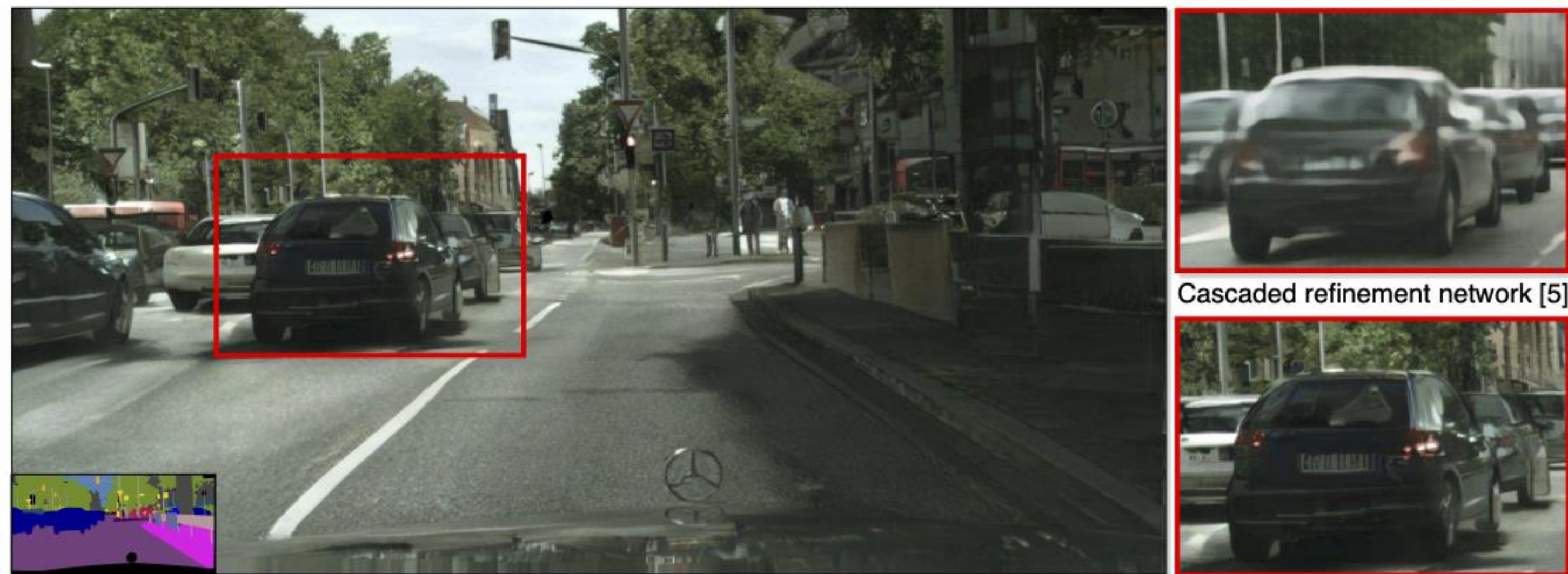


horse → zebra

CycleGAN: Limitations

- Cannot handle shape changes (e.g., dog to cat)
- Can get confused on images outside of the training domains (e.g., horse with rider)
- Cannot close the gap with paired translation methods

High-resolution, high-quality pix2pix



(a) Synthesized result

Our result



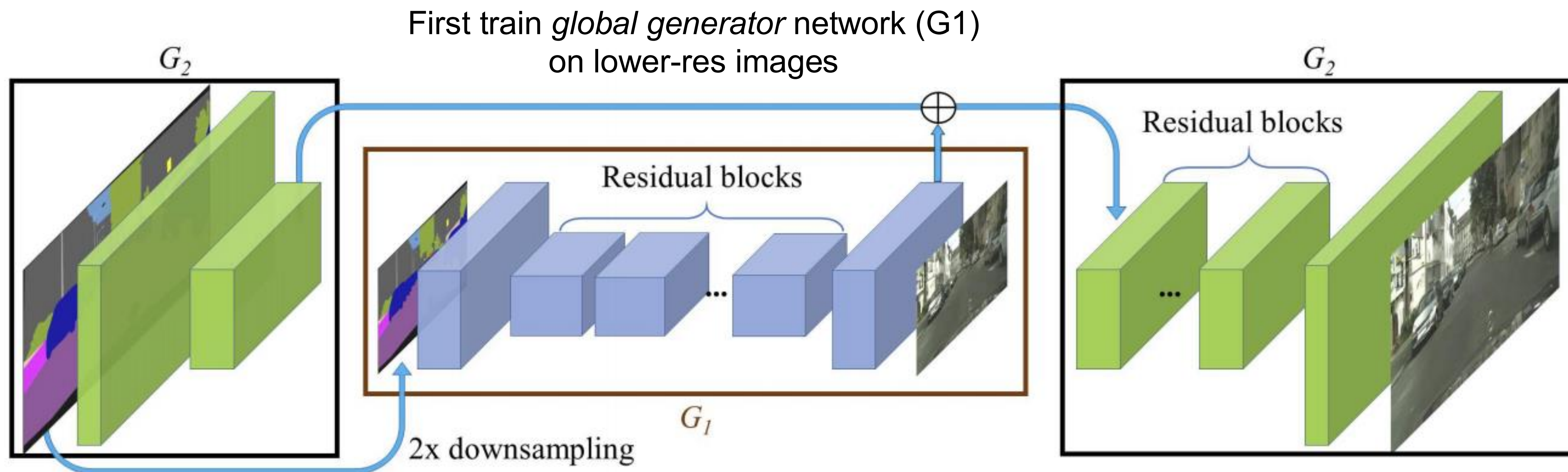
(b) Application: Change label types



(c) Application: Edit object appearance

High-resolution, high-quality pix2pix

- Two-scale generator architecture (up to 2048 x 1024 resolution)



Then append higher-res *enhancer network* (G_2) blocks and train G_1 and G_2 jointly



Human generation conditioned on pose

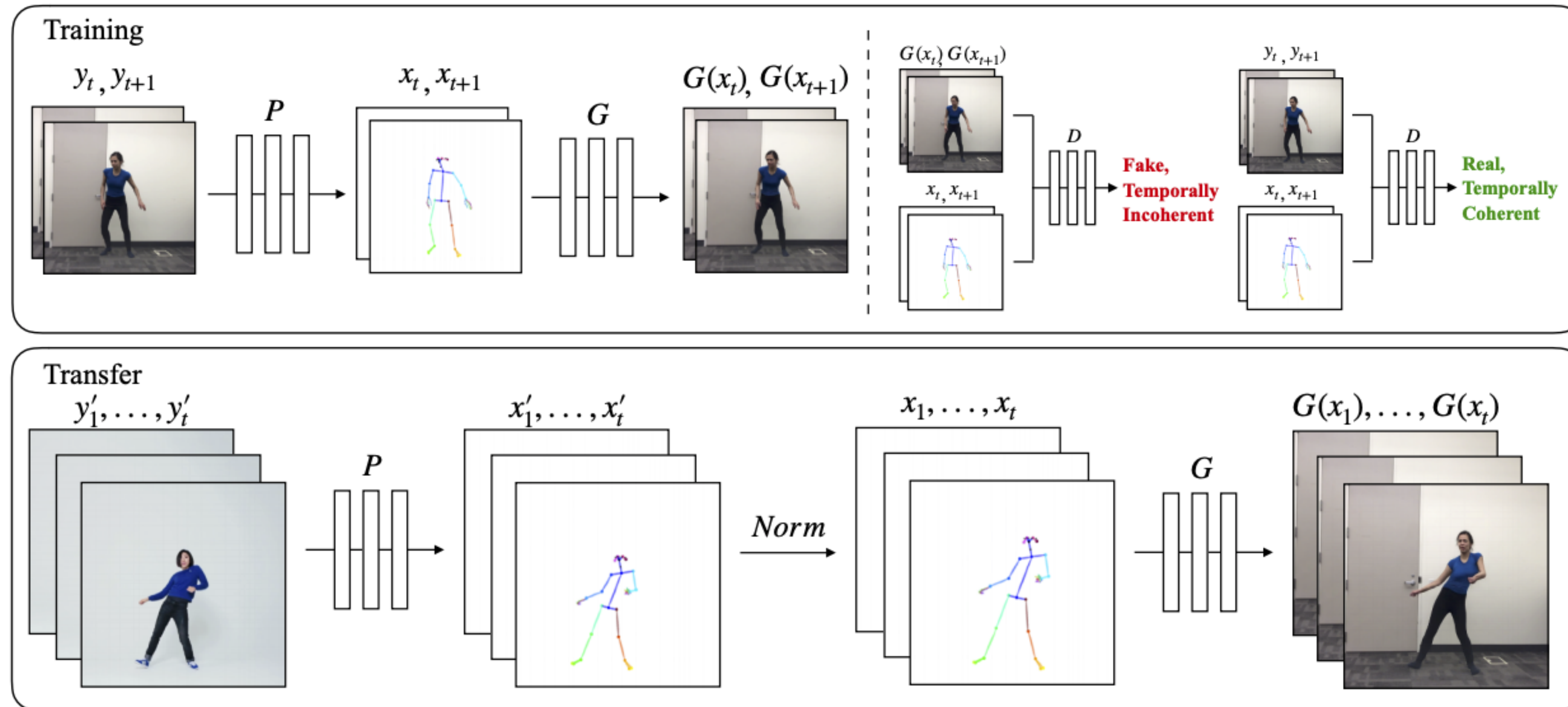
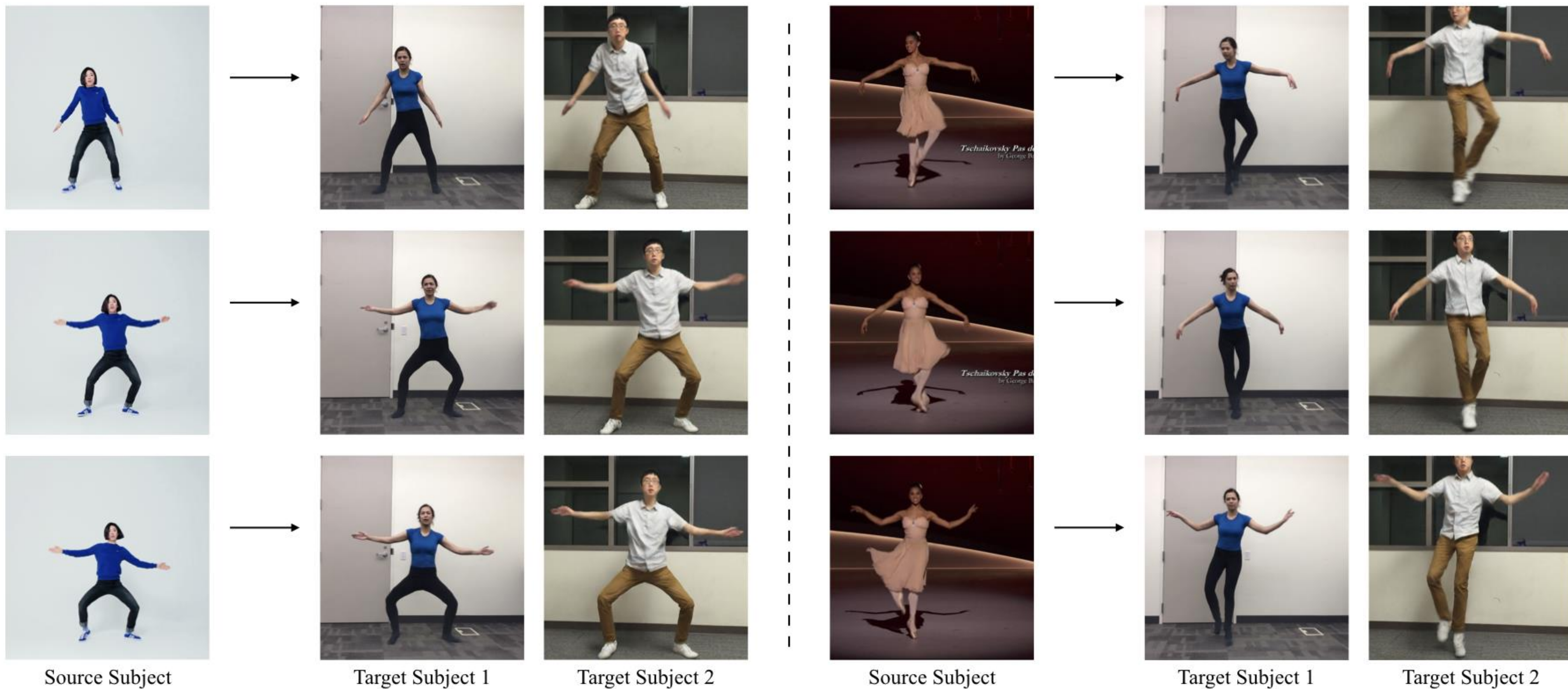


Figure 3: (Top) **Training:** Our model uses a pose detector P to create pose stick figures from video frames of the target subject. We learn the mapping G alongside an adversarial discriminator D which attempts to distinguish between the “real” correspondences $(x_t, x_{t+1}), (y_t, y_{t+1})$ and the “fake” sequence $(x_t, x_{t+1}), (G(x_t), G(x_{t+1}))$. (Bottom) **Transfer:** We use a pose detector P to obtain pose joints for the source person that are transformed by our normalization process $Norm$ into joints for the target person for which pose stick figures are created. Then we apply the trained mapping G .



https://carolineec.github.io/everybody_dance_now/

C. Chan, S. Ginosar, T. Zhou, A. Efros. [Everybody Dance Now](#). ICCV 2019

Source Video



Distorted
Copy



Source for Target 1: Source

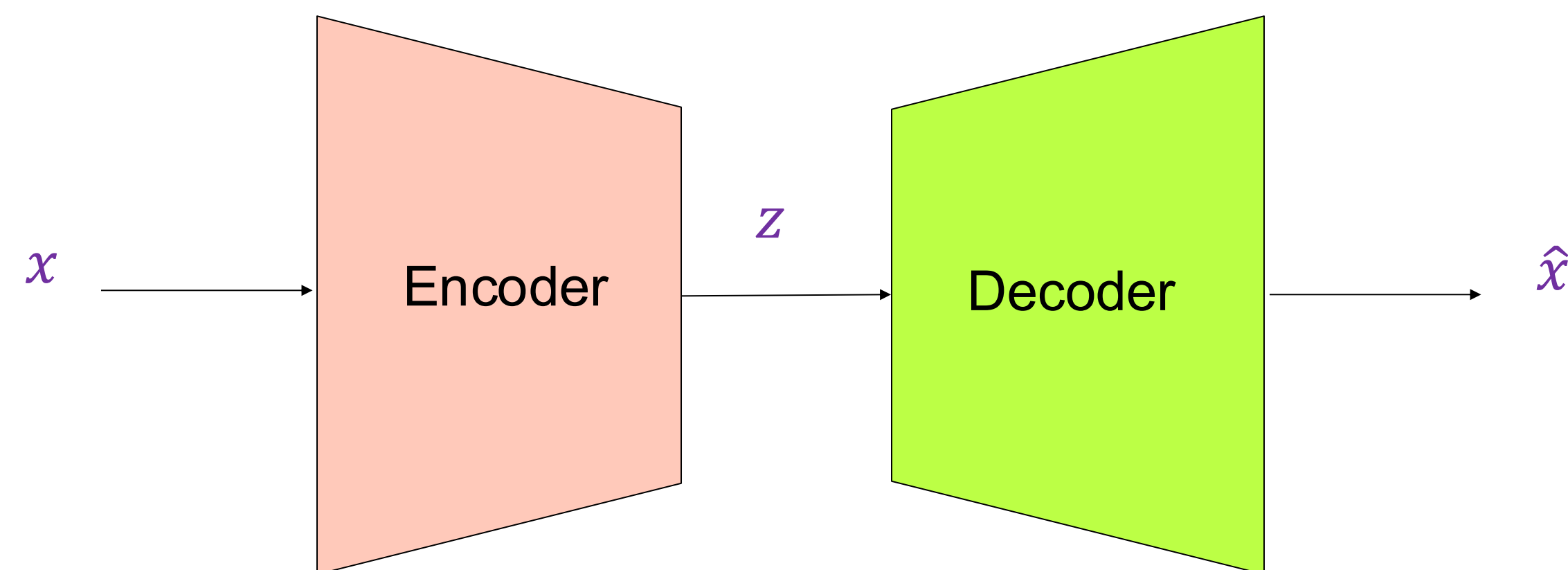


Source for Target 2: Source

Variational Autoencoder (VAE)

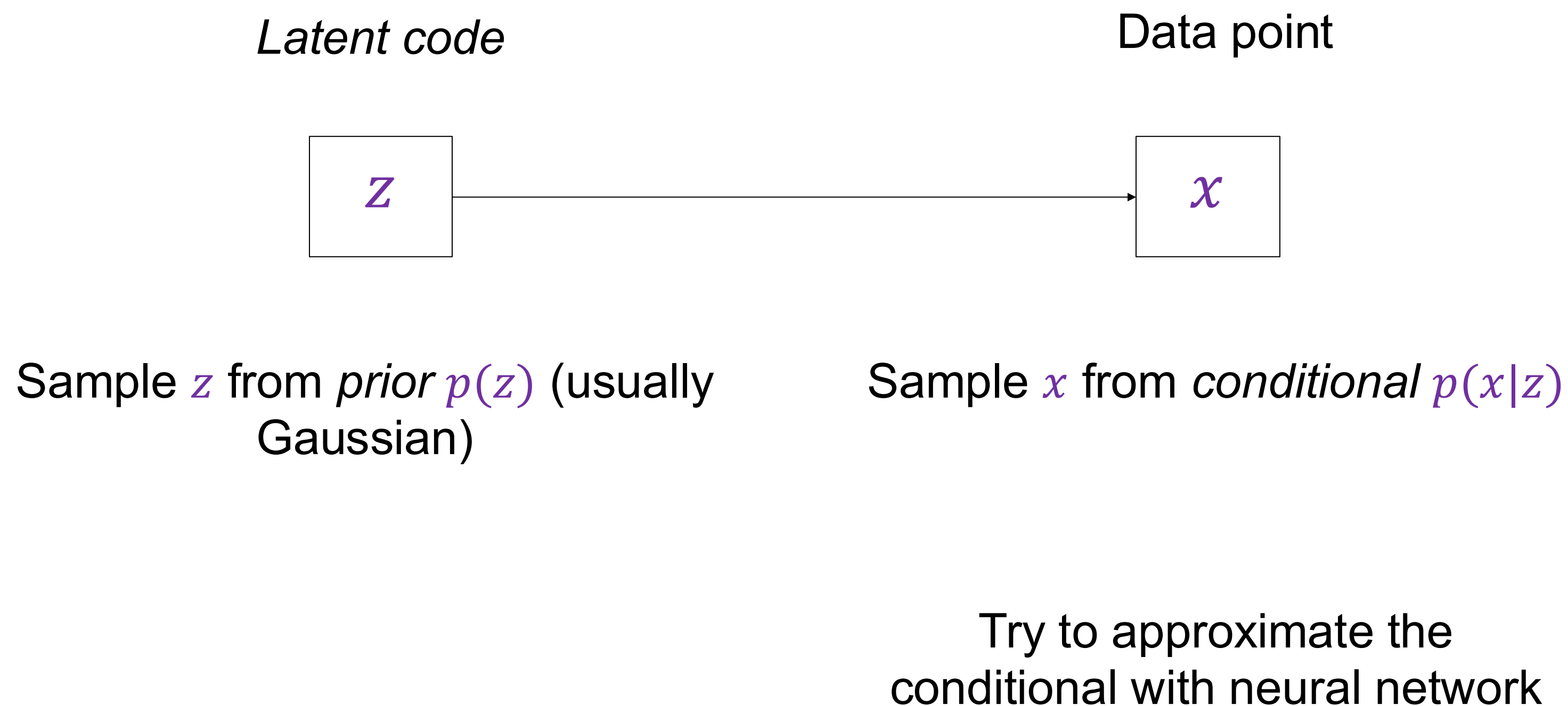
Variational autoencoders: Overview

- Probabilistic formulation based on *variational Bayes* framework
- At training time, jointly learn *encoder* and *decoder* by maximizing (a bound on) the data likelihood
- At test time, discard encoder and use decoder to sample from the learned distribution



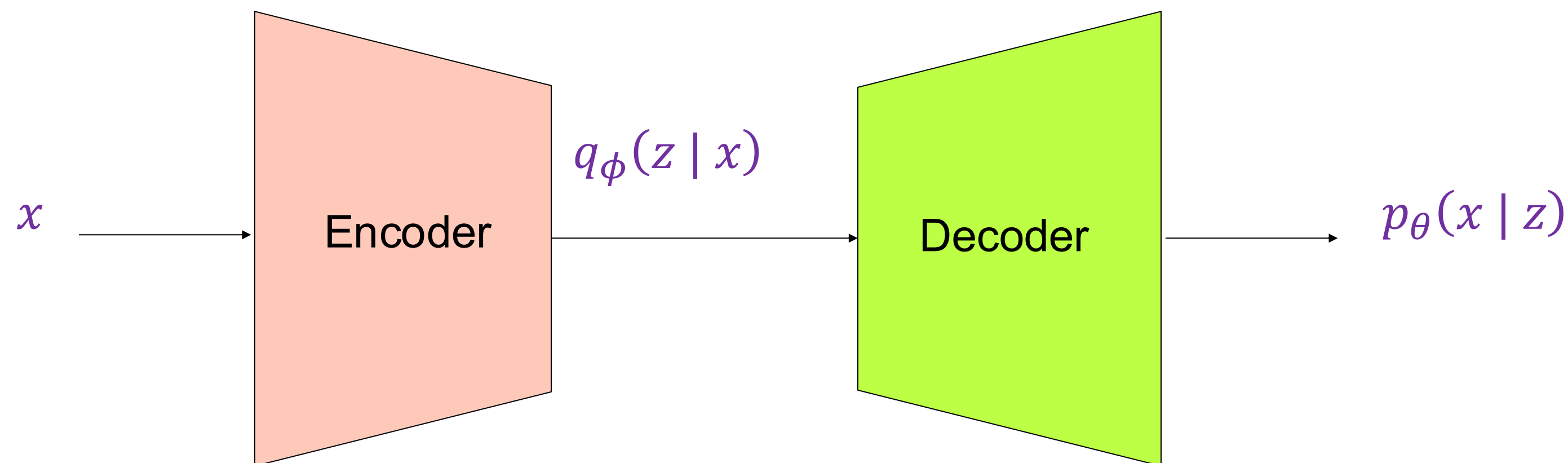
Variational autoencoders: Overview

- Probabilistic generative model of the data distribution:



Variational autoencoders: Training

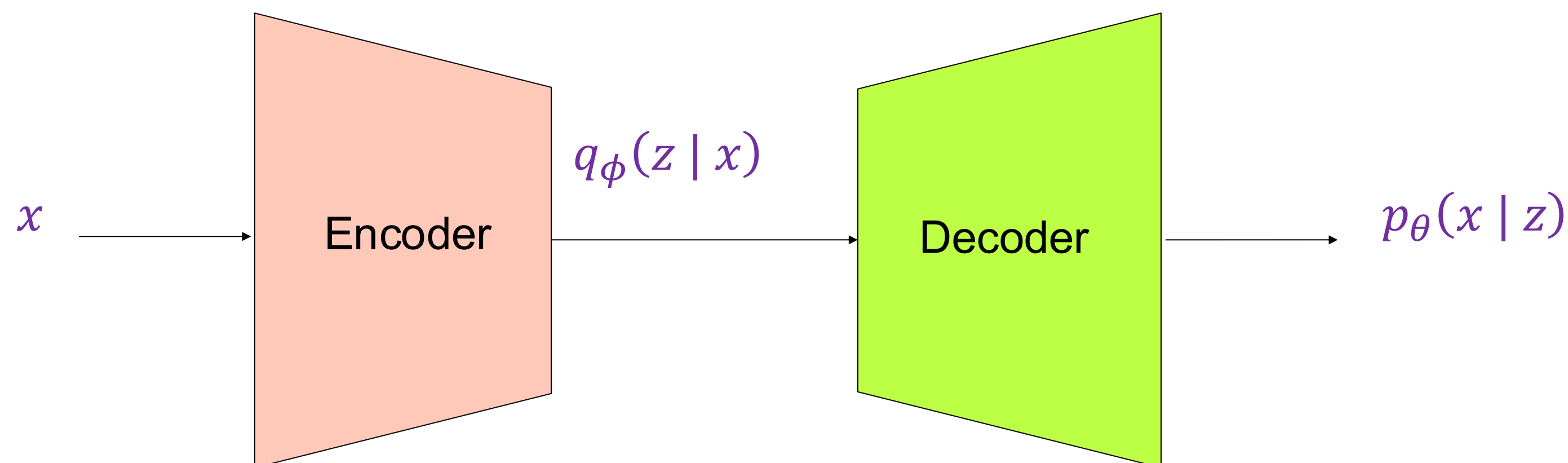
- **Encoder:** given inputs x , output $q_\phi(z | x)$
 - Specifically, output mean and (diagonal) covariance, or $\mu_{z|x}$ and $\Sigma_{z|x}$, so that $q_\phi(z | x) = N(\mu_{z|x}, \Sigma_{z|x})$
 - Approximate $q_\phi(z | x)$ to $N(0, I)$
- **Decoder:** given z , which is sampled from $q_\phi(z | x)$, then output $p_\theta(x | z)$



Variational autoencoders: Training

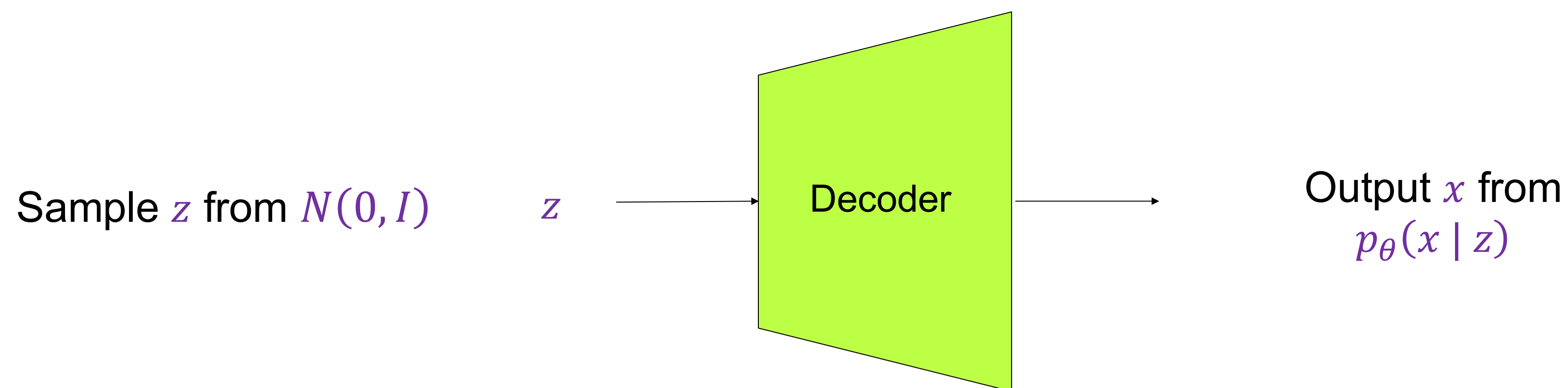
- Objective: maximize the *variational lower bound* on the data likelihood:

$$\log p_{\theta}(x) \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - D_{KL} \left(q_{\phi}(z|x) \parallel N(0, I) \right)$$

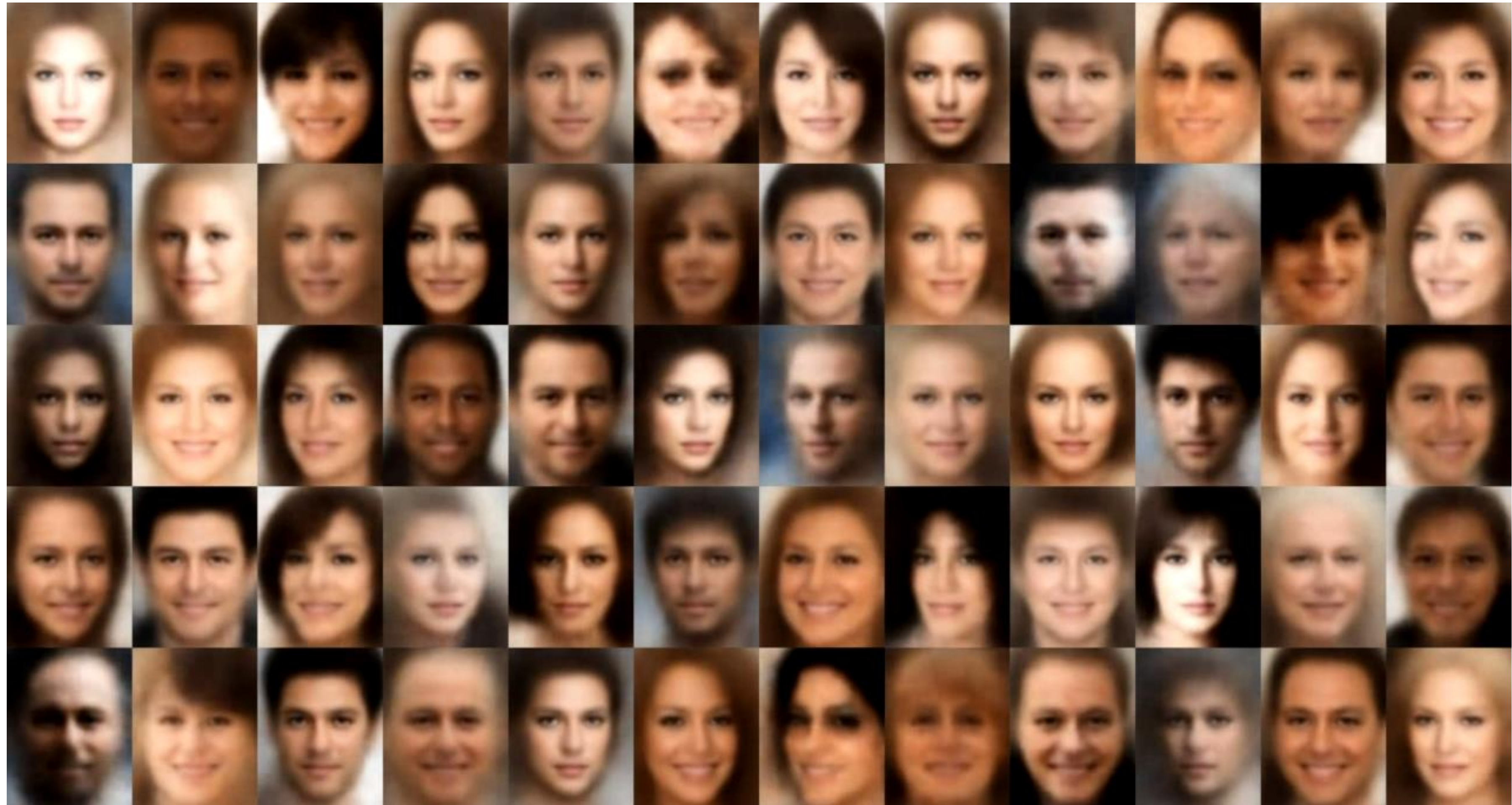


Variational autoencoders: Testing

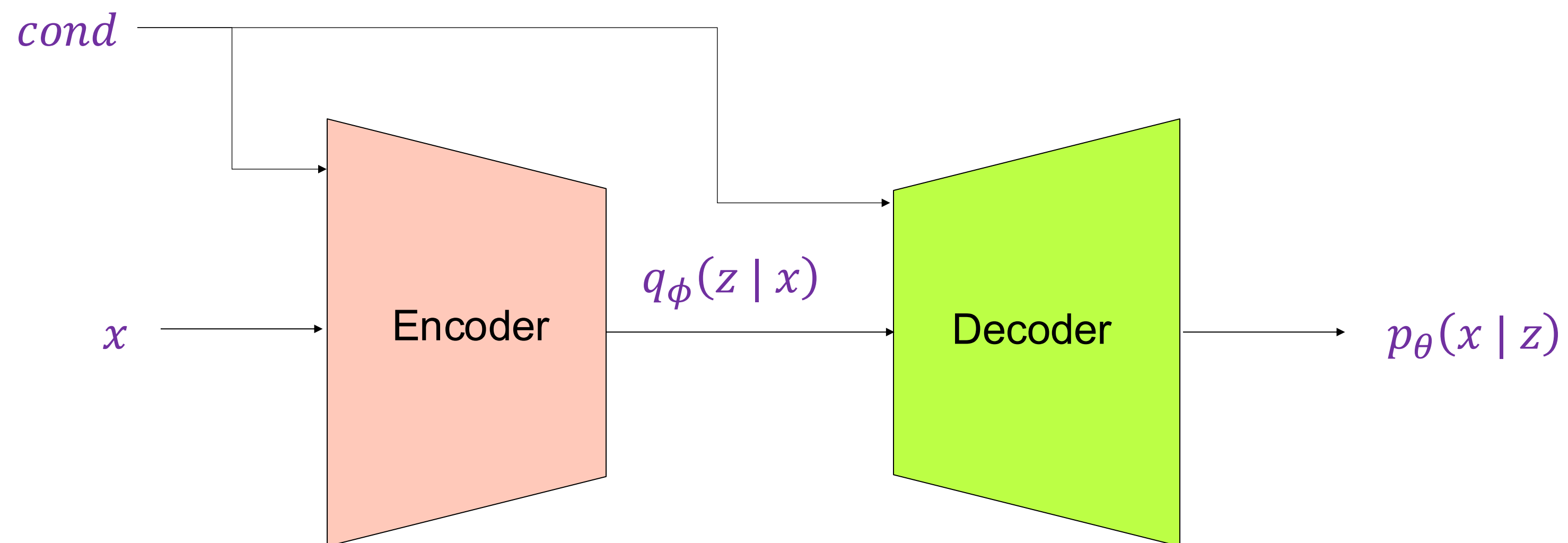
- At test time, discard encoder and use decoder to sample z from $N(0, I)$ and obtain output $p_{\theta}(x | z)$



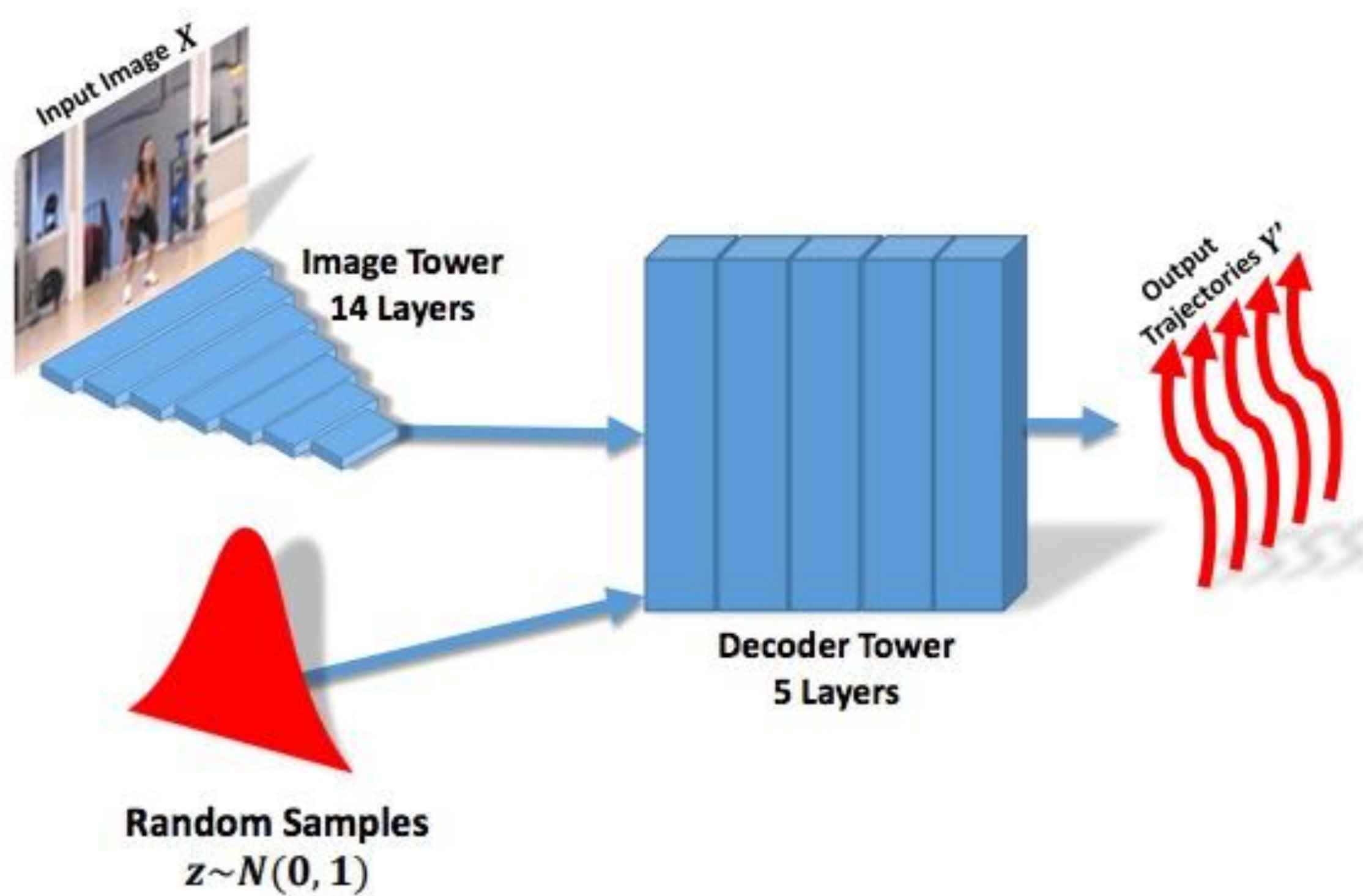
Variational autoencoders: Generating data



Conditional VAE

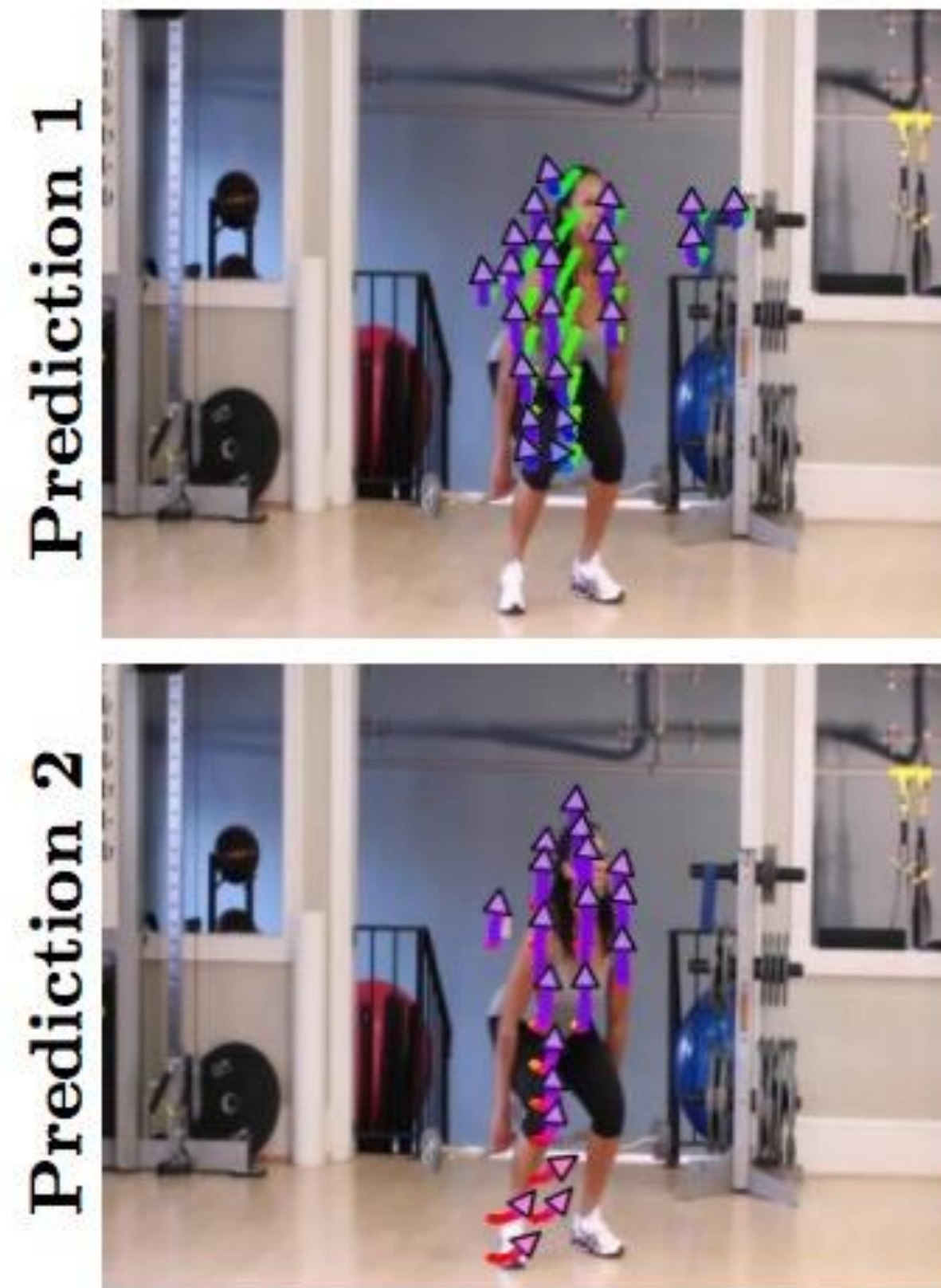


Conditional VAE for Video Prediction

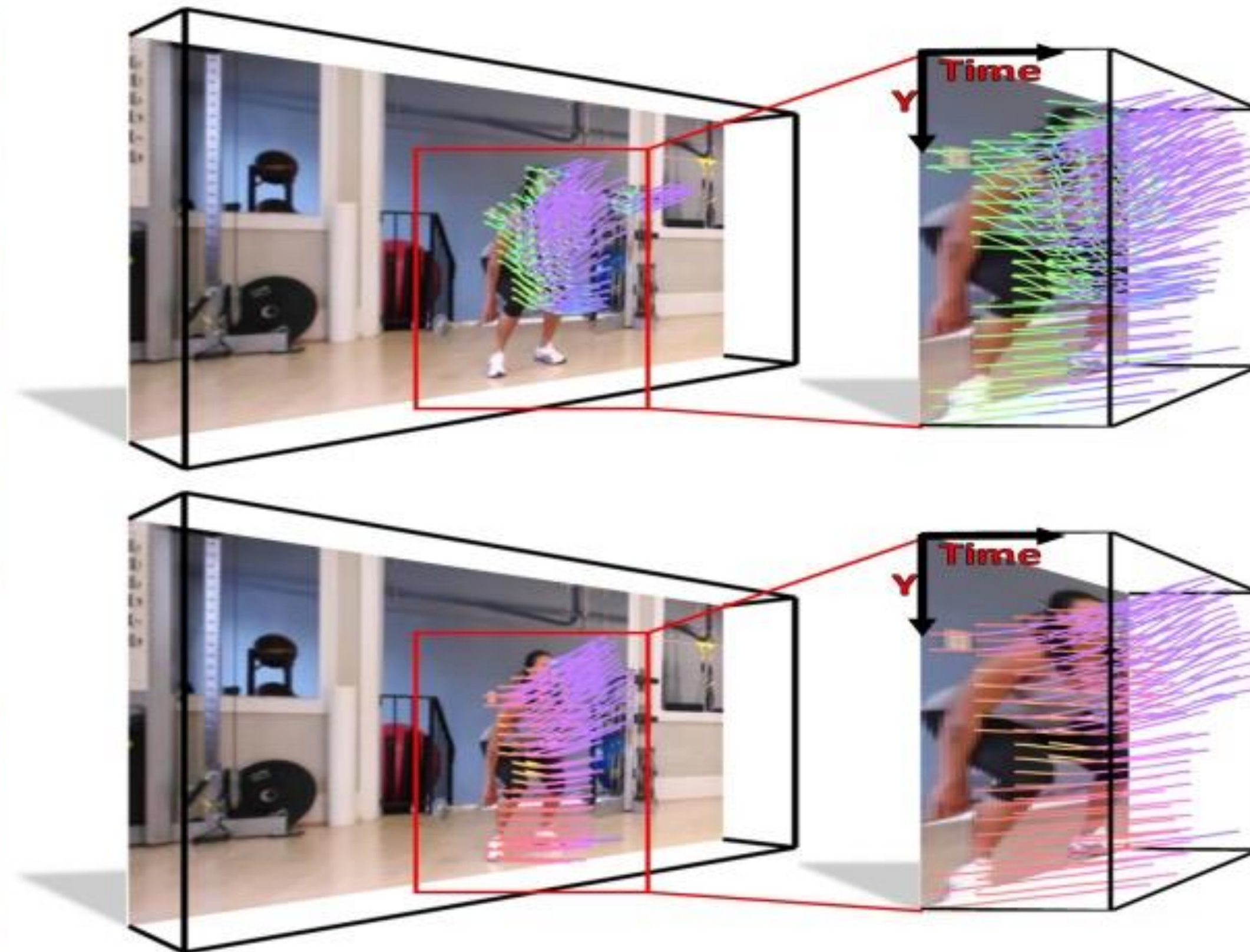


(a) Testing Architecture

Conditional VAE for Video Prediction



(a) Trajectories on Image



(b) Trajectories in Space-Time

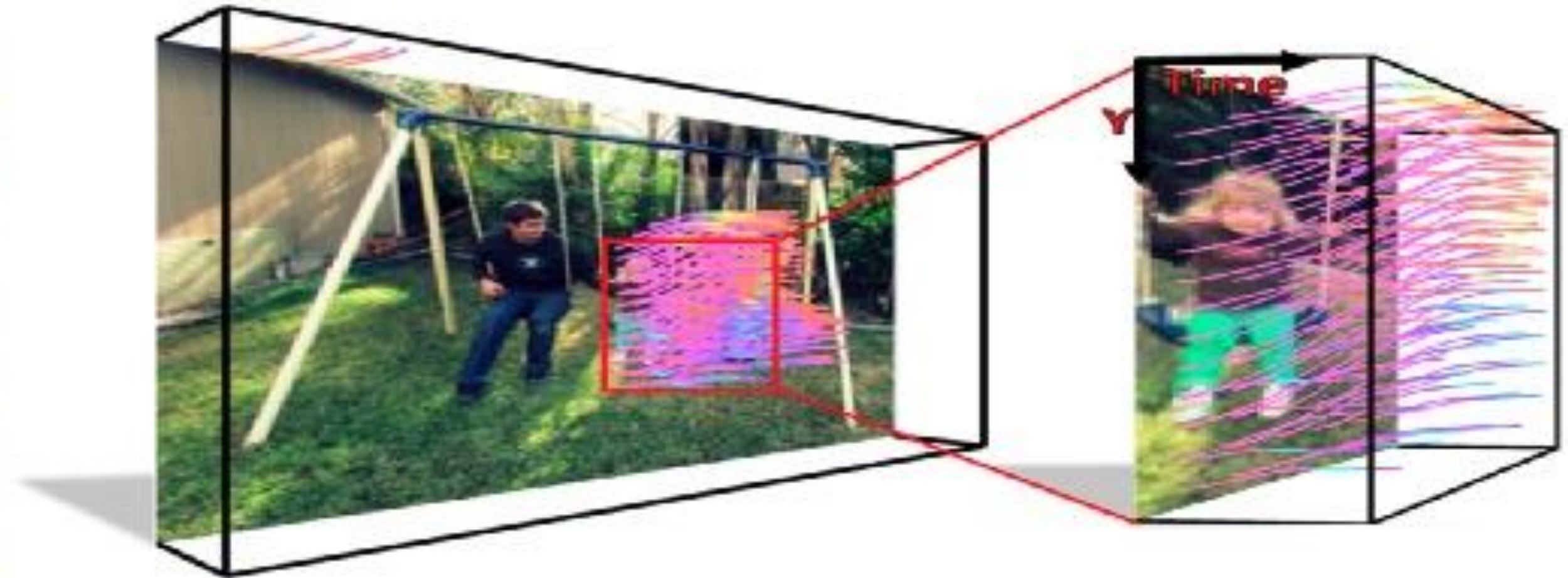
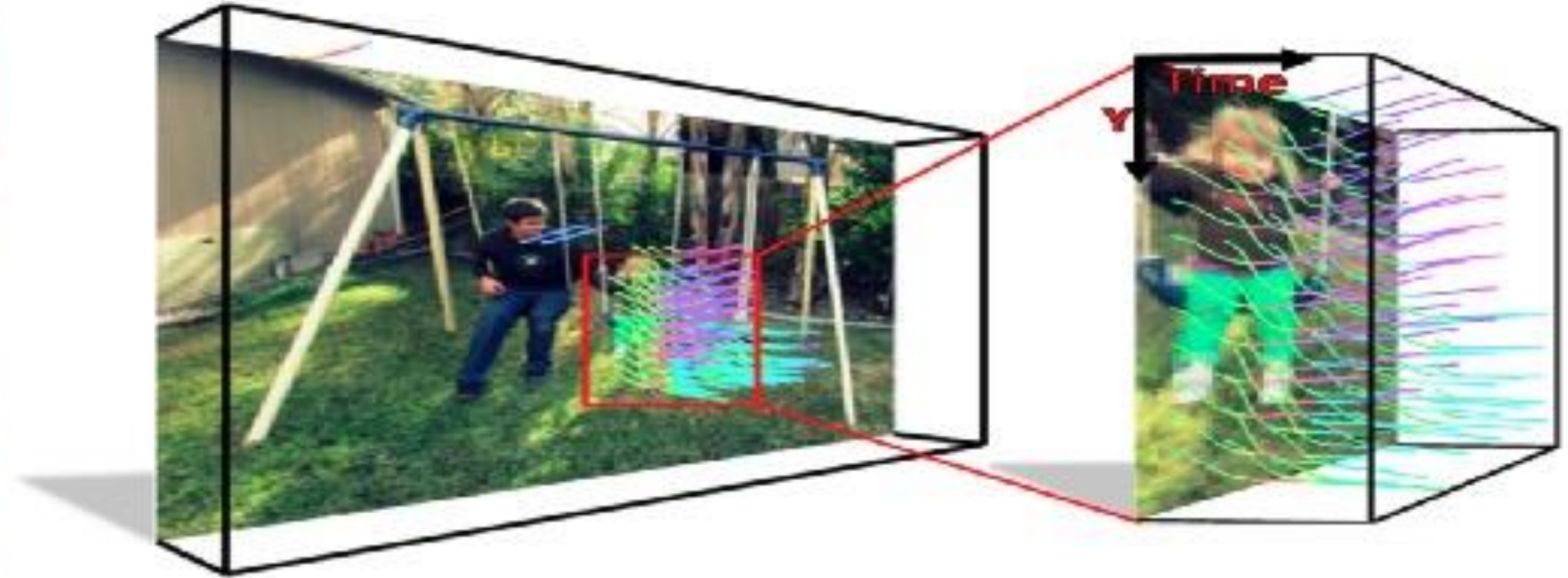


Conditional VAE for Video Prediction

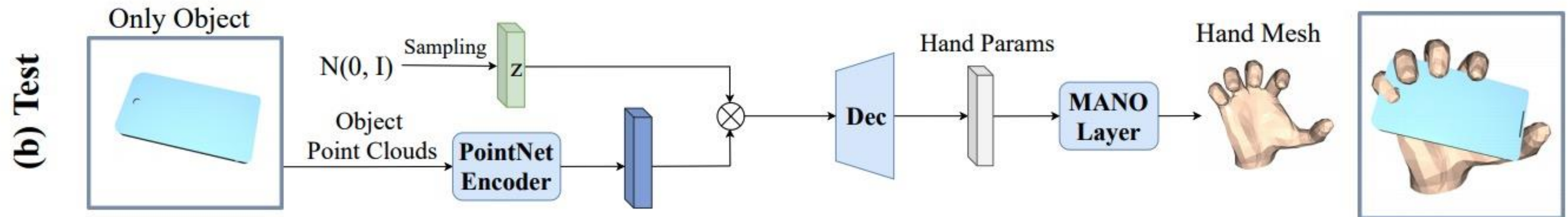
Prediction 1



Prediction 2

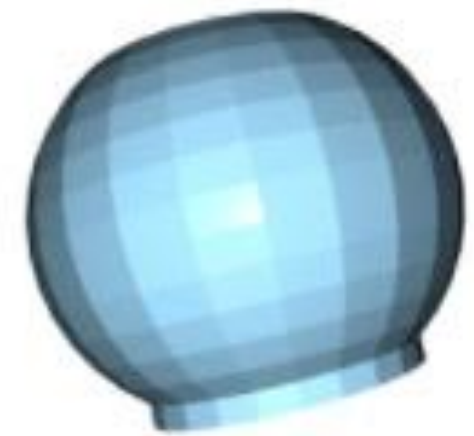


Conditional VAE for Grasp Generation



Conditional VAE for Grasp Generation

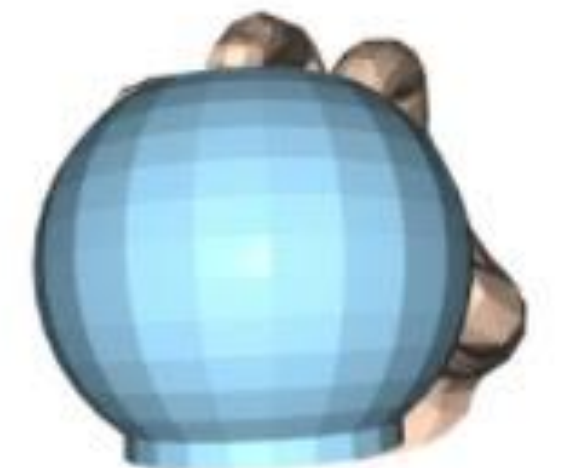
Input



Output



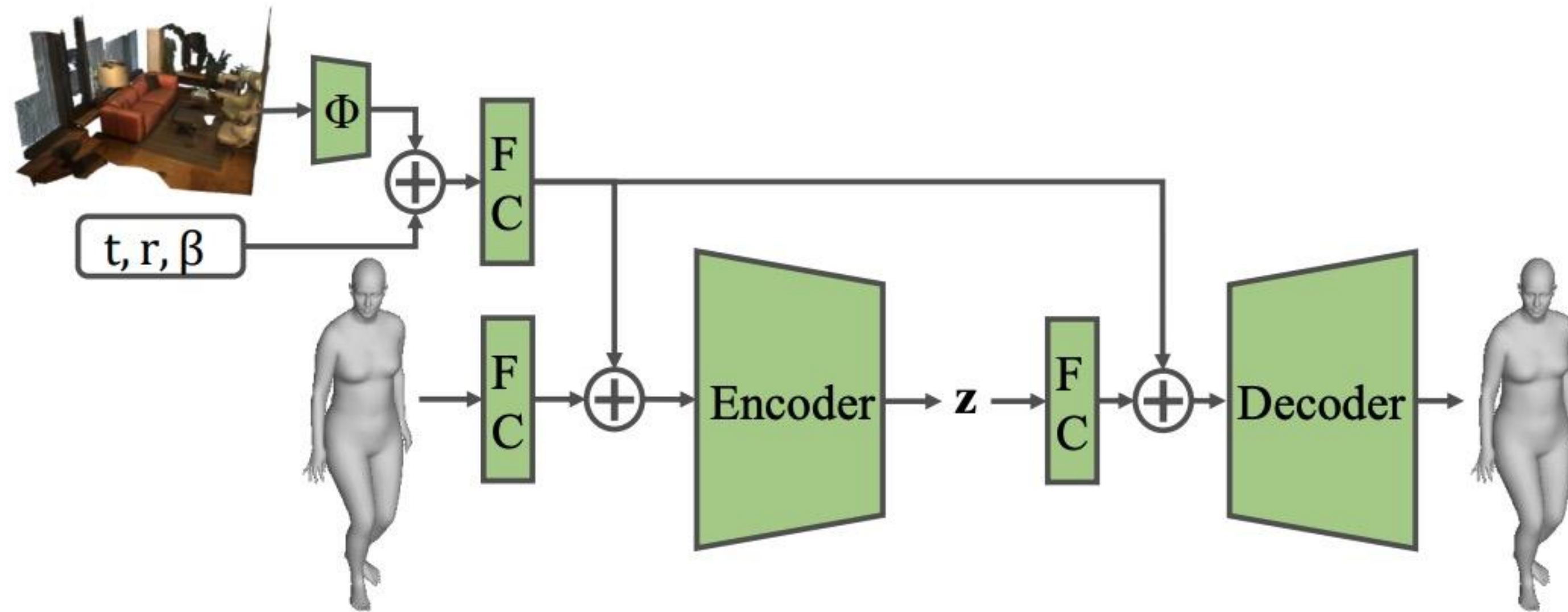
Input + Output (3 Views)



Conditional VAE for Grasp Generation



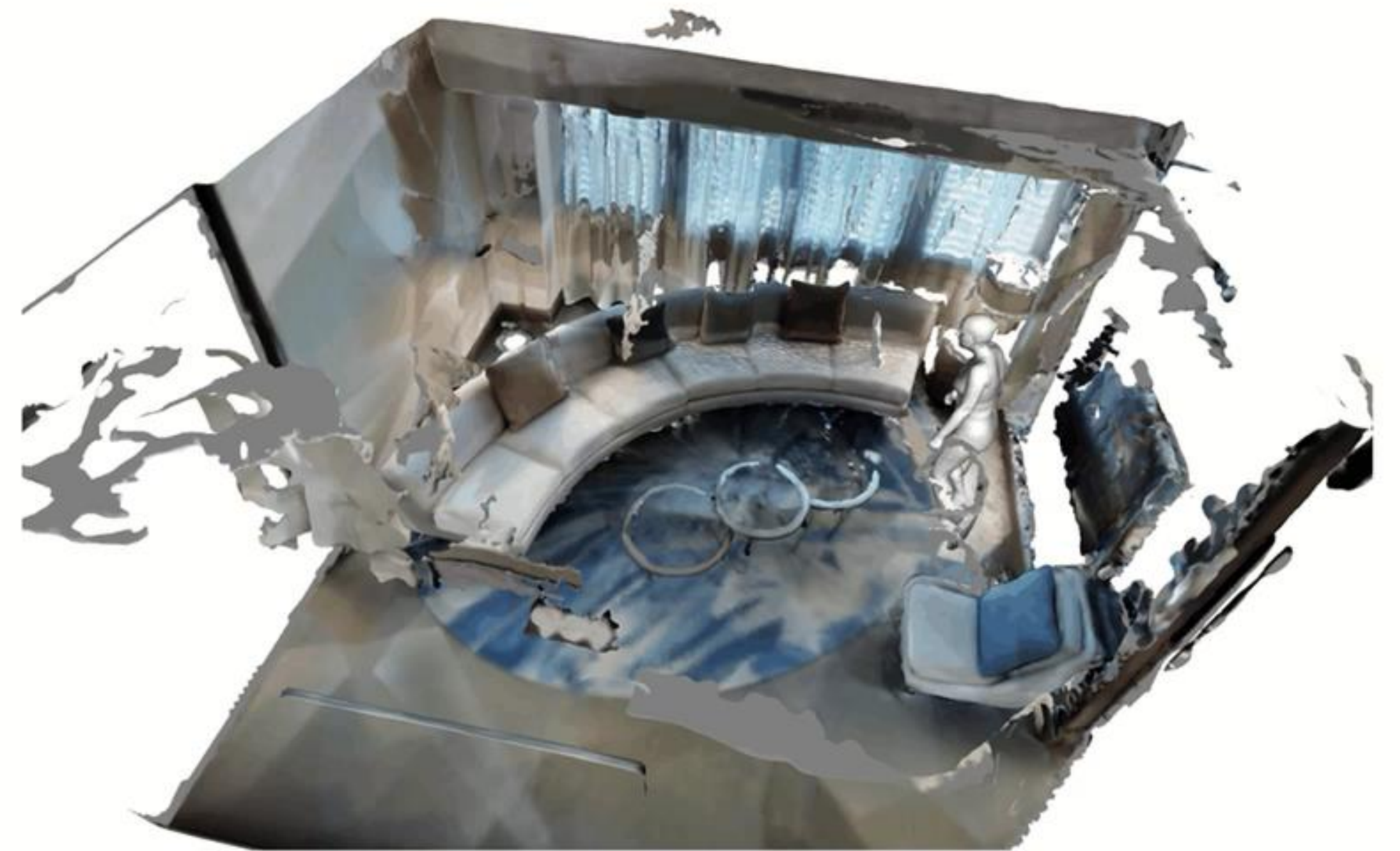
Conditional VAE for Human Motion Synthesis



Conditional VAE for Human Motion Synthesis



Conditional VAE for Human Motion Synthesis



Summary

- Image-to-Image Translation: pix2pix
- Unpaired Image-to-Image Translation: CycleGAN
- Variational Autoencoder (VAE)