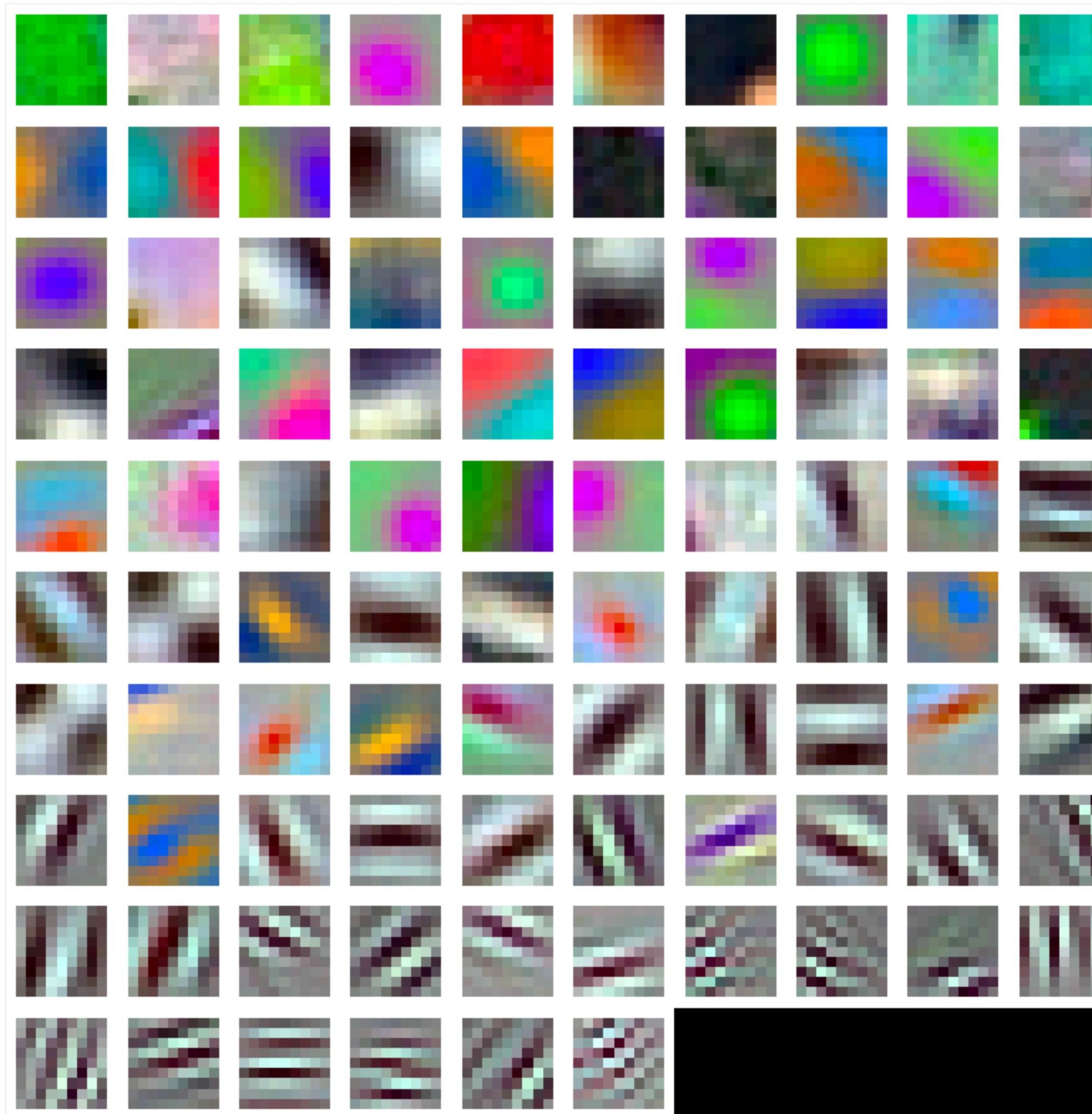


Lecture 13: GANs and autoregressive models

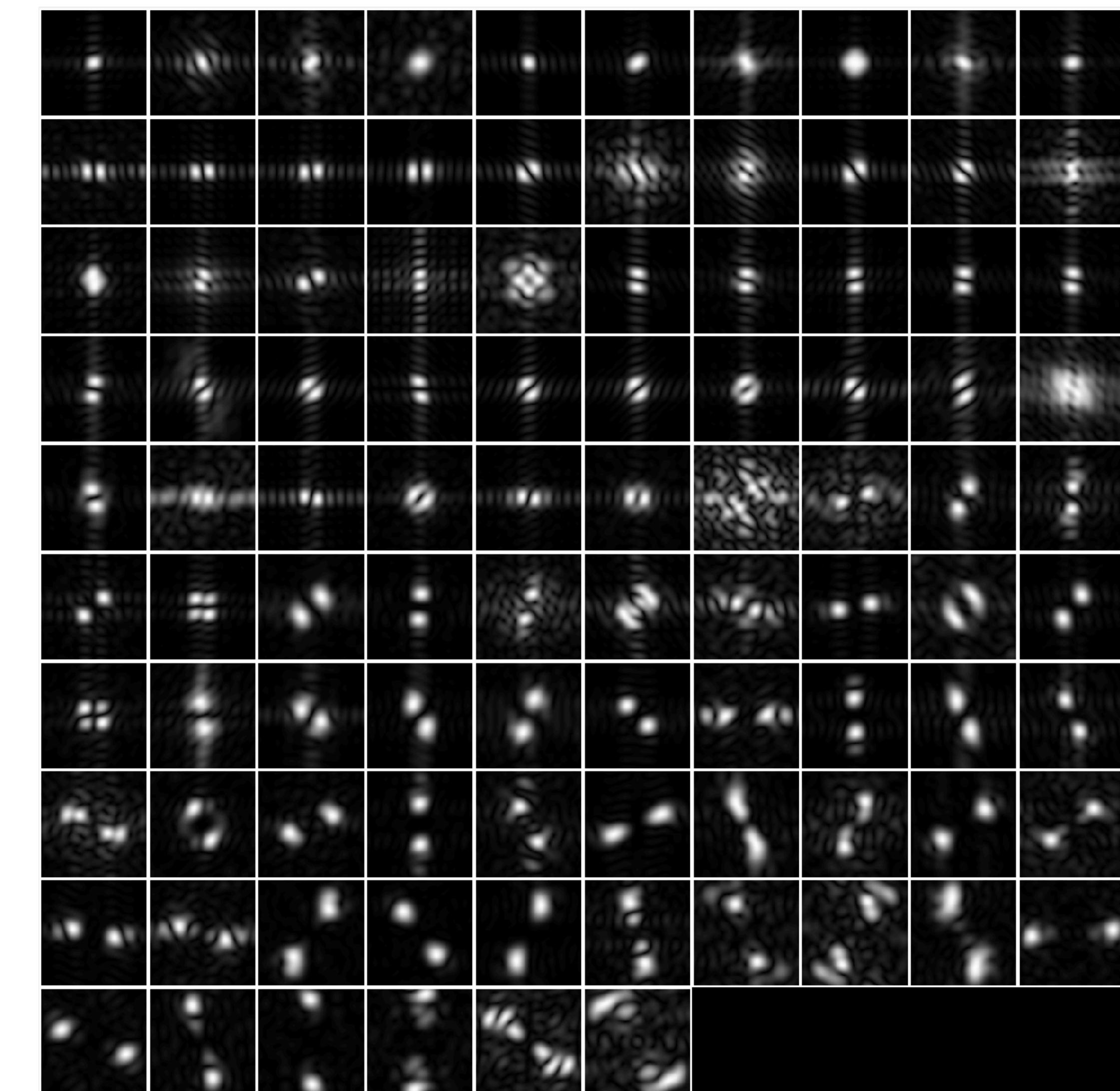
Today

- Texture synthesis
- Generative adversarial nets (GANs)
- Autoregressive models

Recall: neural net feature visualization



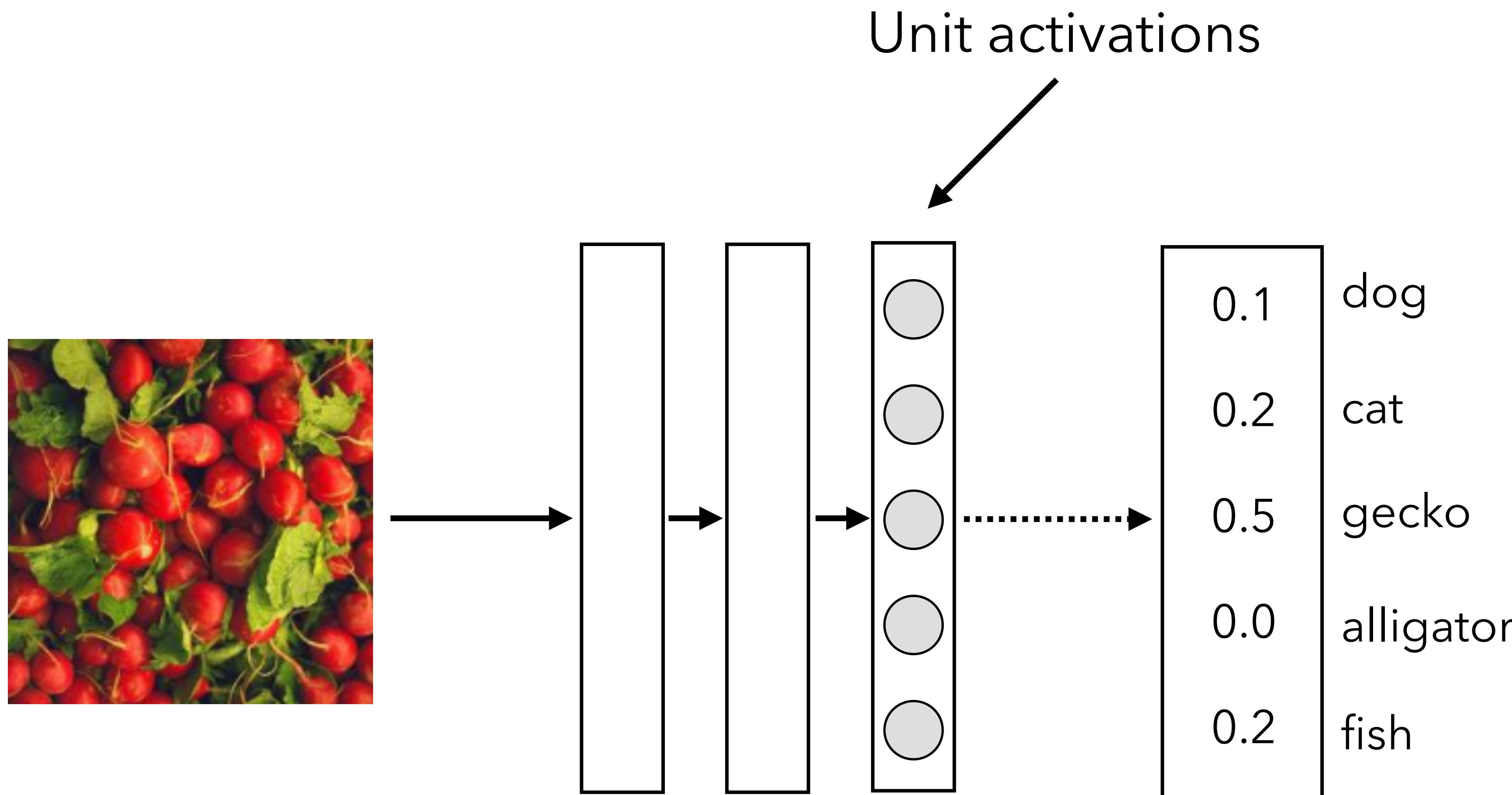
96 Units in conv1



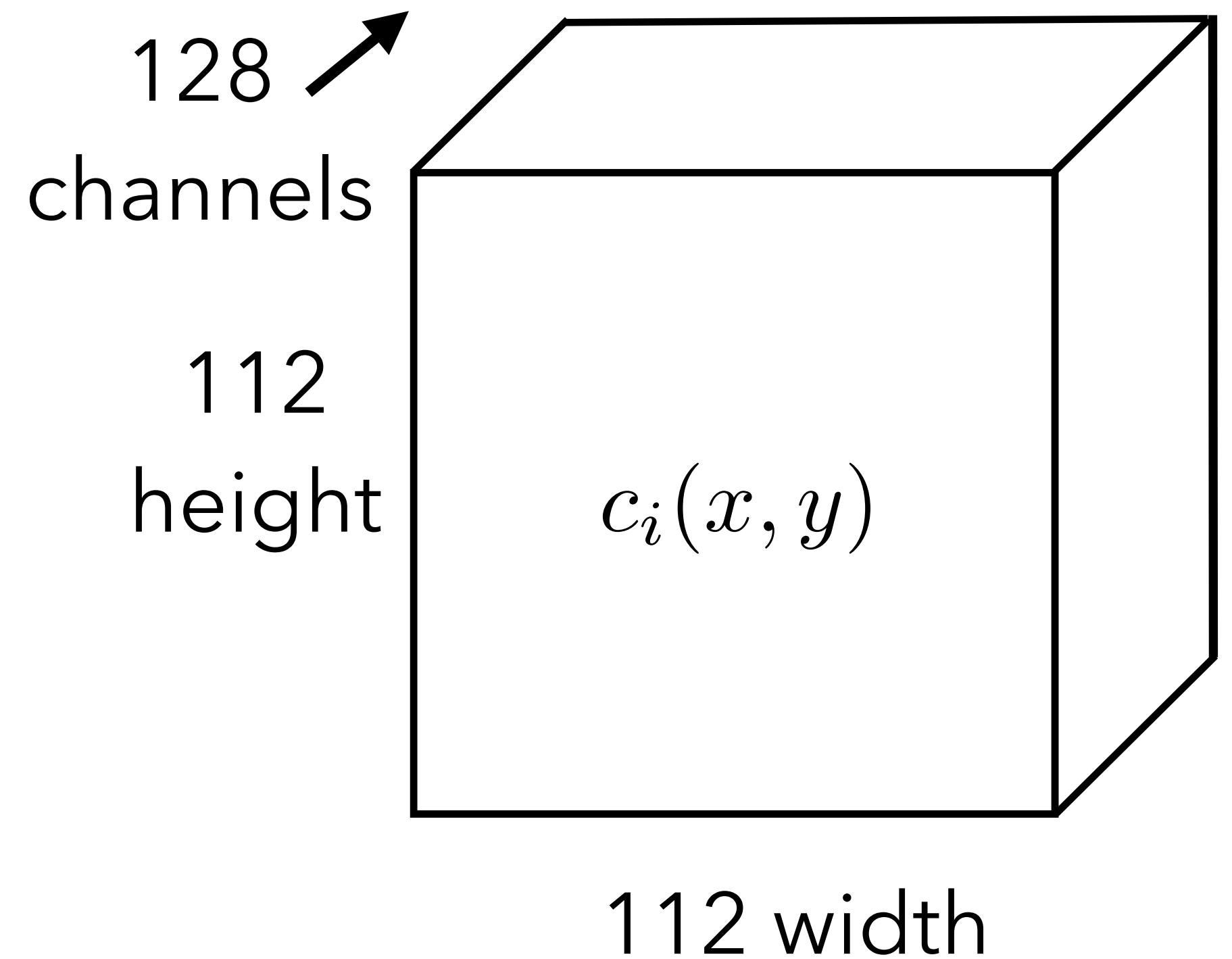
3

Source: Isola, Torralba, Freeman

Extracting features from a trained network



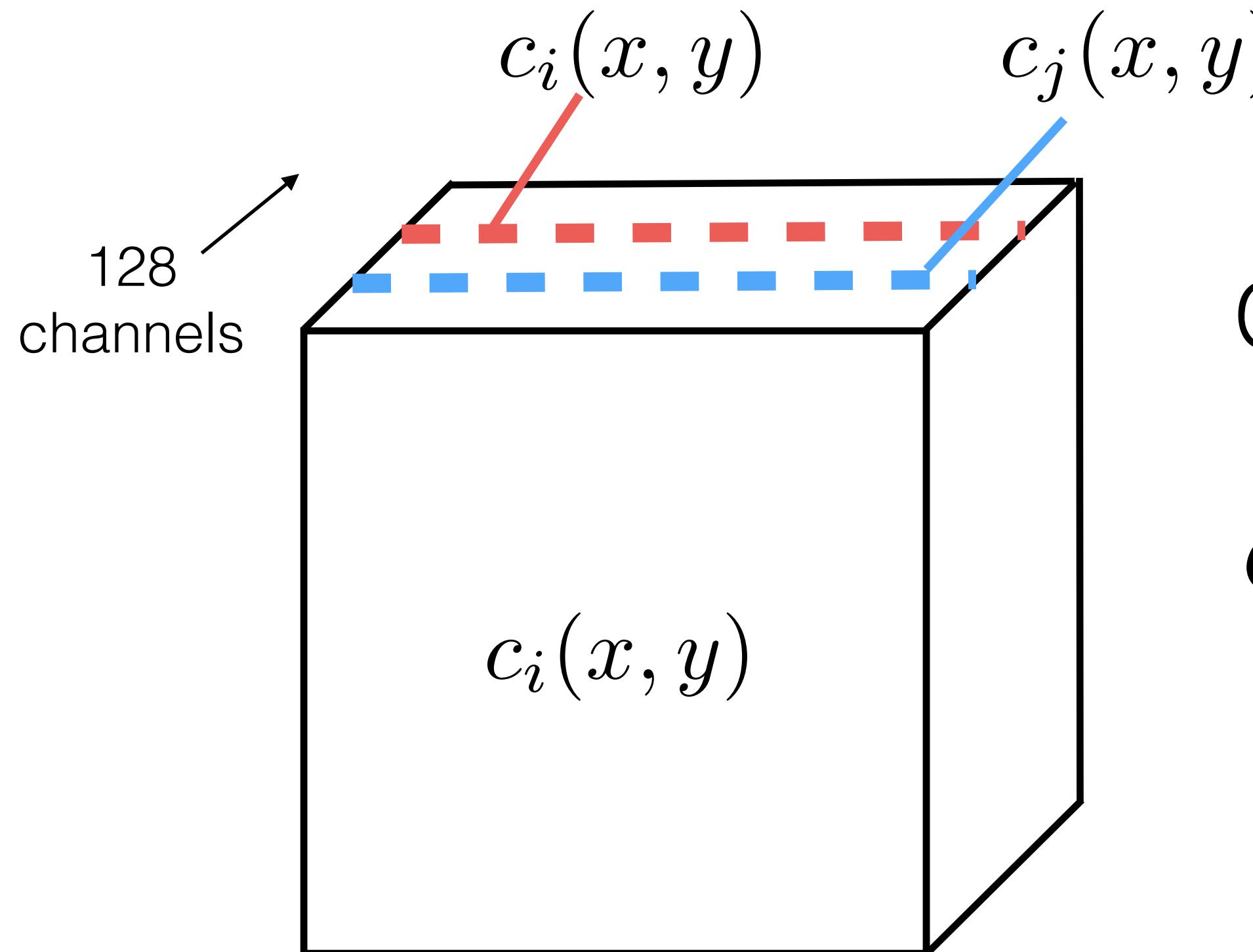
Extracting neural net features



conv2 feature activations

$(112 \times 112) \times 128$ for conv2 of VGG19

Capturing feature correlations



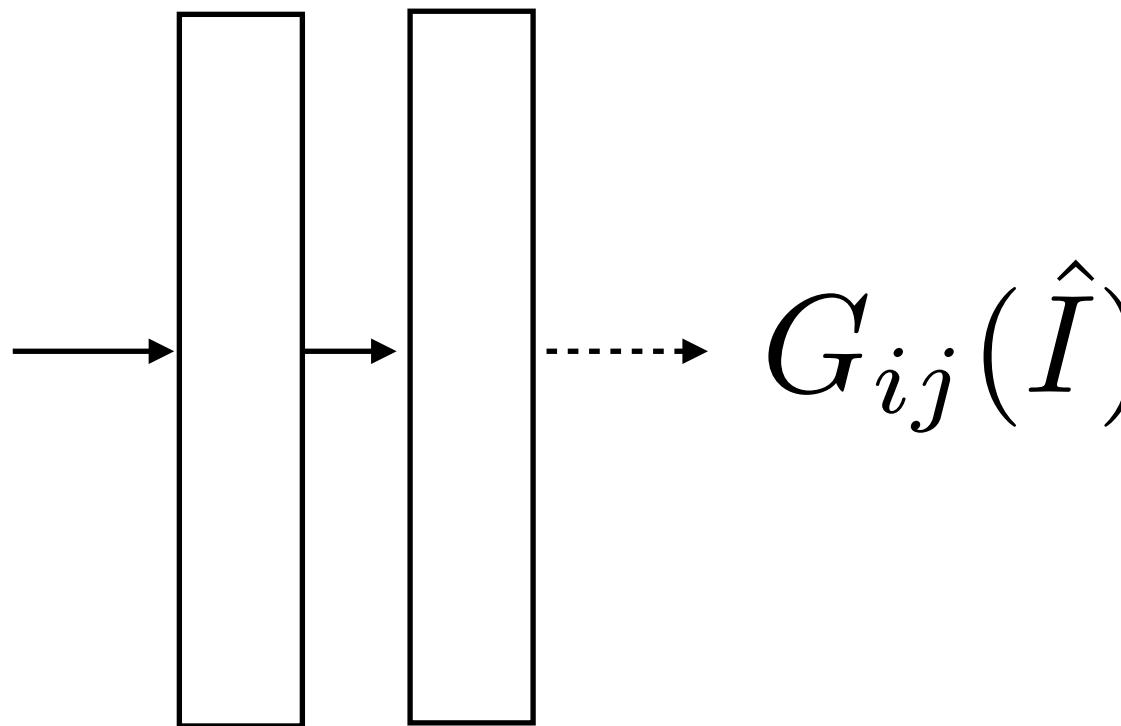
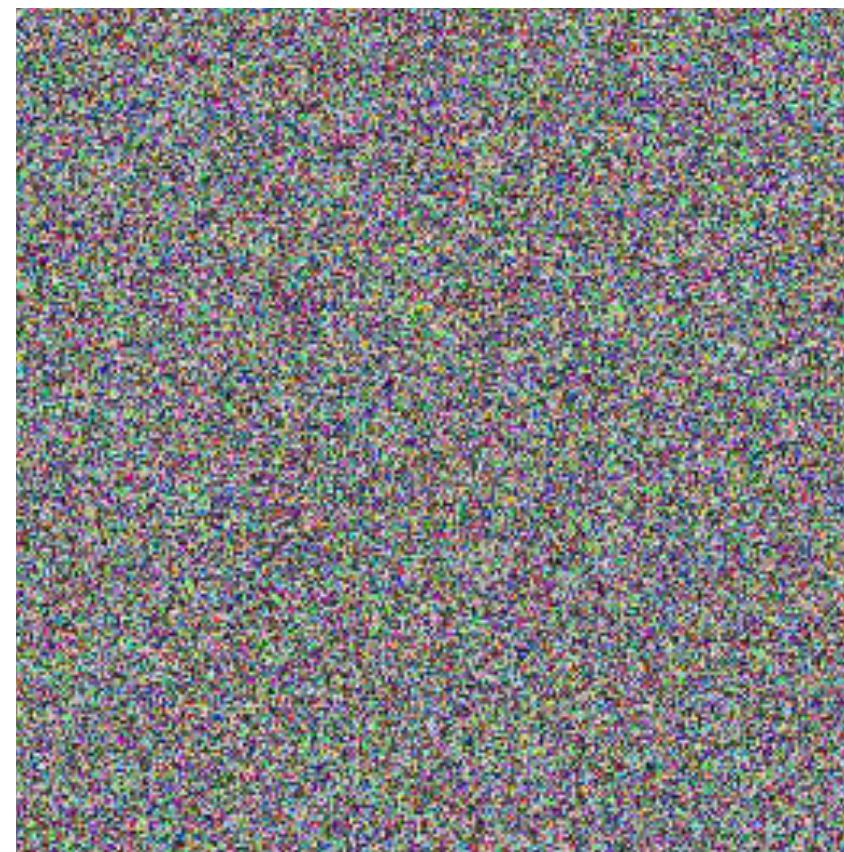
Gram (\approx covariance) matrix:

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

[Gatys et al. 2016]

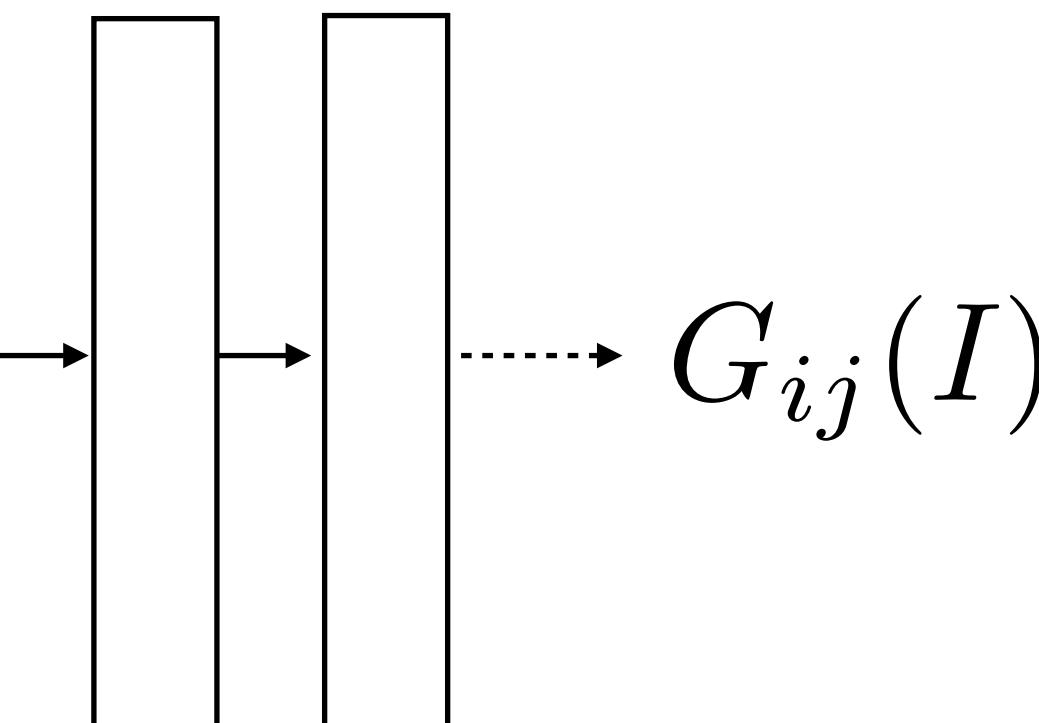
Idea: correlations between unit activations convey texture.
Discard global spatial information.

Matching image statistics



Find \hat{I} by minimizing:

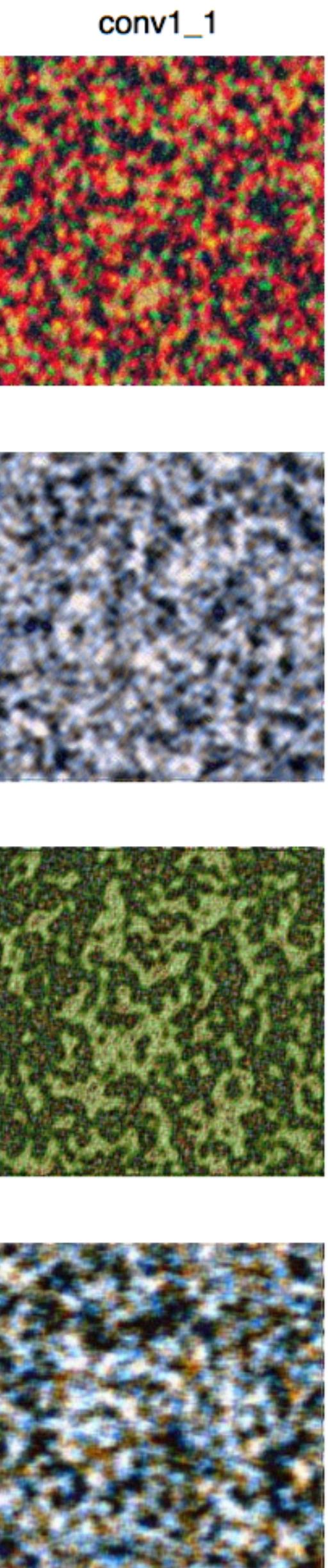
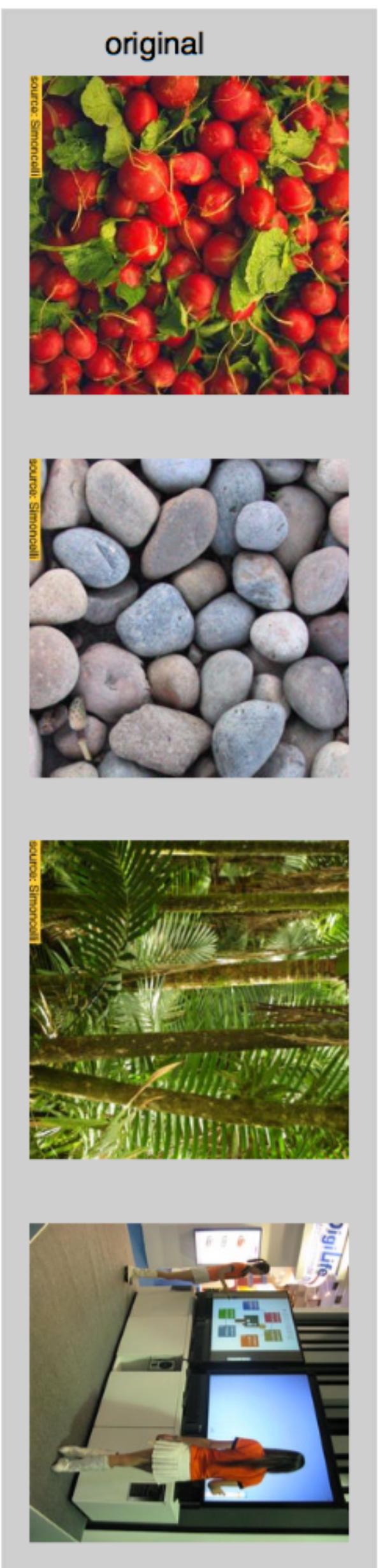
$$\sum_{i=1}^{128} \sum_{j=1}^{128} (G_{ij}(I) - G_{ij}(\hat{I}))^2$$



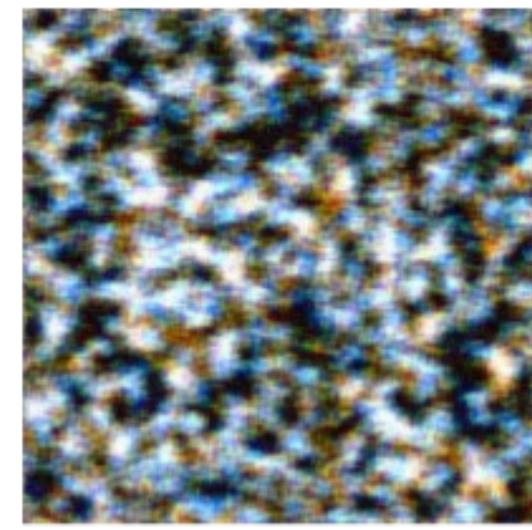
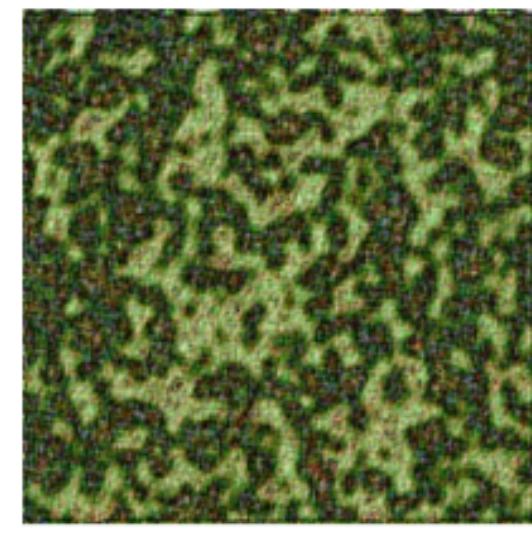
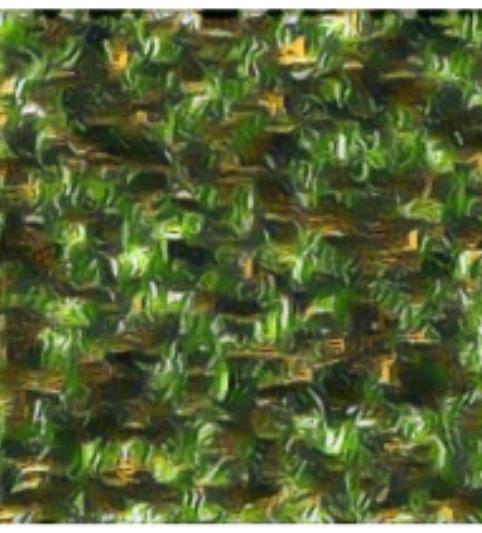
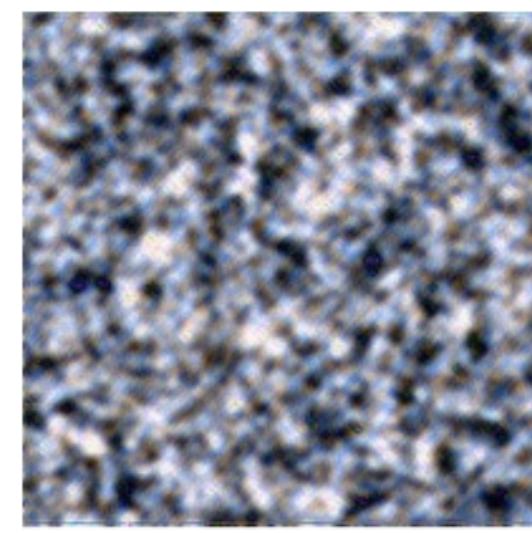
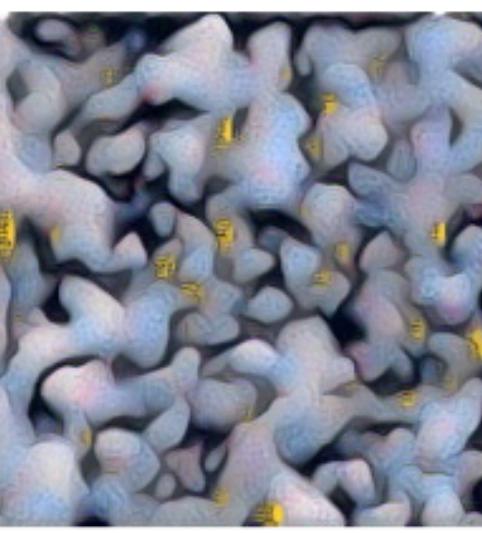
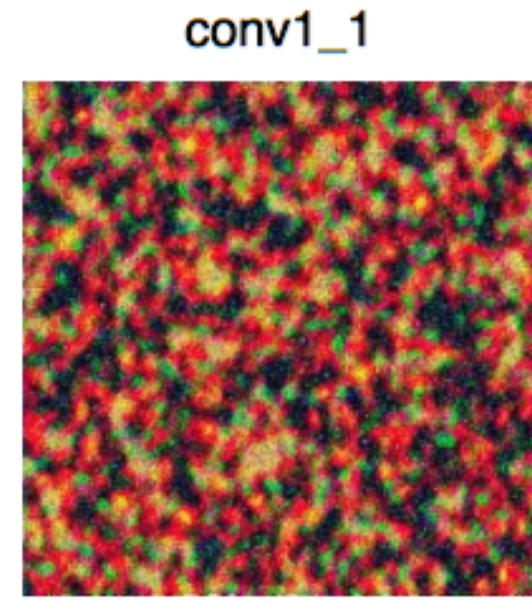
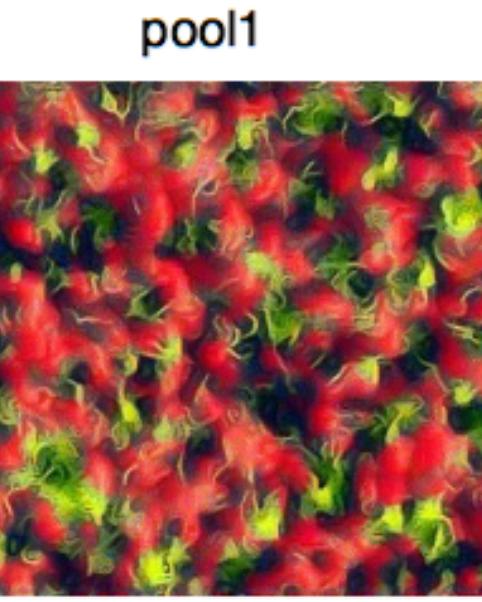
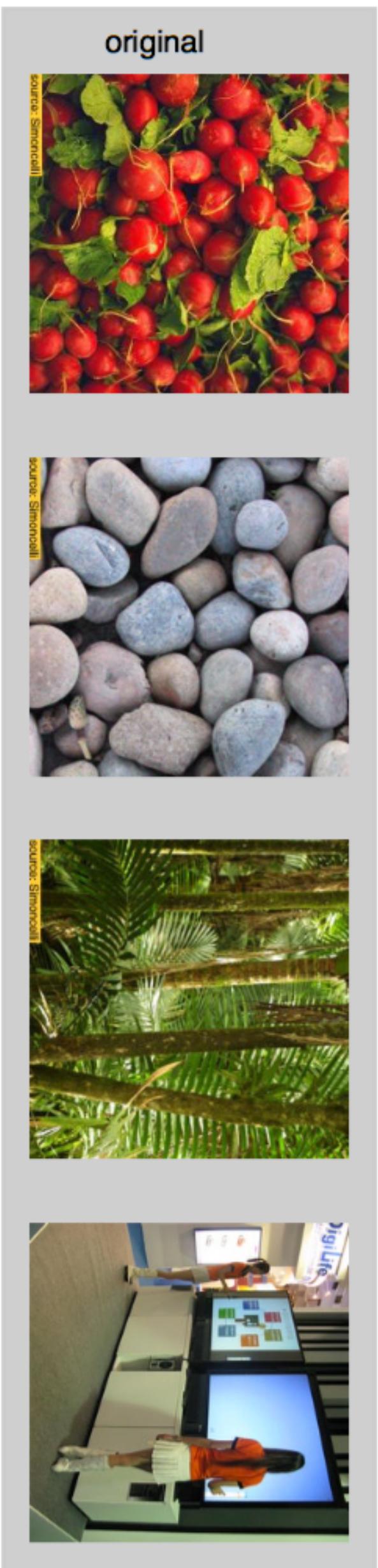
Implementation details:

- Minimize with gradient descent
- Use many layers of network

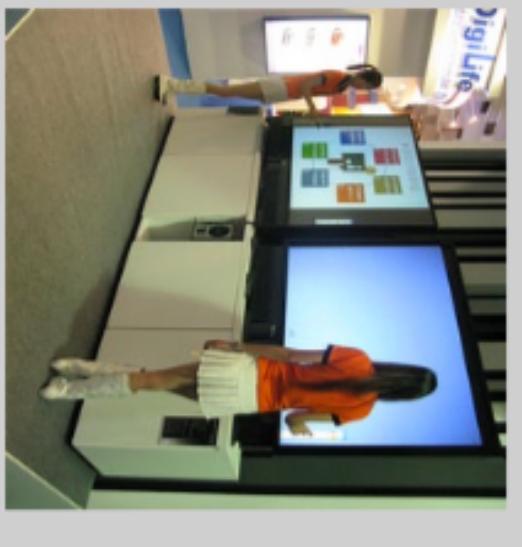
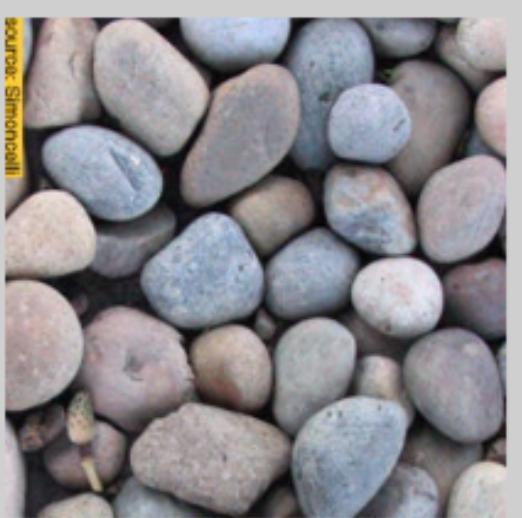
High → Low



High → Low



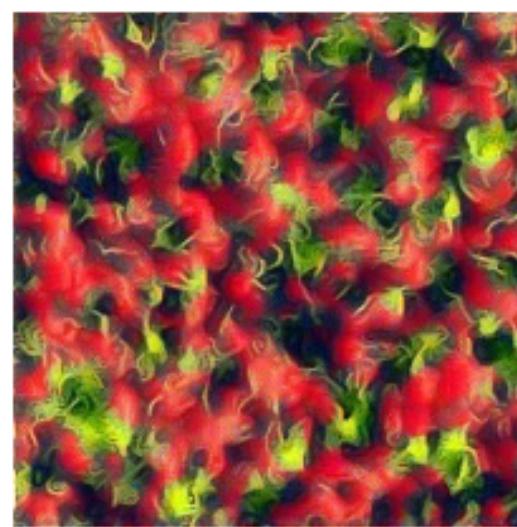
High → Low



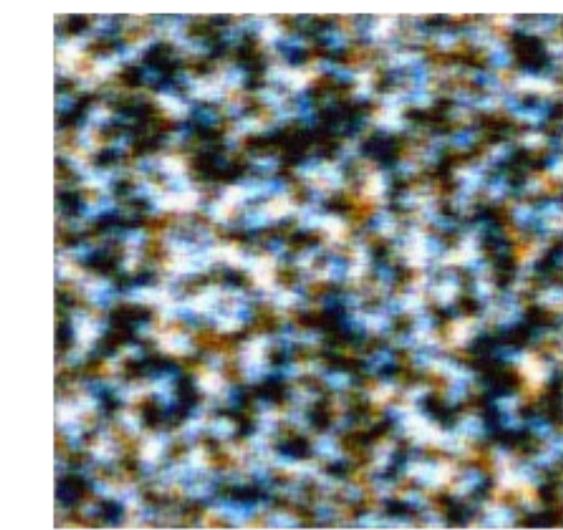
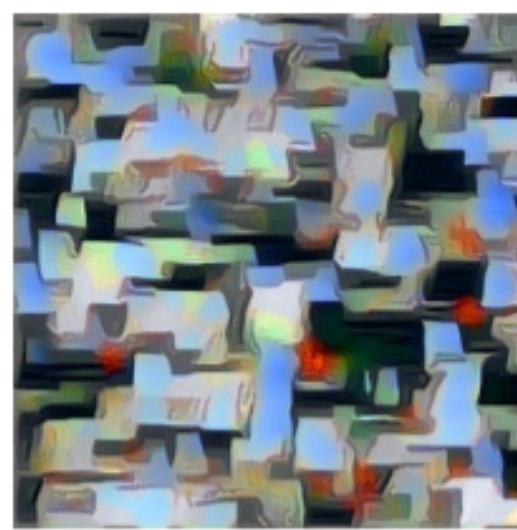
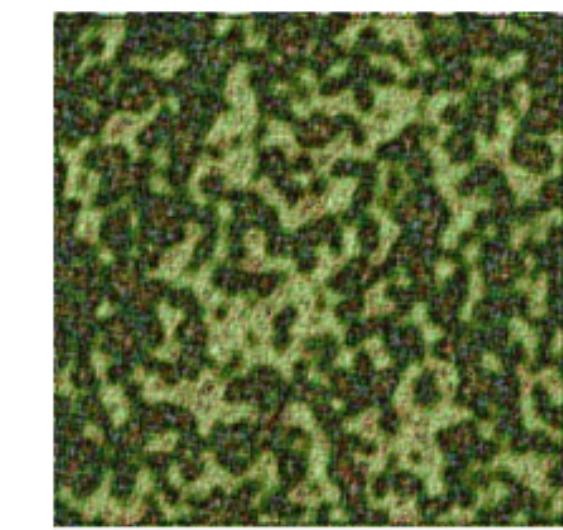
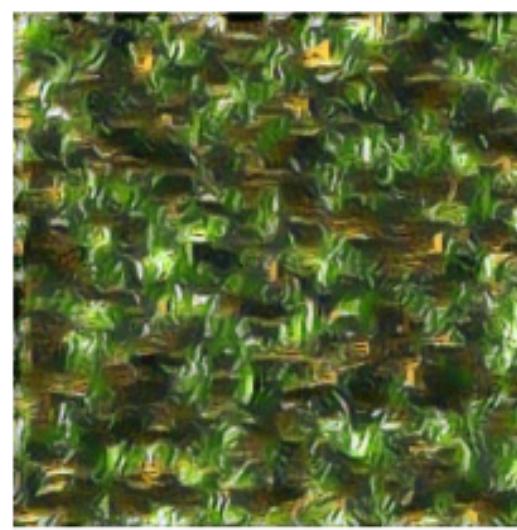
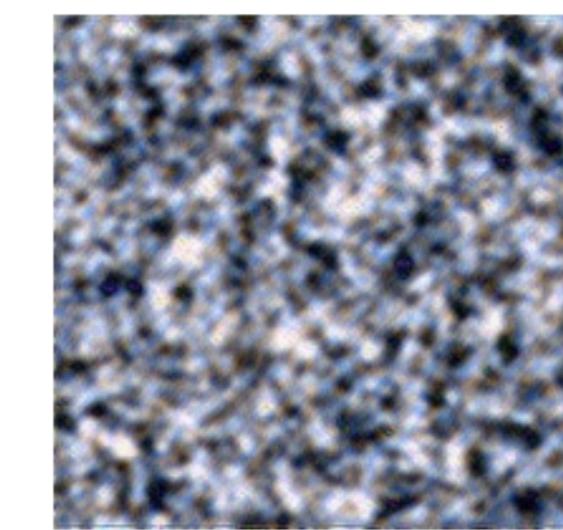
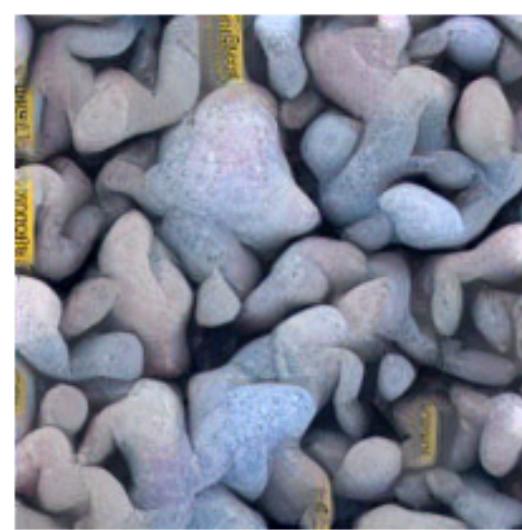
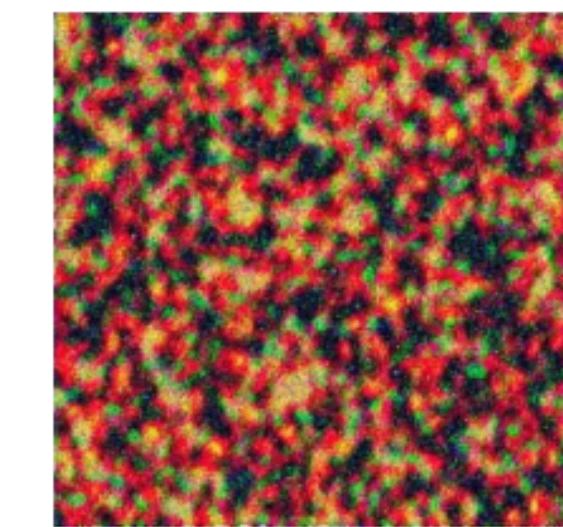
pool2



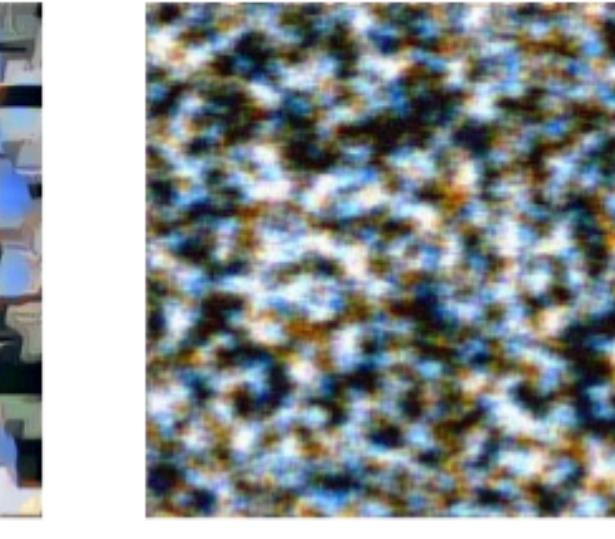
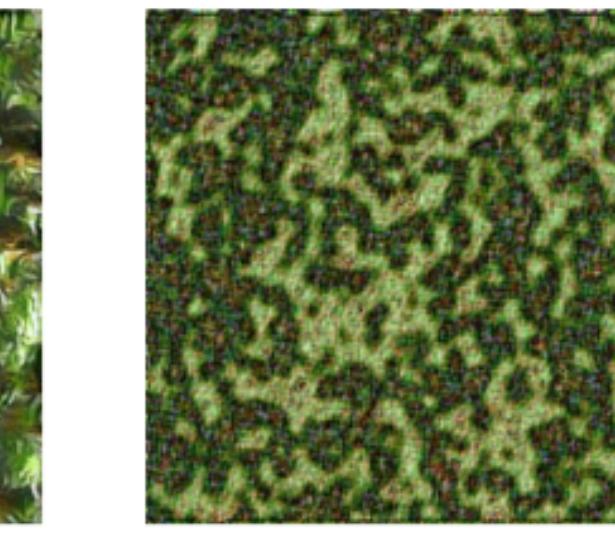
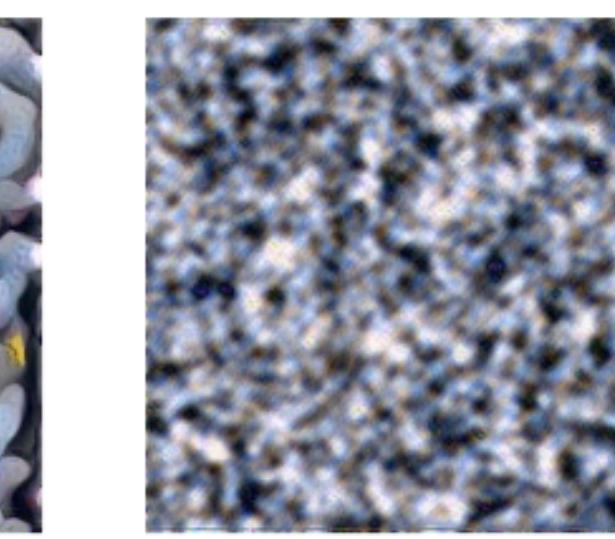
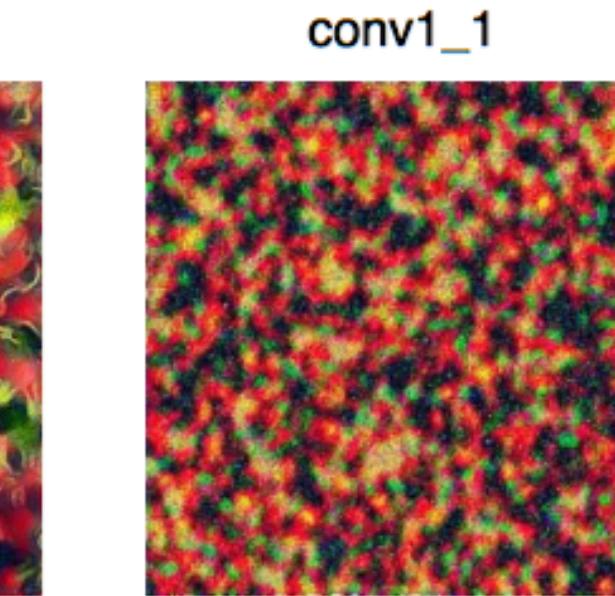
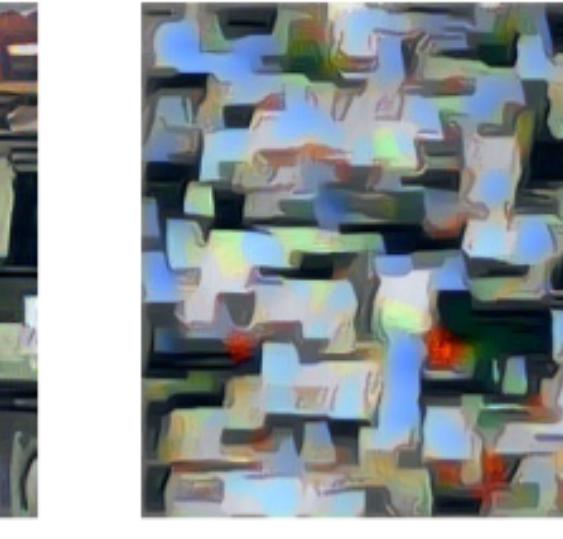
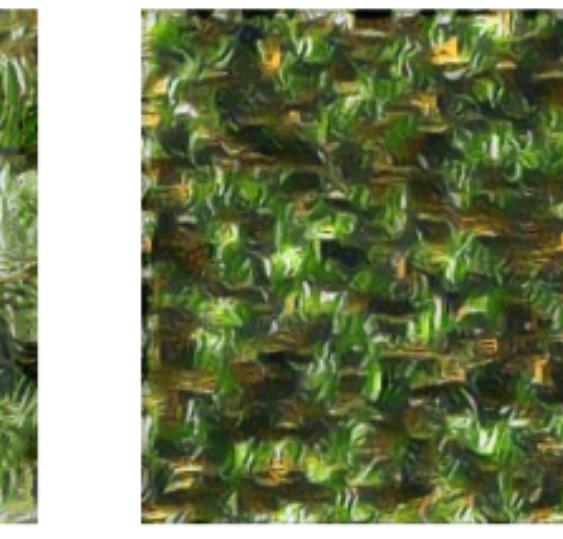
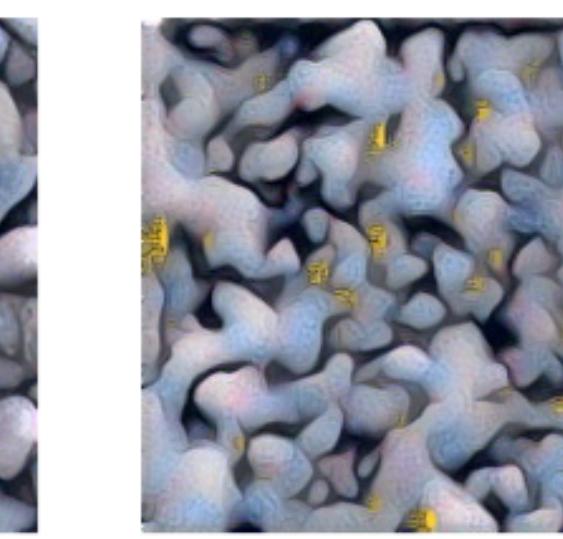
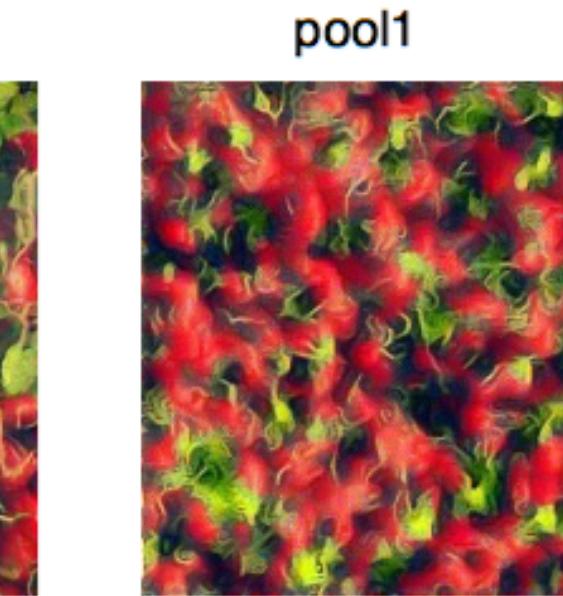
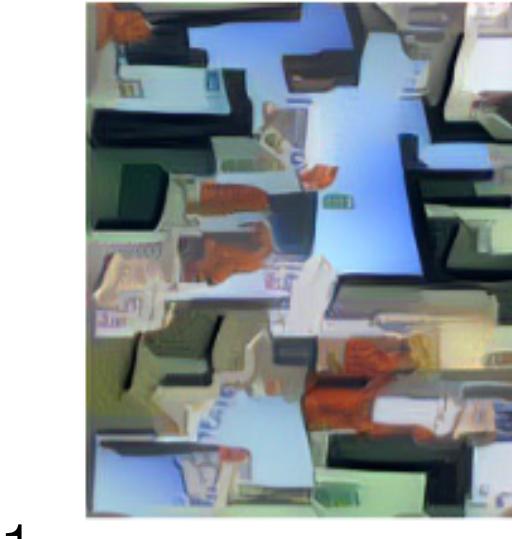
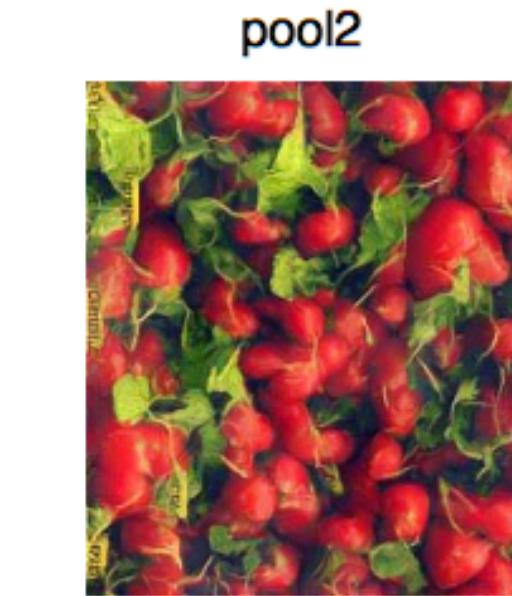
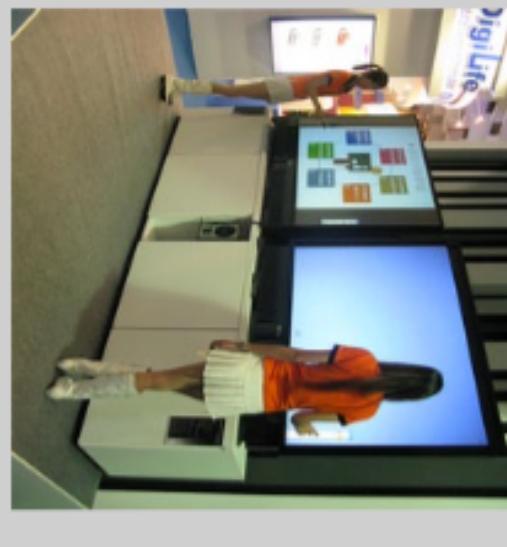
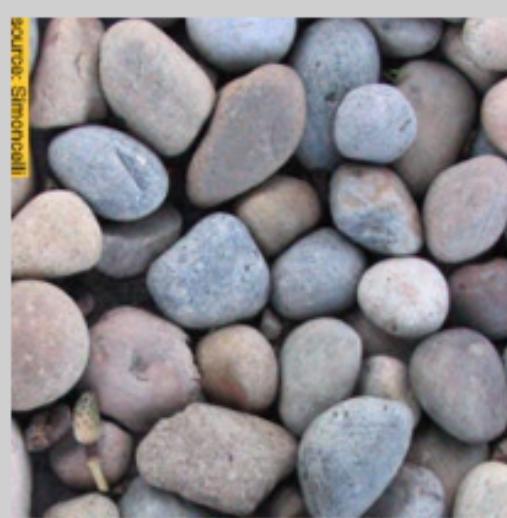
pool1



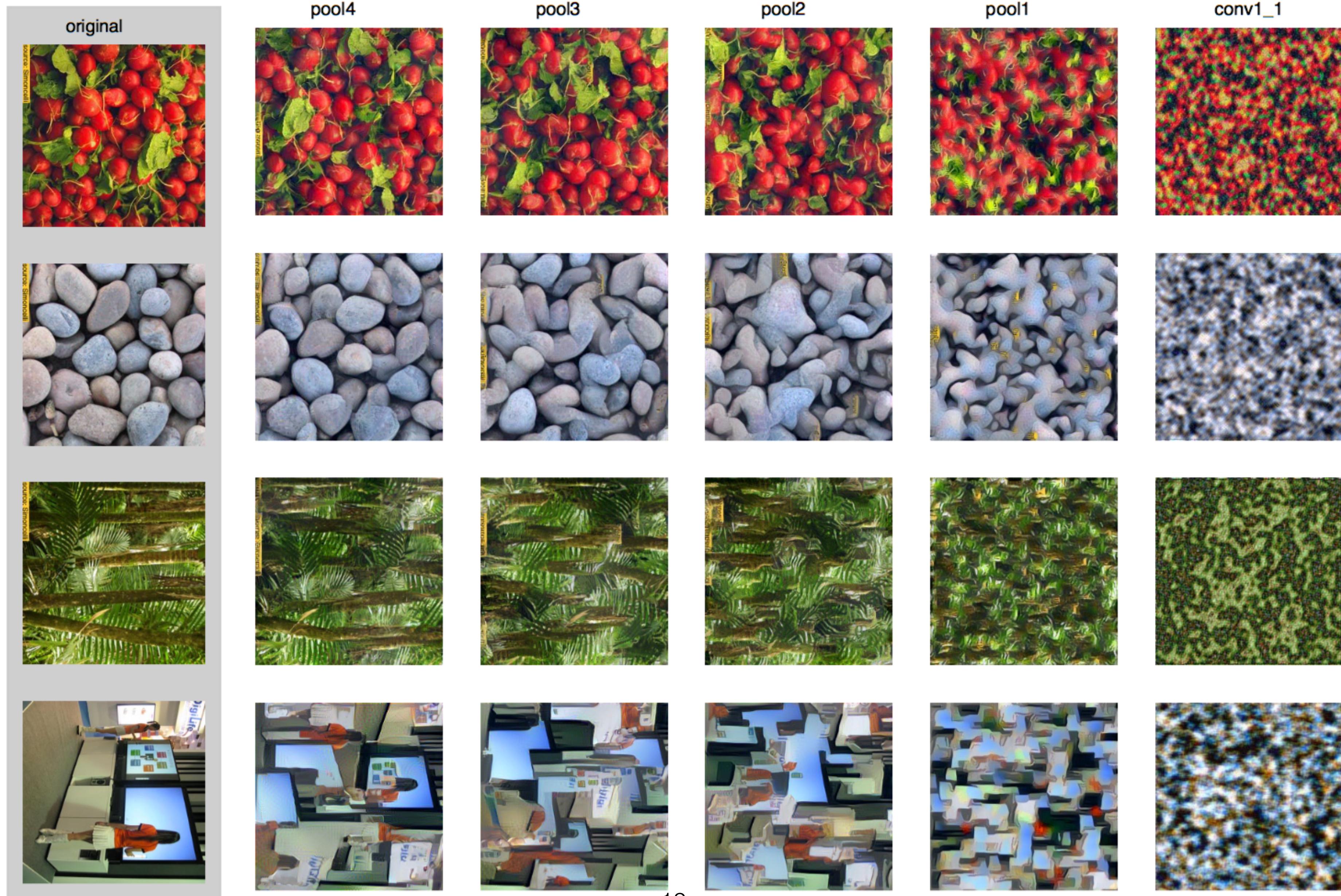
conv1_1

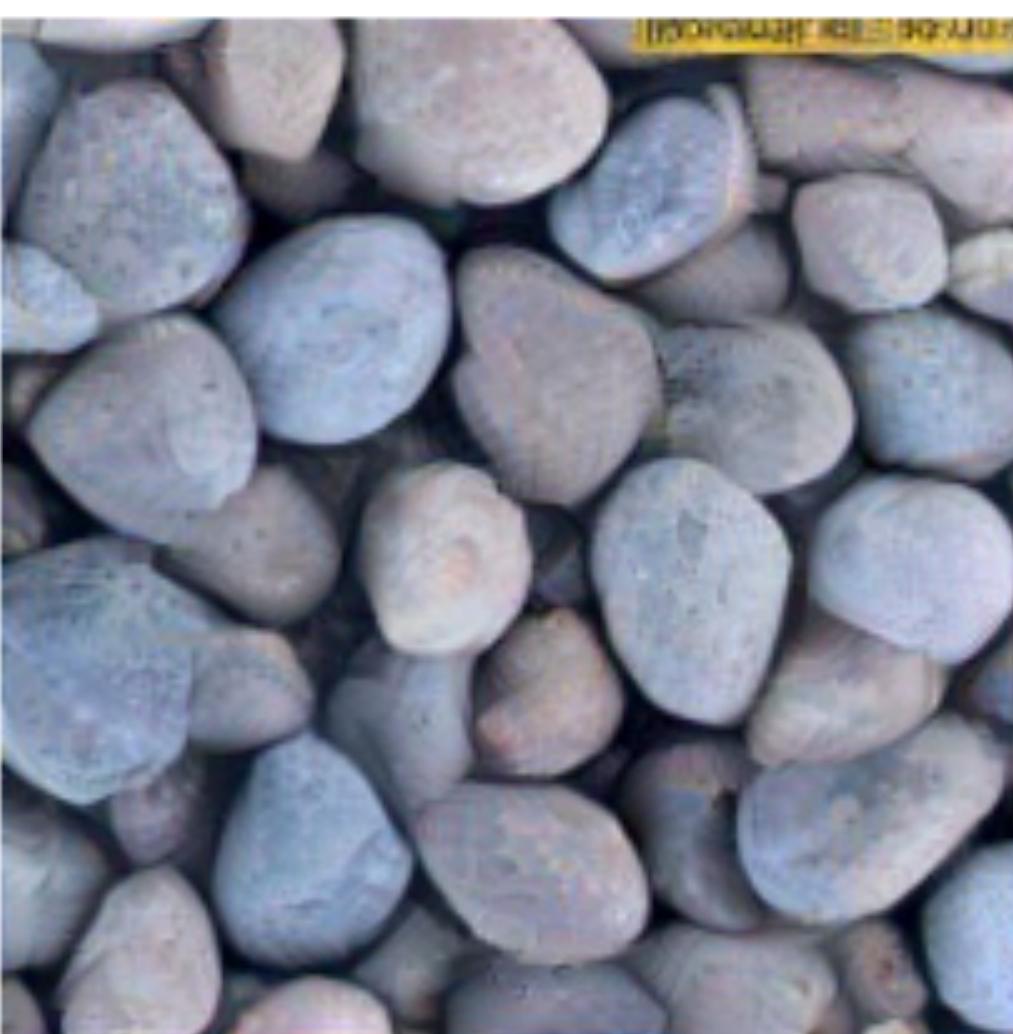
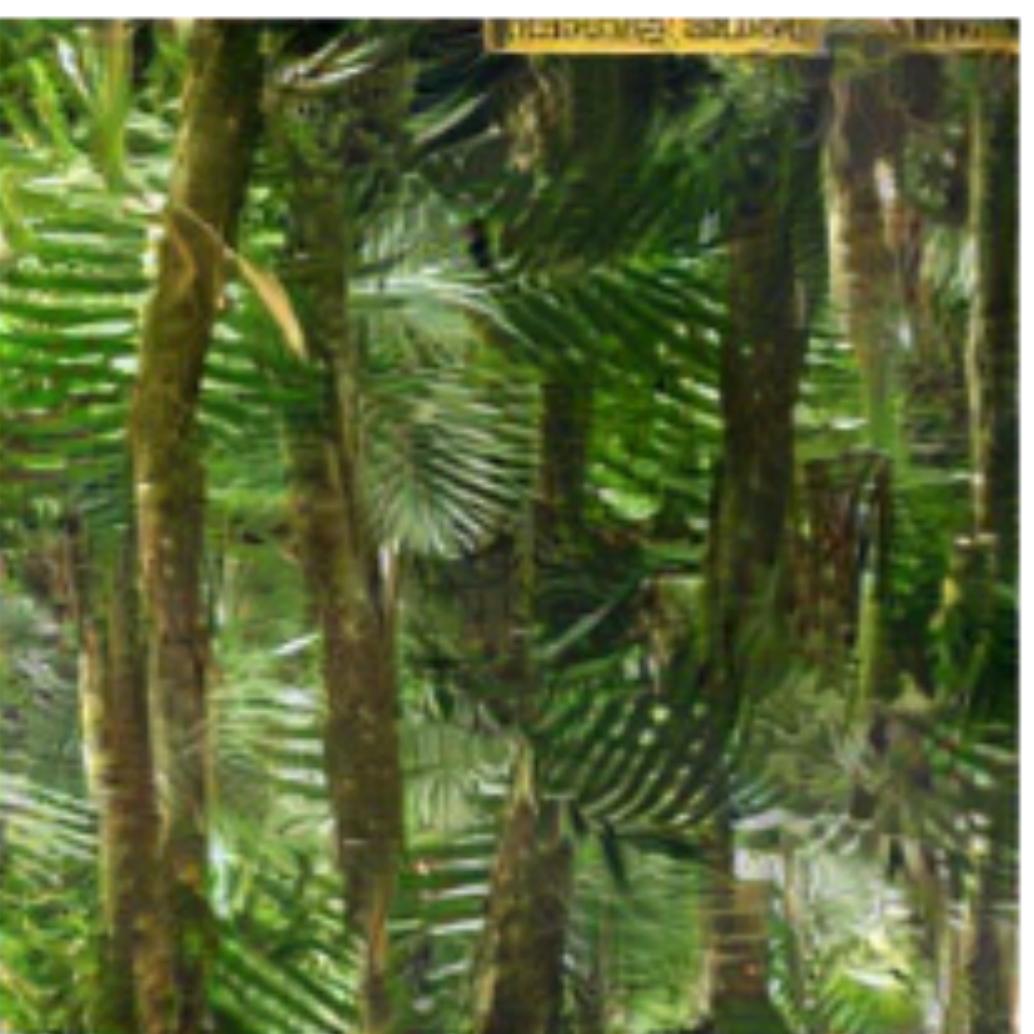


High → Low

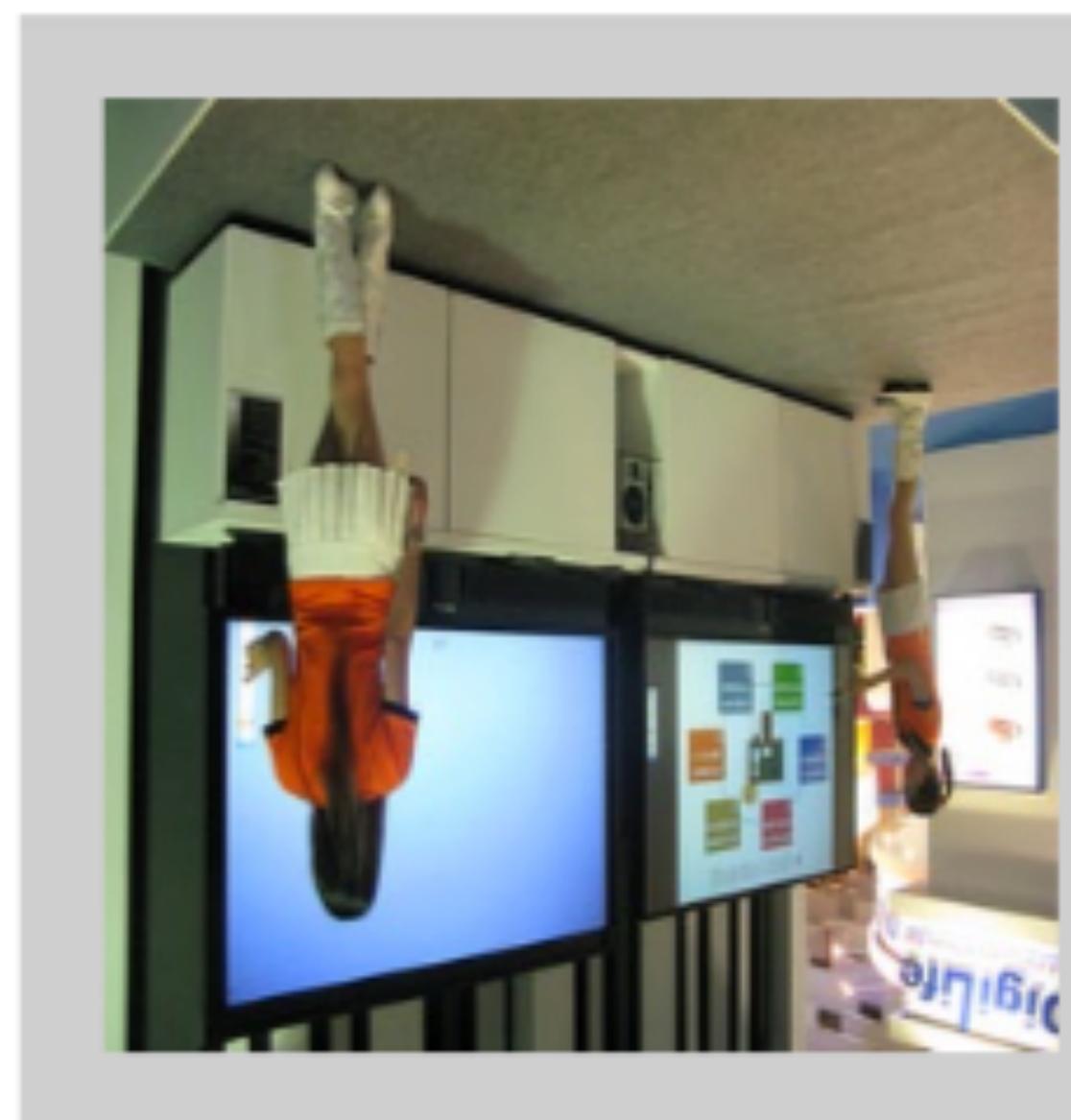


High → Low





pool4



original

Texture captures artistic style

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

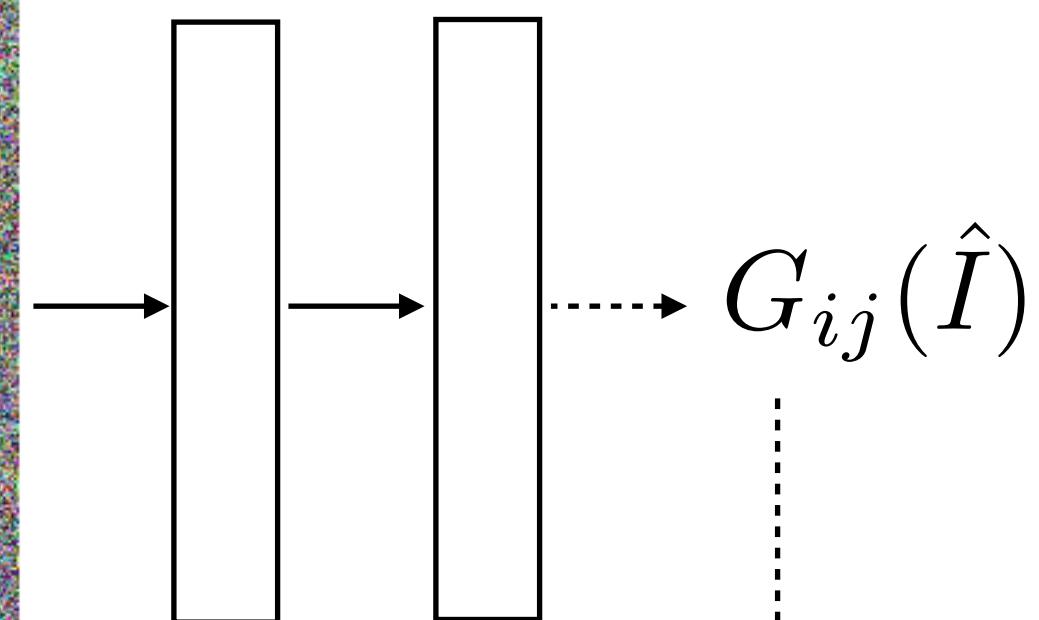


[Gatys et al. 2016]

Match the **style** of the painting.



Synthesized image



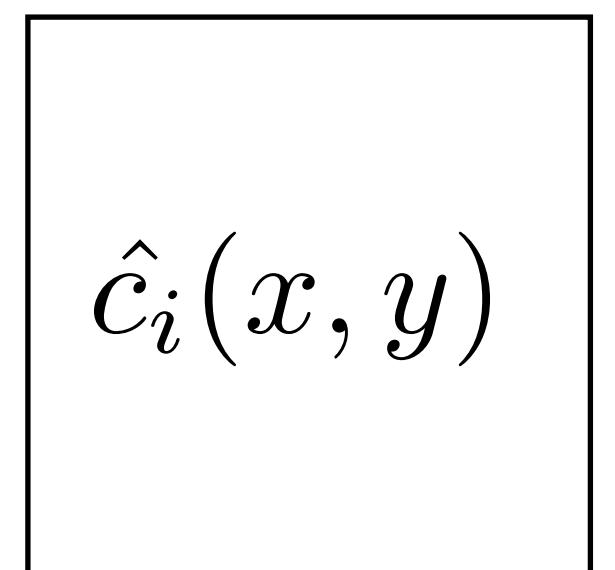
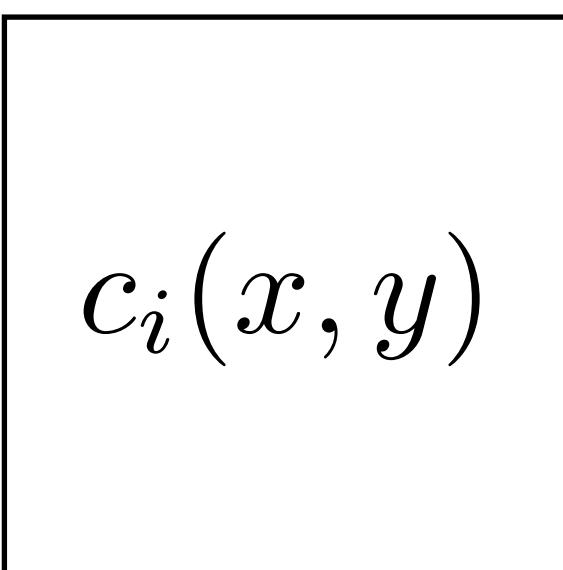
$$\sum_{i=1}^{128} \sum_{j=1}^{128} (G_{ij}(\hat{I}) - G_{ij}(I))^2$$

Perceptual loss:

usually distance in feature space

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

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













London during the day.

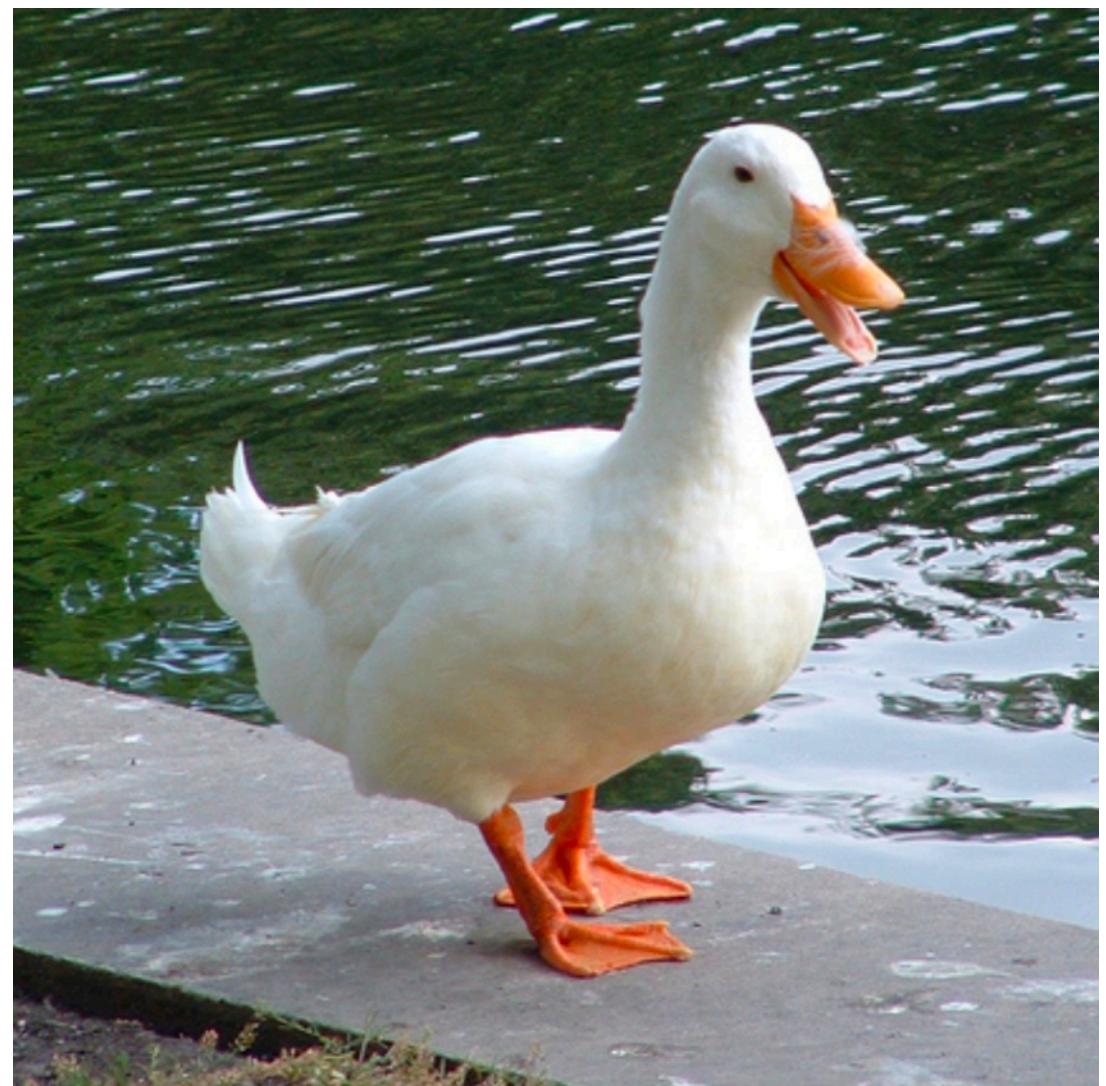


New York at night.



Neural networks that generate images

Image classification



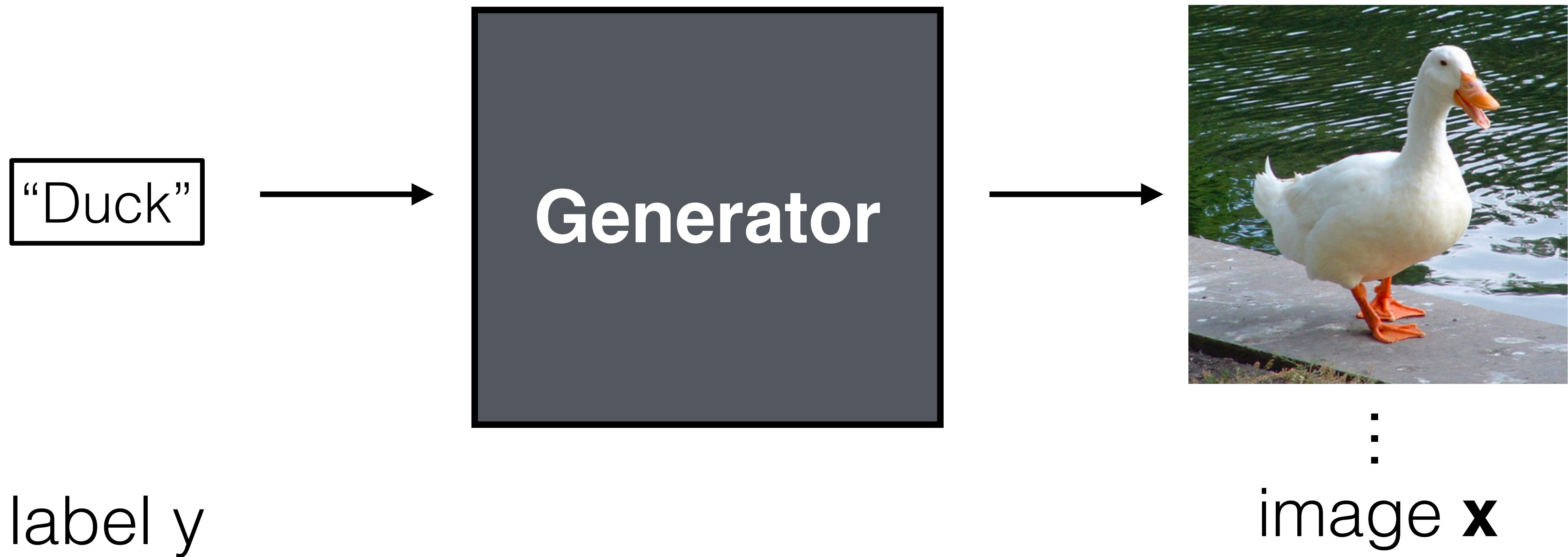
“Duck”

:

image **x**

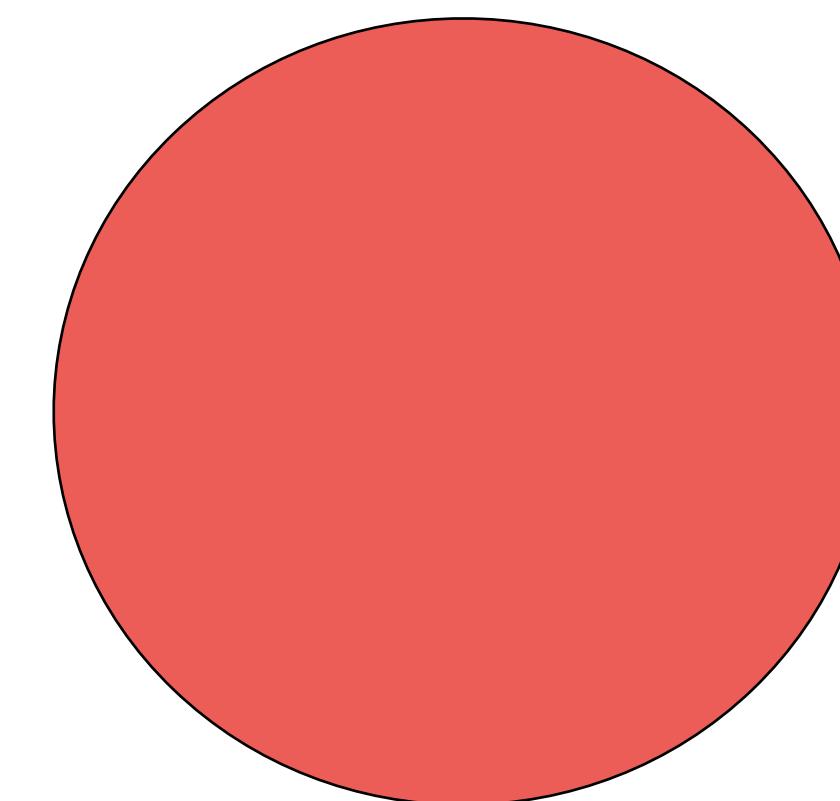
label **y**

Image synthesis



Neural networks as distribution transformers

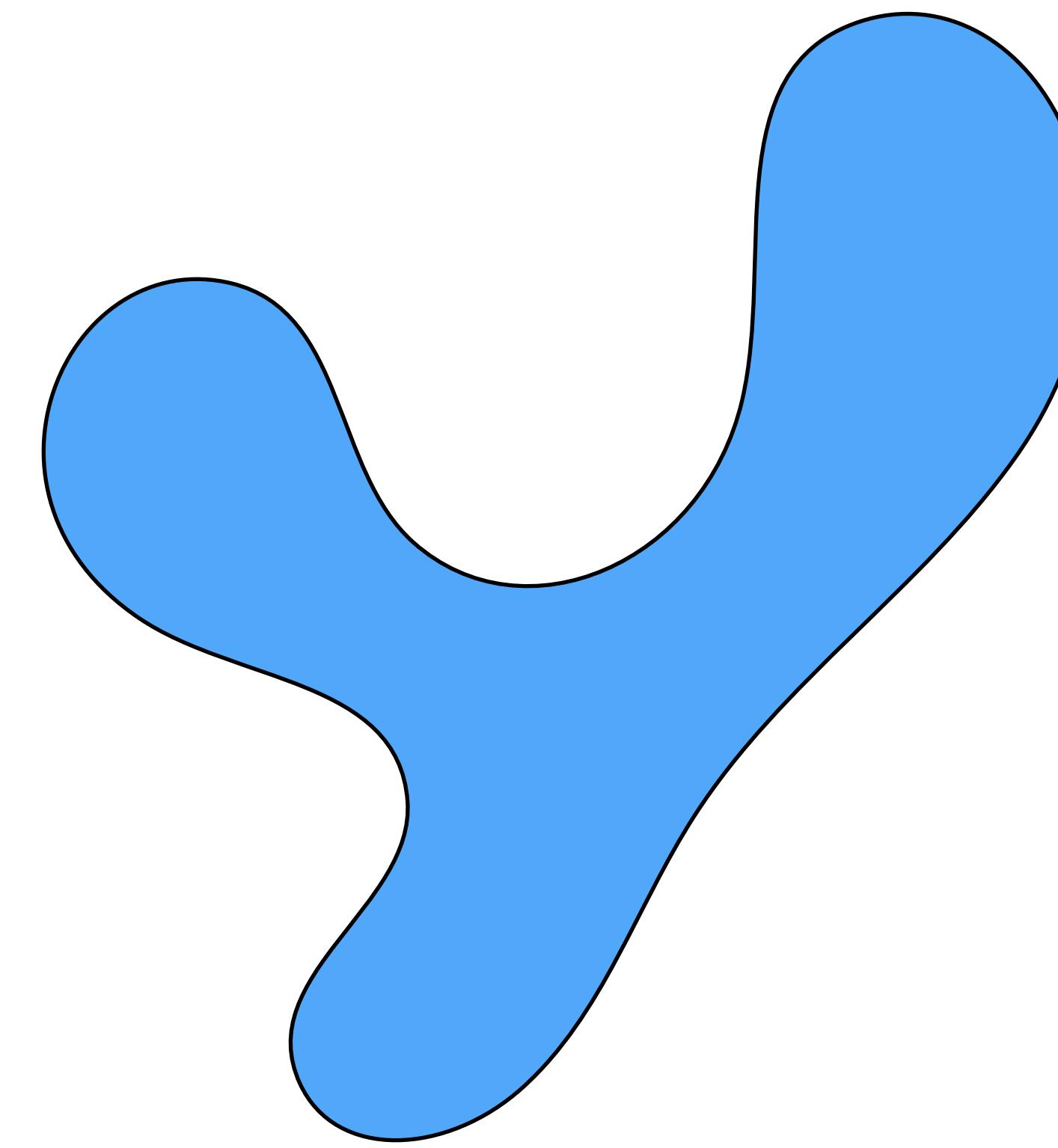
Source distribution



$$p(z)$$

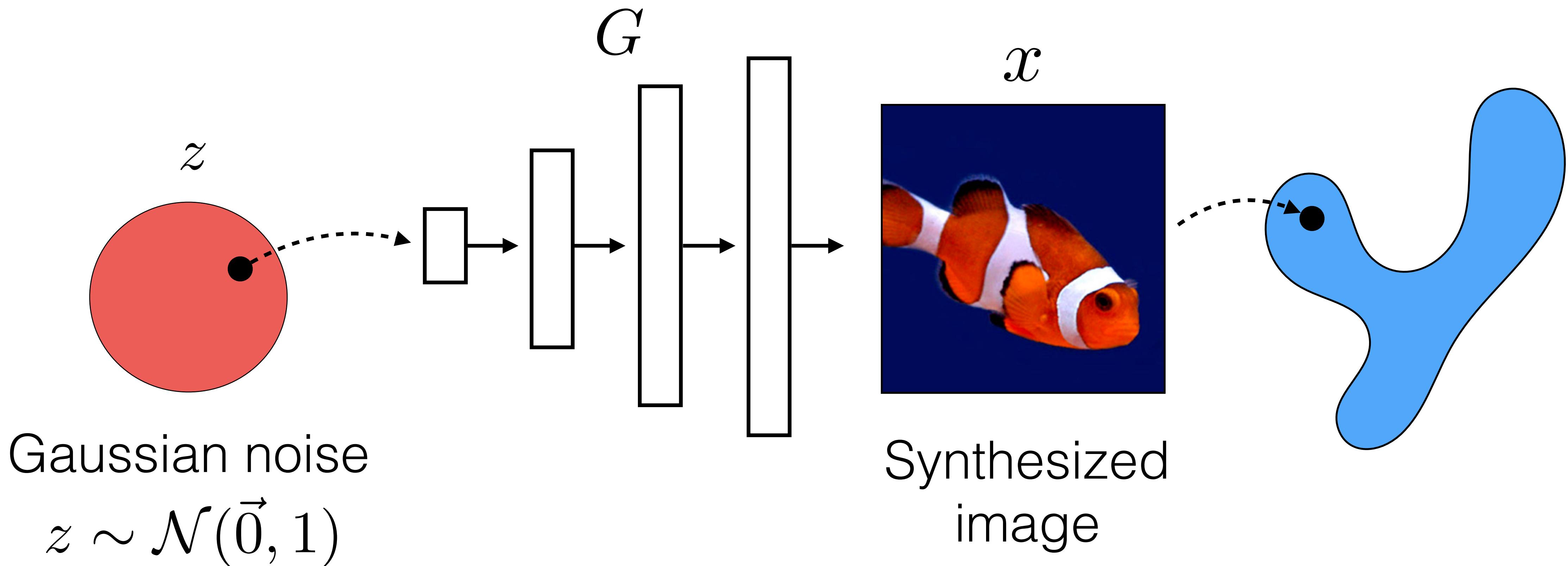


Target distribution

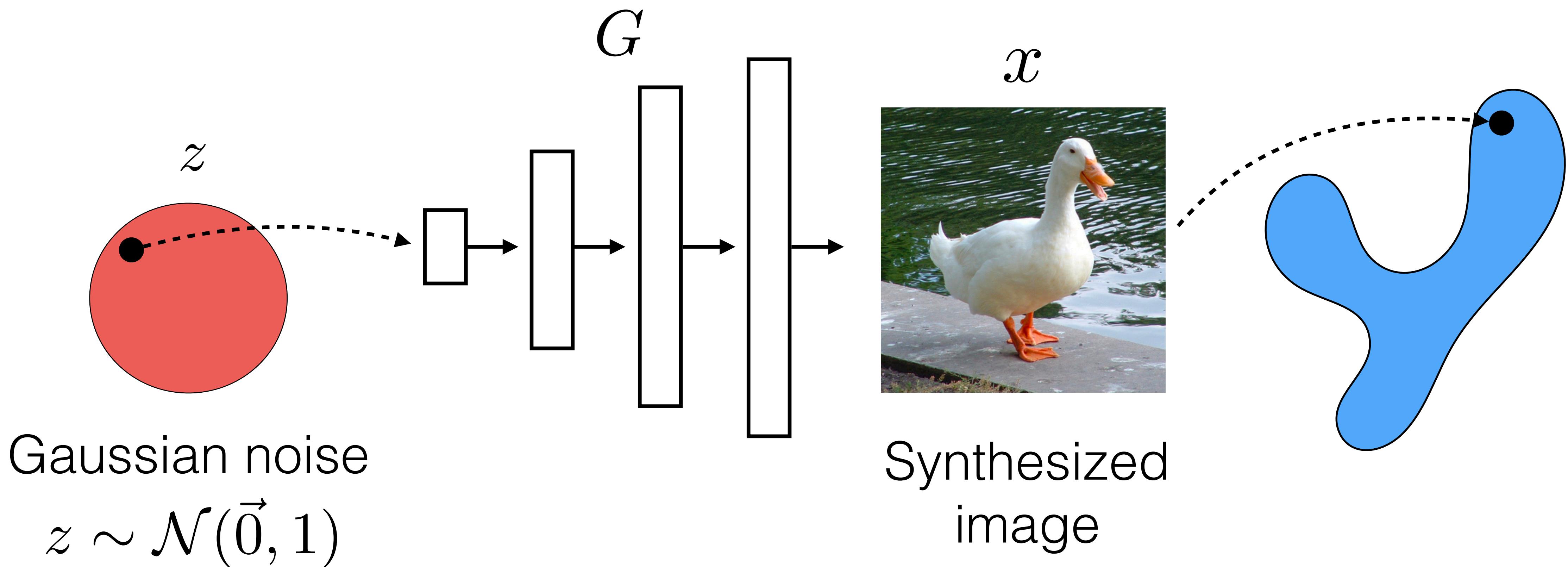


$$p(x)$$

Neural networks as distribution transformers

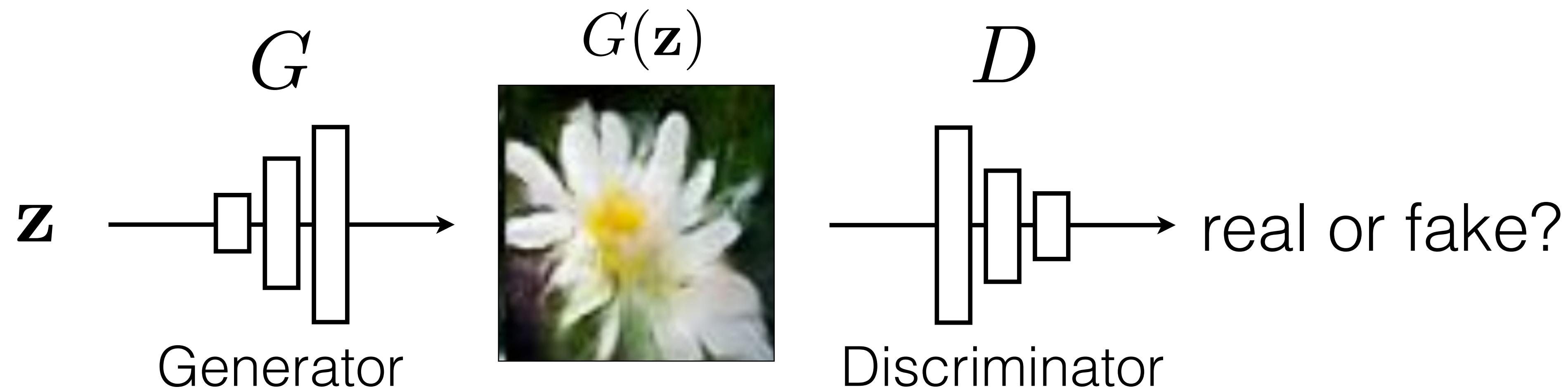


Neural networks as distribution transformers



Today

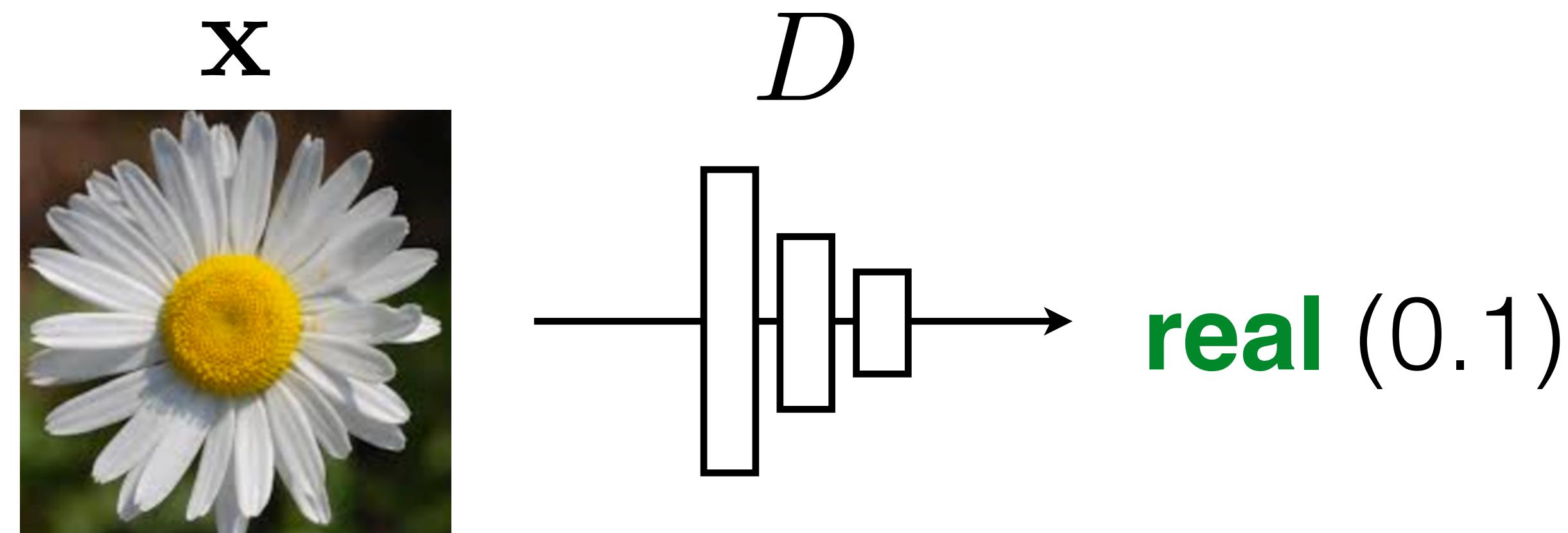
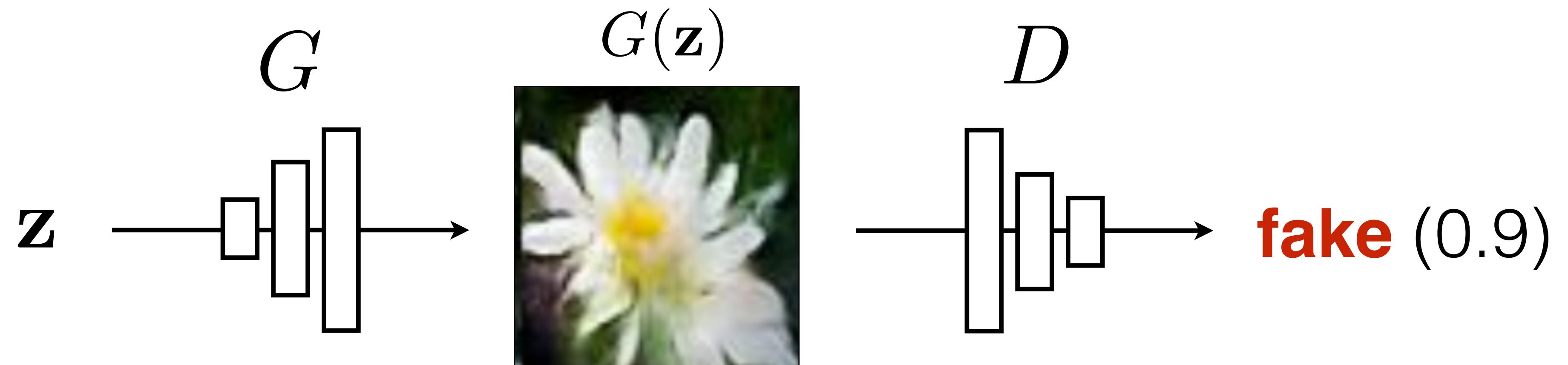
- Texture synthesis
- **Generative adversarial nets (GANs)**
- Autoregressive models



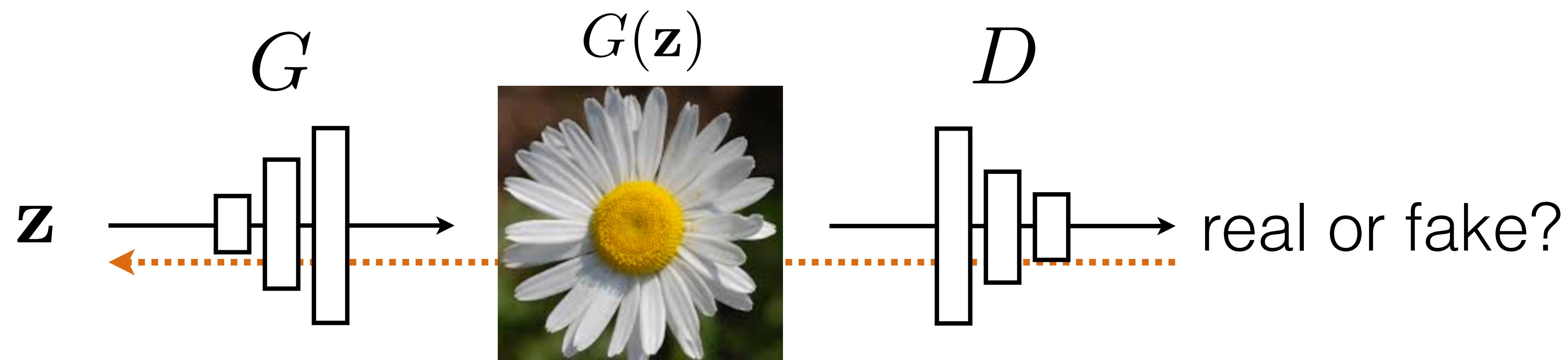
G tries to synthesize fake images that fool **D**

D tries to identify the fakes

[Goodfellow et al., 2014]



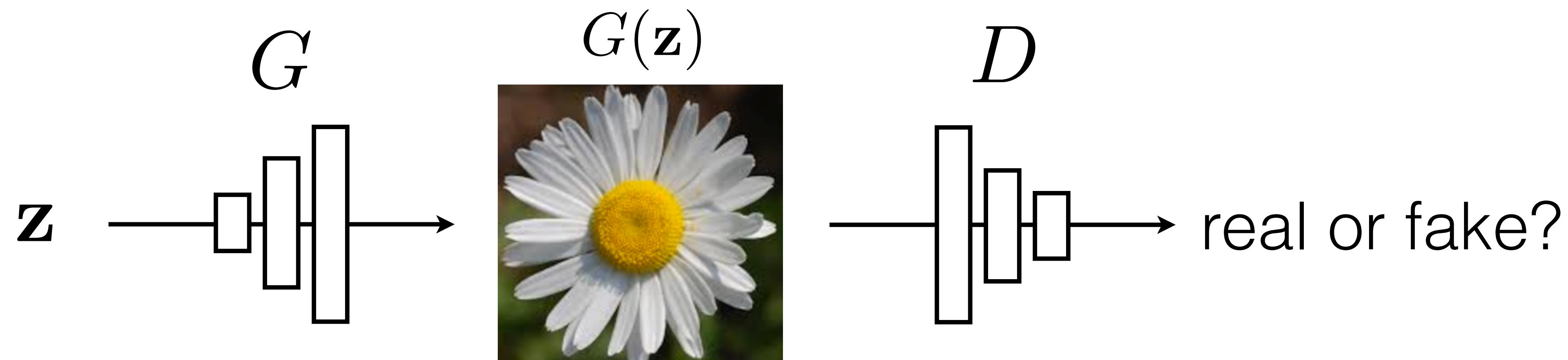
$$\arg \max_D \mathbb{E}_{\mathbf{z}, \mathbf{x}} [\boxed{\log D(G(\mathbf{z}))} + \boxed{\log (1 - D(\mathbf{x}))}]$$



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]

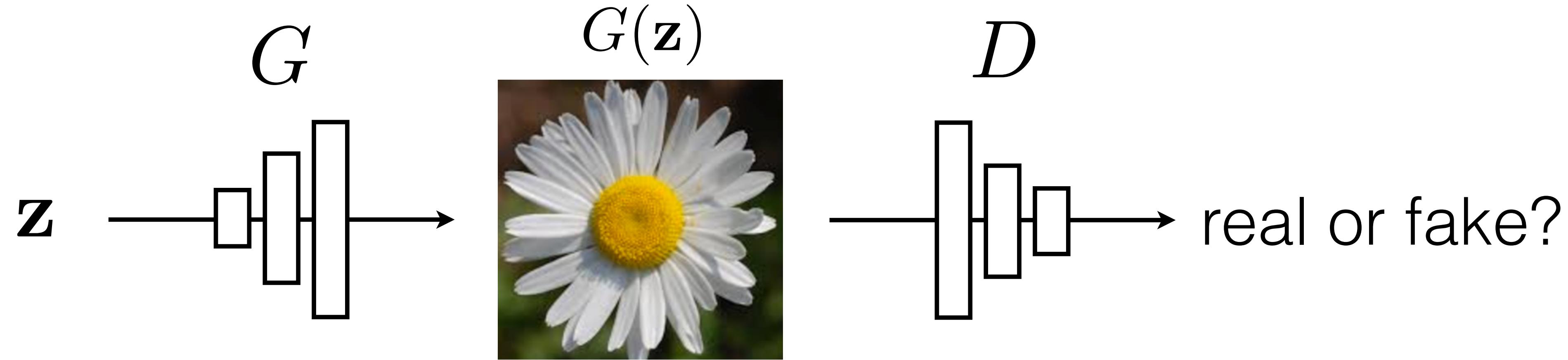


G tries to synthesize fake images that **fool** the **best** **D**:

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

[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 (see book)
- Note: when training G, we take derivatives through discriminator!

[Goodfellow et al., 2014]

Samples from BigGAN

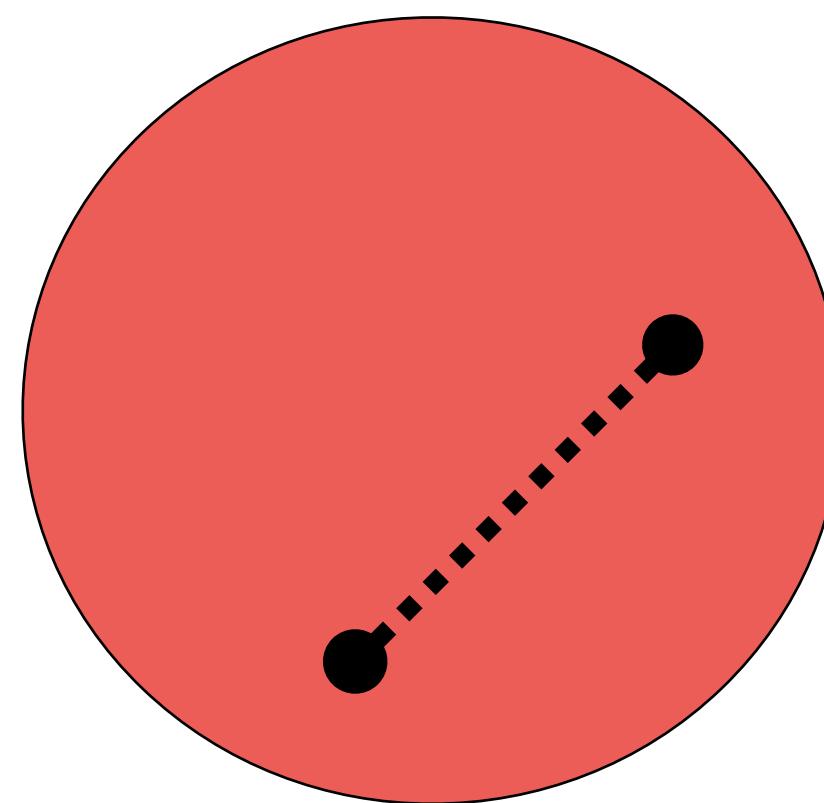
[Brock et al. 2018]



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

Latent space
(Gaussian)

z



Data space
(Natural image manifold)

X



[BigGAN, Brock et al. 2018]





<https://github.com/NVlabs/stylegan3> [Karras et al., "Alias-Free Generative Adversarial Networks", 2021]

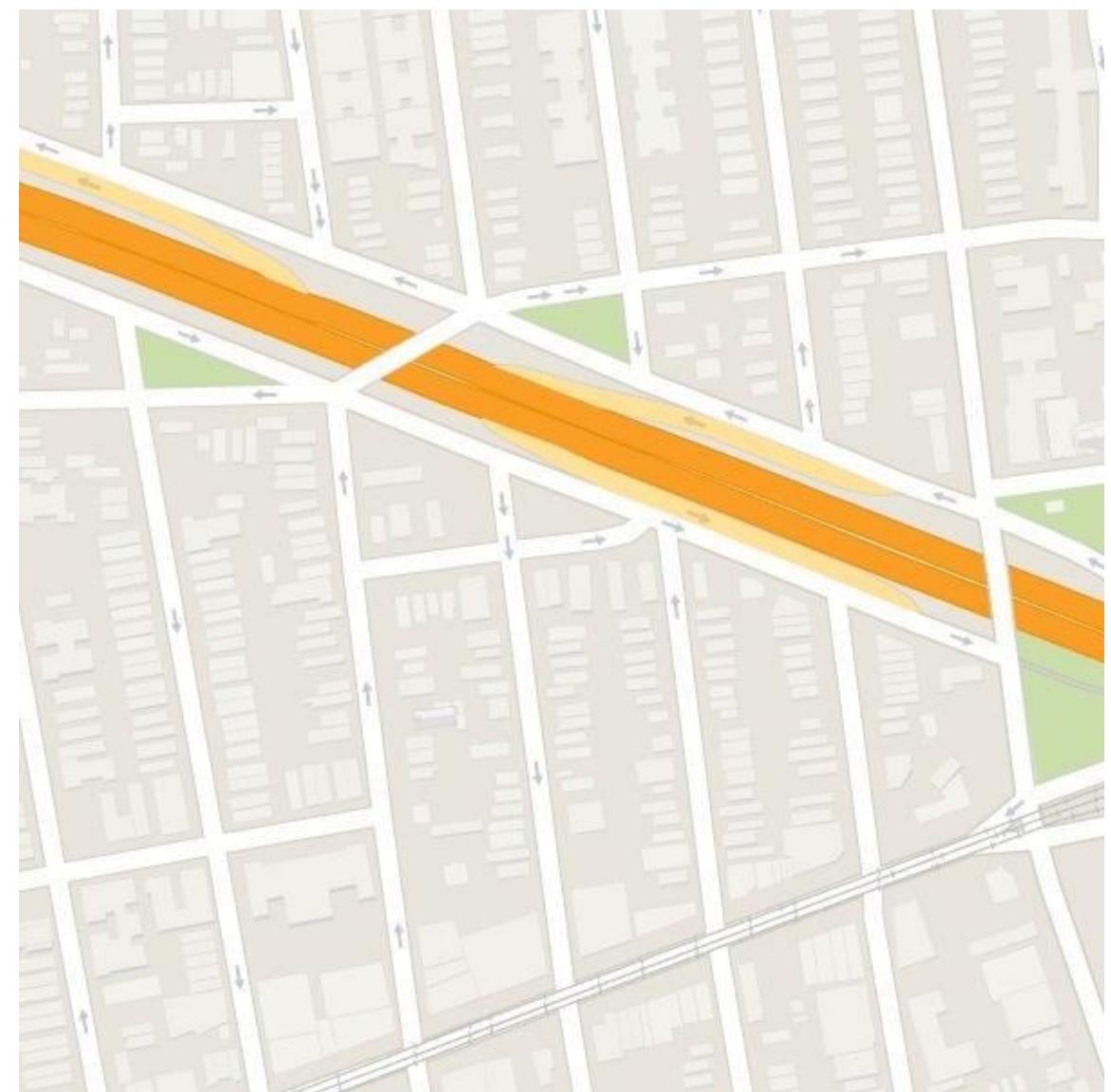


[Karras et al., "Alias-Free Generative Adversarial Networks", 2021]

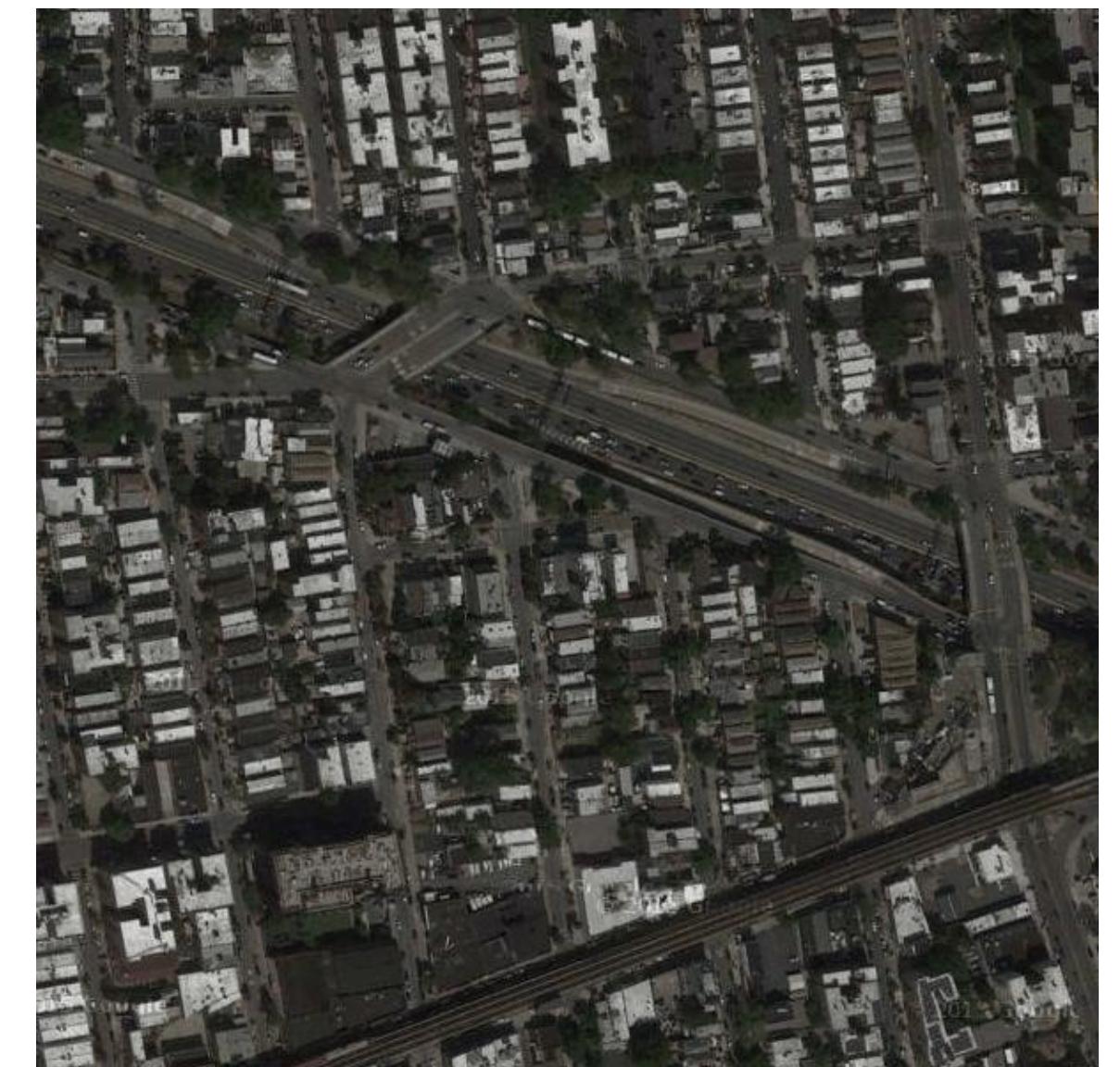


[Karras et al., "Alias-Free Generative Adversarial Networks", 2021]

Image translation



Google Map



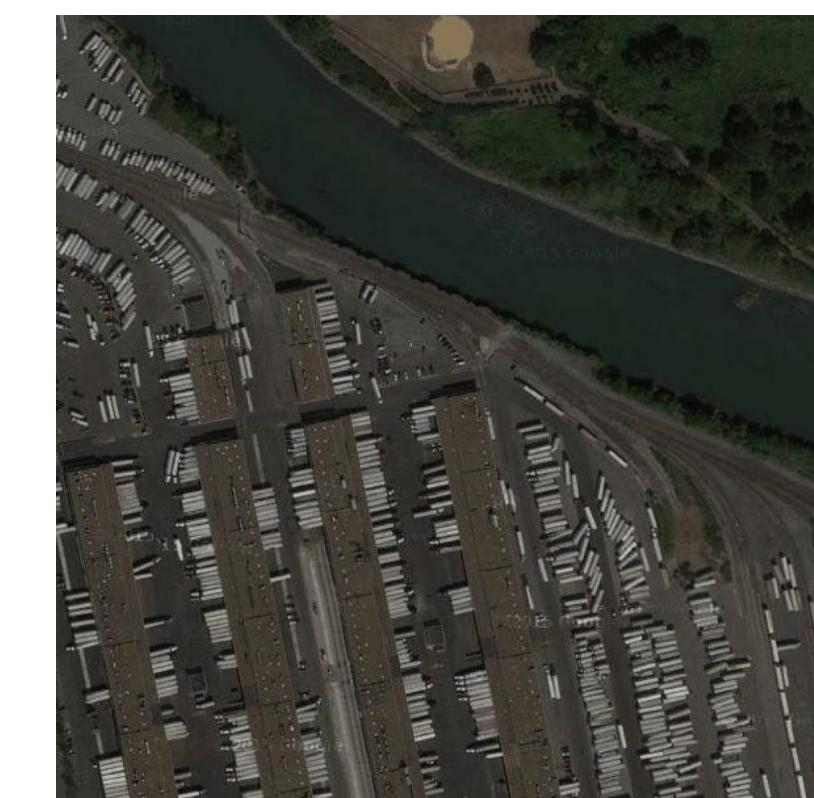
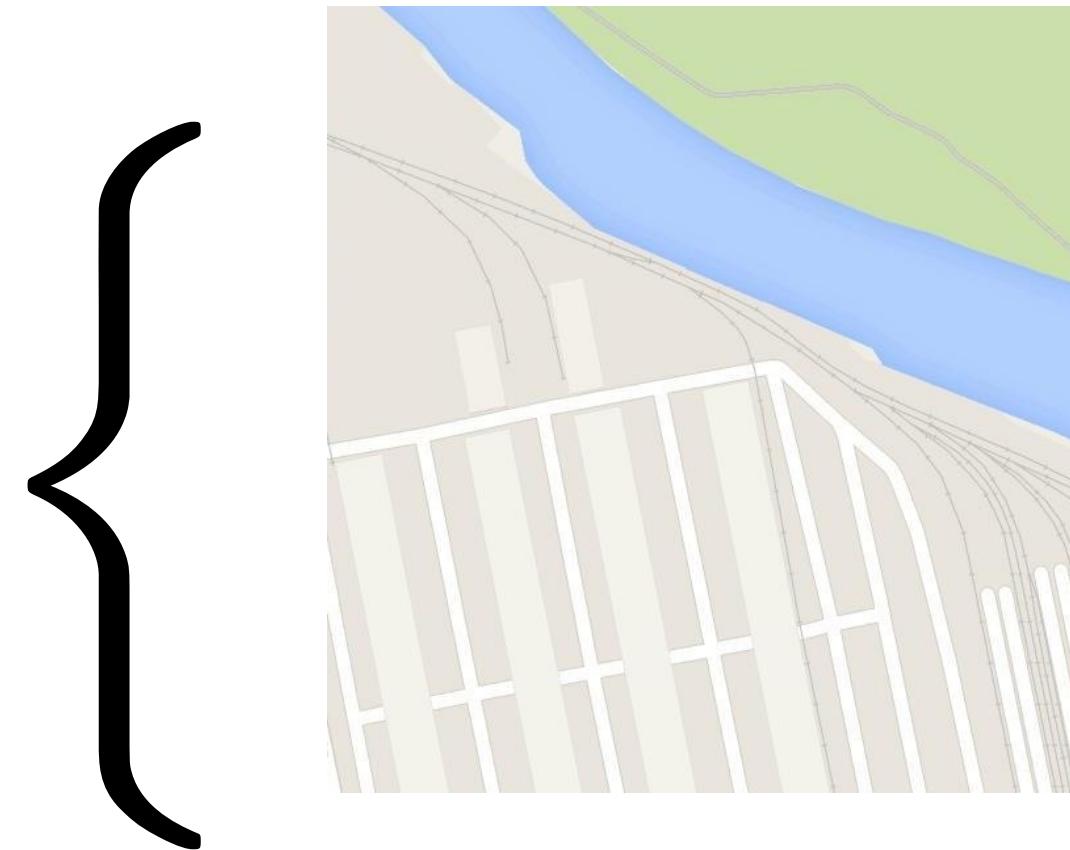
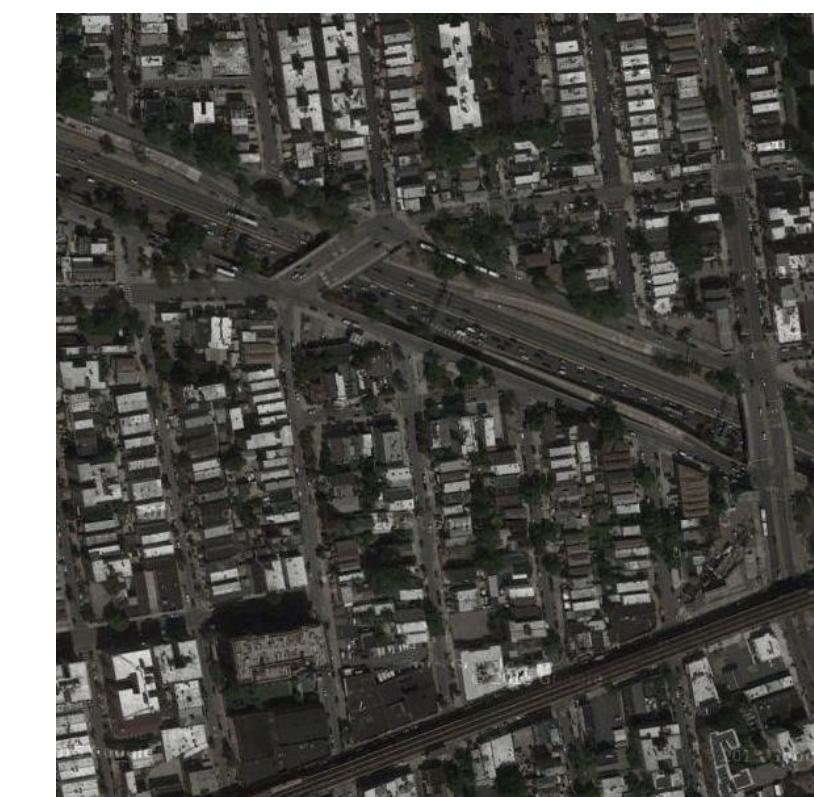
Satellite photo

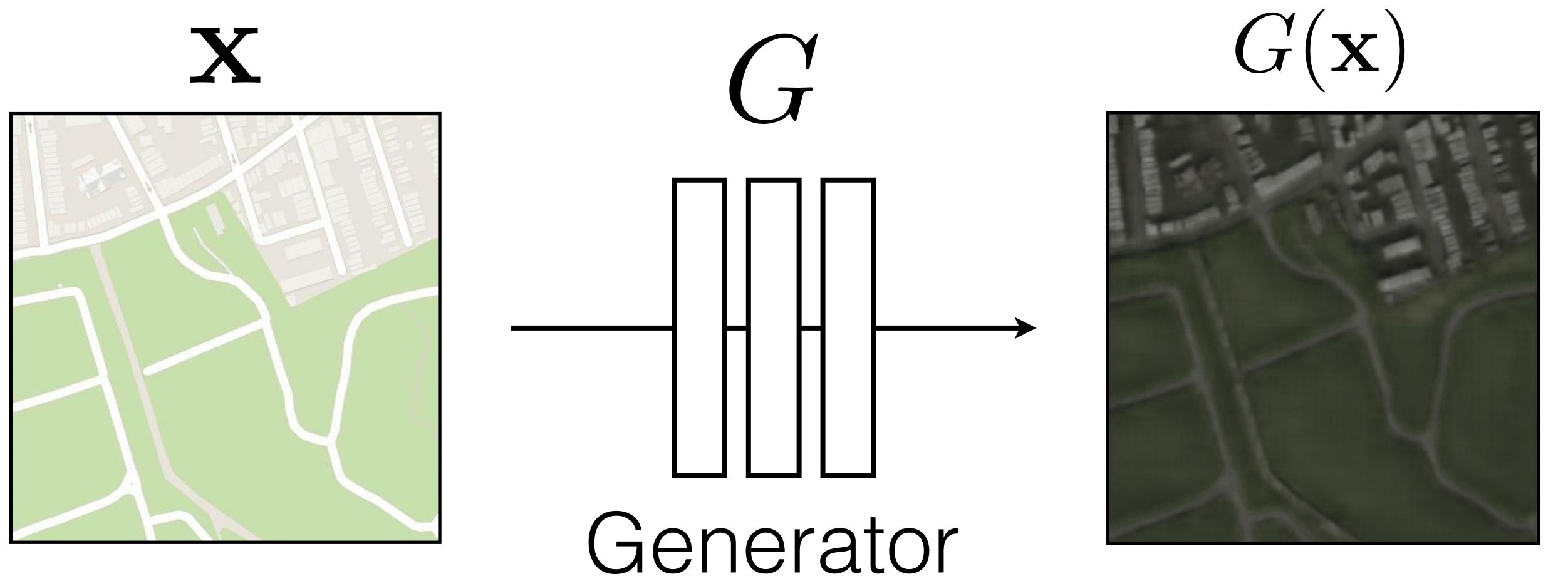
Map2Sat

x



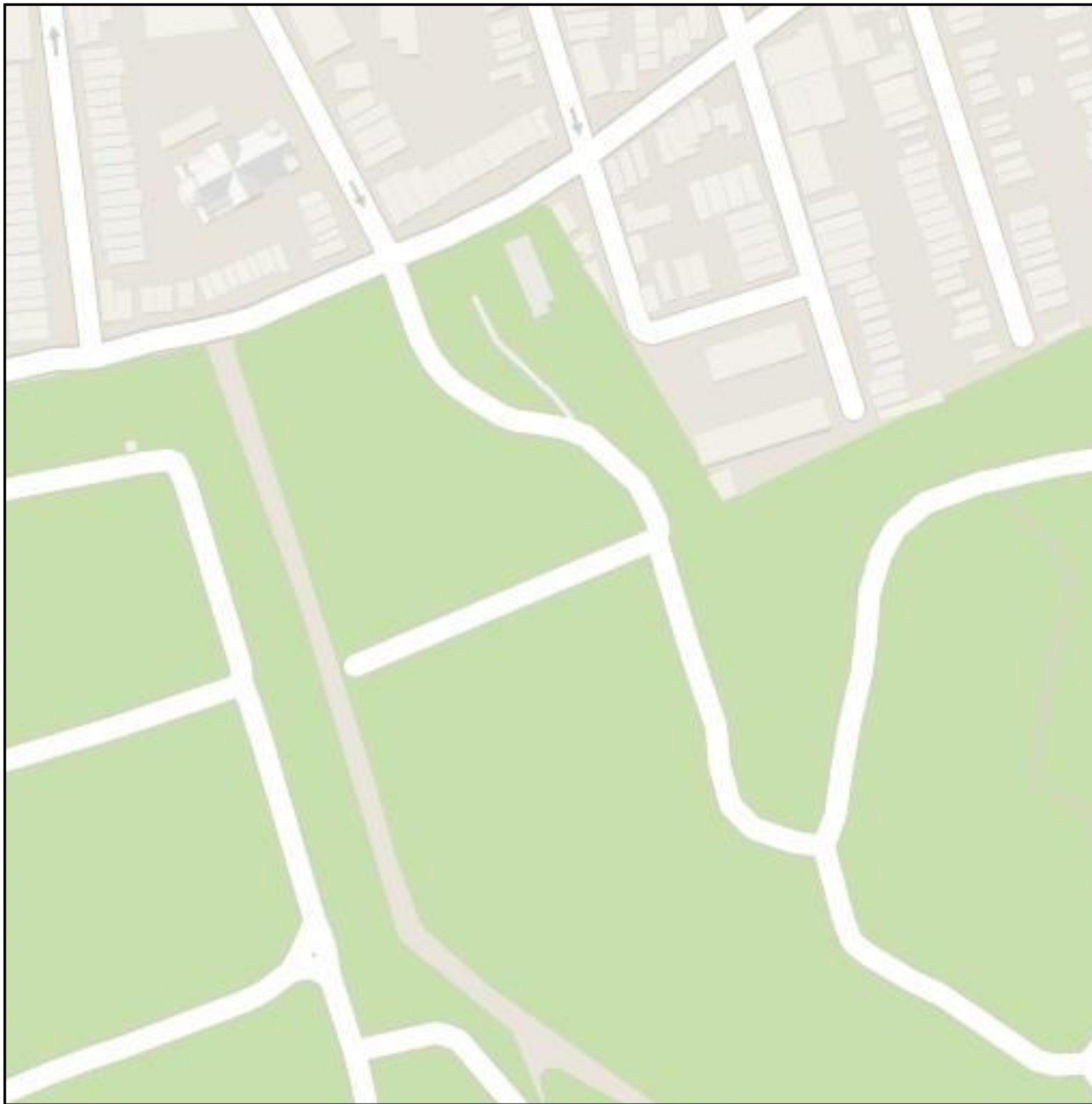
y





Idea: L1 loss
 $\|G(\mathbf{x}) - \mathbf{y}\|_1$

Input



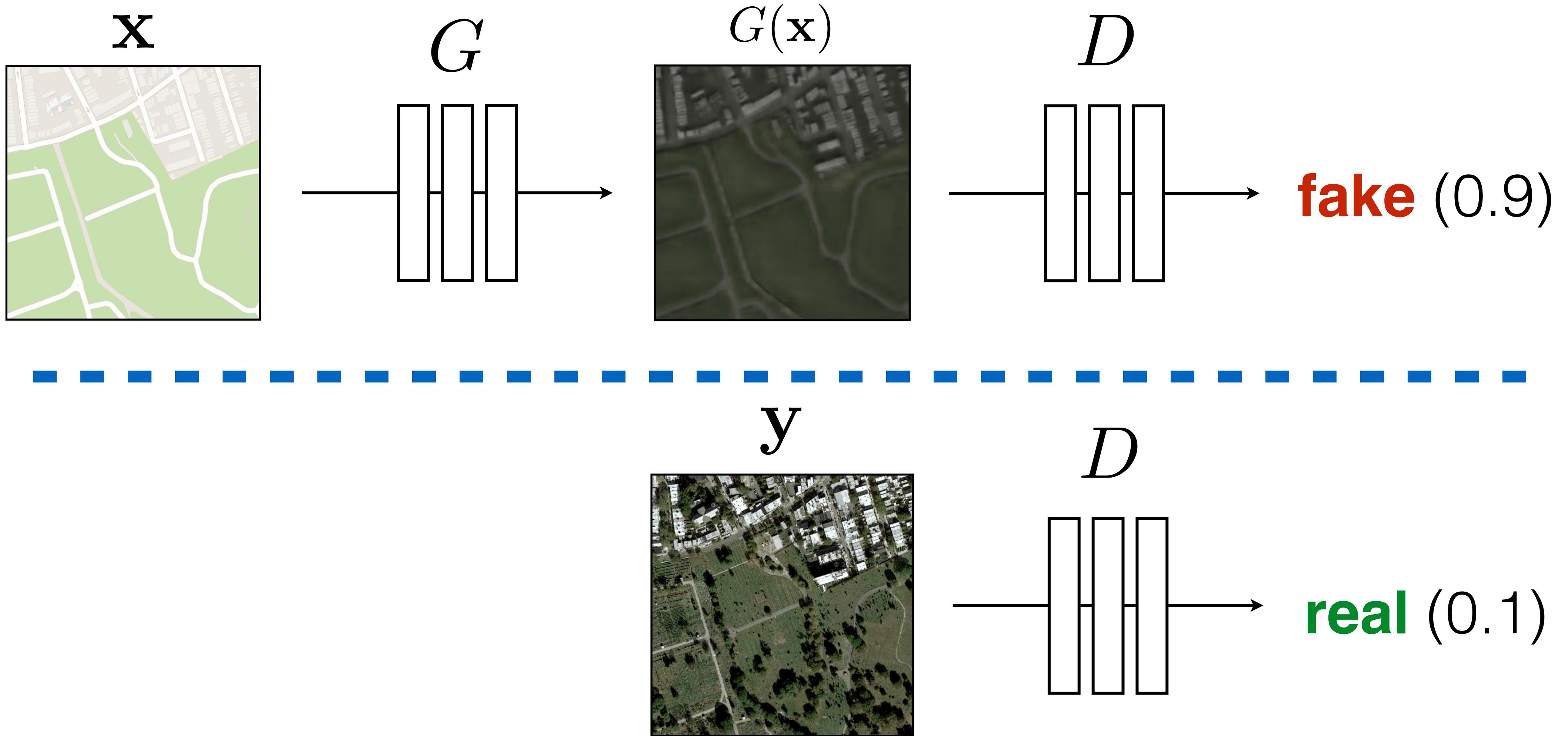
L1 loss



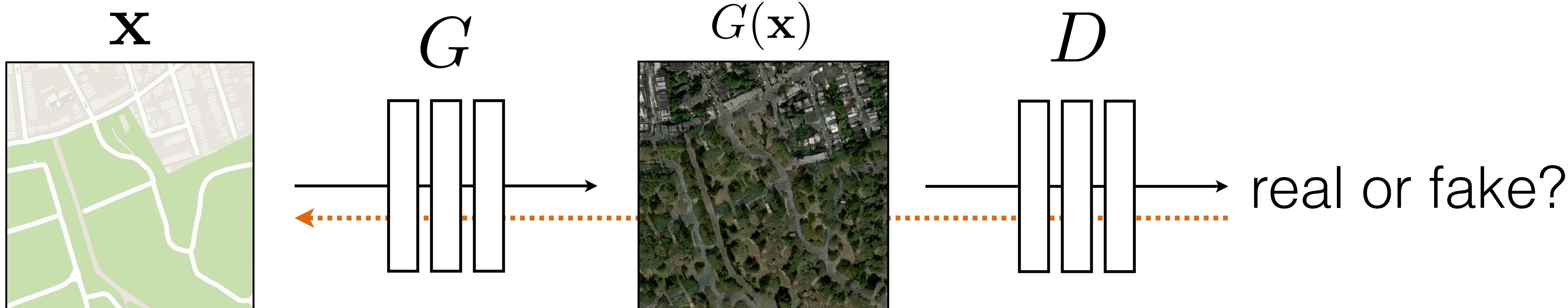


G tries to synthesize fake images that fool **D**

D tries to identify the fakes

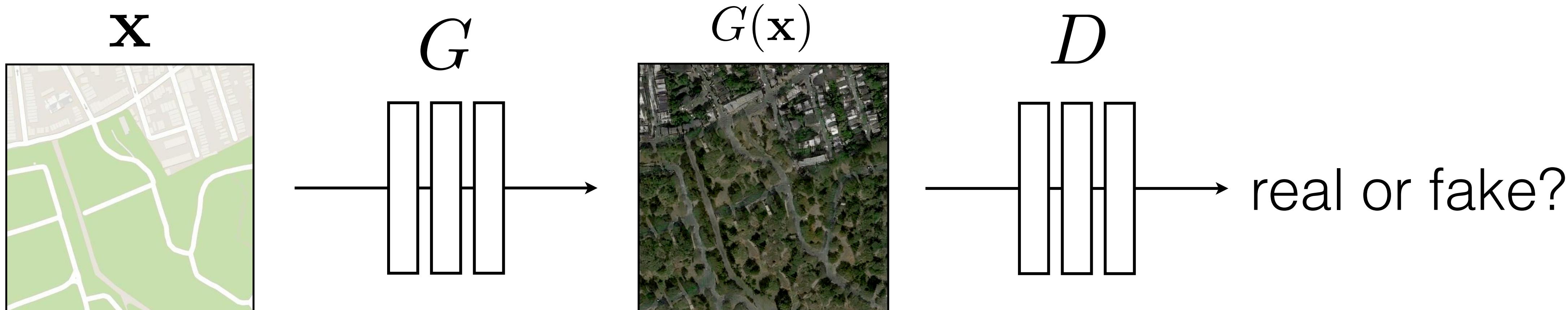


$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\boxed{\log D(G(\mathbf{x}))} + \boxed{\log(1 - D(\mathbf{y}))}]$$



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

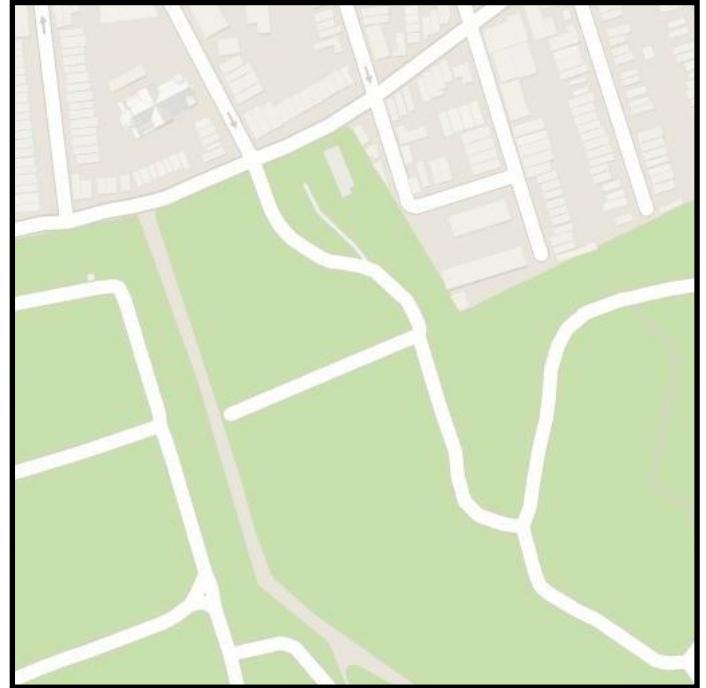
$$\arg \min_G \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$



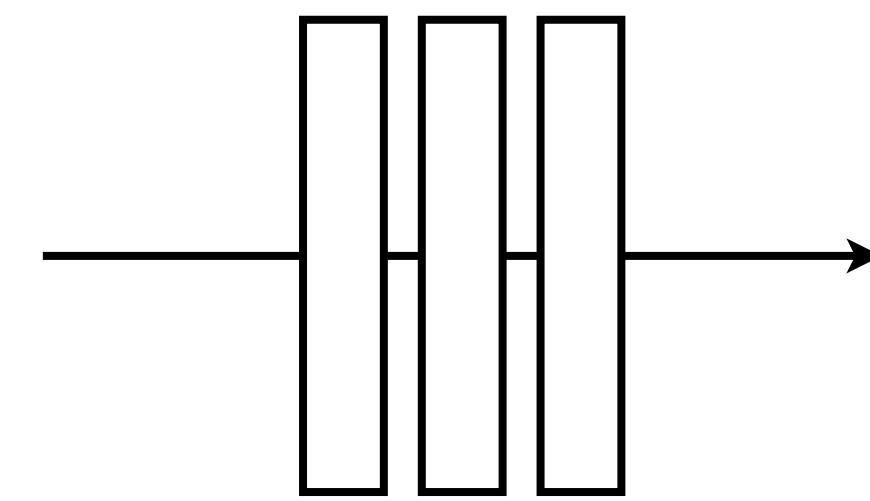
G tries to synthesize fake images that **fool** the **best** **D**:

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

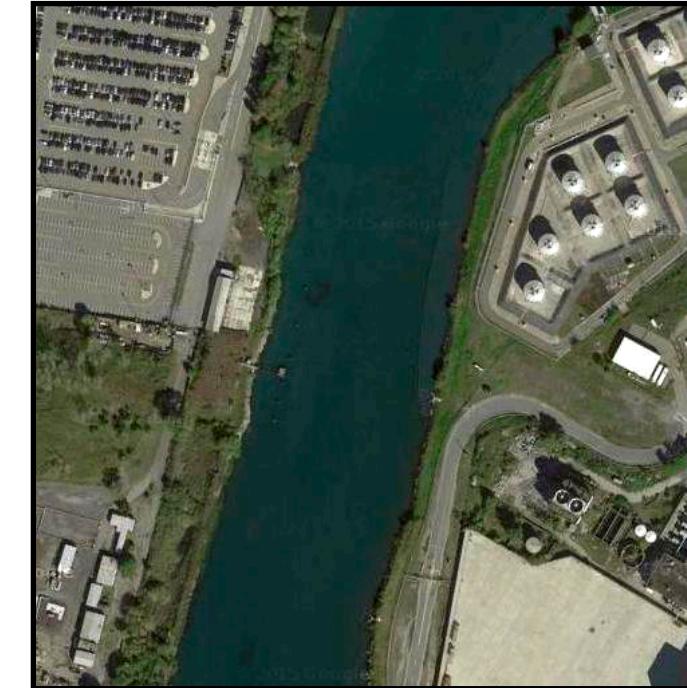
x



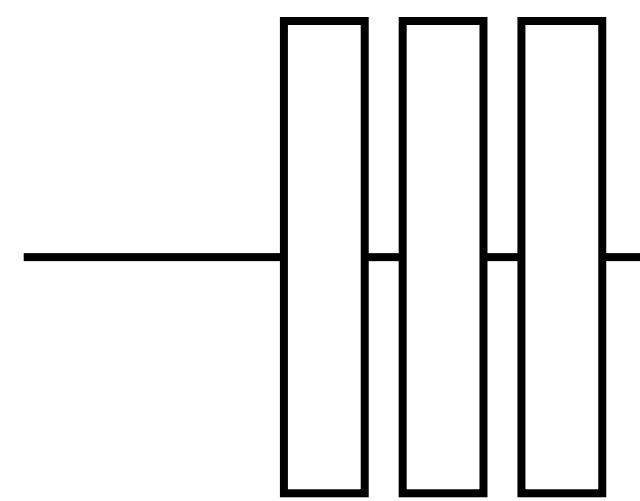
G



G(x)



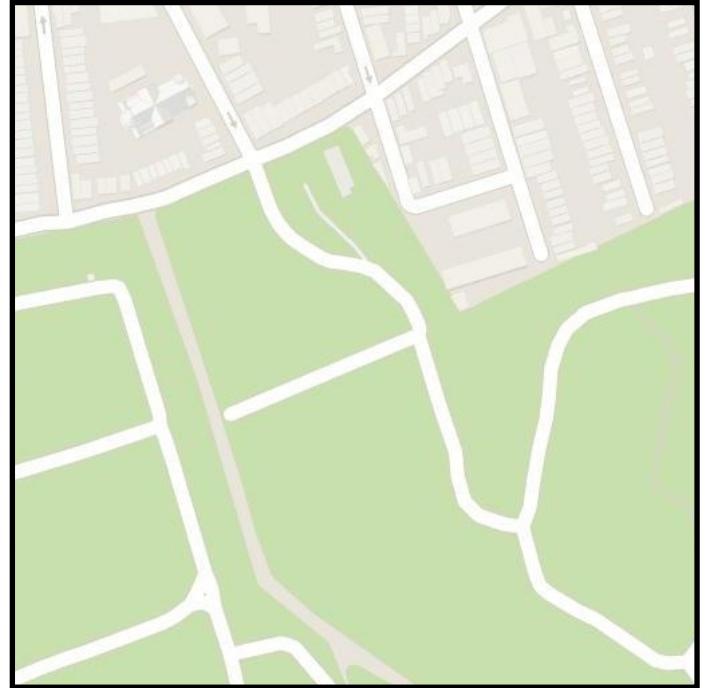
D



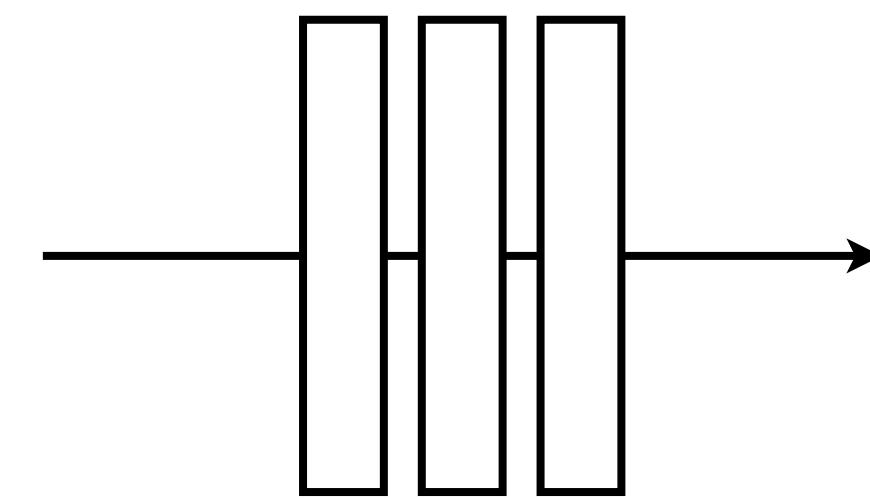
real or fake?

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

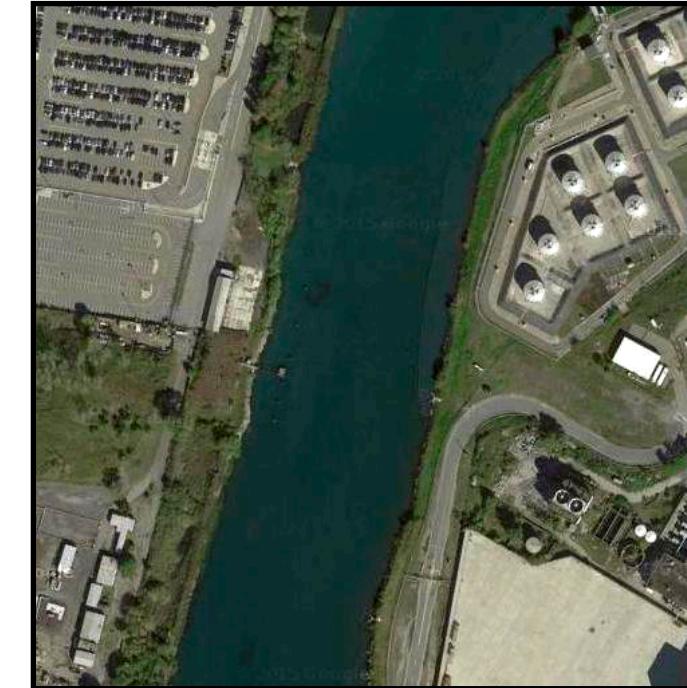
x



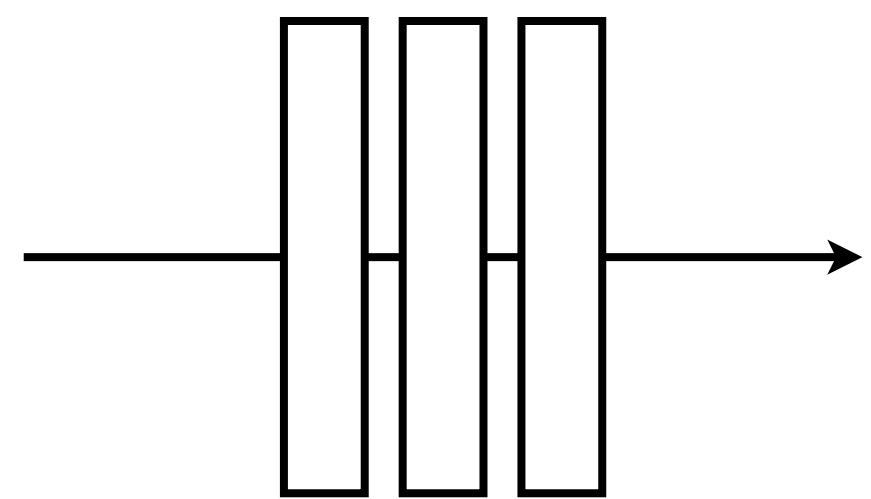
G



G(x)

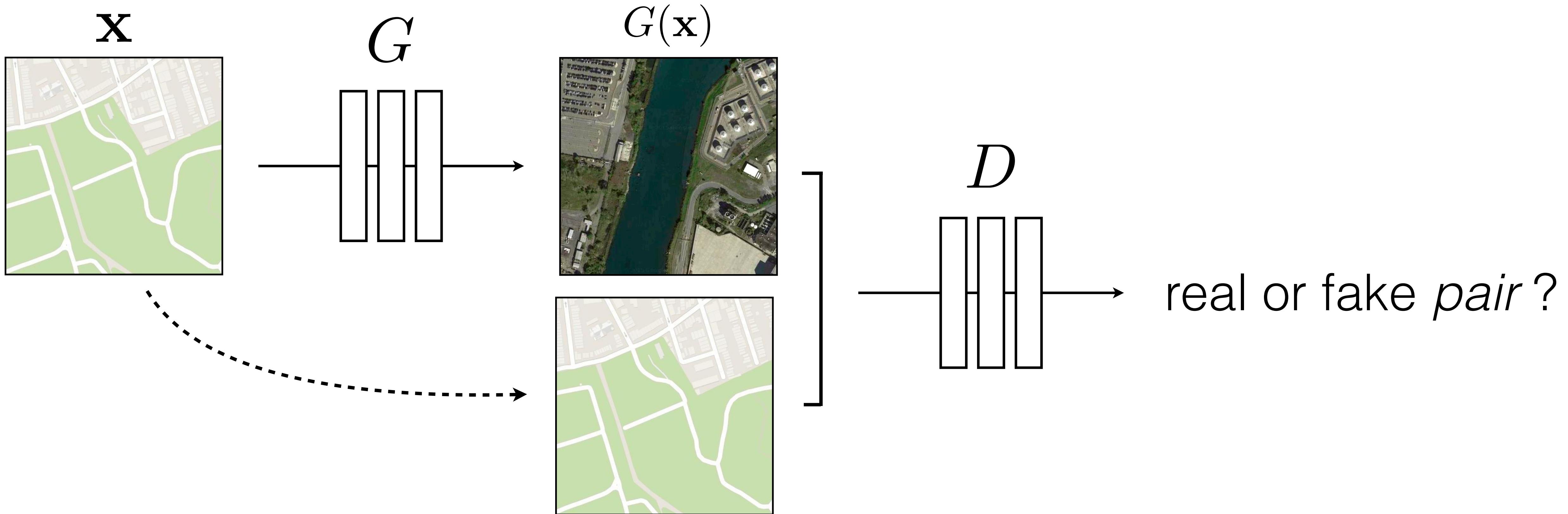


D

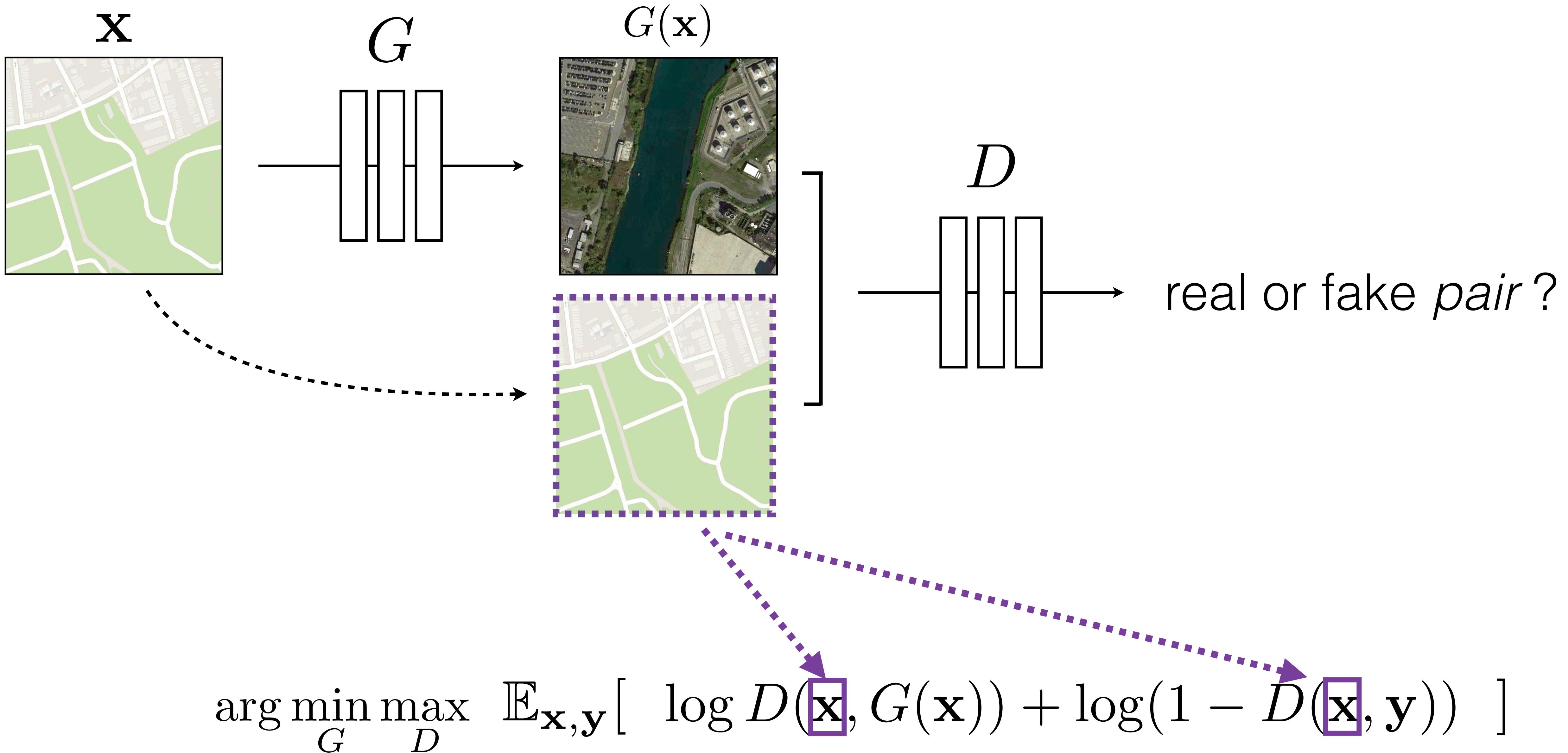


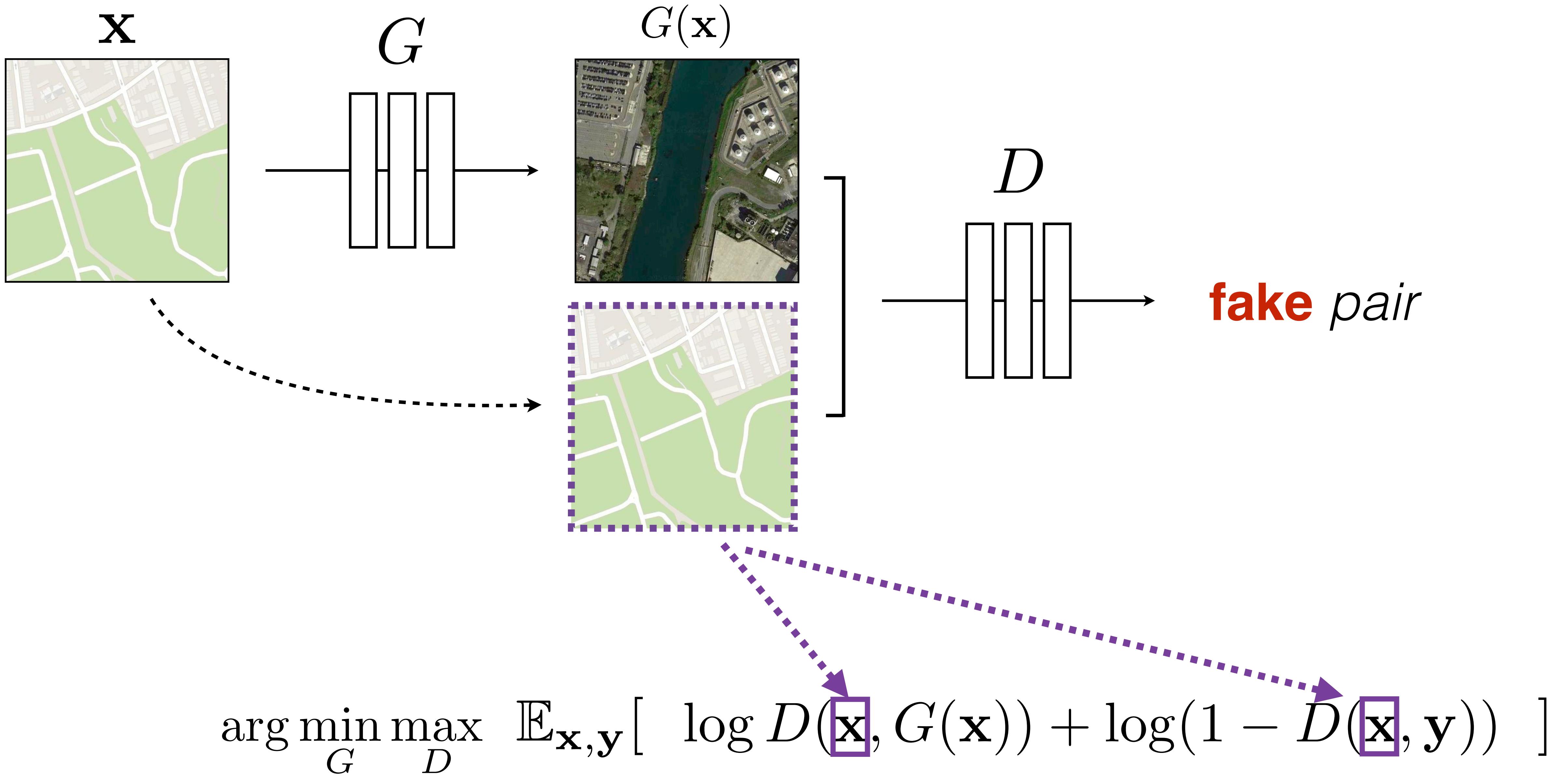
real!

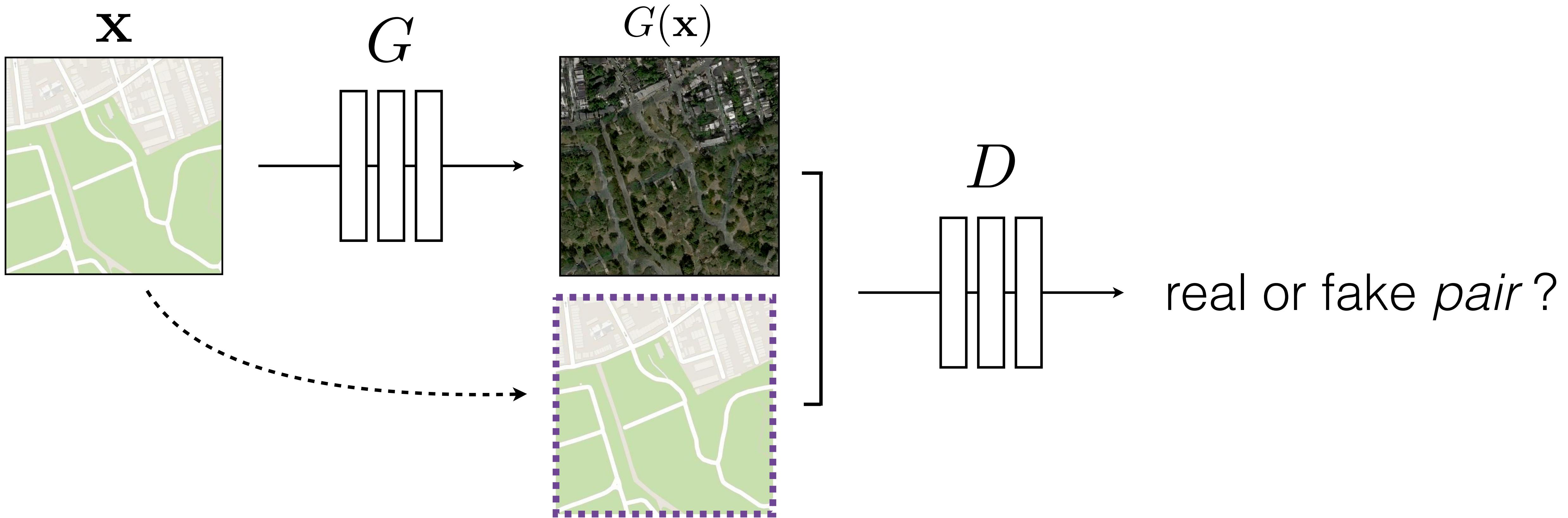
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$



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





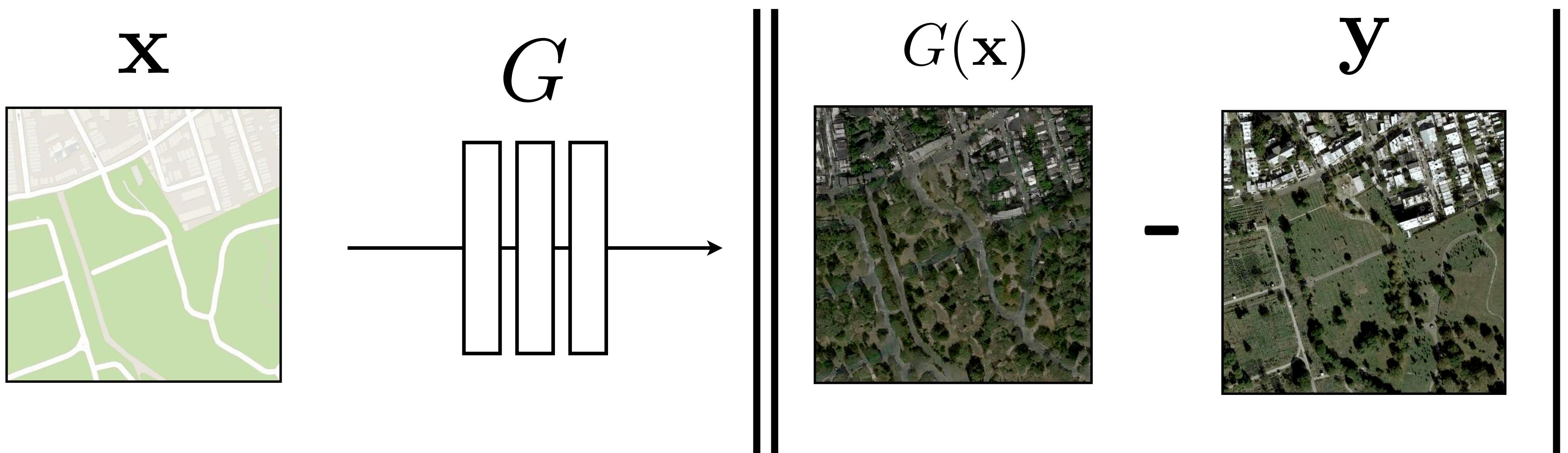


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))]$$

Training Details: Loss function

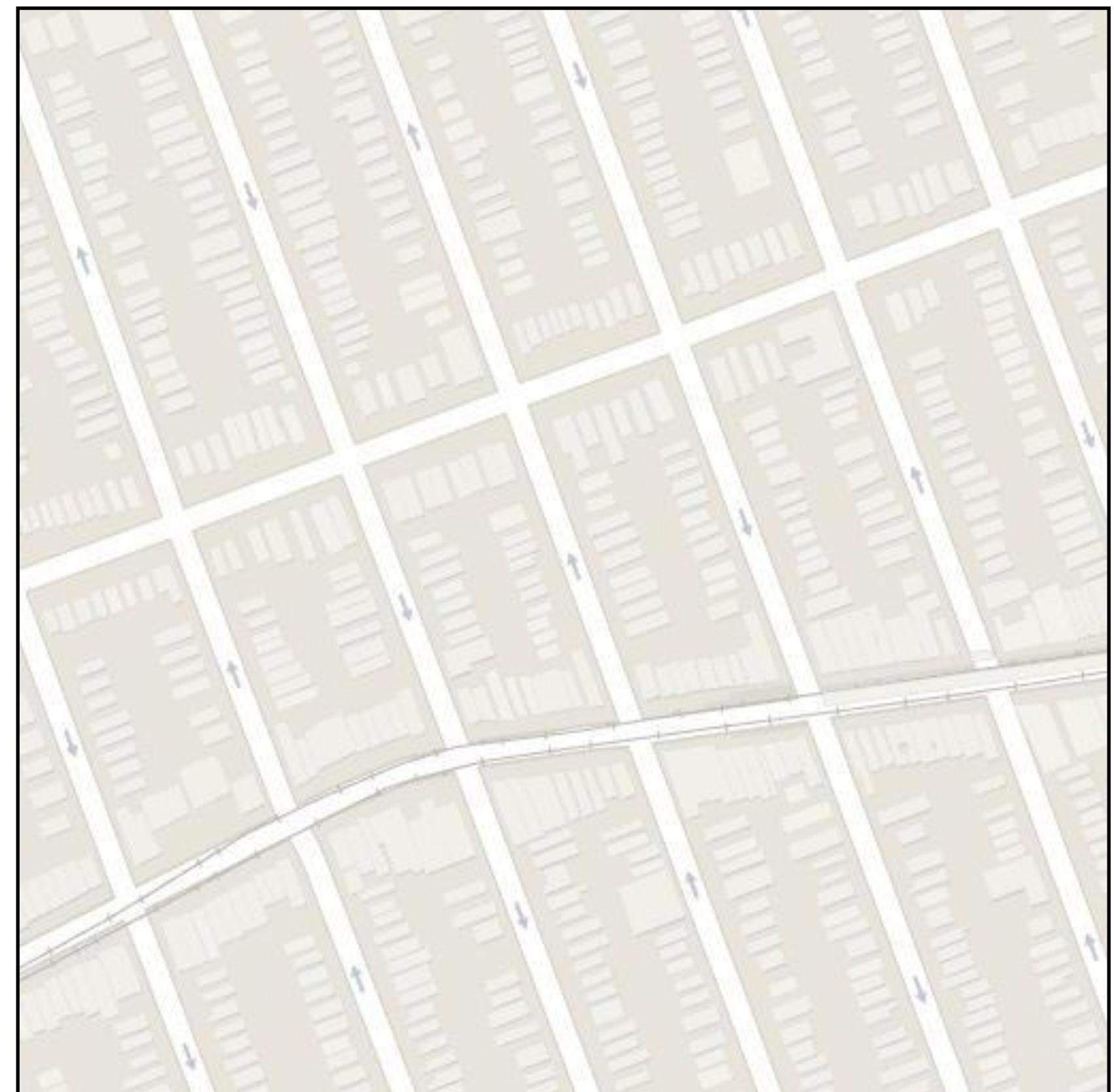
Conditional GAN

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

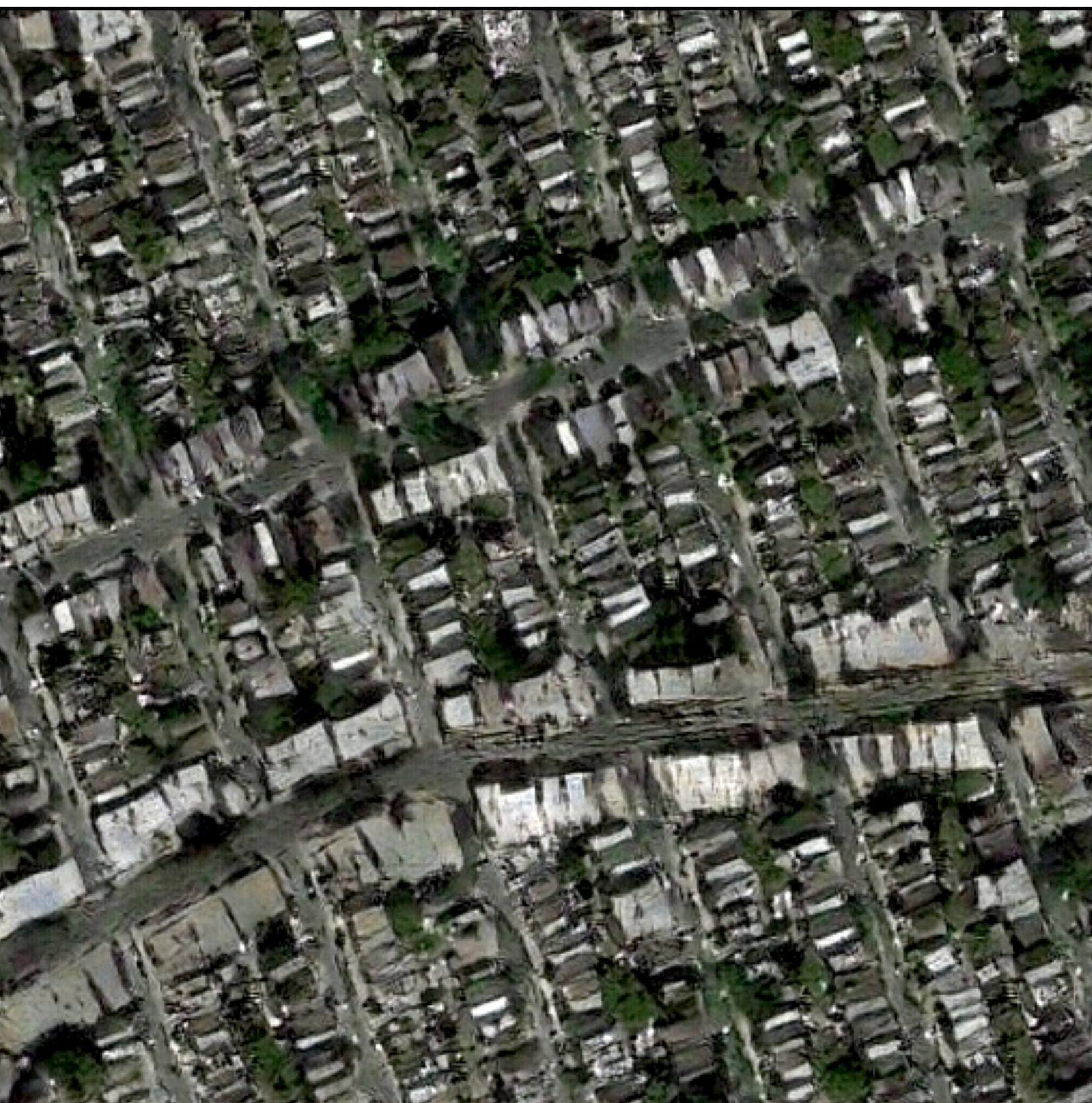


Pixel-wise loss helps stabilize training + faster convergence

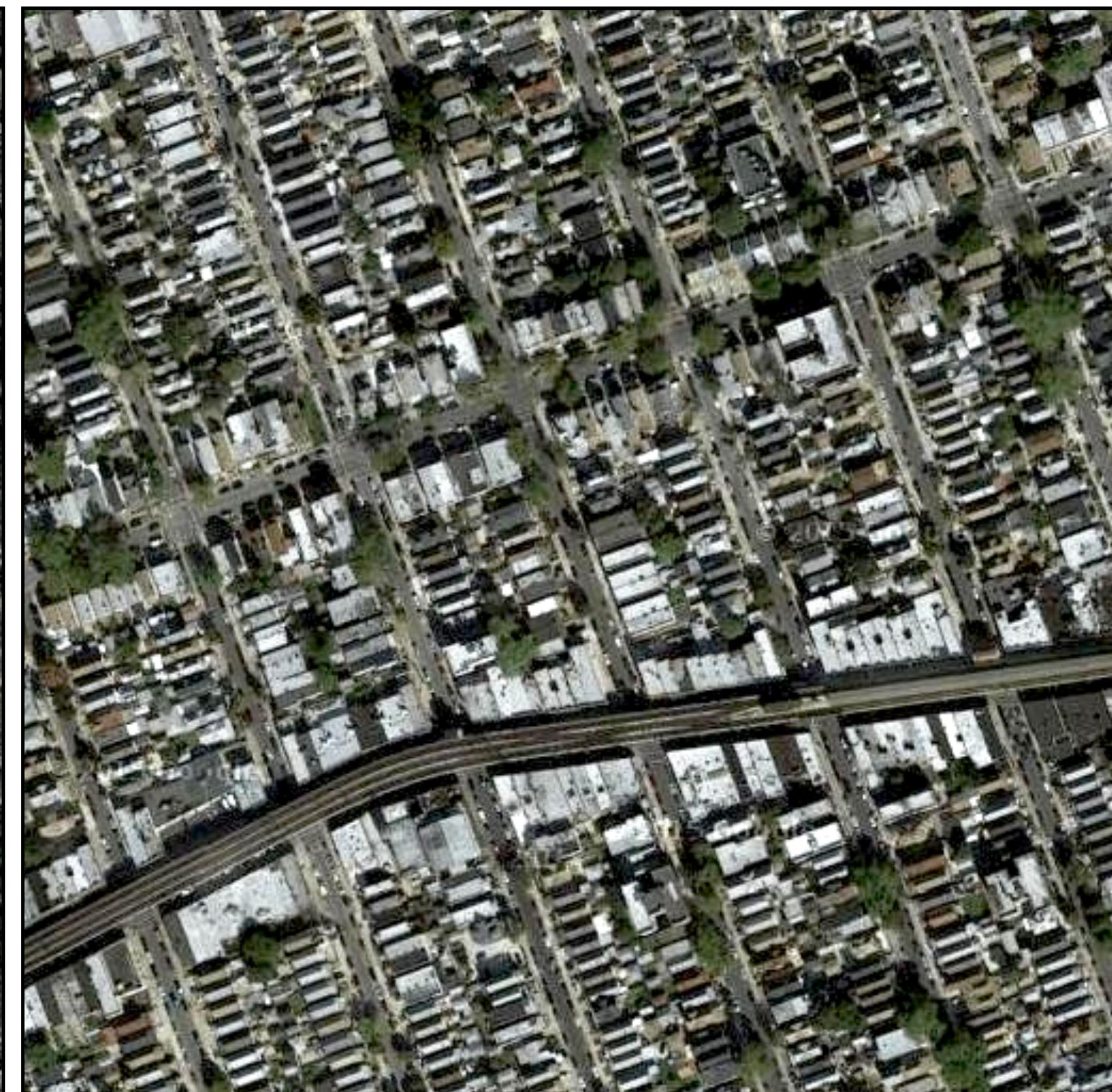
Input



Output



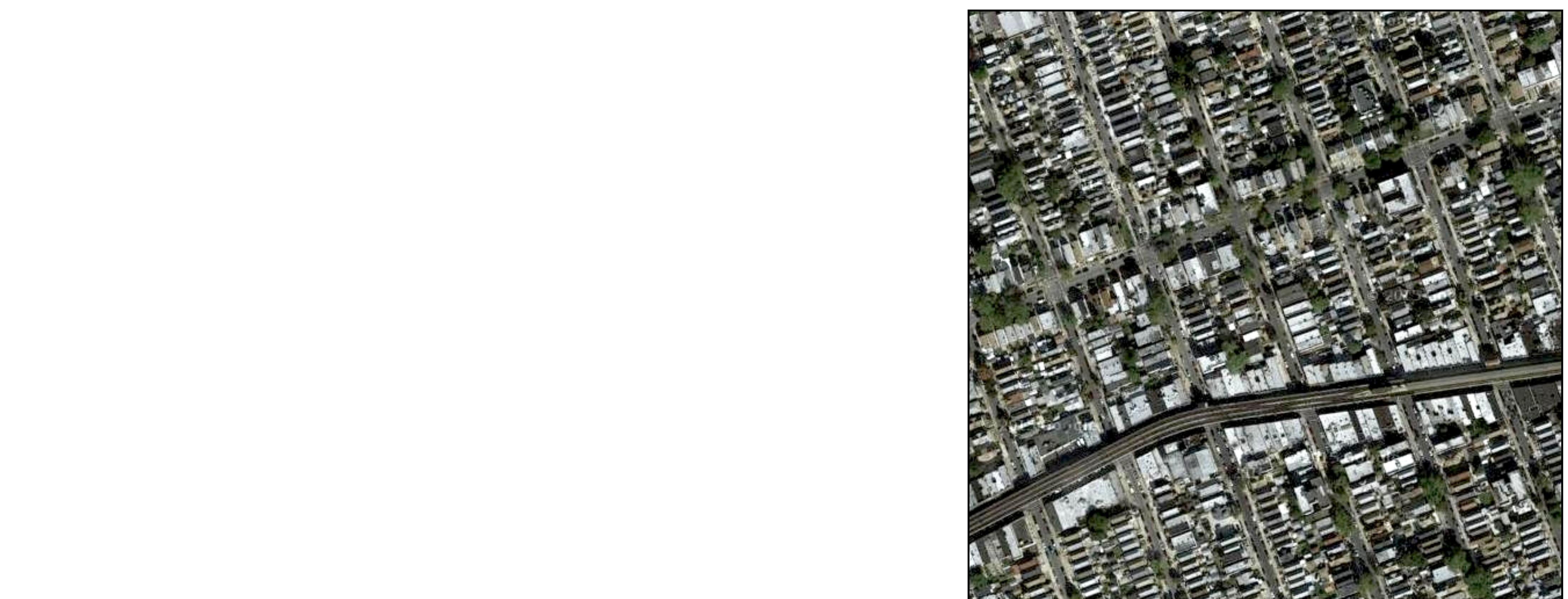
Groundtruth



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



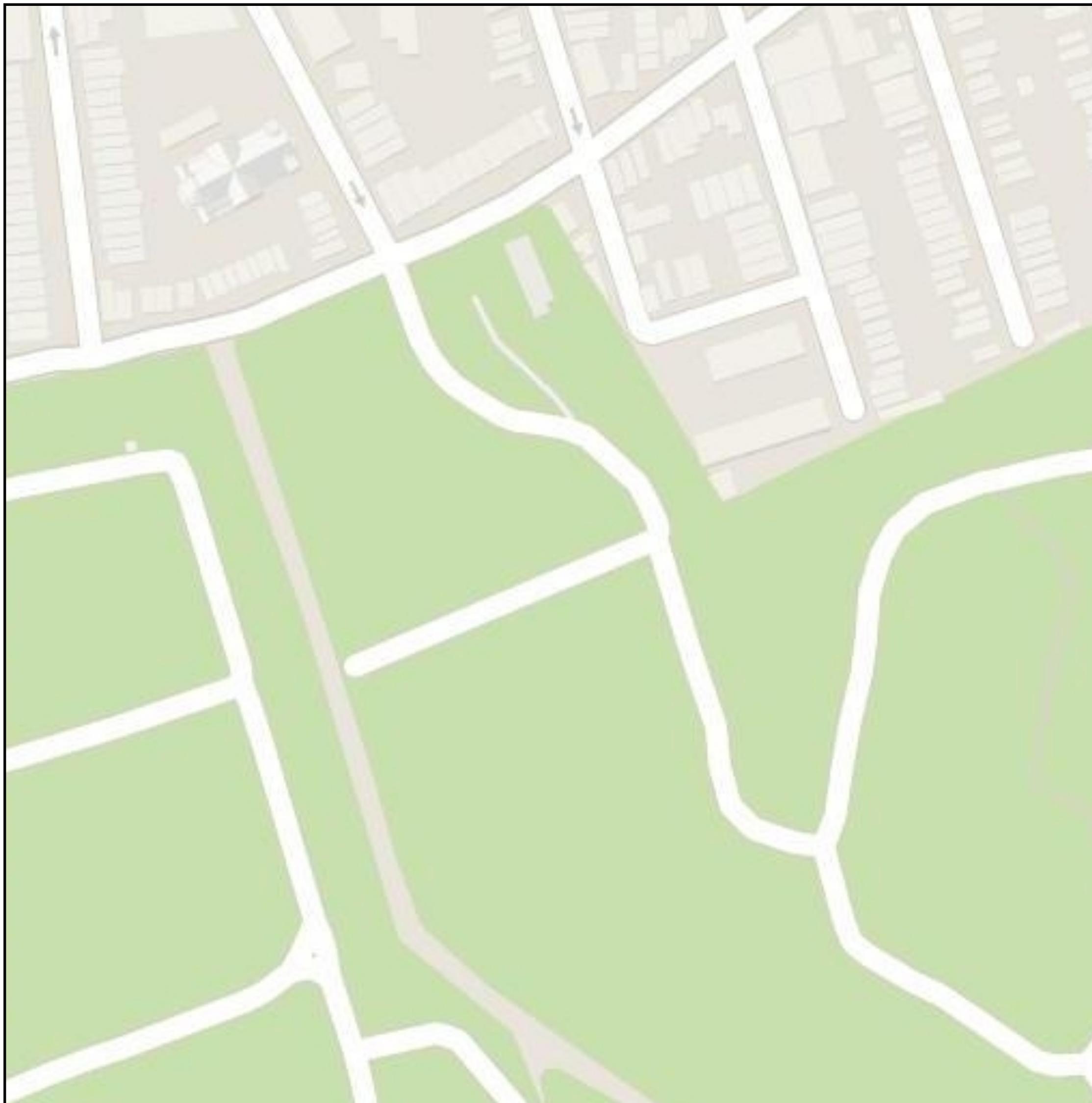
Input



Output

Groundtruth

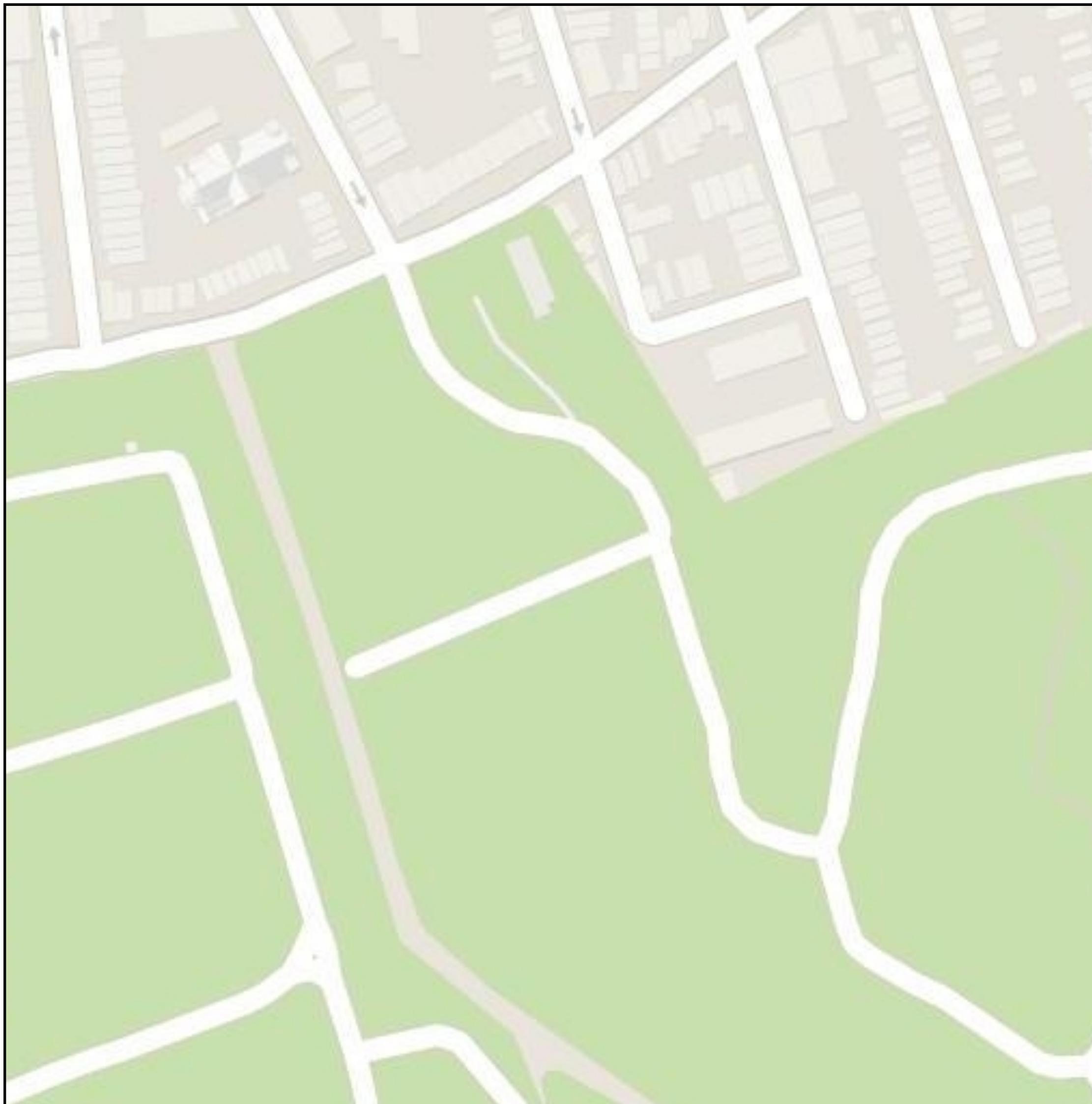
Input



L1 loss only



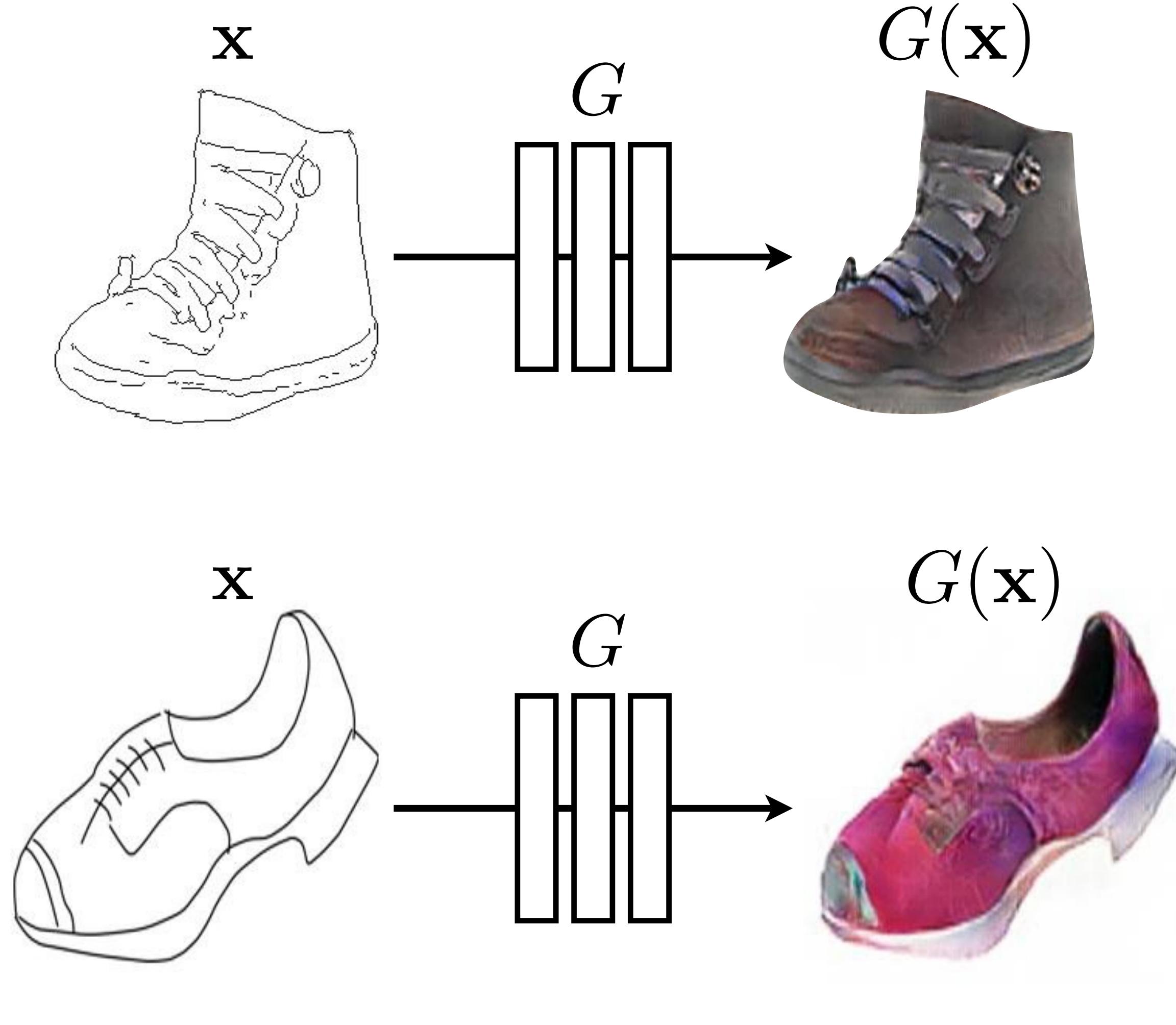
Input



L1 loss + discriminator

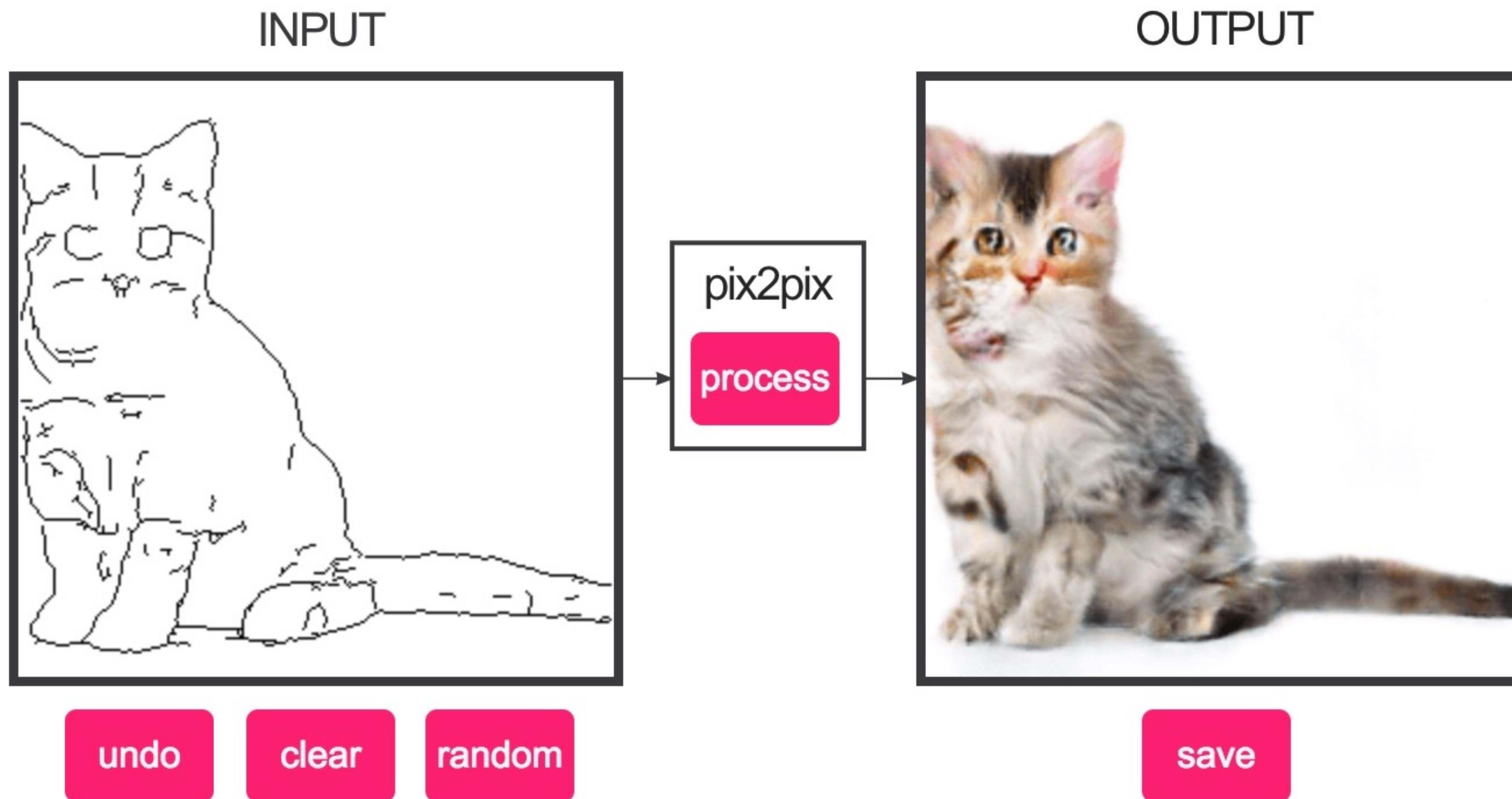


Training data



[HED, Xie & Tu, 2015]

edges2cats

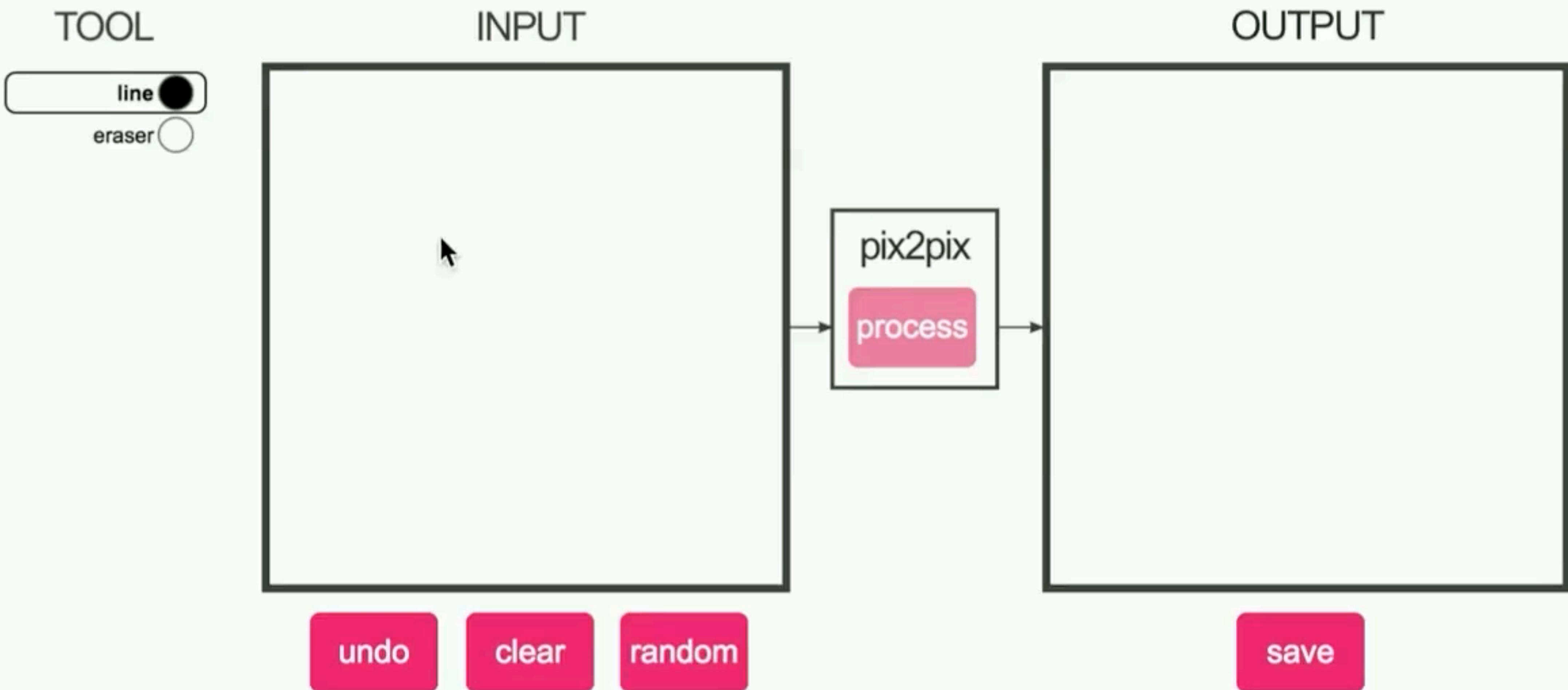


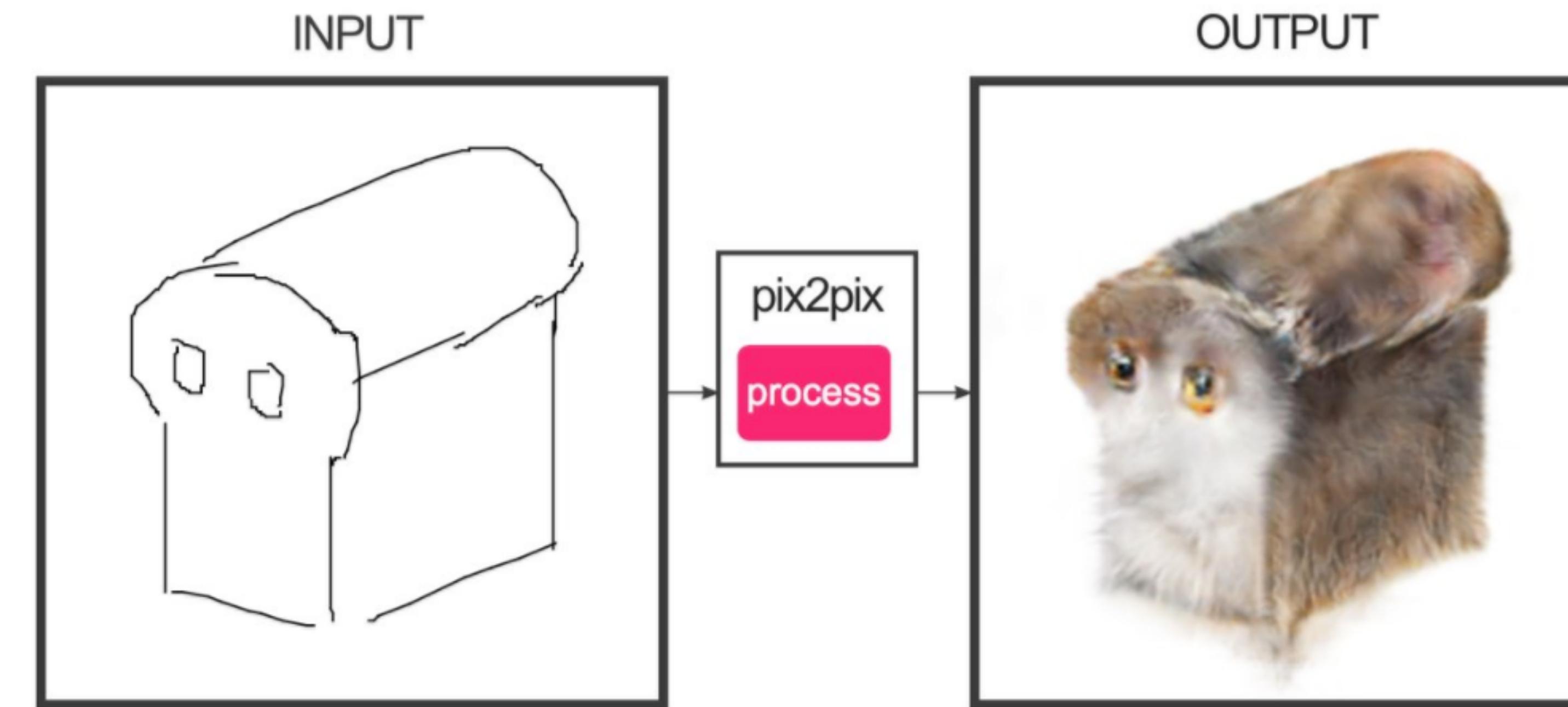
[Chris Hess, [edges2cats](#)]

60

Source: Isola, Freeman, Torralba

edges2cats





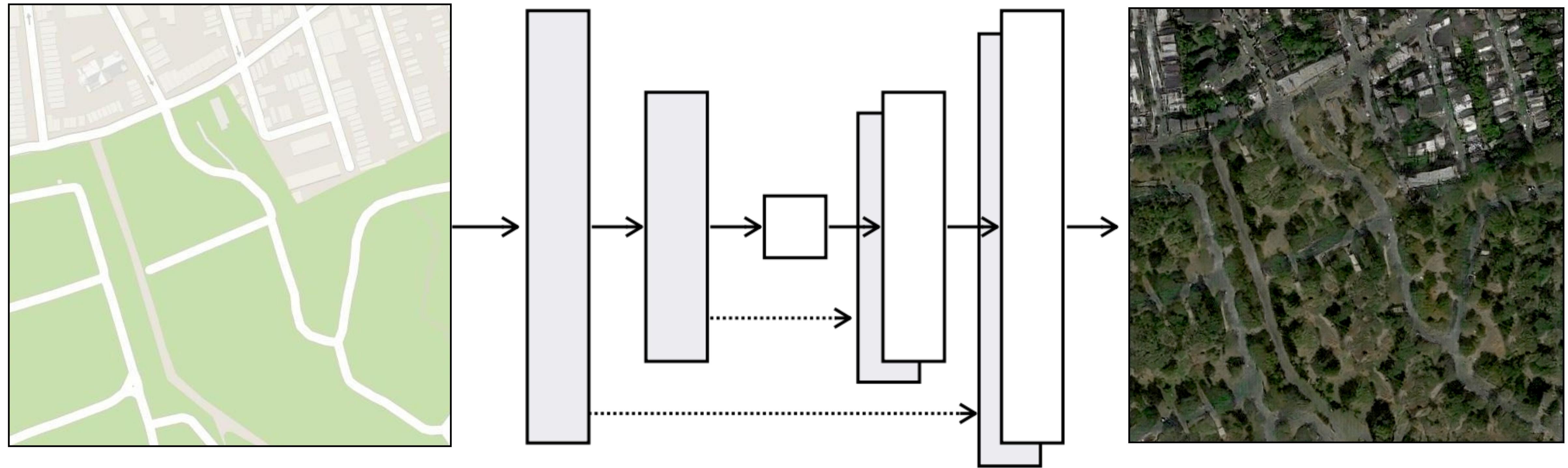
Ivy Tasi @ivymyt



Vitaly Vidmirov @vvid

Architectures

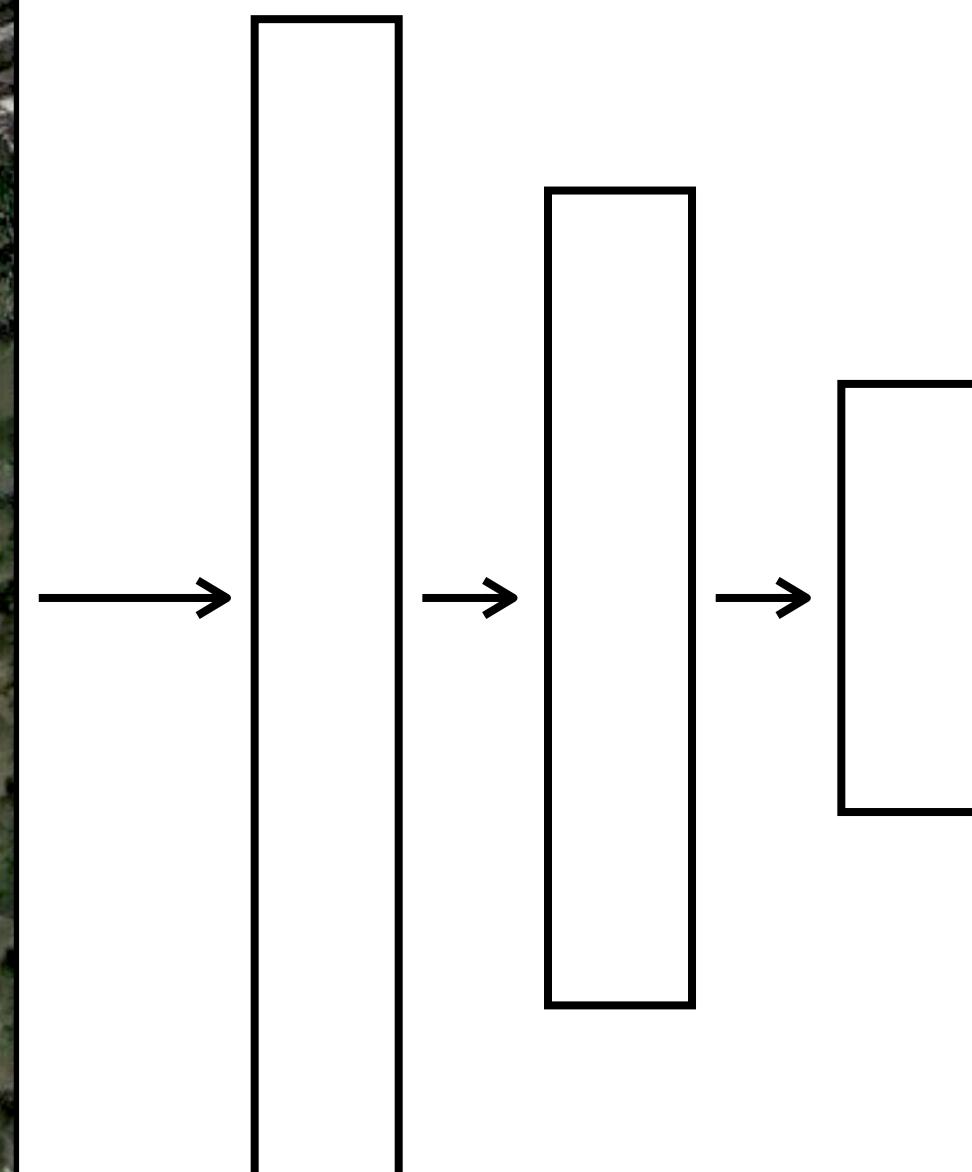
Generator: U-Net



Skip connections between encoder and decoder layers

Architectures

Discriminator: fully convolutional network

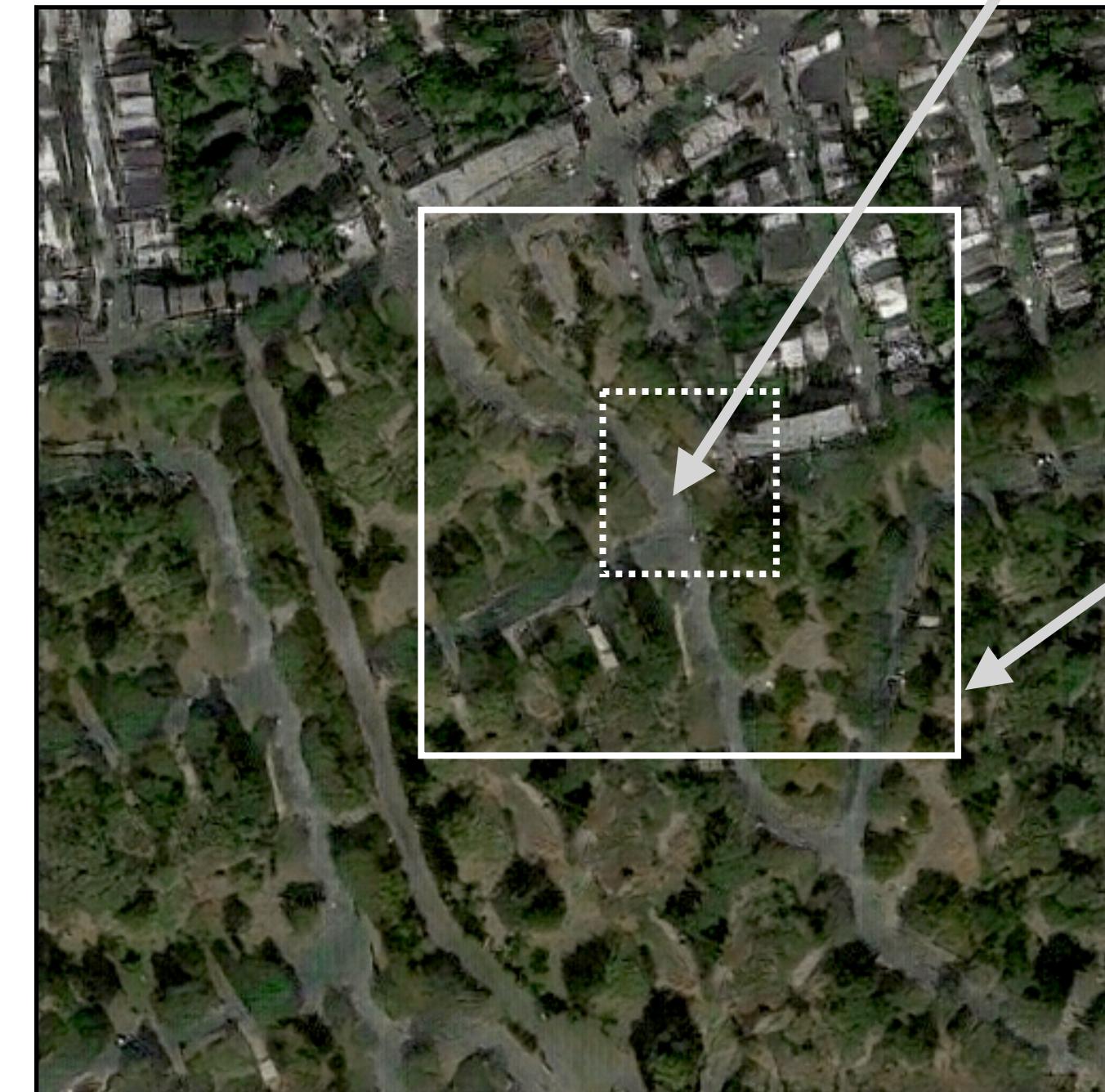
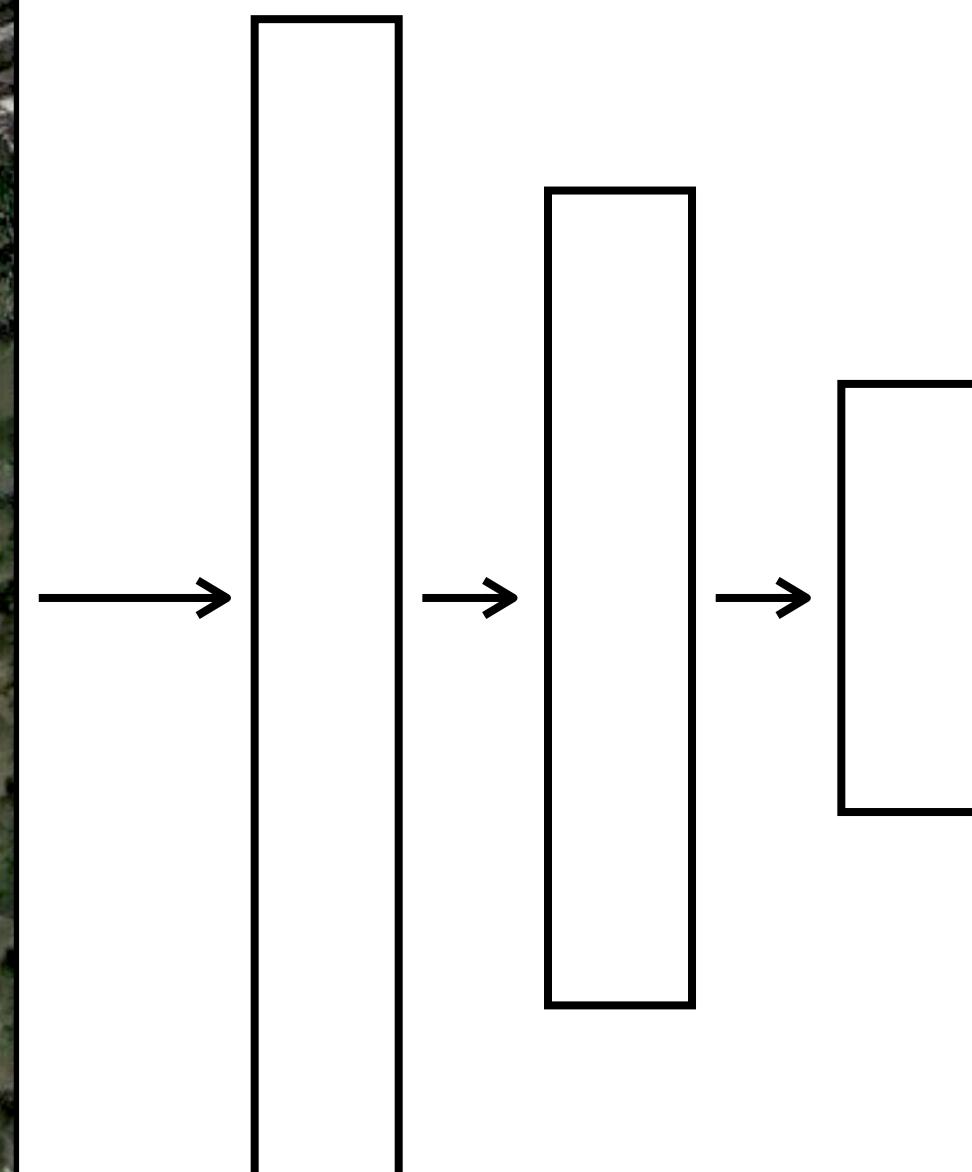


$n \times n$ output map (last conv. layer)

Sequence of strided convolutions

Architectures

Discriminator: fully convolutional network

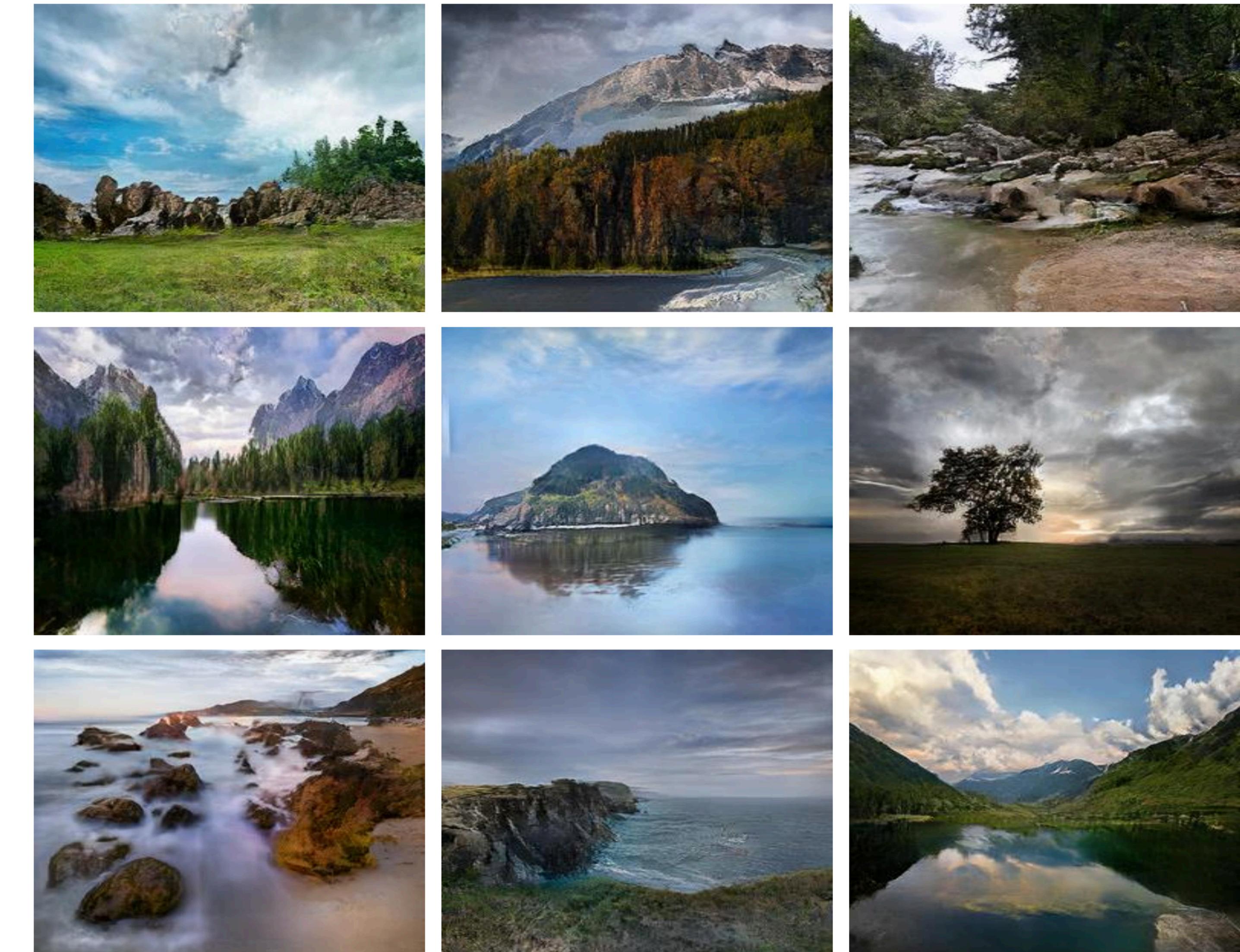
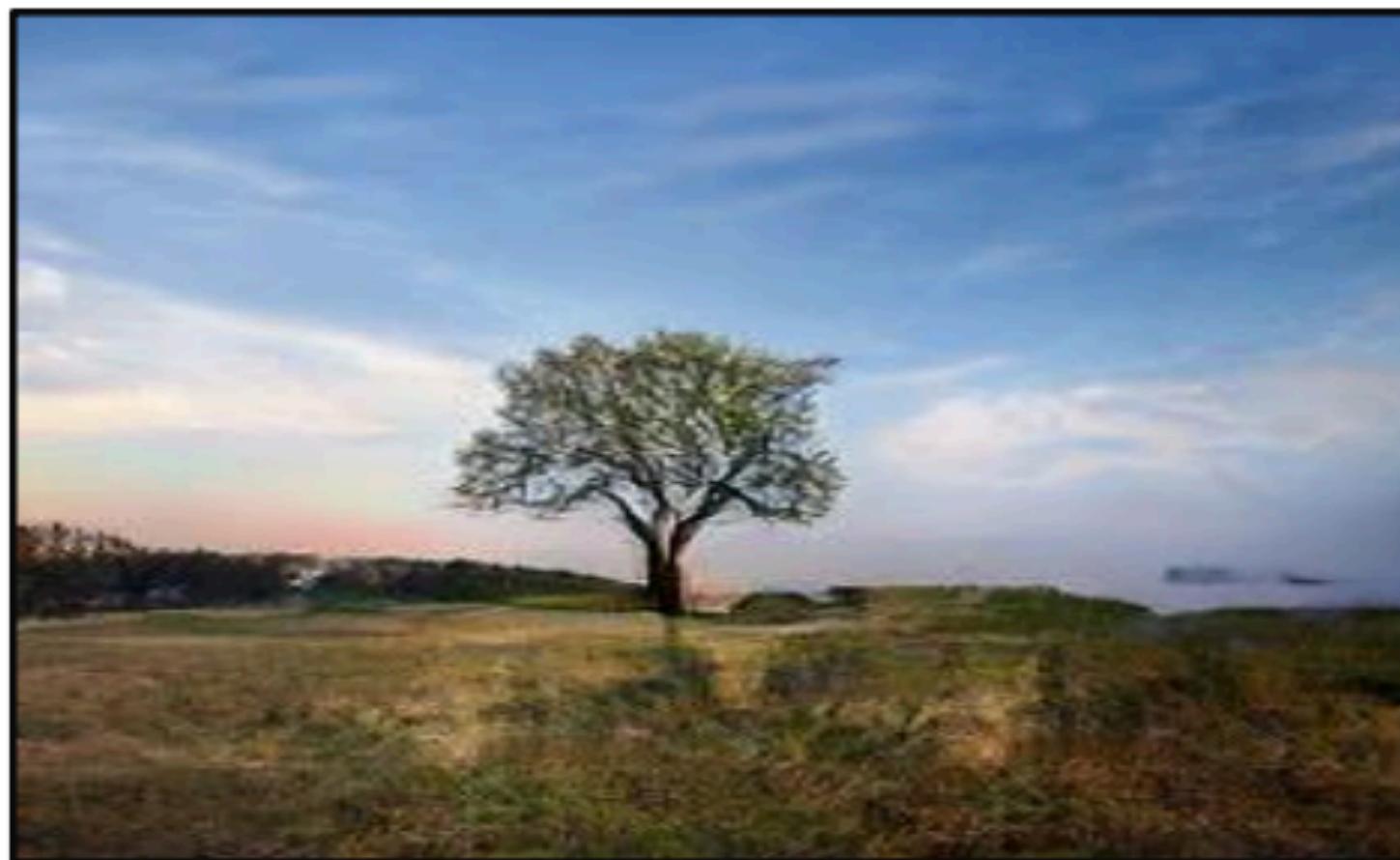
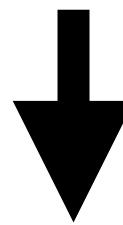


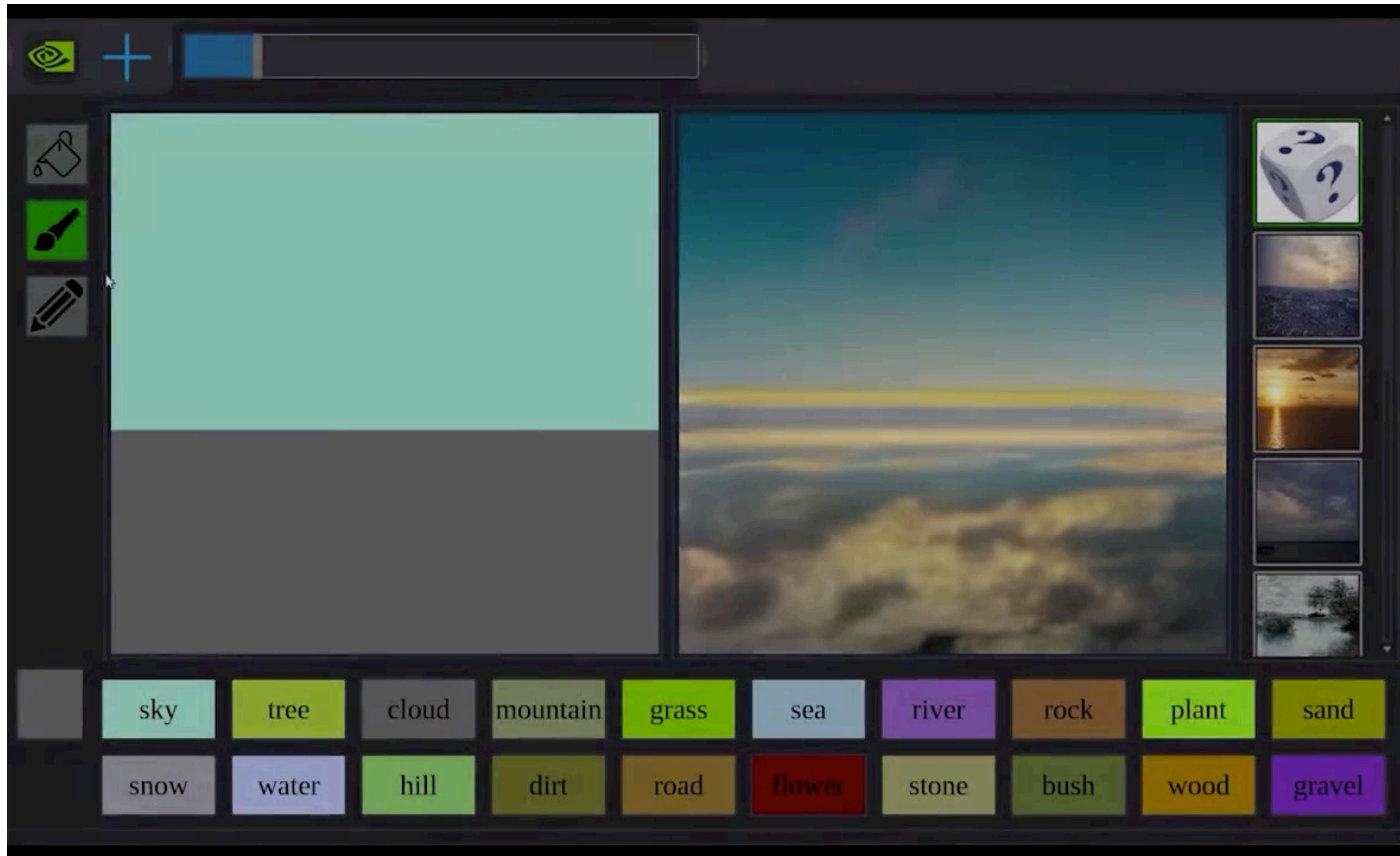
real or fake?

receptive field

Sometimes called a Patch GAN, since it effectively only looks at patches

More recent architectures





What's wrong with GANs?



Category: fish

[Brock et al., "BigGAN", 2018]

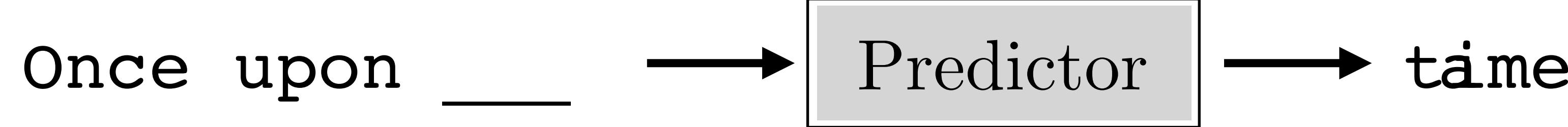
- Optimization challenges
- Lack of diversity
- Mode collapse: only models one “mode” of the probability distribution

Figure source: [Razavi et al., 2019]

Today

- Texture synthesis
- Generative adversarial nets (GANs)
- **Autoregressive models**

Recall: autoregressive models for text

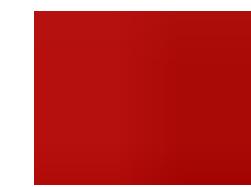


Can we do something analogous for images?

Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time



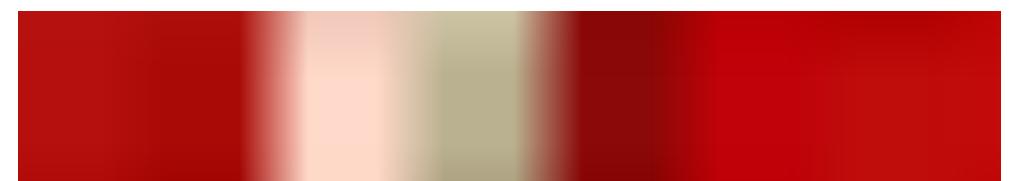
Generating images one pixel at a time



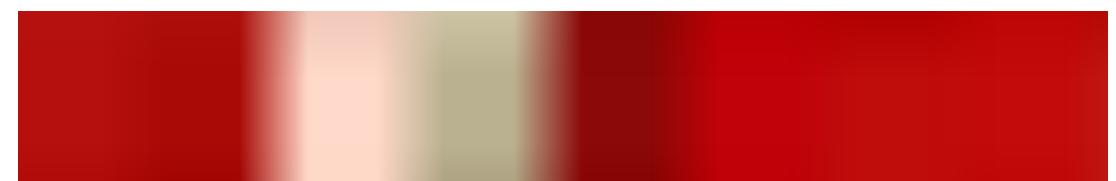
Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time



Generating images one pixel at a time

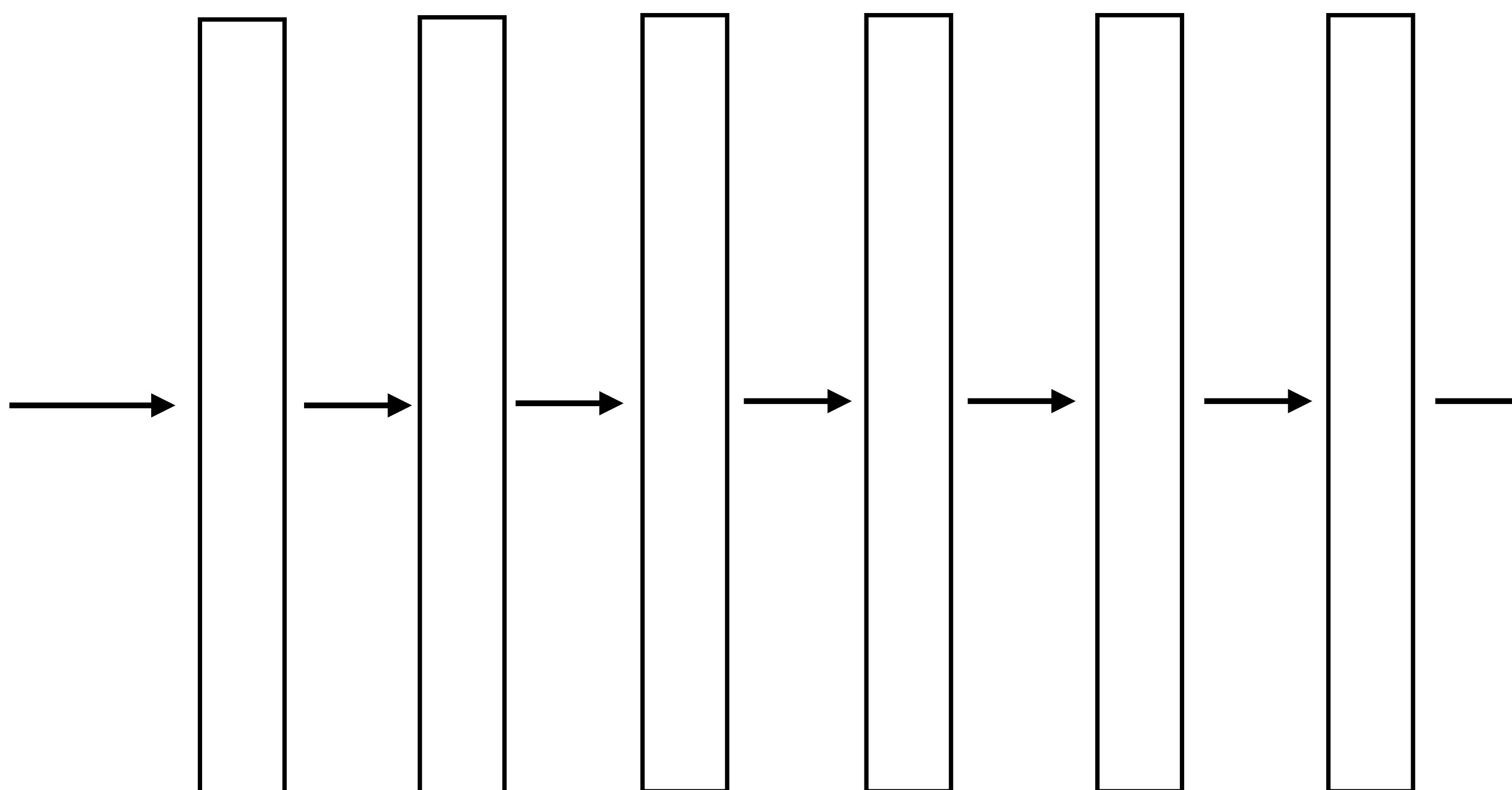
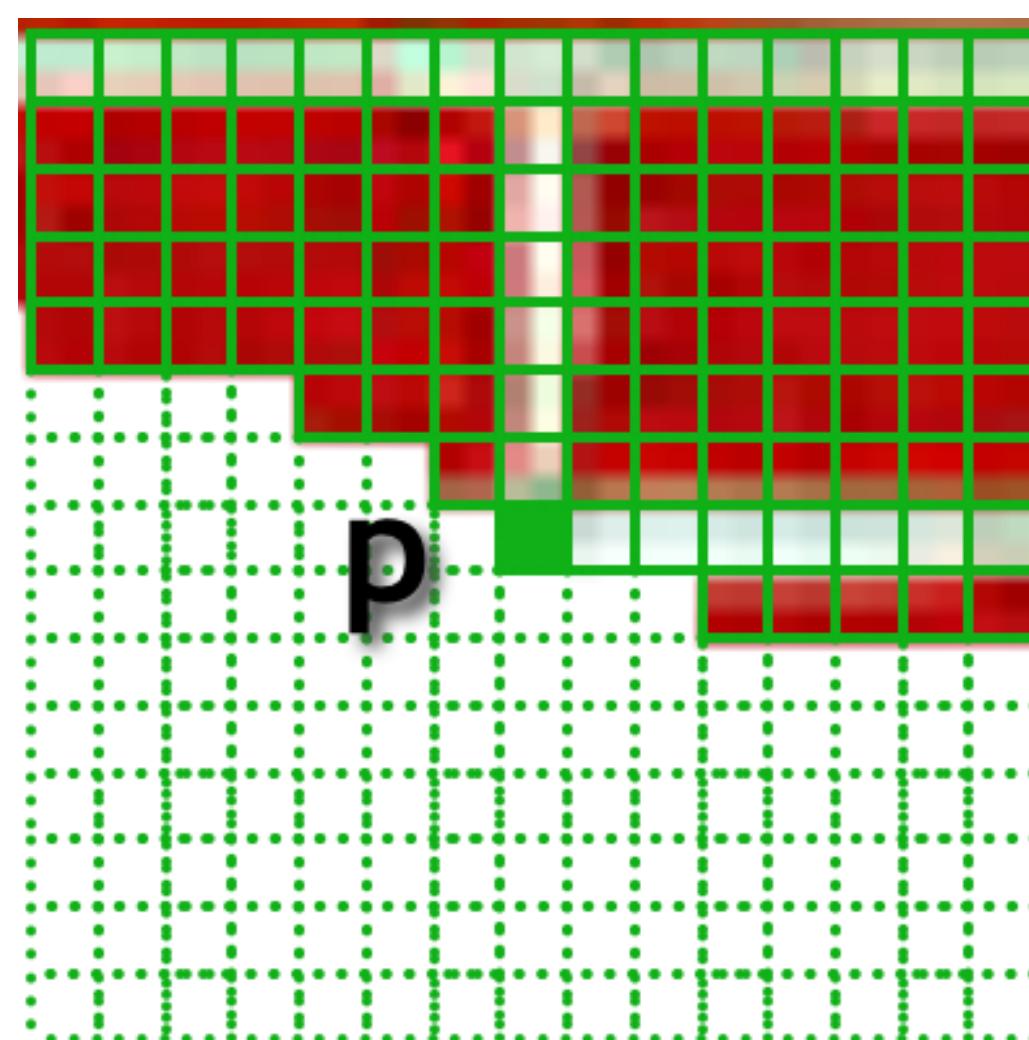


Generating images one pixel at a time



Texture synthesis with a neural network

Input partial
image

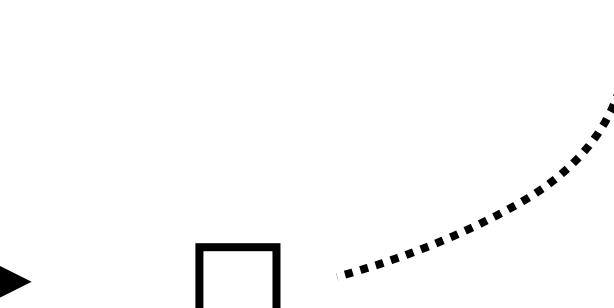
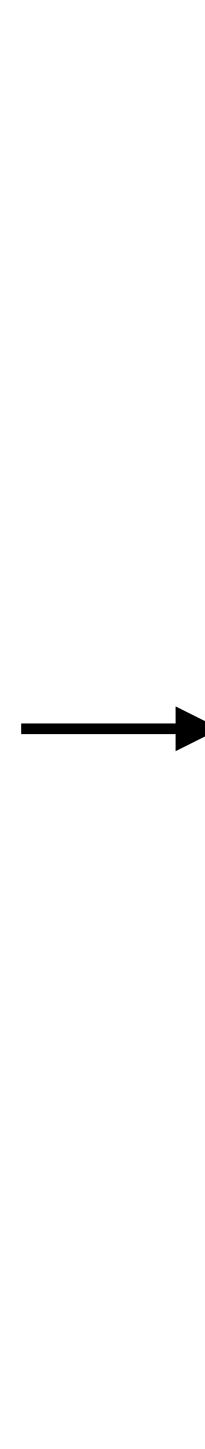
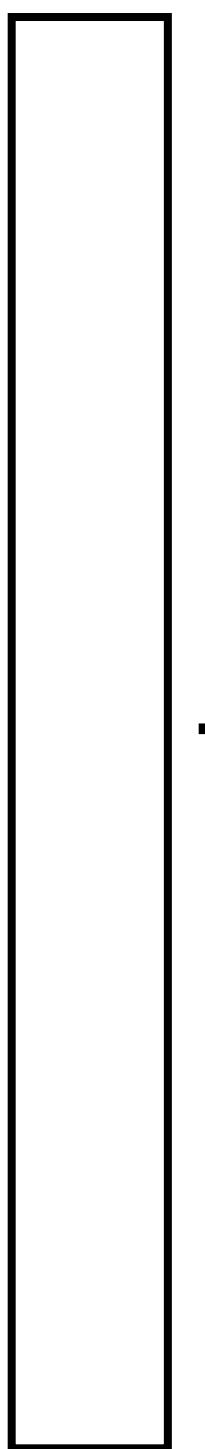
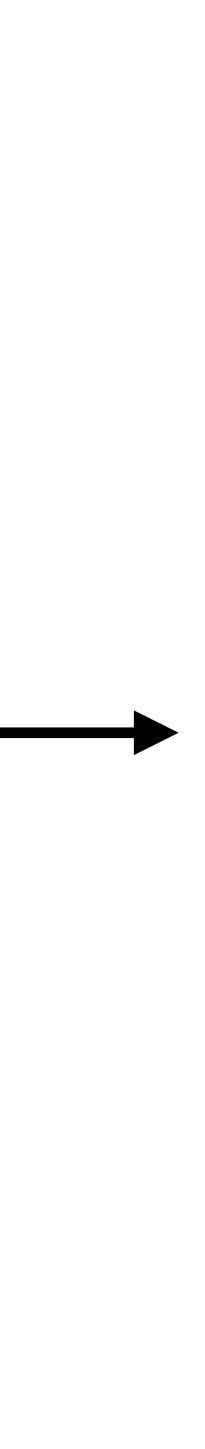
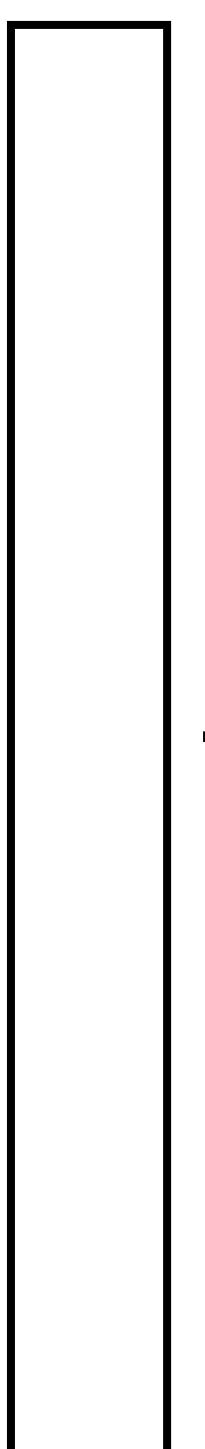
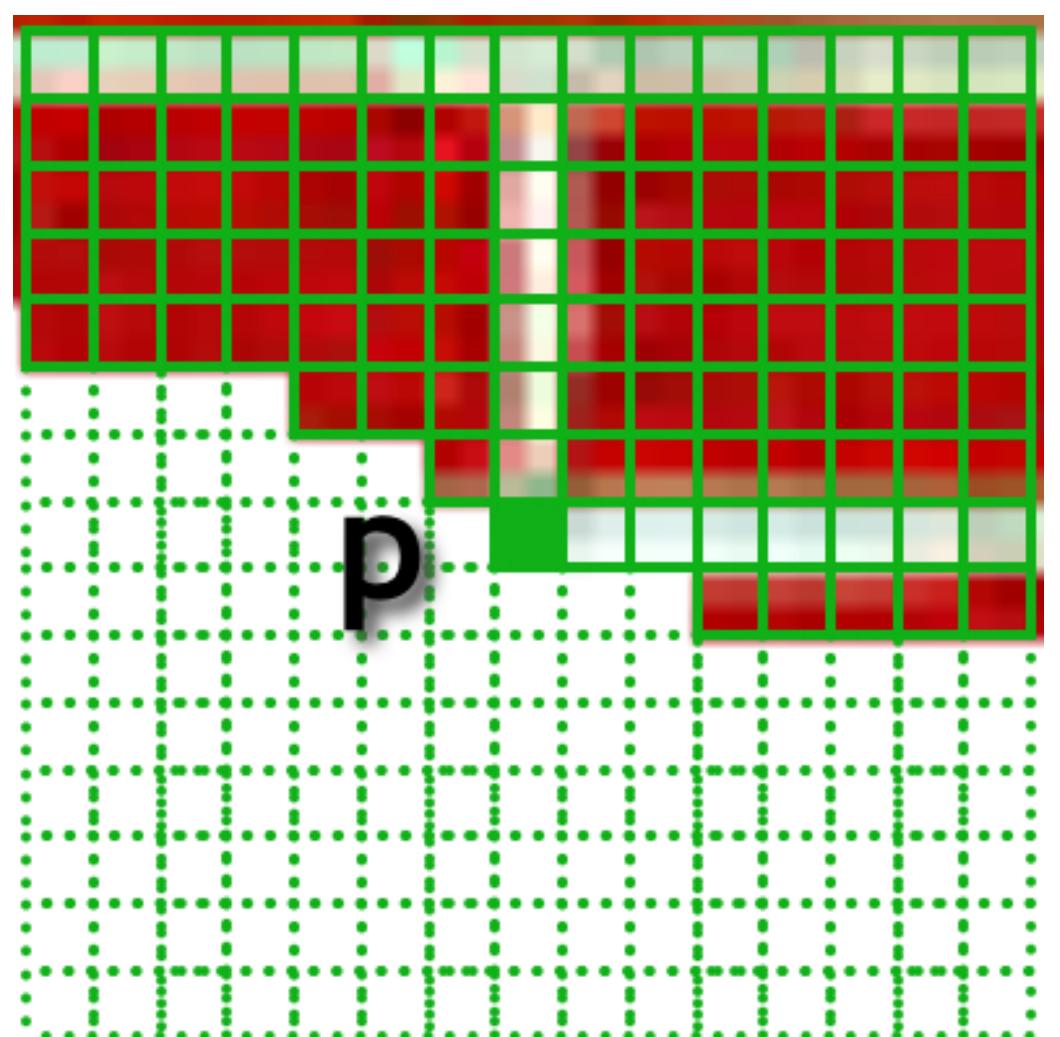


Predicted color
of next pixel

"white"

[PixelRNN, PixelCNN, van der Oord et al. 2016]

Input partial
image



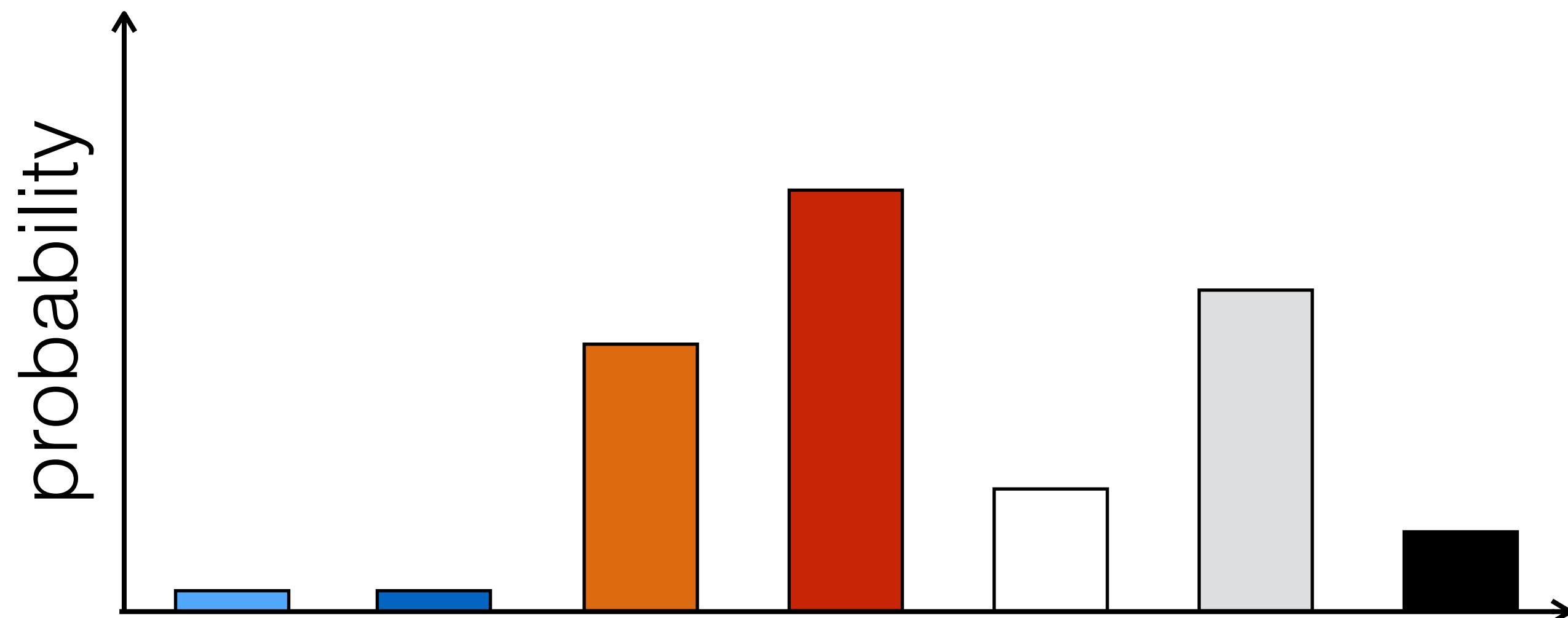
Predicted color
of next pixel

"white"

[PixelRNN, PixelCNN, van der Oord et al. 2016]

Source: Torralba, Freeman, Isola

Represent colors as discrete categories



- One class for each intensity value, e.g. 256 classes
- One label per color channel
- Newer methods use mixture density networks

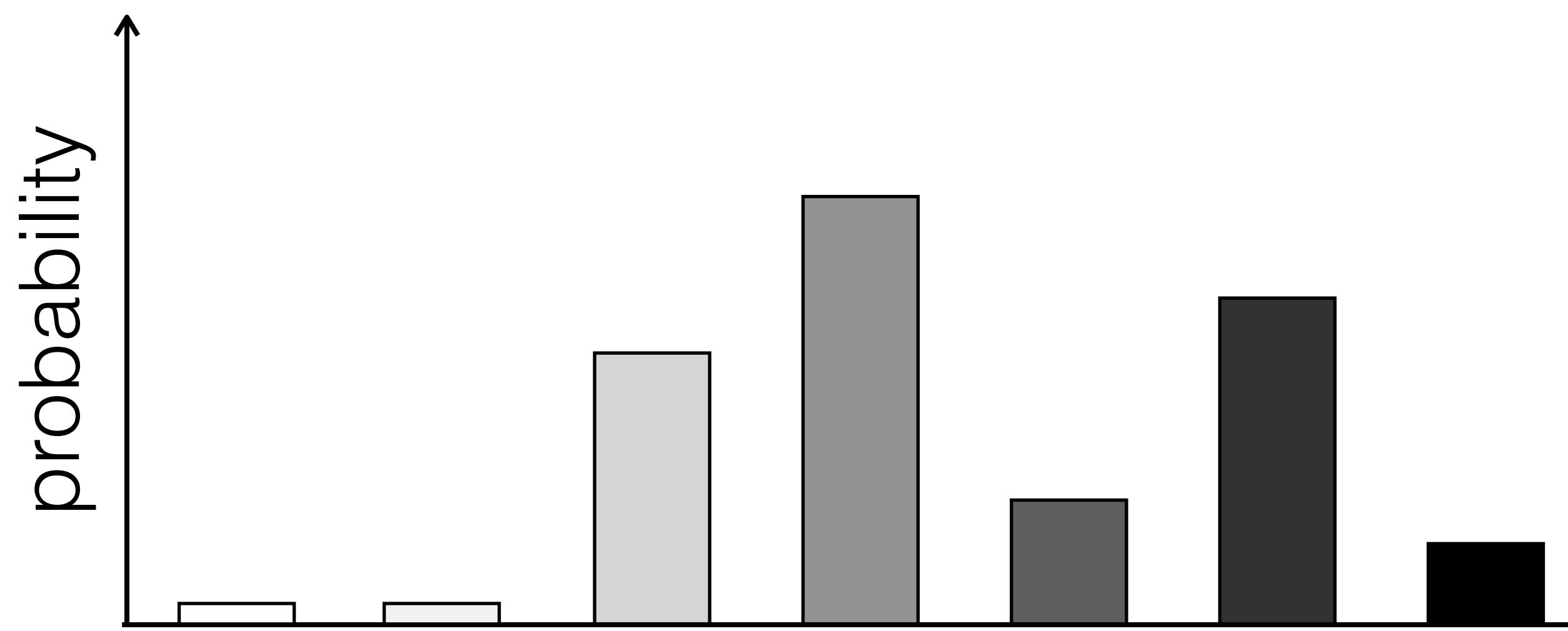
Represent colors as discrete categories

The model is learning: $P(\text{next pixel} \mid \text{previous pixels})$:

$$\hat{\mathbf{y}} \equiv [P_\theta(Y = 1|X = \mathbf{x}), \dots, P_\theta(Y = K|X = \mathbf{x})] \leftarrow \text{predicted probability of each class given input } \mathbf{x}$$

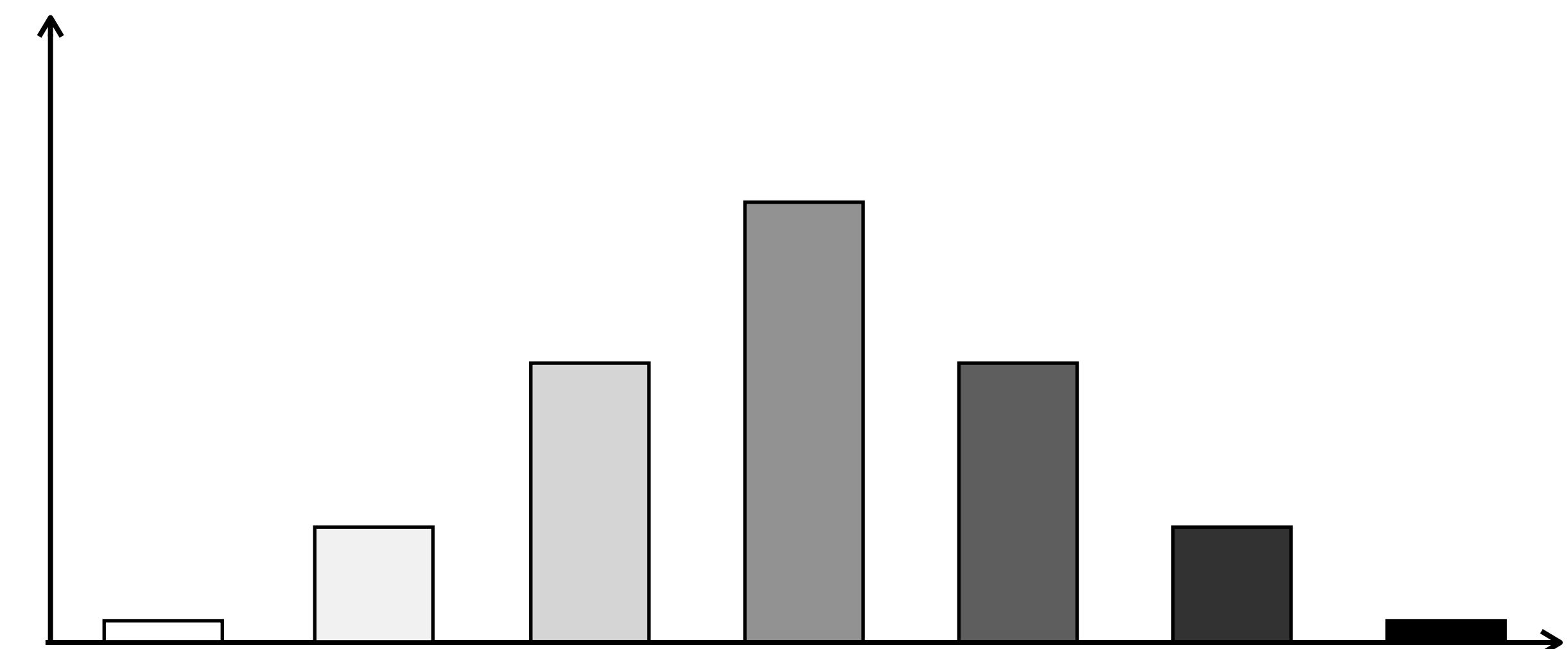
$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k \leftarrow \text{picks out the -log likelihood of the ground truth class } \mathbf{y \text{ under the model prediction } \hat{\mathbf{y}}}$$

Representing distributions



Distribution has multiple modes

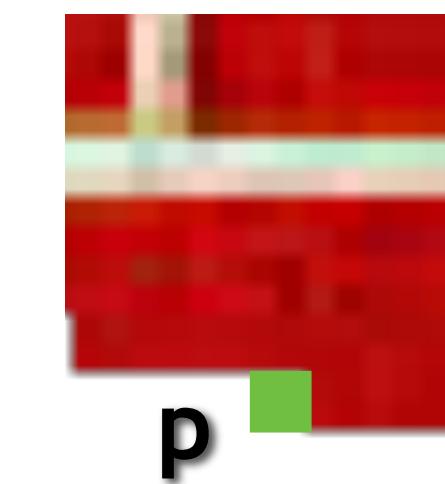
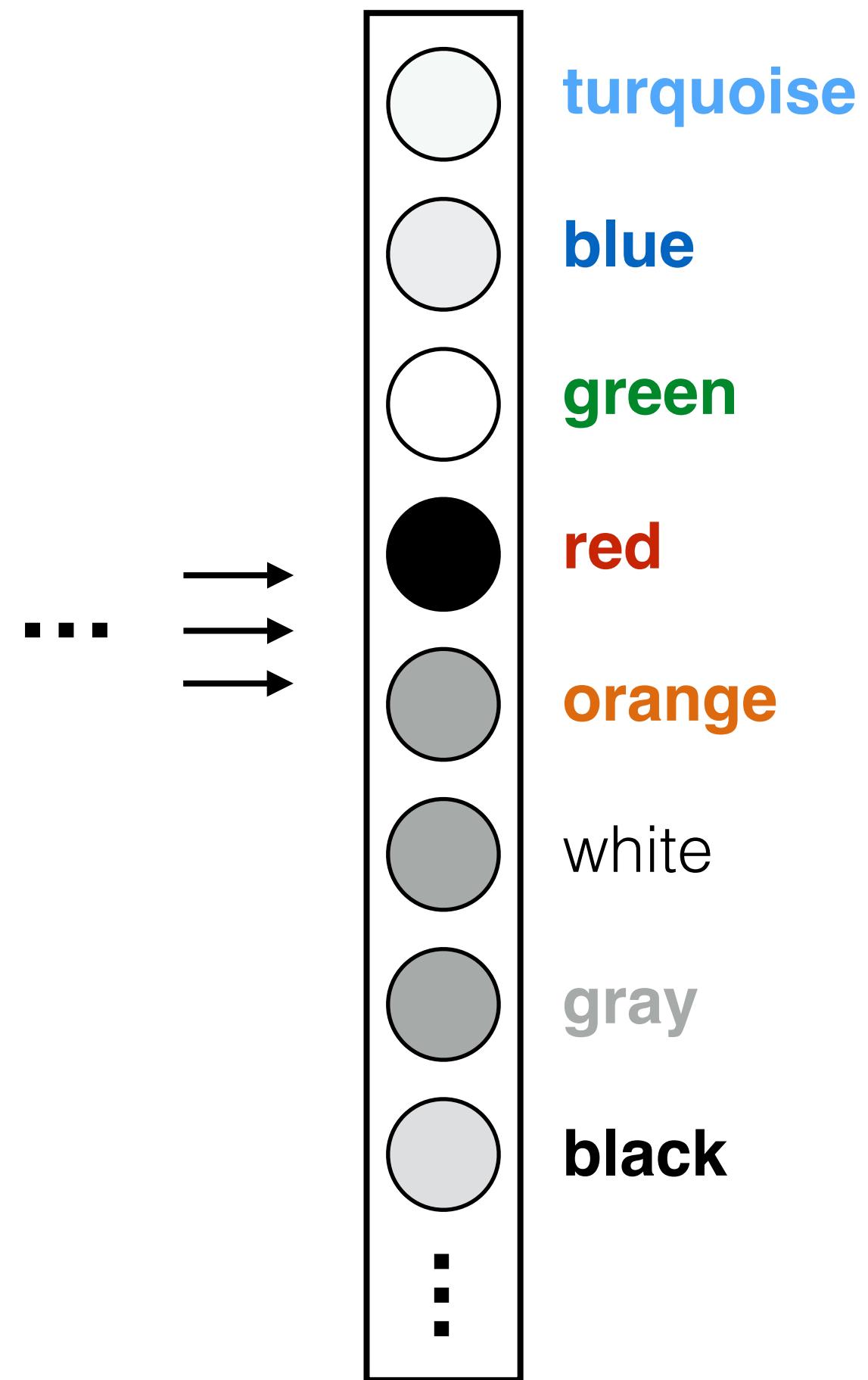
$$H(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{k=1}^K y_k \log \hat{y}_k$$



Single mode

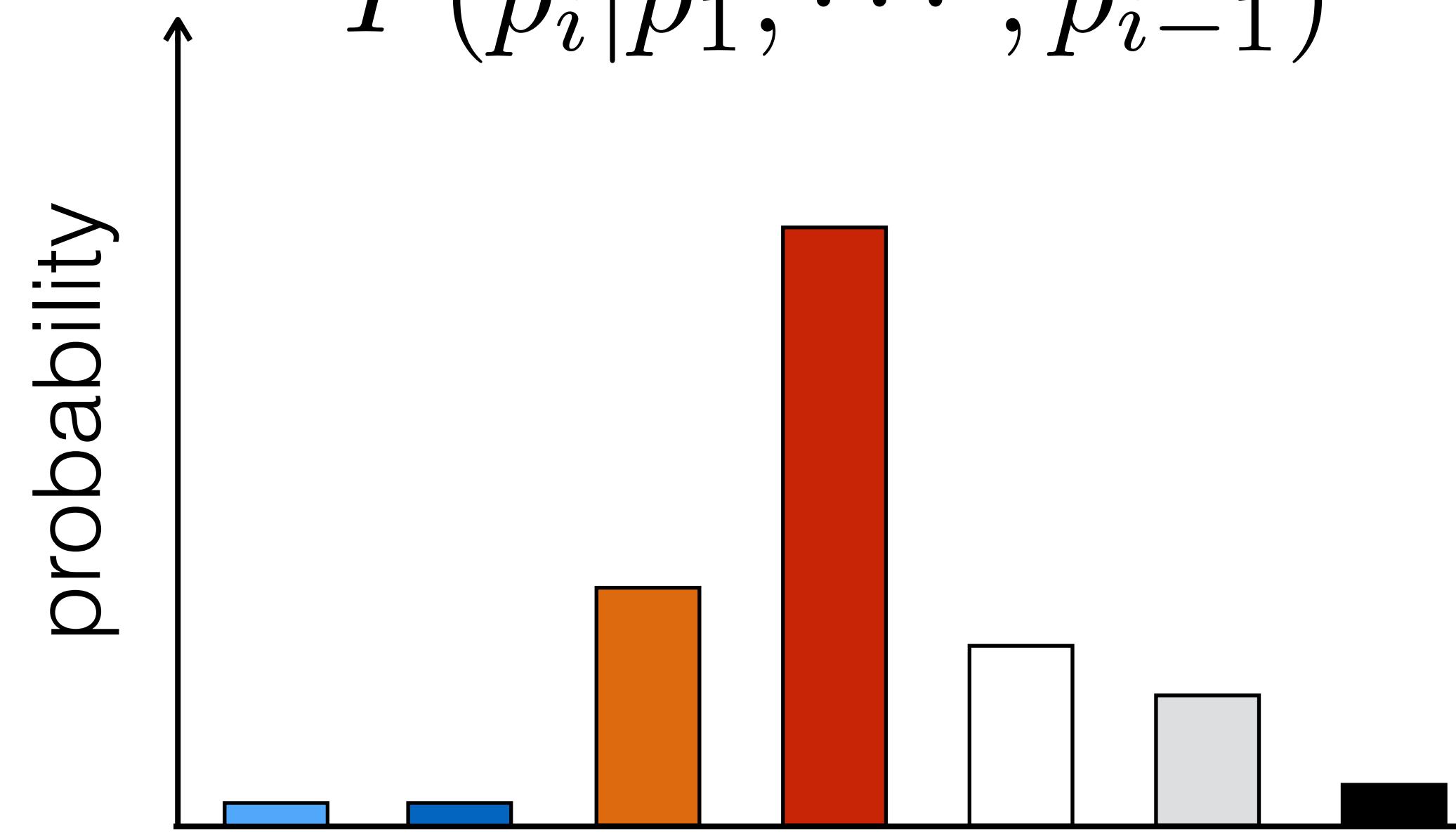
$$L_2(y, \hat{y}) = -\log(P_{\text{gauss}}(\hat{y}|y)) \propto (y - \hat{y})^2$$

Network output

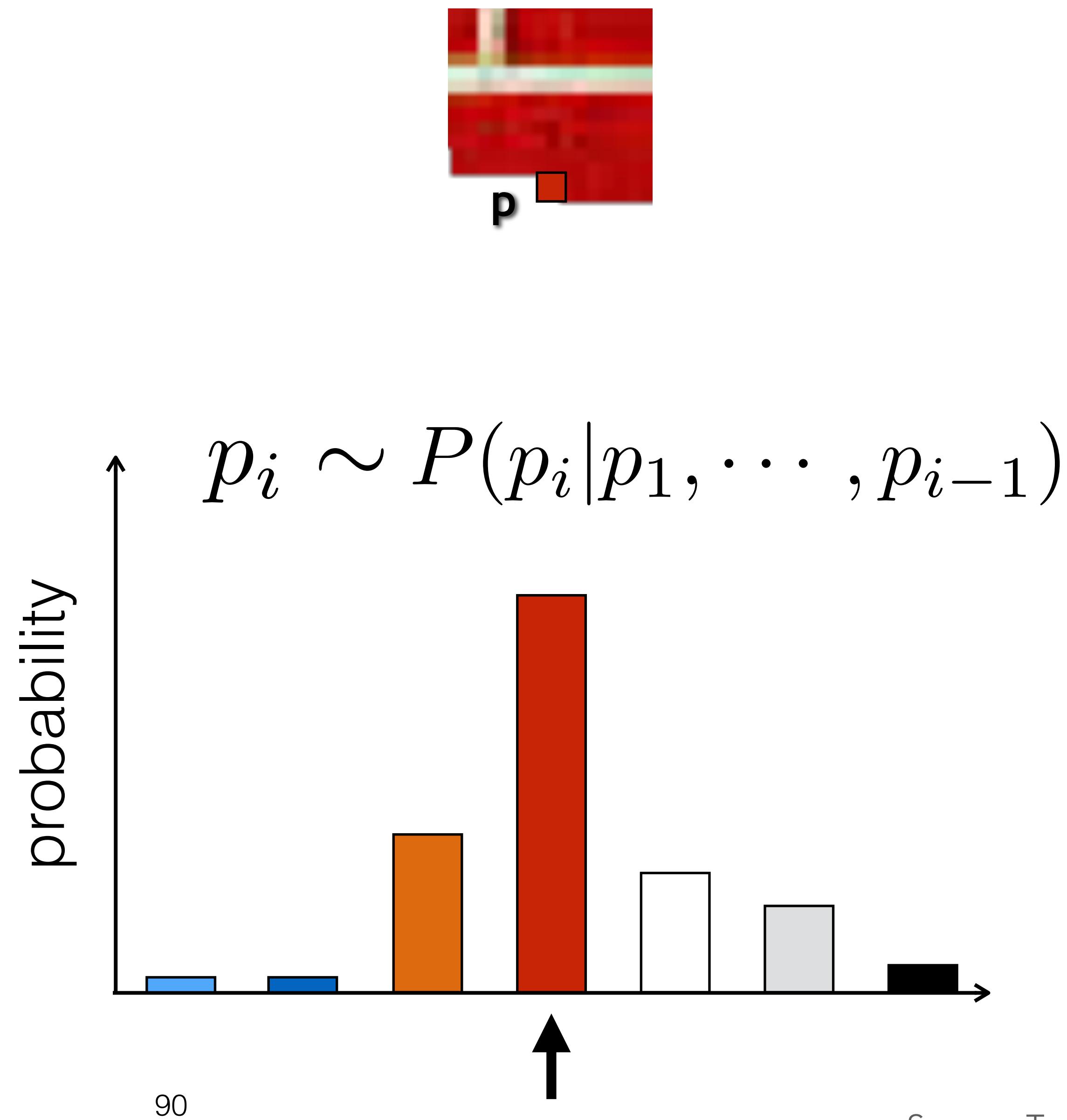
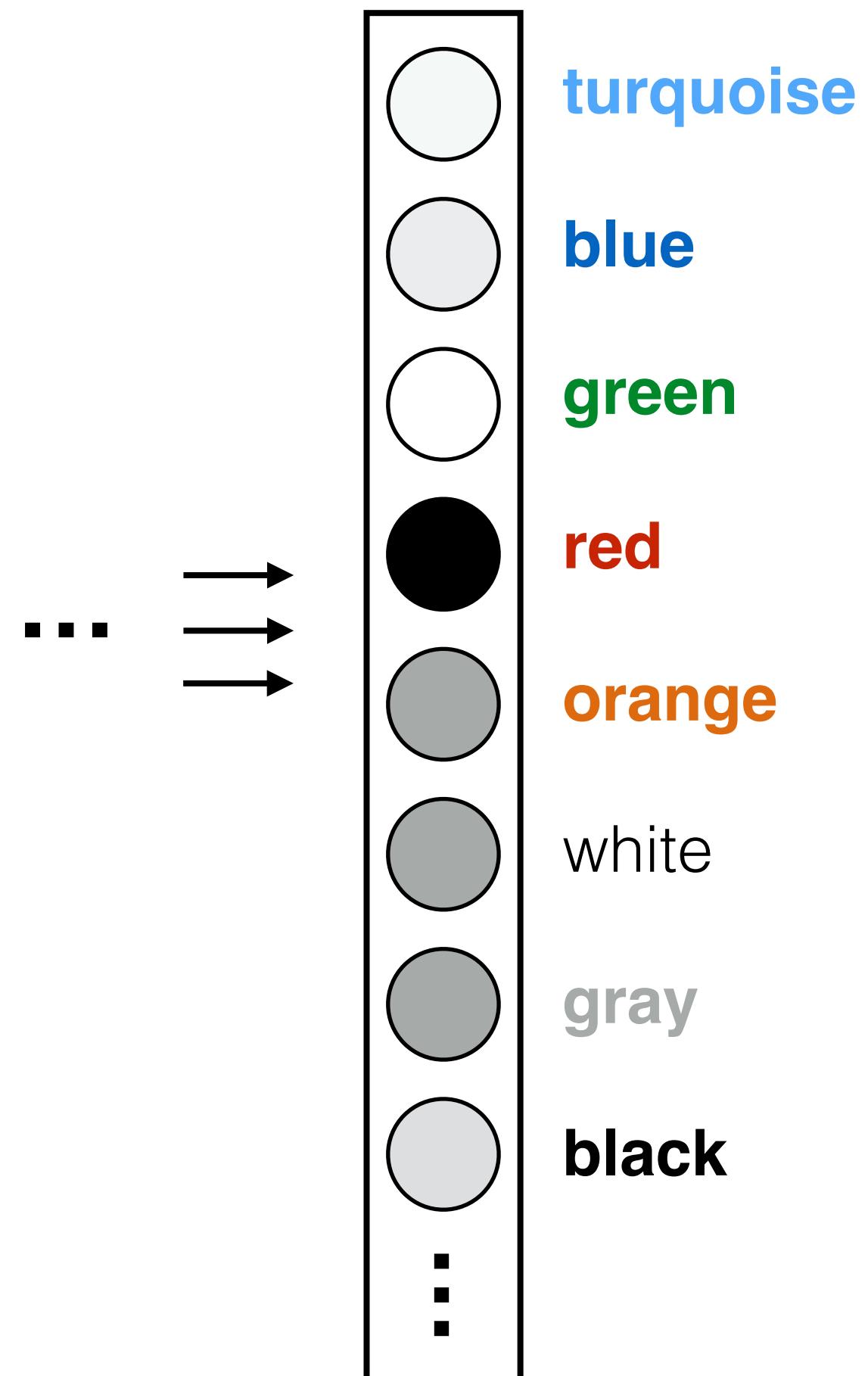


$P(\text{next pixel} \mid \text{previous pixels})$

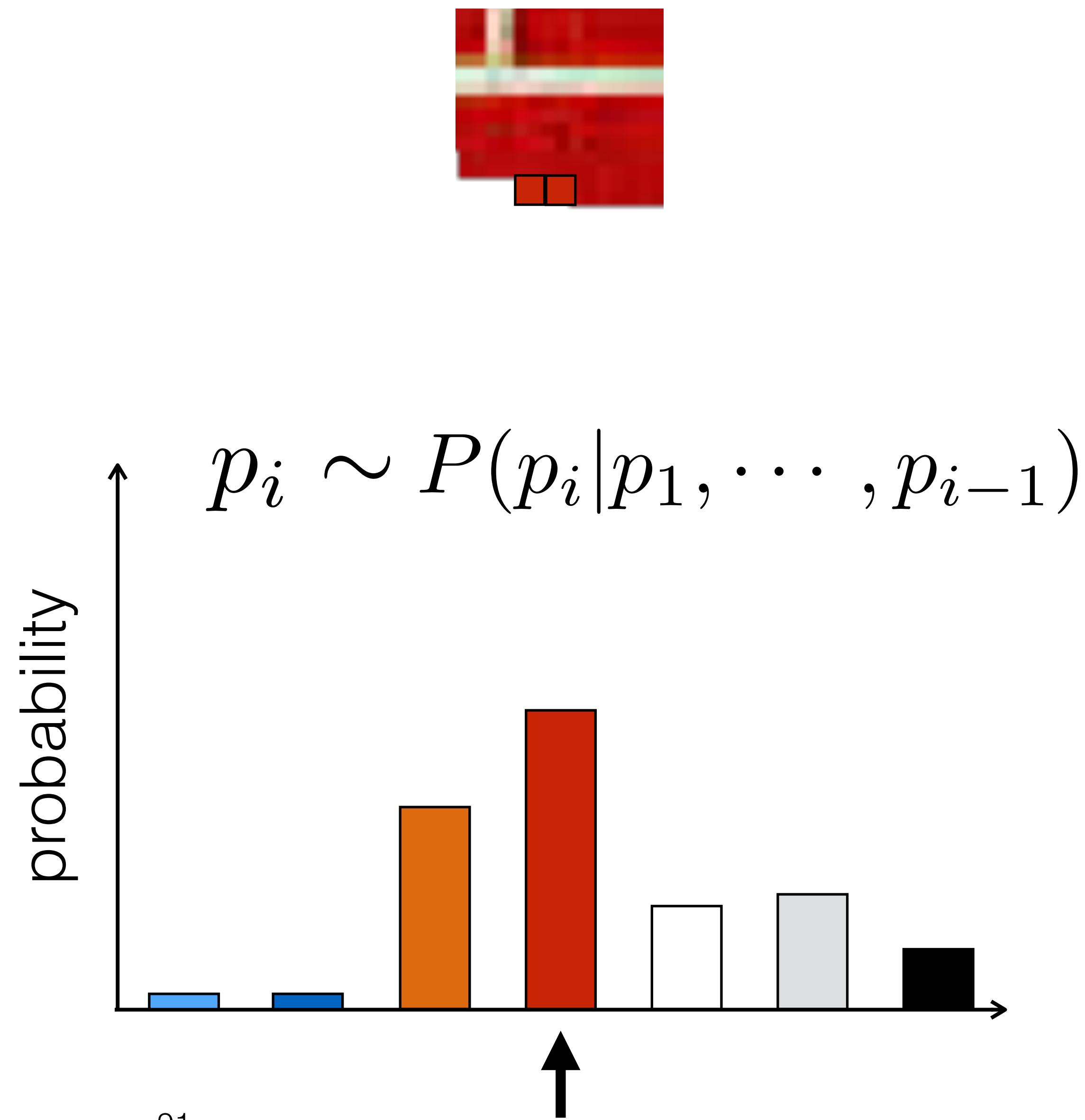
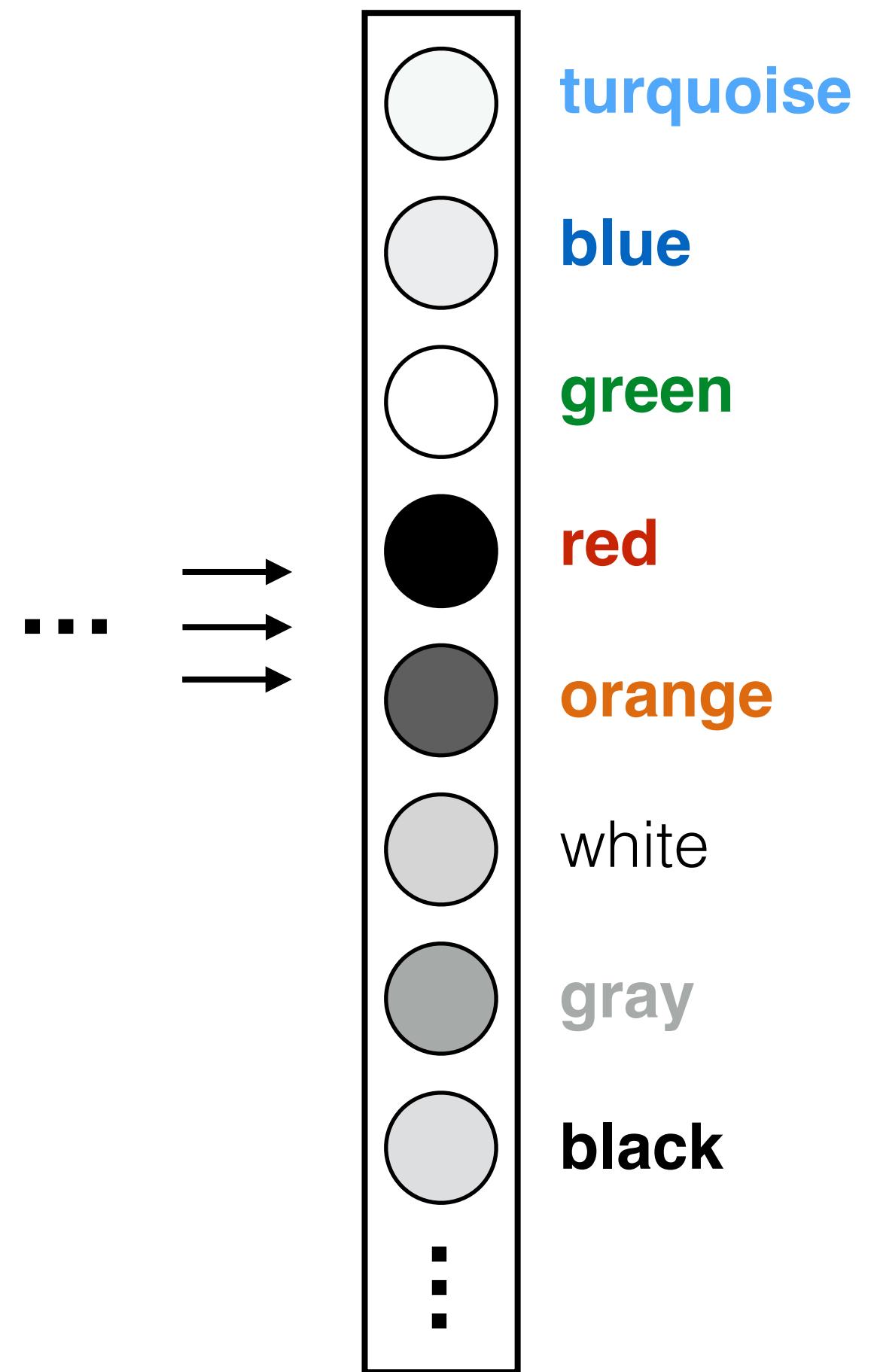
$$P(p_i \mid p_1, \dots, p_{i-1})$$



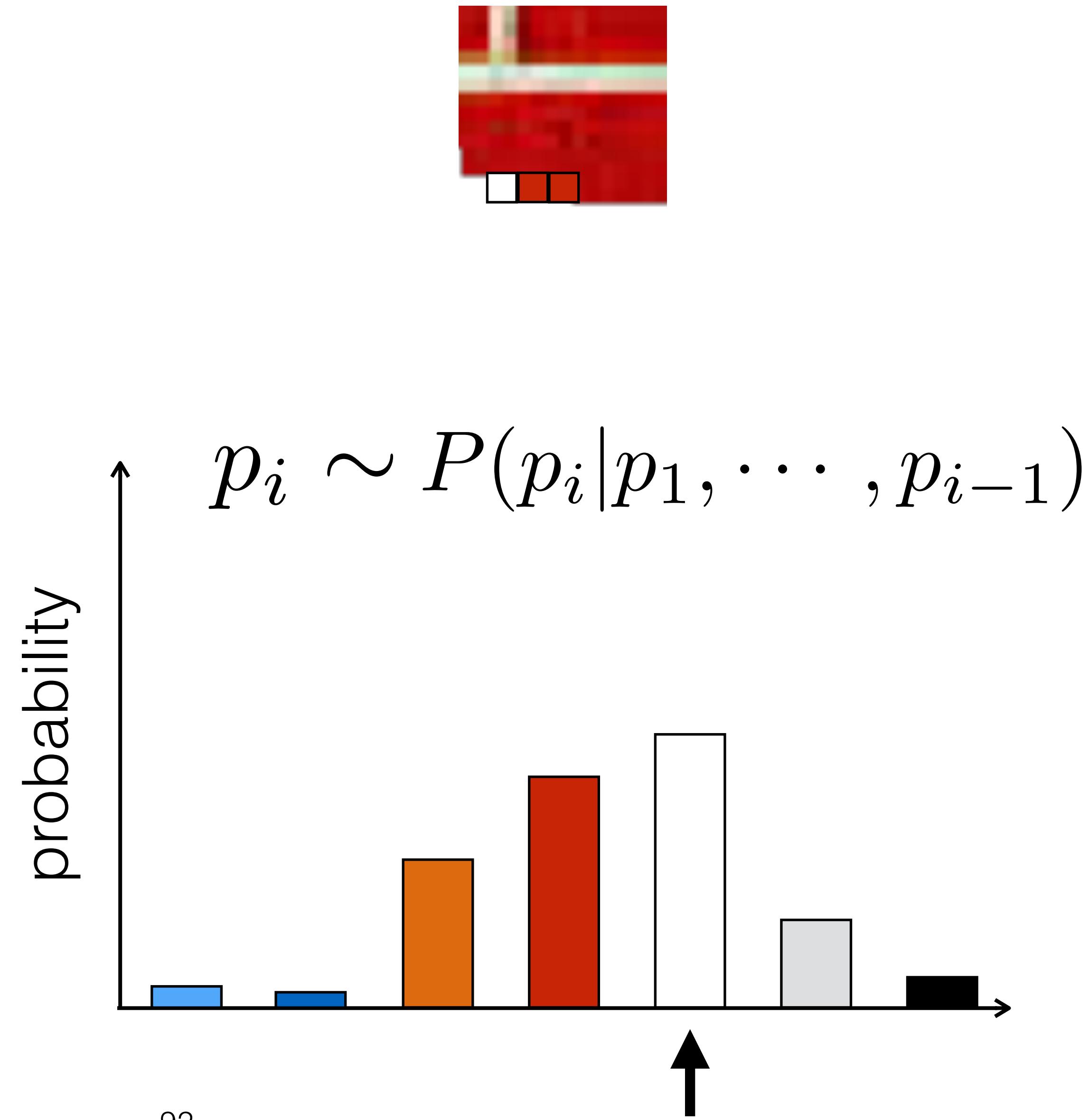
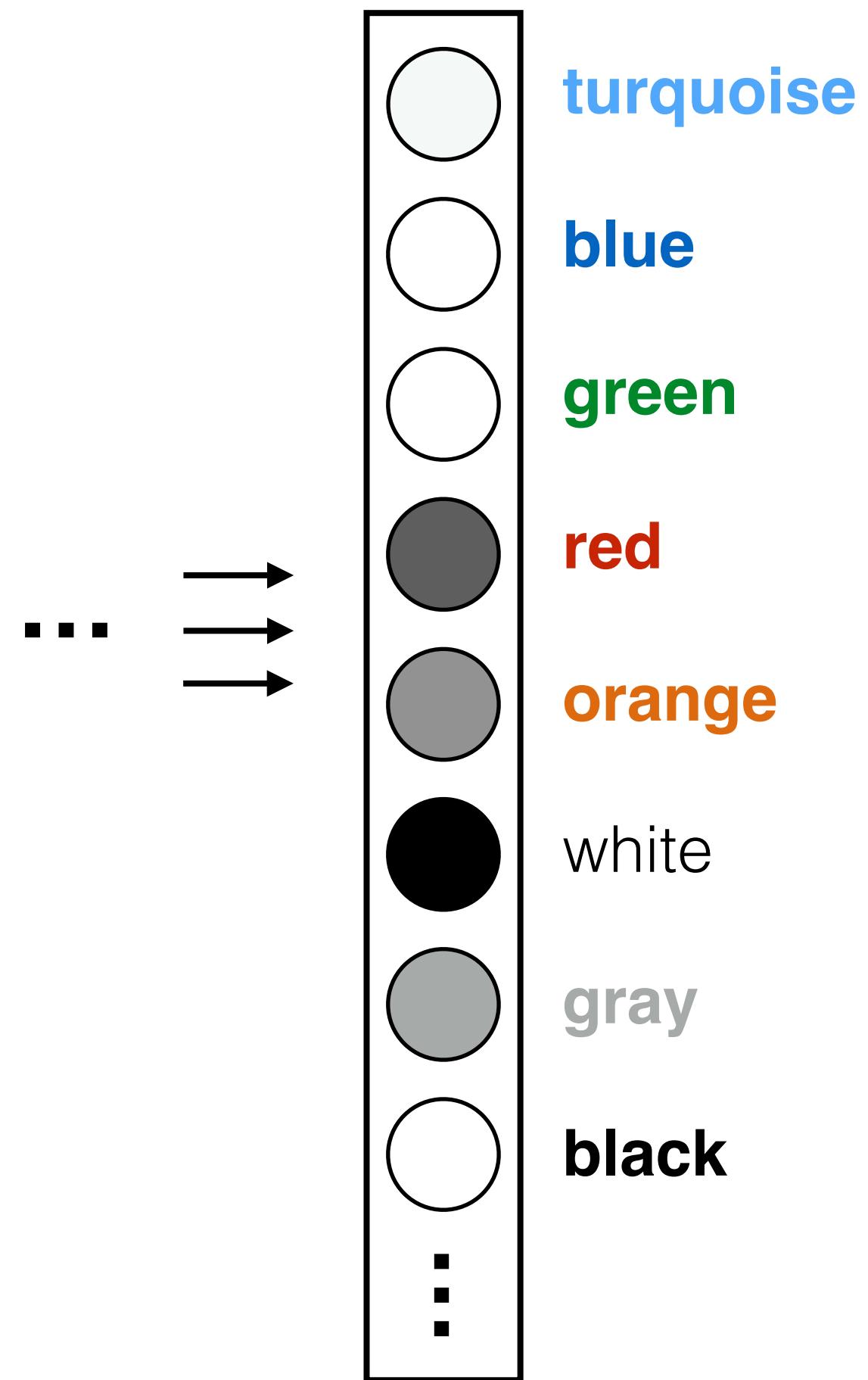
Network output



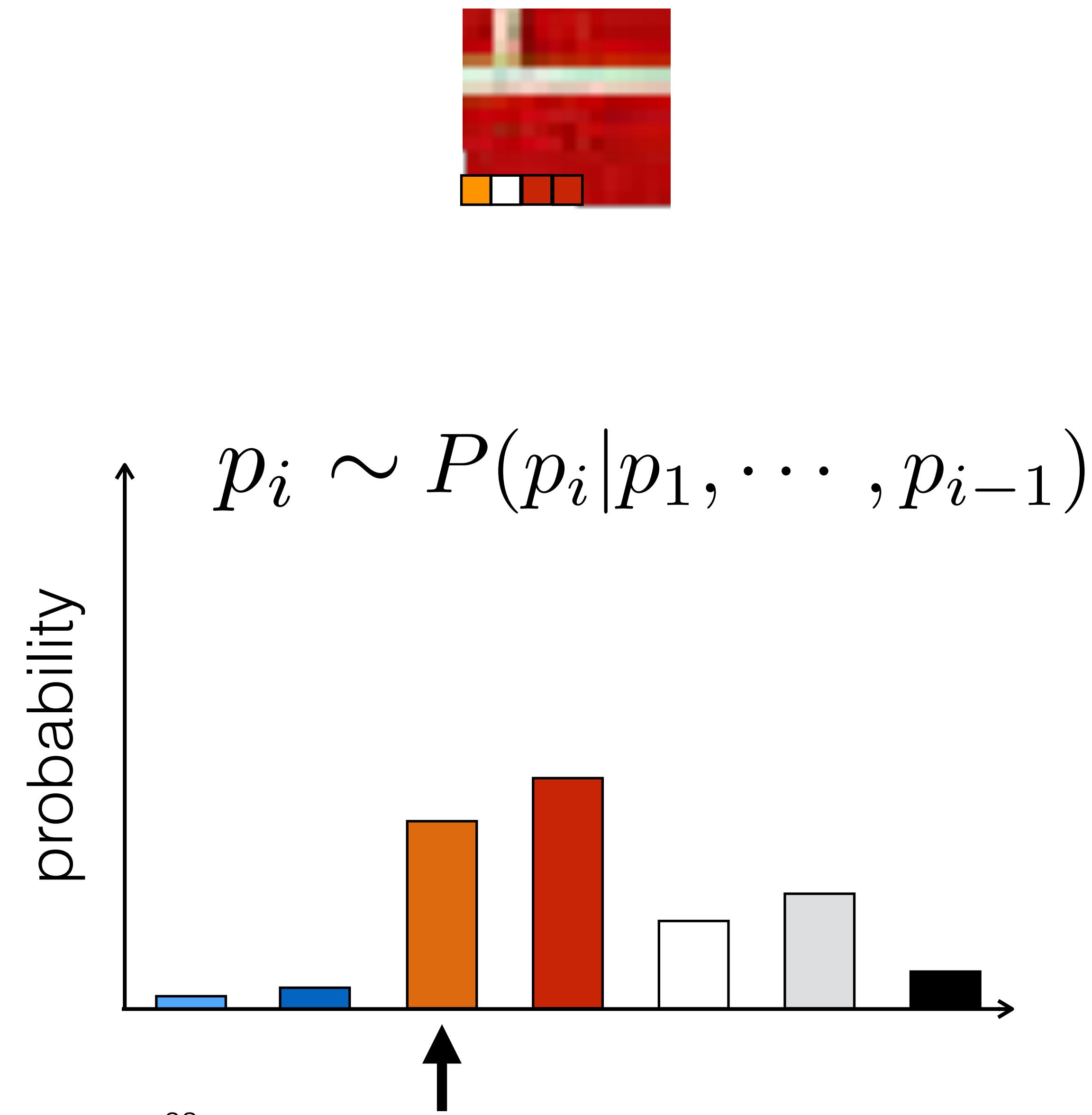
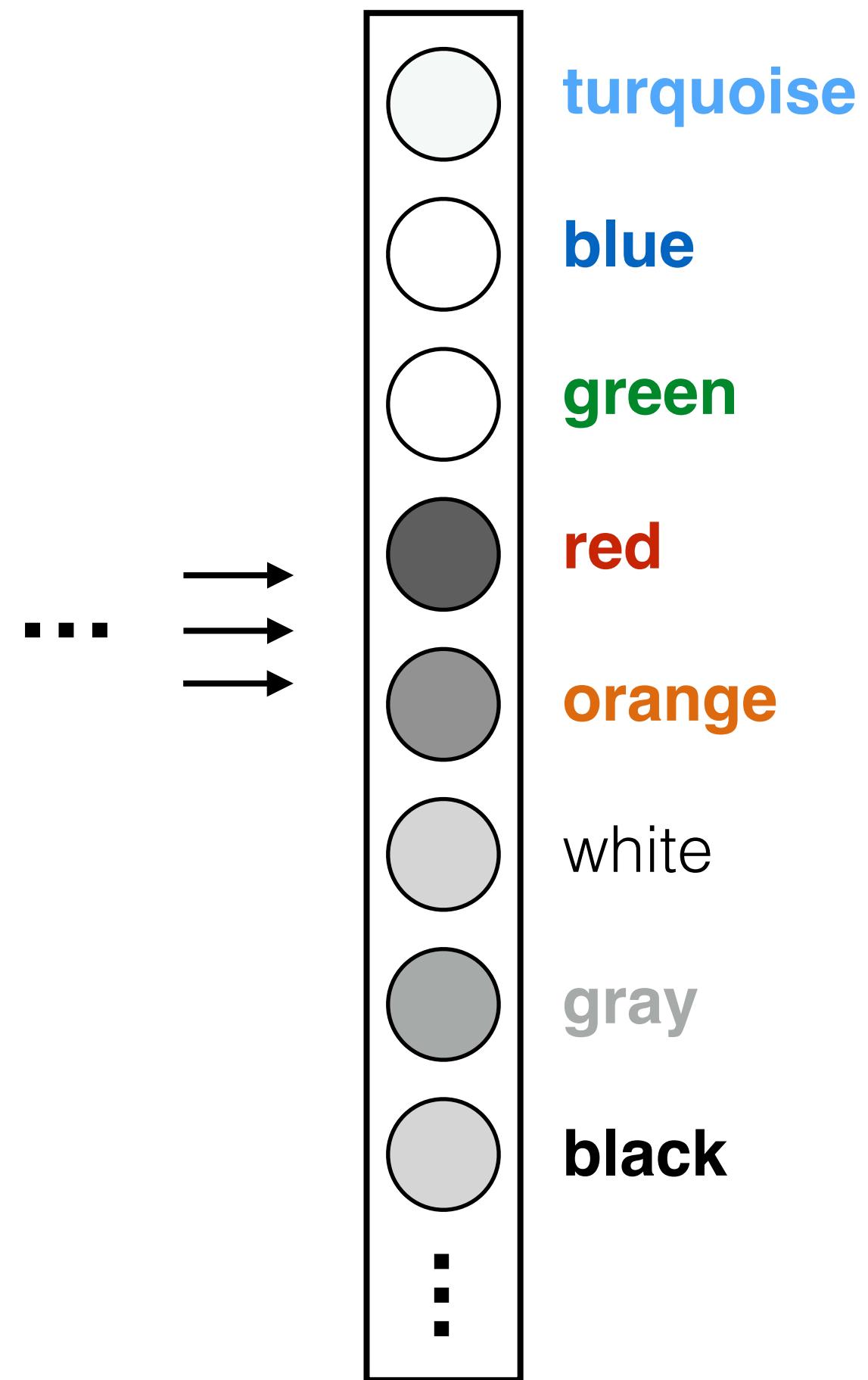
Network output



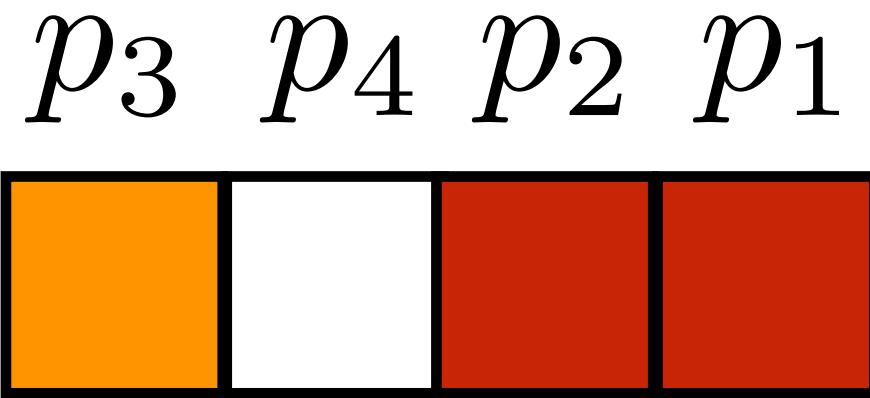
Network output



Network output



$$p_1 \sim P(p_1)$$



$$p_2 \sim P(p_2|p_1)$$

$$p_3 \sim P(p_3|p_1, p_2)$$

$$p_4 \sim P(p_4|p_1, p_2, p_3)$$

$$\{p_1, p_2, p_3, p_4\} \sim P(p_4|p_1, p_2, p_3)P(p_3|p_1, p_2)P(p_2|p_1)P(p_1)$$

$$p_i \sim P(p_i|p_1, \dots, p_{i-1})$$

$$\mathbf{p} \sim \prod_{i=1}^N P(p_i|p_1, \dots, p_{i-1})$$

Autoregressive probability model

$$\mathbf{p} \sim \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1})$$

$$P(\mathbf{p}) = \prod_{i=1}^N P(p_i | p_1, \dots, p_{i-1}) \quad \leftarrow \textbf{General product rule}$$

The sampling procedure we defined above takes exact samples from the learned probability distribution (pdf).

Multiplying all conditionals evaluates the probability of a full joint configuration of pixels.

Samples from PixelRNN

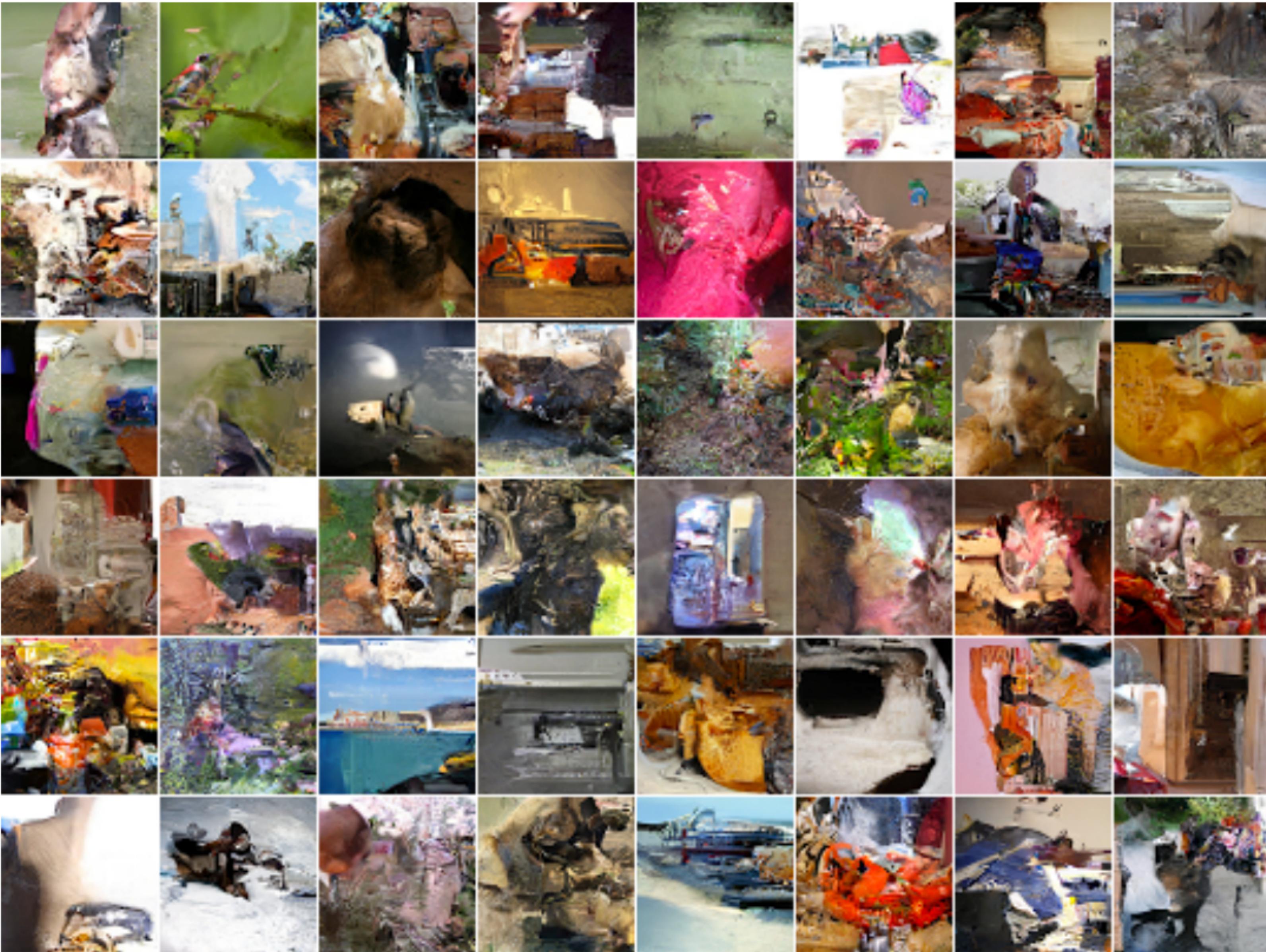


Image completions (conditional samples) from PixelRNN

occluded

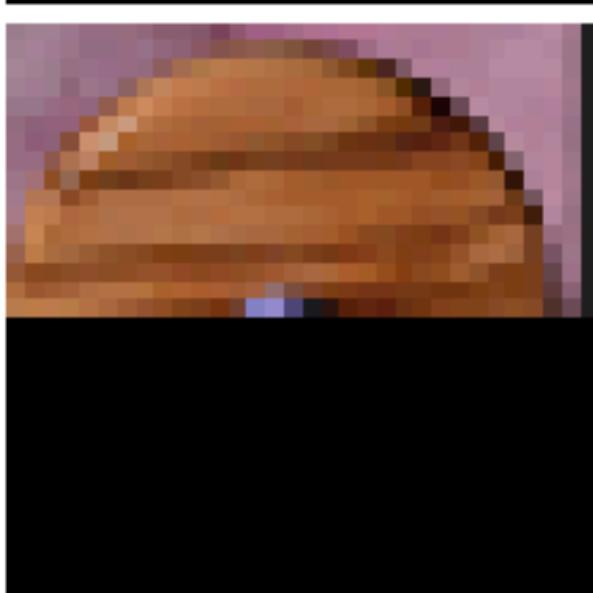
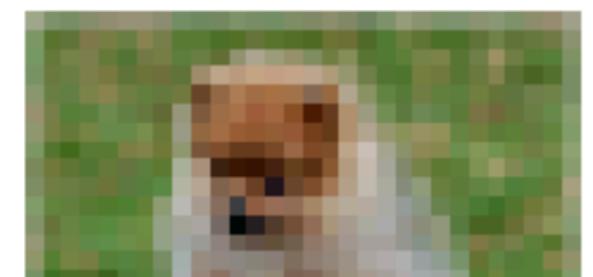
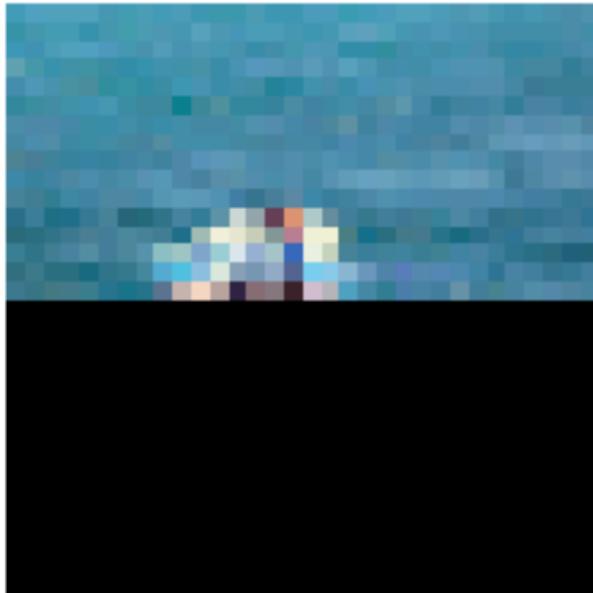


Image completions (conditional samples) from PixelRNN

occluded

completions

original



Image completions (conditional samples) from PixelRNN

occluded



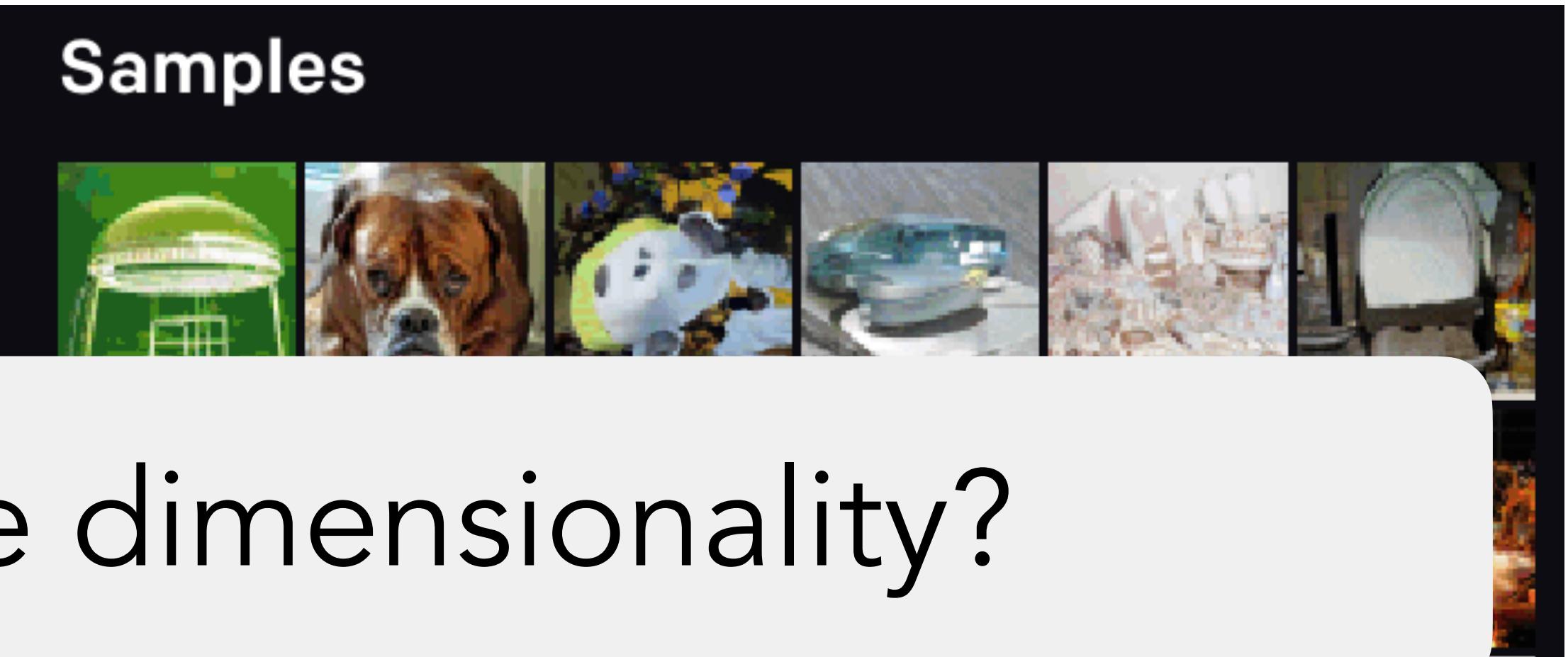
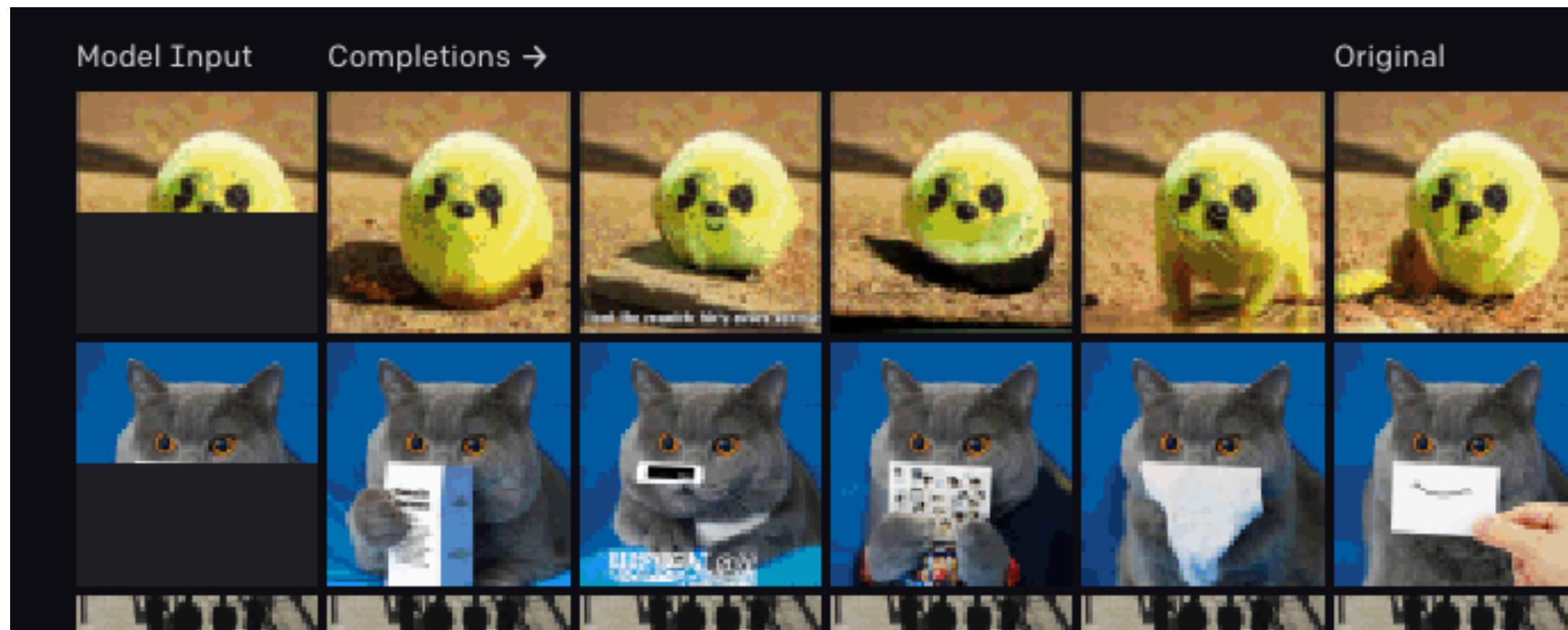
completions



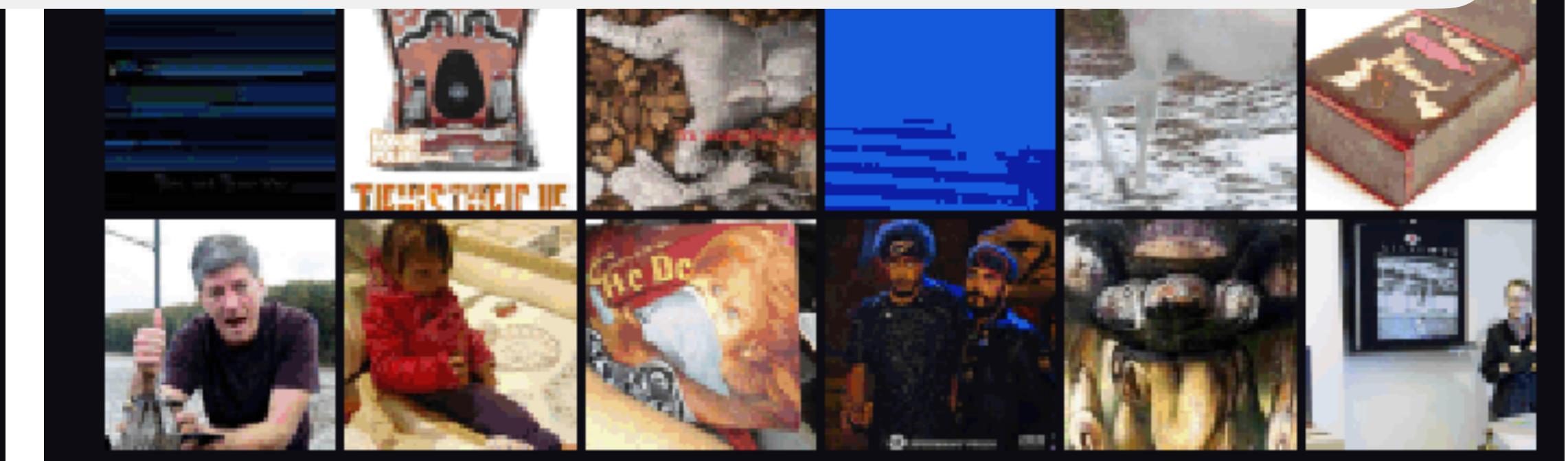
original



Hard to scale up, though...

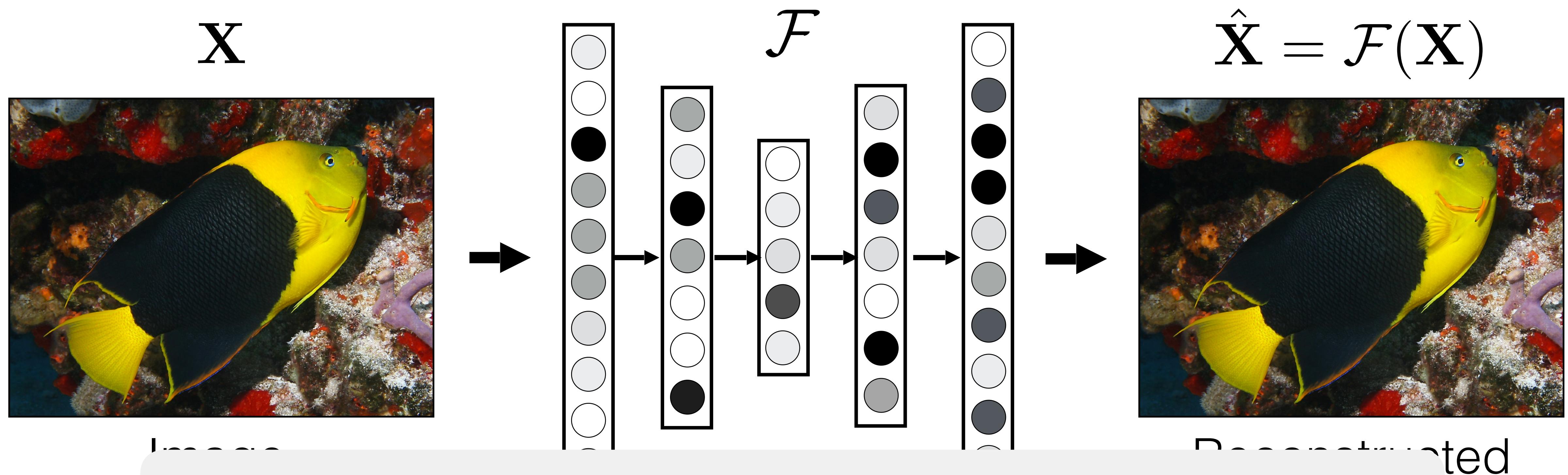


Can we lower the dimensionality?



<https://openai.com/blog/image-gpt/>

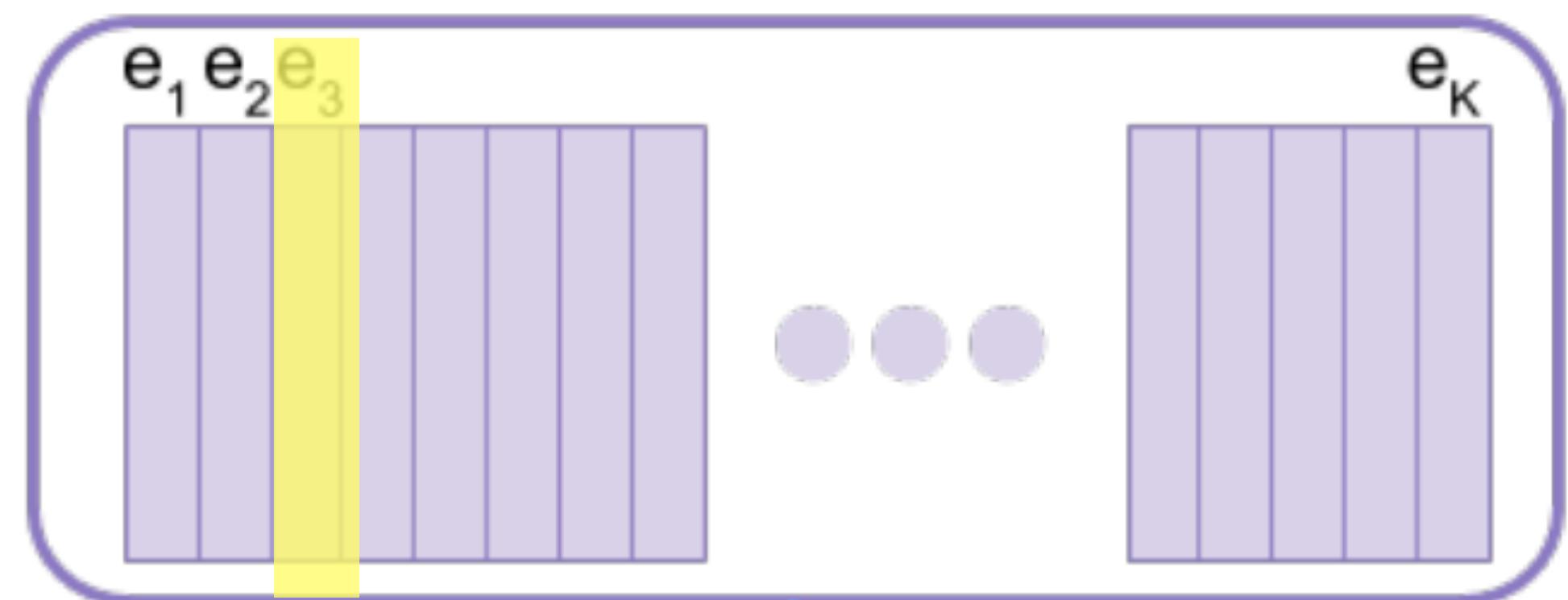
Lower dimensionality using an autoencoder



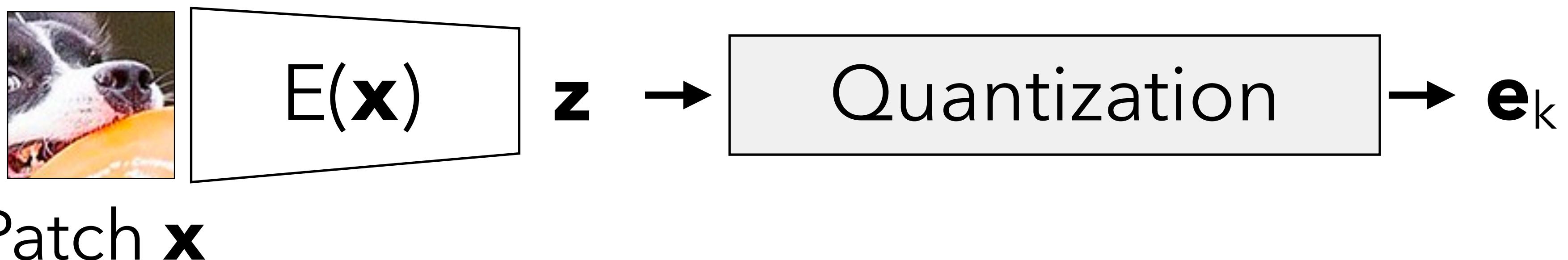
How can we make the output discrete?

Vector quantization

Predict a real-valued vector, then “snap” it to a nearest neighbor from a codebook.



Codebook



$$\text{Quantize}(E(\mathbf{x})) = \mathbf{e}_k \quad \text{where } k = \arg \min_j ||E(\mathbf{x}) - \mathbf{e}_j||$$

Vector quantized (variational) autoencoder

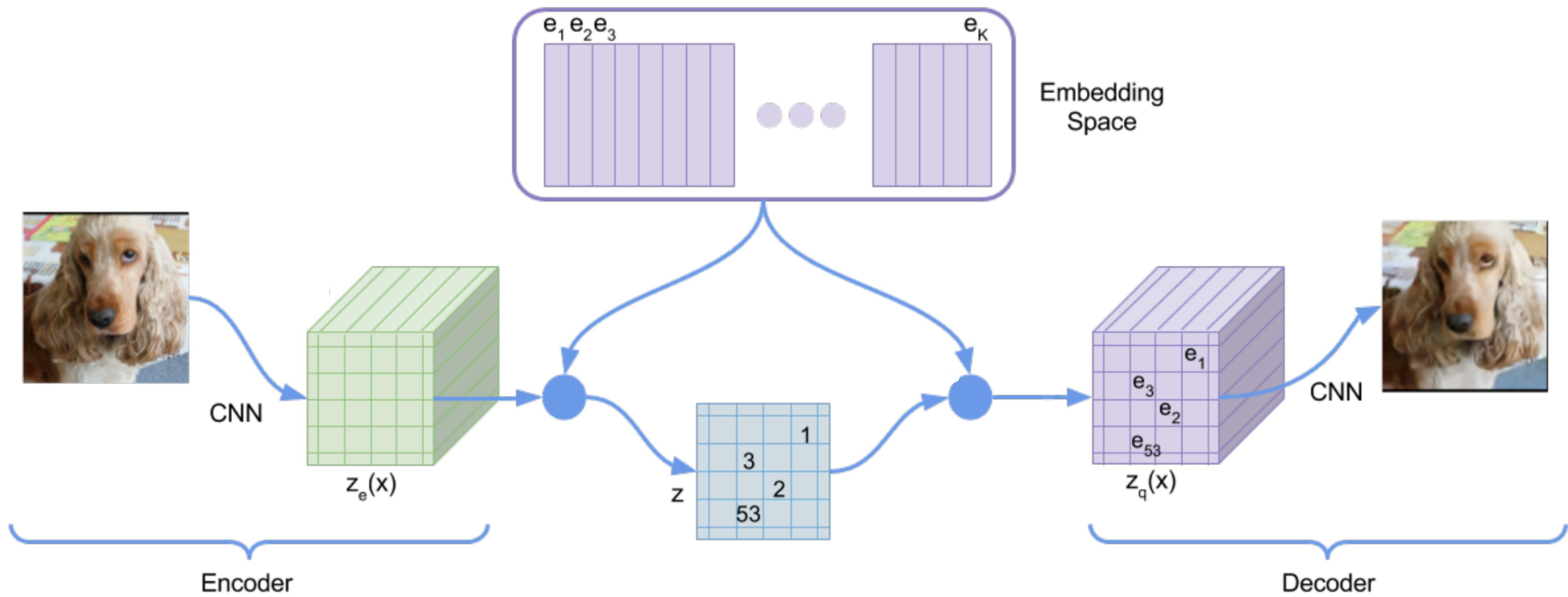
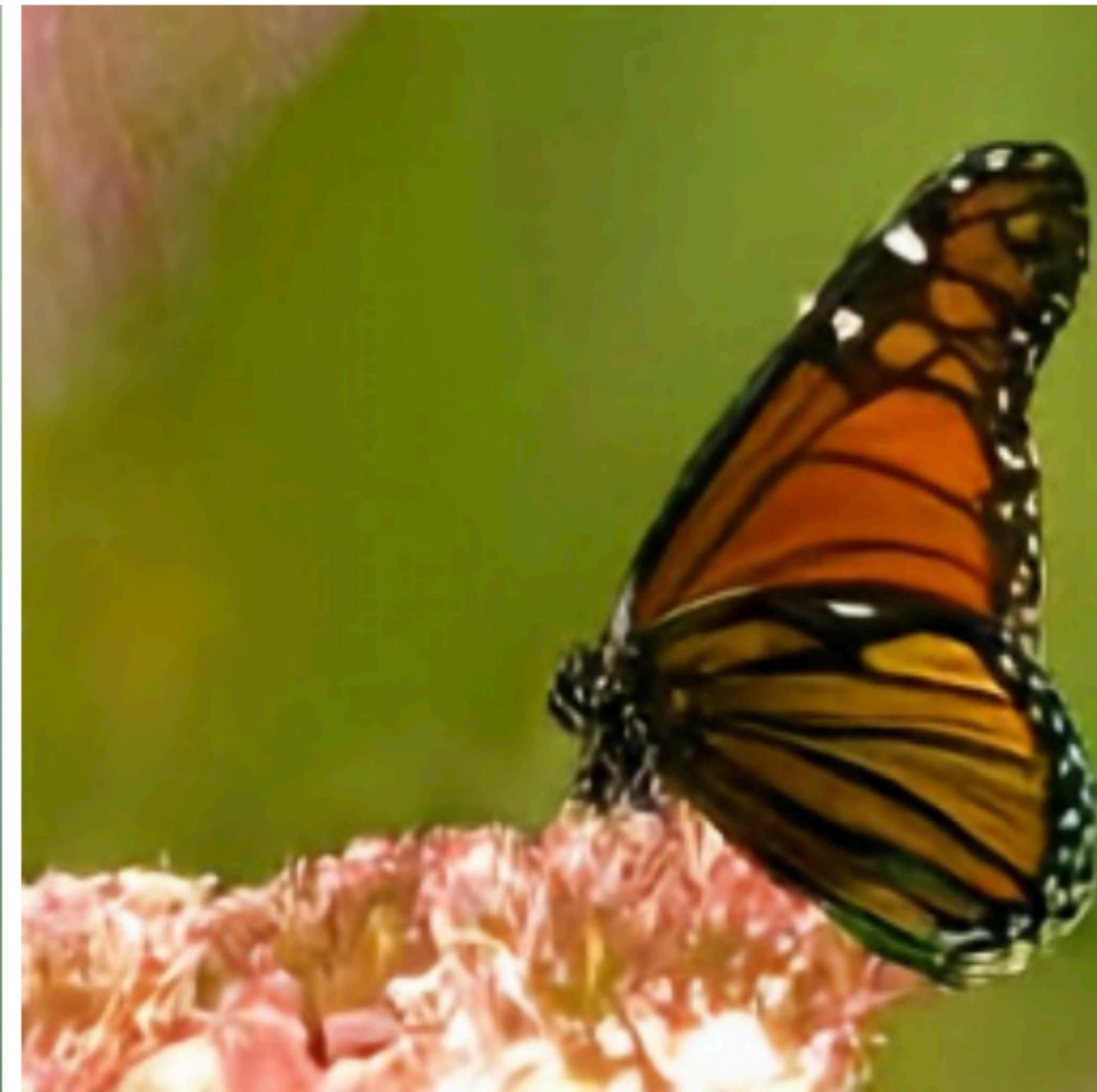


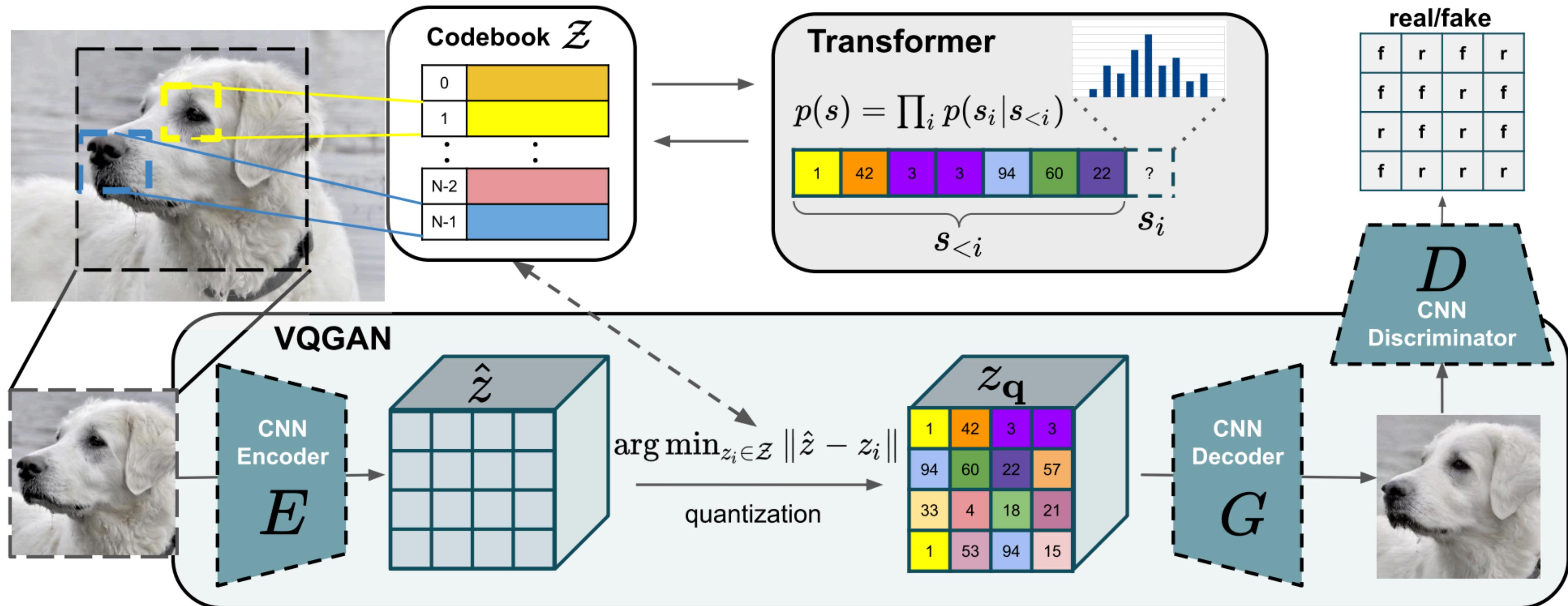
Figure source: [van den Oord et al., "VQ-VAE", 2017]



Generated images (256×256)

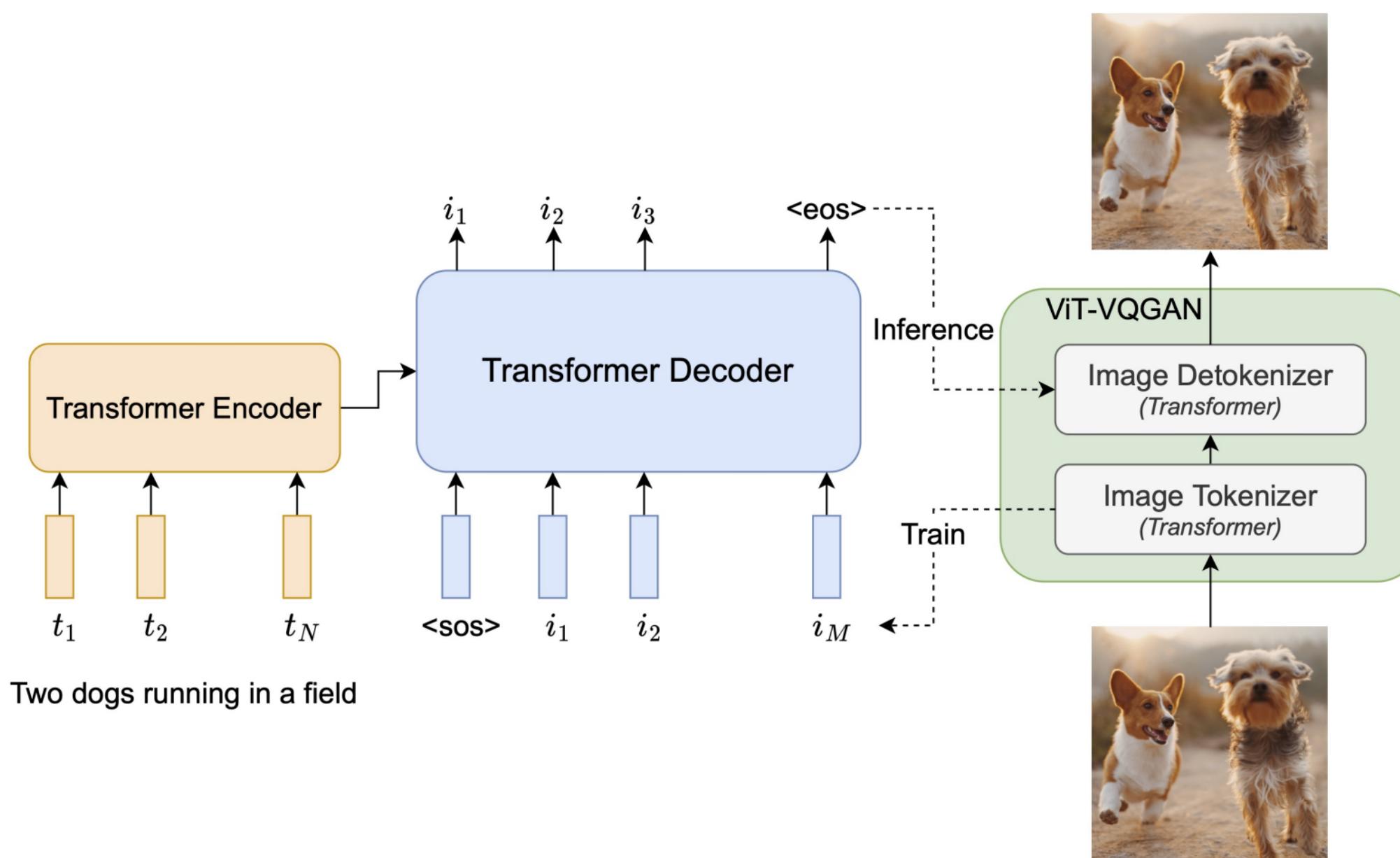
Figure source: [van den Oord et al., "VQ-VAE 2", 2017]

Vector quantized GAN (VQ-GAN)



Add a GAN loss to get crisper images

Autoregressive text-to-image generation



[Yu et al., "Parti", 2022]

Next time: more image synthesis