

# CMP784

## DEEP LEARNING

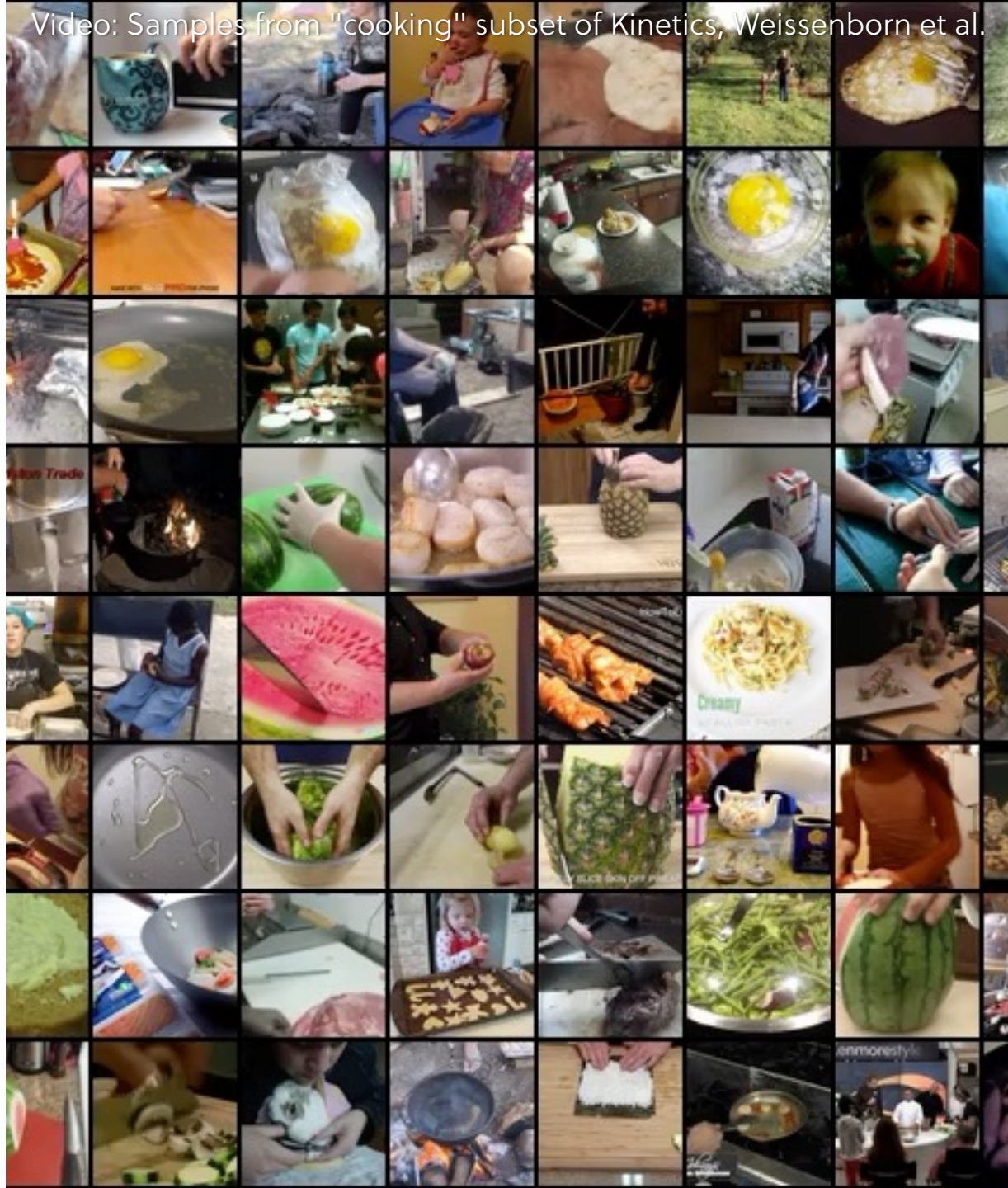
### Lecture #10 – Deep Generative Models – Part 2



HACETTEPE  
UNIVERSITY  
COMPUTER  
VISION LAB

# Previously on CMP784

- Supervised vs. Unsupervised Learning
- Generative Modeling
- Basic Foundations
  - Sparse Coding
  - Autoencoders
- Autoregressive Generative Models



# Lecture overview

- Generative Adversarial Networks (GANs)

**Disclaimer:** Some of the material and slides for this lecture were borrowed from

—Ian Goodfellow’s tutorial on “Generative Adversarial Networks”

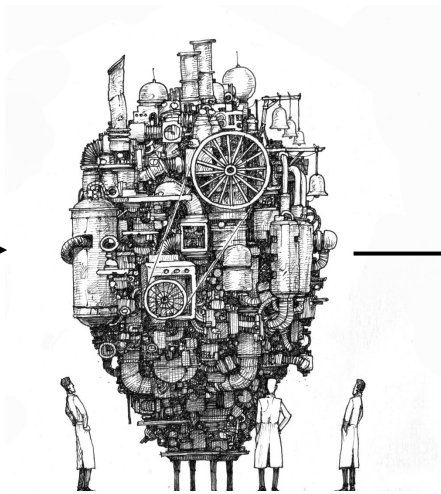
—Aaron Courville’s IFT6135 class

—Bill Freeman, Antonio Torralba and Phillip Isola’s MIT 6.869 class

# Discriminative vs. Generative Models

$$p(y|x)$$

$$p(x|y)$$



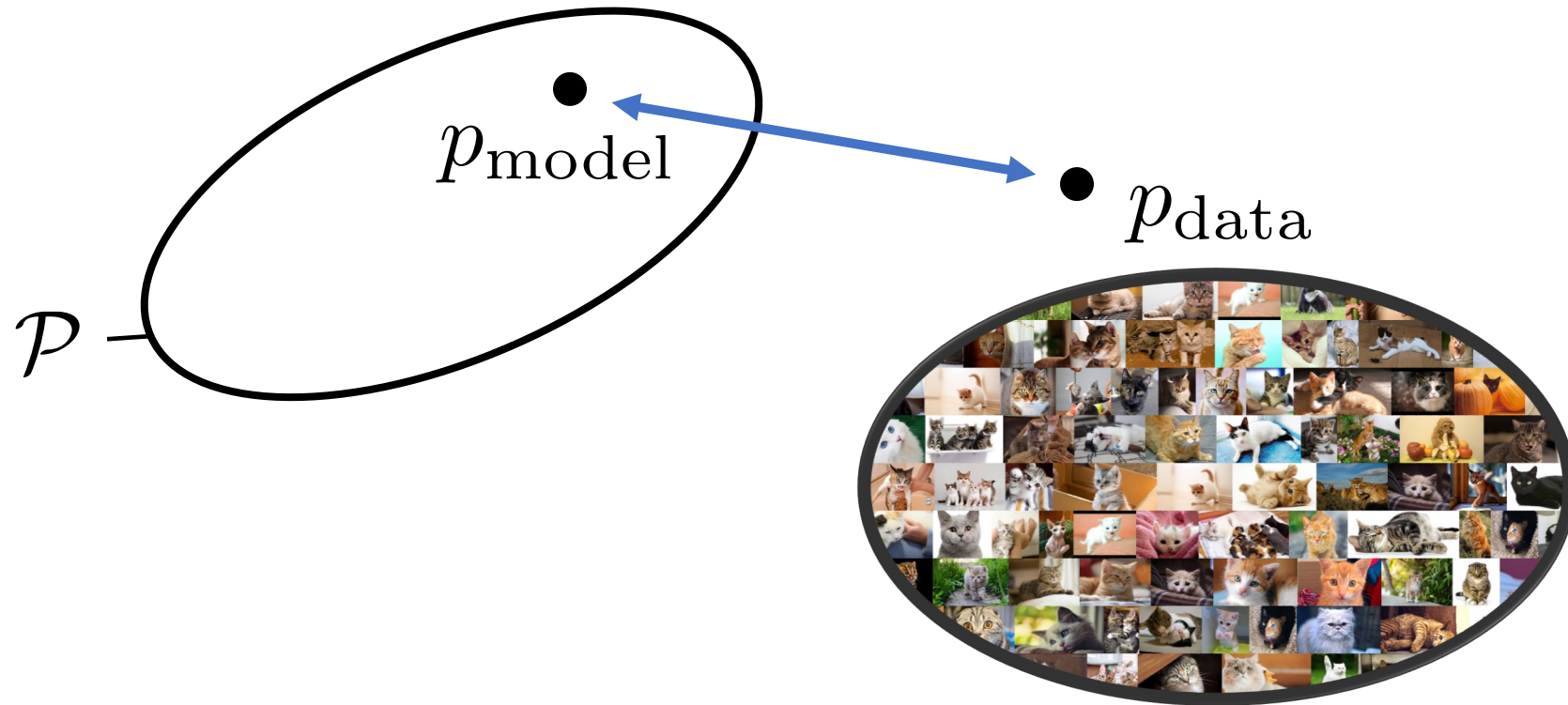
"Cat"



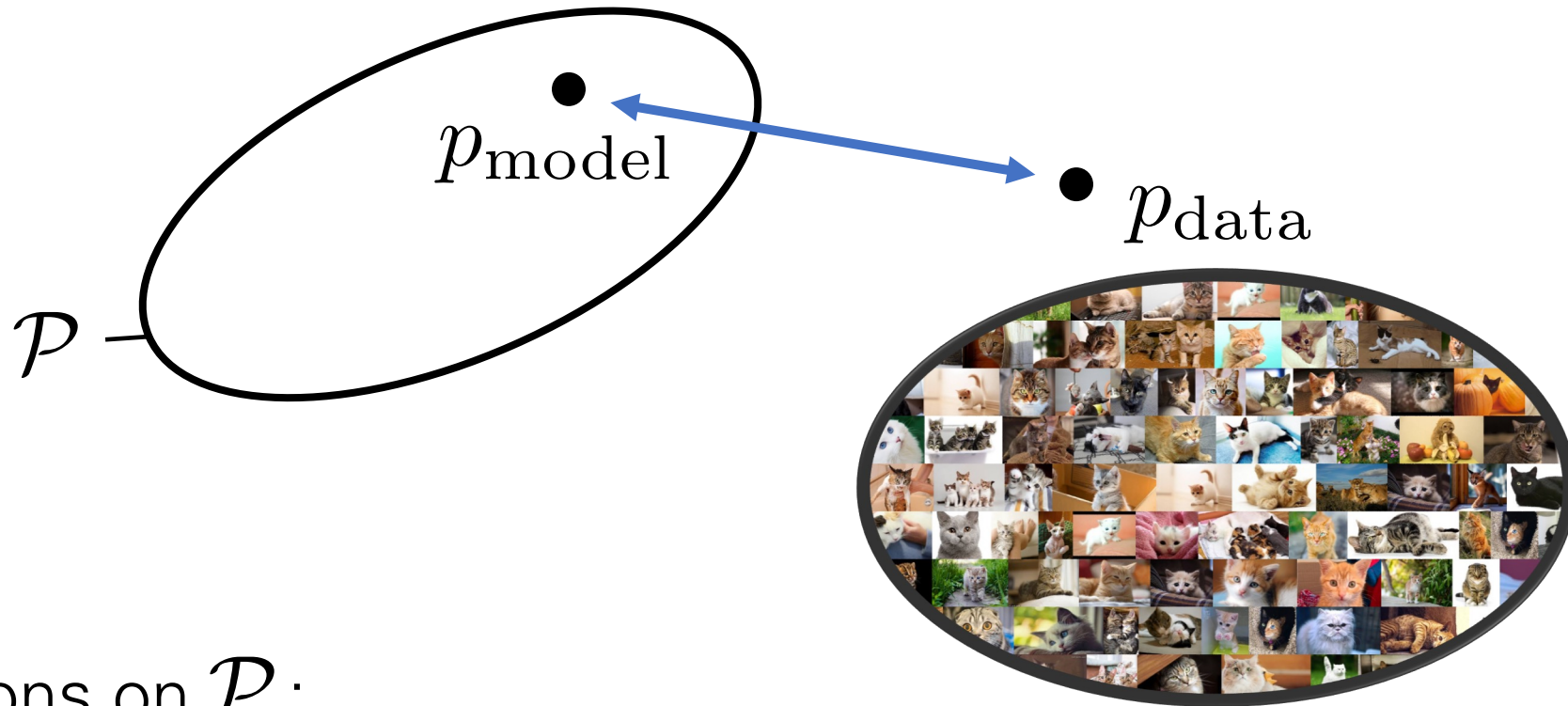
Discriminative models

Generative models

# Generative Modeling

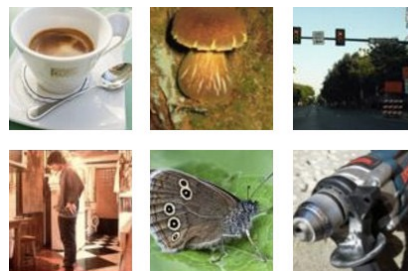


# Generative Modeling

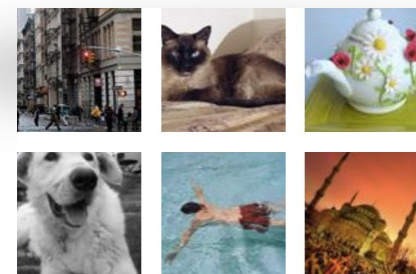


Assumptions on  $\mathcal{P}$ :

- tractable sampling

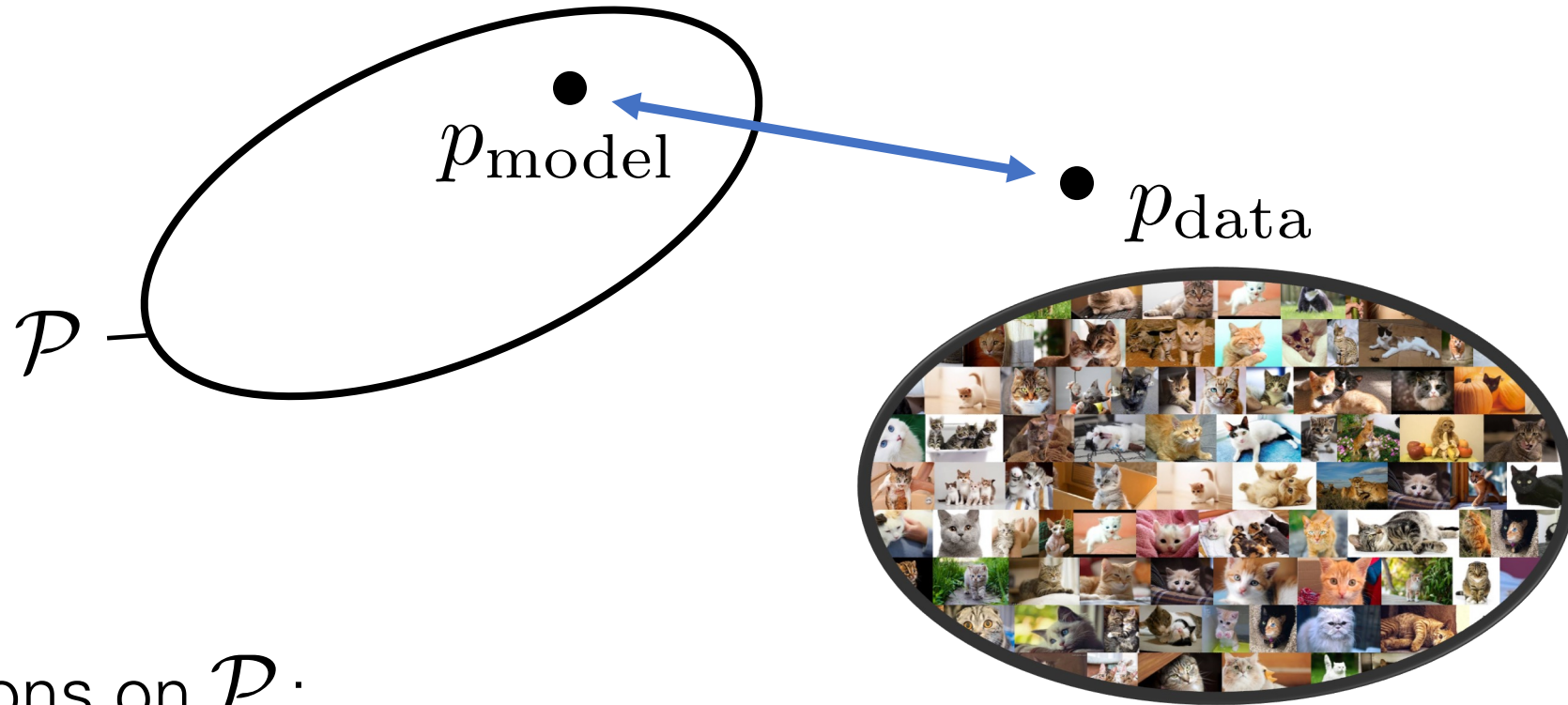


Training examples



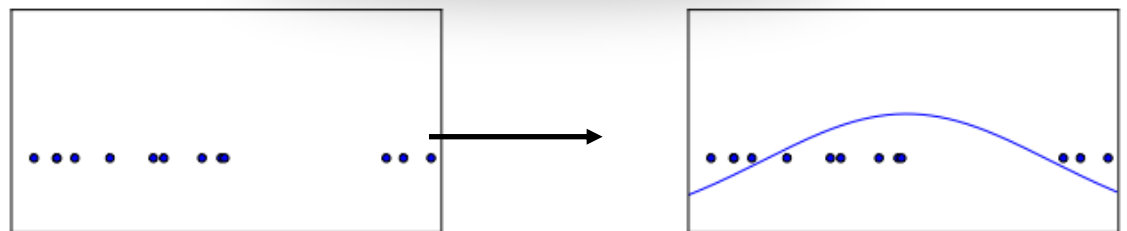
Model samples

# Generative Modeling



Assumptions on  $\mathcal{P}$ :

- tractable sampling
- tractable likelihood function



# Broad Categories of Generative Models

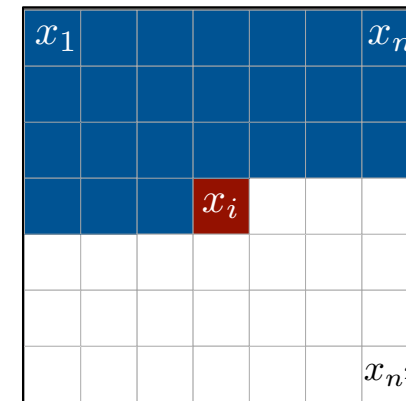
- Autoregressive Models
- Generative Adversarial Networks (GANs)
- Flow-based Models
- Variational Autoencoders
- Energy-based Models



# Autoregressive Models

- Explicitly model conditional probabilities:

$$p_{\text{model}}(\mathbf{x}) = p_{\text{model}}(x_1) \prod_{i=2}^n p_{\text{model}}(x_i \mid x_1, \dots, x_{i-1})$$



Each conditional can be a complicated neural net

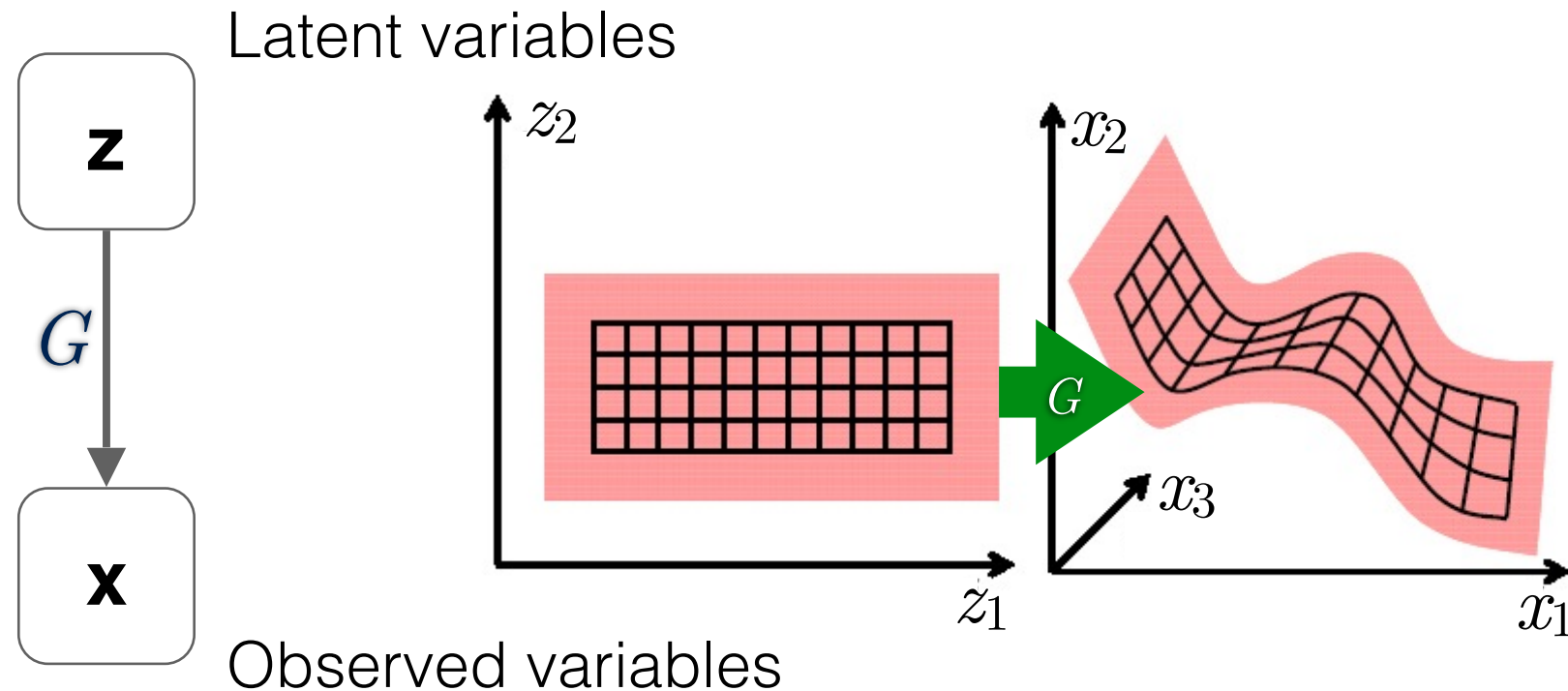
## Disadvantages:

- Generation can be too costly
- Generation can not be controlled by a latent code

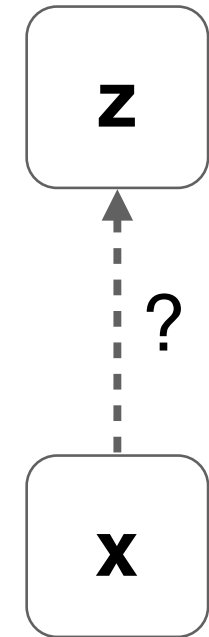


PixelCNN elephants  
(van den Ord et al. 2016)

# Another way to train a latent variable model



inference



# Generative Adversarial Networks

# Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)



Noise  
(random input)

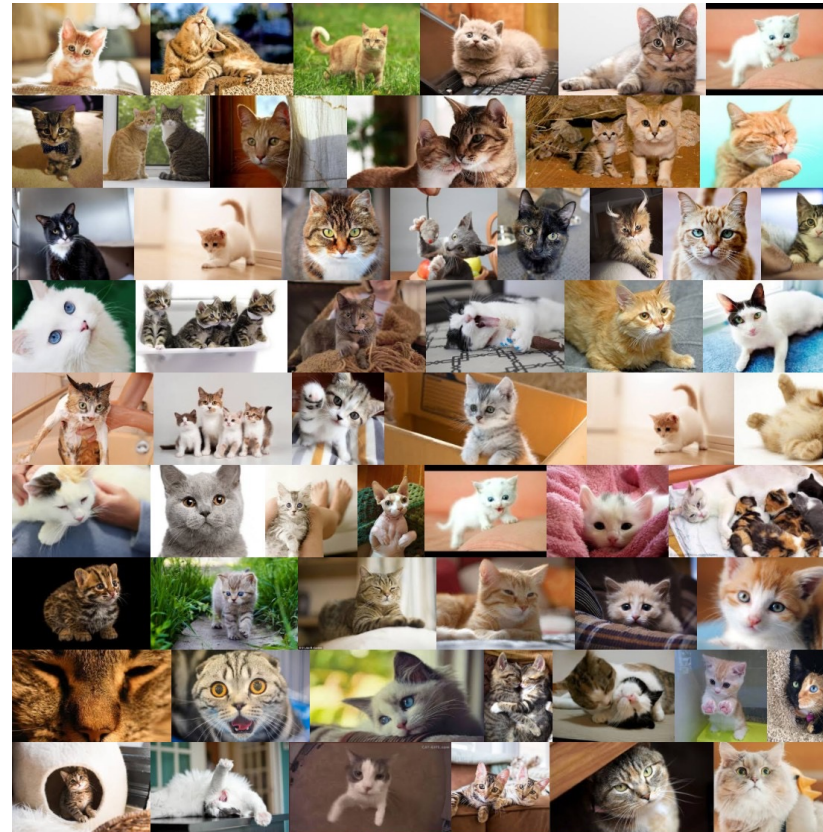


Generative  
Model



$z \sim \text{Uniform}_{100}$

think of this as  
a transformation



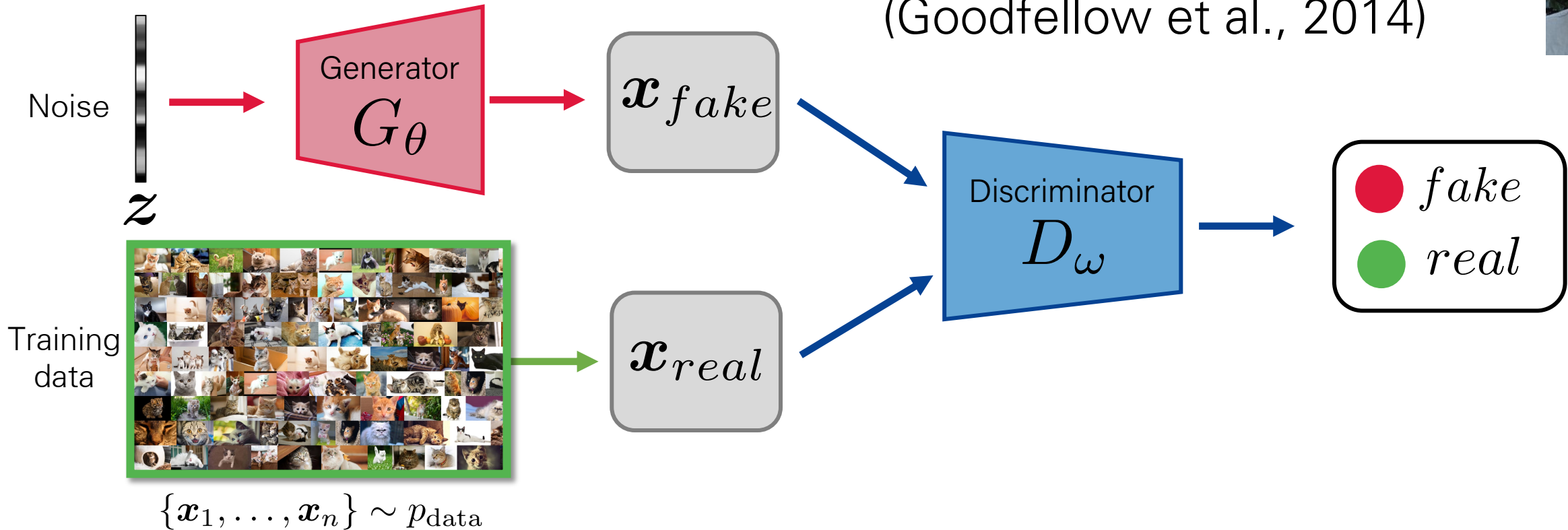
- A game-theoretic likelihood free model

## Advantages:

- Uses a latent code
- No Markov chains needed
- Produces the best looking samples

# Generative Adversarial Networks (GANs)

(Goodfellow et al., 2014)



- A game between a generator  $G_\theta(z)$  and a discriminator  $D_\omega(x)$ 
  - Generator tries to fool discriminator (i.e. generate realistic samples)
  - Discriminator tries to distinguish fake from real samples

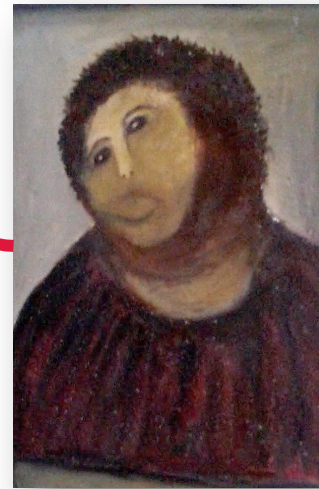
# Intuition behind GANs



$D_\omega$ : Discriminator (Art Critic)



$x_{real}$



$x_{fake}$

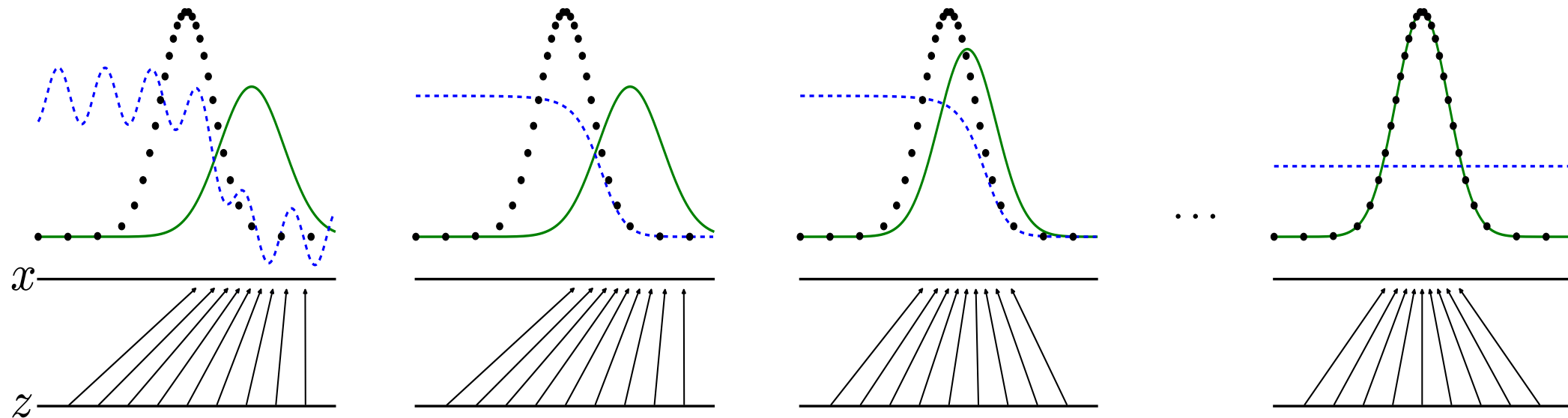


$G_\theta$ : Generator (Forger)

# Training Procedure

(Goodfellow et al., 2014)

- Use SGD on two minibatches simultaneously:
  - A minibatch of training examples
  - A minibatch of generated samples



# GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$

Real data

Noise vector used  
to generate data

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} \log D(\mathbf{x}) - \frac{1}{2} \mathbb{E}_{\mathbf{z}} \log (1 - D(G(\mathbf{z})))$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z}} \log D(G(\mathbf{z}))$$

Cross-entropy  
loss for binary  
classification

Generator maximizes the log-probability  
of the discriminator being mistaken

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence between  $p_{\text{data}}$  and  $p_{\mathbf{x}}$



# GAN Training: Minimax Game (Goodfellow et al., 2014)

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$

Real data

Noise vector used

$$J^{(D)} = -\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log D_{\omega}(\mathbf{x})]$$

$$J^{(G)} = -\frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [\log (1 - D_{\omega}(G_{\theta}(\mathbf{z})))]$$

Important question is  
"Does this converge??"

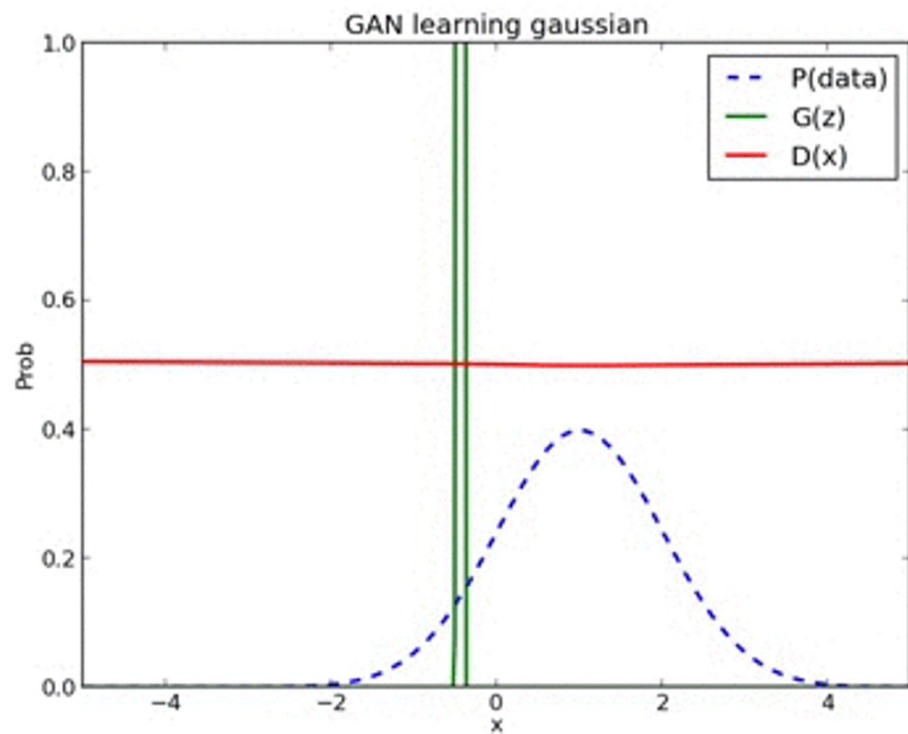
Cross-entropy loss for binary classification

log-probability of the discriminator being mistaken

- Equilibrium of the game
- Minimizes the Jensen-Shannon divergence

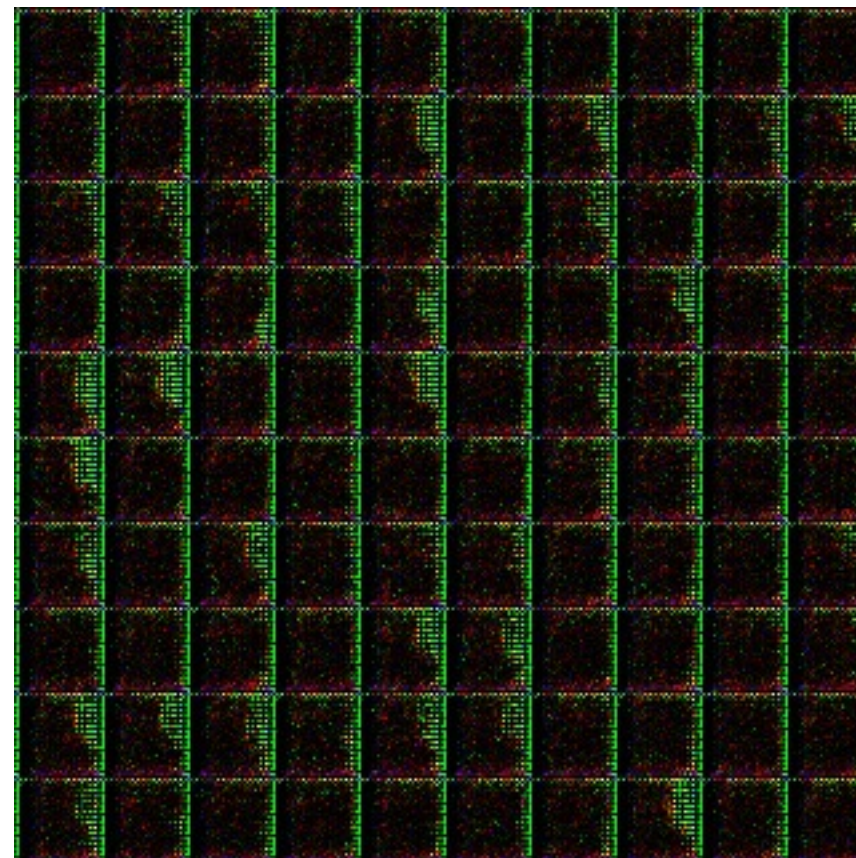
# Training Procedure

(Goodfellow et al., 2014)



Source: Alec Radford

Generating 1D points



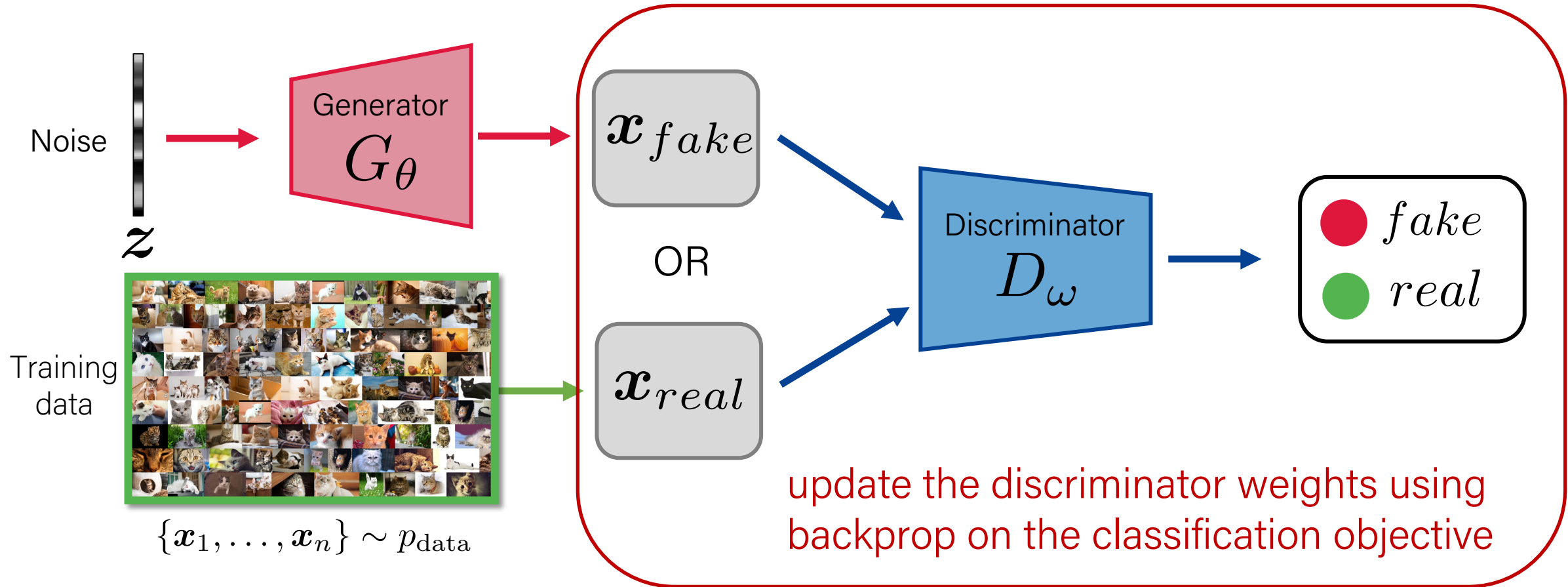
Source: OpenAI blog

Generating images



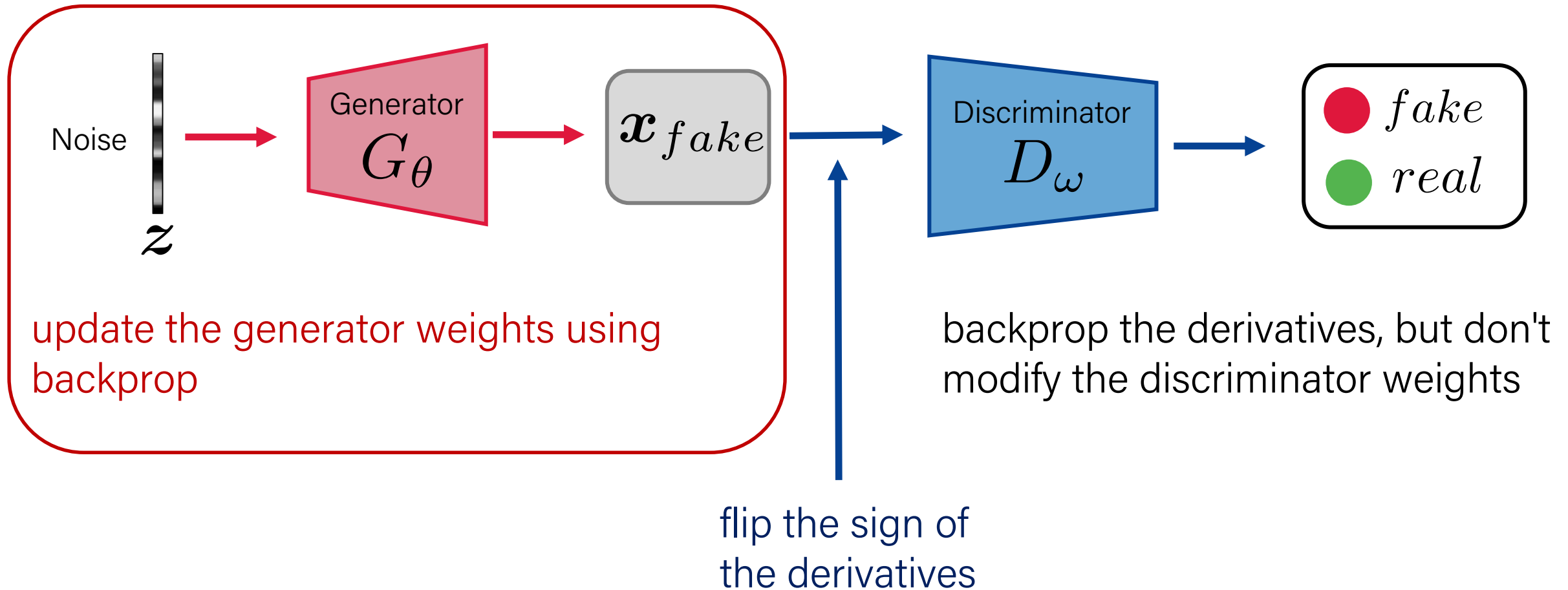
# Training Procedure

- Updating the discriminator:



# Training Procedure

- Updating the generator:



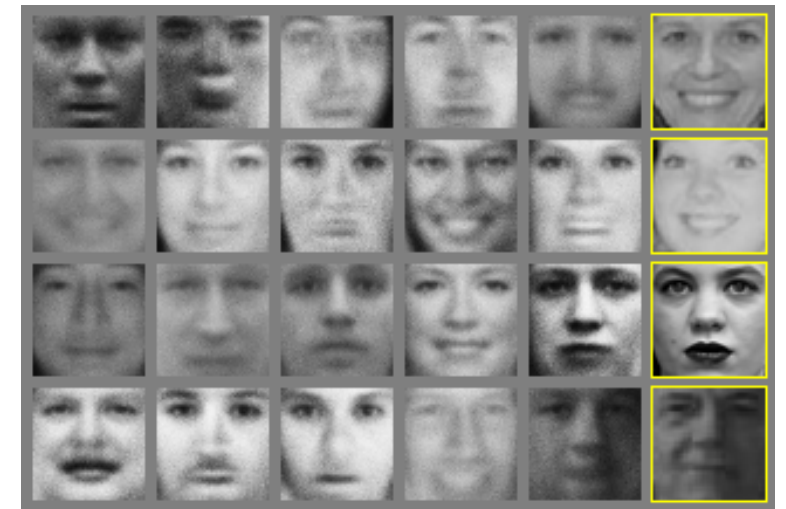
# Results

(Goodfellow et al., 2014)

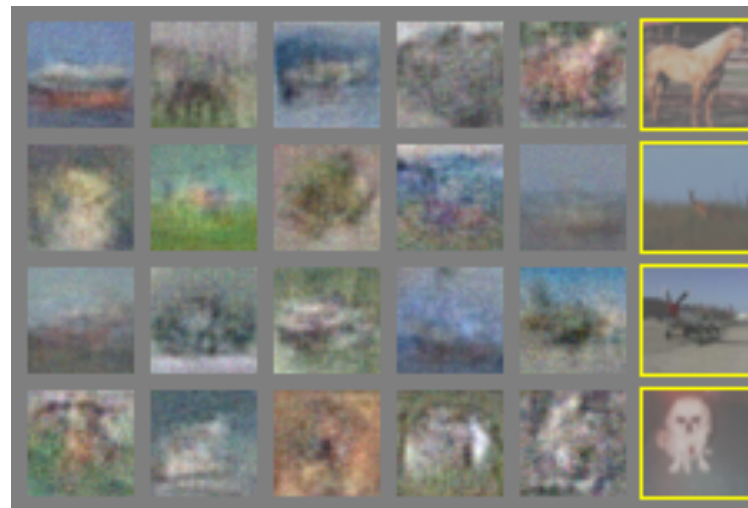
- The generator uses a mixture of rectifier linear activations and/or sigmoid activations
- The discriminator net used maxout activations.



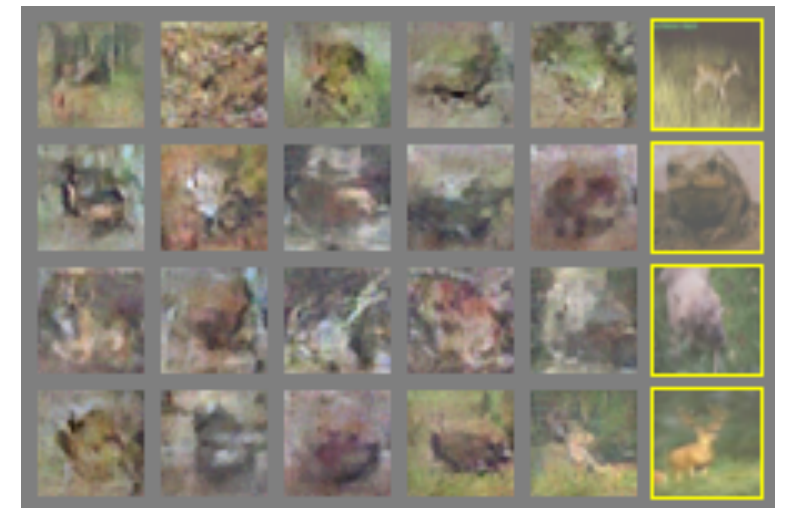
MNIST samples



TFD samples



CIFAR10 samples  
(fully-connected model)



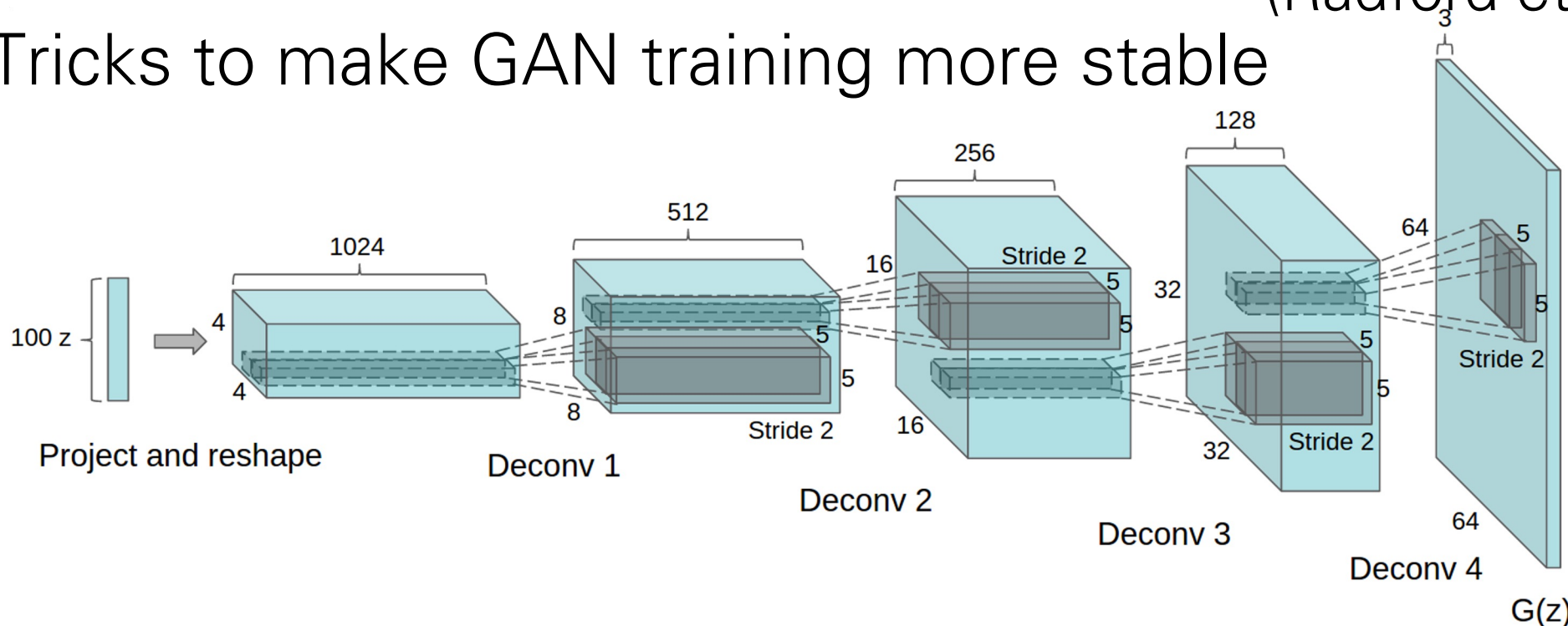
CIFAR10 samples  
(convolutional discriminator,  
deconvolutional generator)

# Deep Convolutional GANs (DCGAN)



(Radford et al., 2015)

- Idea: Tricks to make GAN training more stable

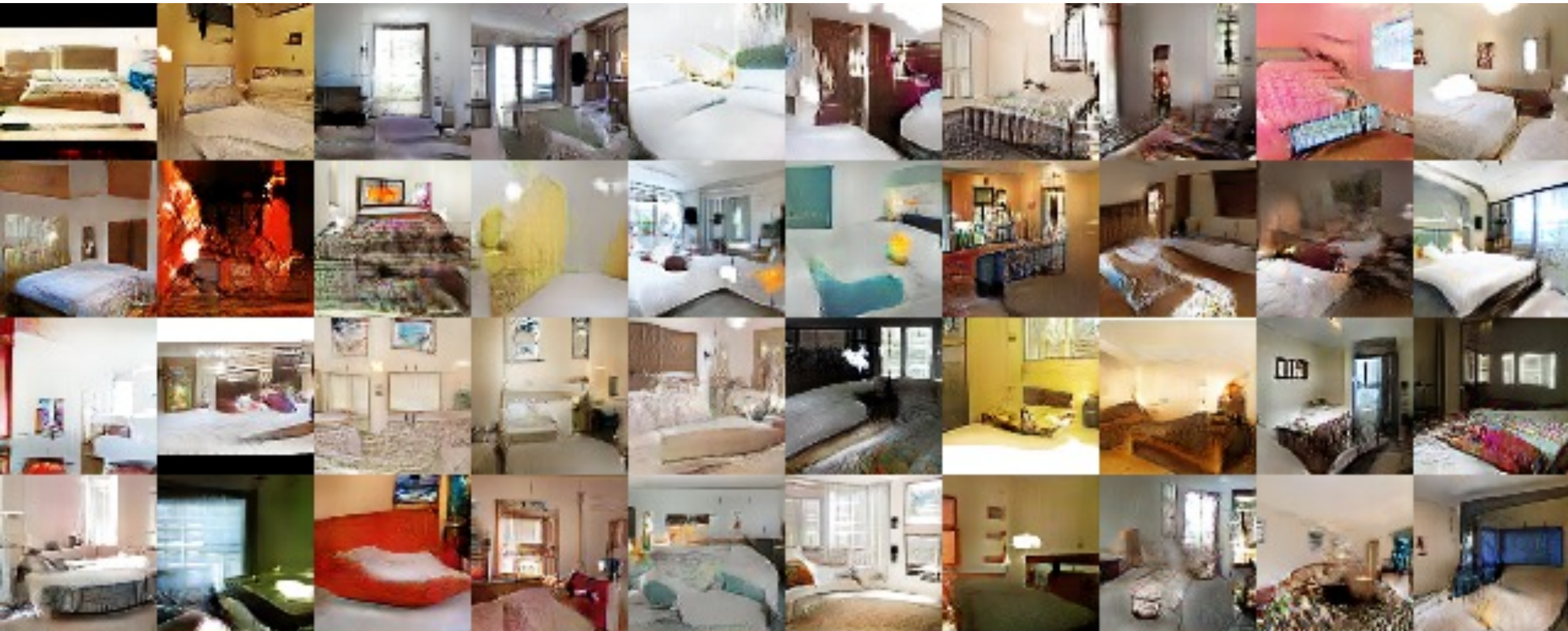


- No fully connected layers
- Batch Normalization (Ioffe and Szegedy, 2015)
- Leaky Rectifier in D
- Use Adam (Kingma and Ba, 2015)
- Tweak Adam hyperparameters a bit ( $\beta_1=0.0002$ ,  $\beta_2=0.5$ )

# DCGAN for LSUN Bedrooms

64×64 pixels  
~3M images

(Radford et al.,  
2015)

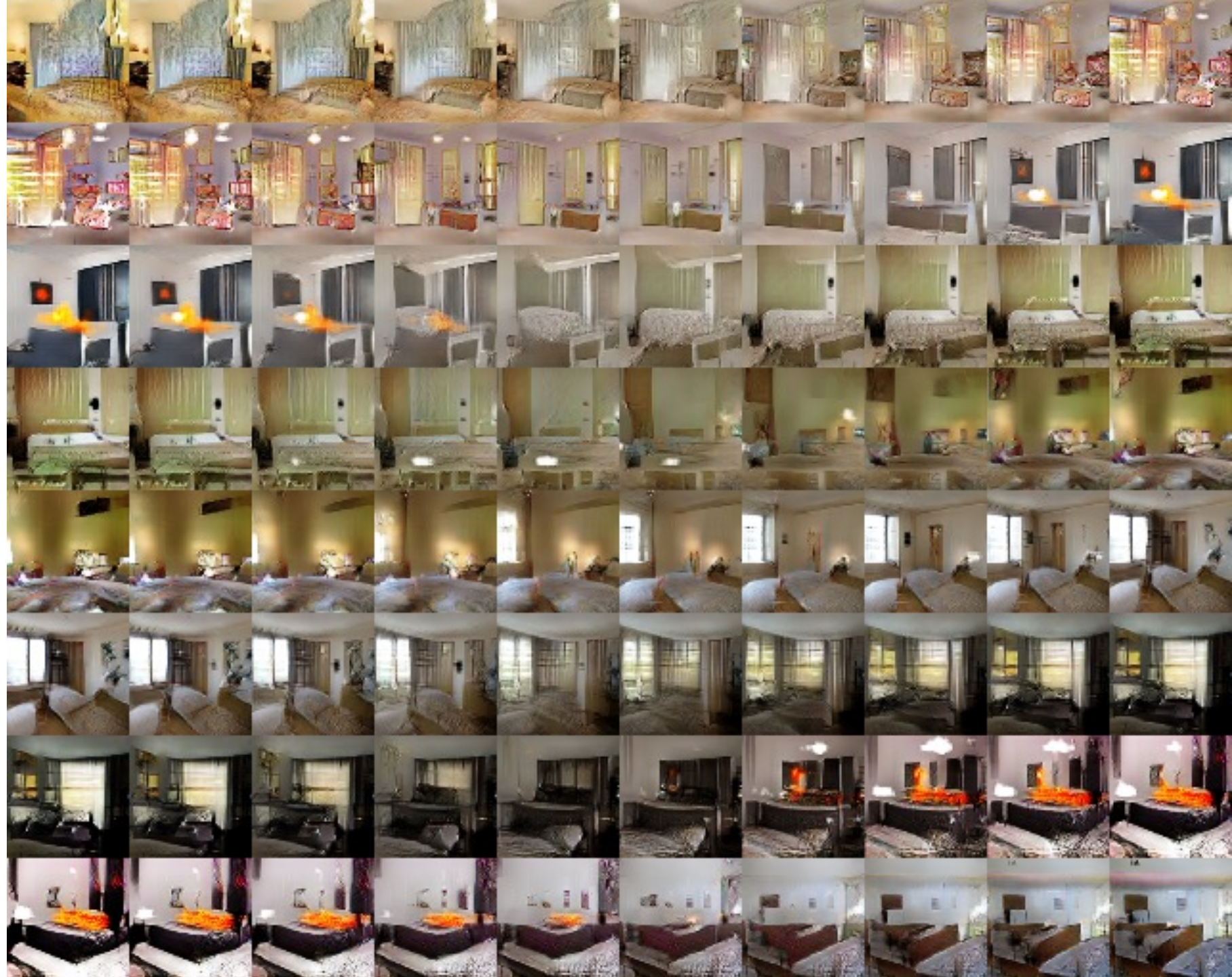




# Walking over the latent space

(Radford et al., 2015)

- Interpolation suggests non-overfitting behavior



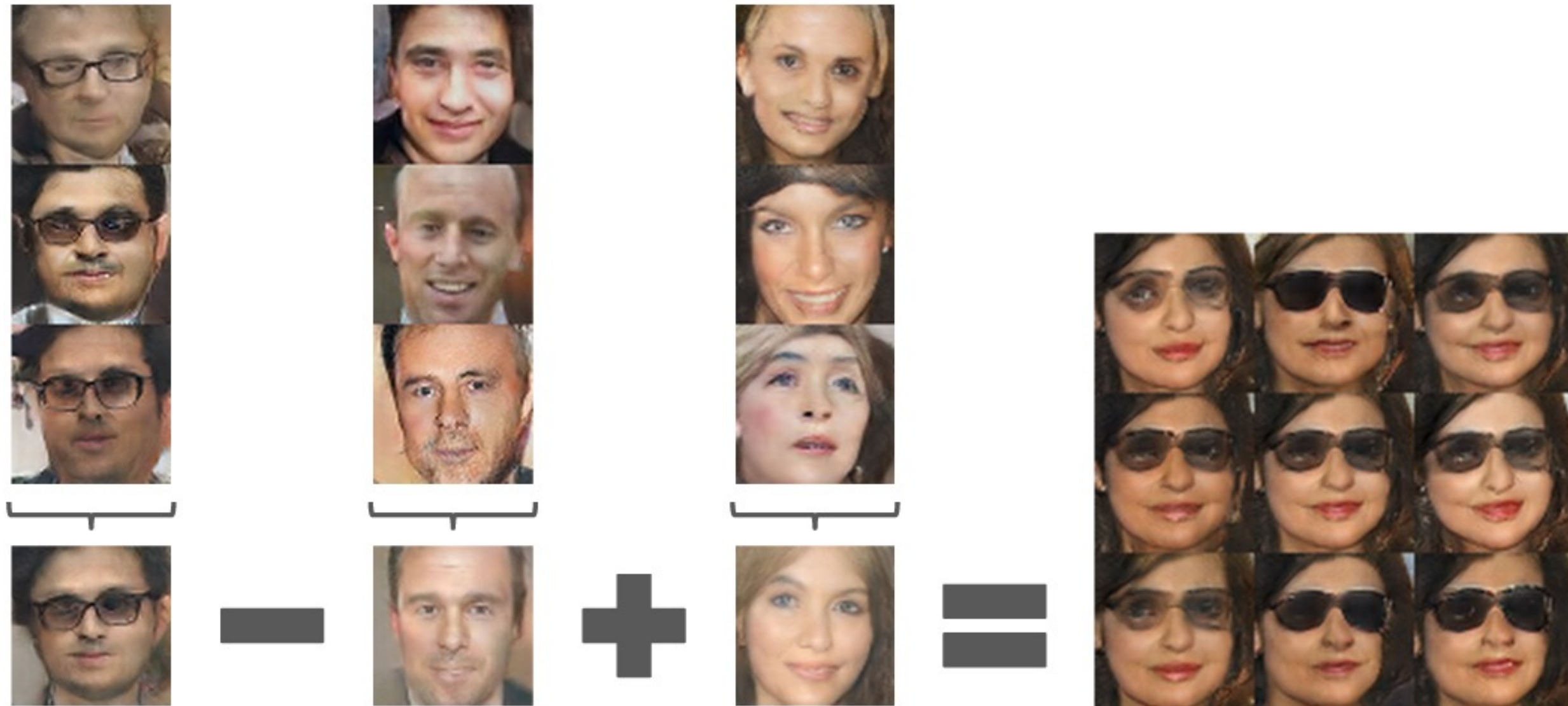
# Walking over the latent space

(Radford et al., 2015)



# Vector Space Arithmetic

(Radford et al., 2015)



man  
with glasses

man  
without glasses

woman  
without glasses

woman with glasses

# Vector Space Arithmetic

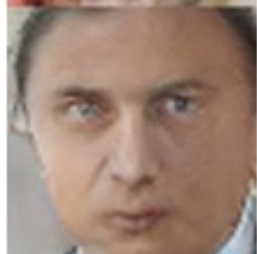
(Radford et al., 2015)



smiling woman



neutral woman

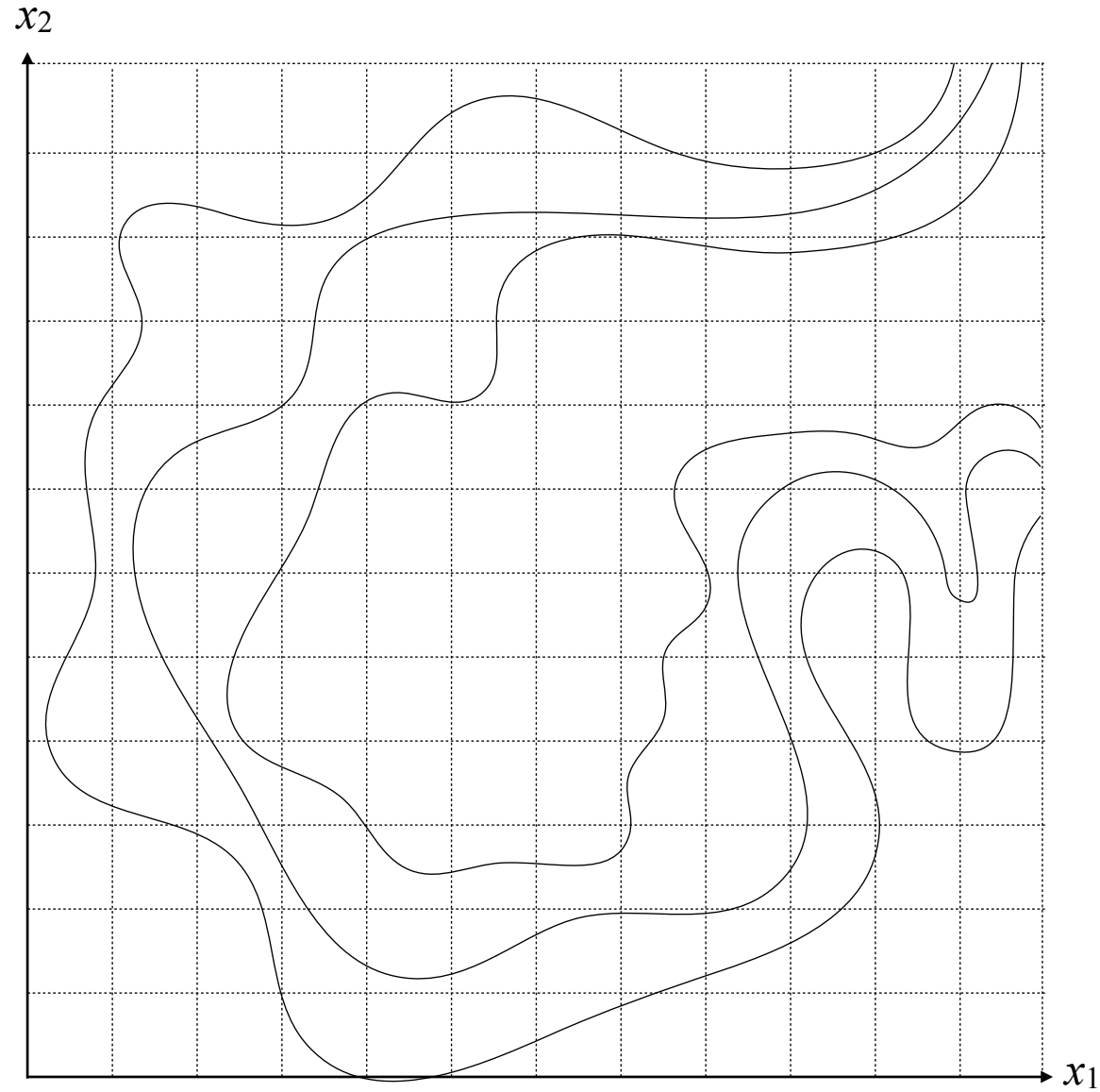


neutral man

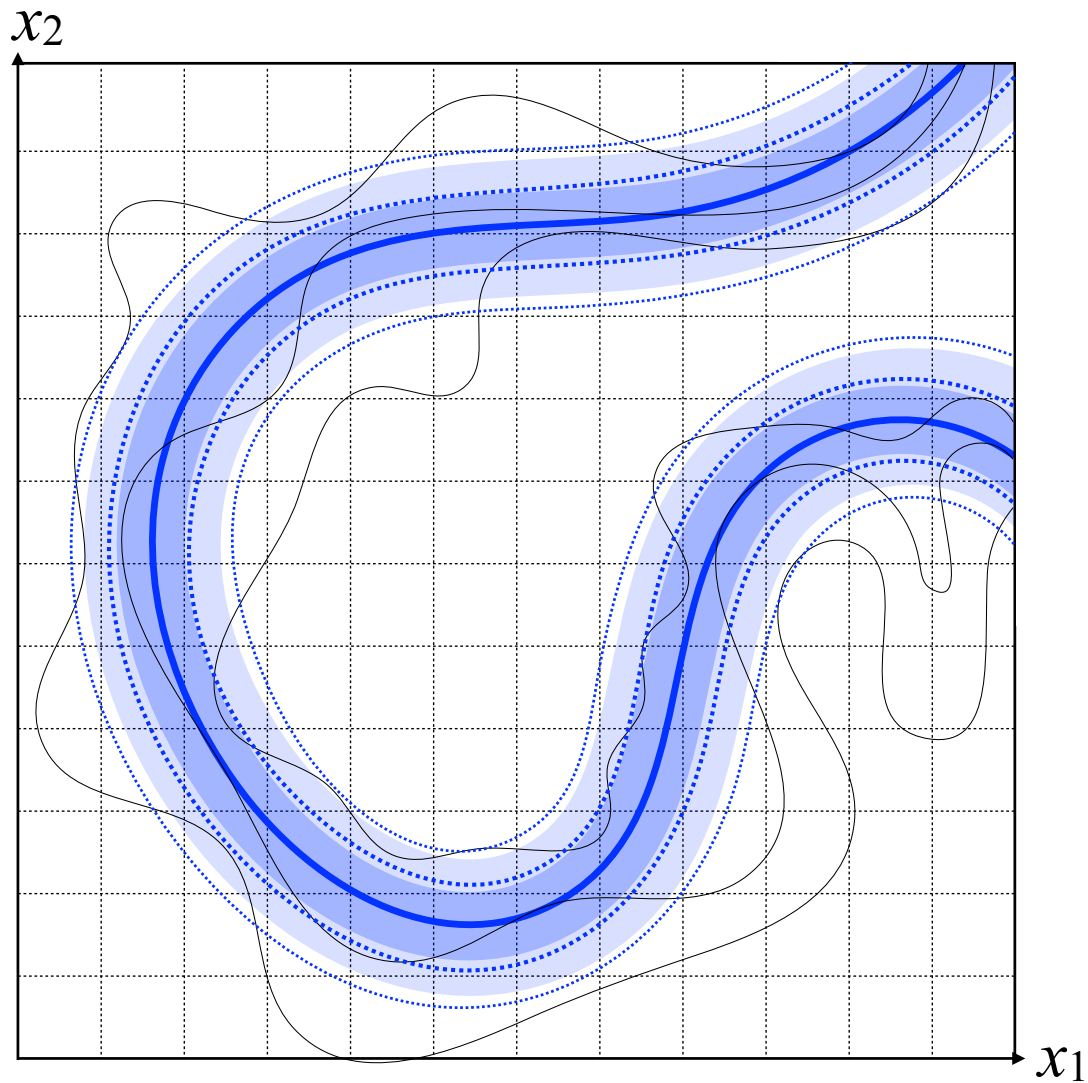


smiling man

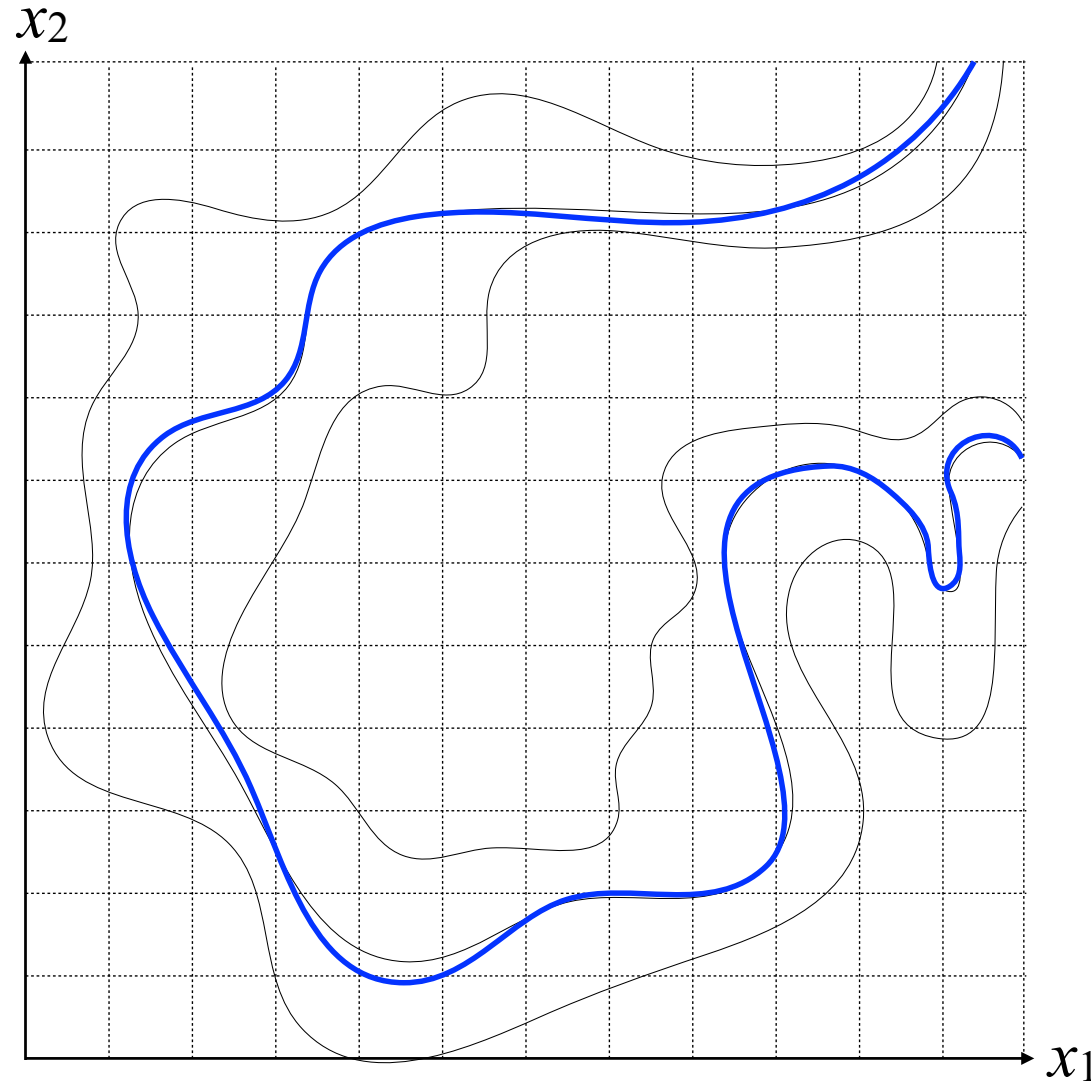
# Cartoon of the Image manifold



# What makes GANs special?



more traditional max-likelihood approach

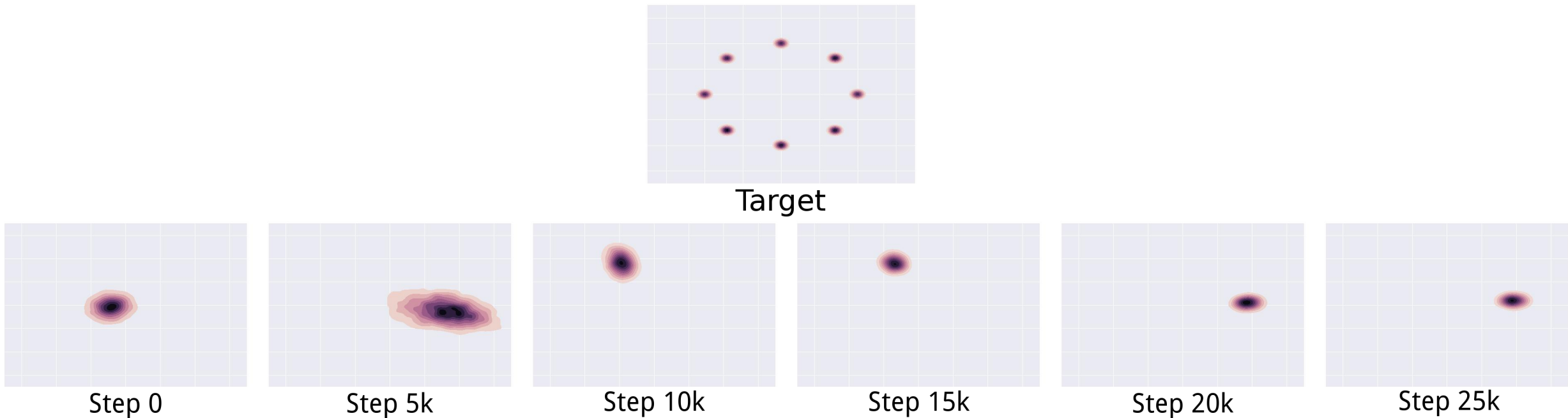


GAN

# GAN Failures: Mode Collapse

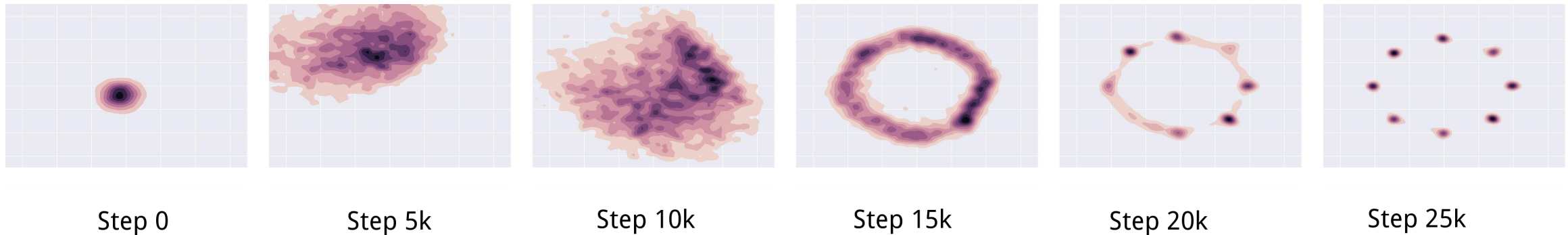
$$\min_G \max_D V(G, D) \neq \max_D \min_G V(G, D)$$

- $D$  in inner loop: convergence to correct distribution
- $G$  in inner loop: place all mass on most likely point



# Mode Collapse: Solutions

- **Unrolled GANs** (Metz et al 2016): Prevents mode collapse by backproping through a set of (k) updates of the discriminator to update generator parameters

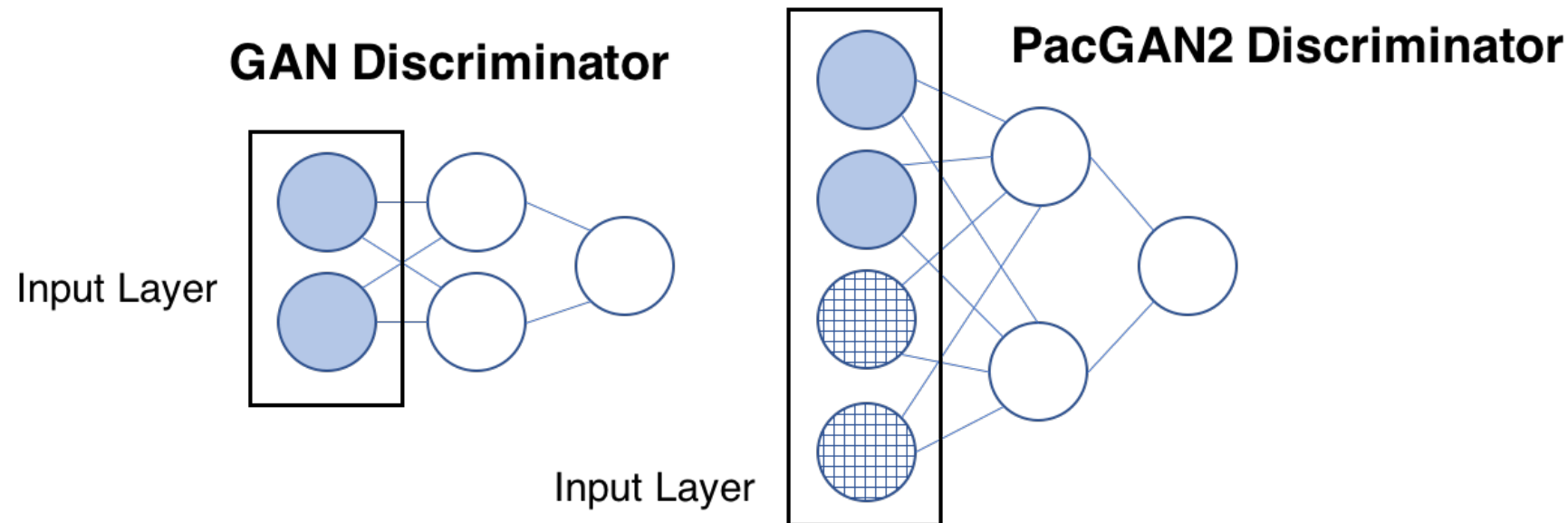


- **VEEGAN** (Srivastava et al 2017): Introduce a reconstructor network which is learned both to map the true data distribution  $p(x)$  to a Gaussian and to approximately invert the generator network.



# Mode Collapse: Solutions

- **Minibatch Discrimination** (Salimans et al 2016): Add minibatch features that classify each example by comparing it to other members of the minibatch (Salimans et al 2016)
- **PacGAN**: The power of two samples in generative adversarial networks (Lin et al 2017): Also uses multisample discrimination.



# Mode Collapse: Solutions

- **PacGAN:** The power of two samples in generative adversarial networks (Lin et al 2017): Also uses multisample discrimination.

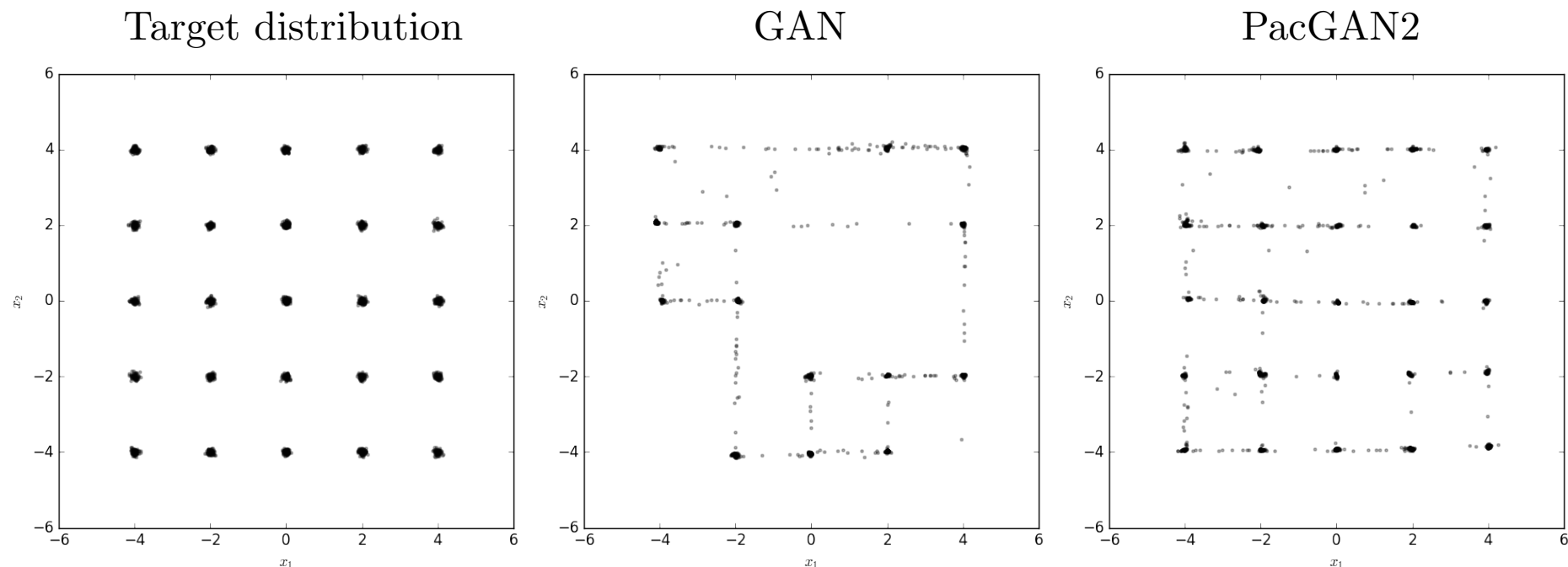


Figure 2: Scatter plot of the 2D samples from the true distribution (left) of 2D-grid and the learned generators using GAN (middle) and PacGAN2 (right). PacGAN2 captures all of the 25 modes.

# GAN Evaluation

- Quantitatively evaluating GANs is not straightforward:
  - Max Likelihood is a poor indication of sample quality
- Some evaluation metrics

- **Inception Score (IS):**

$y$  = labels given gen. image.  $p(y|x)$  is from classifier - InceptionNet

$$\text{IS}(\mathbb{P}_g) = e^{\mathbb{E}_{\mathbf{x} \sim \mathbb{P}_g} [KL(p_{\mathcal{M}}(y|\mathbf{x}) || p_{\mathcal{M}}(y))]}$$

- **Fréchet inception distance (FID):** (Currently most popular)

Estimate mean  $m$  and covariance  $C$  from classifier output - InceptionNet

$$d^2((m, C), (m_w, C_w)) = \|m - m_w\|_2^2 + \text{Tr}(C + C_w - 2(CC_w)^{1/2})$$

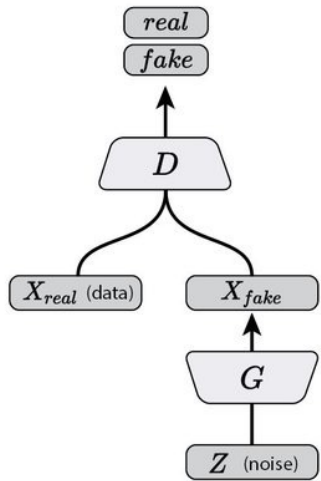
- **Kernel MMD** (Maximum Mean Discrepancy):

$$\text{MMD}(\mathbb{P}_r, \mathbb{P}_g) = \left( \mathbb{E}_{\substack{\mathbf{x}_r, \mathbf{x}'_r \sim \mathbb{P}_r, \\ \mathbf{x}_g, \mathbf{x}'_g \sim \mathbb{P}_g}} \left[ k(\mathbf{x}_r, \mathbf{x}'_r) - 2k(\mathbf{x}_r, \mathbf{x}_g) + k(\mathbf{x}_g, \mathbf{x}'_g) \right] \right)^{\frac{1}{2}}$$

# Subclasses of GANs

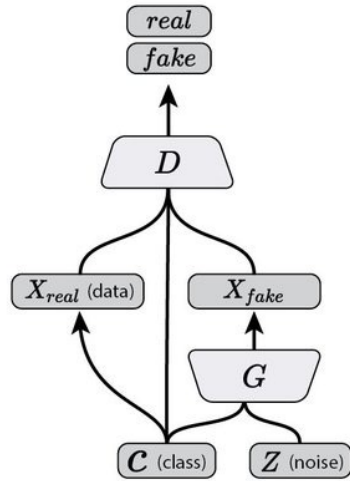
## Vanilla GAN

**Vanilla GAN**  
(Goodfellow, et al., 2014)

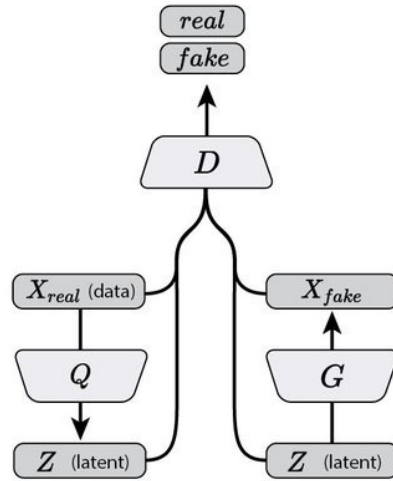


## Discriminator Looks at Latent Variables

**Conditional GAN**  
(Mirza & Osindero, 2014)

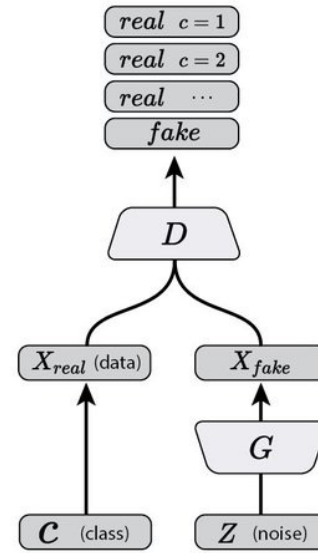


**Bidirectional GAN**  
(Donahue, et al., 2016; Dumoulin, et al., 2016)

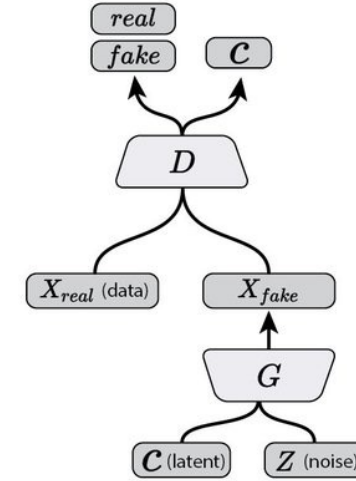


## Discriminator Predicts Latent Variables

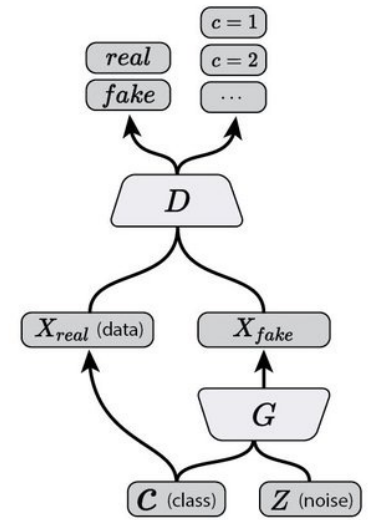
**Semi-Supervised GAN**  
(Odena, 2016; Salimans, et al., 2016)



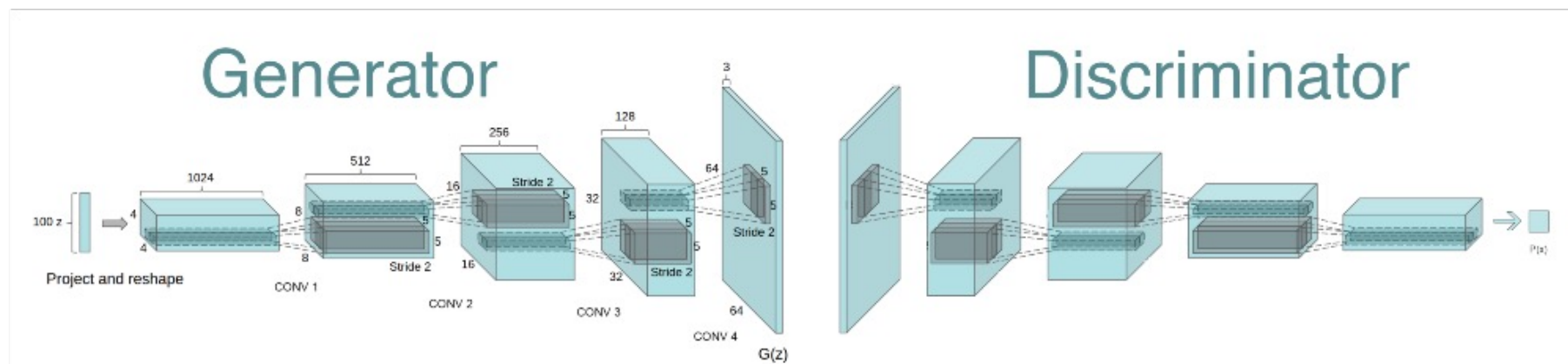
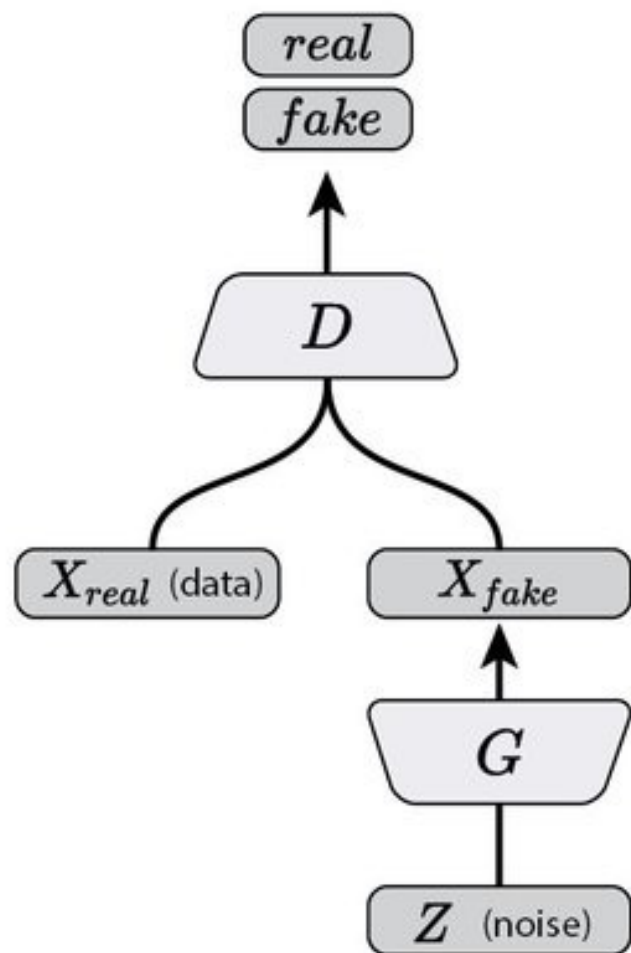
**InfoGAN**  
(Chen, et al., 2016)



**Auxiliary Classifier GAN**  
(Odena, et al., 2016)



# Vanilla GAN (Goodfellow et al., 2014)

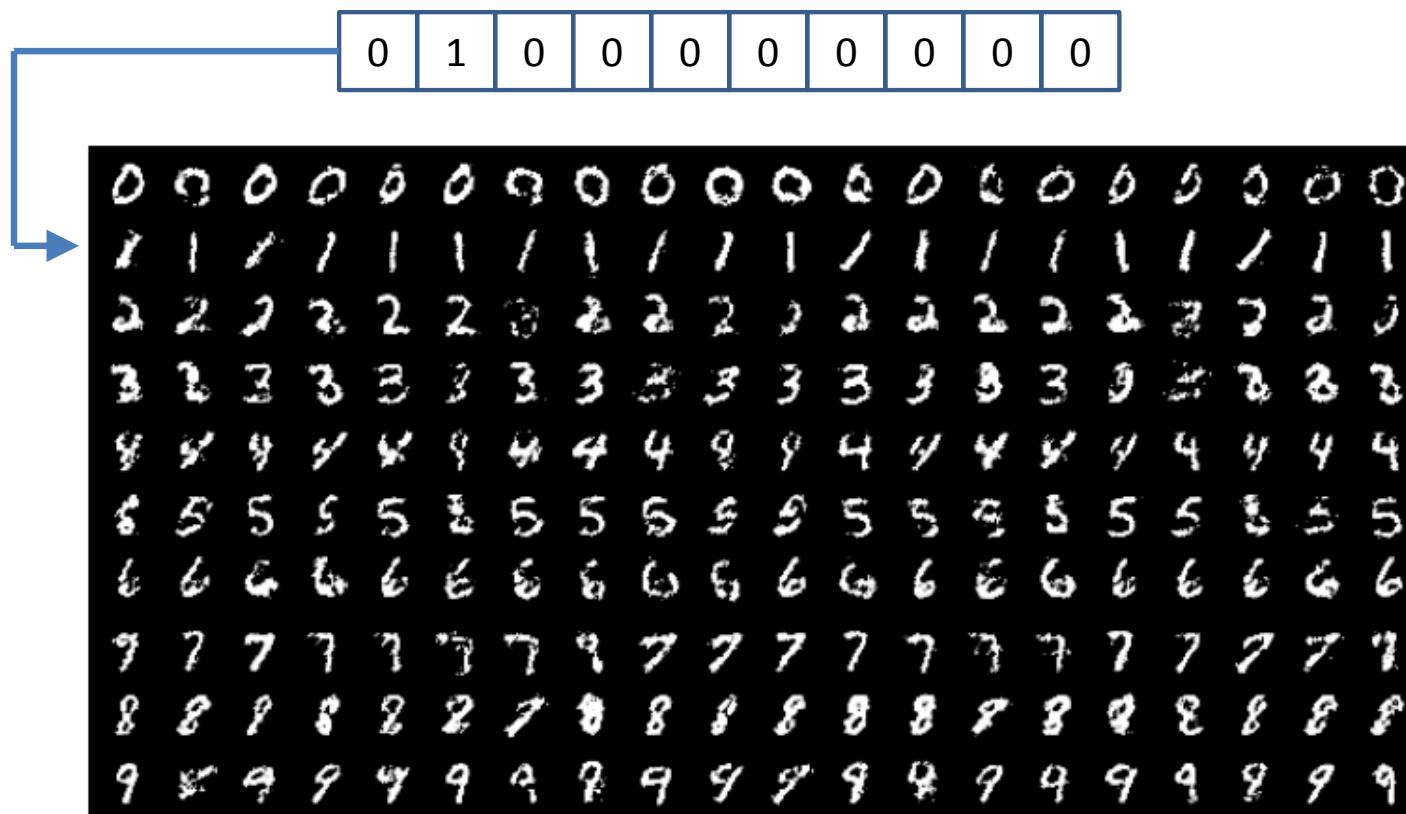
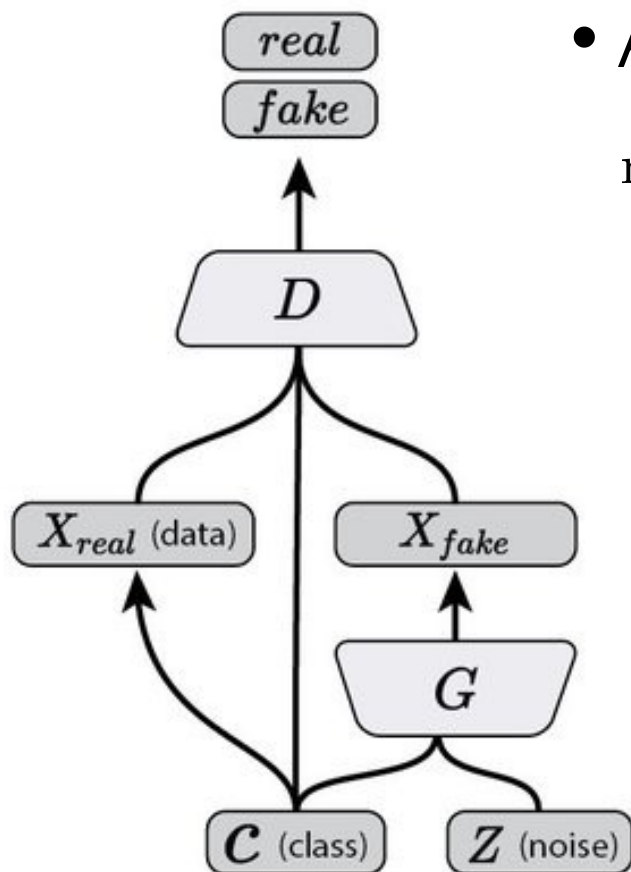


DCGAN (Radford et al., 2015)

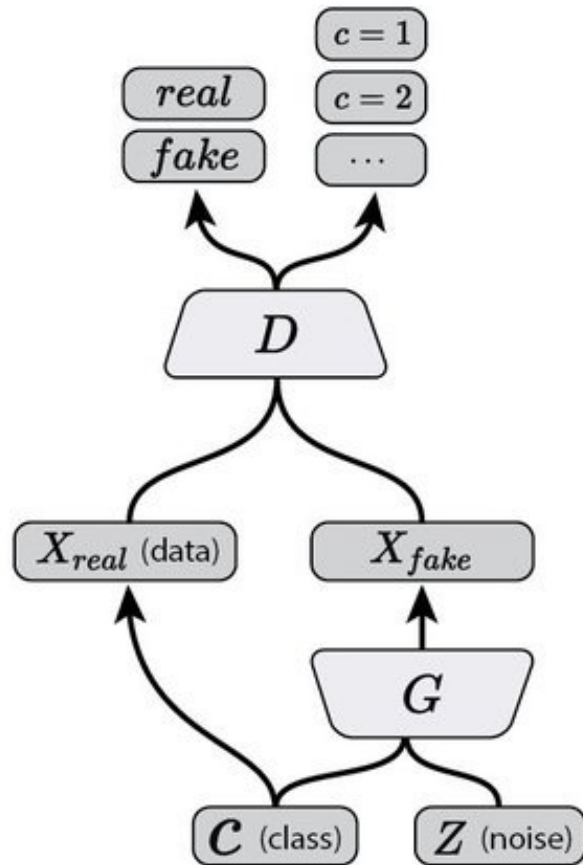
# Conditional GAN (Mirza and Osindero, 2014)

- Add conditional variables  $\mathbf{y}$  into  $G$  and  $D$

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x} | \mathbf{y})] + \mathbb{E}_{\mathbf{z} \sim p_z(\mathbf{z})} [\log(1 - D(G(\mathbf{z} | \mathbf{y})))]$$



# Auxiliary Classifier GAN (Odena et al., 2016)



- Every generated sample has a corresponding class label

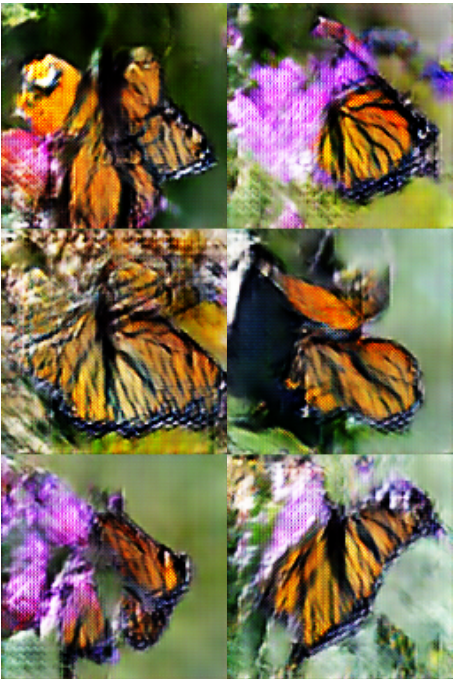
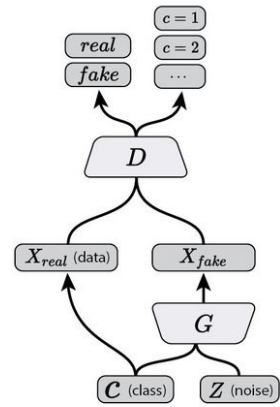
$$L_S = E[\log P(S = real \mid X_{real})] + E[\log P(S = fake \mid X_{fake})]$$

$$L_C = E[\log P(C = c \mid X_{real})] + E[\log P(C = c \mid X_{fake})]$$

- $D$  is trained to maximize  $L_S + L_C$
- $G$  is trained to maximize  $L_C - L_S$
- Learns a representation for  $z$  that is independent of class label

# Auxiliary Classifier GAN (Odena et al., 2016)

128×128 resolution samples from 5 classes taken from an AC-GAN trained on the ImageNet



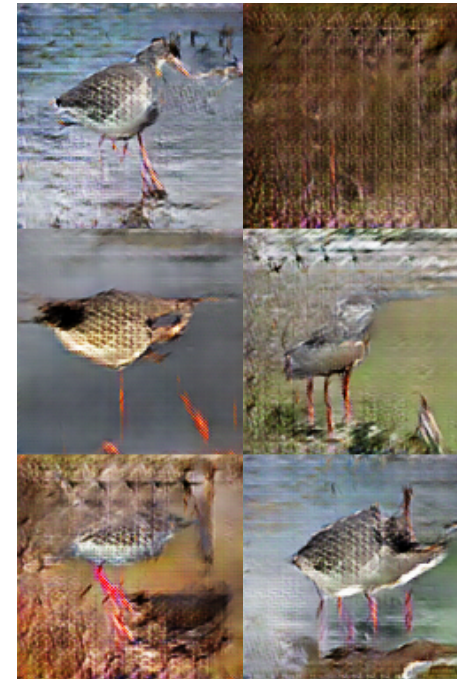
monarch butterfly



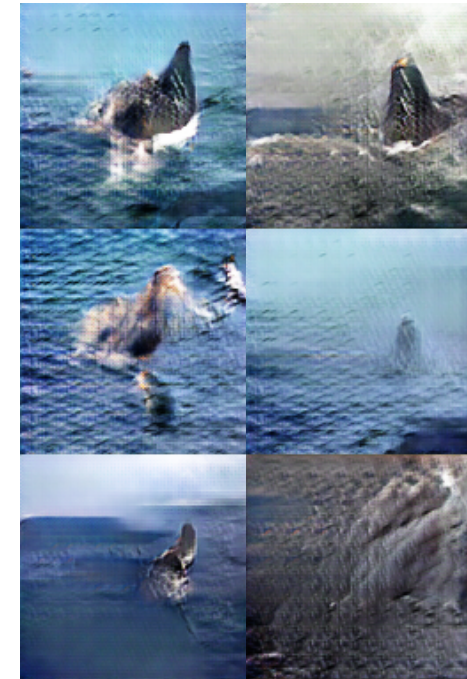
goldfinch



daisy



redshank



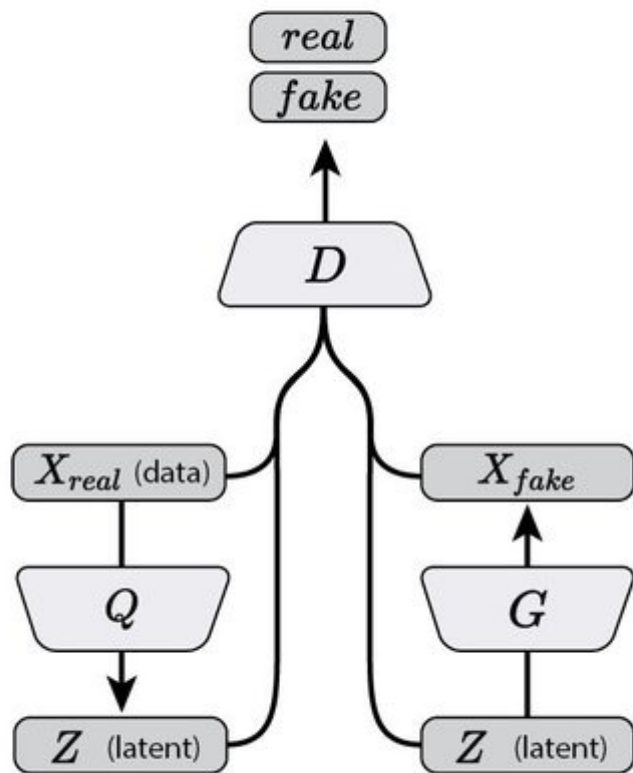
grey whale



# Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)

- Jointly learns a generator network and an inference network using an adversarial process.

$$\begin{aligned} \min_G \max_D V(D, G) &= \mathbb{E}_{q(\mathbf{x})} [\log(D(\mathbf{x}, G_z(\mathbf{x})))] + \mathbb{E}_{p(\mathbf{z})} [\log(1 - D(G_x(\mathbf{z}), \mathbf{z}))] \\ &= \iint q(\mathbf{x})q(\mathbf{z} | \mathbf{x}) \log(D(\mathbf{x}, \mathbf{z})) d\mathbf{x}d\mathbf{z} \\ &+ \iint p(\mathbf{z})p(\mathbf{x} | \mathbf{z}) \log(1 - D(\mathbf{x}, \mathbf{z})) d\mathbf{x}d\mathbf{z}. \end{aligned}$$

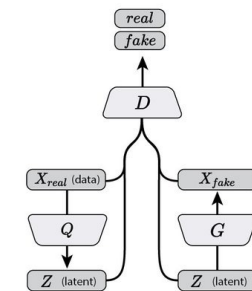


CelebA reconstructions



SVNH reconstructions

# Bidirectional GAN (Donahue et al., 2016; Dumoulin et al., 2016)



LSUN bedrooms



Tiny ImageNet

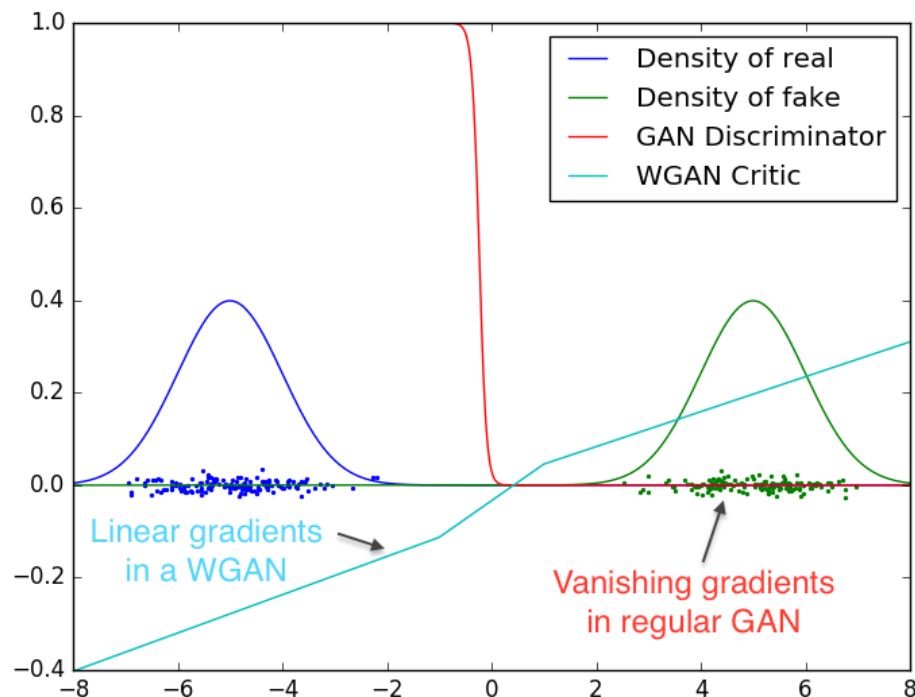


# Wasserstein GAN (Arjovsky et al., 2016)

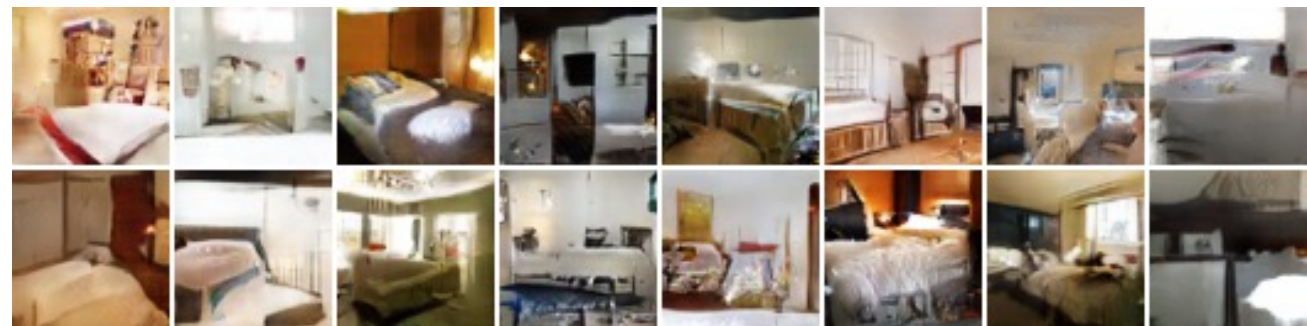
- Objective based on Earth-Mover or Wasserstein distance:

$$\min_{\theta} \max_{\omega} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [D_{\omega}(\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}} [D_{\omega}(G_{\theta}(\mathbf{z}))]$$

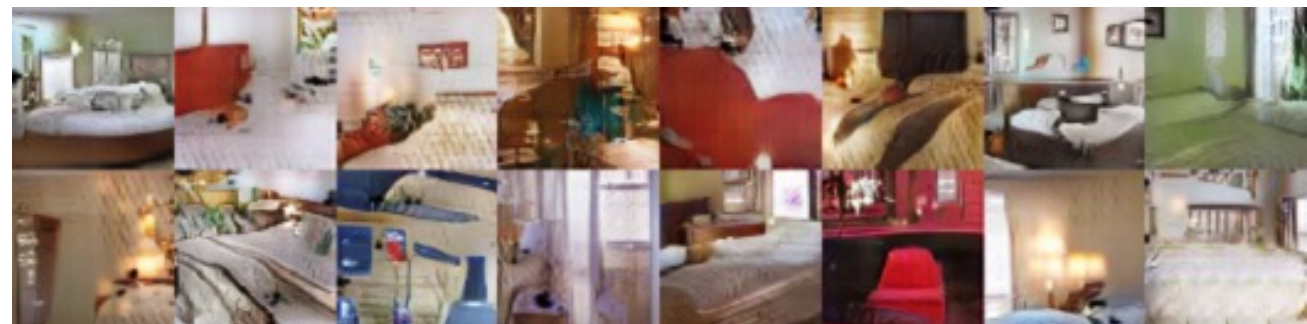
- Provides nice gradients over real and fake samples



WGAN

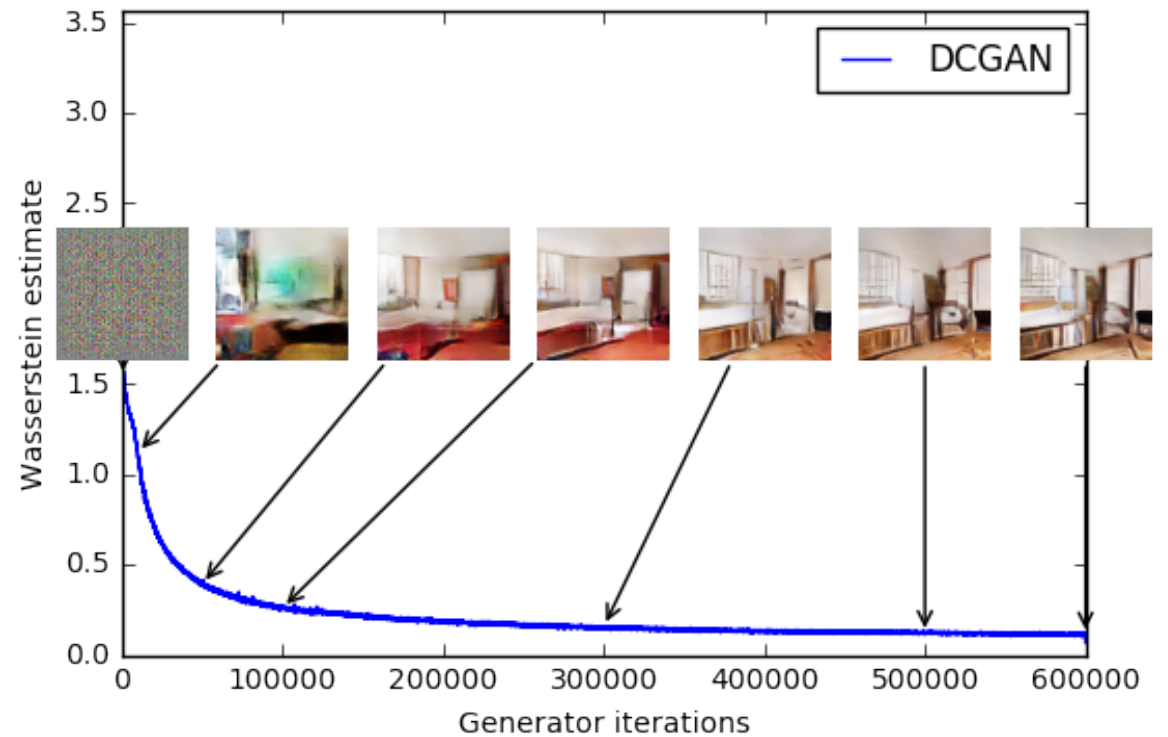
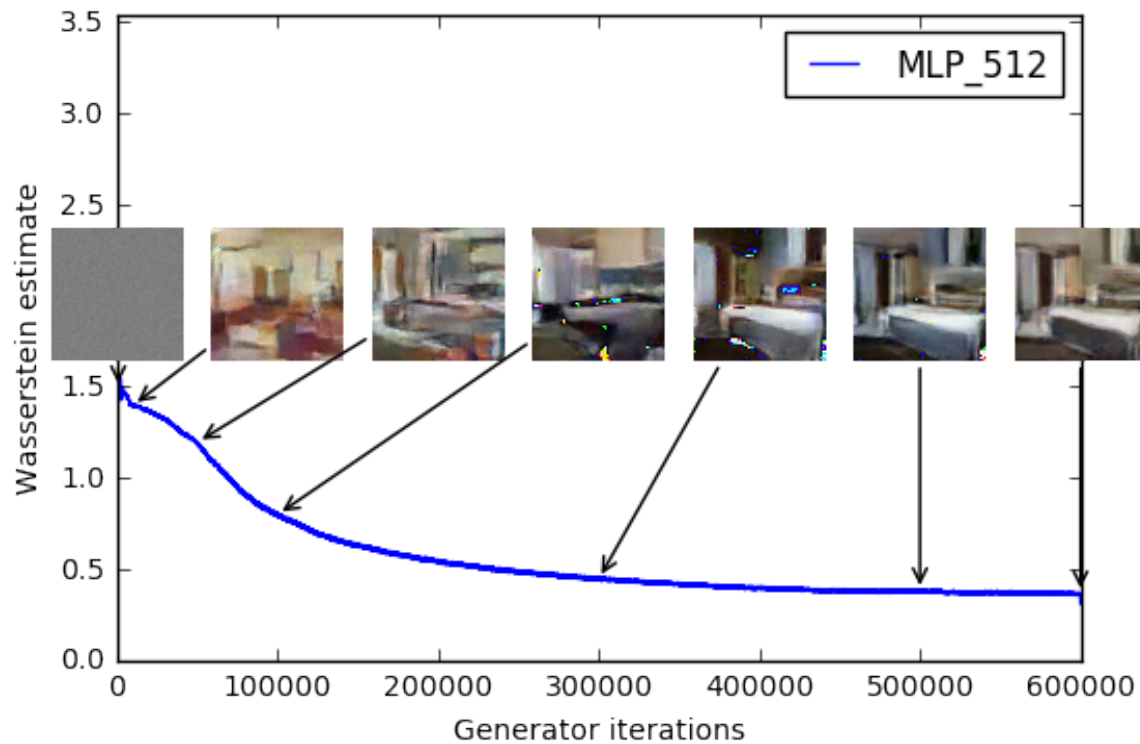


DCGAN



# Wasserstein GAN (Arjovsky et al., 2016)

- Wasserstein loss seems to correlate well with image quality.



# WGAN with gradient penalty (Gulraani et al., 2017)

$$L = \underbrace{\mathbb{E}_{\tilde{\mathbf{x}} \sim \mathbb{P}_g} [D(\tilde{\mathbf{x}})] - \mathbb{E}_{\mathbf{x} \sim \mathbb{P}_r} [D(\mathbf{x})]}_{\text{Original critic loss}} + \lambda \underbrace{\mathbb{E}_{\hat{\mathbf{x}} \sim \mathbb{P}_{\hat{\mathbf{x}}}} [(\|\nabla_{\hat{\mathbf{x}}} D(\hat{\mathbf{x}})\|_2 - 1)^2]}_{\text{Our gradient penalty}}$$

- Faster convergence and higher-quality samples than WGAN with weight clipping
- Train a wide variety of GAN architectures with almost no hyperparameter tuning, including discrete models

Samples from a character-level GAN language model on Google Billion Word

---

## WGAN with gradient penalty

Busino game camperate spent odea  
 In the bankaway of smarling the  
 SingersMay , who kill that invic  
 Keray Pents of the same Reagun D  
 Manging include a tudancs shat "  
 His Zuith Dudget , the Denmbern  
 In during the Uitational questio  
 Divos from The ' noth ronkies of  
 She like Monday , of macunsuer S  
 The investor used ty the present  
 A papees are country congress oo  
 A few year inom the group that s  
 He said this syenn said they wan  
 As a world 1 88 ,for Autouries  
 Foand , th Word people car , Il  
 High of the upseader homing pull  
 The guipe is worly move dogsfor  
 The 1874 incidested he could be  
 The allo tooks to security and c

Solice Norkedin pring in since  
 ThiS record ( 31. ) UBS ) and Ch  
 It was not the annuas were plogr  
 This will be us , the ect of DAN  
 These leaded as most-worsd p2 a0  
 The time I paid0a South Cubry i  
 Dour Fraps higs it was these del  
 This year out howneed allowed lo  
 Kaulna Seto consficutes to repor  
 A can teal , he was schoon news  
 In th 200. Pesish picriers rega  
 Konney Panice rimimber the teami  
 The new centuct cut Denester of  
 The near , had been one injustie  
 The incestion to week to shorted  
 The company the high product of  
 20 - The time of accomplete , wh  
 John WVuderenson seqiivic spends  
 A ceetens in indestedly the Wat

---

## Standard GAN objective

dddddddddddddddddddddddddddddd  
 ddddddddddddddddddddddddddd

dddddddddddddddddddddddddd  
 ddddddddddddddddddddddd

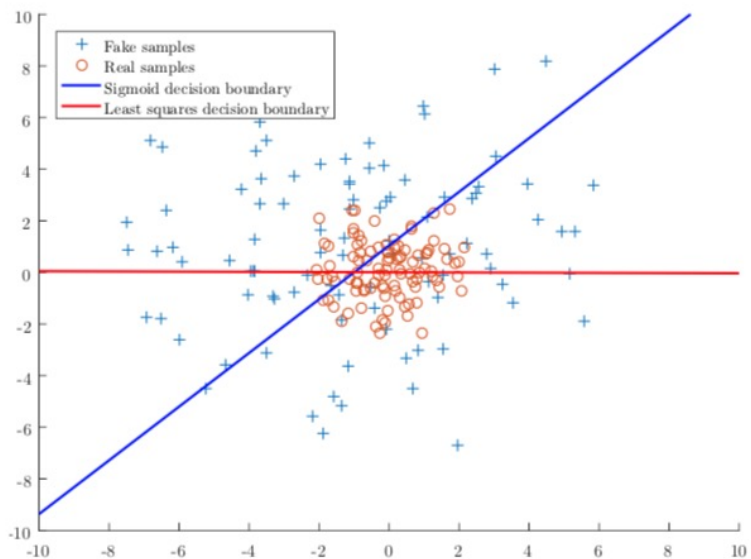
---

# Least Squares GAN (LSGAN) (Mao et al., 2017)

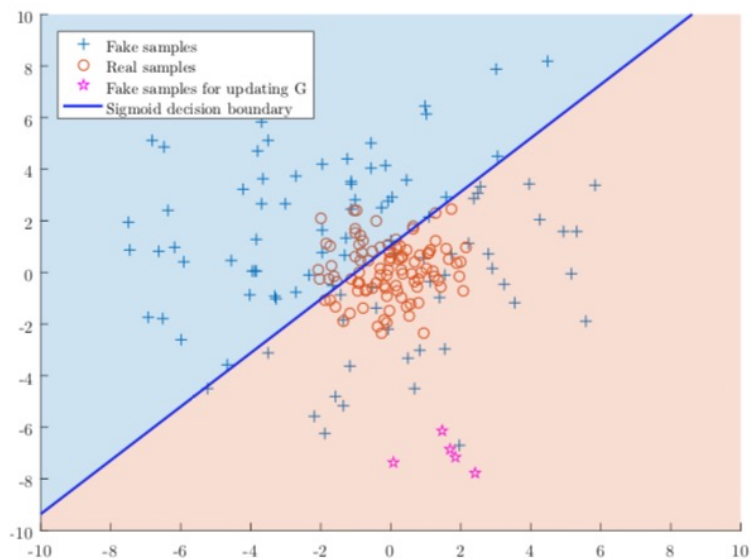
- Use a loss function that provides smooth and non-saturating gradient in discriminator D

$$\min_D V_{\text{LSGAN}}(D) = \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [(D(\mathbf{x}) - b)^2] + \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z}))) - a]^2]$$

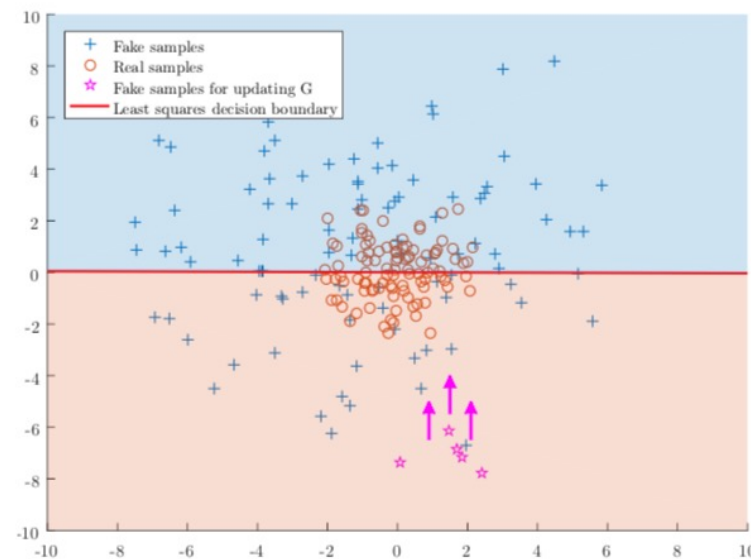
$$\min_G V_{\text{LSGAN}}(G) = \frac{1}{2} \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [(D(G(\mathbf{z}))) - c]^2,$$



Decision boundaries of Sigmoid & Least Squares loss functions



Sigmoid decision boundary

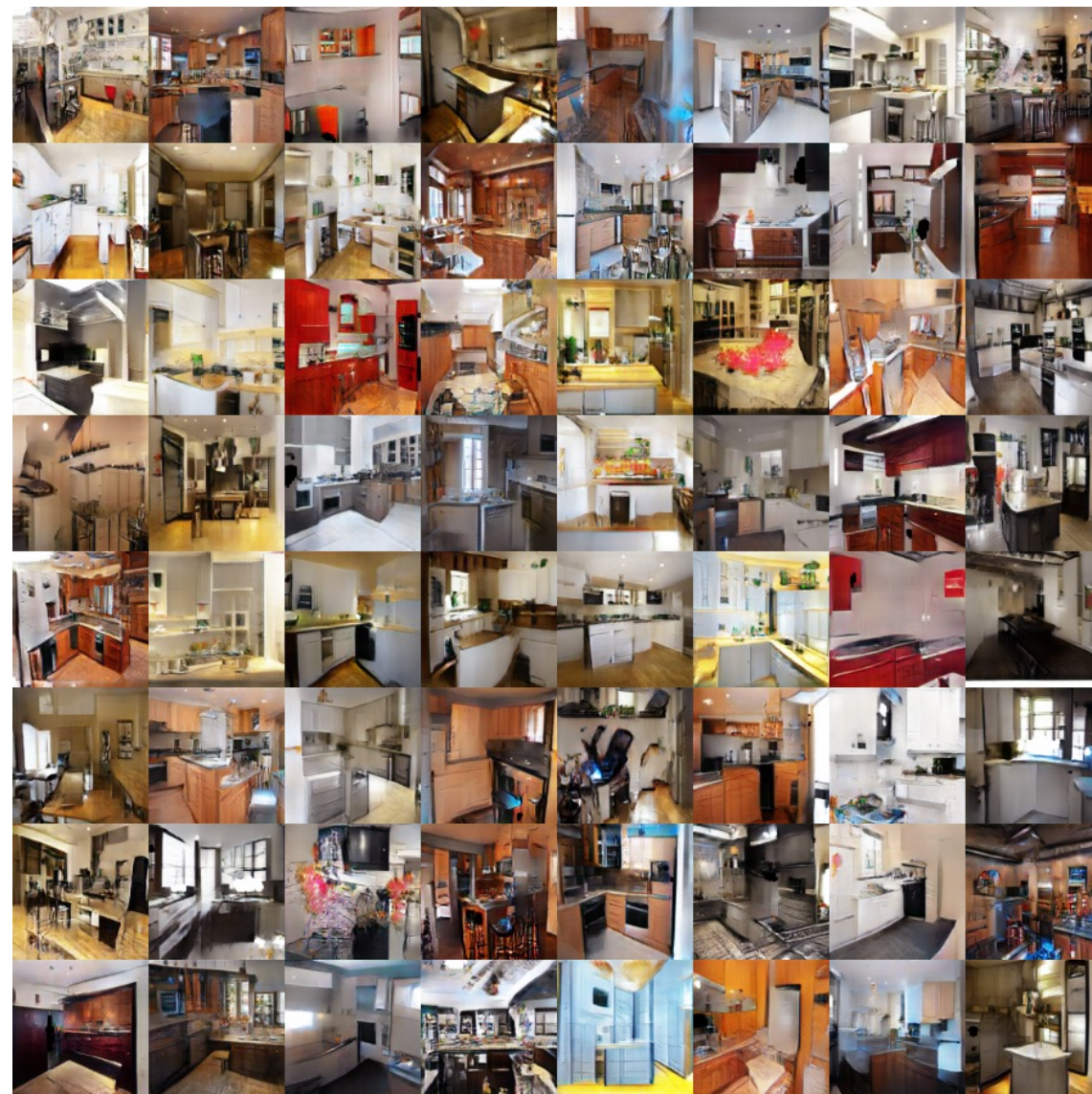


Least Squares decision boundary

# Least Squares GAN (LSGAN) (Mao et al., 2017)



Church



Kitchen

# Boundary Equilibrium GAN (BEGAN)

(Berthelot et al., 2017)

- A loss derived from the Wasserstein distance for training auto-encoder based GANs

$$\mathcal{L}(v) = |v - D(v)|^\eta \text{ where } \begin{cases} D : \mathbb{R}^{N_x} \mapsto \mathbb{R}^{N_x} & \text{is the autoencoder function.} \\ \eta \in \{1, 2\} & \text{is the target norm.} \\ v \in \mathbb{R}^{N_x} & \text{is a sample of dimension } N_x. \end{cases}$$

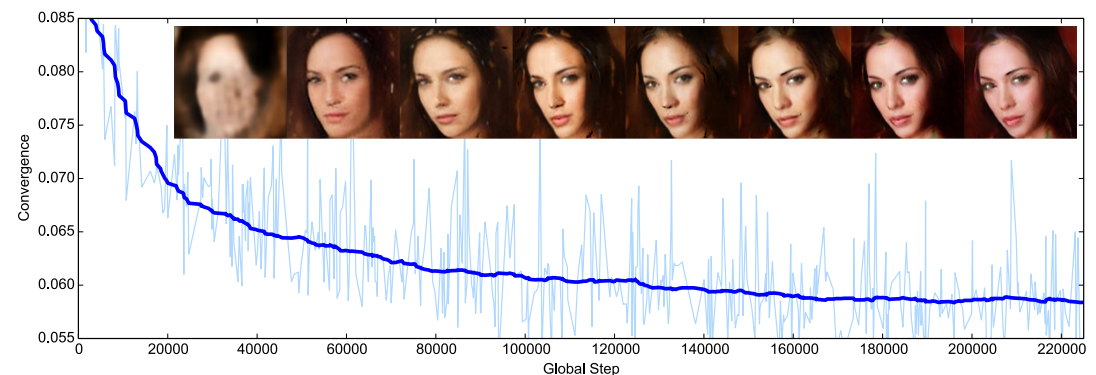
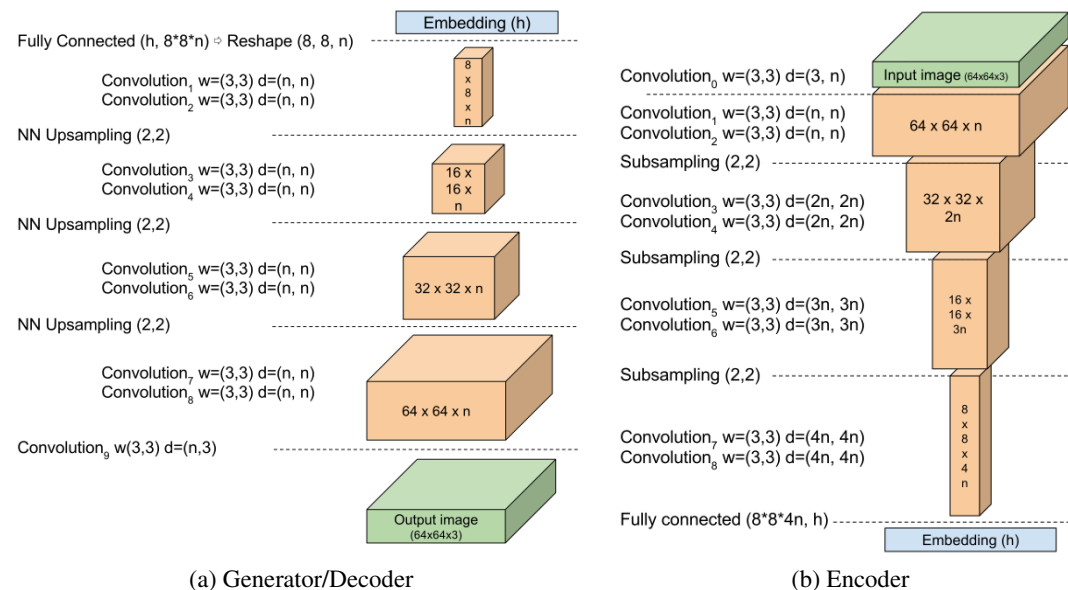
- Wasserstein distance btw. the reconstruction losses of real and generated data

- Convergence measure:

$$\mathcal{M}_{global} = \mathcal{L}(x) + |\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))|$$

- Objective:

$$\begin{cases} \mathcal{L}_D = \mathcal{L}(x) - k_t \cdot \mathcal{L}(G(z_D)) & \text{for } \theta_D \\ \mathcal{L}_G = \mathcal{L}(G(z_G)) & \text{for } \theta_G \\ k_{t+1} = k_t + \lambda_k (\gamma \mathcal{L}(x) - \mathcal{L}(G(z_G))) & \text{for each training step } t \end{cases}$$





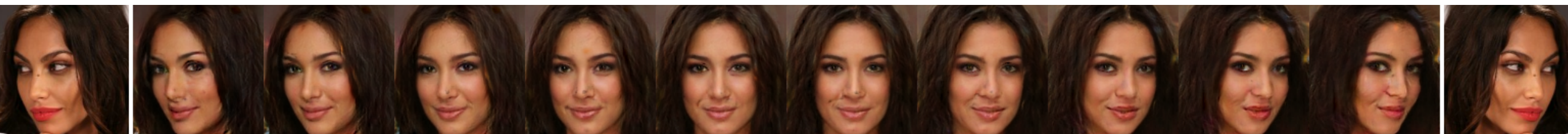
# BEGANs for CelebA

360K celebrity face images  
128x128 with 128 filters

(Berthelot et al., 2017)



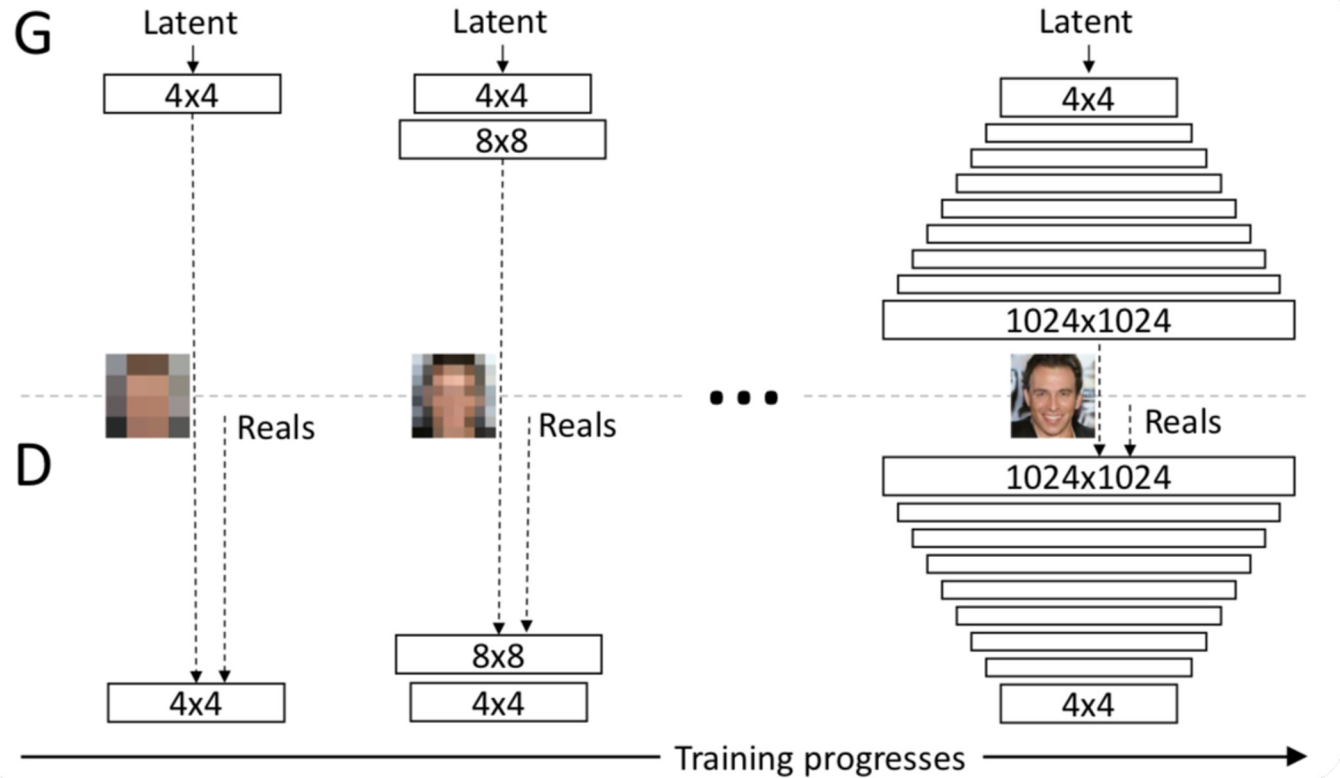
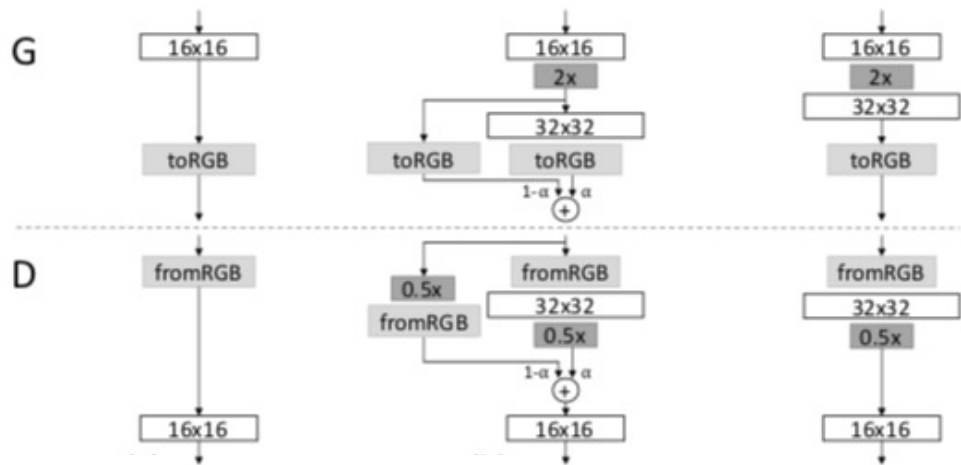
Interpolations in the latent space



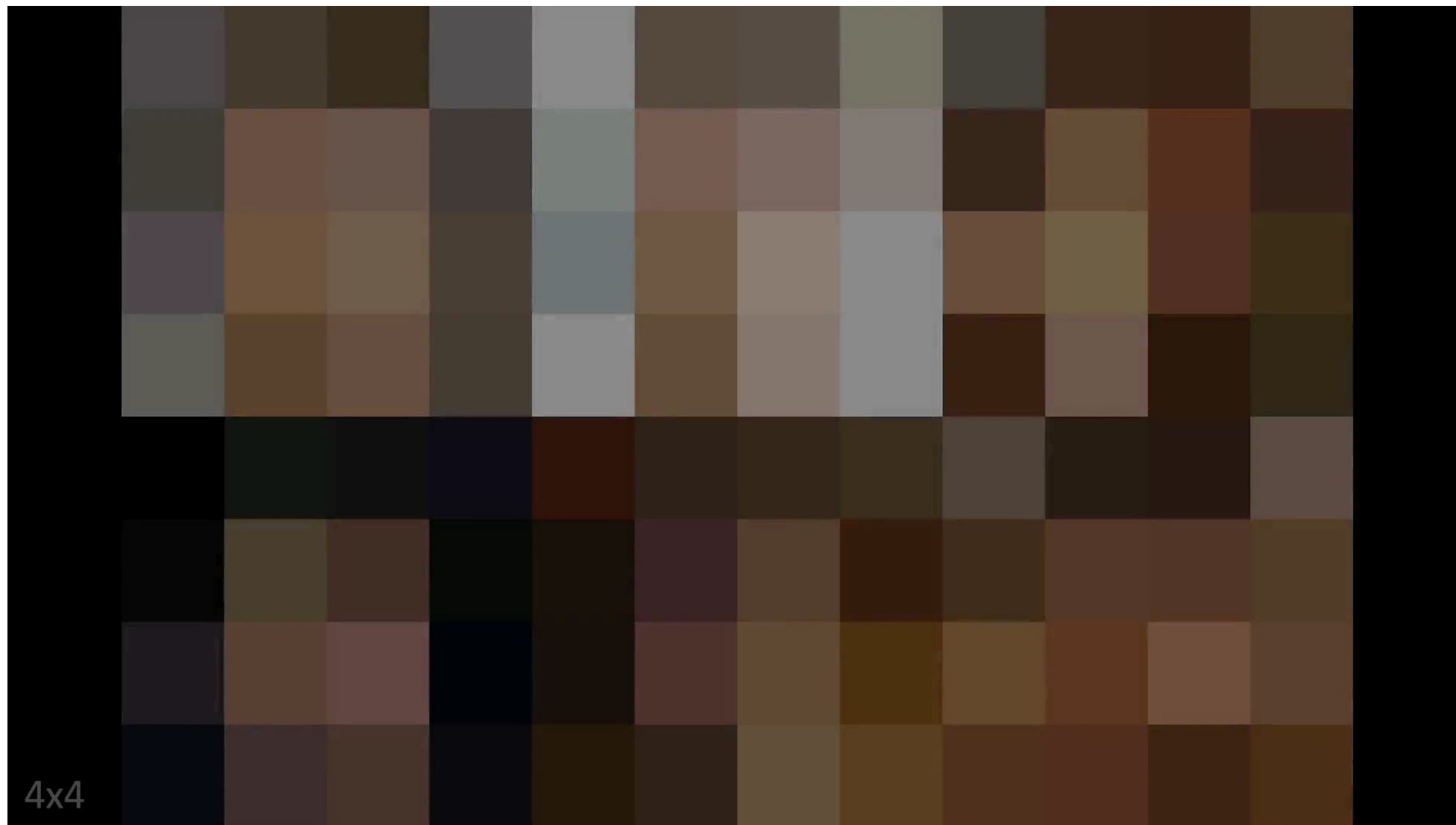
Mirror interpolation example

# Progressive GANs (Karras et al., 2018)

- Progressively generate high-res images
- Multi-step training from low to high resolutions



# Progressive GANs (Karras et al., 2018)



- Training process

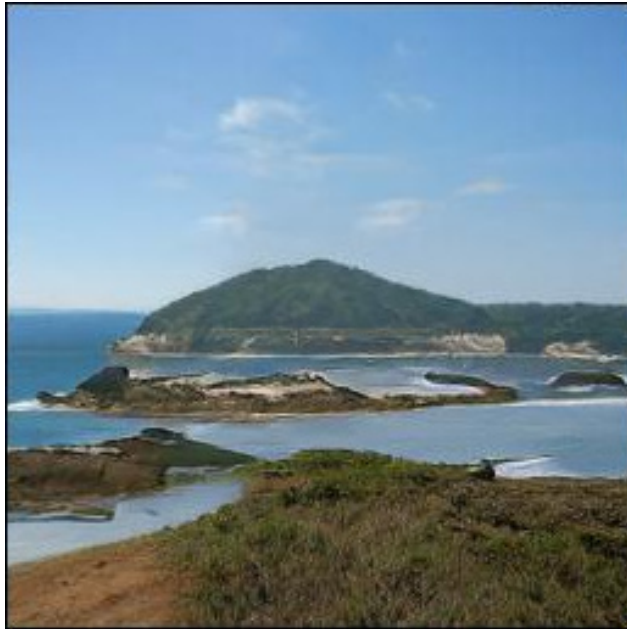
# Progressive GANs (Karras et al., 2018)



CelebA-HQ  
random interpolations

# BigGANs (Brock et al., 2019)

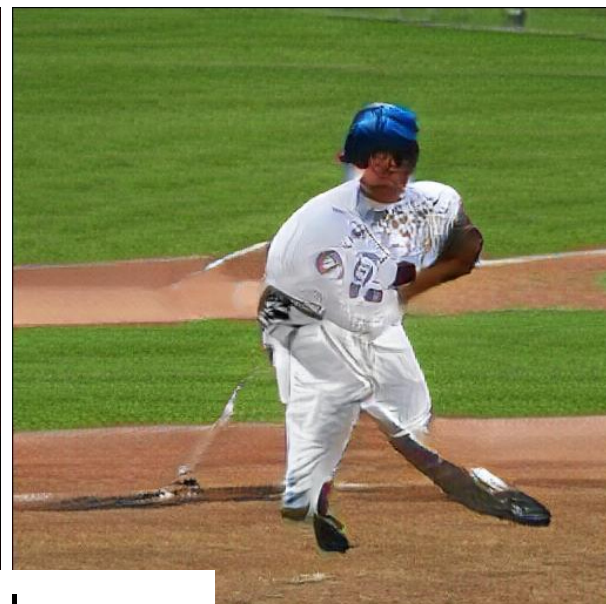
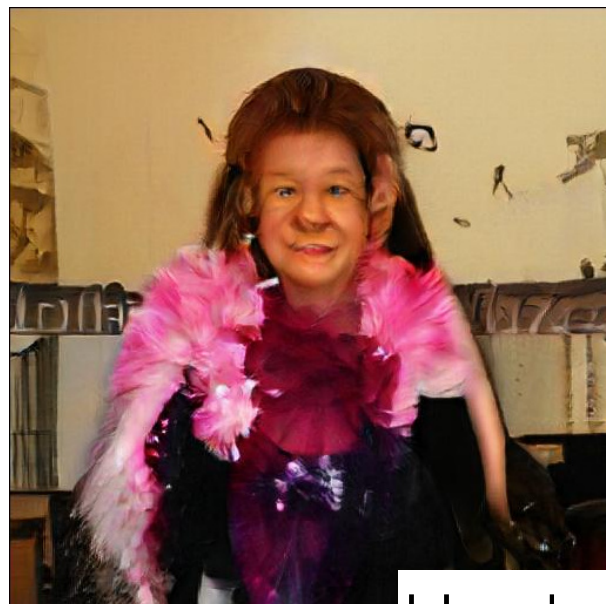
High resolution, class-conditional samples generated by the model



- BigGANs trained with 2-4x as many parameters and 8x the batch size compared to prior art.
- Uses Gaussian truncation to sample  $z$  (avoid sampling from the tail of the Gaussian distribution)
- Uses multiple other tricks including multiple regularizations including a Gradient penalty regularization and an Orthogonal Regularization:

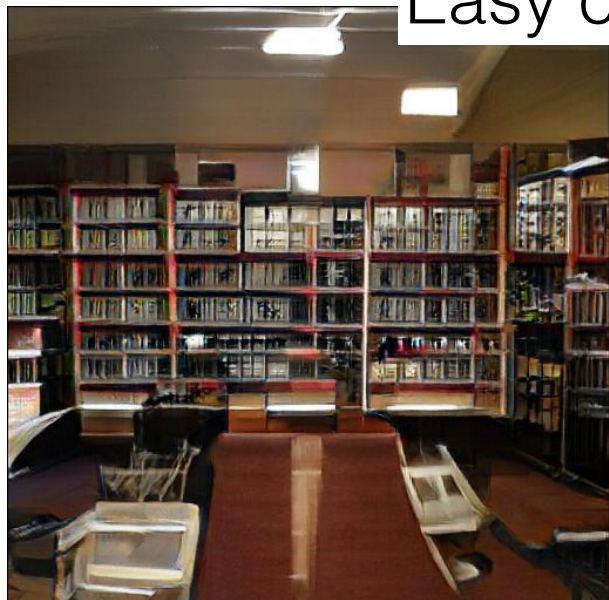
$$R_{\beta}(W) = \beta \|W^T W \odot (\mathbf{1} - I)\|_F^2,$$

# BigGANs (Brock et al., 2019)



Easy classes

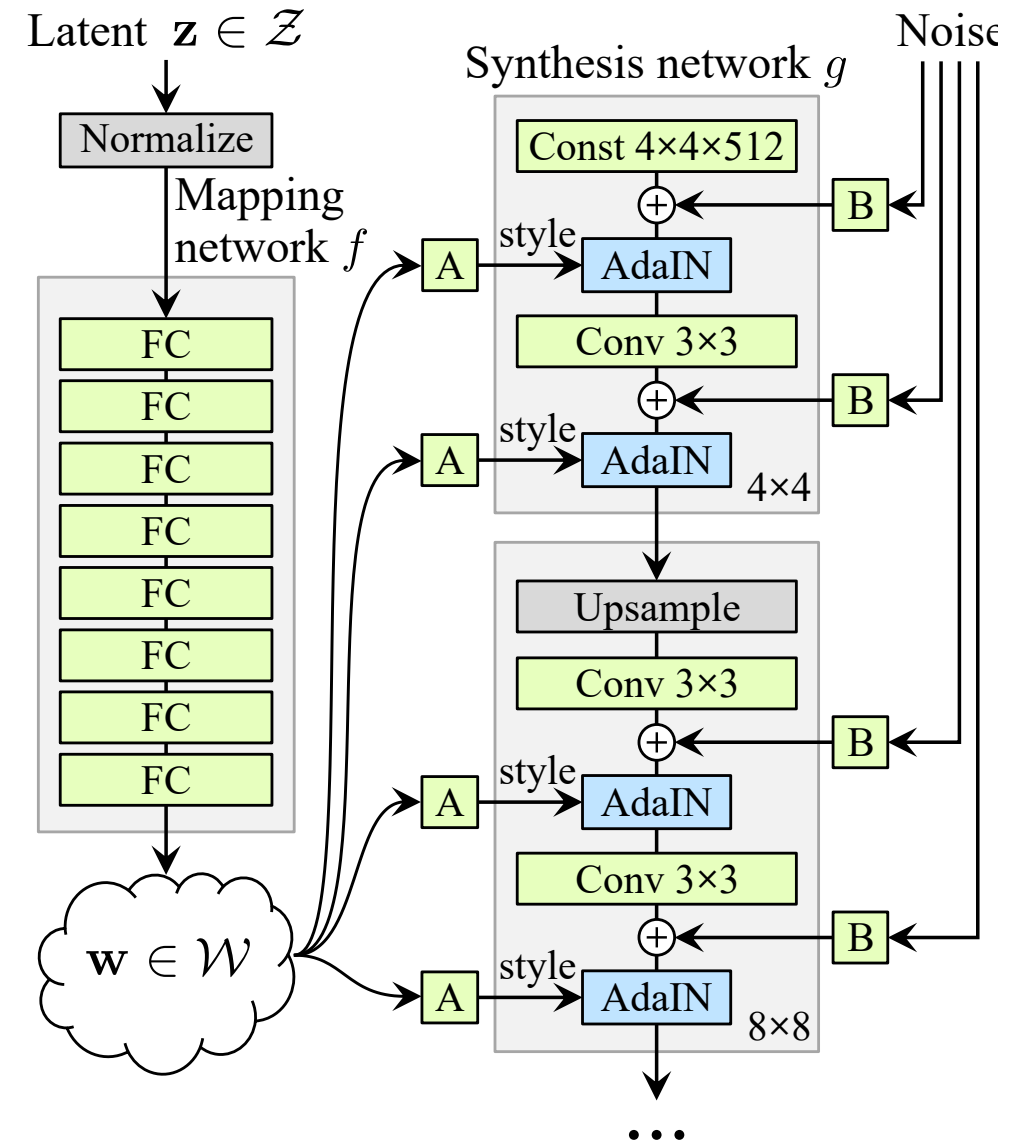
Hard classes



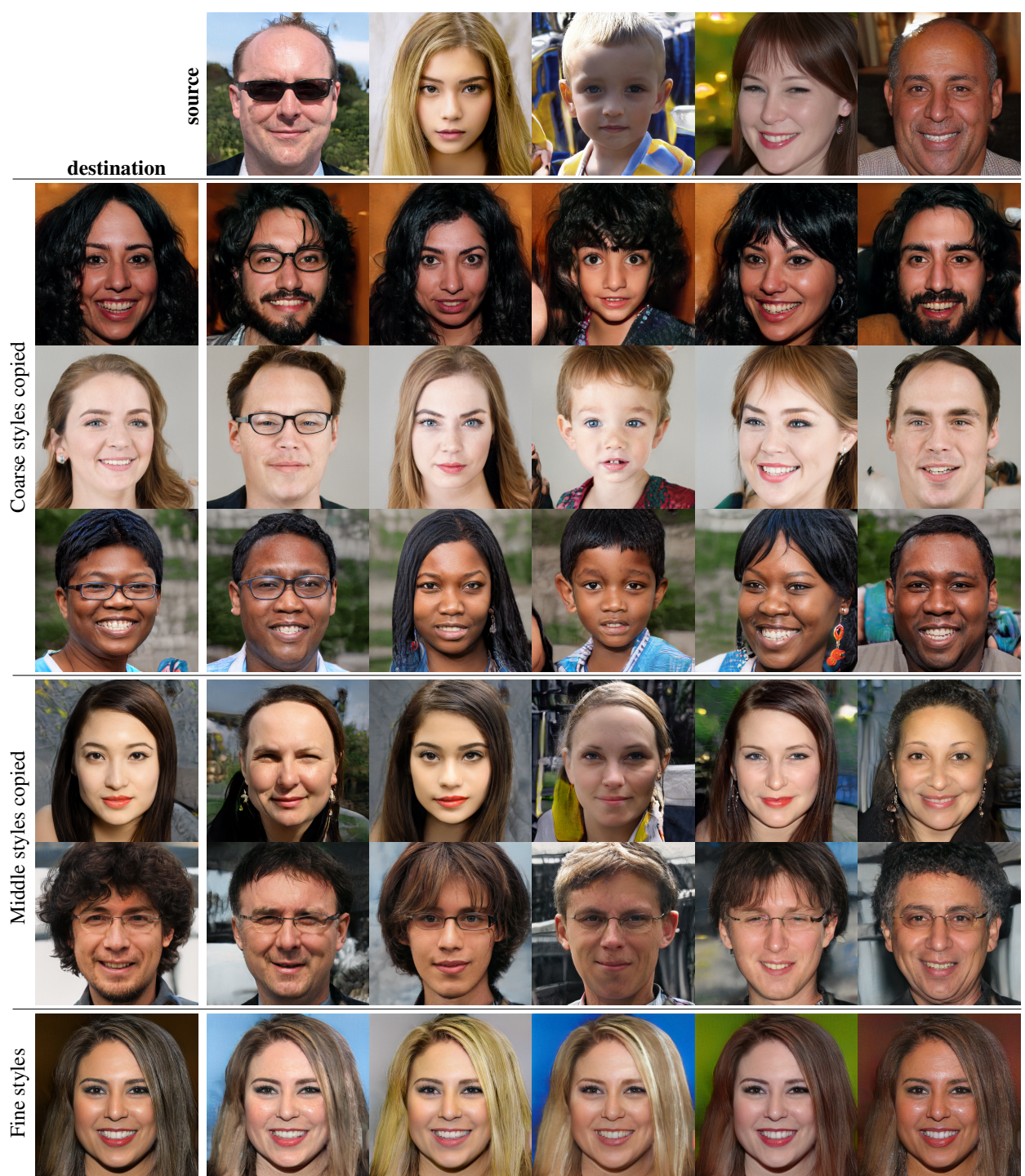
Resolution: 512x512

# StyleGANs (Karras et al., 2019)

- A new architecture motivated by the style transfer networks
- allows unsupervised separation of high-level attributes and stochastic variation in the generated images



# StyleGANs (Karras et al., 2018)





# Some Applications of GANs

# Semi-supervised Classification

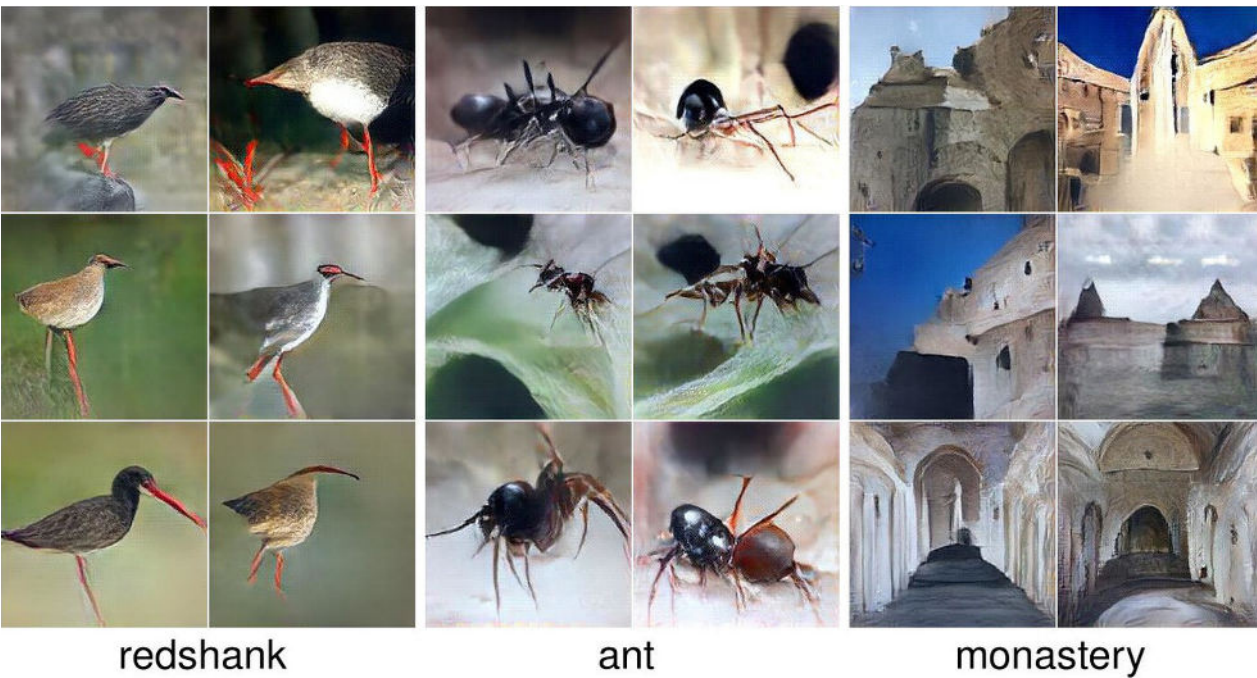
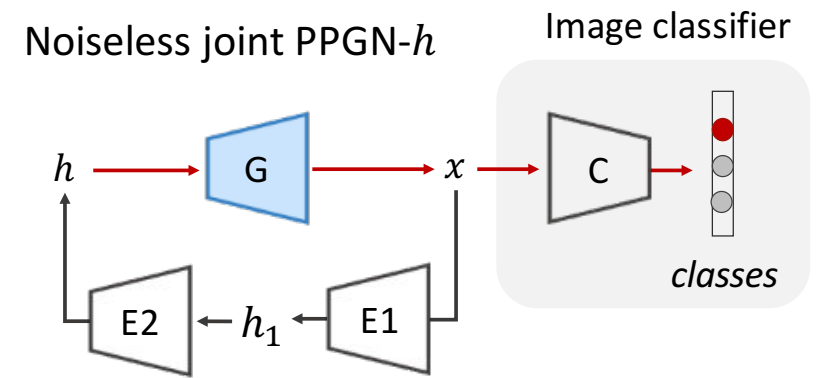
(Salimans et al., 2016;  
Dumoulin et al., 2016)

## SVNH

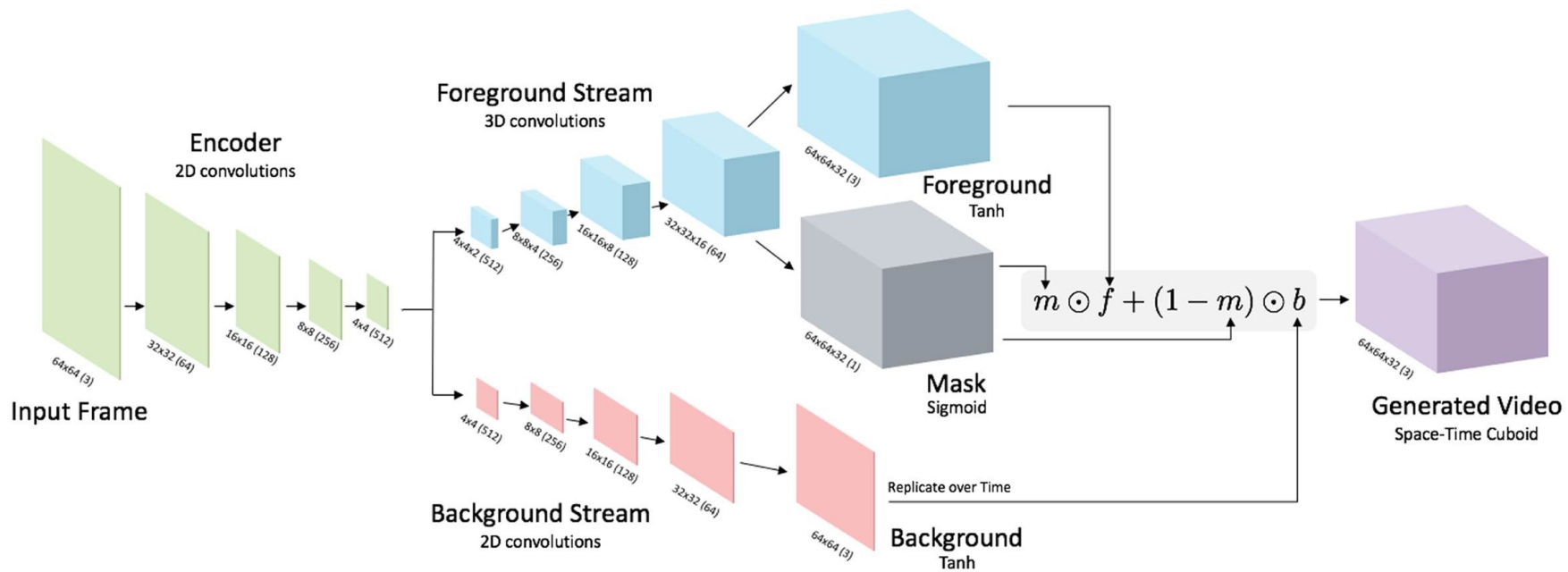
Model	Misclassification rate
VAE (M1 + M2) (Kingma et al., 2014)	36.02
SWWAE with dropout (Zhao et al., 2015)	23.56
DCGAN + L2-SVM (Radford et al., 2015)	22.18
SDGM (Maaløe et al., 2016)	16.61
<b>GAN (feature matching) (Salimans et al., 2016)</b>	<b>8.11 ± 1.3</b>
ALI (ours, L2-SVM)	19.14 ± 0.50
<b>ALI (ours, no feature matching)</b>	<b>7.42 ± 0.65</b>

# Class-specific Image Generation (Nguyen et al., 2016)

- Generates  $227 \times 227$  realistic images from all ImageNet classes
- Combines adversarial training, moment matching, denoising autoencoders, and Langevin sampling



# Video Generation (Vondrick et al., 2016)



Beach



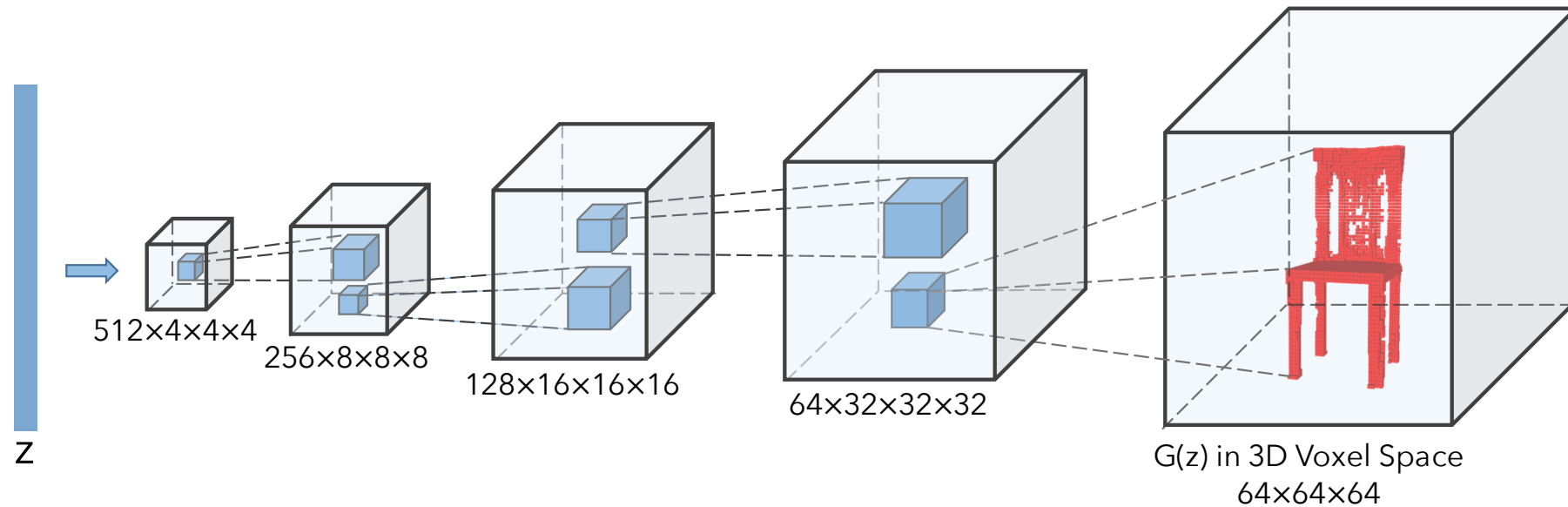
Golf



Train Station



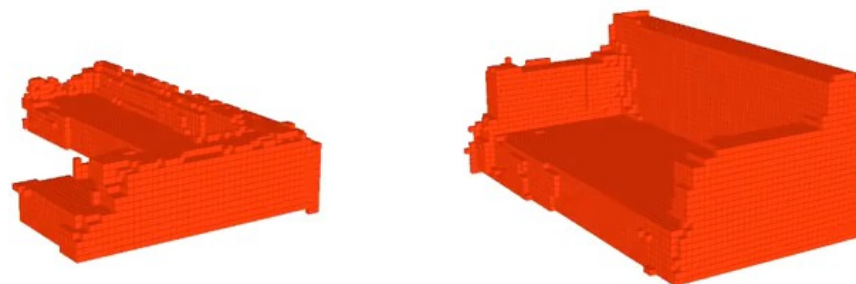
# Generative Shape Modeling (Wu et al., 2016)



Chairs



Sofas



# Text-to-Image Synthesis (Zhang et al., 2016)

The small bird has a red head with feathers that fade from red to gray from head to tail



The petals of this flower are white with a large stigma



A unique yellow flower with no visible pistils protruding from the center



This flower is pink and yellow in color, with petals that are oddly shaped



This is a light colored flower with many different petals on a green stem



This flower is yellow and green in color, with petals that are ruffled



The flower have large petals that are pink with yellow on some of the petals



A flower that has white petals with some tones of yellow and green filaments



# Text-to-Image Synthesis (Zhu et al., 2019)

This bird has a white throat and a dark yellow bill and grey wings.



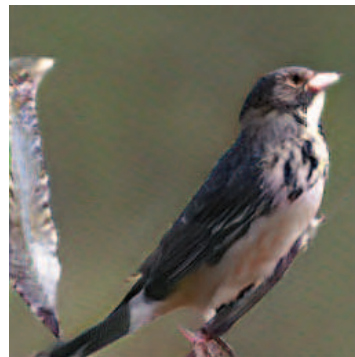
This bird has wings that are grey and has a white belly.



This particular bird has a belly that is yellow and brown.



This bird has wings that are black and has a white belly.



This bird is a lime green with greyish wings and long legs.



This is a grey bird with a brown wing and a small orange beak.



This yellow bird has a thin beak and jet black eyes and thin feet.



This bird has a short brown bill, a white eyering, and a medium brown crown.



# Single Image Super-Resolution (Ledig et al., 2016)

- Combine content loss with adversarial loss

bicubic



SRResNet



SRGAN

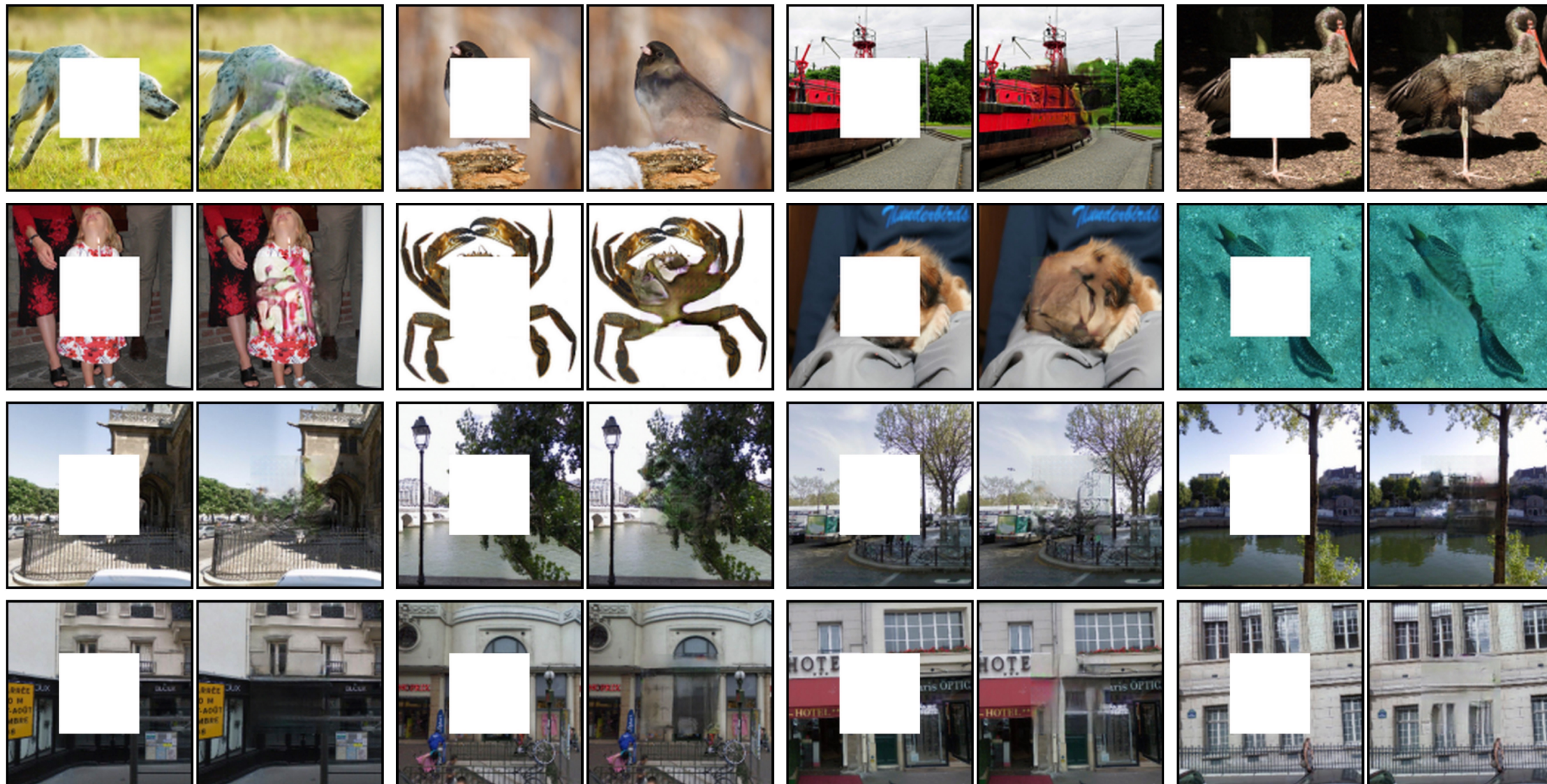


original

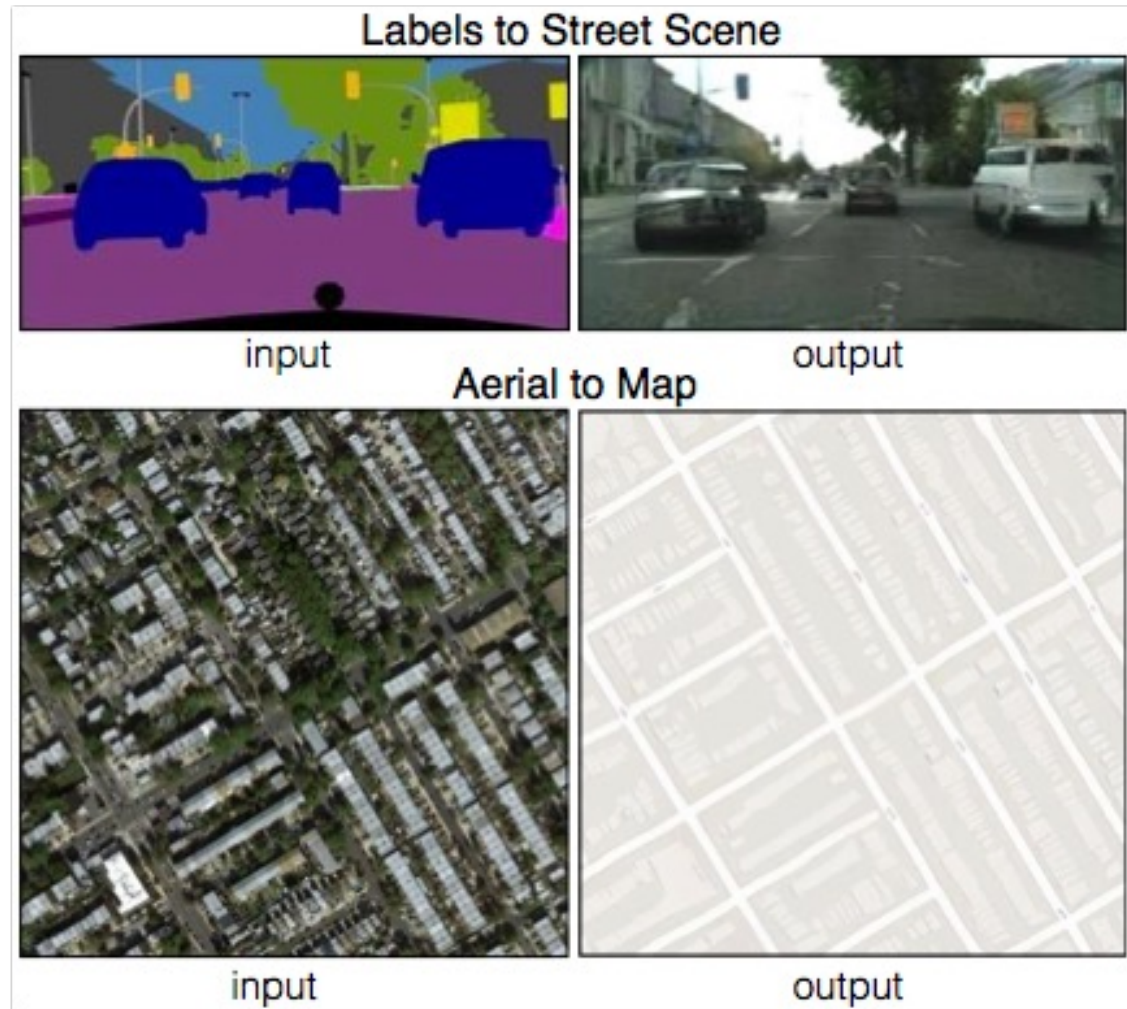




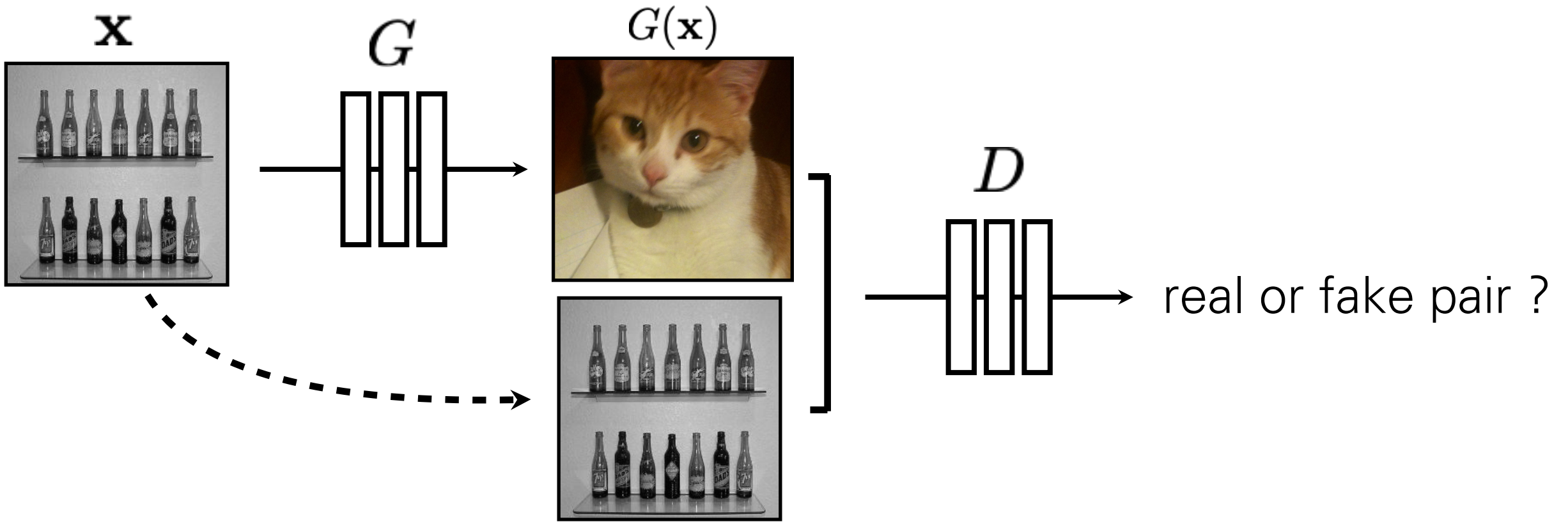
# Image Inpainting (Pathak et al., 2016)



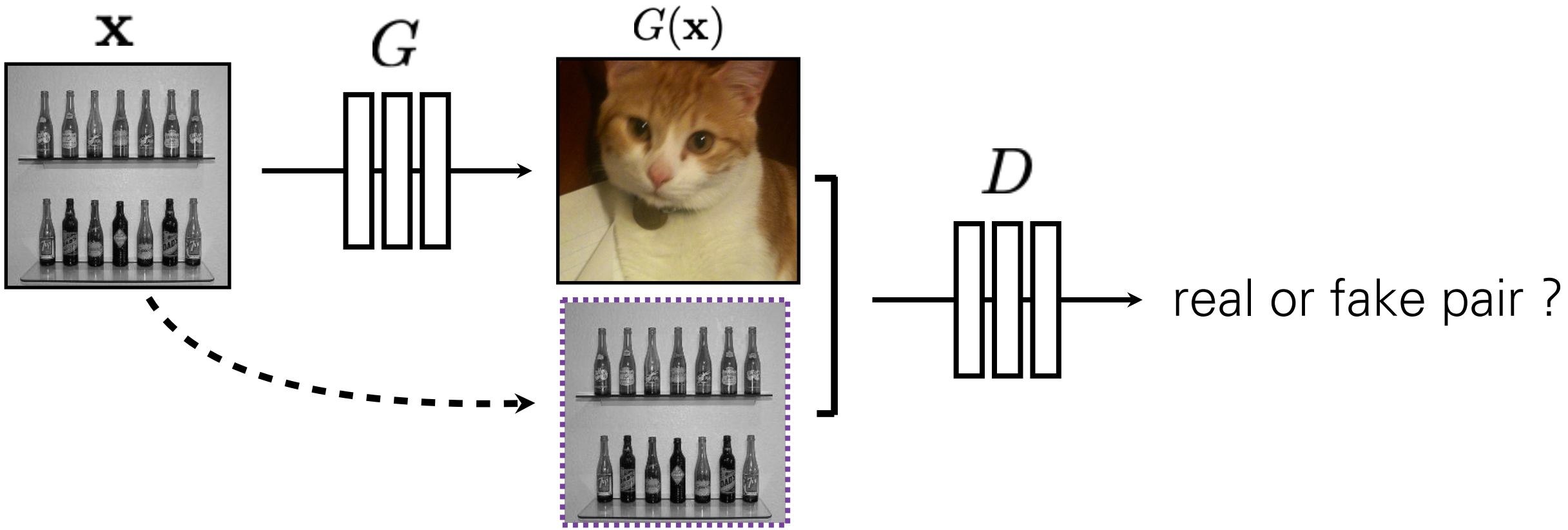
# Image to Image Translation (Pix2Pix)



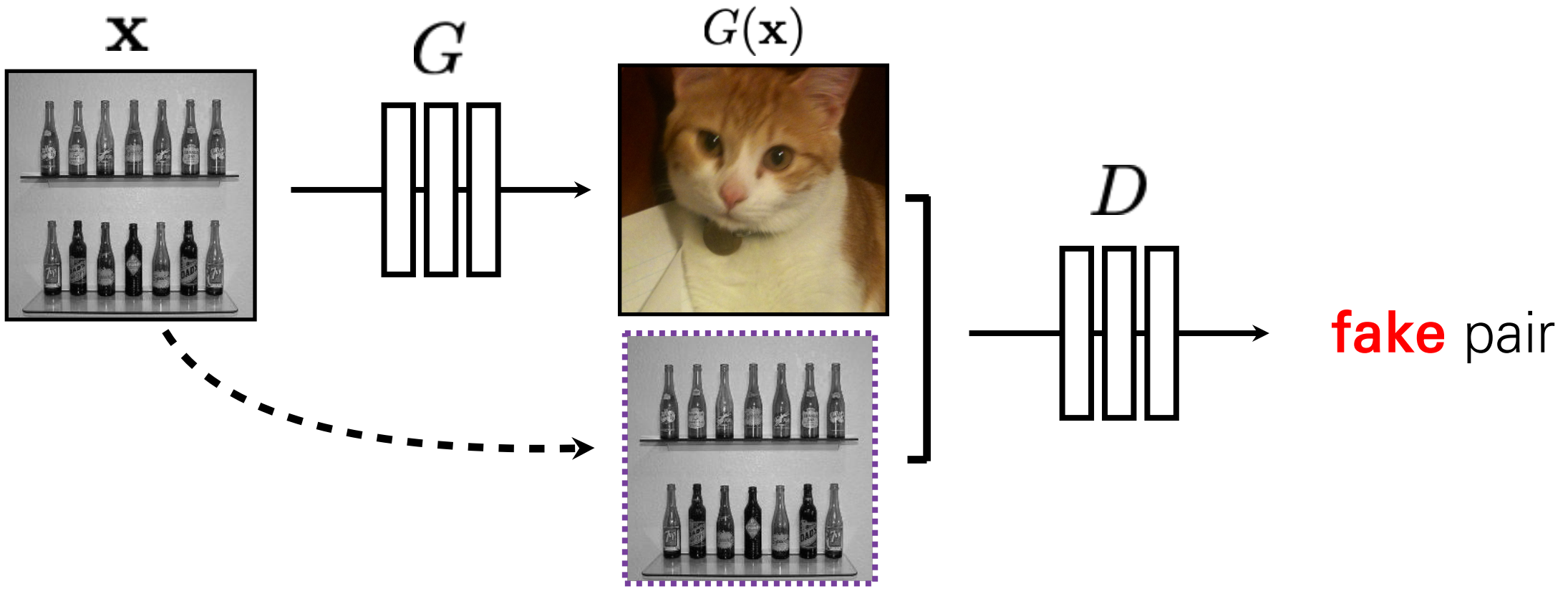
(Isola et al. 2016)



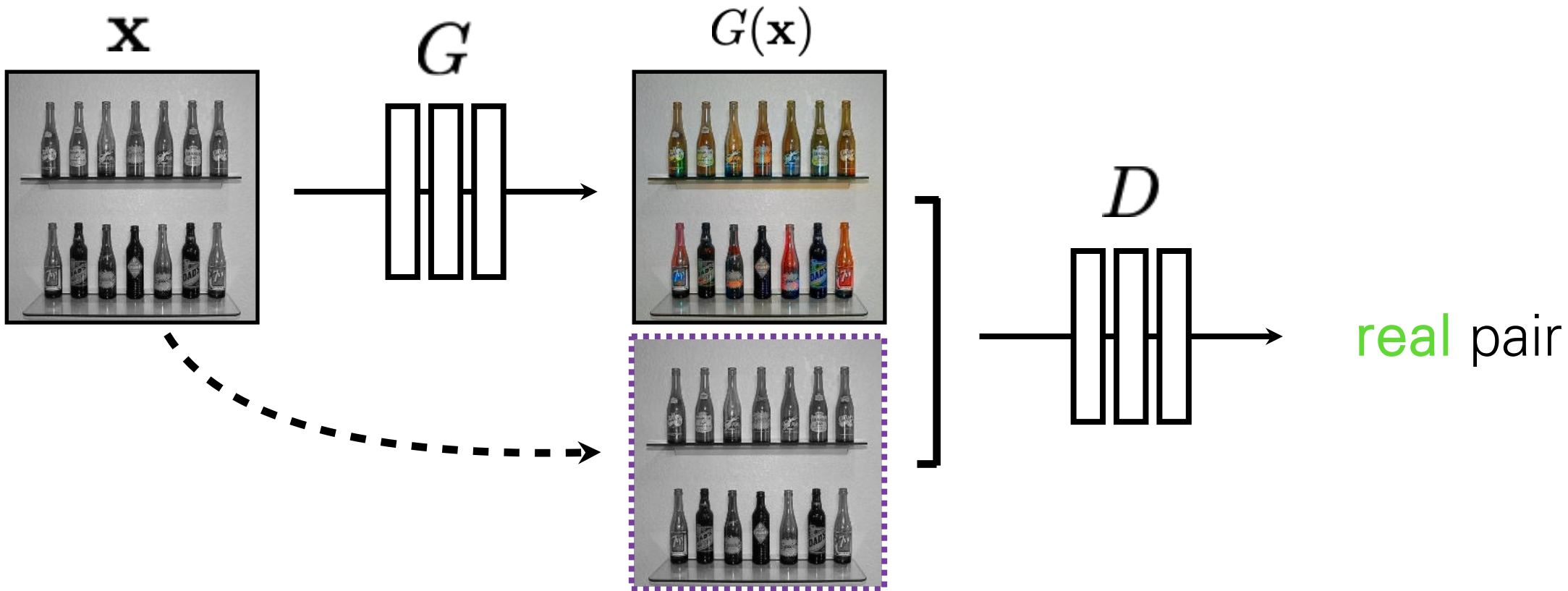
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



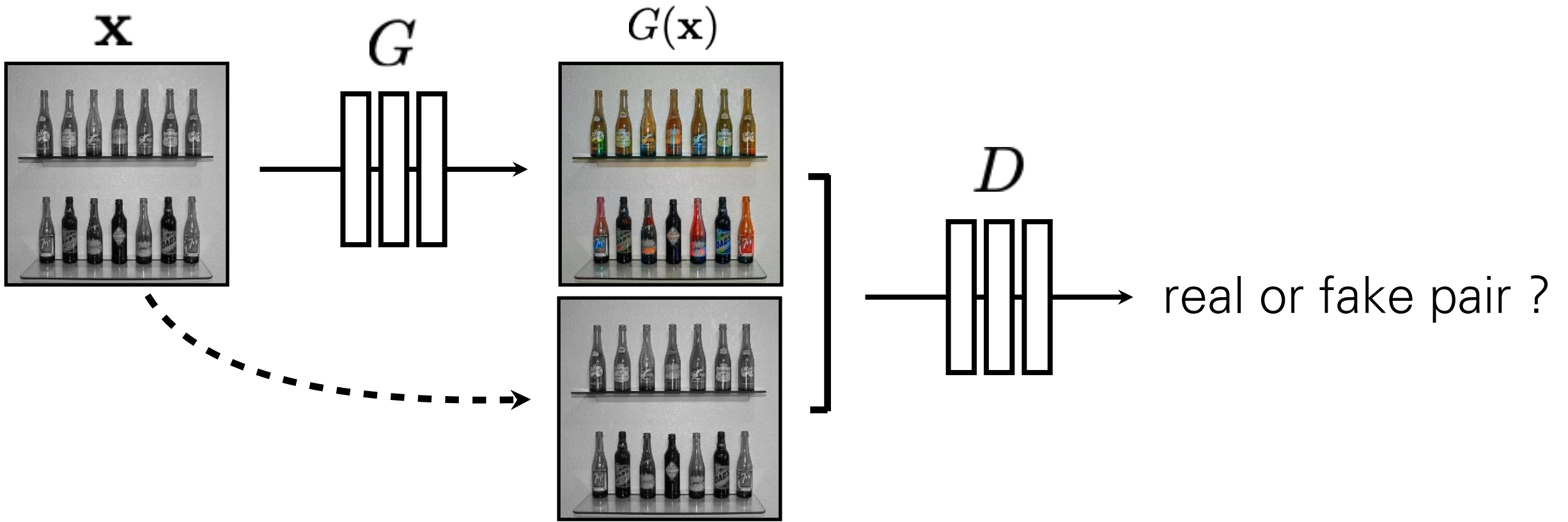
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$



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



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



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

# BW → Color

Input

Output



Input

Output



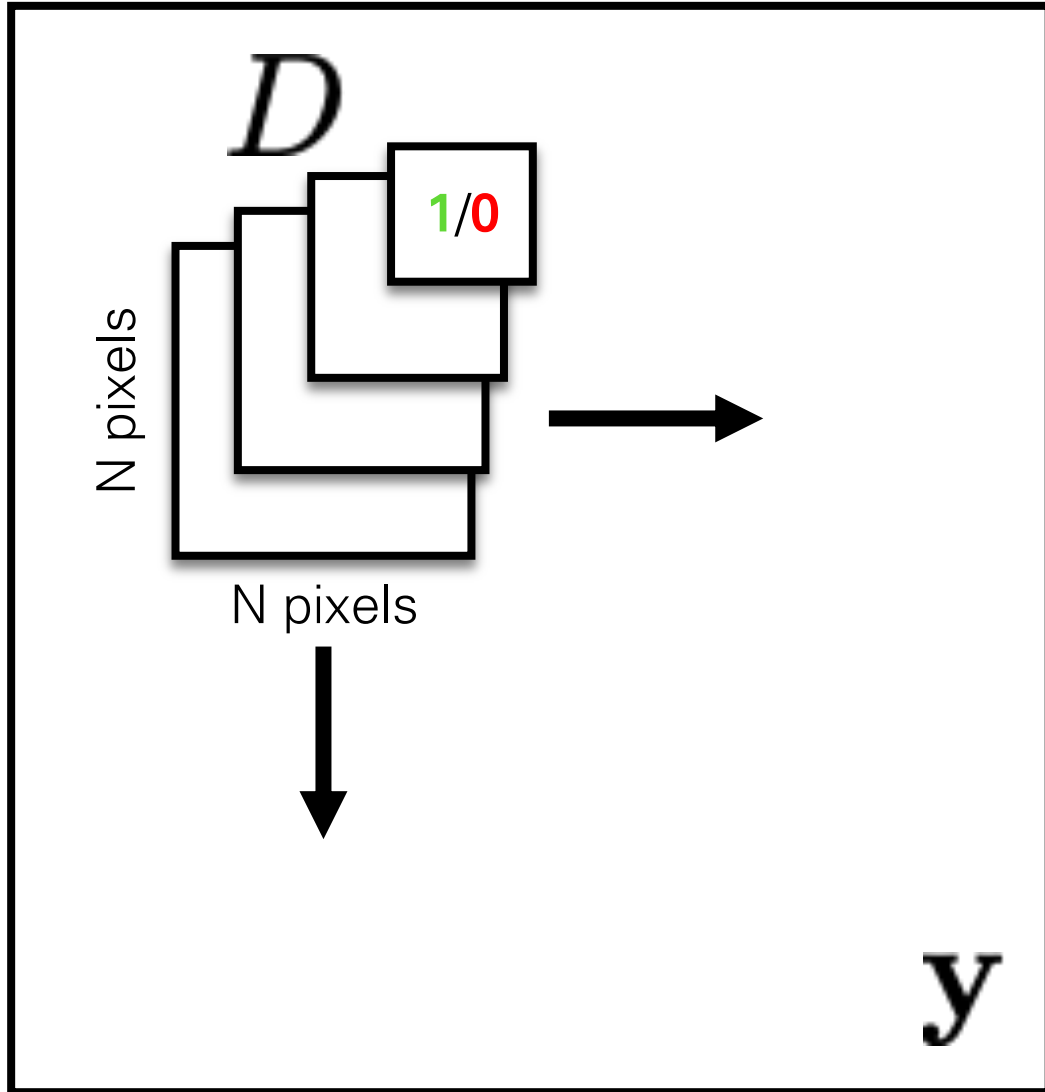
Input

Output





# Shrinking the capacity: Patch Discriminator



Rather than penalizing if output image looks fake, penalize if each overlapping patch in output looks fake

- Faster, fewer parameters
- More supervised observations
- Applies to arbitrarily large images

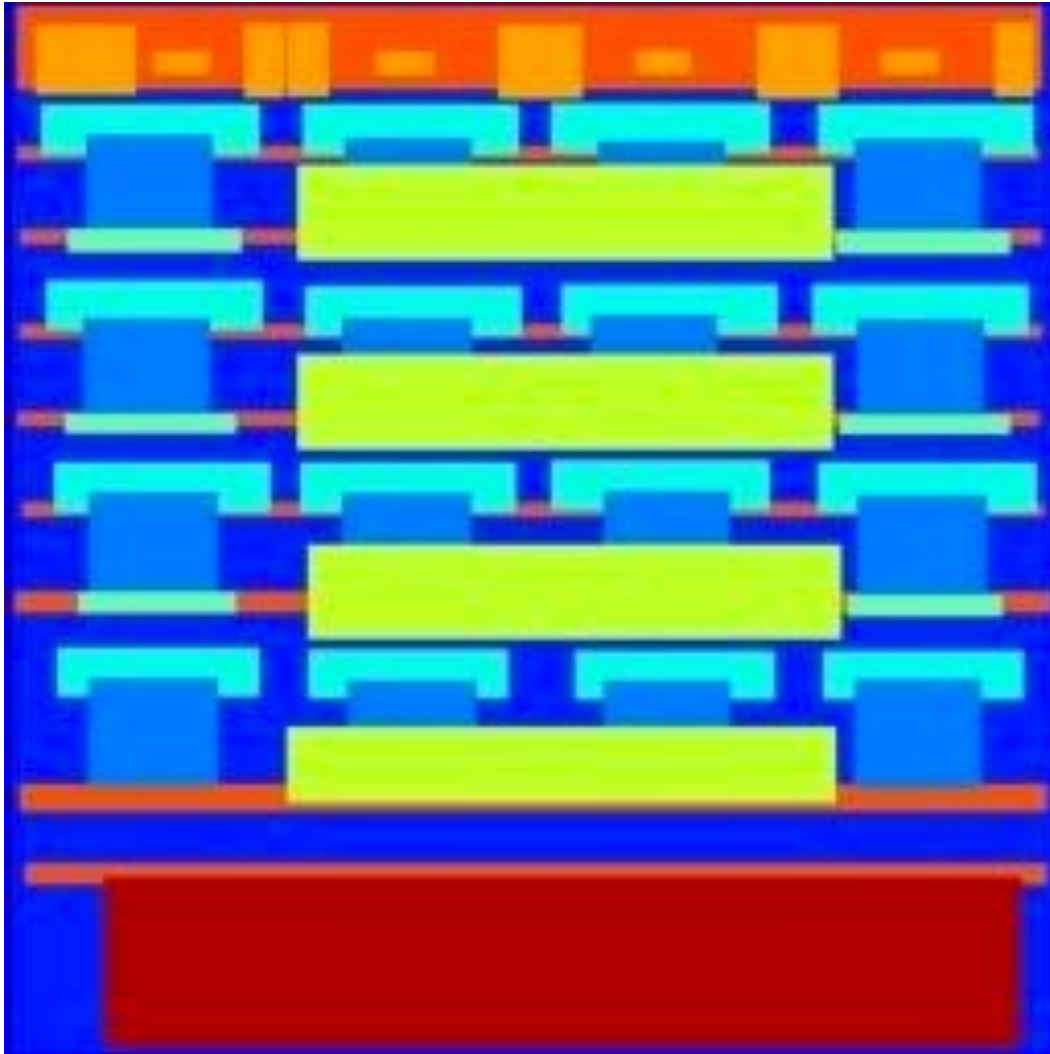
[Li & Wand 2016]

[Shrivastava et al. 2017]

[Isola et al. 2017]

# Labels $\rightarrow$ Facades

Input

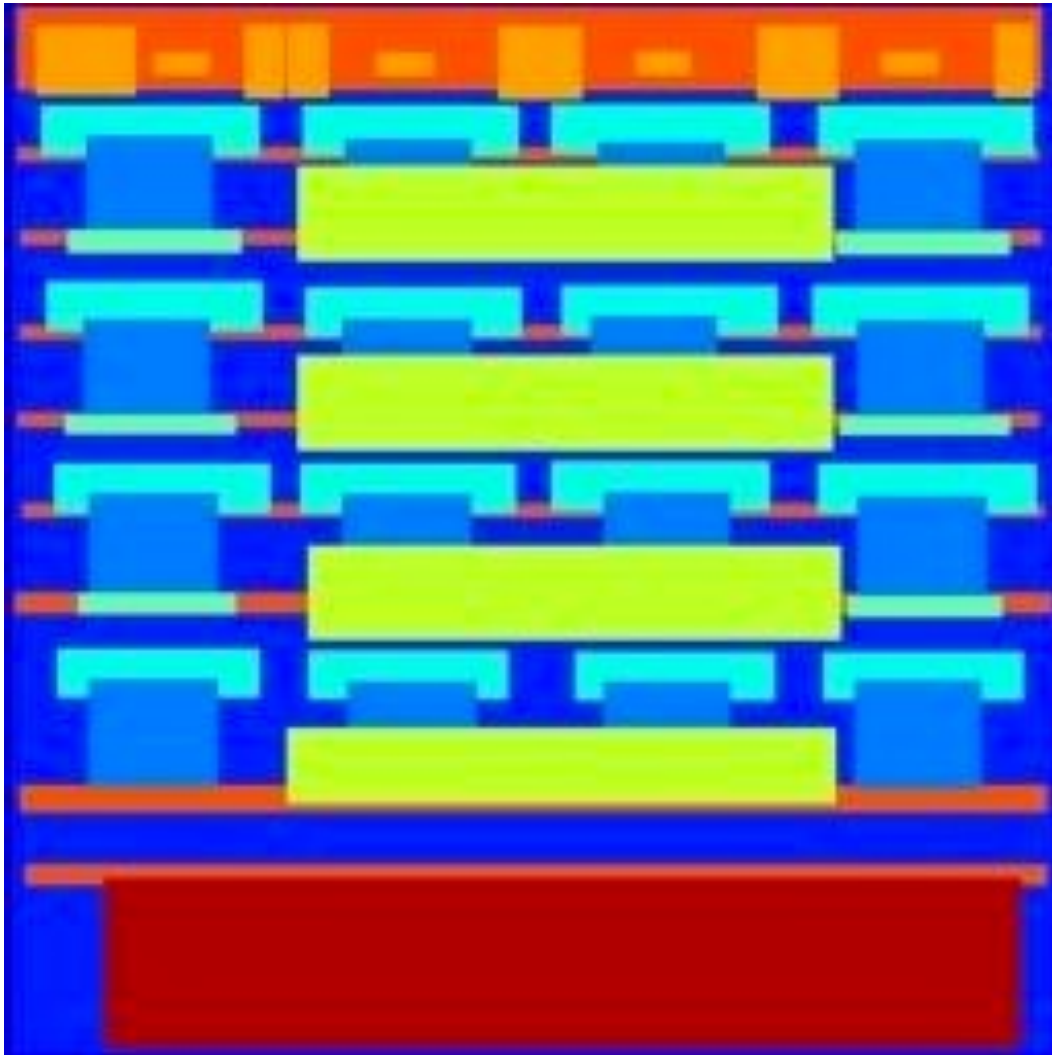


1x1 Discriminator



# Labels → Facades

Input

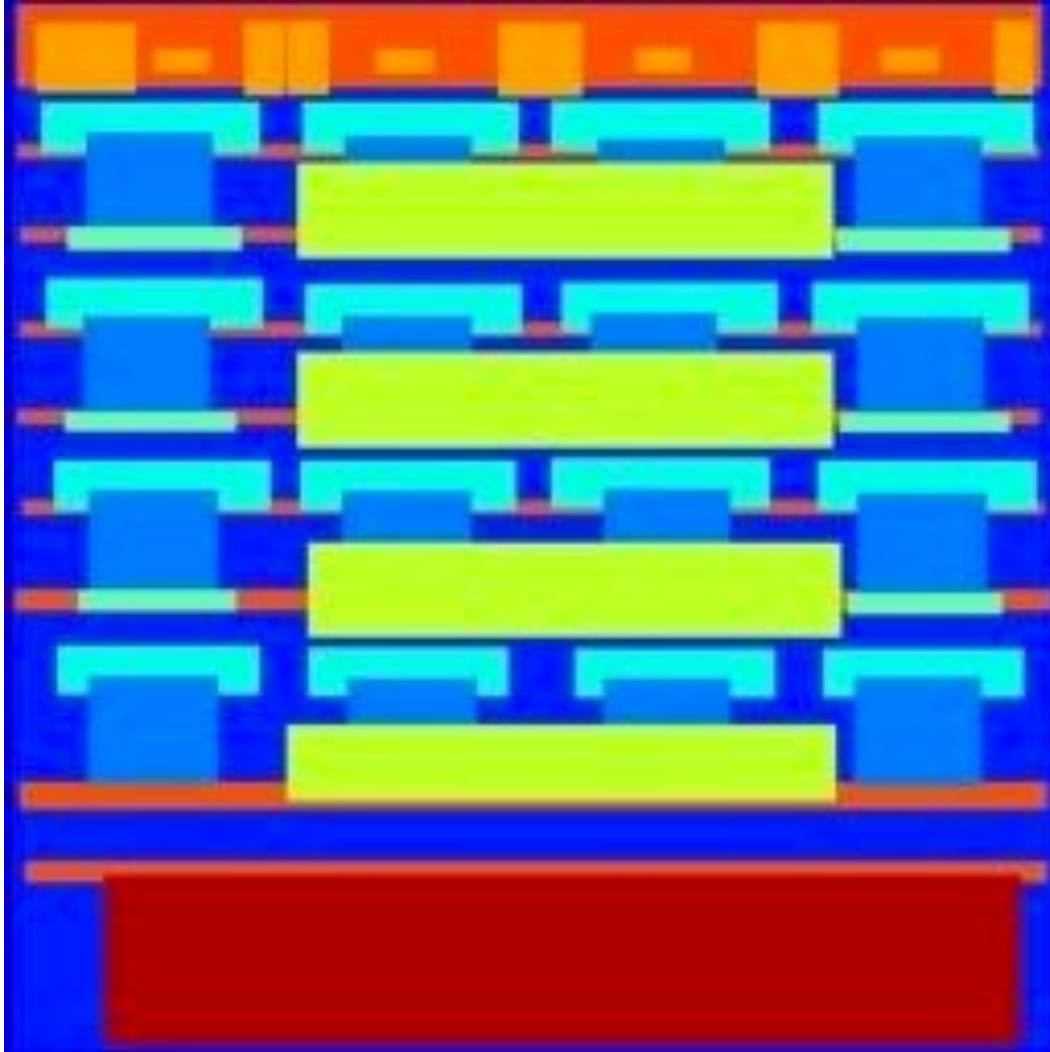


16x16 Discriminator



# Labels → Facades

Input

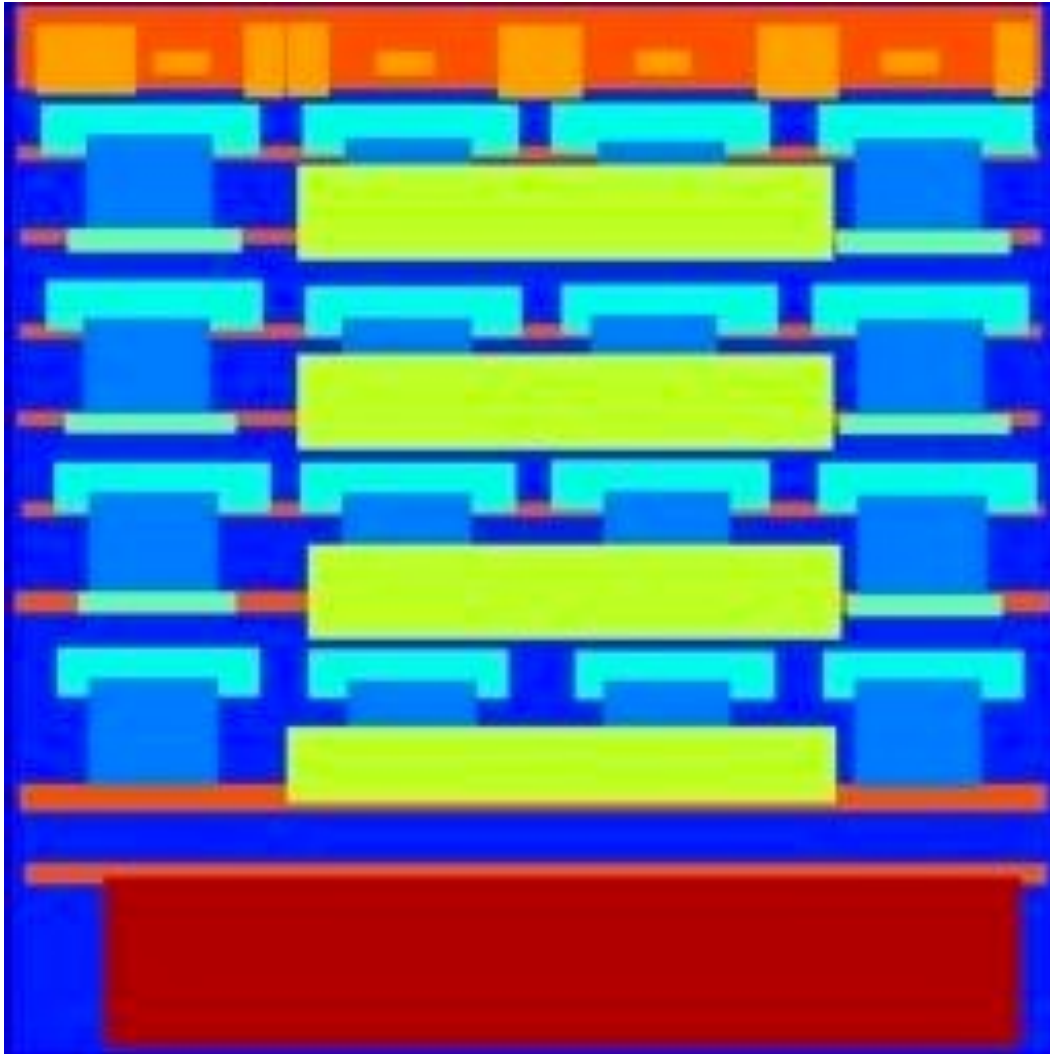


70x70 Discriminator



# Labels → Facades

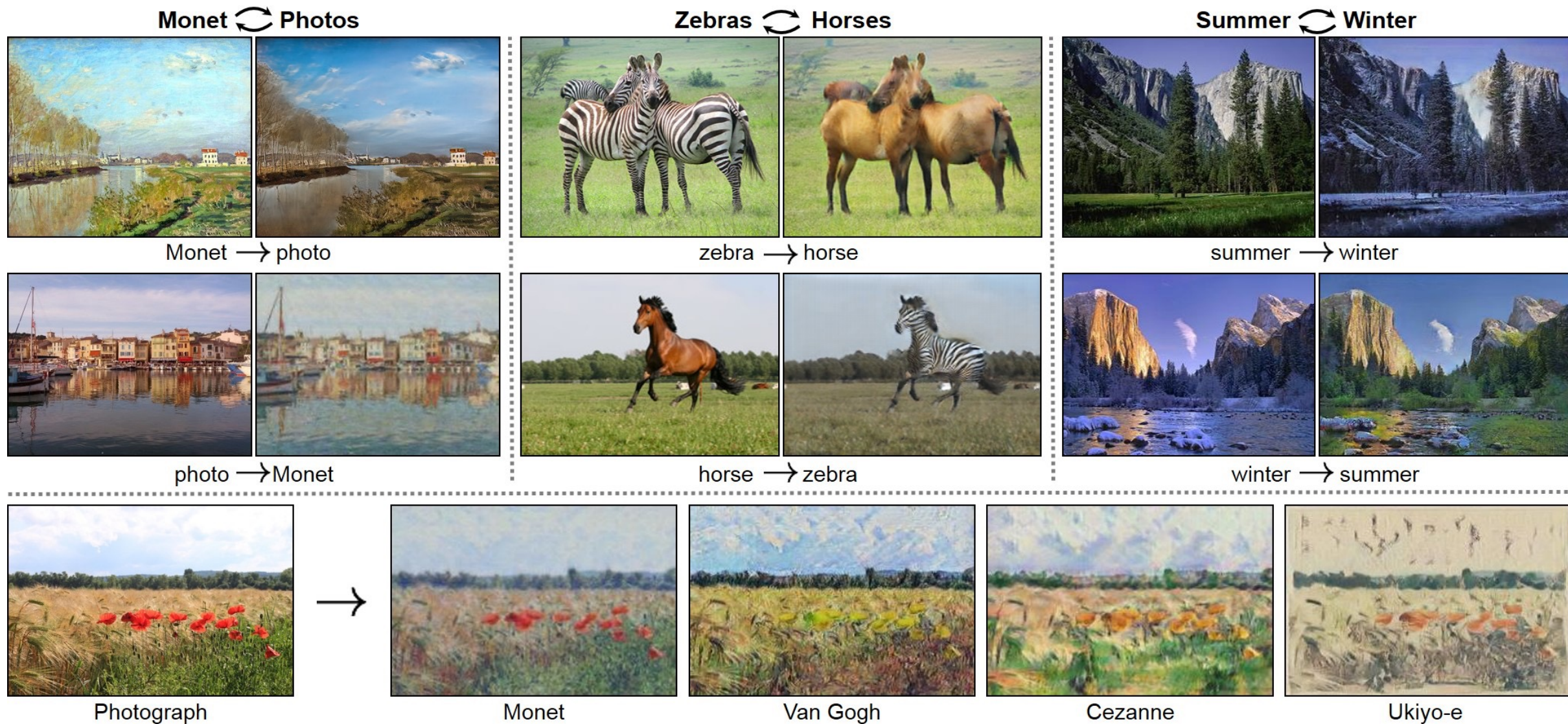
Input



Full image Discriminator



# CycleGAN: Pix2Pix w/o input-output pairs

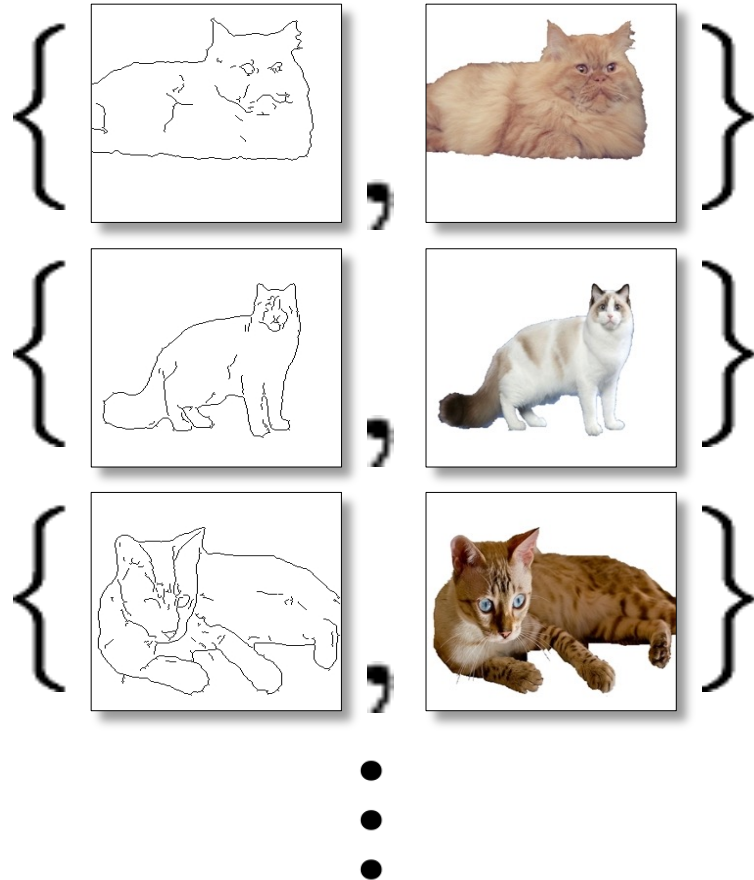


(Zhu et al. 2017)

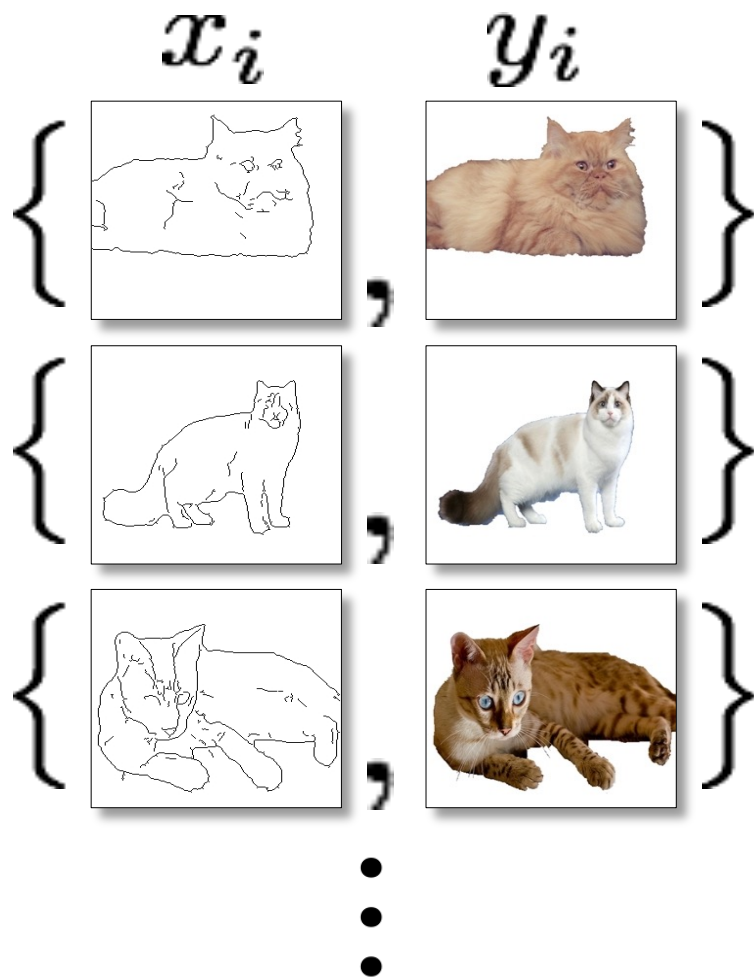
# Paired data

$x_i$

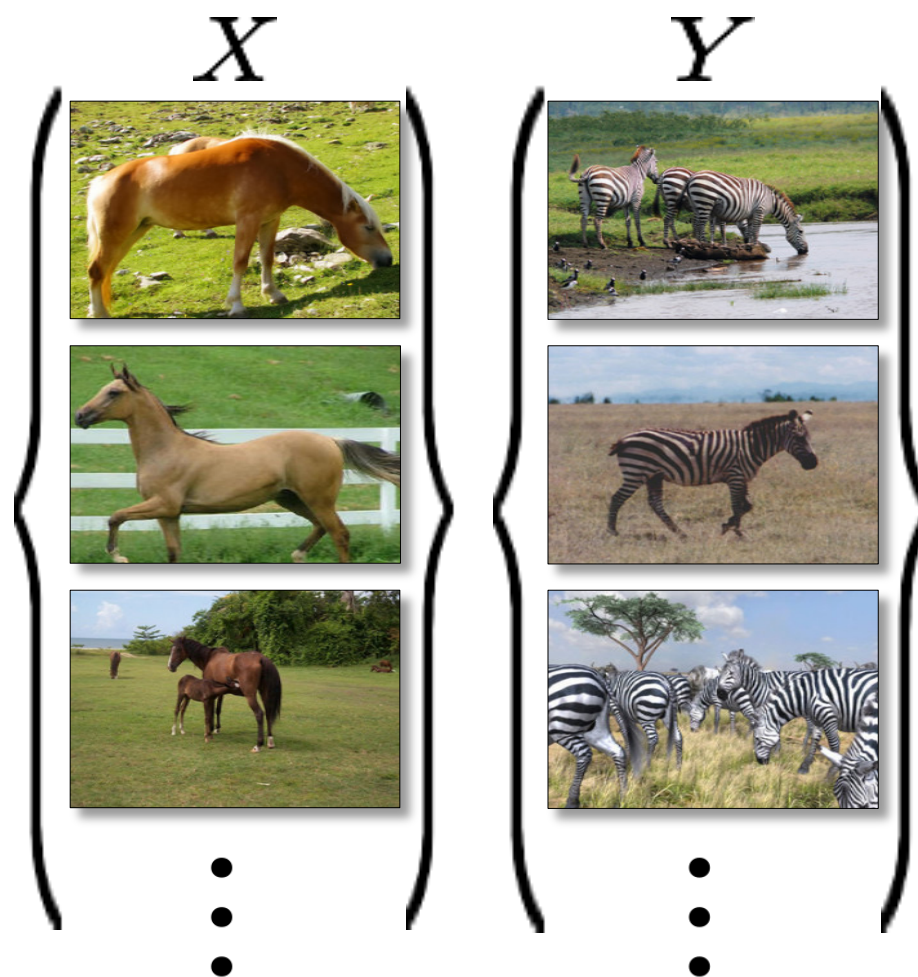
$y_i$



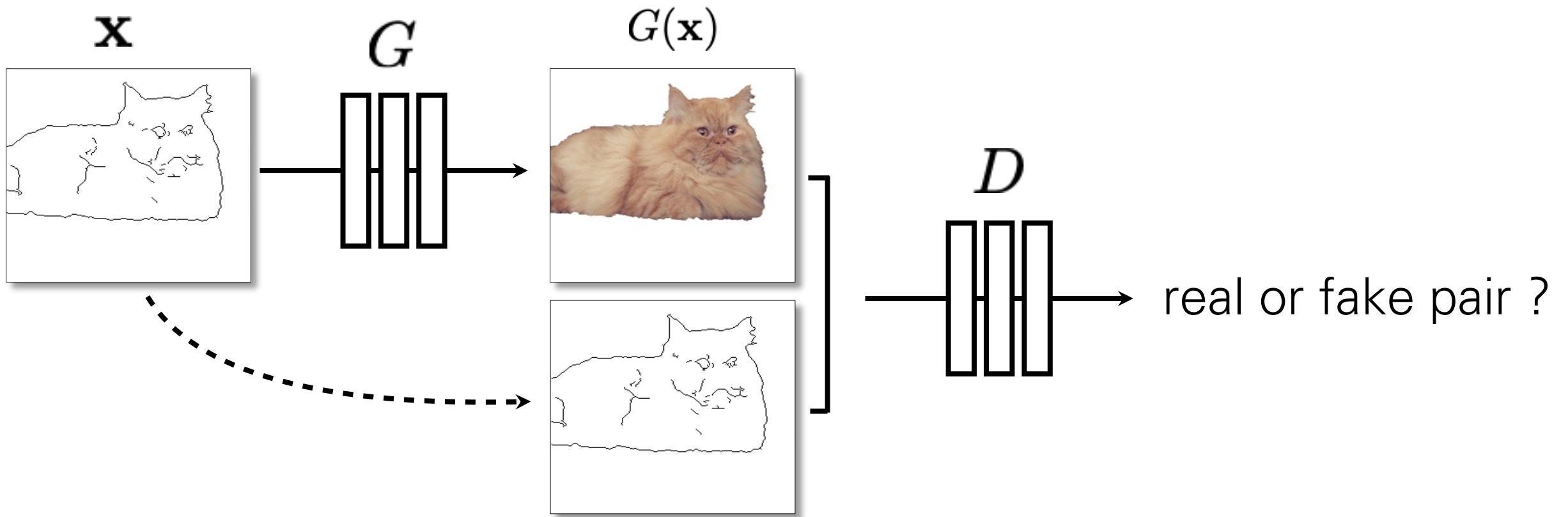
# Paired data



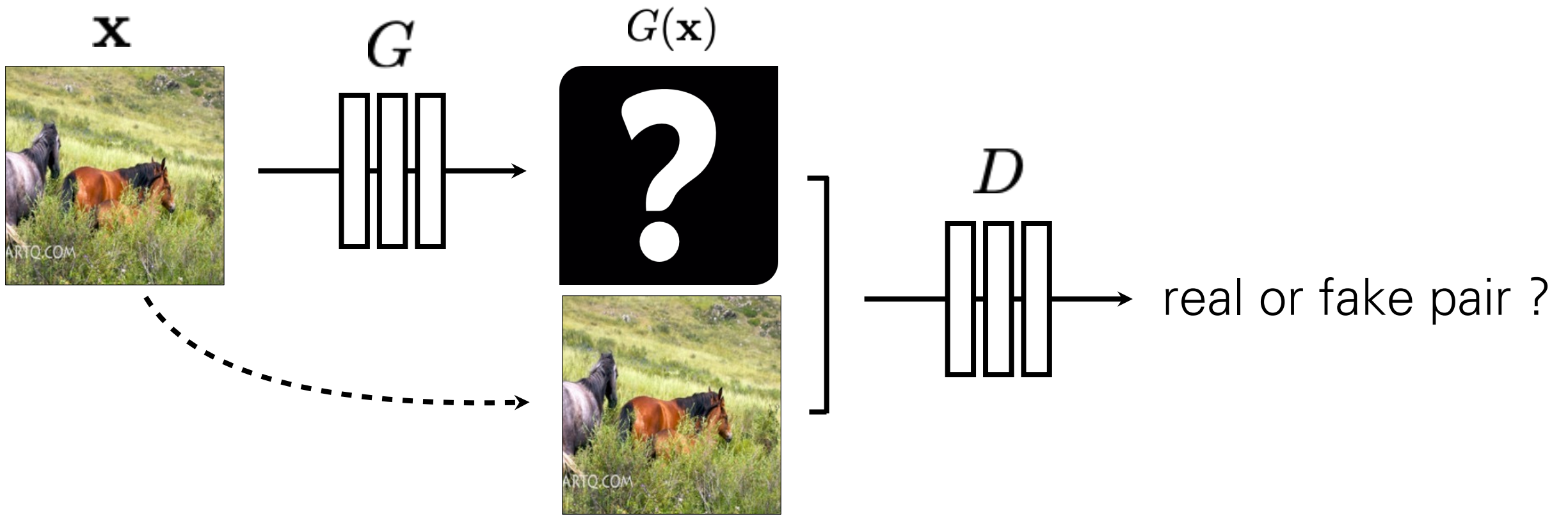
# Unpaired data





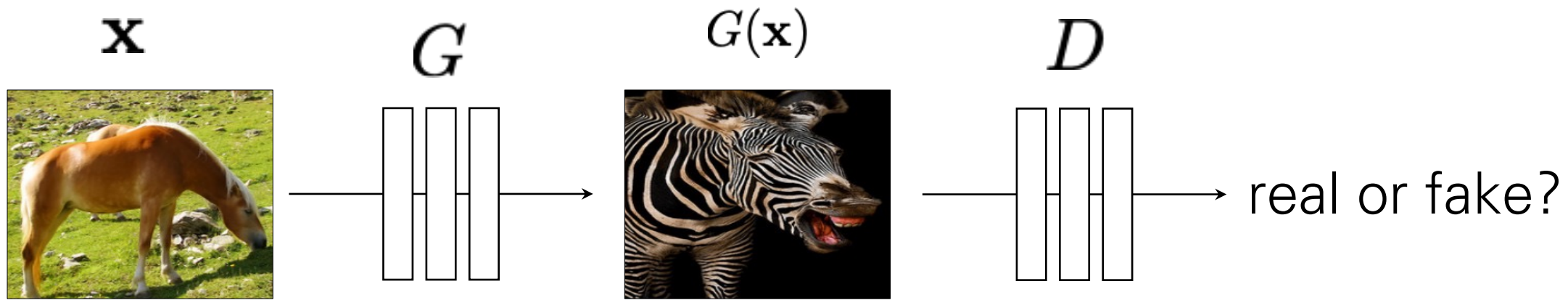


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$



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

No input-output pairs!

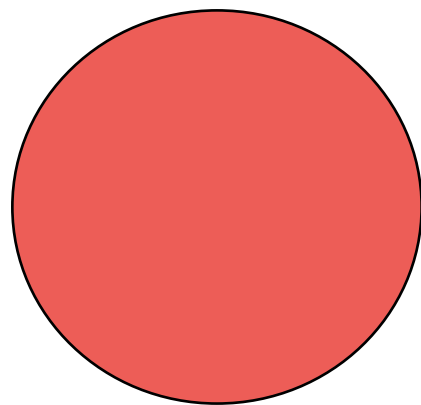


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

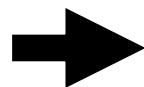
Usually loss functions check if output matches a target instance

GAN loss checks if output is part of an admissible set

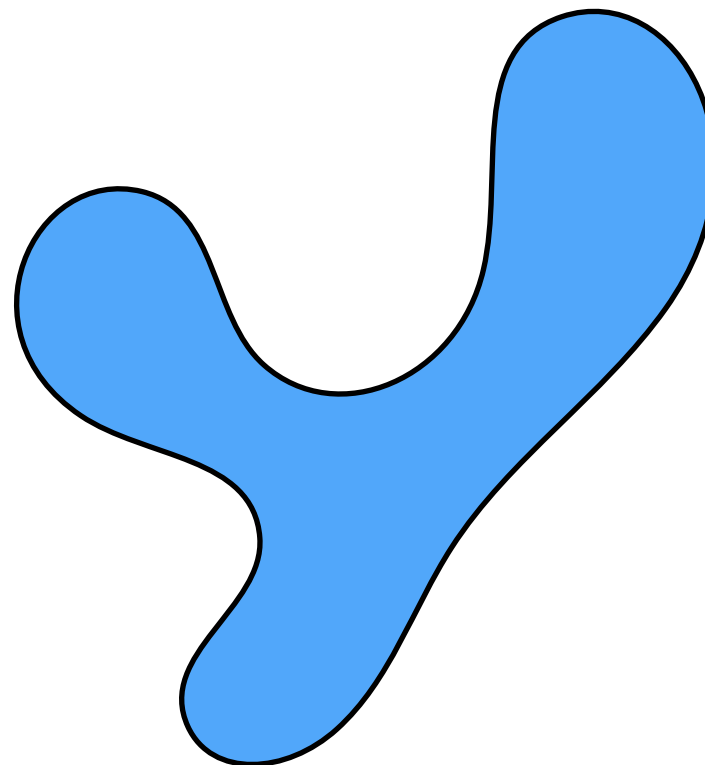
Gaussian



**Z**

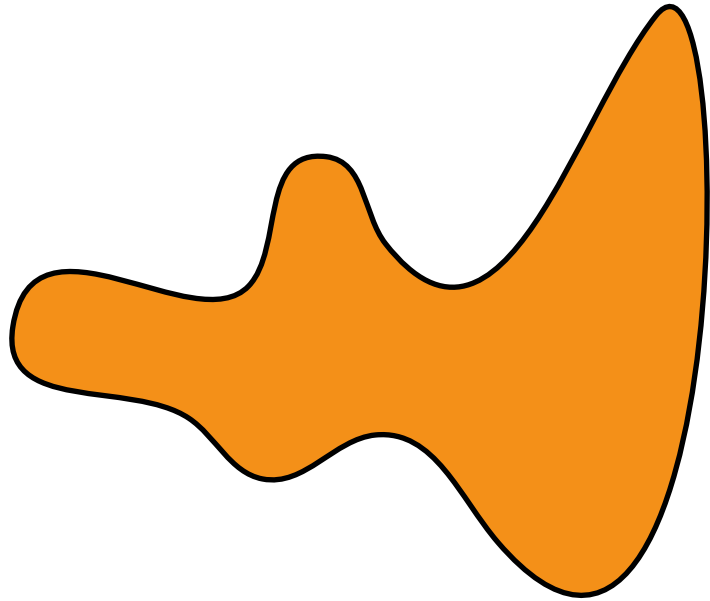


Target distribution

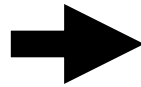


**Y**

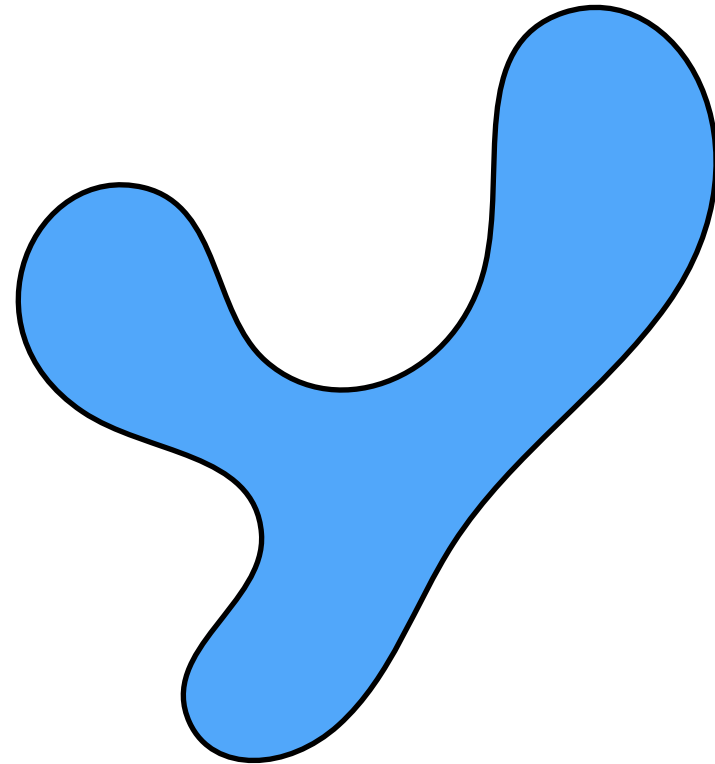
Horses



**X**



Zebras

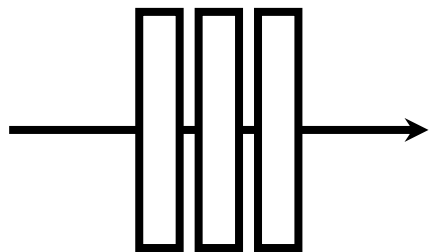


**Y**

$\mathbf{x}$



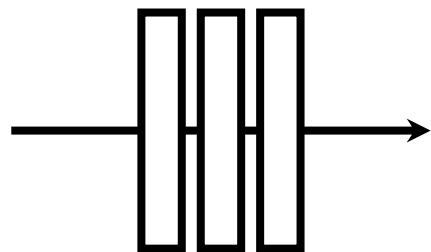
$G$



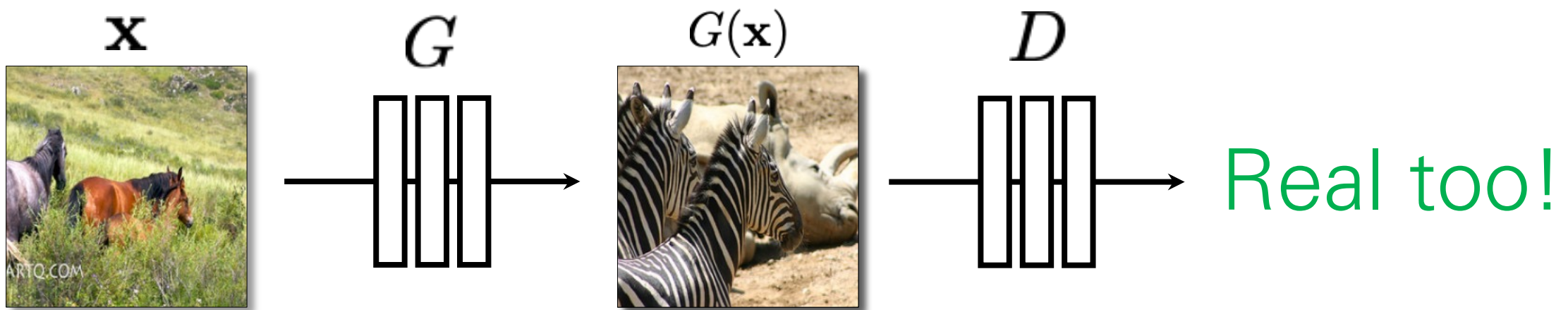
$G(\mathbf{x})$



$D$

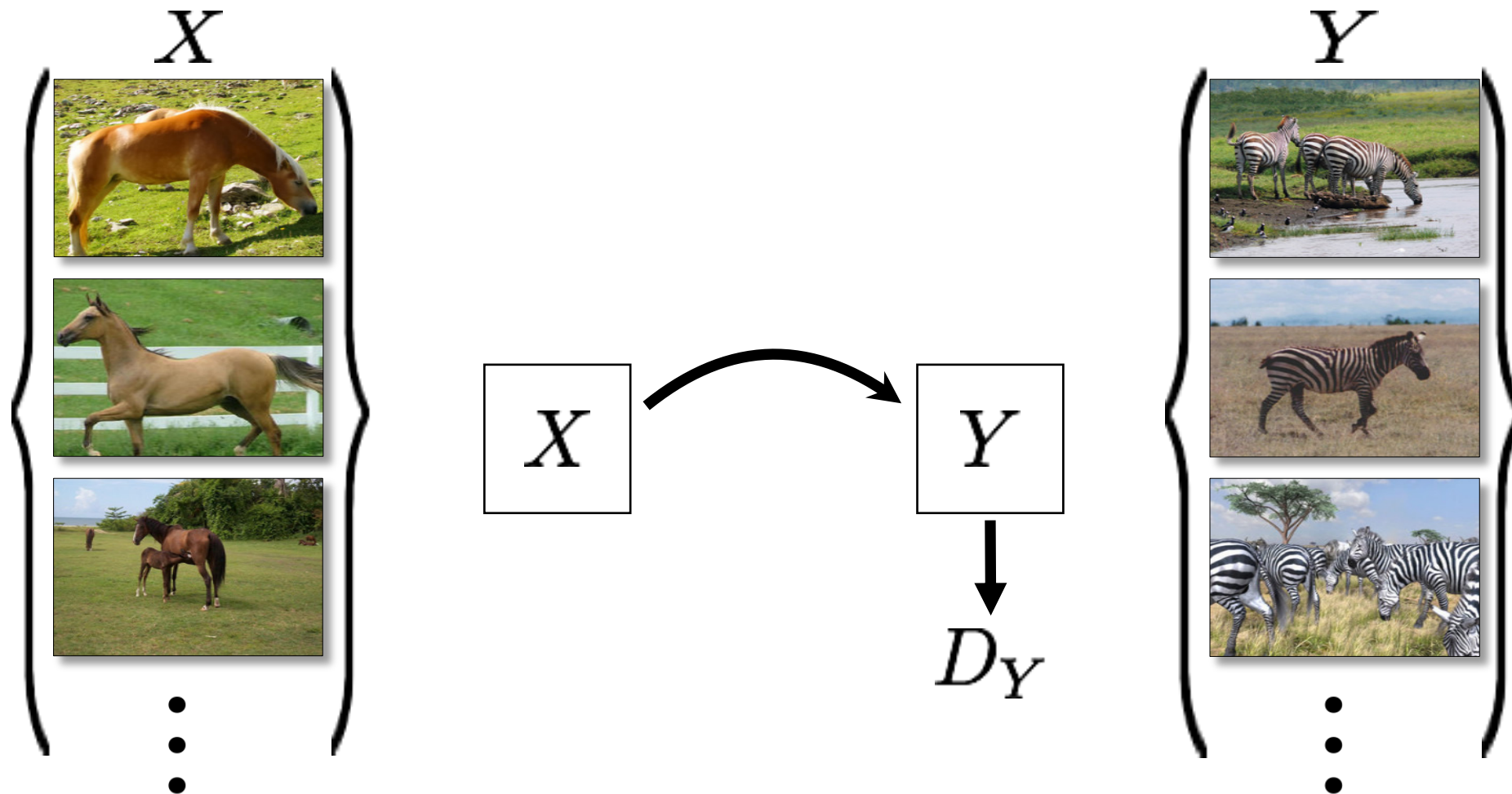


Real!



Nothing to force output to correspond to input

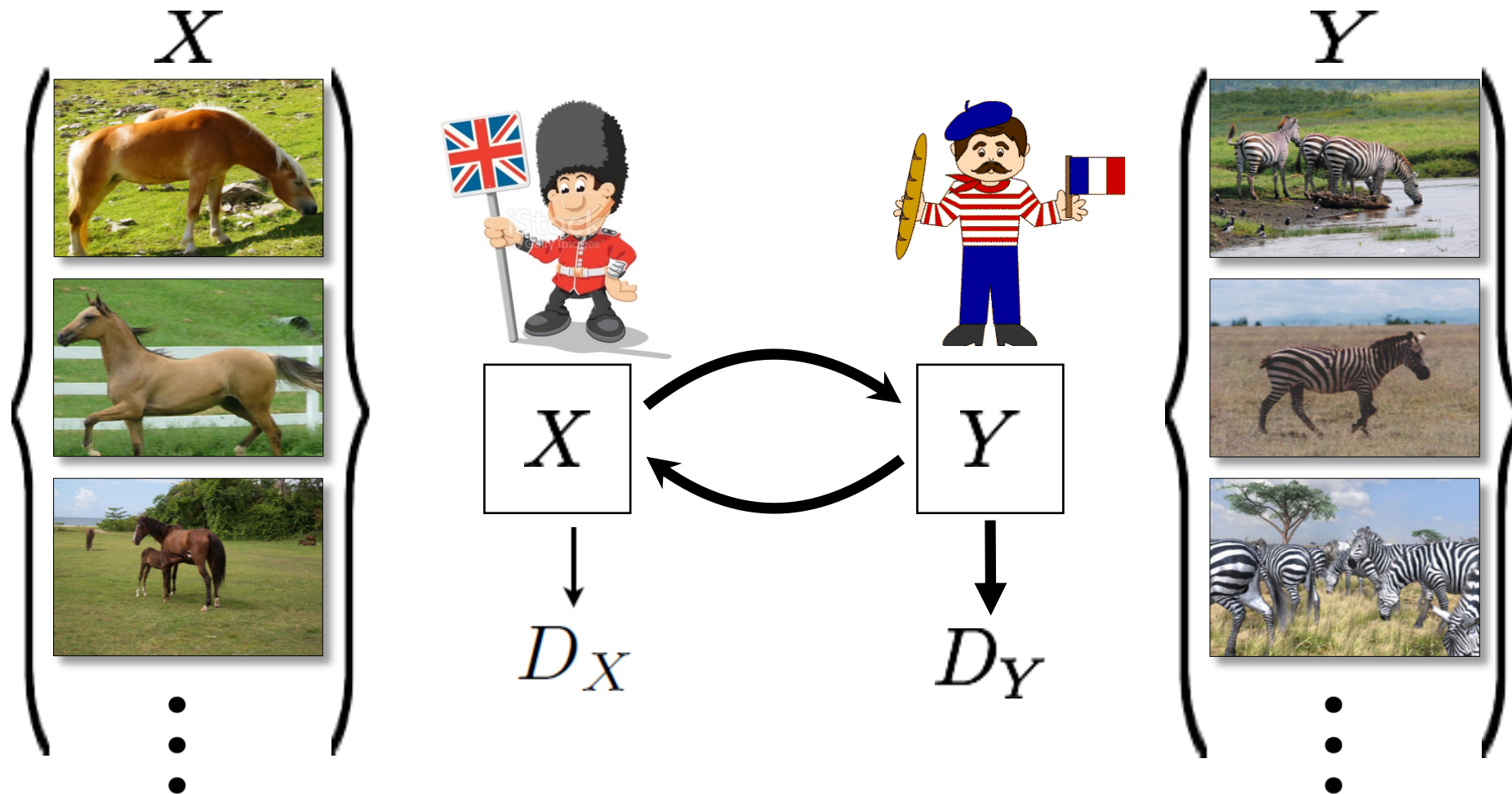
# Cycle-Consistent Adversarial Networks



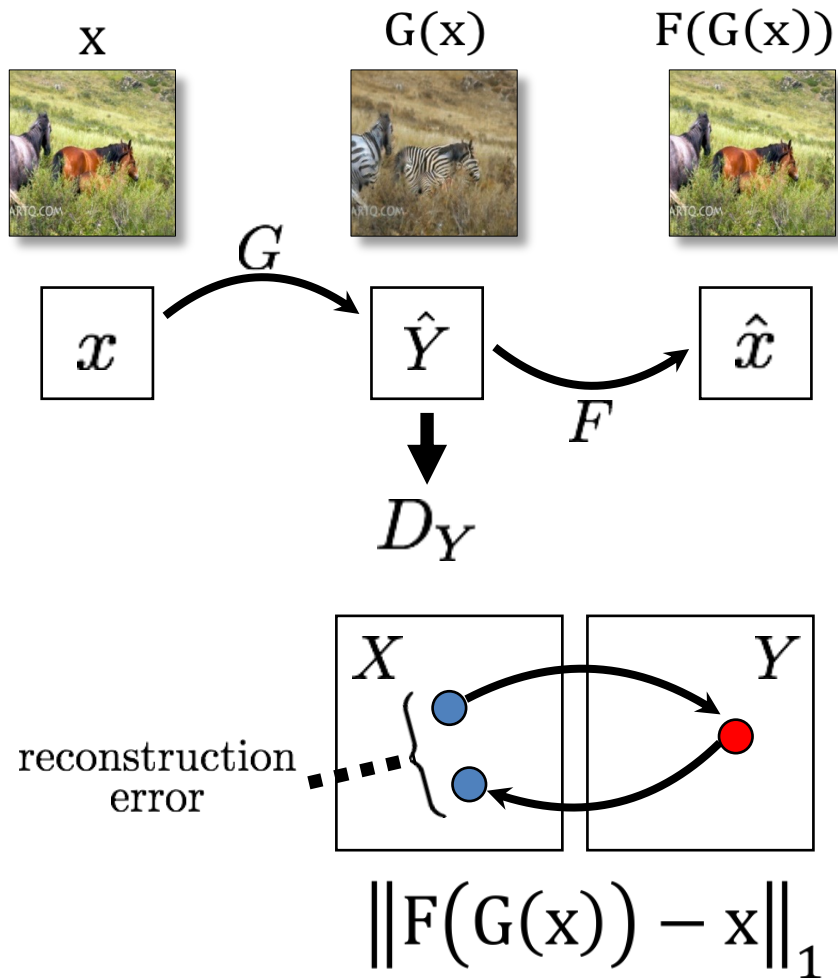
[Zhu et al. 2017], [Yi et al. 2017], [Kim et al. 2017]



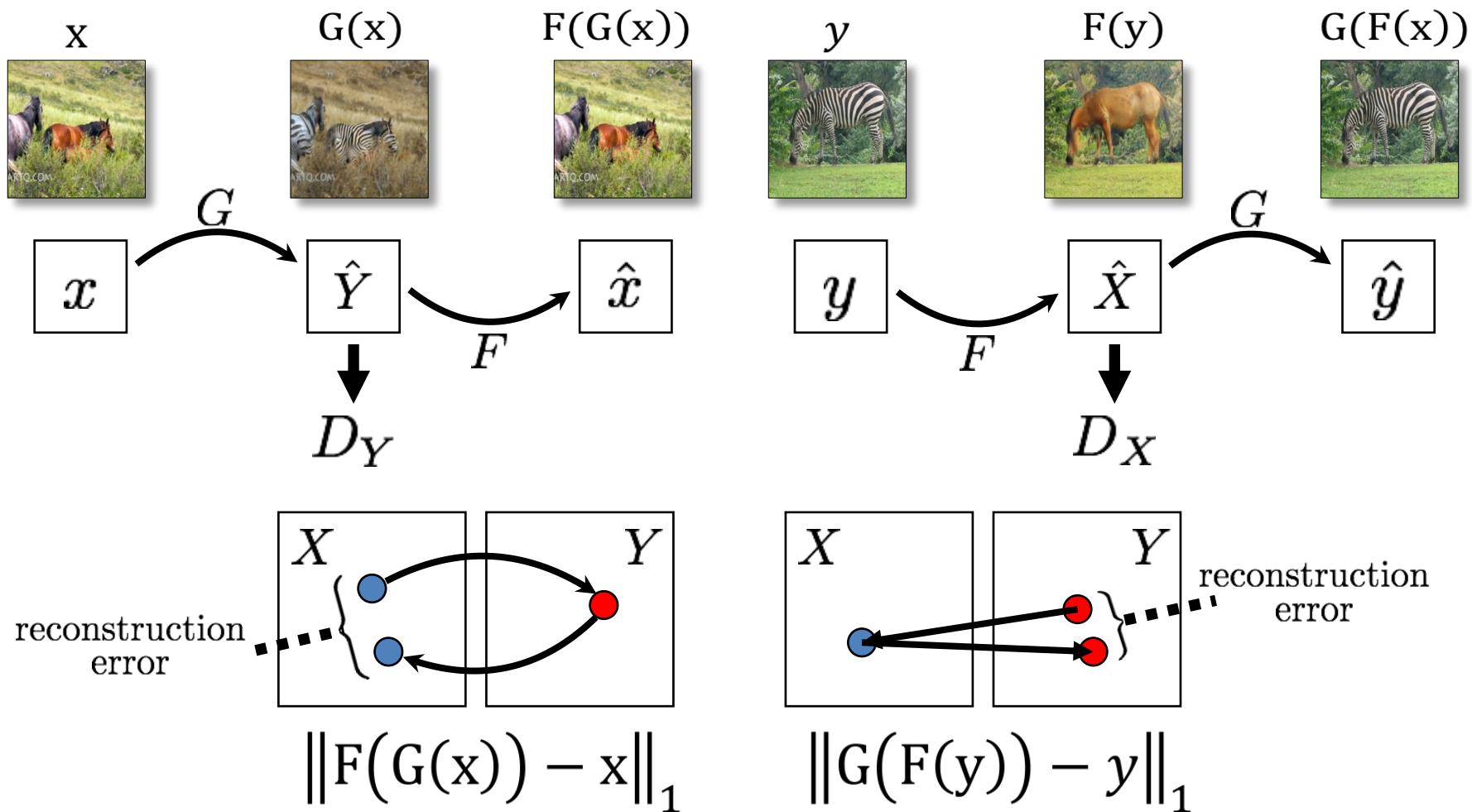
# Cycle-Consistent Adversarial Networks

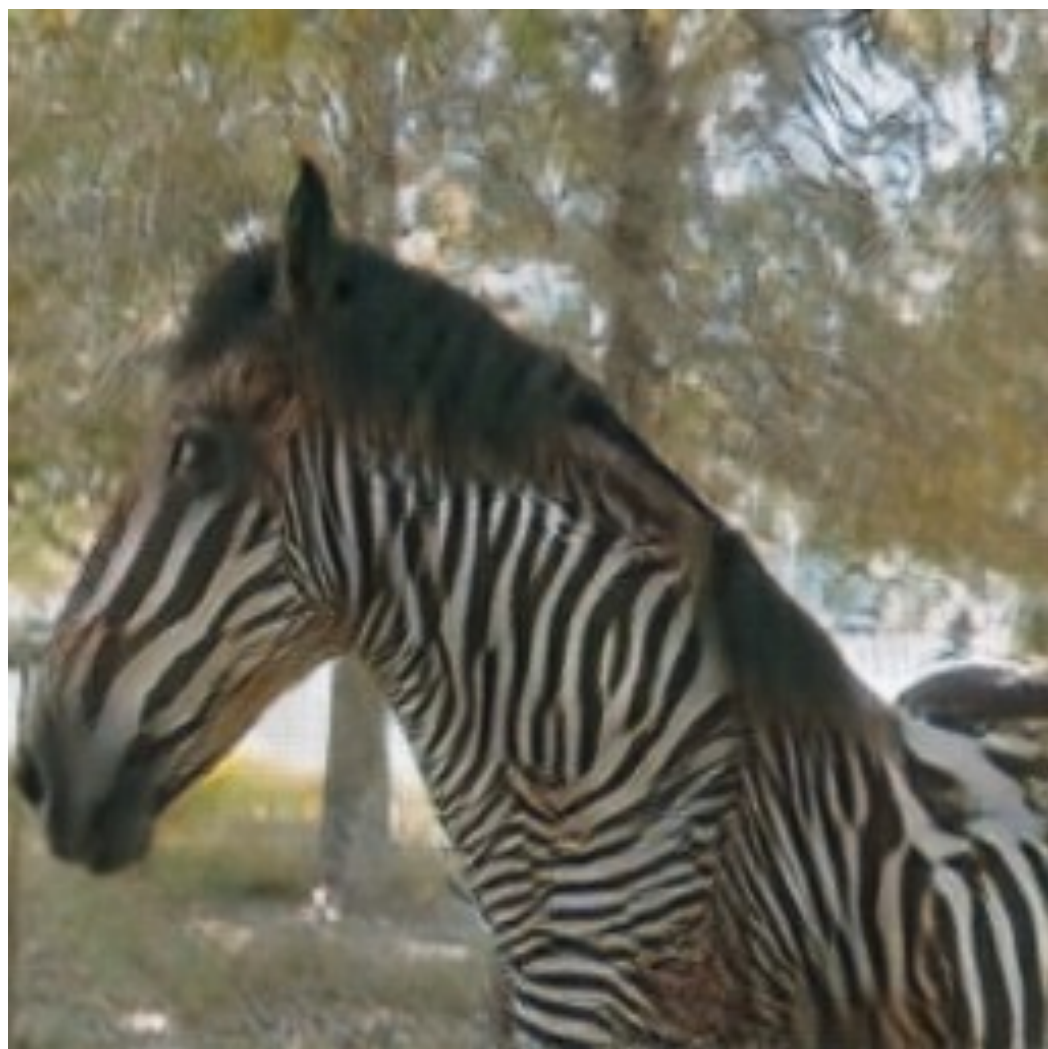


# Cycle Consistency Loss



# Cycle Consistency Loss







# Collection Style Transfer



Photograph  
@ Alexei Efros



Monet



Van Gogh



Cezanne



Ukiyo-e

Input



Monet



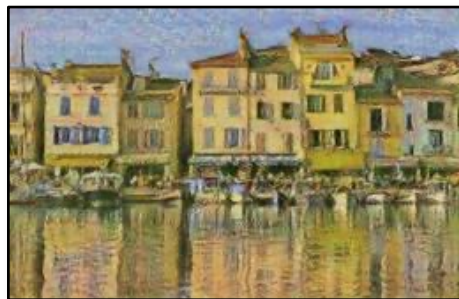
Van Gogh



Cezanne



Ukiyo-e



# Monet's paintings → photos





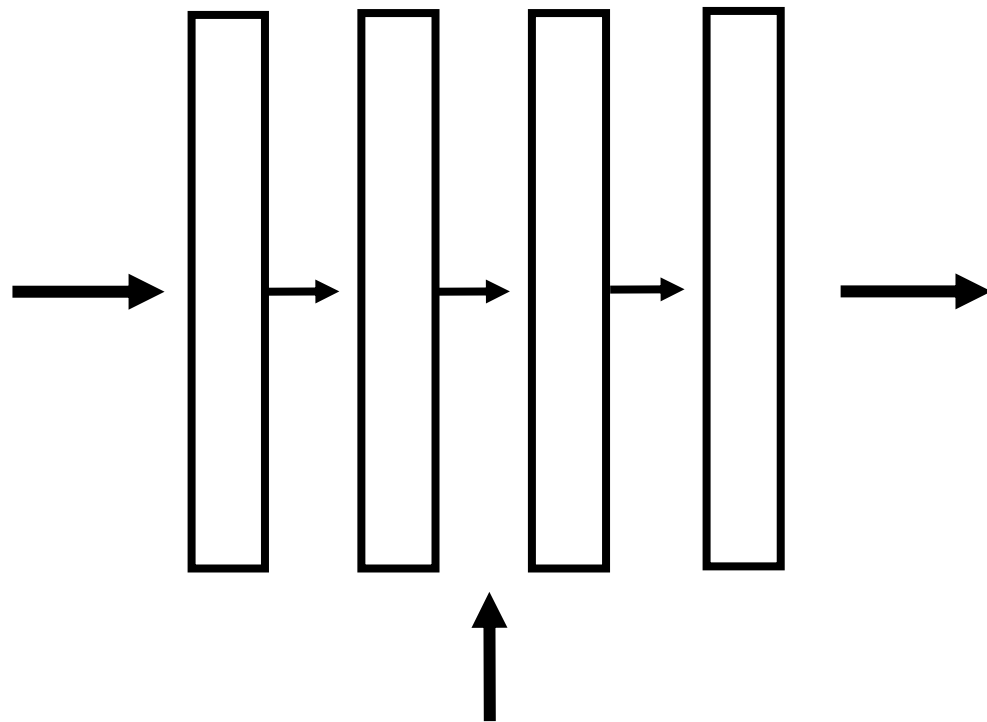
# Monet's paintings → photos



# Semantic Image Synthesis



semantic layout  $\mathbf{x}$



synthesized image  $\mathbf{y}$

- **Input:** input layout + target style image
- Output:** synthesized image

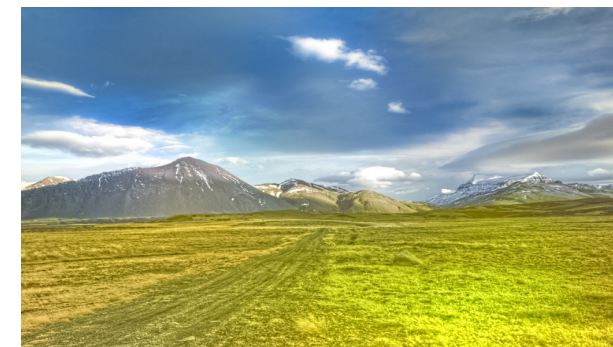
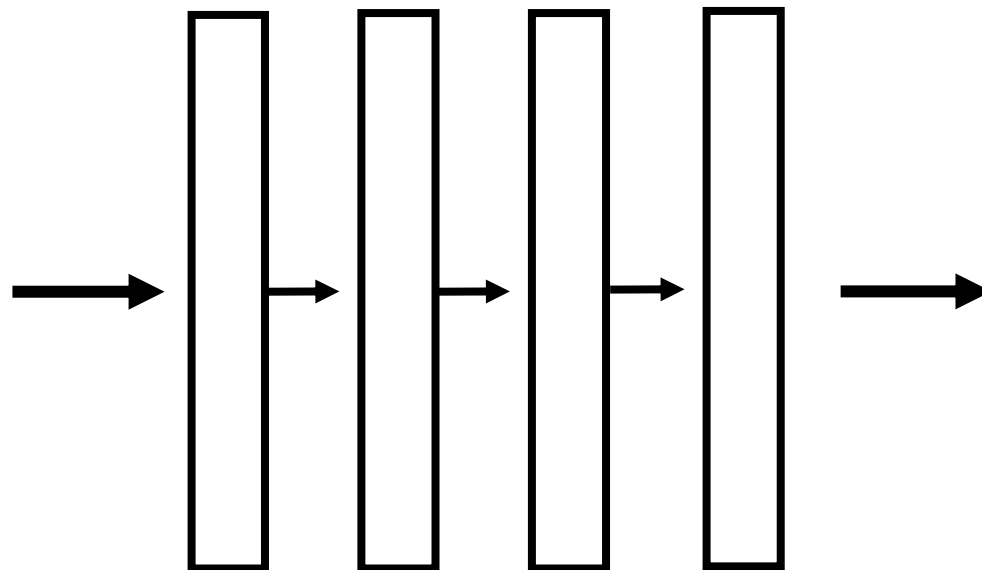


Target style image  $\mathbf{t}$

# Semantic Image Editing



image  $\mathbf{x}$



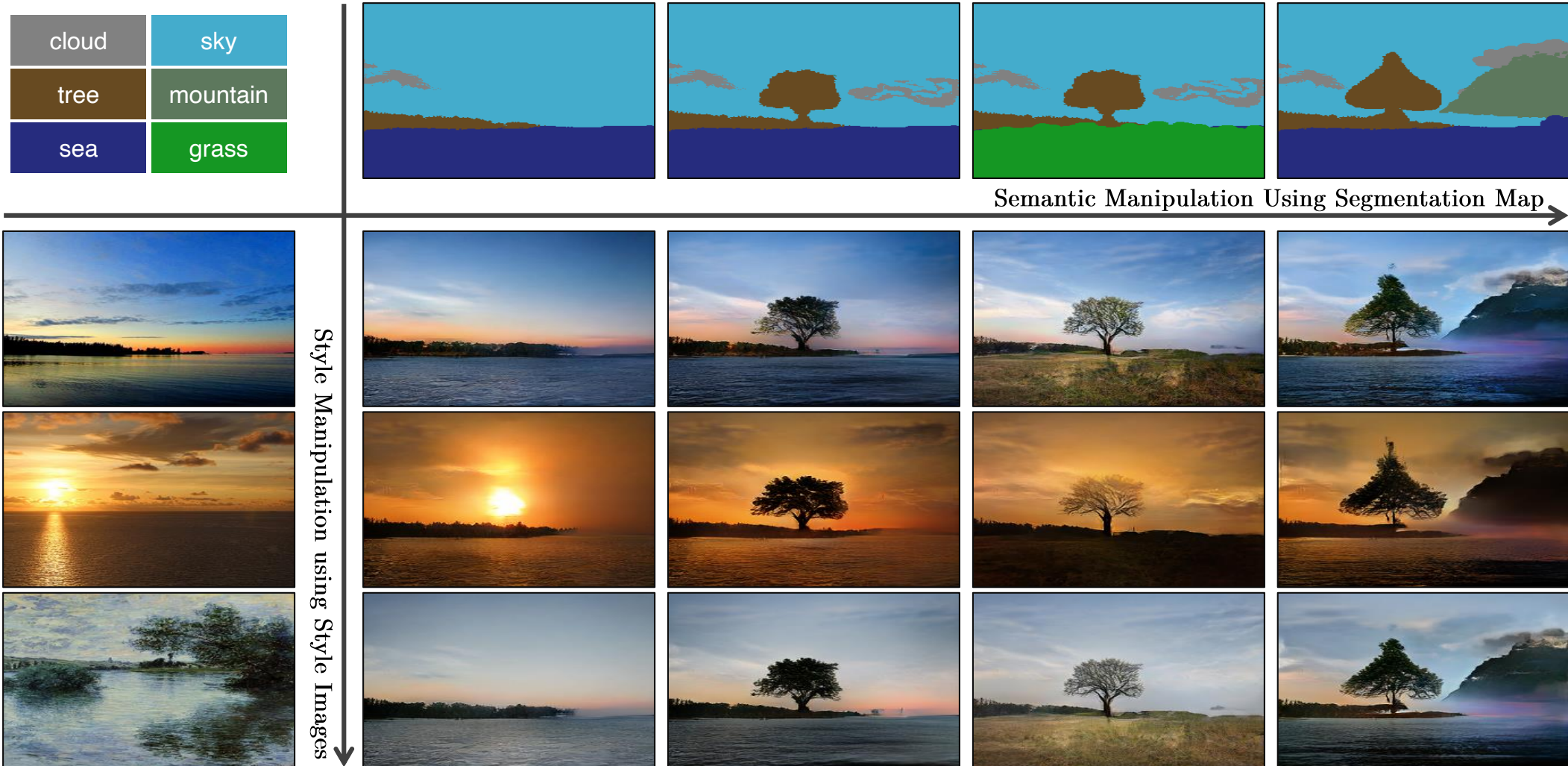
manipulated  
image  $\mathbf{y}$

- **Input:** input image  
(+ semantic layout)  
+ target attribute
- Output:** manipulated image

“more flowers”  
“more cloudy”  
Target scene  
attribute(s)  $\mathbf{t}$

# Semantic Image Synthesis (SPADE) (Park et al., 2019)

- Image generation conditioned on semantic layouts



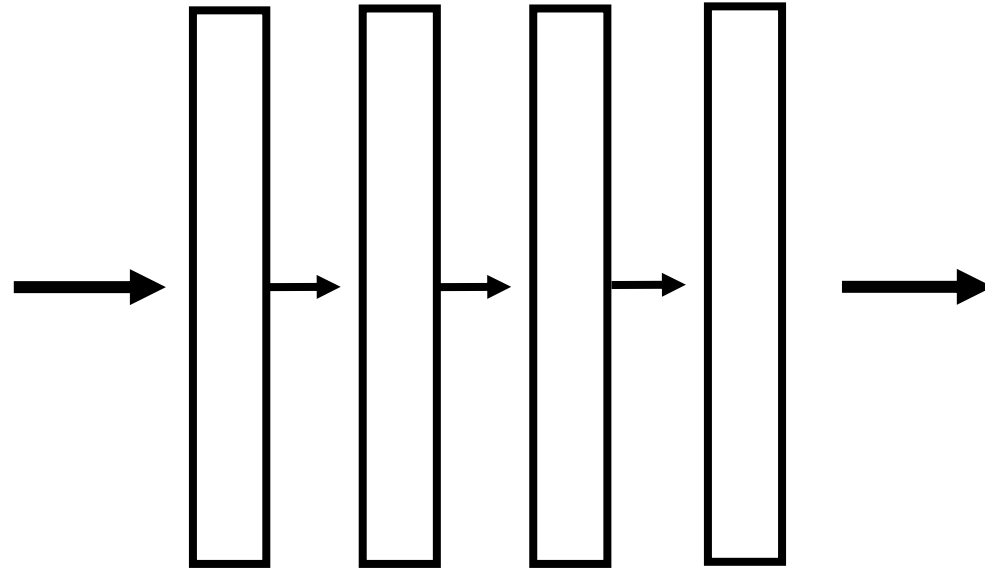
# Imagine this scene in a snowy winter day...







image  $\mathbf{x}$



“more flowers”  
“more cloudy”

Target transient  
scene attribute  $\mathbf{t}$



manipulated  
image  $\mathbf{y}$

• **Input:** input image  
+ target attribute

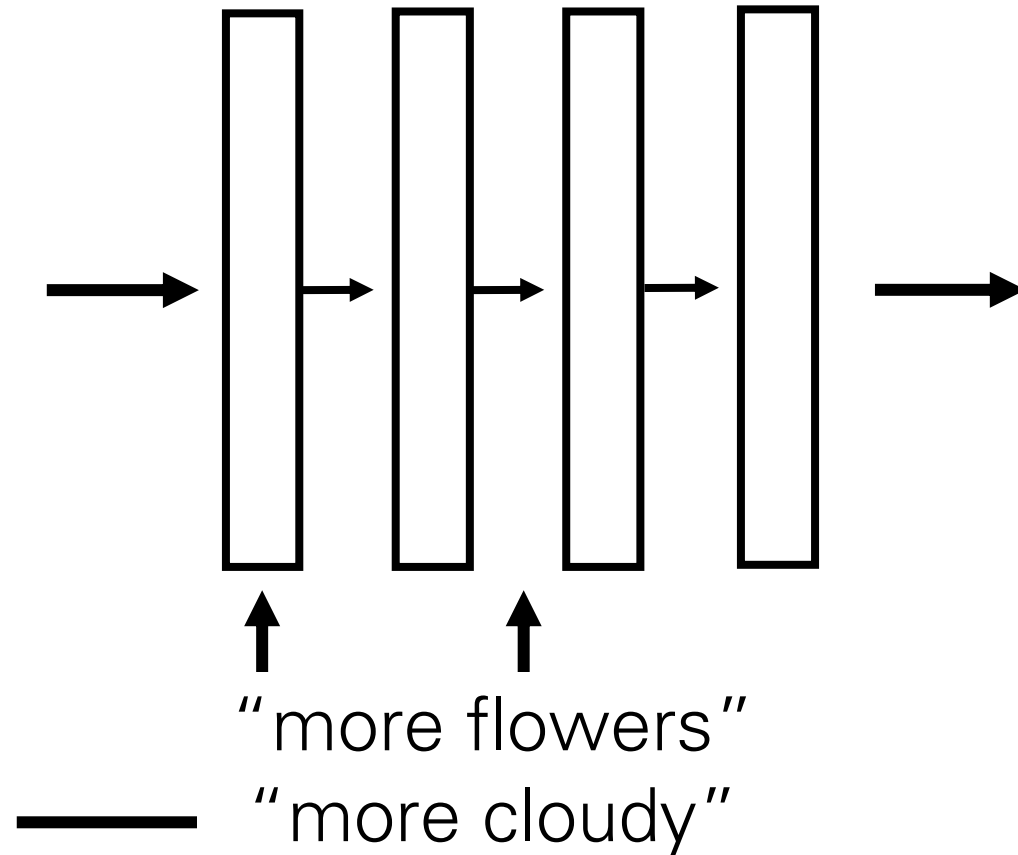
**Output:** manipulated image



image  $\mathbf{x}$



semantic layout  $\mathbf{l}$



Target transient scene attribute  $\mathbf{t}$



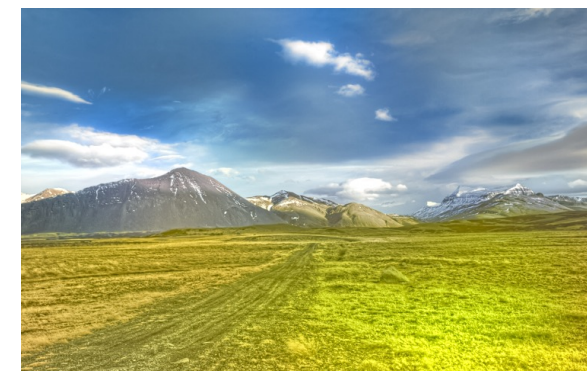
manipulated image  $\mathbf{y}$

- **Input:** input image + semantic layout + target attribute
- Output:** manipulated image

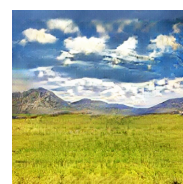
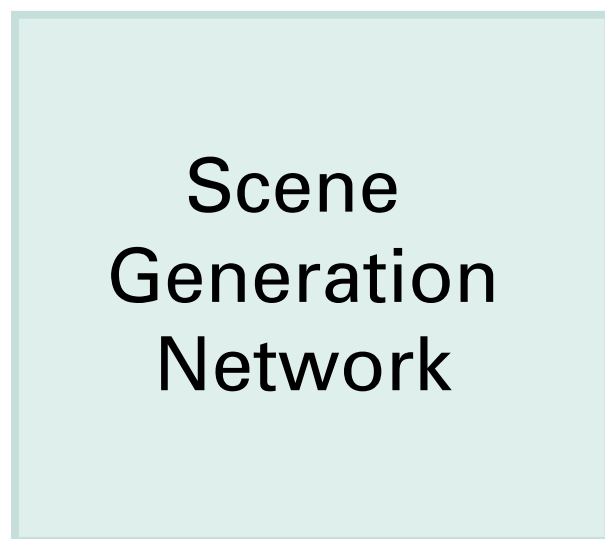




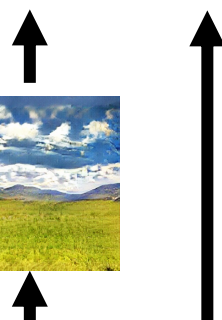
image  $\mathbf{x}$  and semantic layout  $\mathbf{l}$



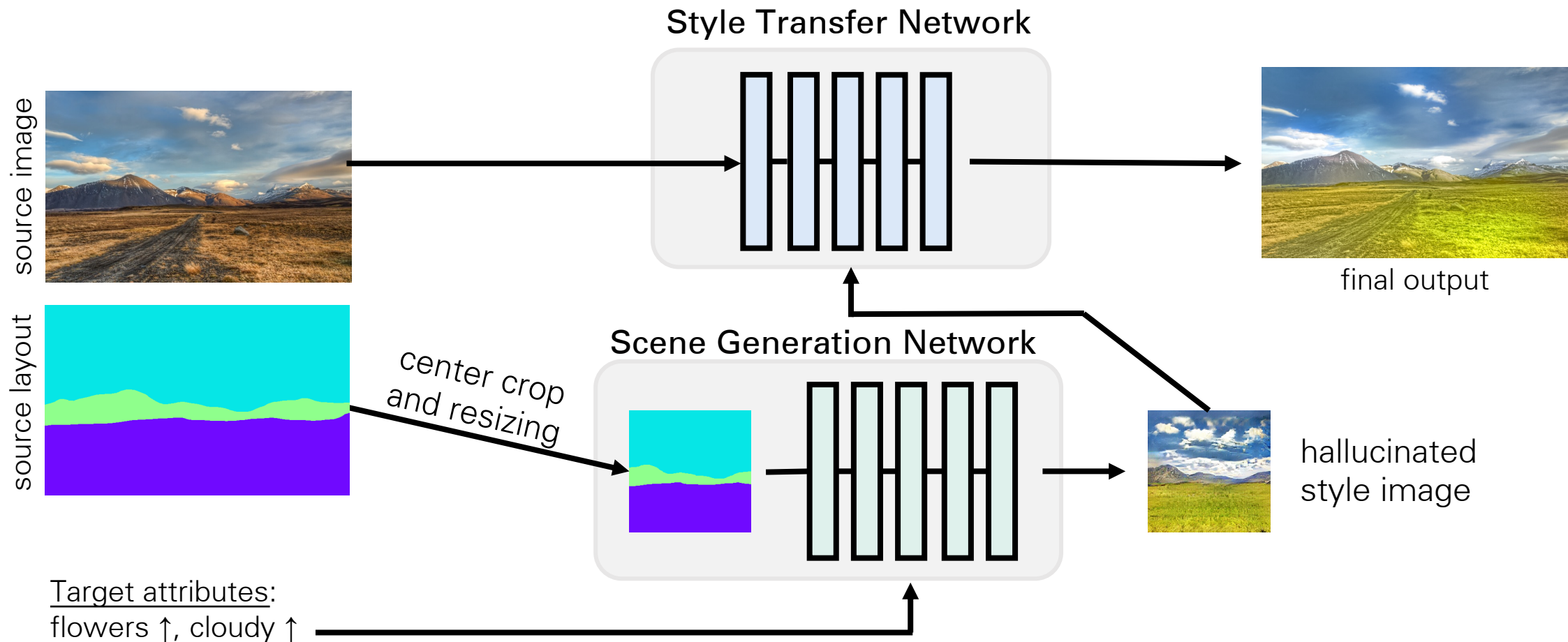
manipulated image  $\mathbf{y}$



"more flowers" target transient  
"more cloudy" scene attribute  $\mathbf{t}$



**Idea:** Hallucinate alternative version of the input scene consistent with target attributes and use this image as the style image in the photo style transfer.



- **Scene Generation Network**

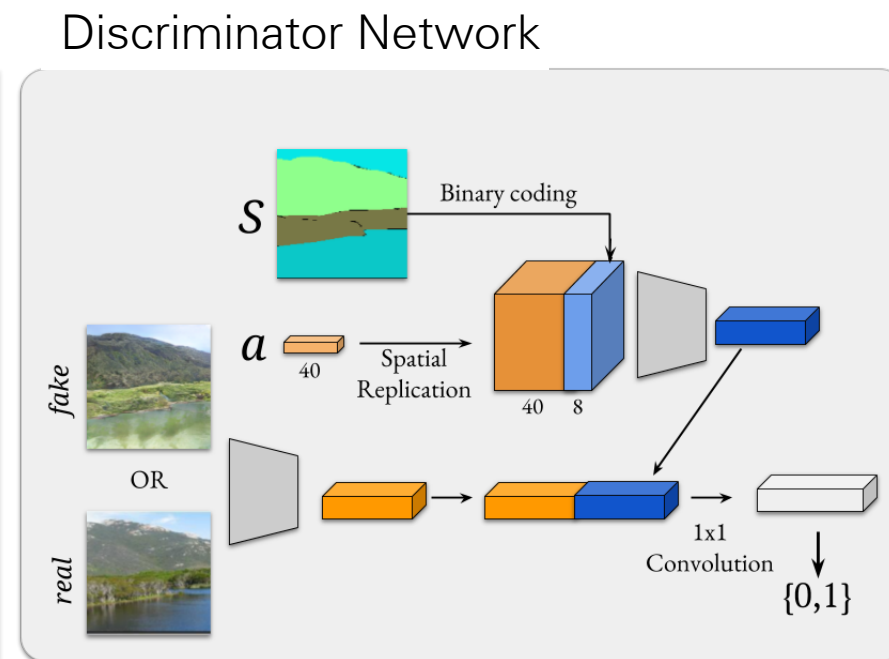
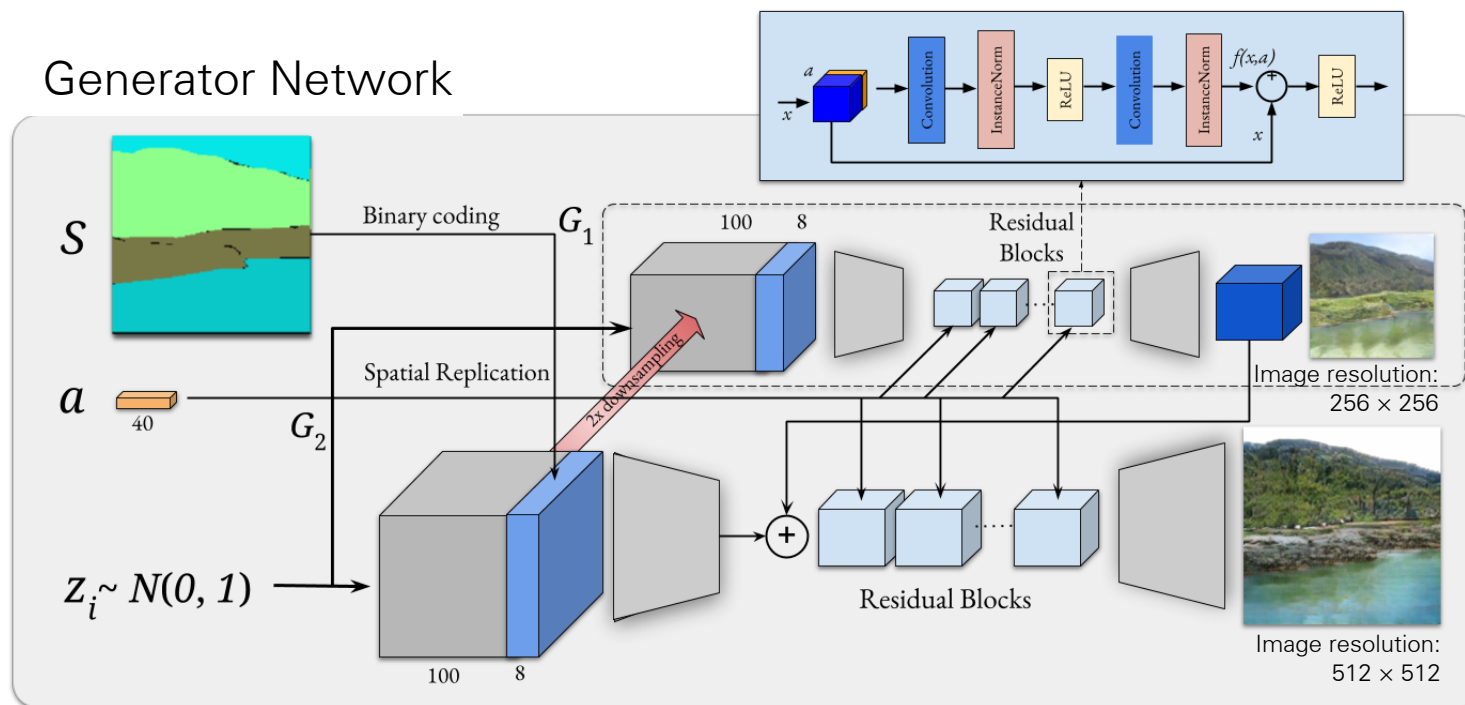
- A conditioned GAN model with two conditions:
  - (1) semantic layout,
  - (2) target attributes

- **Style Transfer Network**

- A deep photo style transfer network that modifies the look of the source image based on the hallucinated style image

# Scene Generation Network

- The semantic layout categories are encoded into 8-bit binary codes
- The transient attributes are represented by a 40-d vector.



- An architecture similar to Pix2pixHD model (Wang et al. 2018)
- **Generator network:** A coarse-to-fine model with 2 generator networks
- **Discriminator network:** A combination of three different discriminator networks operating at an image pyramid of 3 scales





Manipulating Attributes of Natural Scenes via Hallucination [Karacan et al., 2020]



Manipulating Attributes of Natural Scenes via Hallucination [Karacan et al., 2020]



night

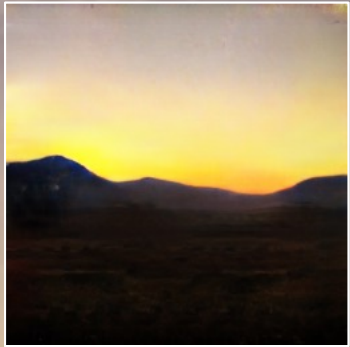


prediction

Manipulating Attributes of Natural Scenes via Hallucination [Karacan et al., 2020]



sunset

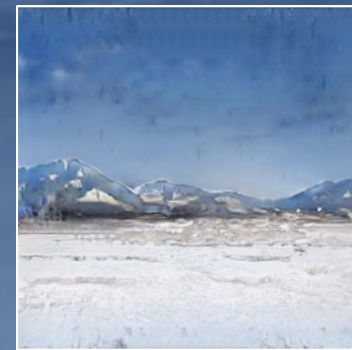


prediction





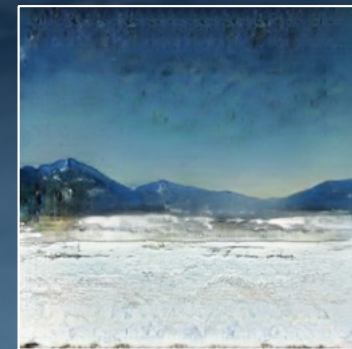
snow



prediction

Manipulating Attributes of Natural Scenes via Hallucination [Karacan et al., 2020]

winter



prediction



Manipulating Attributes of Natural Scenes via Hallucination [Karacan et al., 2020]

# Spring and clouds



prediction



Moist, rain and fog



prediction

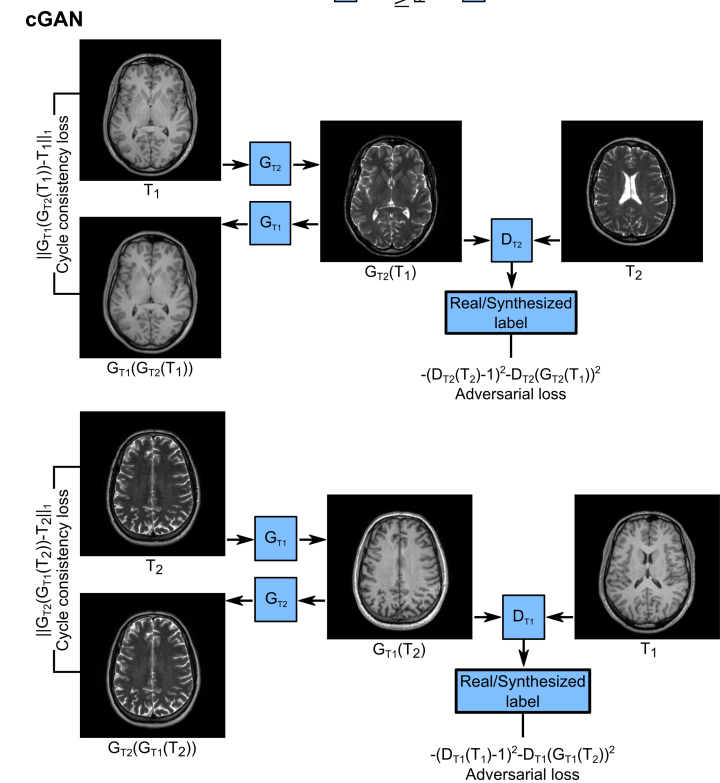
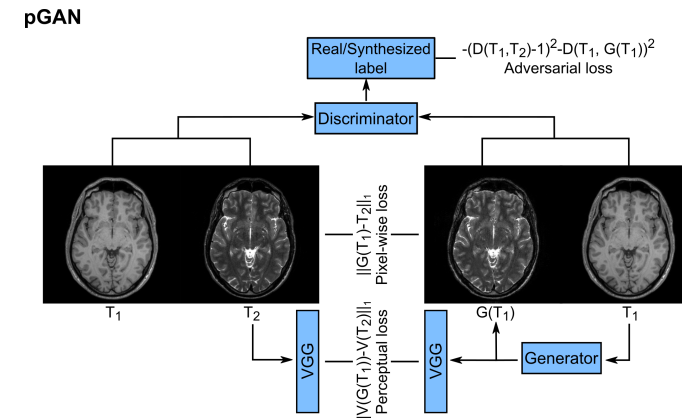
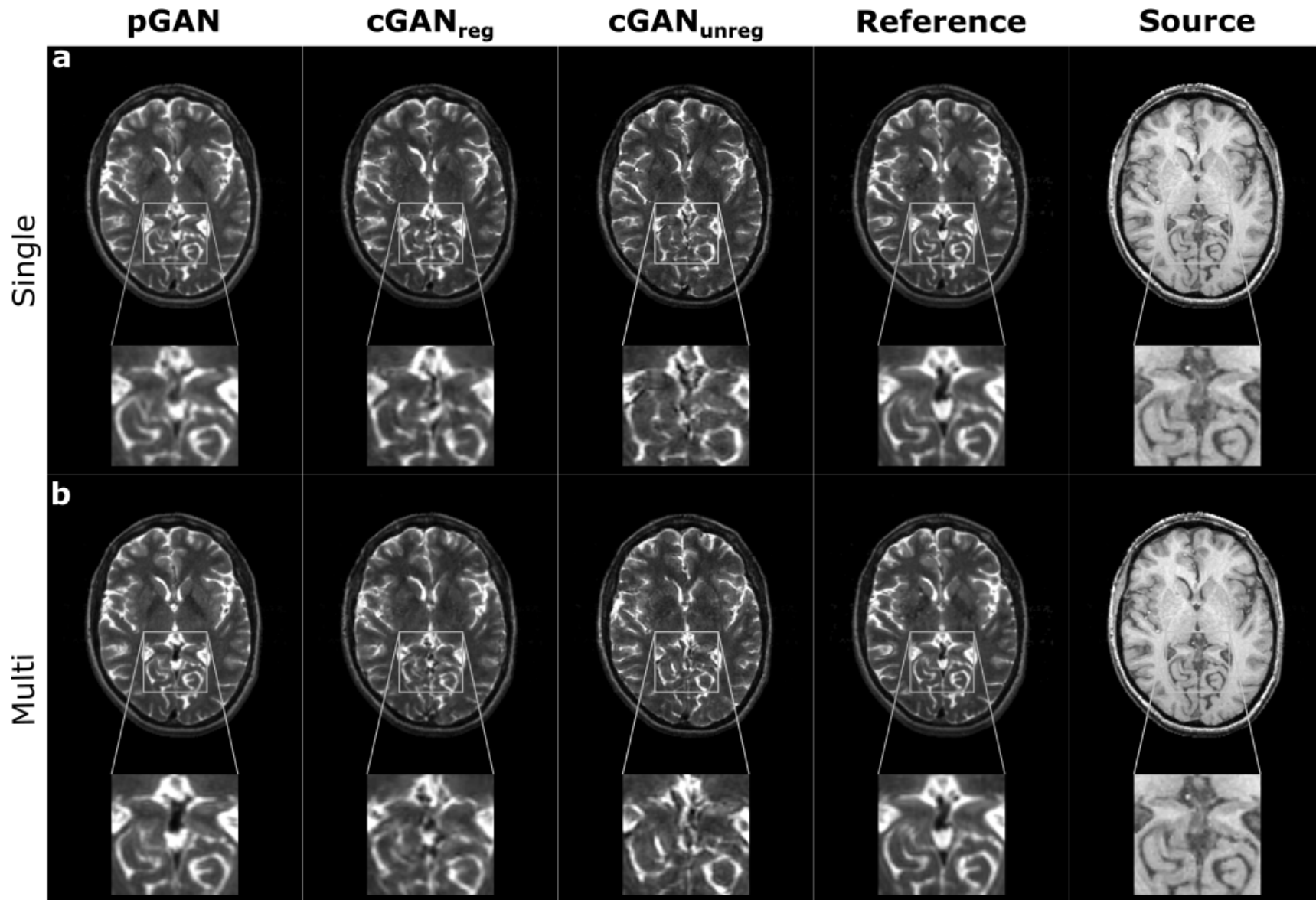


flowers

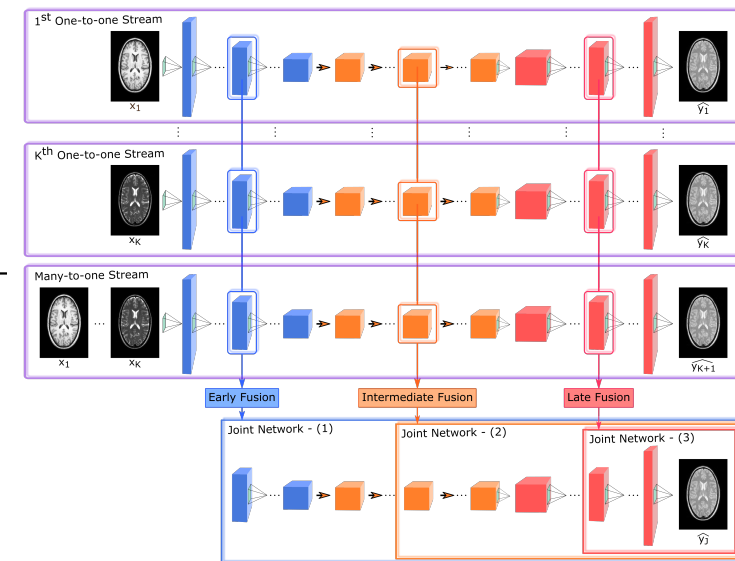
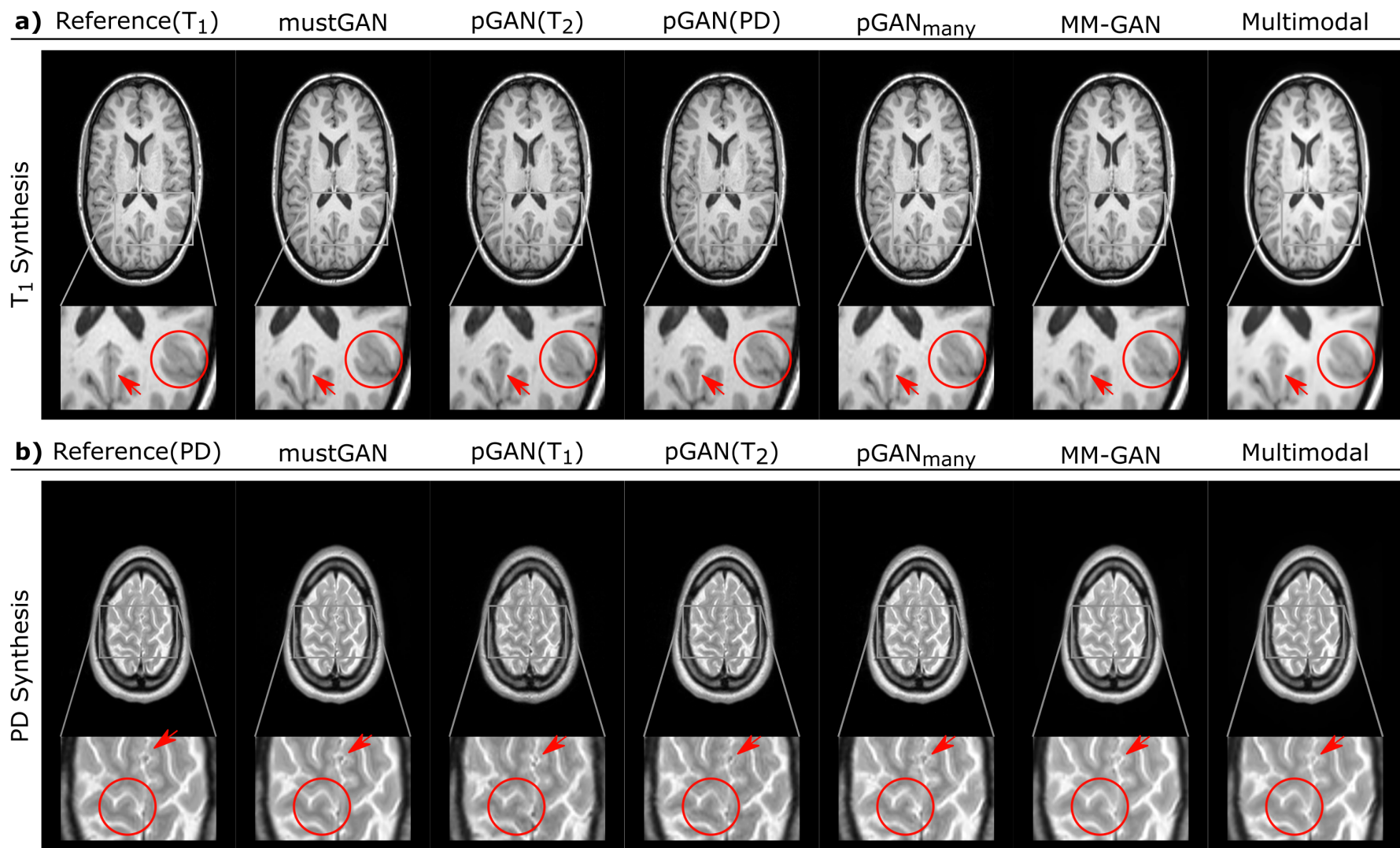


prediction





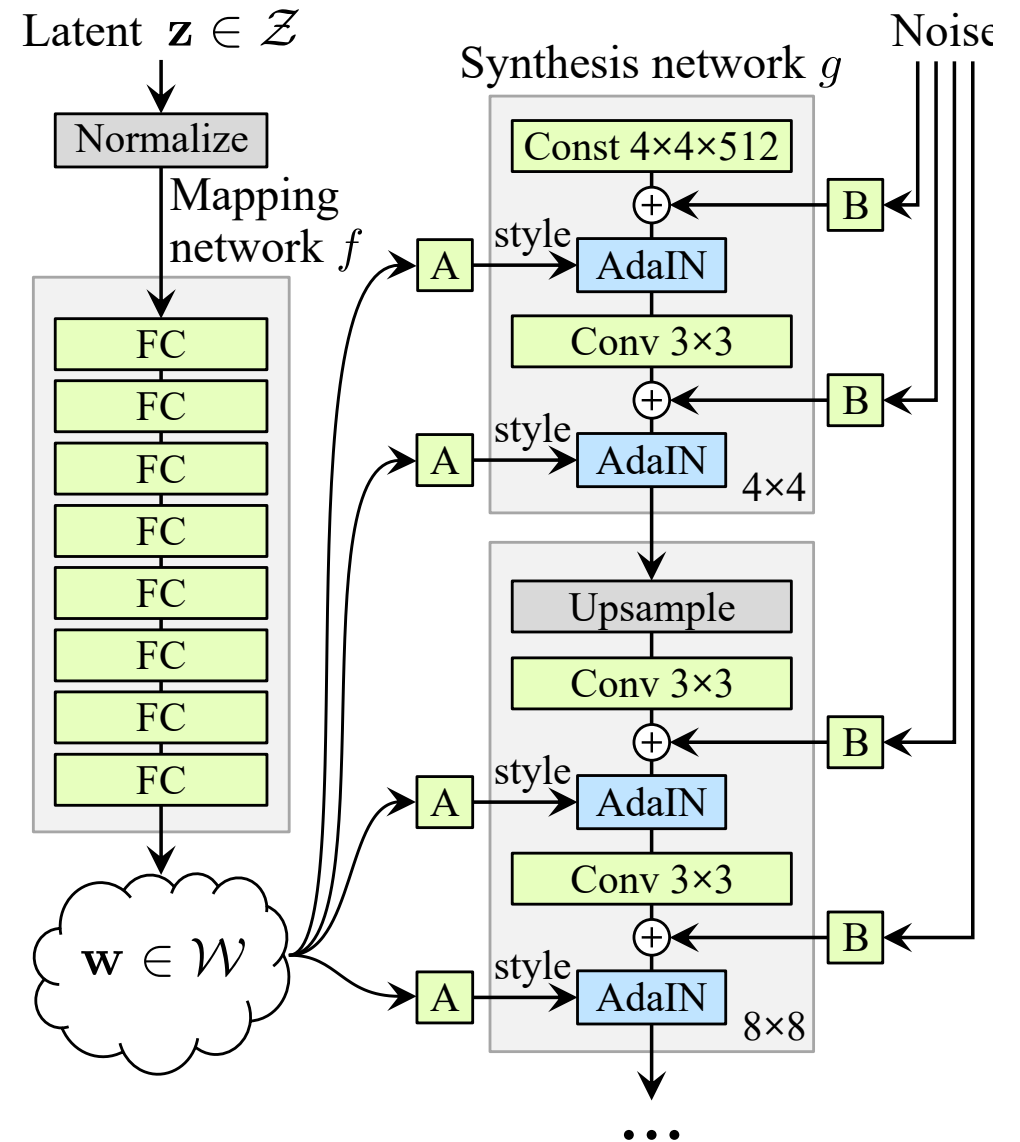
- Image Synthesis in Multi-Contrast MRI [Ul Hassan Dar et al. 2019]



- Image Synthesis in Multi-Contrast MRI [Mahmut Yurt et al. 2021]

# Recall StyleGANs (Karras et al., 2019)

- A new architecture motivated by the style transfer networks
- allows unsupervised separation of high-level attributes and stochastic variation in the generated images

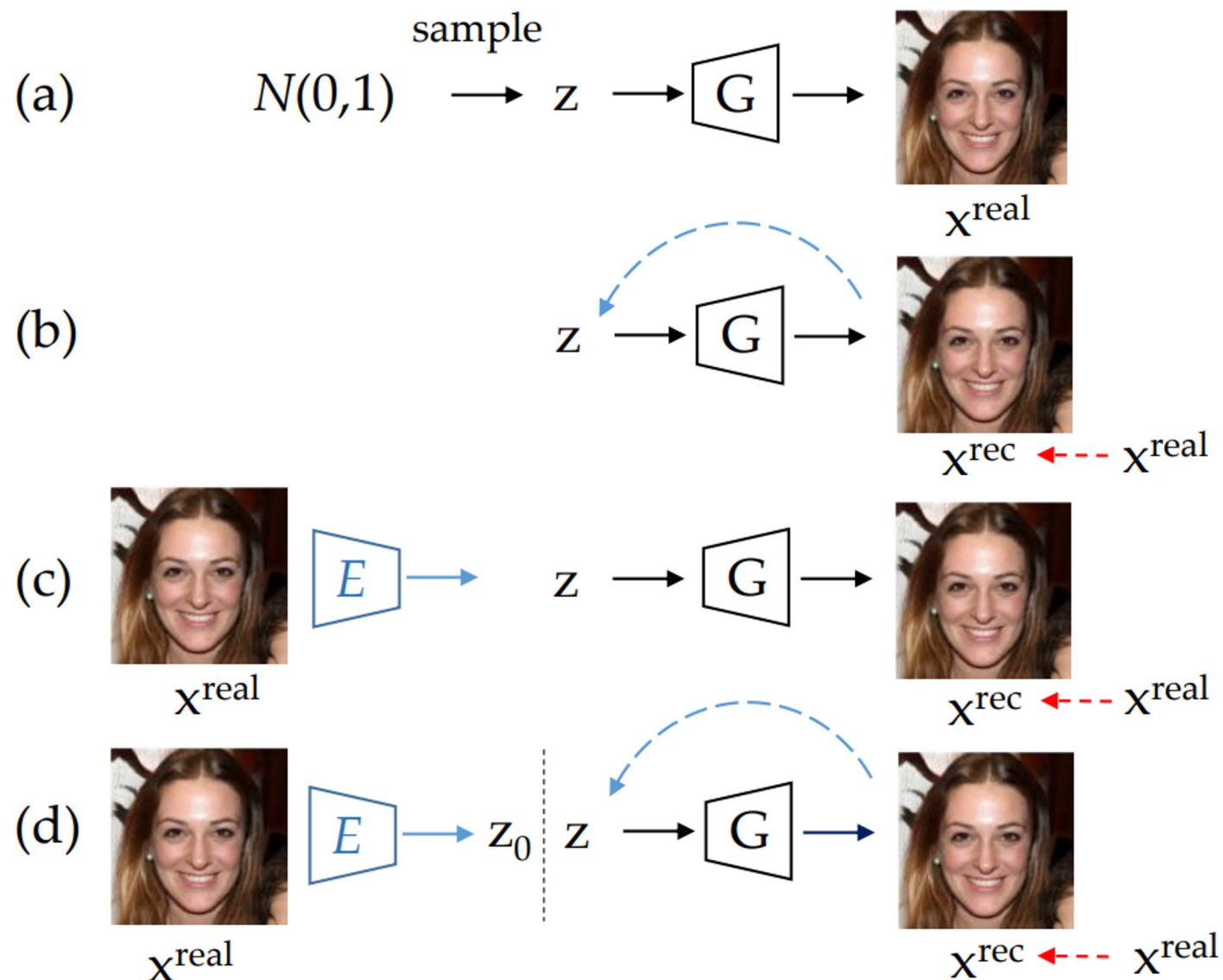




# GAN Inversion

3 methods of inversion:

- Optimization-based (b)
- Learning-based (c)
- Hybrid (d)



# Semantic Image Editing

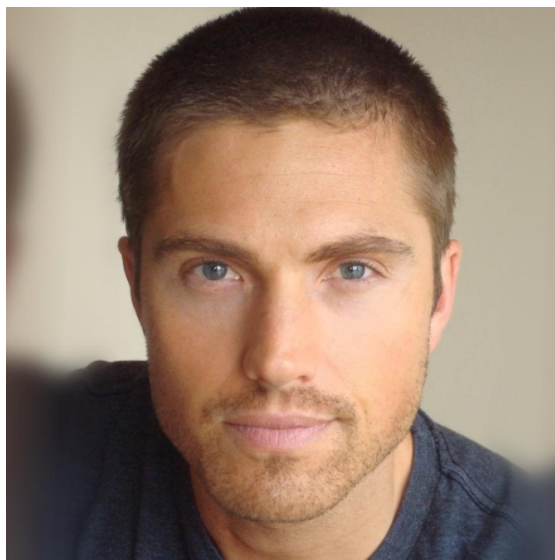
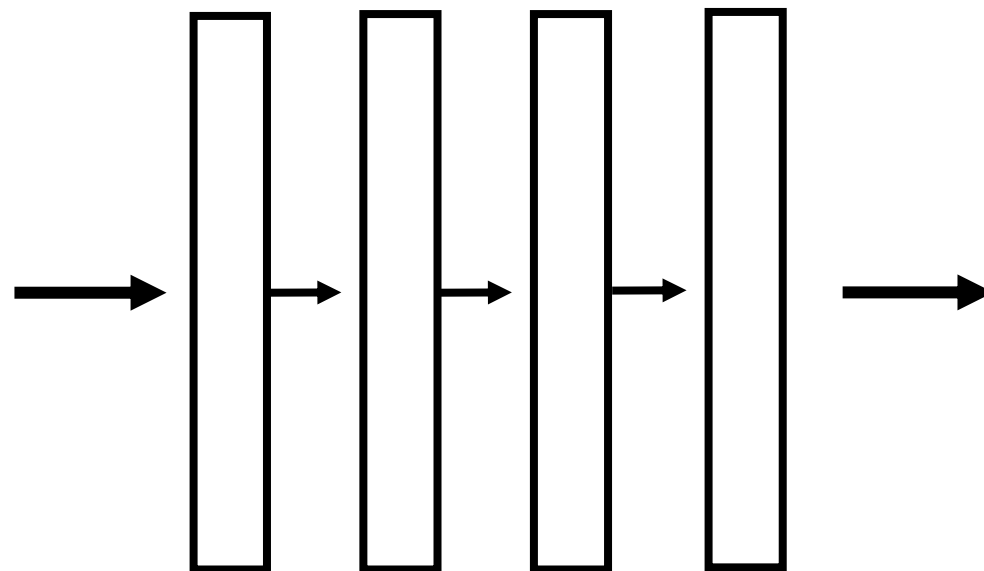
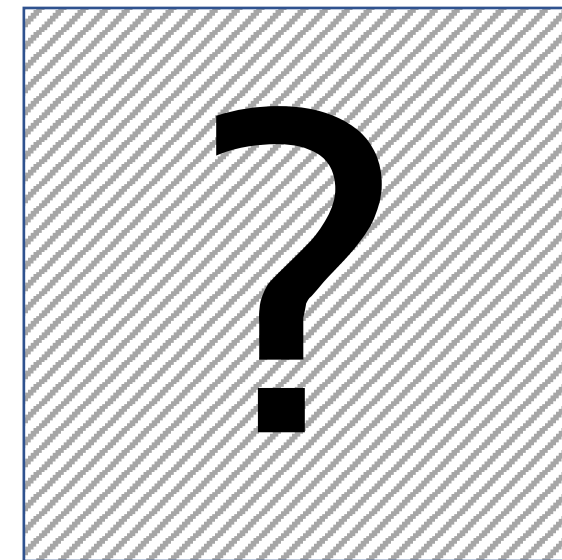


image  $\mathbf{x}$



"he is tanned"



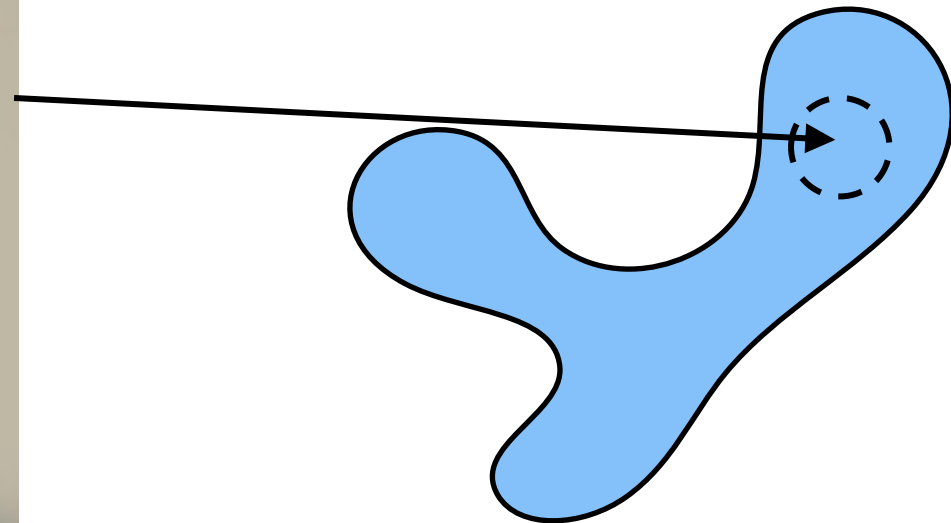
manipulated  
image  $\mathbf{y}$

- **Input:** input image  
+ target description
- Output:** manipulated image

Target description  $\mathbf{t}$

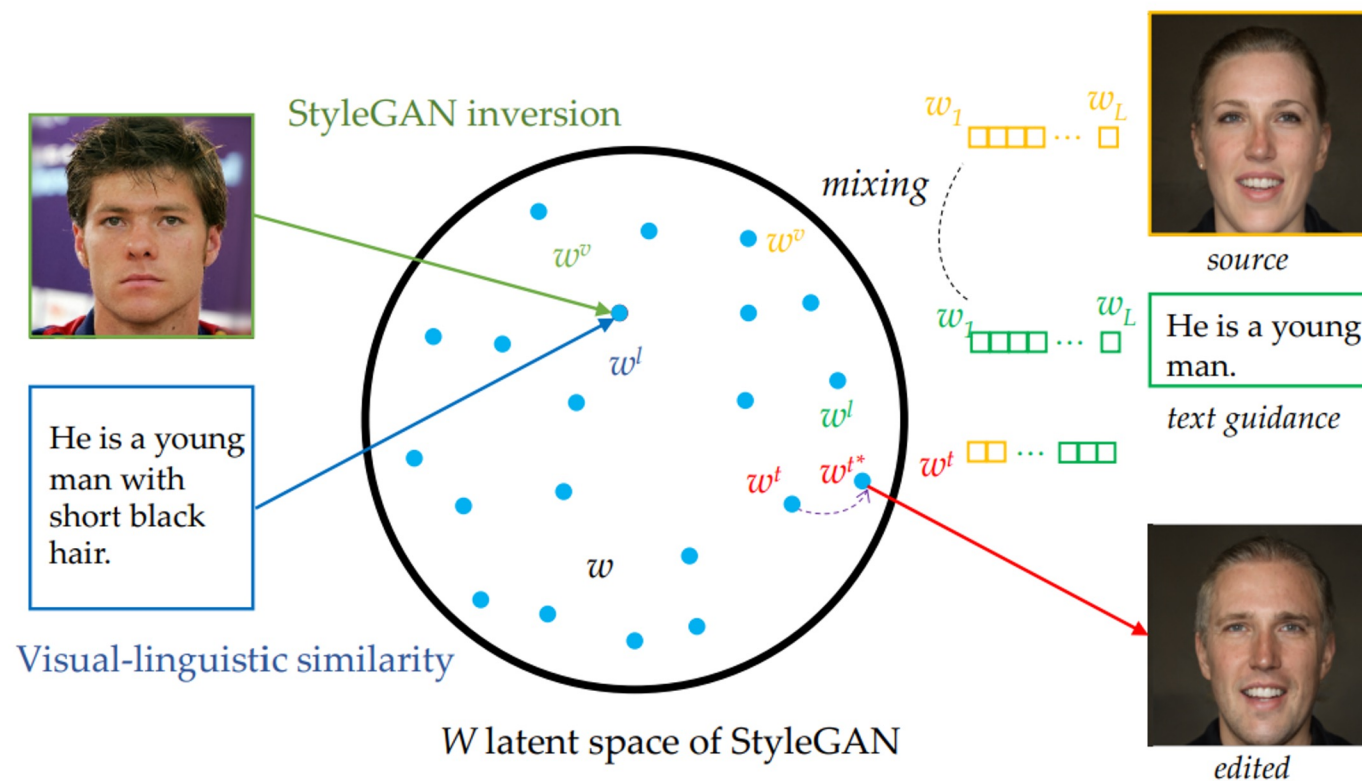


- Utilize pretrained StyleGAN model as a natural prior for face images
- Project input image to the latent space and perform edit in the vicinity of that point.



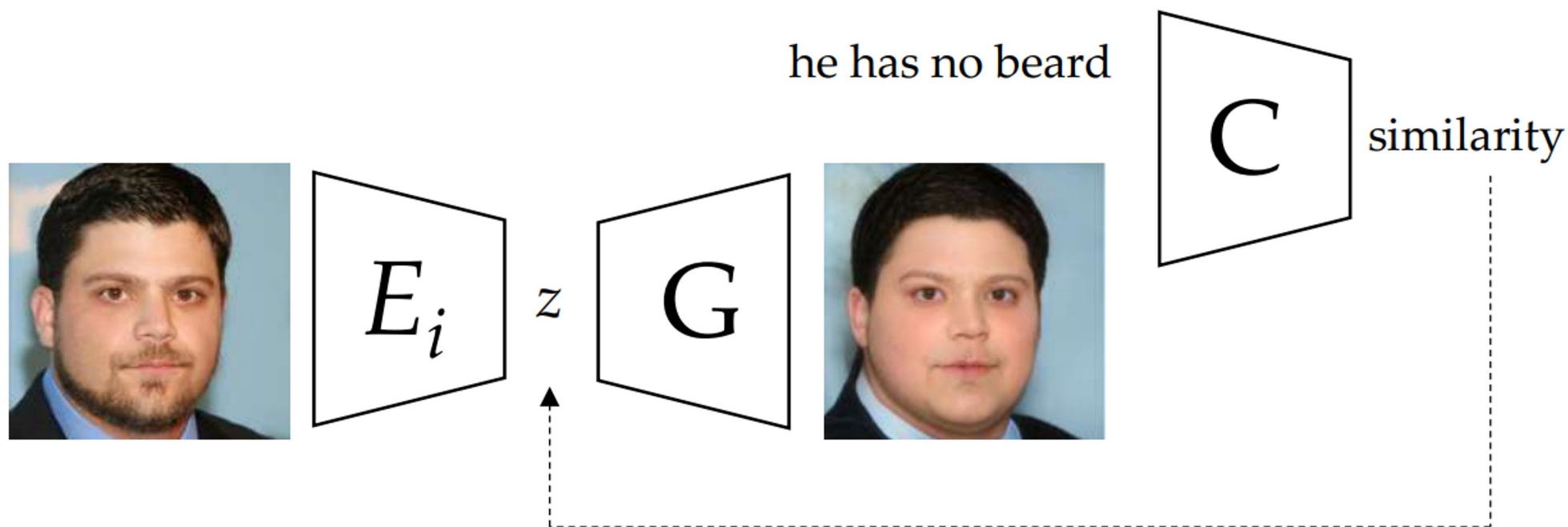
# TediGAN

- A recent inversion based manipulation model.
- They proposed two different approaches.
- In their first approach, they learn a common visual-linguistic semantic space.
- They train a text encoder such that both modalities are embedded to the same space



# TediGAN

- Their second method is optimization based.
- They use CLIP to provide a signal to directly optimize the latent code



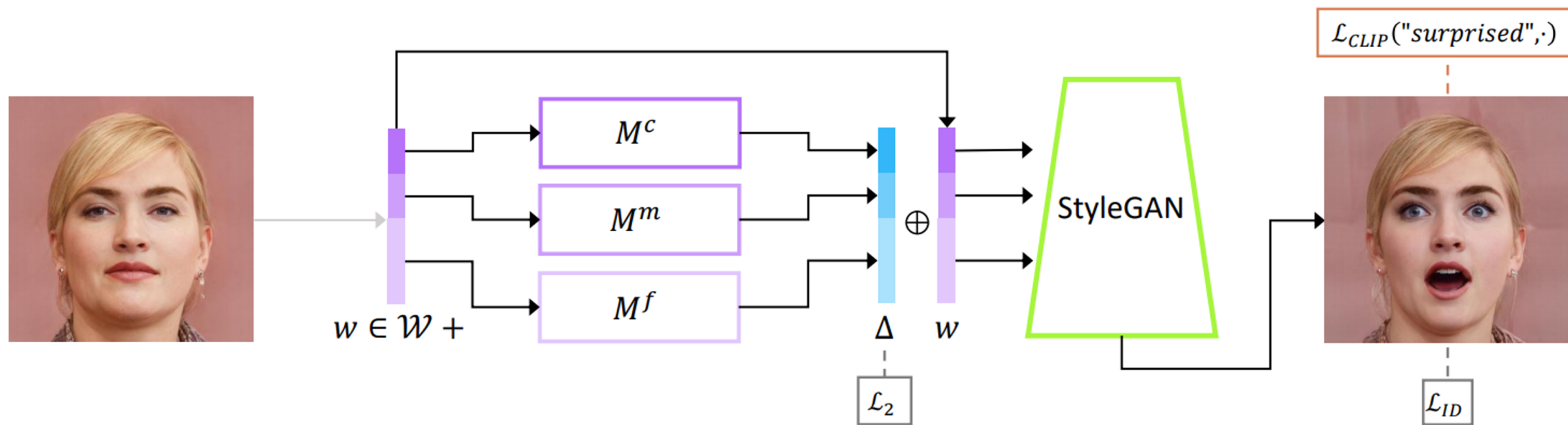
# StyleCLIP

- StyleCLIP is another inversion based manipulation model
- Direct optimization
- Similar to TediGAN, they use CLIP to optimize the latent code directly
- They include a term to preserve identity

$$\arg \min_{w \in \mathcal{W}^+} D_{\text{CLIP}}(G(w), t) + \lambda_{\text{L2}} \|w - w_s\|_2 + \lambda_{\text{ID}} \mathcal{L}_{\text{ID}}(w)$$

# StyleCLIP

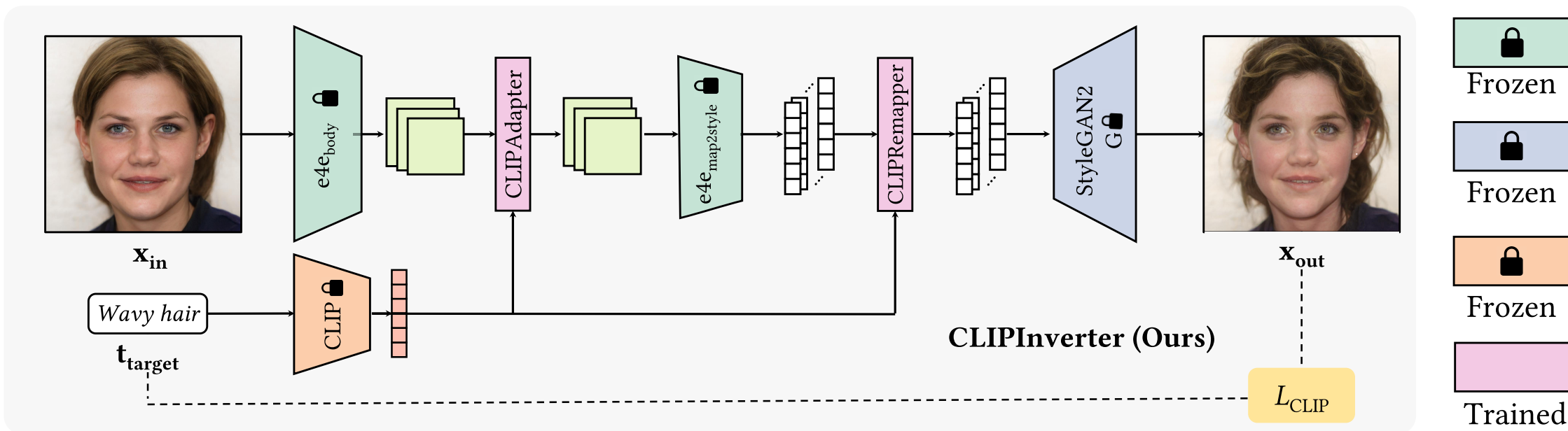
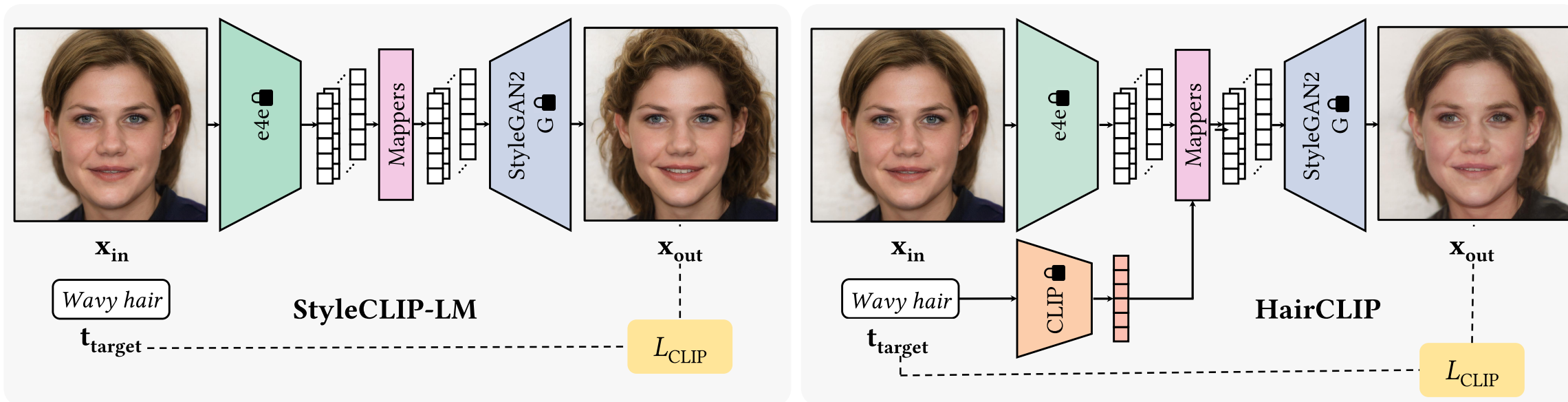
- Images are inverted to latent code and residuals (delta in the figure).
- Textual input is not used during inversion, it is only used in the loss
- Required to train different mappers for different text prompts.



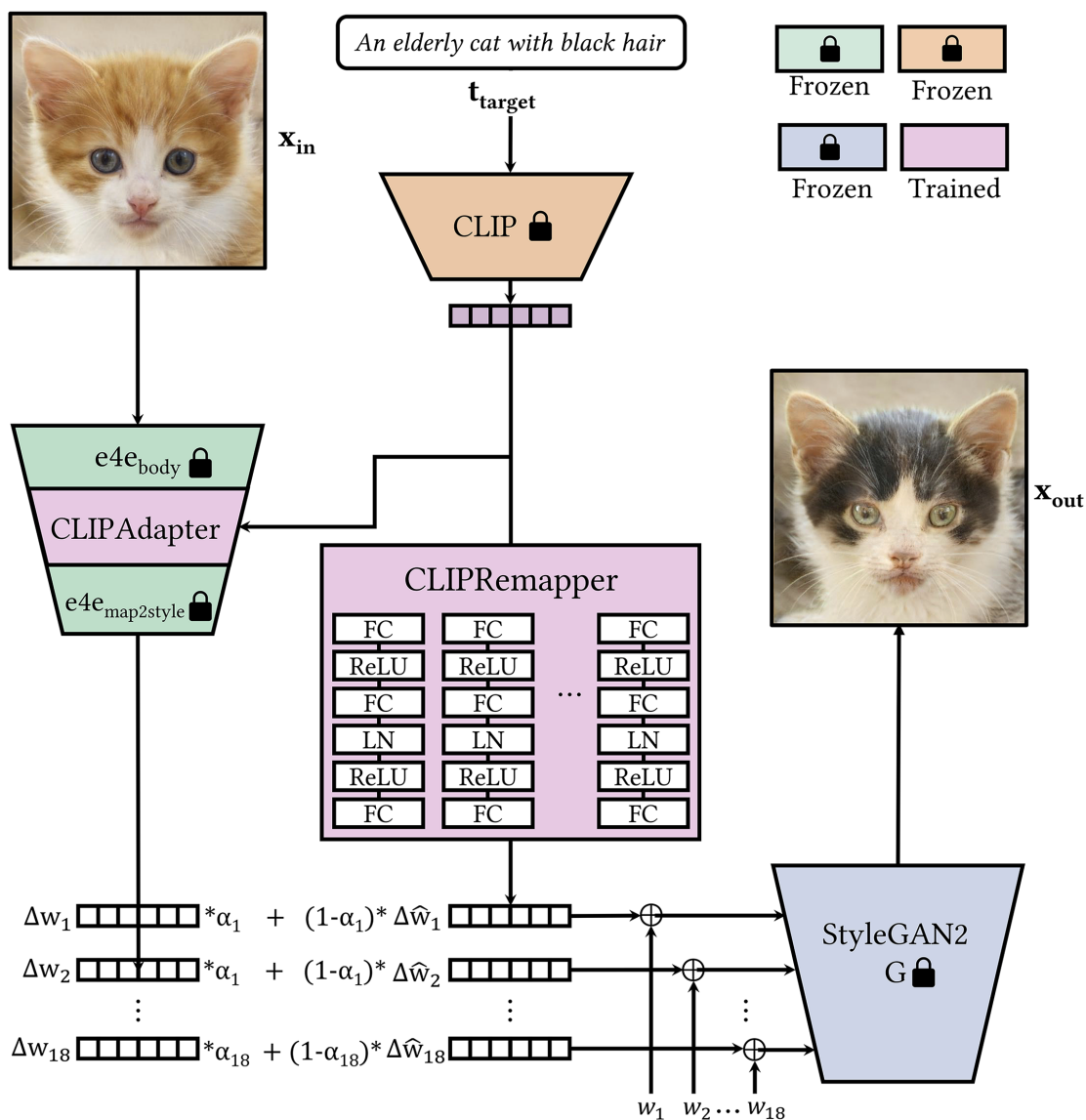
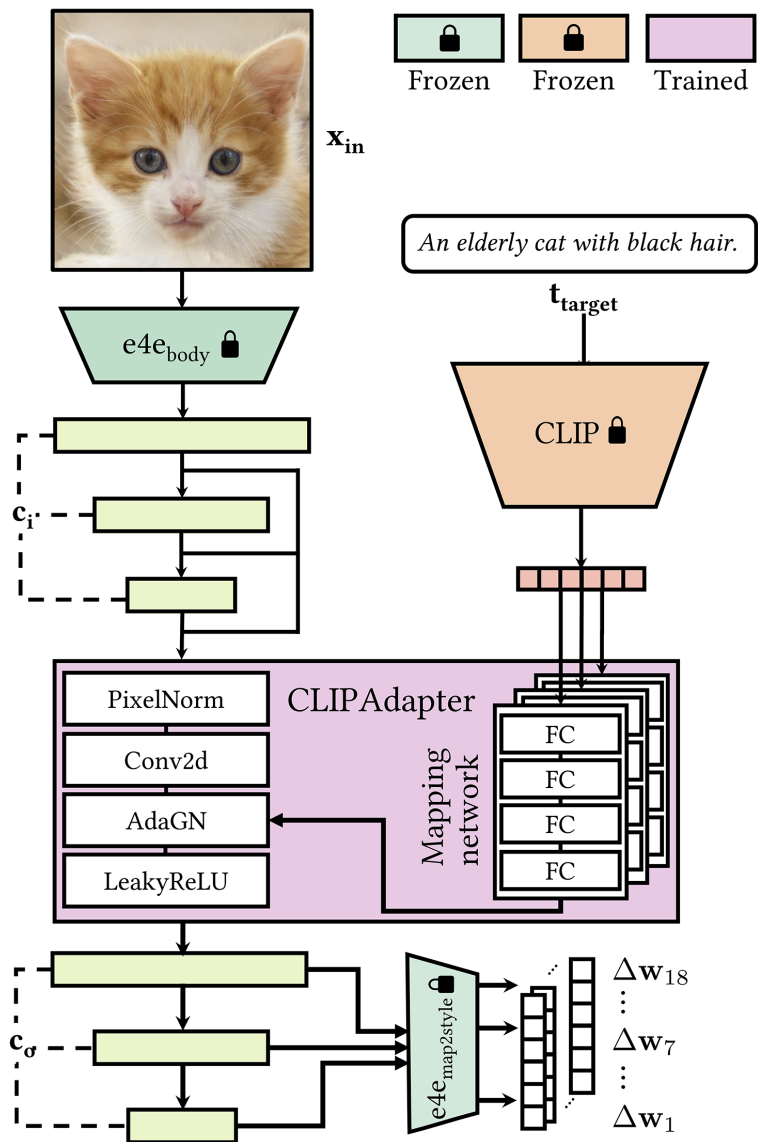




# CLIPInverter



# CLIPInverter



# CLIPInverter



He is smiling and has gray hair, high cheekbones, eyeglasses, double chin, and bags under eyes.



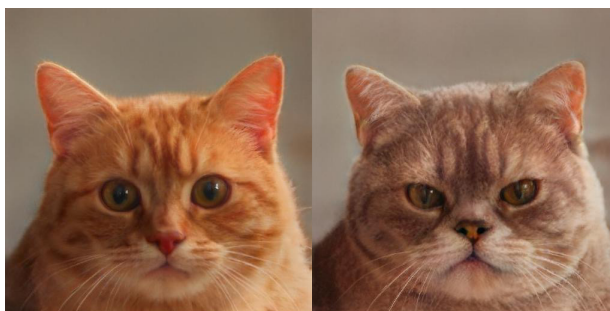
An old norwegian forest cat



This smaller bird is bright red and has black wings.



She has bangs. She is young and wears lipstick.



A grumpy elderly british shorthair cat



This is a small bird with green, yellow and blue on the breast, cheek patches, and crown.



This woman is attractive and has arched eyebrows, straight hair, and blond hair.



A fearful cat with grey hair



This bird is brown with yellow and has a long, pointy beak.

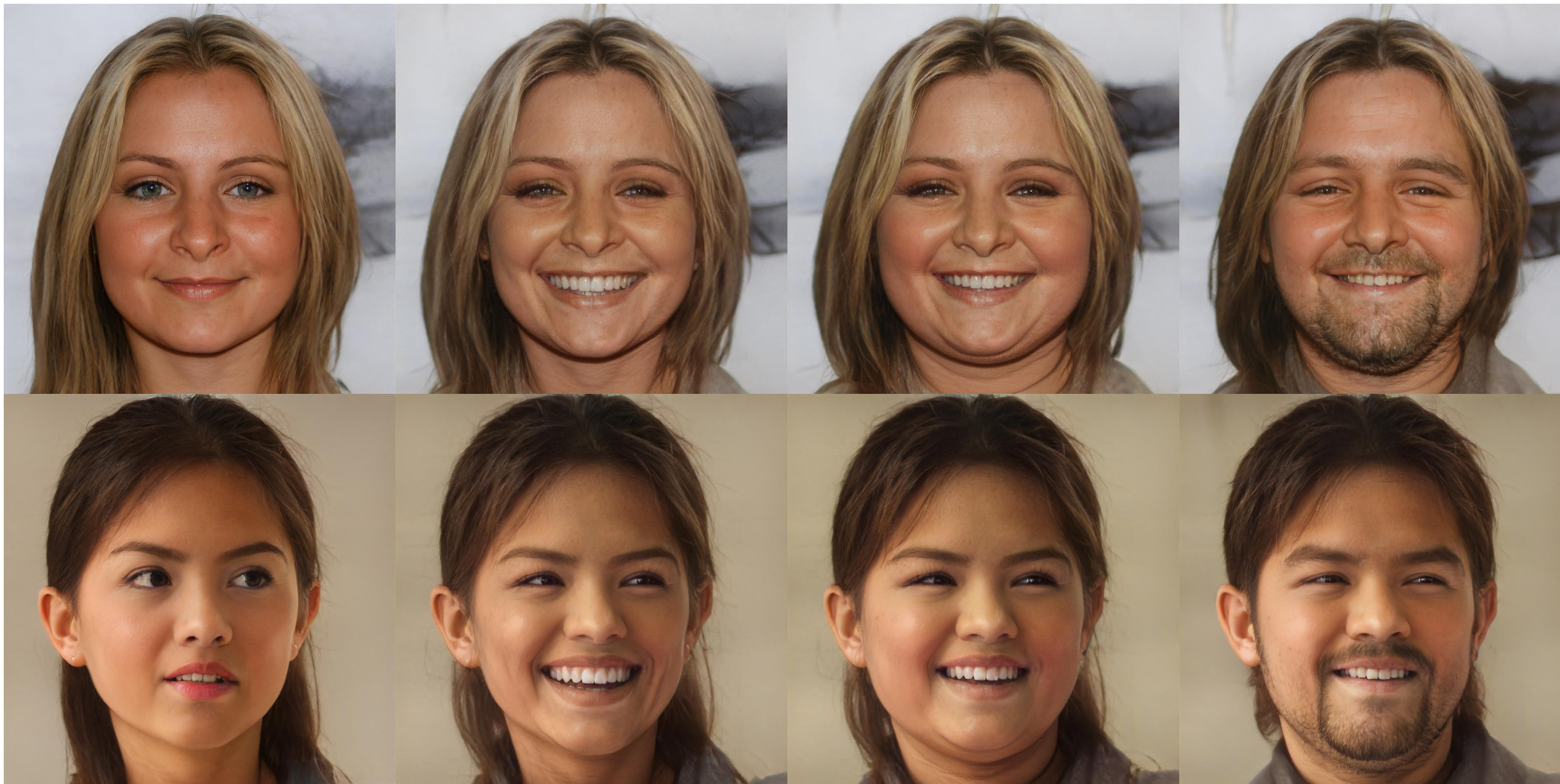
# CLIPInverter

Original

Smiling

Smiling+Chubby

Smiling+Chubby+Beard



# CLIPInverter



The person has bags under eyes. ←

→ She has wavy hair. She is wearing lipstick. She is young.



This person is young and has brown hair. ←

→ This person has mustache.

# CLIPInverter



An elderly cat with grey hair. ←

→ A british shorthair kitty.



An old cat. ←

→ A cat with ginger hair.

# CLIPInverter

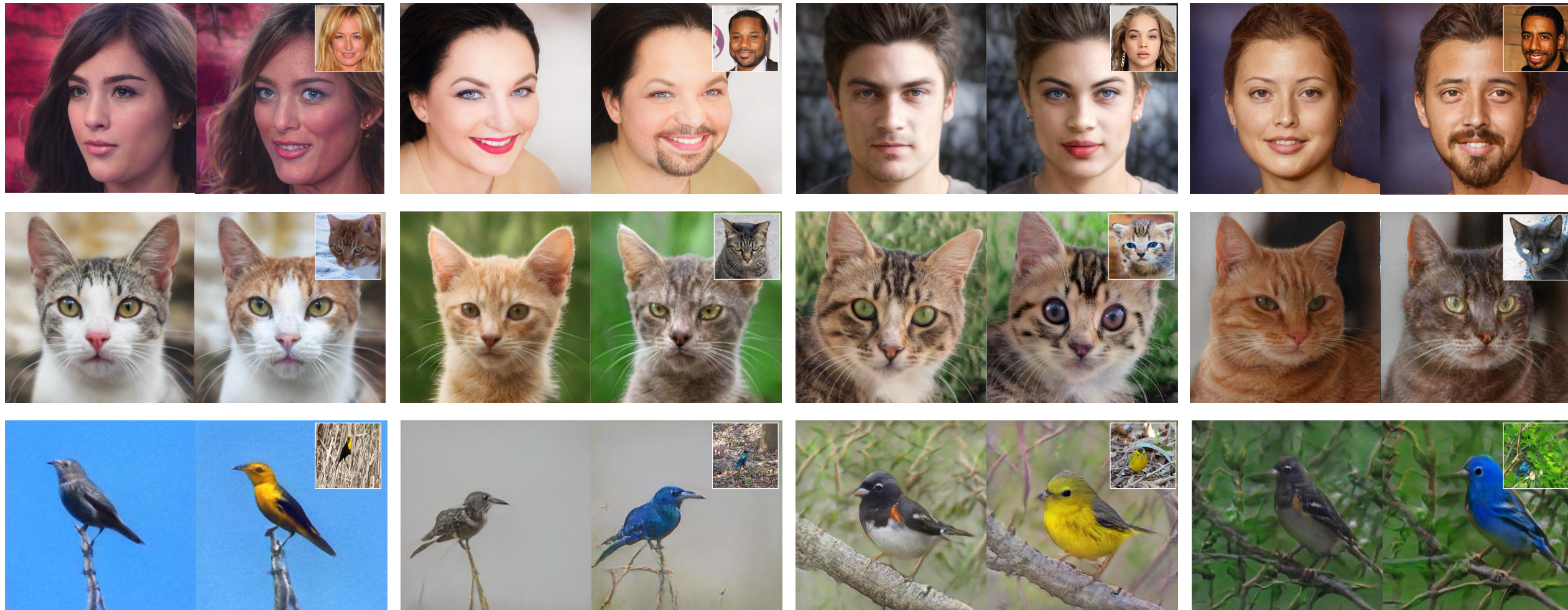


This small bird has a blue top and tail, and a small, straight beak that is pointed. ← → This bird is gray, yellow, and orange in color, with a light colored beak.



This bird has a small head, a light grey belly and brown wings with long skinny black tarsus. ← → This bird has wings that are blue and has a white belly.

# CLIPInverter





Original

TediGAN-B

StyleCLIP-LO

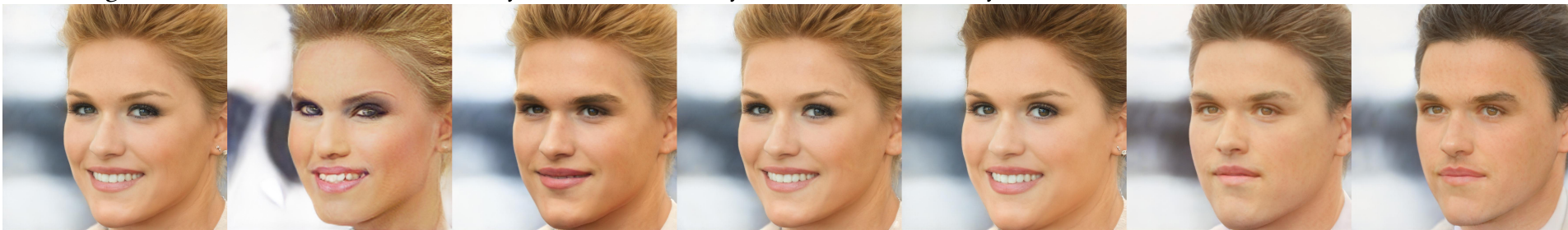
StyleCLIP-GD

StyleMC

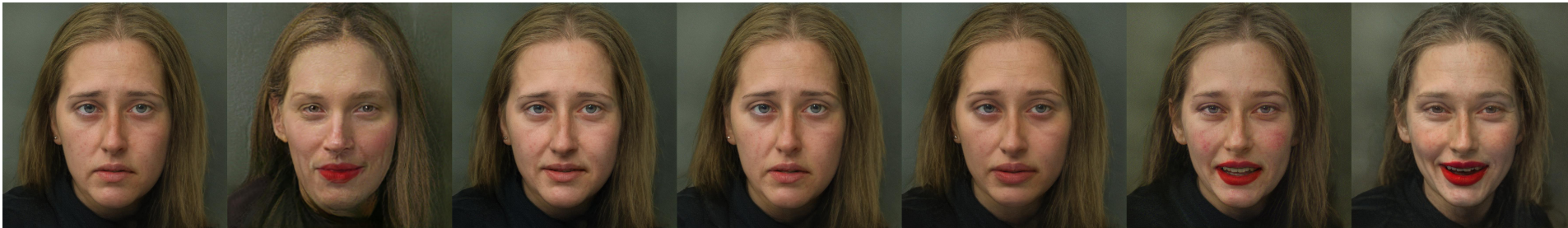
HairCLIP

Ours

# CLIPInverter



This person has bushy eyebrows, pointy nose, black hair, bags under eyes, and big lips. He wears necktie.



She has mouth slightly open, and high cheekbones. She is smiling and is wearing lipstick.



The woman has arched eyebrows, and blond hair and is wearing heavy makeup, and lipstick. She is attractive, and young.



He wears necktie. He has bushy eyebrows, mouth slightly open, bags under eyes, big nose, and high cheekbones. He is smiling. (Baykal et al., ACM TOG 2023) 143

# CLIPInverter

Original

TediGAN-B

StyleCLIP-LO

HairCLIP

Ours



This particular bird has a belly that is gray and white.



This bird has a yellow head with brown and cream colored body.



This is a small yellow bird with greenish wings and a small pointed beak. (Baykal et al., ACM TOG 2023)

# CLIPInverter

Original

TediGAN-B

StyleCLIP-LO

HairCLIP

Ours



A fearful elderly cat.

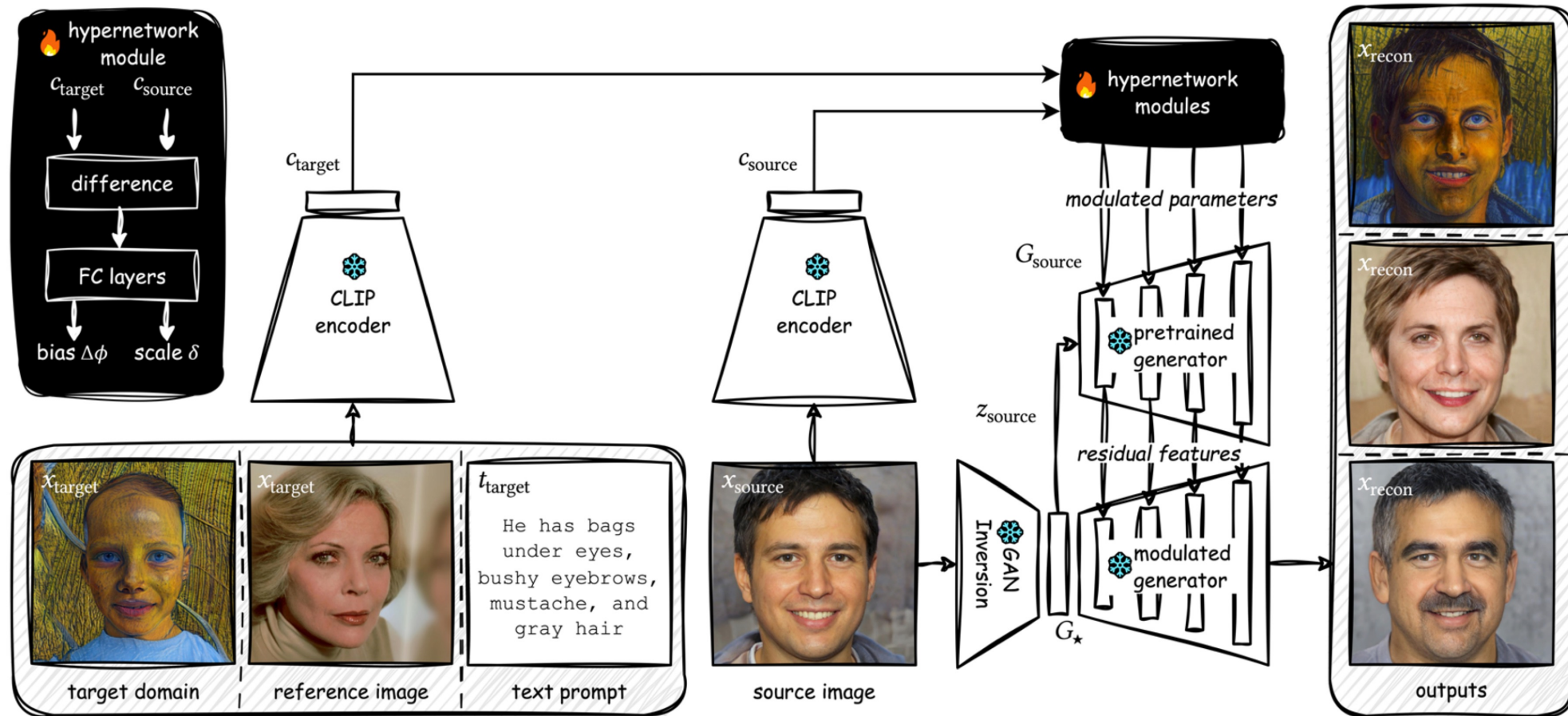


A cat with ginger hair.



A kitten with white hair.

# HyperGAN-CLIP

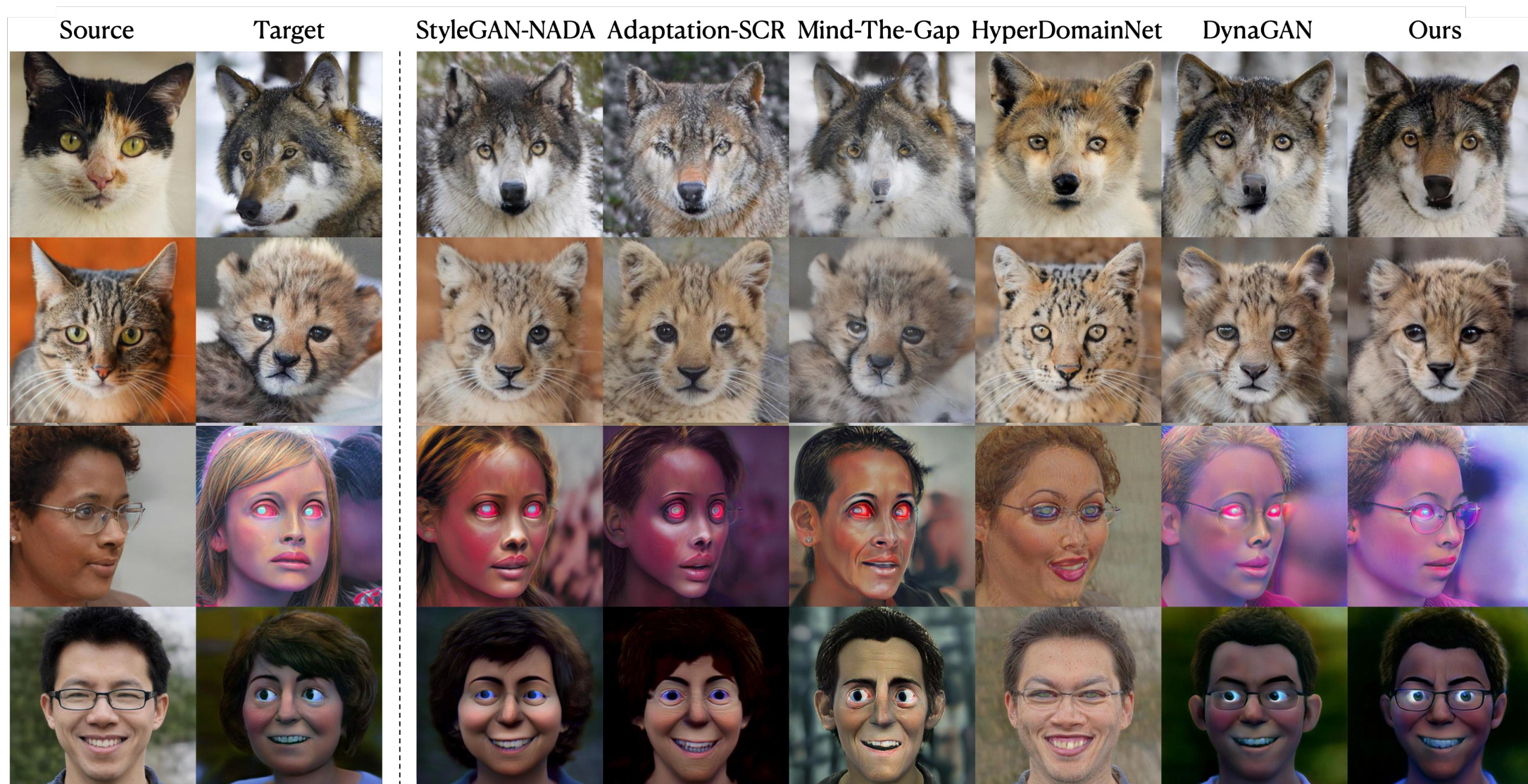


- a flexible framework that is capable of handling domain adaptation, reference-guided image synthesis and text-guided image manipulation.

# HyperGAN-CLIP – Domain Adaptation



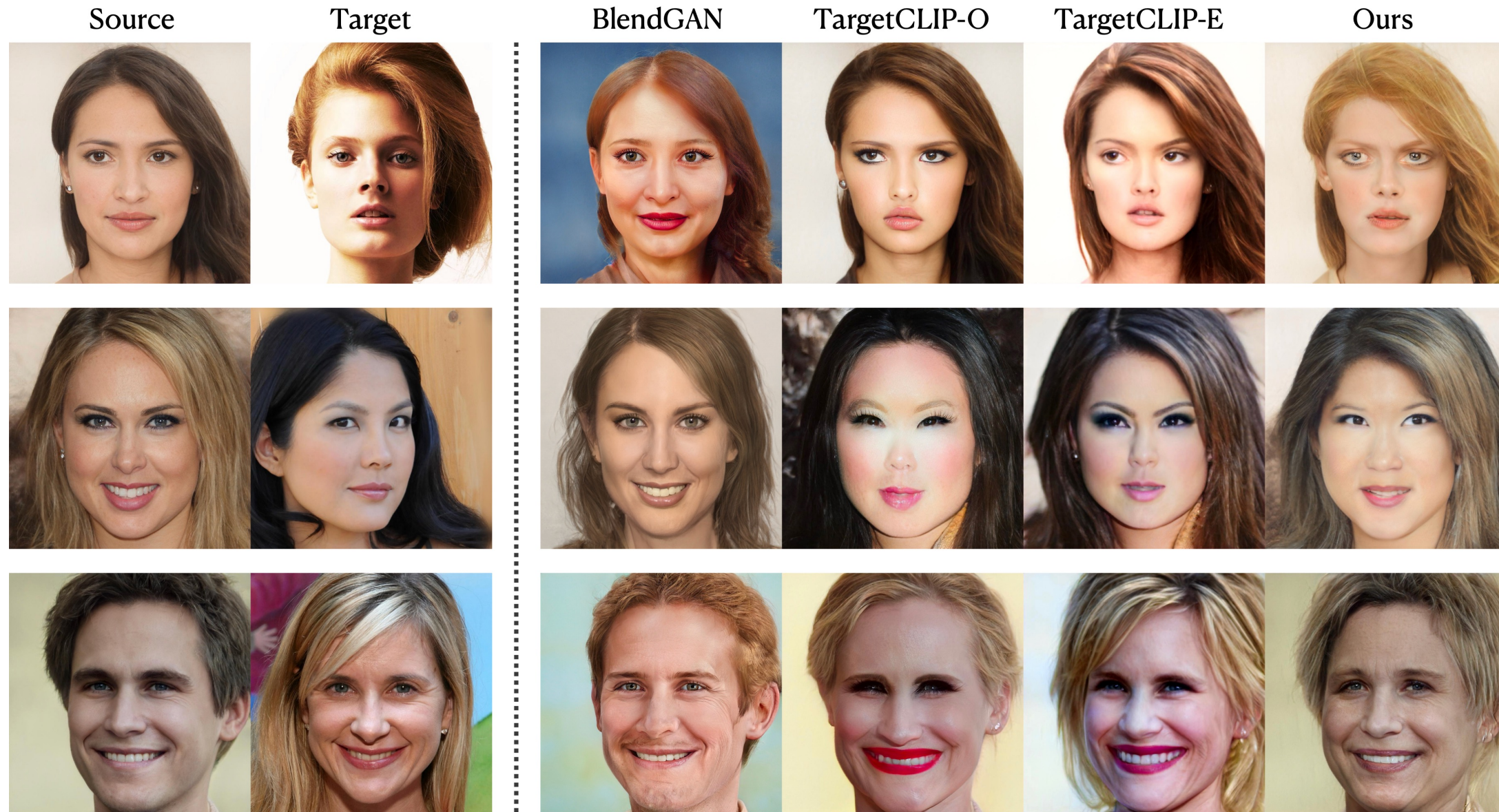
# HyperGAN-CLIP – Domain Adaptation



# HyperGAN-CLIP – Reference-Guided Image Synthesis

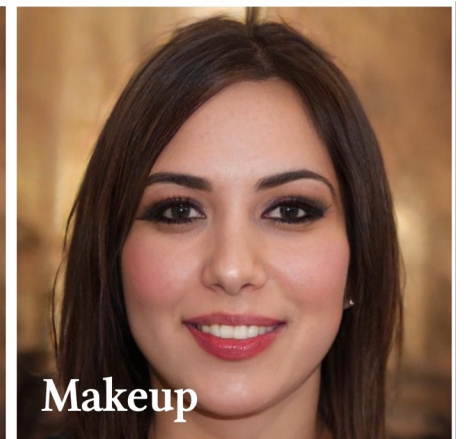
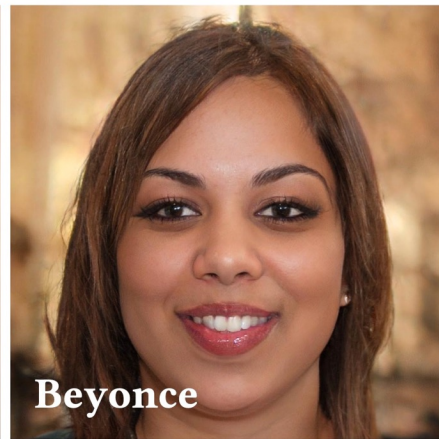
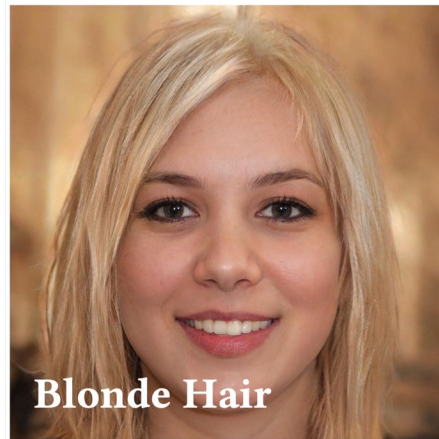


# HyperGAN-CLIP – Reference-Guided Image Synthesis





# HyperGAN-CLIP – Text-Guided Image Manipulation



# HyperGAN-CLIP – Text-Guided Image Manipulation

Source    TediGAN-B    StyleCLIP-LO    StyleCLIP-GD    HairCLIP    CLIPInverter    DiffusionCLIP    Plug-and-play    DeltaEdit    Ours



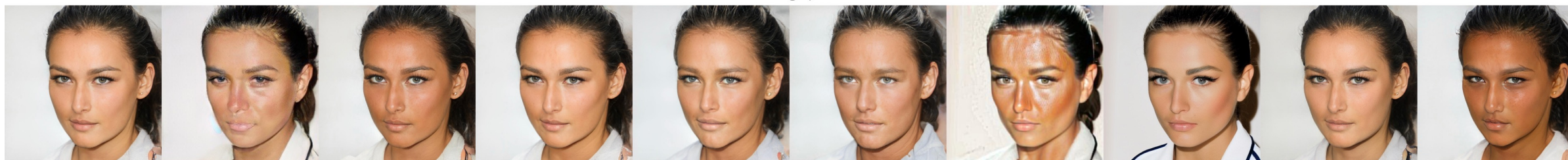
*She wears earrings, lipstick. She has high cheekbones, big lips, and bangs. She is smiling.*



*The person has black hair, and wavy hair.*



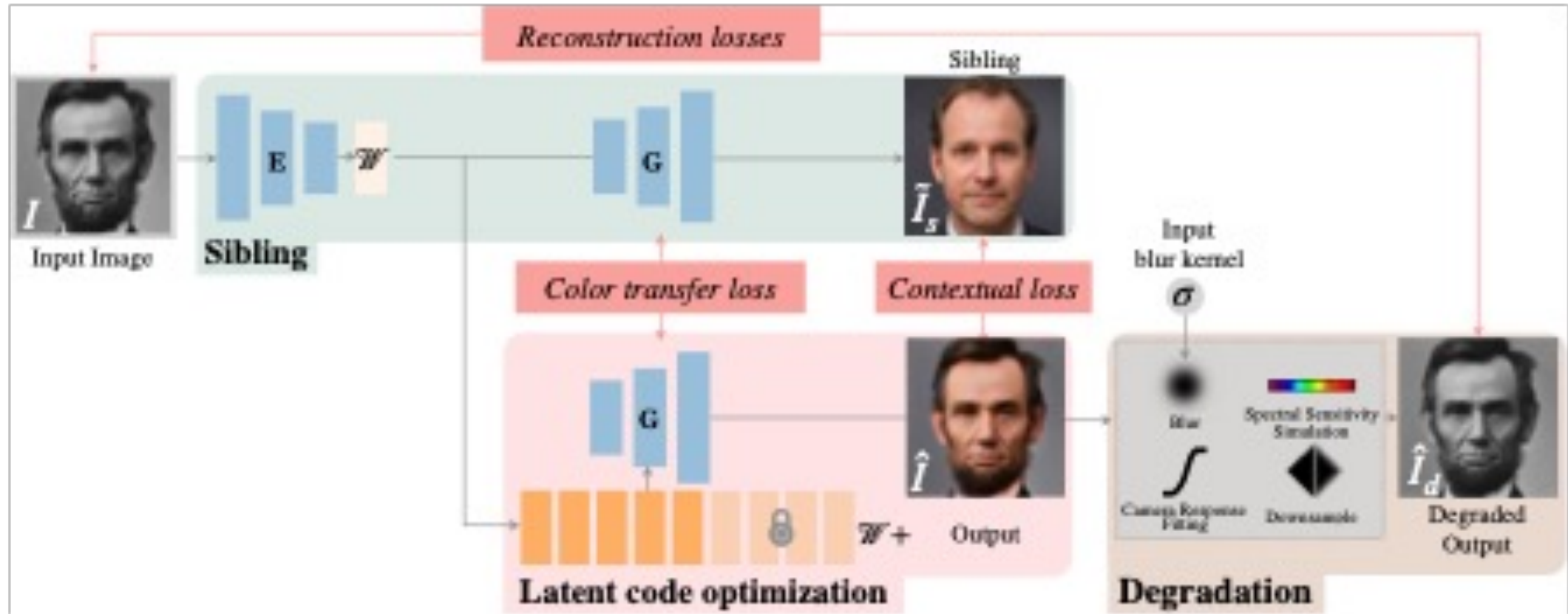
*Angry*



*Tanned*

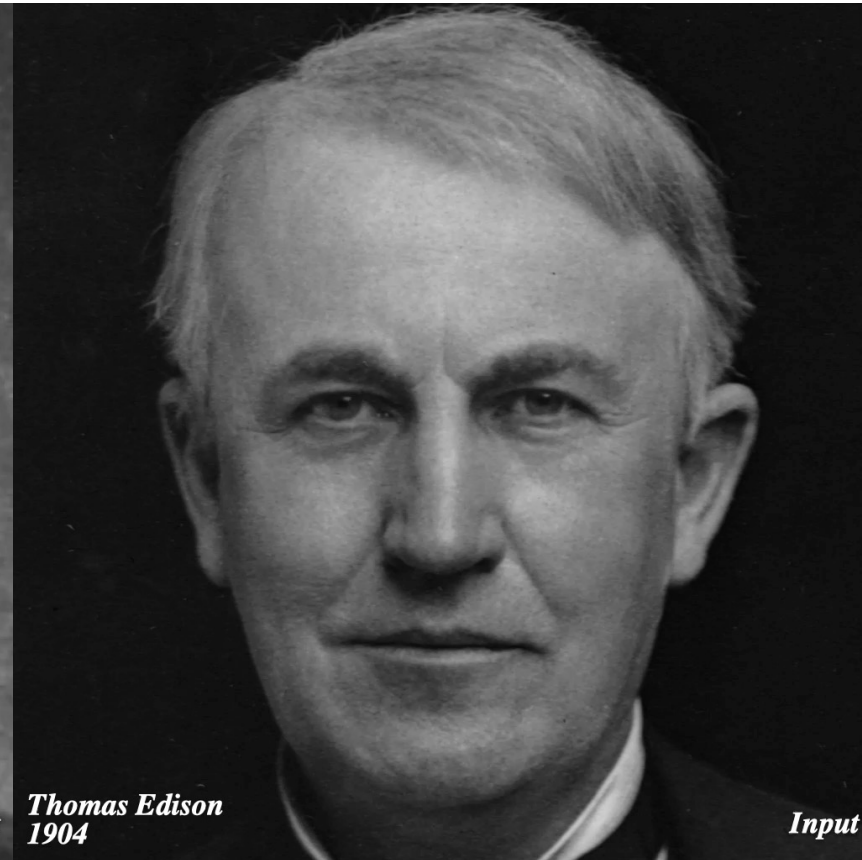
(Anees et al., SIGGRAPH Asia 2024)

# Time-travel Rephotography



- Key idea: Use the StyleGAN2 framework to project old photos into the space of modern high-resolution photos for enhancing their quality.

# Time-travel Rephotography



# Time-travel Rephotography



Input

Ours

DeOldify

InstColorization

Zhang

Zhang (FFHQ)

**Next Lecture:**  
Deep Generative Models  
Part 3