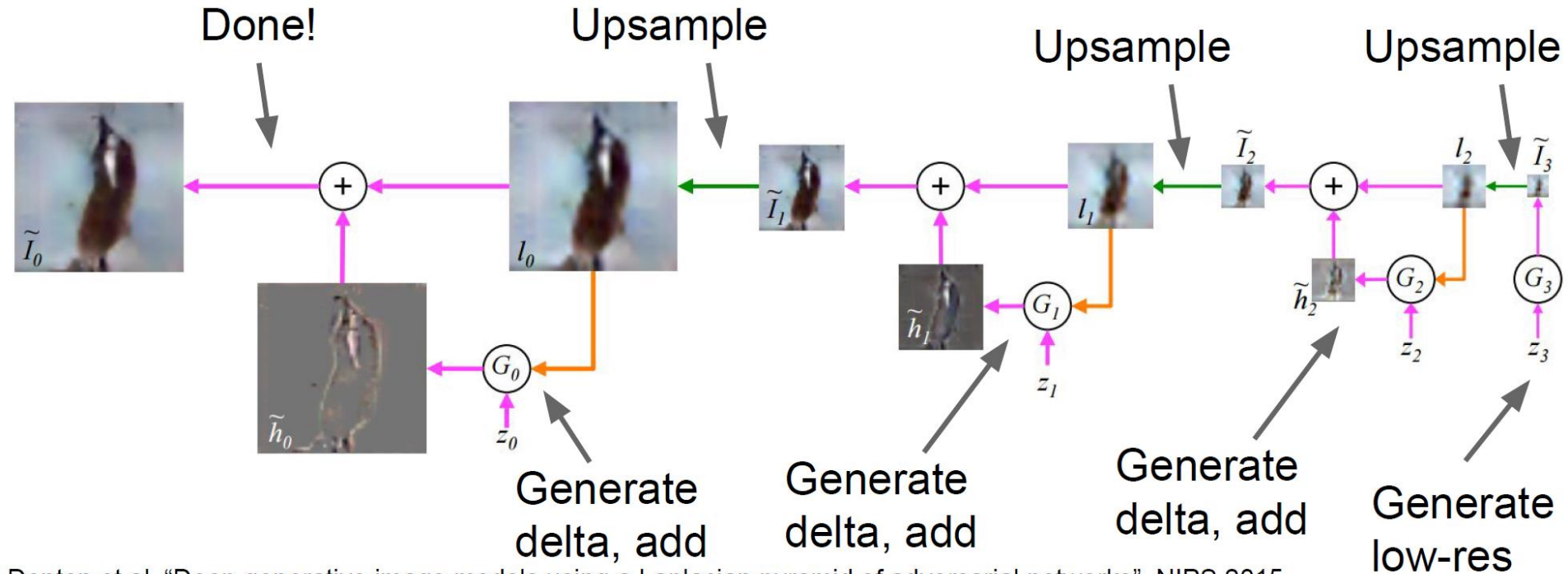


# GAN Architectures and Conditional GANs

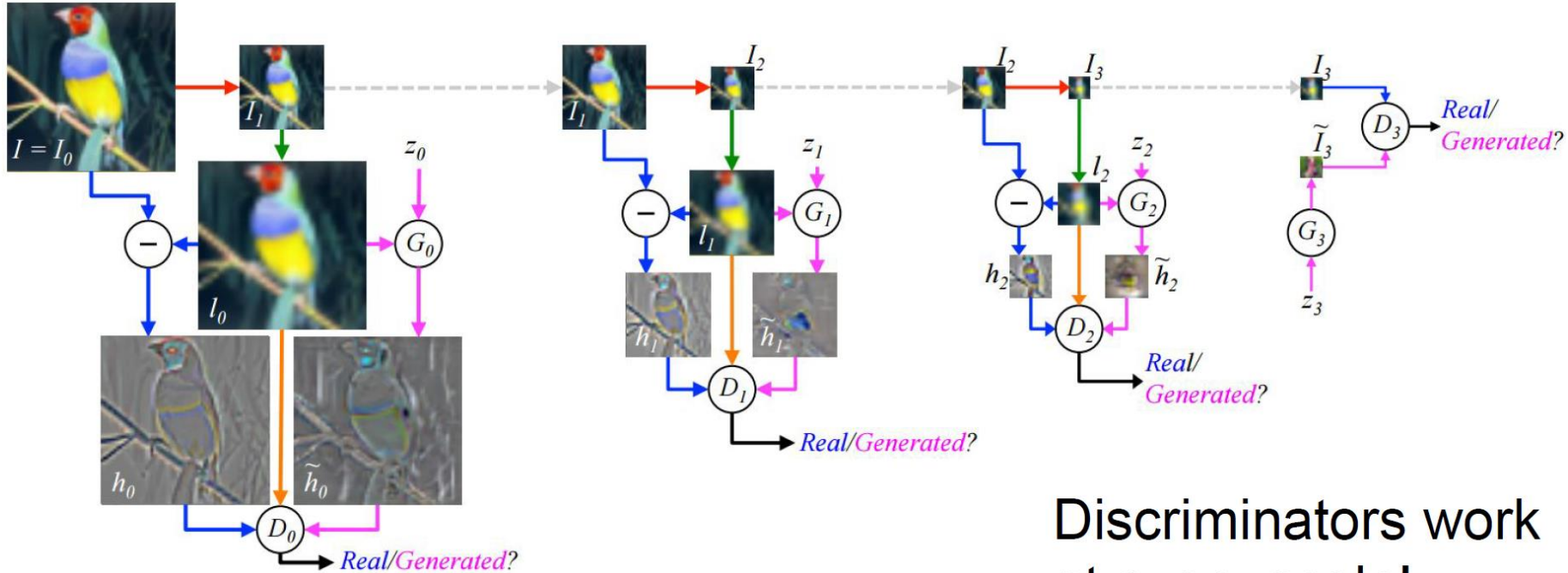
# GAN Architectures

# Multiscale GANs



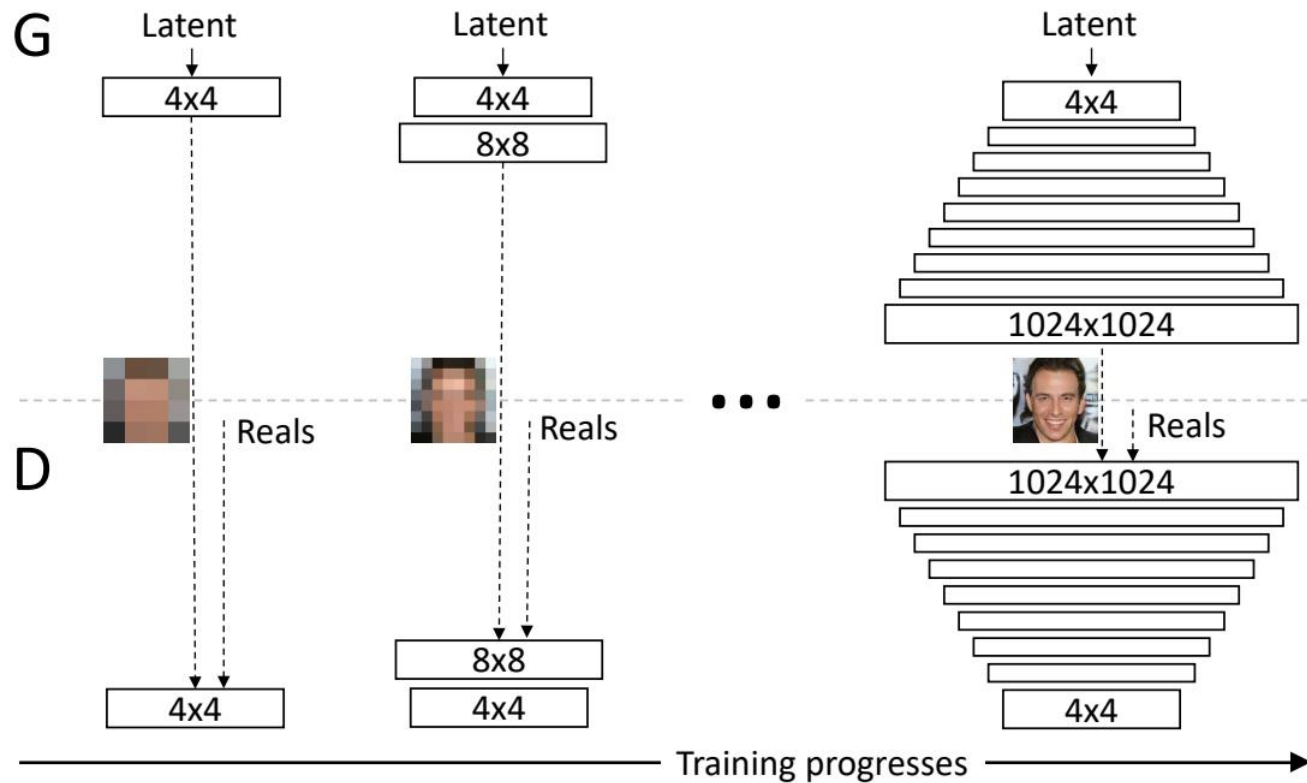
Denton et al, "Deep generative image models using a Laplacian pyramid of adversarial networks", NIPS 2015

# Multiscale GANs



Denton et al, NIPS 2015

# Progressive Growing GANs



G

4×4
4×4

2x
8×8
8×8

2x
16×16
16×16

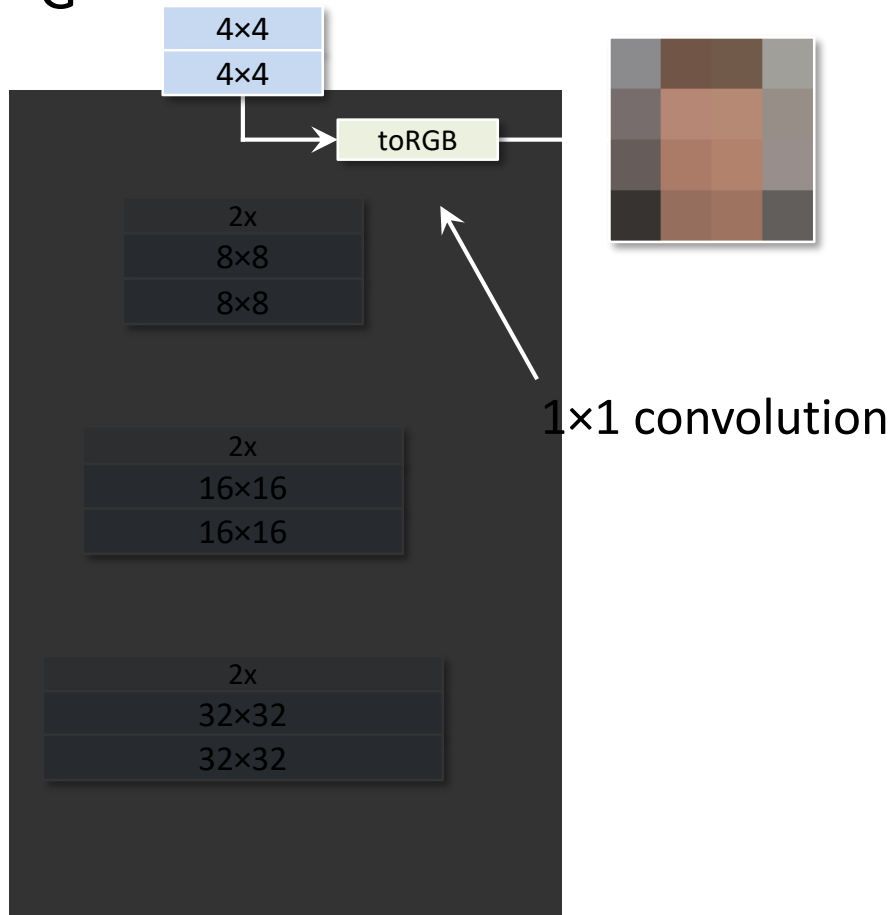
2x
32×32
32×32

Replicated block

Nearest-neighbor upsampling

3×3 convolution

G



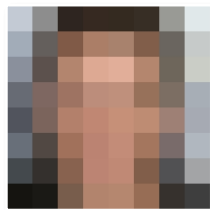
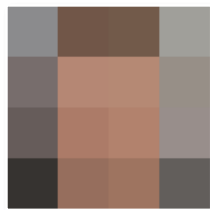
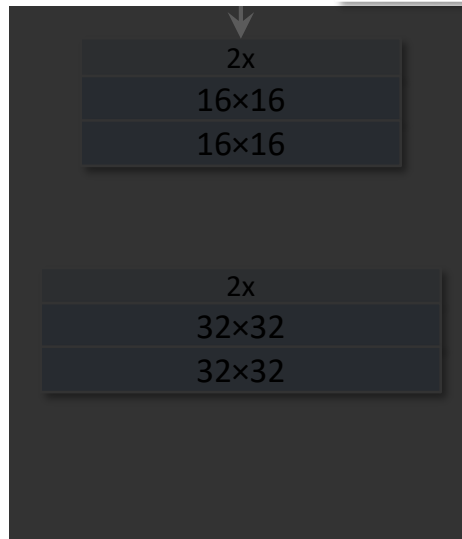
G

4×4  
4×4

toRGB

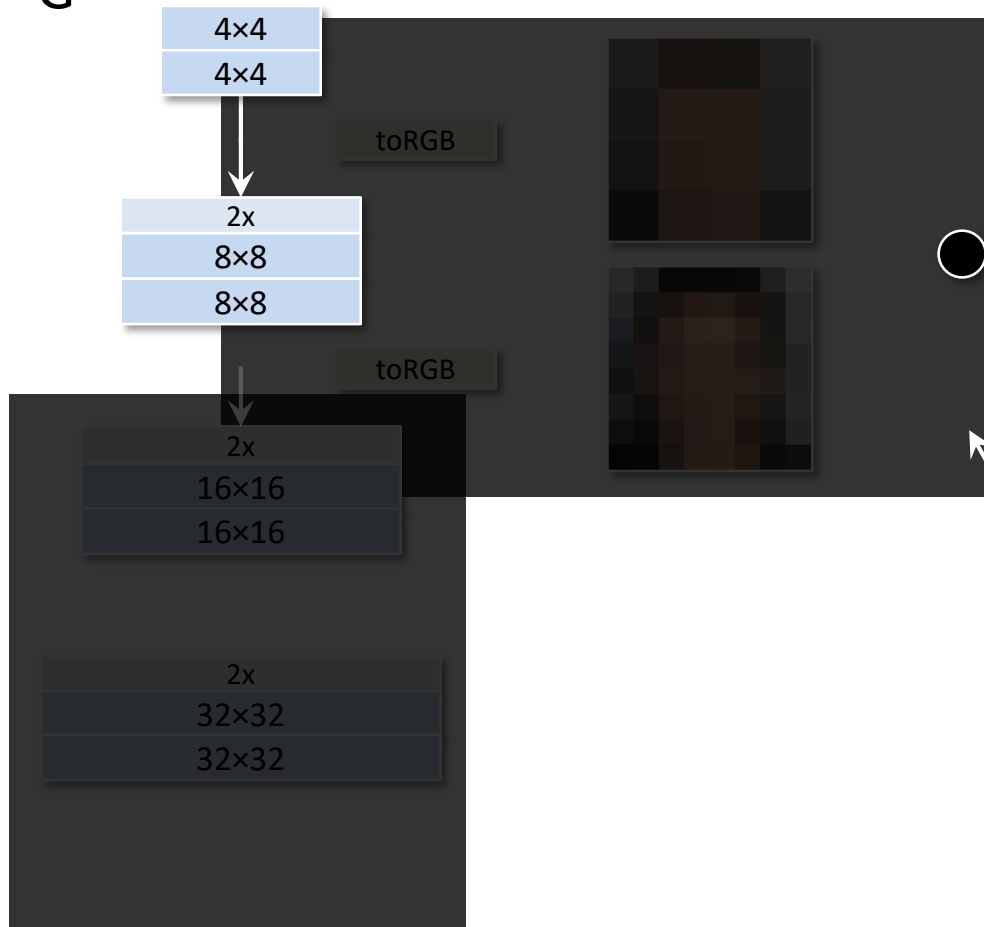
2x  
8×8  
8×8

toRGB



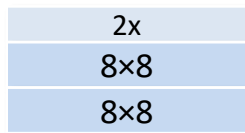
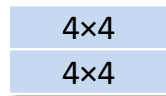


G

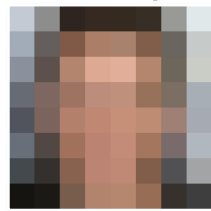


Linear crossfade

G



toRGB

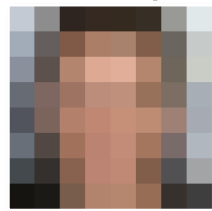
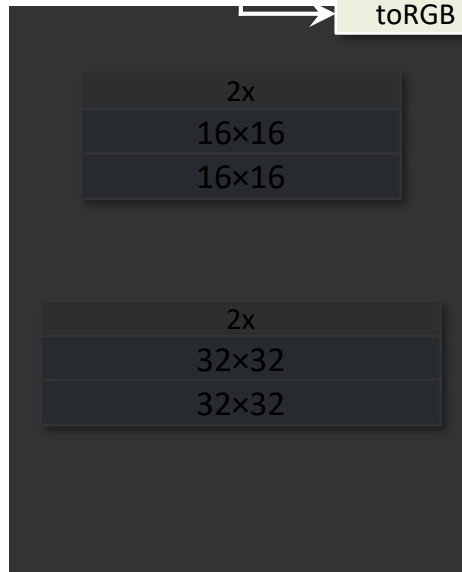


G

4×4
4×4

2x
8×8
8×8

toRGB



D

32×32
32×32
0.5x

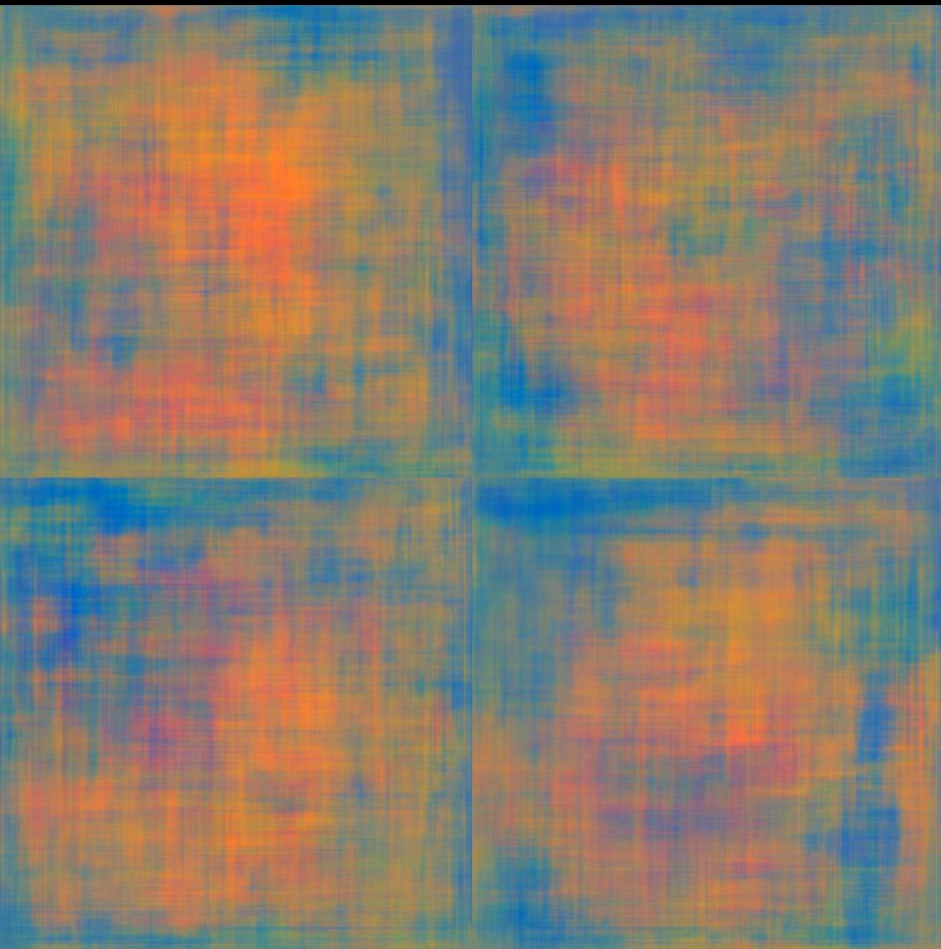
16×16
16×16
0.5x

fromRGB

8×8
8×8
0.5x

4×4
4×4

5 min 00 sec



Fixed resolution



Progressive growing

# Progressive Growing GANs

CelebA-HQ

$1024 \times 1024$

Latent space interpolations

# Lots of GAN Variations

- Hundreds of GAN papers in the last two years
  - > Mostly with different losses
  - > Extremely hard to train and evaluate

## Are GANs Created Equal? A Large-Scale Study

Mario Lucic\*   Karol Kurach\*   Marcin Michalski   Sylvain Gelly   Olivier Bousquet  
Google Brain

### Abstract

*Generative adversarial networks (GAN) are a powerful subclass of generative models. Despite a very rich research*

GAN algorithm(s) perform objectively better than the others. That's partially due to the lack of robust and consistent metric, as well as limited comparisons which put all algorithms on equal footage, including the computational

# Conditional Generative Adversarial Networks (cGANs)

# Conditional GANs (cGANs)

- Gain control of output
- Modeling (e.g., sketch-based modeling, etc.)
  - Add semantic meaning to latent space manifold
- Domain transfer
  - Labels on A -> transfer to B, train network on 'B', test on B
  - More later



# GAN Manifold



## Train Data



Sampled Data  $\rightarrow G(z)$

# GAN Manifold

a



b

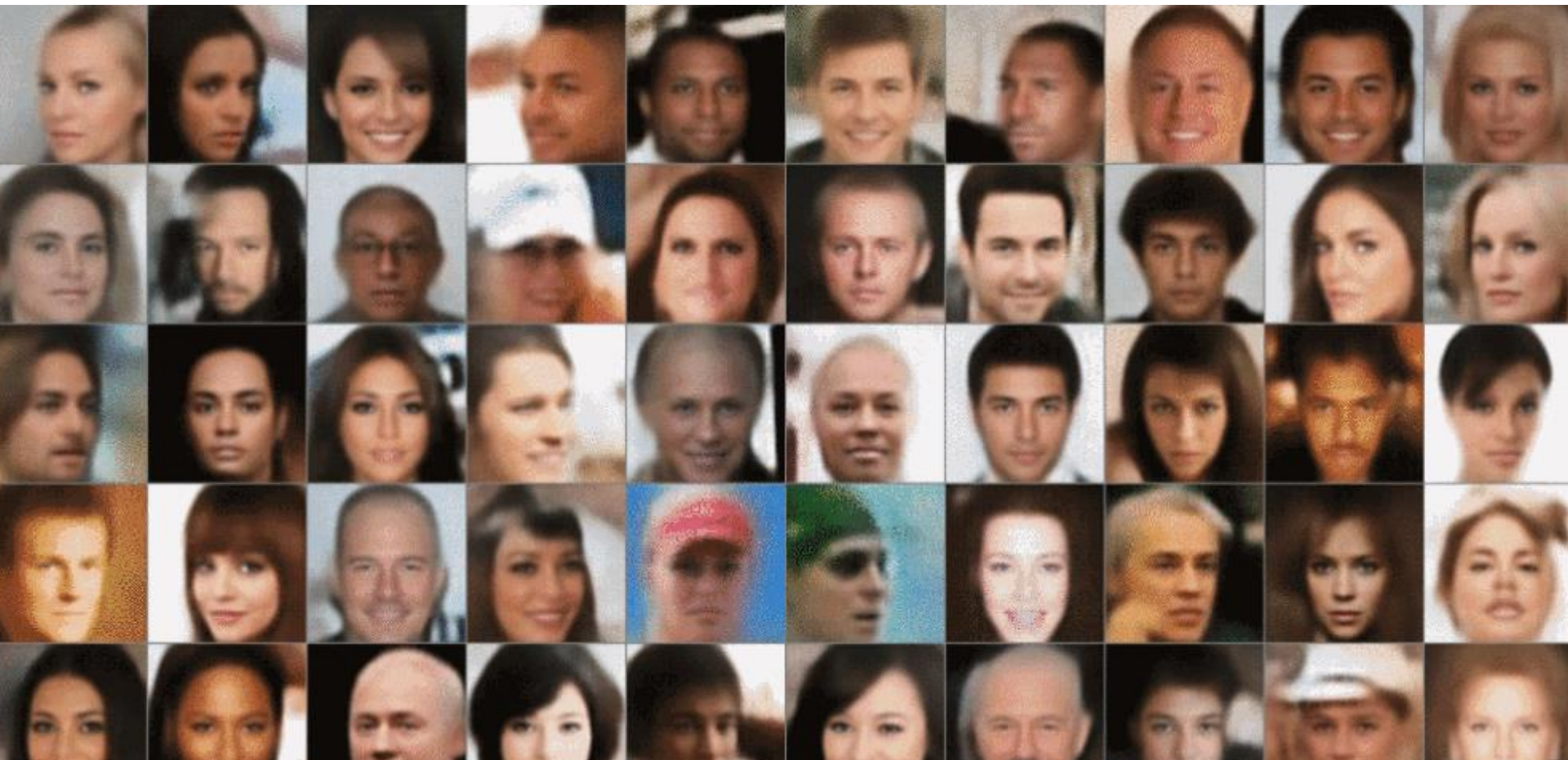


c



$a - b + c$

# GAN Manifold





# GAN Manifold

$G(z_0)$



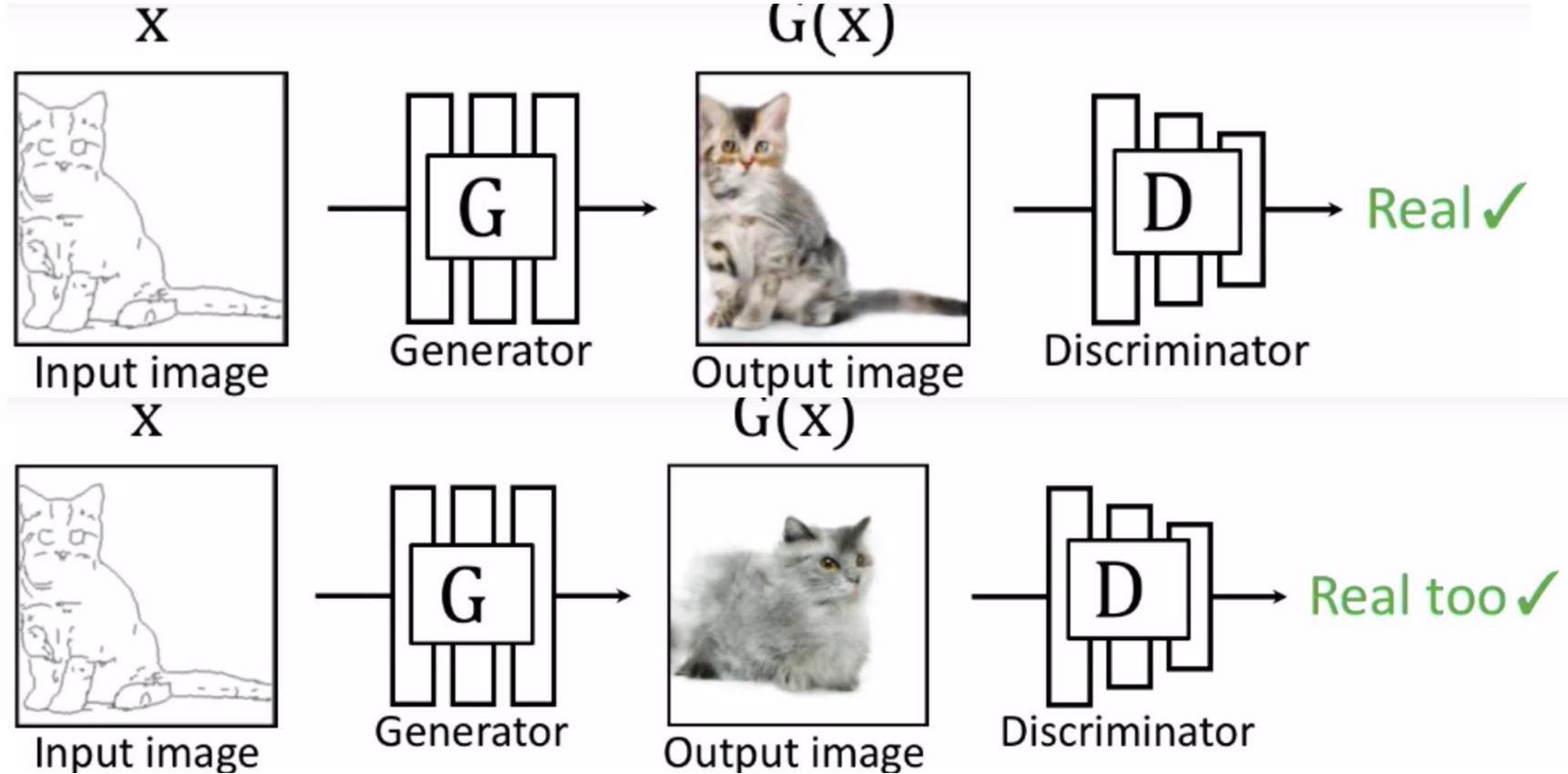
Linear interpolation in  $z$  space:  $G(z_0 + t \cdot (z_1 - z_0))$



$G(z_1)$



# Conditional GANs (cGANs)



# iGANs: Overview



original photo



projection on manifold



Editing UI



different degree of image manipulation



Edit Transfer



transition between the original and edited projection

# iGANs: Overview



original photo



projection on manifold



different degree of image manipulation



Edit Transfer



transition between the original and edited projection

# iGANs: Projecting an Image onto the Manifold

Input: real image  $x^R$   
Output: latent vector  $z$

**Optimization**

$$z^* = \arg \min \mathcal{L}(G(z), x^R)$$

Reconstruction loss  $L$

Generative model  $G(z)$



0.196



0.238



0.332



# iGANs: Projecting an Image onto the Manifold

Input: real image  $x^R$   
Output: latent vector  $z$

## Optimization

$$z^* = \arg \min \mathcal{L}(G(z), x^R)$$

Inverting Network  $z = P(x)$

$$\theta_P^* = \arg \min_{\theta_P} \sum_{x_n^R} \mathcal{L}(\underbrace{G(P(x_n^R; \theta_P))}_{\text{Auto-encoder}}, x_n^R)$$

Auto-encoder  
*with a fixed decoder G*



# iGANs: Projecting an Image onto the Manifold

Input: real image  $x^R$

Output: latent vector  $z$

## Optimization

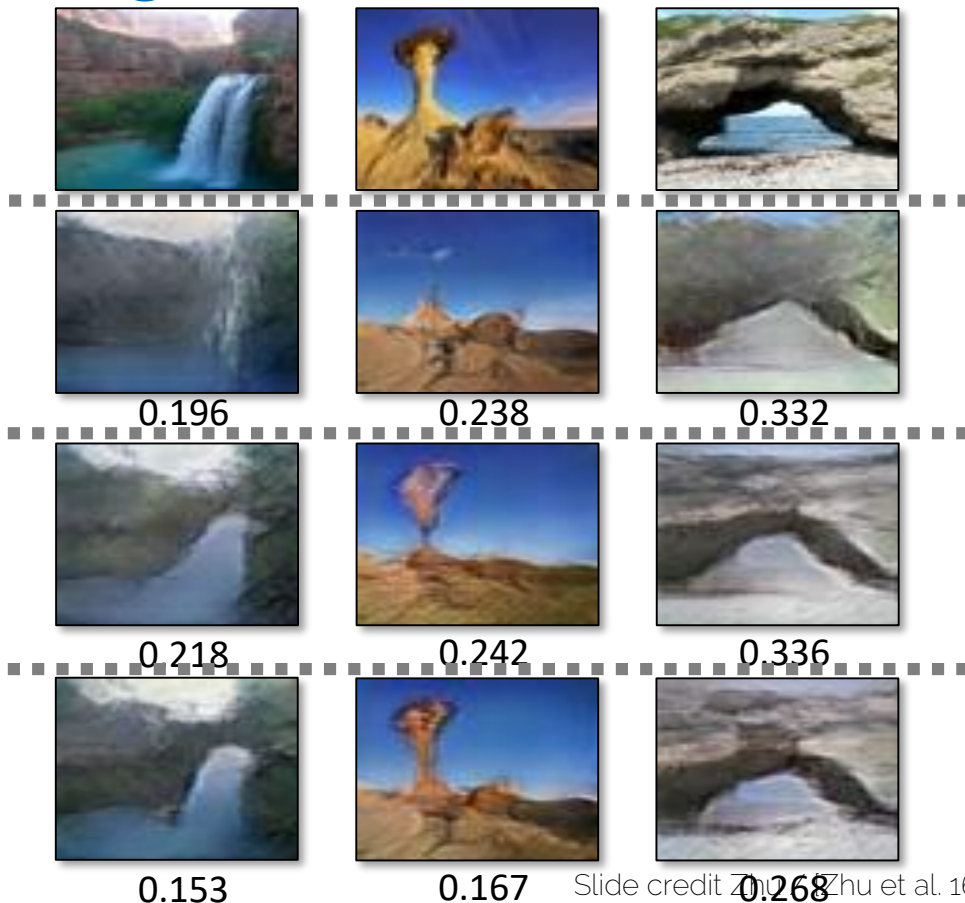
$$z^* = \arg \min \mathcal{L}(G(z), x^R)$$

Inverting Network  $z = P(x)$

$$\theta_P^* = \arg \min_{\theta_P} \sum_{x_n^R} \mathcal{L}(G(P(x_n^R; \theta_P)), x_n^R)$$

## Hybrid Method

Use the **network** as initialization  
for the **optimization** problem



# iGANs: Overview



original photo



projection on manifold



Editing UI



different degree of image manipulation



Edit Transfer



transition between the original and edited projection

# iGANs: Manipulating the Latent Vector

constraint violation loss  $L_g$

user guidance image

Objective:  $z^* = \arg \min_{z \in \mathbb{Z}} \left\{ \underbrace{\sum_g (\mathcal{L}_g(G(z)) \underbrace{v_g}_{\text{data term}})}_{\text{data term}} + \underbrace{\lambda_s \cdot \|z - z_0\|_2^2}_{\text{manifold smoothness}} \right\}.$

Guidance

$v_g$



$z_0$

# iGANs: Overview



original photo



projection on manifold



Editing UI



different degree of image manipulation



Edit Transfer

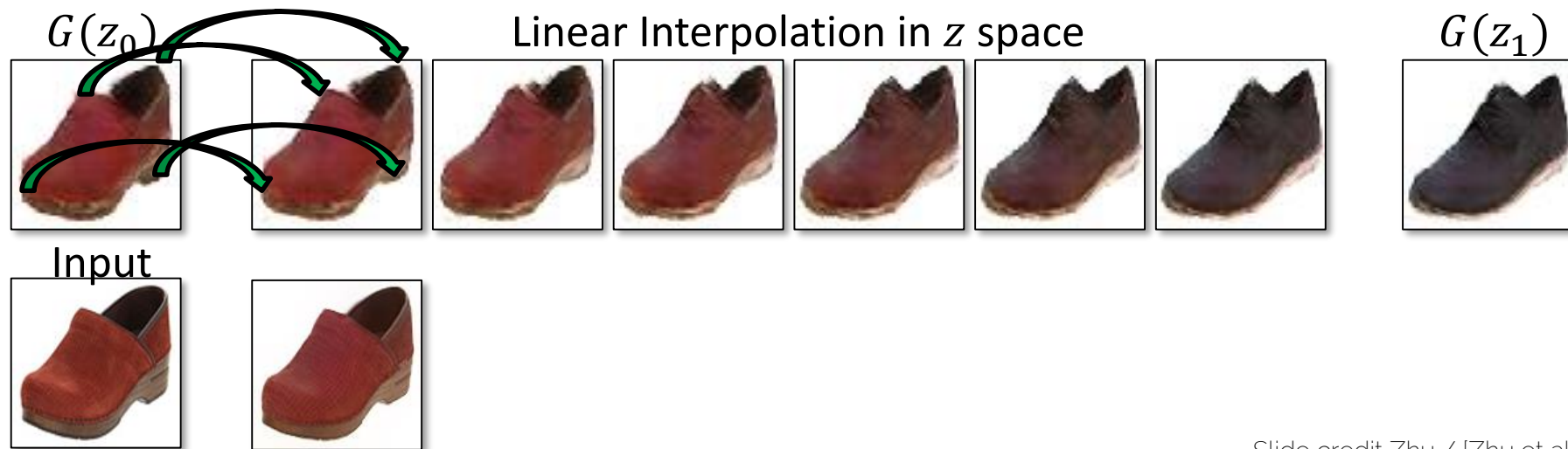


transition between the original and edited projection

# iGANs: Edit Transfer

**Motion** ( $\mathbf{u}, \mathbf{v}$ ) + **Color** ( $A_{3 \times 4}$ ): estimate per-pixel geometric and color variation

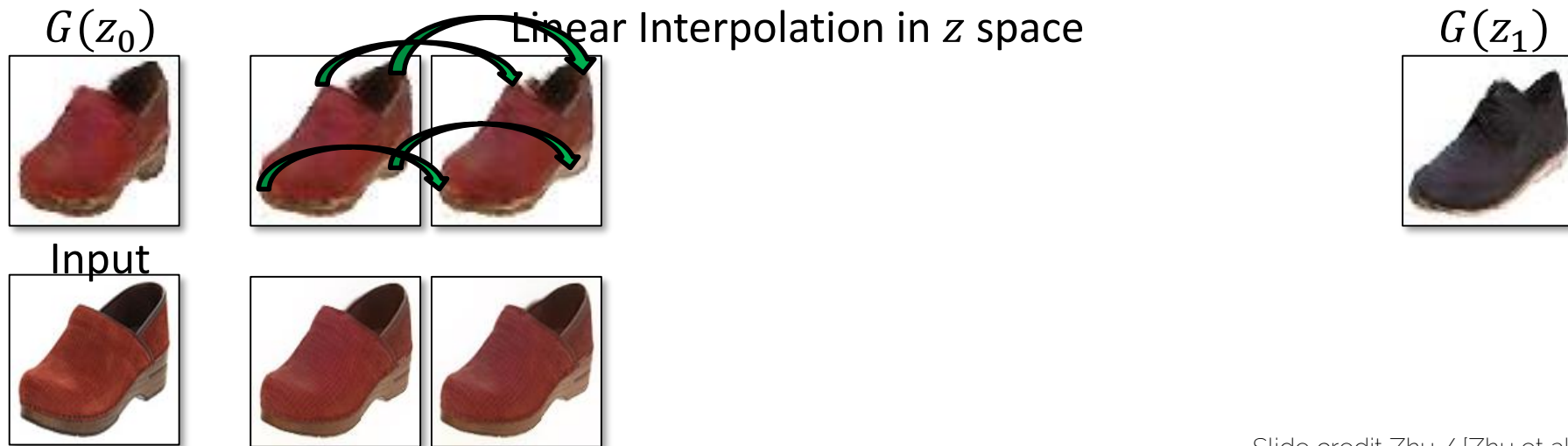
$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x+u, y+v, t+1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$



# iGANs: Edit Transfer

**Motion** ( $\mathbf{u}, \mathbf{v}$ ) + **Color** ( $A_{3 \times 4}$ ): estimate per-pixel geometric and color variation

$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x+u, y+v, t+1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$





# iGANs: Edit Transfer

**Motion (u, v) + Color ( $A_{3 \times 4}$ ):** estimate per-pixel geometric and color variation

$$\iint \underbrace{\|I(x, y, t) - A \cdot I(x+u, y+v, t+1)\|^2}_{\text{data term}} + \underbrace{\sigma_s (\|\nabla u\|^2 + \|\nabla v\|^2)}_{\text{spatial reg}} + \underbrace{\sigma_c \|\nabla A\|^2}_{\text{color reg}} dx dy$$

$G(z_0)$



Linear Interpolation in z space



$G(z_1)$



Input

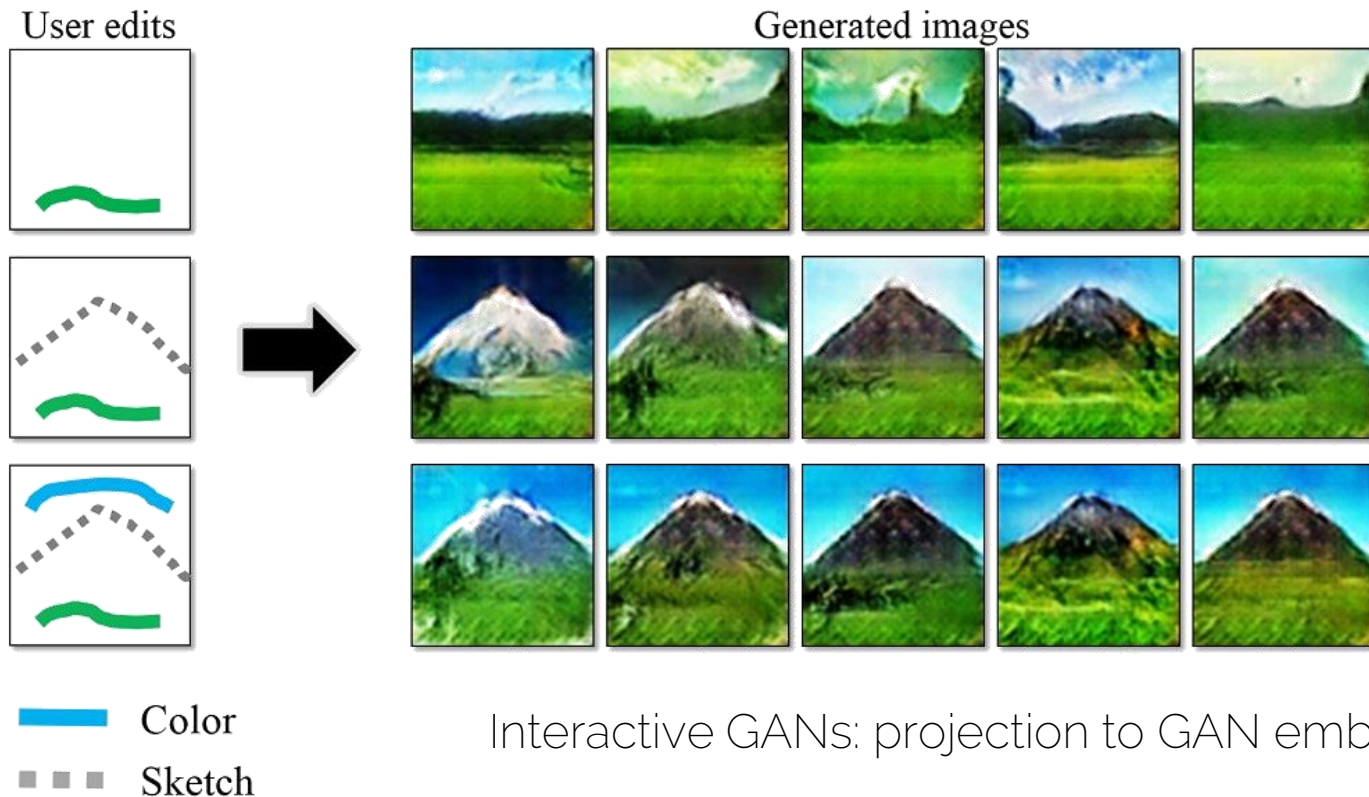


Result





# cGANs: Interactive GANs



# cGANs: Interactive GANs

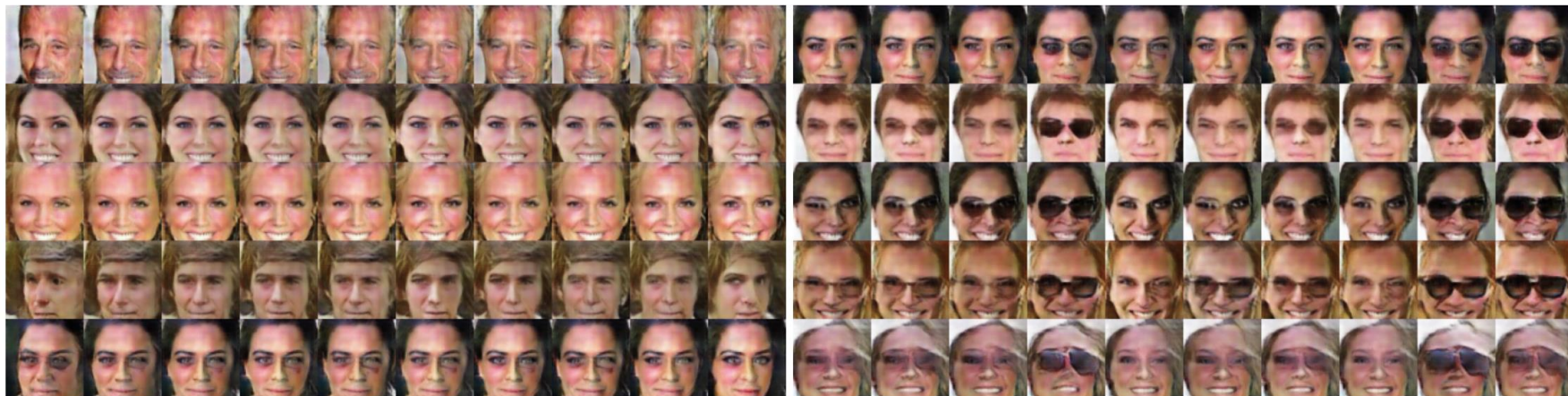
Original photos										
Reconstruction via Optimization										
	0.165	0.164	0.370	0.279	0.350	0.249	0.437	0.255	0.178	0.227
Reconstruction via Network										
	0.198	0.190	0.382	0.302	0.251	0.339	0.482	0.270	0.248	0.263
Reconstruction via Hybrid Method										
	0.133	0.141	0.298	0.218	0.160	0.204	0.318	0.185	0.183	0.190

# cGANs: Interactive GANs



# Mapping in Latent Space is Difficult!

- Semantics are missing
- In most cases, no labels available
- Ideally, need some unsupervised disentangled rep.

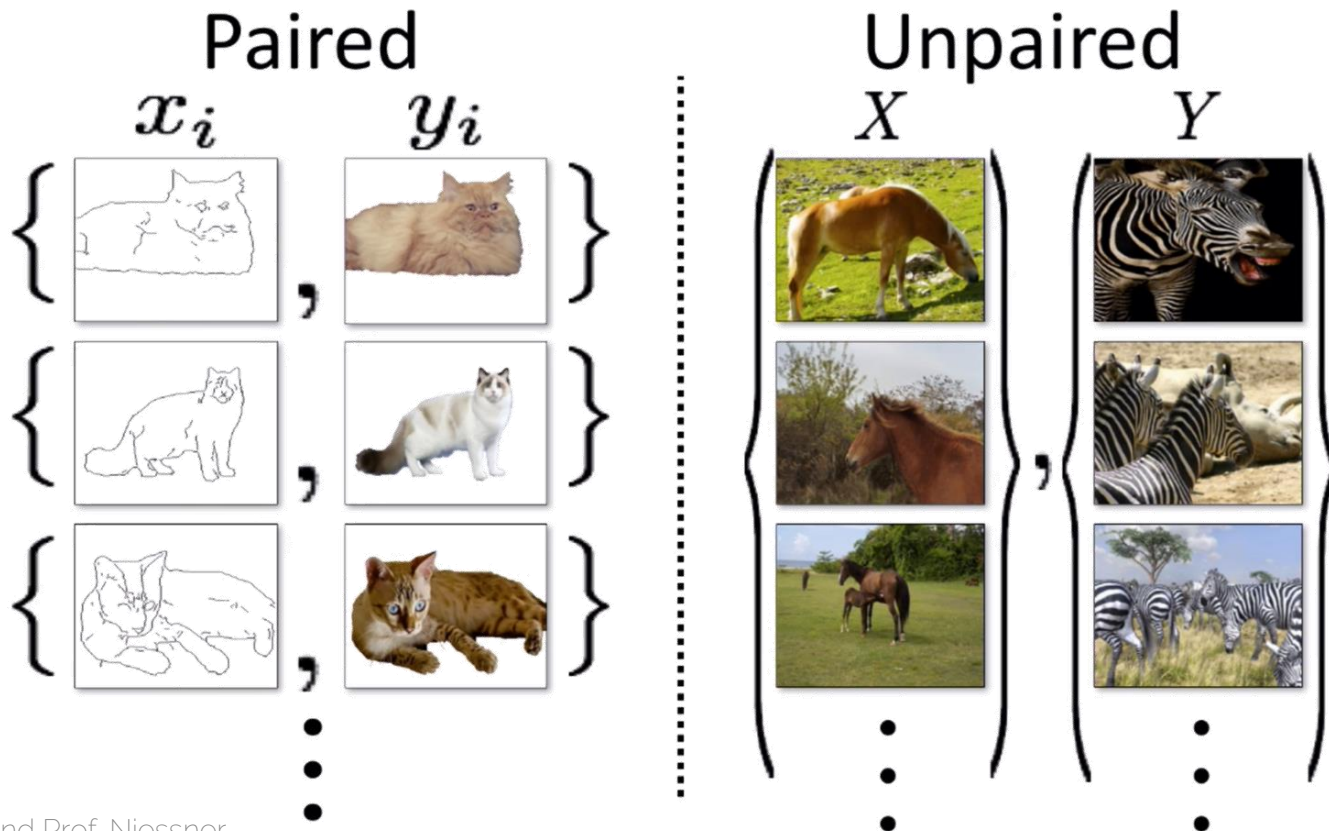


(a) Azimuth (pose)

(b) Presence or absence of glasses



# Paired vs Unpaired Setting



# pix2pix: Image-to-Image Translation

Labels to Street Scene

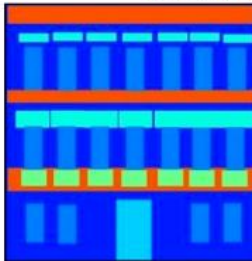


input



output

Labels to Facade



input



output

BW to Color



input



output

Aerial to Map



input

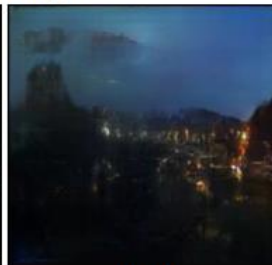


output

Day to Night



input



output

Edges to Photo



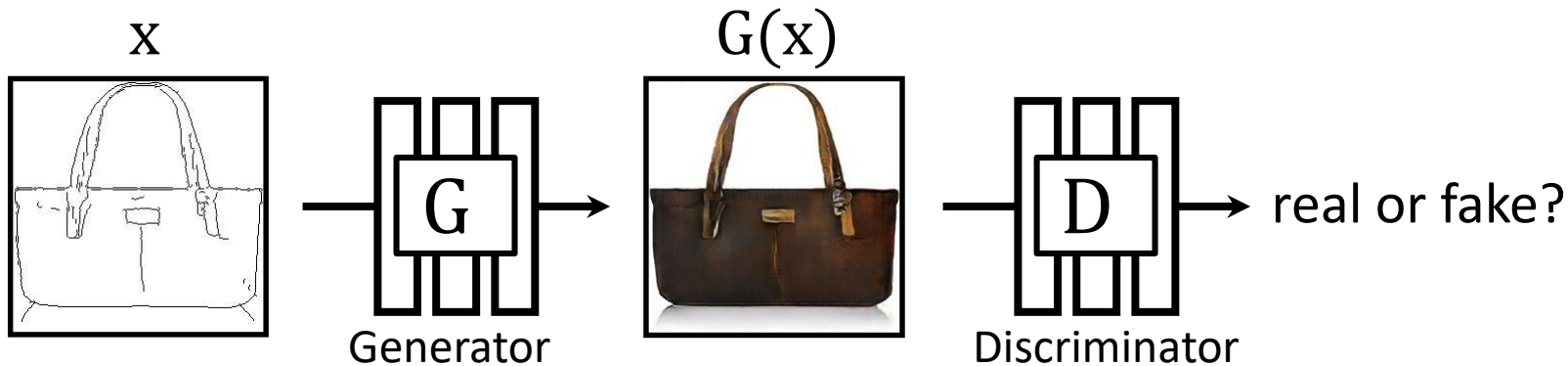
input



output

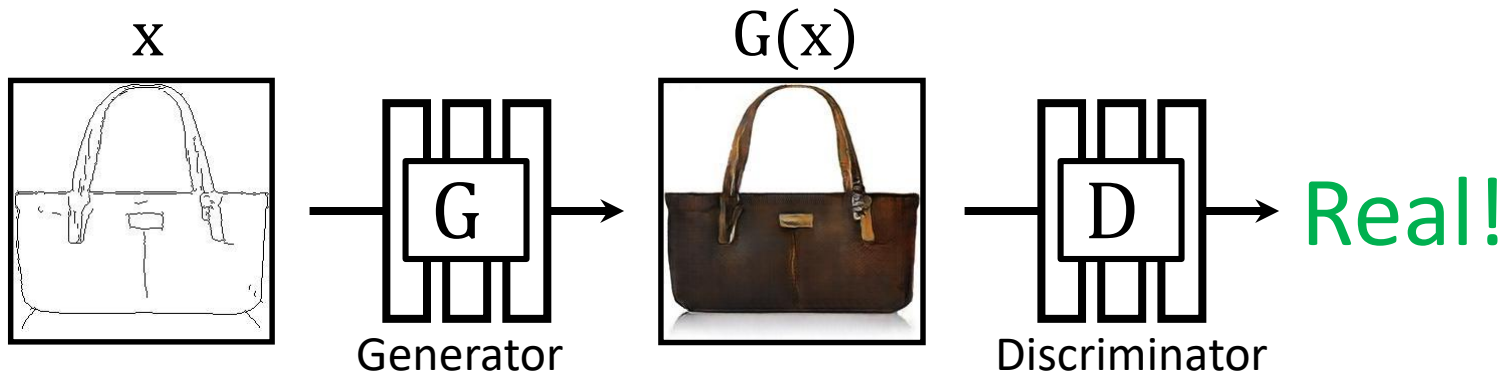


$$\min_G \max_D \mathbb{E}_{z,x} [\log D(G(z)) + \log(1 - D(x))]$$



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

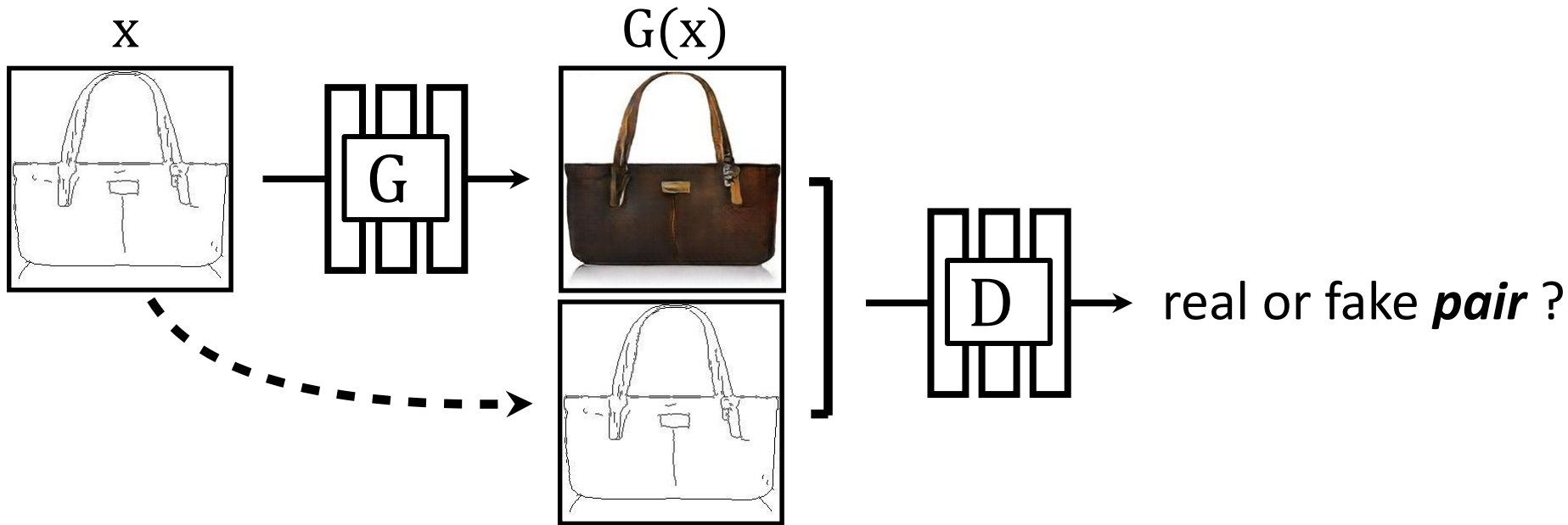




$$\min_G \max_D \mathbb{E}_{x,y} [\log D(G(x)) + \log(1 - D(y))]$$



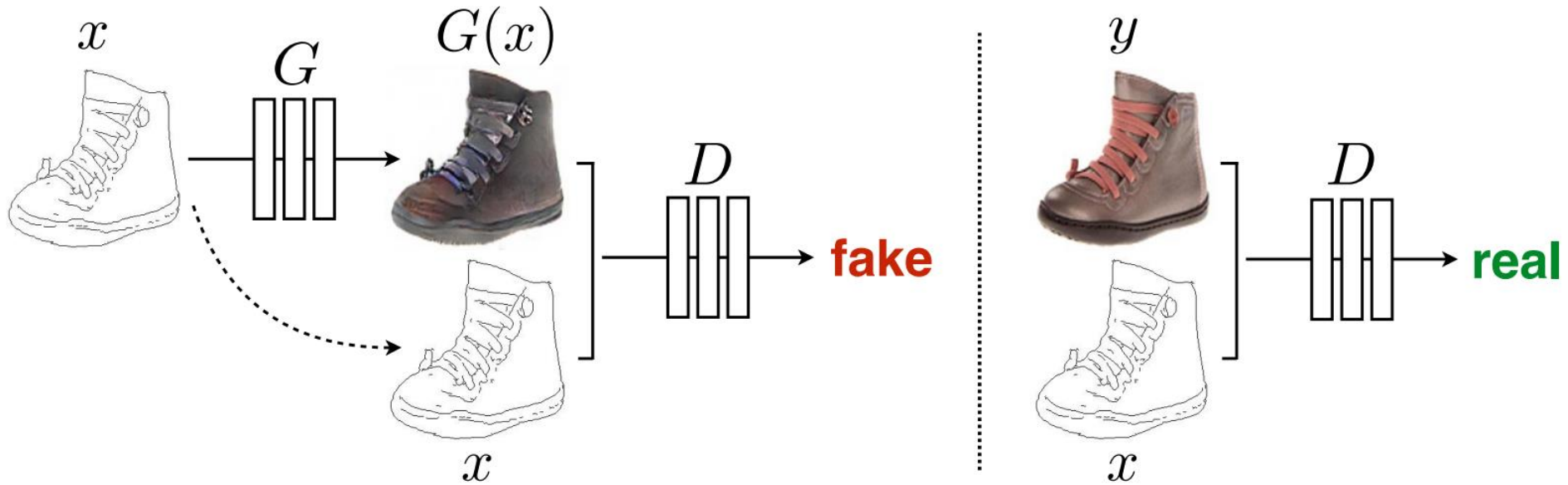
$$\min_G \max_D \mathbb{E}_{x,y} [\log D(G(x)) + \log(1 - D(y))]$$



$$\min_G \max_D \mathbb{E}_{x,y} [\log \underbrace{D(x, G(x))}_{\text{fake pair}} + \log(1 - \underbrace{D(x, y)}_{\text{real pair}})]$$

match joint distribution  $p(G(x), y) \sim p(x, y)$

# pix2pix



# Edges → Images

Input

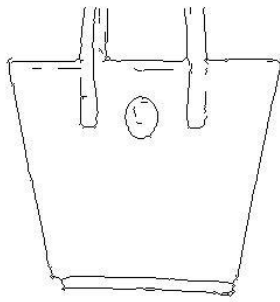
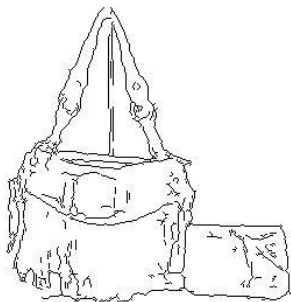
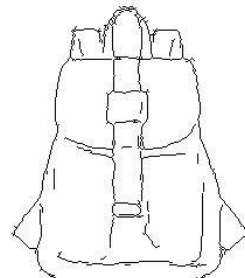
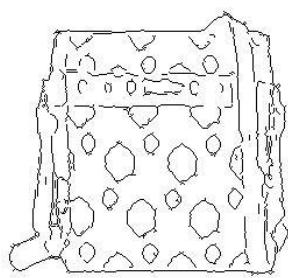
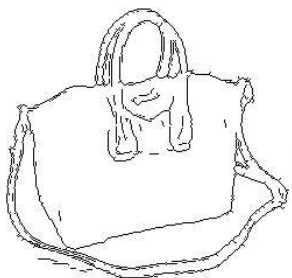
Output

Input

Output

Input

Output

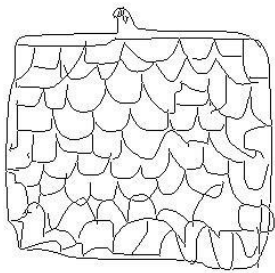


# pix2pix: Paired Setting

- Great when we have 'free' training data
- Often called self-supervised
- Think about these settings 😊

# Sketches $\rightarrow$ Images

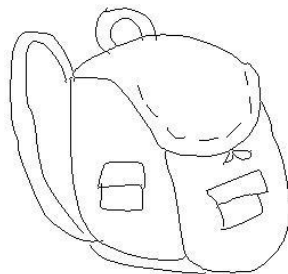
Input



Output



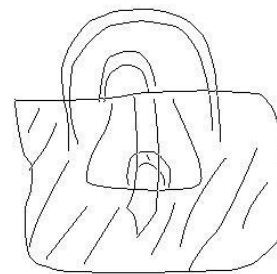
Input



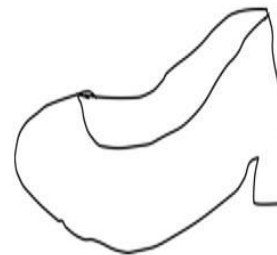
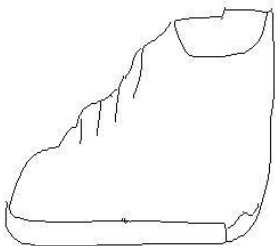
Output



Input

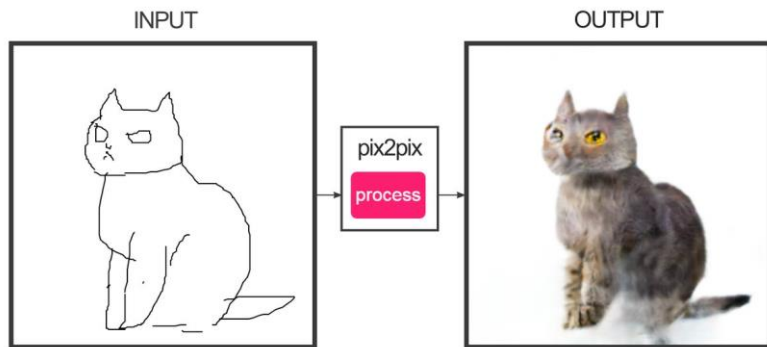


Output

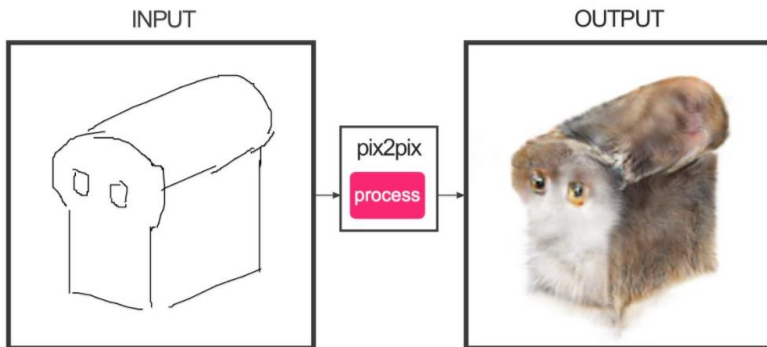


Trained on Edges  $\rightarrow$  Images

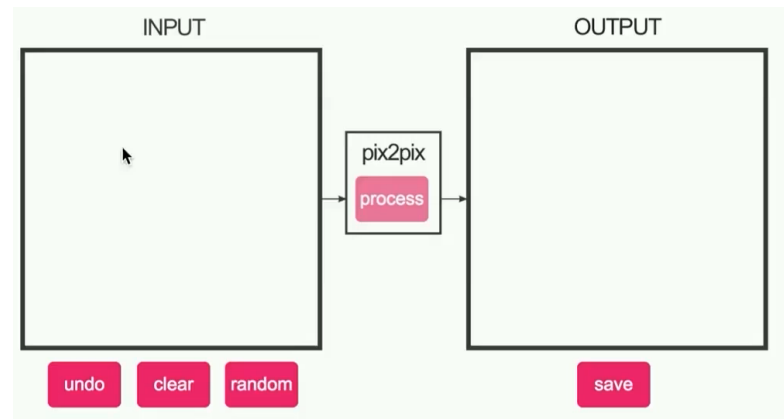
## #edges2cats [Christopher Hesse]



@gods\_tail



Ivy Tasi @ivymyt



@matthematician



Vitaly Vidmirov @vvid

<https://affinelayer.com/pixsrv/>



Input



Output



Groundtruth



Data from  
[[maps.google.com](https://maps.google.com)]



slides credit: Isola / Zhu

# BW $\rightarrow$ Color

Input

Output



Input

Output



Input

Output



# Ideas behind Pix2Pix

- $L = L_{GAN} + \lambda L_1$  (makes it more constraint)
- Unet / skip connections for preserving structure
- Noise only through dropout
  - cGANs tend to learn to ignore the random vector  $z$
  - Still want probabilistic model

# Ideas behind Pix2Pix

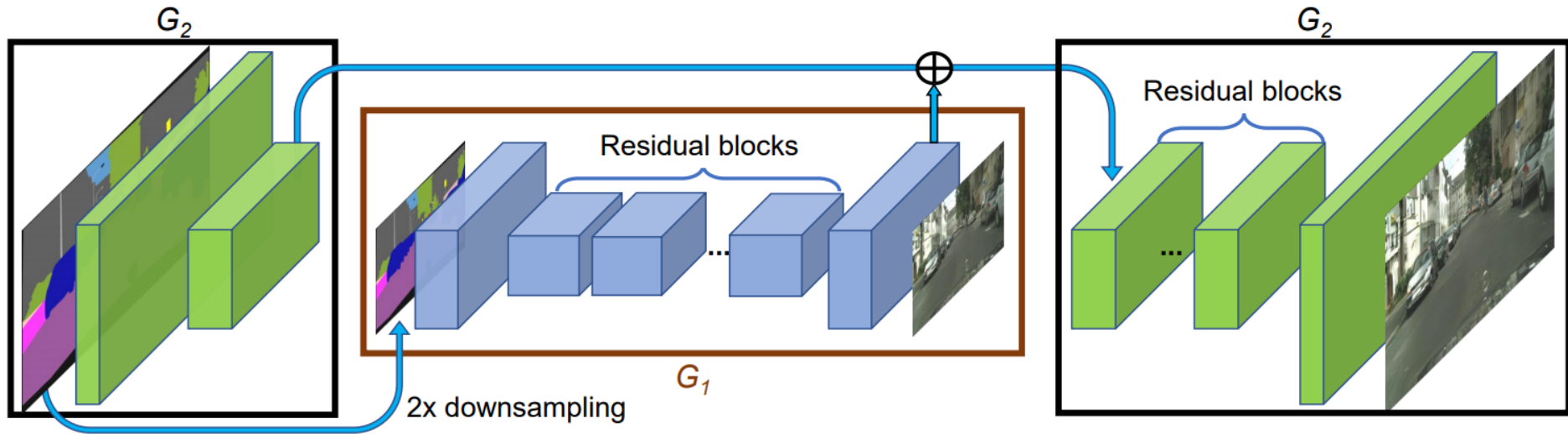
- L1 or L2 loss for low frequency details
  - GAN discriminator for high frequency details
- > PatchGAN
- GAN discriminator applied only to local patches
  - It's fully-convolutional; i.e., can run on arbitrary image sizes

# Pix2PixHD

- Expand the pix2pix idea to multi-scale
- Coarse-to-fine generator + discriminator
- G's and D's are the same but since they operate on different resolutions, they have effectively a larger receptive field



# Pix2PixHD

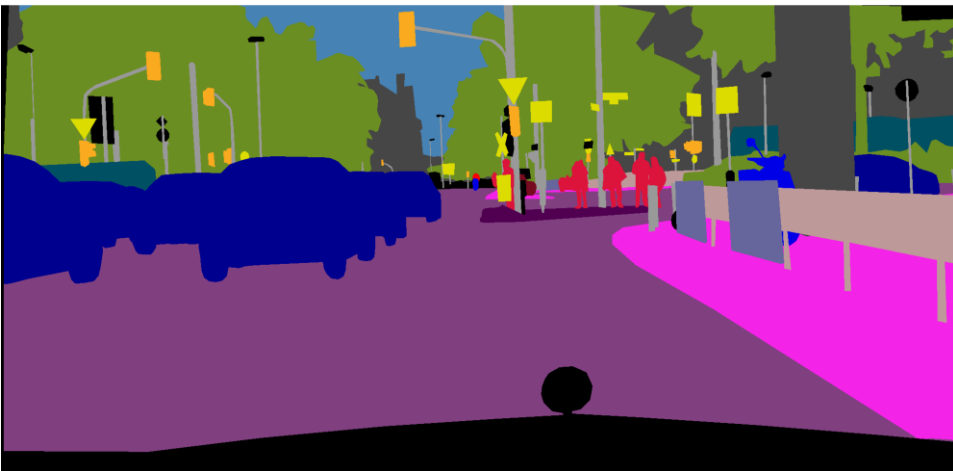


# Pix2PixHD

- Use of multi-scale discriminators
- $\min_G \max_{D_1, D_2, D_3} \sum_{k=1,2,3} L_{GAN}(G, D_k)$
- Can make various combinations of stacking discriminator and generator
  - E.g., have a single G and downsample generated and real images – or have intermediate real images (cf. ProGAN)

# Pix2PixHD

Input labels

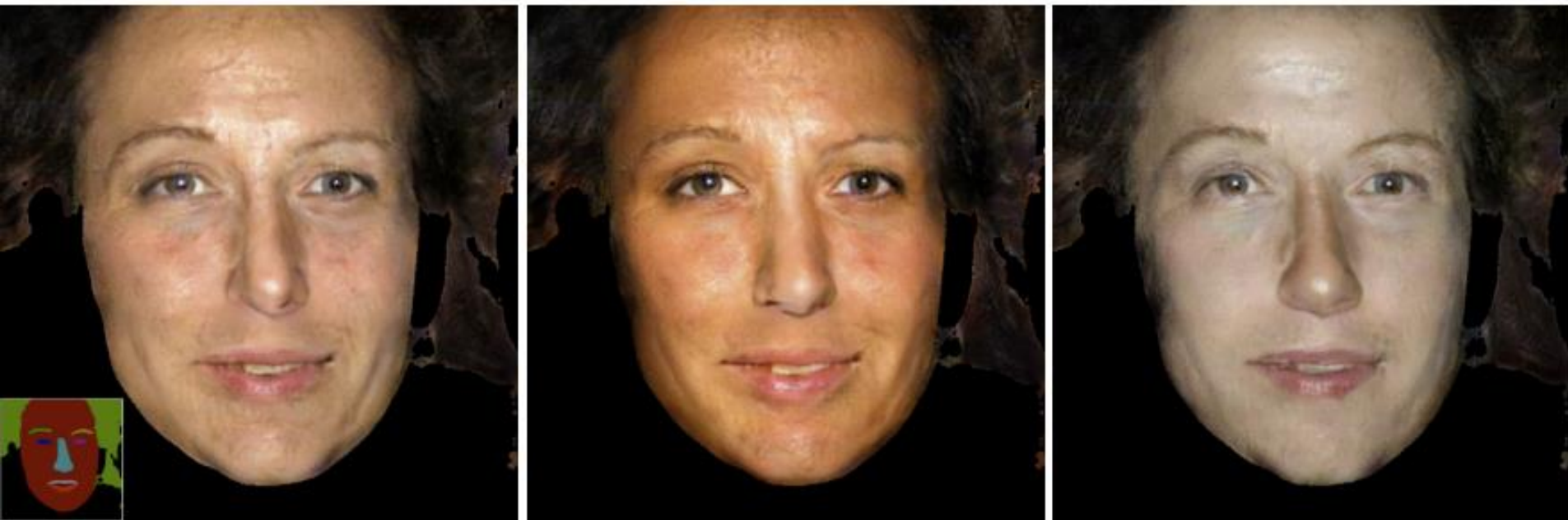


Synthesized image





# Pix2PixHD



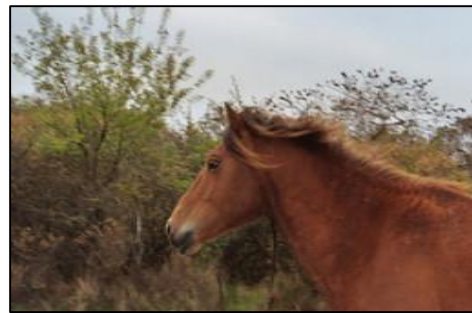
# Pix2PixHD (Interactive Results)



# Paired



Label  $\leftrightarrow$  photo: per-pixel labeling



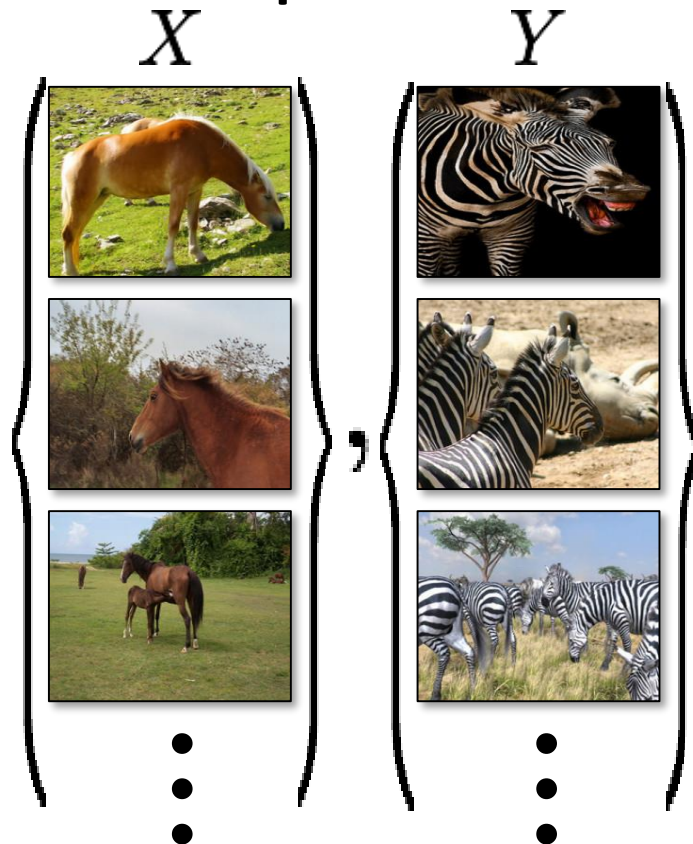
Horse  $\leftrightarrow$  zebra: how to get zebras?

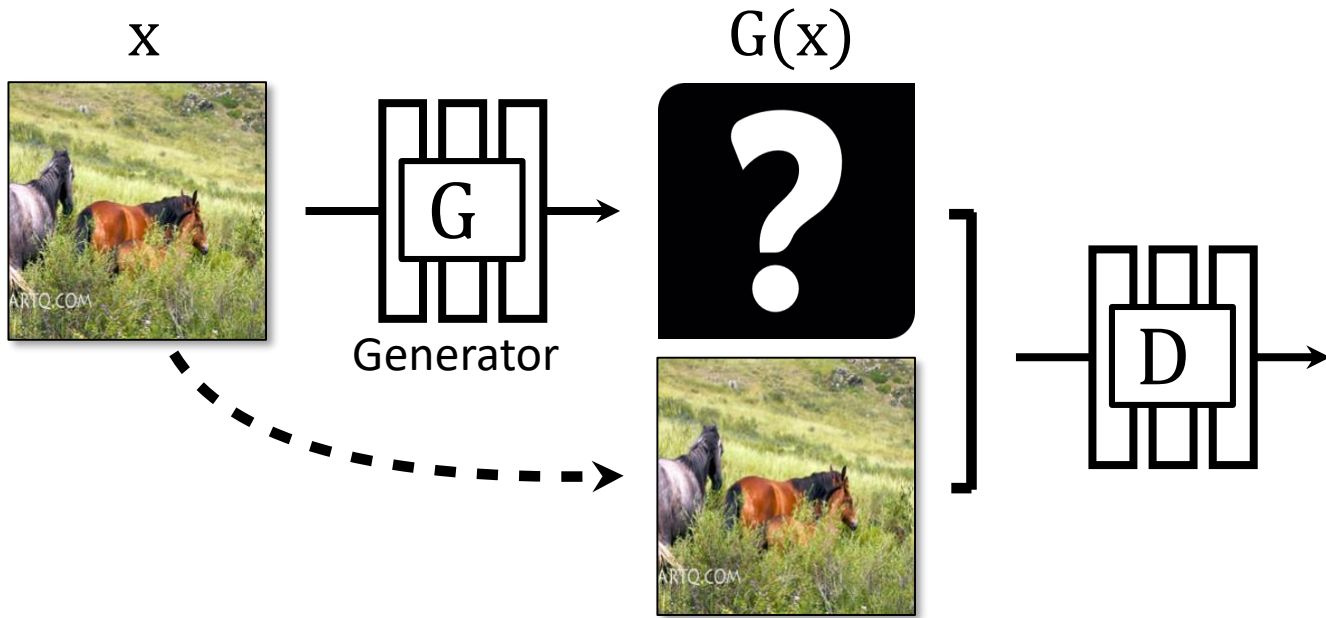
- Expensive to collect pairs.
- Impossible in many scenarios.

# Paired



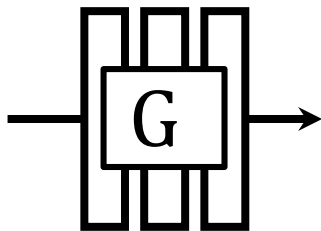
# Unpaired





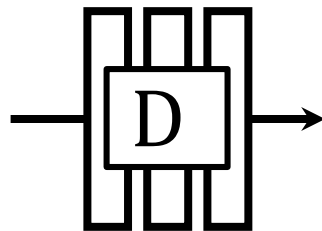
No input-output pairs!

X



Generator

$G(X)$

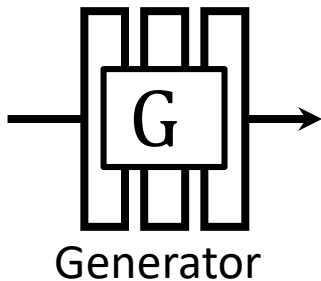


Discriminator

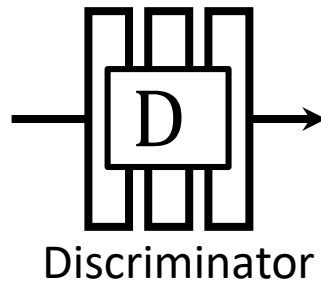
Real!



X



$G(x)$



Real too!

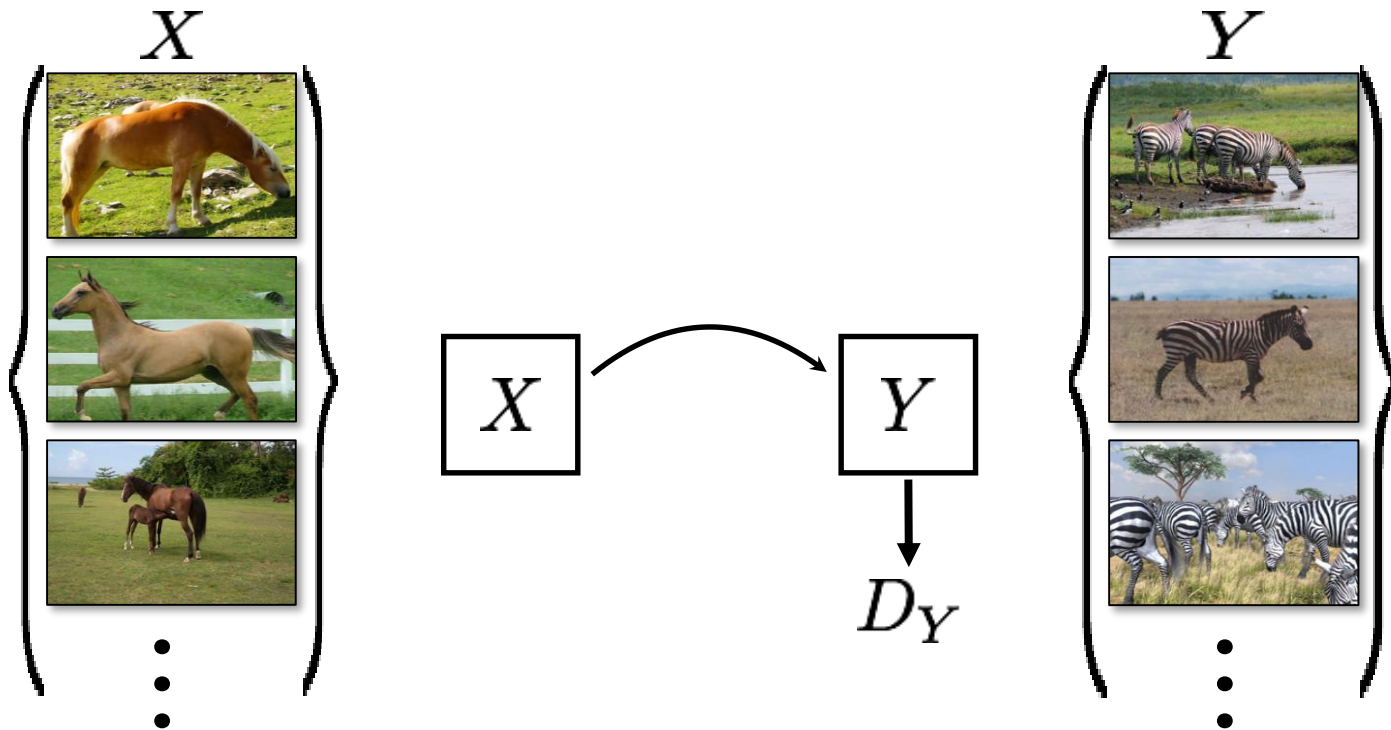
GANs doesn't force output to correspond to input



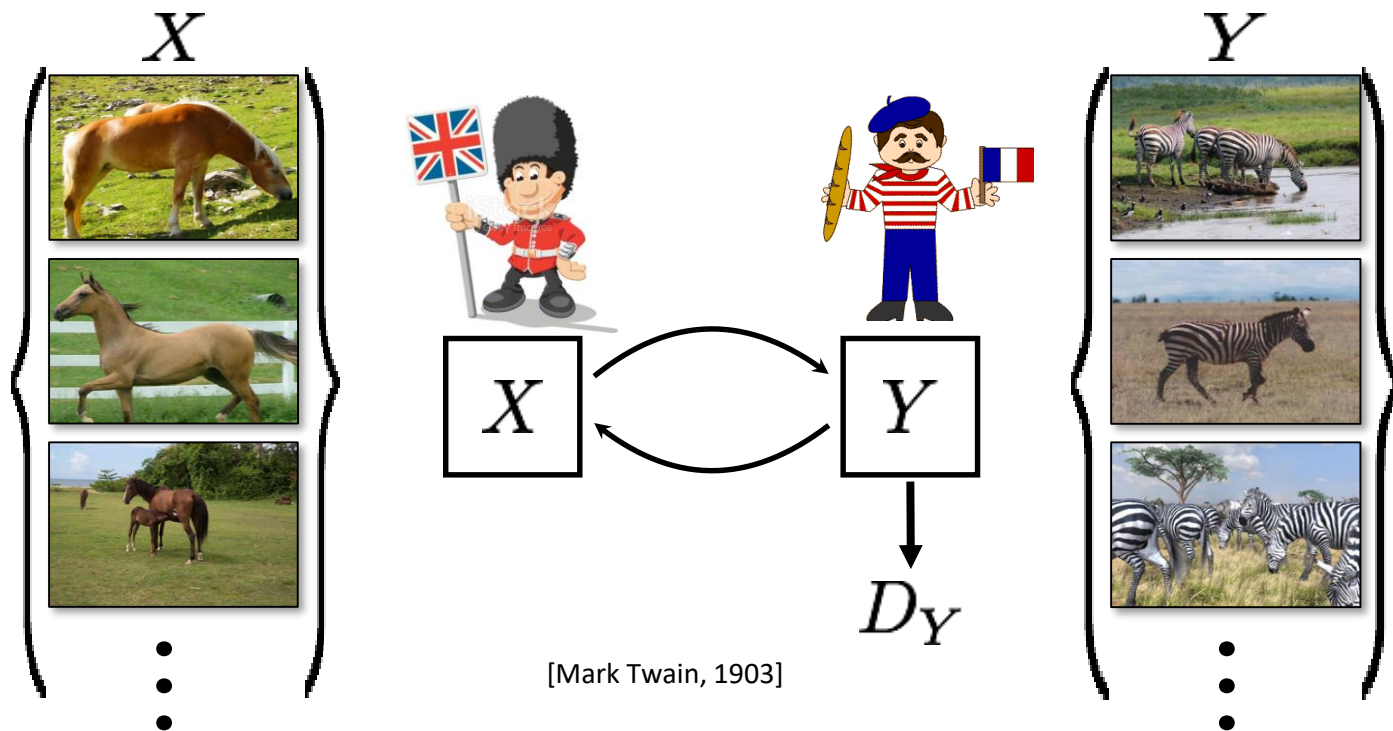
mode collapse!



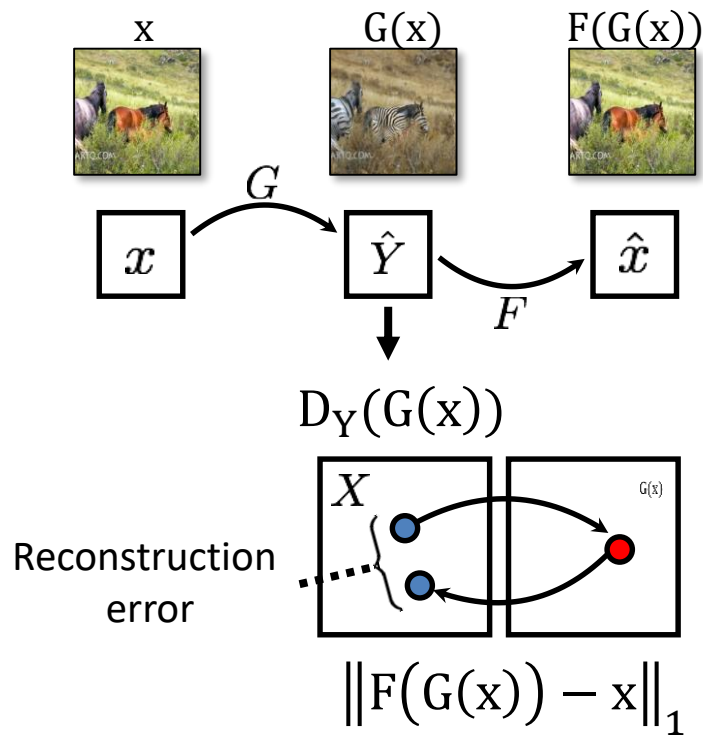
# Cycle-Consistent Adversarial Networks



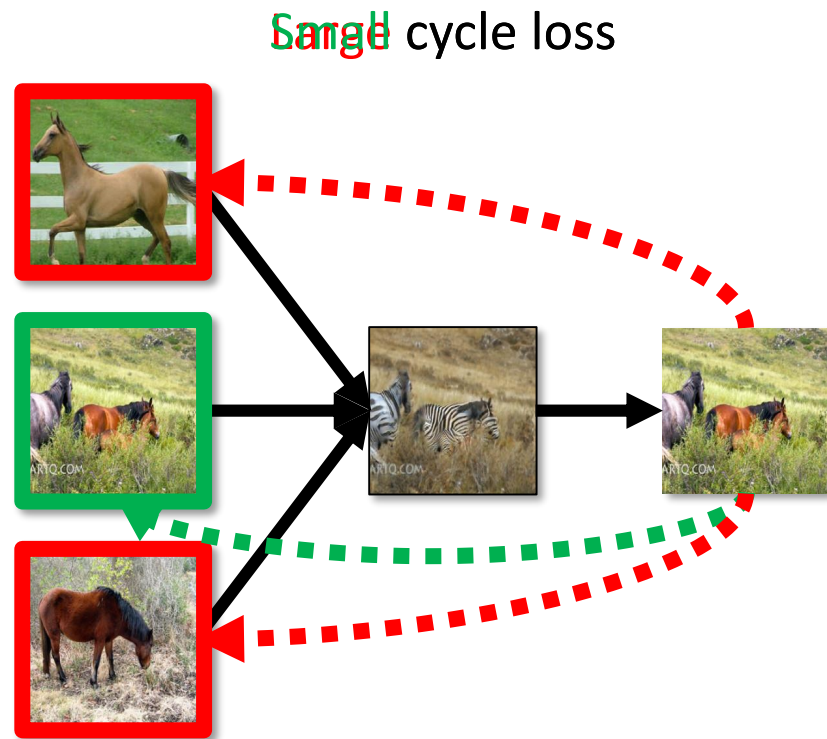
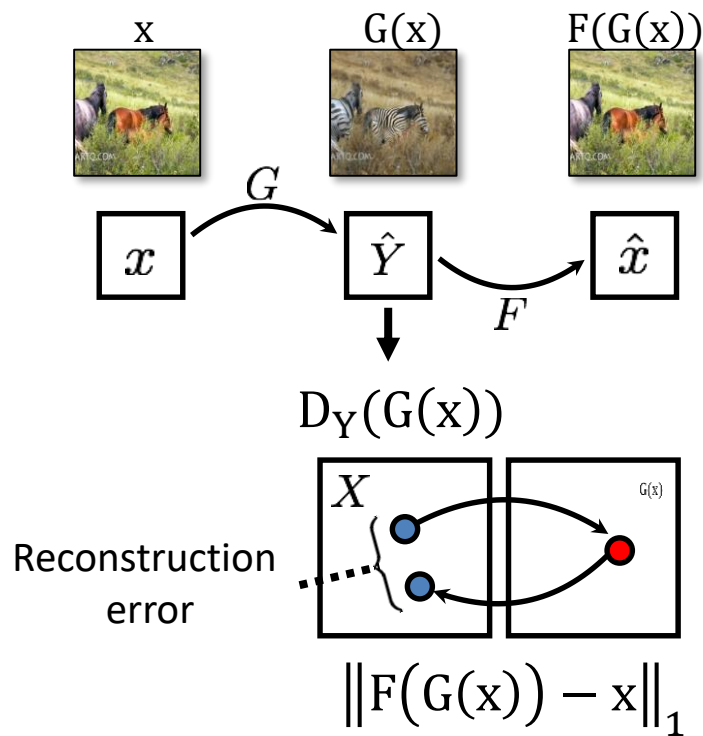
# Cycle-Consistent Adversarial Networks



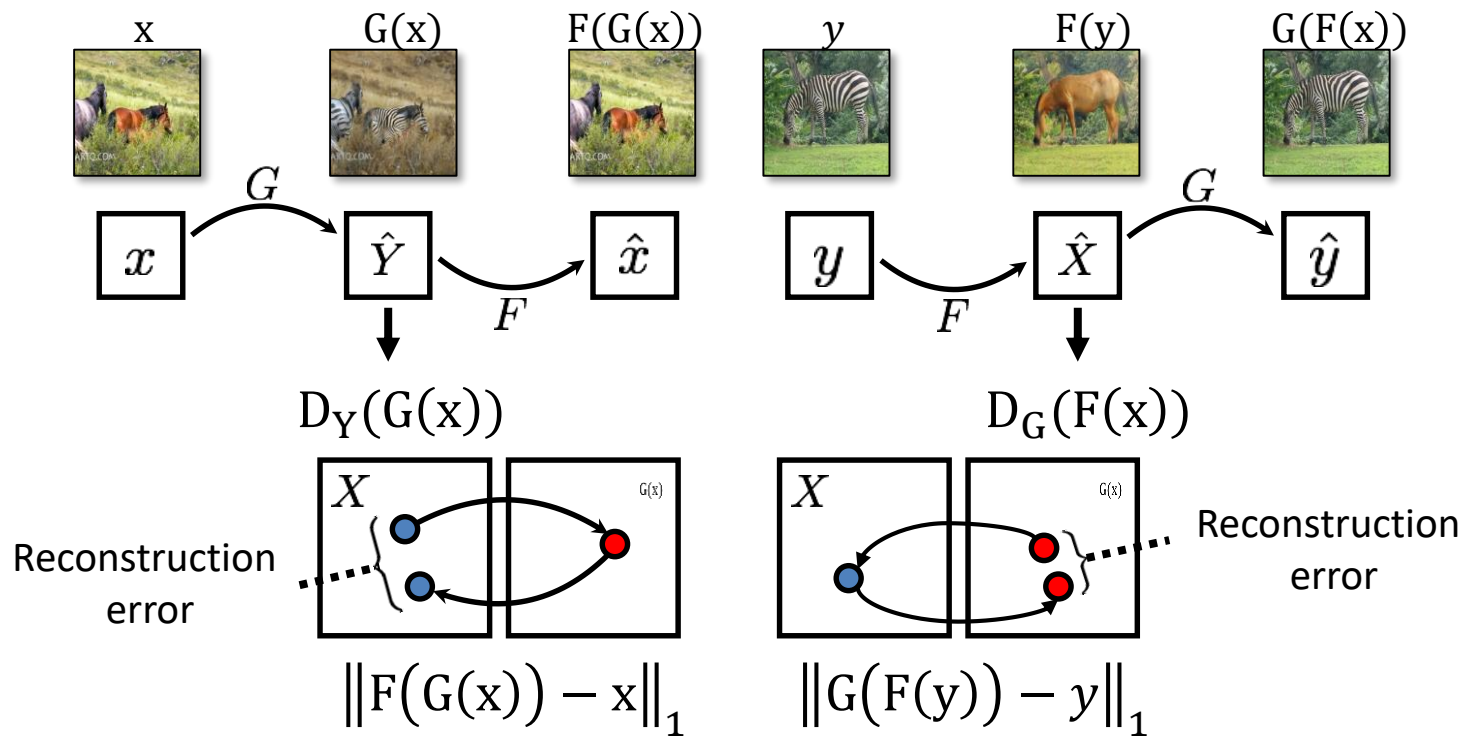
# Cycle Consistency Loss



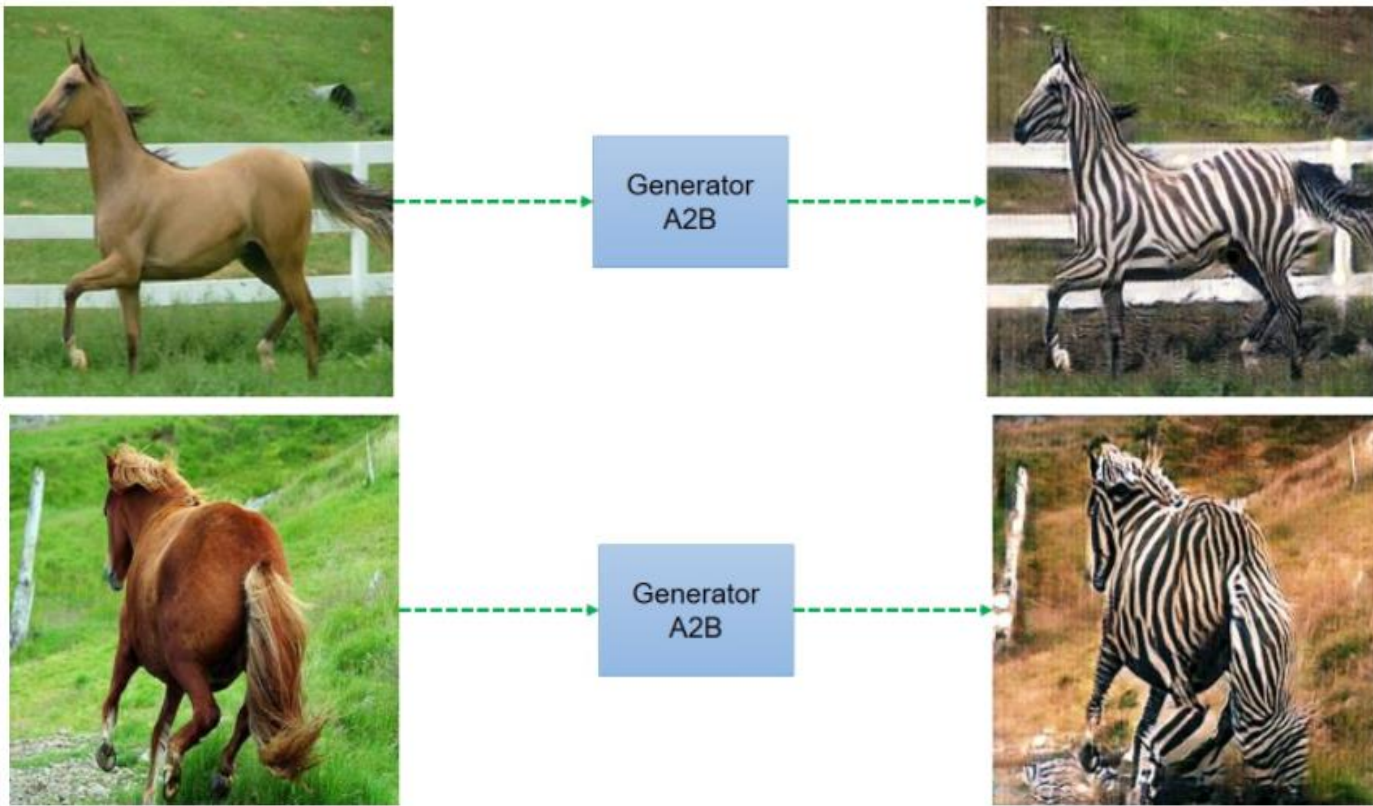
# Cycle Consistency Loss



# Cycle Consistency Loss



# Cycle GAN - Overview





# Monet's paintings → photos









# Next Lectures

- Next Lectures:
  - Videos
  - Neural Rendering
  - 3D Deep Learning
- Keep working on the projects!

See you next week 😊