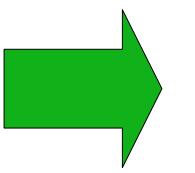


Lecture 13: Image synthesis

Synthesizing textures

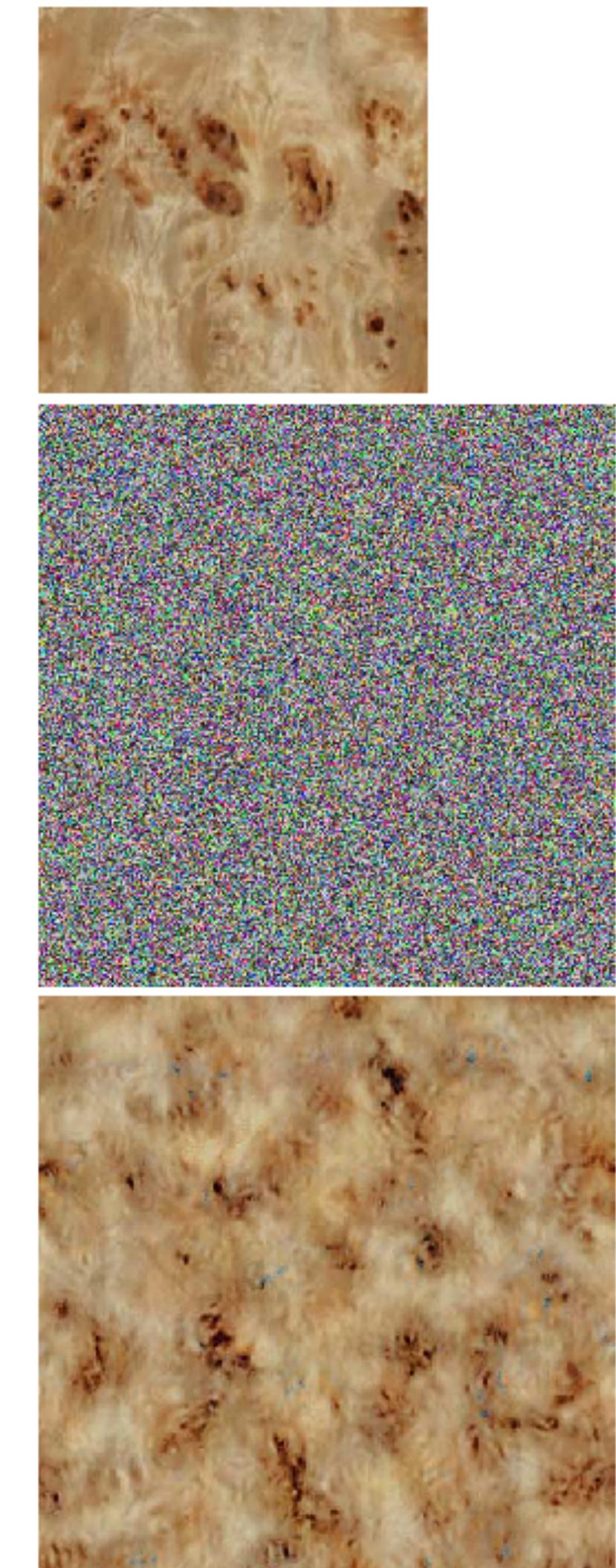


Recall: parametric texture synthesis

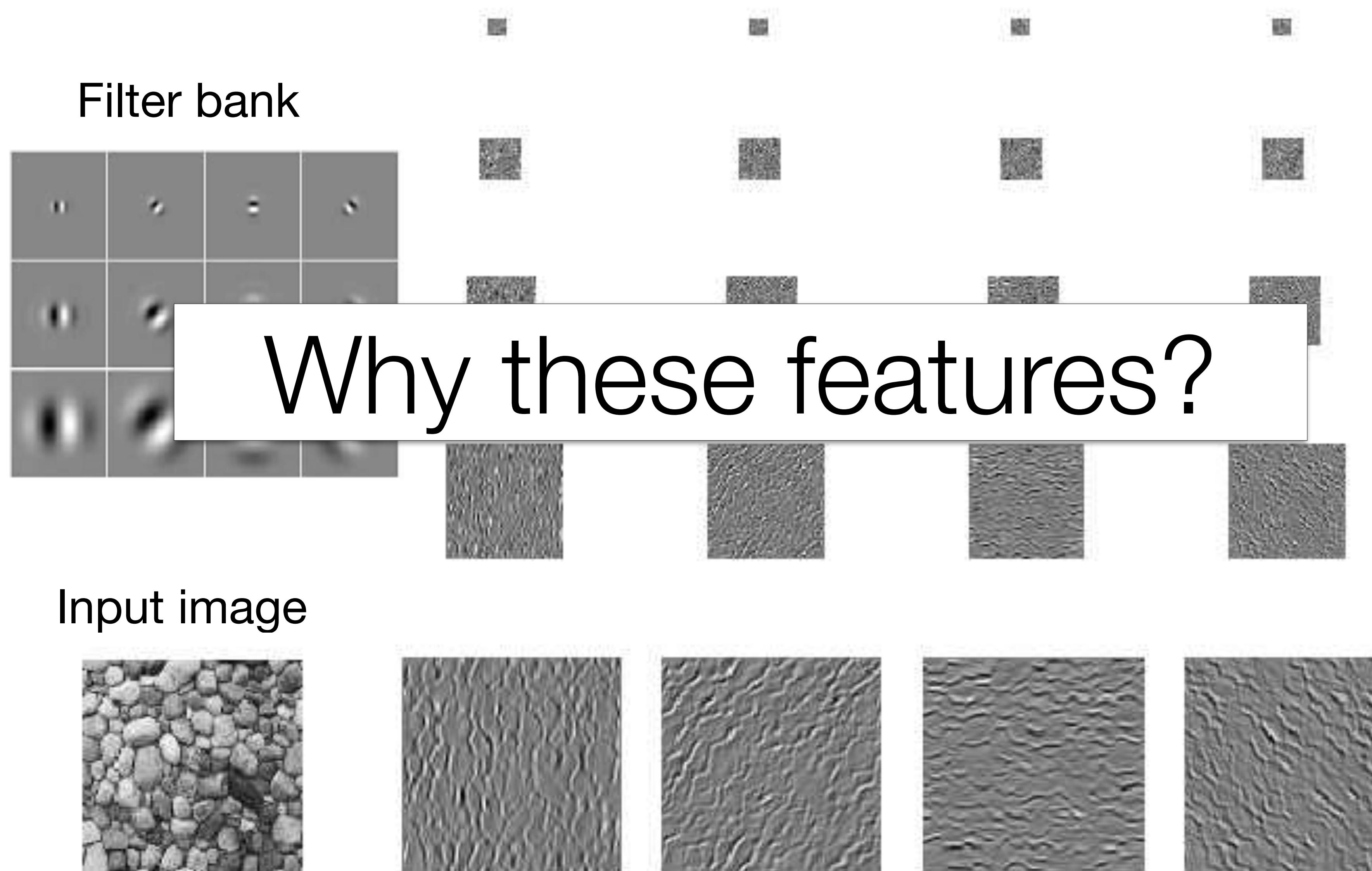
Start with a noise image as output.

Main loop:

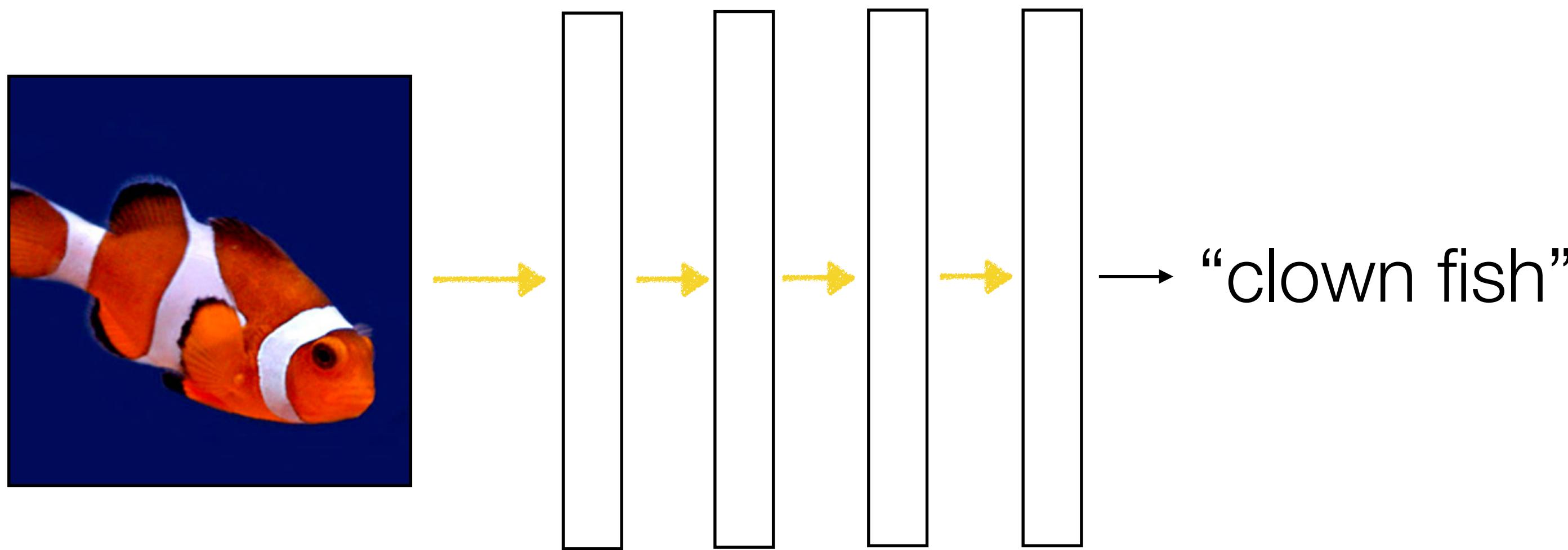
- Match pixel histogram of output image to input
- Decompose input/output images using a Steerable Pyramid
- Match subband histograms of input and output pyramids
- Reconstruct input and output images (collapse the pyramids)



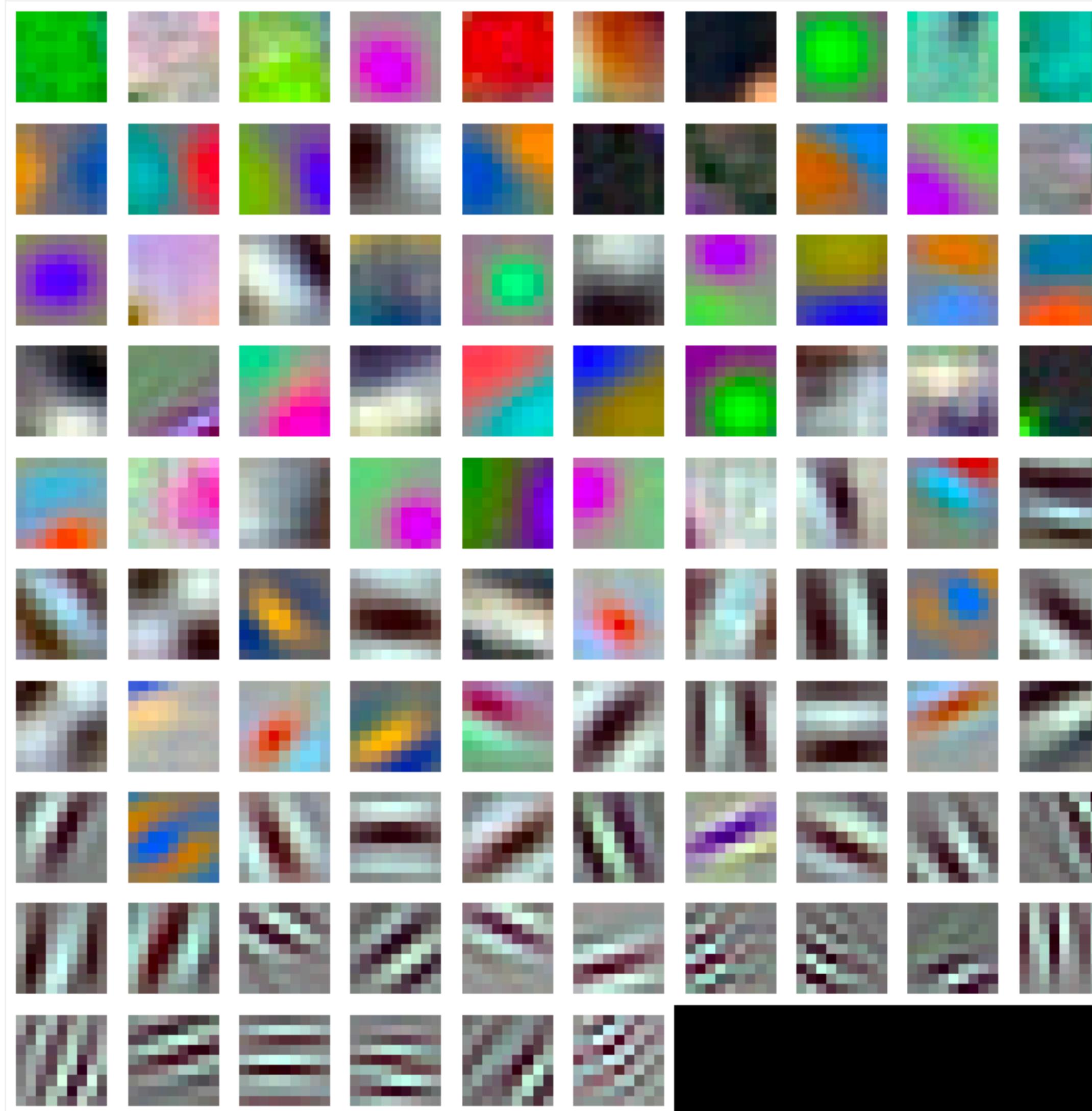
Recall: steerable filter features



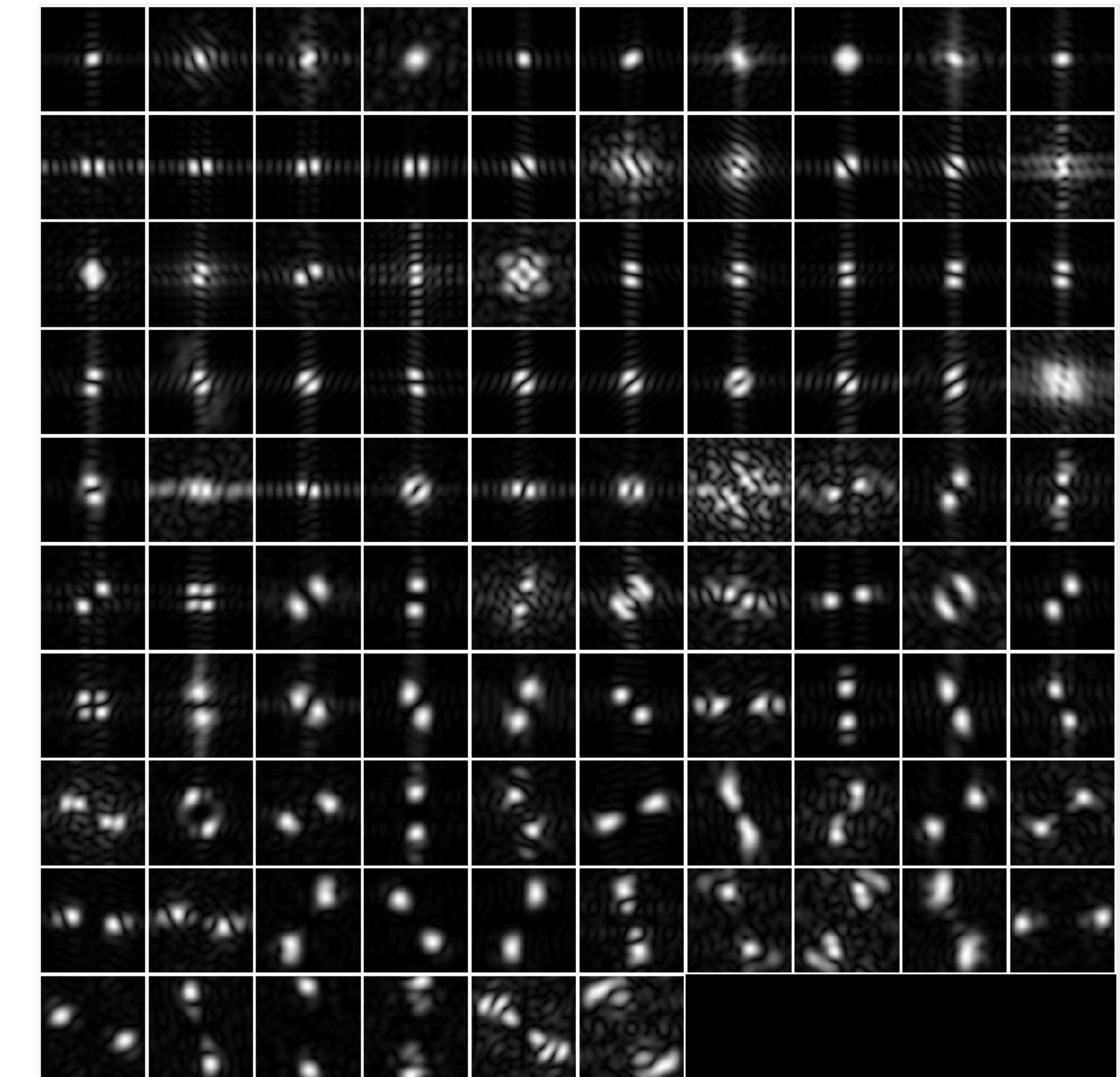
Use CNN features instead!



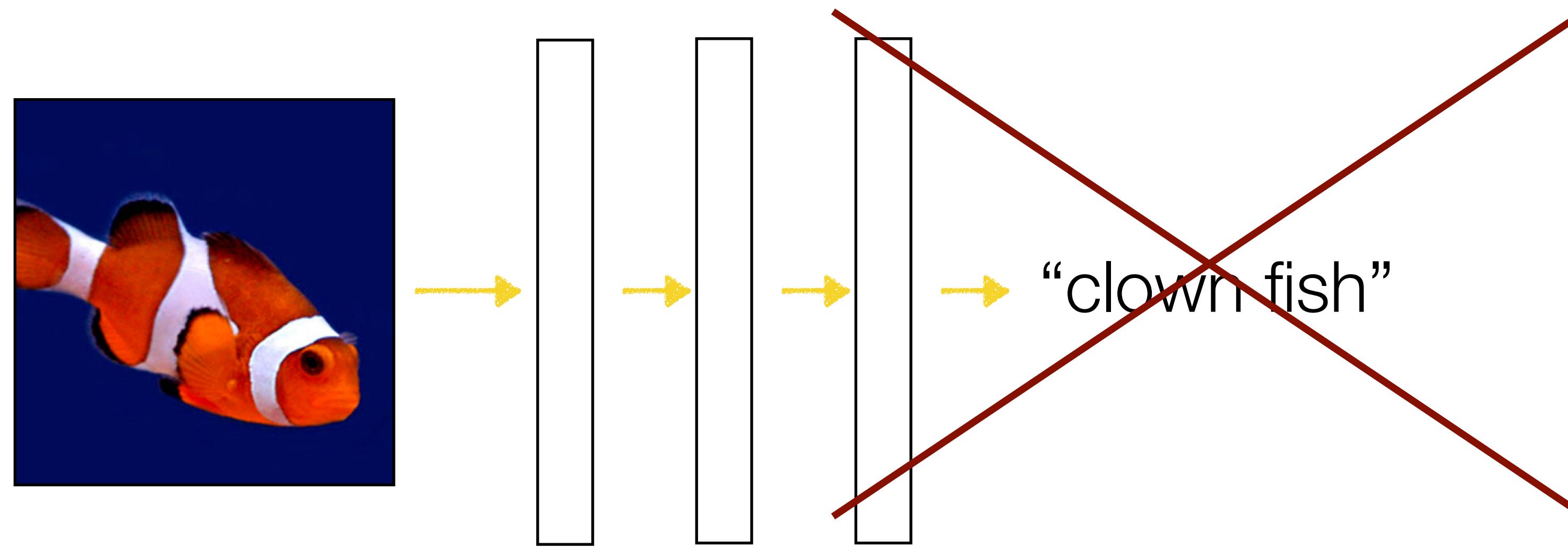
Recall: AlexNet units



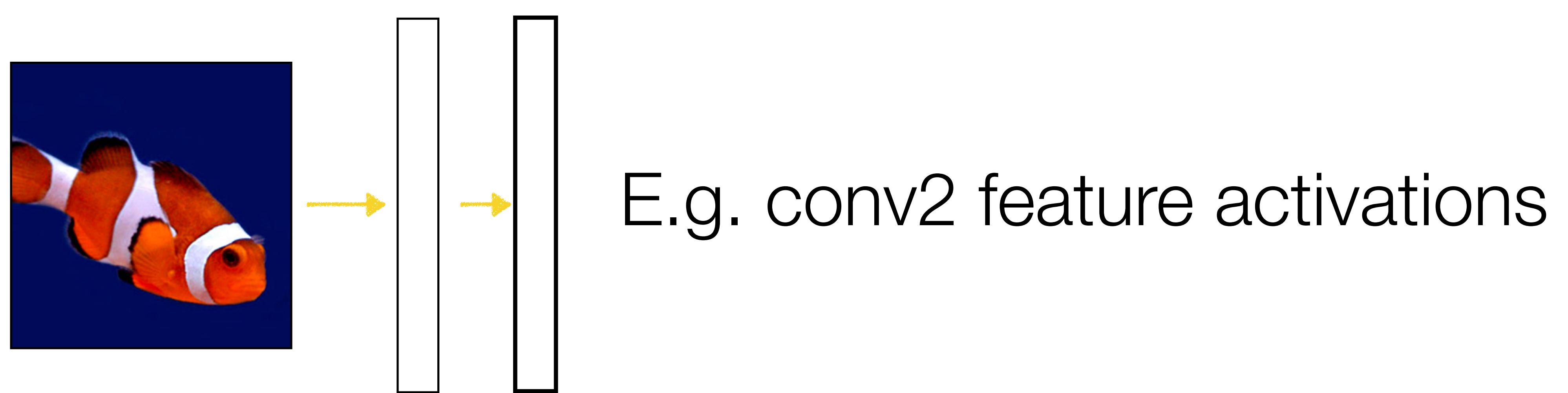
96 Units in conv1



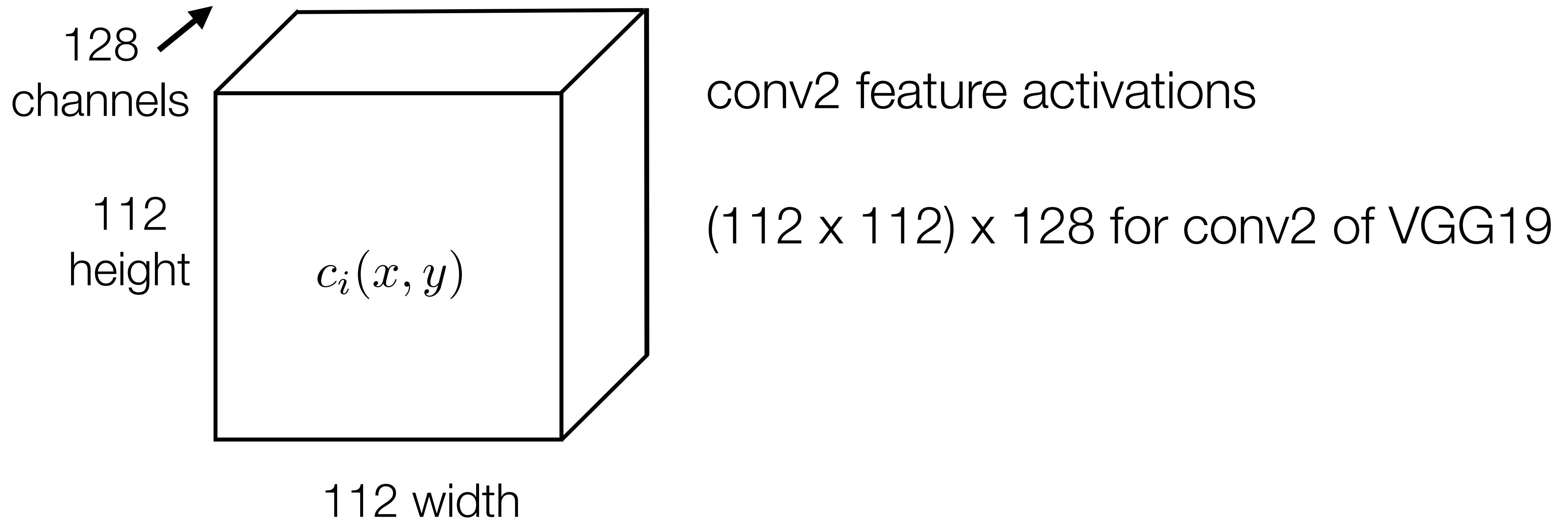
Extracting neural net features



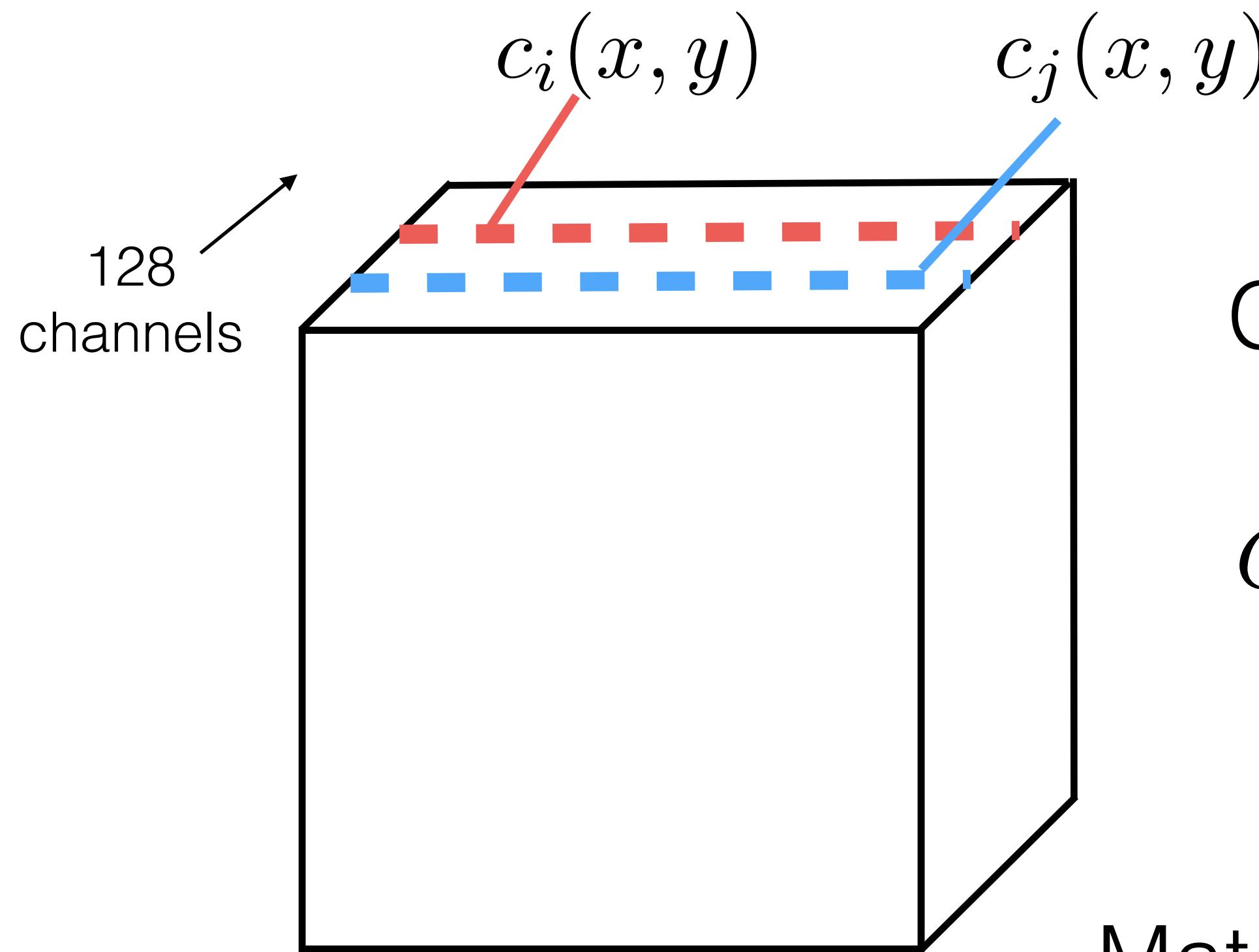
Extracting neural net features



Extracting neural net features



Texture synthesis

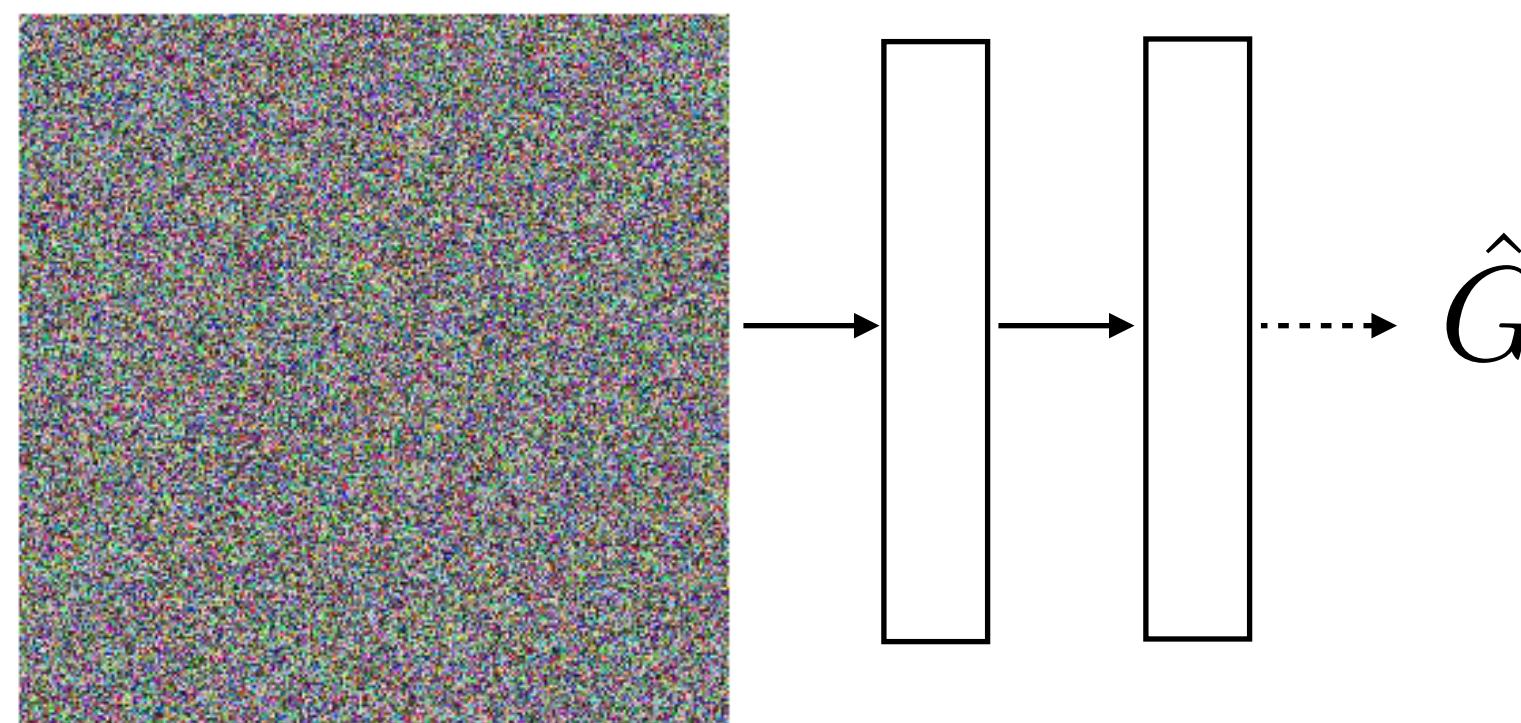


Gram (\approx covariance) matrix:

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

[Gatys et al. 2016]

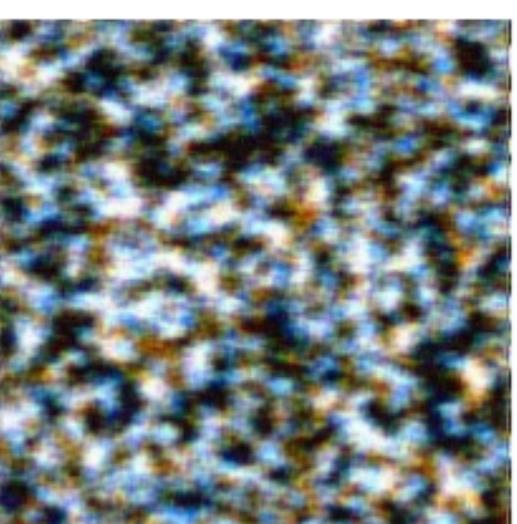
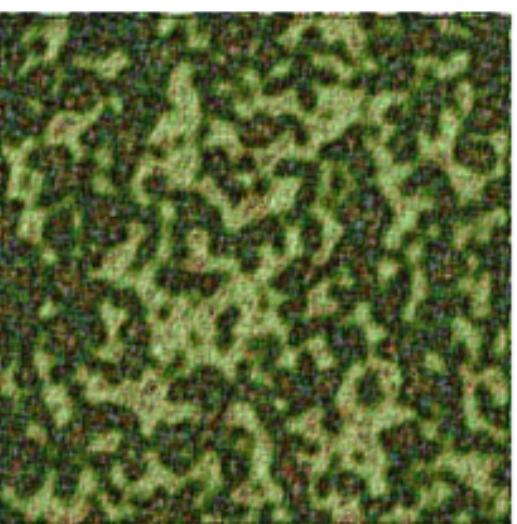
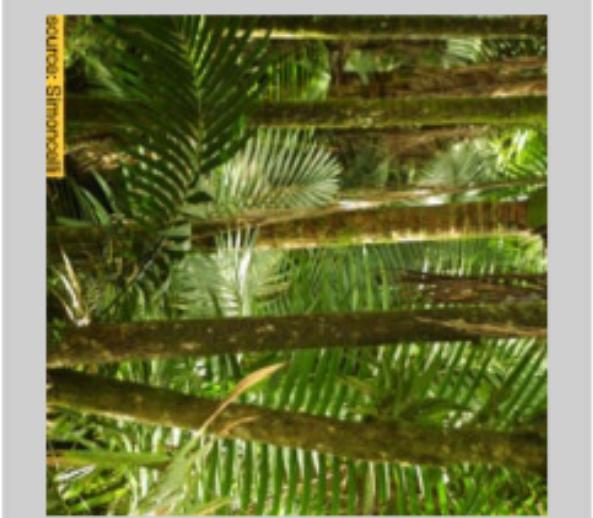
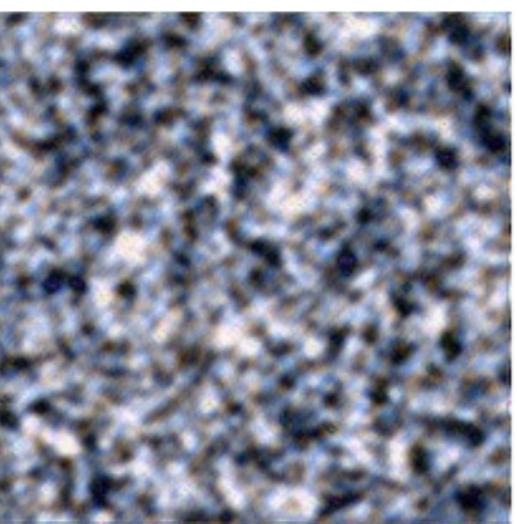
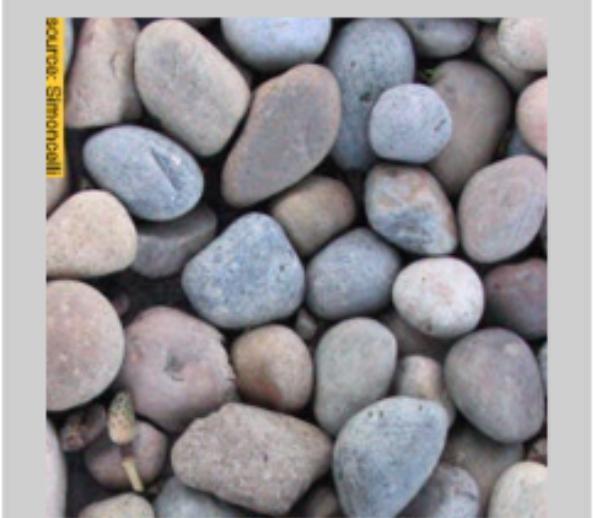
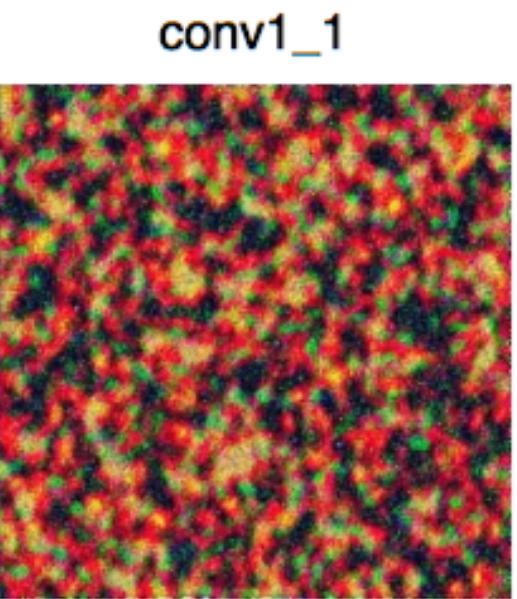
Match target image stats! Minimize:



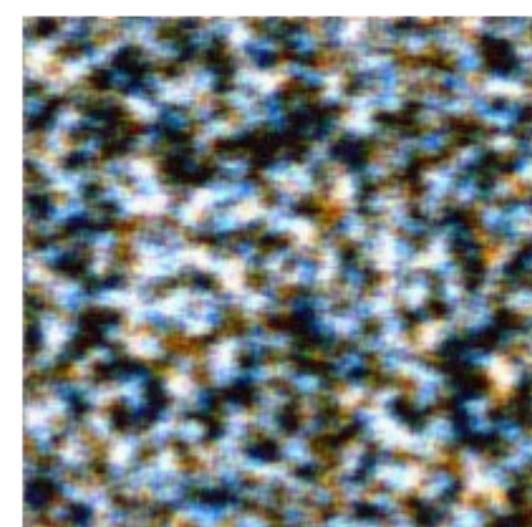
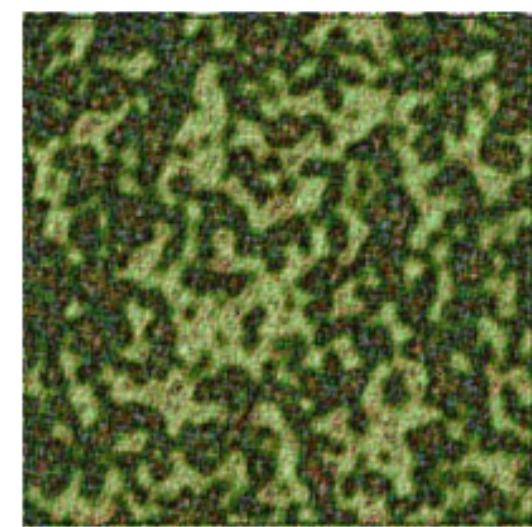
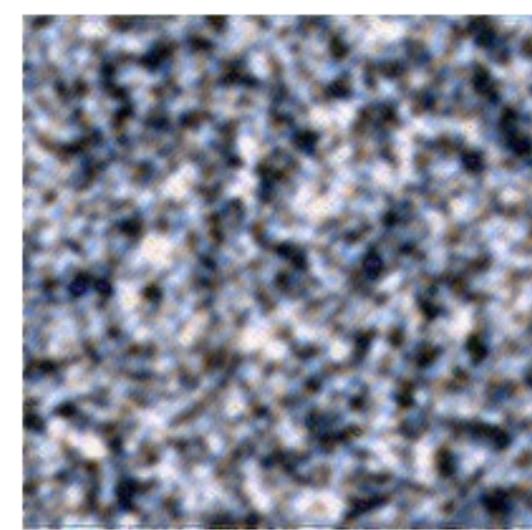
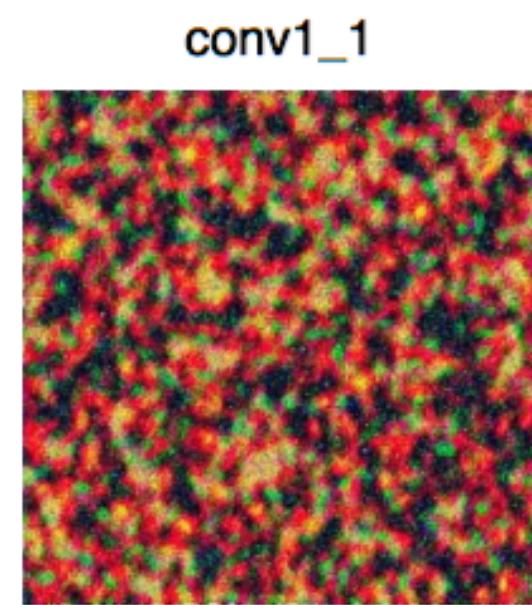
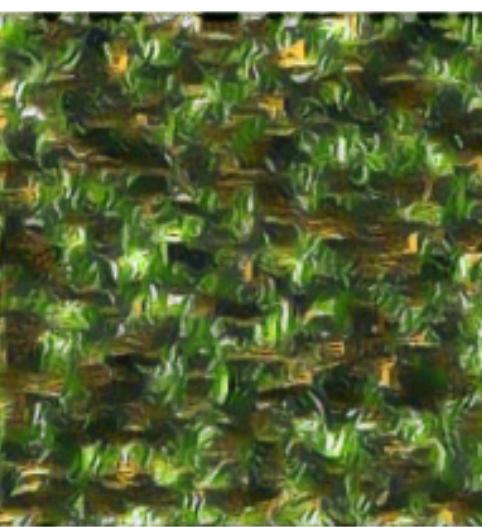
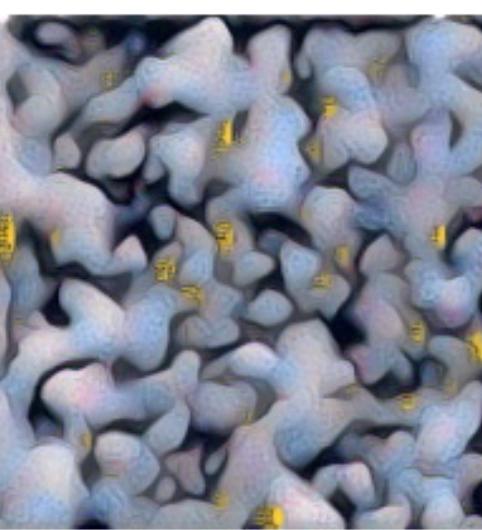
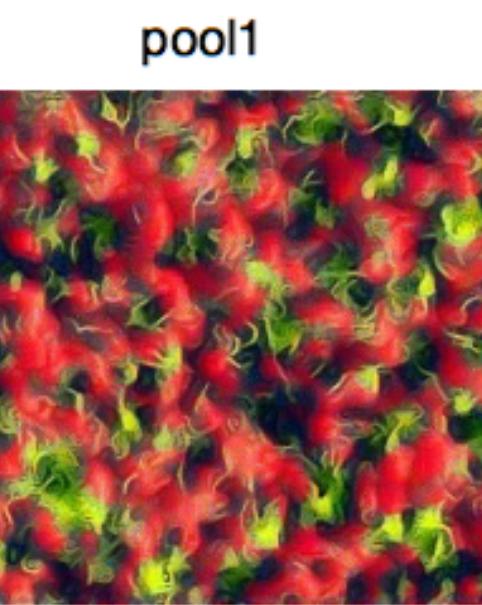
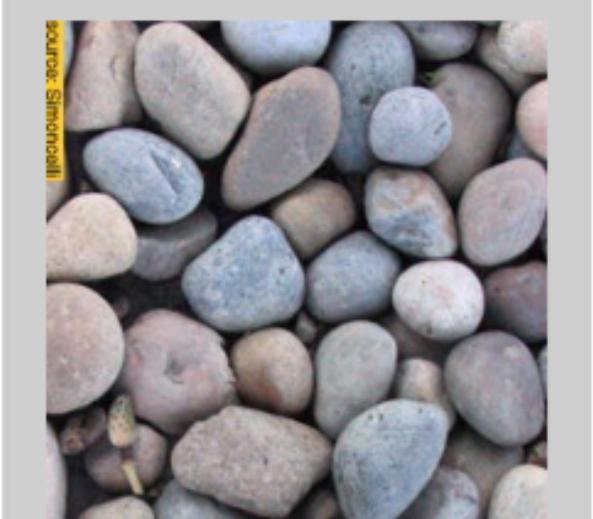
$$\sum_{i,j} (G_{ij} - \hat{G}_{ij})^2$$

- Use all layers.
- Minimize with gradient descent!

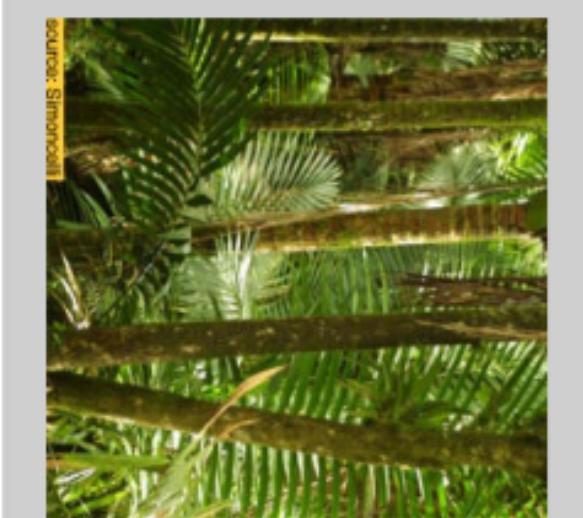
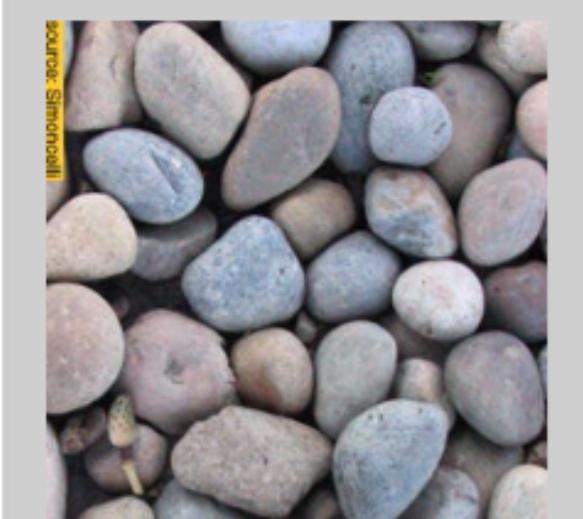
High → Low



High → Low



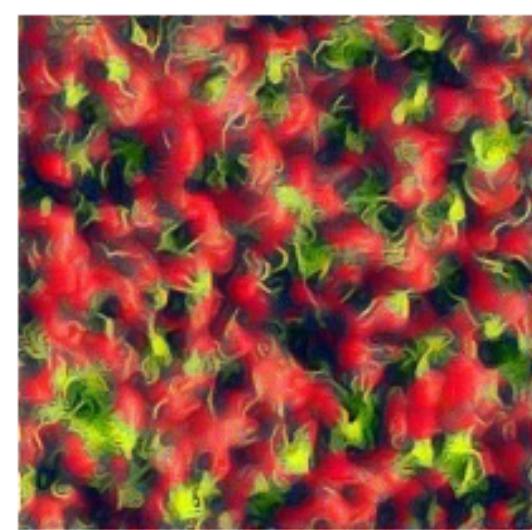
High → Low



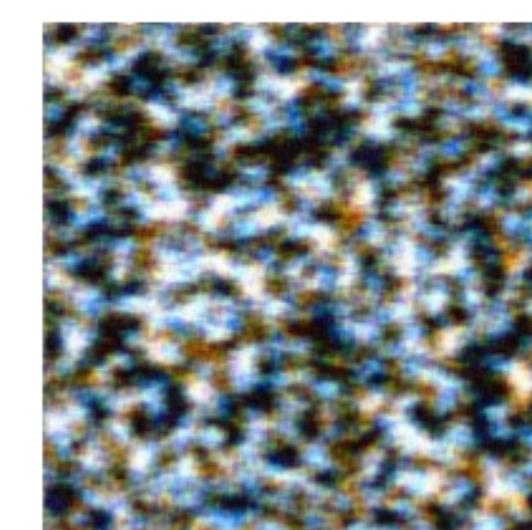
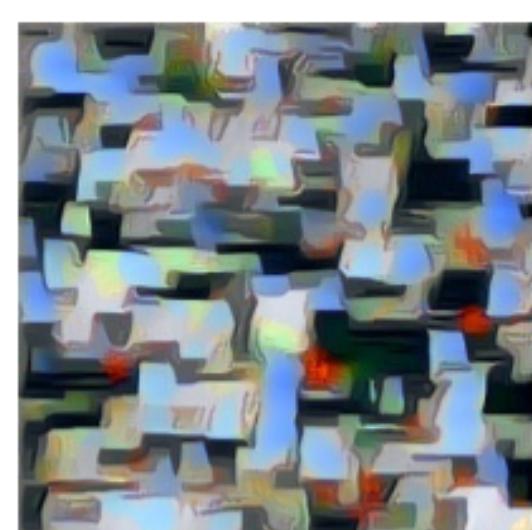
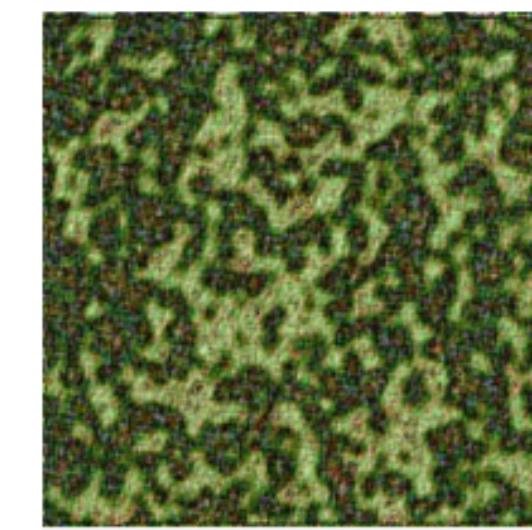
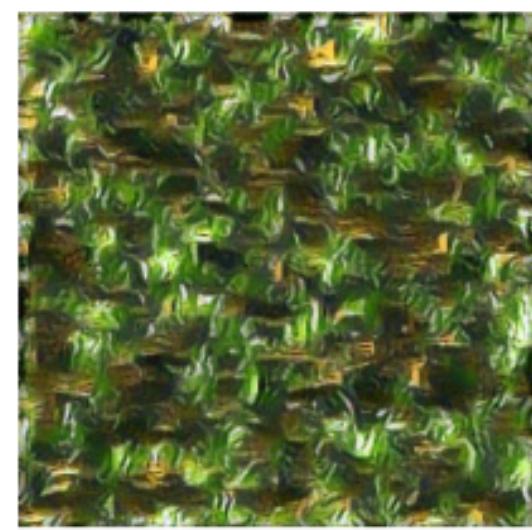
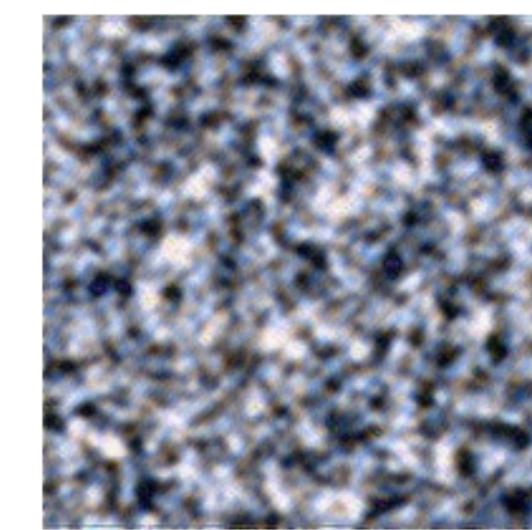
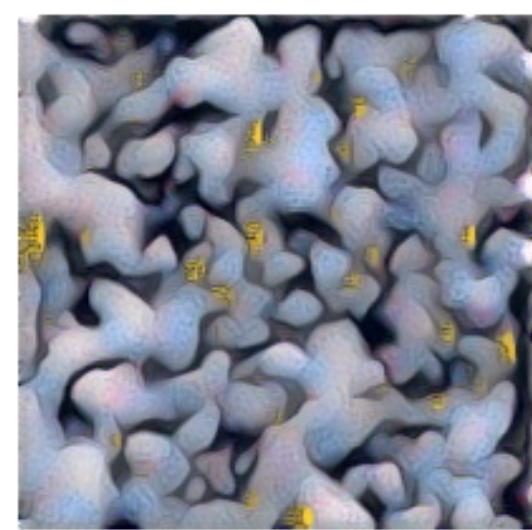
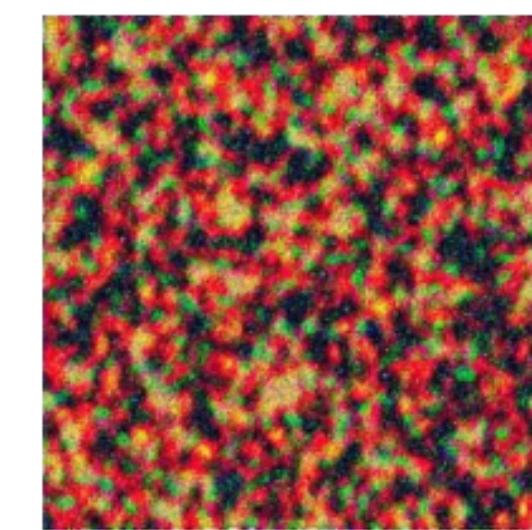
pool2



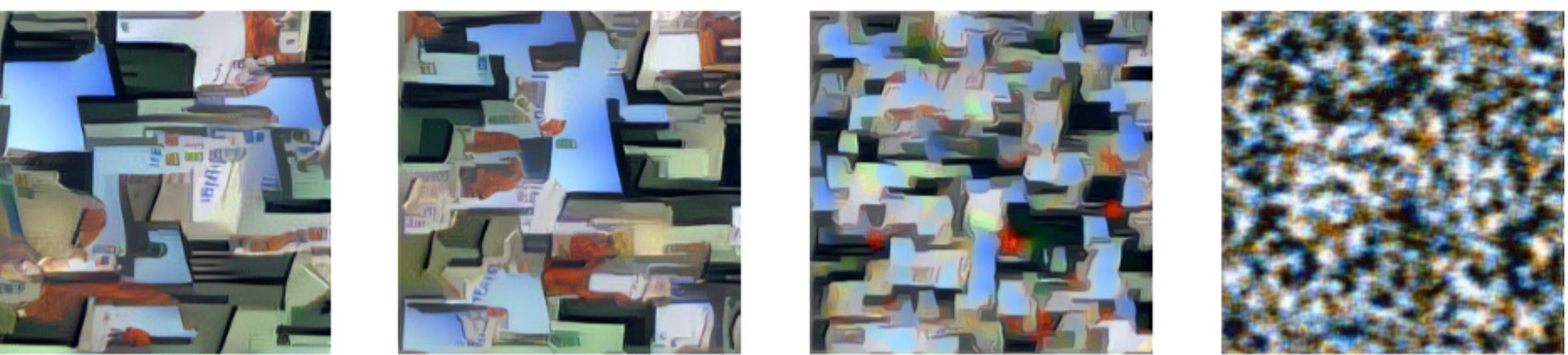
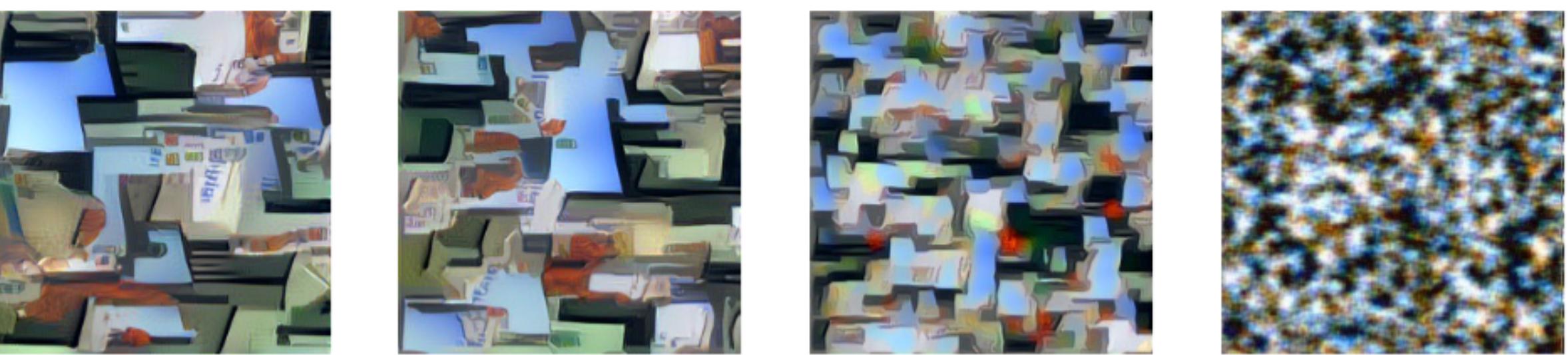
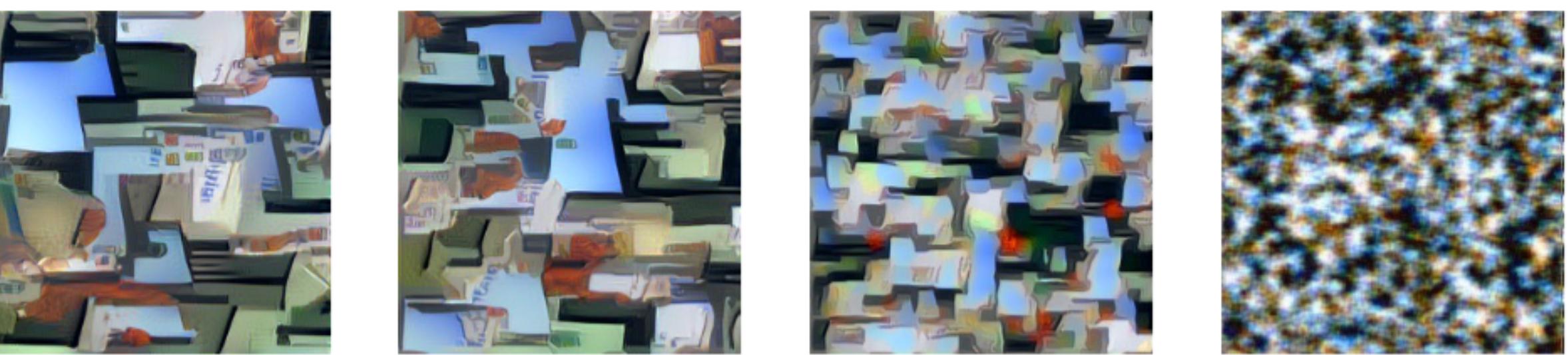
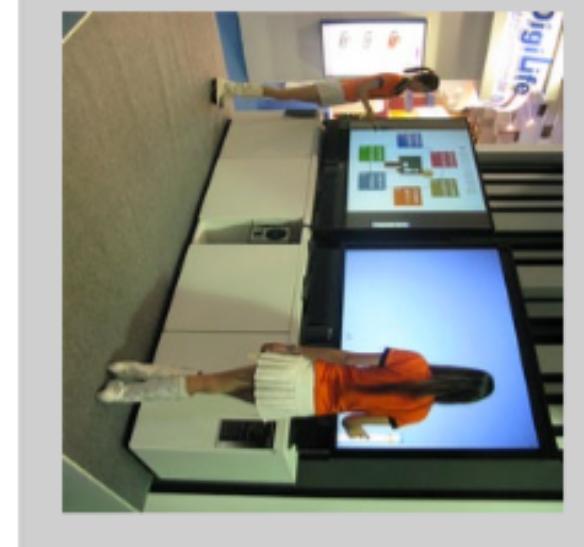
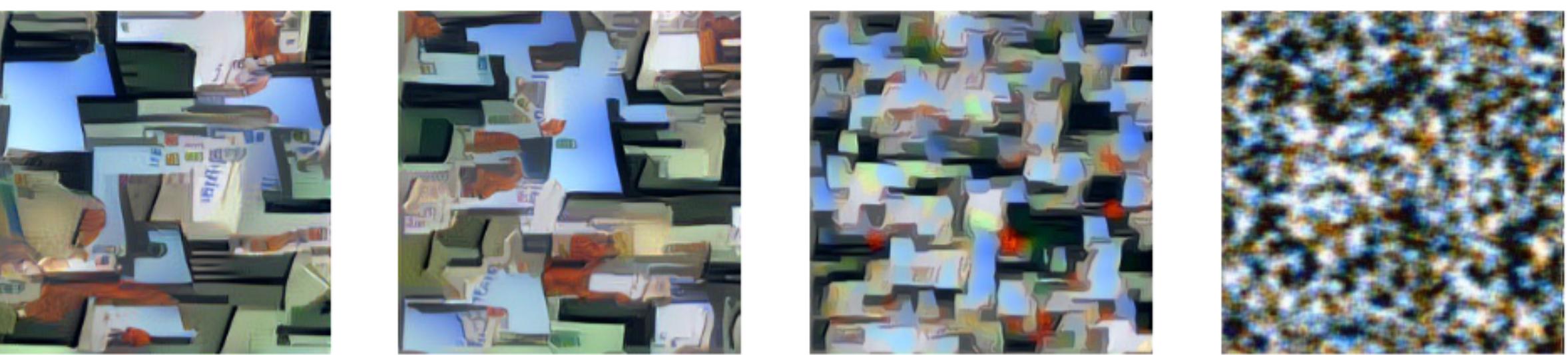
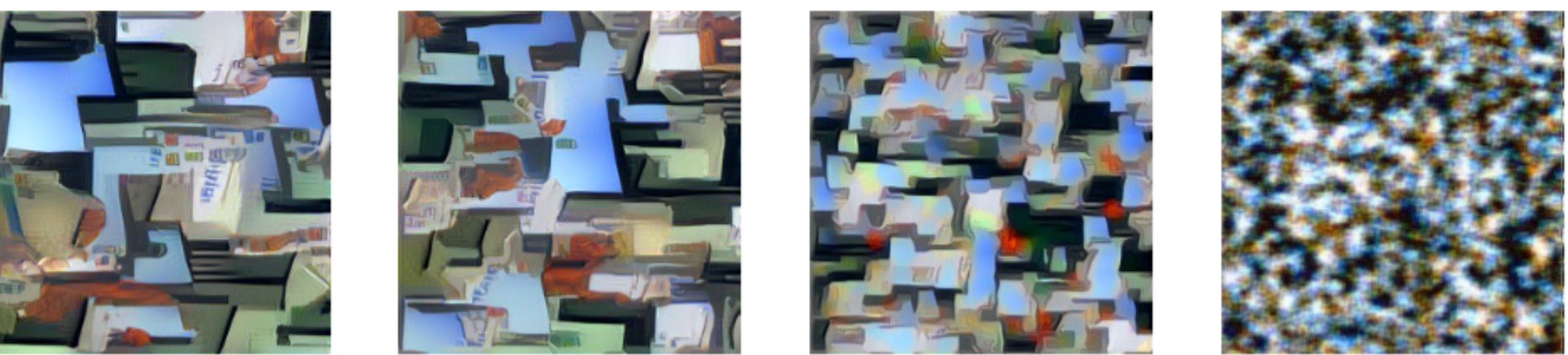
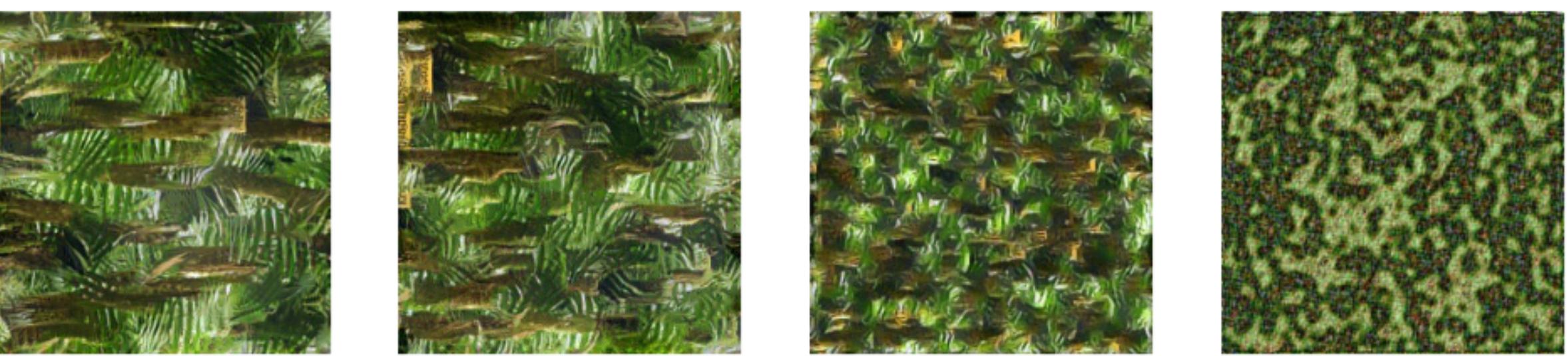
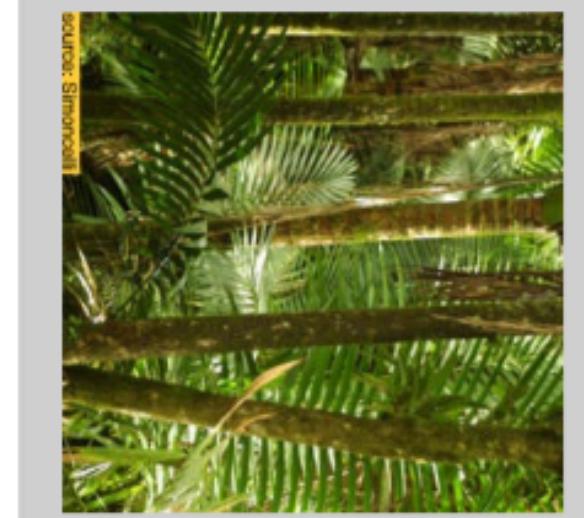
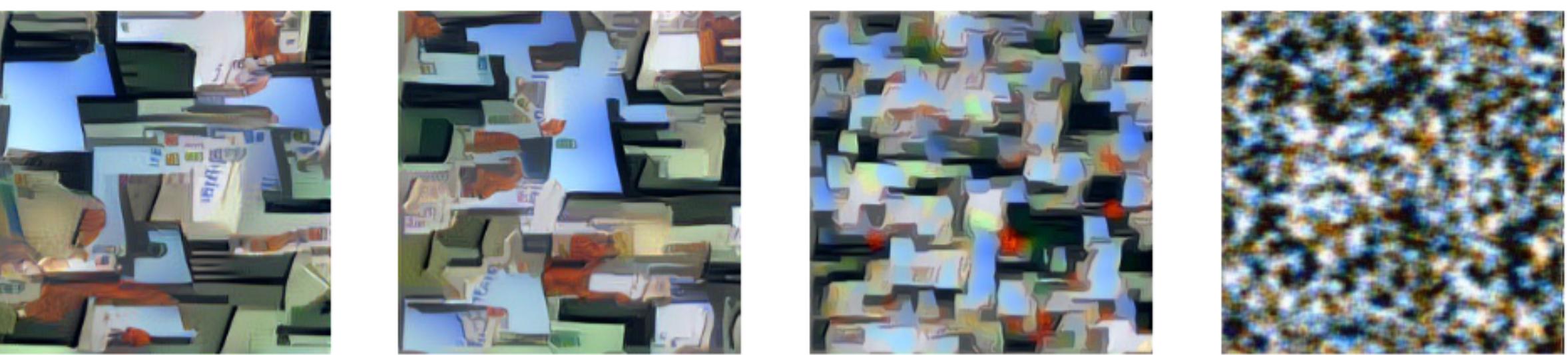
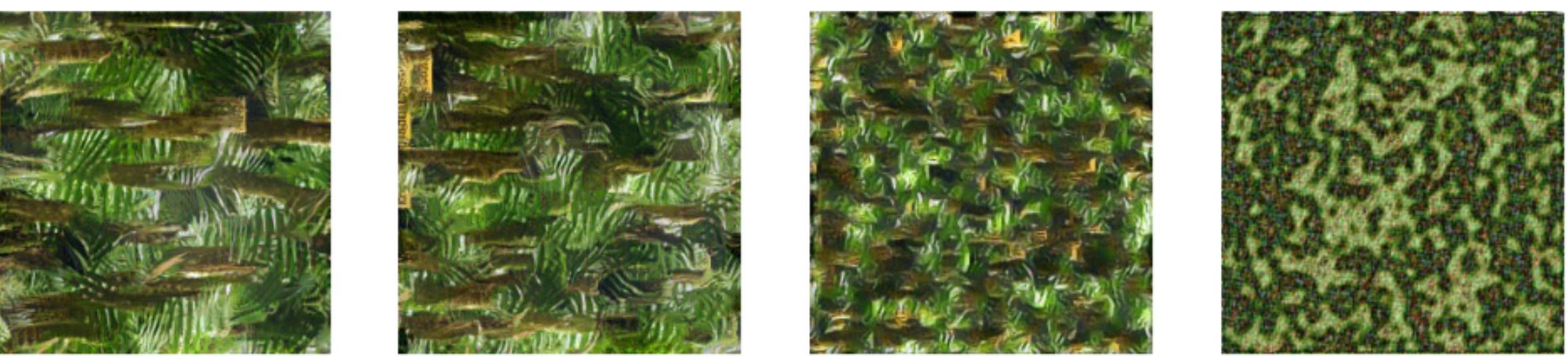
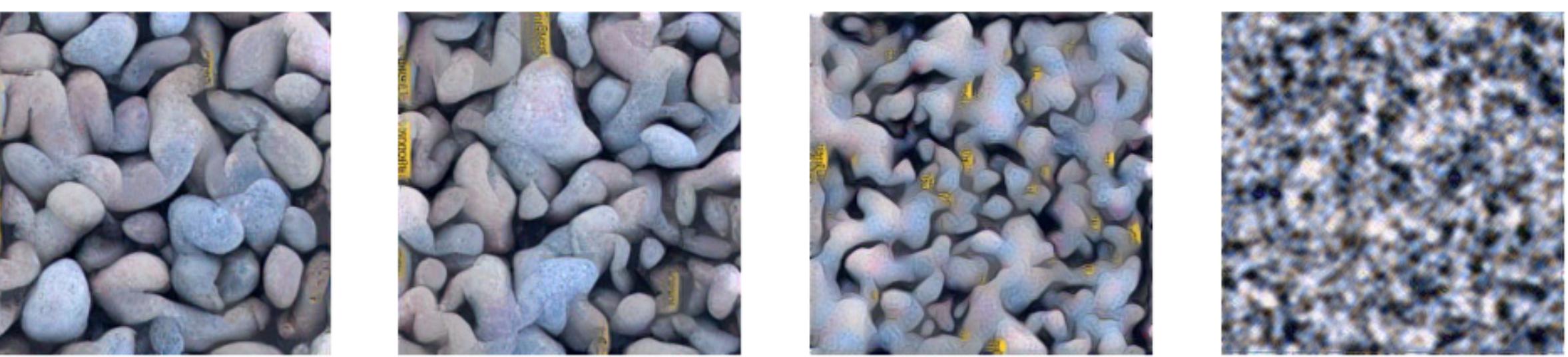
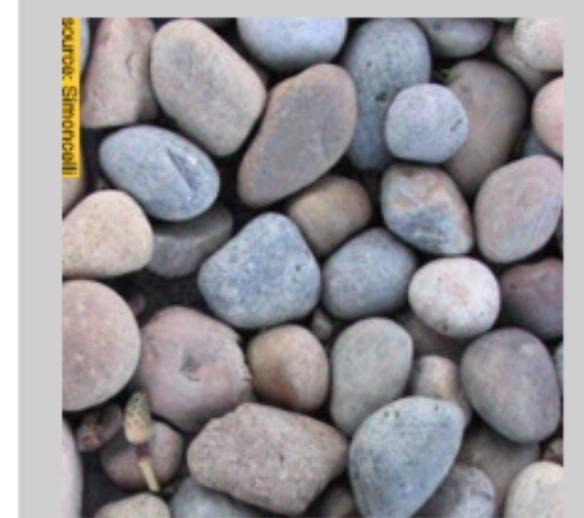
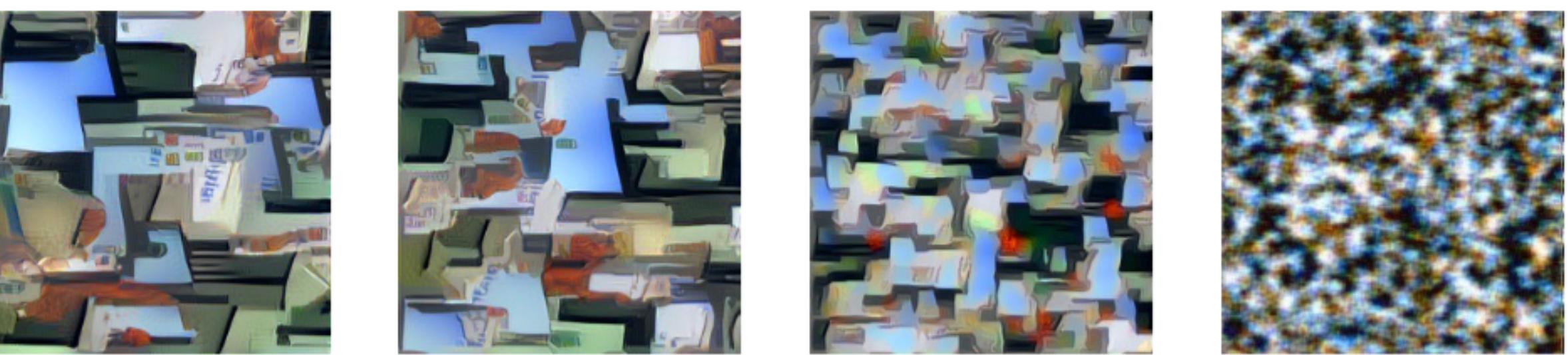
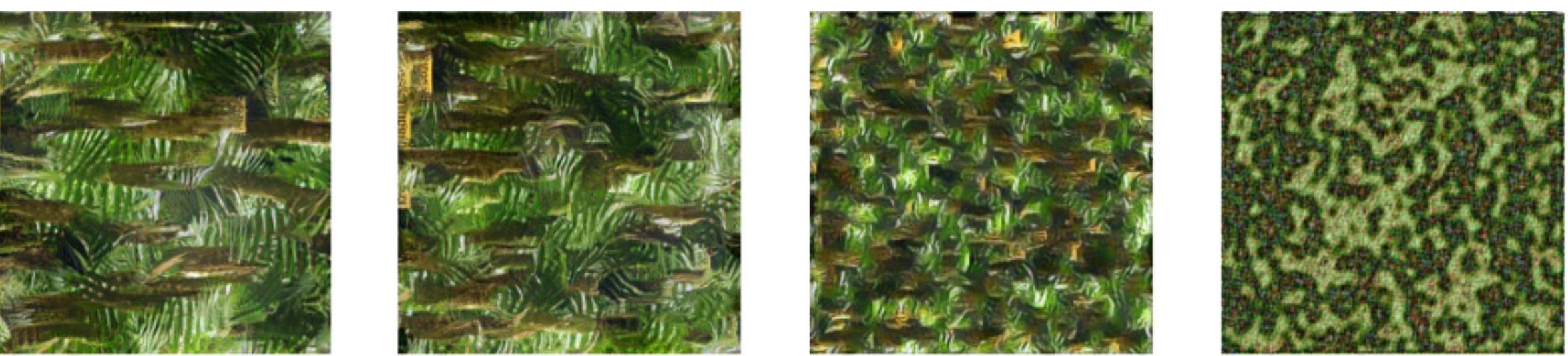
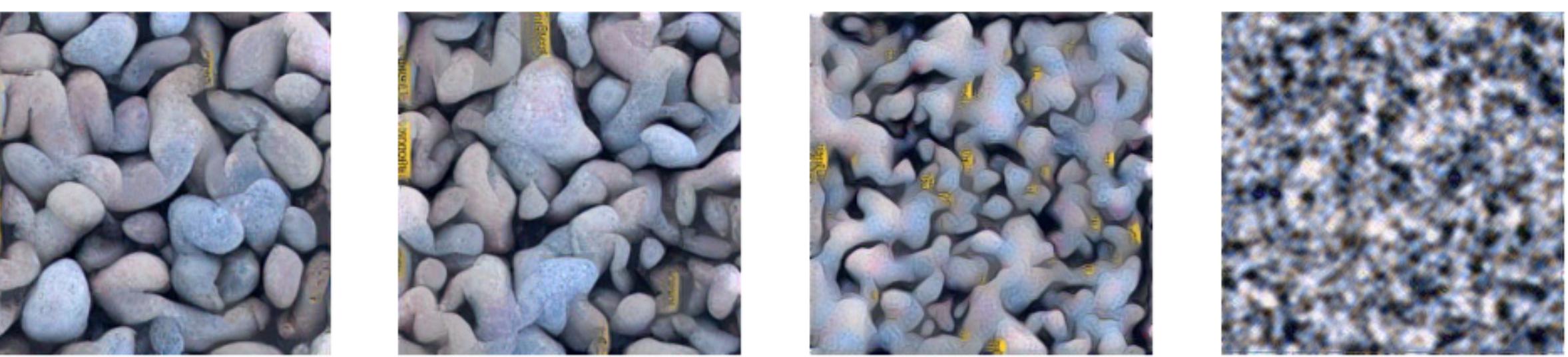
pool1



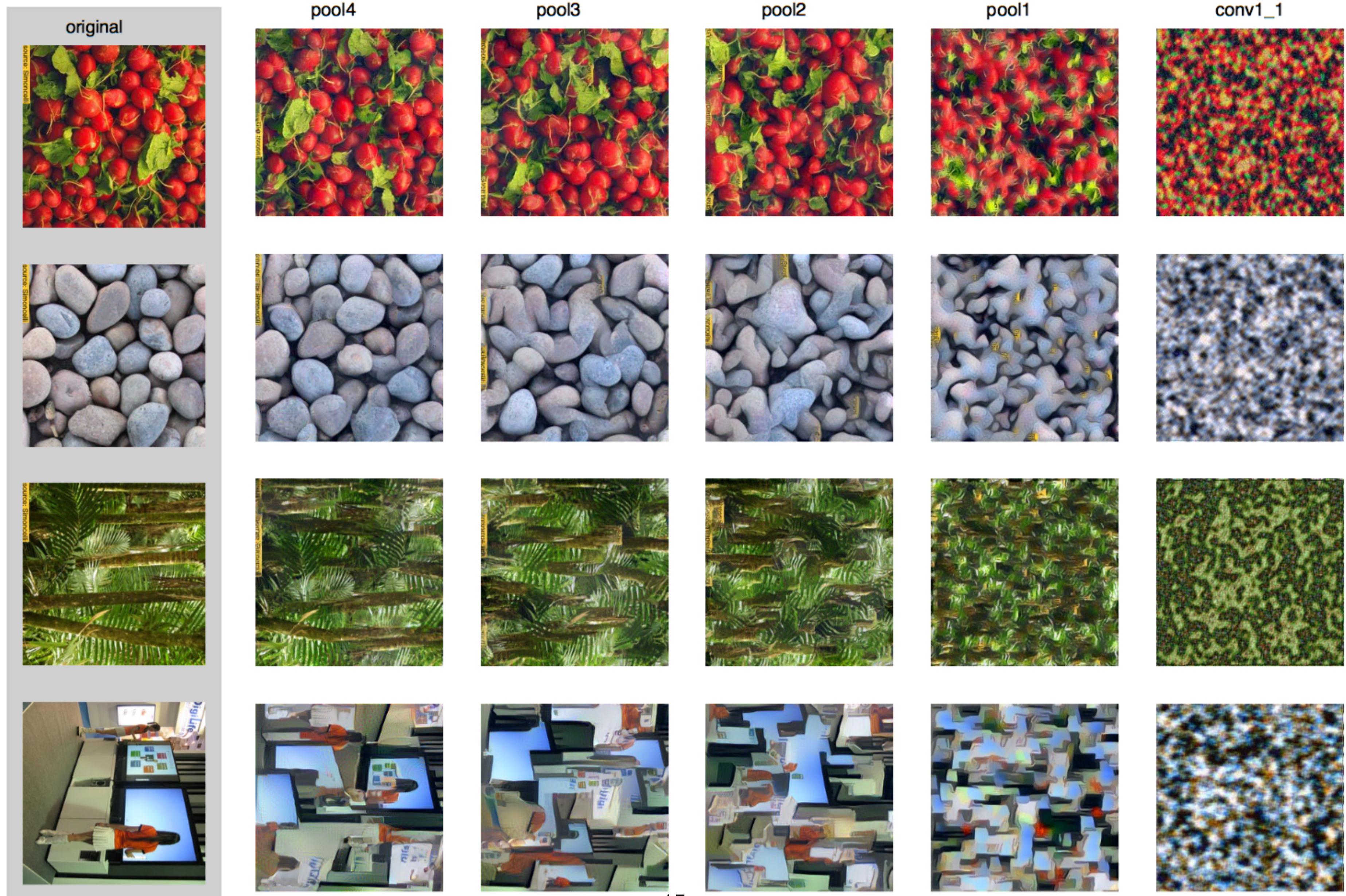
conv1_1



High → Low



High → Low



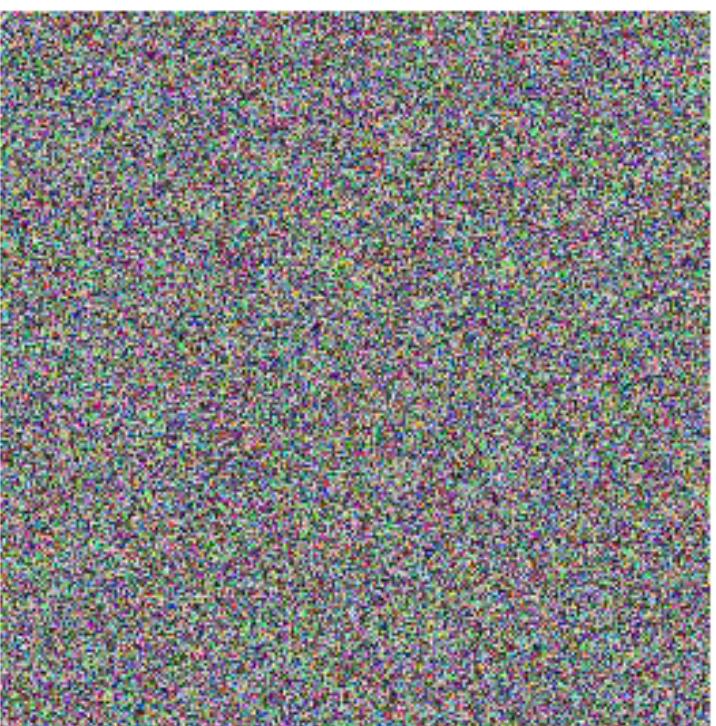
Texture captures artistic style

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

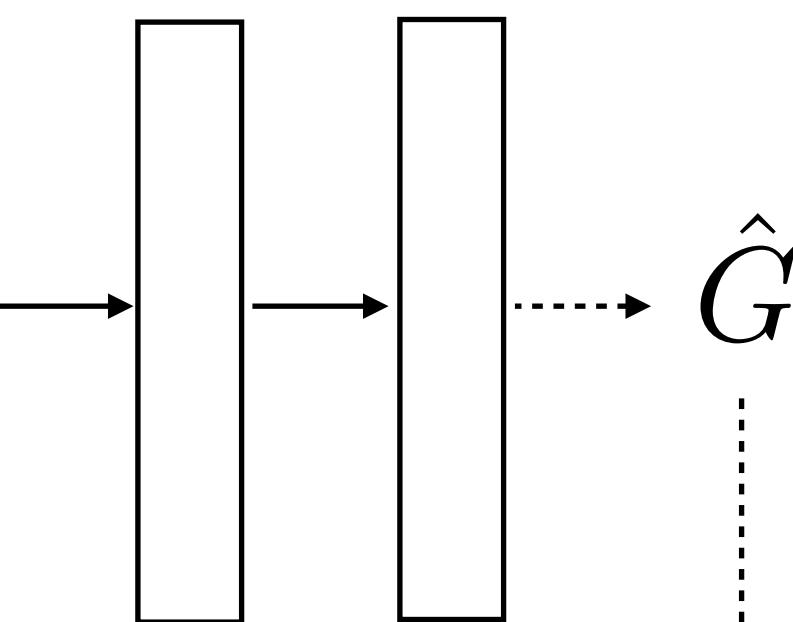


[Gatys et al. 2016]

Match the **style** of the painting.



Synthesized image



“perceptual loss”

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

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



$$c_i(x, y)$$

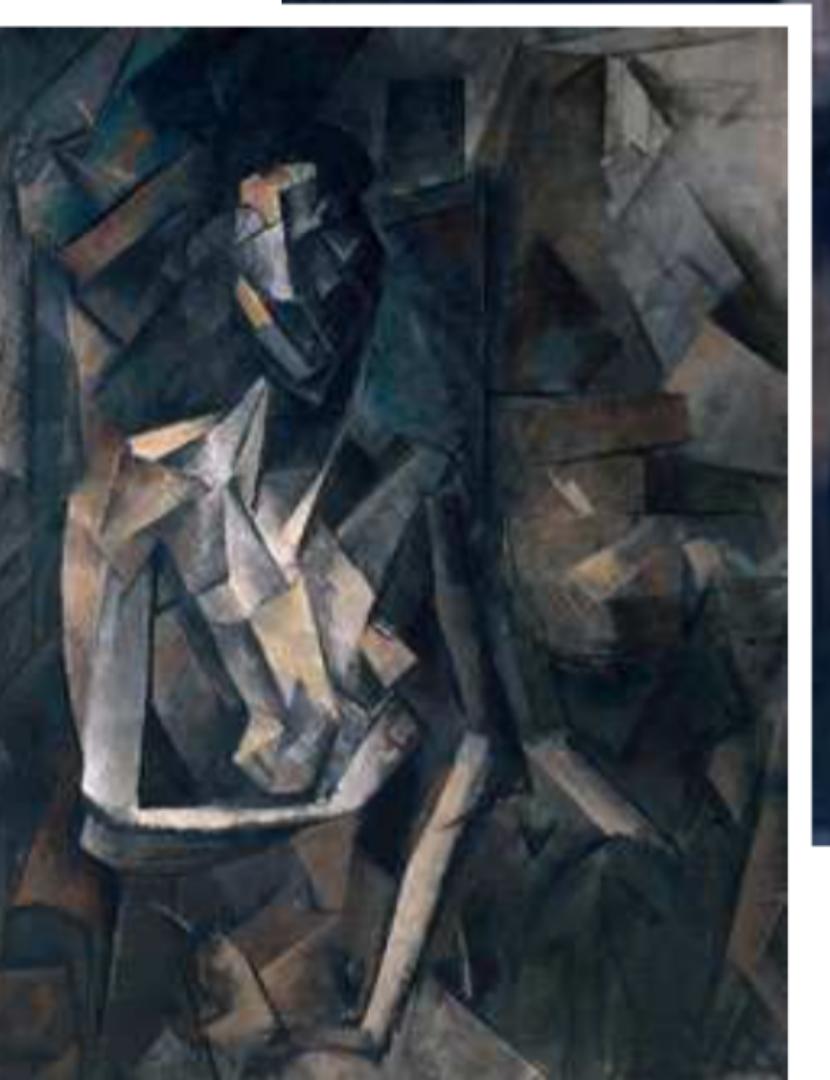
$$\hat{c}_i(x, y)$$











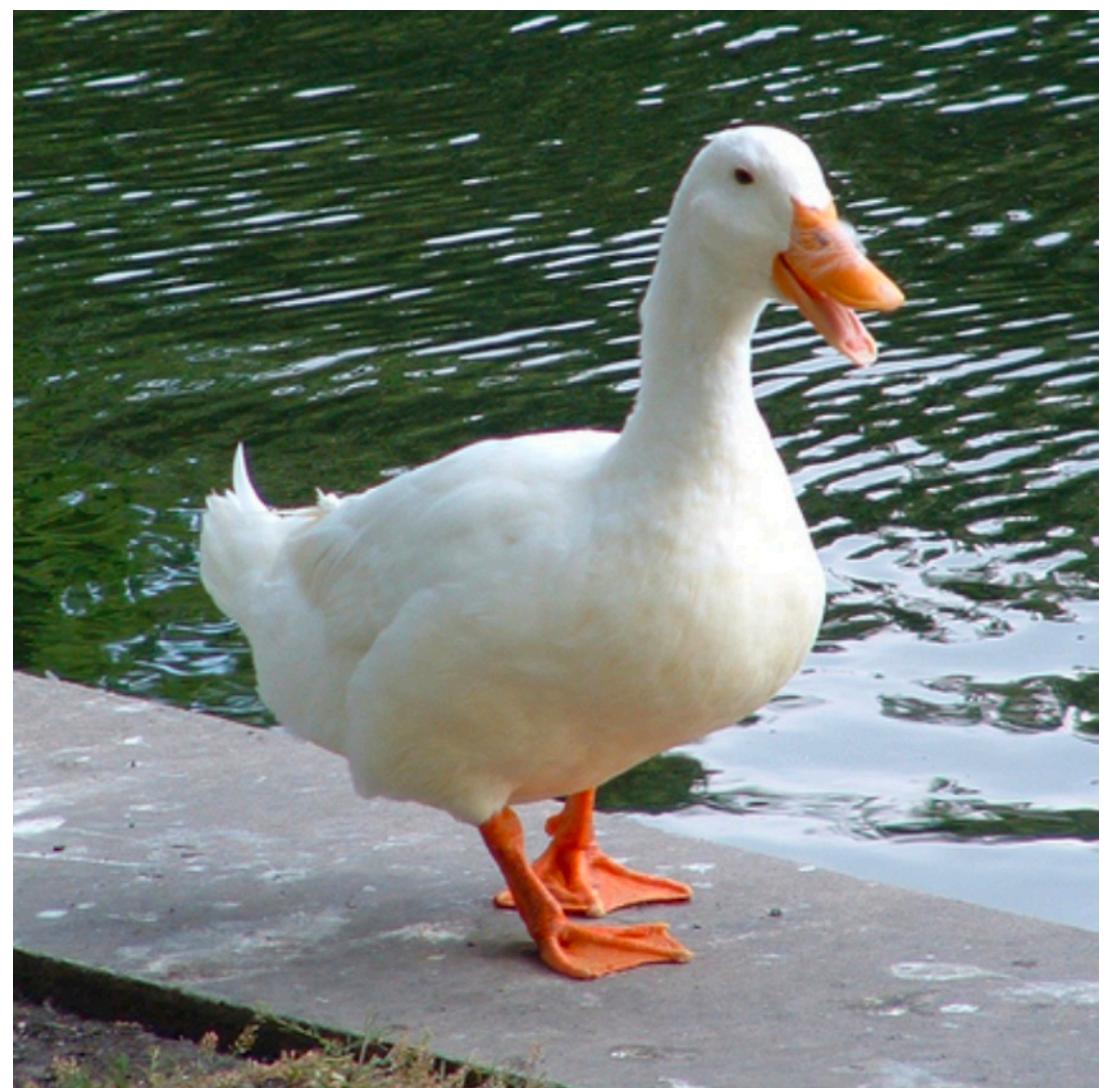
London during the day.



New York at night.



Image classification



“Duck”

:

image x

label y

Image synthesis

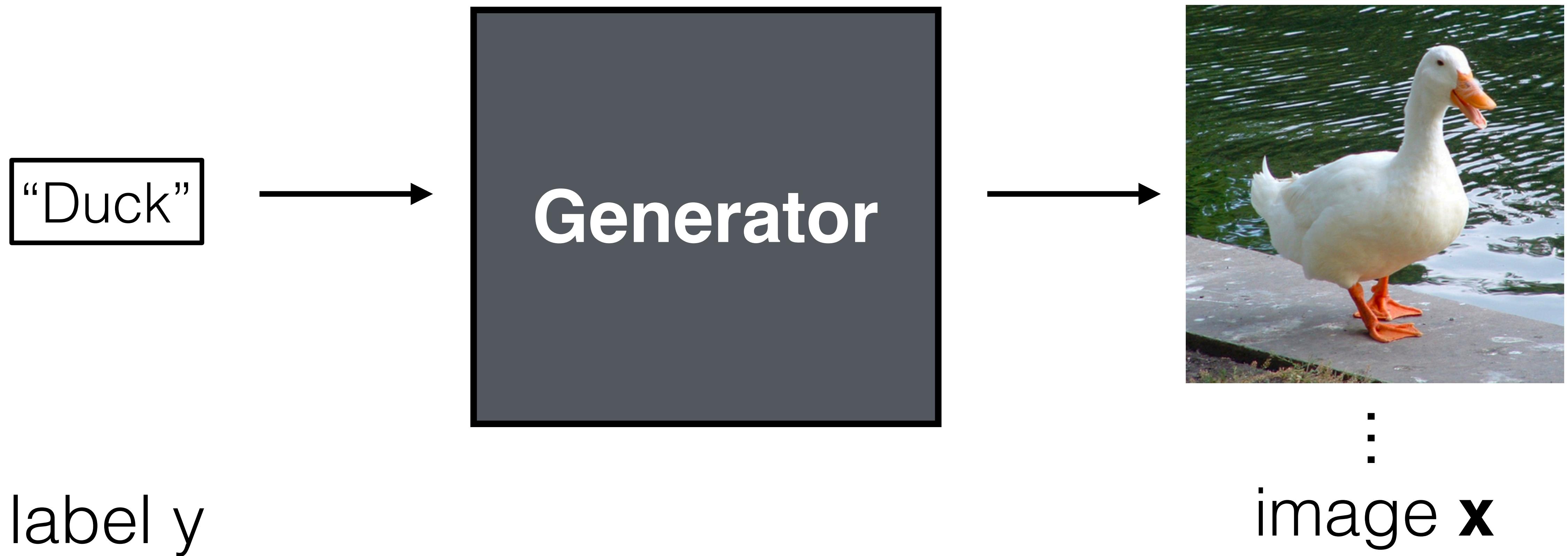


Image translation

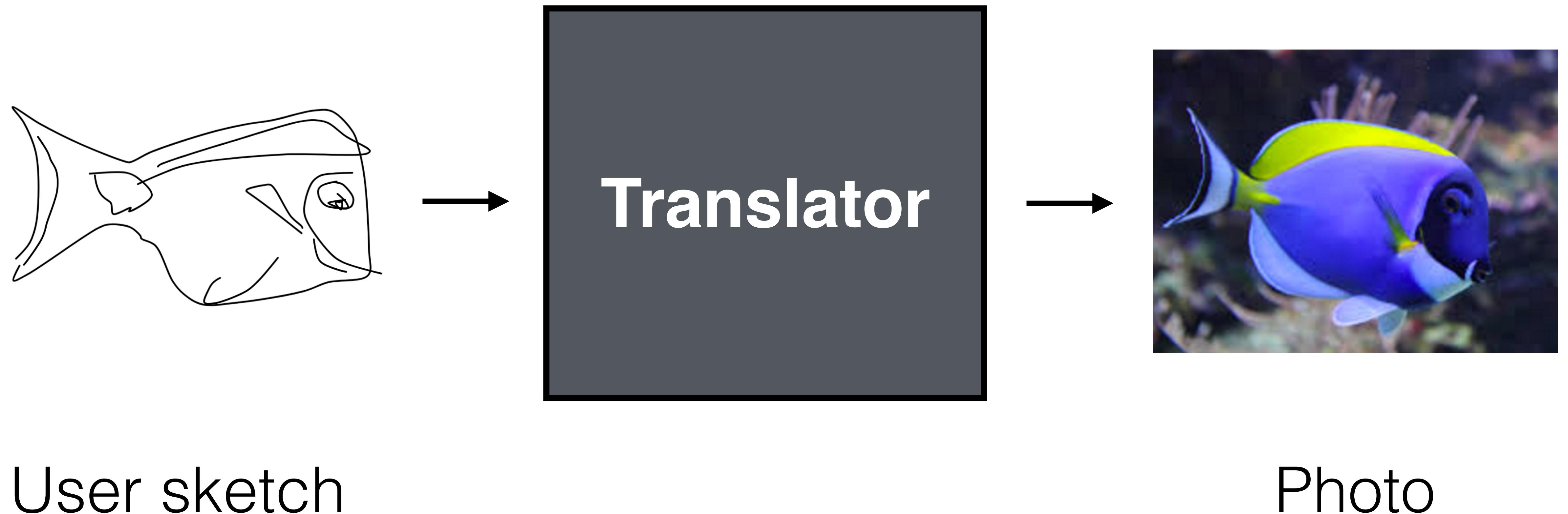
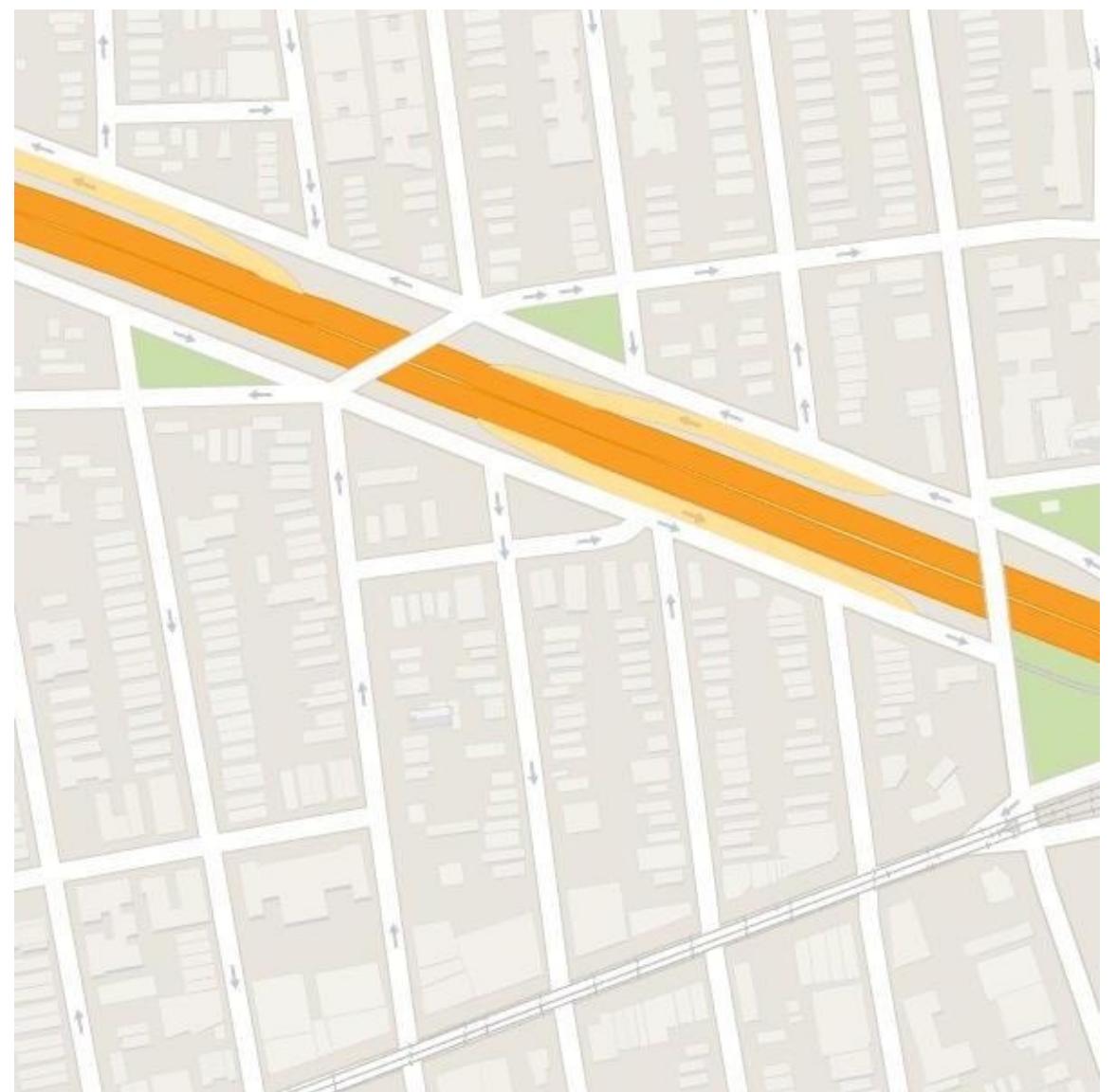
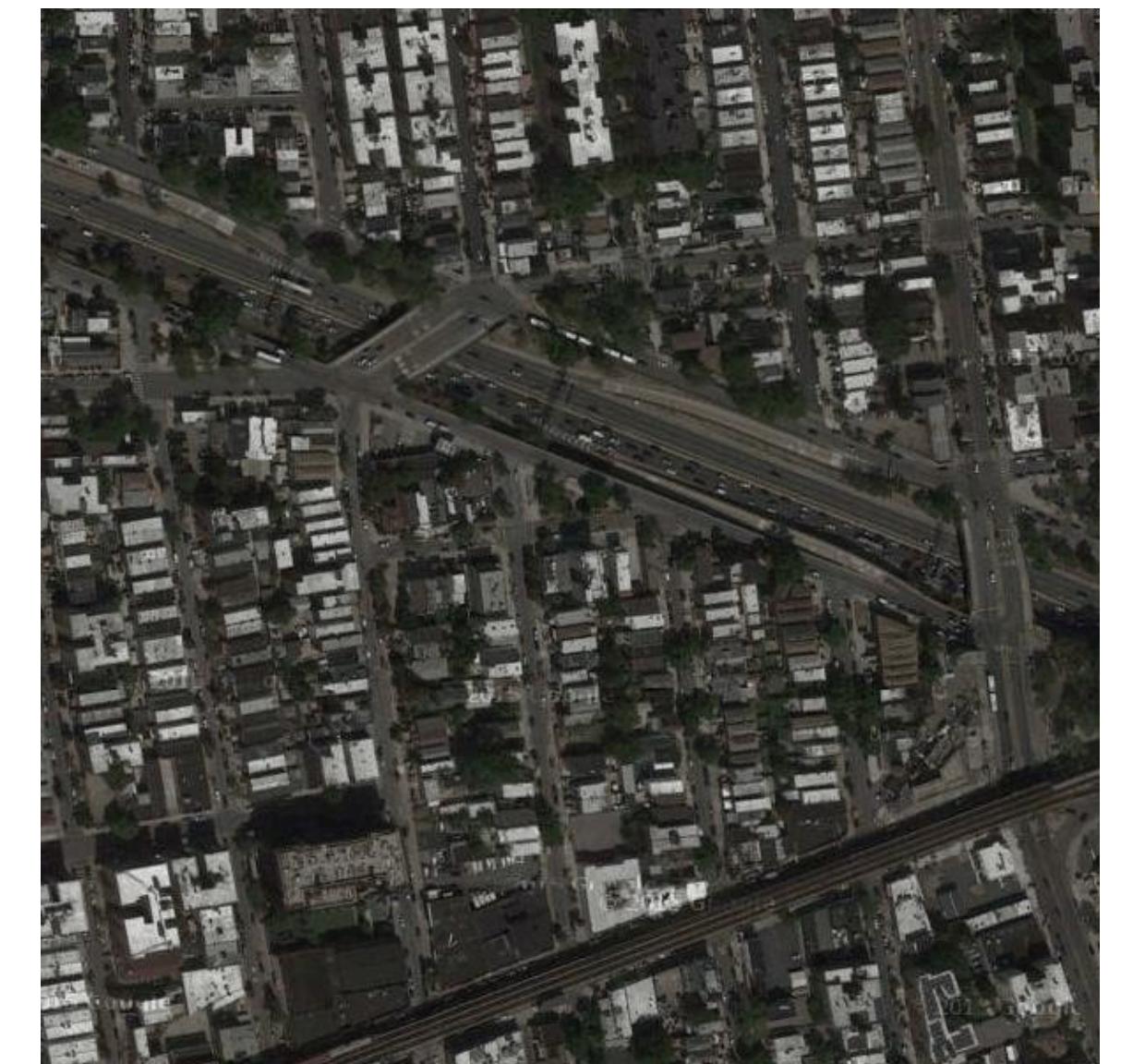


Image translation



Google Map

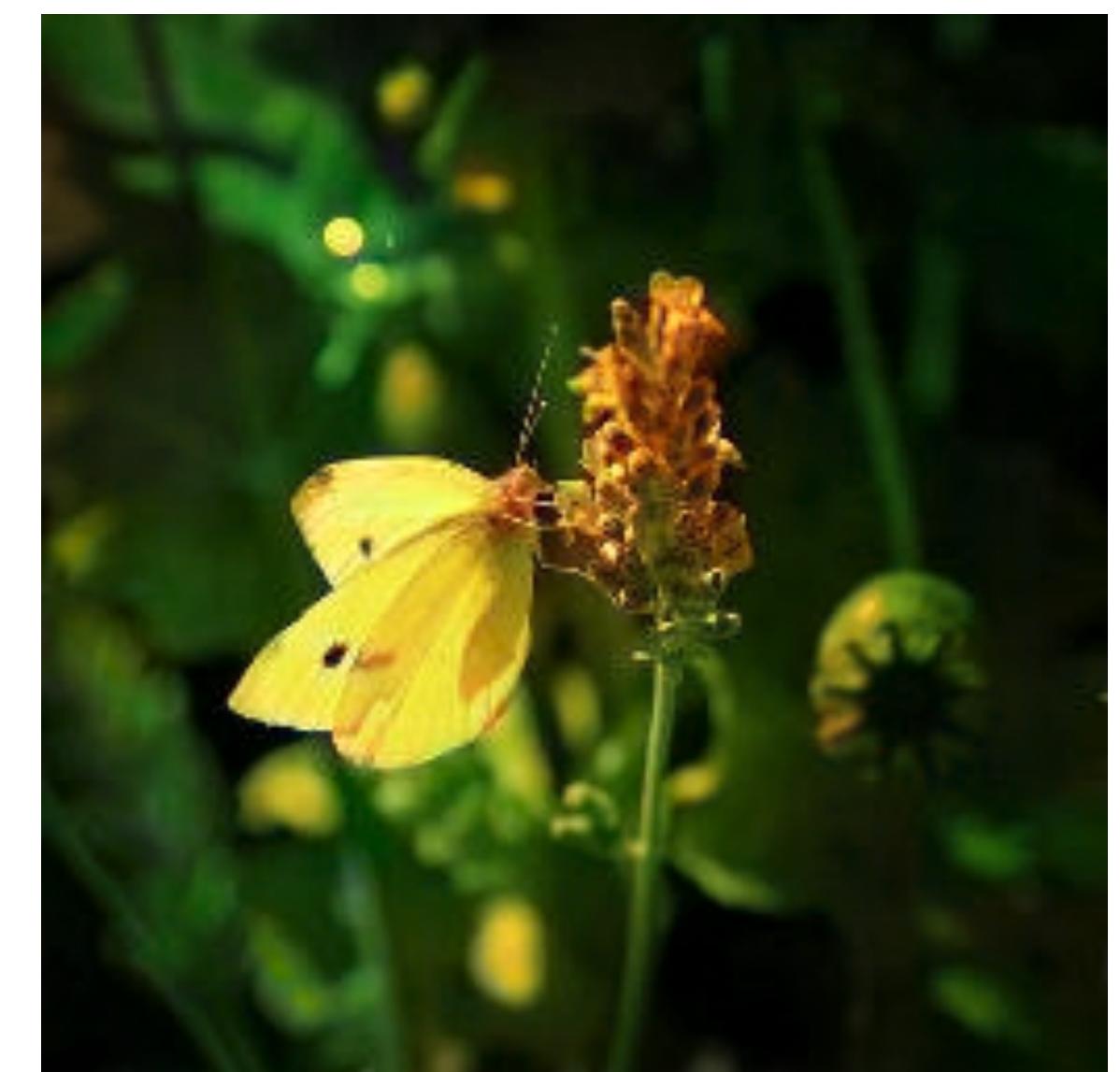


Satellite photo

Image translation



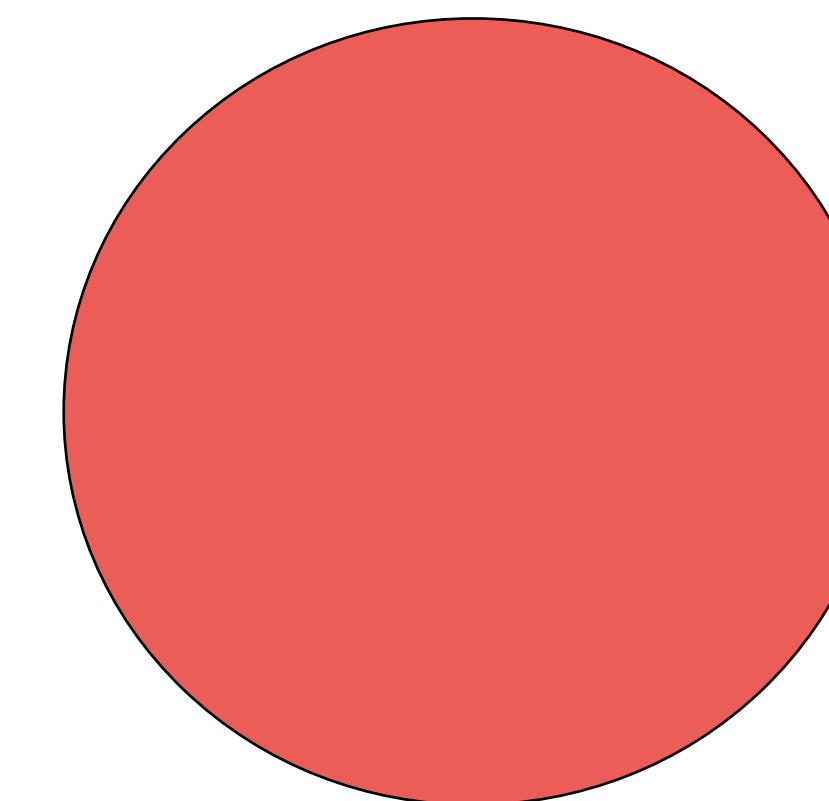
BW image



Color image

Deep generative models as distribution transformers

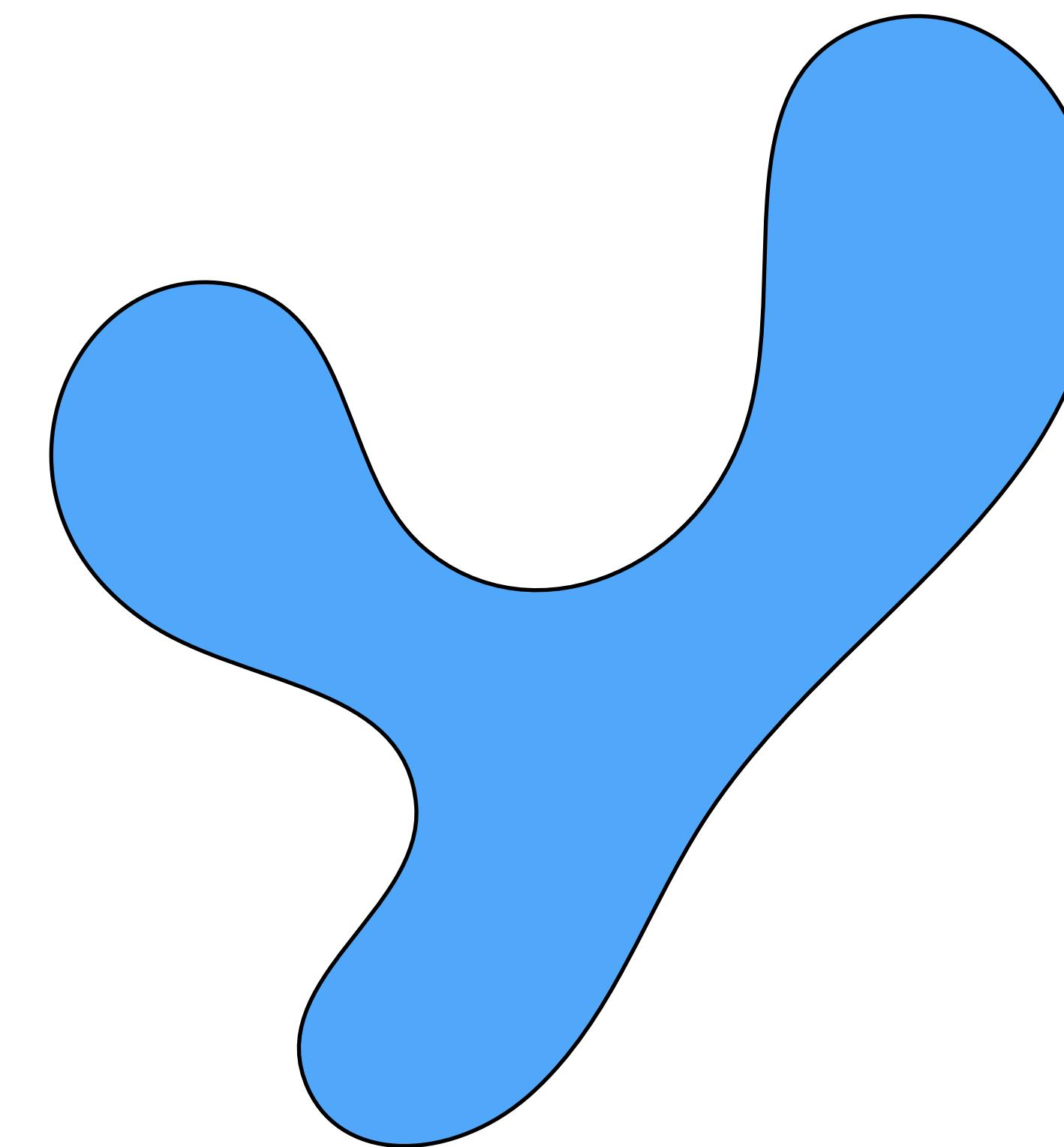
Source distribution



$$p(z)$$

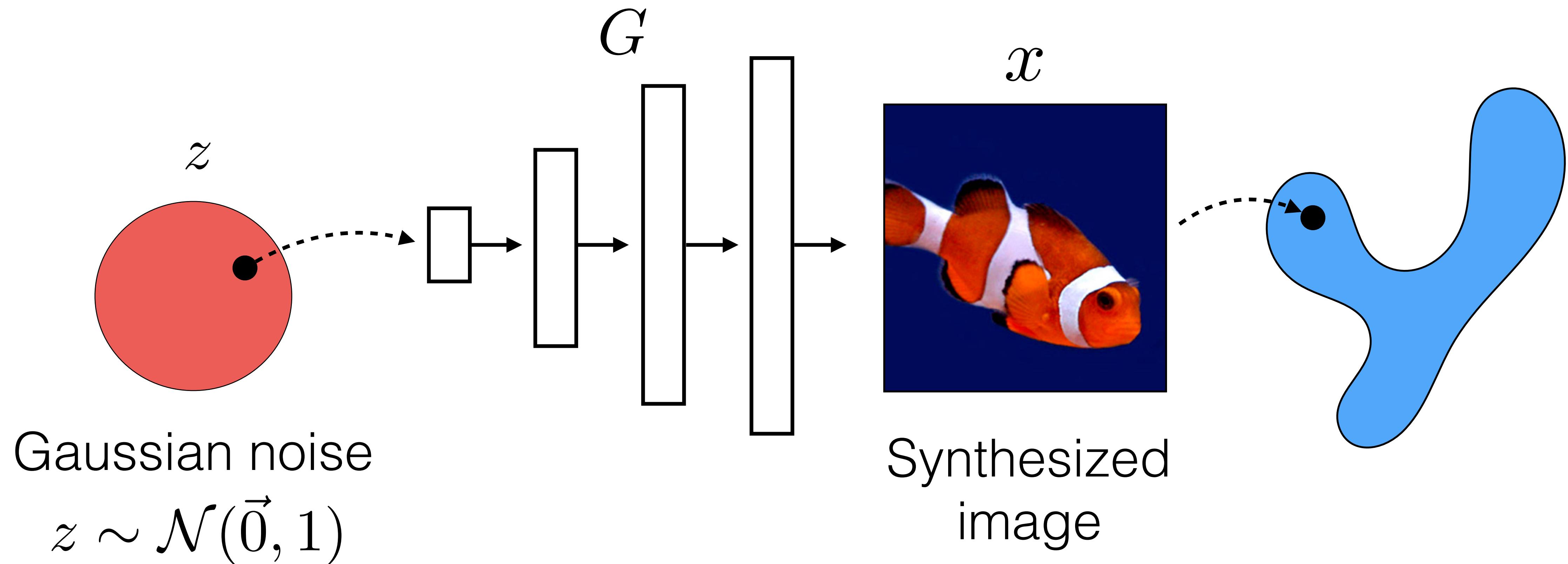


Target distribution

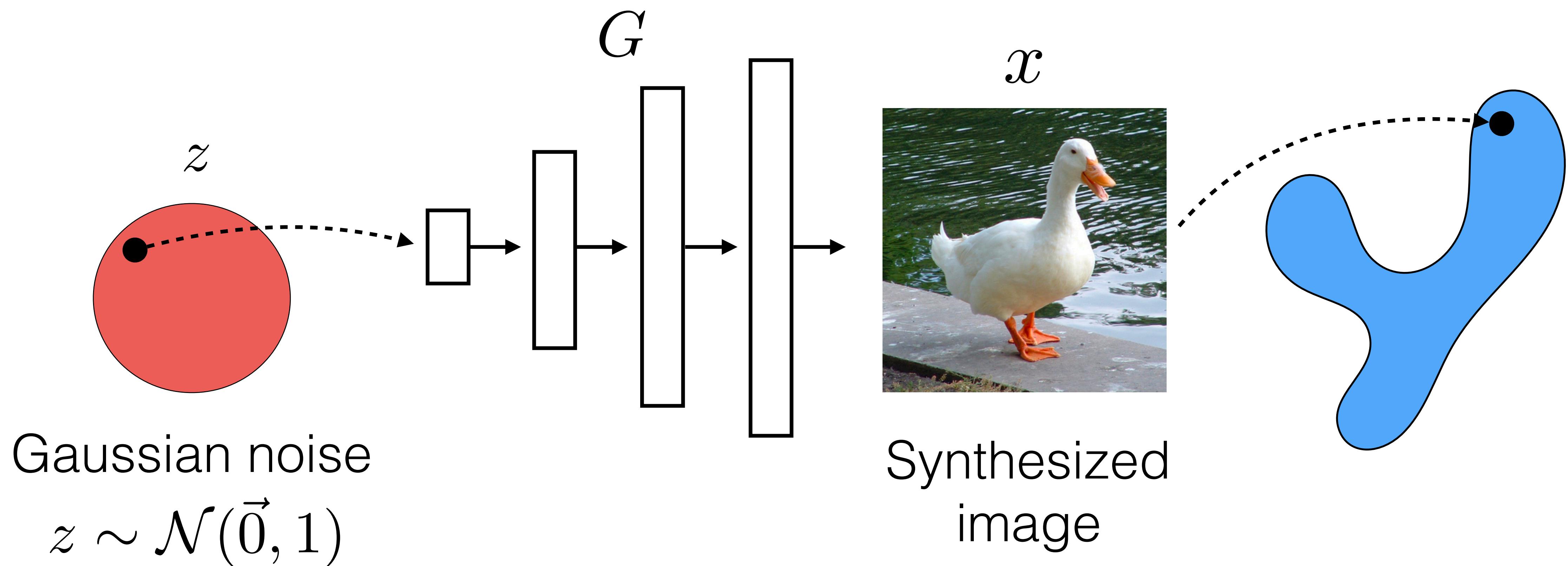


$$p(x)$$

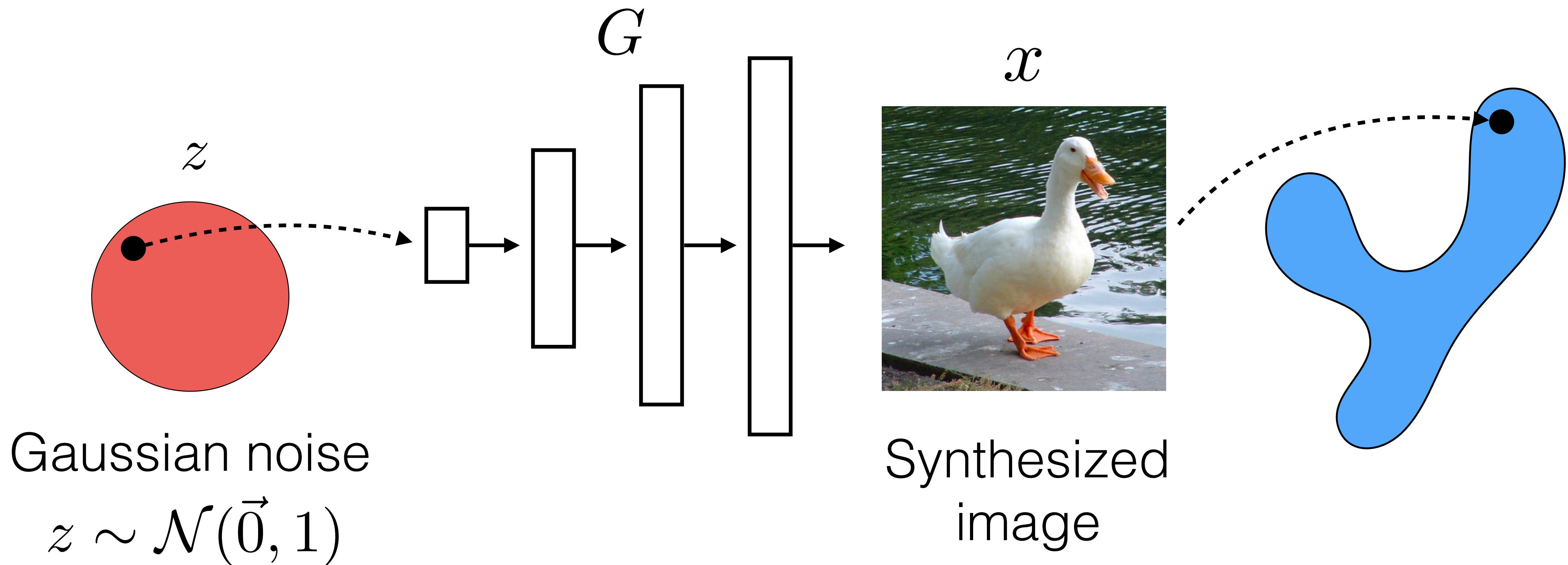
Deep generative models as distribution transformers

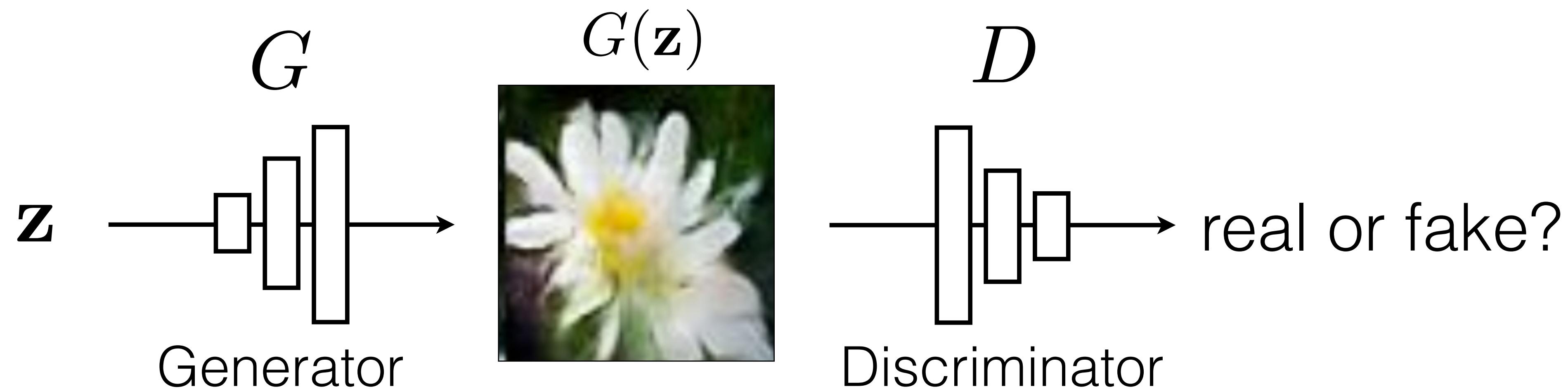


Deep generative models as distribution transformers



Generative Adversarial Networks (GANs)

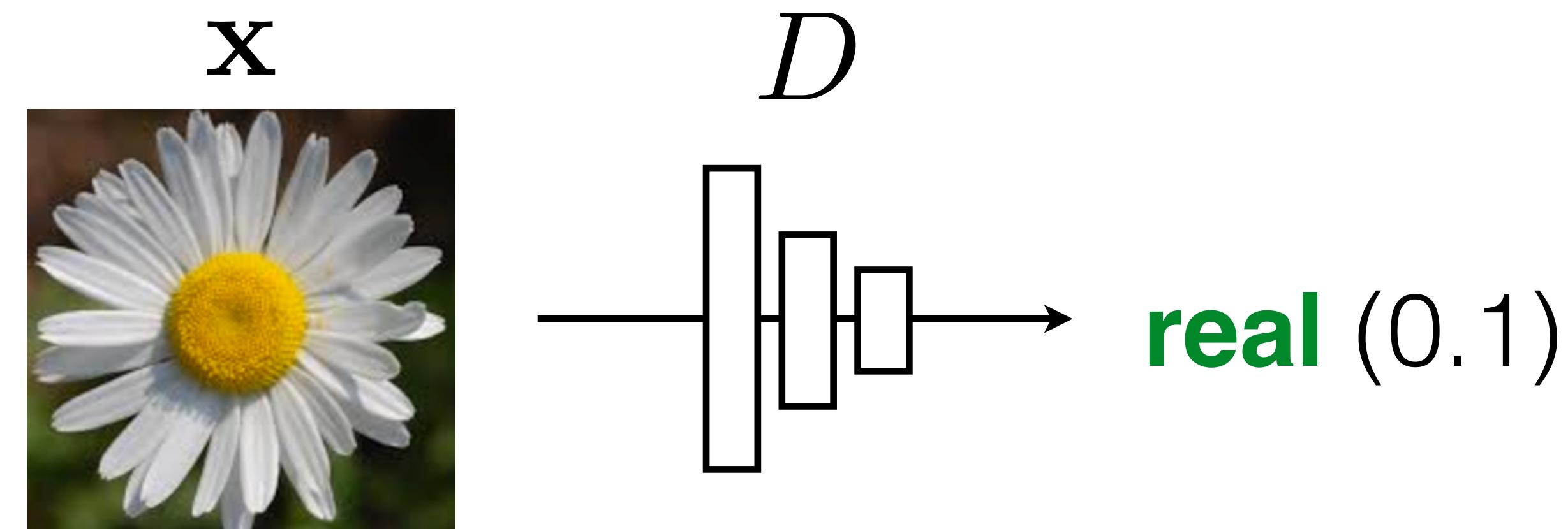
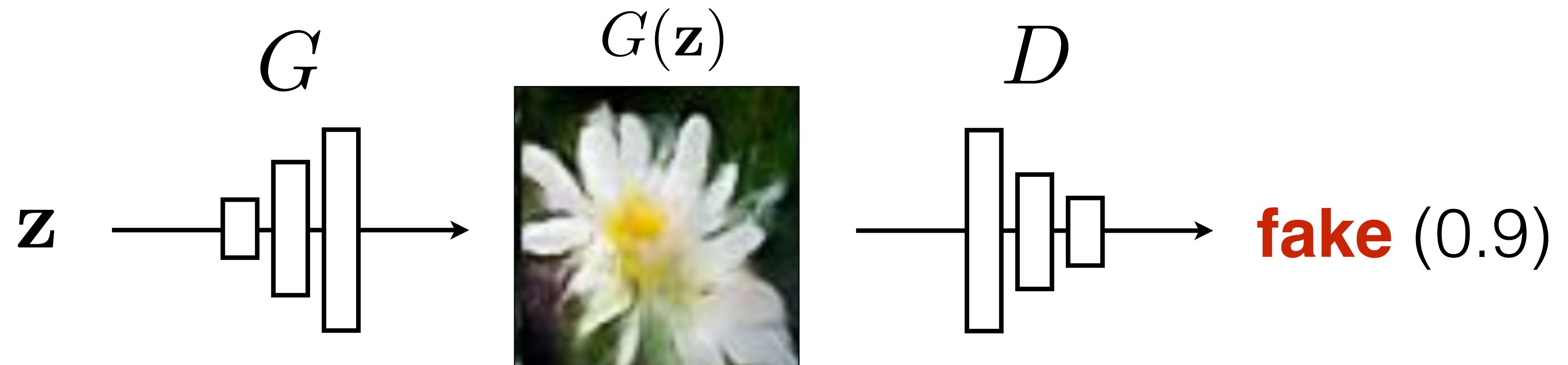




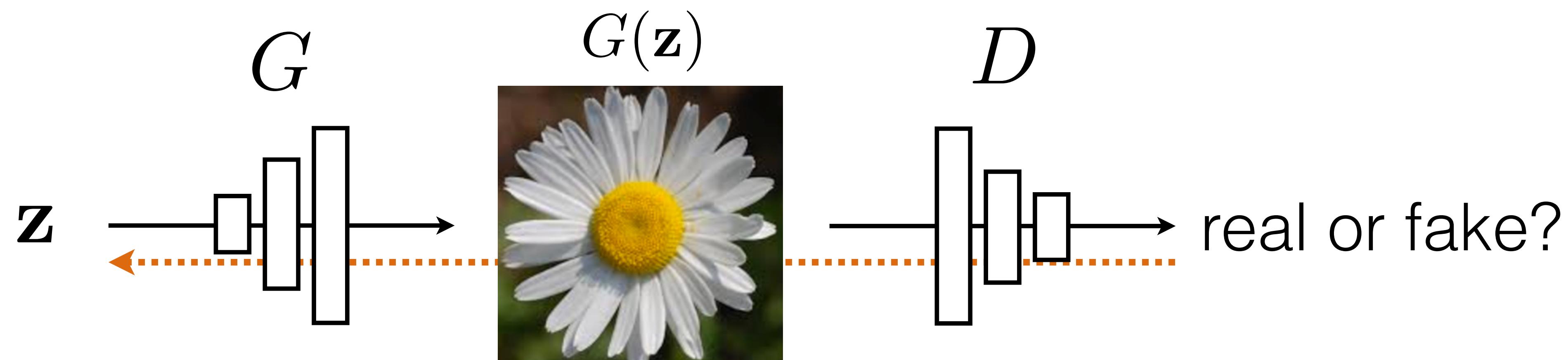
G tries to synthesize fake images that fool **D**

D tries to identify the fakes

[Goodfellow et al., 2014]



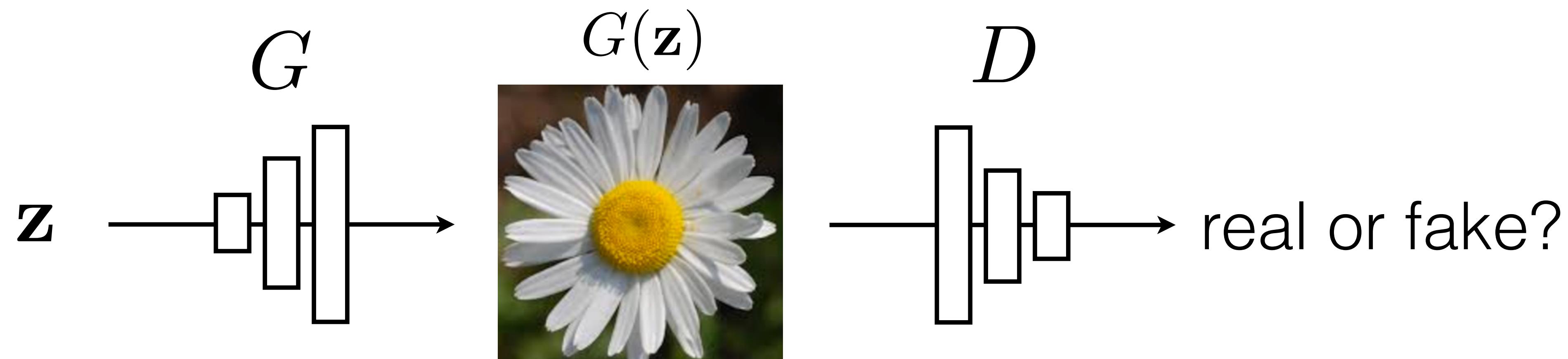
$$\arg \max_D \mathbb{E}_{\mathbf{z}, \mathbf{x}} [\boxed{\log D(G(\mathbf{z}))} + \boxed{\log (1 - D(\mathbf{x}))}]$$



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

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

[Goodfellow et al., 2014]

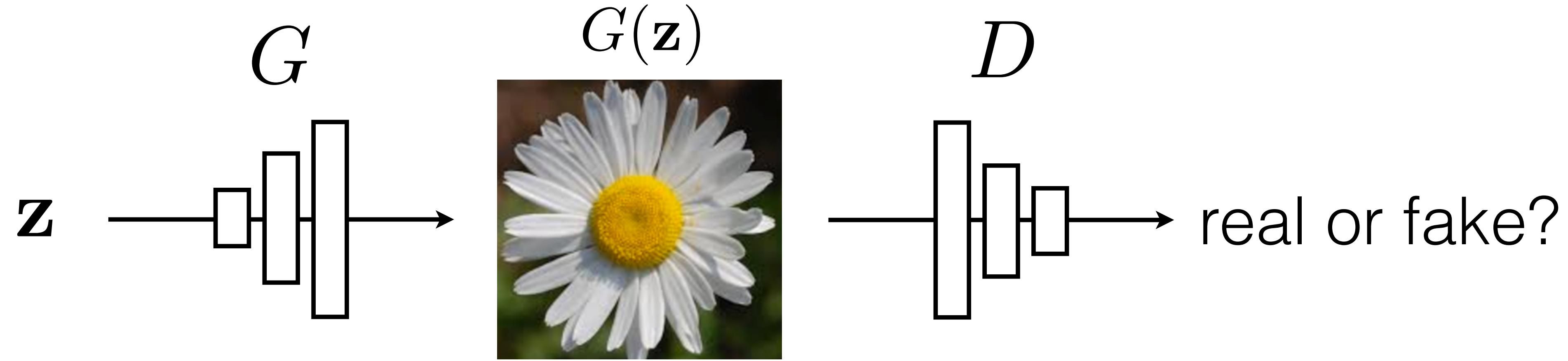


G tries to synthesize fake images that **fool** the **best** **D**:

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

[Goodfellow et al., 2014]

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.

[Goodfellow et al., 2014]

$p_g = p_{data}$ is the unique global minimizer of the GAN objective.

Proof sketch:

Can show this is
optimal D, given G

$$\begin{aligned} C(G) &= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_g} [\log(1 - D_G^*(\mathbf{x}))] \\ &= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \frac{p_{data}(\mathbf{x})}{P_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_g} \left[\log \frac{p_g(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \right] \end{aligned}$$

$$C(G) = -\log(4) + KL \left(p_{data} \middle\| \frac{p_{data} + p_g}{2} \right) + KL \left(p_g \middle\| \frac{p_{data} + p_g}{2} \right)$$

$$\begin{aligned} C(G) &= -\log(4) + 2 \cdot JSD(p_{data} \| p_g) \\ &\underbrace{\qquad\qquad\qquad}_{\geq 0, \quad 0} \stackrel{39}{\iff} p_g = p_{data} \quad \square \end{aligned}$$

Samples from BigGAN

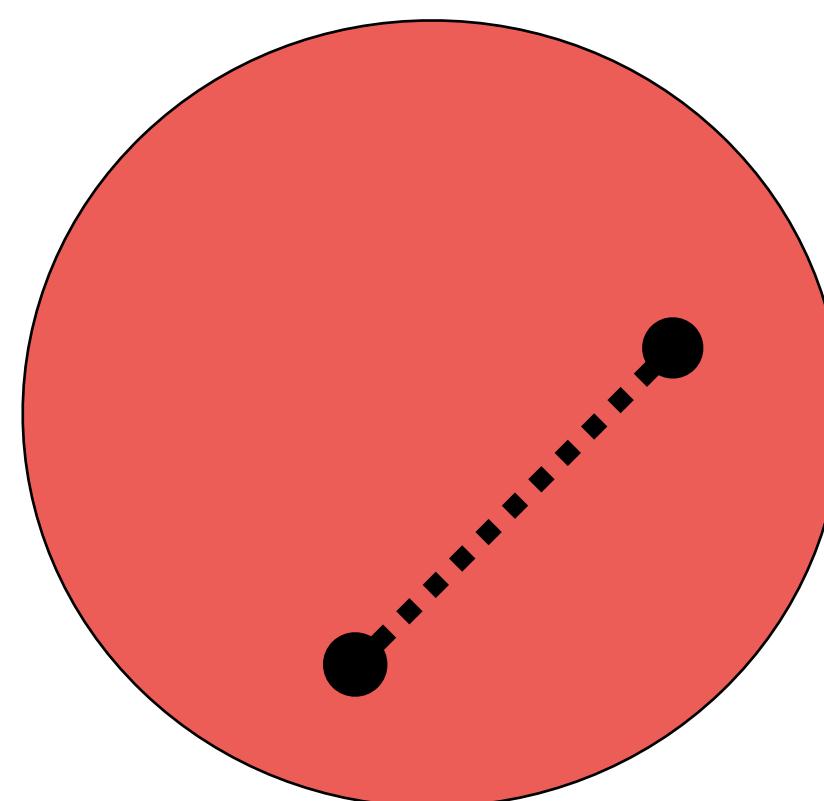
[Brock et al. 2018]



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

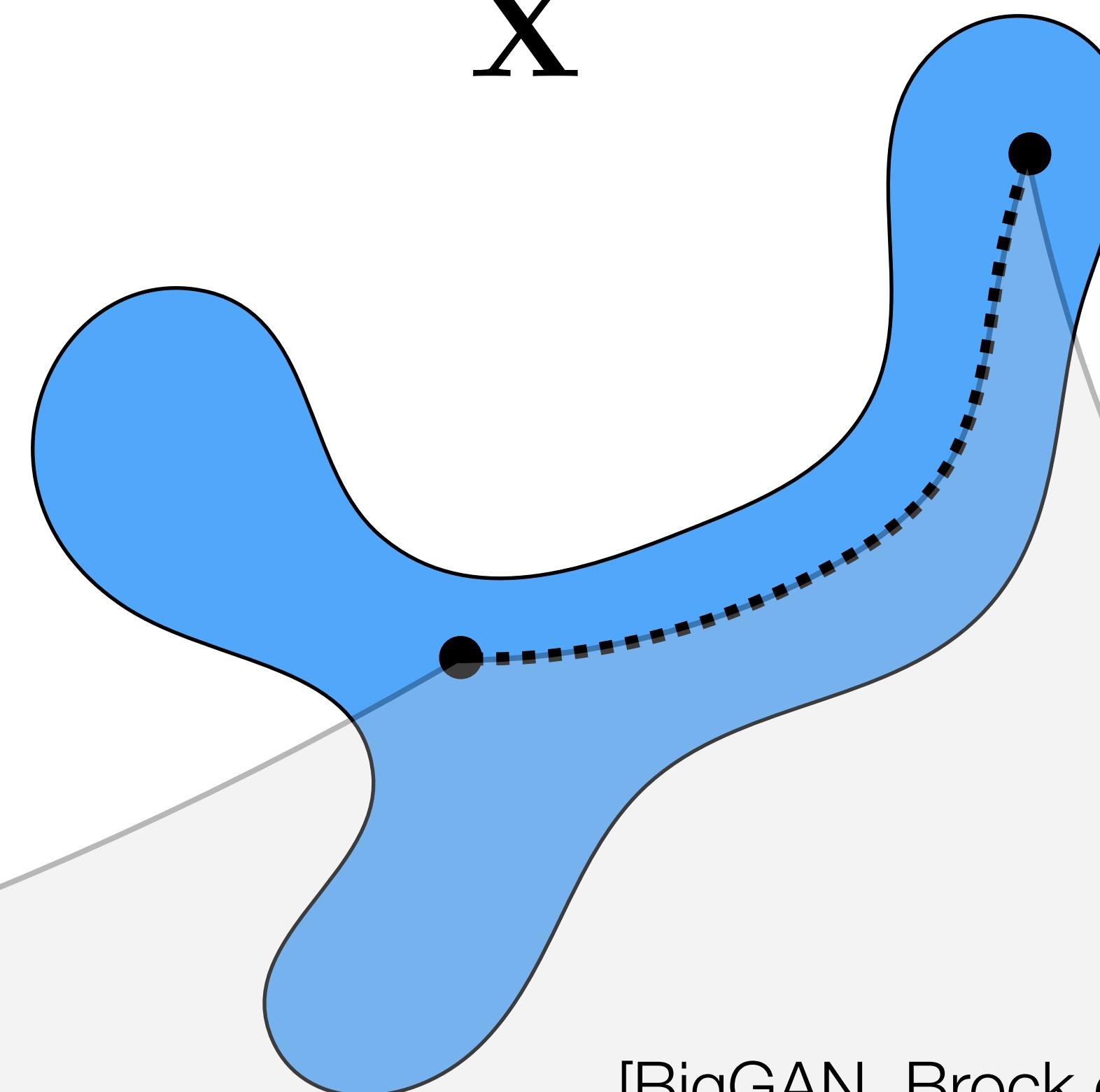
Latent space
(Gaussian)

z



Data space
(Natural image manifold)

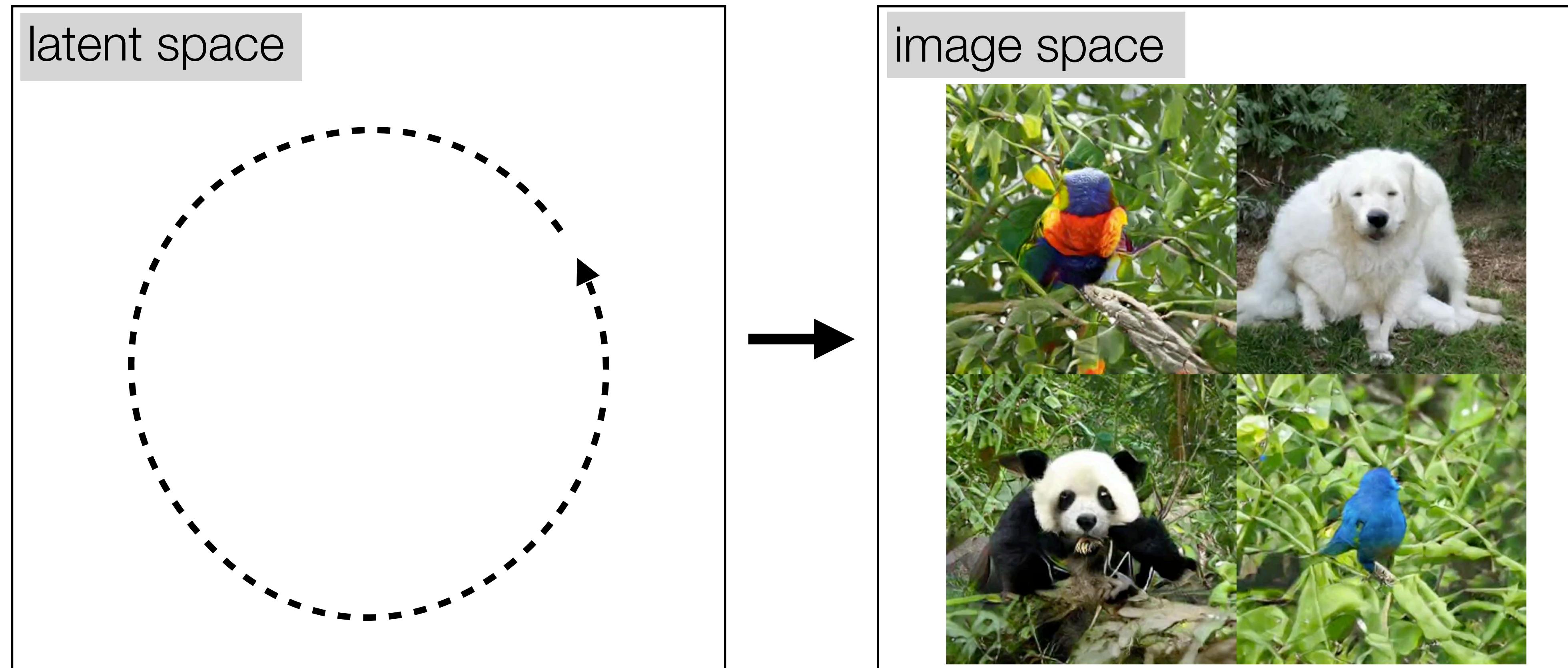
X



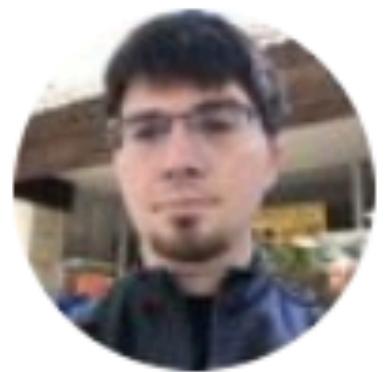
[BigGAN, Brock et al. 2018]



Generative models organize the manifold of natural images



What has driven GAN progress?



Ian Goodfellow @goodfellow_ian · Jan 14

▼

4.5 years of **GAN progress** on face generation. arxiv.org/abs/1406.2661

arxiv.org/abs/1511.06434 arxiv.org/abs/1606.07536 arxiv.org/abs/1710.10196

arxiv.org/abs/1812.04948



Better objectives? Optimization?

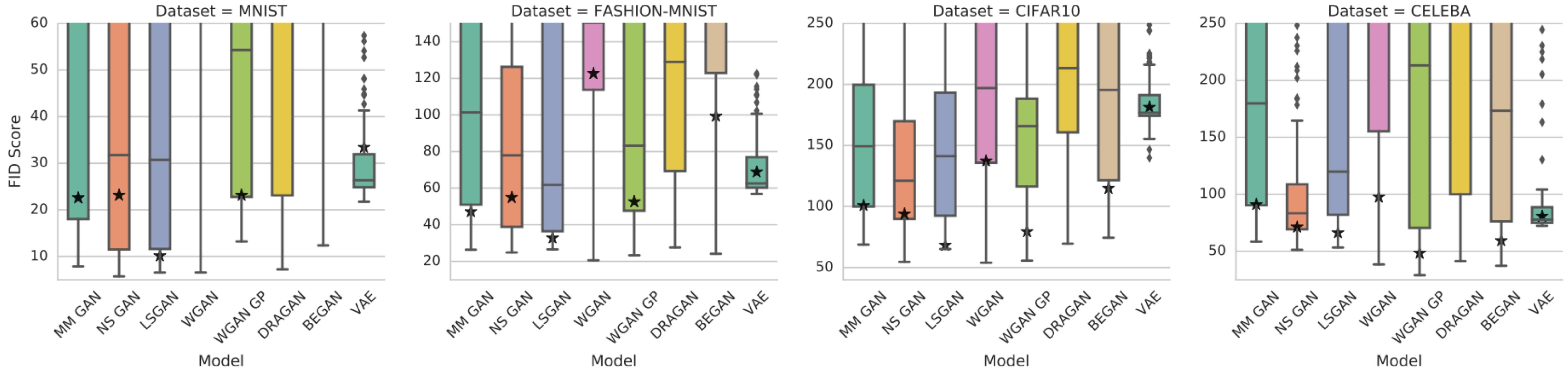
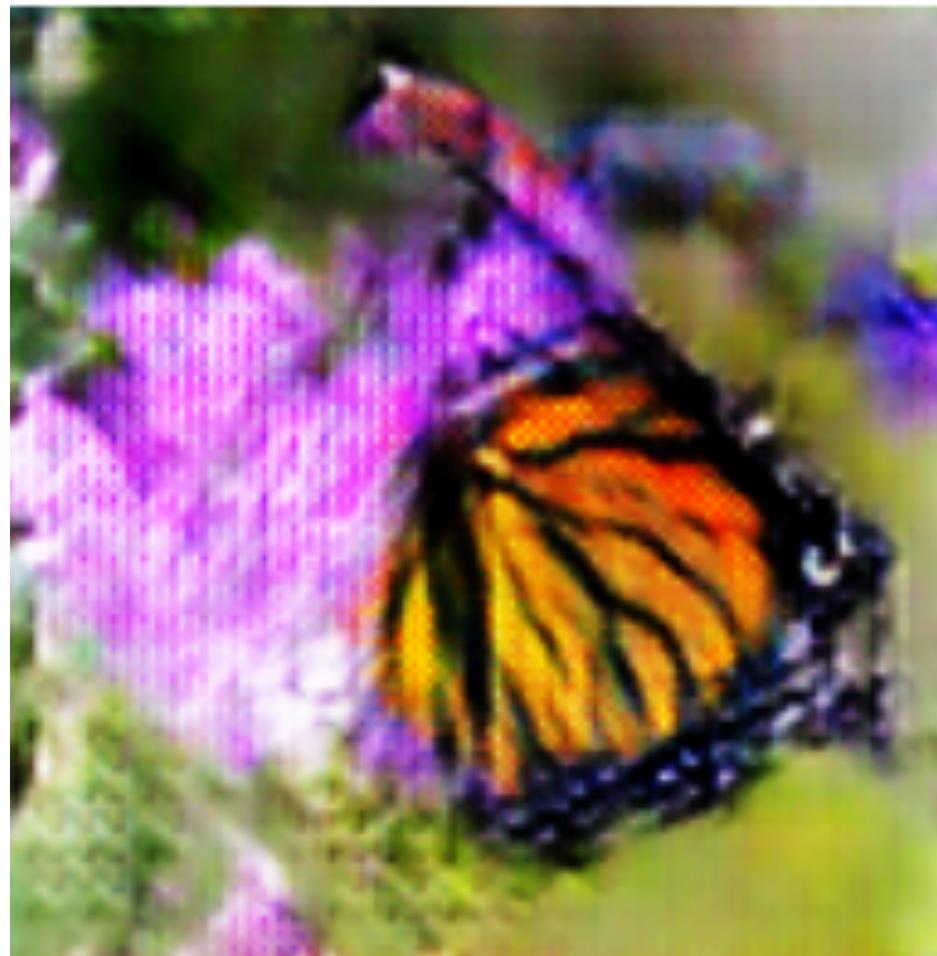


Figure 4: A *wide range* hyperparameter search (100 hyperparameter samples per model). Black stars indicate the performance of suggested hyperparameter settings. We observe that GAN training is extremely sensitive to hyperparameter settings and there is no model which is significantly more stable than others.

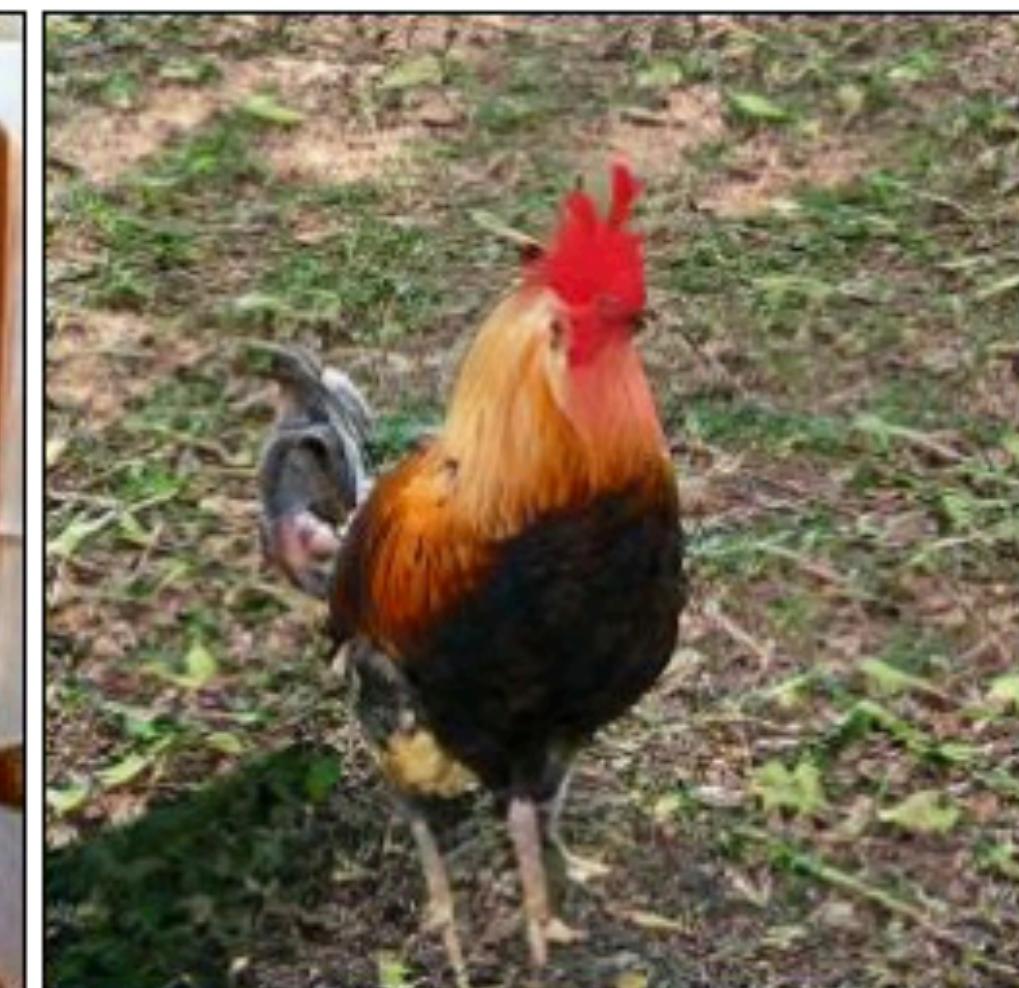
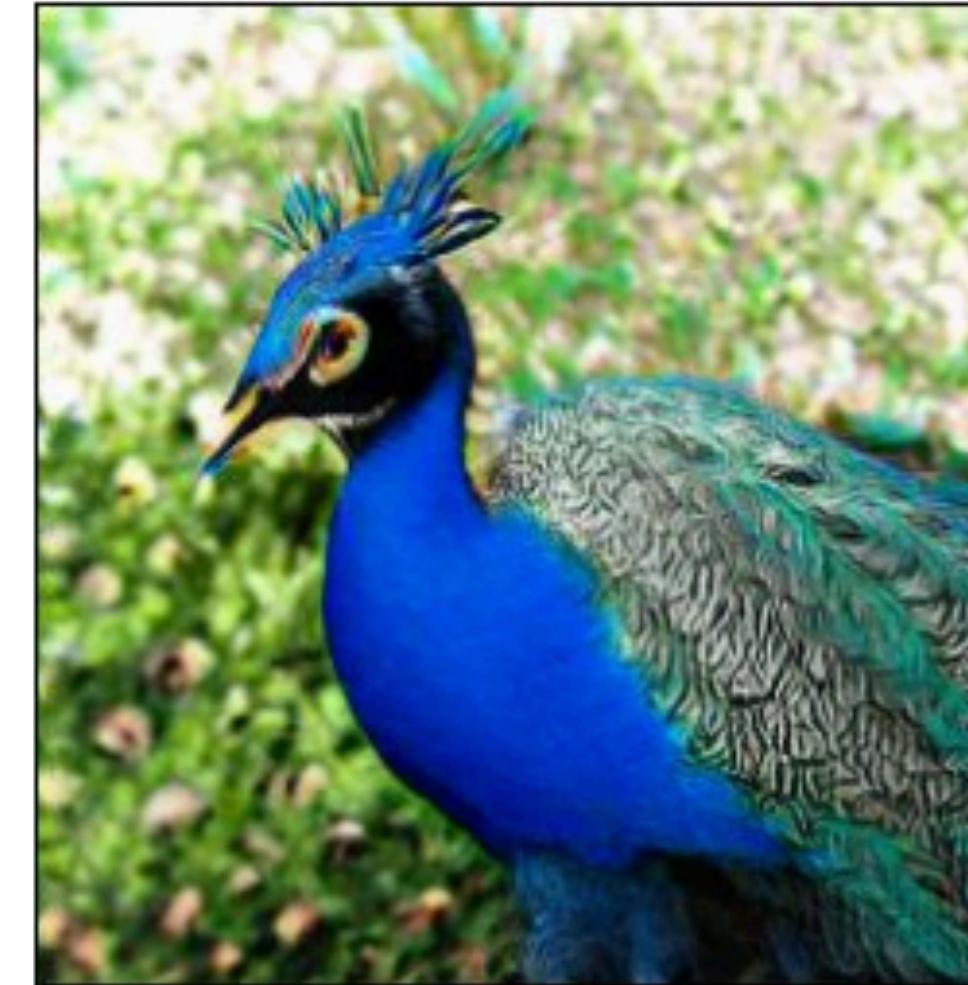
[“Are all GANs Created Equal⁴⁴?”, Lucic*, Kurach*, et al. 2018]

More data?

ACGAN [Odena et al. 2016]



BigGAN [Brock et al. 2018]



Both trained on Imagenet

Architectures

DCGAN

[Radford, Metz, Chintala 2016]



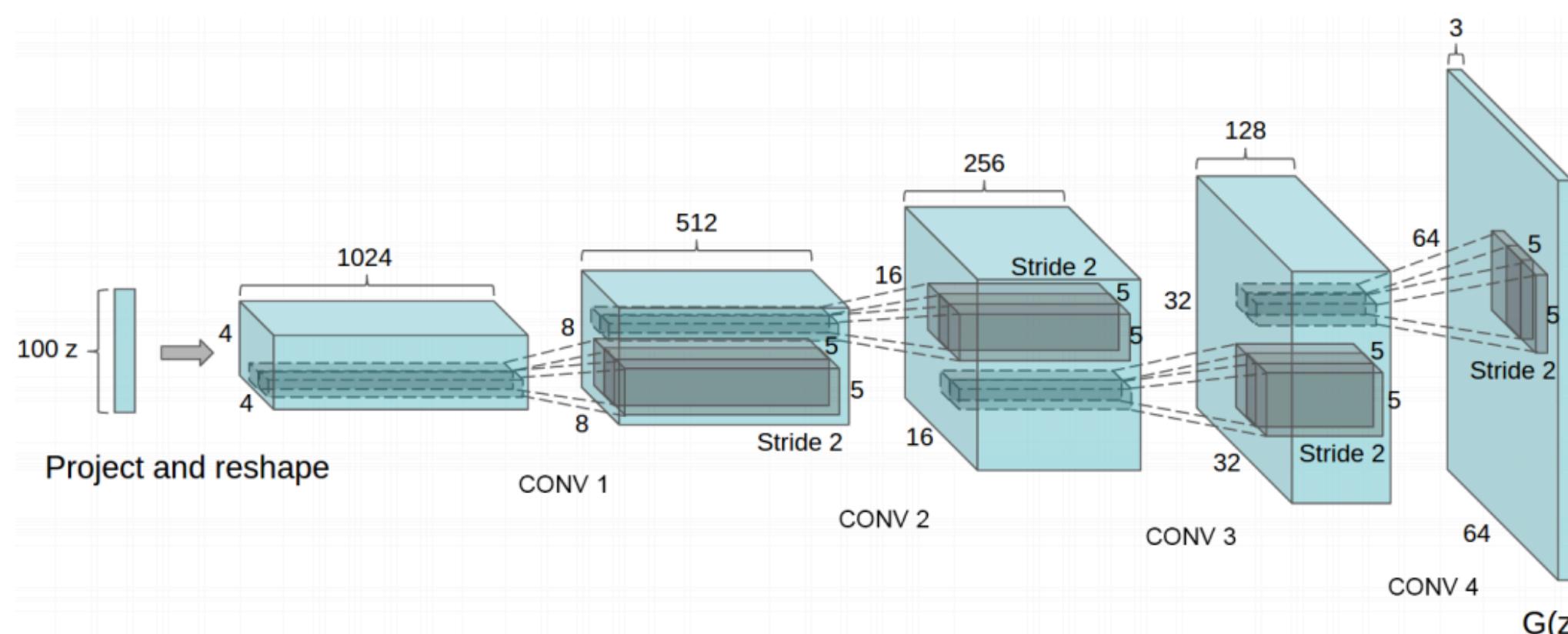
StyleGAN

[Karras, Laine, Aila 2019]



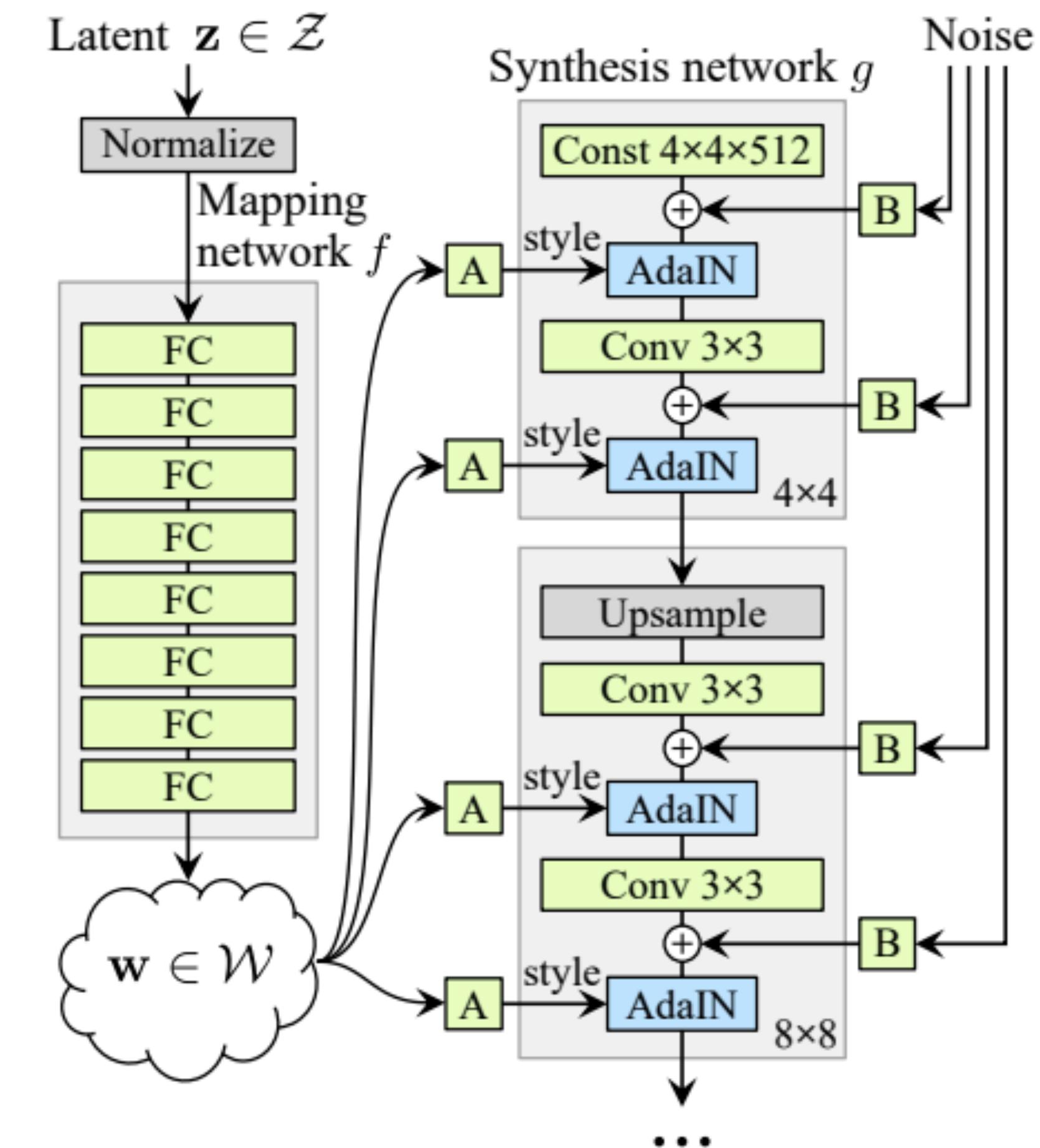
Architectures

DCGAN
[Radford, Metz, Chintala 2016]



47

StyleGAN
[Karras, Laine, Aila 2019]



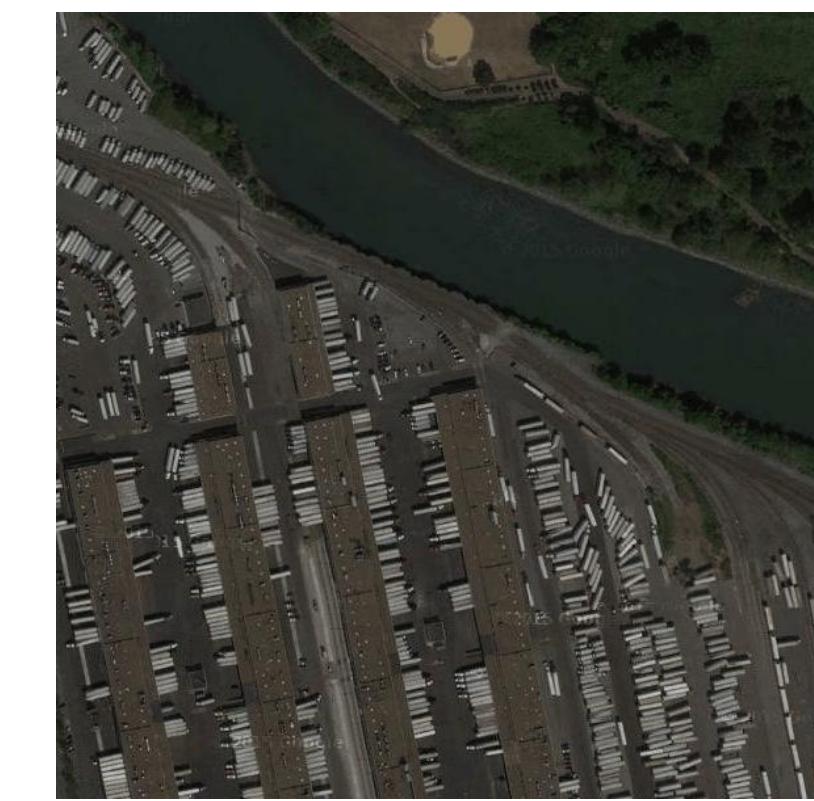
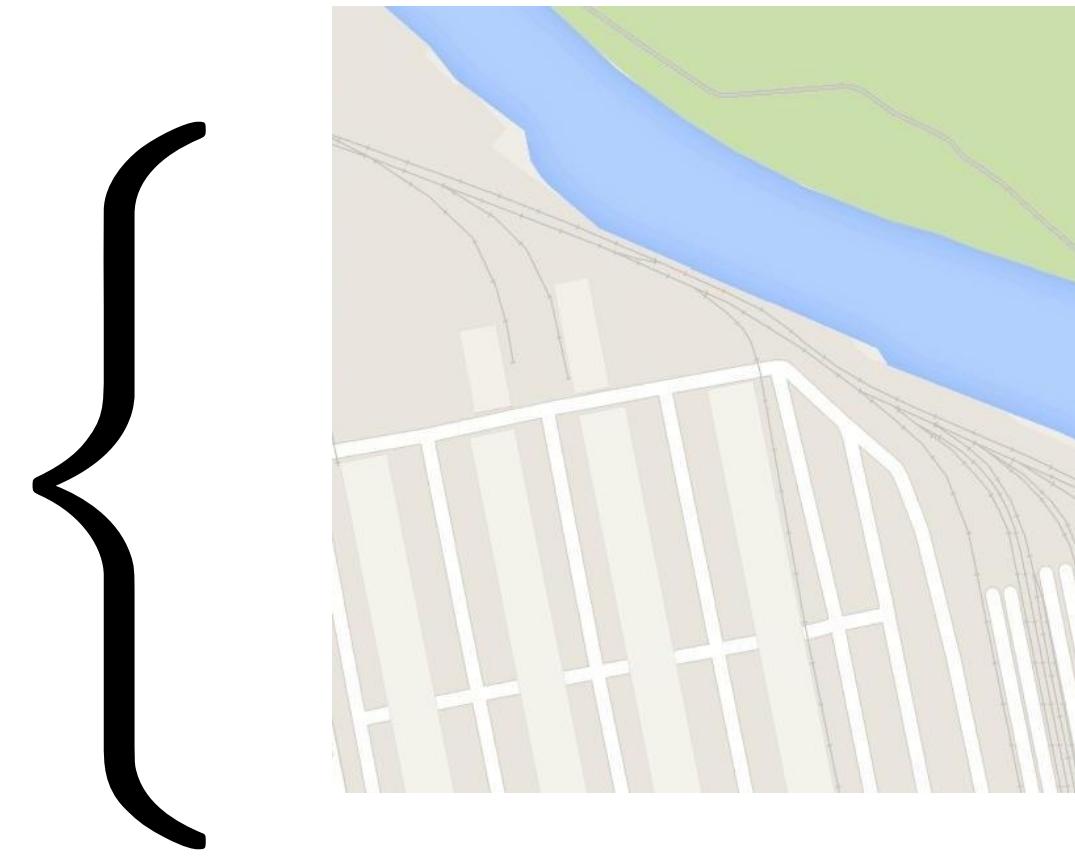
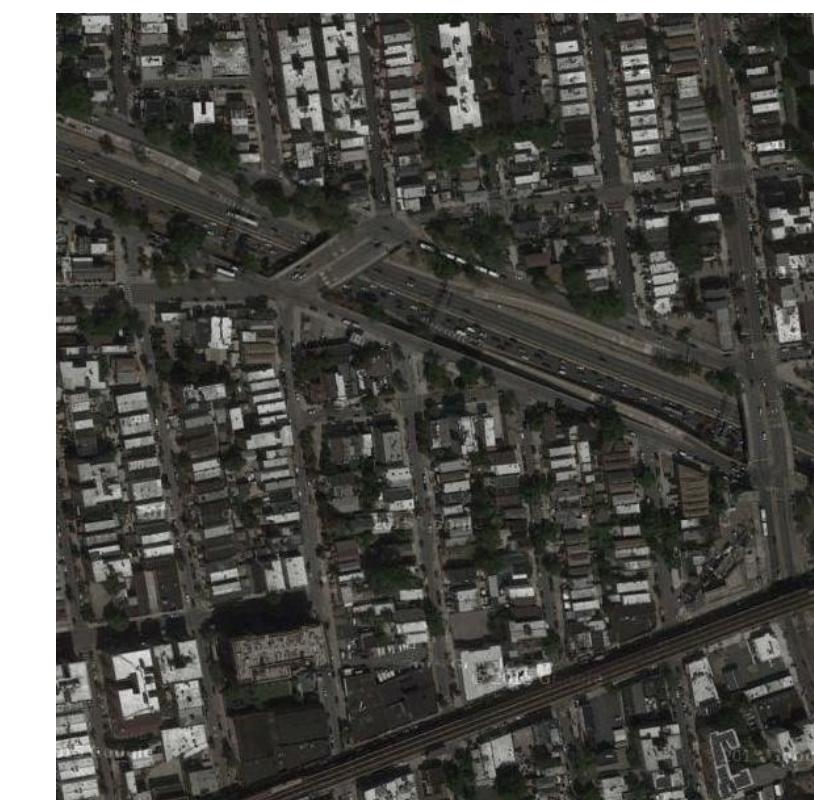
Source: Isola, Freeman, Torralba

Map2Sat

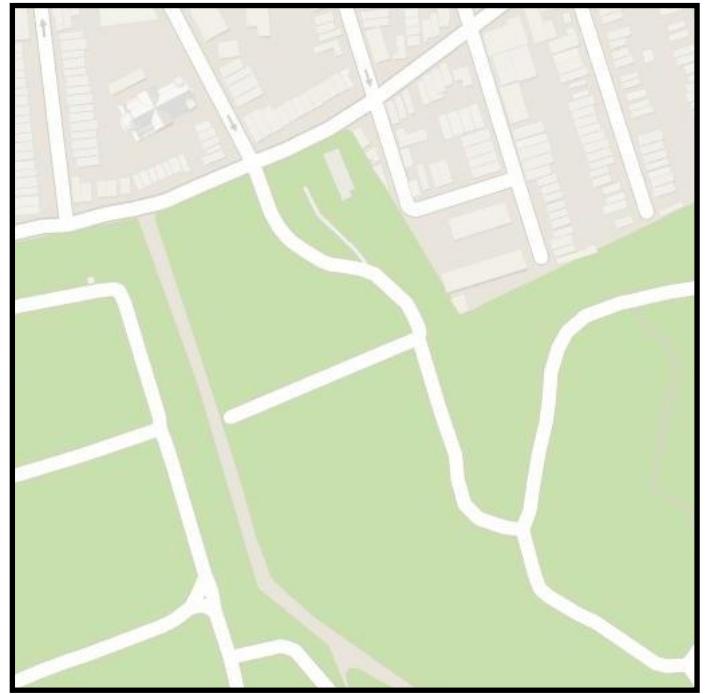
x



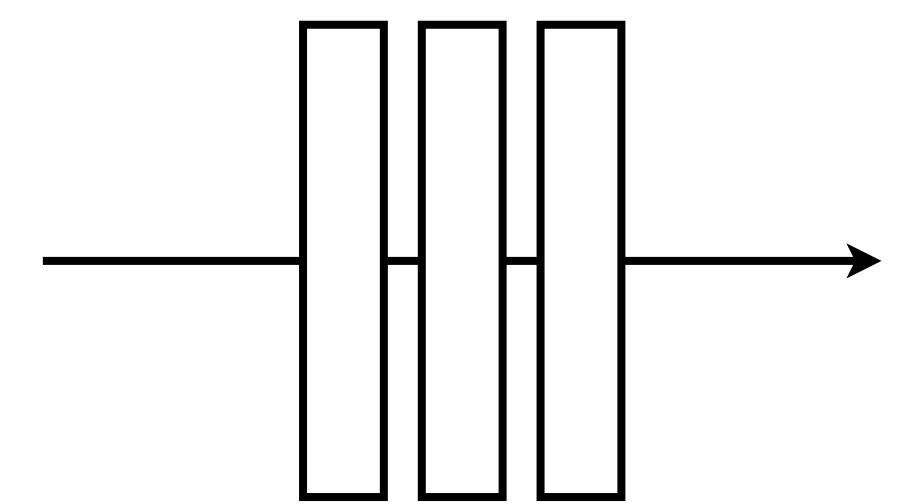
y



x



G



Generator

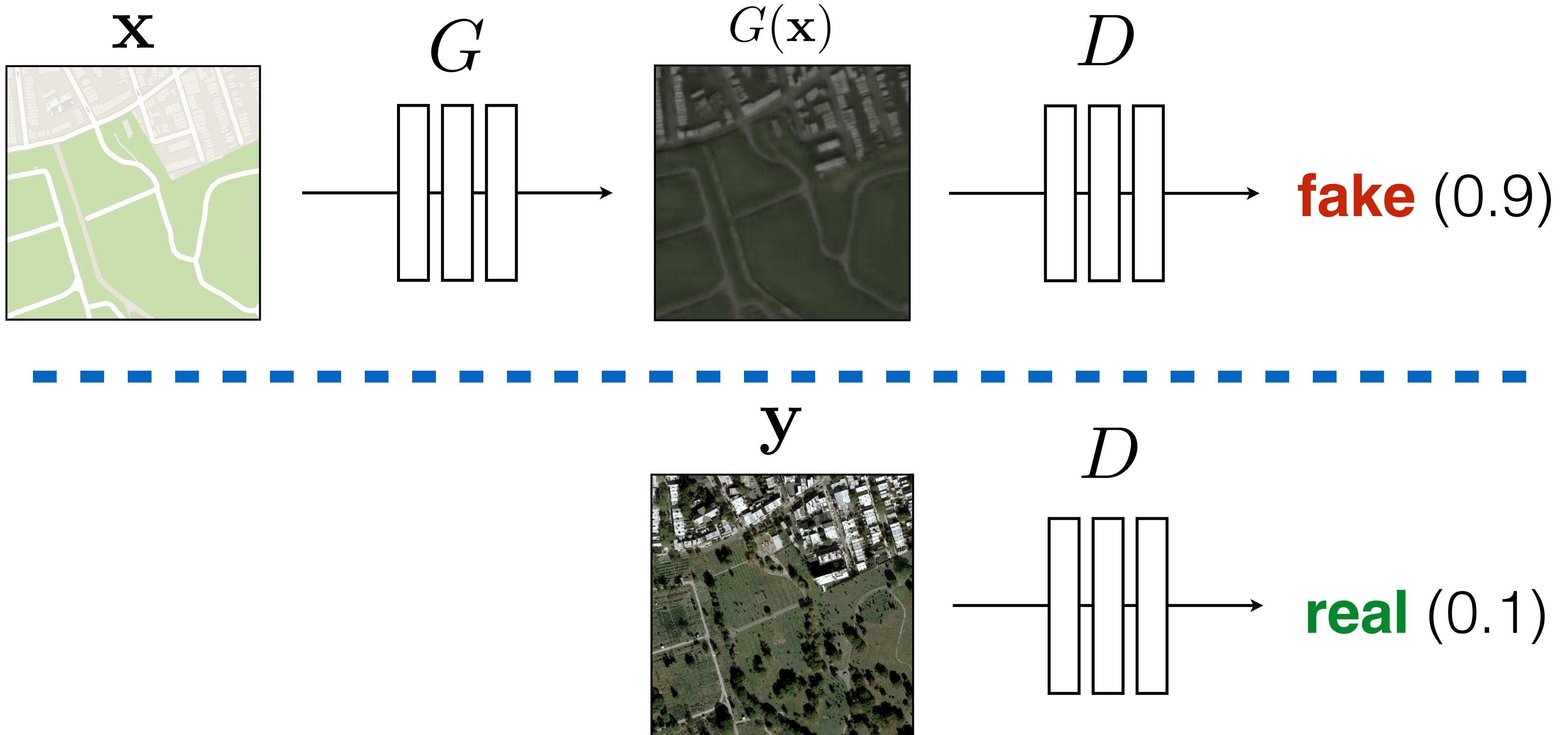
$G(x)$



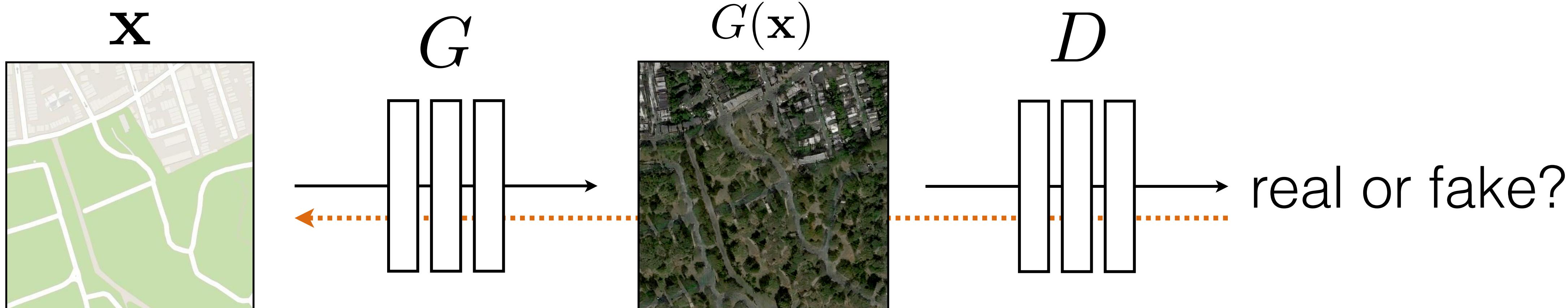


G tries to synthesize fake images that fool **D**

D tries to identify the fakes

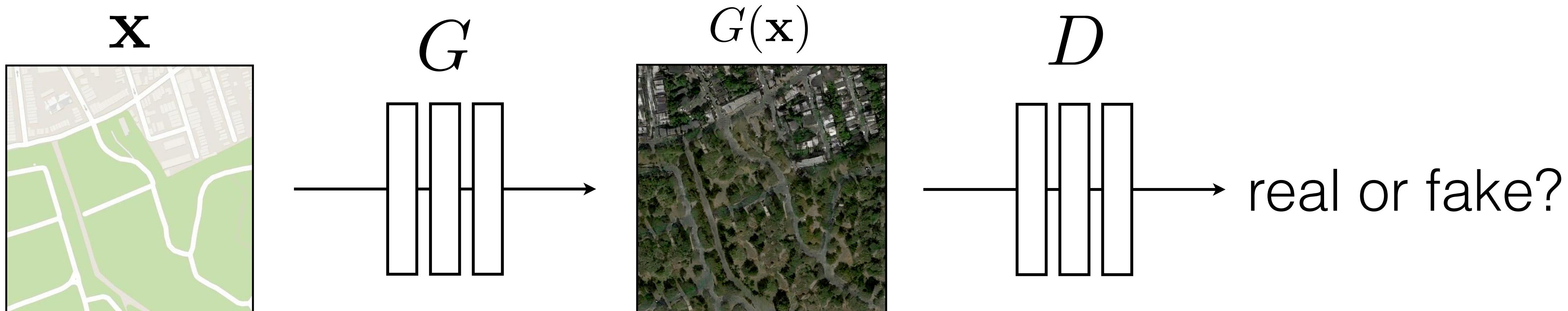


$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\boxed{\log D(G(\mathbf{x}))} + \boxed{\log(1 - D(\mathbf{y}))}]$$



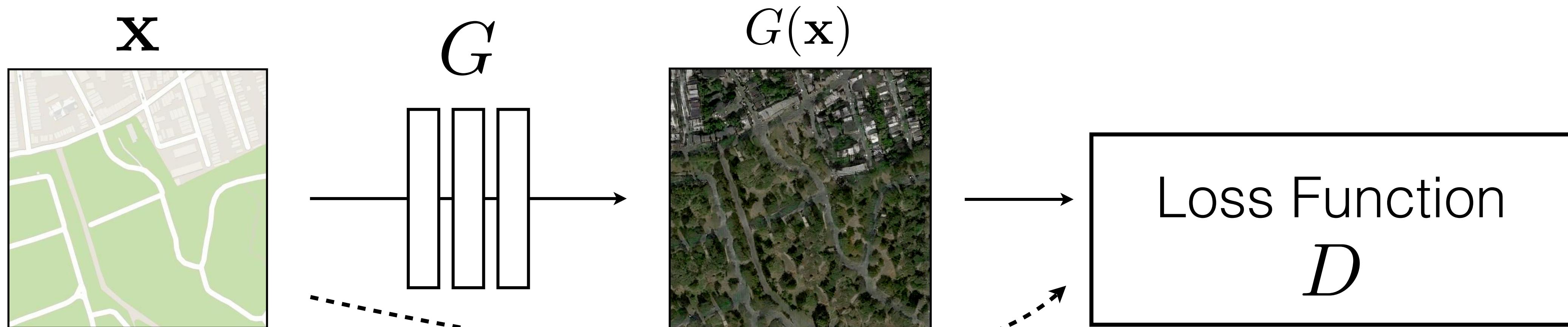
G tries to synthesize fake images that **fool** **D**:

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



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

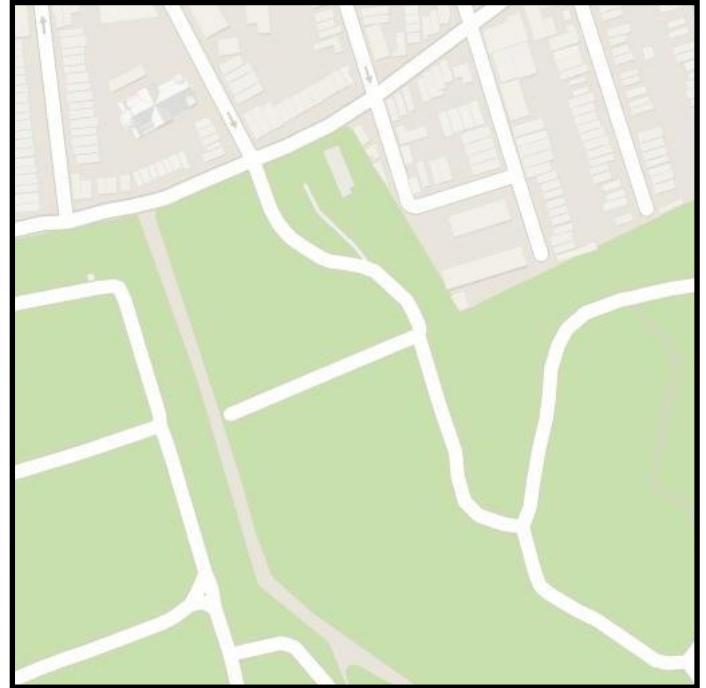
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$



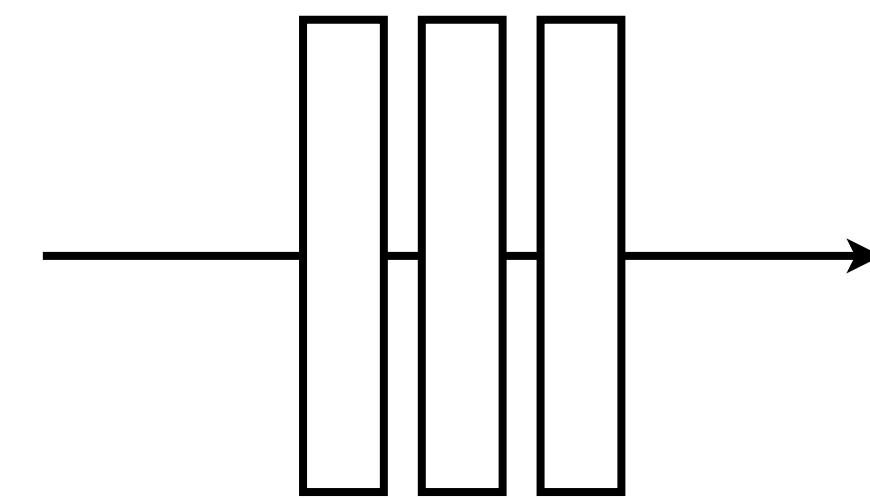
G's perspective: **D** is a loss function.

Rather than being hand-designed, it is *learned* and *highly structured*.

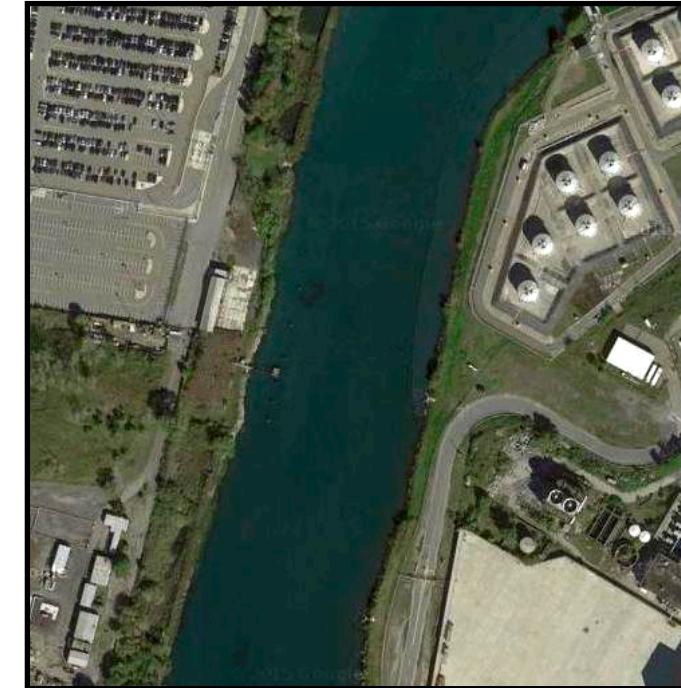
x



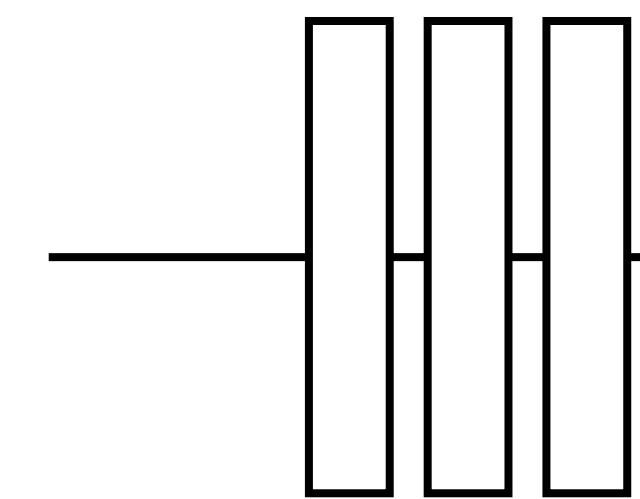
G



G(x)



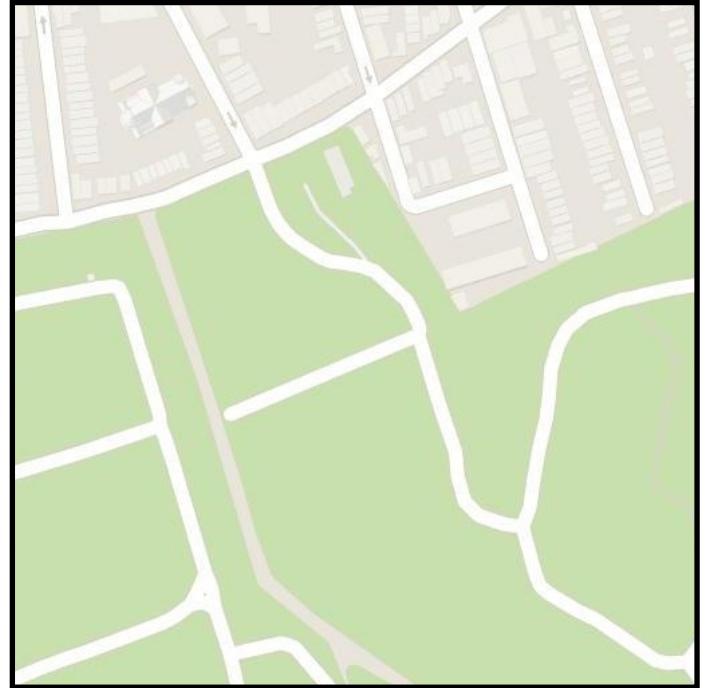
D



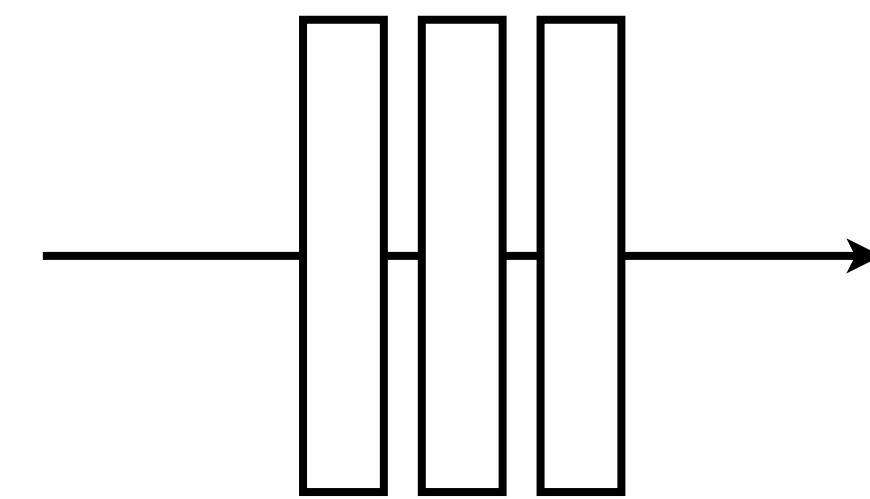
real or fake?

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

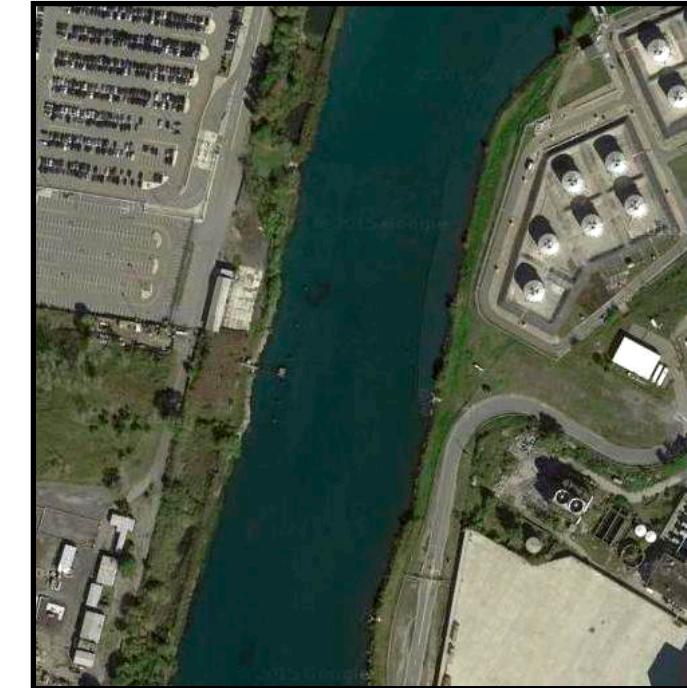
x



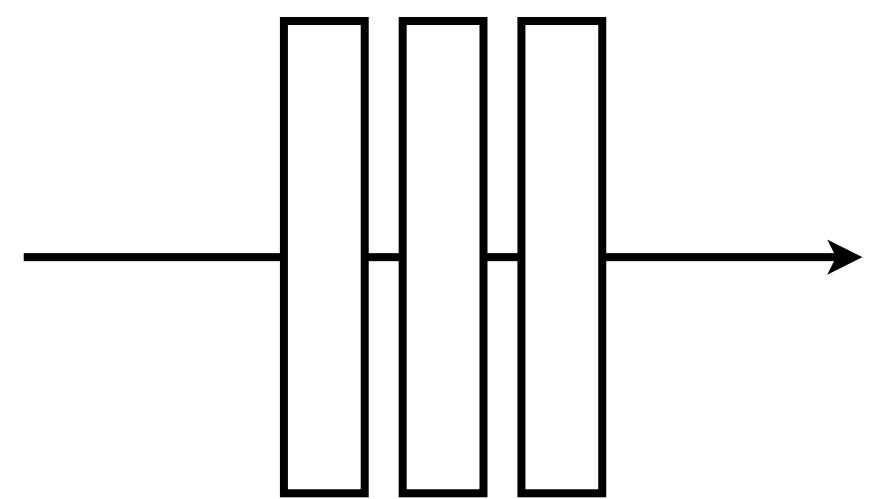
G



G(x)

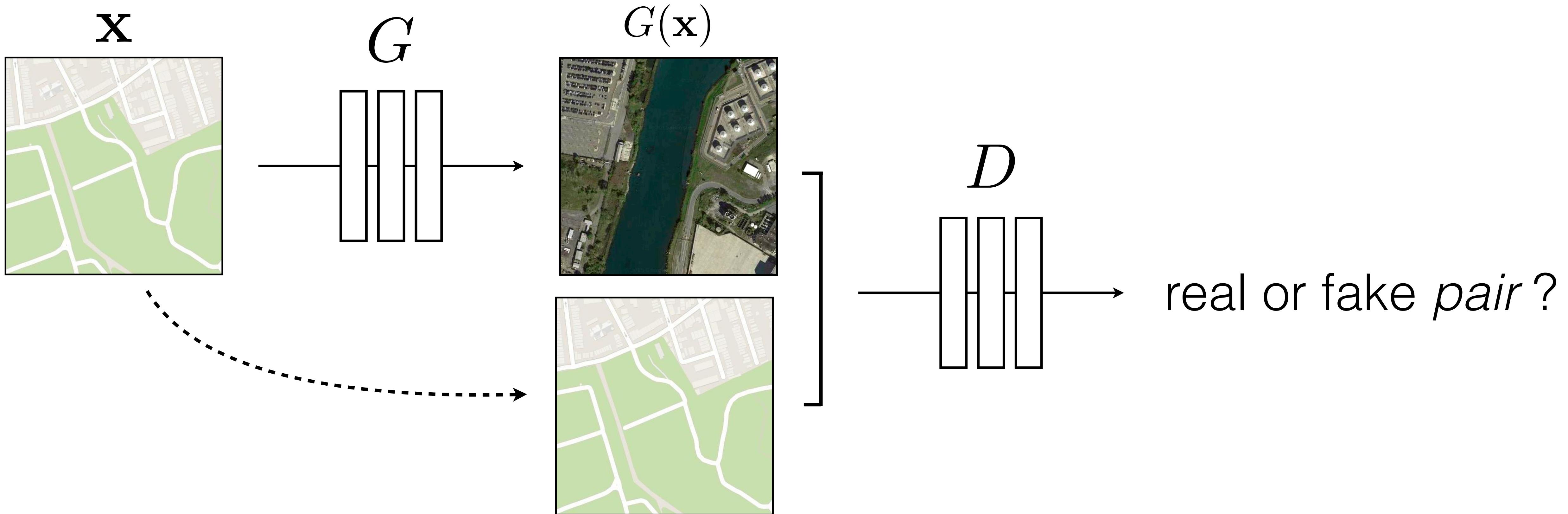


D

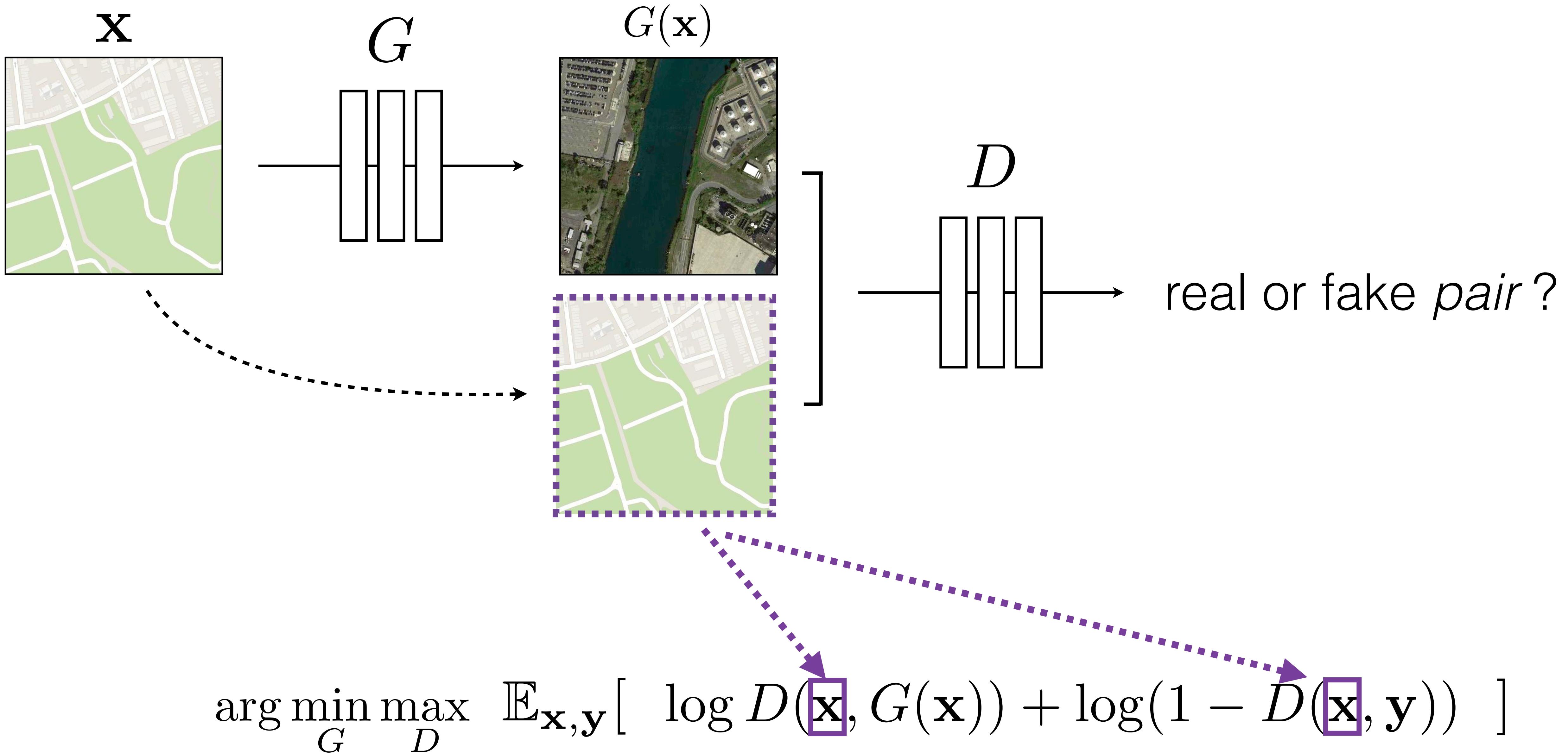


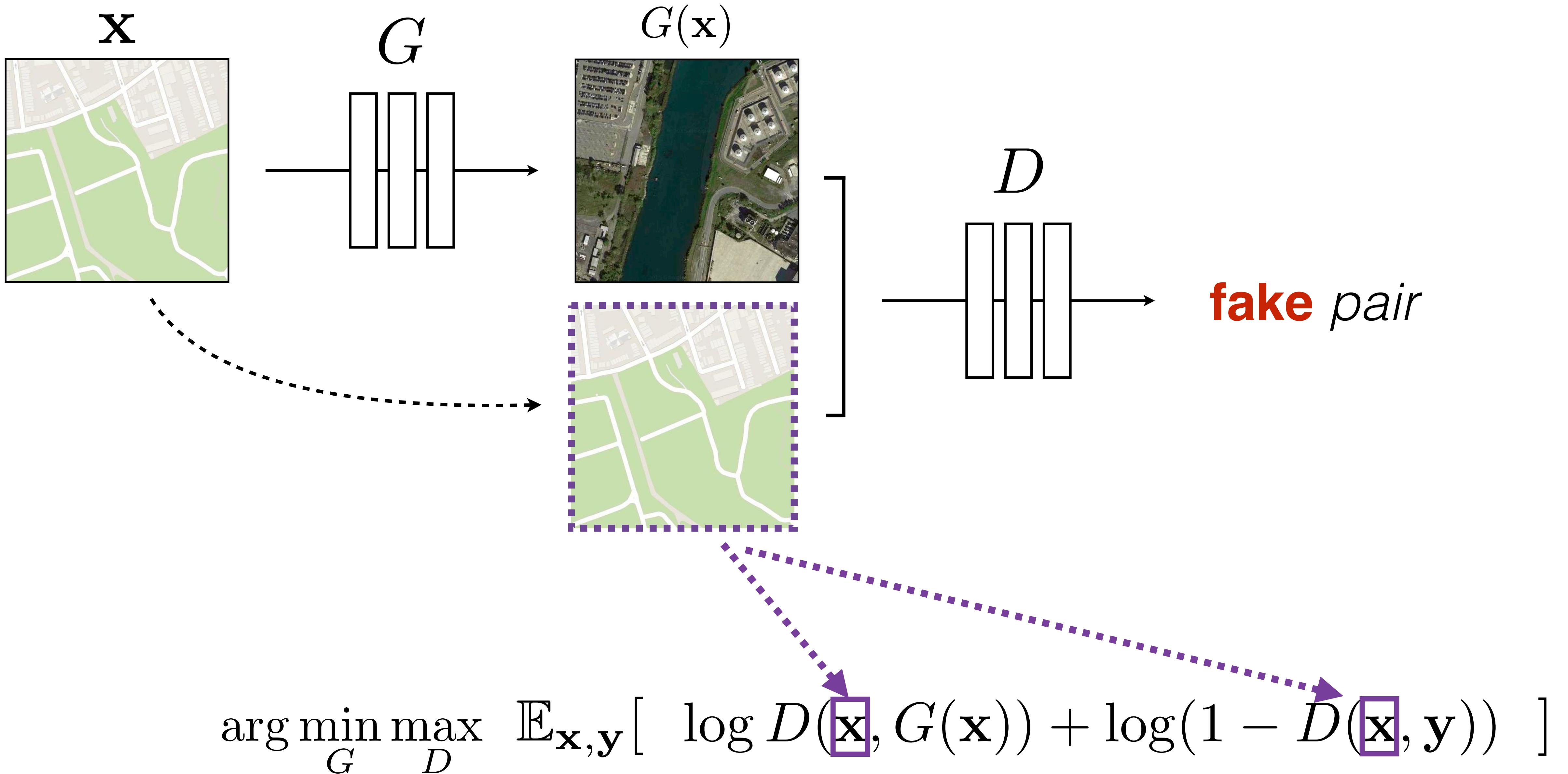
real!

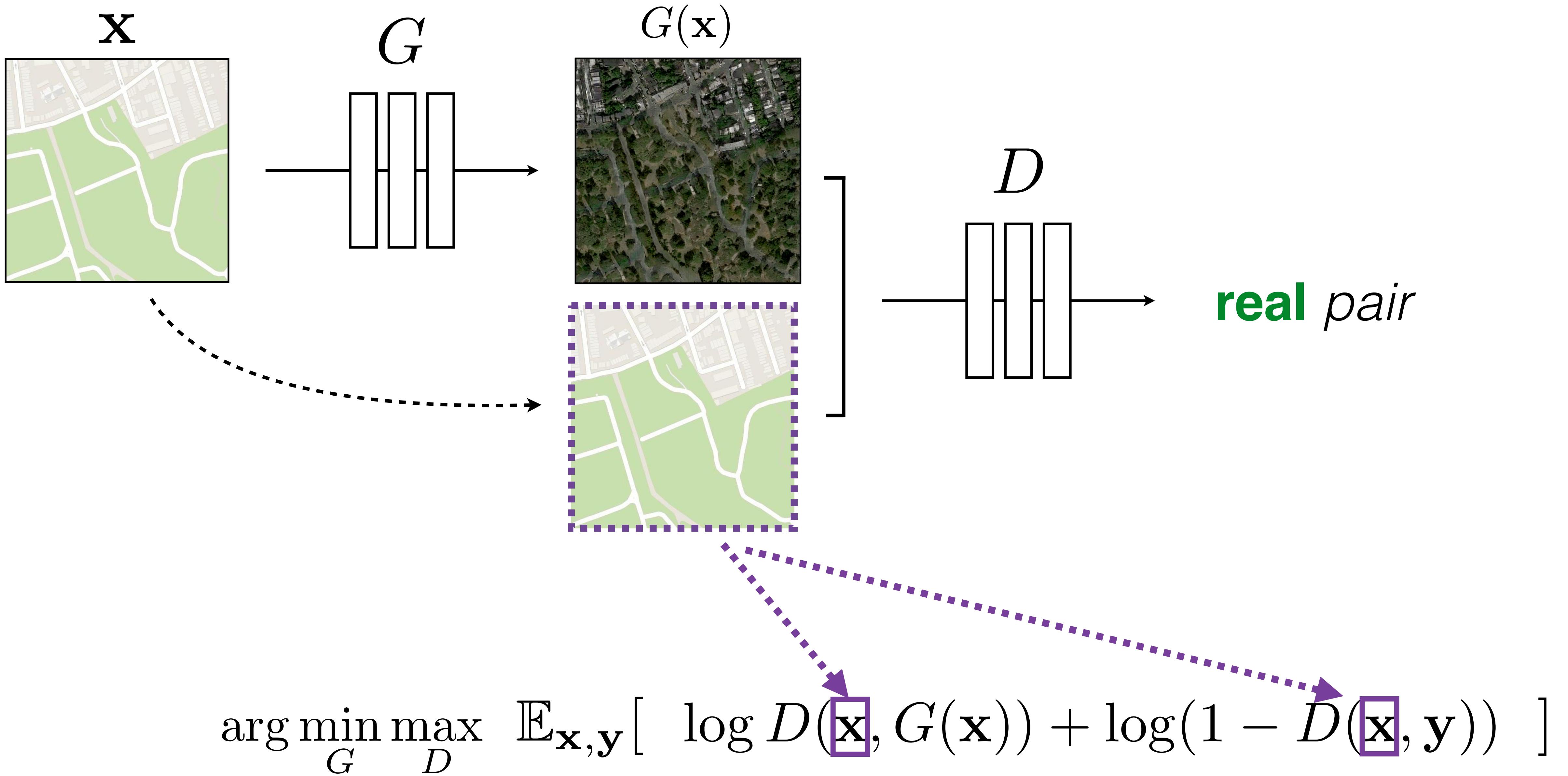
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$

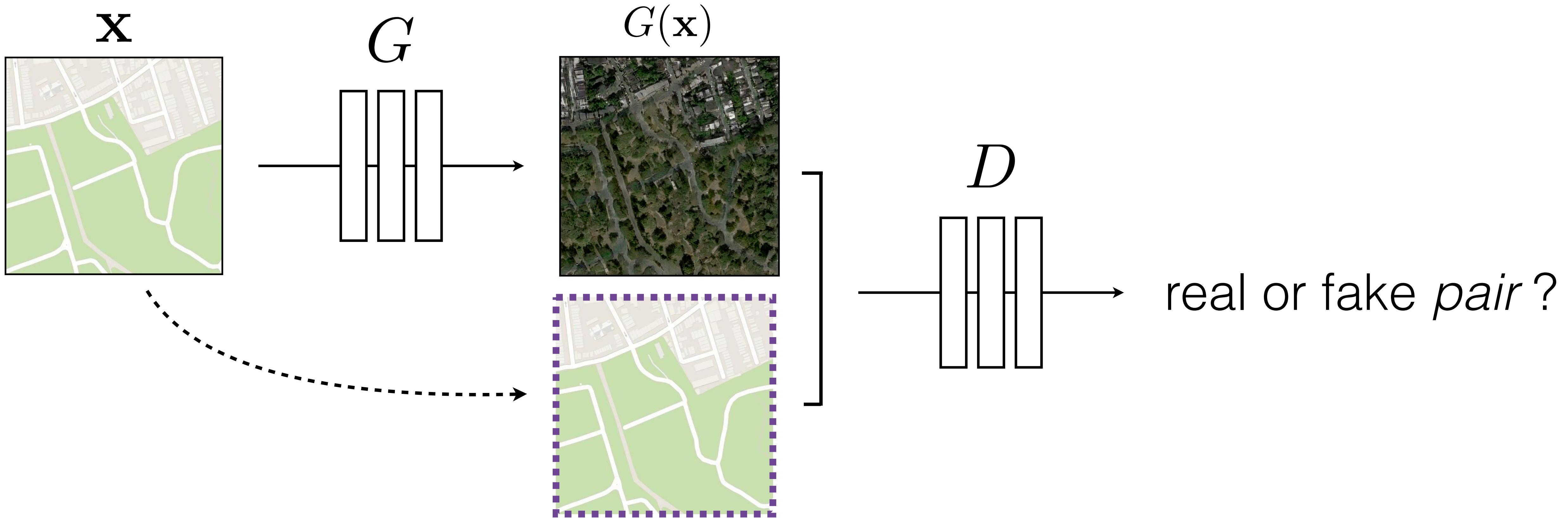


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$









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

Training Details: Loss function

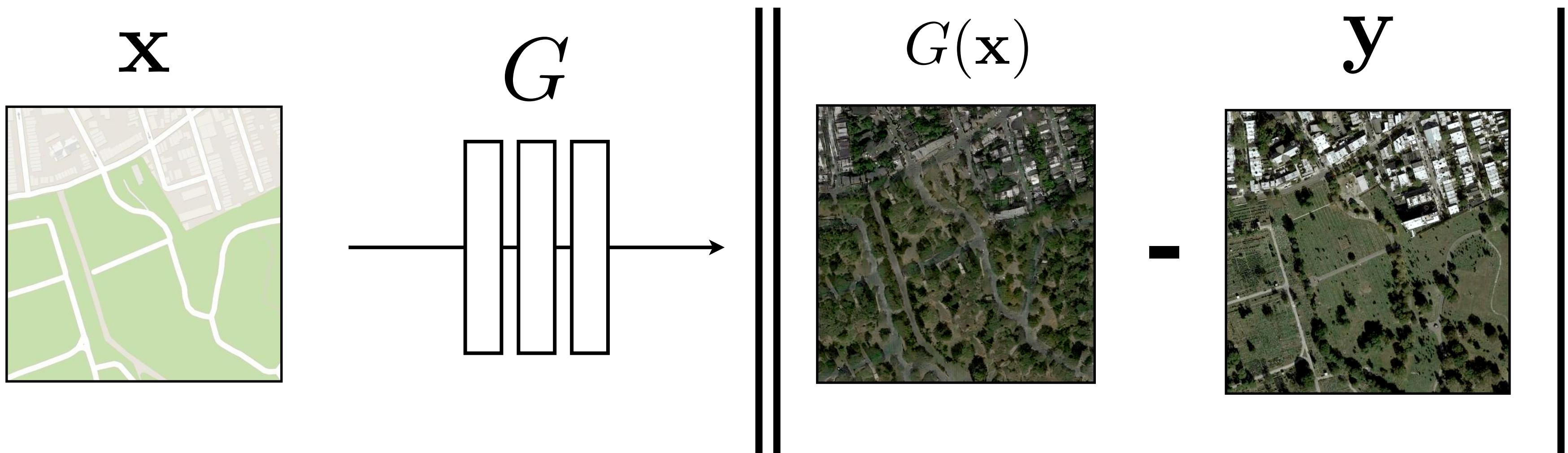
Conditional GAN

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

Training Details: Loss function

Conditional GAN

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



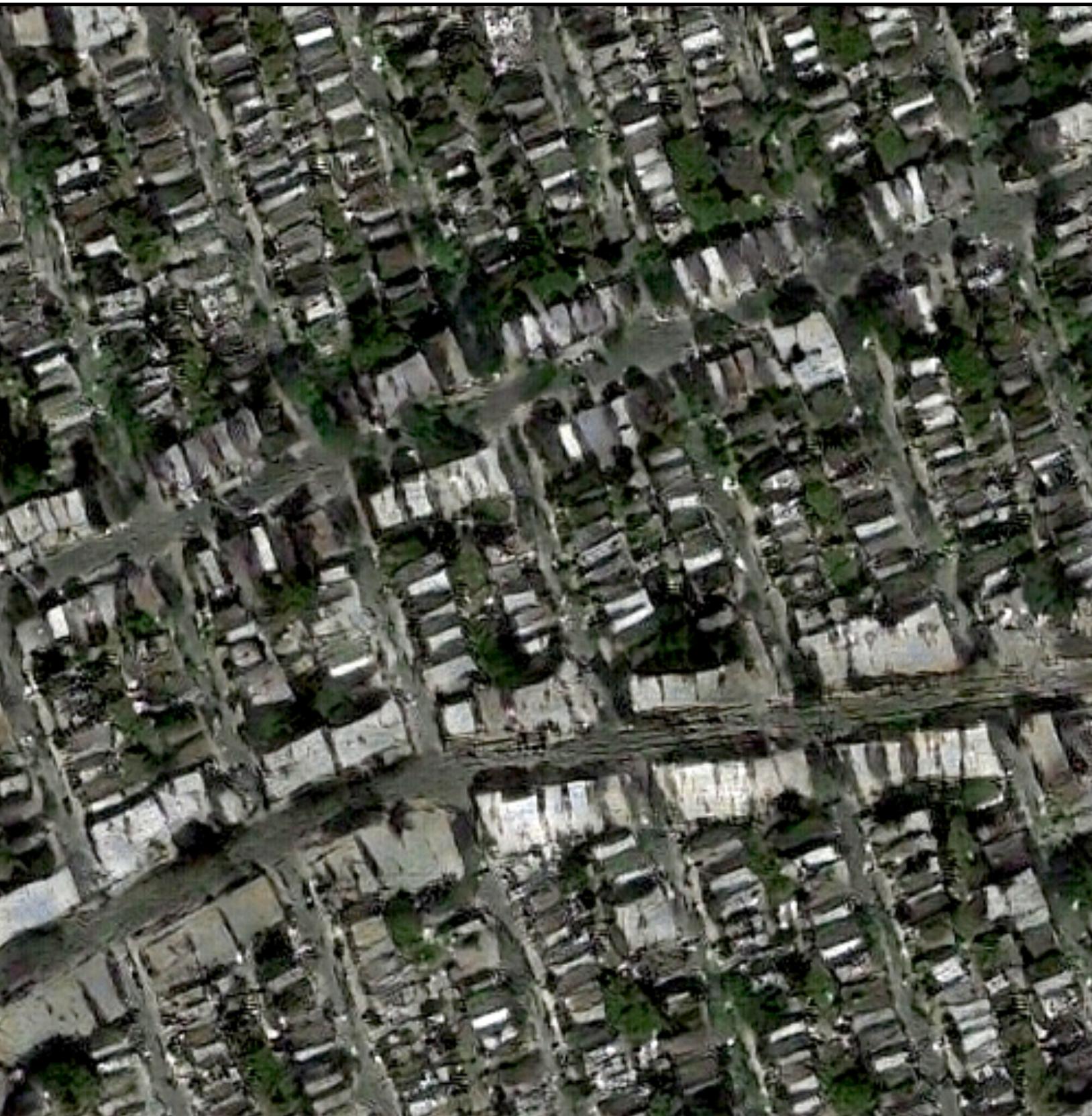
Stable training + fast convergence

[c.f. Pathak et al. CVPR 2016]

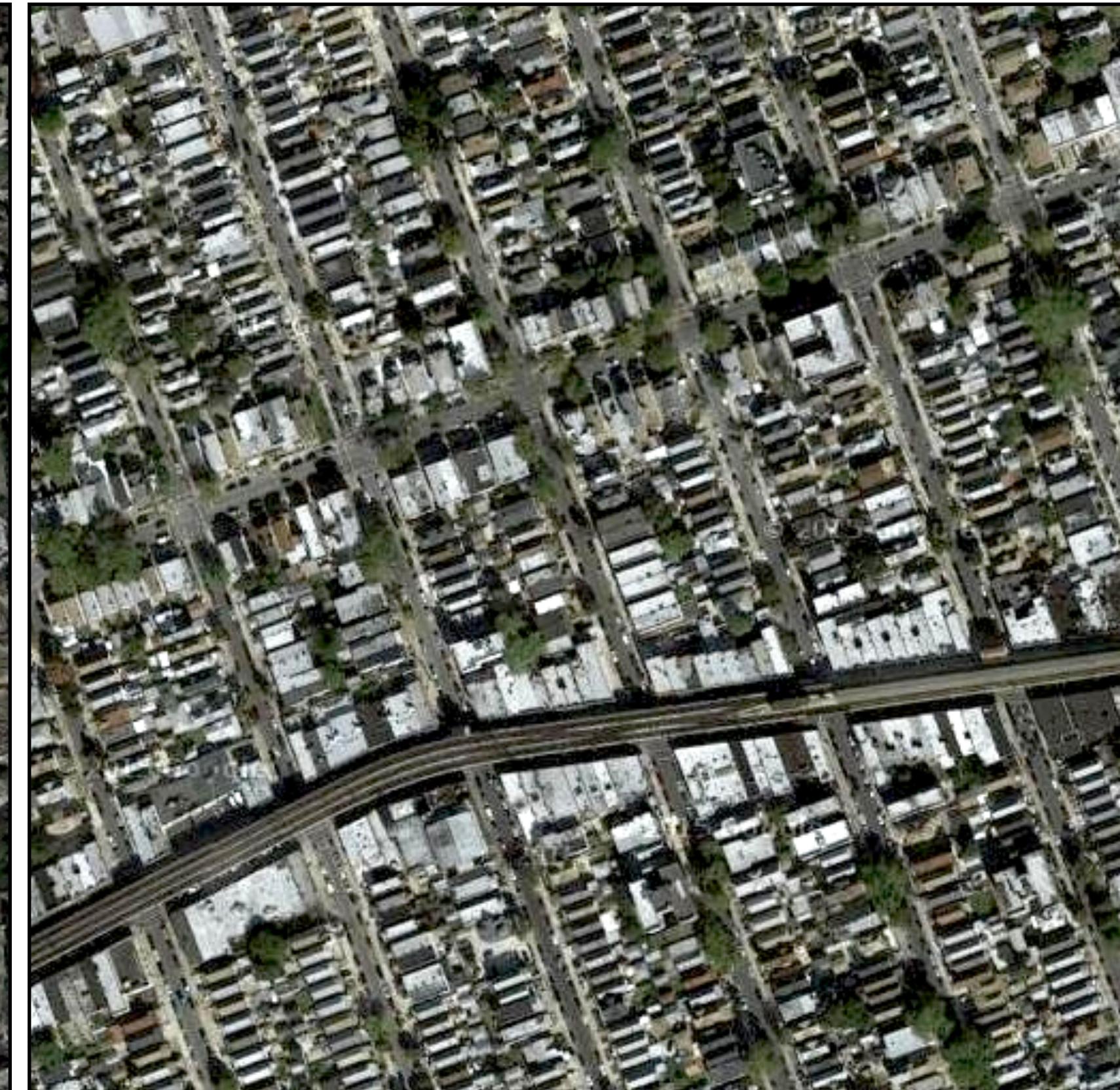
Input



Output



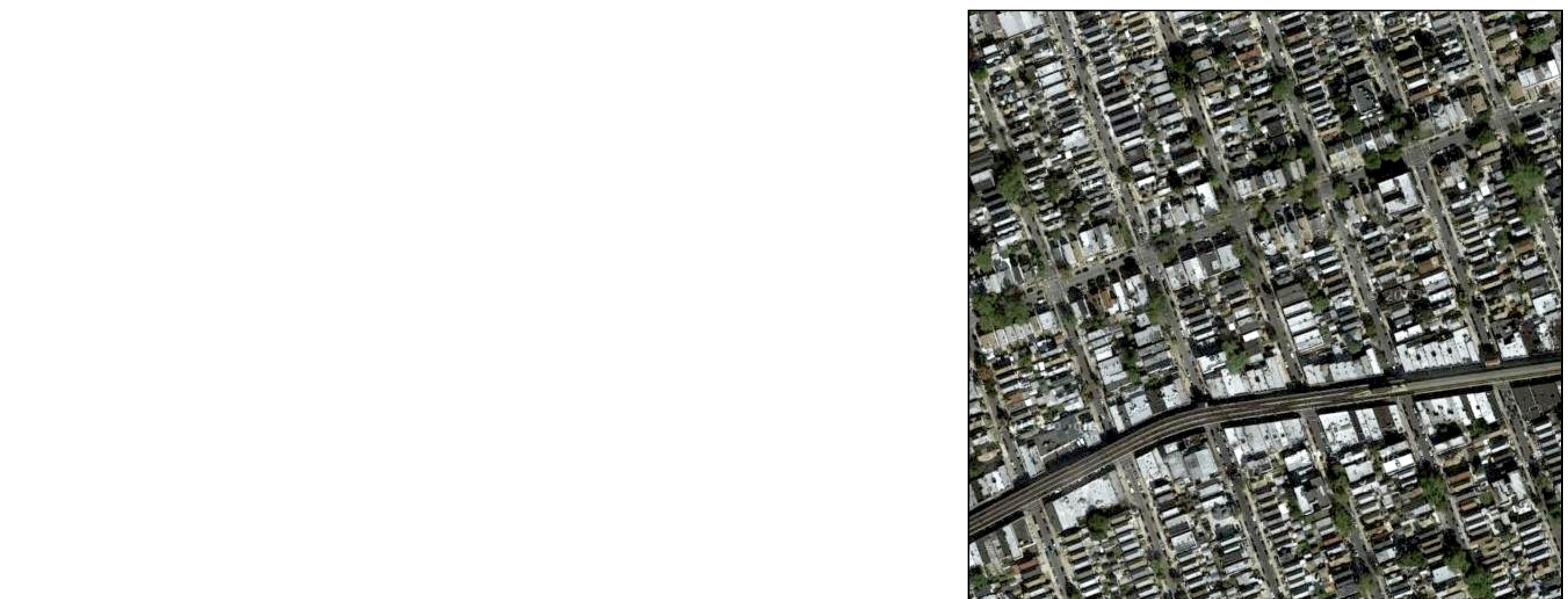
Groundtruth



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

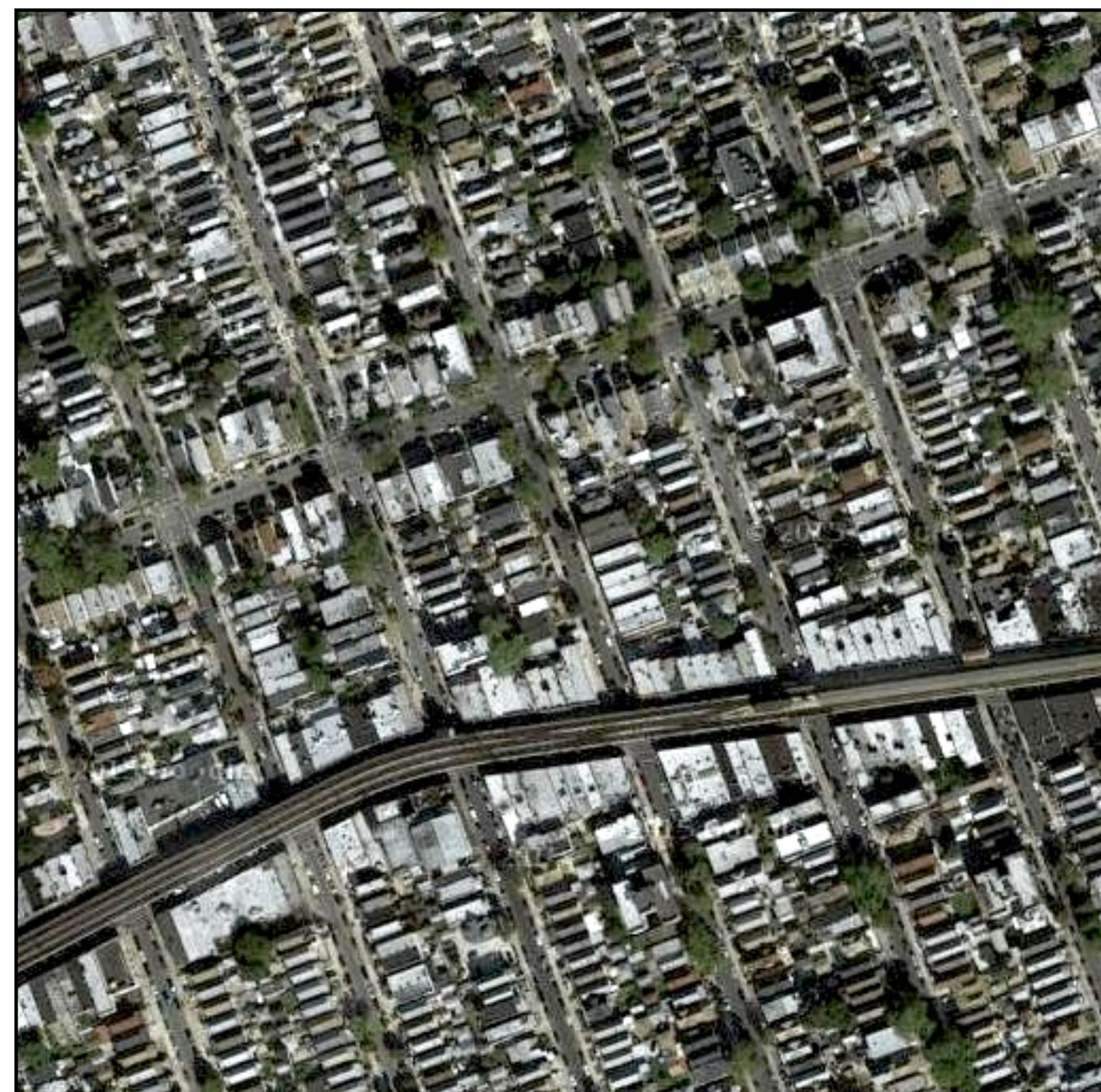


Input



Output

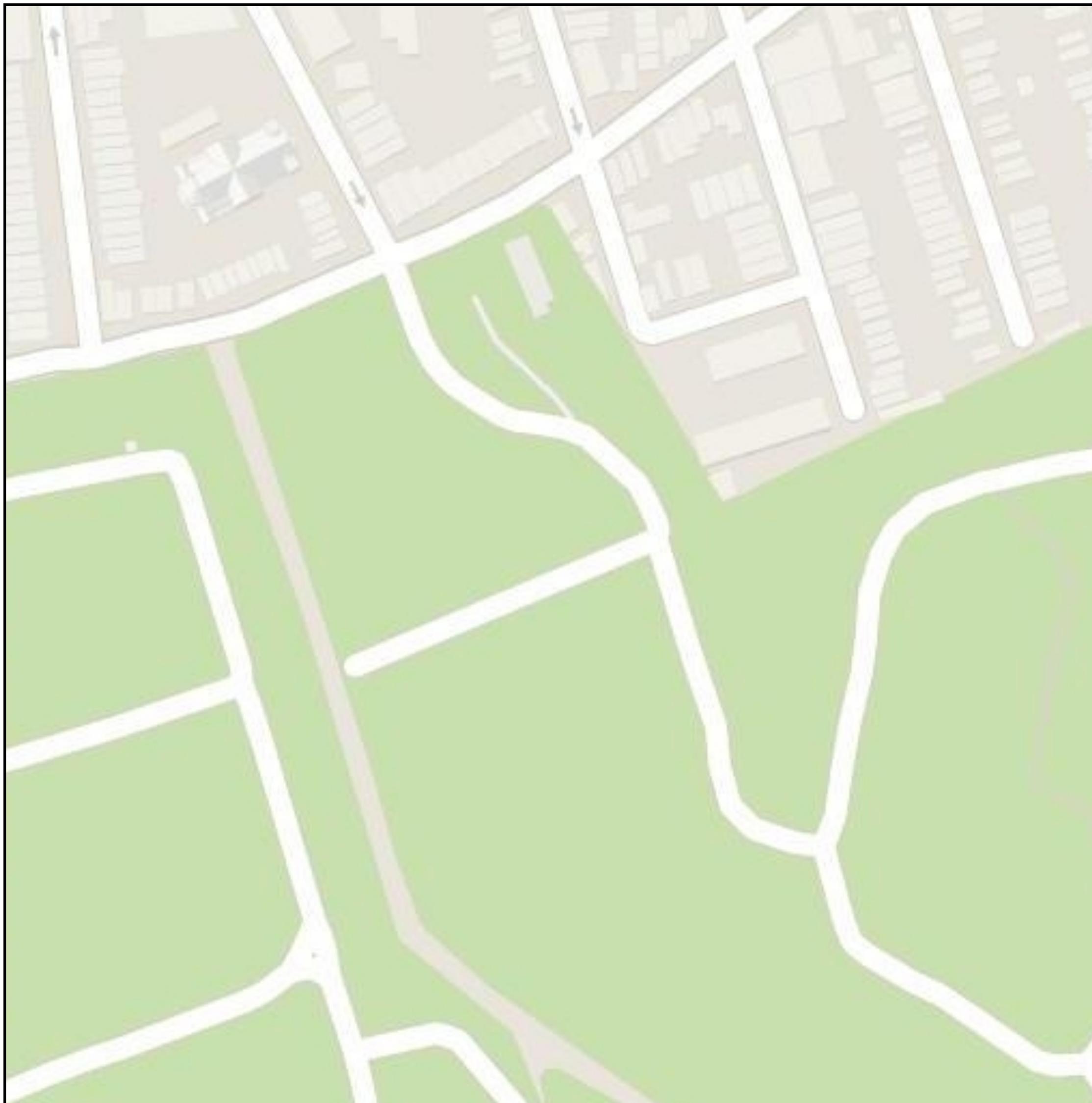
Groundtruth



Data from [maps.google.com]⁶⁵

Source: Isola, Freeman, Torralba

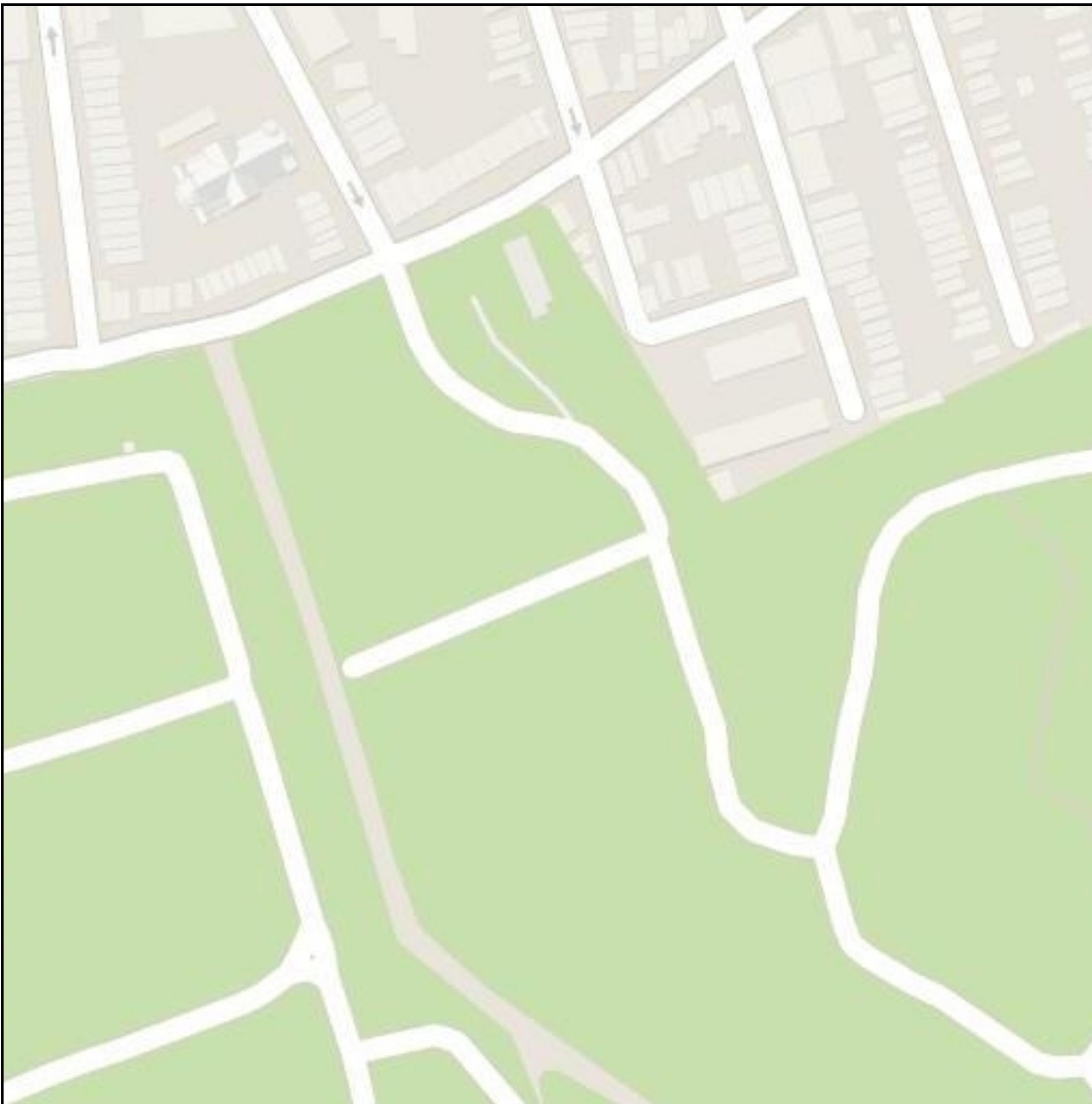
Input



Unstructured prediction (L1)



Input



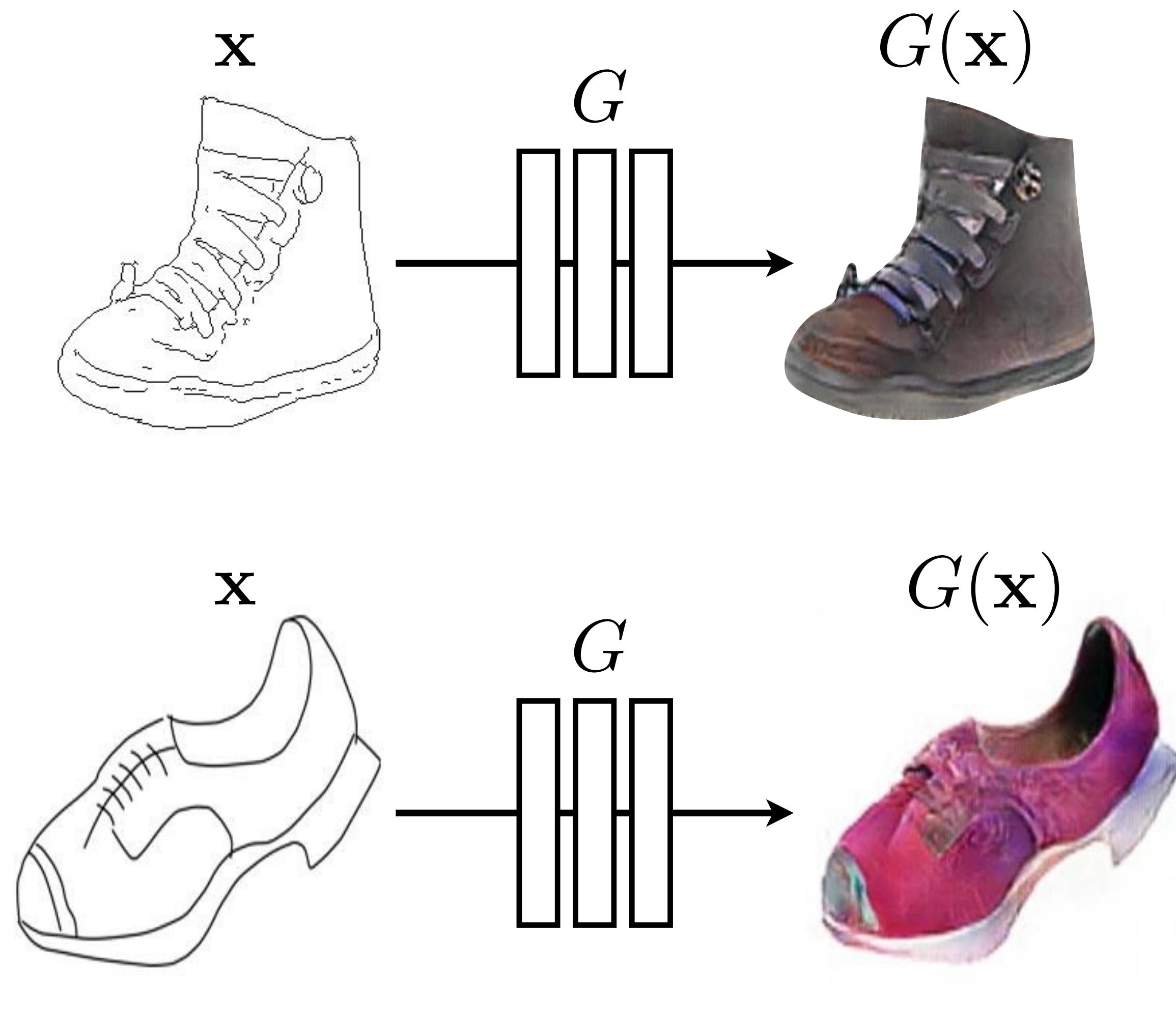
Structured Prediction (cGAN)



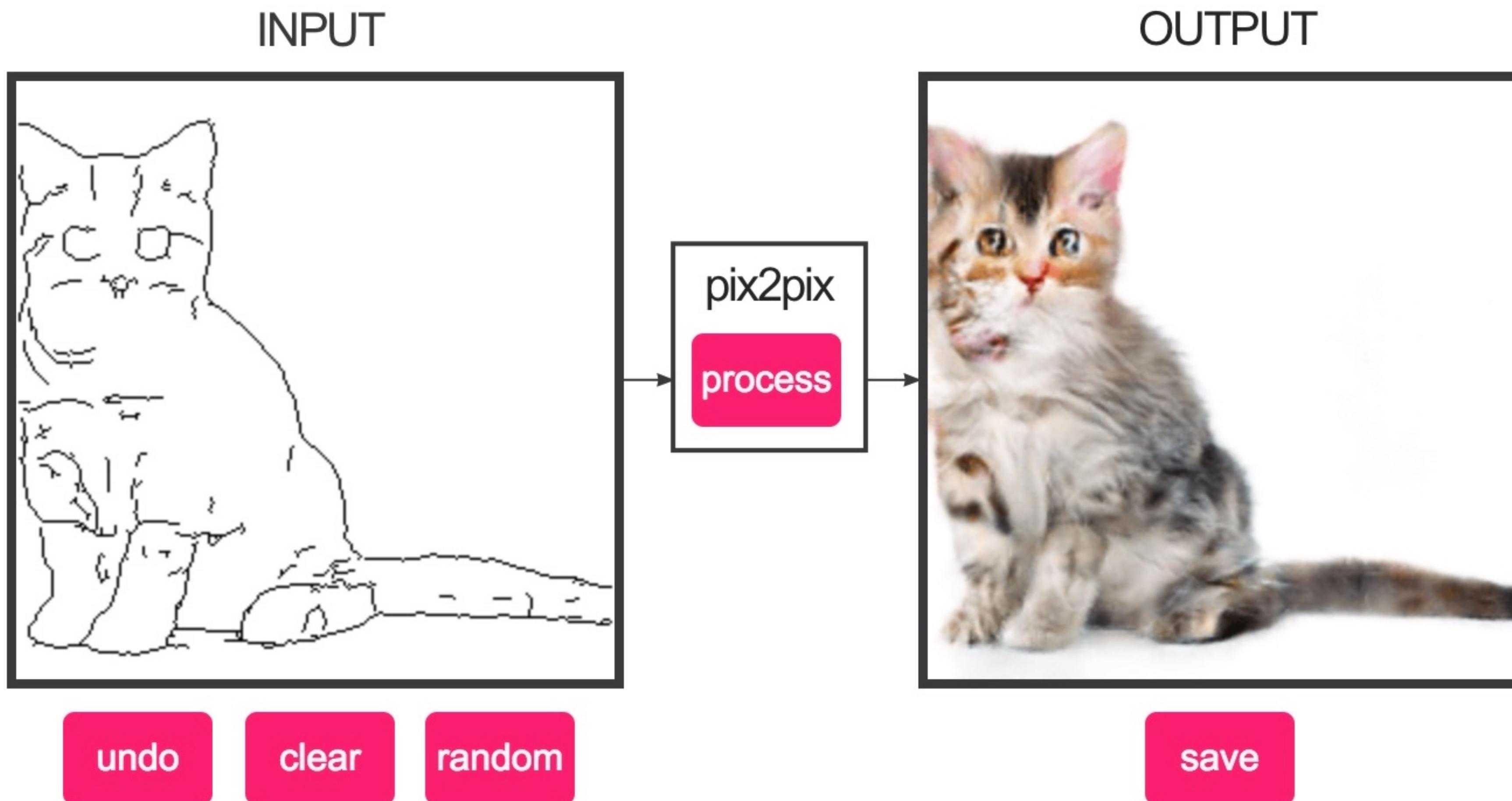
Training data



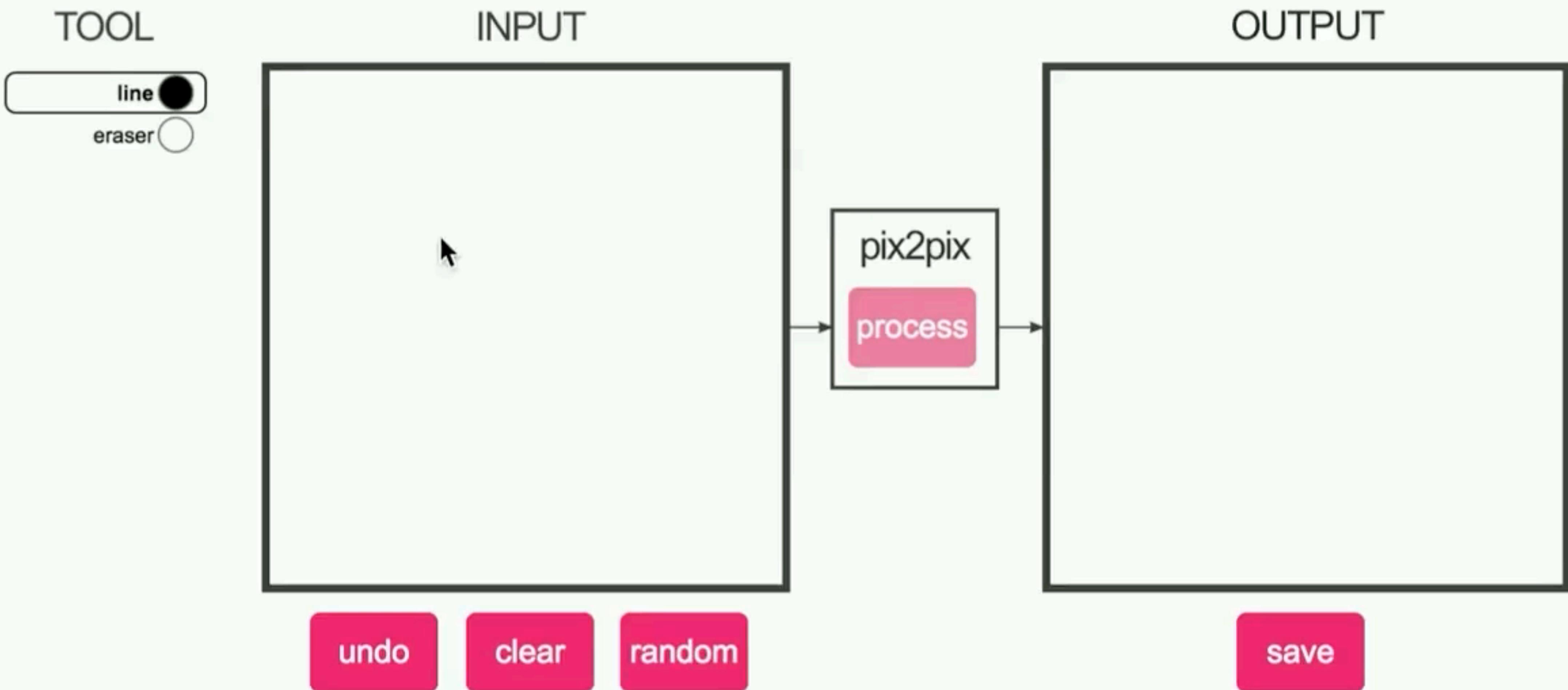
[HED, Xie & Tu, 2015]

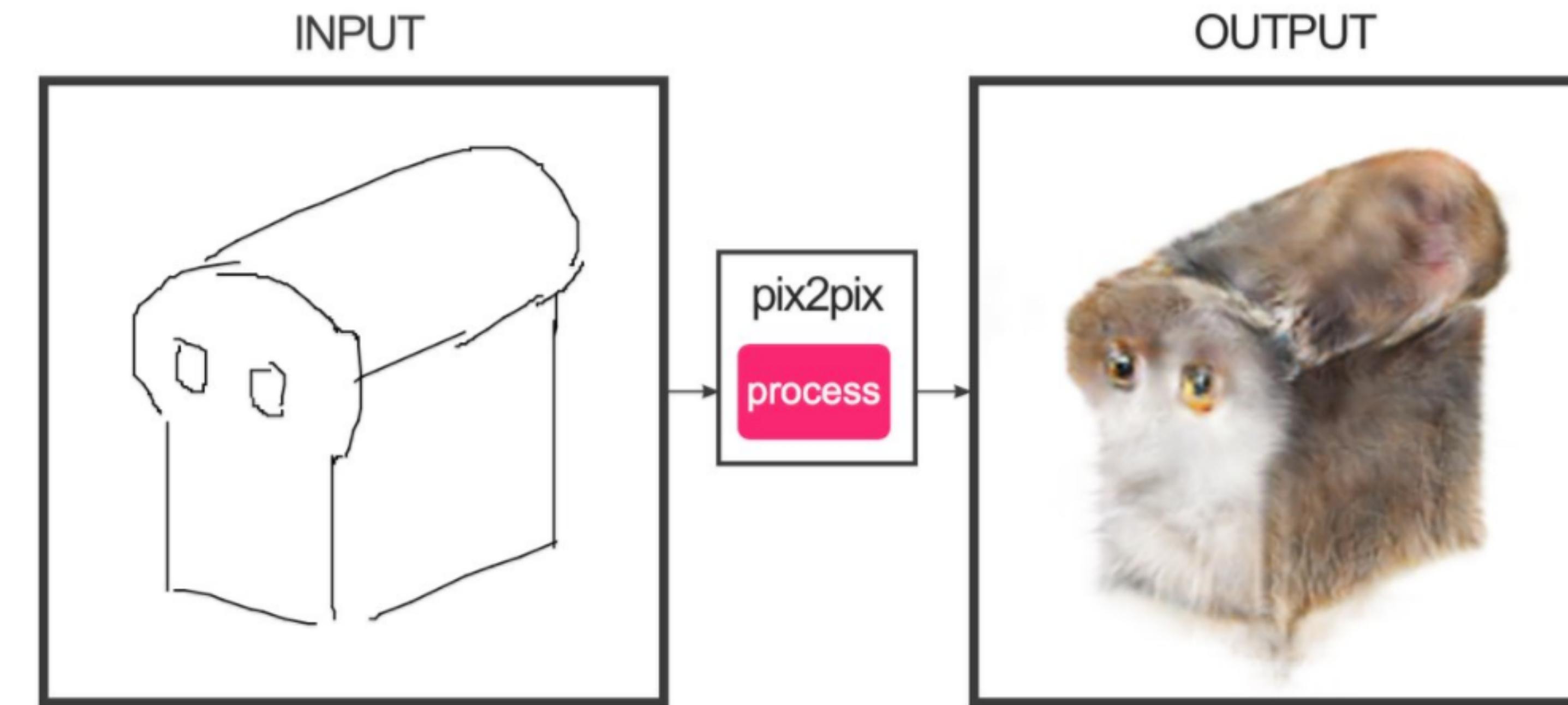


#edges2cats [Chris Hesse]



edges2cats

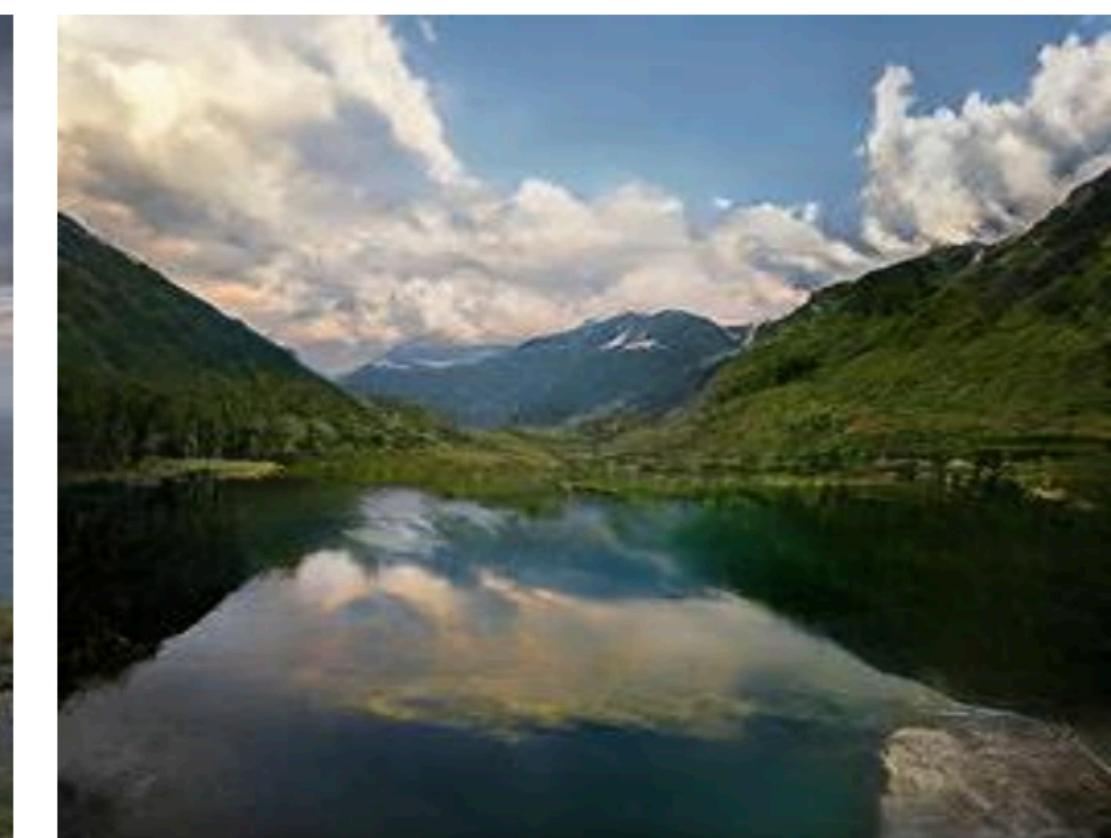
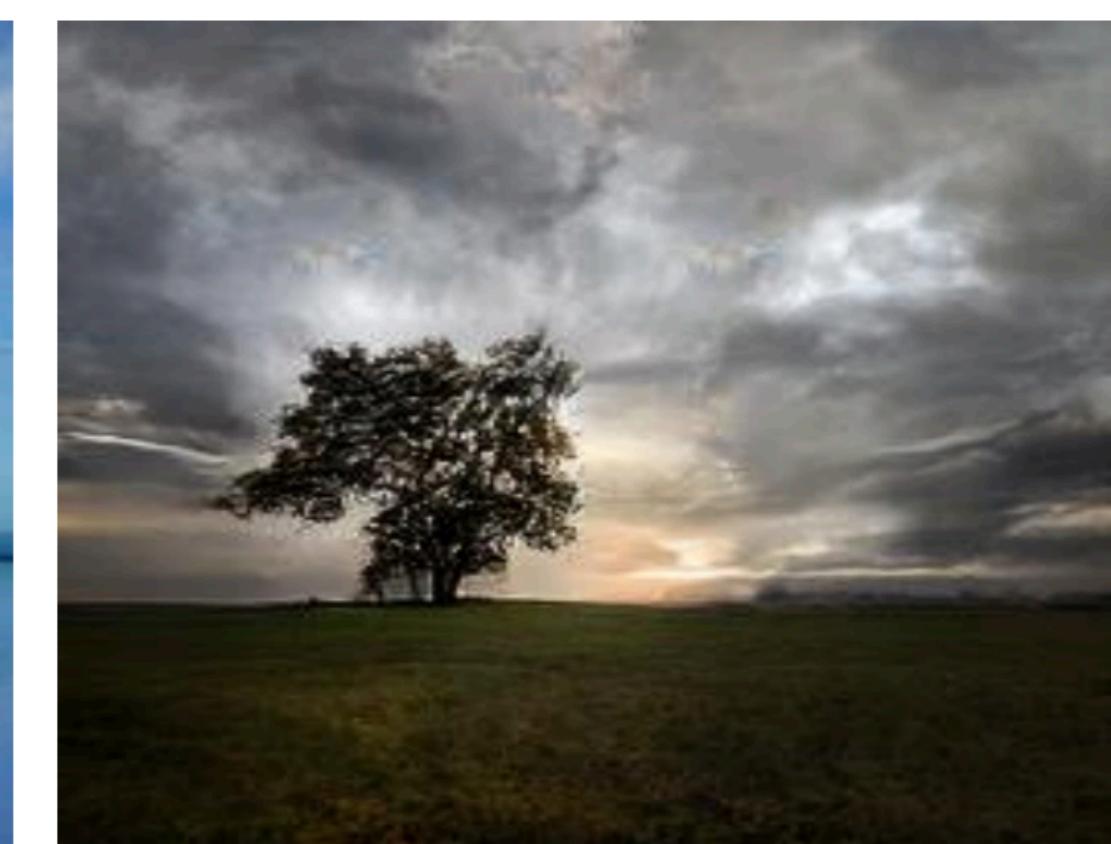
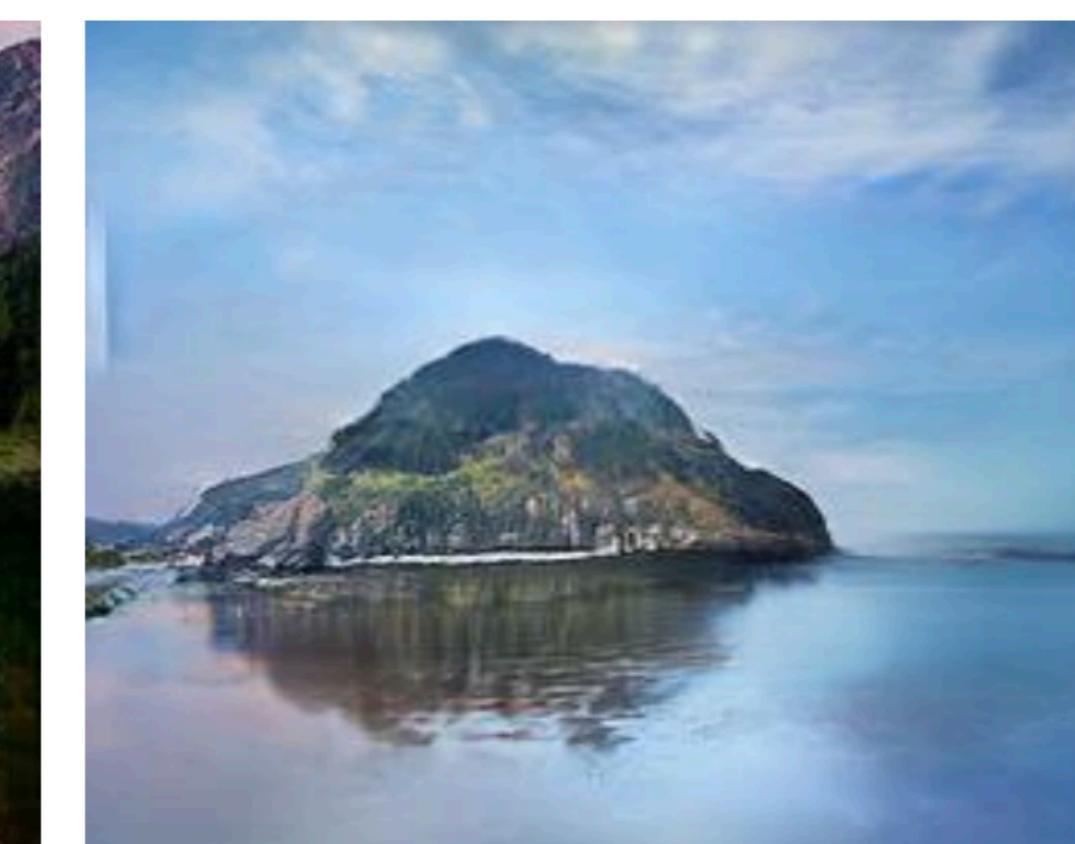
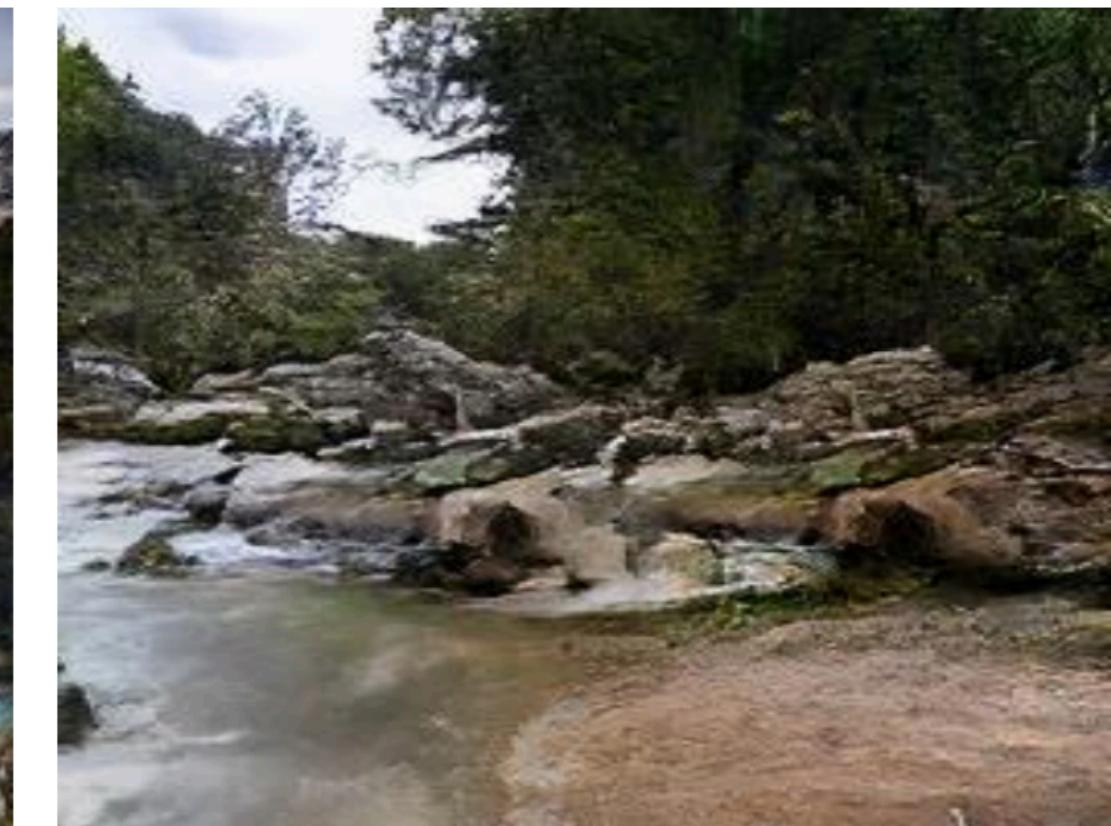
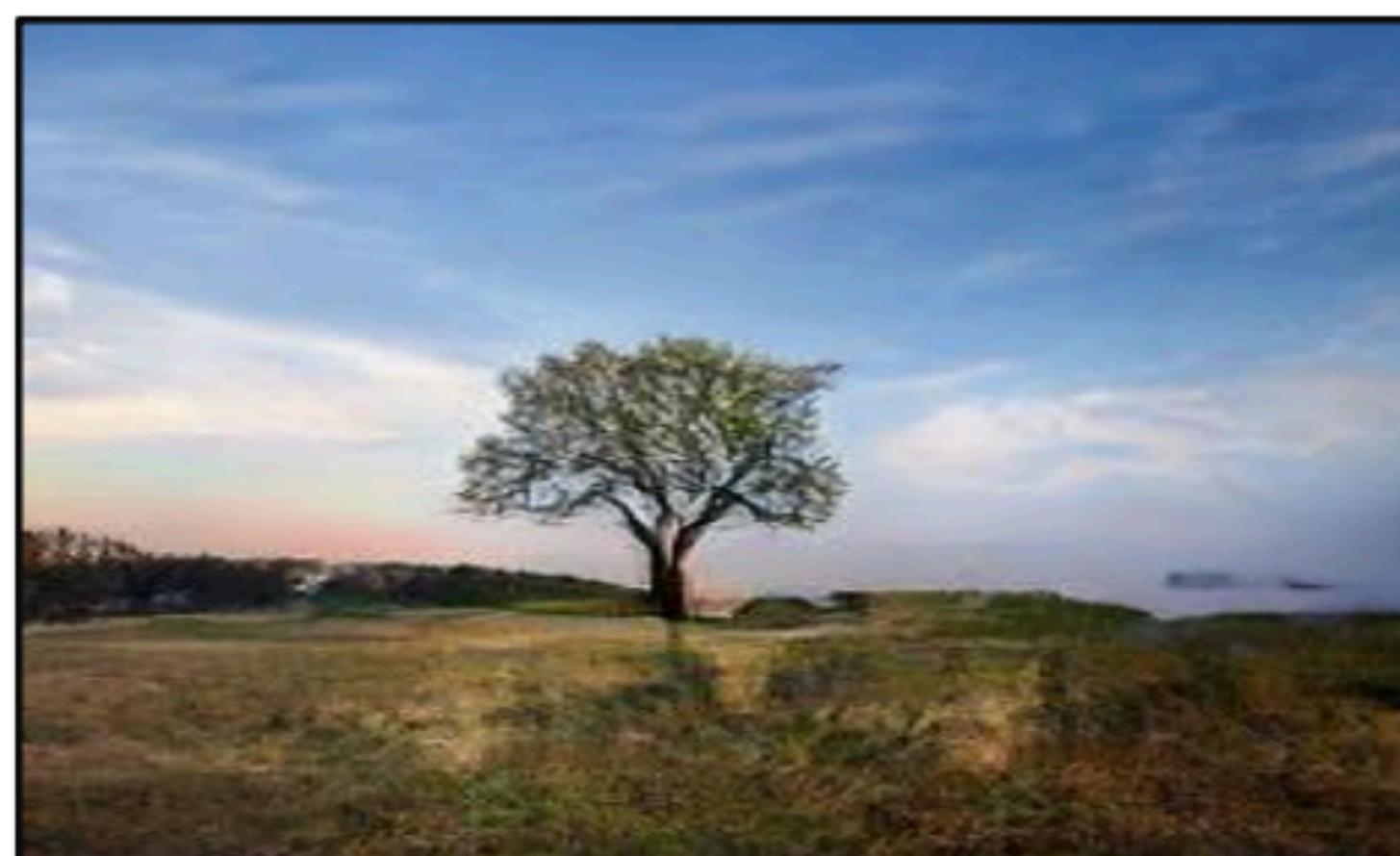
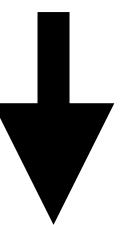


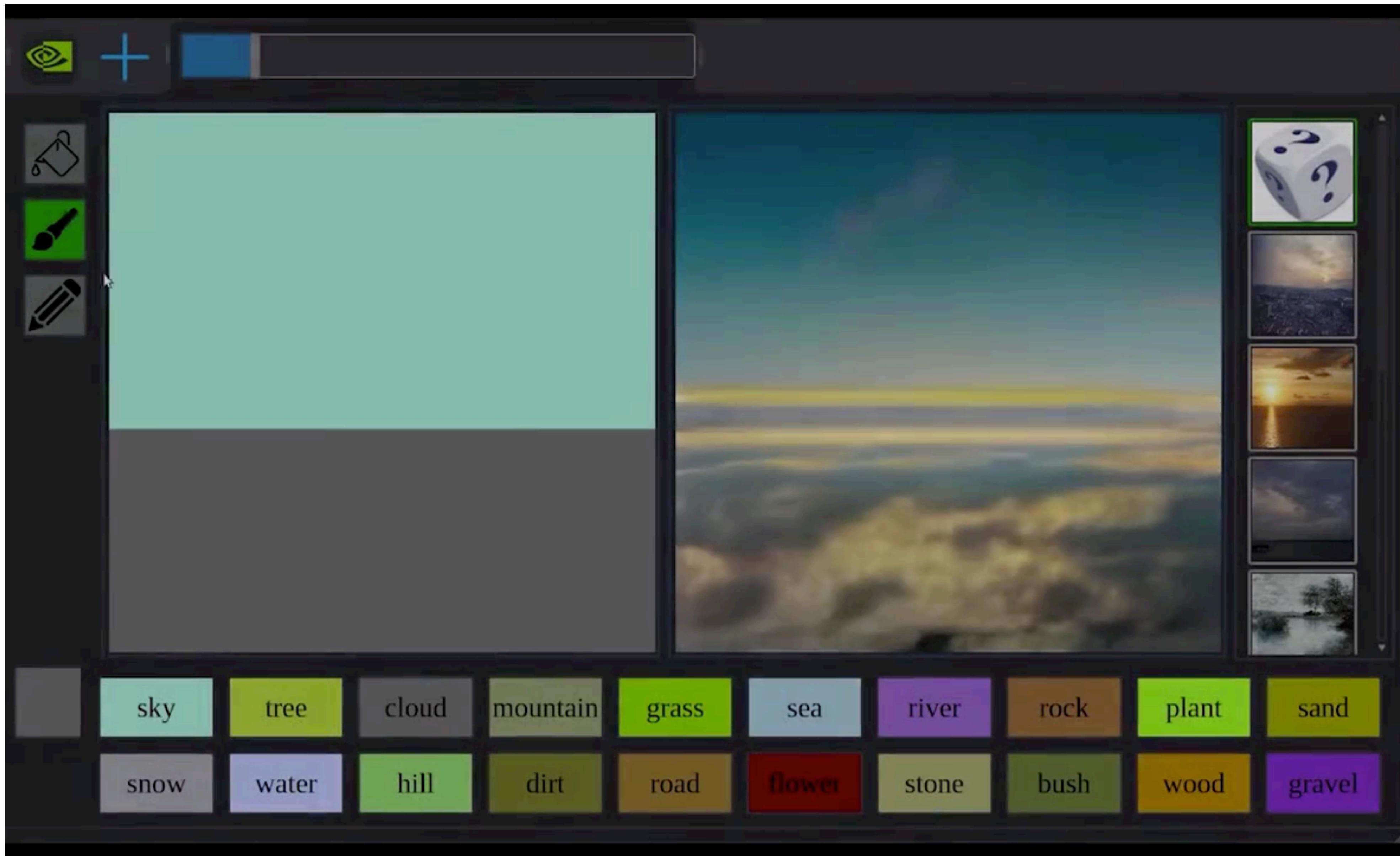


Ivy Tasi @ivymyt



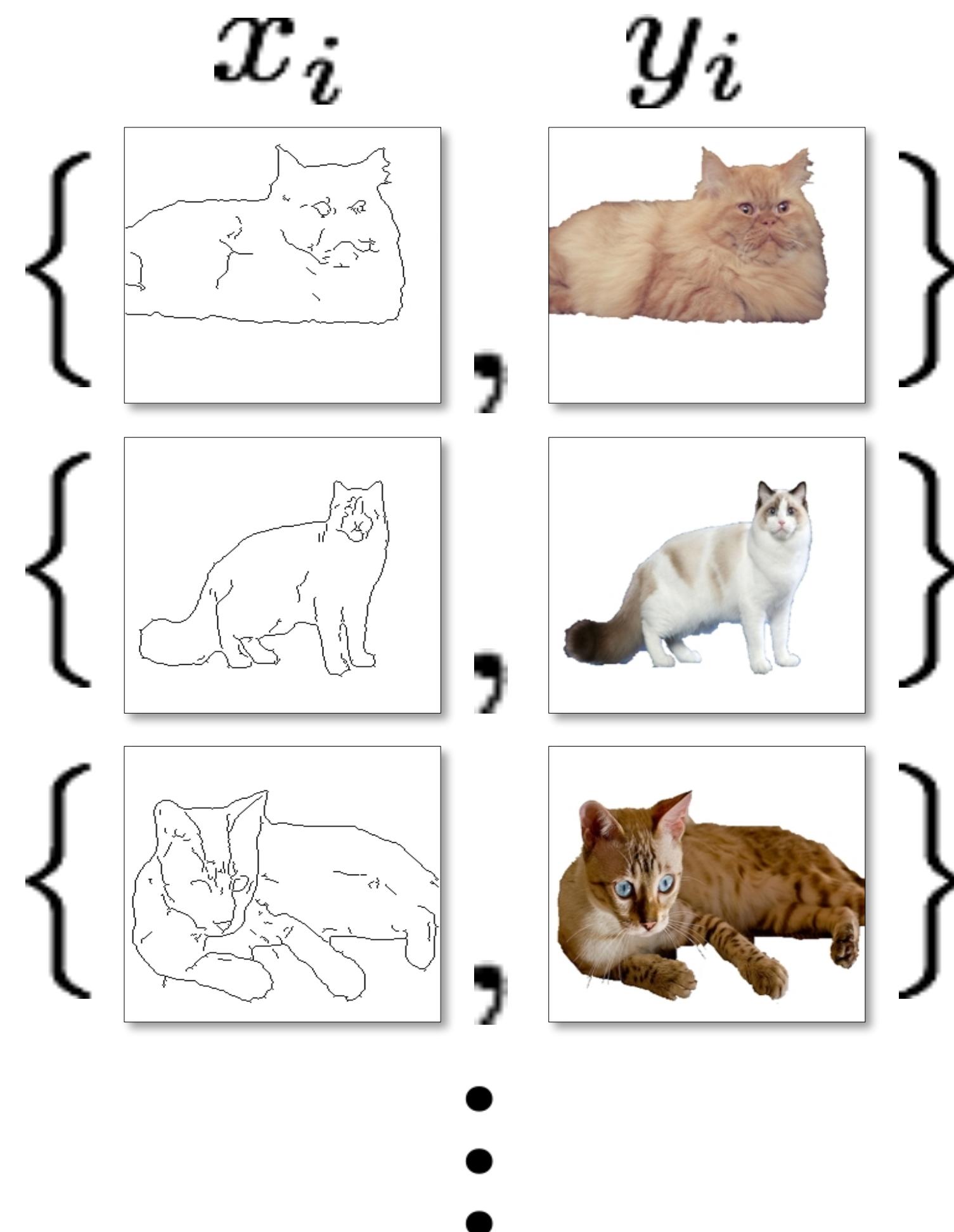
Vitaly Vidmirov @vvid



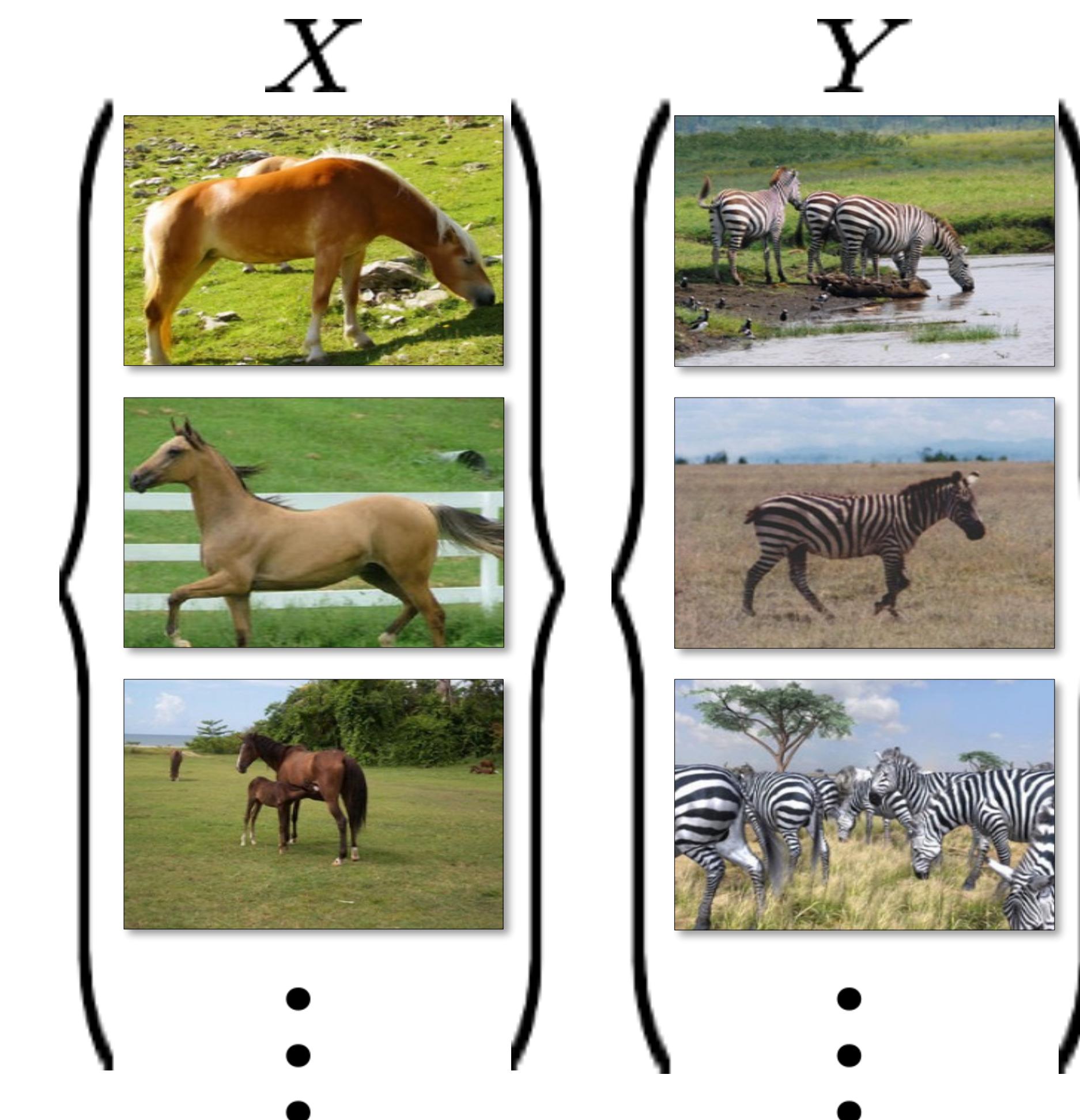


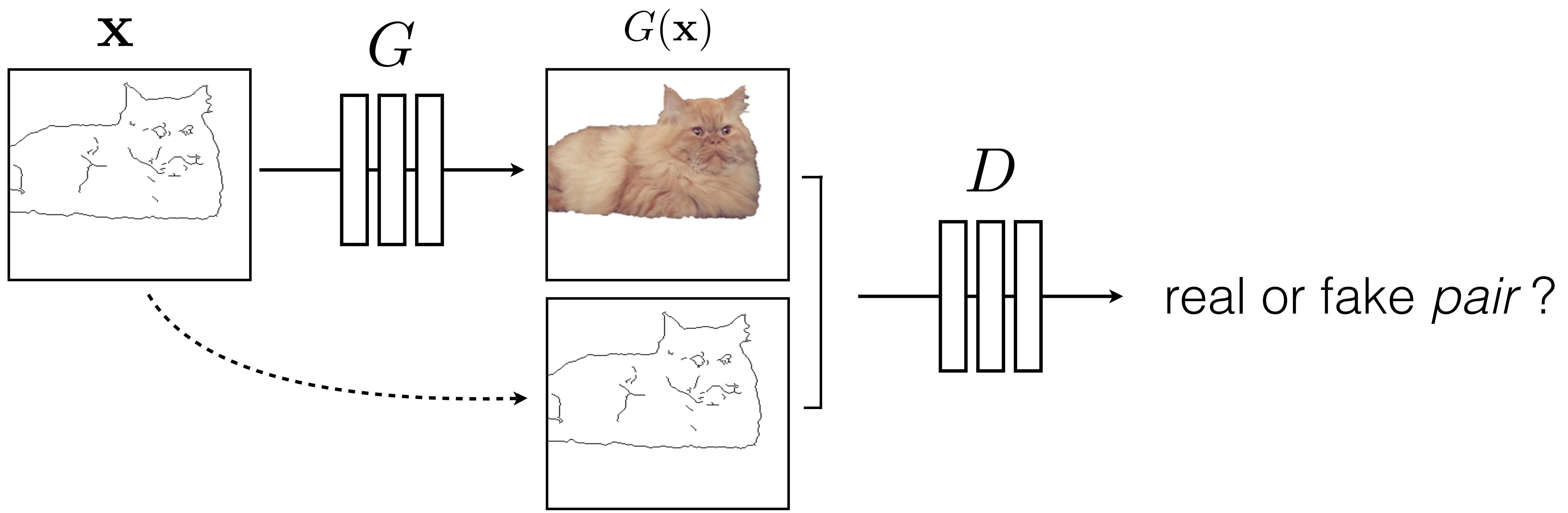
Handling unpaired data

Paired data

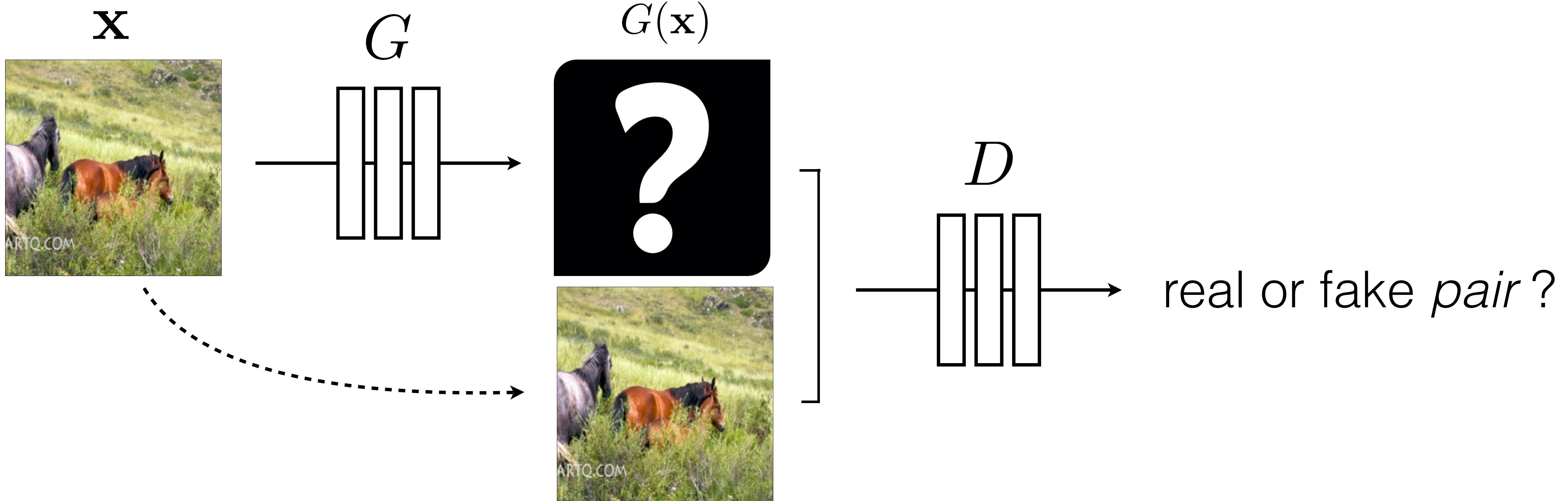


Unpaired data



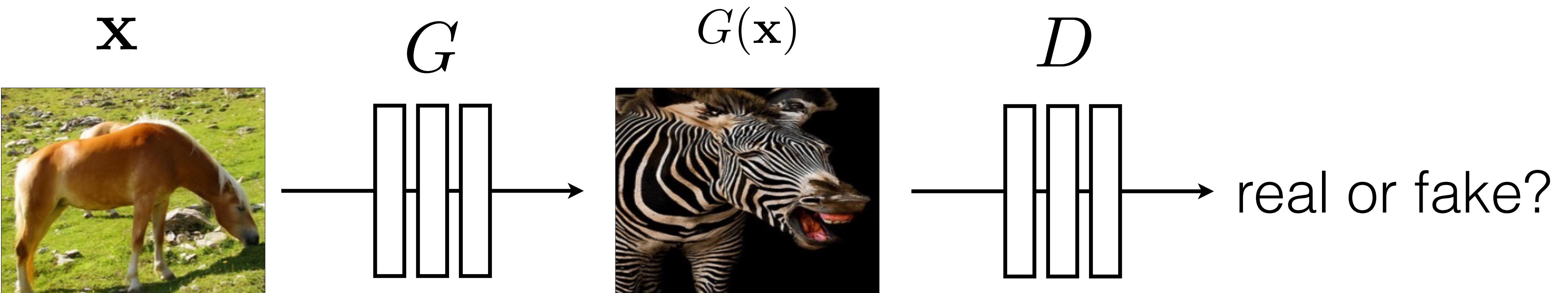


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))]$$



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

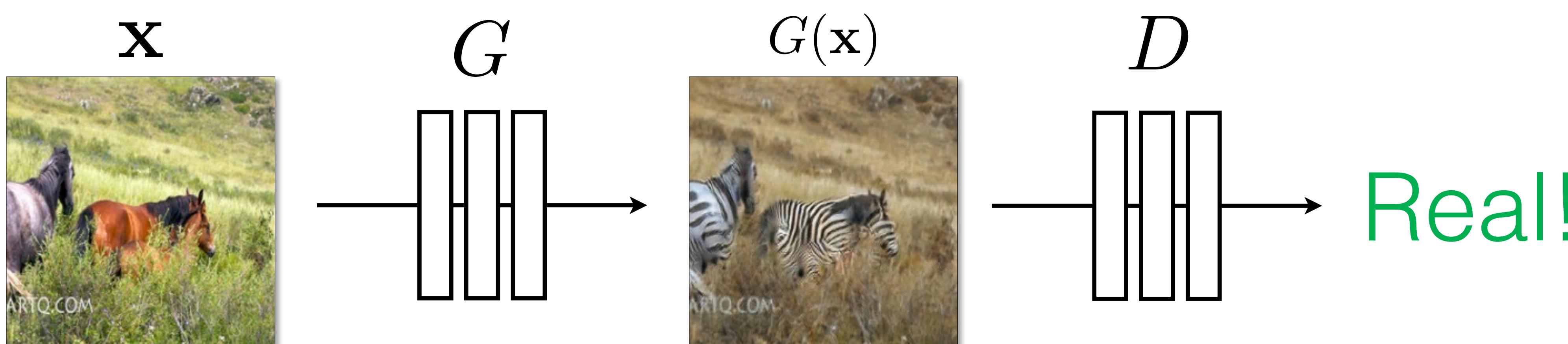
No input-output pairs!



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

Usually loss functions check if output matches a target *instance*

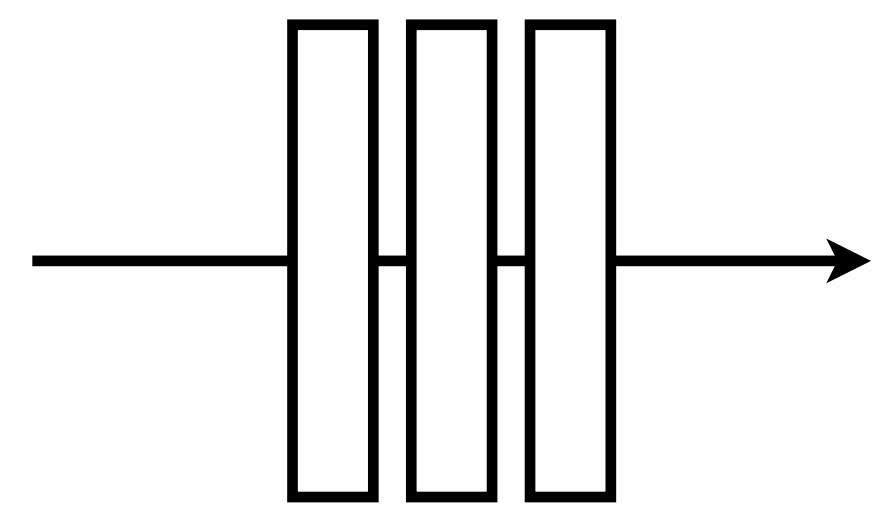
GAN loss checks if output is part of an admissible set



\mathbf{x}



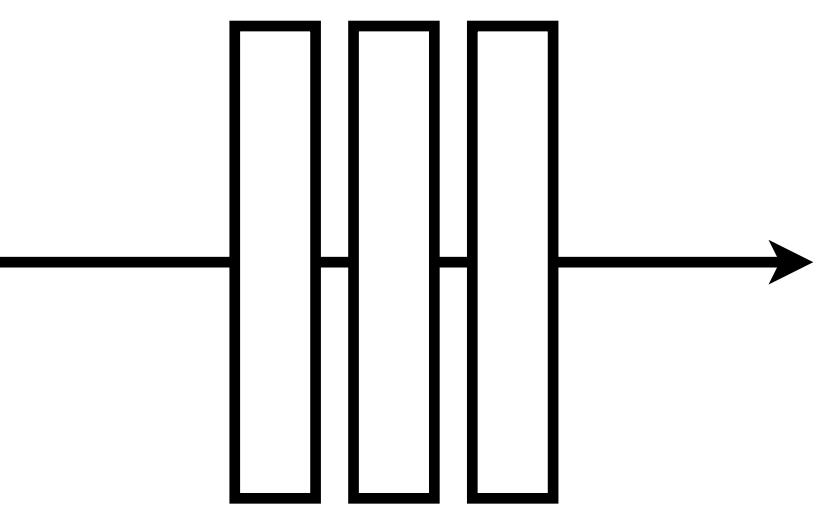
G



$G(\mathbf{x})$



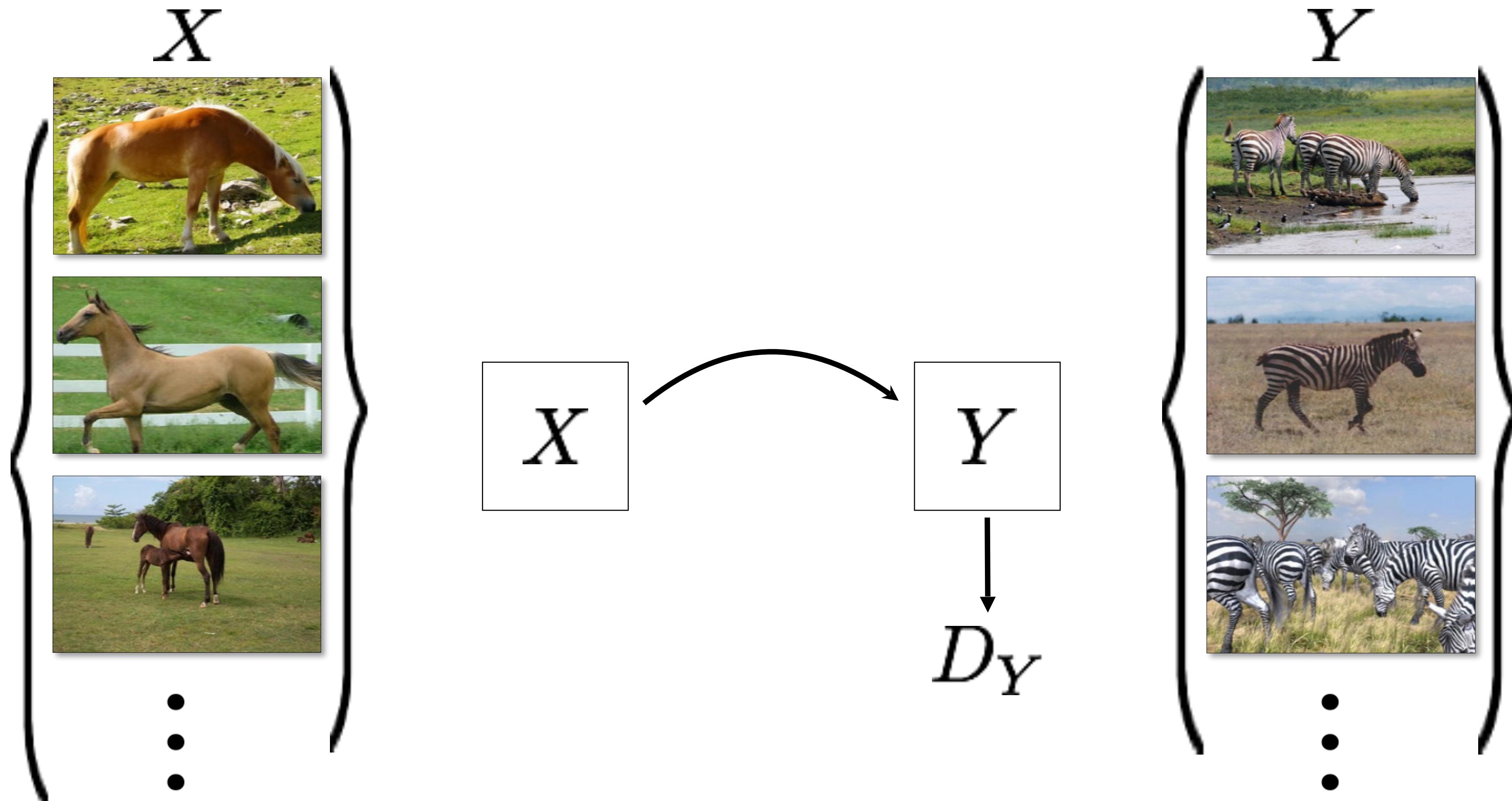
D



Real too!

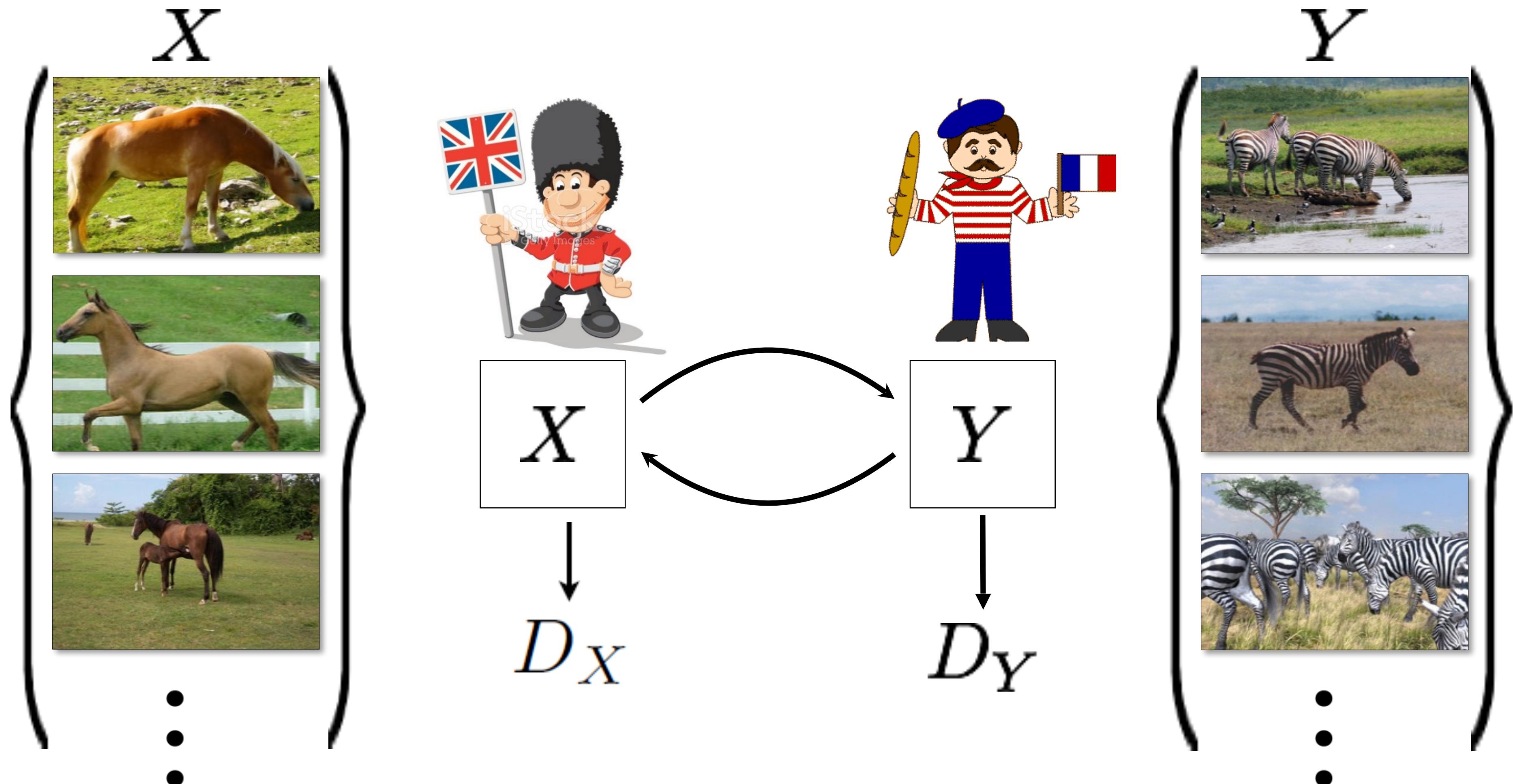
Nothing to force output to correspond to input

CycleGAN

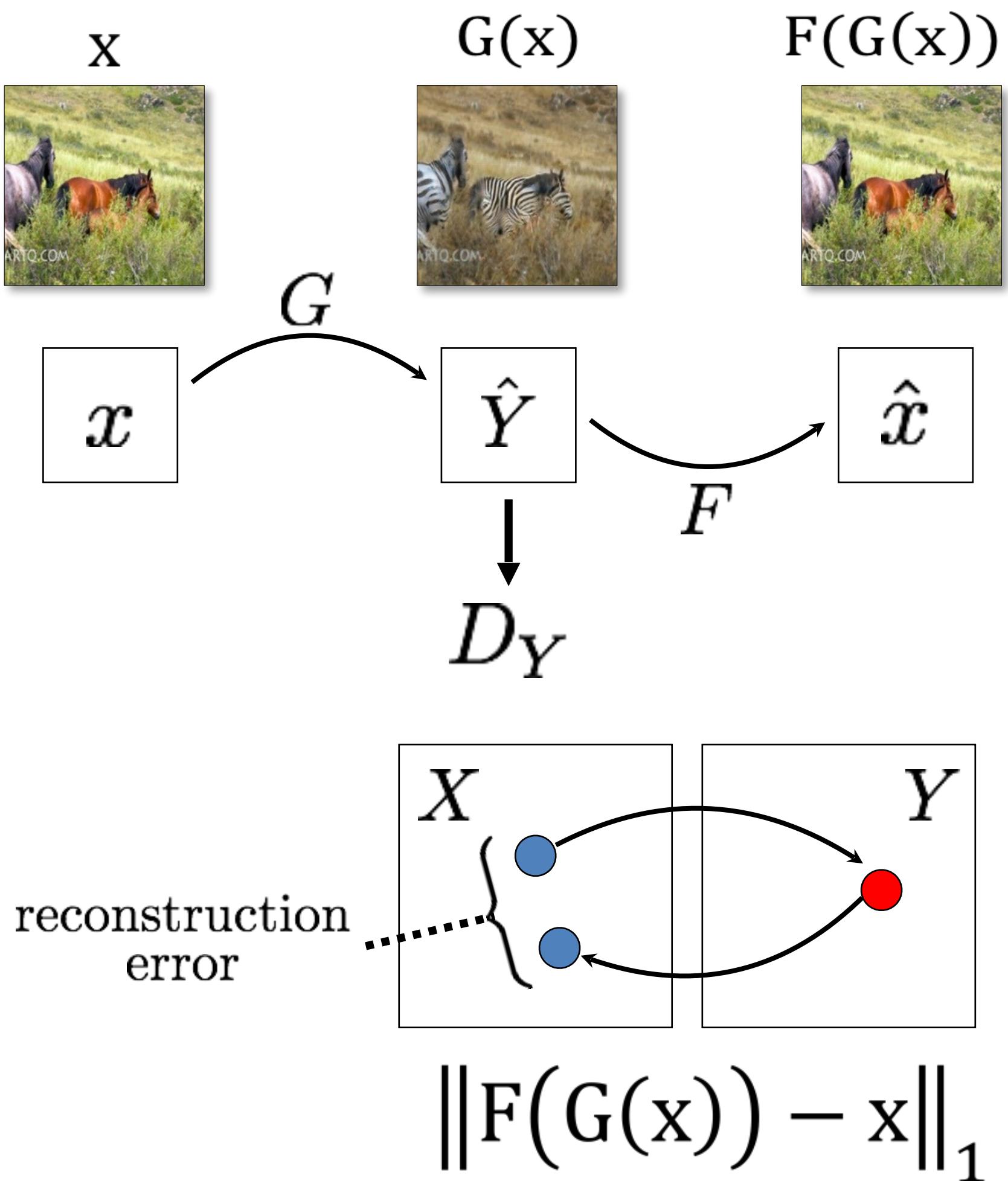


[Zhu*, Park* et al. 2017], [Yi et al. 2017], [Kim et al. 2017]

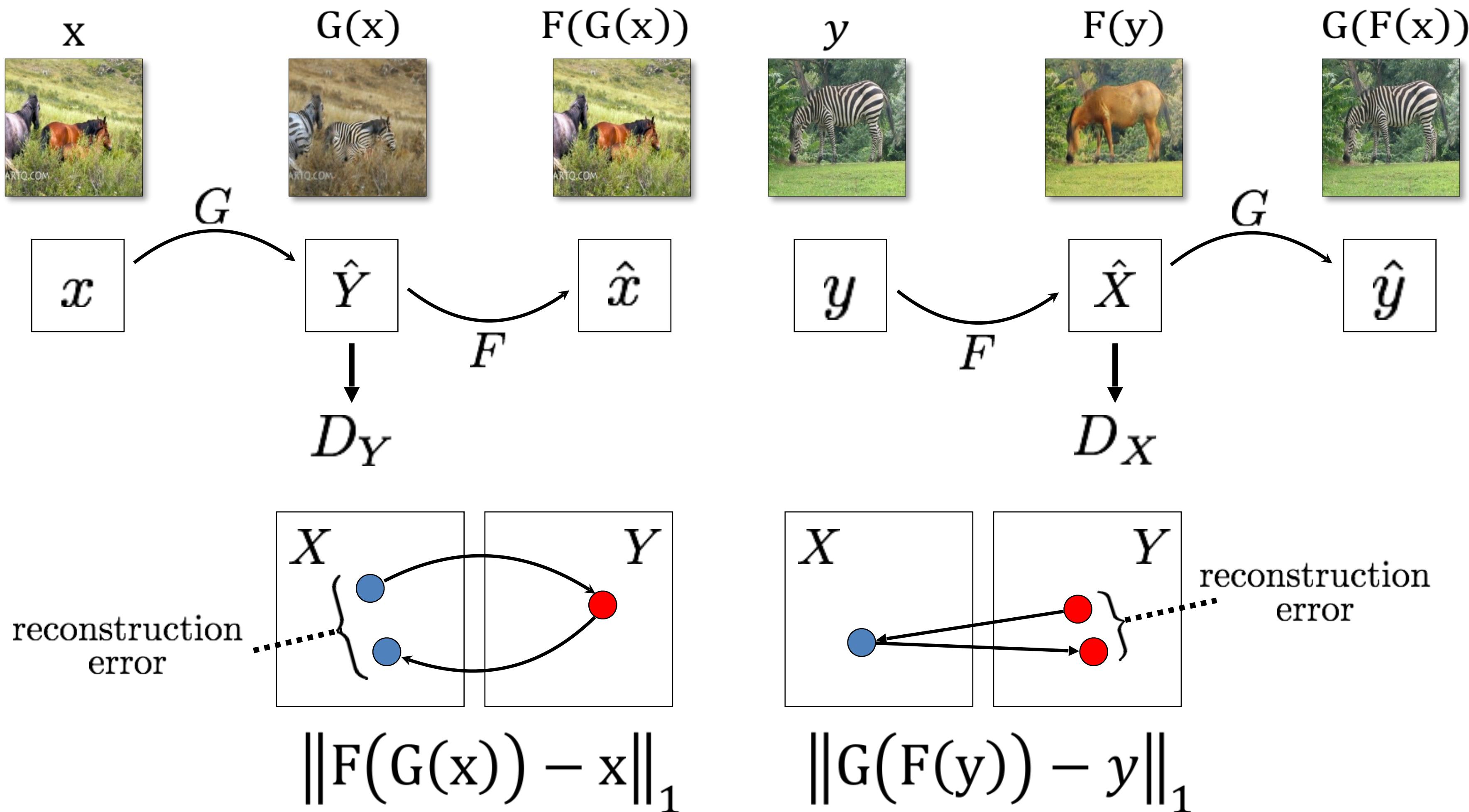
CycleGAN

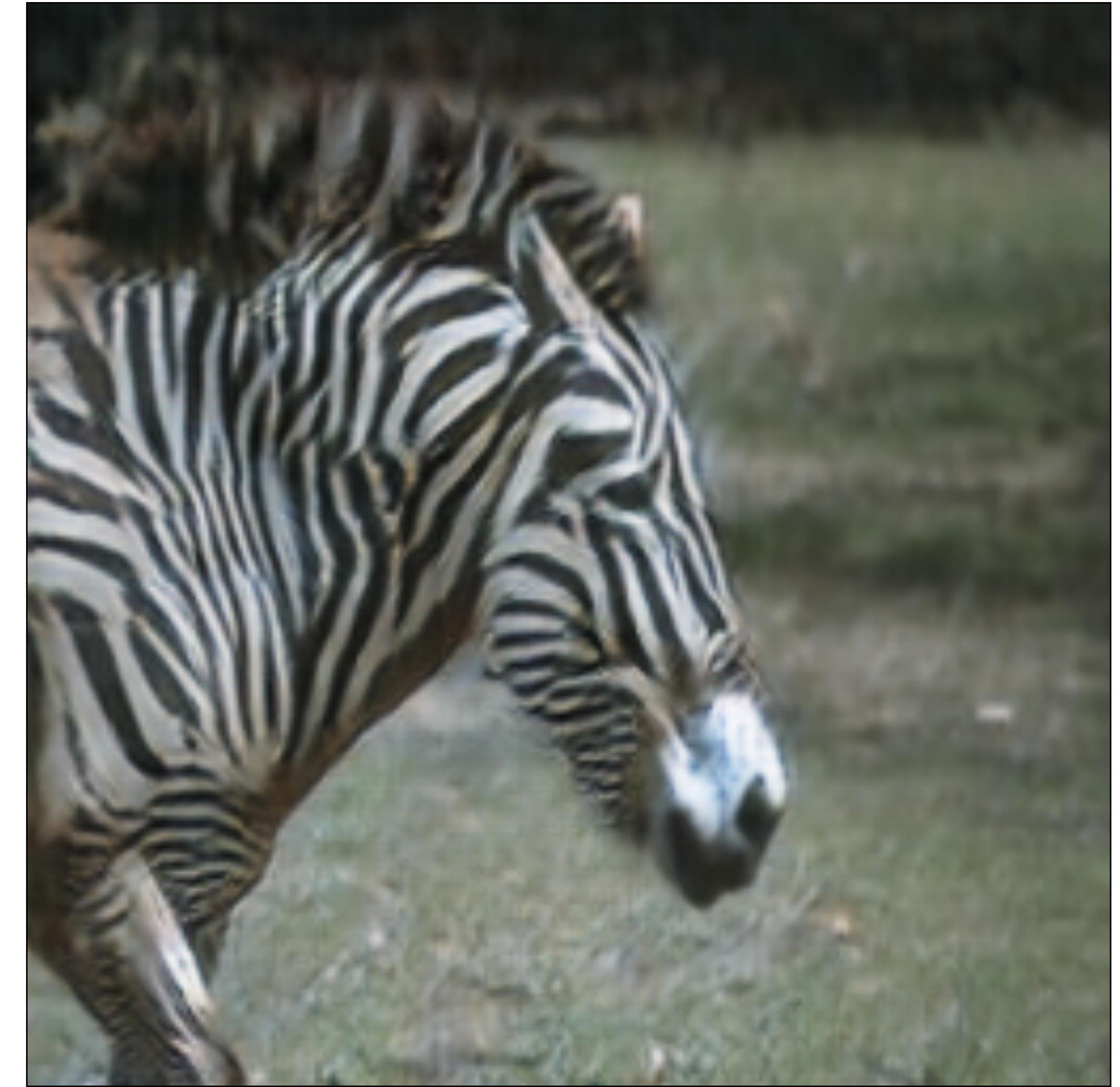


Cycle Consistency Loss



Cycle Consistency Loss







Input



Monet



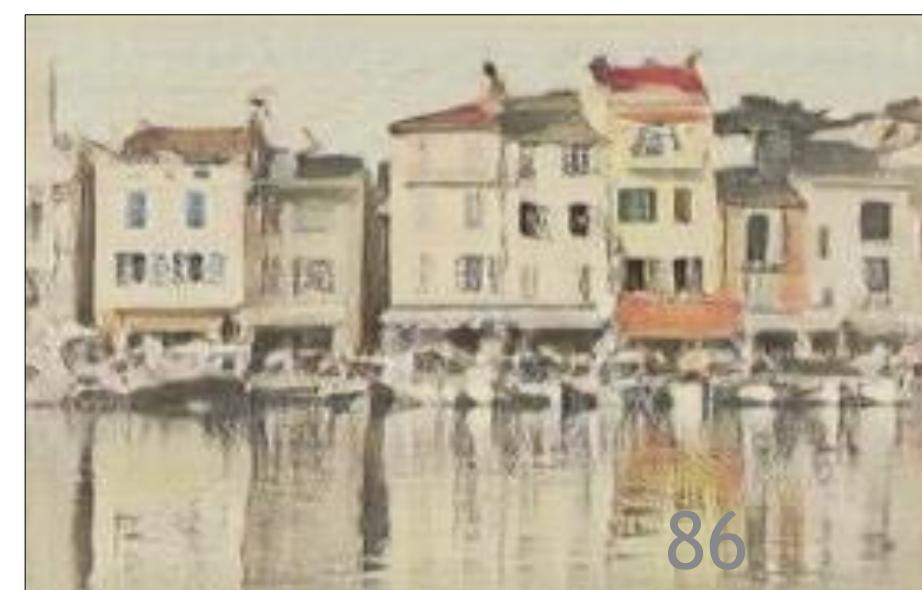
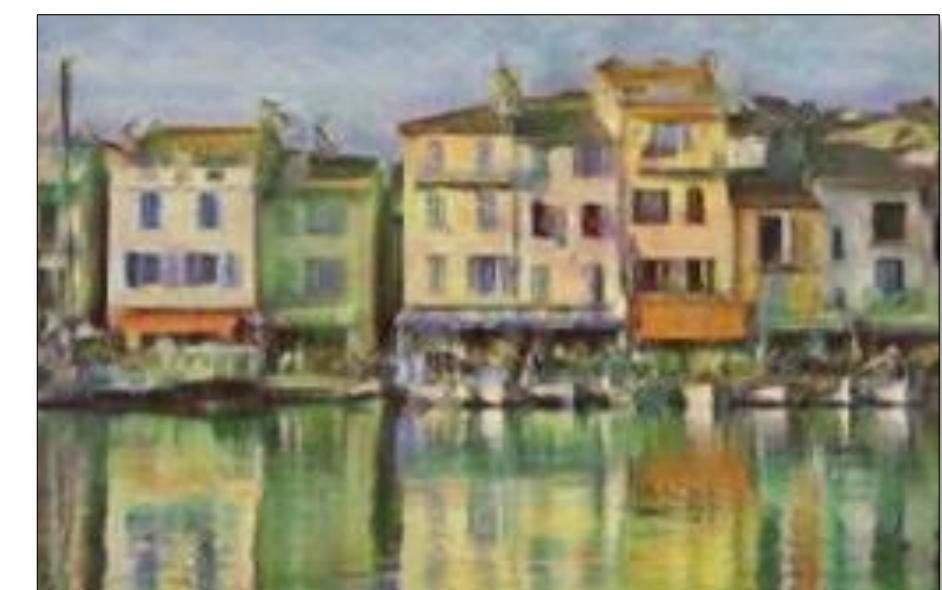
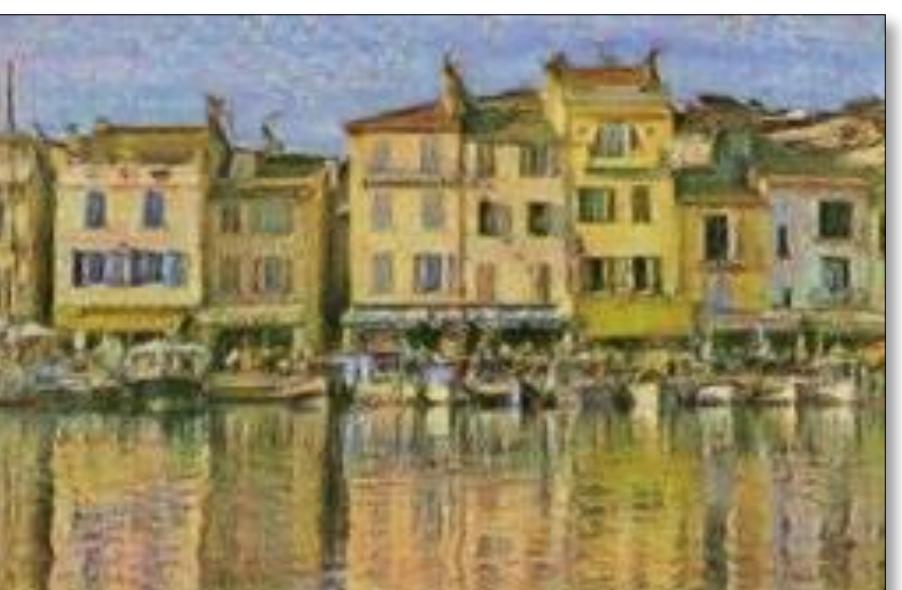
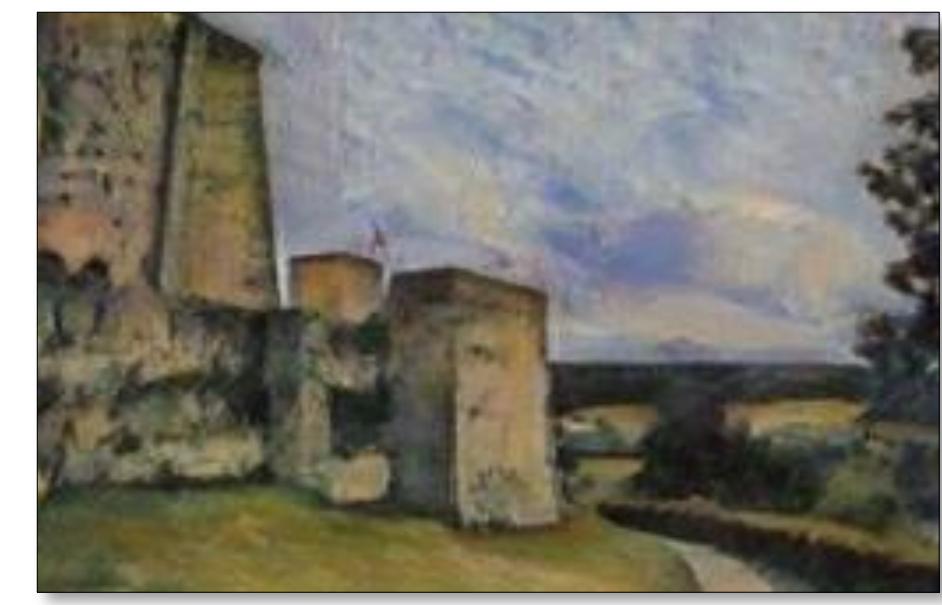
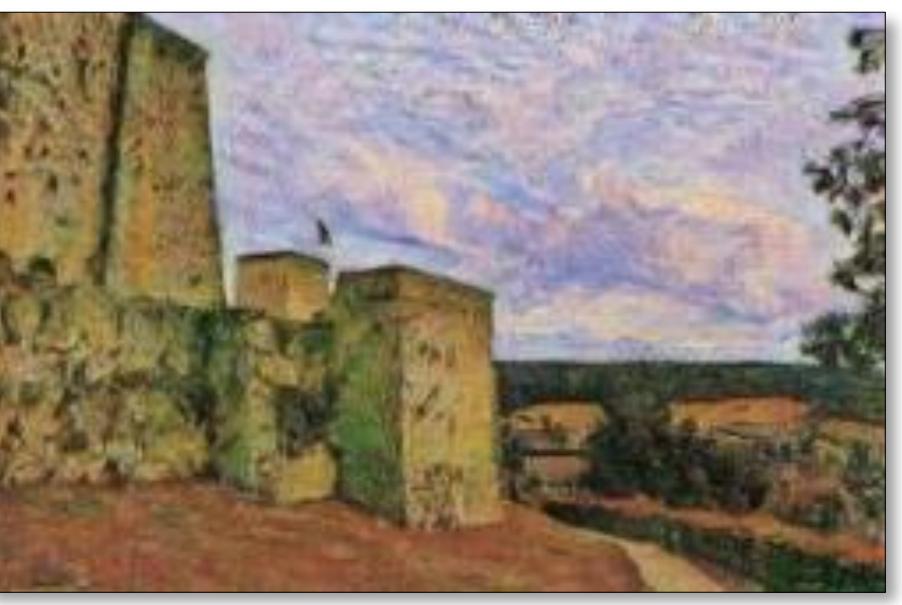
Van Gogh



Cezanne



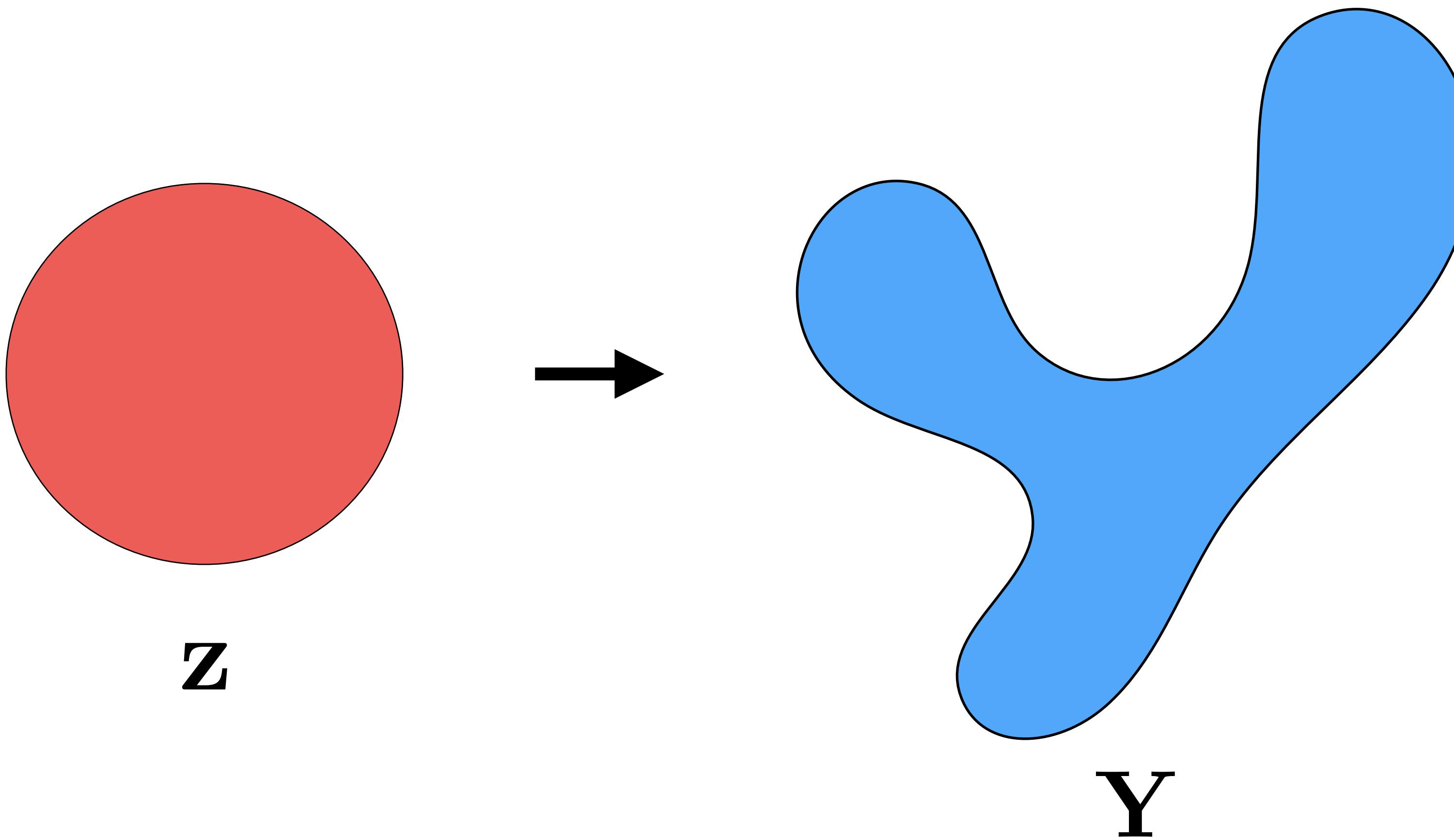
Ukiyo-e



GANs

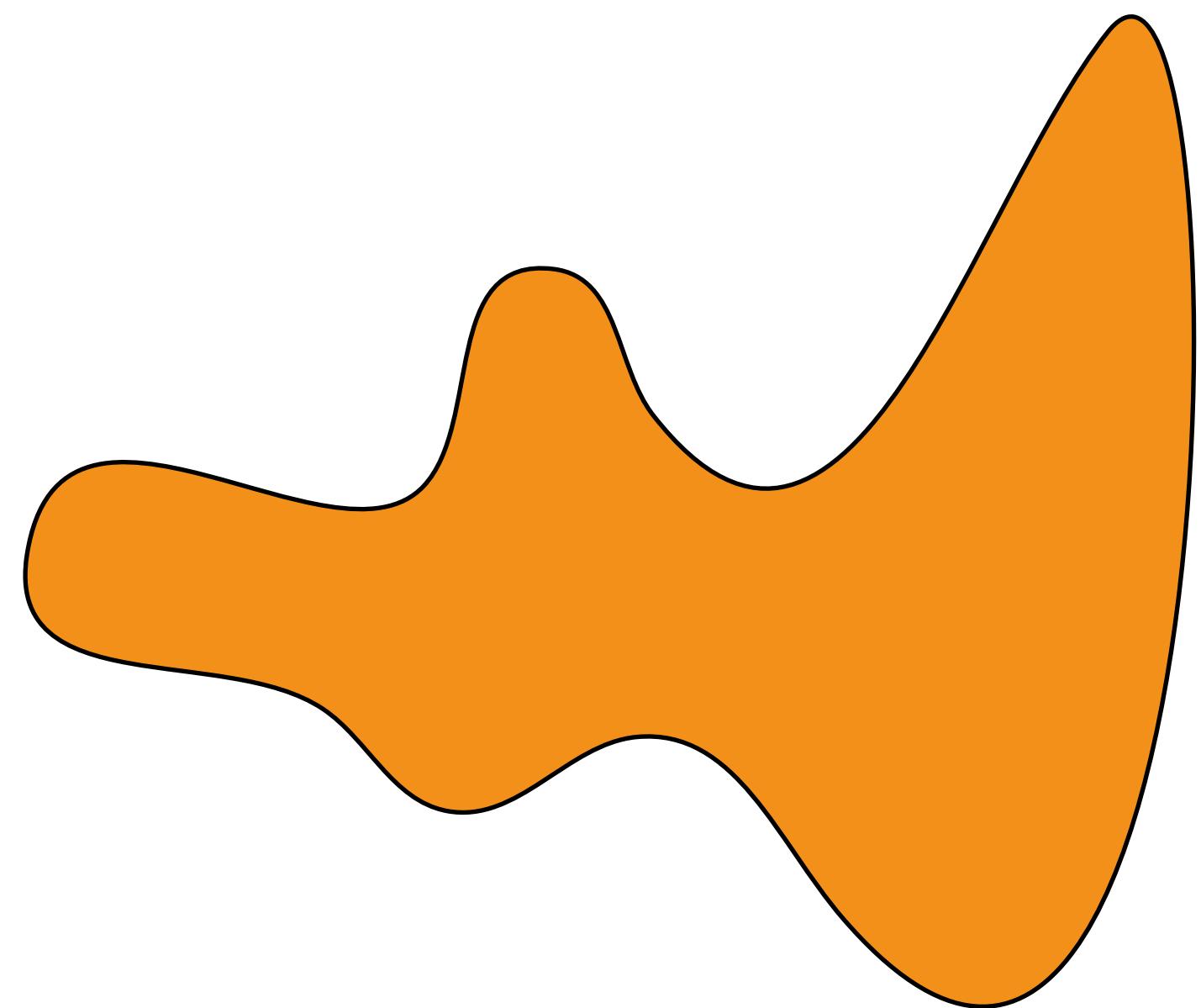
Gaussian

Target distribution



CycleGAN

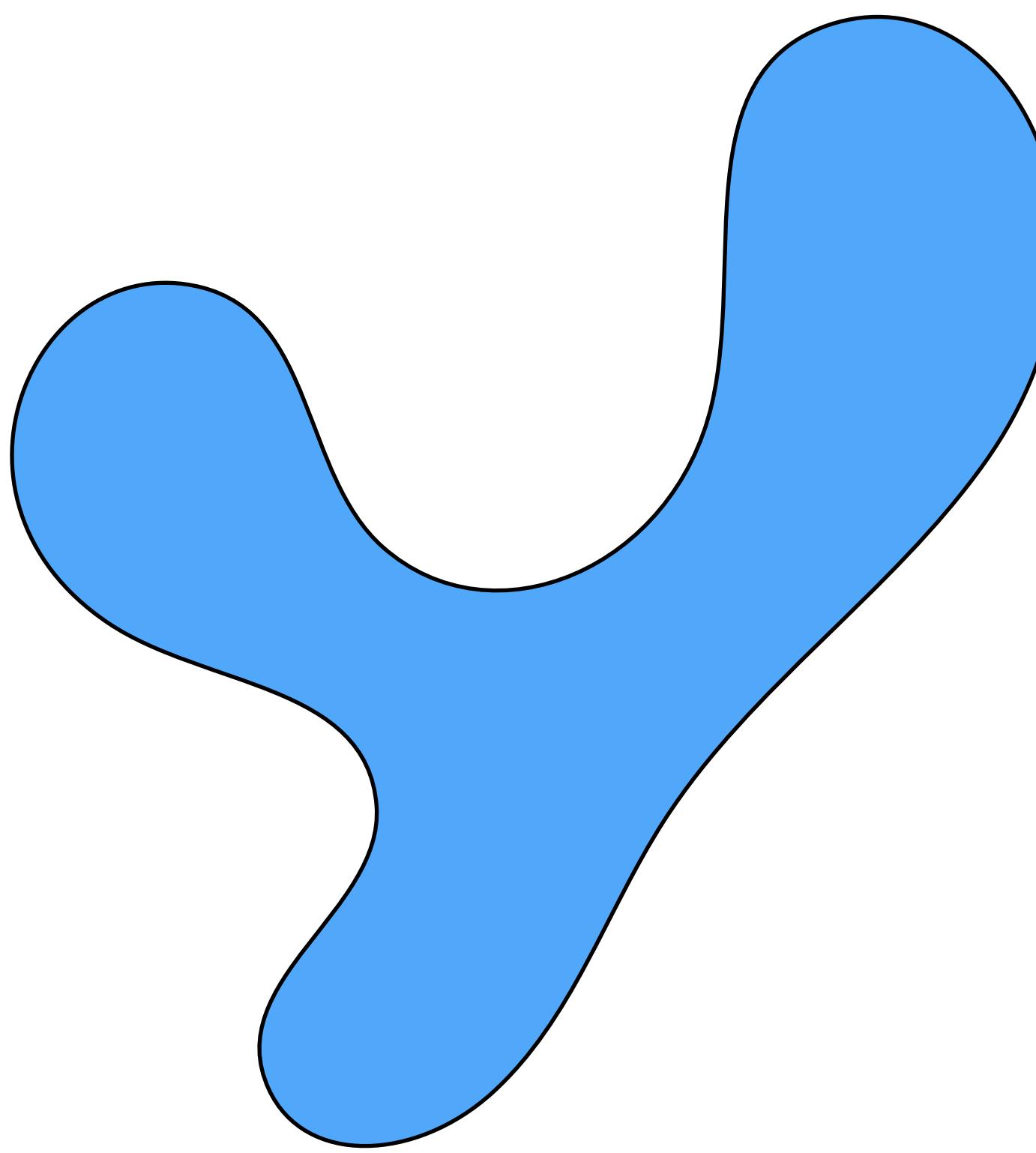
Horses



X

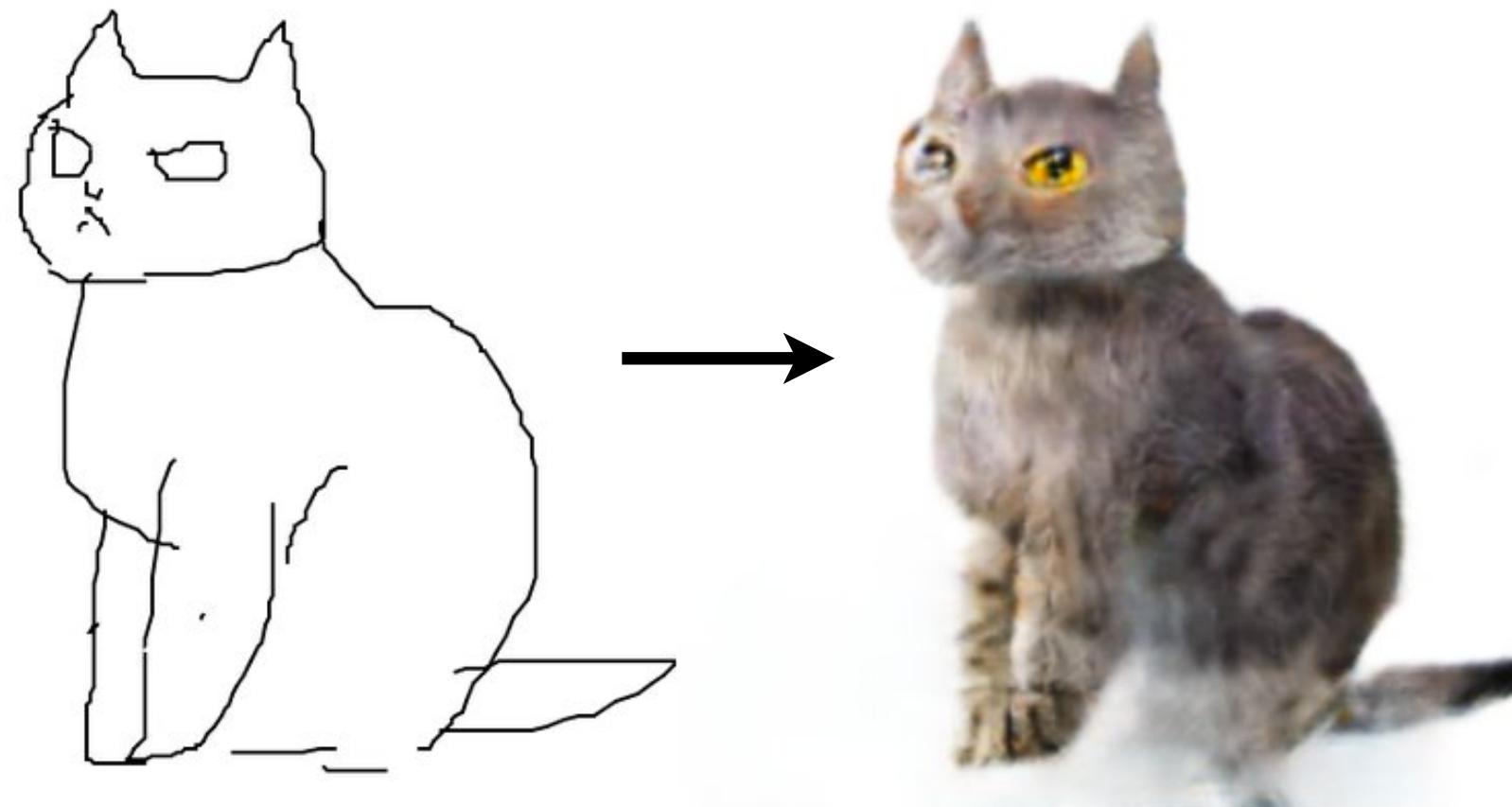


Zebras



Y

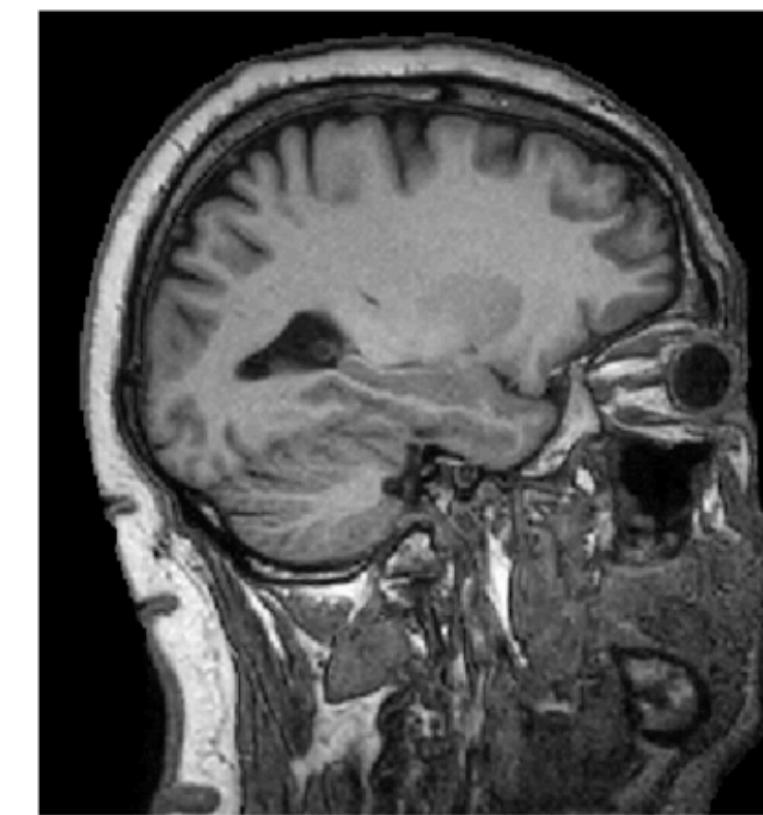
What would it look like if...?



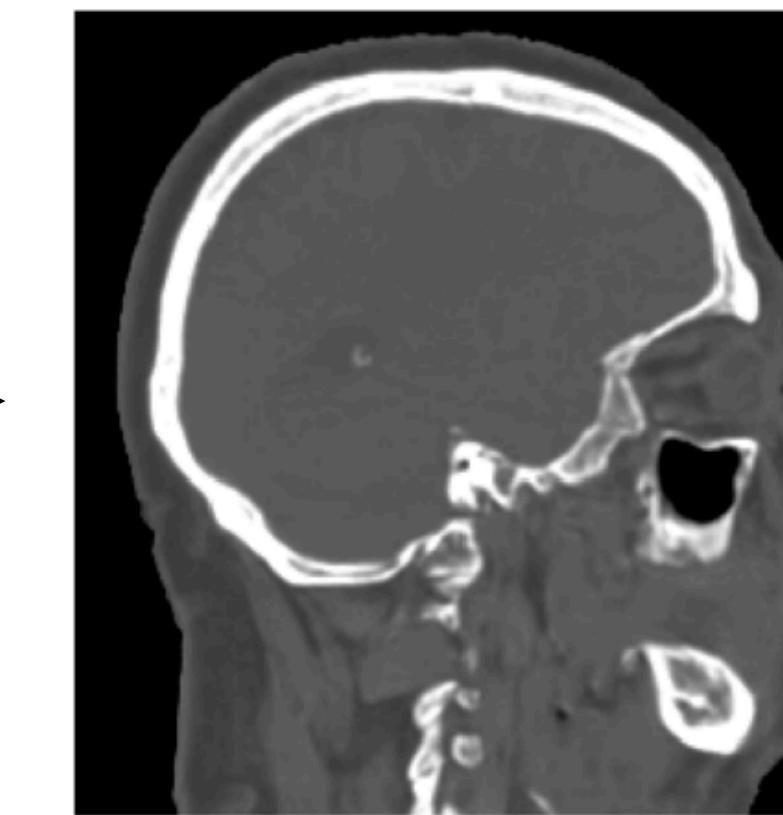
89

What would it look like if...?

MRI



CT



[Wolterink et al, 2017]

Sim



“Real”



90

[Hoffman et al, 2018]