

# Lecture 10: Image Synthesis

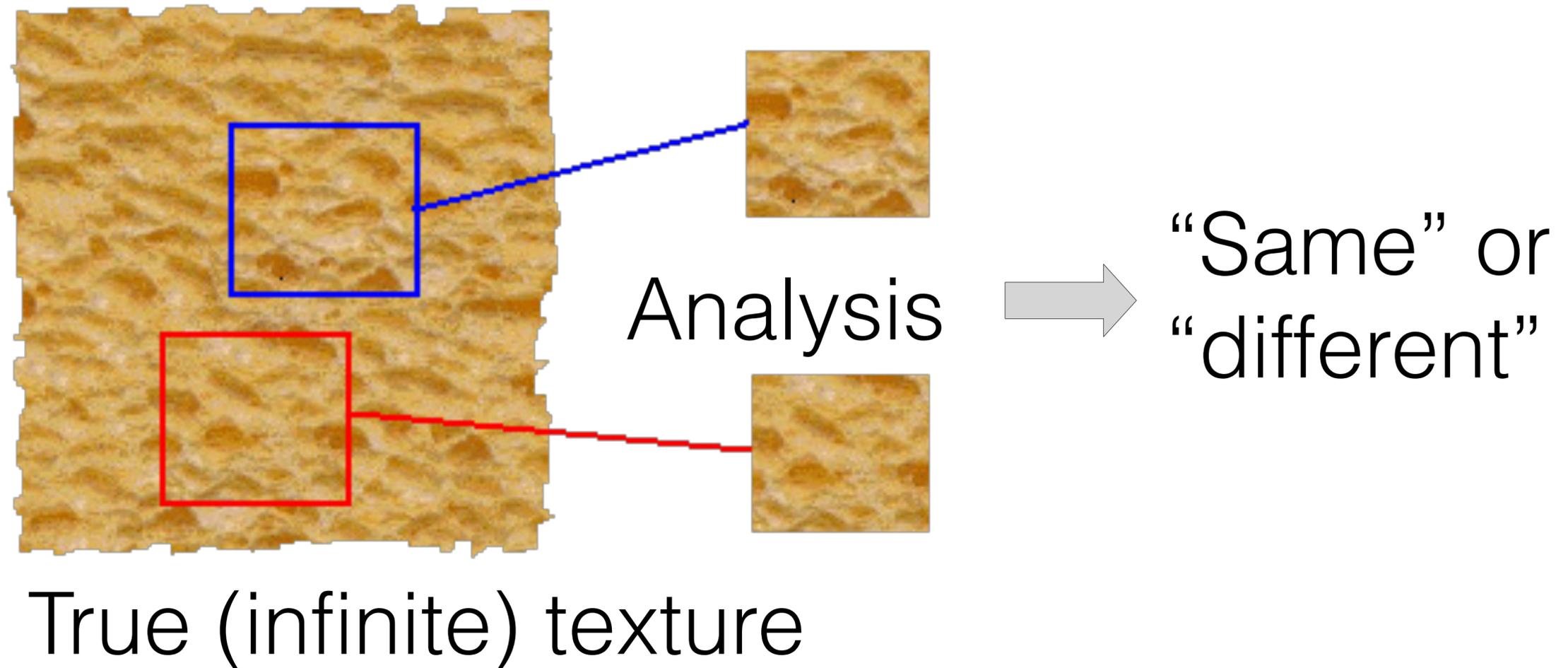
# Announcements

- Problem Set 4 is due on Wednesday
- Need to submit to both Canvas and Gradescope

# Today

- Texture synthesis
- Image stylization
- Generative adversarial nets (GANs)
- Conditional GANs

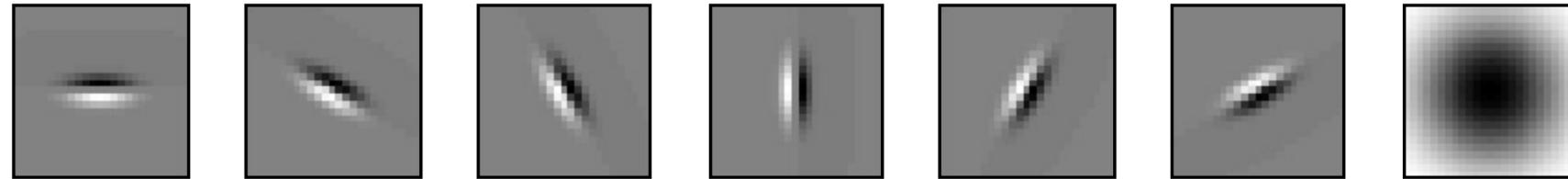
# Texture analysis



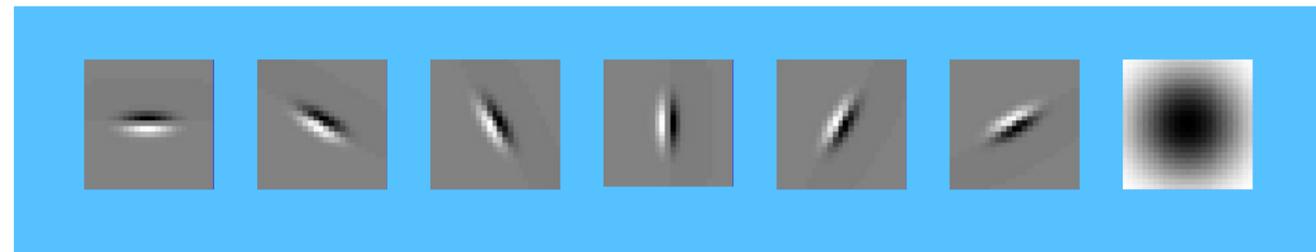
What we'd like: are they made of the same “stuff”. Are these textures similar?

# How can we represent texture in natural images?

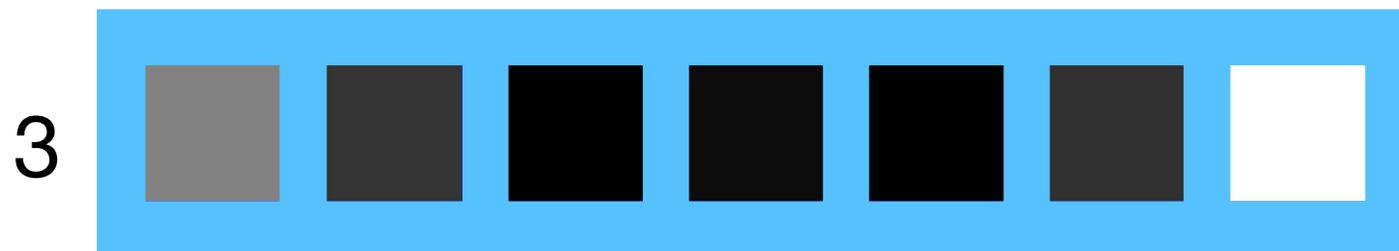
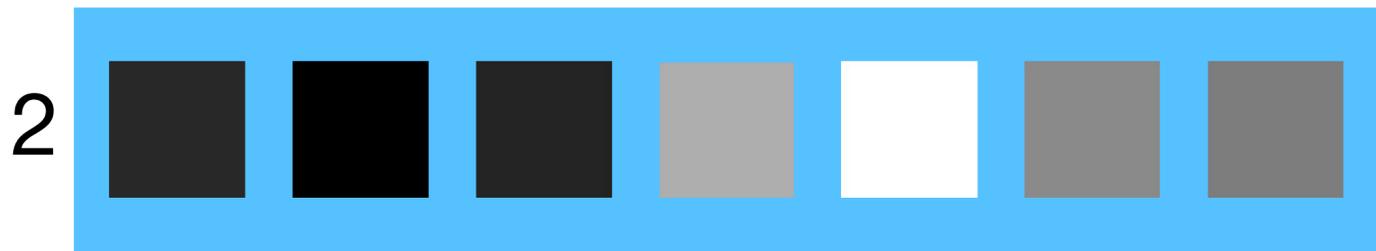
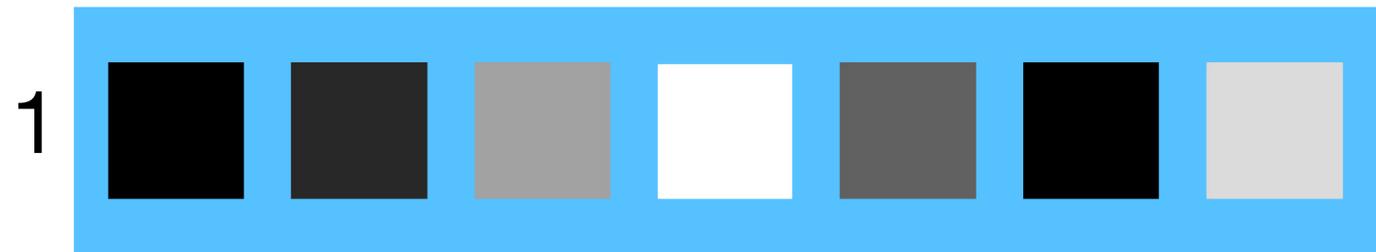
Idea #1: Record simple statistics (e.g., mean, std.) of absolute filter responses



# Can you match the texture to the response?

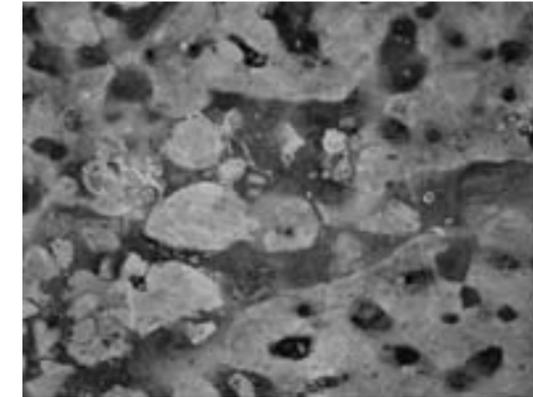


Filters



Mean abs. responses

A



B



C

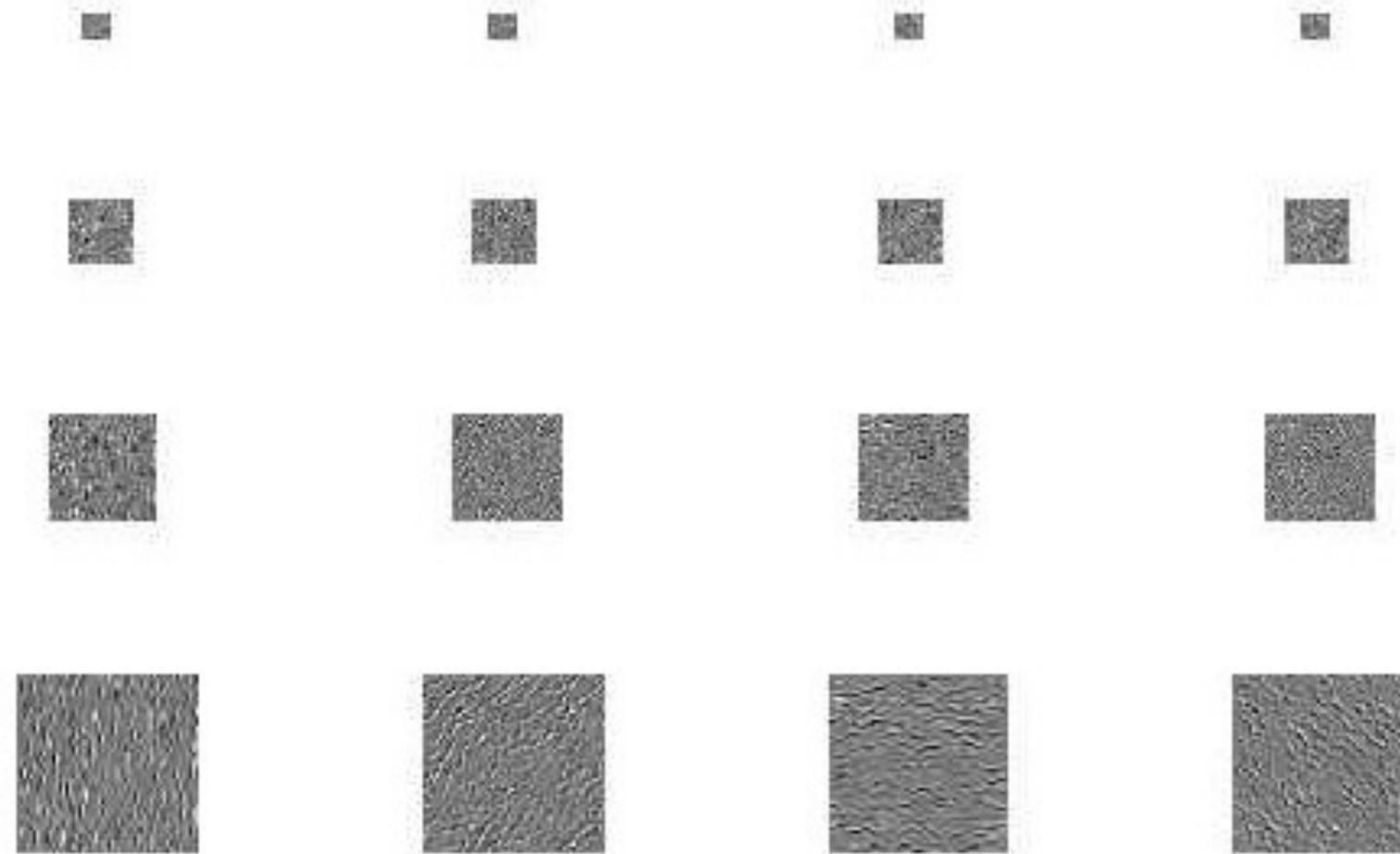
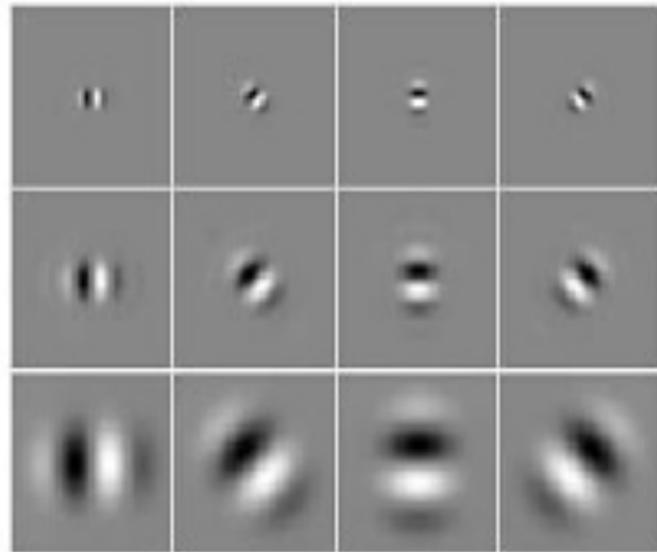


# How can we represent texture?

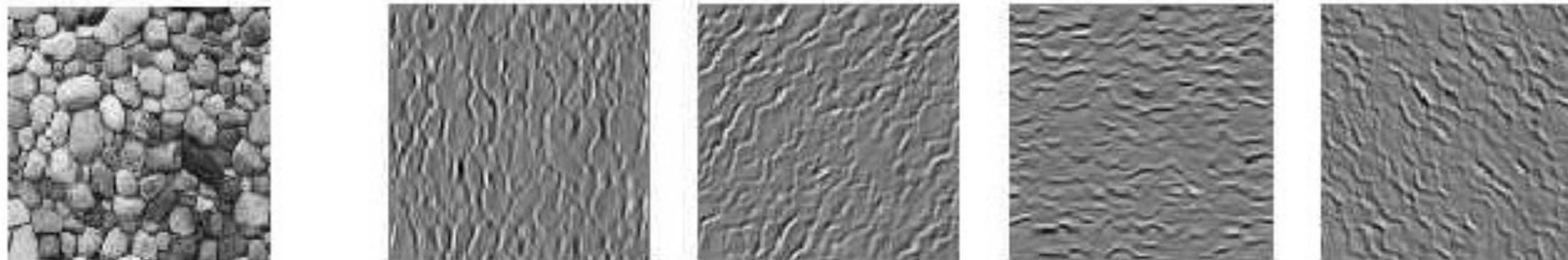
- Generalize this to “orientation histogram”
- Idea #2: Histograms of filter responses
  - One histogram per filter

# Pyramid of filter responses

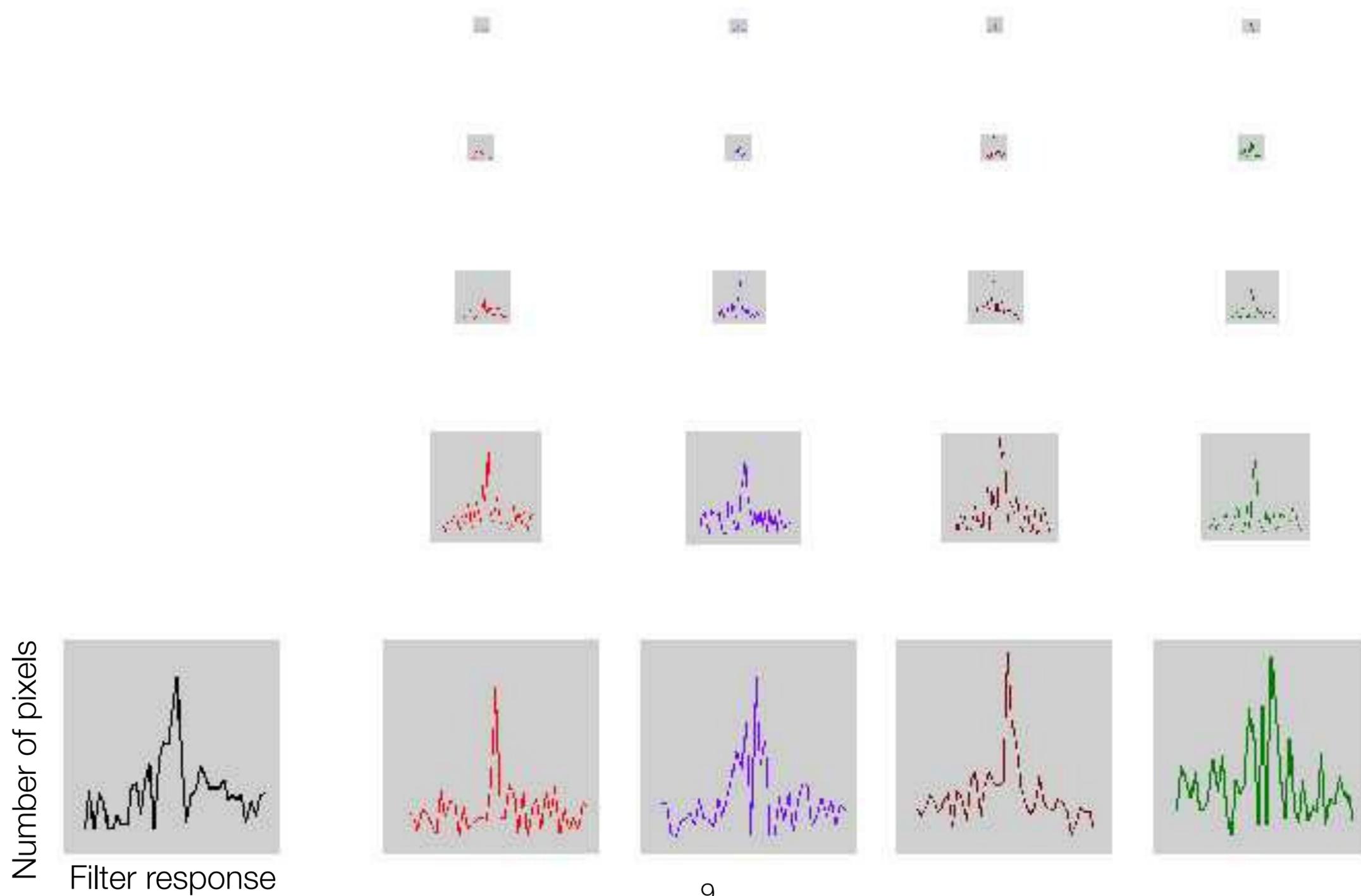
Filter bank



Input image



# Filter response histograms

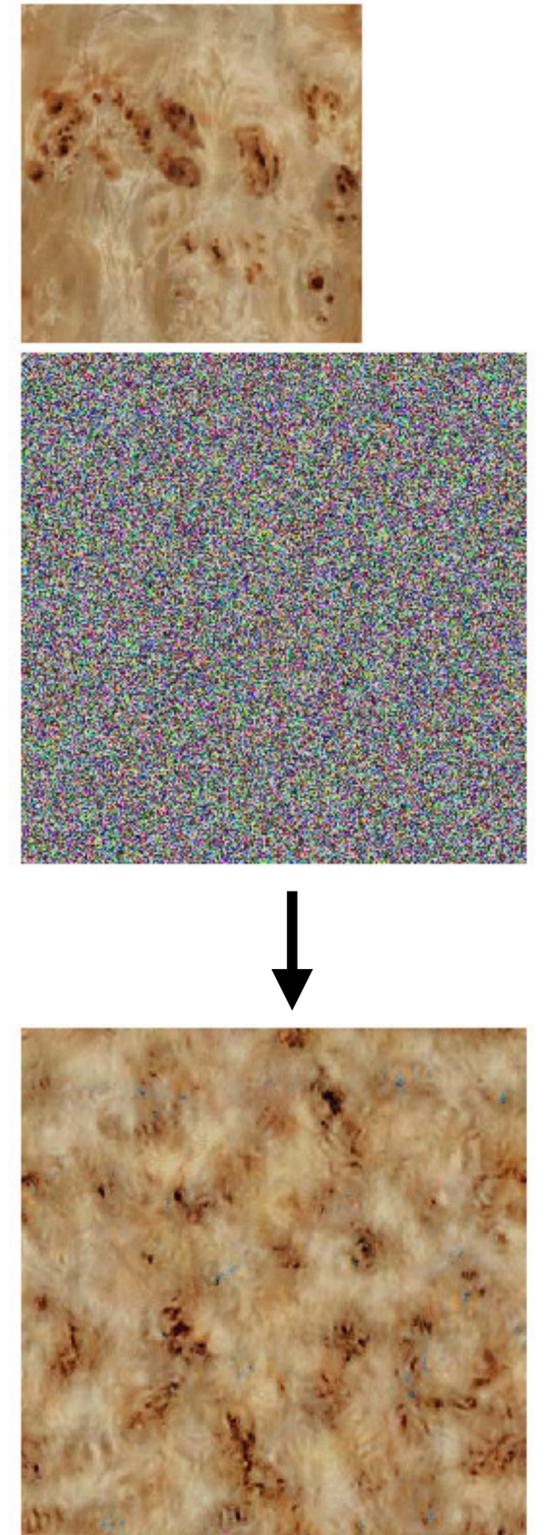


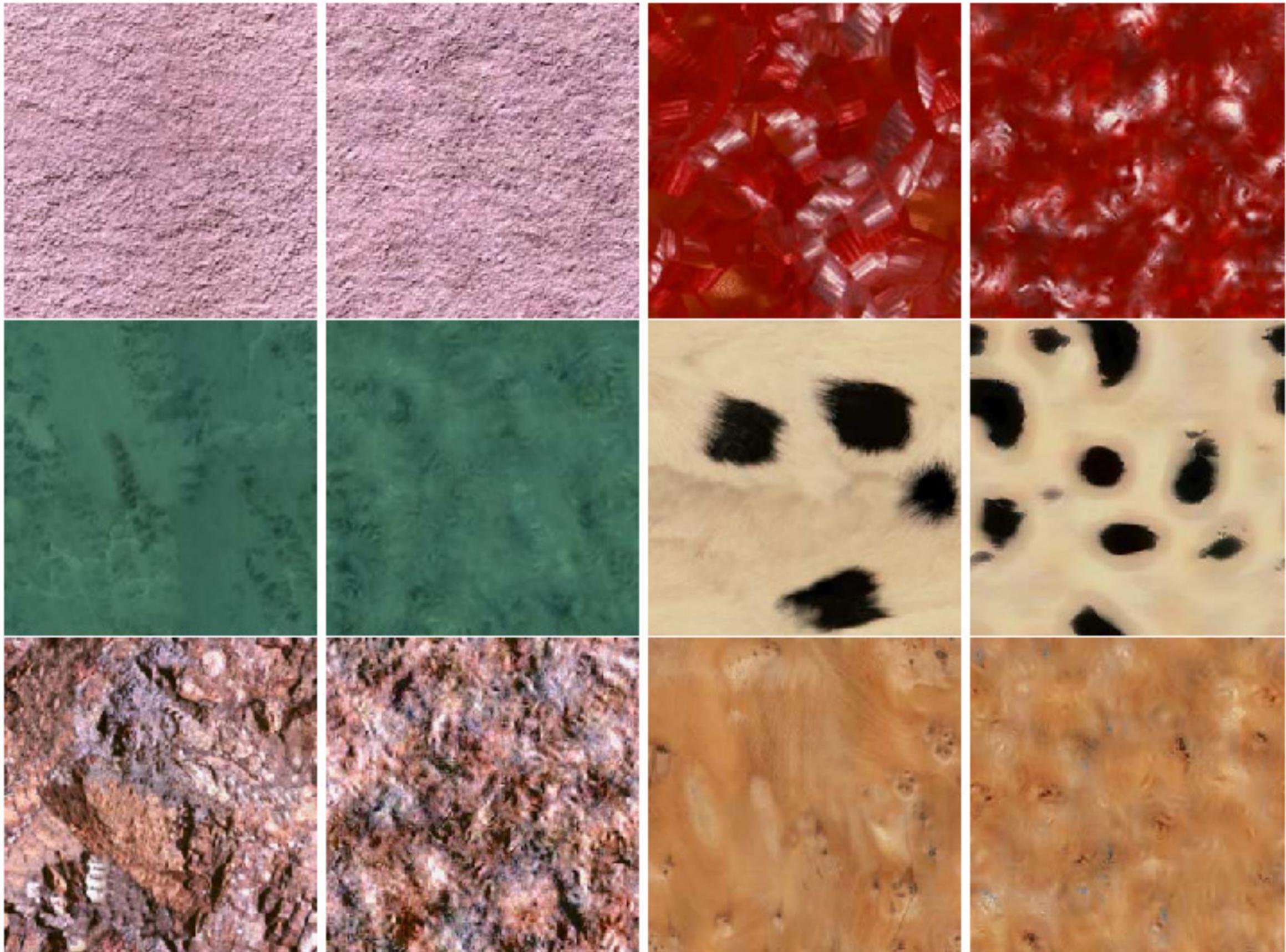
# Texture synthesis

Start with a noise image as output.

**Iterative algorithm** [Heeger & Bergen, 95]:

- Match pixel histogram of output image to input
- Decompose input/output images using a Steerable Pyramid
- Match histograms of input and output pyramids
- Reconstruct image and repeat







# Failure cases

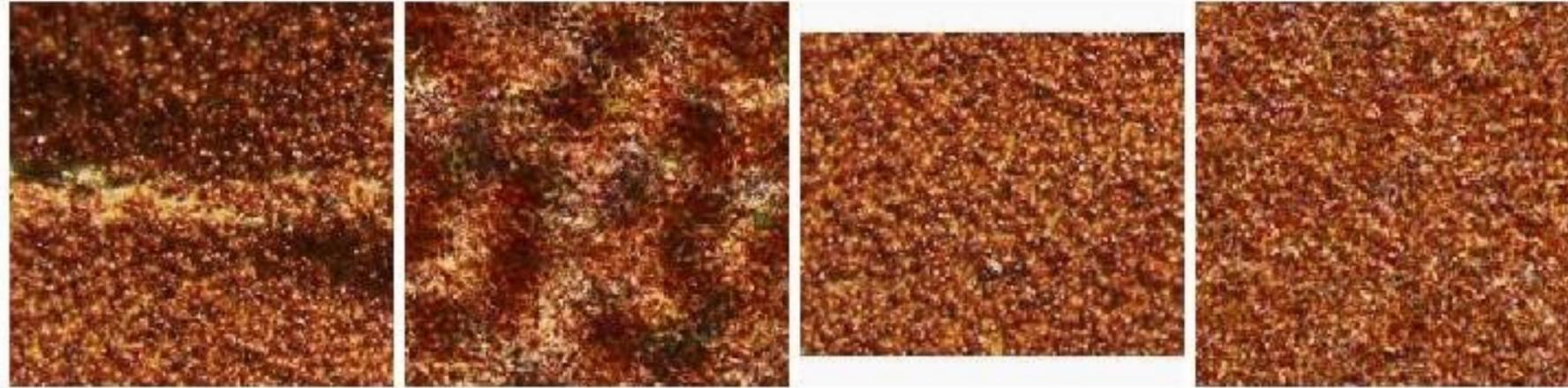


Figure 7: (Left pair) Inhomogeneous input texture produces blotchy synthetic texture. (Right pair) Homogeneous input.

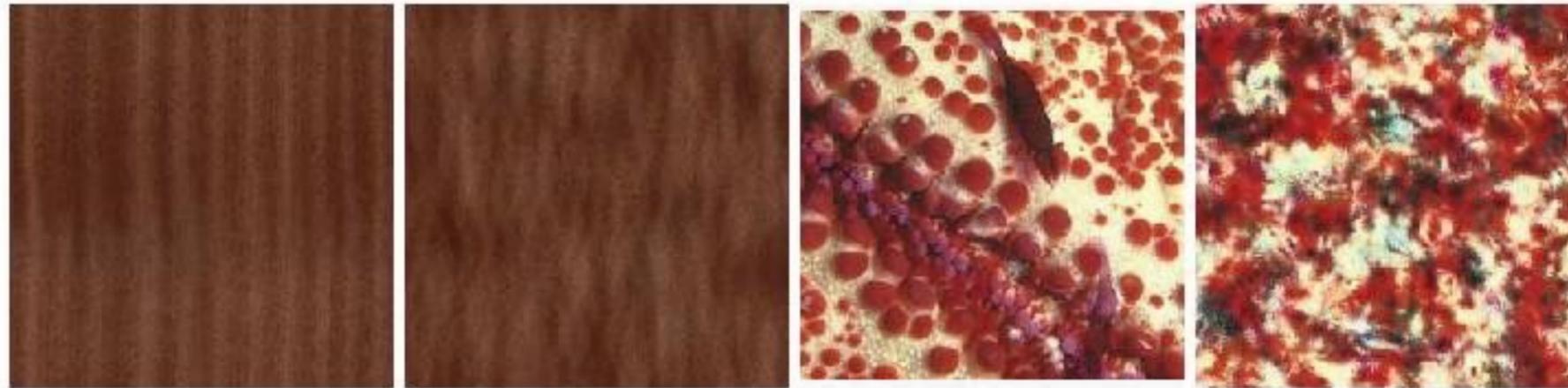


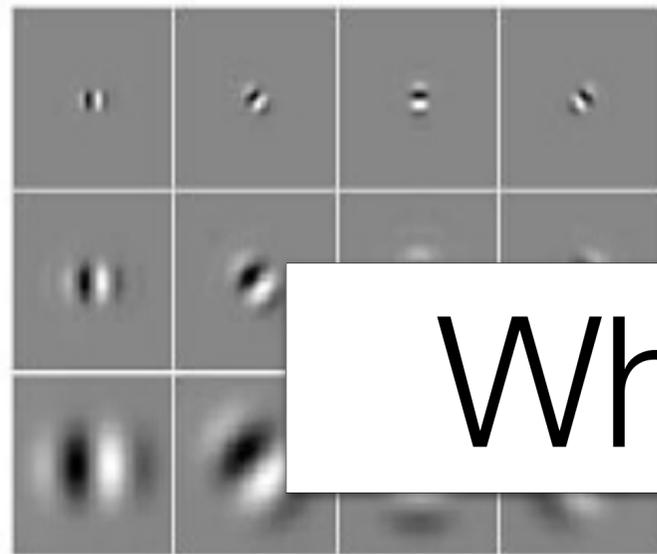
Figure 8: Examples of failures: wood grain and red coral.



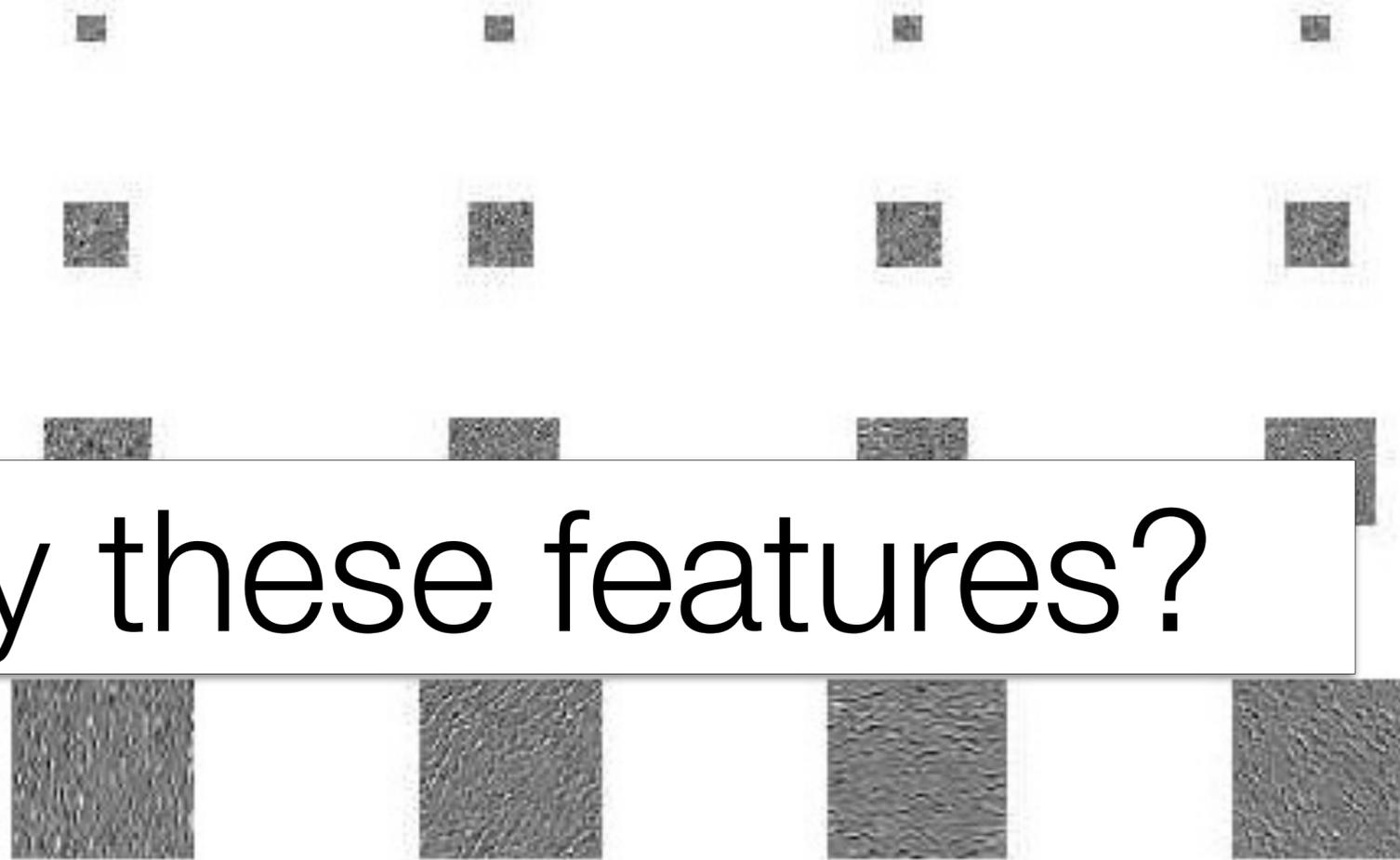
Figure 9: More failures: hay and marble.

# Pyramid of filter responses

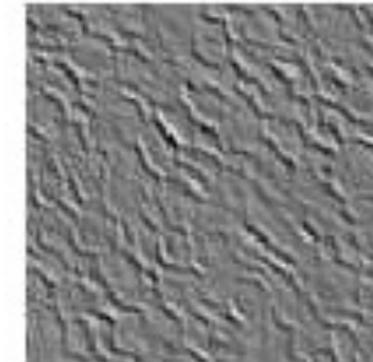
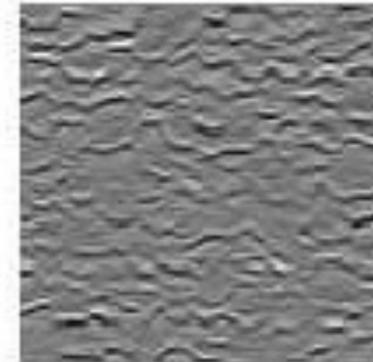
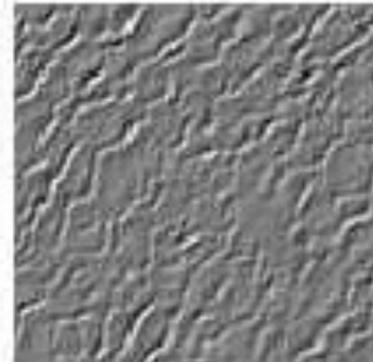
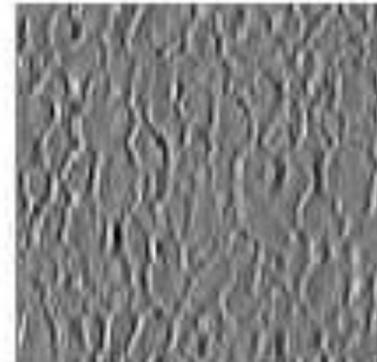
Filter bank



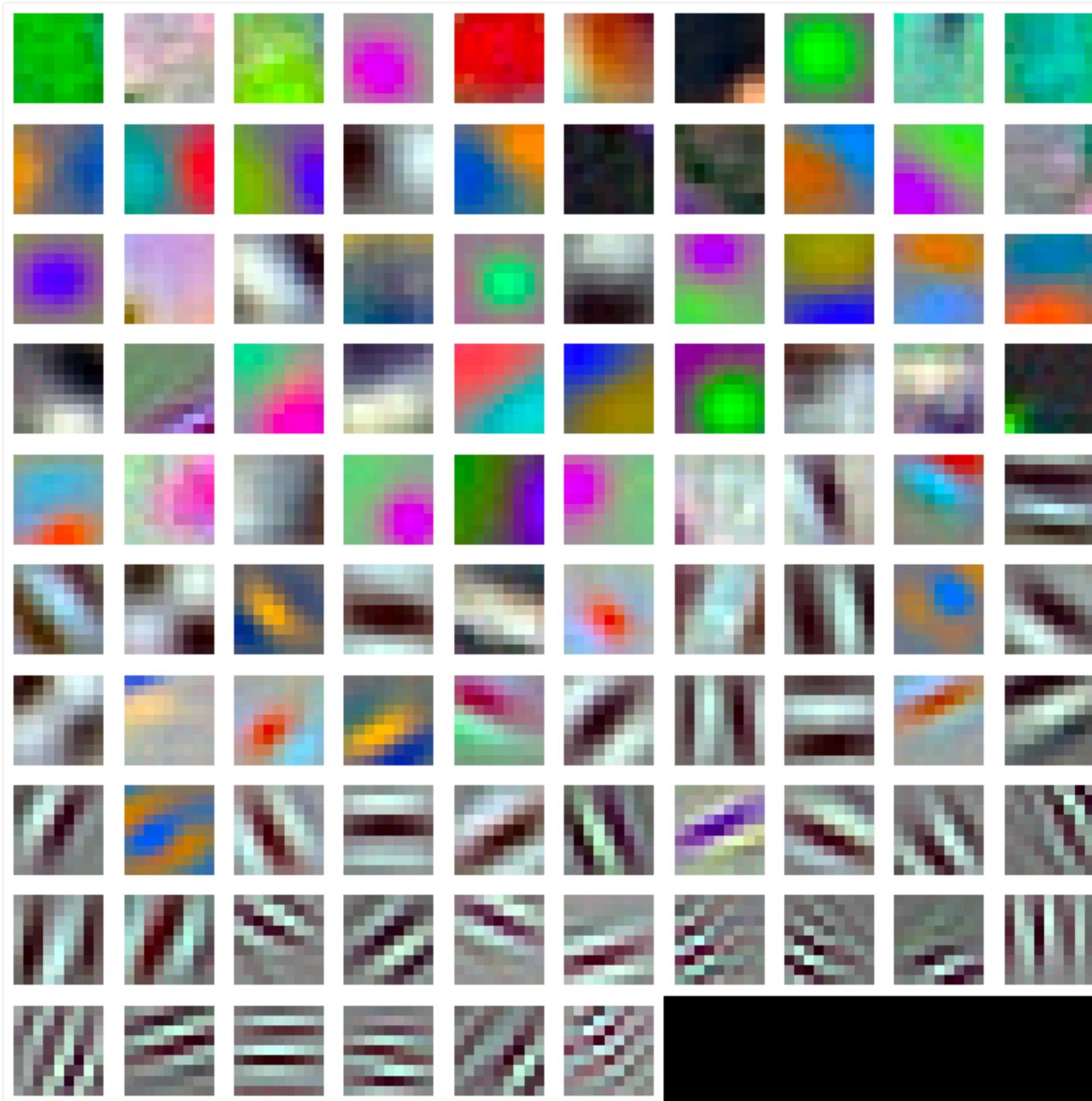
Why these features?



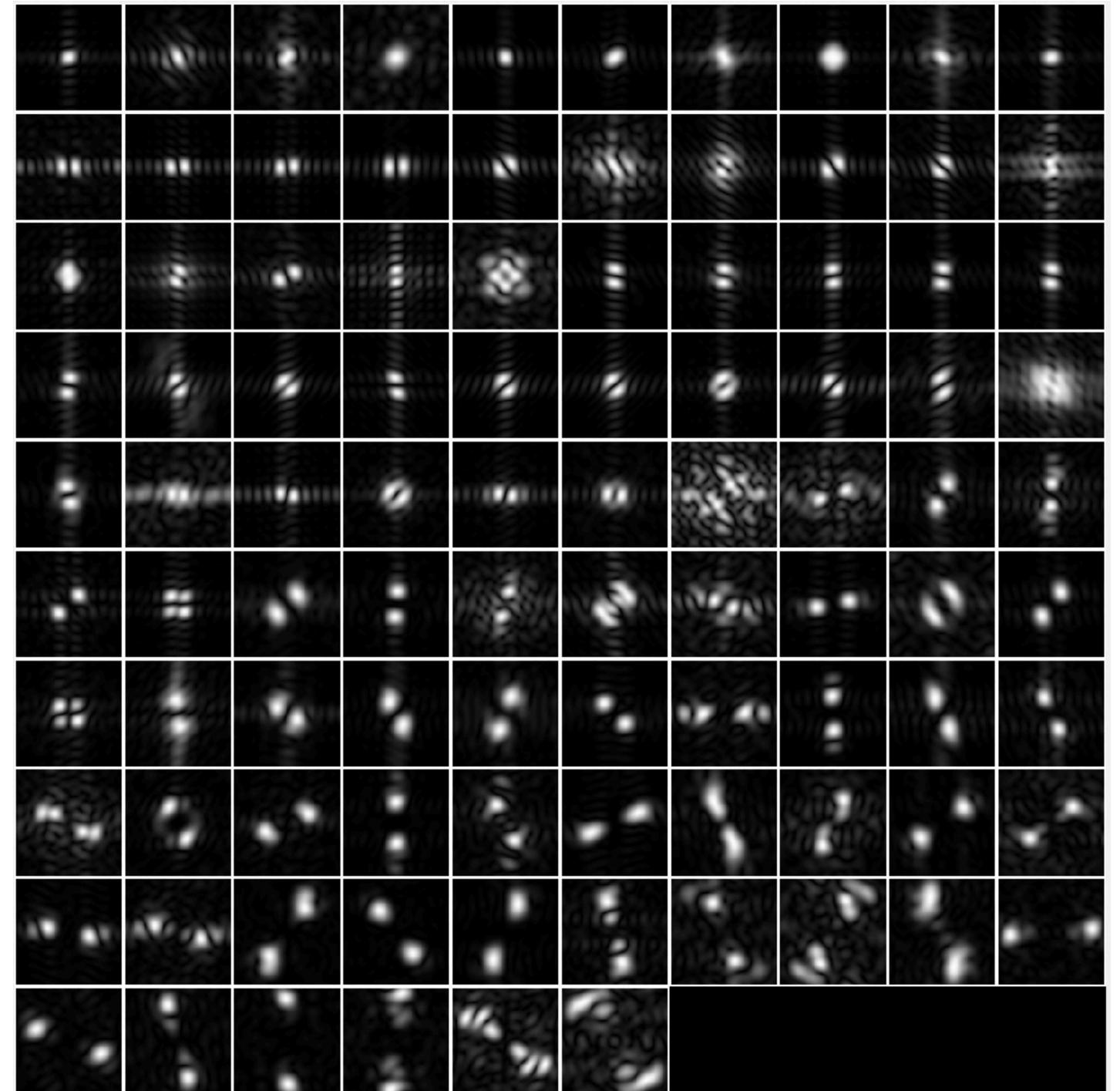
Input image



# Recall: neural net feature visualization

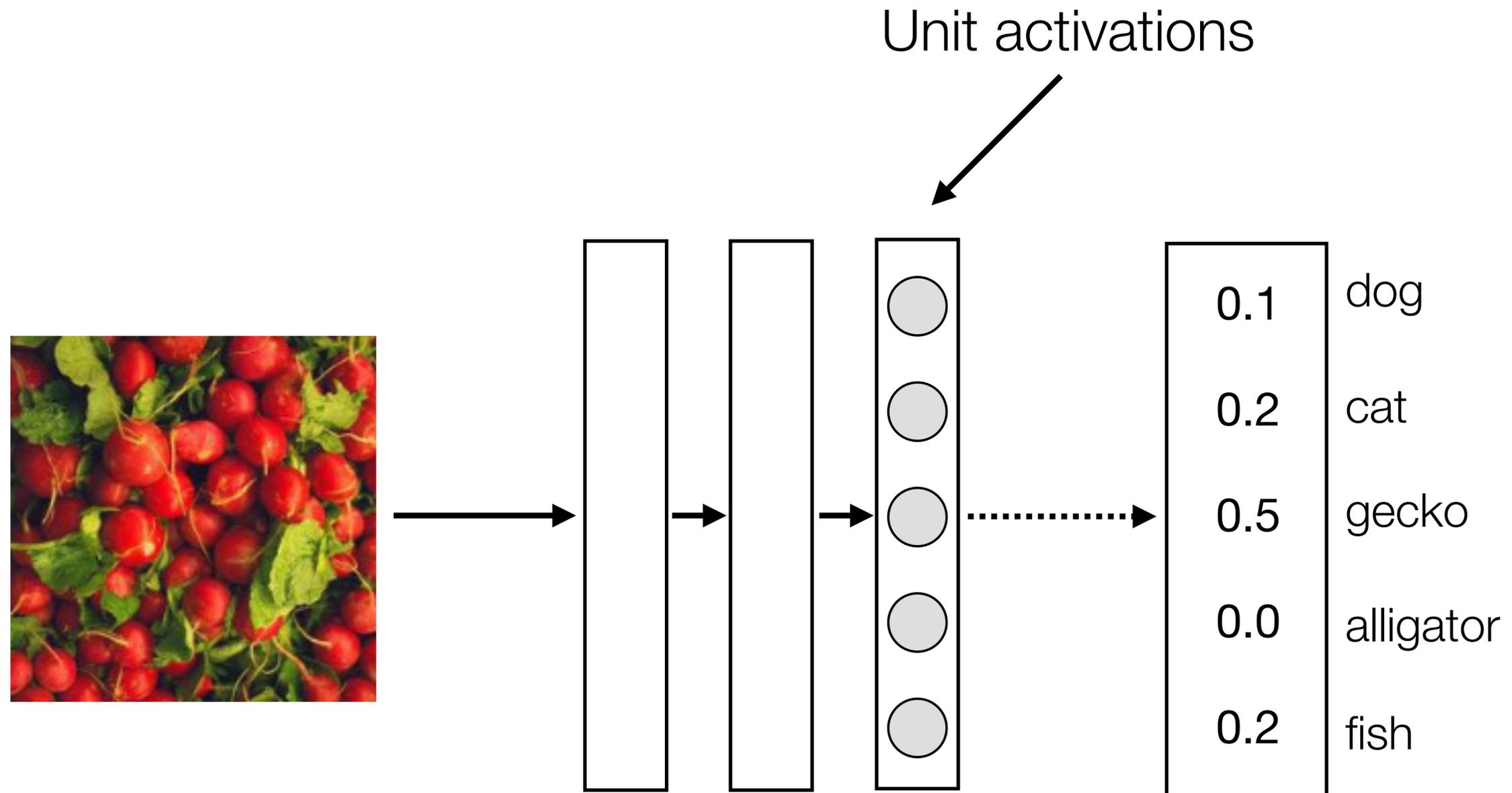


96 Units in conv1

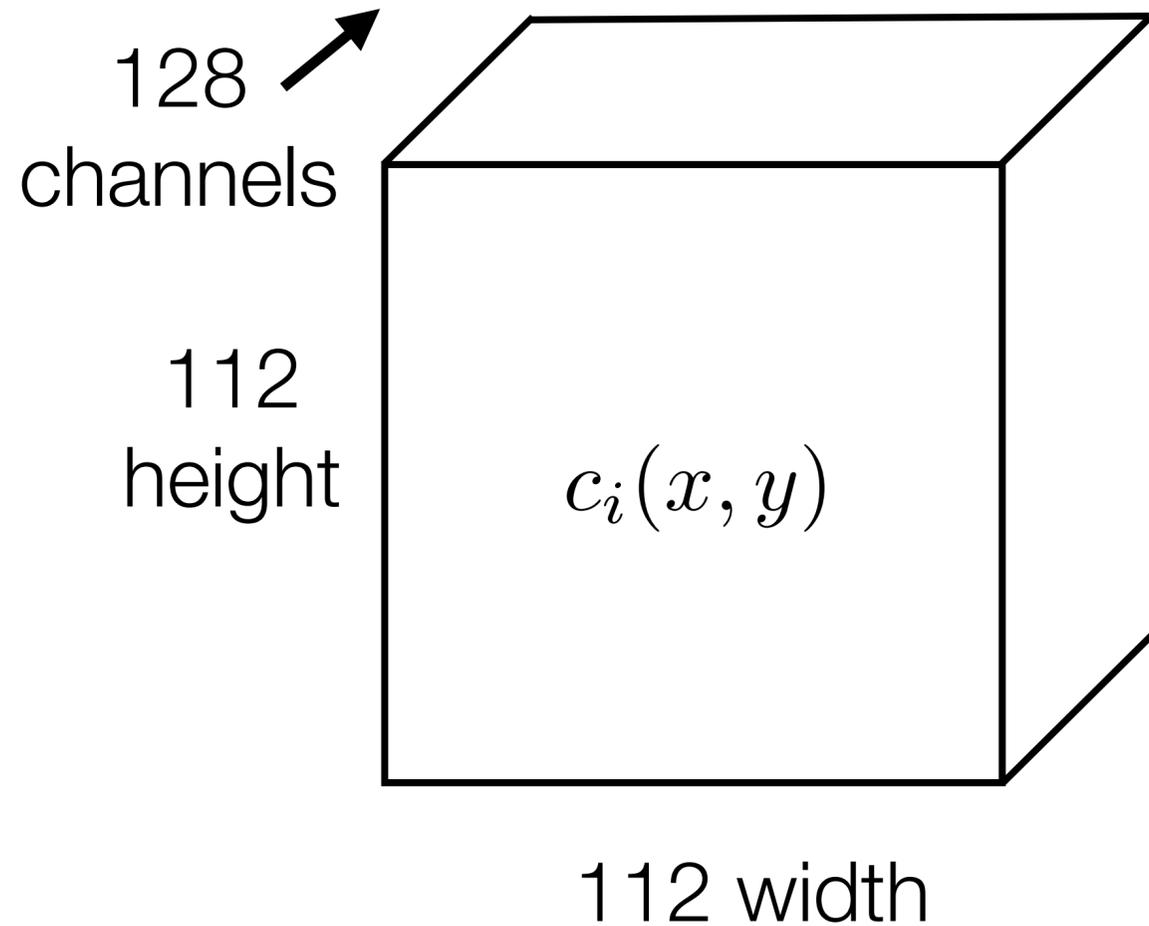


15

# Extracting features from a trained network

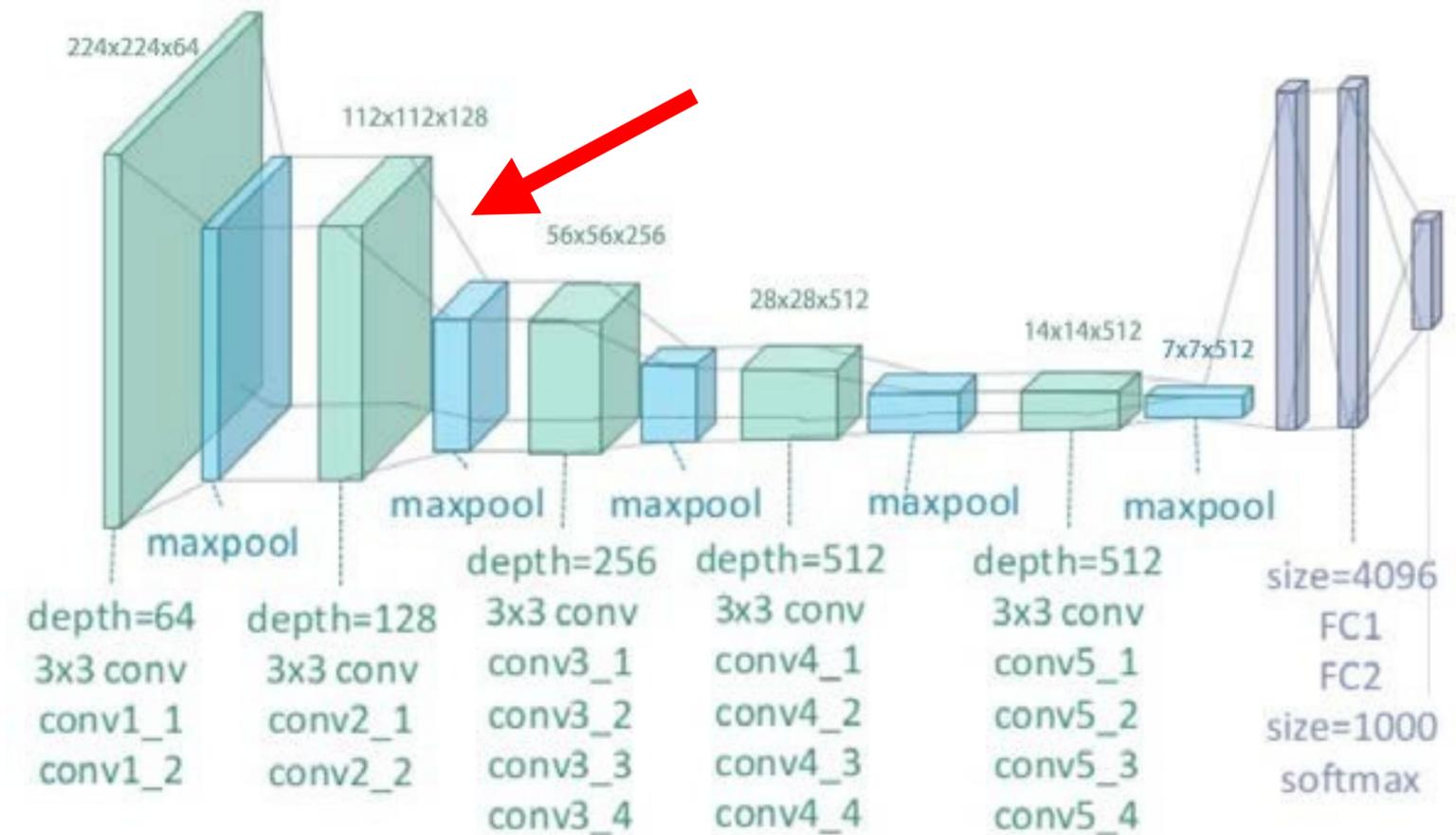


# Extracting neural net features

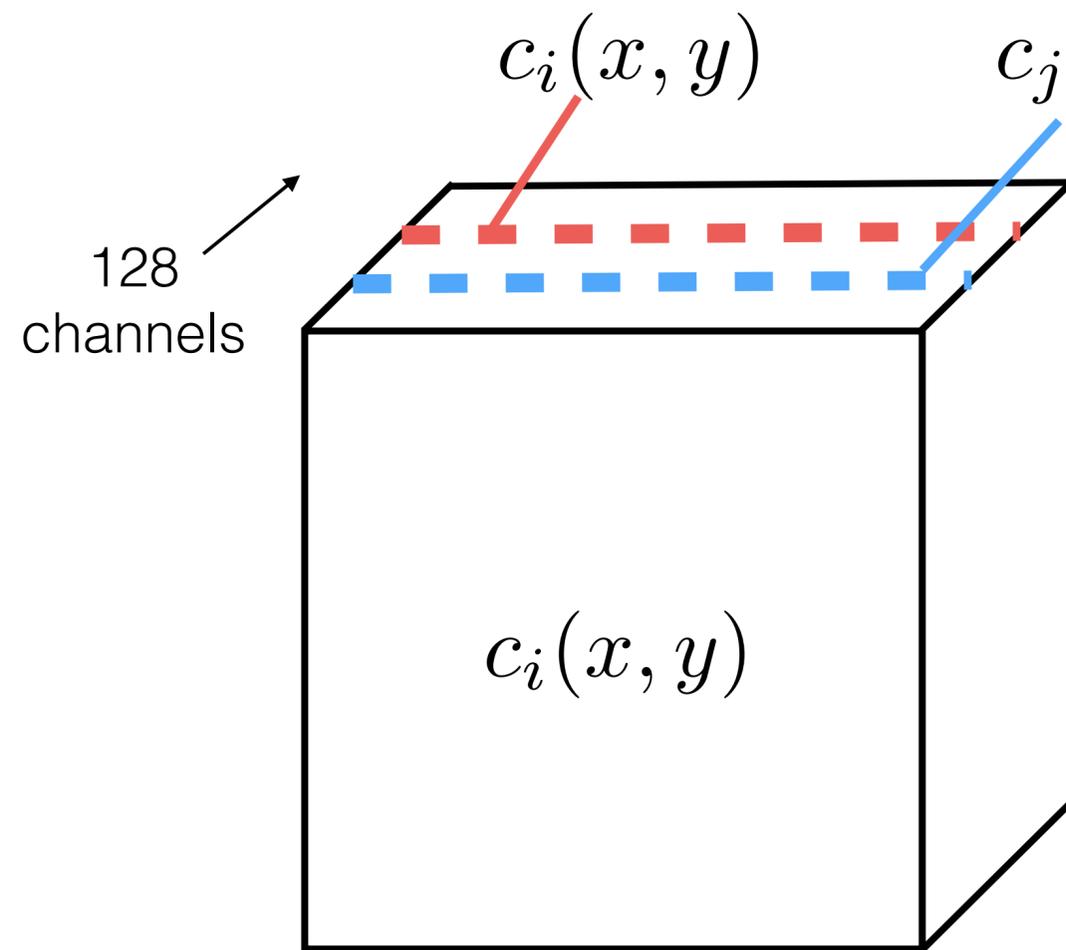


conv2 feature activations

(112 x 112) x 128 for conv2 of VGG19



# Capturing feature correlations



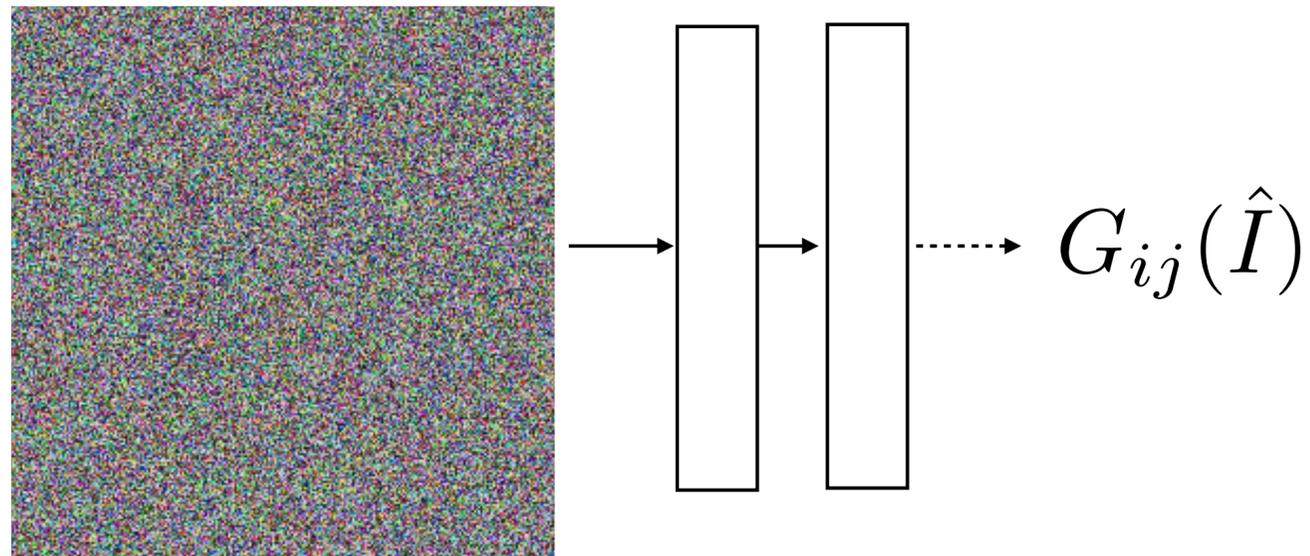
Gram ( $\approx$  covariance) matrix:

$$G_{ij} = \sum_{x=1}^w \sum_{y=1}^h c_i(x, y) c_j(x, y)$$

[Gatys et al. 2016]

Idea: correlations between unit activations convey texture.  
Discard global spatial information.

# Matching image statistics

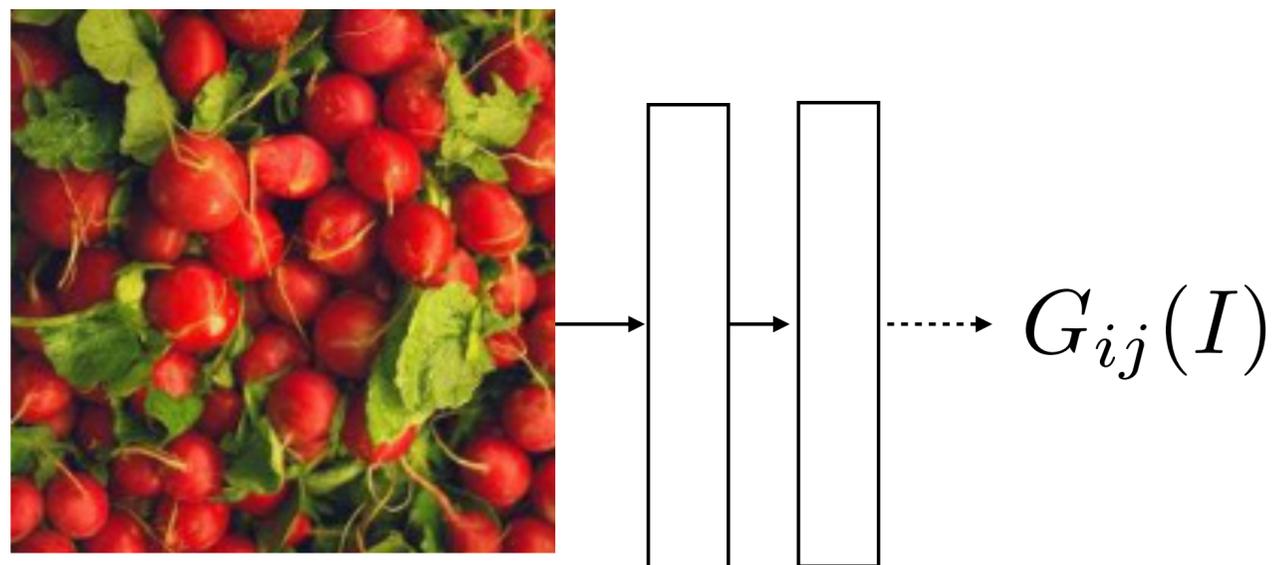


Find  $\hat{I}$  by minimizing:

$$\sum_{i=1}^{128} \sum_{j=1}^{128} (G_{ij}(I) - G_{ij}(\hat{I}))^2$$

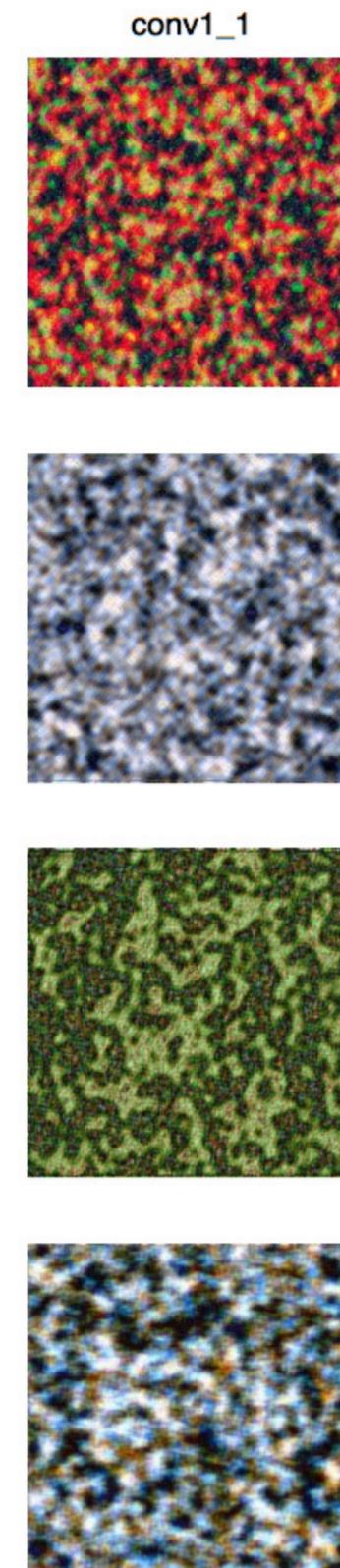
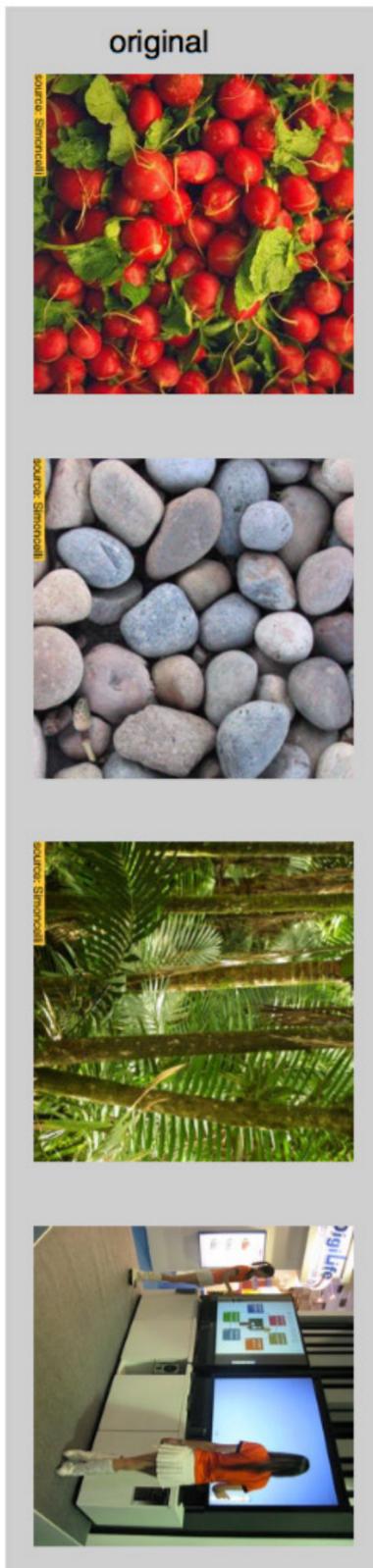
Implementation details:

- Minimize with gradient descent
- How do we compute gradients? Backprop!
- Use many layers of network.

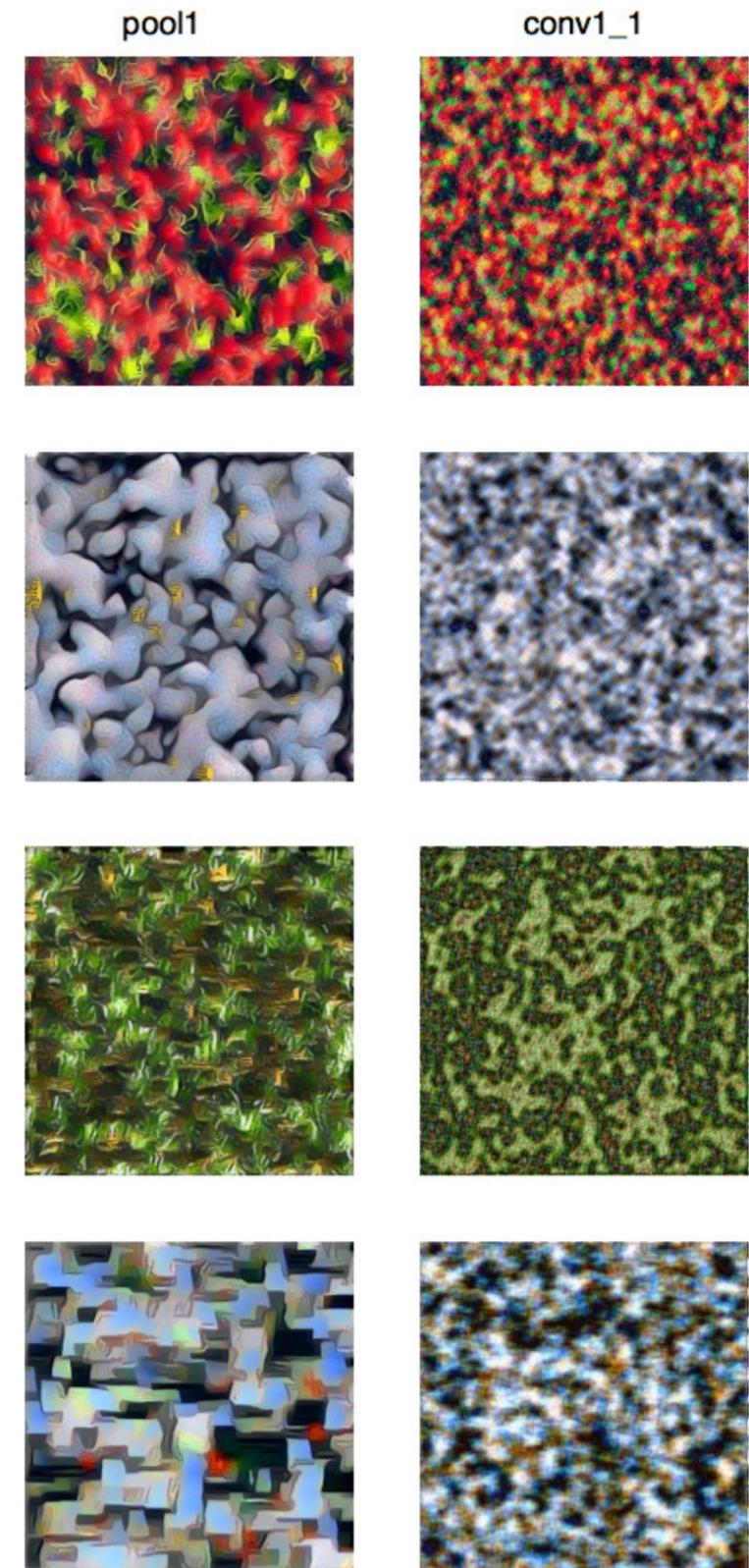
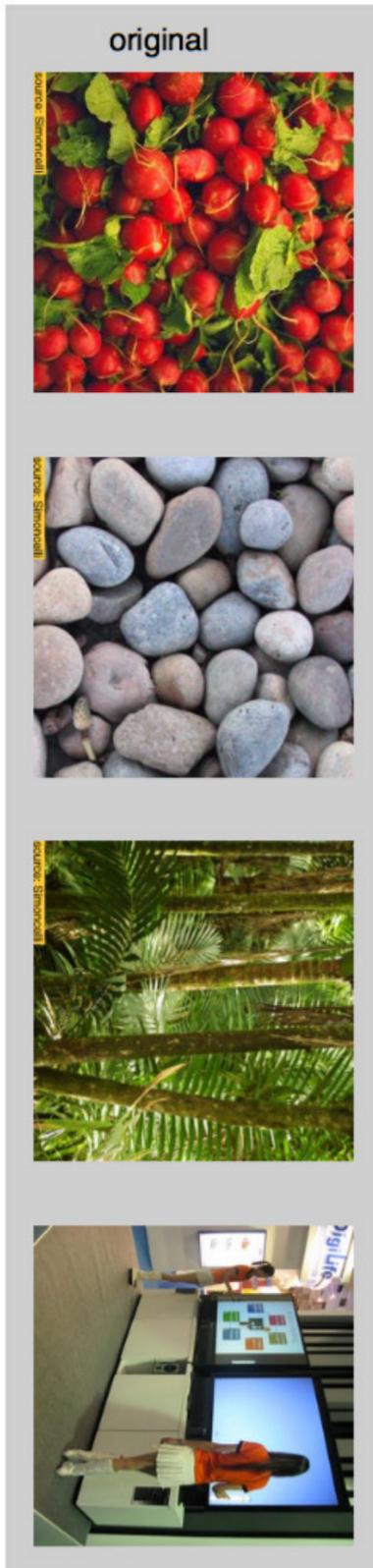


Target

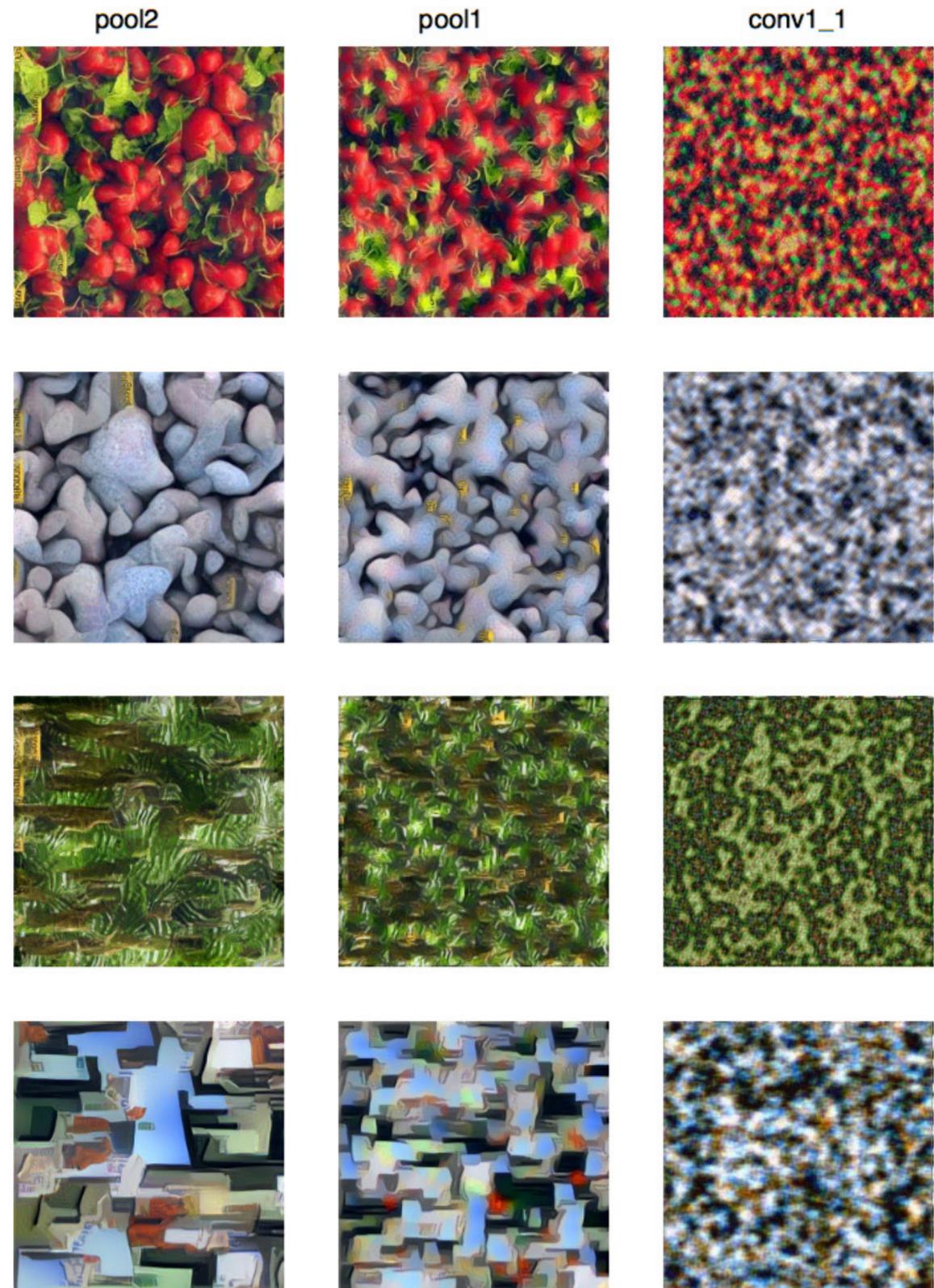
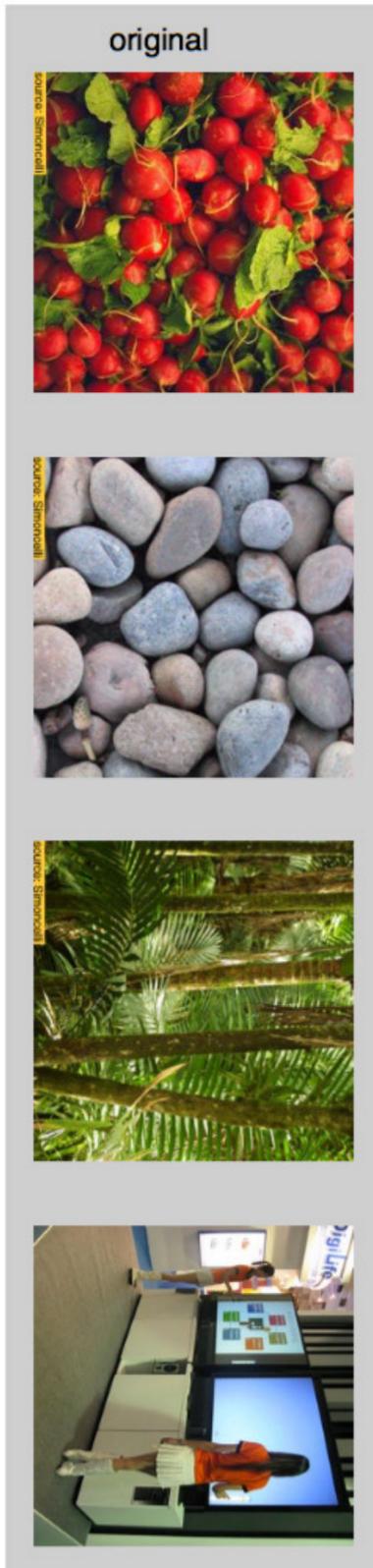
High  $\longrightarrow$  Low



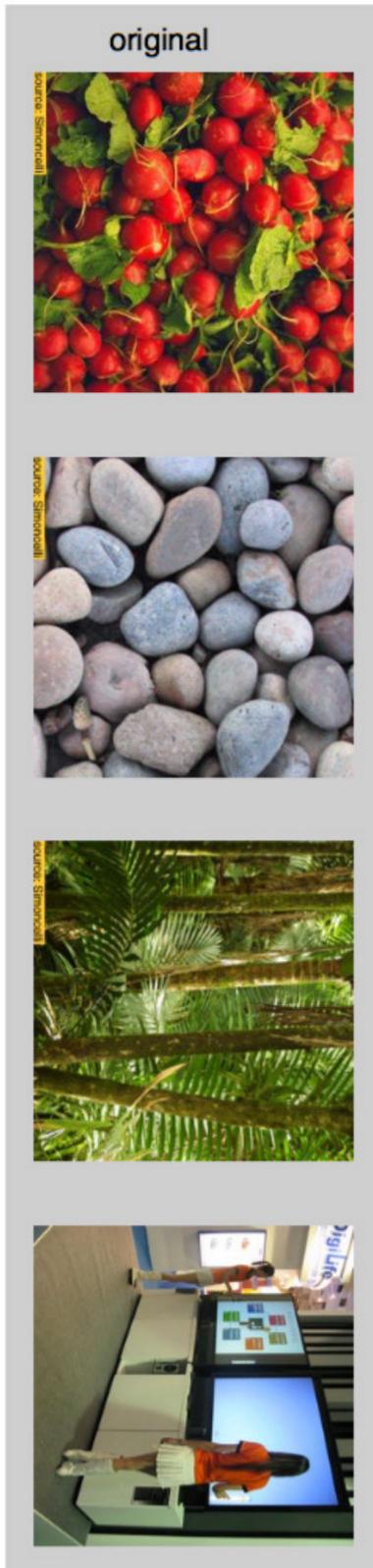
High  $\longrightarrow$  Low



High  $\longrightarrow$  Low



High  $\longrightarrow$  Low



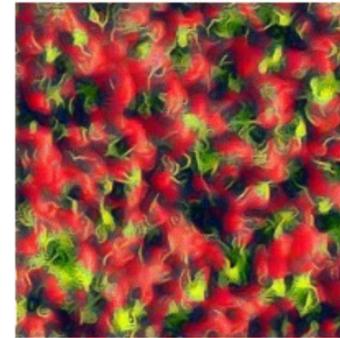
pool3



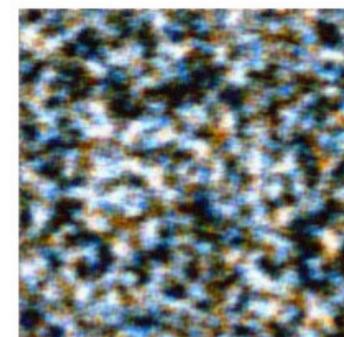
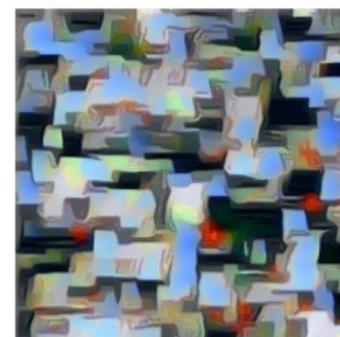
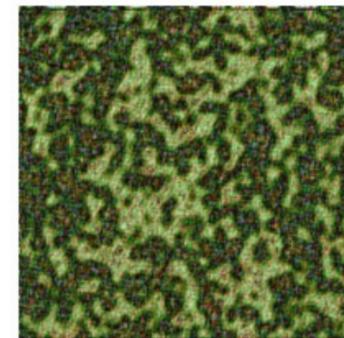
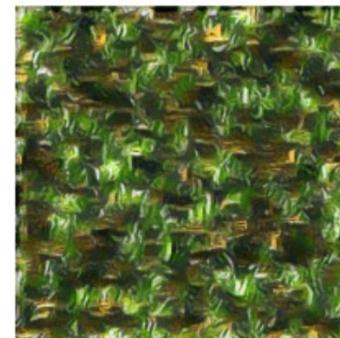
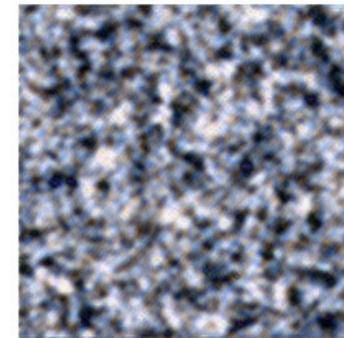
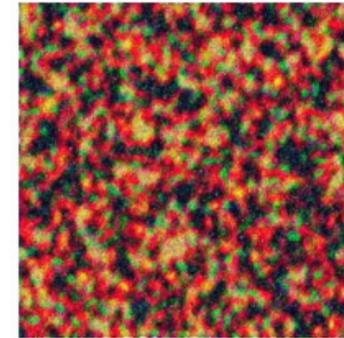
pool2



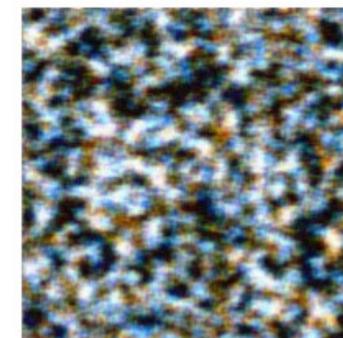
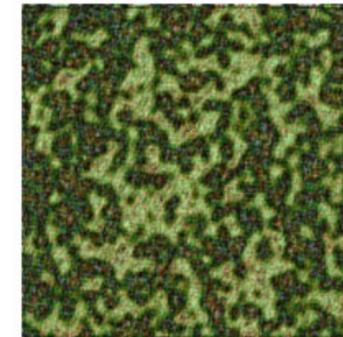
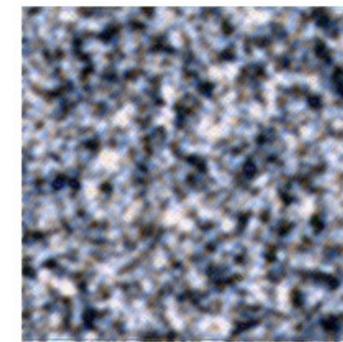
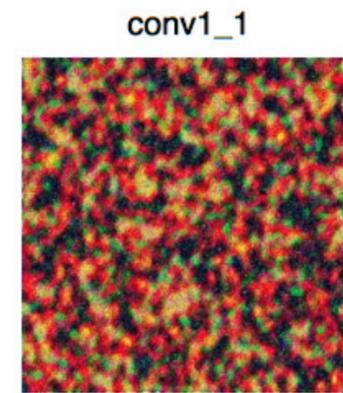
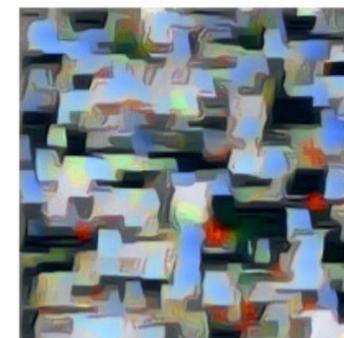
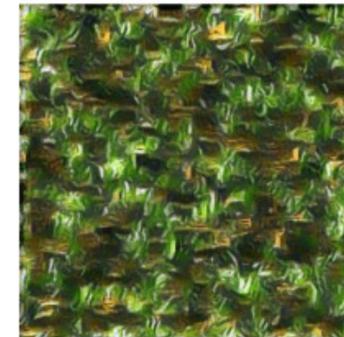
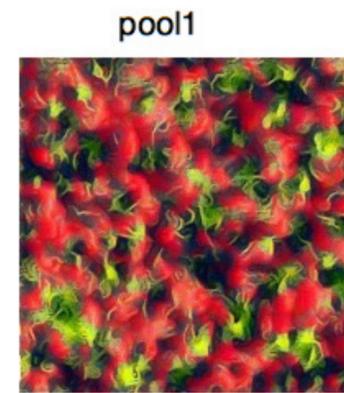
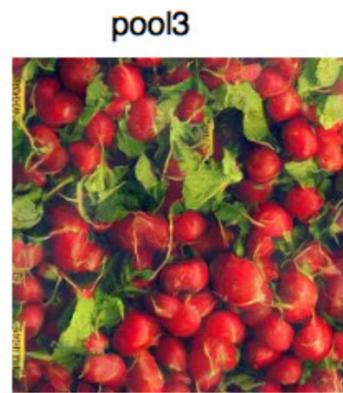
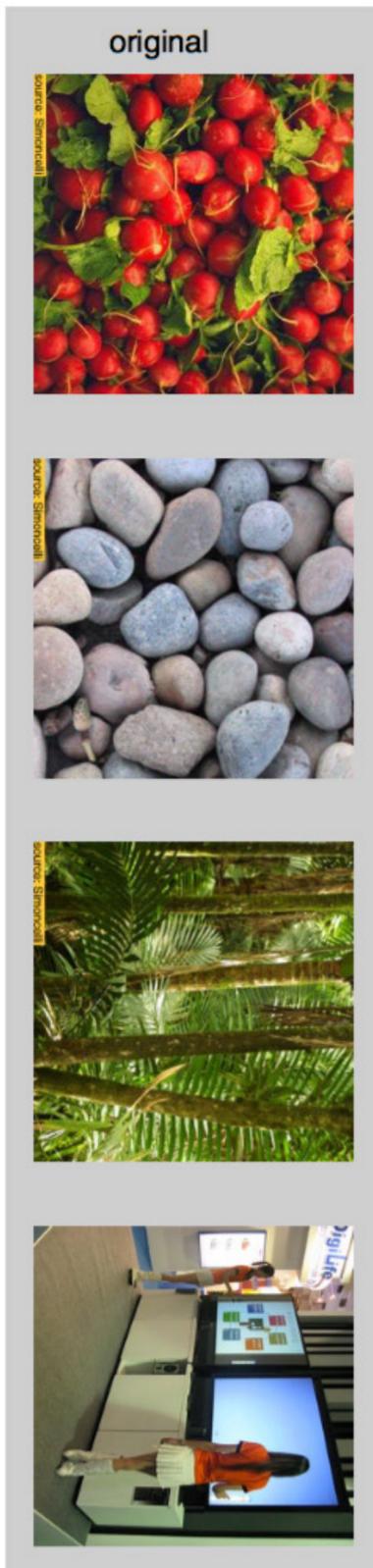
pool1



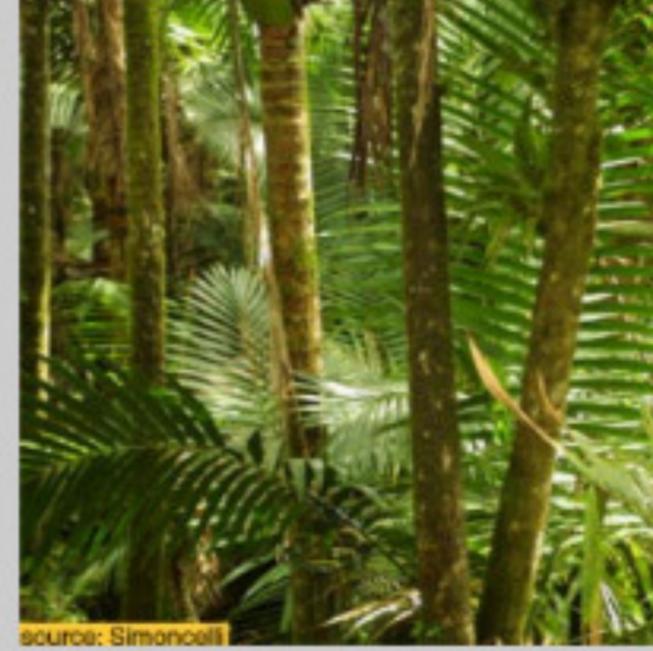
conv1\_1



High  $\longrightarrow$  Low



original



pool4



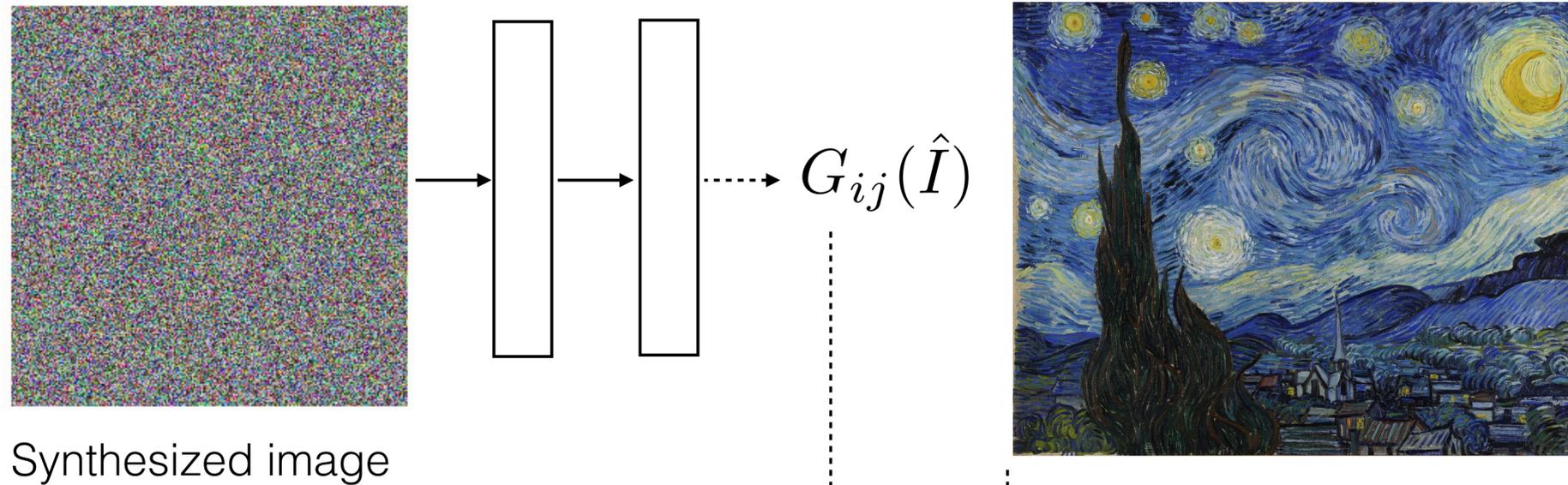
# Texture captures artistic style

Can we transfer the style of a painting to a photo?



[Gatys et al. 2016]

Match the **style** of the painting.

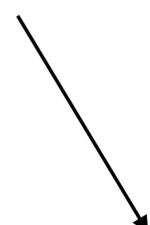


Synthesized image

$$\sum_{i=1}^{128} \sum_{j=1}^{128} (G_{ij}(\hat{I}) - G_{ij}(I))^2$$

**Perceptual loss:**

usually distance in feature space



... and the **content** of the photo.

$$\sum_i \sum_{x,y} (c_i(x,y) - \hat{c}_i(x,y))^2$$

$$c_i(x,y)$$

$$\hat{c}_i(x,y)$$











London during the day.

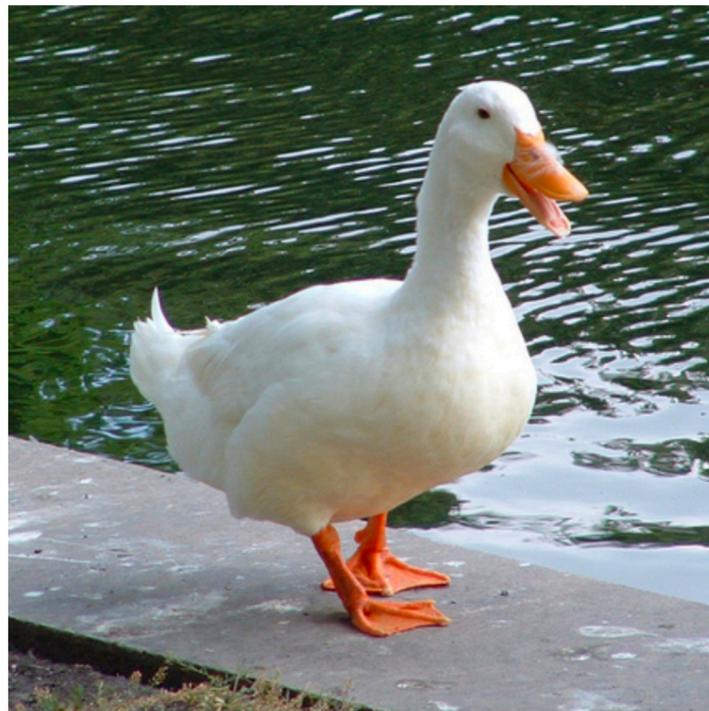


New York at night.



# Neural networks that generate images

# Image classification



⋮

image  $\mathbf{x}$

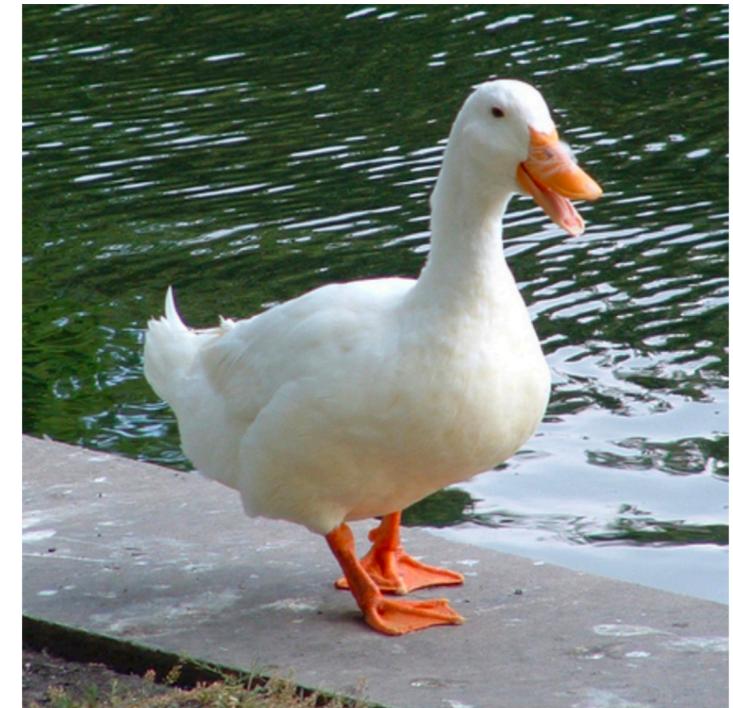
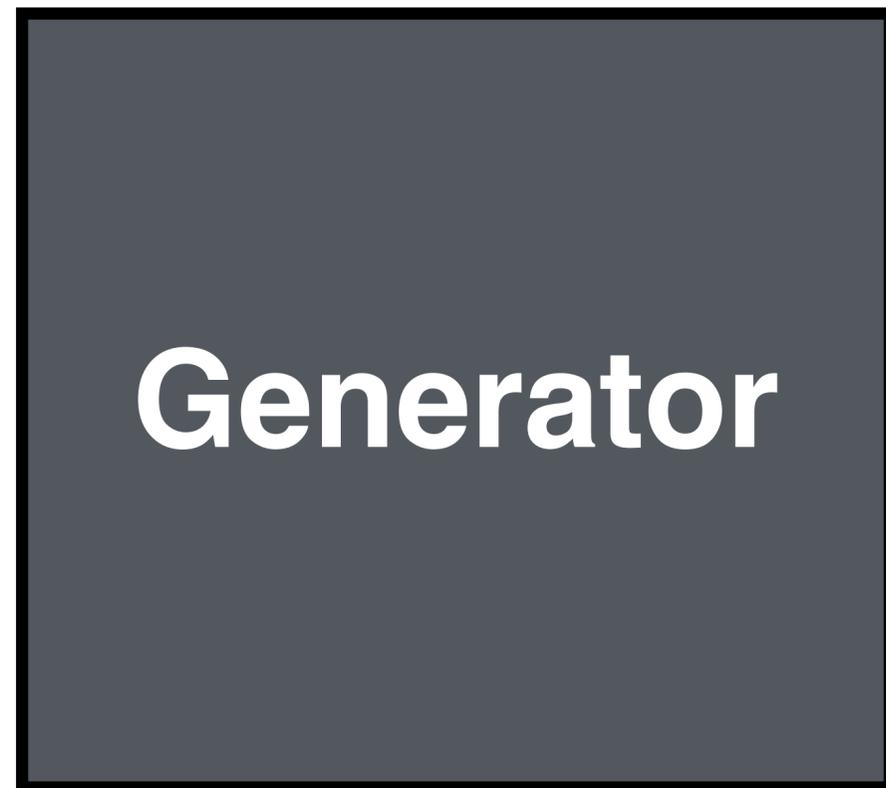


“Duck”

label  $y$

# Image synthesis

“Duck”



⋮

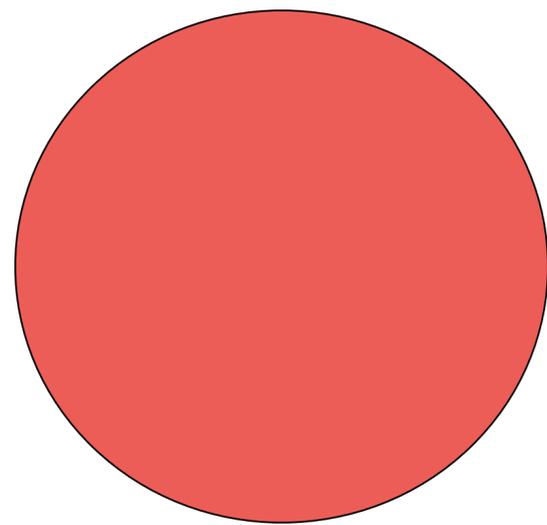
label  $y$

image  $x$

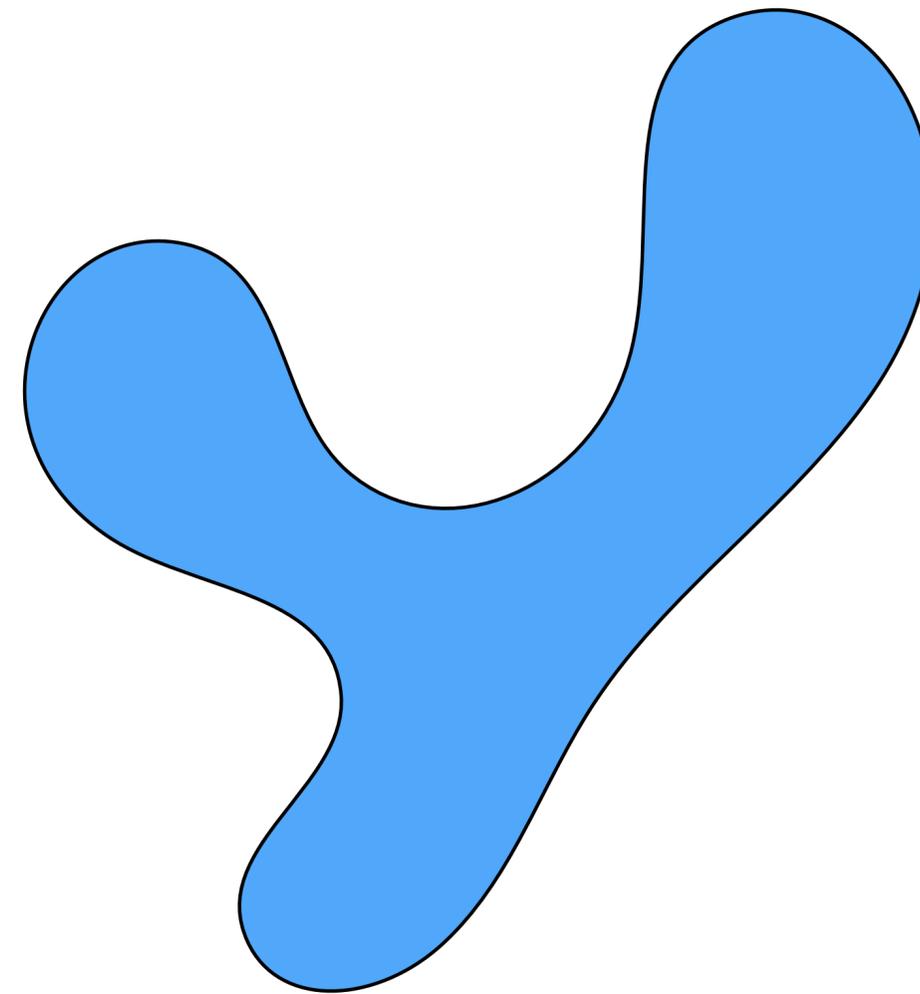
# Neural networks as distribution transformers

Source distribution

Target distribution

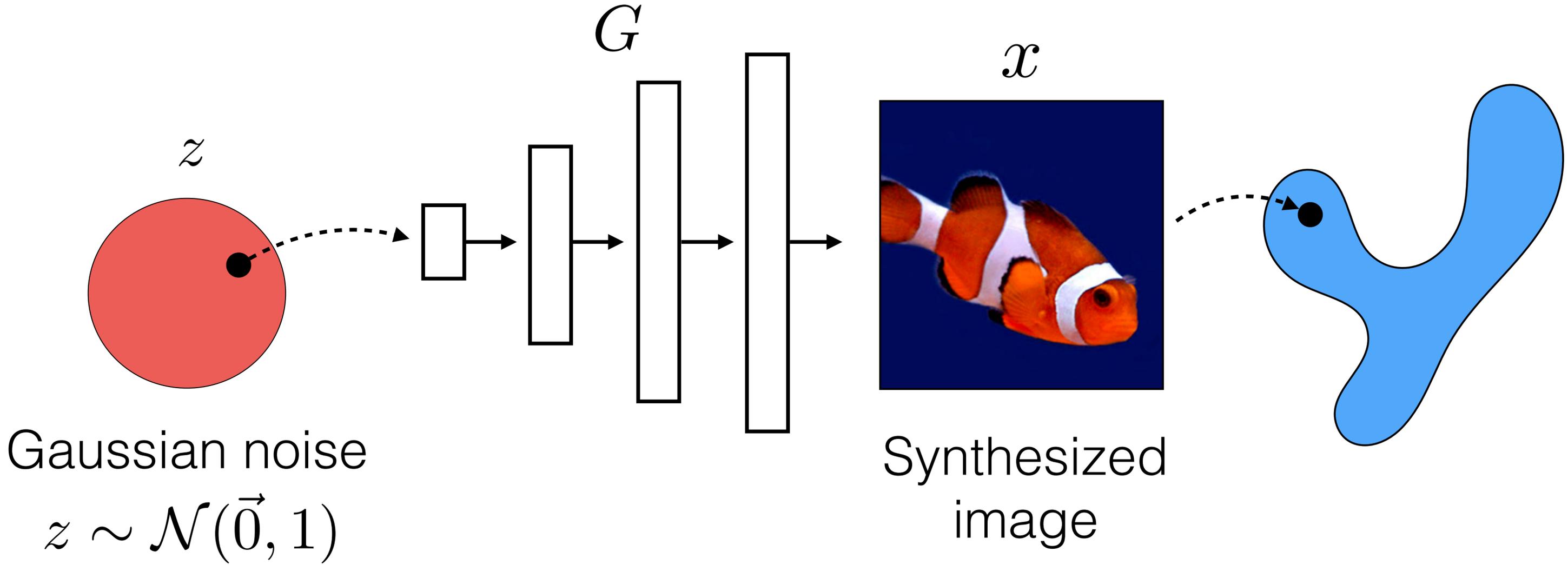


$p(z)$

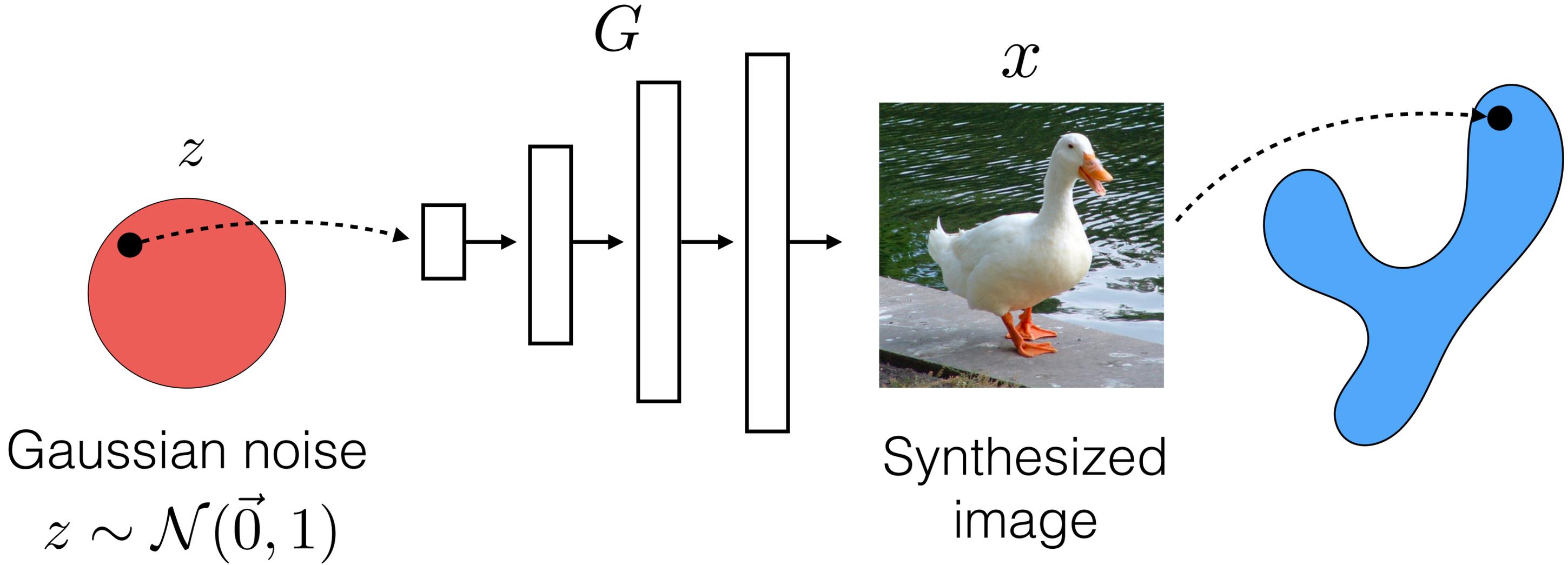


$p(x)$

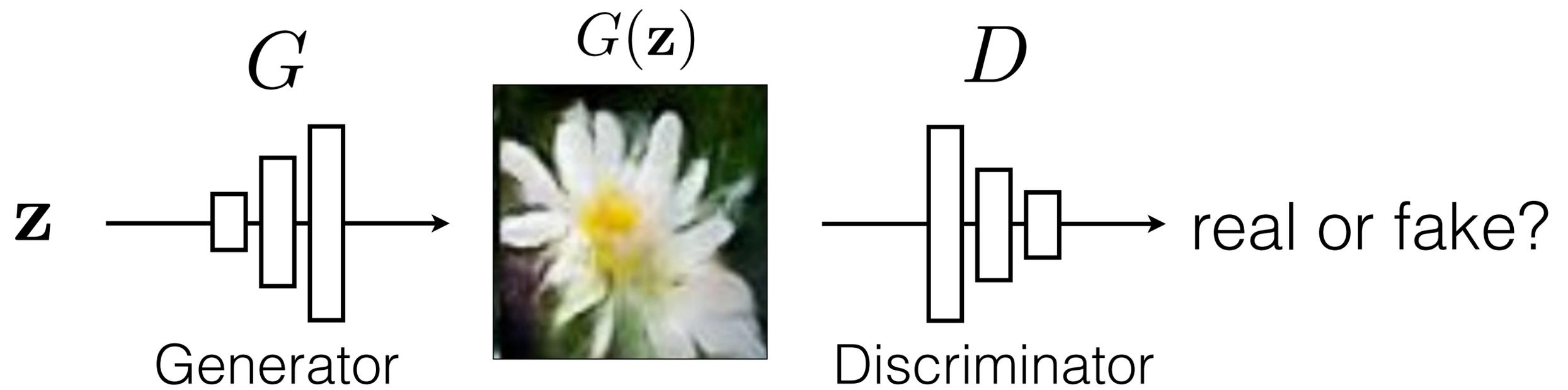
# Neural networks as distribution transformers



# Neural networks as distribution transformers

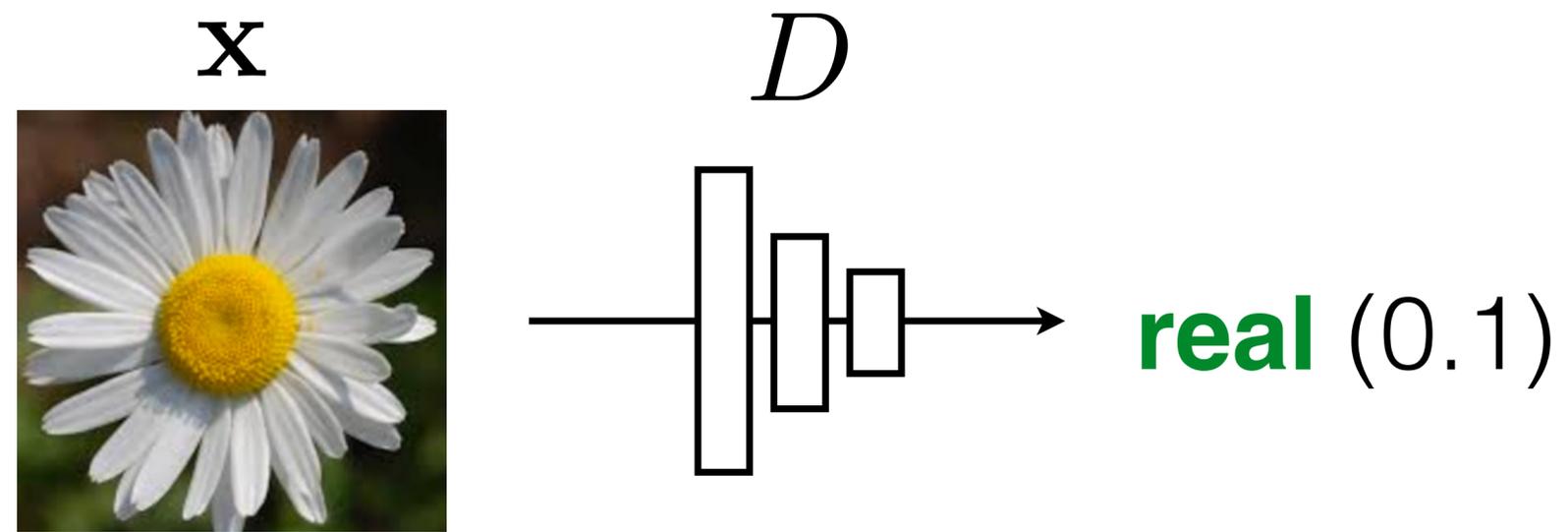
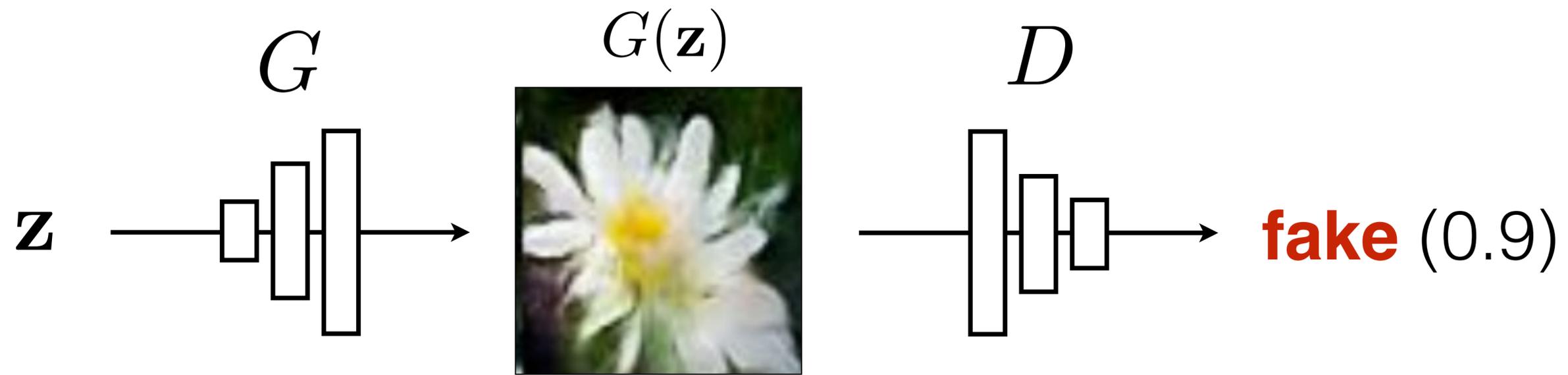


# Generative adversarial networks (GANs)

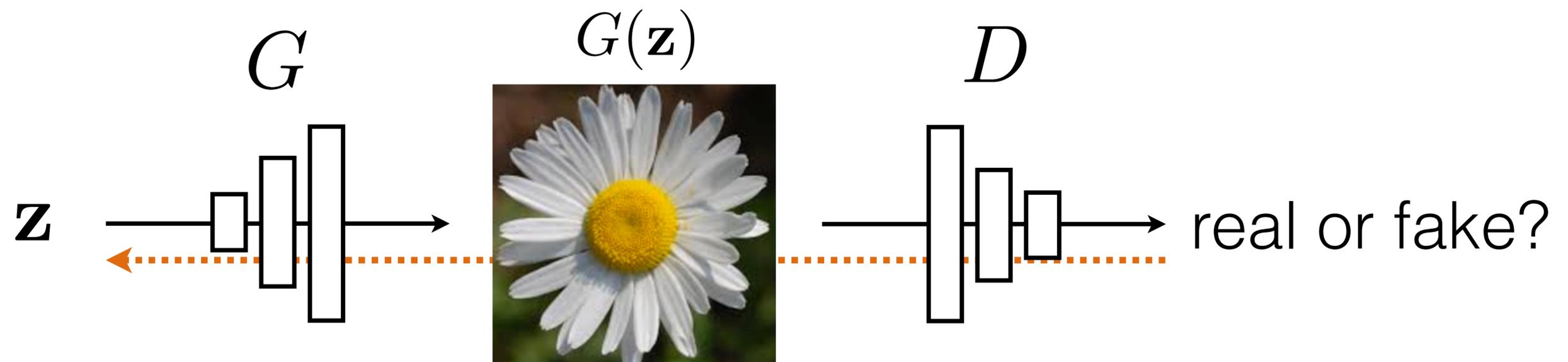


**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes

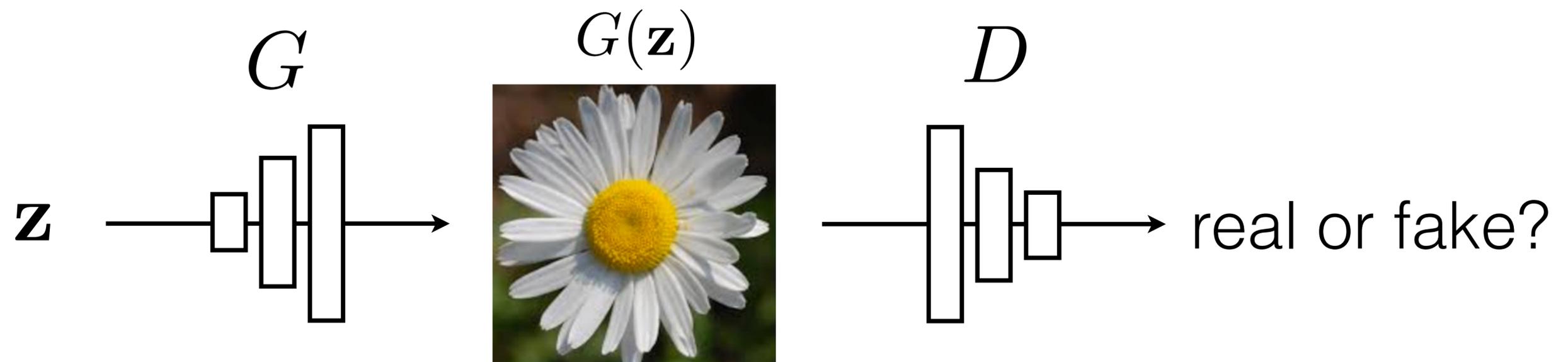


$$\arg \max_D \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[ \log D(G(\mathbf{z})) + \log (1 - D(\mathbf{x})) \right]$$



**G** tries to synthesize fake images that **fool D**:

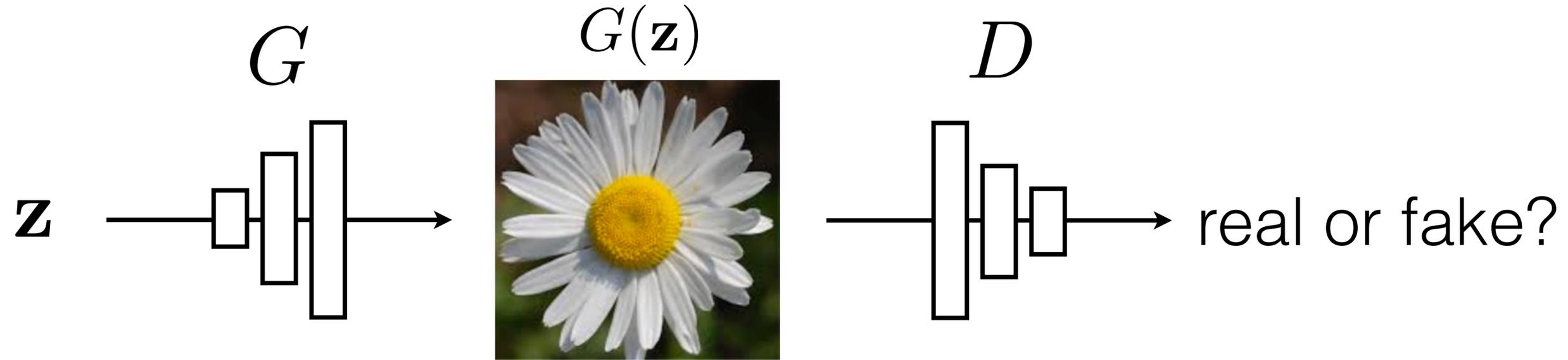
$$\arg \min_G \mathbb{E}_{\mathbf{z}, \mathbf{x}} [ \log D(G(\mathbf{z})) + \log (1 - D(\mathbf{x})) ]$$



**G** tries to synthesize fake images that *fool* the *best* **D**:

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

# Training



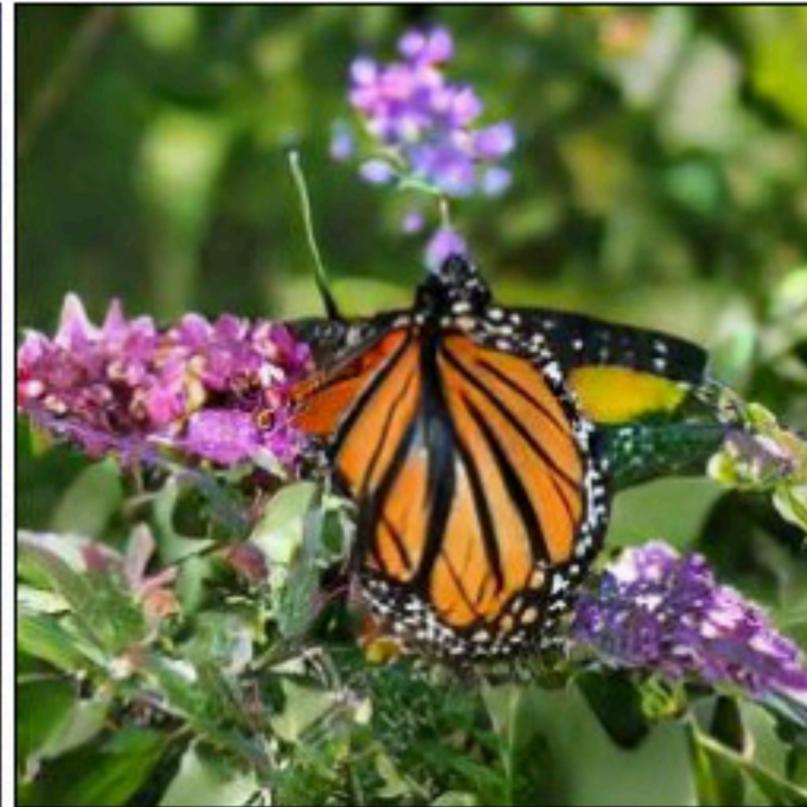
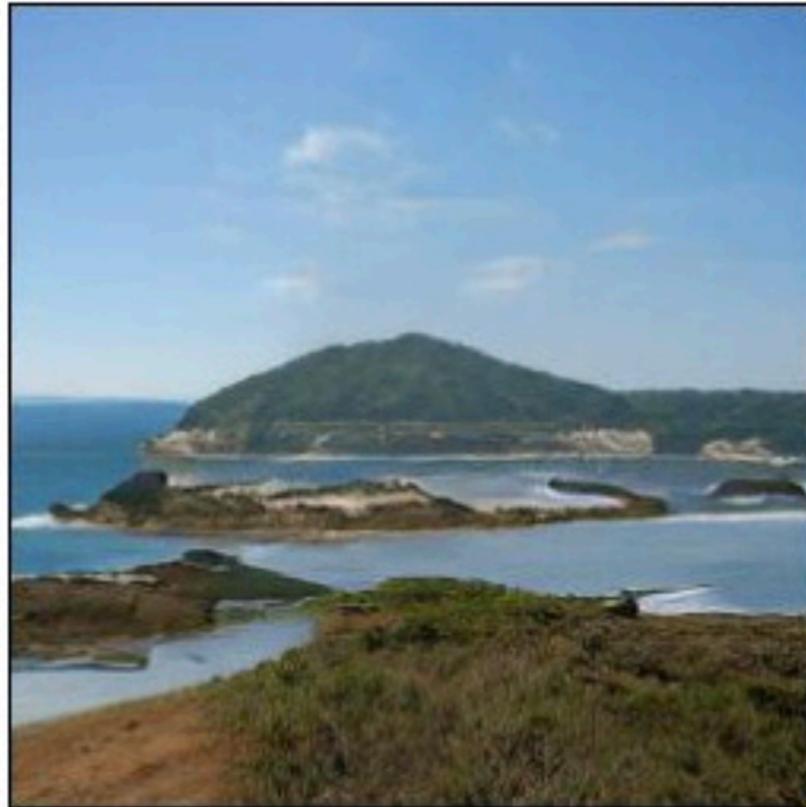
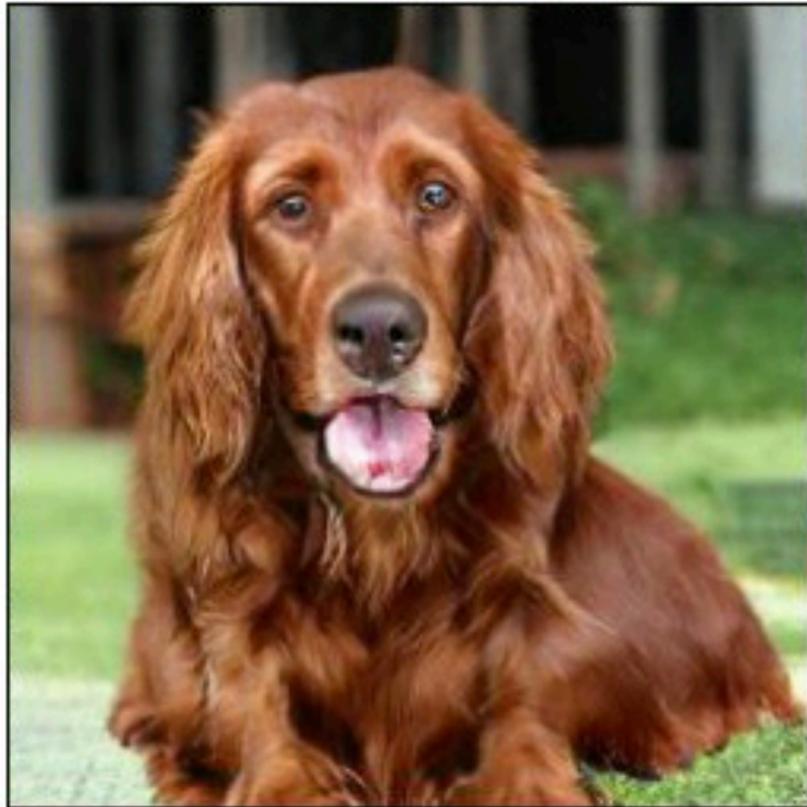
**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes

- Training: iterate between training  $D$  and  $G$  with backprop.
- Global optimum when  $G$  reproduces data distribution (see book)

# Samples from BigGAN

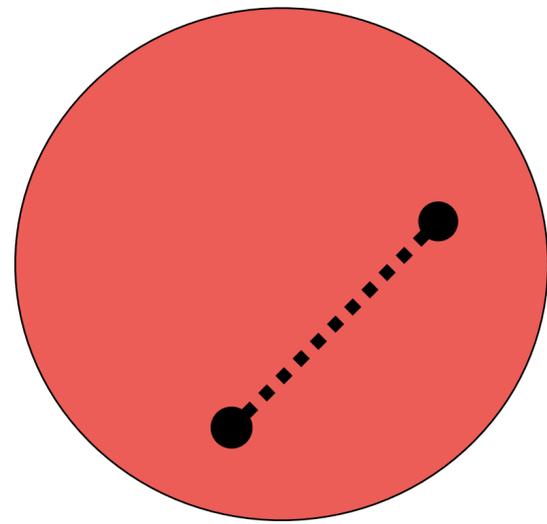
[Brock et al. 2018]



More here: <https://arxiv.org/pdf/1809.11096.pdf>

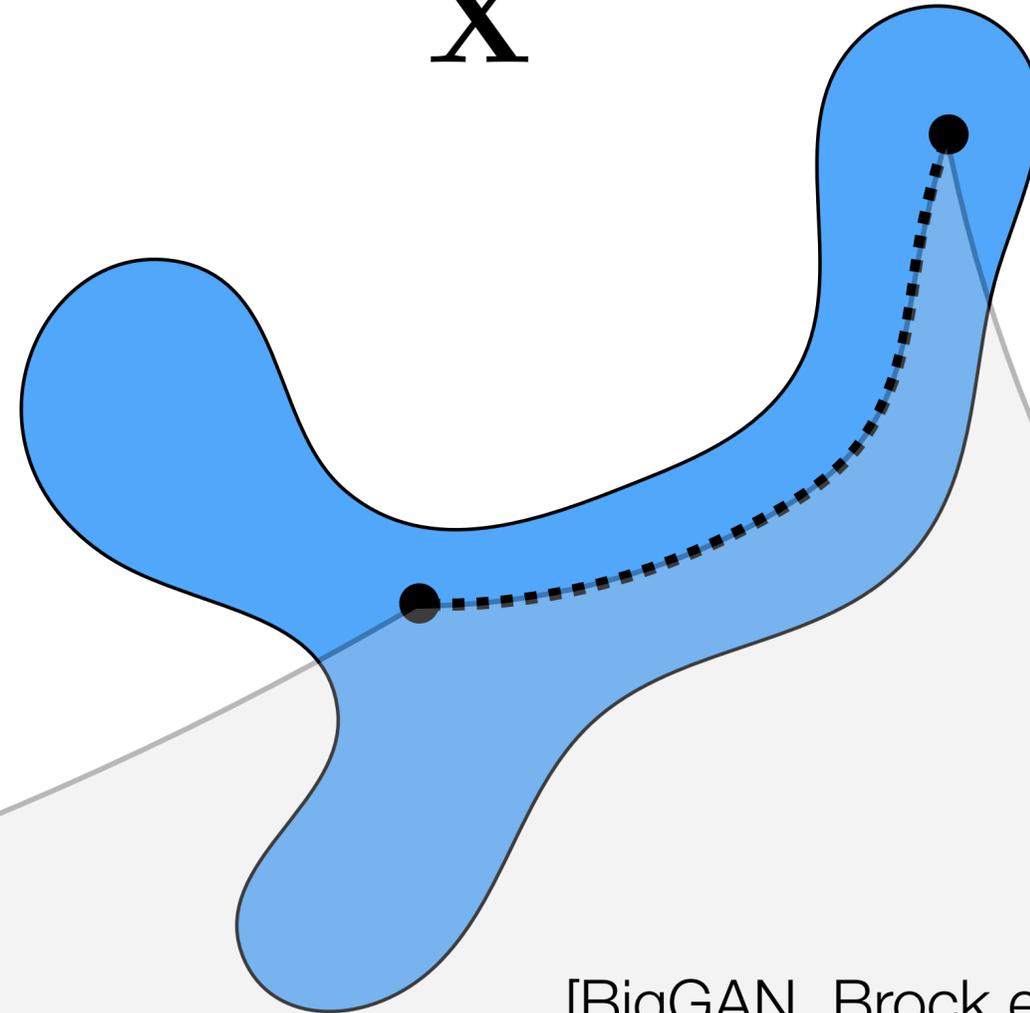
Latent space  
(Gaussian)

$\mathbf{z}$



Data space  
(Natural image manifold)

$\mathbf{X}$



[BigGAN, Brock et al. 2018]





StyleGAN2



StyleGAN3 (Ours)



StyleGAN2



StyleGAN3 (Ours)

[Karras et al., "Alias-Free Generative Adversarial Networks", 2021]



StyleGAN2

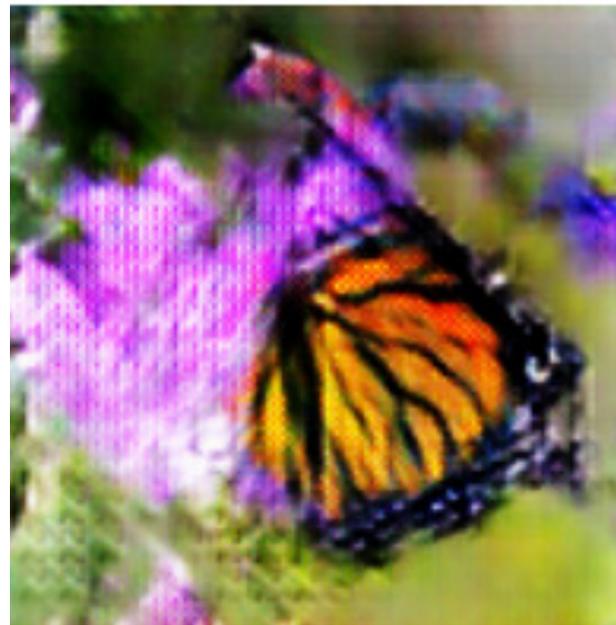


StyleGAN3 (Ours)

[Karras et al., "Alias-Free Generative Adversarial Networks", 2021]

# Rapid progress due mostly to better architectures

ACGAN [Odena et al. 2016]



BigGAN [Brock et al. 2018]



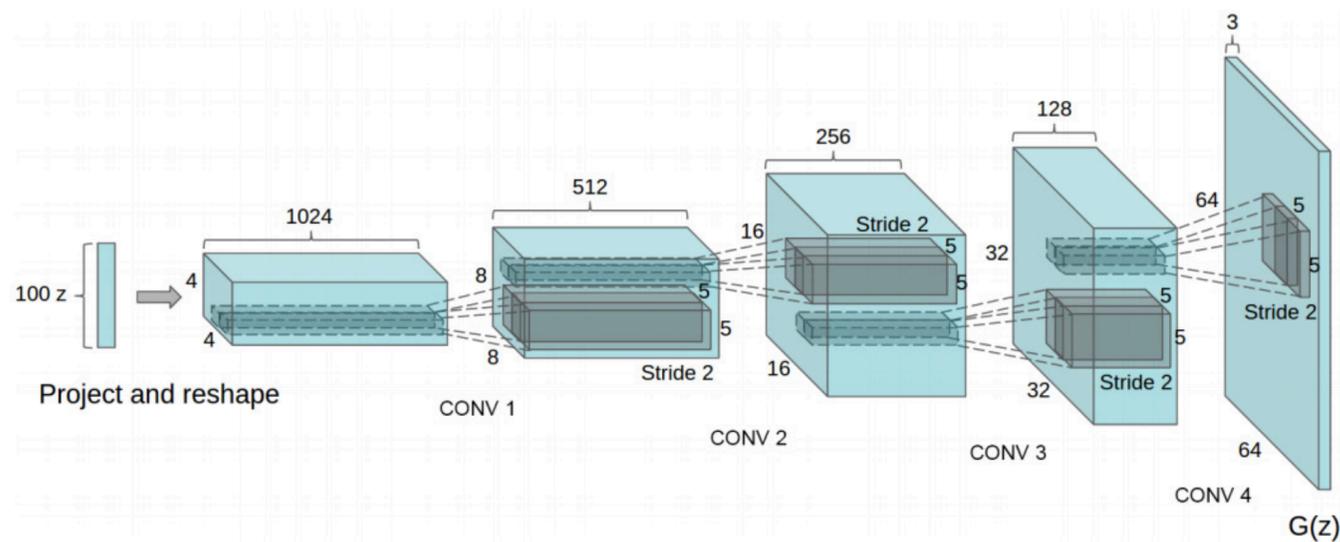
52

Both trained on ImageNet

# Architectures

## DCGAN

[Radford, Metz, Chintala 2016]

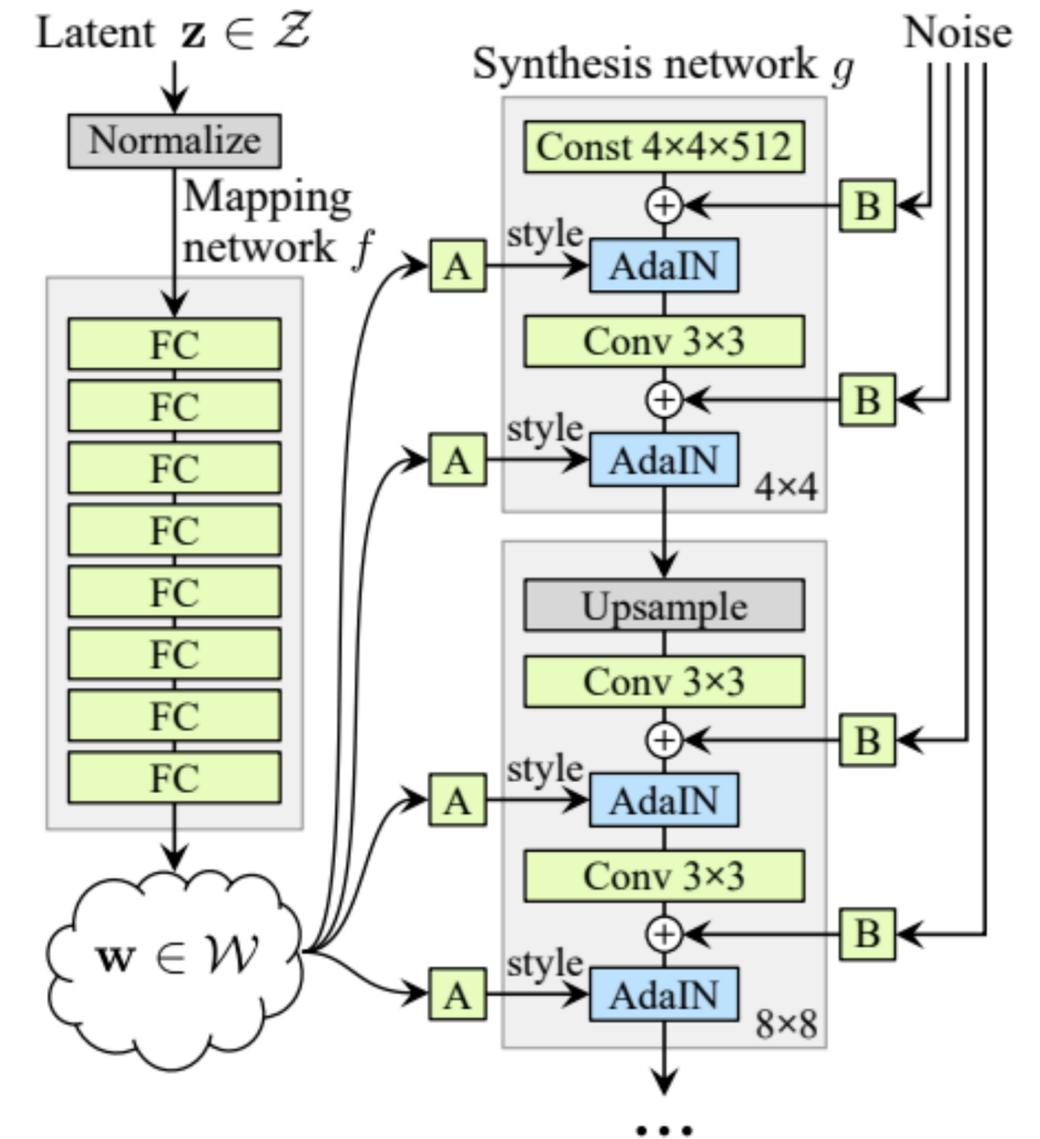


Transpose convolution + batch norm + nonlinearities



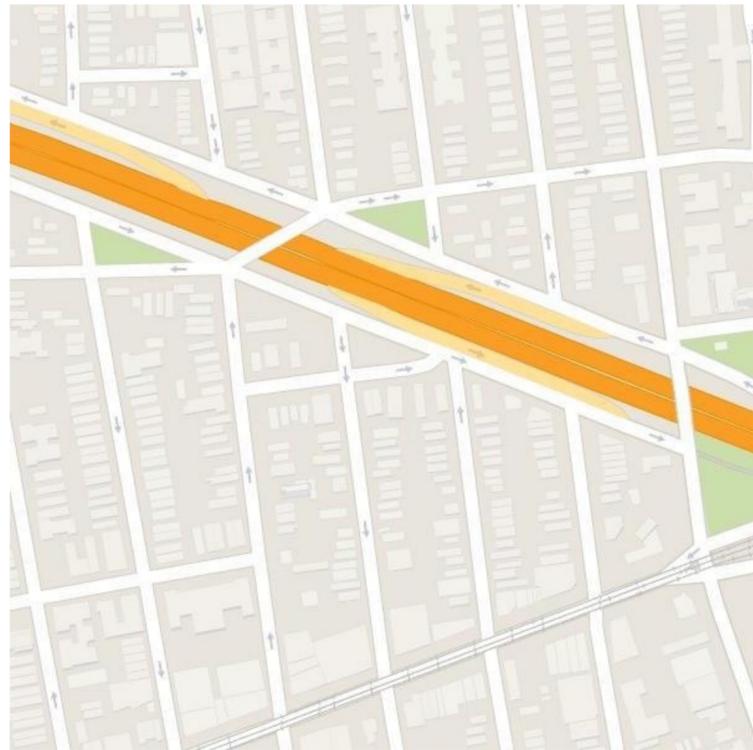
## StyleGAN

[Karras, Laine, Aila 2019]

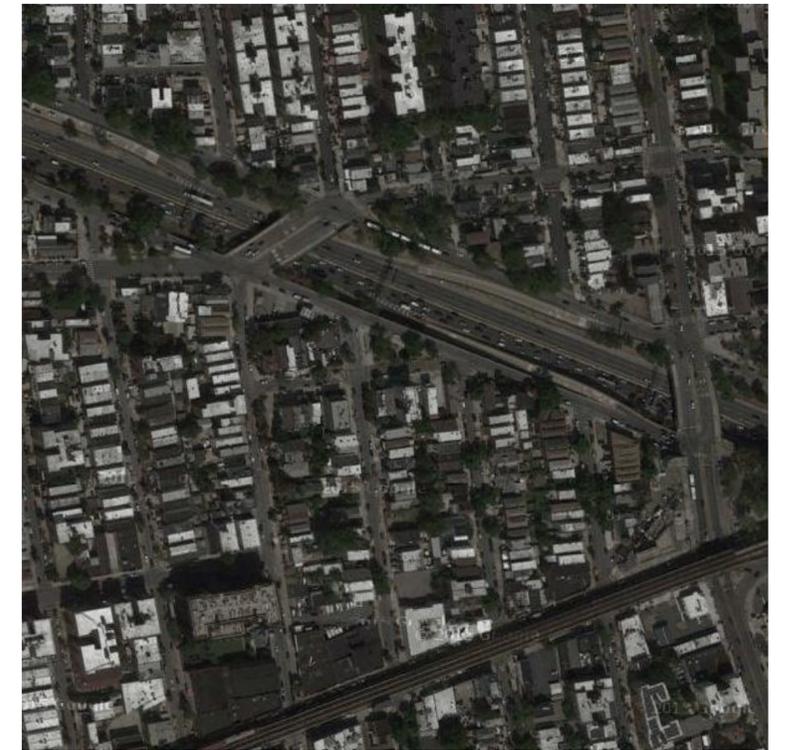


Similar but bigger and lots of engineering details.

# Image translation



Google Map

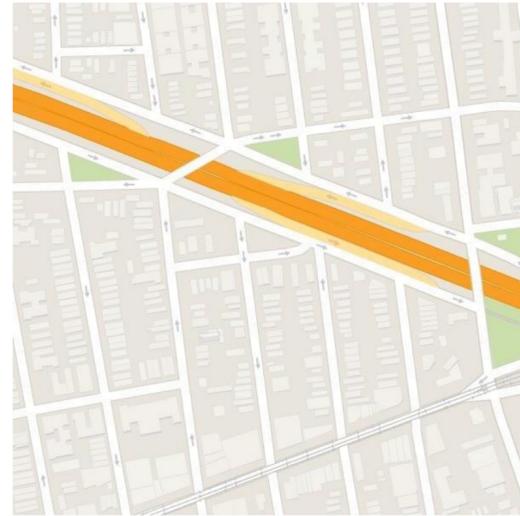
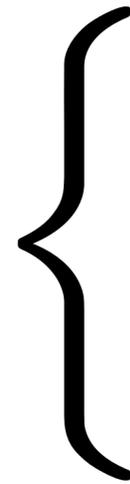


Satellite photo

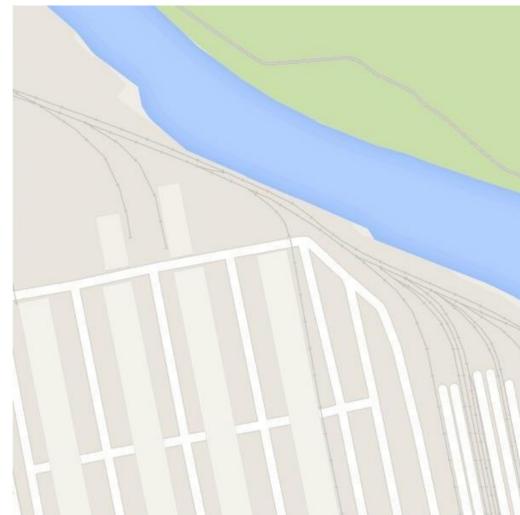
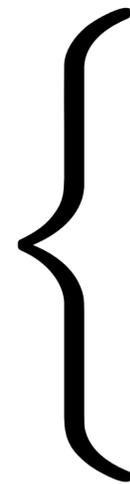
# Map2Sat

**x**

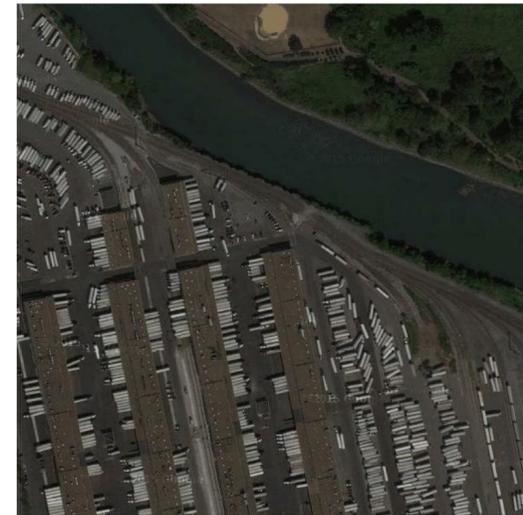
**y**



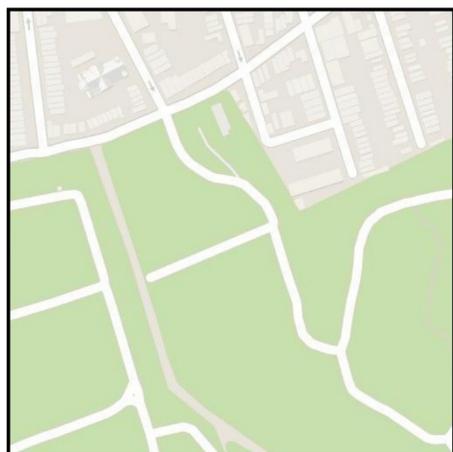
,



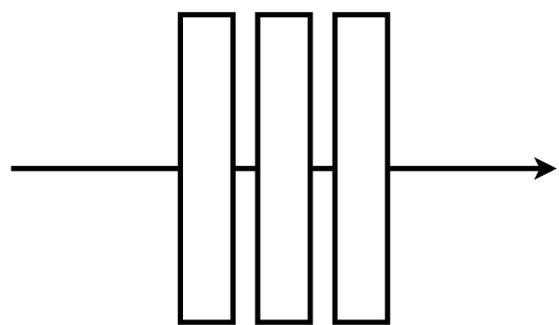
,



$\mathbf{x}$



$G$



Generator

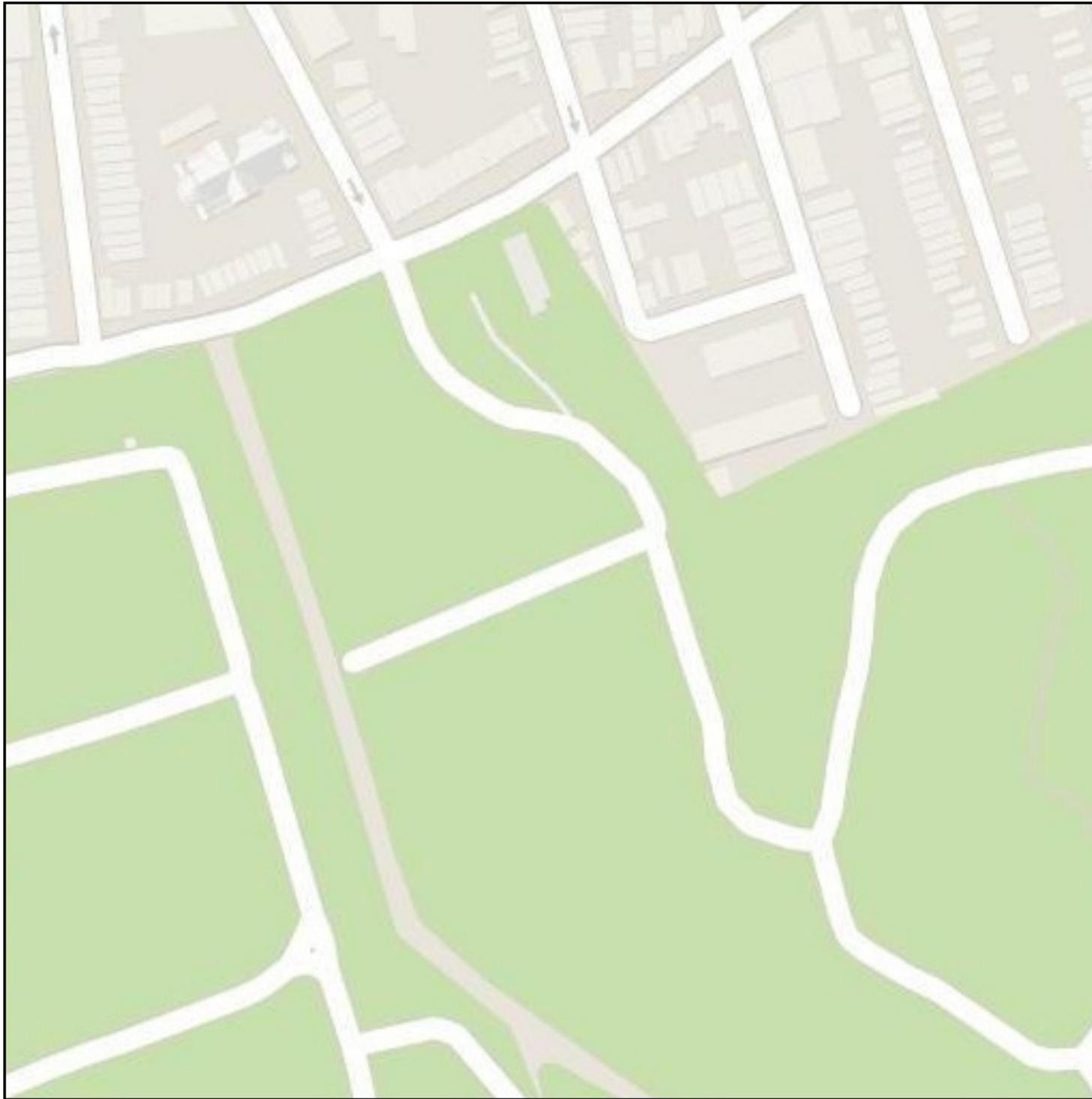
$G(\mathbf{x})$



Idea: L1 loss

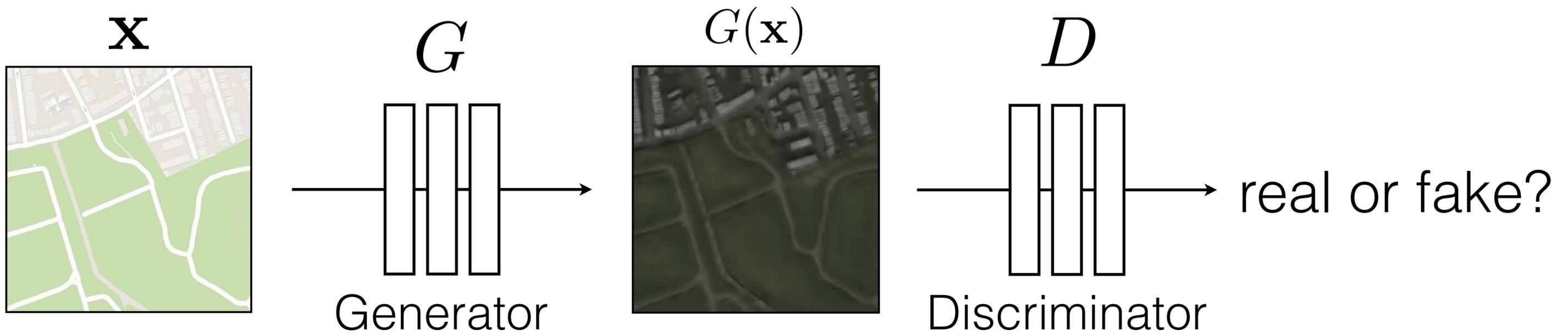
$$\|G(\mathbf{x}) - \mathbf{y}\|_1$$

Input



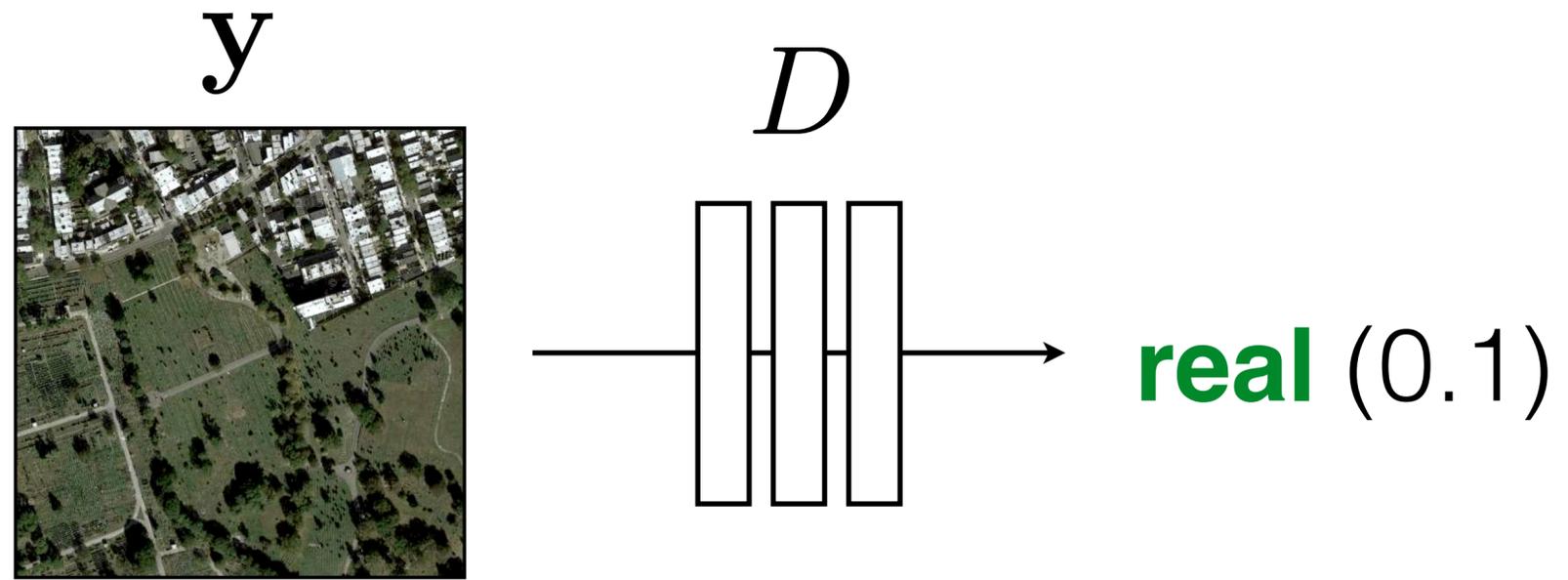
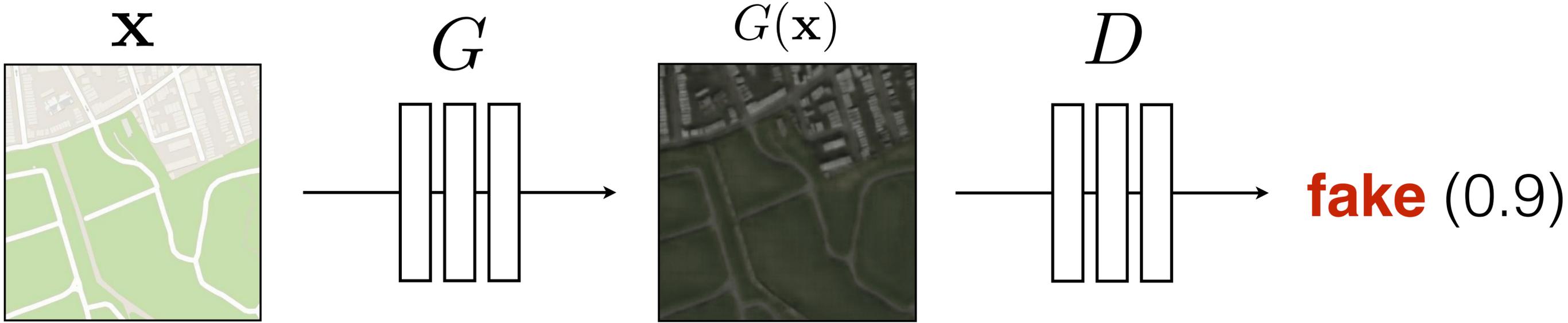
L1 loss



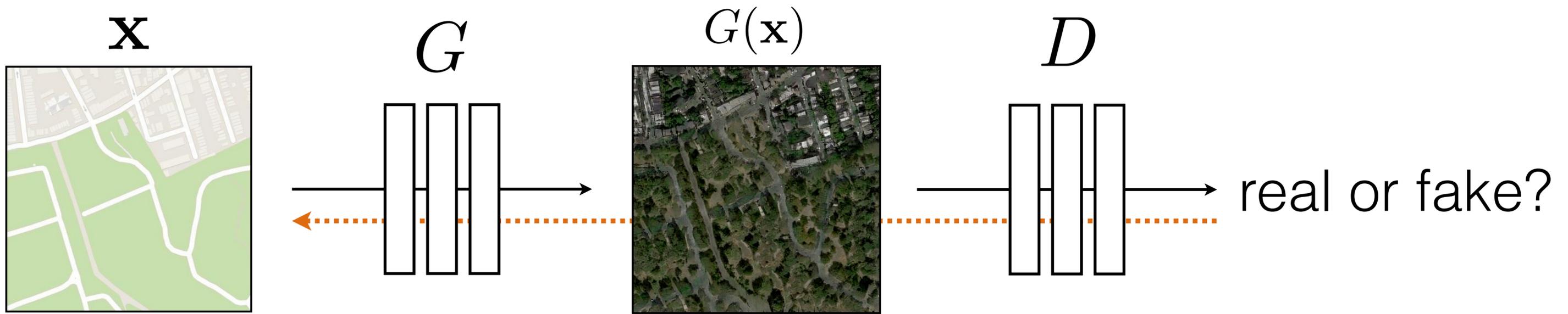


**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes

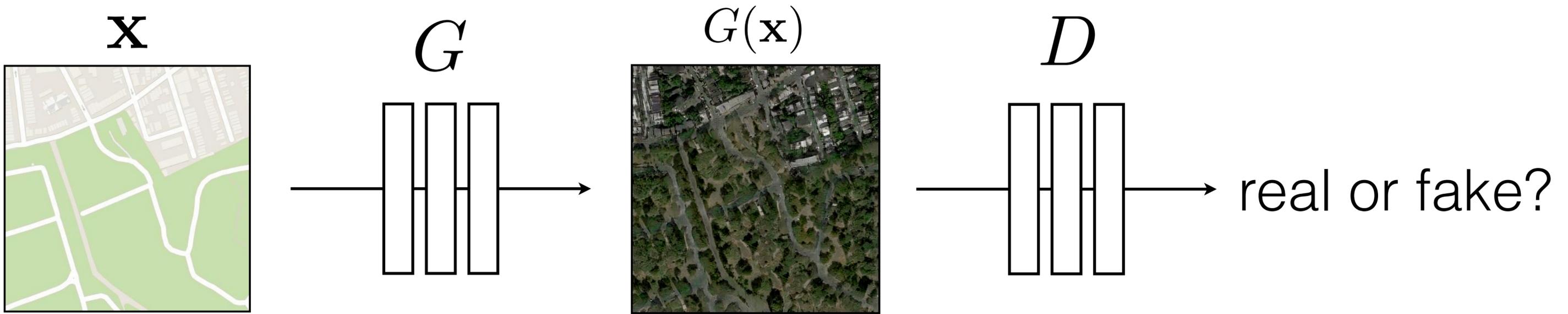


$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) \right]$$



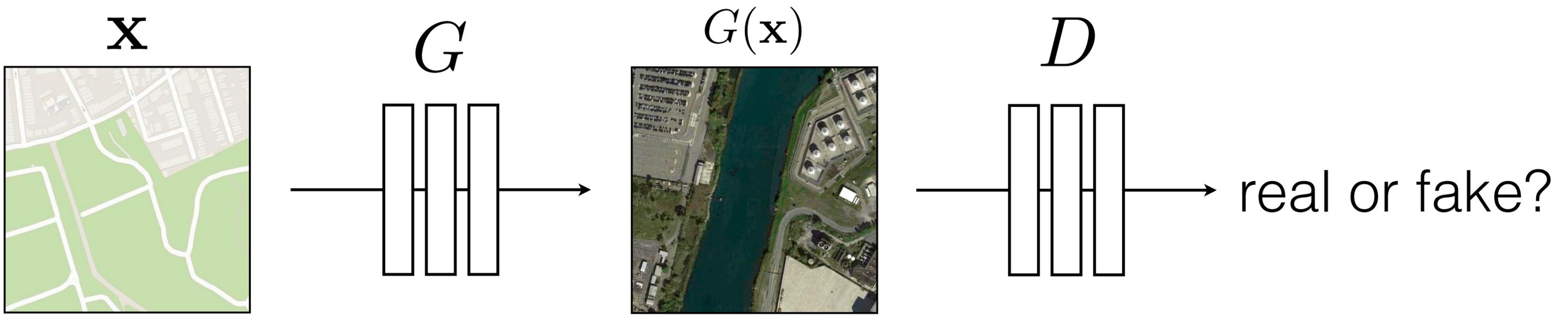
**G** tries to synthesize fake images that **fool D**:

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

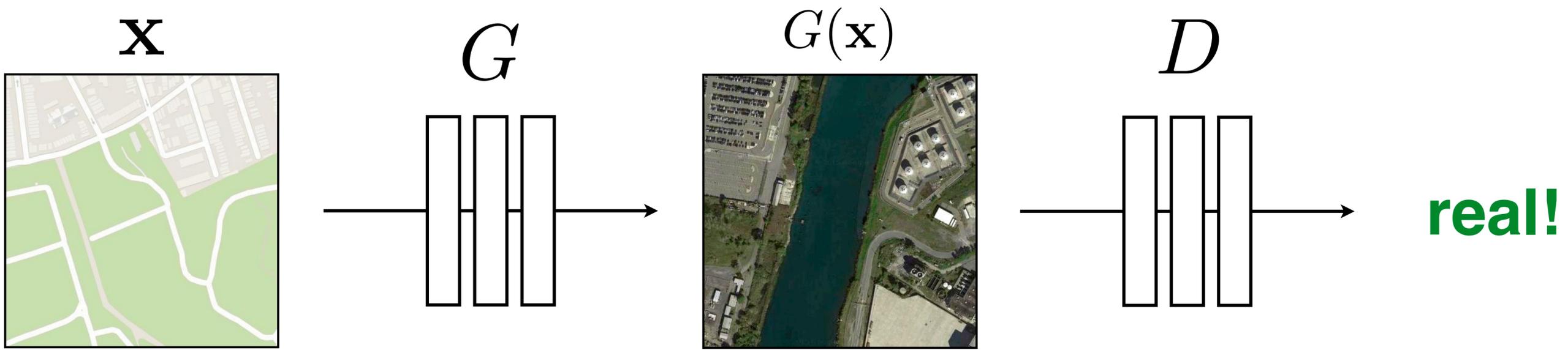


**G** tries to synthesize fake images that *fool* the *best* **D**:

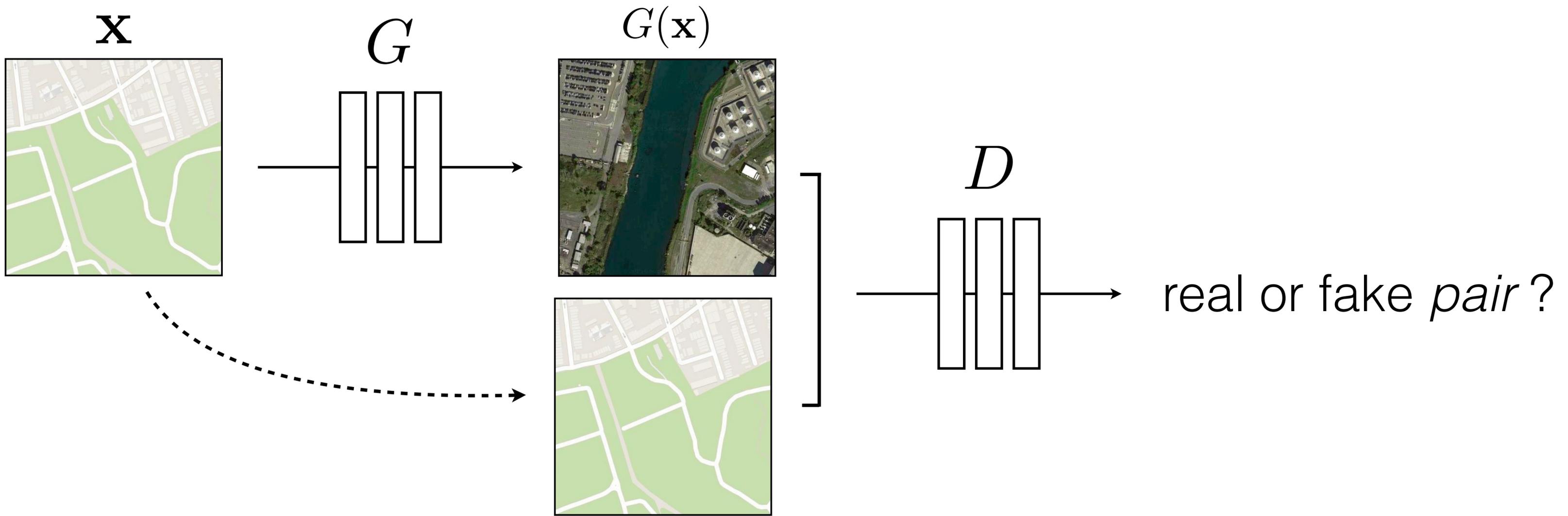
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



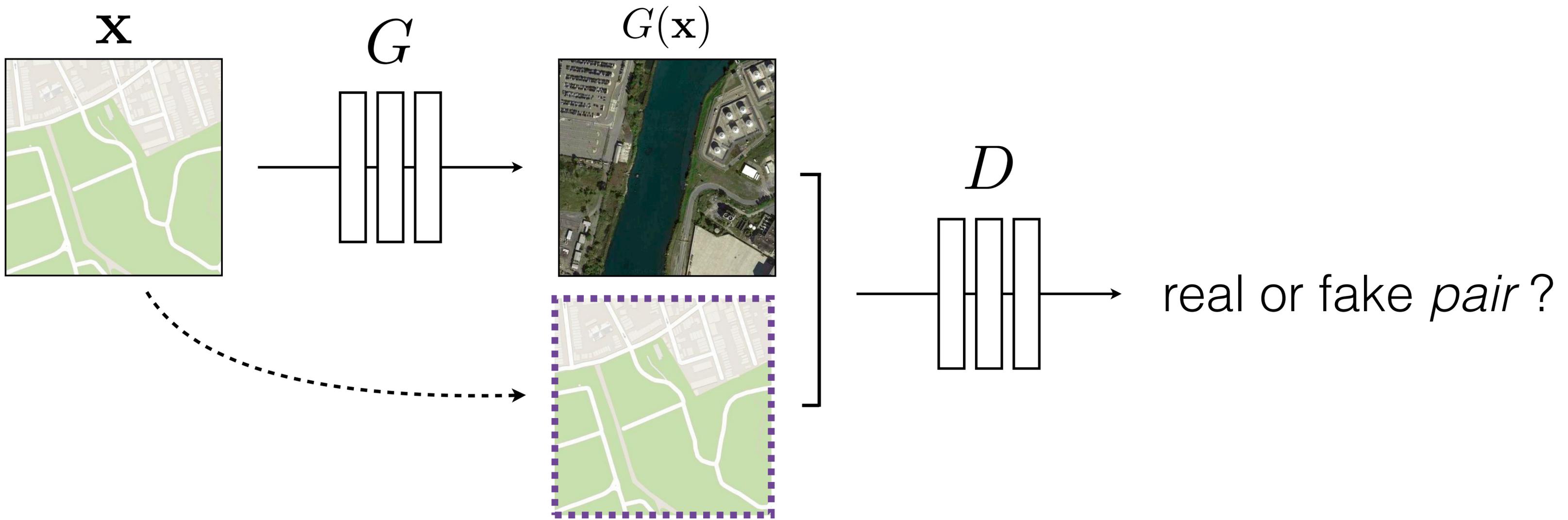
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



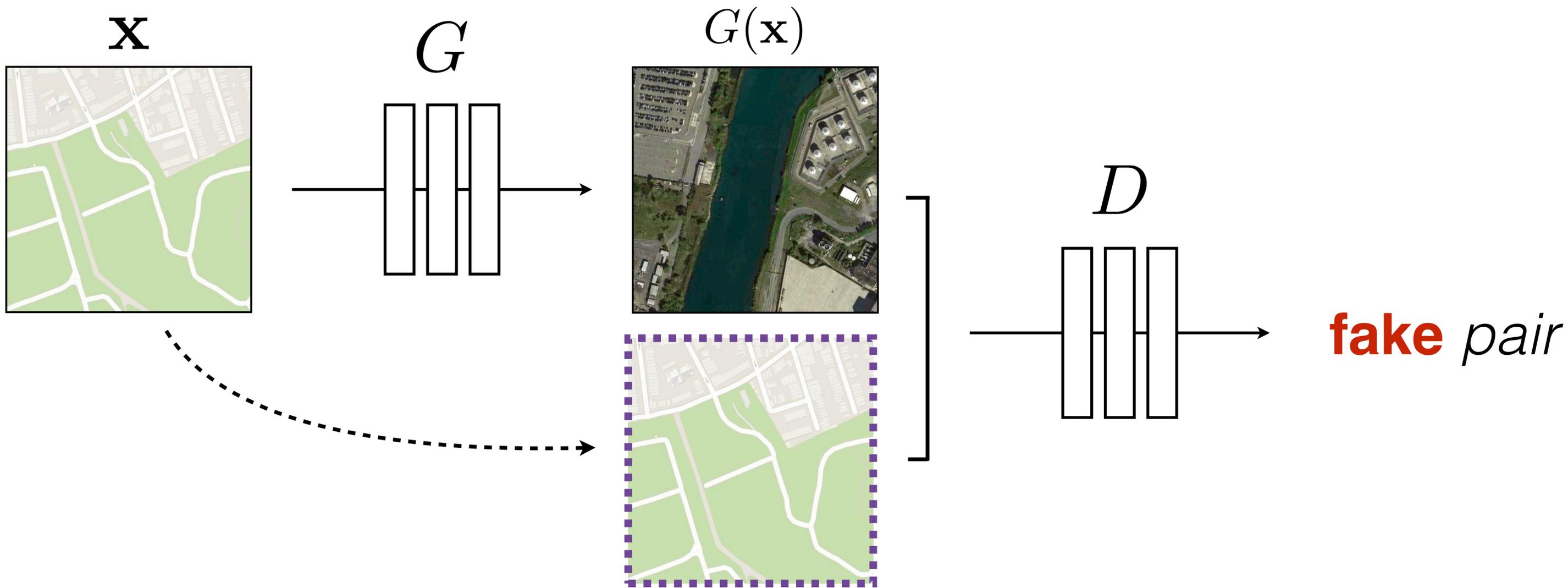
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



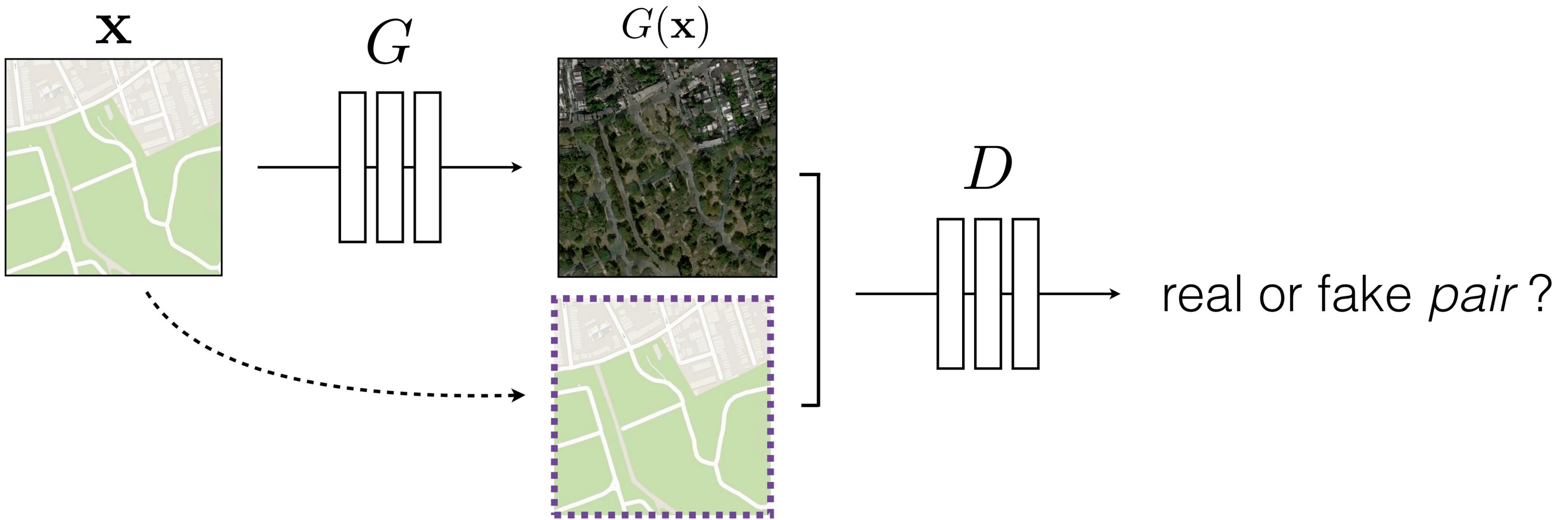
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



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



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

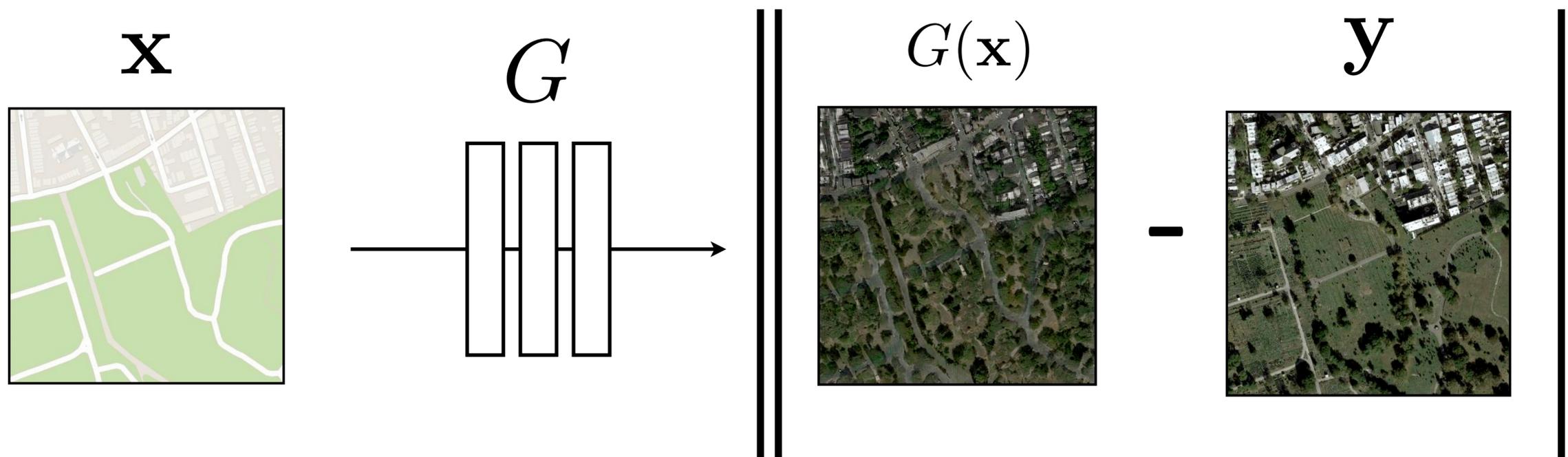


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$

# Training Details: Loss function

## Conditional GAN

$$G^* = \arg \min_G \max_D \mathcal{L}_{cGAN}(G, D) + \lambda \mathcal{L}_{L1}(G).$$

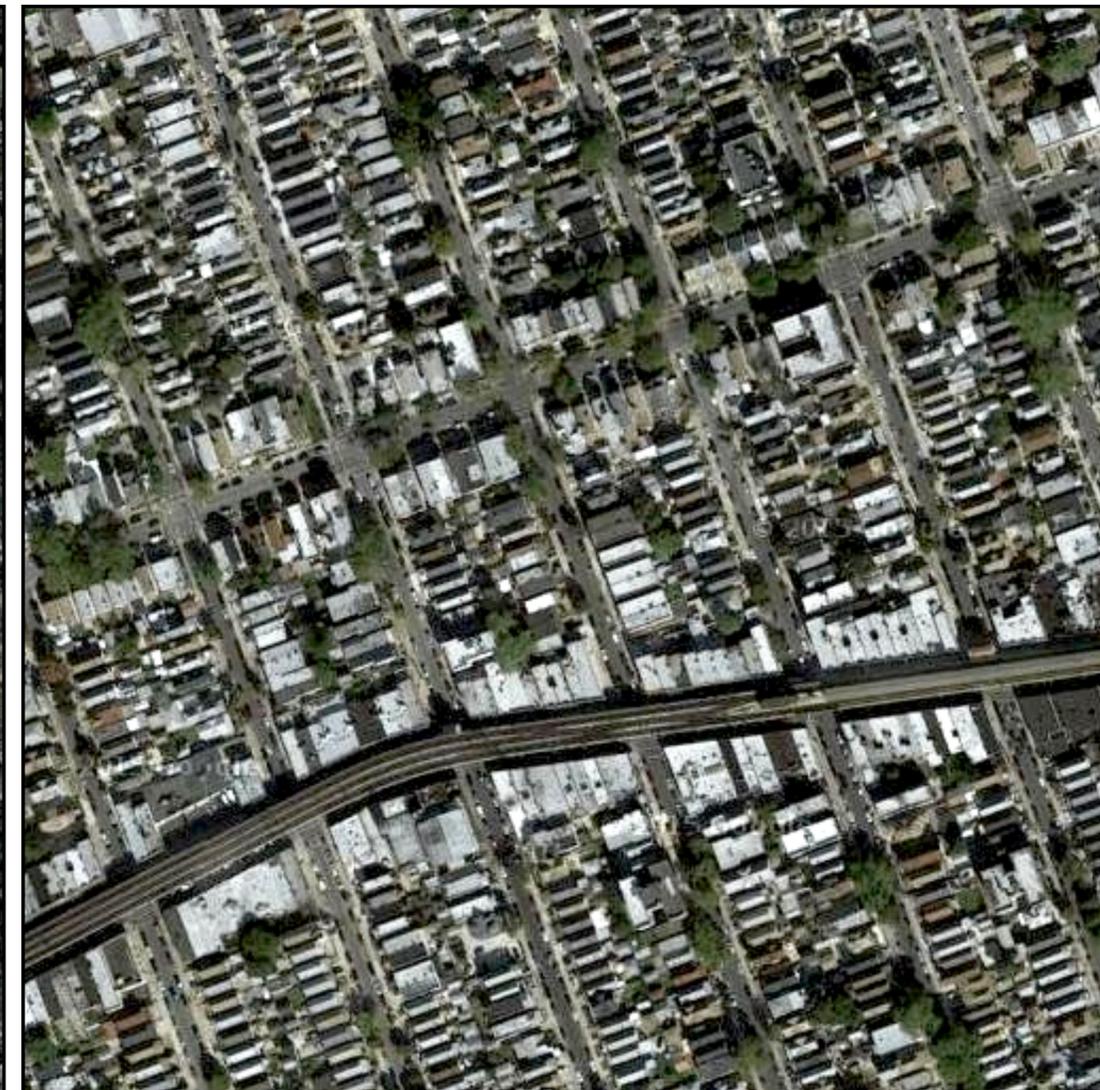
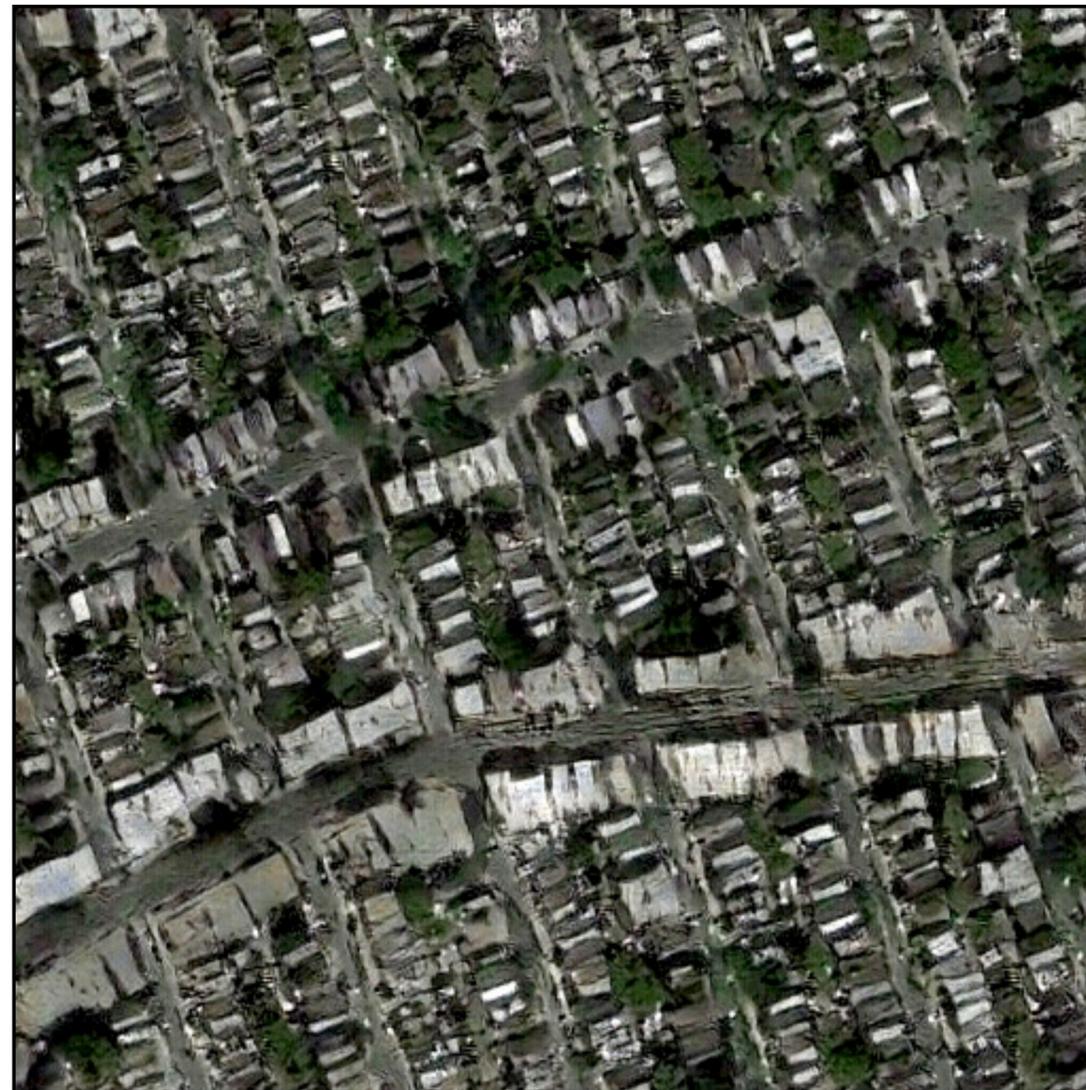


Helps stabilize training + faster convergence

Input

Output

Groundtruth



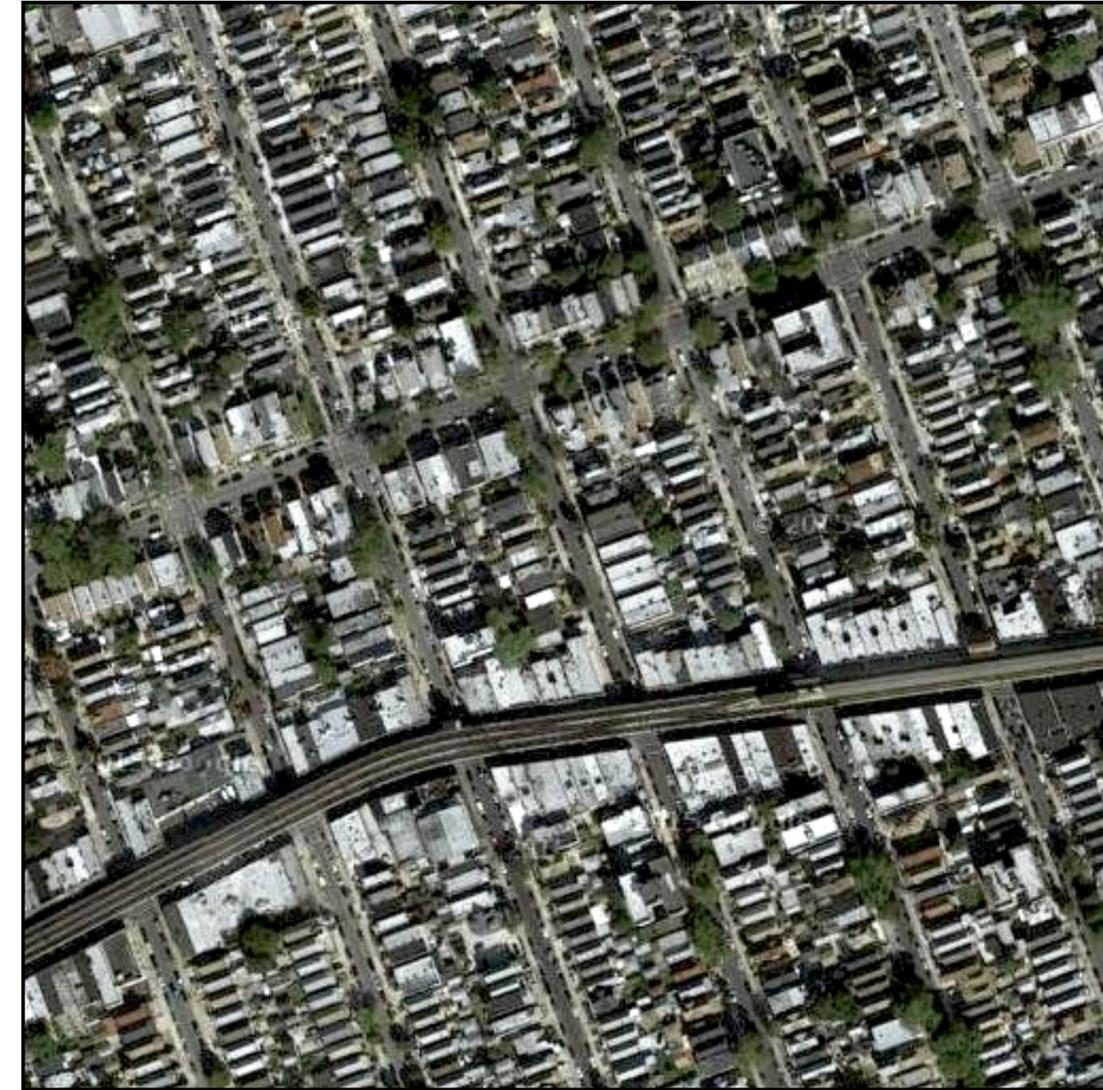
Data from  
[[maps.google.com](https://maps.google.com)]



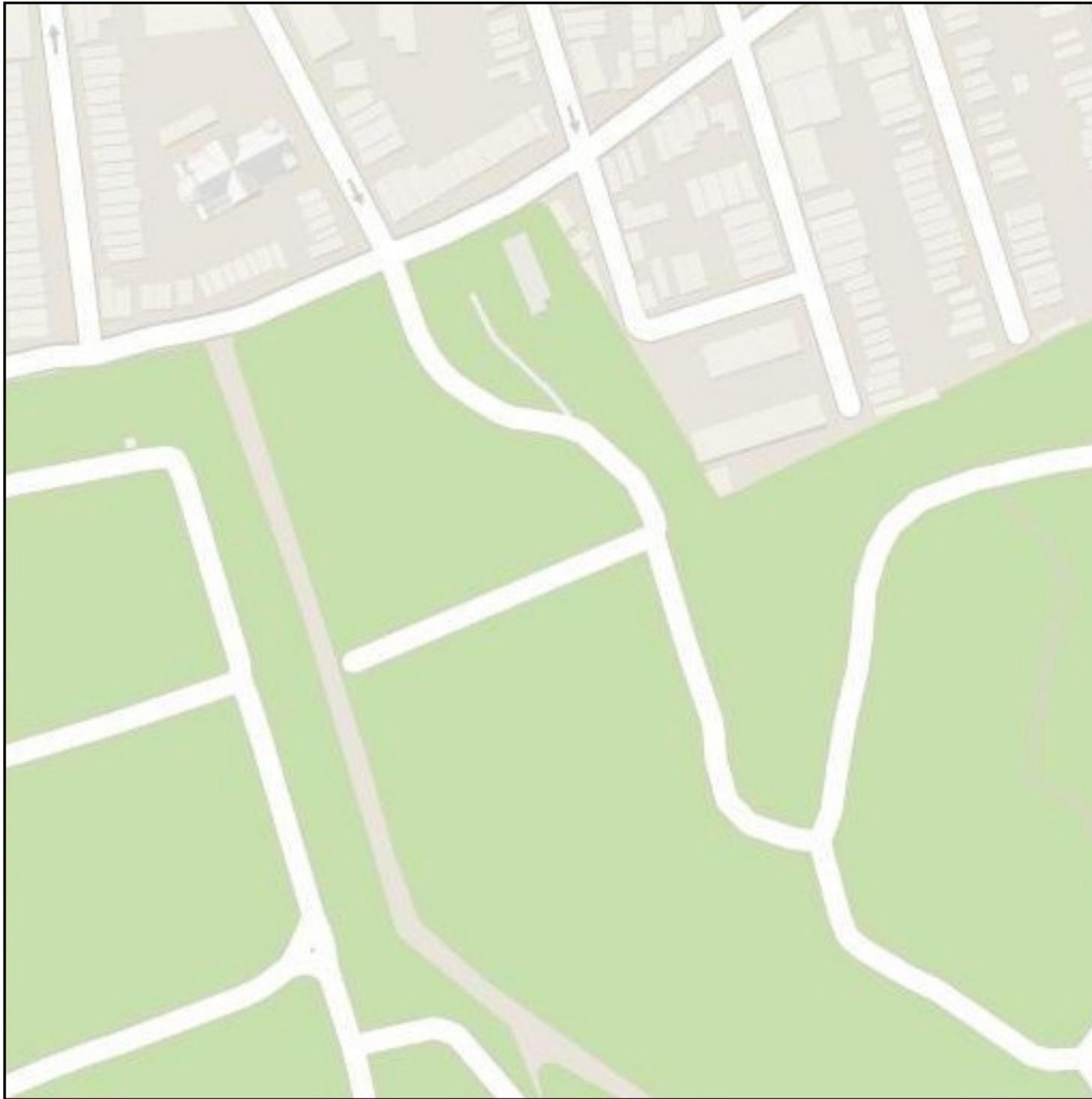
Input

Output

Groundtruth



Input

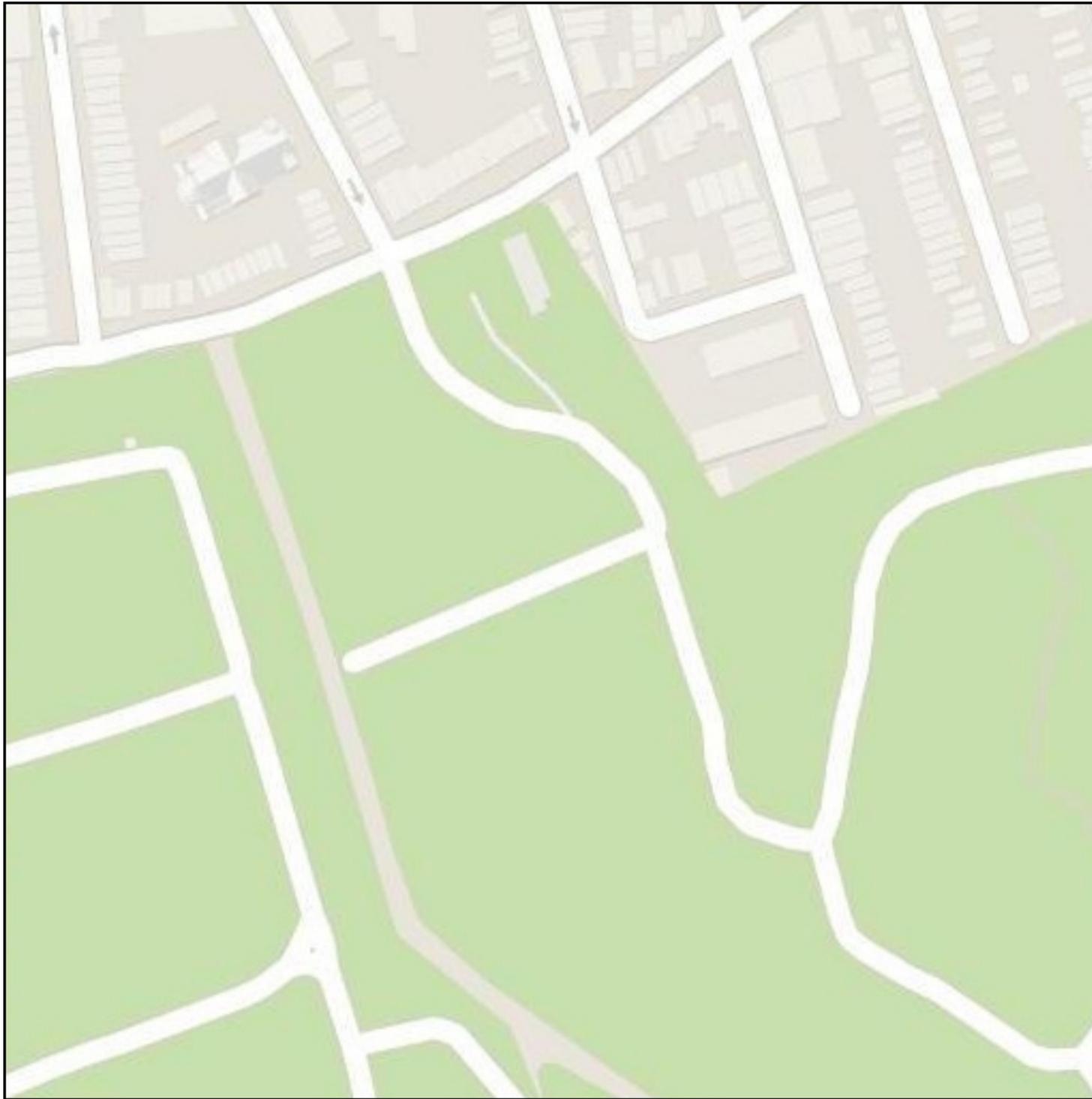


L1 loss only



Input

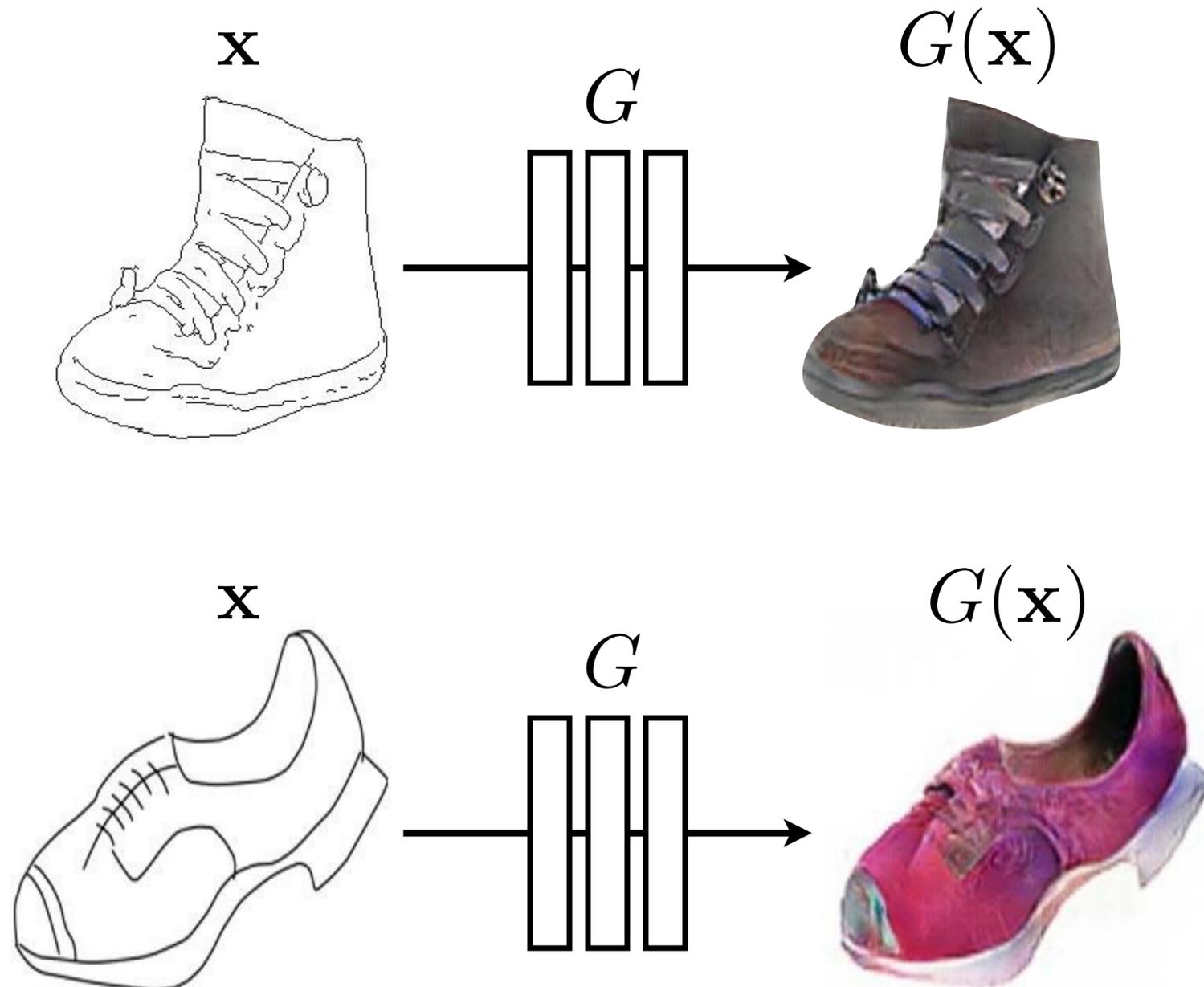
L1 loss + discriminator



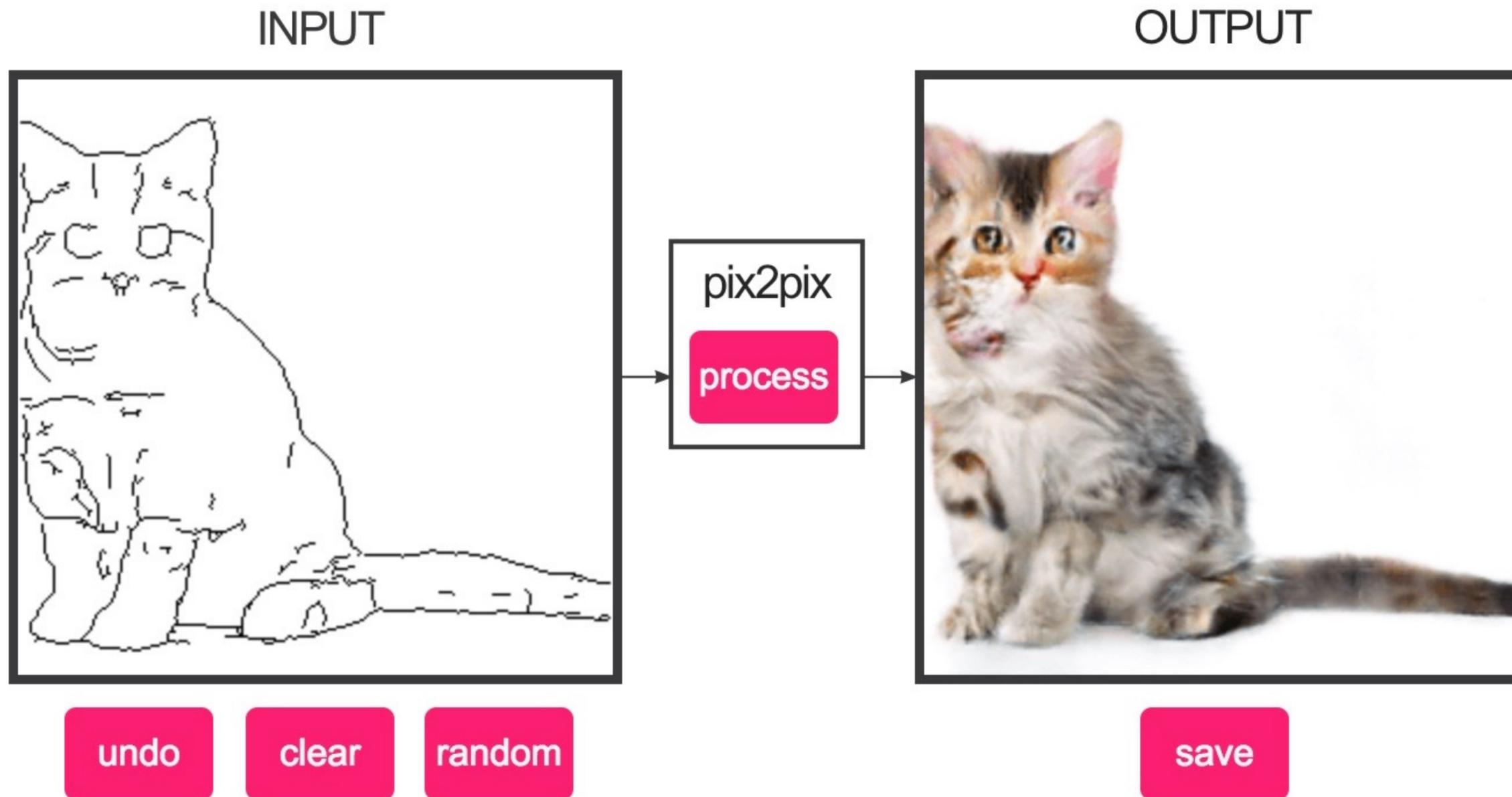
# Training data



[HED, Xie & Tu, 2015]

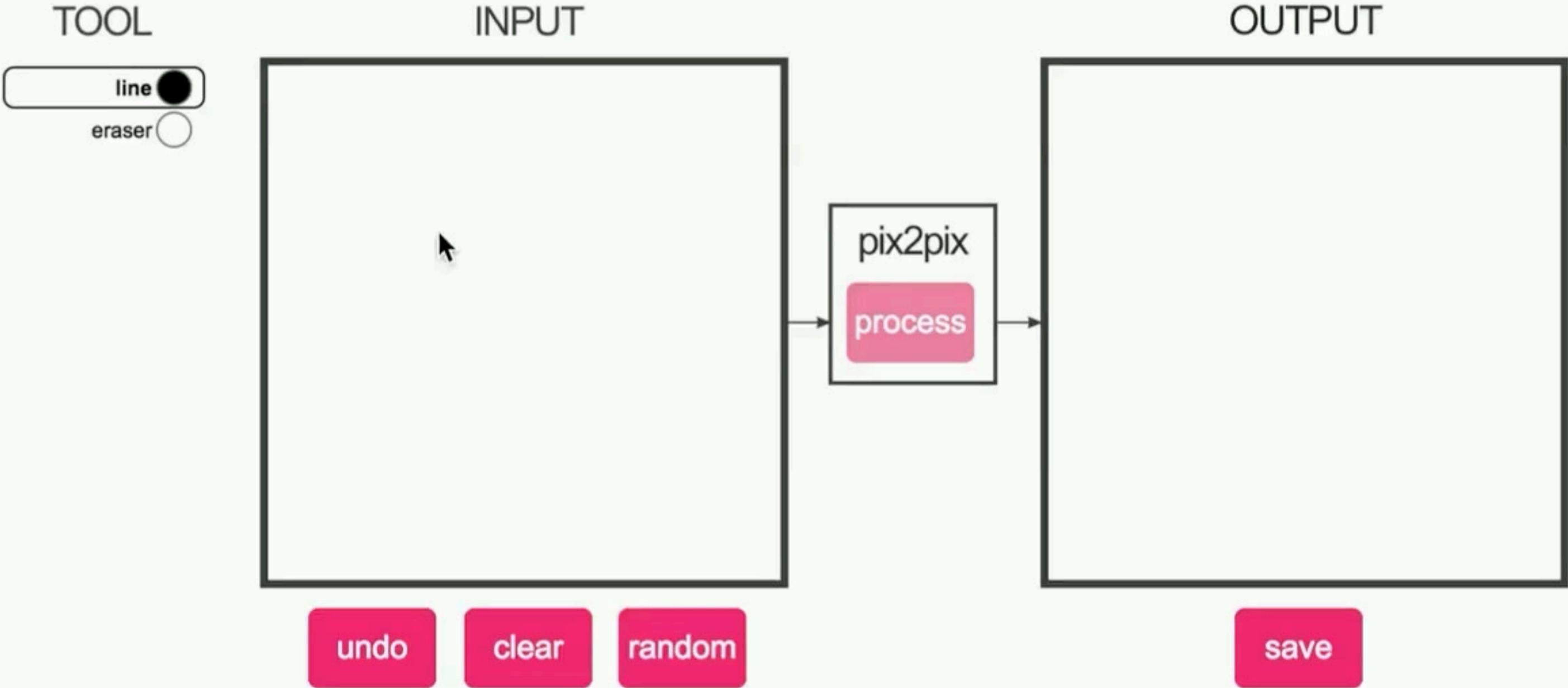


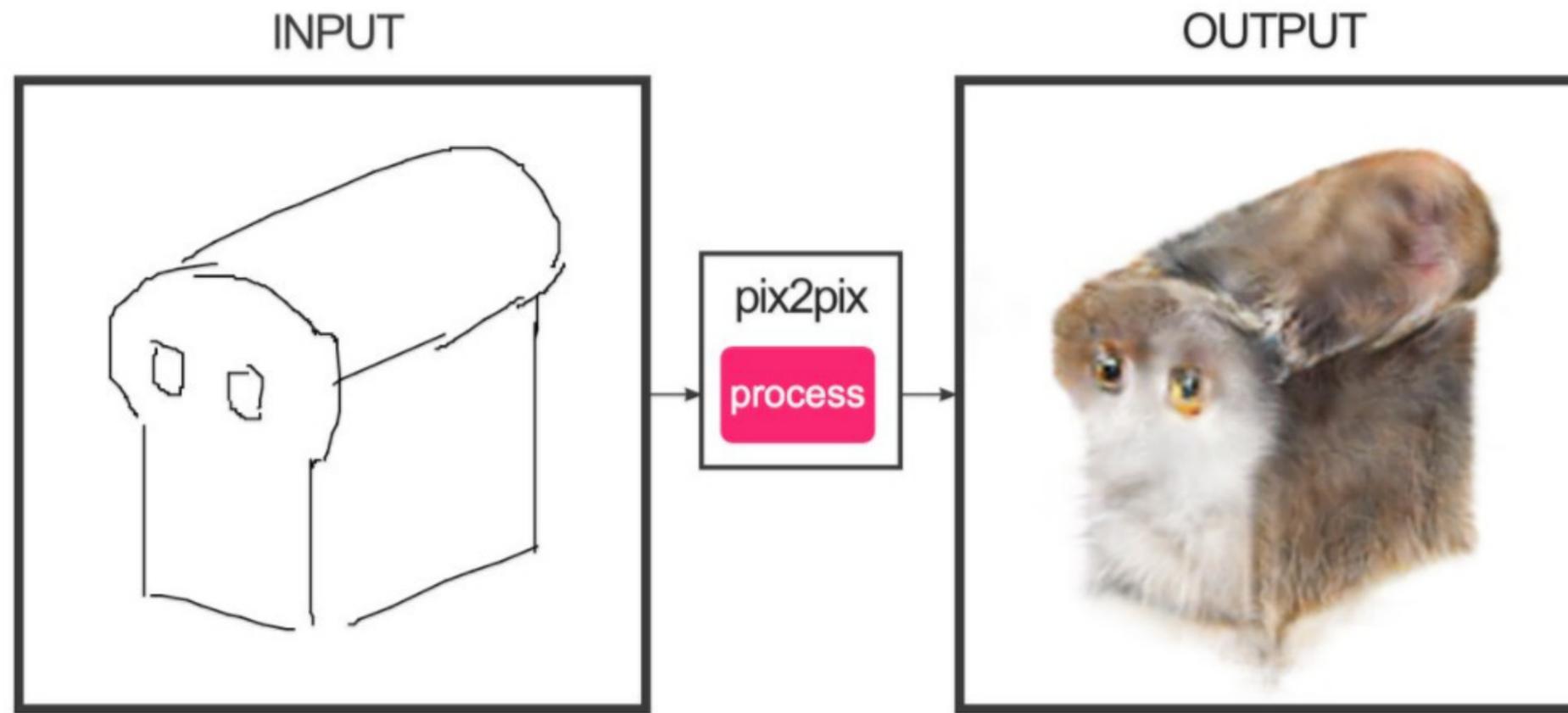
# edges2cats



[Chris Hess, [edges2cats](#)]

# edges2cats





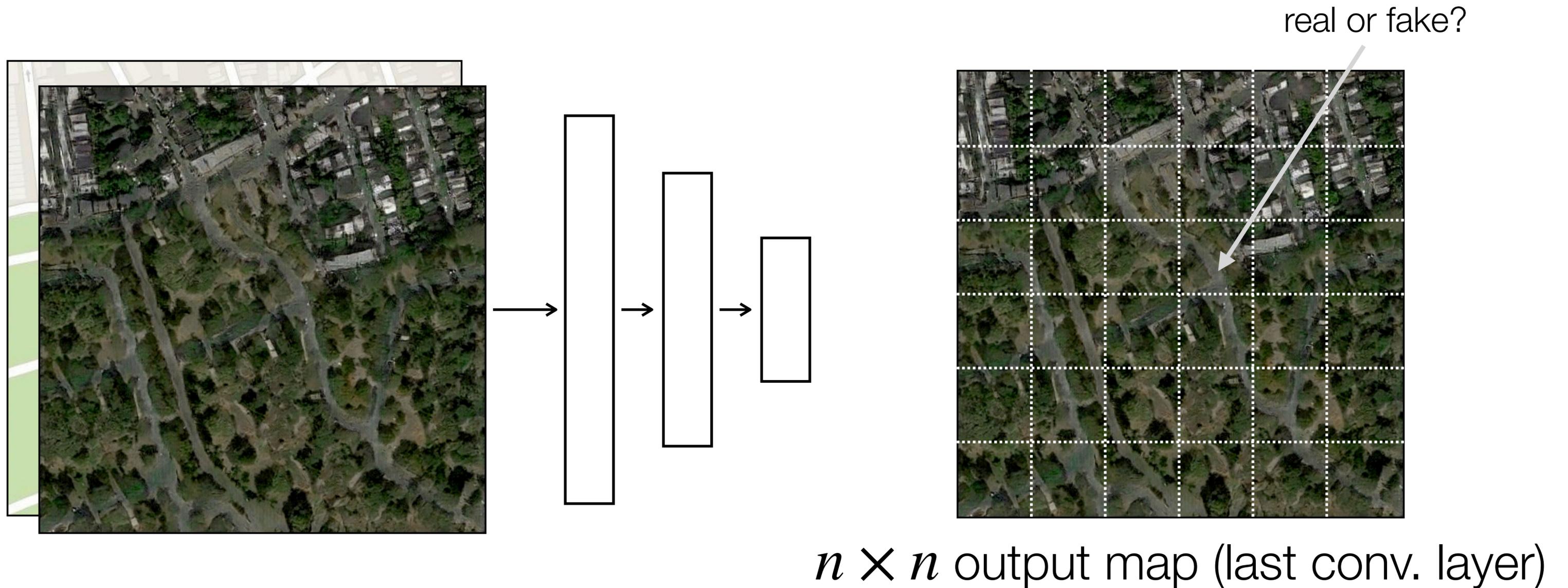
Ivy Tasi @ivymyt



Vitaly Vidmirov @vvid

# Architectures

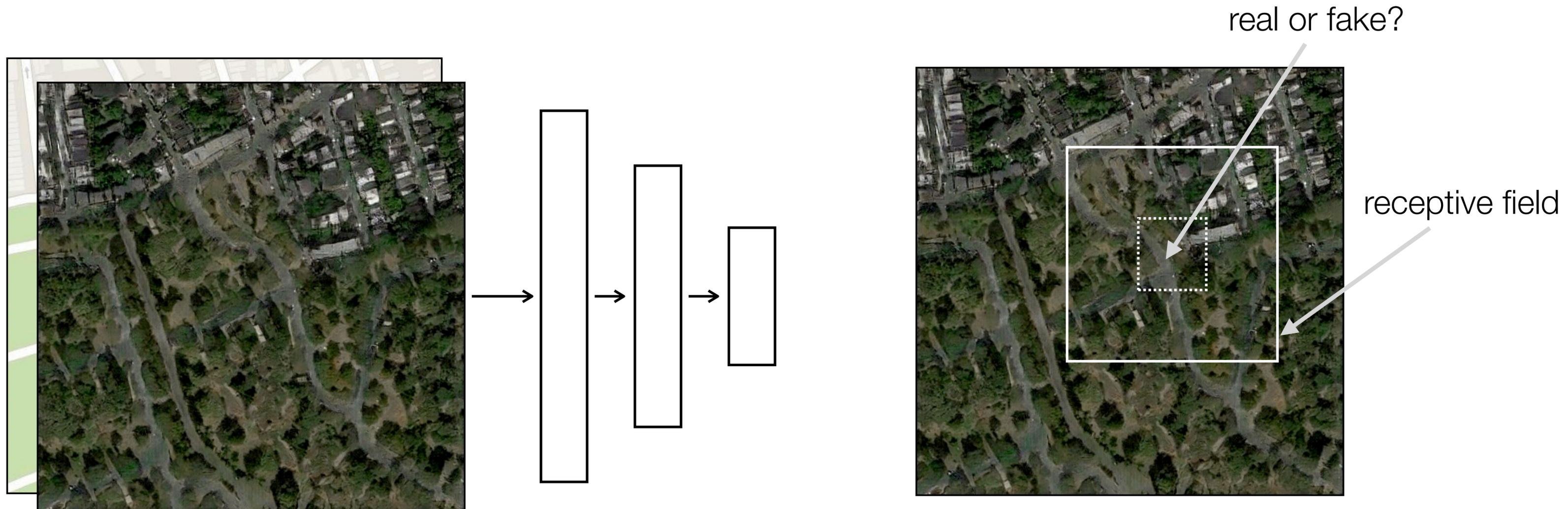
Discriminator: fully convolutional network



Sequence of strided convolutions

# Architectures

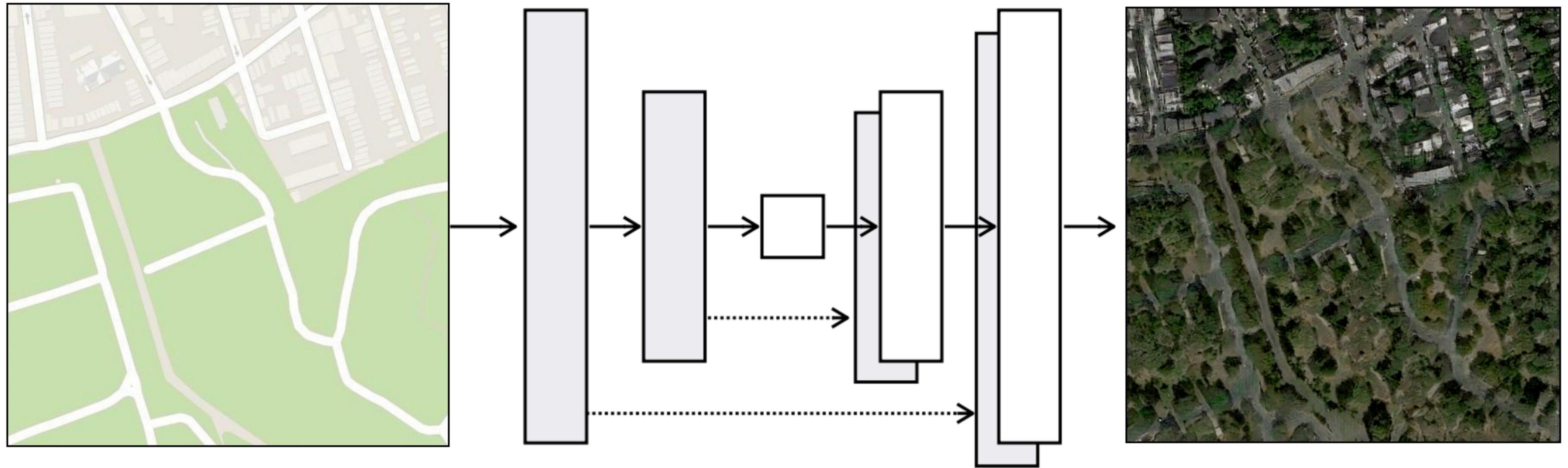
Discriminator: fully convolutional network



Also known as a **Patch GAN**, since it effectively only looks at patches

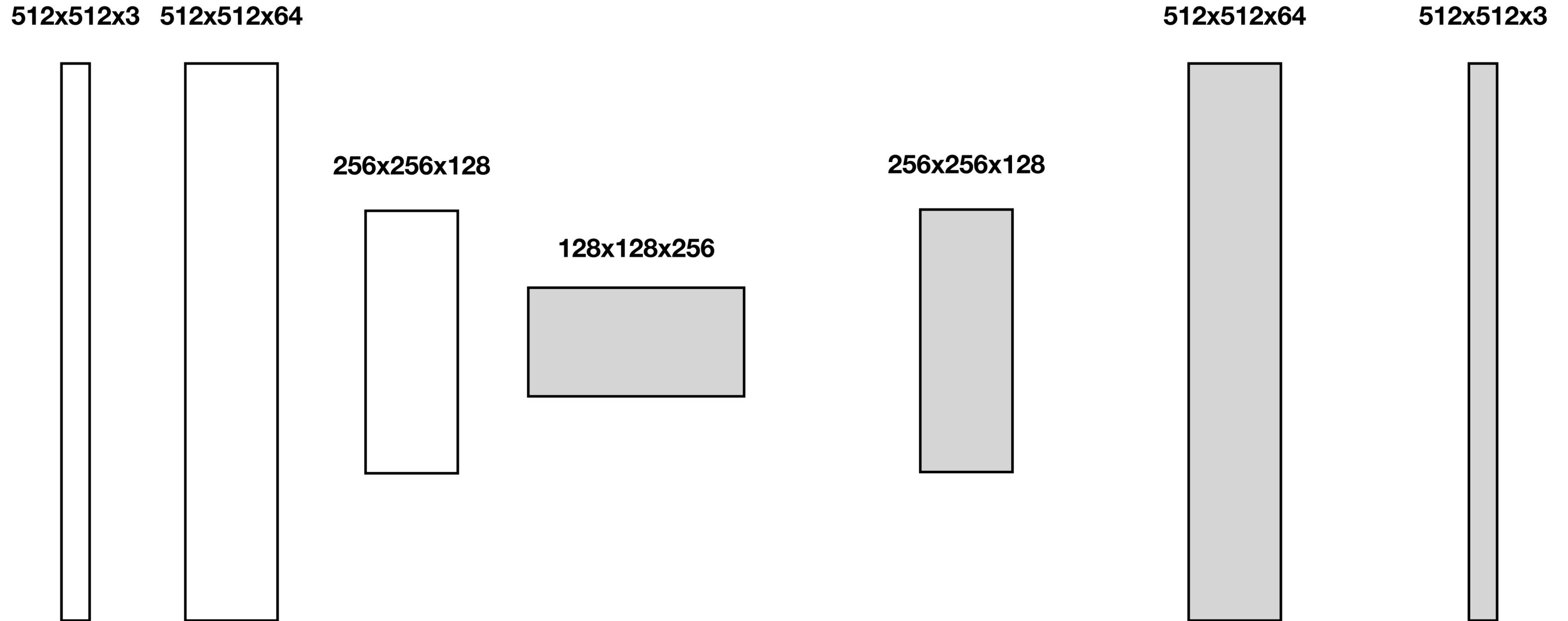
# Architectures

Generator: U-Net

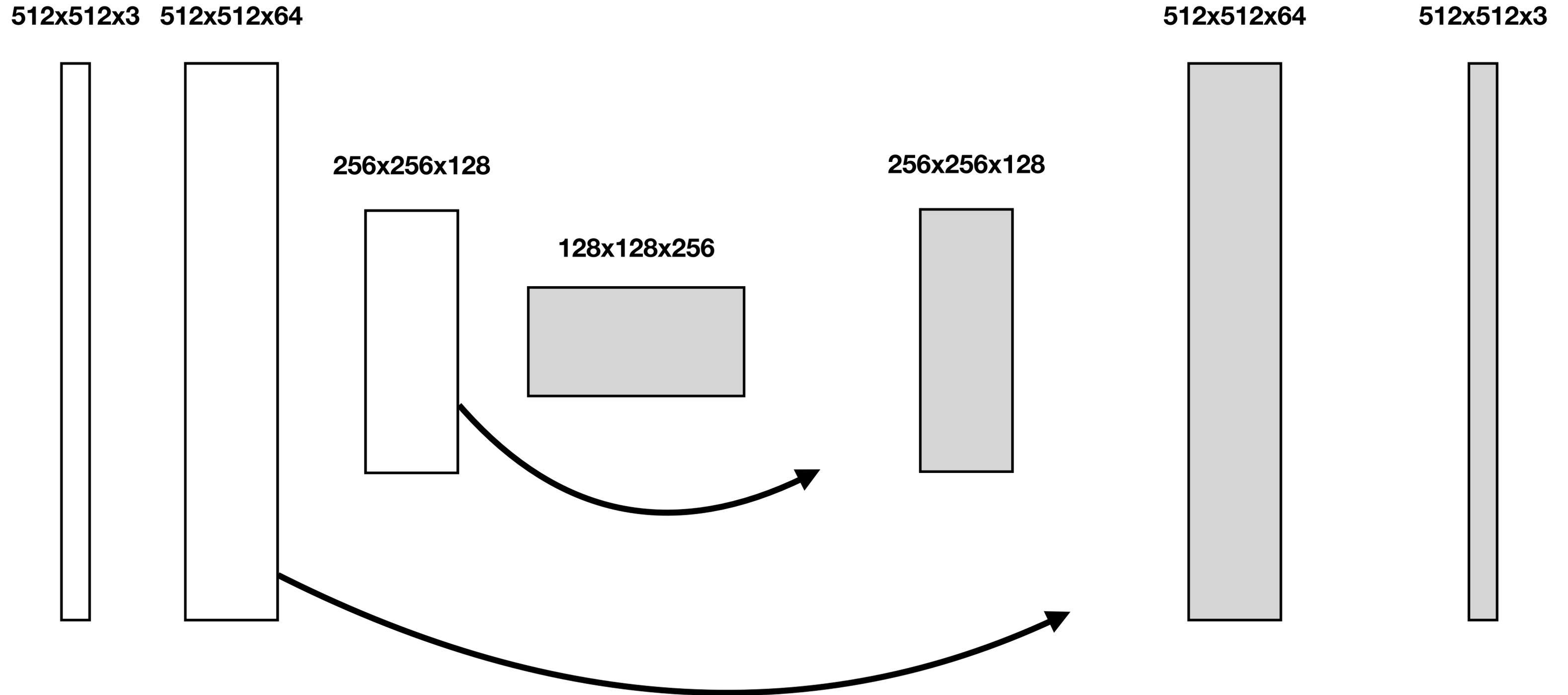


Skip connections between encoder and decoder layers

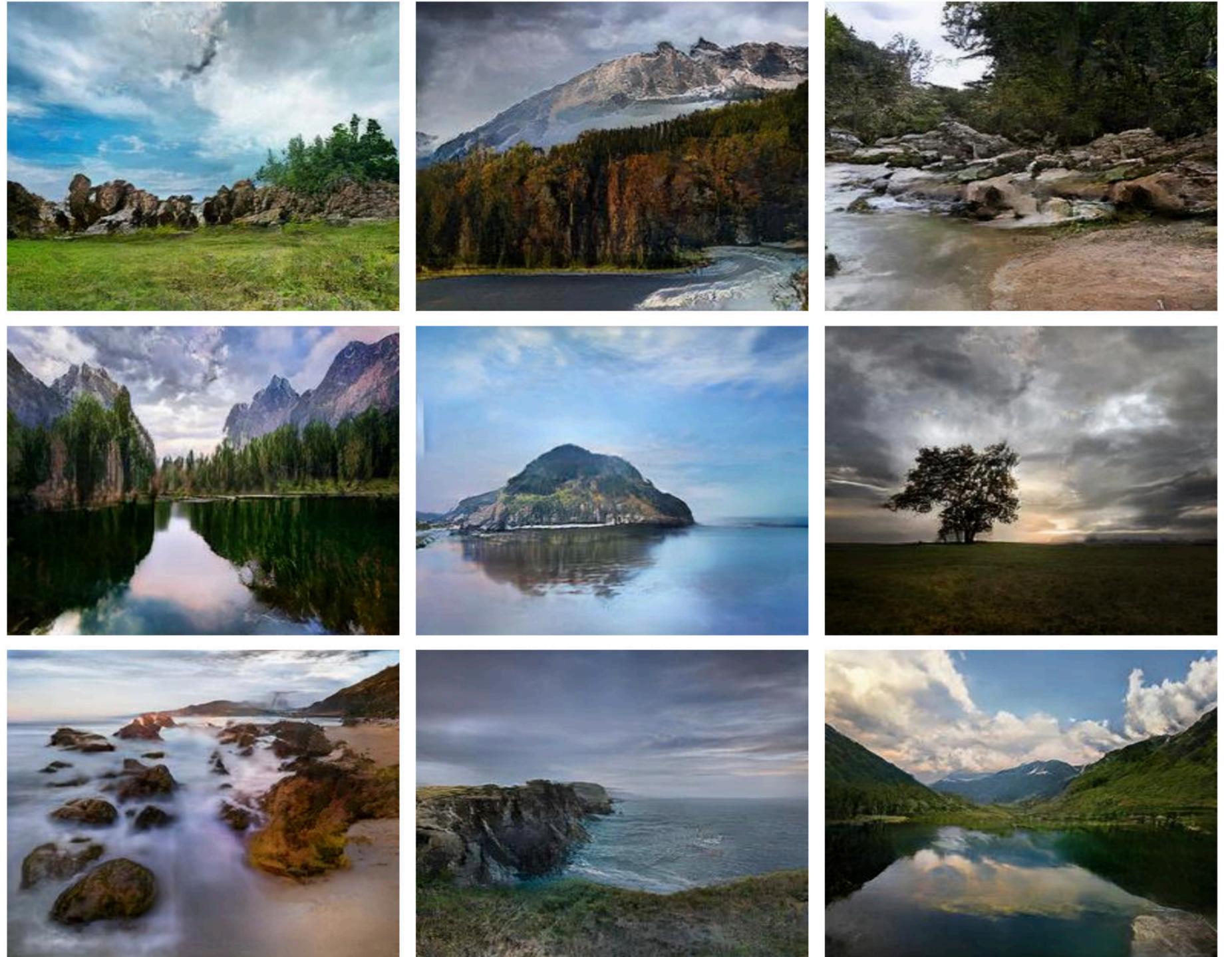
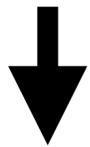
# U-Net

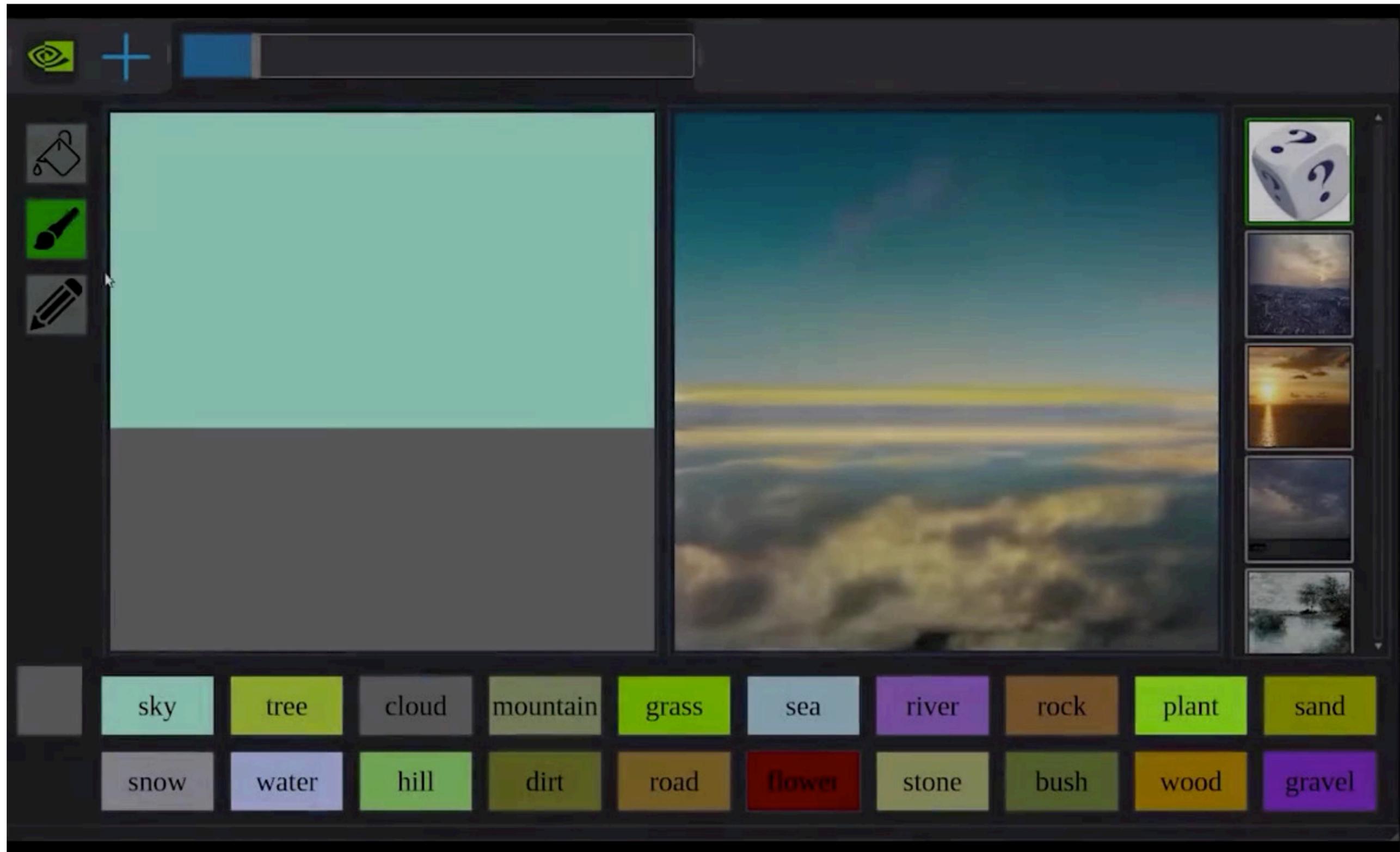


# U-Net



# More recent architectures



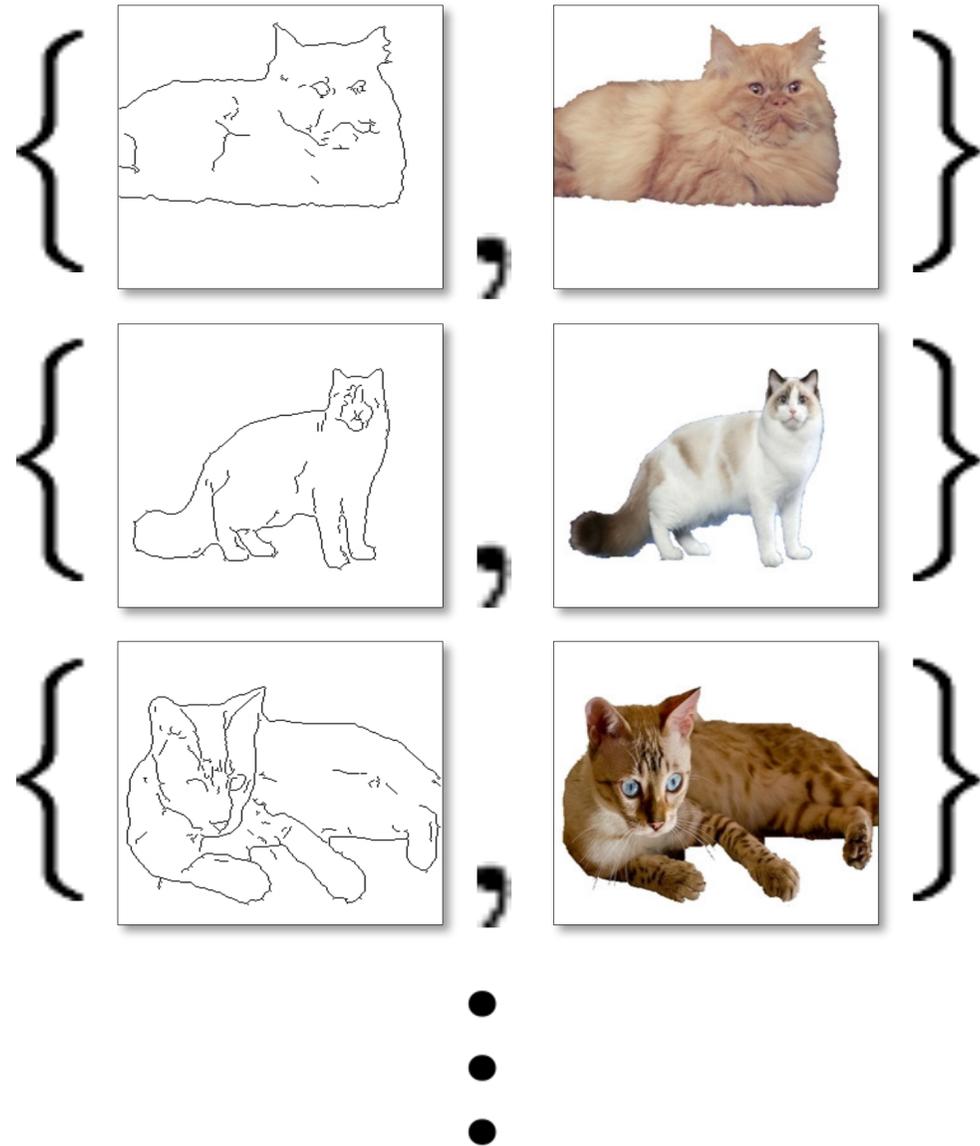


# Handling unpaired data

Paired data

$x_i$

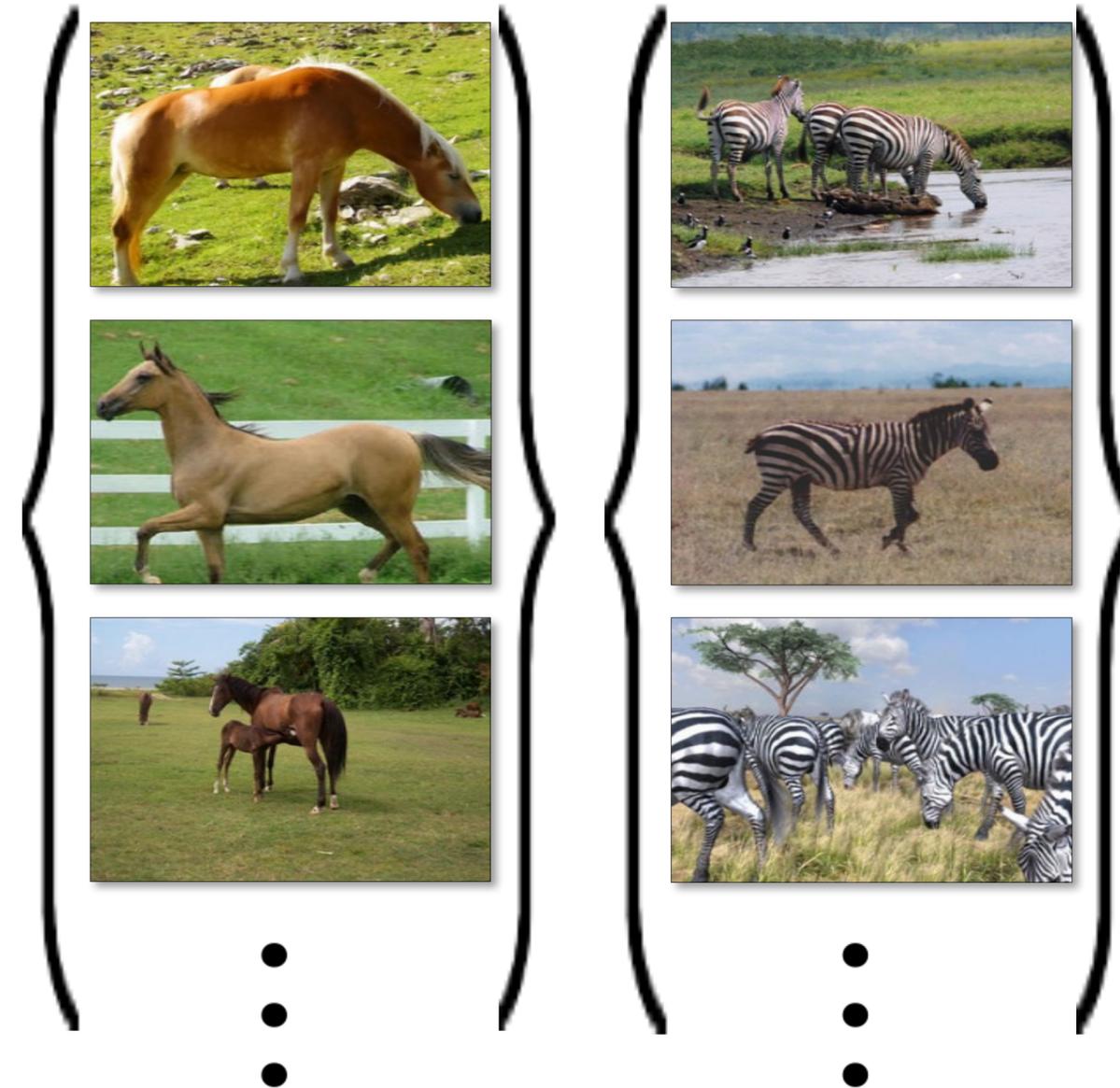
$y_i$

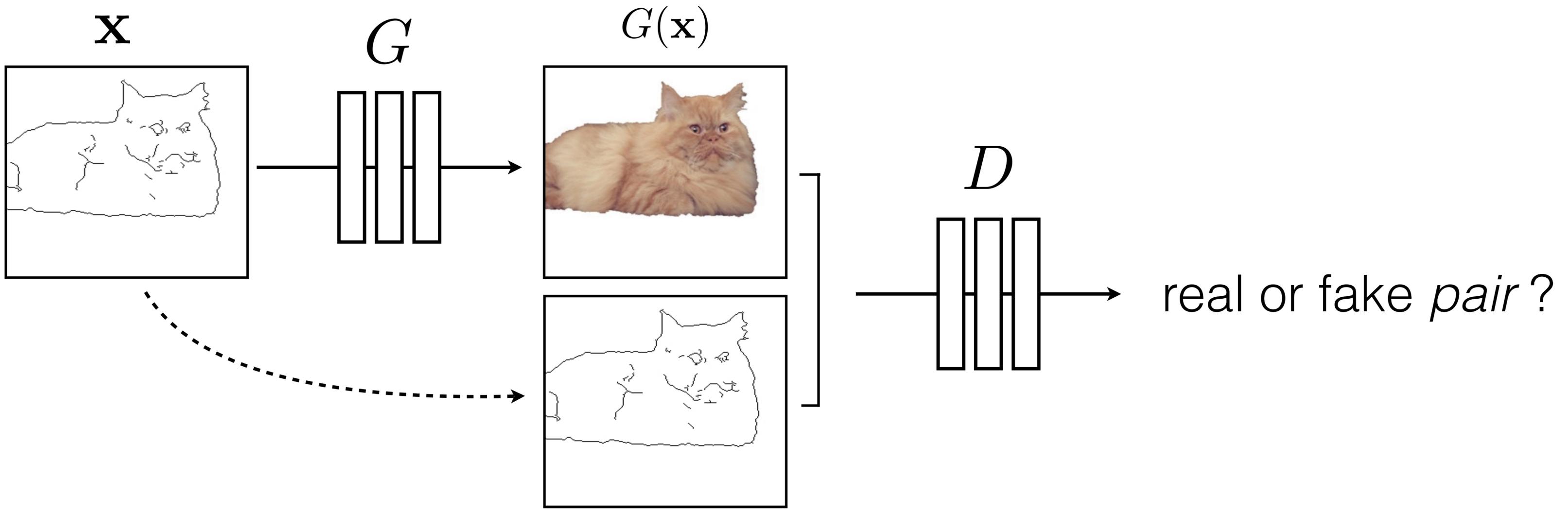


Unpaired data

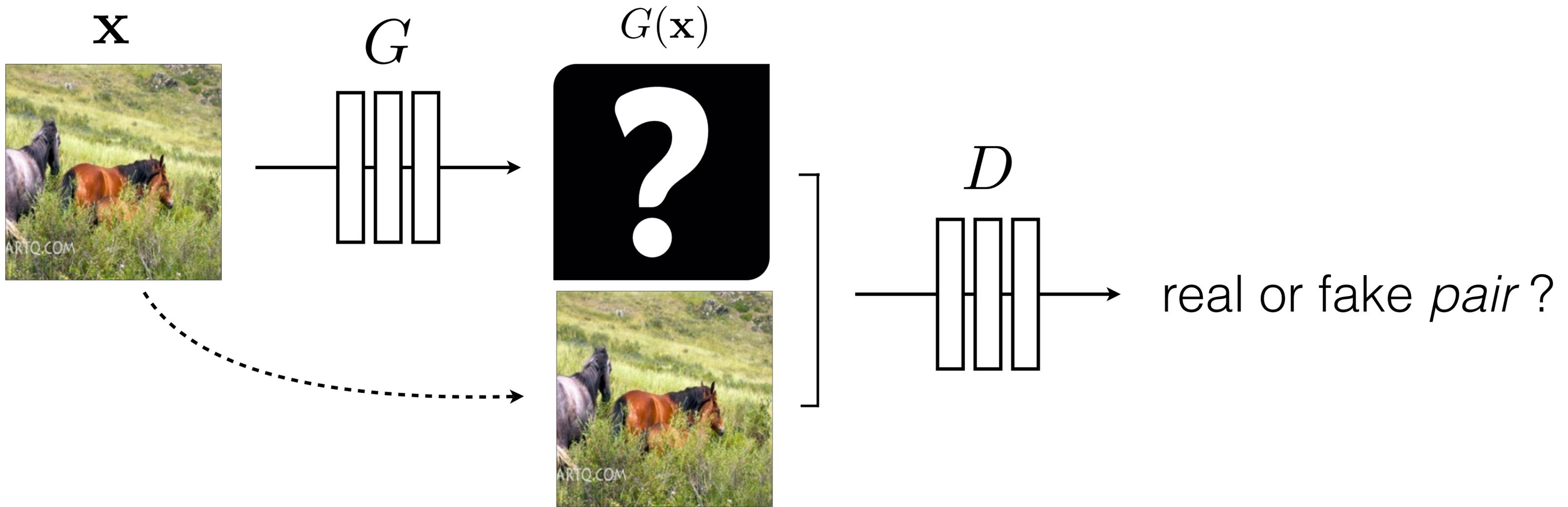
$X$

$Y$





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



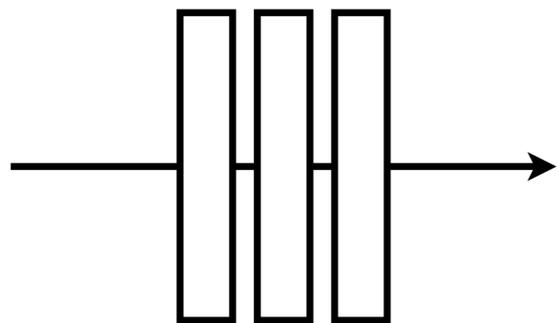
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$

No input-output pairs!

$\mathbf{x}$



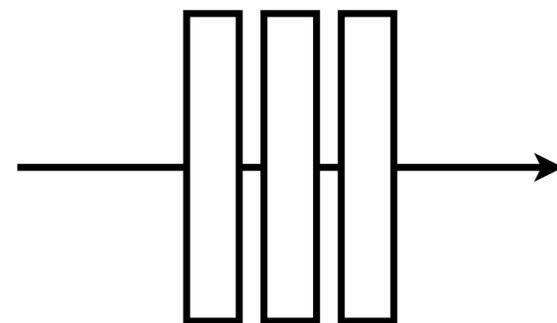
$G$



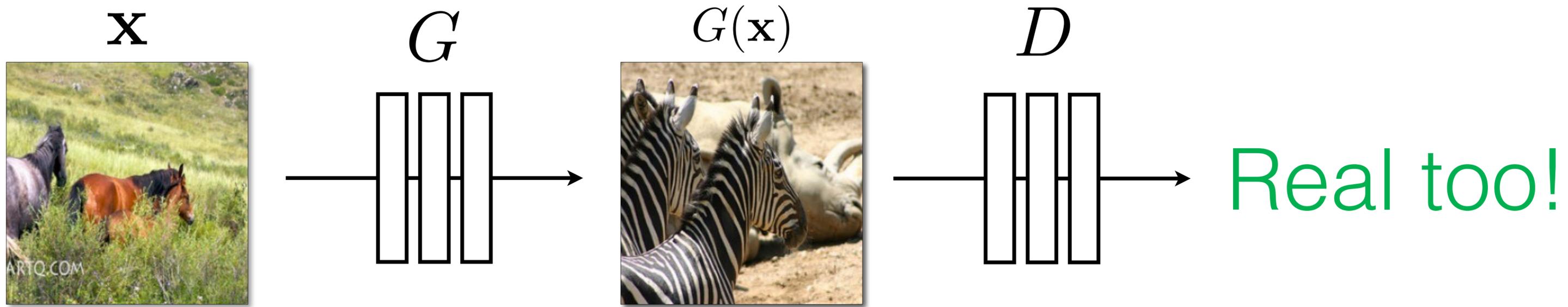
$G(\mathbf{x})$



$D$

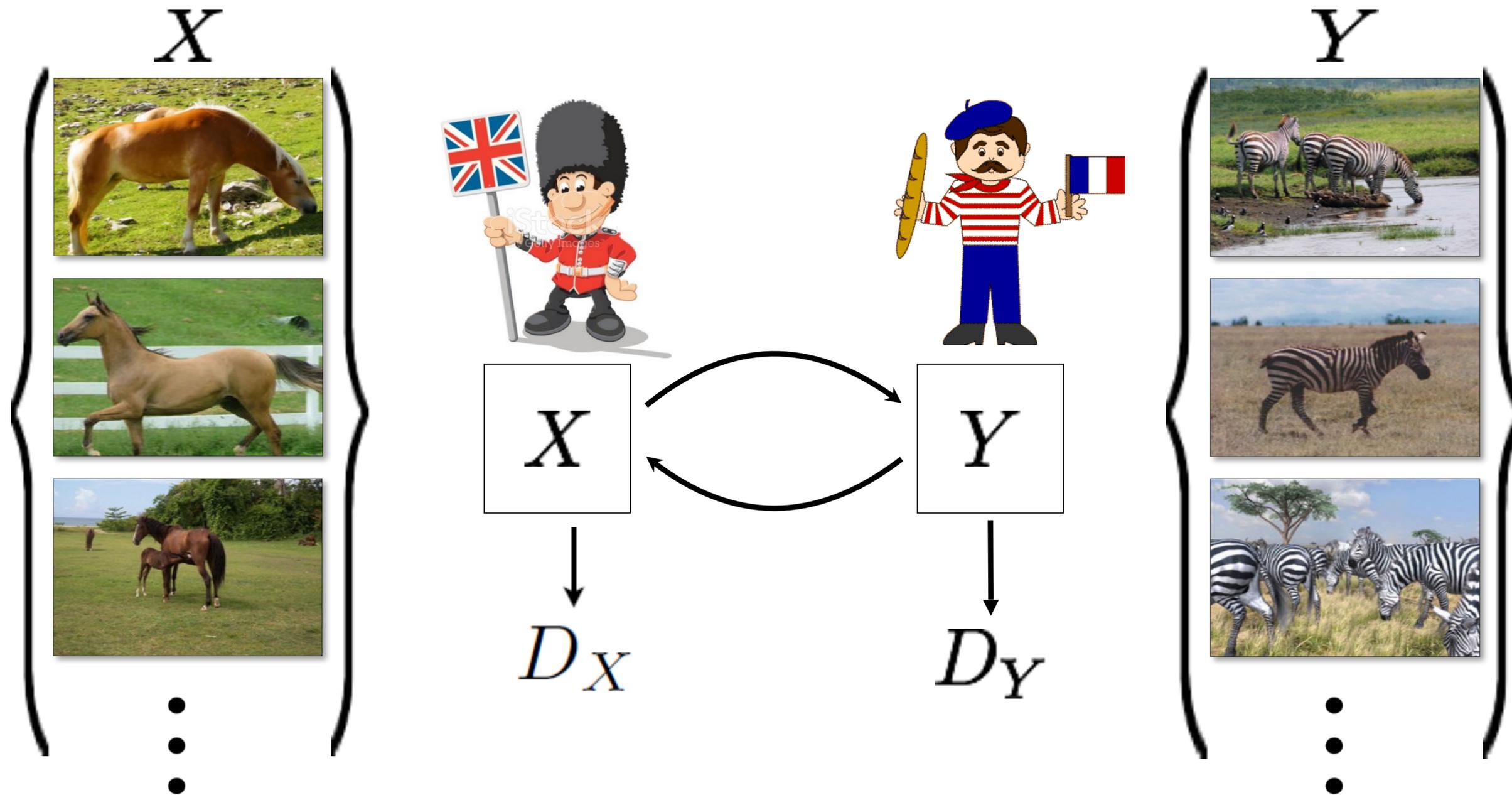


Real!

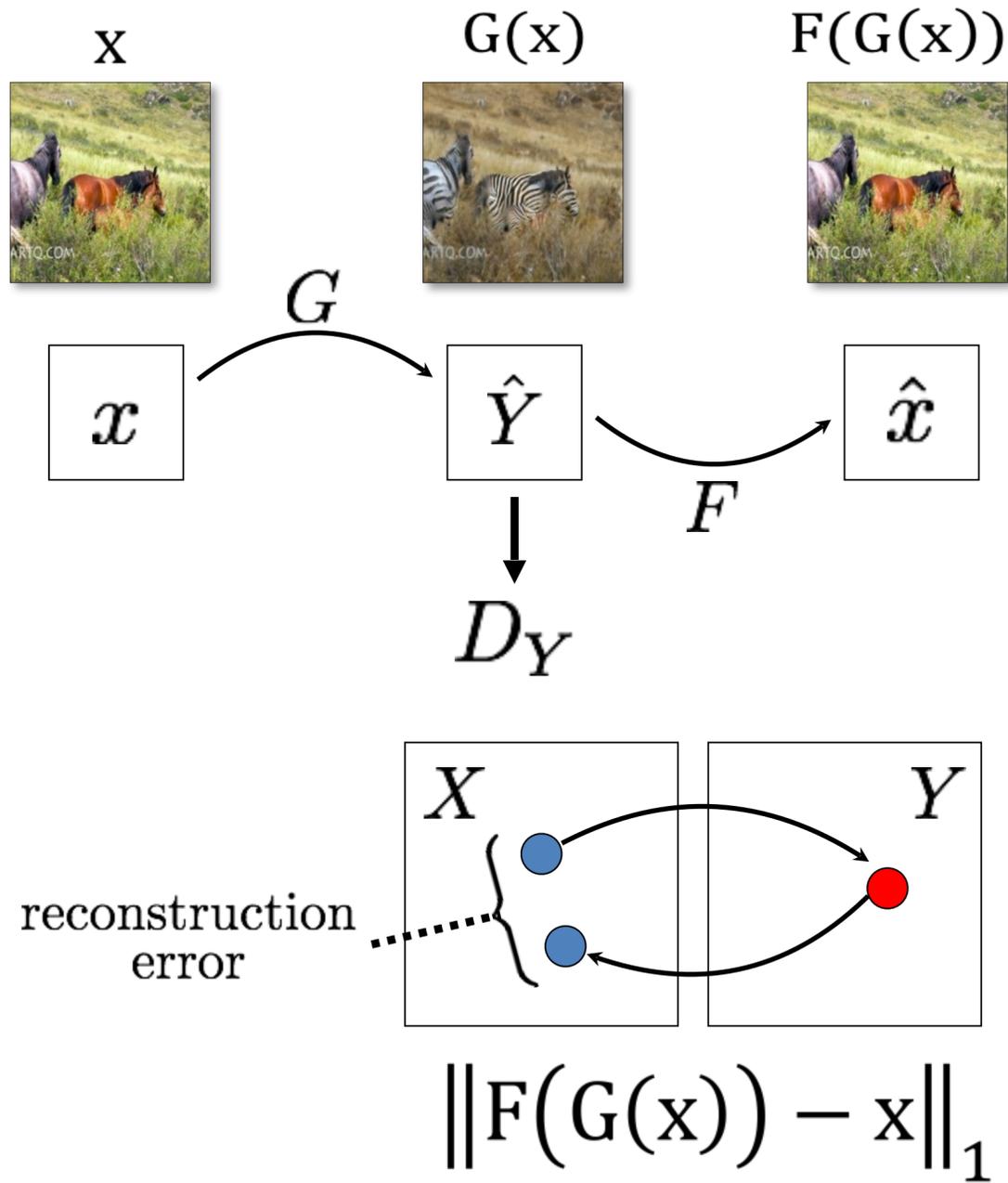


Nothing to force output to correspond to input

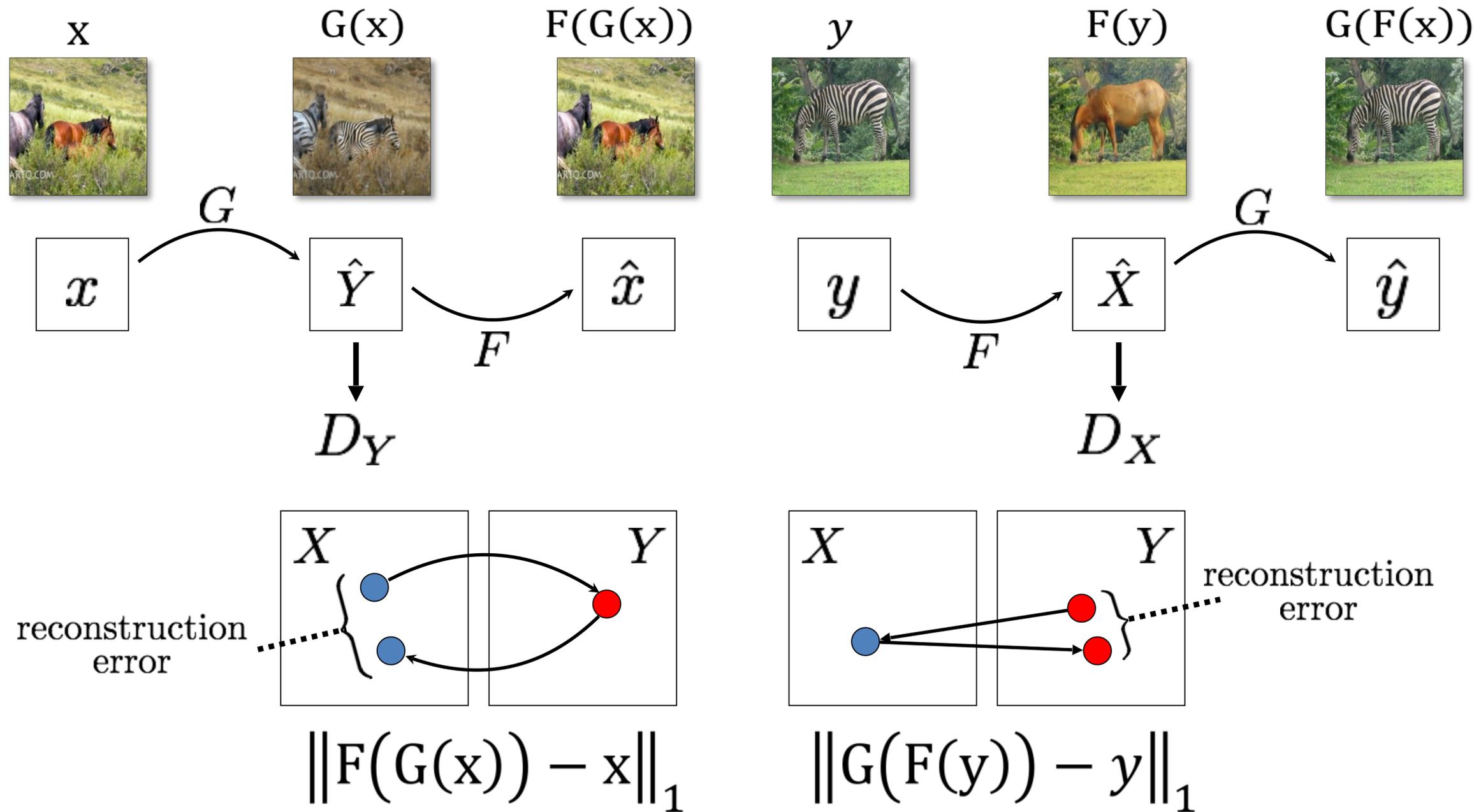
# CycleGAN



# Cycle Consistency Loss



# Cycle Consistency Loss







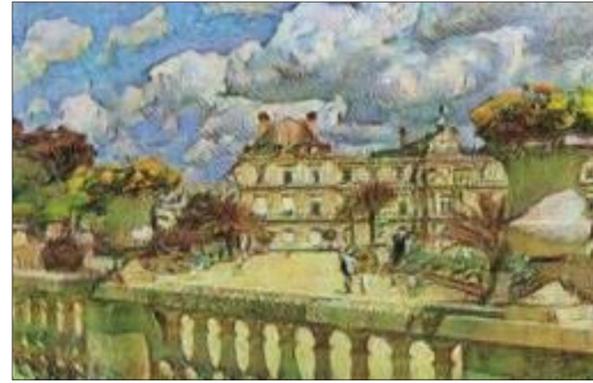
Input



Monet



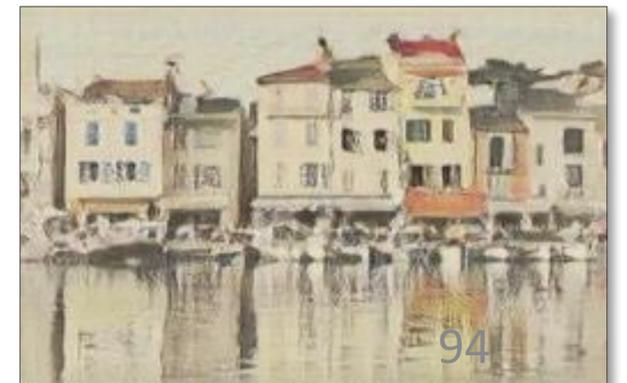
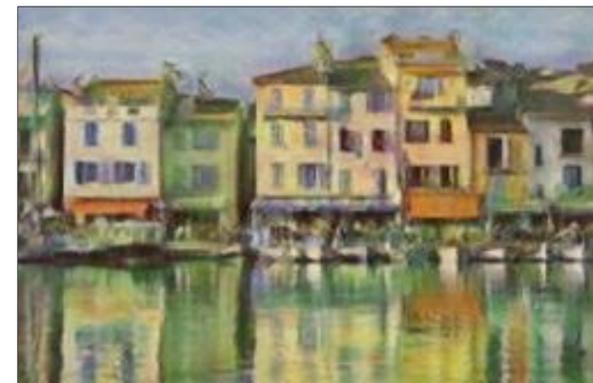
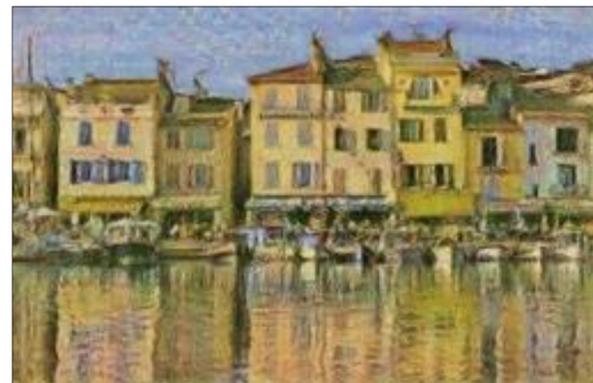
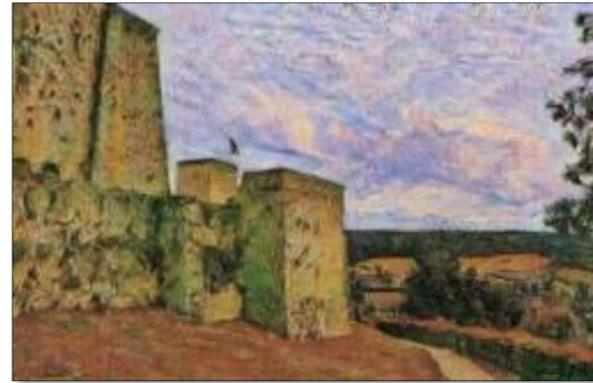
Van Gogh



Cezanne



Ukiyo-e



**Next time:** more image synthesis