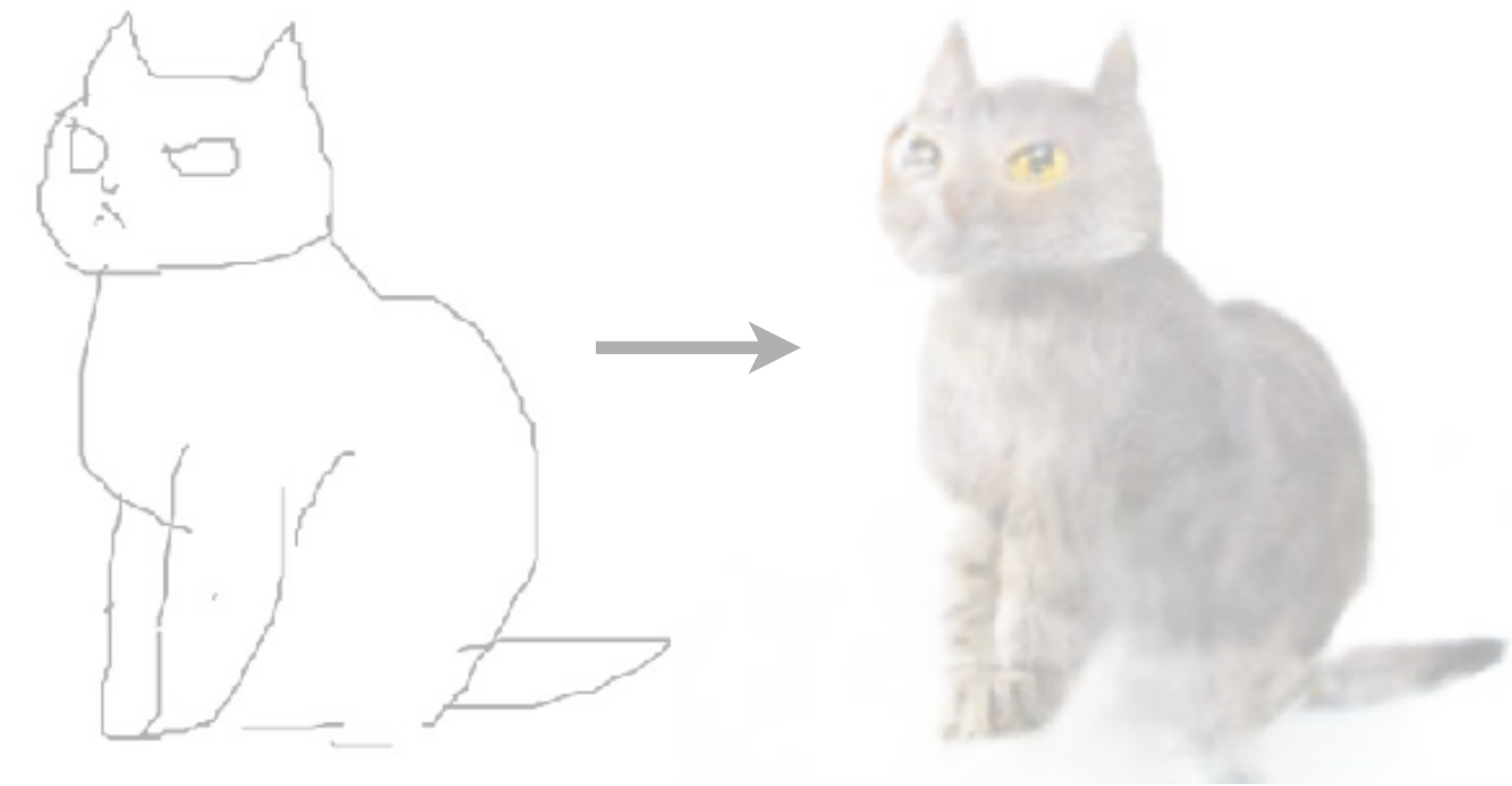
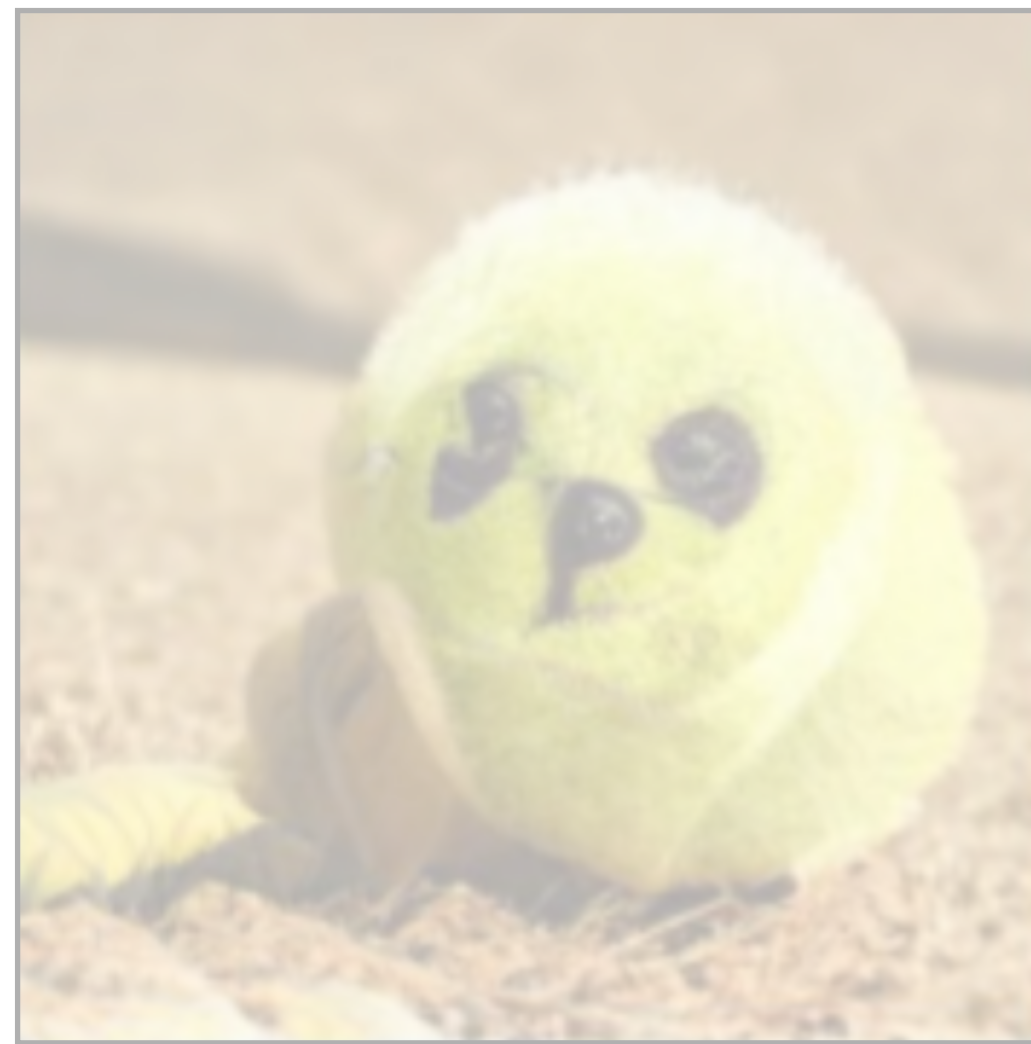


# Lecture 22

## Image Synthesis



# 22. Image Synthesis

- Image synthesis
  - Variational Autoencoders
  - Generative Adversarial Networks
- Structured prediction
  - Image-to-image GANs
- Domain mapping



# Image classification

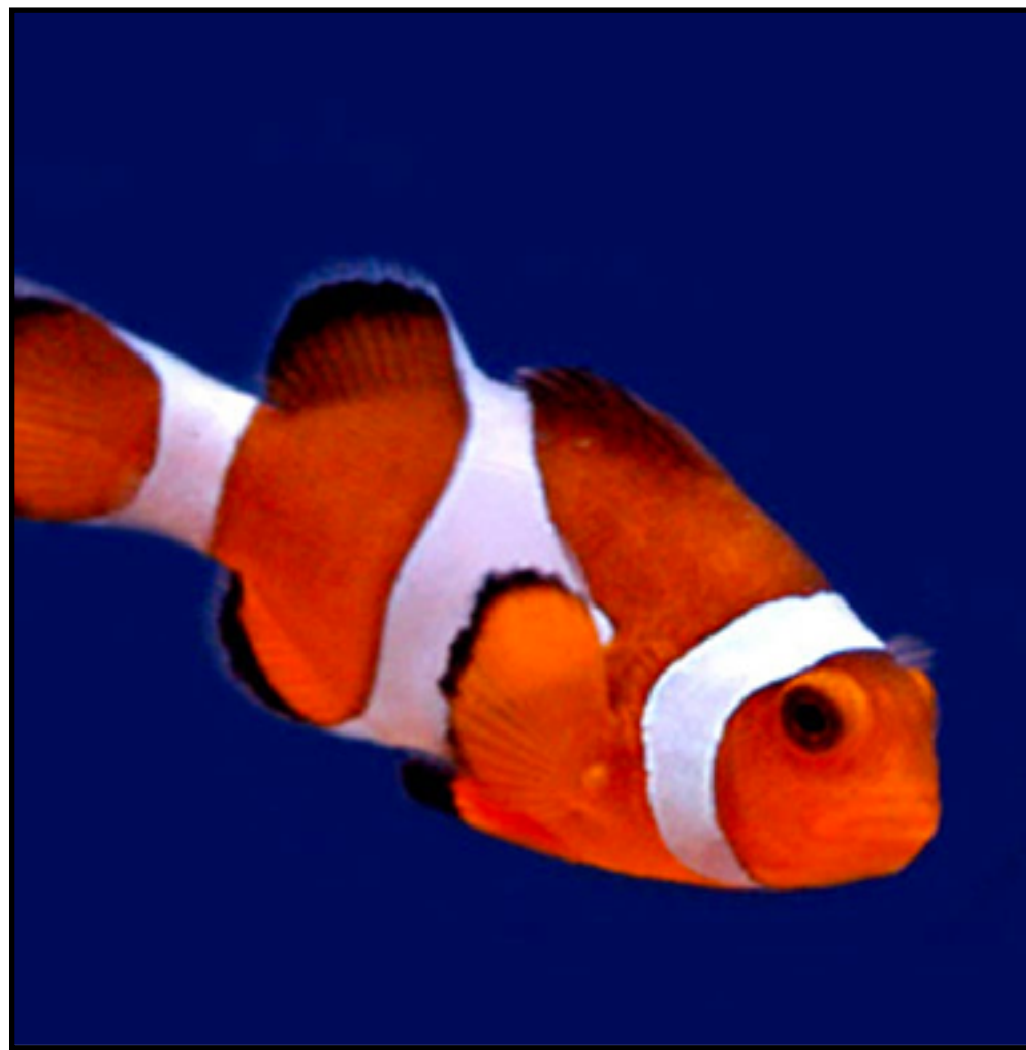


image  $x$



"Fish"

label  $y$

# Image classification



image **x**



"Fish"

label **y**

# Image classification



image **x**

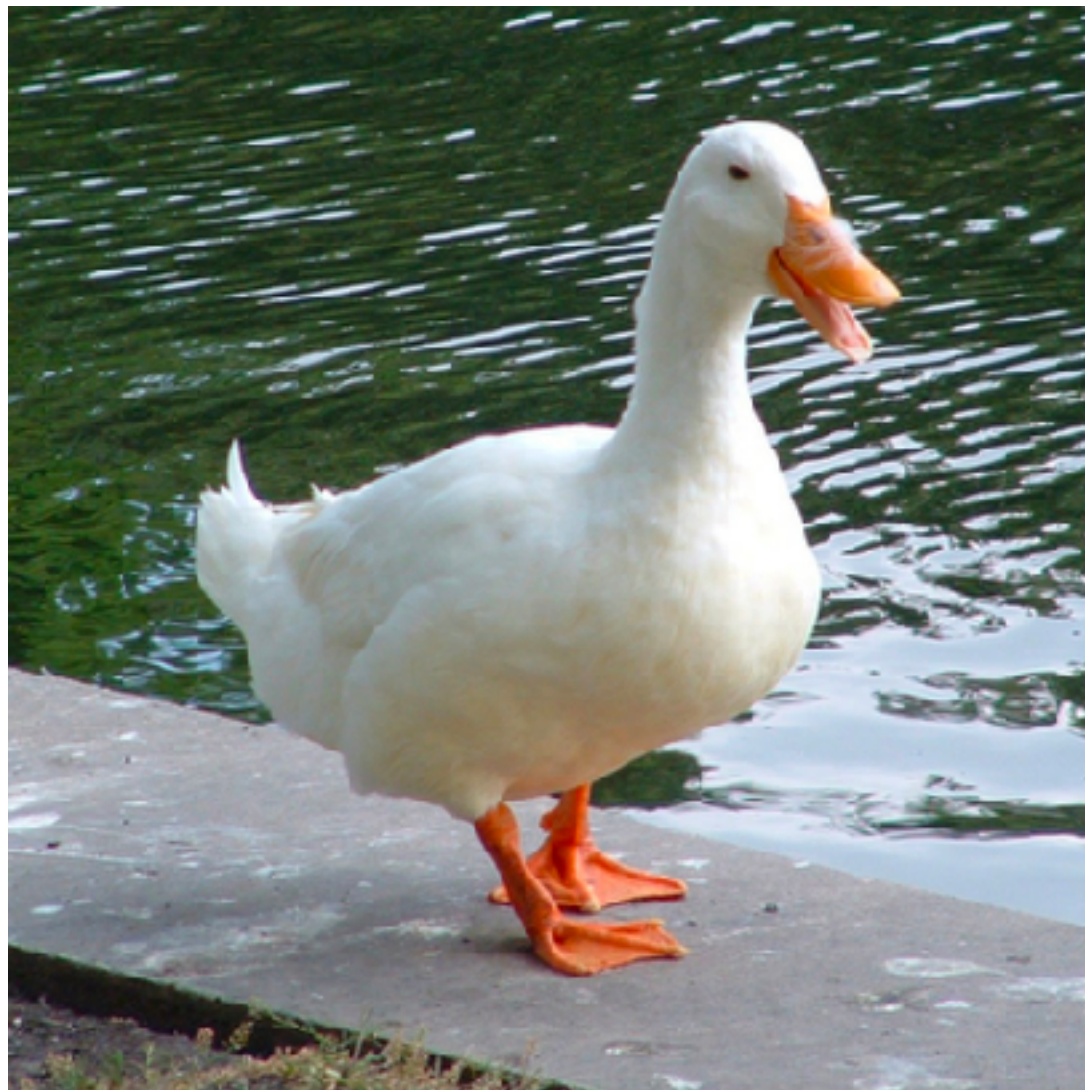


"Fish"

label **y**



# Image classification



⋮  
image  $\mathbf{x}$



“Duck”

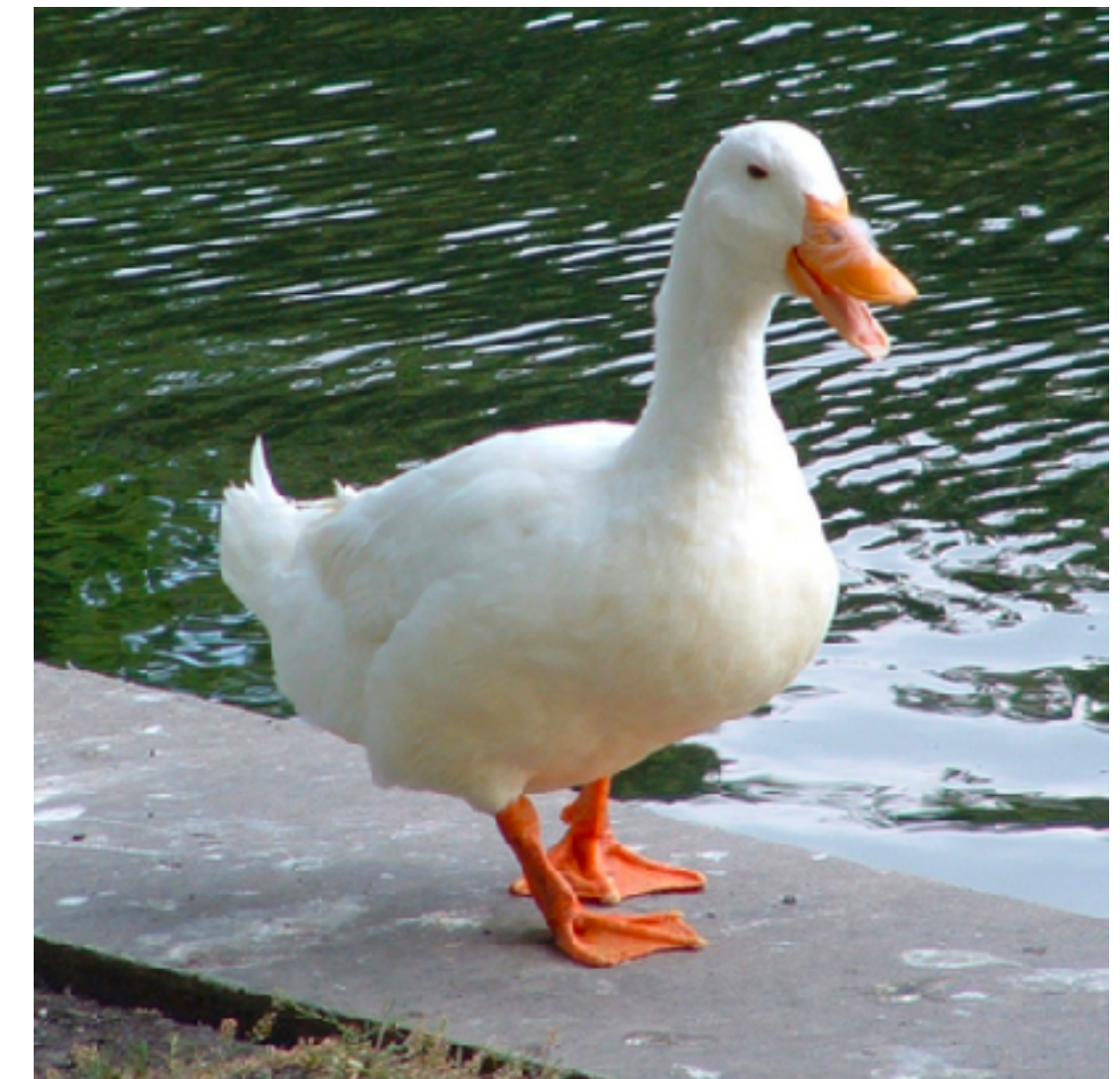
label  $y$

# Image synthesis

“Duck”



Generator



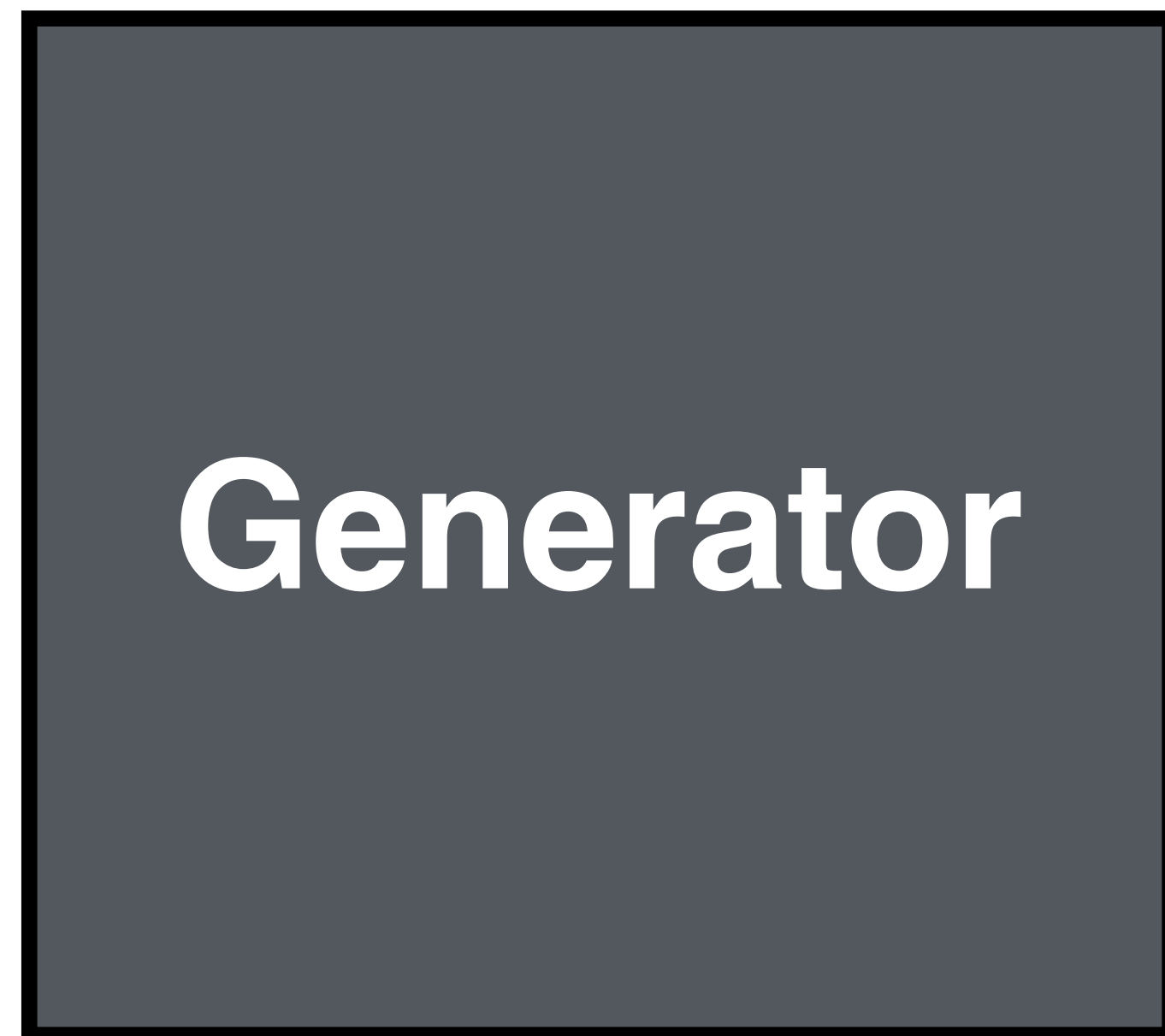
⋮

label  $y$

image  $x$

# Image synthesis

“Fish”

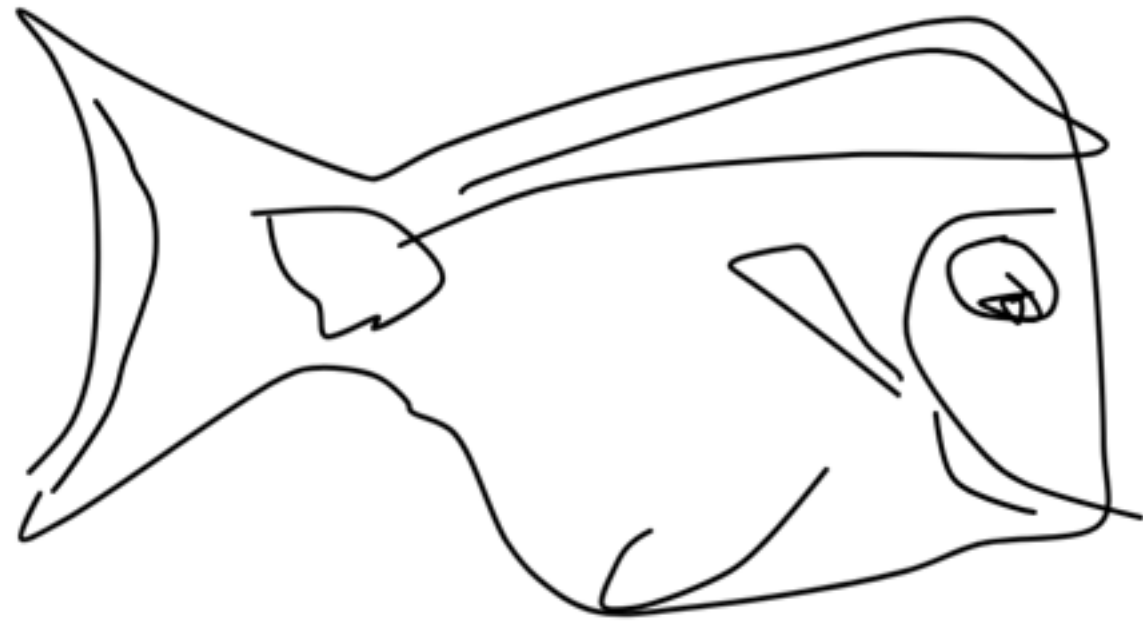


label  $y$

image  $x$



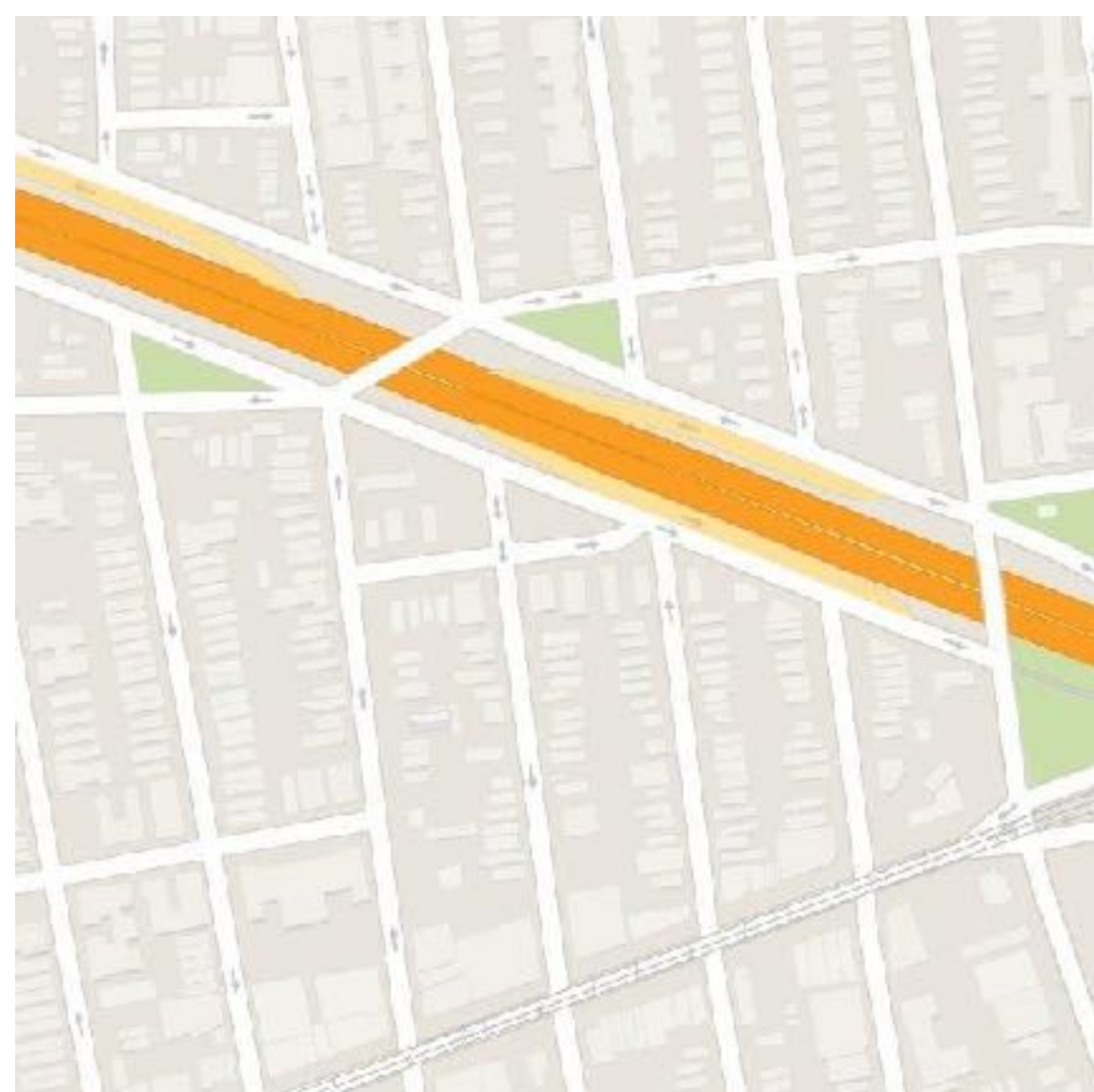
# Image translation



User sketch

Photo

# Image translation



Google Map



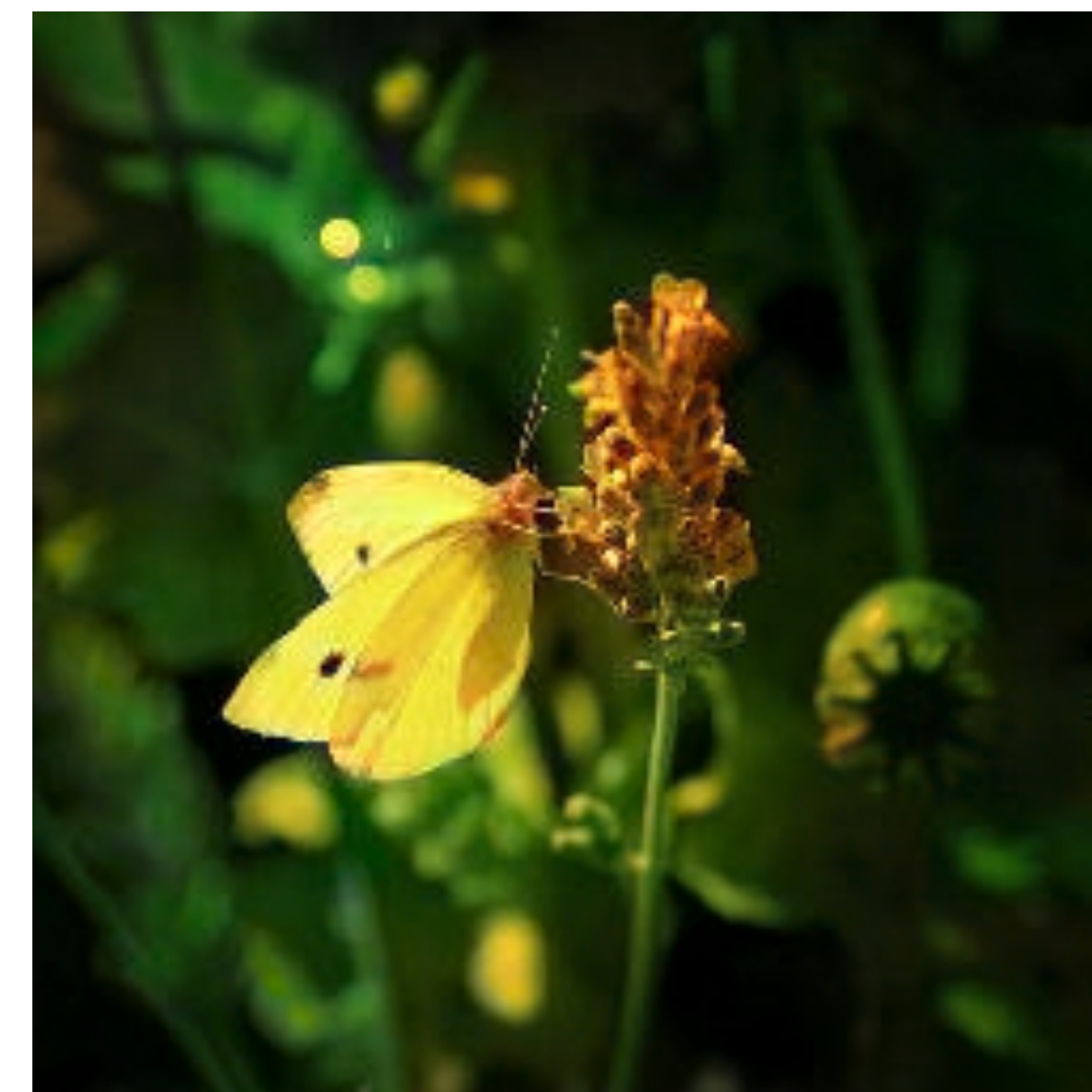
Satellite photo



# Image translation



BW image

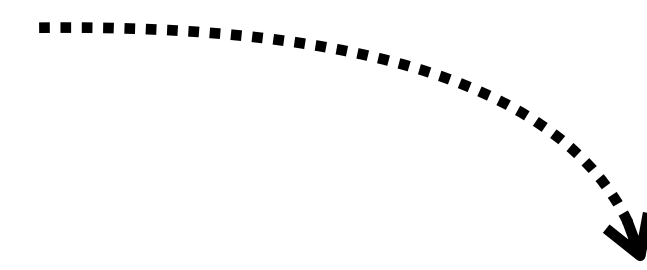


Color image



# Image synthesis via **generative modeling**

**X is high-dimensional!**



Model of high-dimensional structured data  $P(\mathbf{X}|\mathbf{Y} = \mathbf{y})$

In vision, this is usually what we are interested in!

# What can you do with generative models?

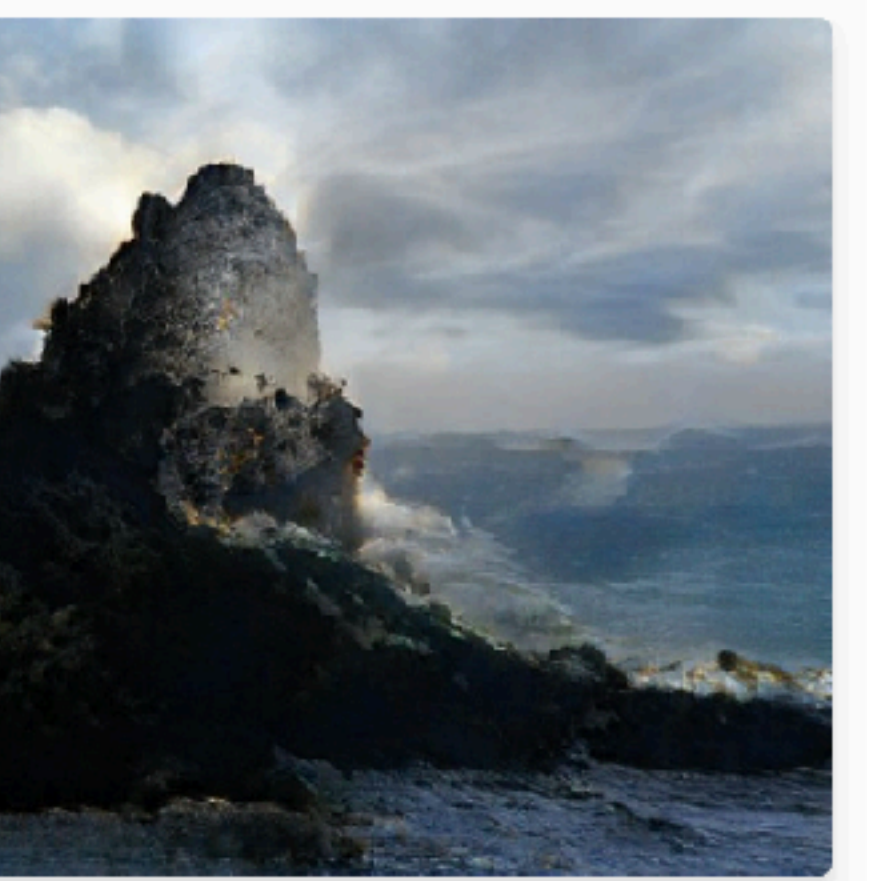
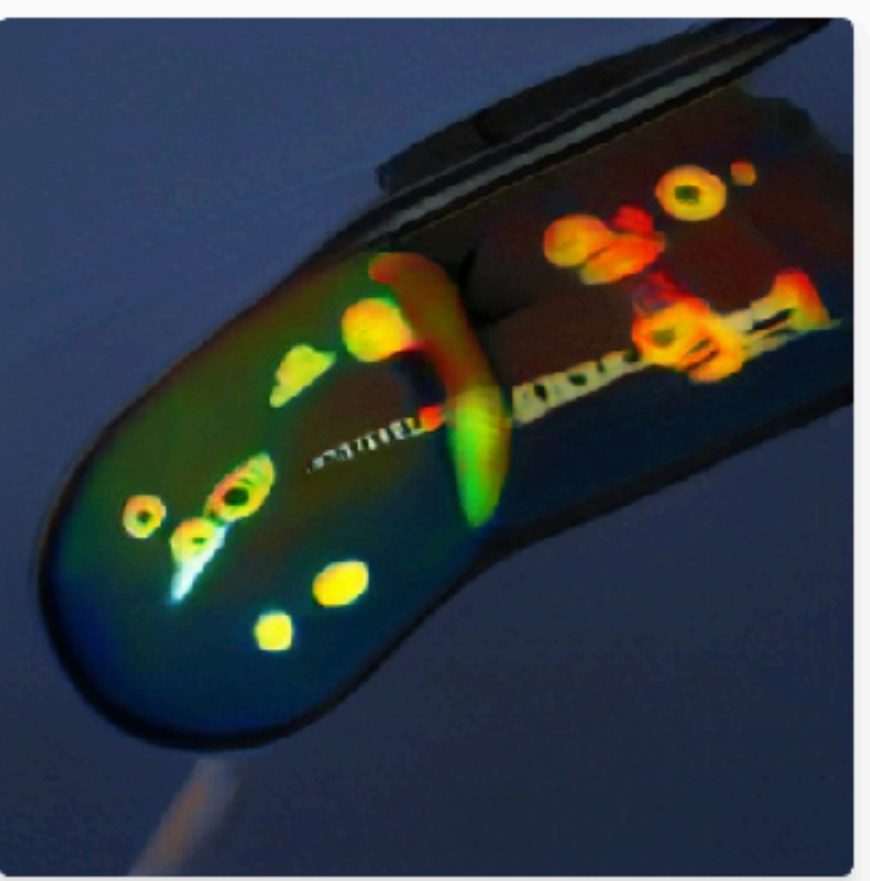
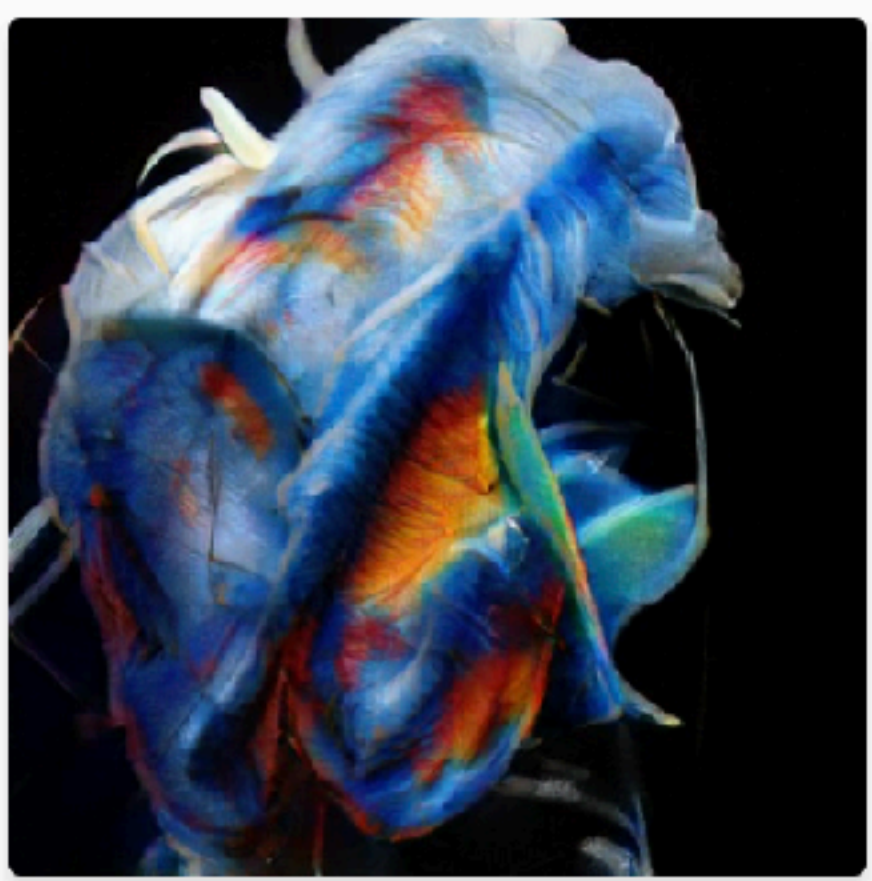
1. Image synthesis
2. Structured prediction
3. Domain mapping
4. (Representation learning)
5. (Model-based intelligence)



**1. Image synthesis**

2. Structured prediction

3. Domain mapping



[Images: <https://ganbreeder.app/>]

# Image synthesis



# Procedural graphics



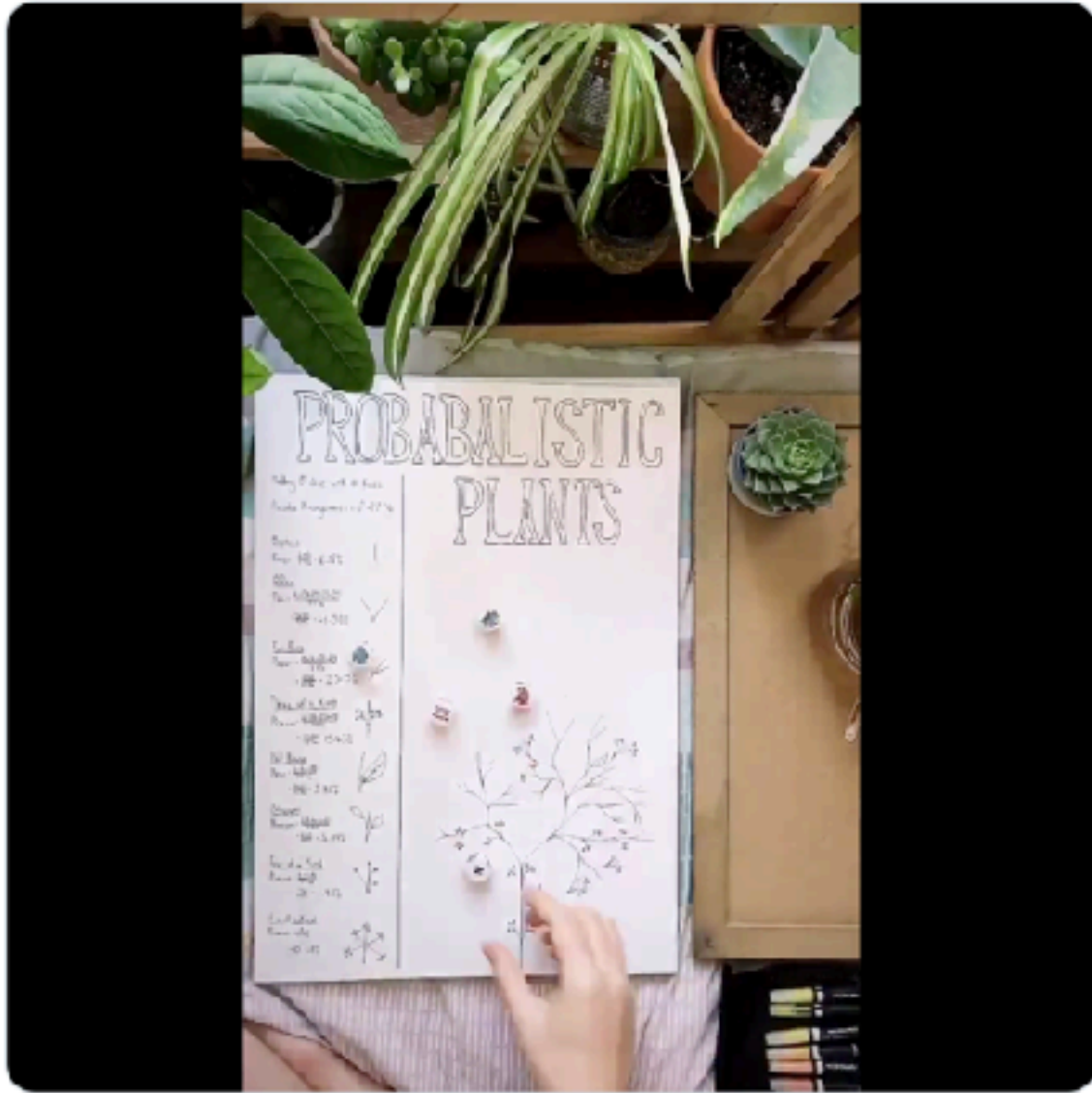
[Anders Scheil]





Aylean @Aylean · Nov 17

Made up a set of rules and rolled some dice to decide how this plant would grow. I never did get that five of a kind, as expected, but I was still hopeful! 🍀🍀🍀



52

1.1K

4.5K



# PROBABALISTIC PLANTS

## PLANTS

Rolling 5 dice with 6 faces

Possible Arrangements =  $6^5 = 7776$

Supernova

Points:  $\frac{1}{7776} = 0.01286\%$



A Pair

Points:  $\frac{6 \cdot (6-1) \cdot (6-1) \cdot (6-1)}{7776} = 0.4630\%$

$\frac{3600}{7776} = 46.30\%$



Two Pairs

Points:  $\frac{(6-1) \cdot (6-1) \cdot (6-1)}{7776} = 0.2315\%$

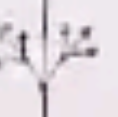
$\frac{1800}{7776} = 23.15\%$



Three of a Kind

Points:  $\frac{6 \cdot (6-1) \cdot (6-1)}{7776} = 0.1543\%$

$\frac{1080}{7776} = 15.43\%$



Full House

Points:  $\frac{6 \cdot (6-1)}{7776} = 0.3842\%$

$\frac{144}{7776} = 3.842\%$



Straight

Points:  $\frac{6 \cdot (6-1) \cdot (6-1) \cdot (6-1)}{7776} = 0.309\%$

$\frac{72}{7776} = 3.09\%$



Four of a Kind

Points:  $\frac{6 \cdot (6-1)}{7776} = 0.193\%$

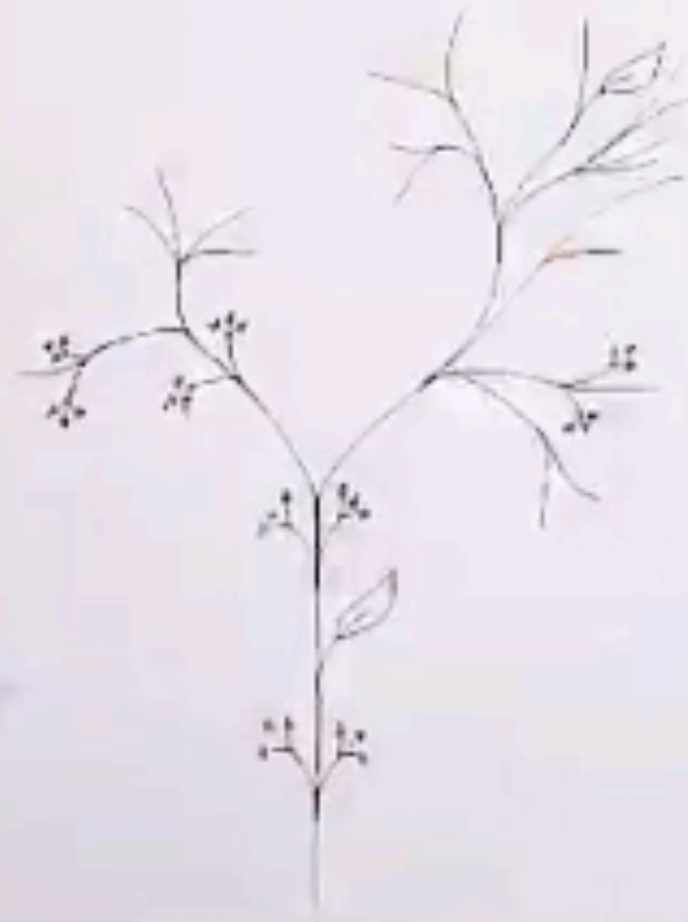
$\frac{36}{7776} = 1.93\%$



Five of a Kind

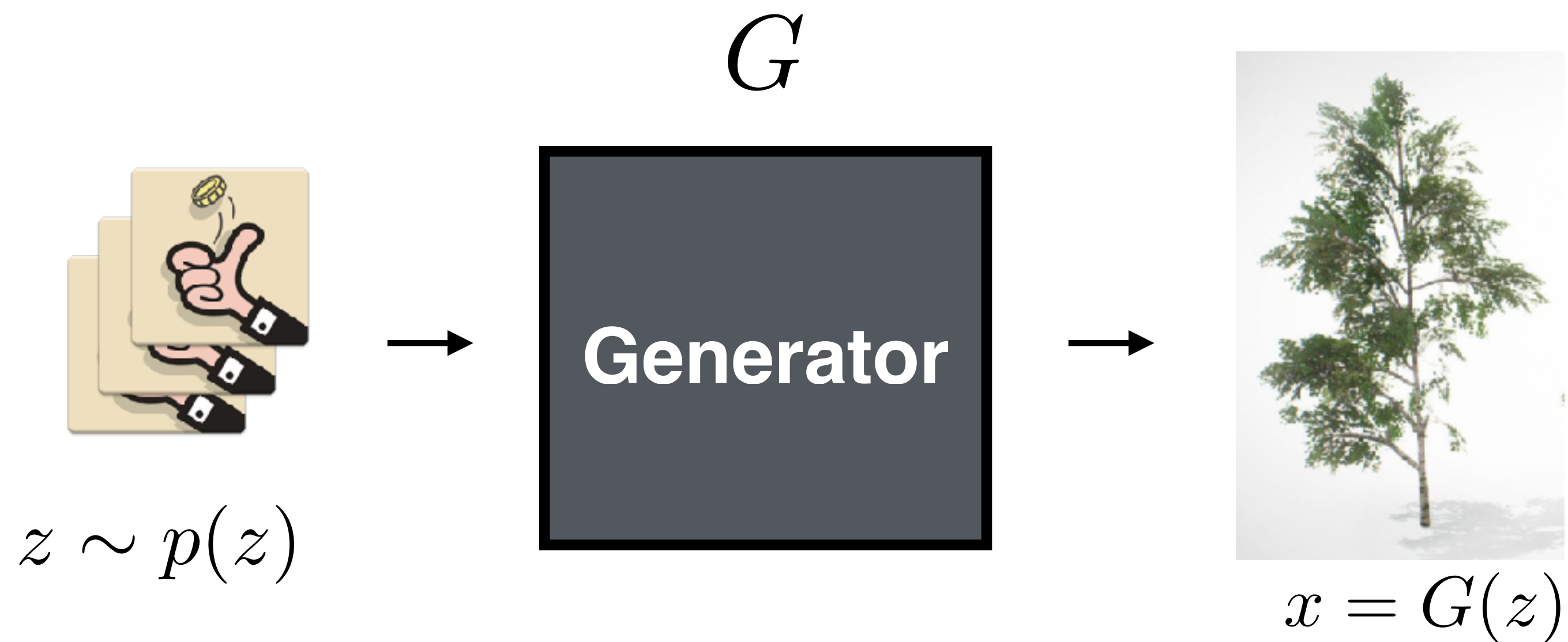
Points:  $\frac{6}{7776} = 0.08\%$

$\frac{6}{7776} = 0.08\%$





# Image synthesis from “noise”



Sampler

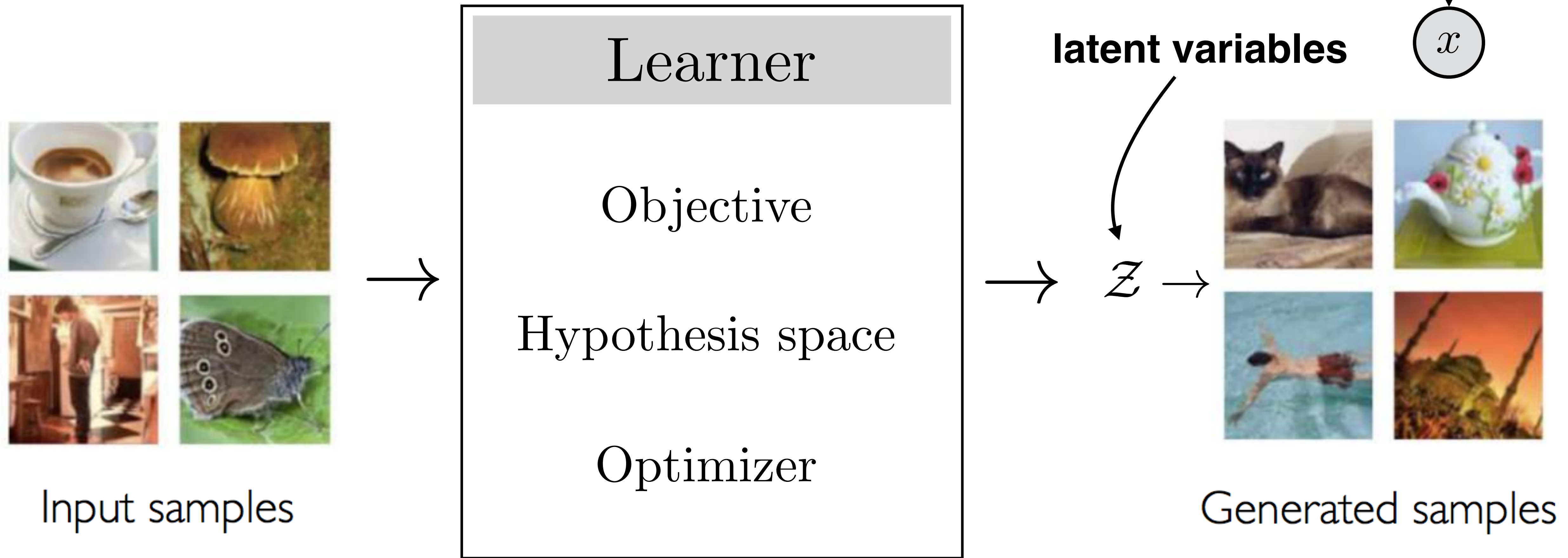
$$G : \mathcal{Z} \rightarrow \mathcal{X}$$

$$z \sim p(z)$$

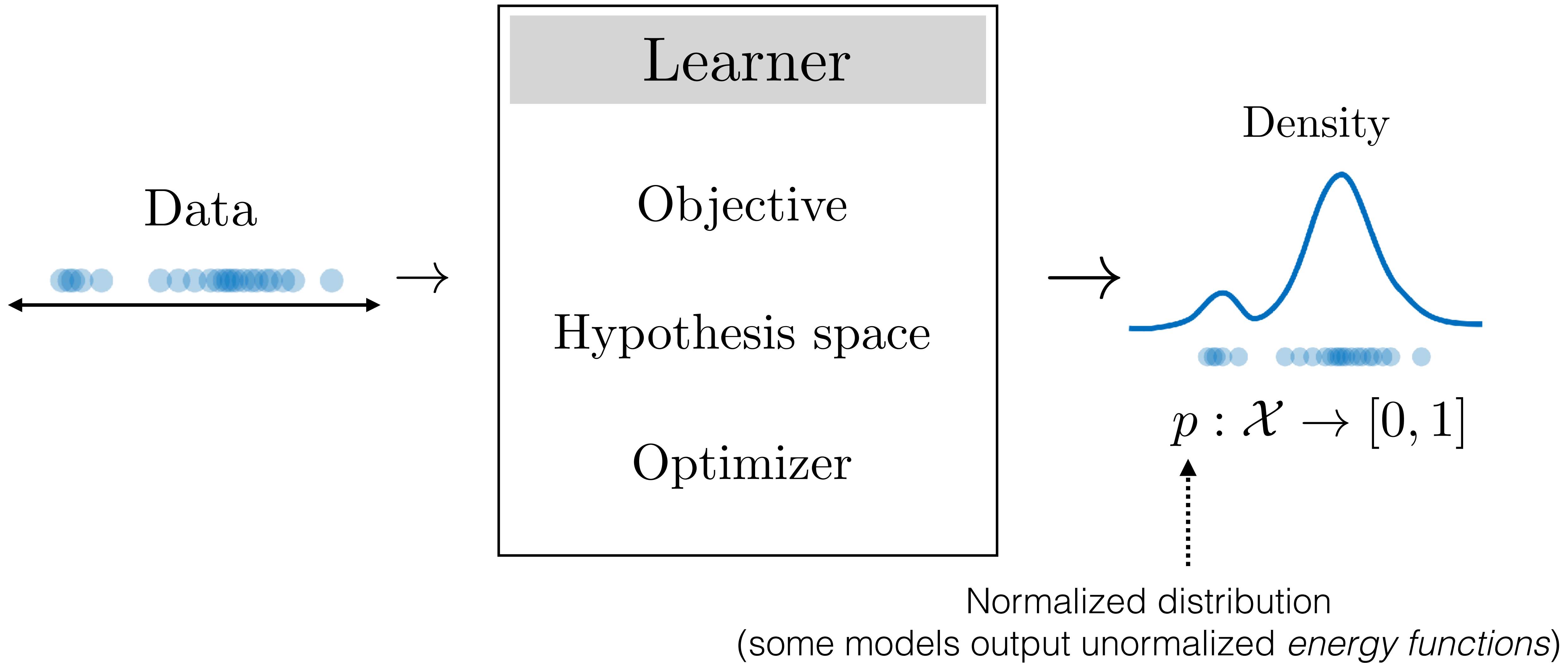
$$x = G(z)$$



# Learning a generative model



# Learning a density model



# Case study #1: Fitting a Gaussian to data

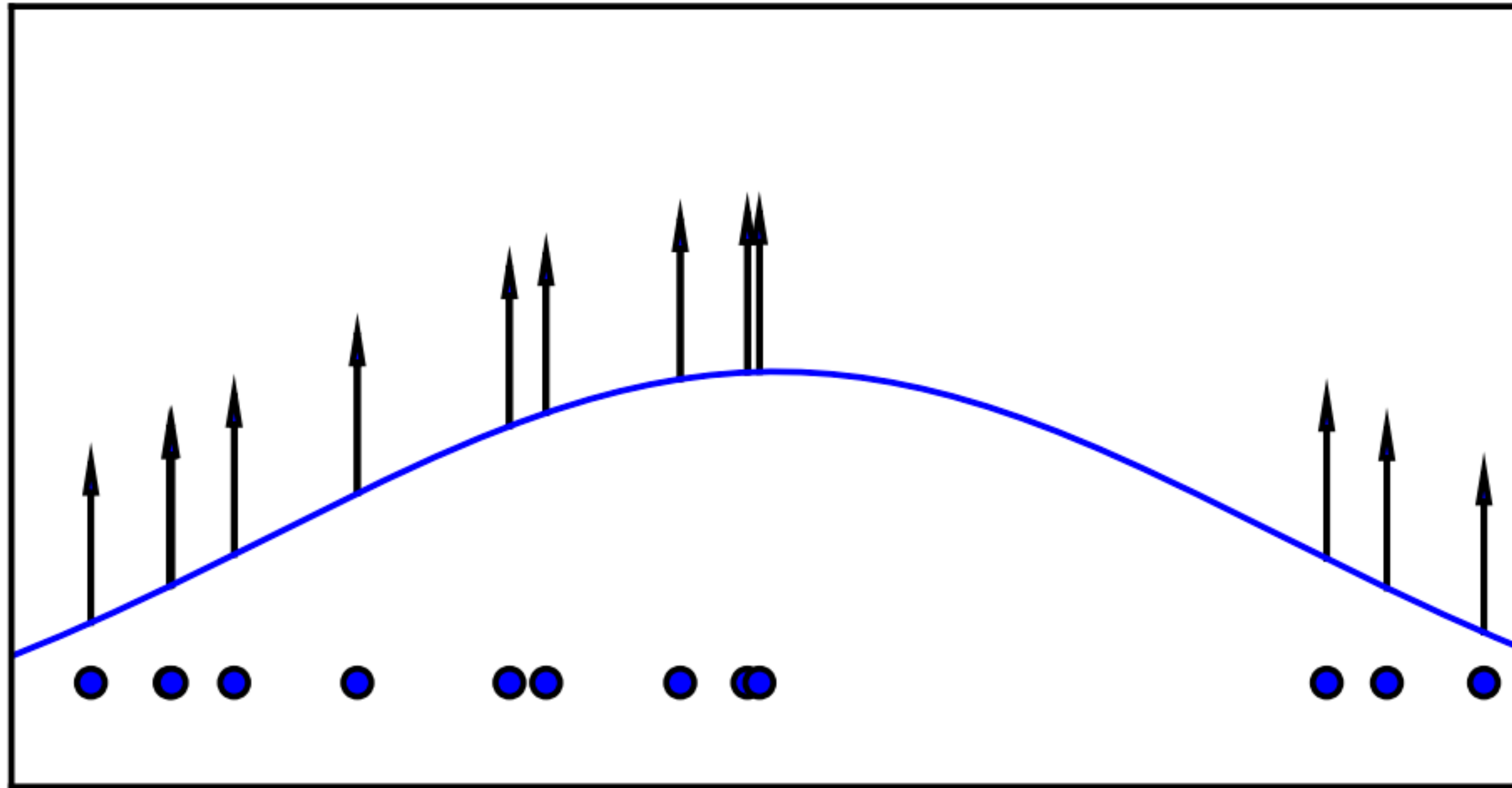


fig from [Goodfellow, 2016]

Max likelihood objective

$$\max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

Considering only Gaussian fits

$$p_{\theta}(x) = \mathcal{N}(x; \mu, \sigma)$$

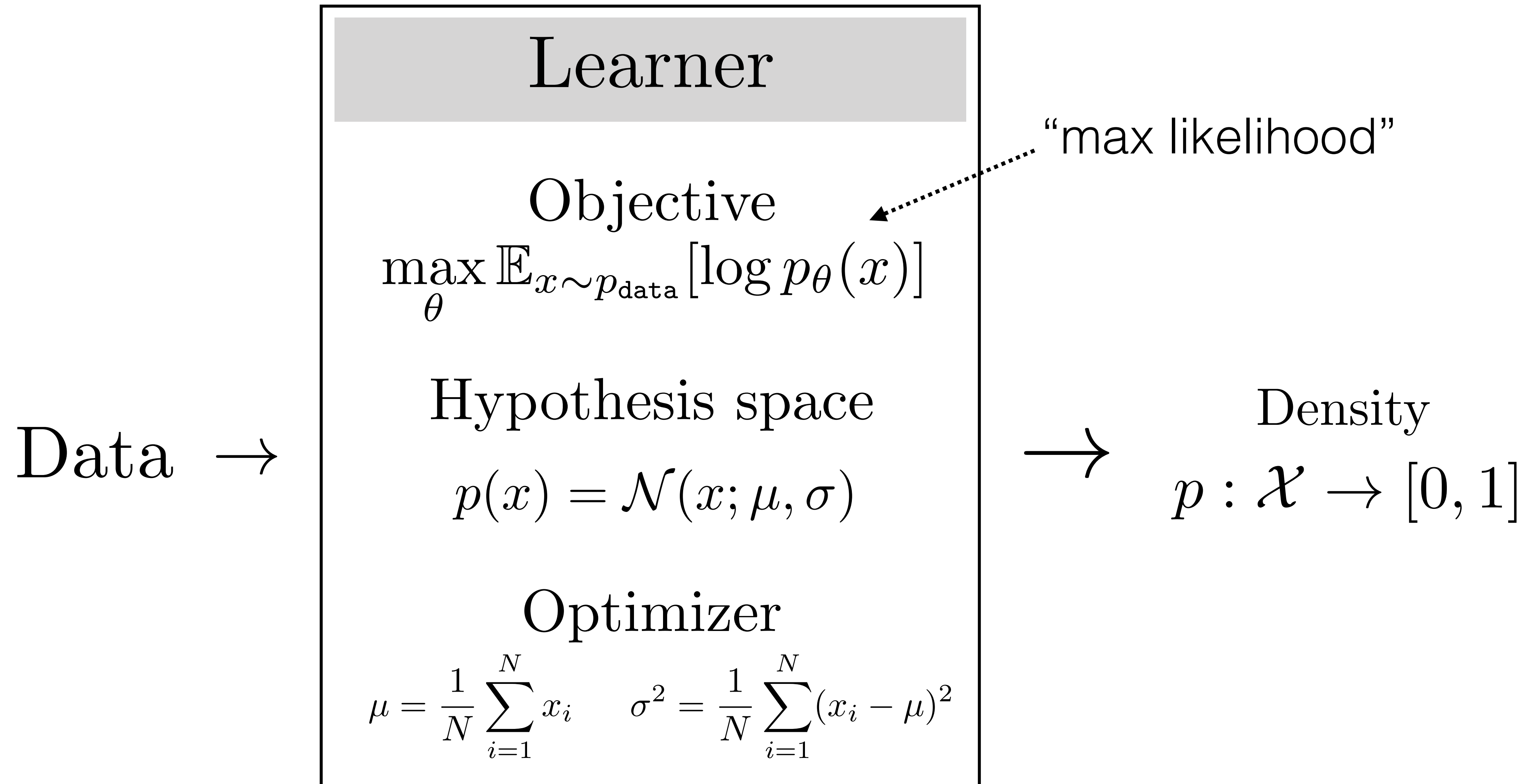
$$\theta = [\mu, \sigma]$$

Closed form optimum:

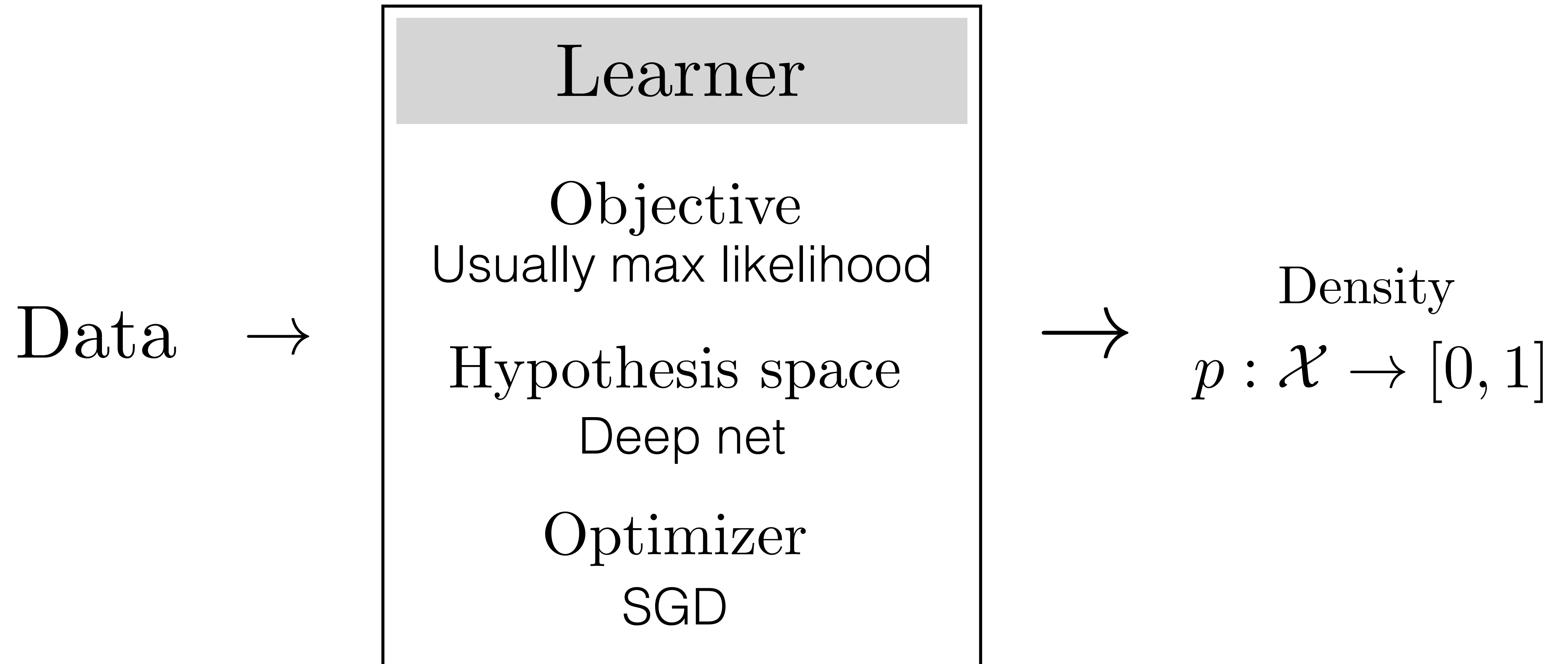
$$\mu = \frac{1}{N} \sum_{i=1}^N x_i \quad \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$$



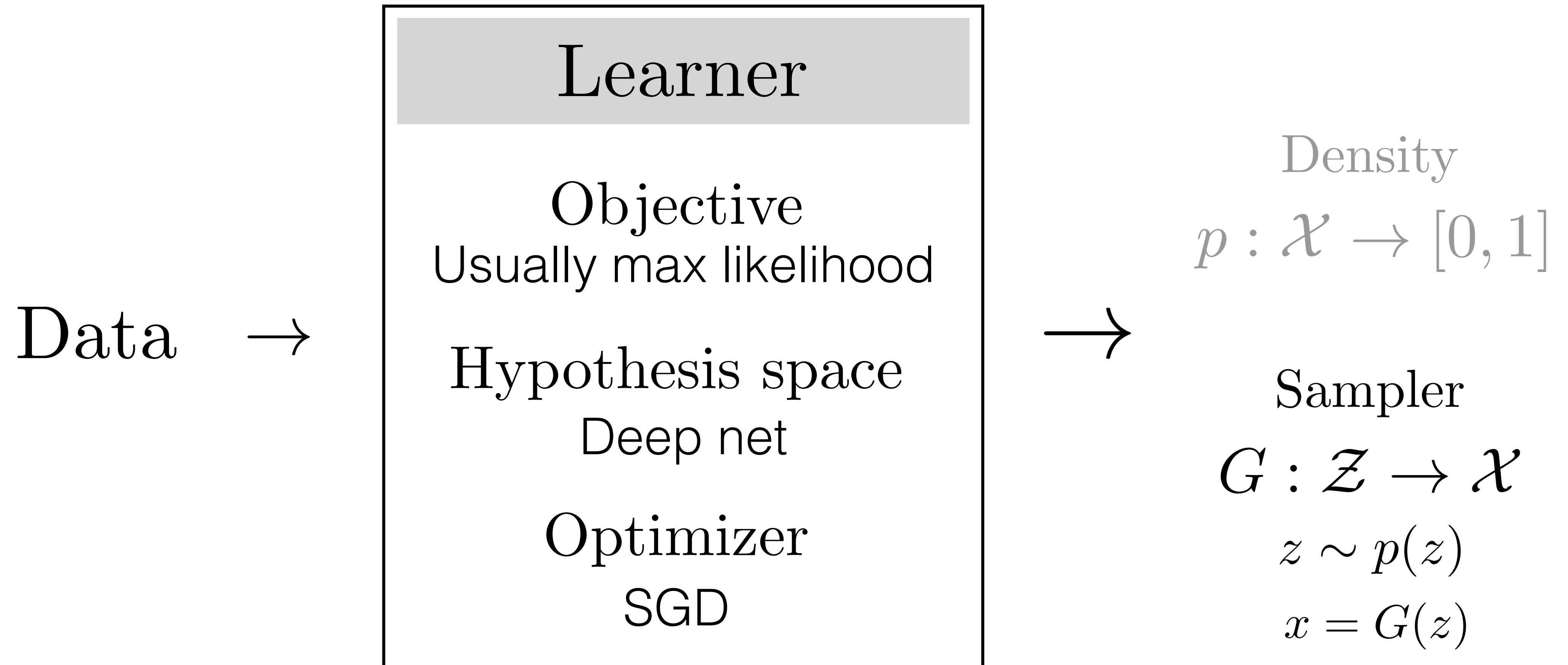
# Case study #1: Fitting a Gaussian to data



# Case study #2: learning a deep generative model



# Case study #2: learning a deep generative model



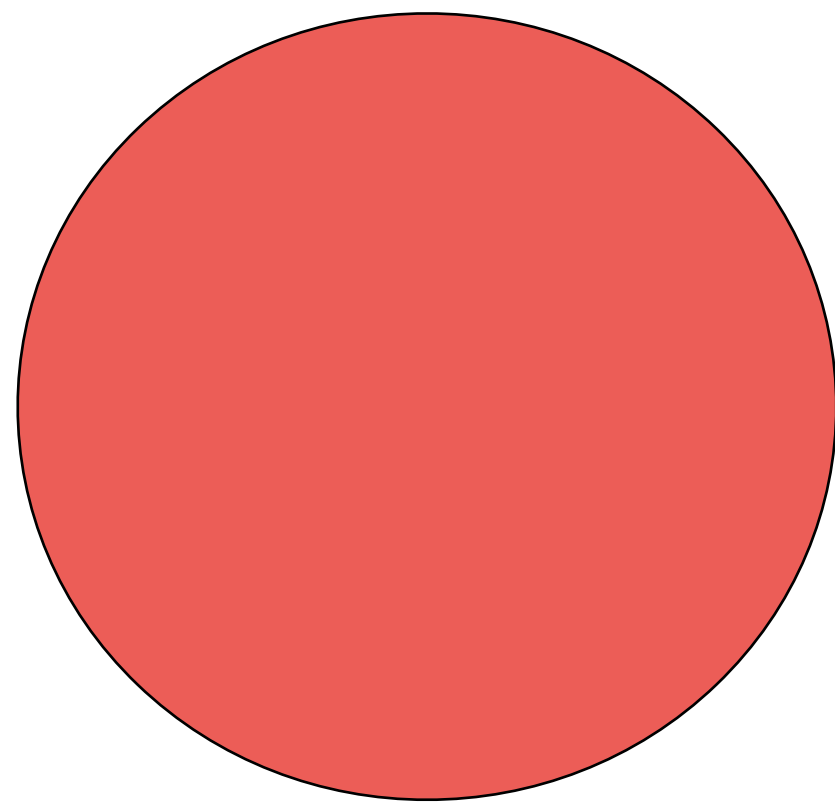
Models that provide a sampler but no density are called **implicit generative models**



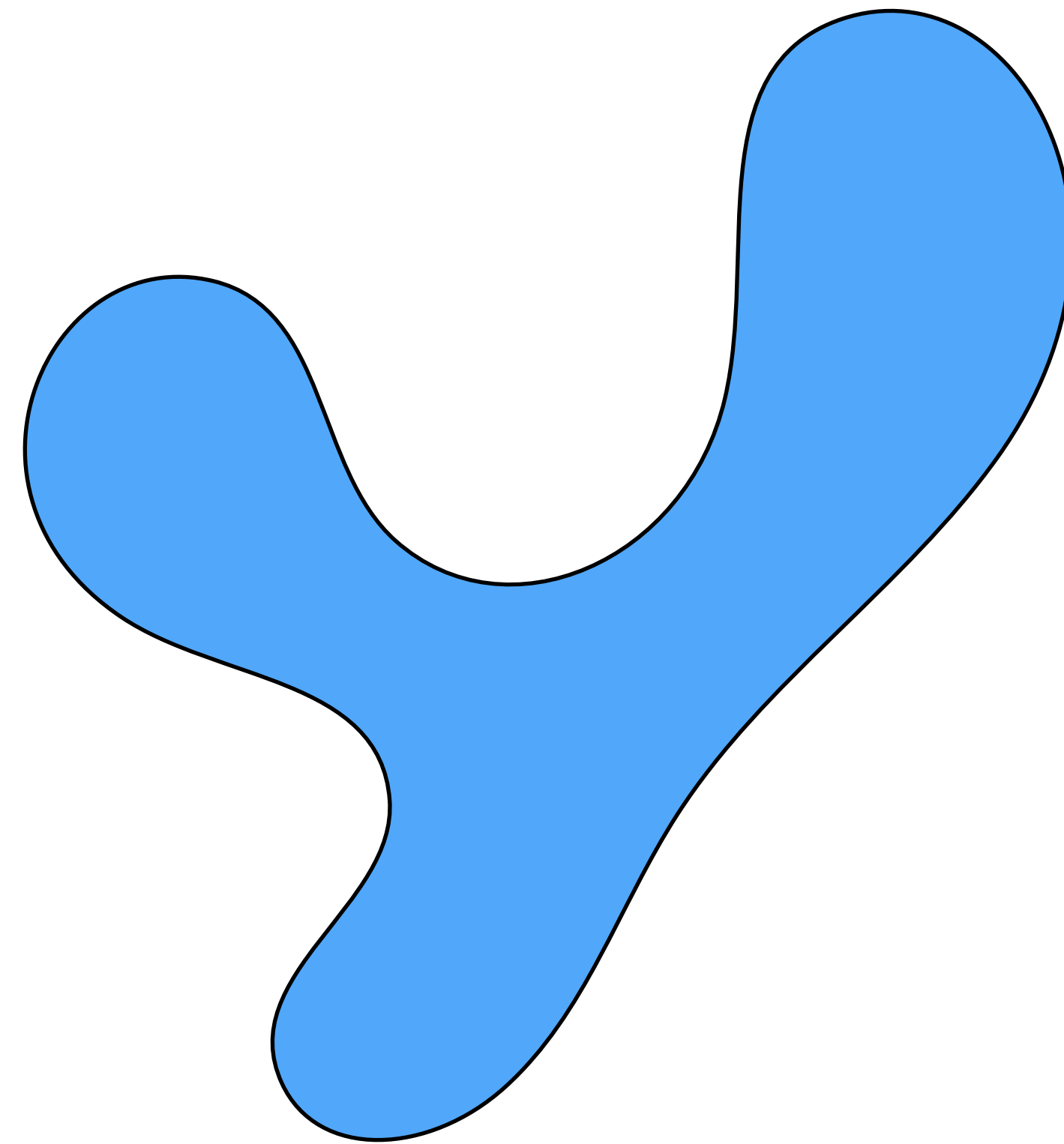
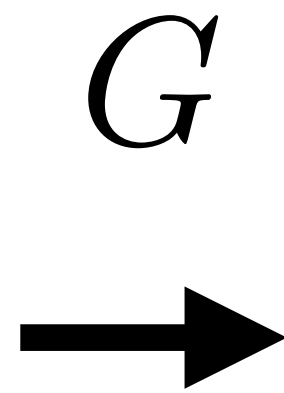
# Deep generative models are distribution transformers

Prior distribution

Target distribution

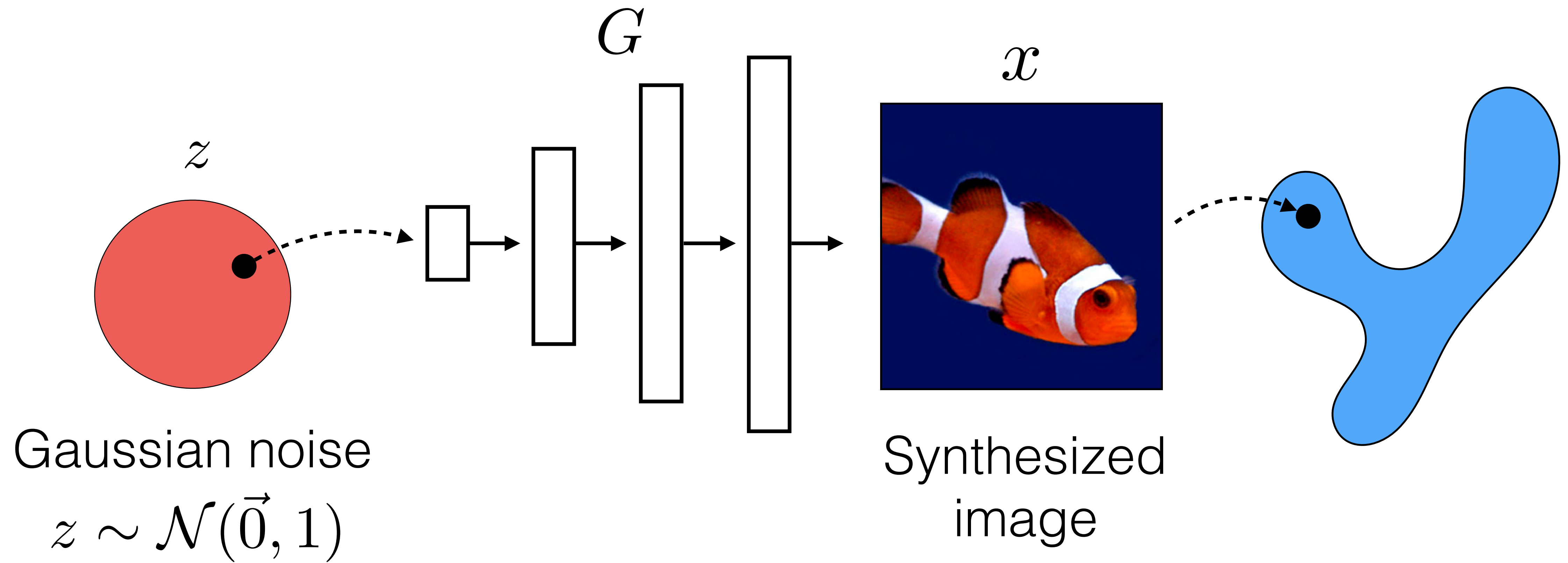


$p(z)$

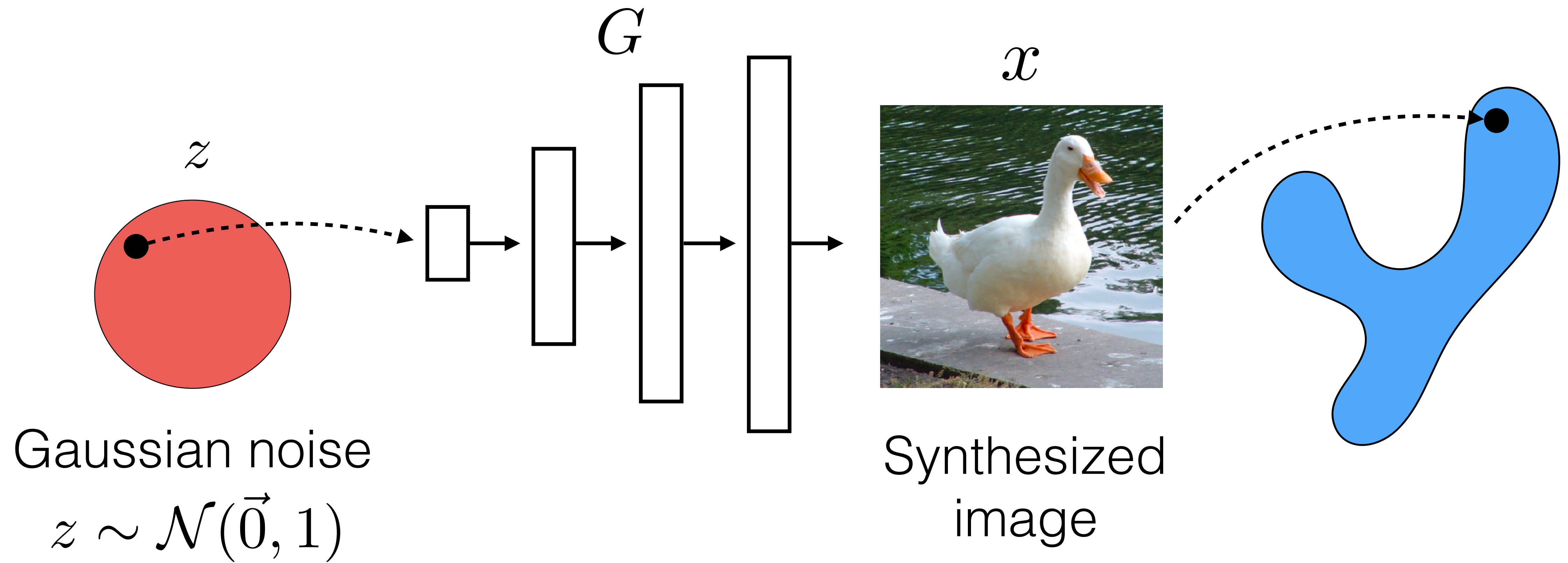


$p(x)$

# Deep generative models are distribution transformers



# Deep generative models are distribution transformers



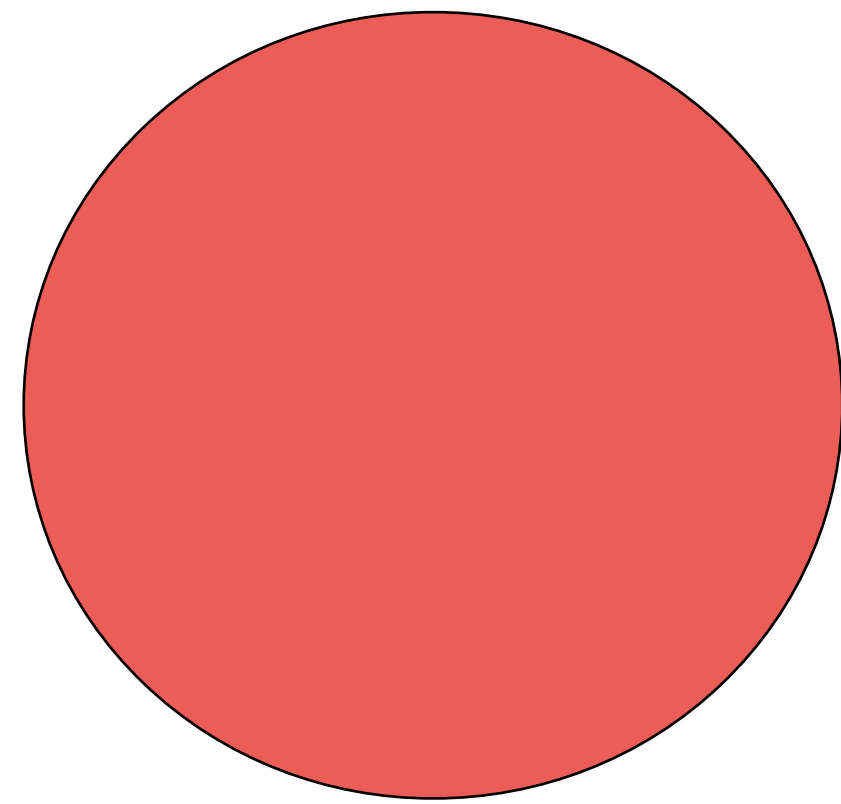


# Variational Autoencoders (VAEs)

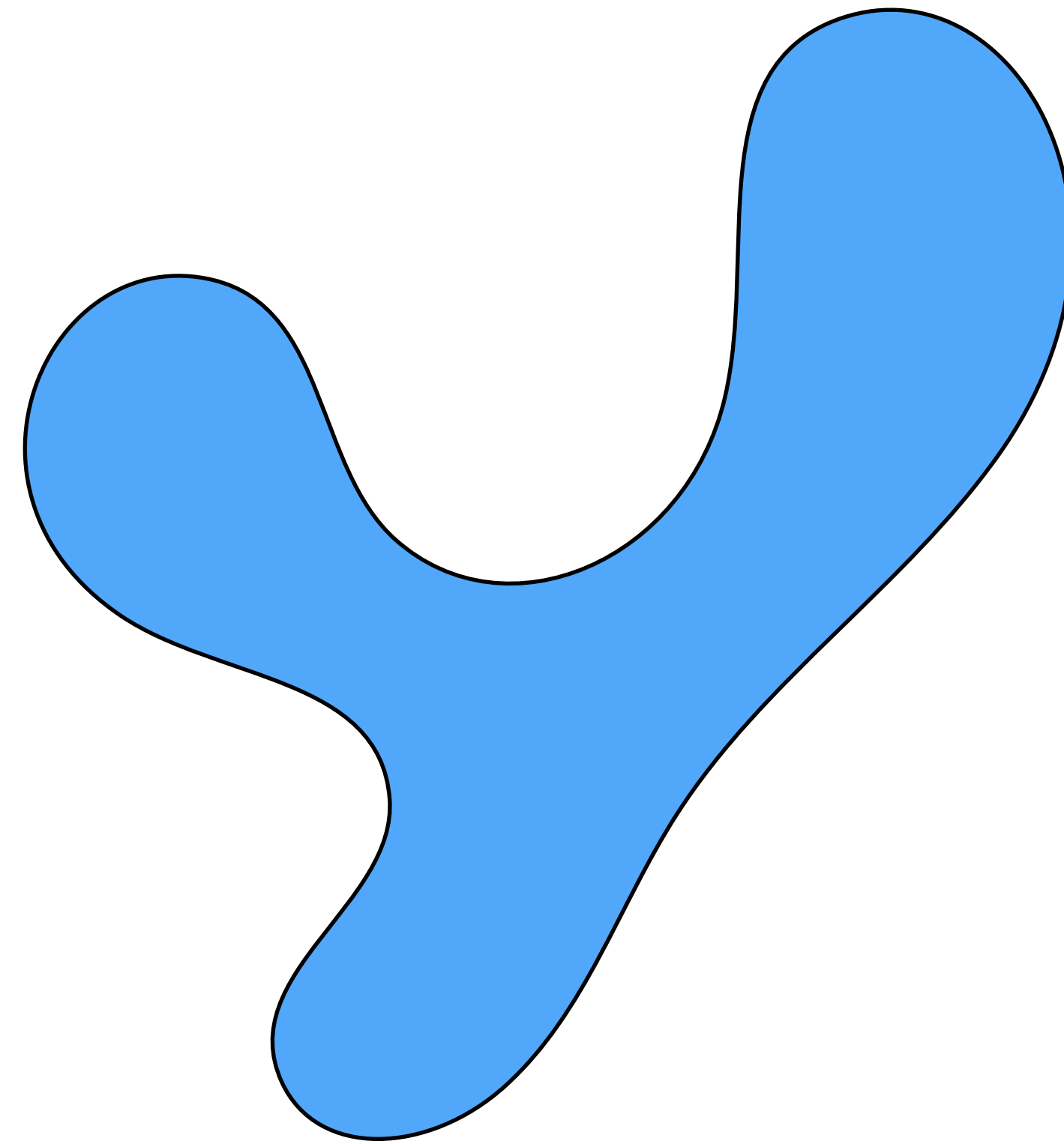
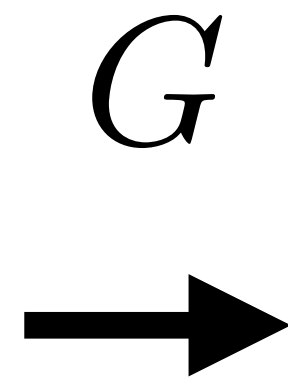
[Kingma & Welling, 2014; Rezende, Mohamed, Wierstra 2014]

Prior distribution

Target distribution

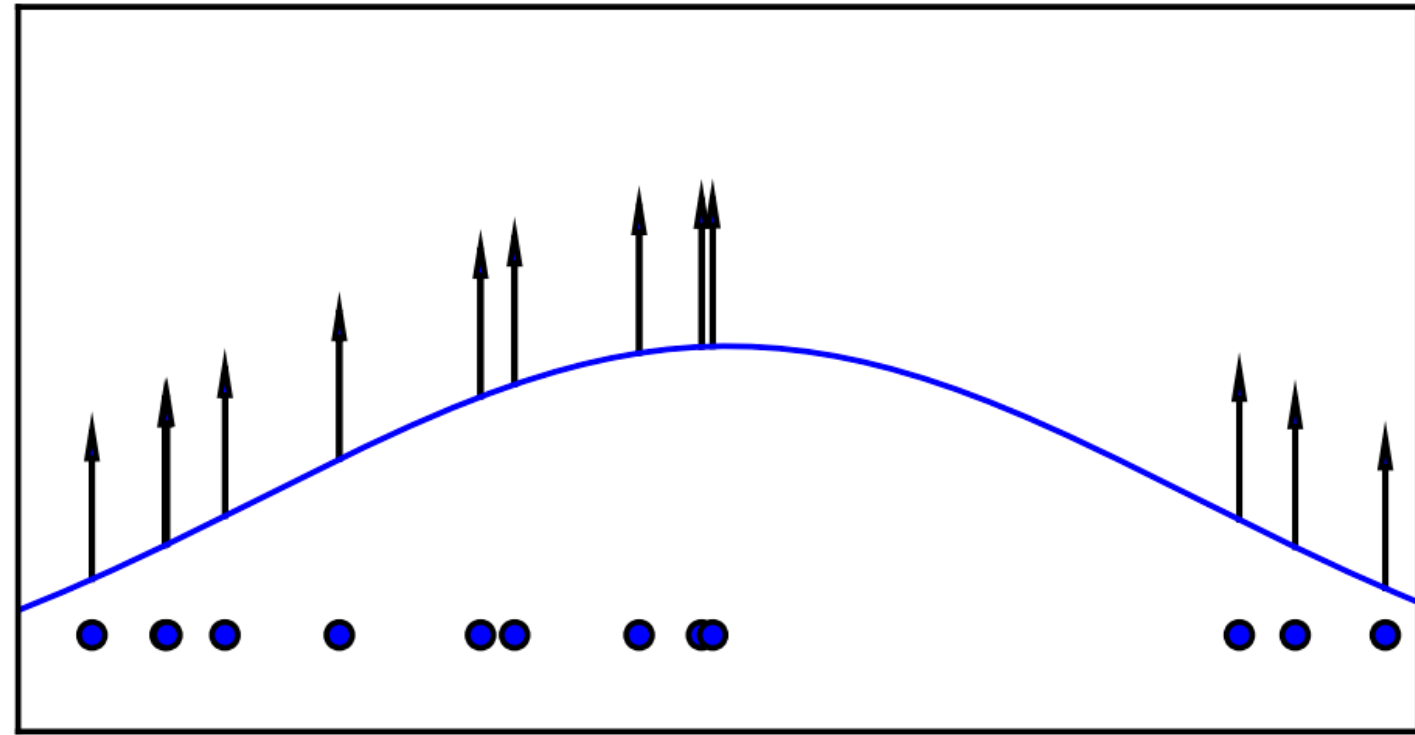


$p(z)$

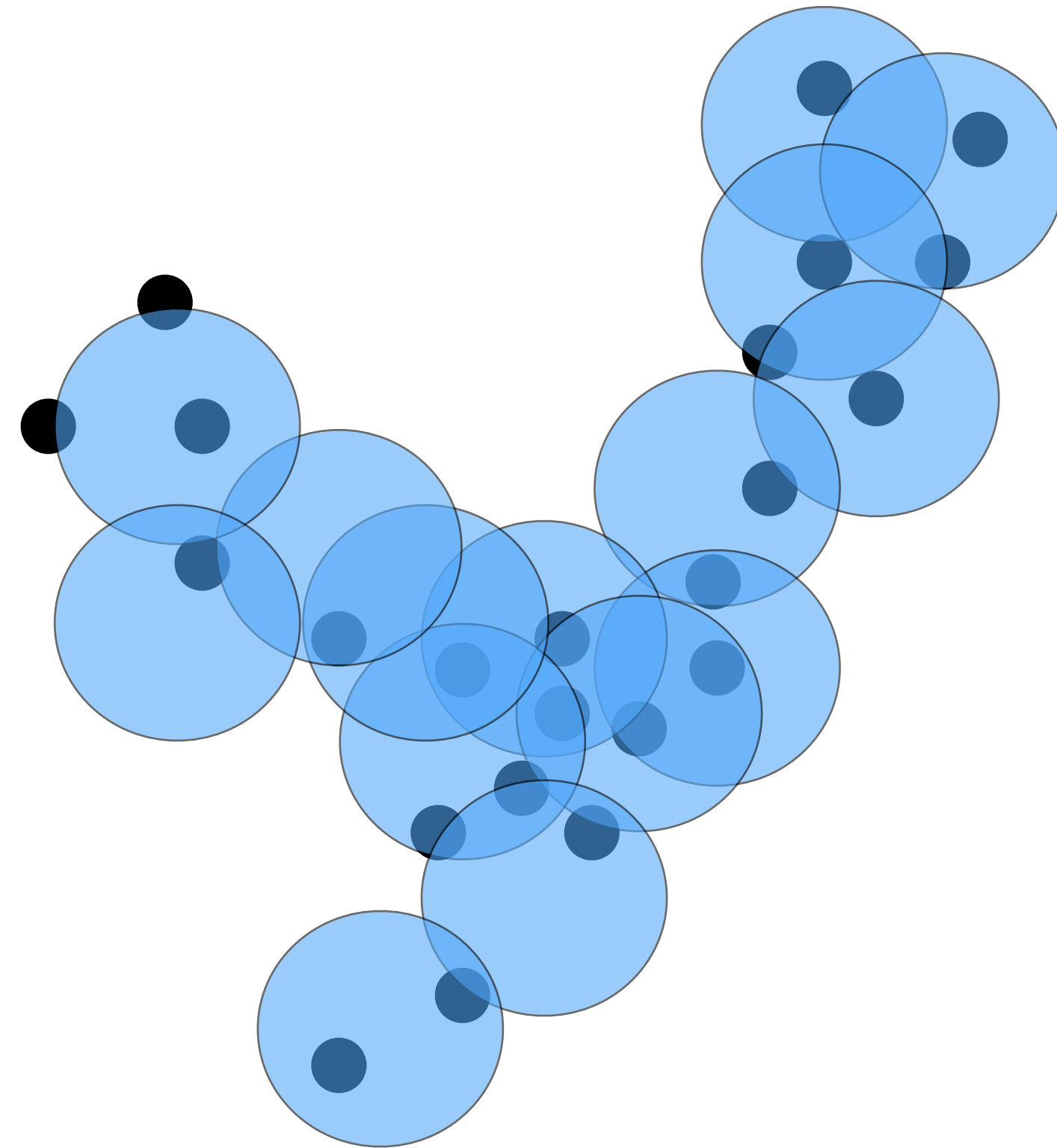


$p(x)$

# Mixture of Gaussians



Target distribution



$x \sim p_{\text{data}}(x)$

$p_{\theta}(x)$

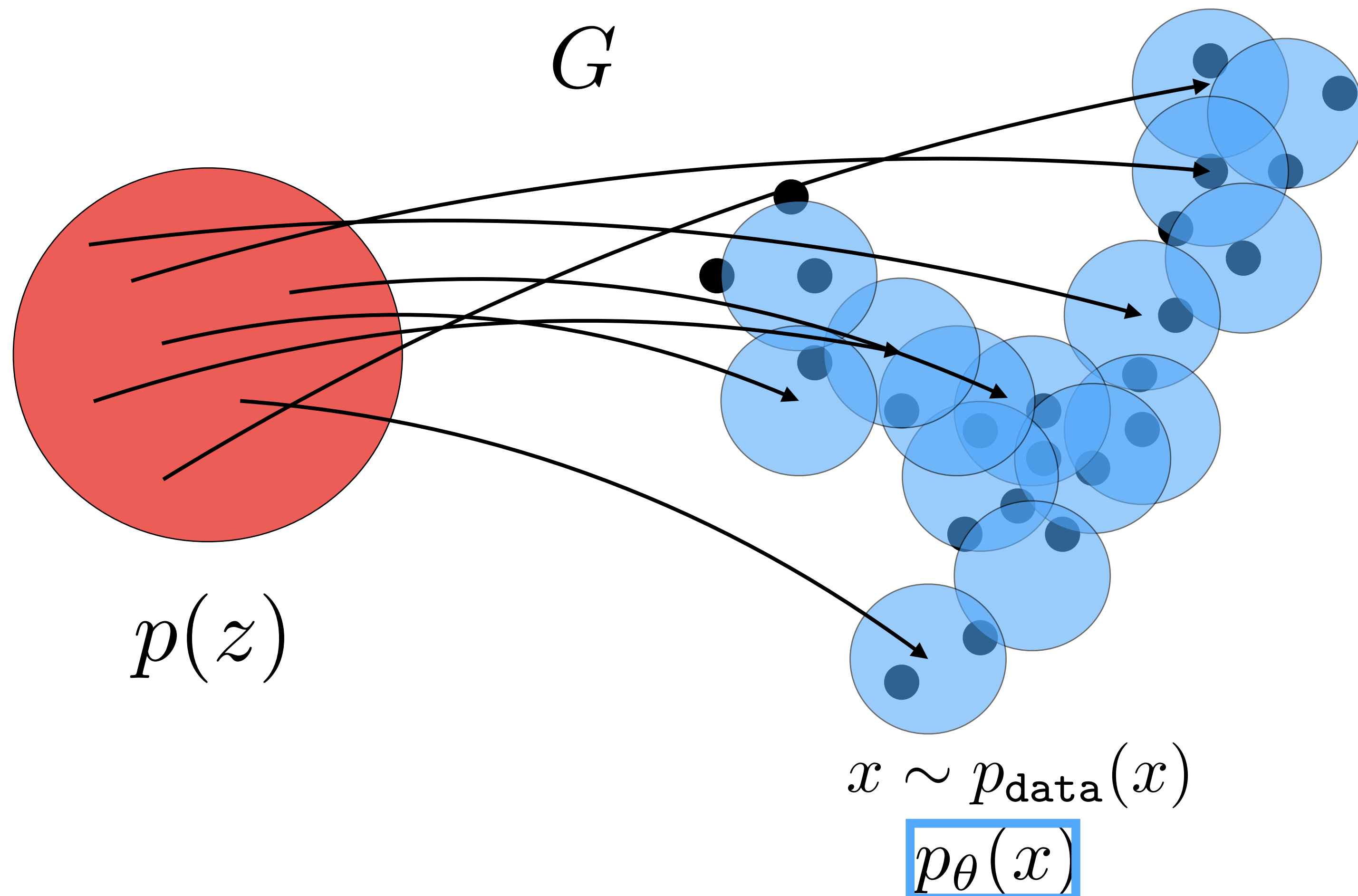
$$p_{\theta}(x) = \sum_{i=1}^k w_i \mathcal{N}(x; \mu_i, \Sigma_i)$$

# Variational Autoencoders (VAEs)

[Kingma & Welling, 2014; Rezende, Mohamed, Wierstra 2014]

Prior distribution

Target distribution



Density model:

$$p_\theta(x) = \int p(x|z; \theta) p(z) dz$$

$$p(x|z; \theta) \sim \mathcal{N}(x; G_\theta^\mu(x), G_\theta^\sigma(x))$$

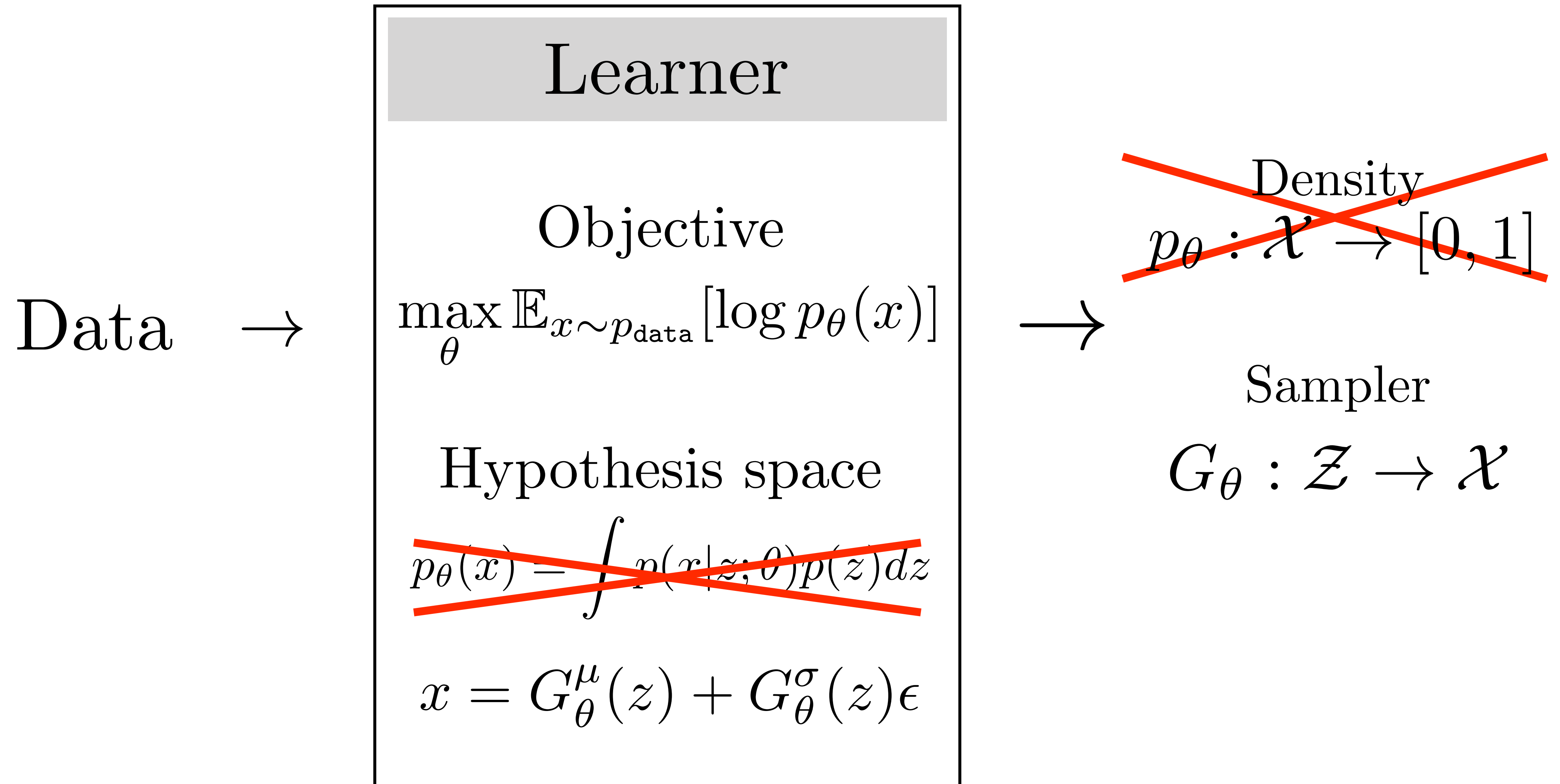
Sampling:

$$z \sim p(z) \quad \epsilon \sim \mathcal{N}(0, 1)$$

$$x = G_\theta^\mu(z) + G_\theta^\sigma(z)\epsilon$$



# Variational Autoencoder (VAE)



# Variational Autoencoders (VAEs) — Training

Fitting a model to data requires computing  $p_{\theta}(x)$

How to compute  $p_{\theta}(x)$  efficiently?

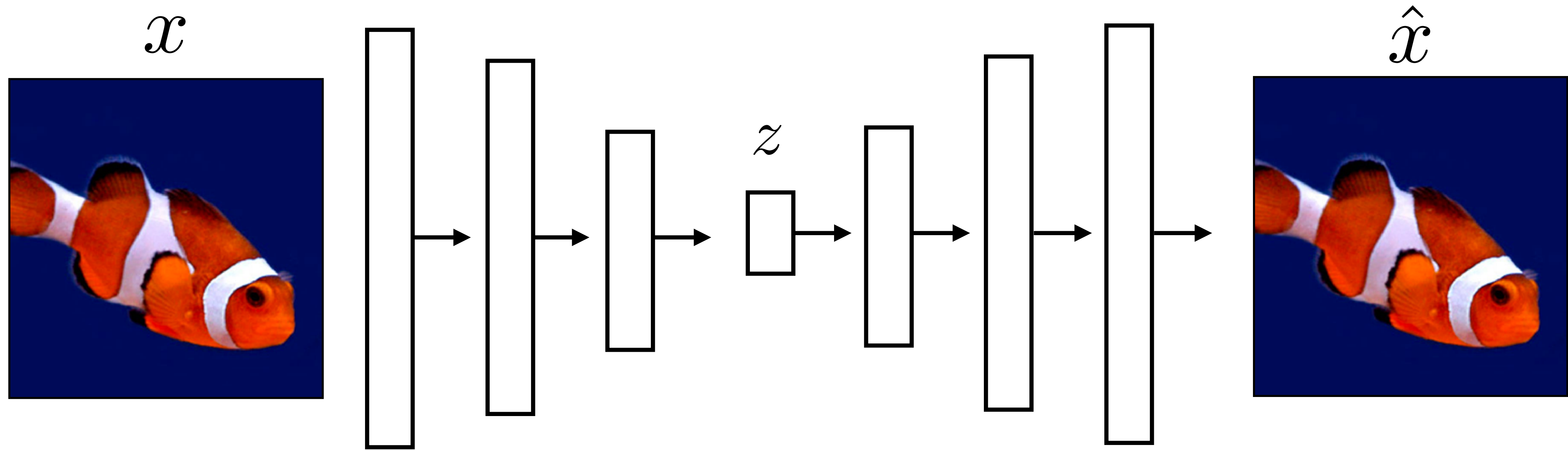
$$p_{\theta}(x) = \int p(x|z; \theta)p(z)dz \quad \longleftarrow \quad \text{almost all terms are near zero}$$

Train “inference network”  $q_{\psi}(z|x)$

to give distribution over the  $z$ 's that are likely to produce  $x$

Approximate  $p_{\theta}(x)$  with  $\mathbb{E}_{q_{\psi}(z|x)}[p_{\theta}(x|z)]$

# Variational Autoencoders (VAEs) — Optimization (aka training)

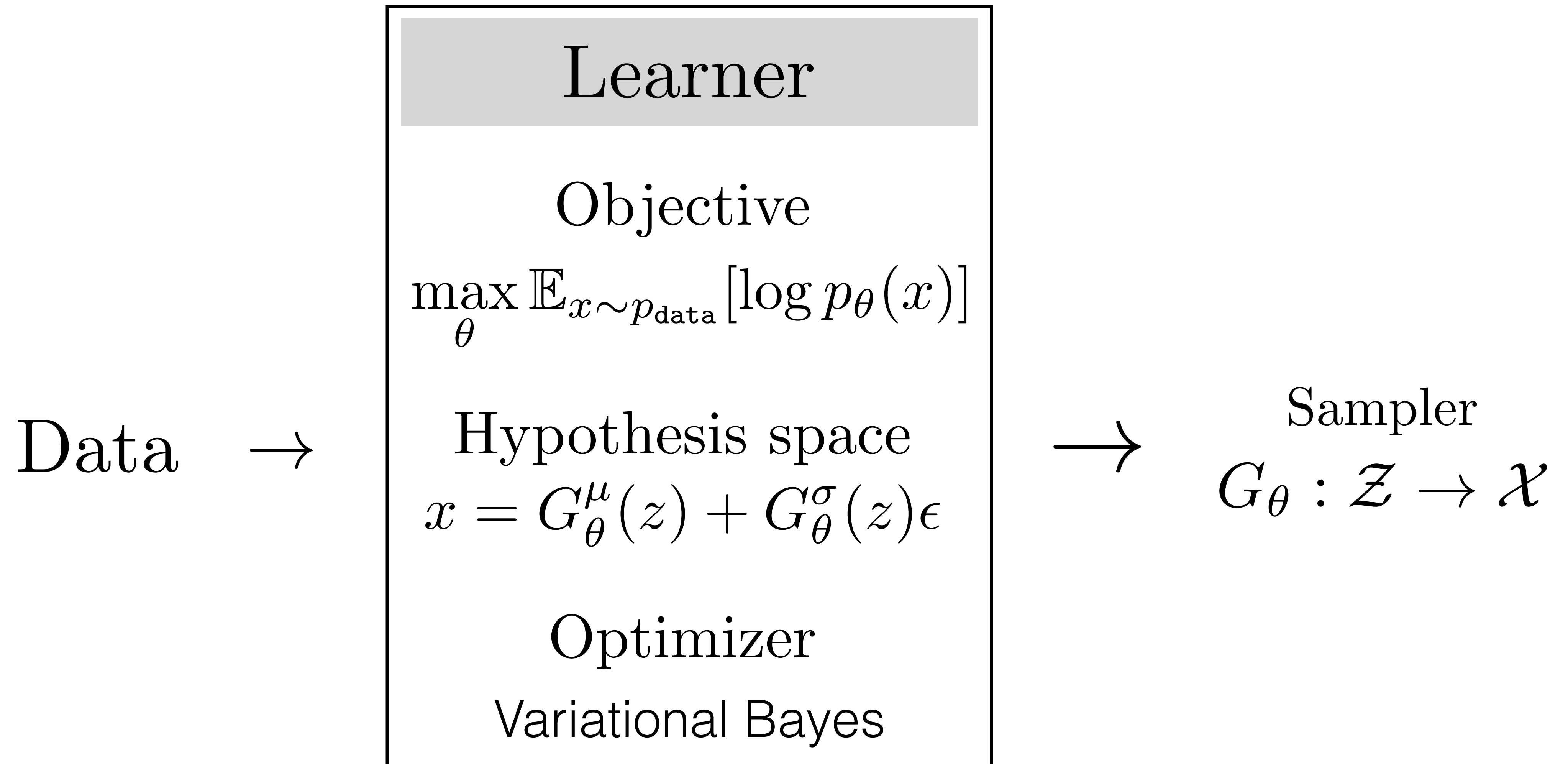


$$\max_{\theta} \mathbb{E}_{x \sim p_{\text{data}}} [\log p_{\theta}(x)]$$

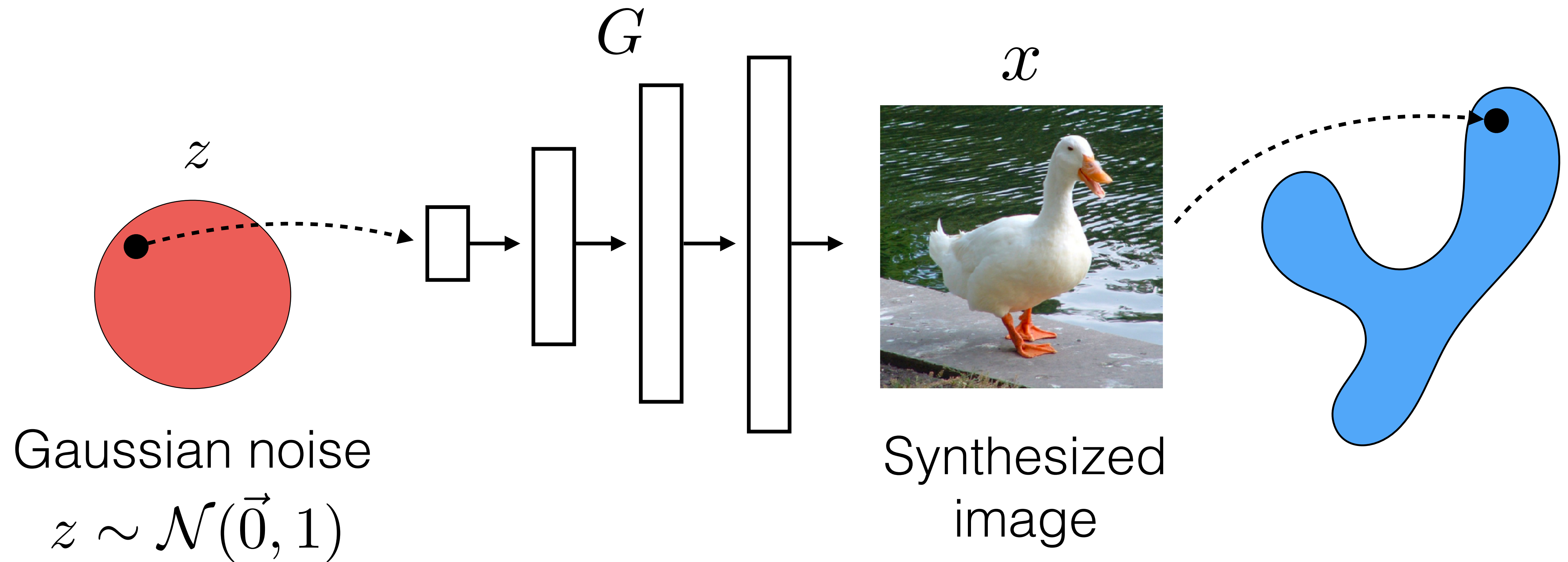
$$\geq \max_{\theta, \psi} \mathbb{E}_{x_i \sim p_{\text{data}}} [\mathbb{E}_{q_{\psi}(z|x_i)} [p_{\theta}(x|z)] - \text{KL}(q_{\psi}(z|x_i) || p(z))]$$

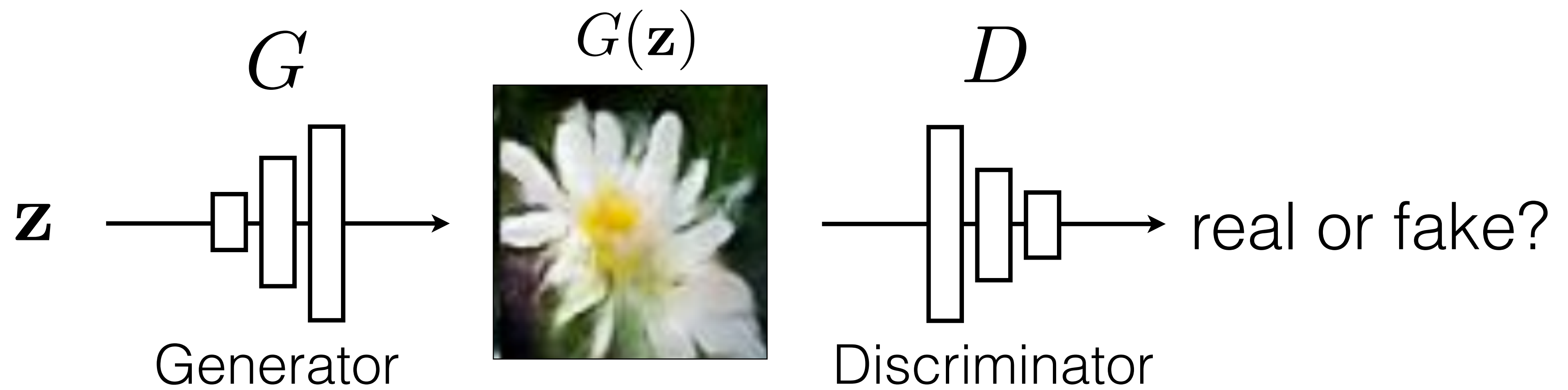


# Variational Autoencoder (VAE)



# Generative Adversarial Networks (GANs)

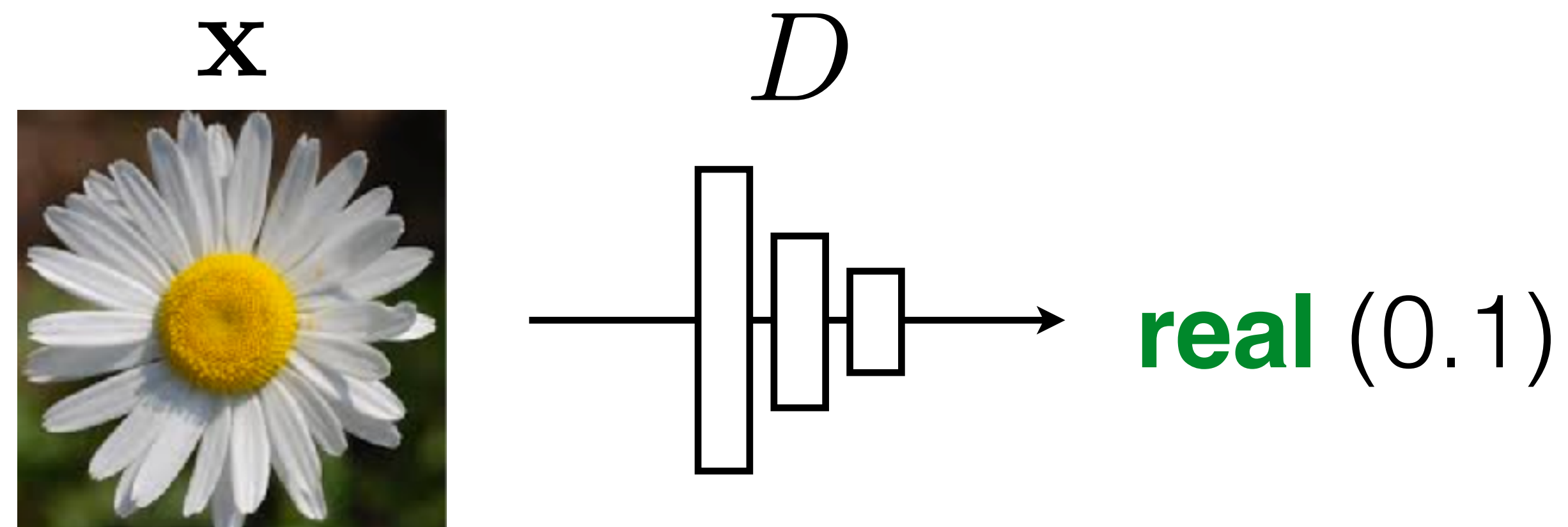
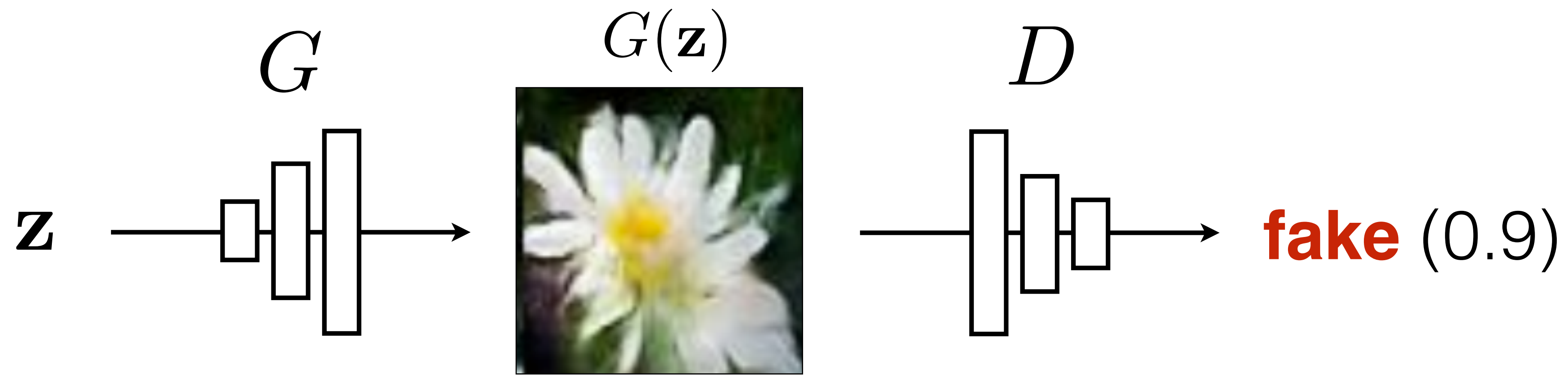




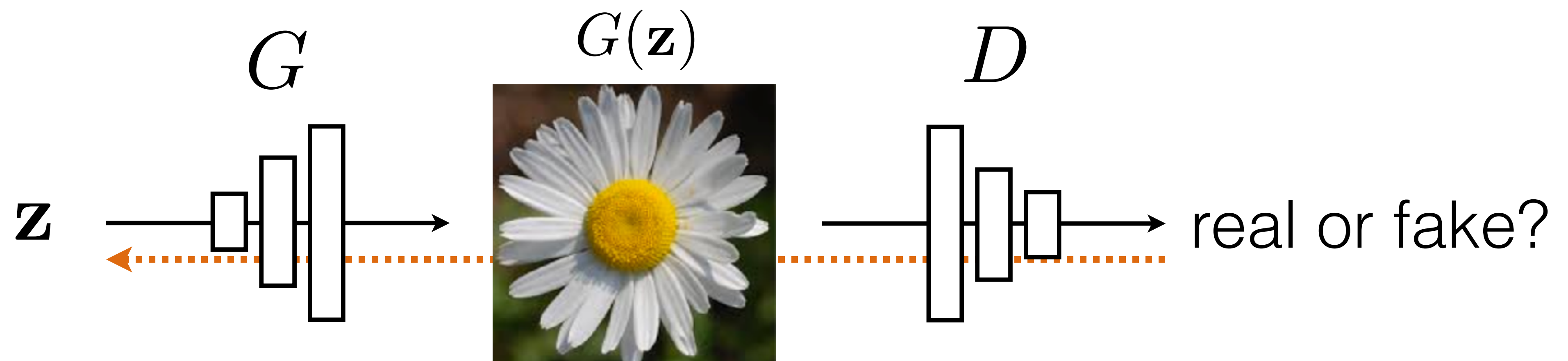
**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes



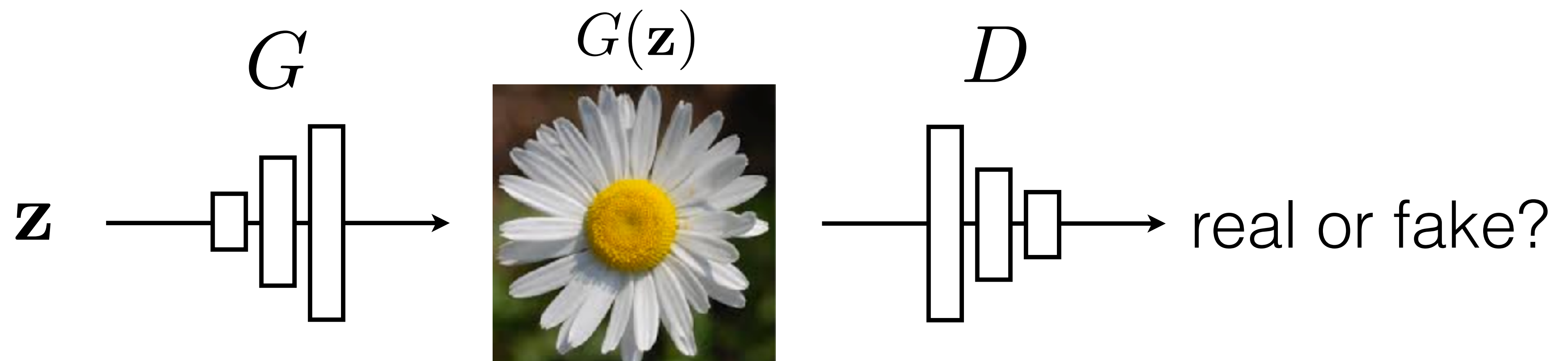


$$\arg \max_D \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[ \log D(G(\mathbf{z})) + \log (1 - D(\mathbf{x})) \right]$$



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

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

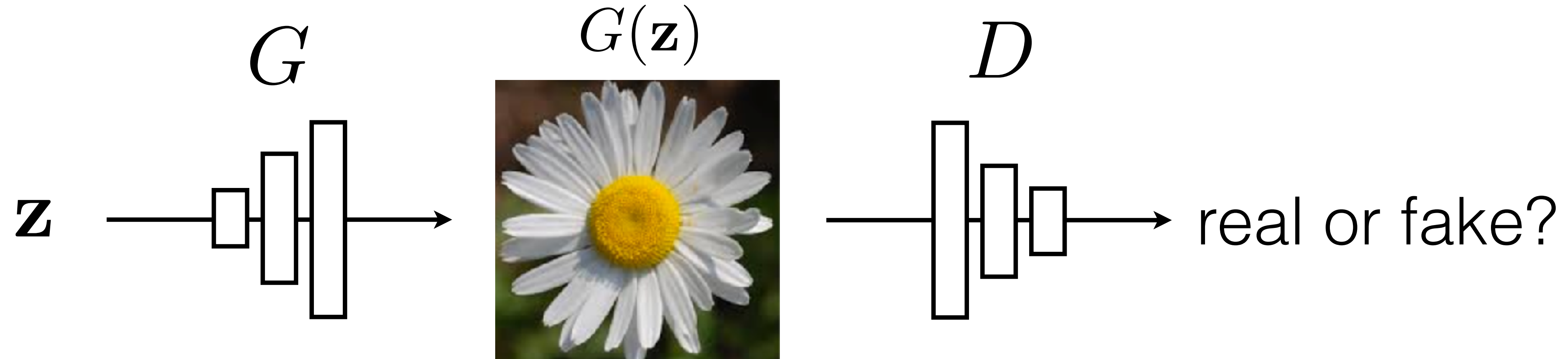


**G** tries to synthesize fake images that *fool* the *best* **D**:

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



# 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.

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

Proof

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

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

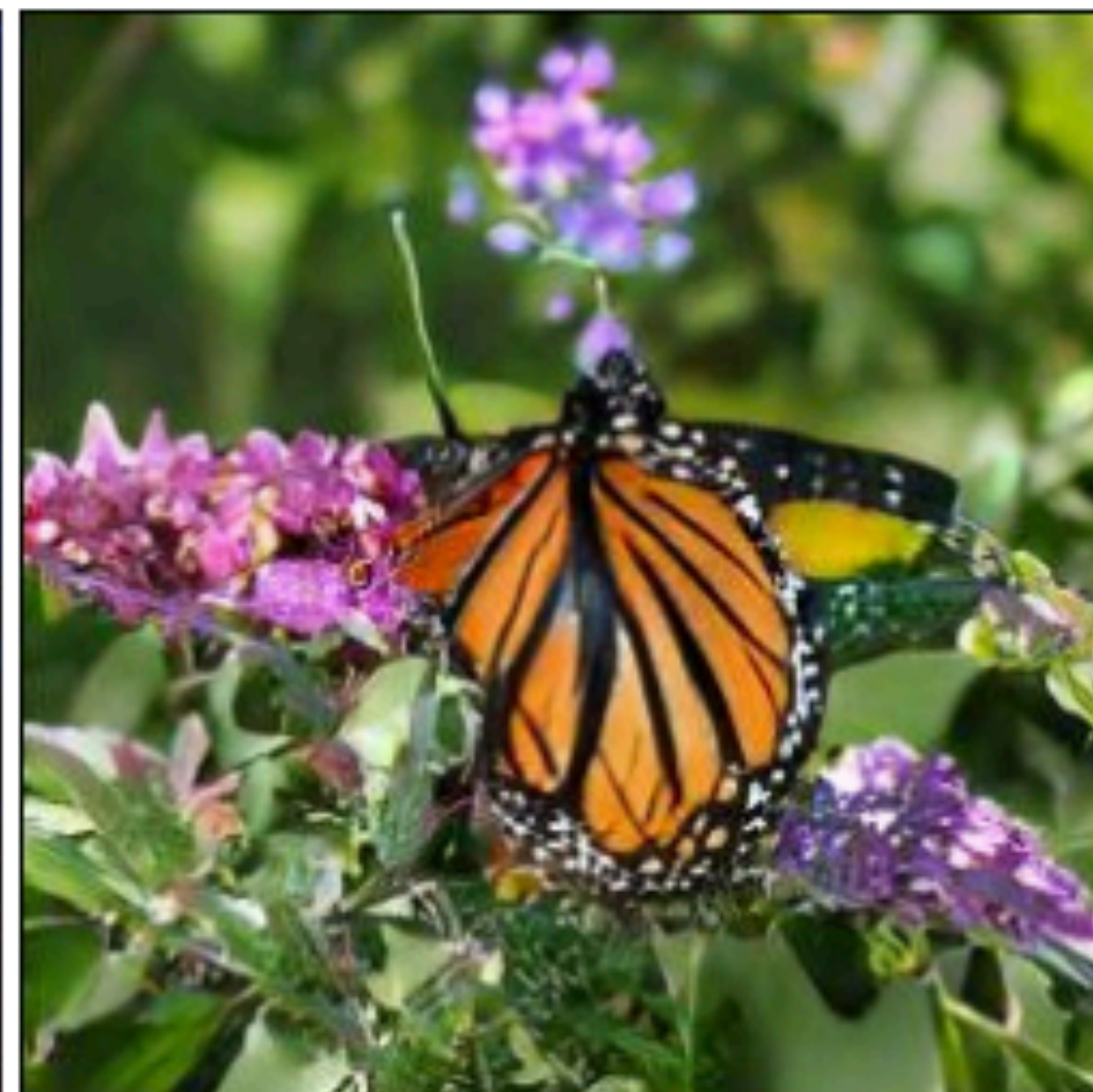
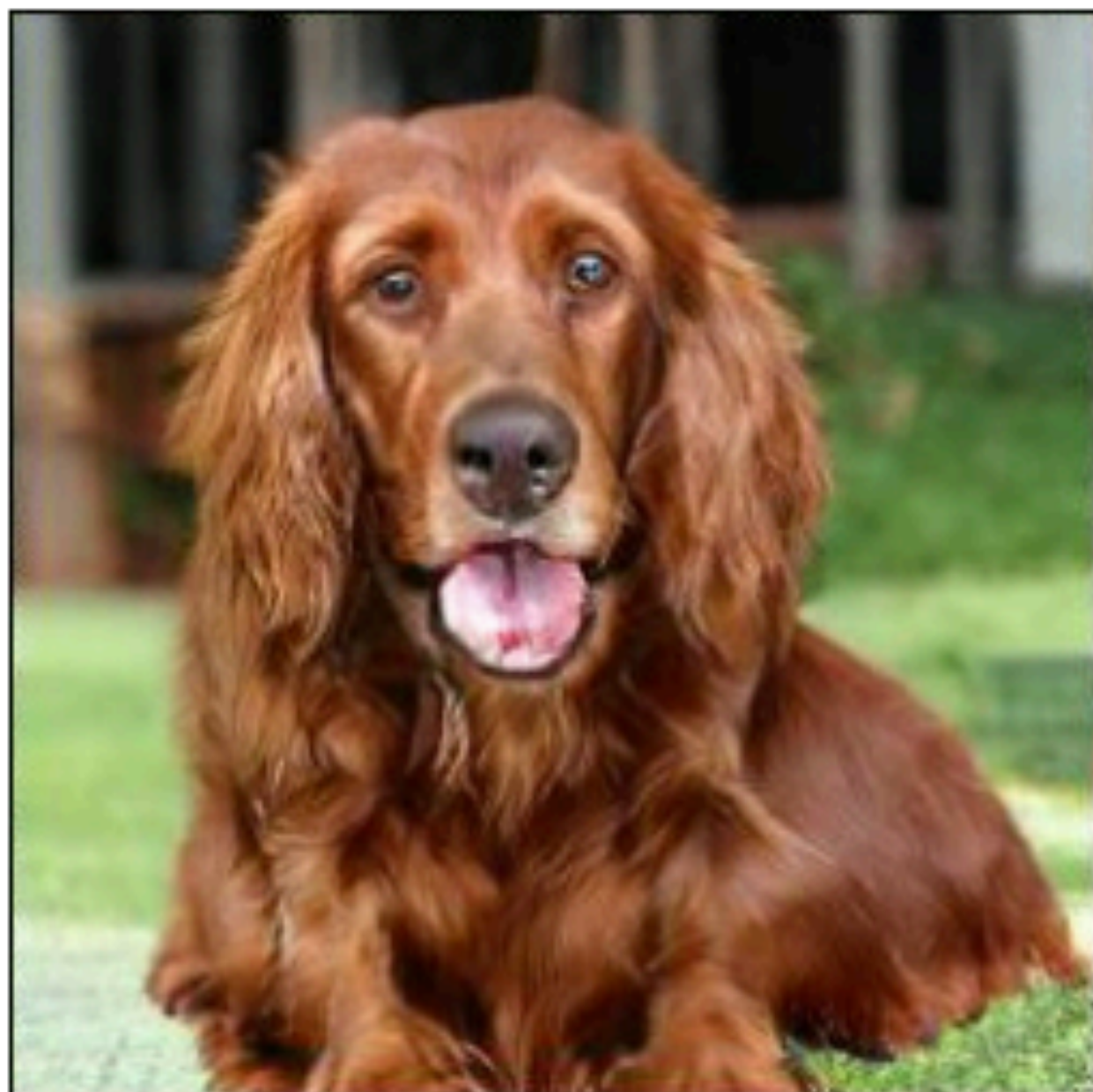
$$C(G) = -\log(4) + 2 \cdot \underbrace{JSD(p_{data} \| p_g)}$$

$$\geq 0, \quad 0 \iff p_g = p_{data} \quad \square$$



# Samples from BigGAN

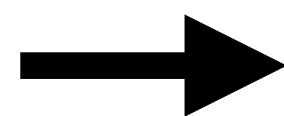
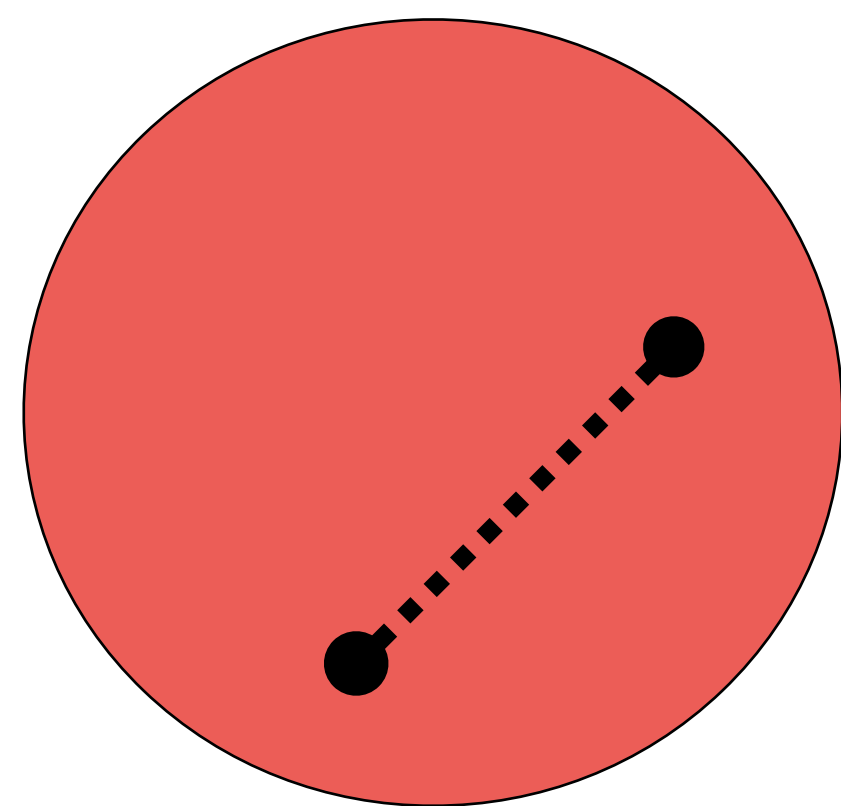
[Brock et al. 2018]





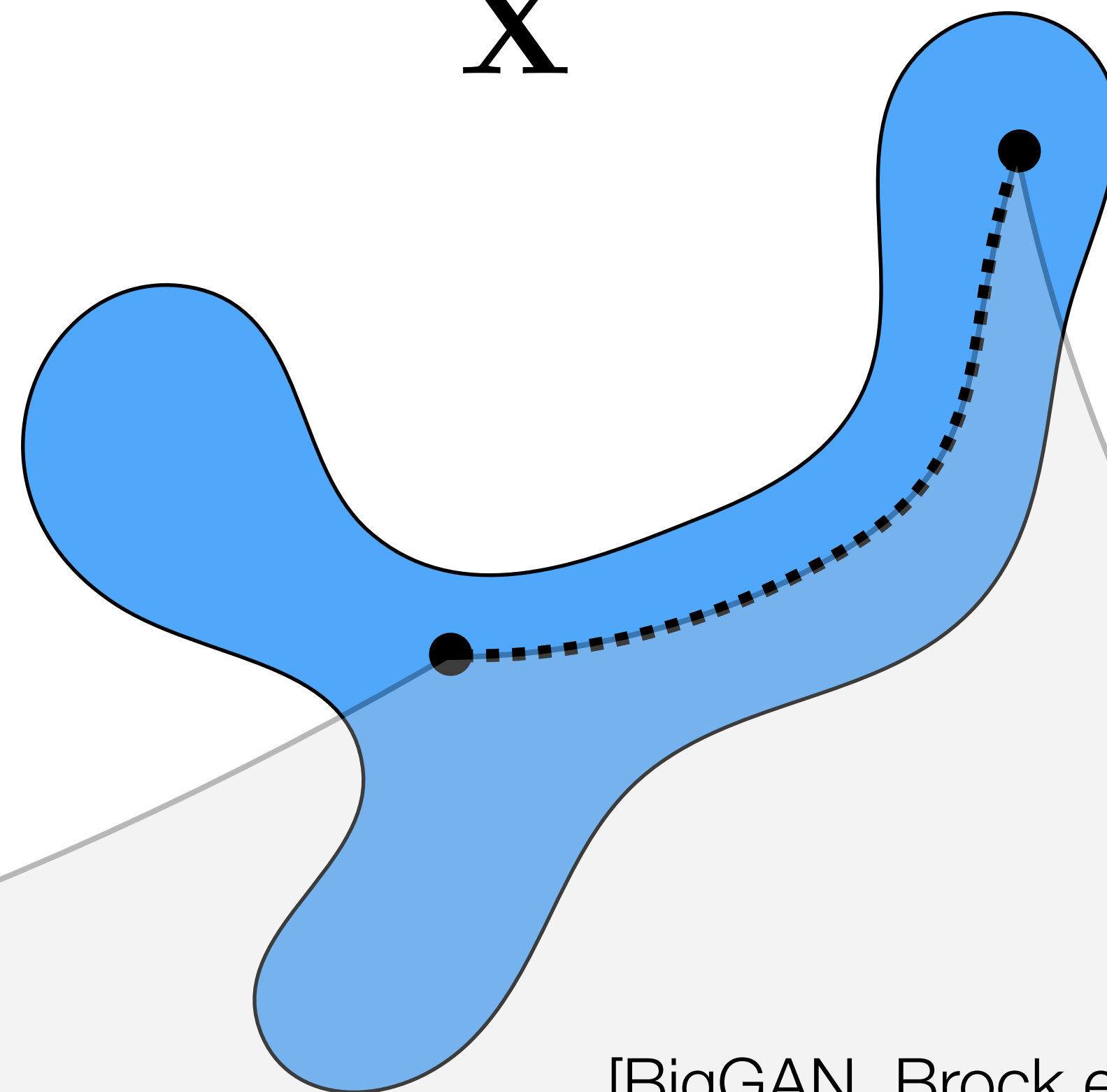
Latent space  
(Gaussian)

$\mathbf{z}$



Data space  
(Natural image manifold)

$\mathbf{X}$

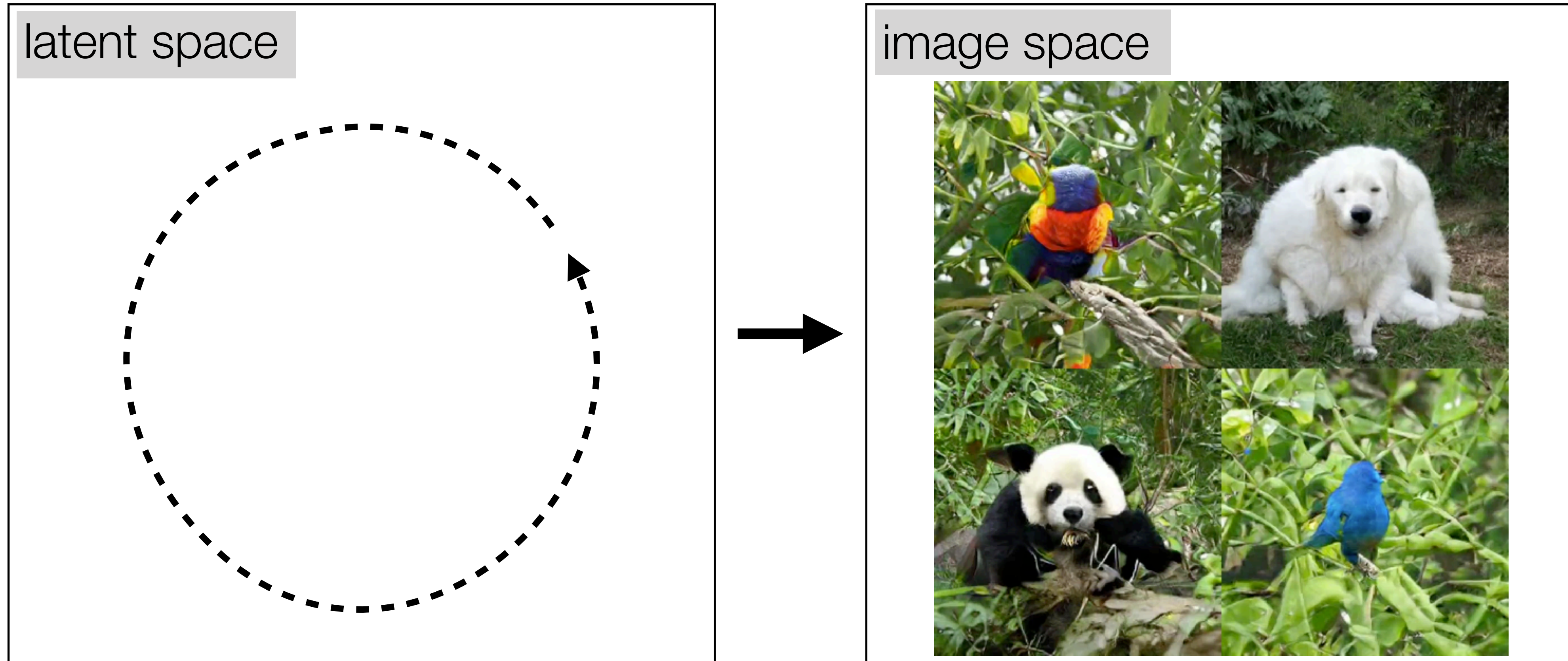


[BigGAN, Brock et al. 2018]

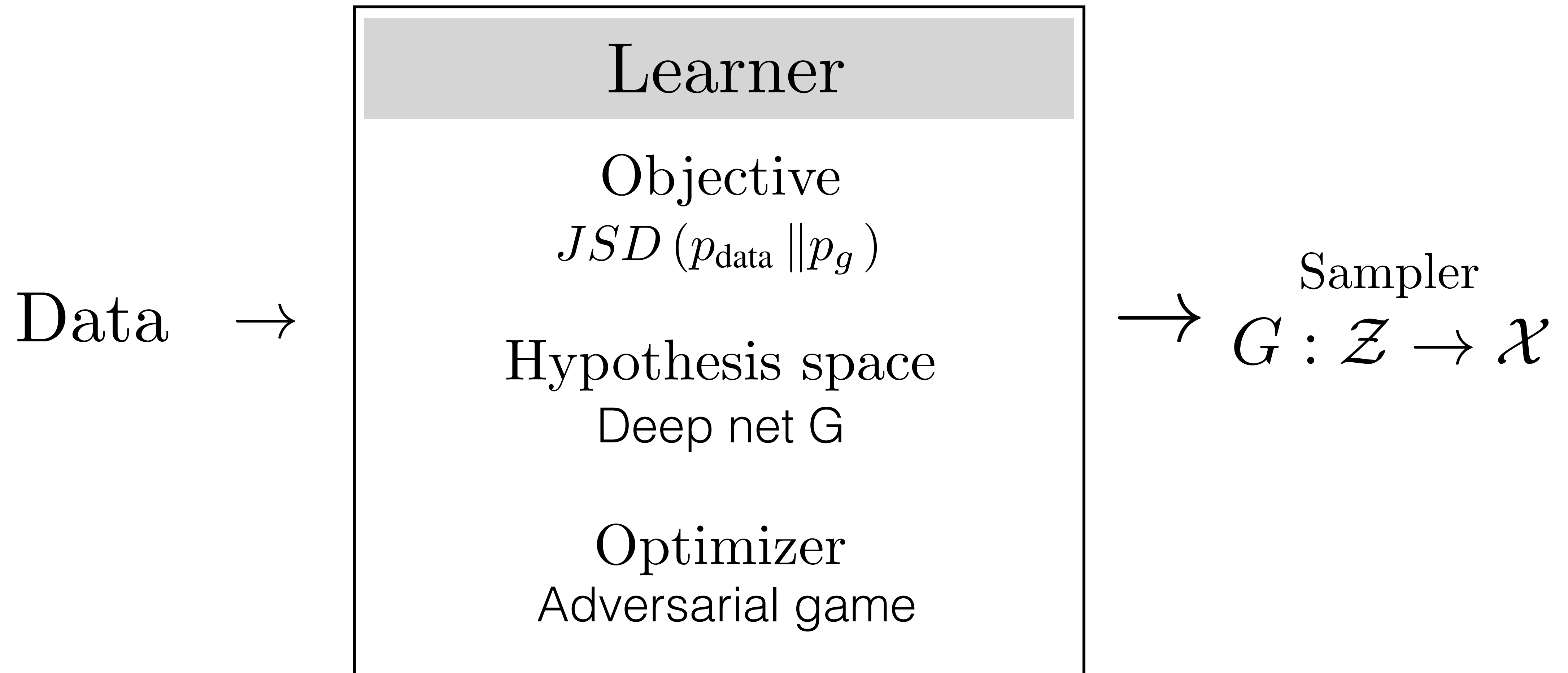




# Generative models organize the manifold of natural images

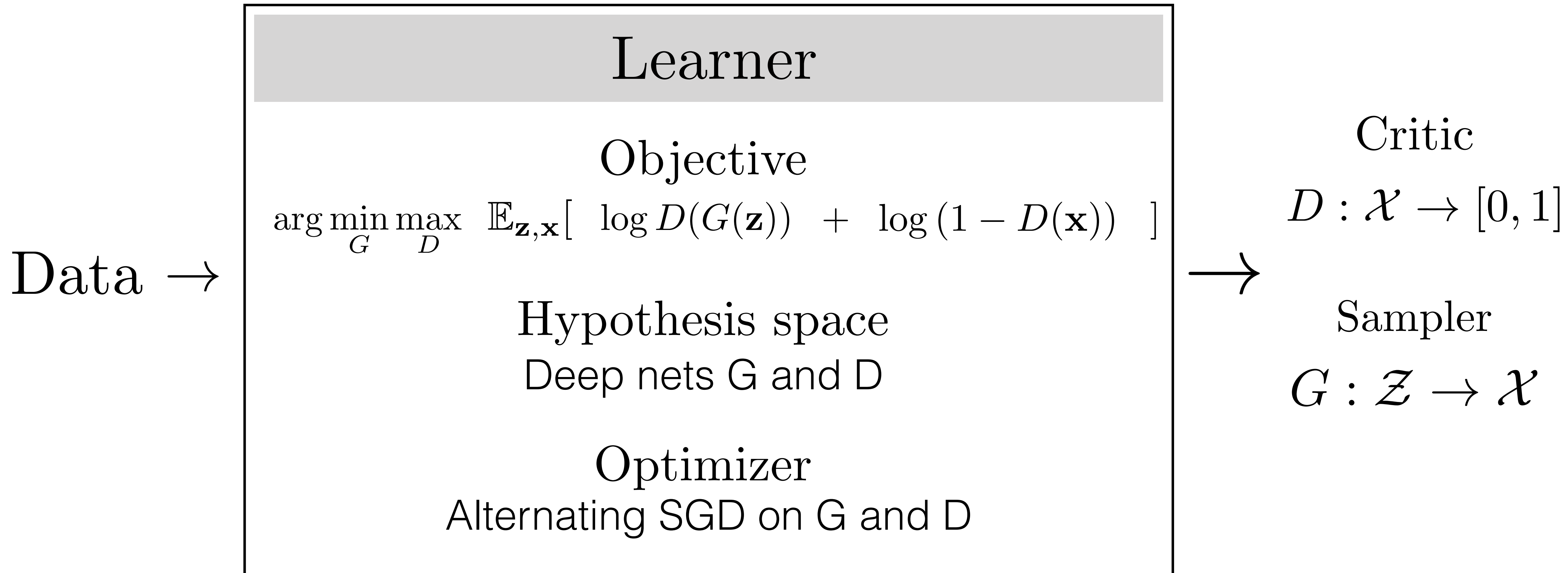


# Generative Adversarial Network





# Generative Adversarial Network



# 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](https://arxiv.org/abs/1406.2661)

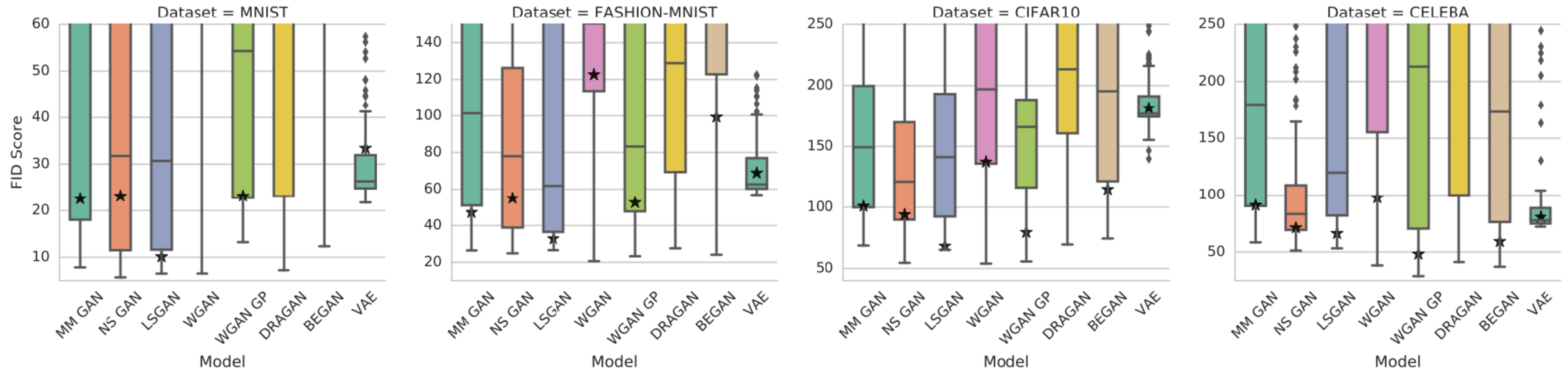
[arxiv.org/abs/1511.06434](https://arxiv.org/abs/1511.06434) [arxiv.org/abs/1606.07536](https://arxiv.org/abs/1606.07536) [arxiv.org/abs/1710.10196](https://arxiv.org/abs/1710.10196)

[arxiv.org/abs/1812.04948](https://arxiv.org/abs/1812.04948)





# Better objectives? optimizers?



**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 Equal?”, Lucic\*, Kurach\*, et al. 2018]

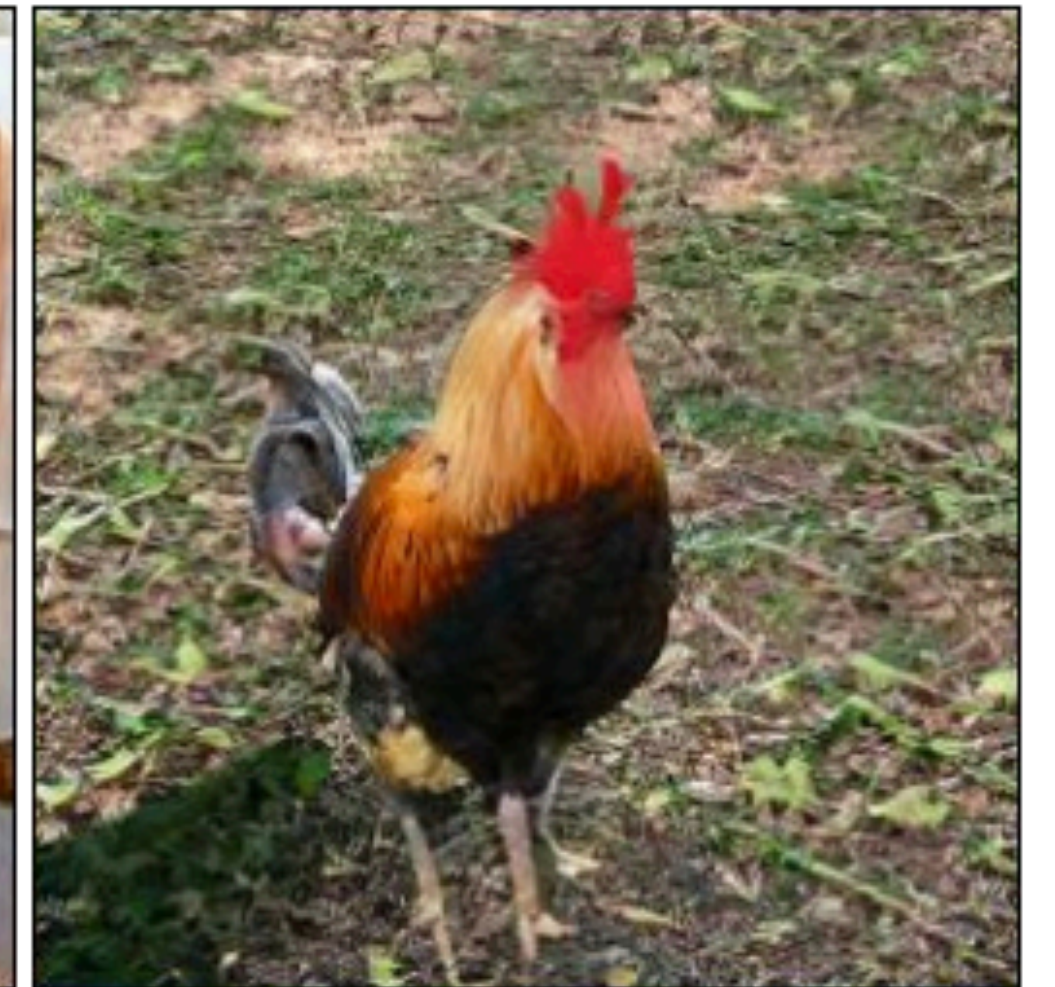


# More data?

ACGAN [Odena et al. 2016]



BigGAN [Brock et al. 2018]



*Both trained on Imagenet*



# Architectures

DCGAN

[Radford, Metz, Chintala 2016]



StyleGAN

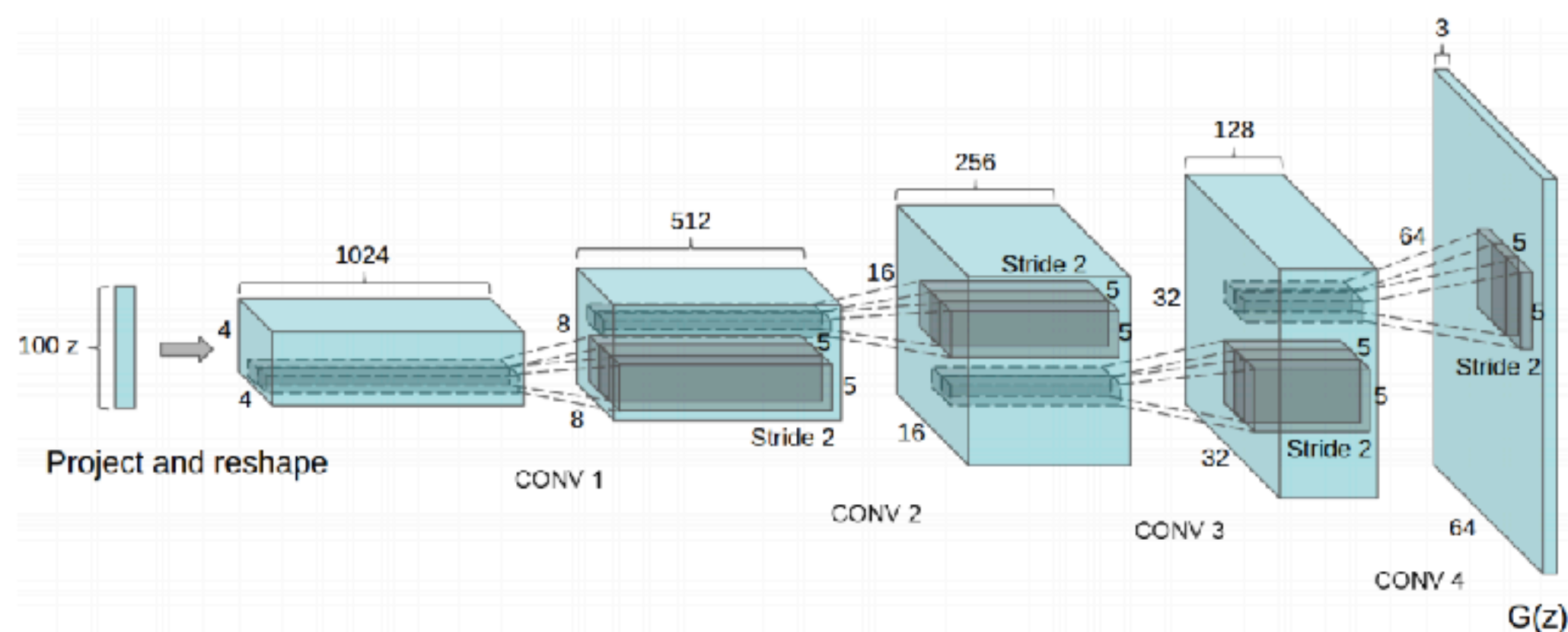
[Karras, Laine, Aila 2019]



# Architectures

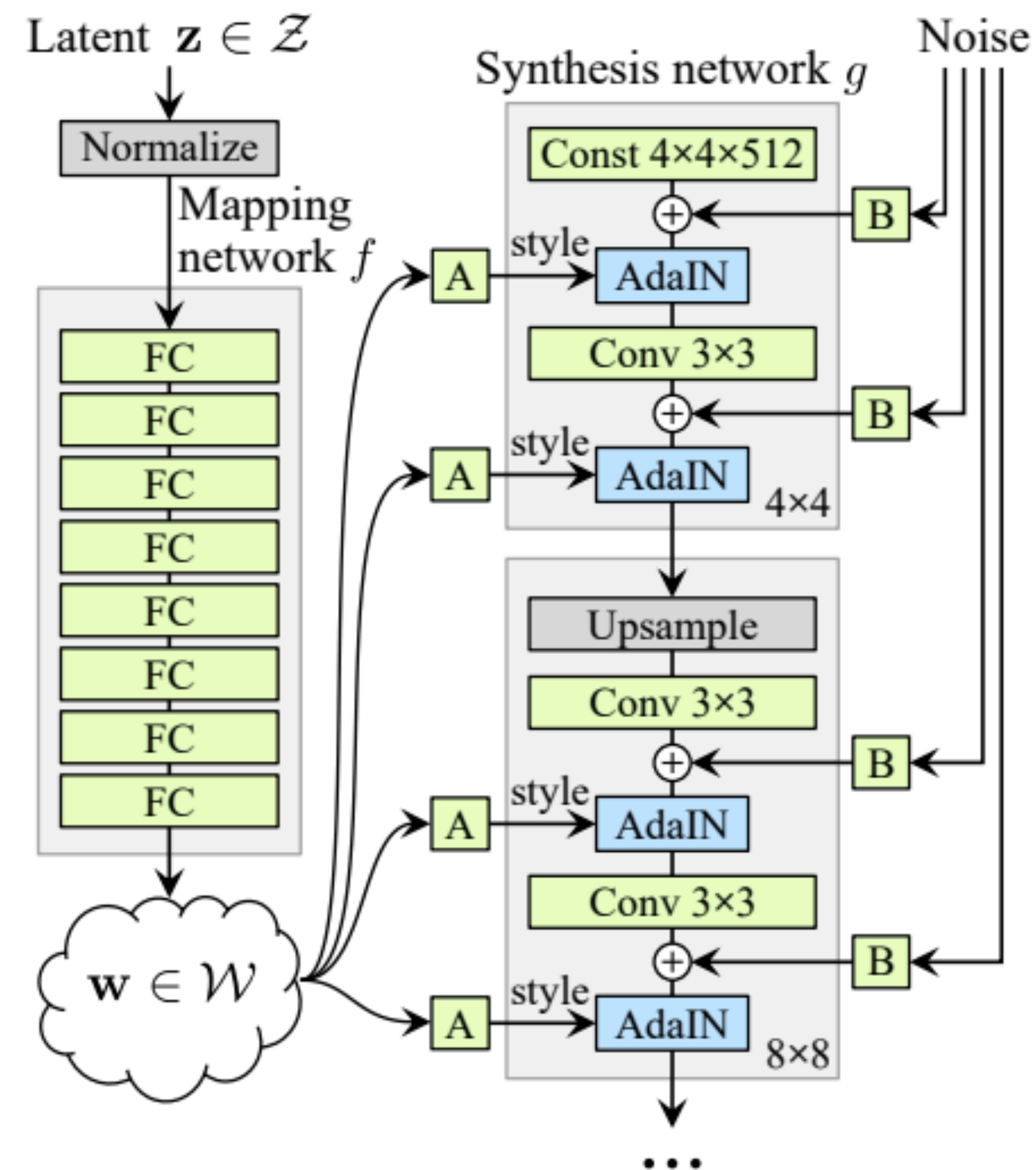
## DCGAN

[Radford, Metz, Chintala 2016]



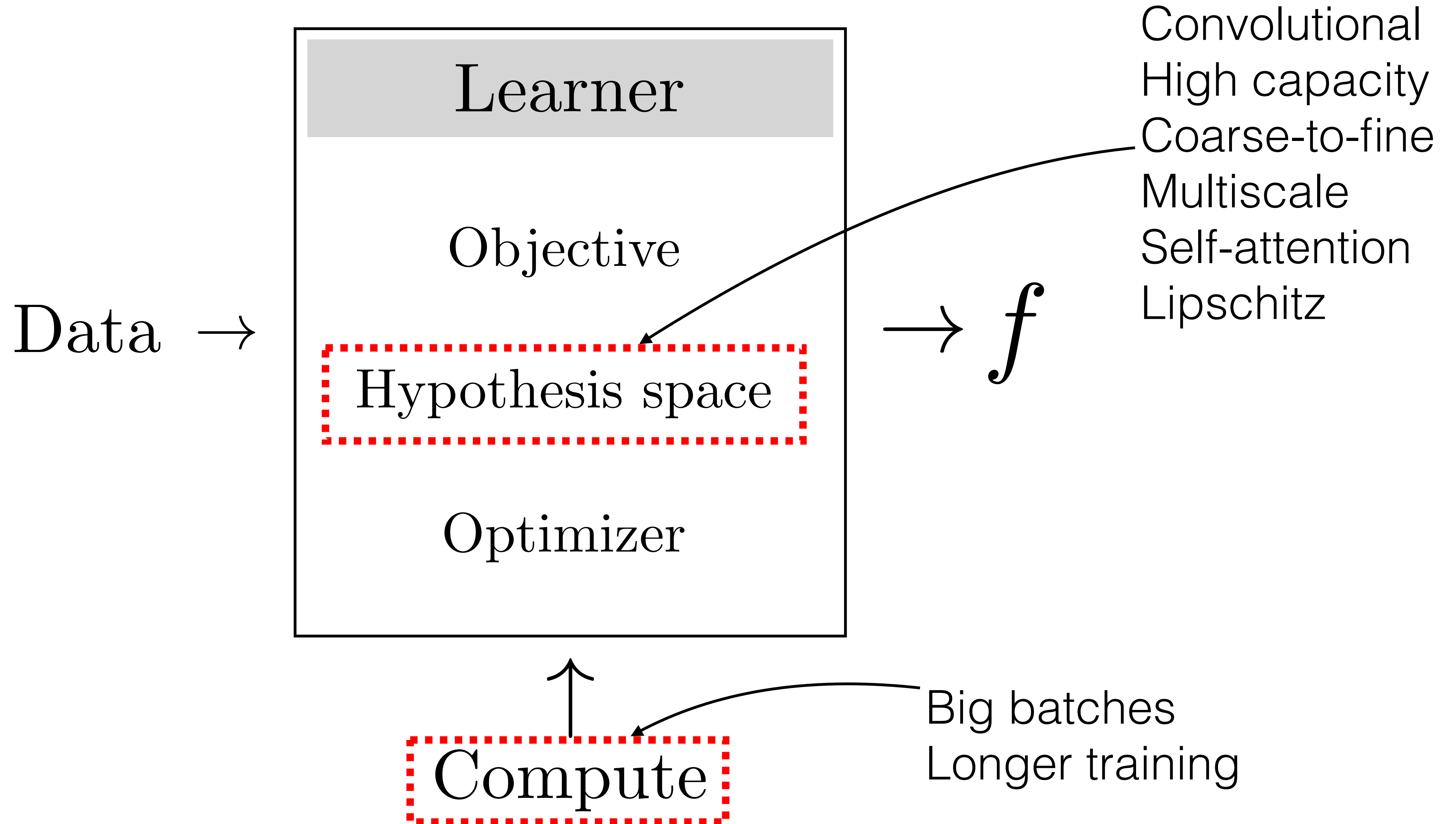
## StyleGAN

[Karras, Laine, Aila 2019]





# What has driven GAN progress?



## **VAEs**

Pros: Cheap to sample, good coverage

Cons: Blurry samples (in practice)

## **GANs**

Pros: Cheap to sample, fast to train, require little data

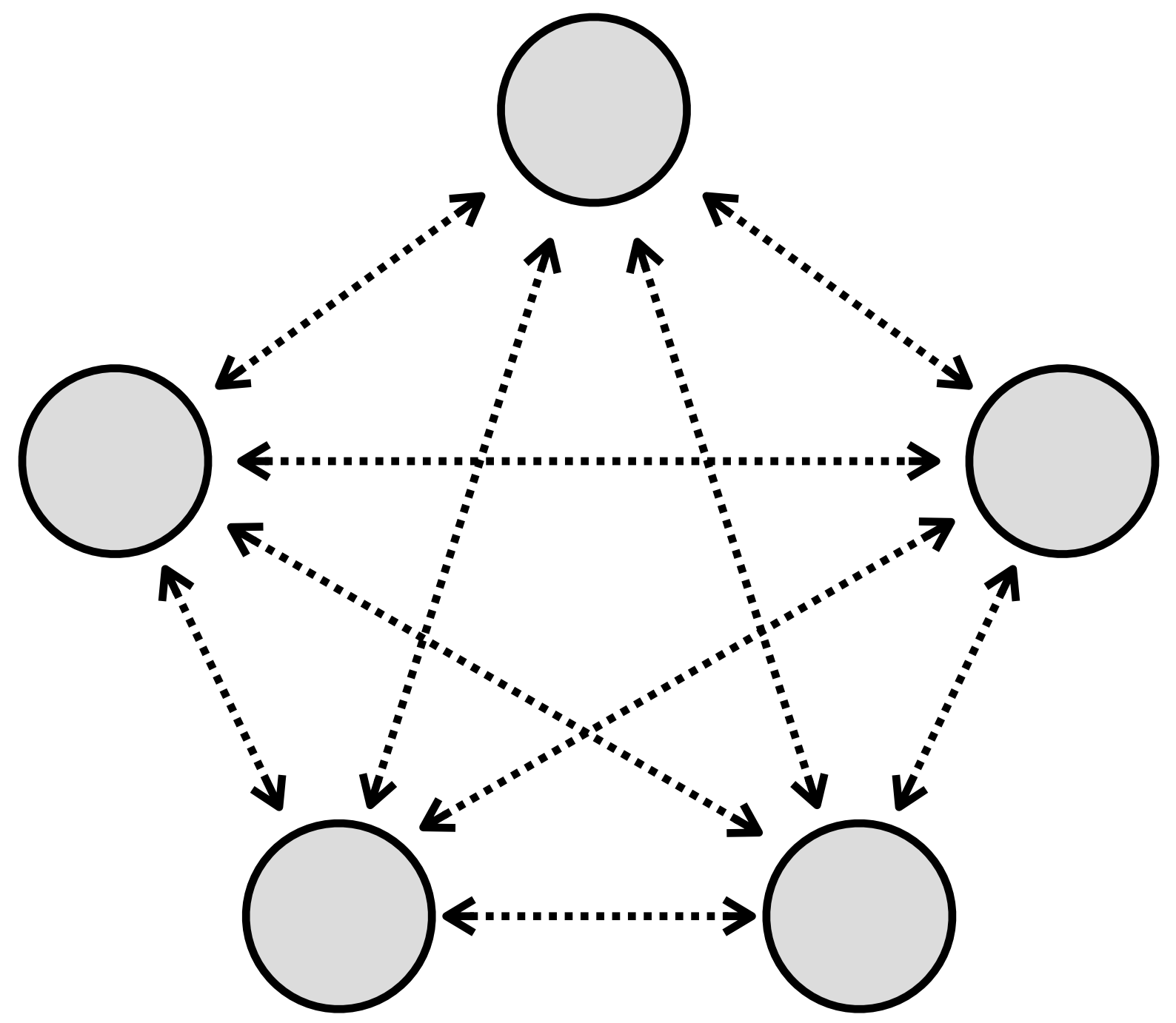
Cons: No likelihoods, bad coverage (mode collapse), finicky to train (minimax)

Other deep generative models:

**Autoregressive models, Normalizing flows, Energy-based models**



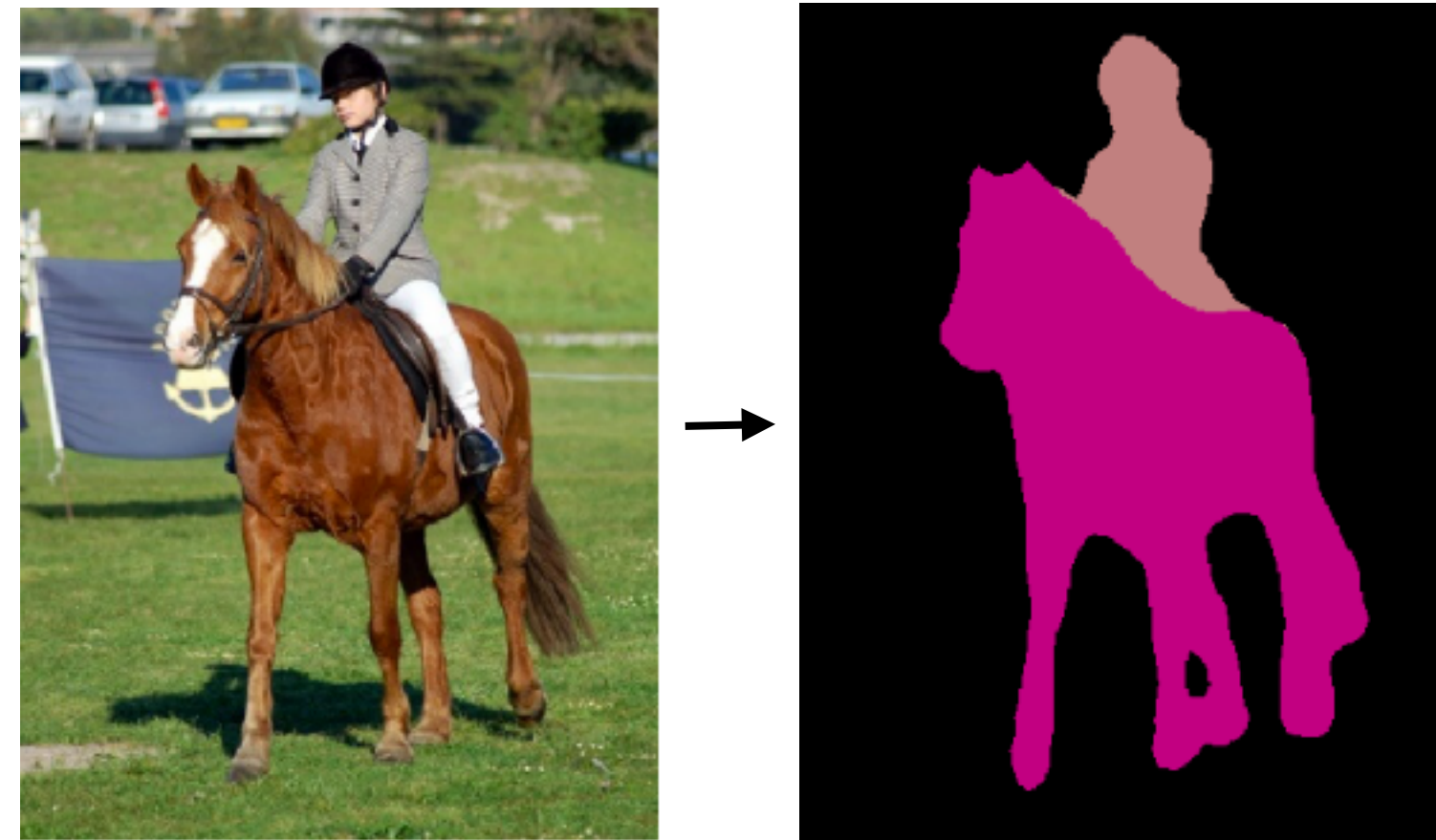
- 1. Image synthesis
- 2. Structured prediction**
- 3. Domain mapping



# Strutured Prediction

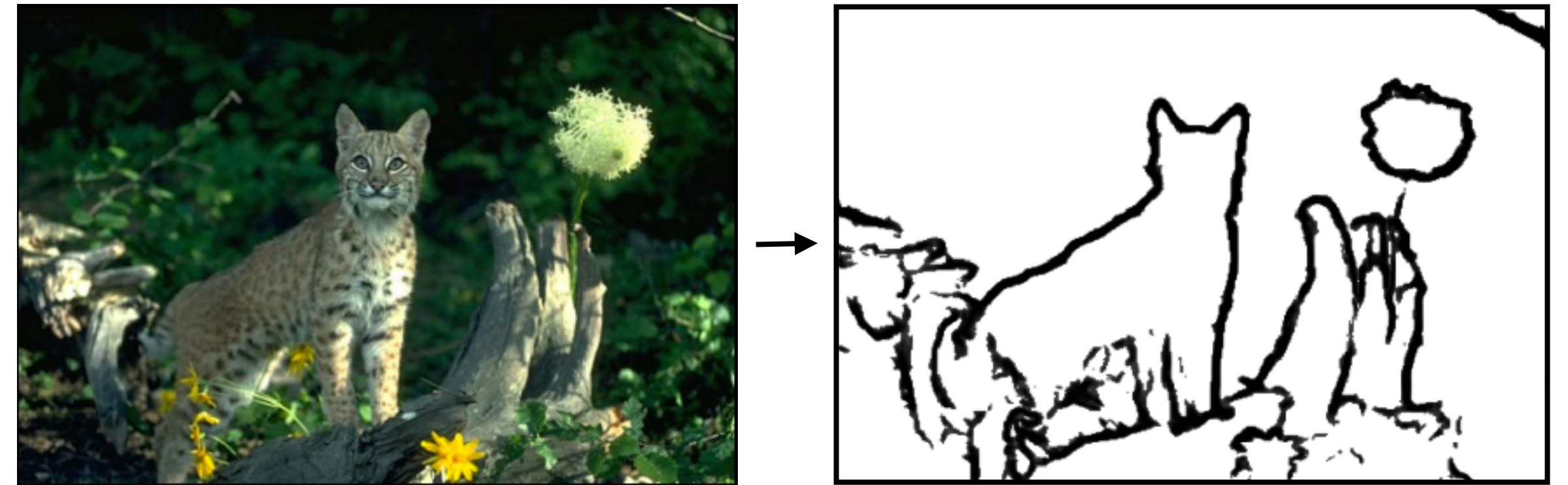
# Data prediction problems (“structured prediction”)

## Semantic segmentation



[Long et al. 2015, ...]

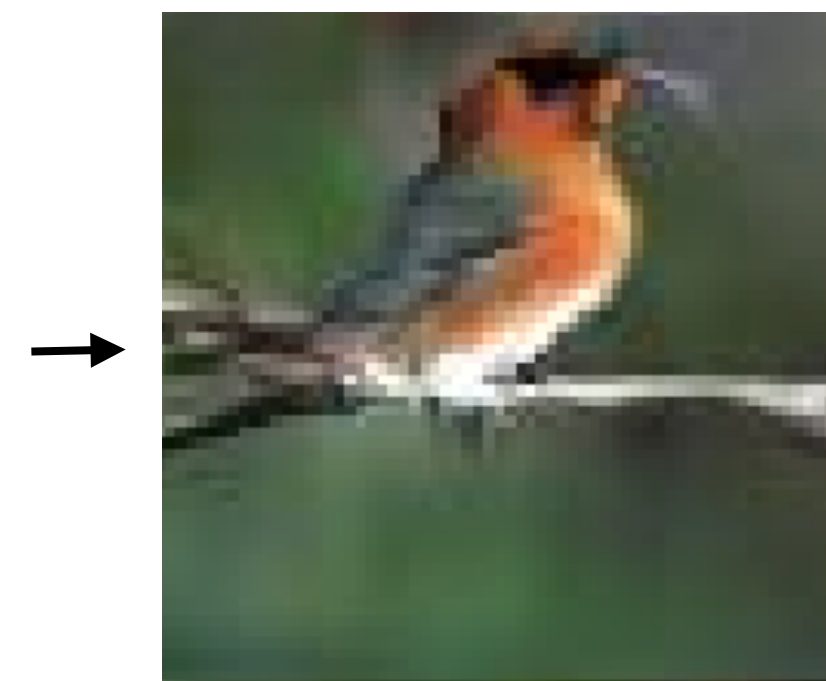
## Edge detection



[Xie et al. 2015, ...]

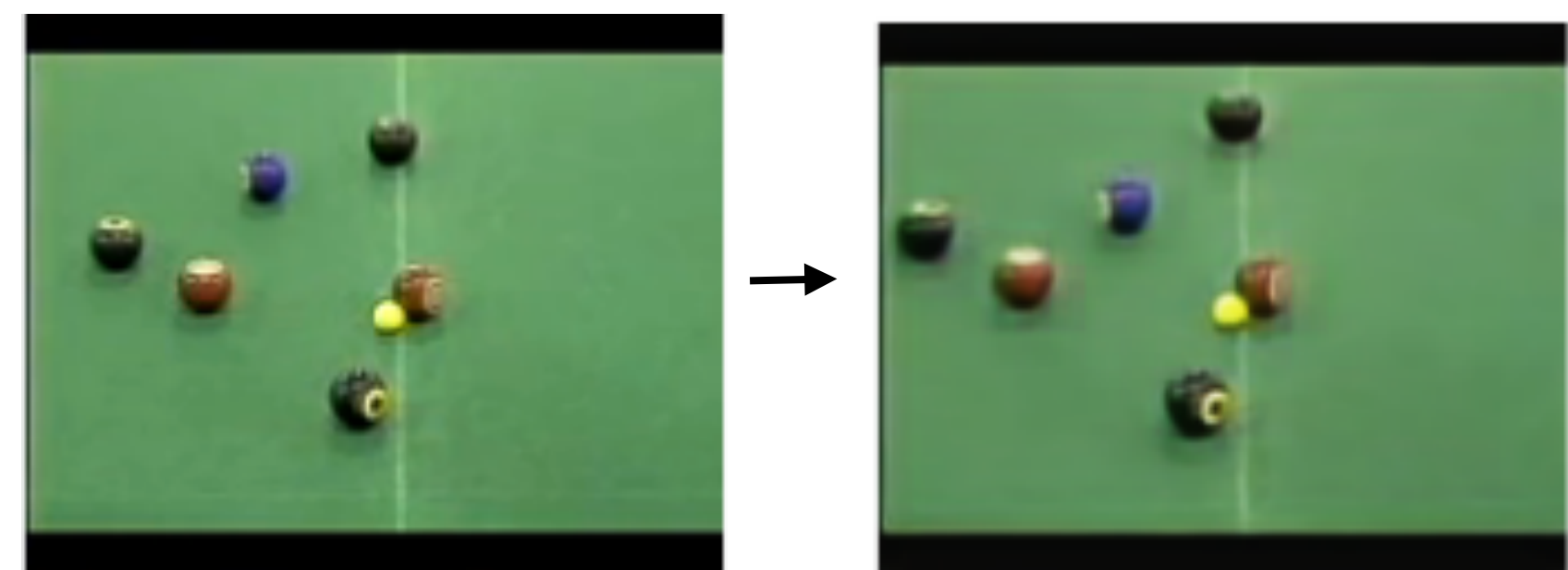
## Text-to-photo

“this small bird has a pink breast and crown...”



[Reed et al. 2014, ...]

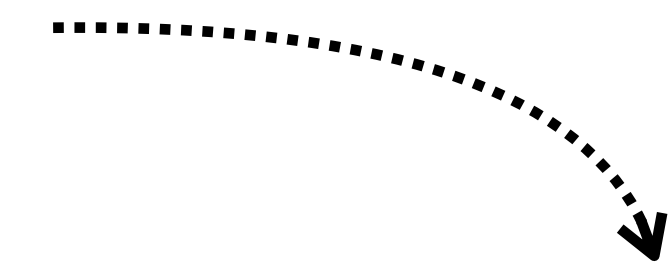
## Future frame prediction



[Mathieu et al. 2016, ...]

# Structured prediction

**X is high-dimensional**



Model *joint* distribution of high-dimensional data  $P(\mathbf{X}|\mathbf{Y} = \mathbf{y})$

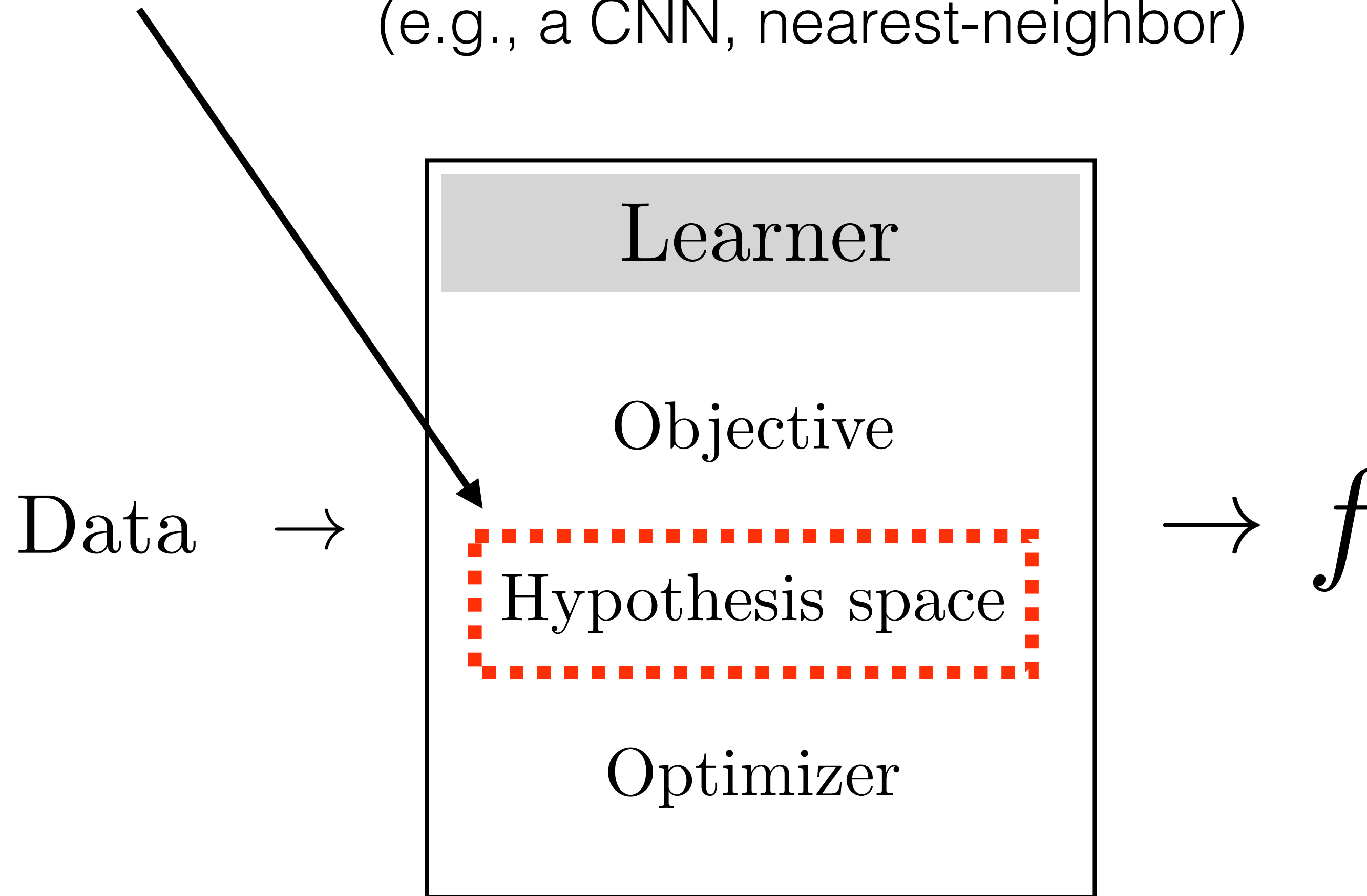
In vision this is usually what we are interested in

Unstructured:  $\prod_i p(X_i|\mathbf{Y} = \mathbf{y})$

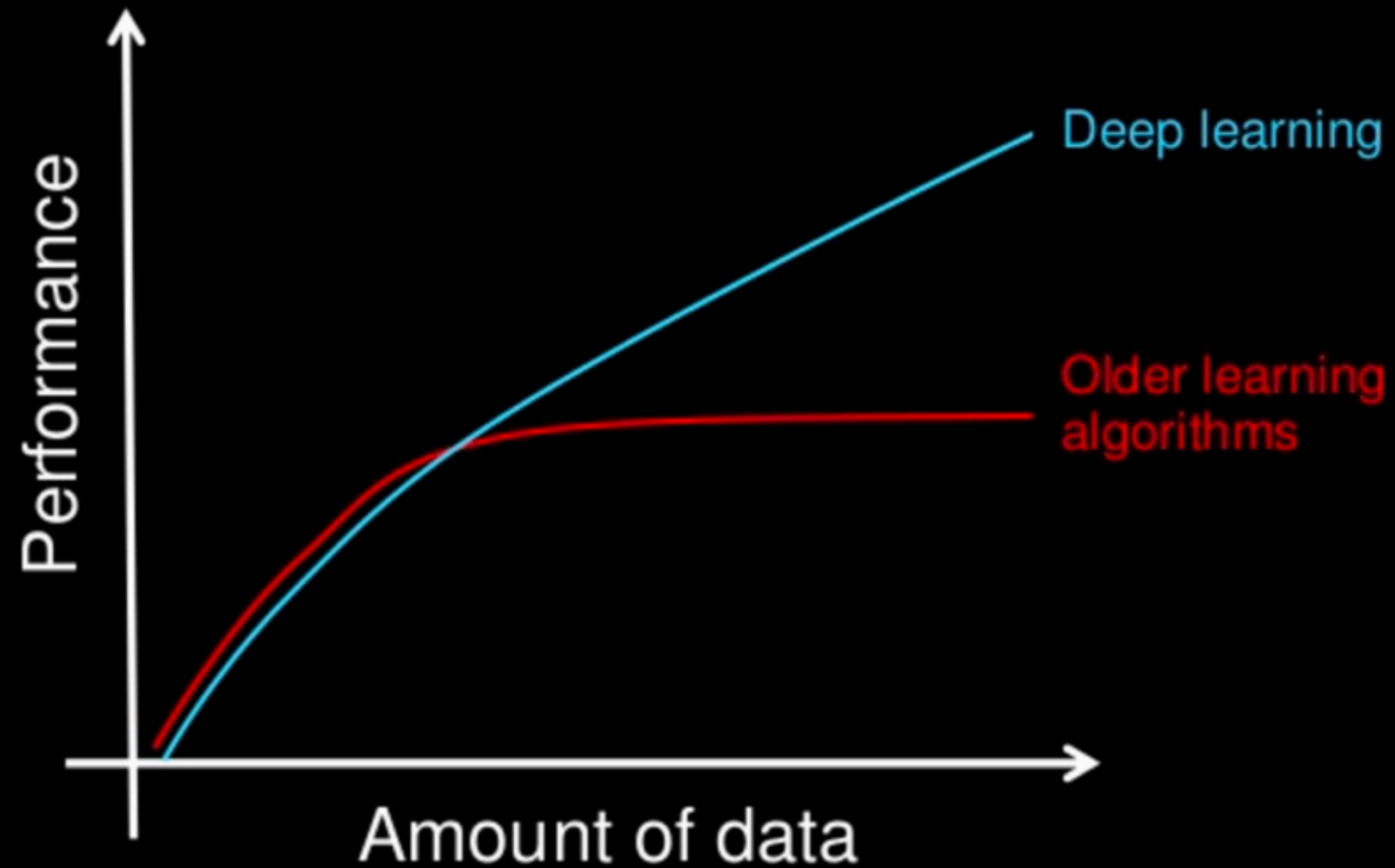


# Deep learning in 2012

Use a **hypothesis space** that can model complex structure  
(e.g., a CNN, nearest-neighbor)



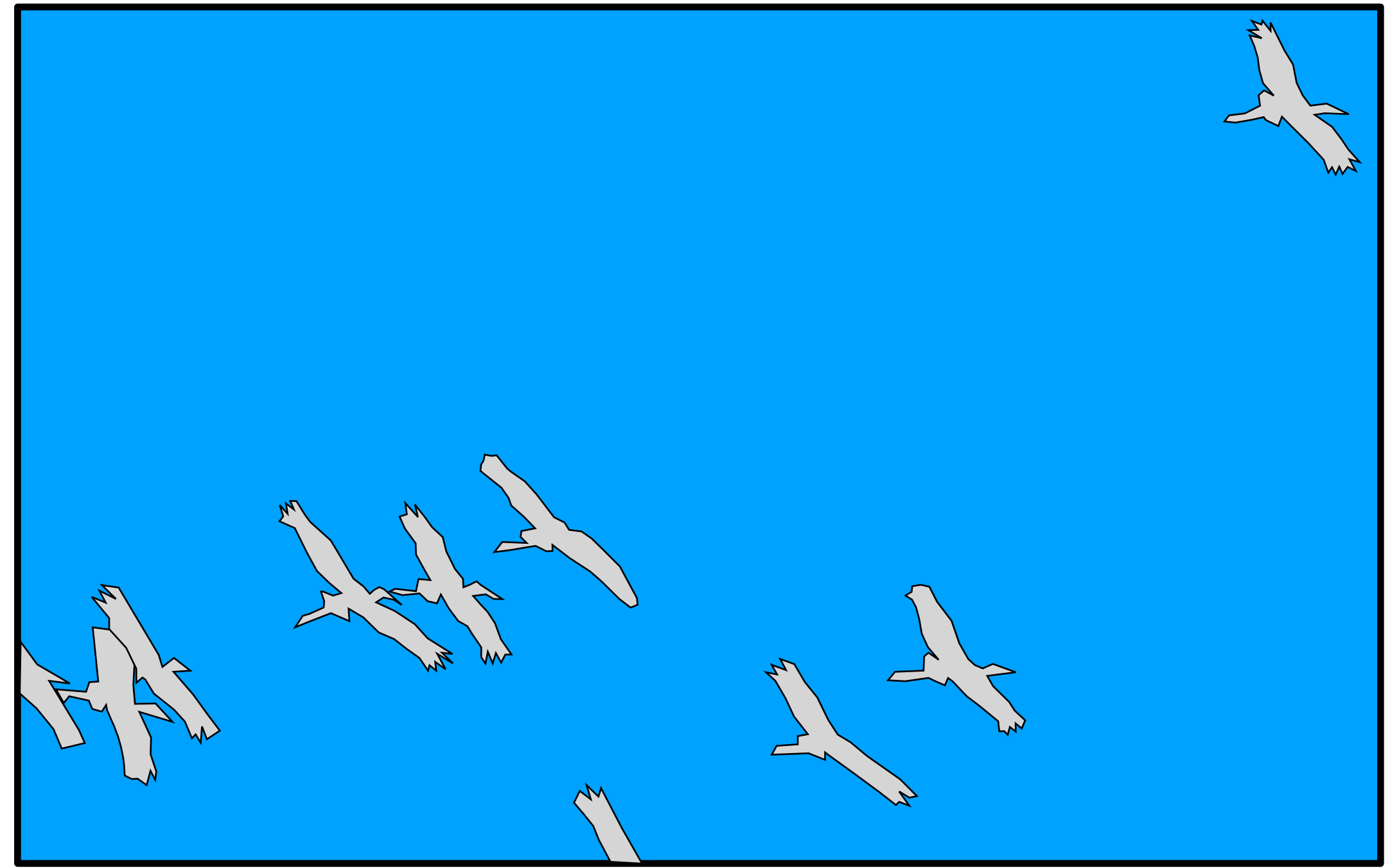
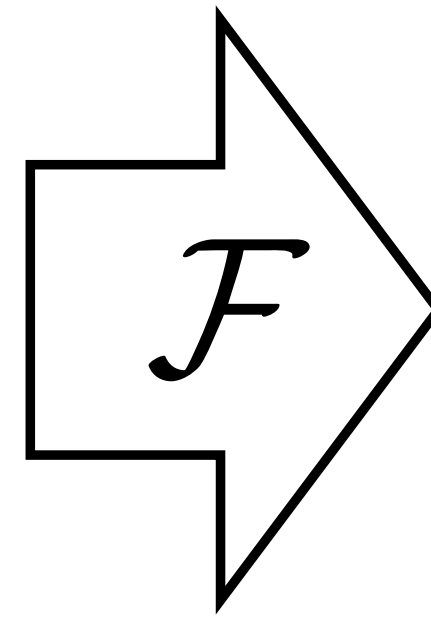
# Why deep learning



How do data science techniques scale with amount of data?



[Photo credit: Fredo Durand]



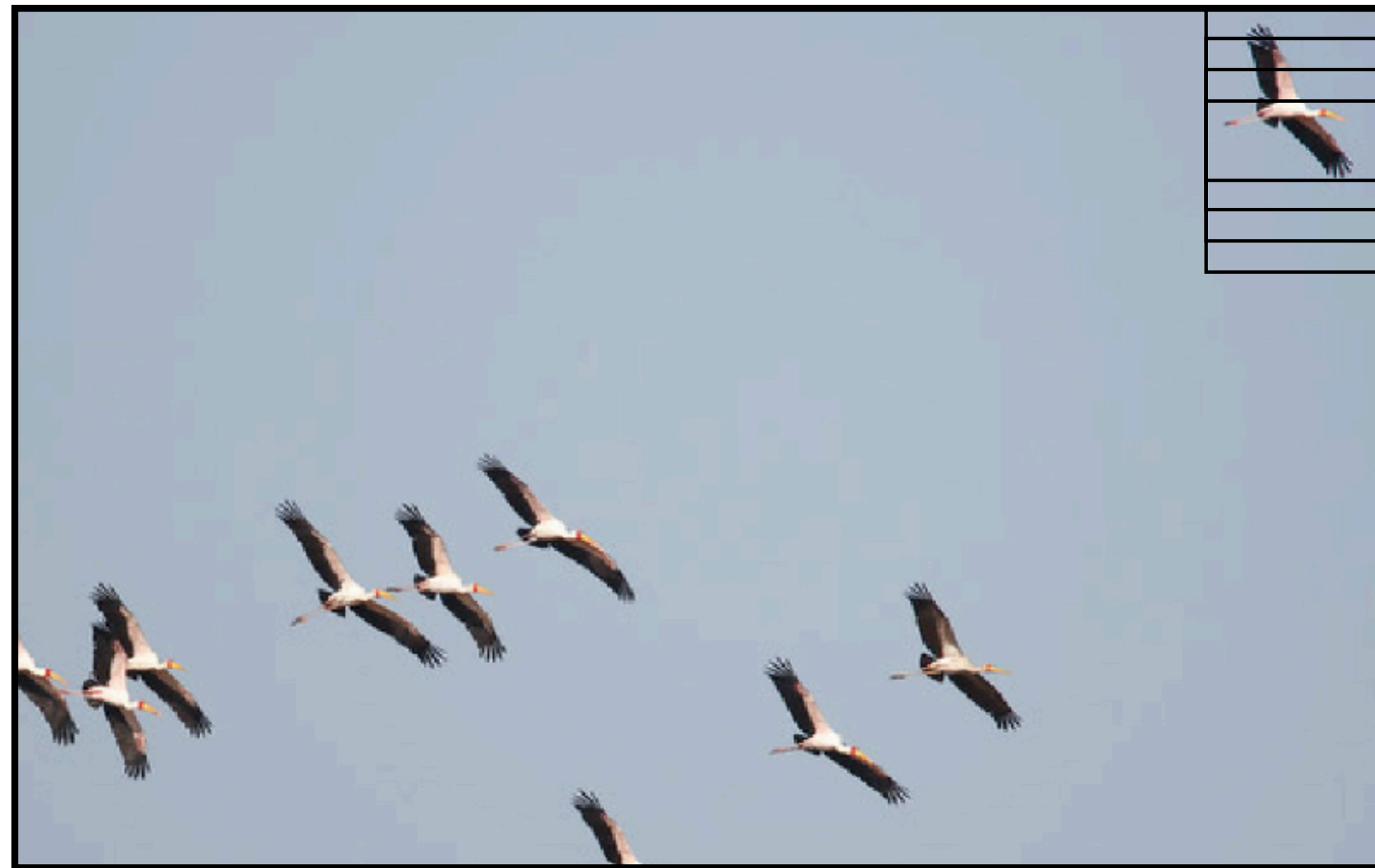
(Colors represent one-hot codes)

$$\arg \min_{\mathcal{F}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [L(\mathcal{F}(\mathbf{x}), \mathbf{y})]$$

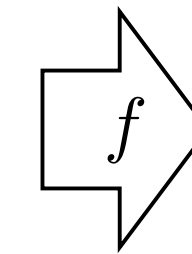
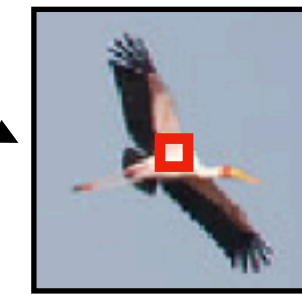
Hypothesis space

Objective function  
(loss)

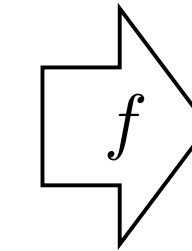




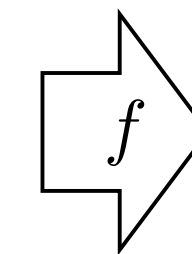
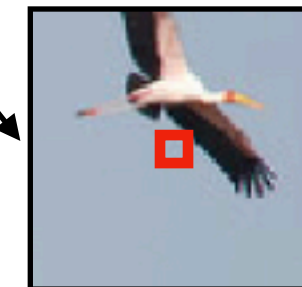
What's the object class of the center pixel?



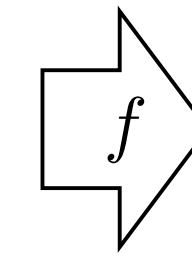
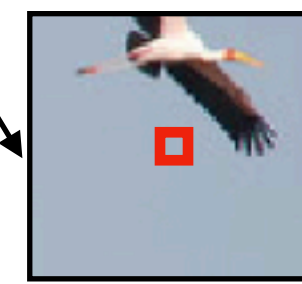
“Bird”



“Bird”



“Sky”

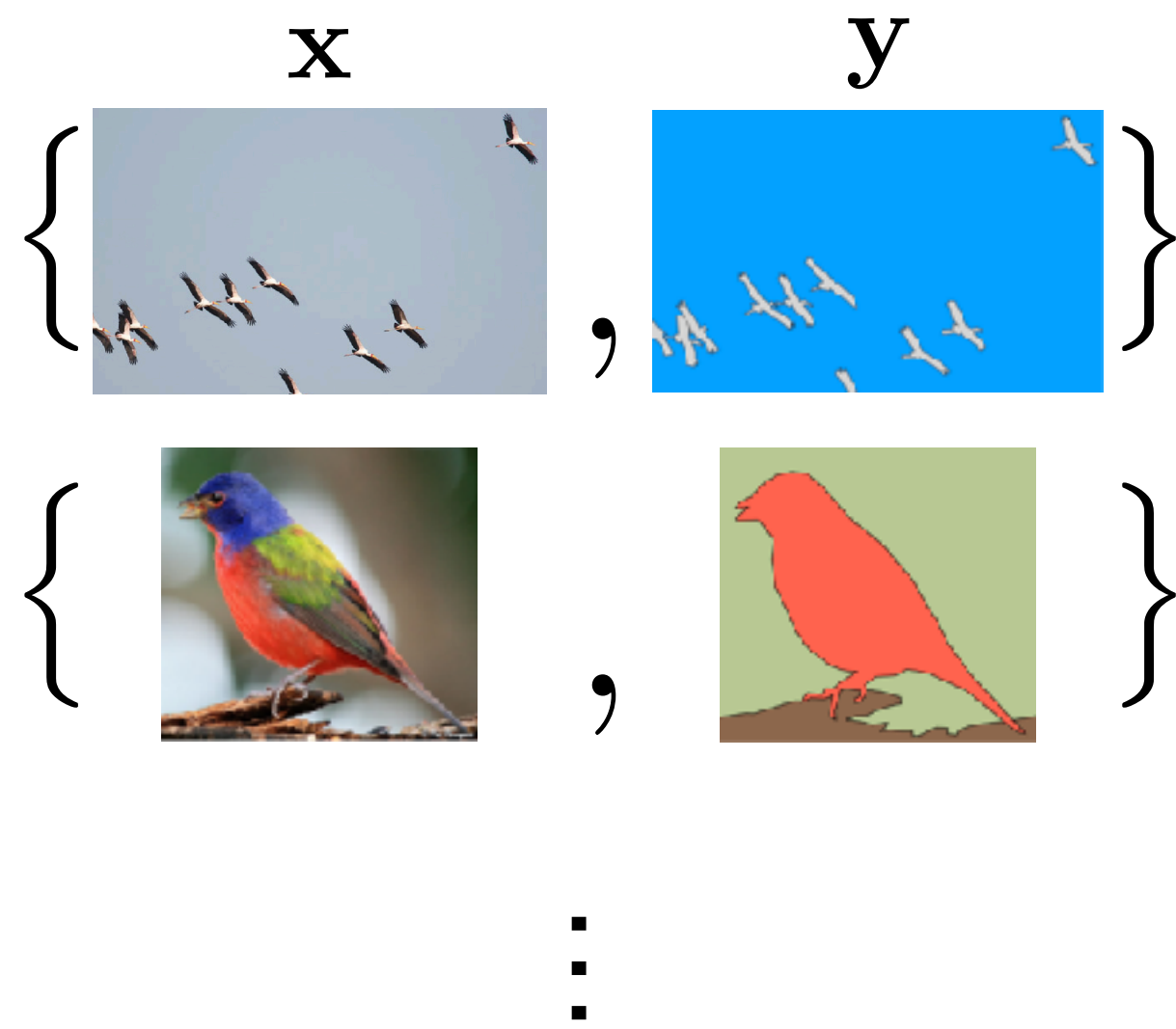


“Sky”

Fully-factored loss: 
$$L(\hat{\mathbf{y}}, \mathbf{y}) = \sum_i \phi_i(\hat{\mathbf{y}}_i, \mathbf{y}_i)$$

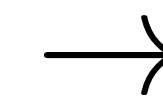
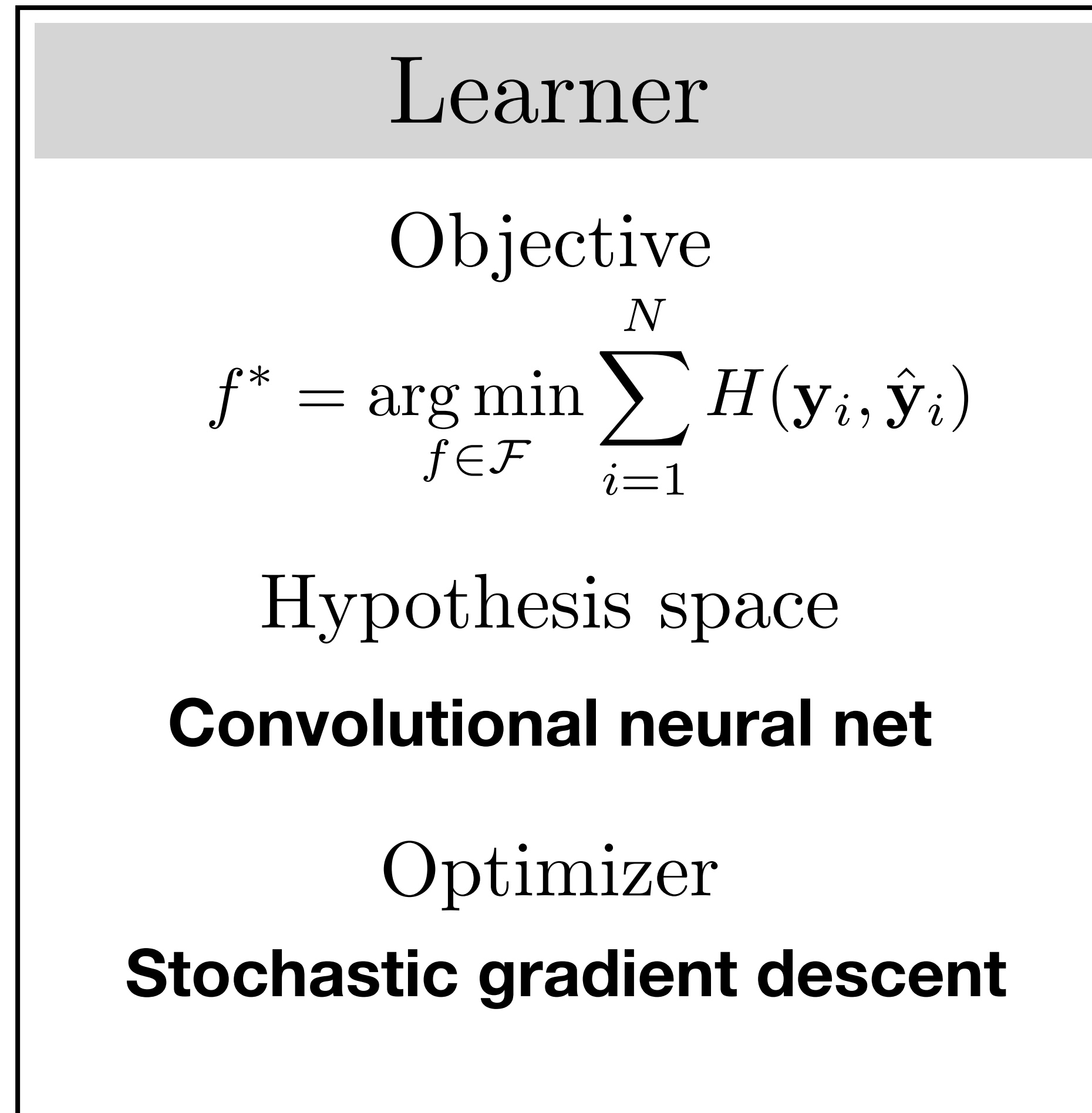
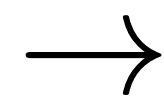
# Semantic Segmentation

Data



$$\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$$

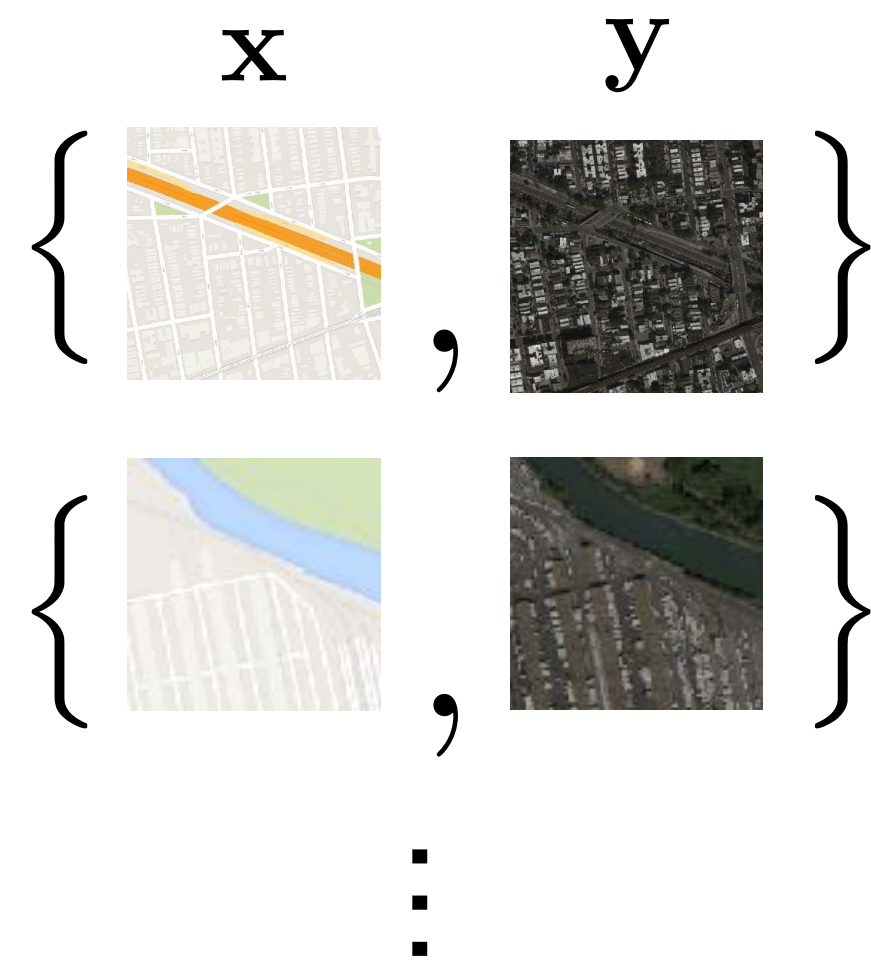
$$\mathbf{y} \in \mathbb{R}^{H \times W \times K}$$



*f*

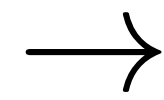
# Sat2Map

Data



$$\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$$

$$\mathbf{y} \in \mathbb{R}^{H \times W \times 3}$$



**Learner**

Objective

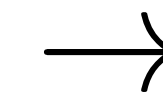
$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (f_{\theta}(\mathbf{x})_i - y_i)^2$$

Hypothesis space

**Convolutional neural net**

Optimizer

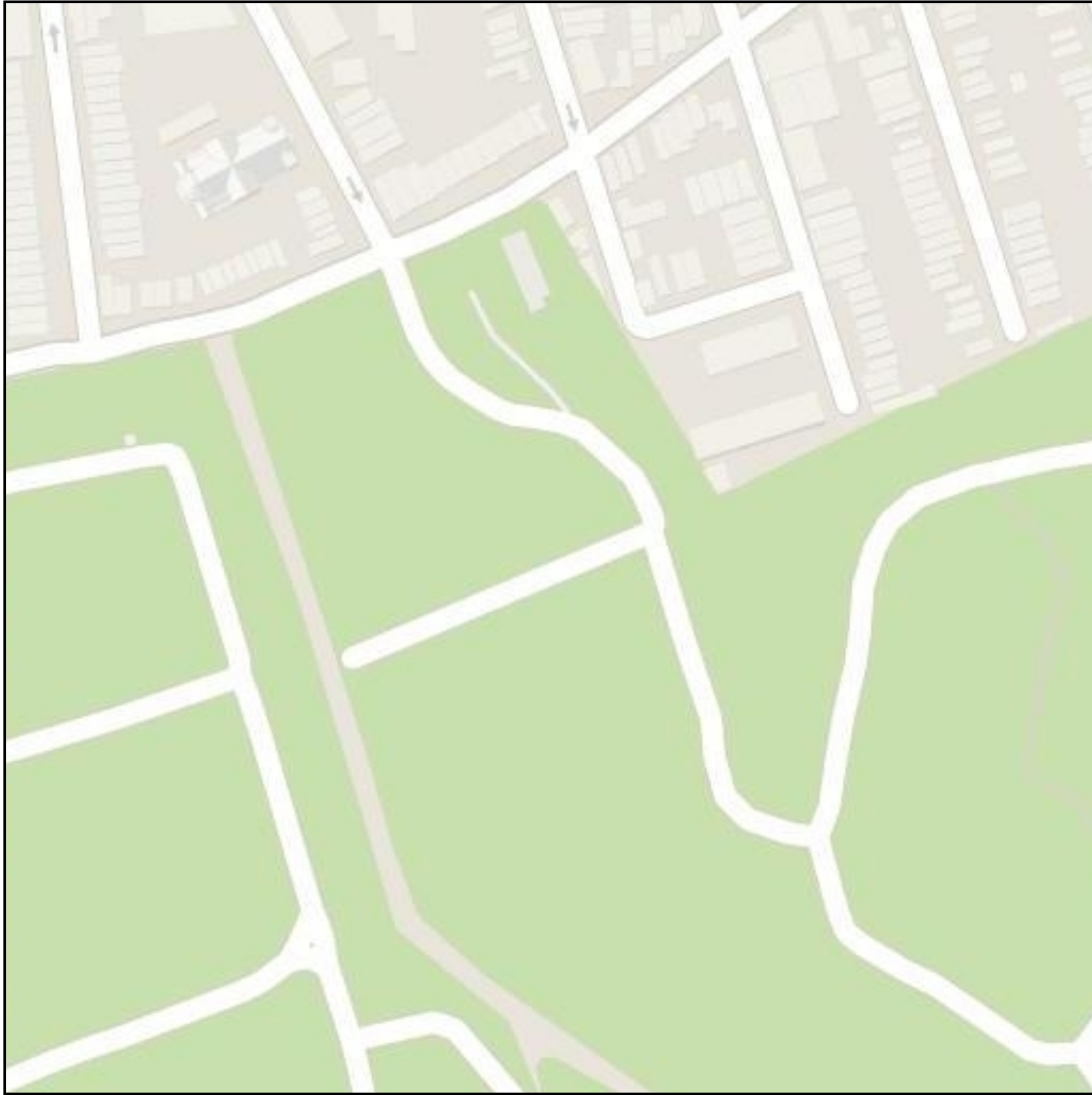
**Stochastic gradient descent**



$f$



Input



Deep net output

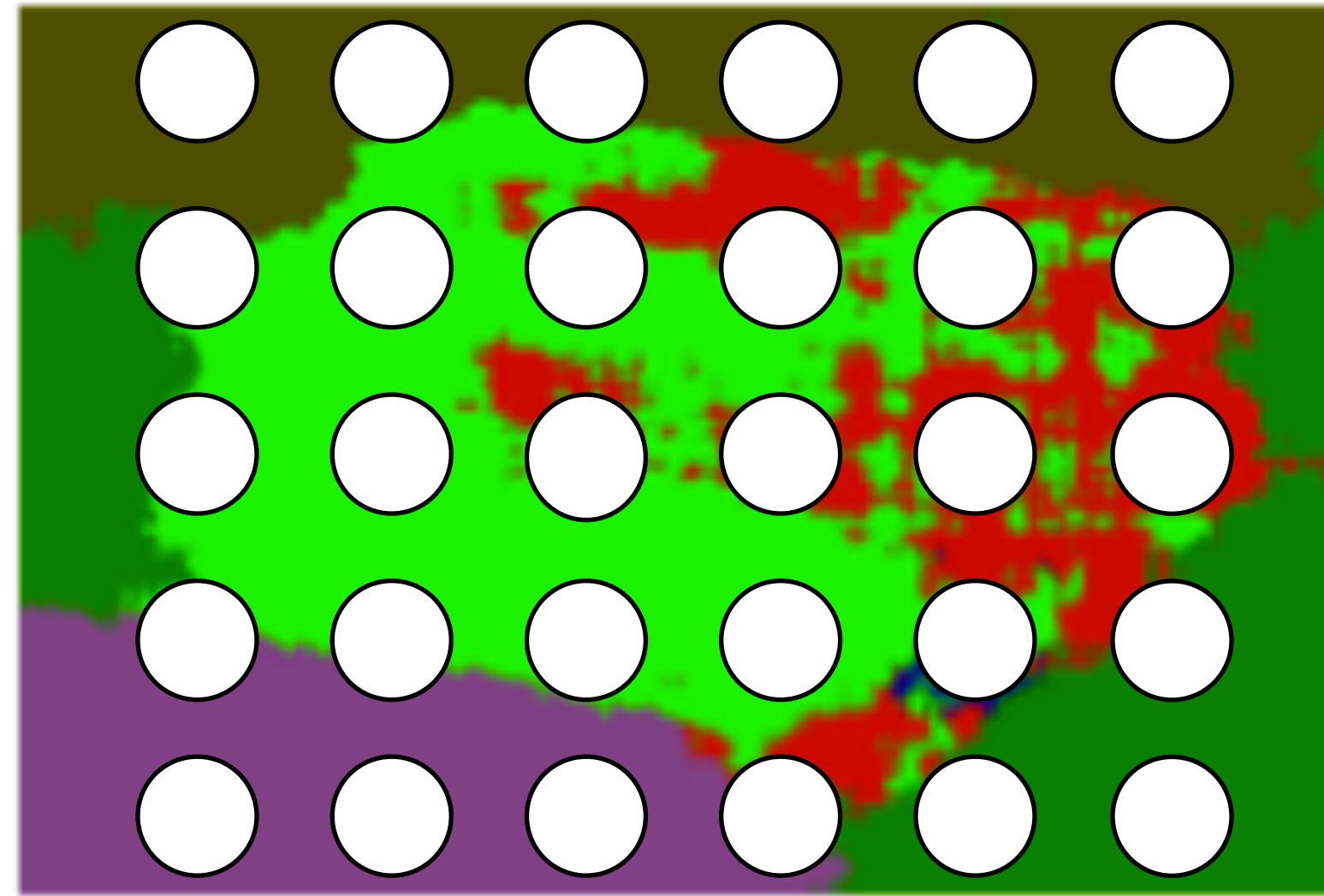




Input

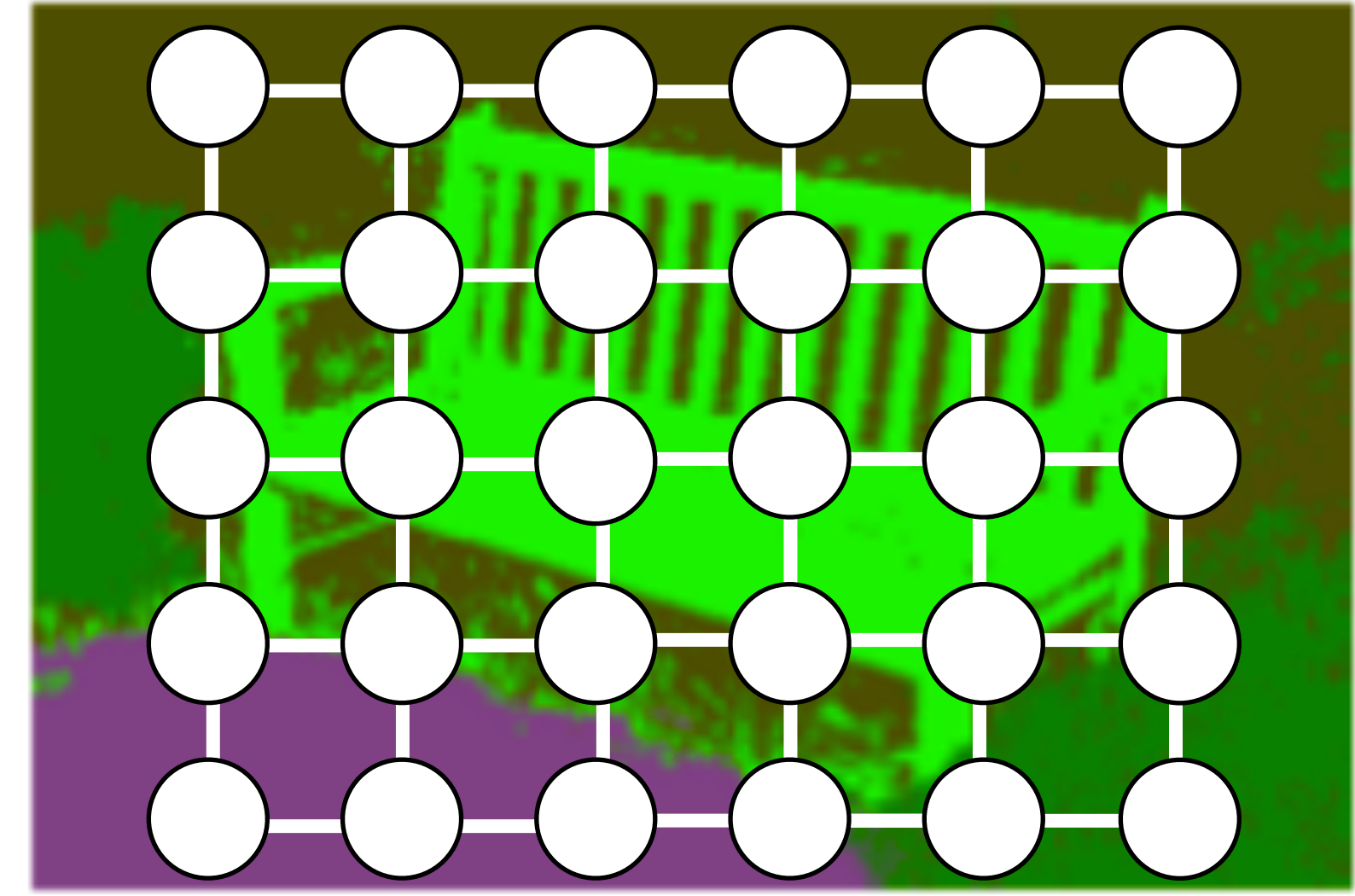


Independent prediction  
per-pixel



$$\max \prod_i p(y_i | \mathbf{x})$$

Find a configuration of  
compatible labels

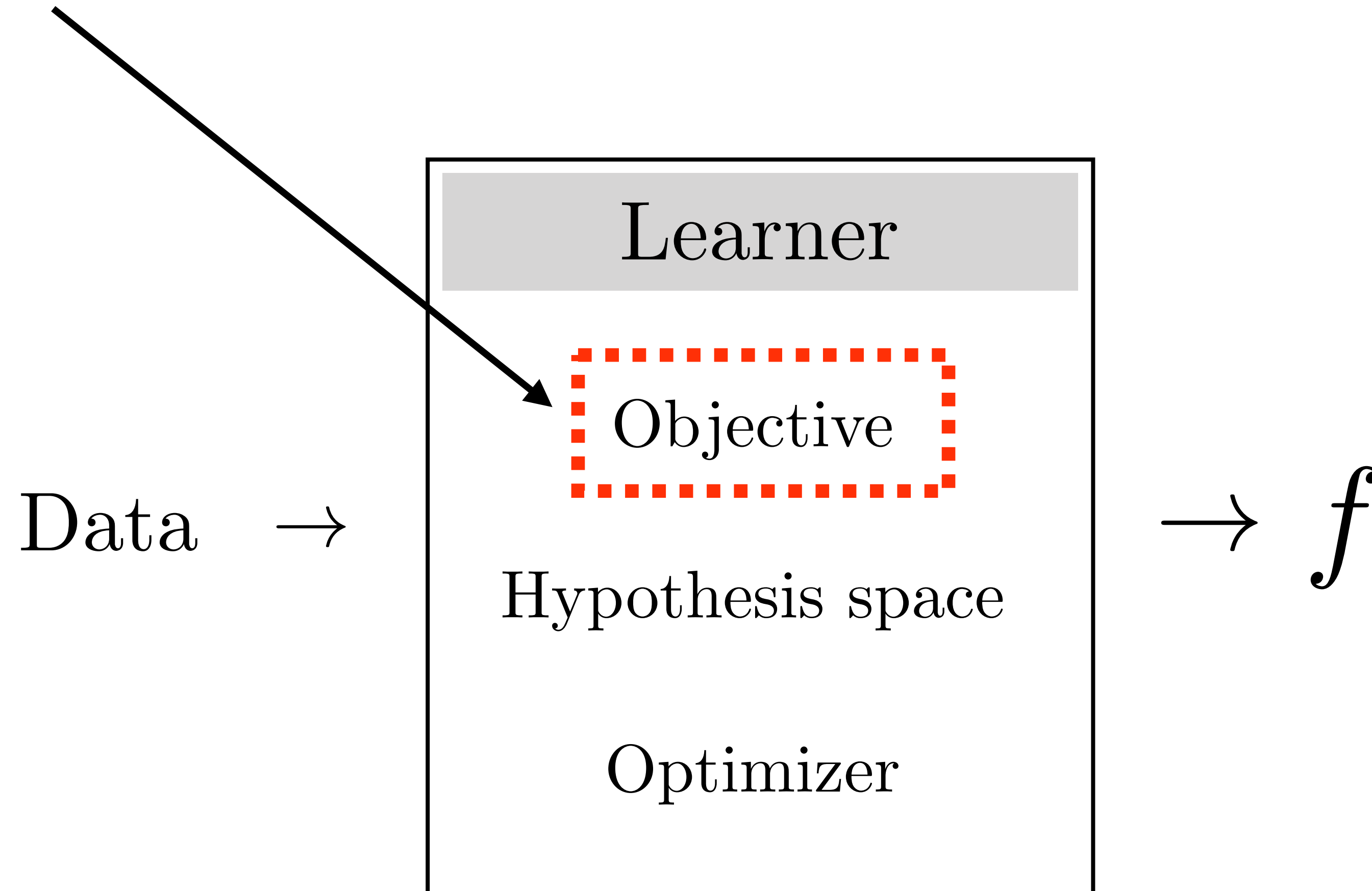


$$\max \frac{1}{Z} \prod_{i,j} p(y_i, y_j | \mathbf{x})$$

[“Efficient Inference in Fully Connected CRFs with Gaussian Edge Potentials”, Krahenbuhl and Koltun, NeurIPS 2011]

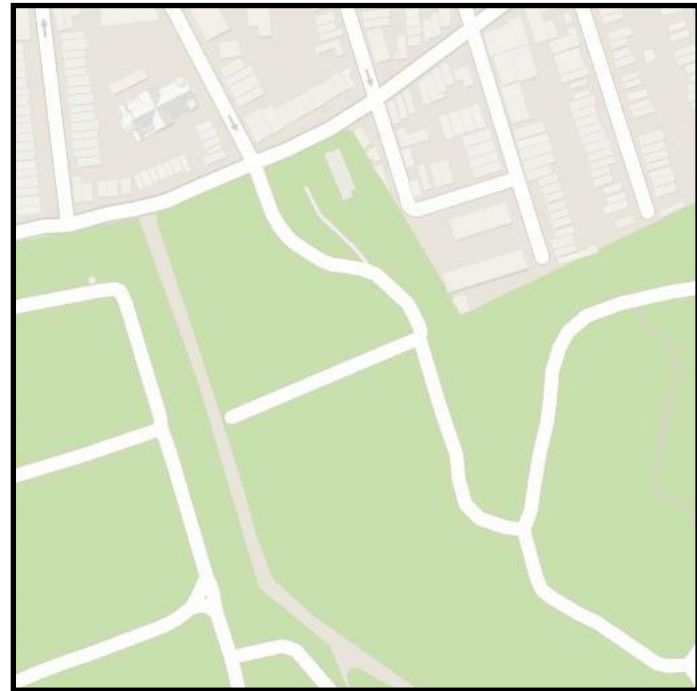
# Structured prediction

Use an **objective** that can model structure! (e.g., a graphical model, a GAN, etc)

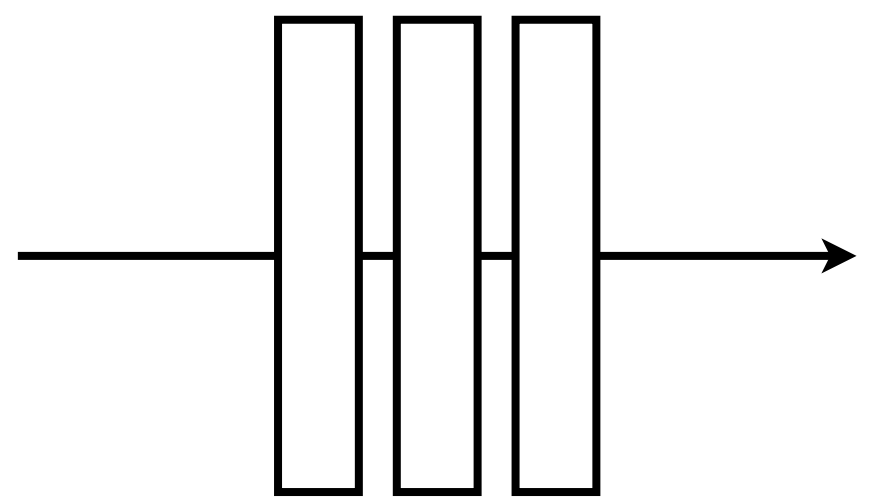




$\mathbf{x}$



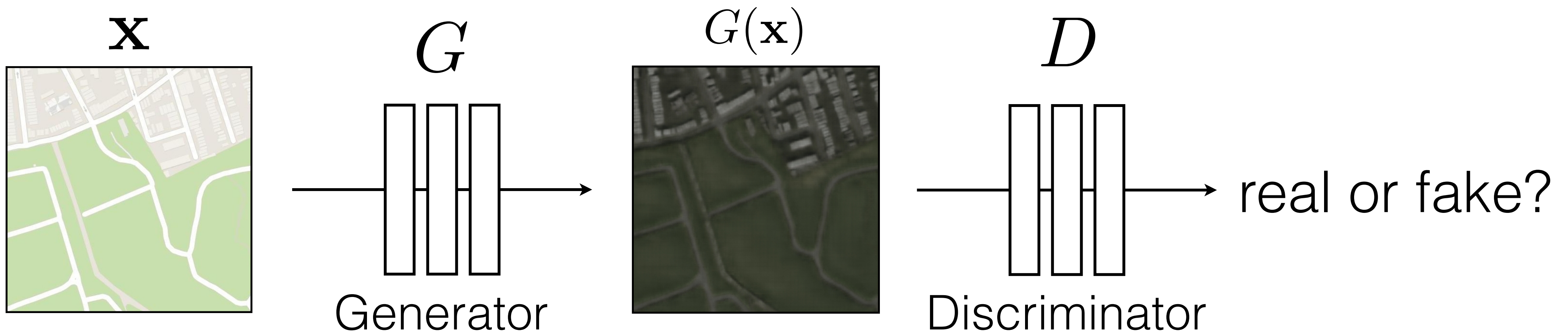
$G$



Generator

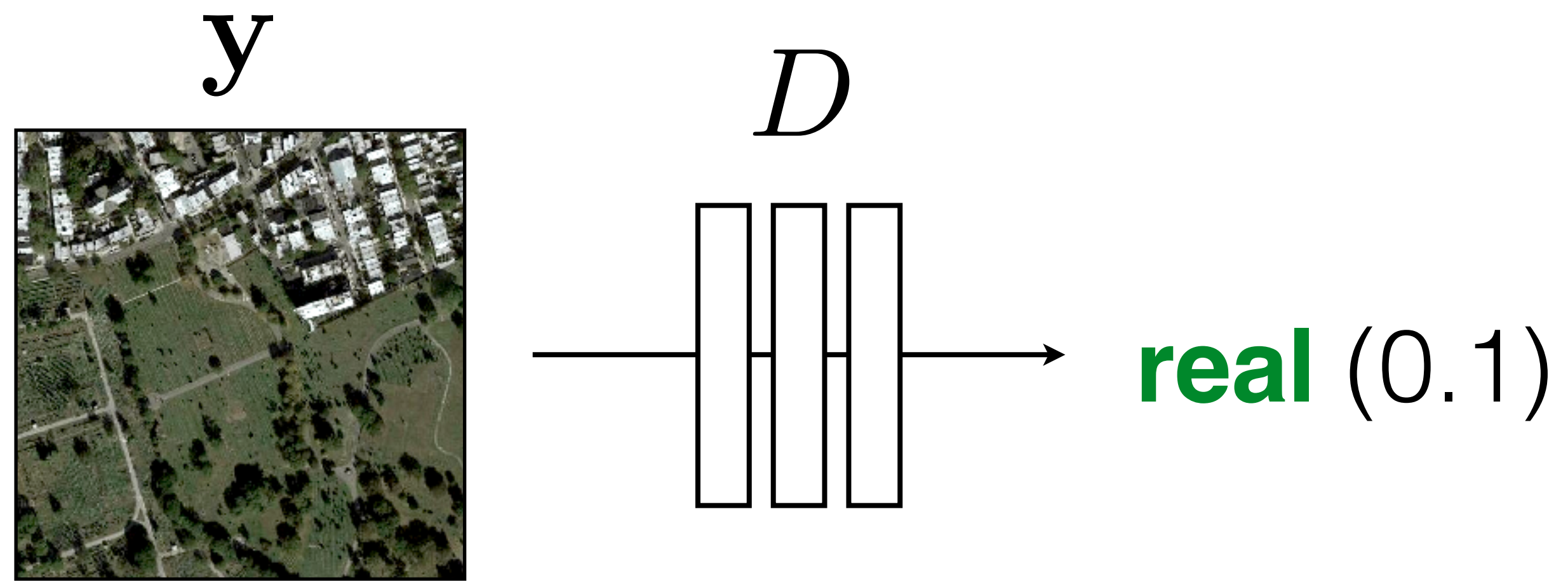
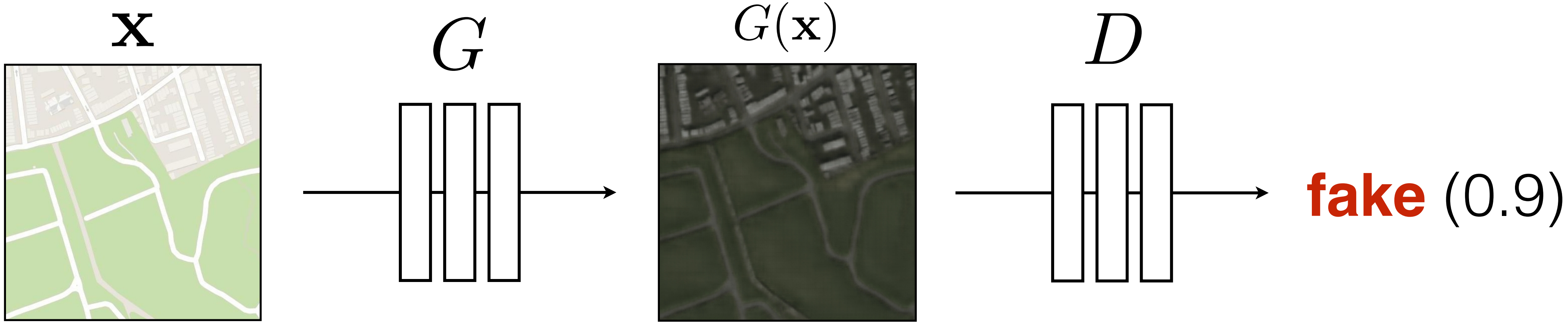
$G(\mathbf{x})$





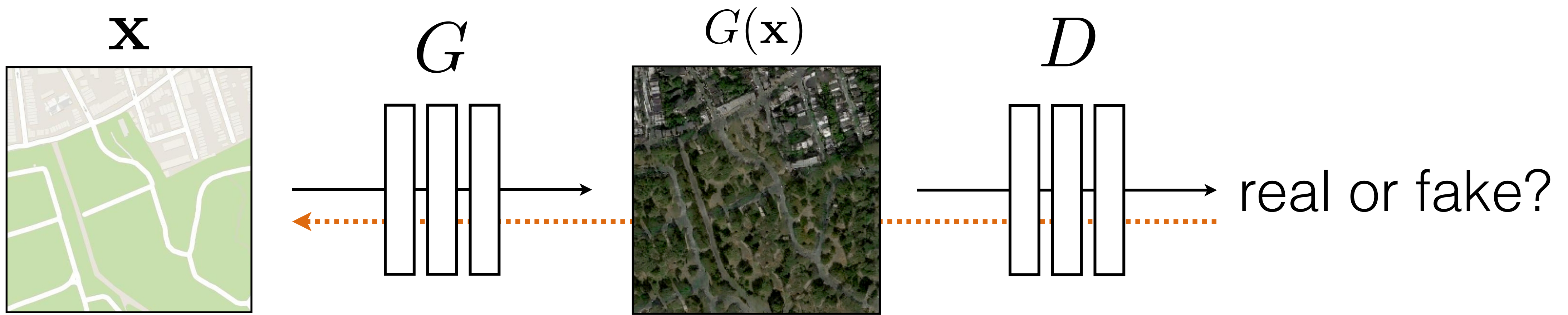
**G** tries to synthesize fake images that fool **D**

**D** tries to identify the fakes



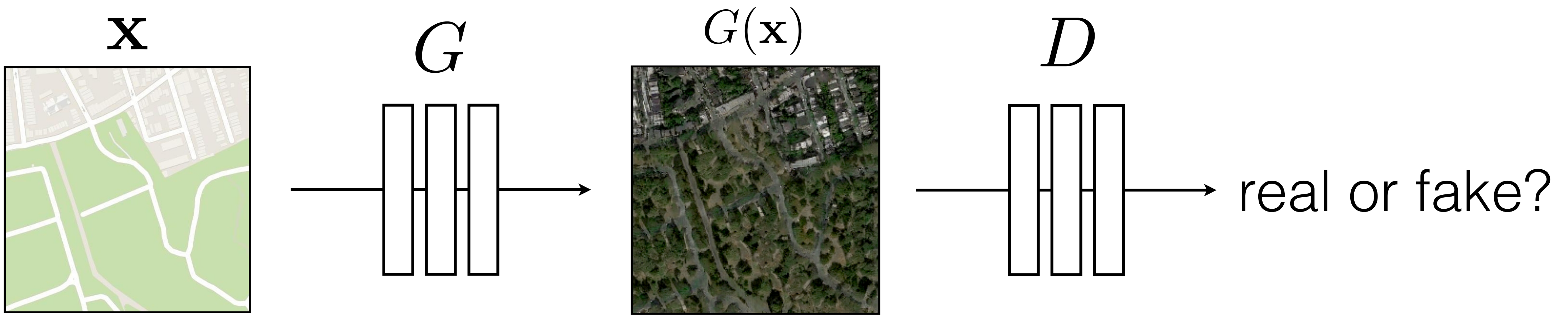
$$\arg \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} \left[ \boxed{\log D(G(\mathbf{x}))} + \boxed{\log(1 - D(\mathbf{y}))} \right]$$





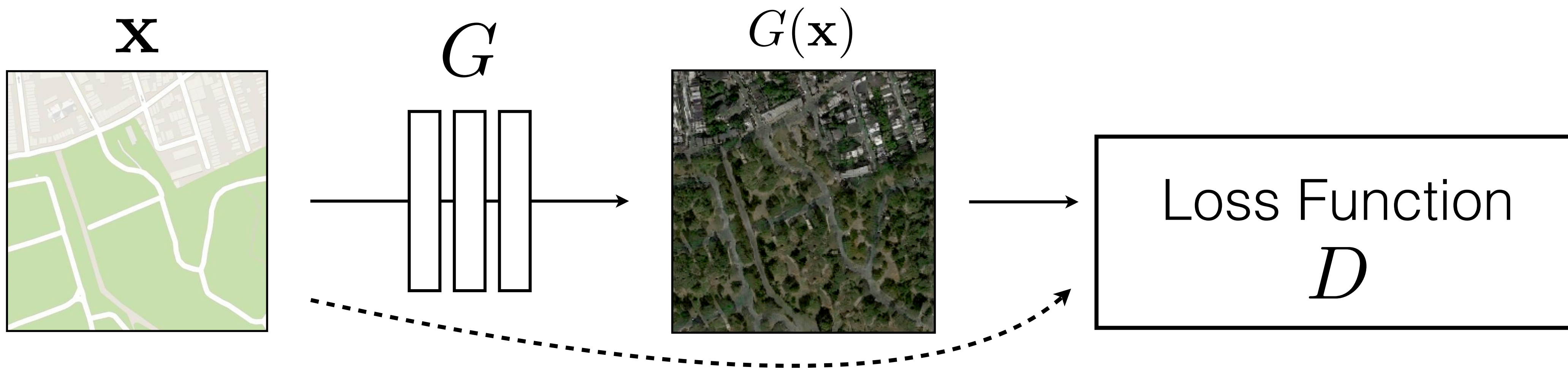
**G** tries to synthesize fake images that *fool* **D**:

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



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

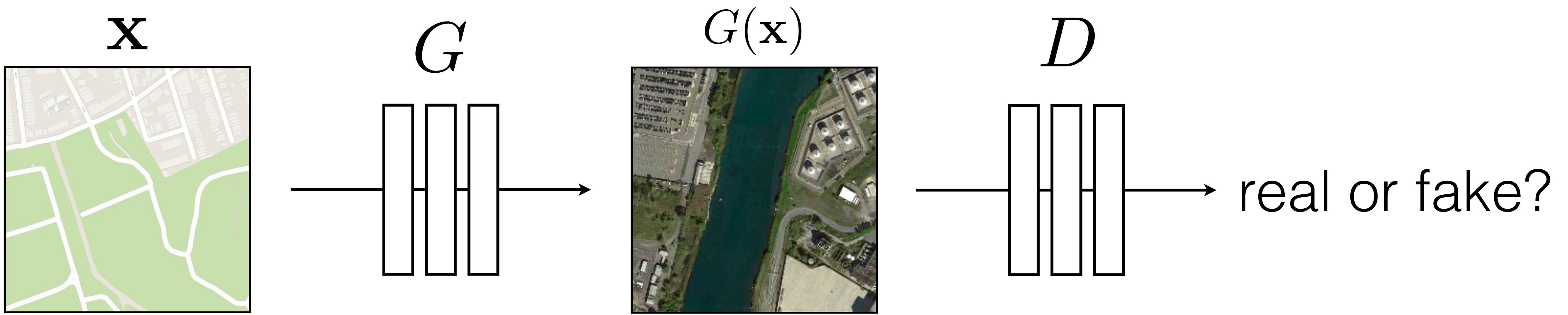
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



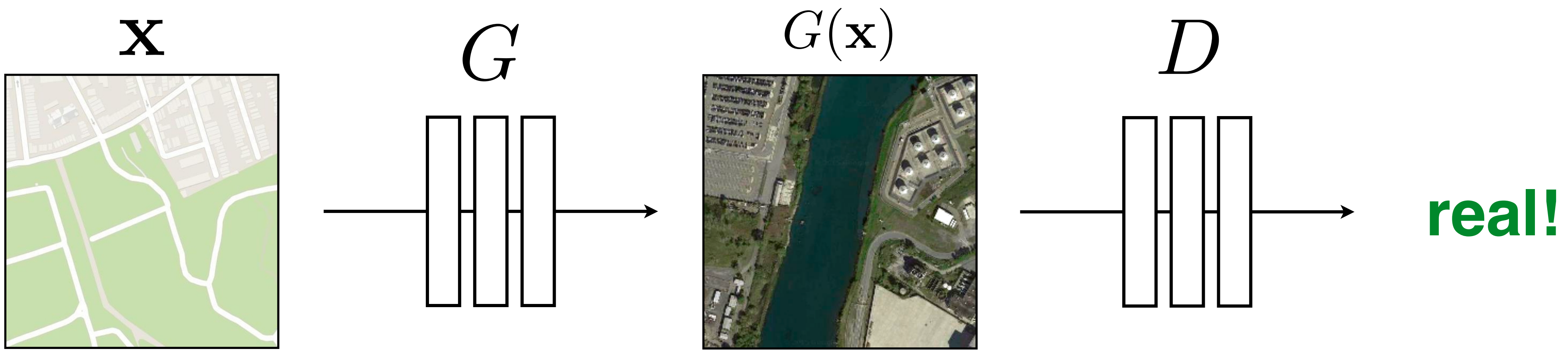
**G**'s perspective: **D** is a loss function.

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

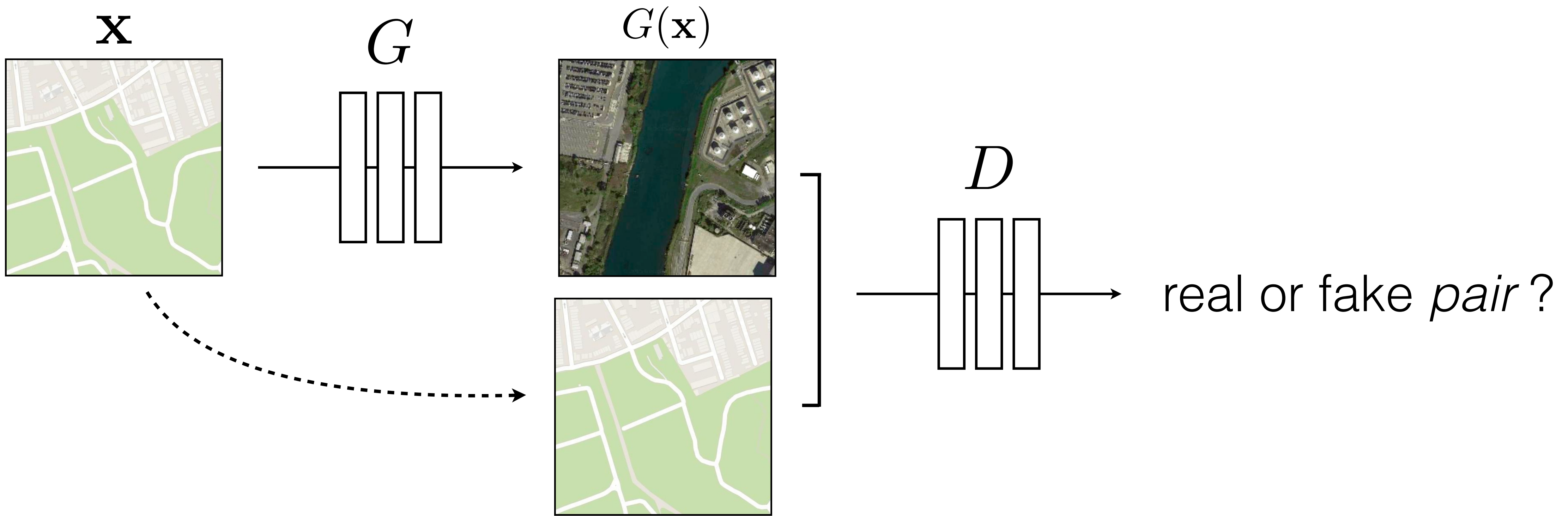




$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

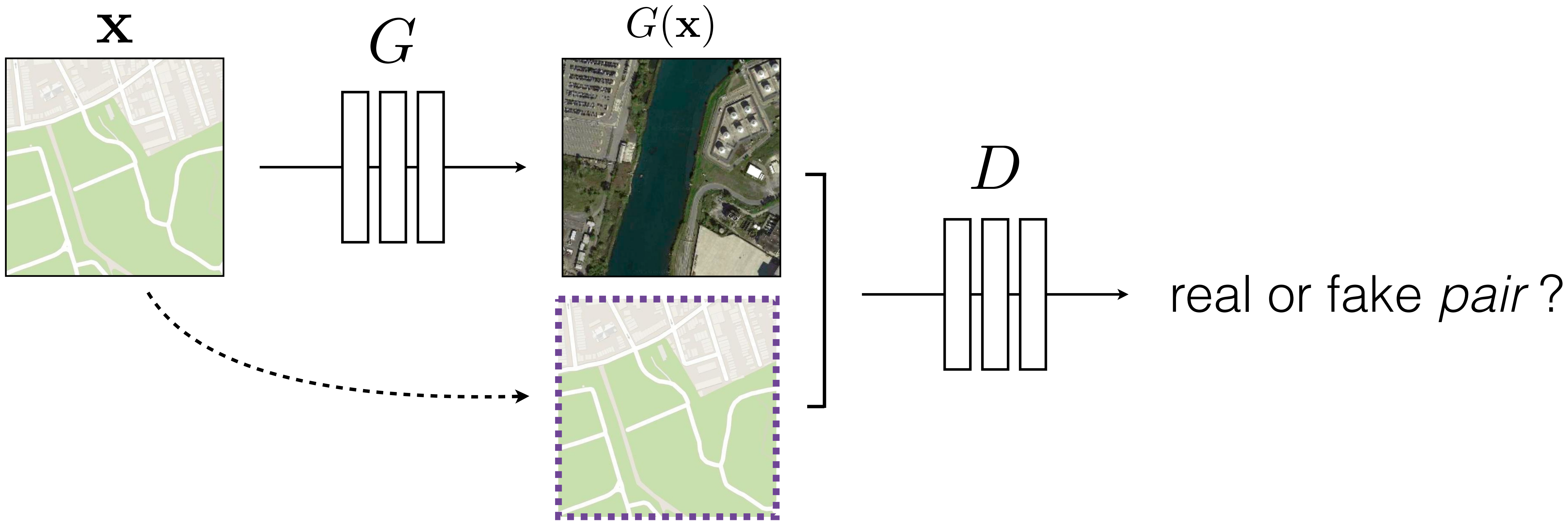


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$

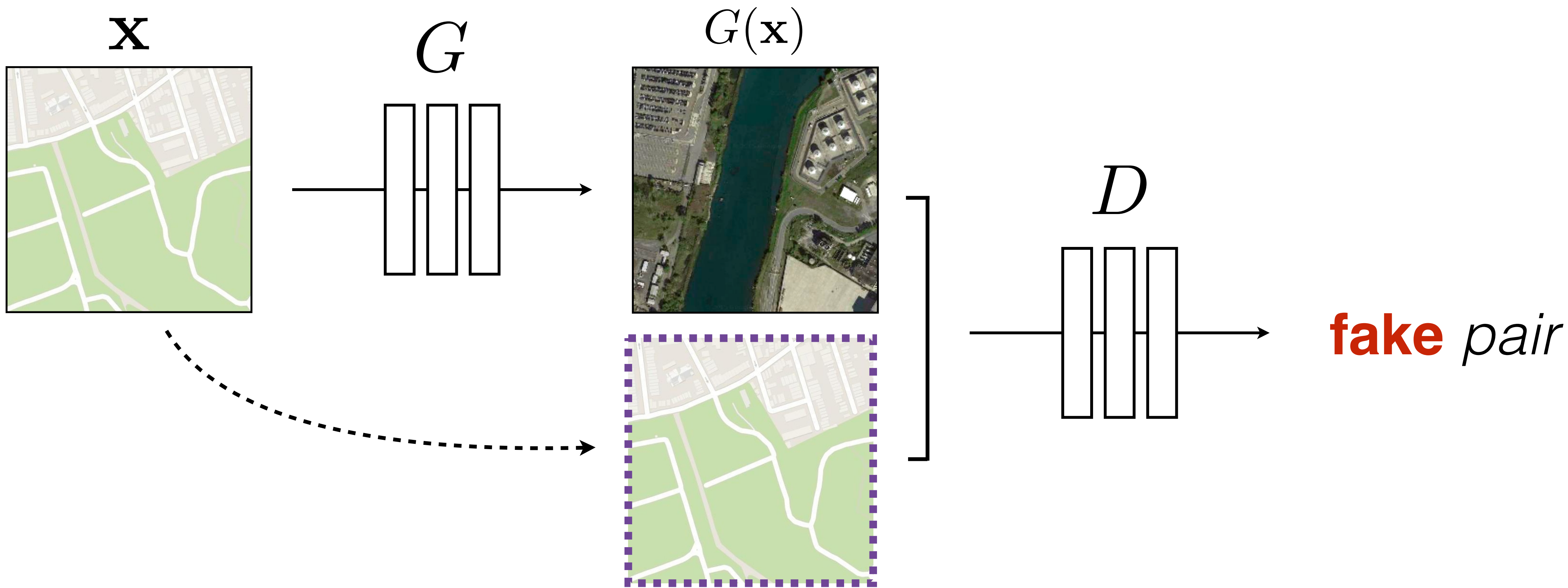


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y})) ]$$



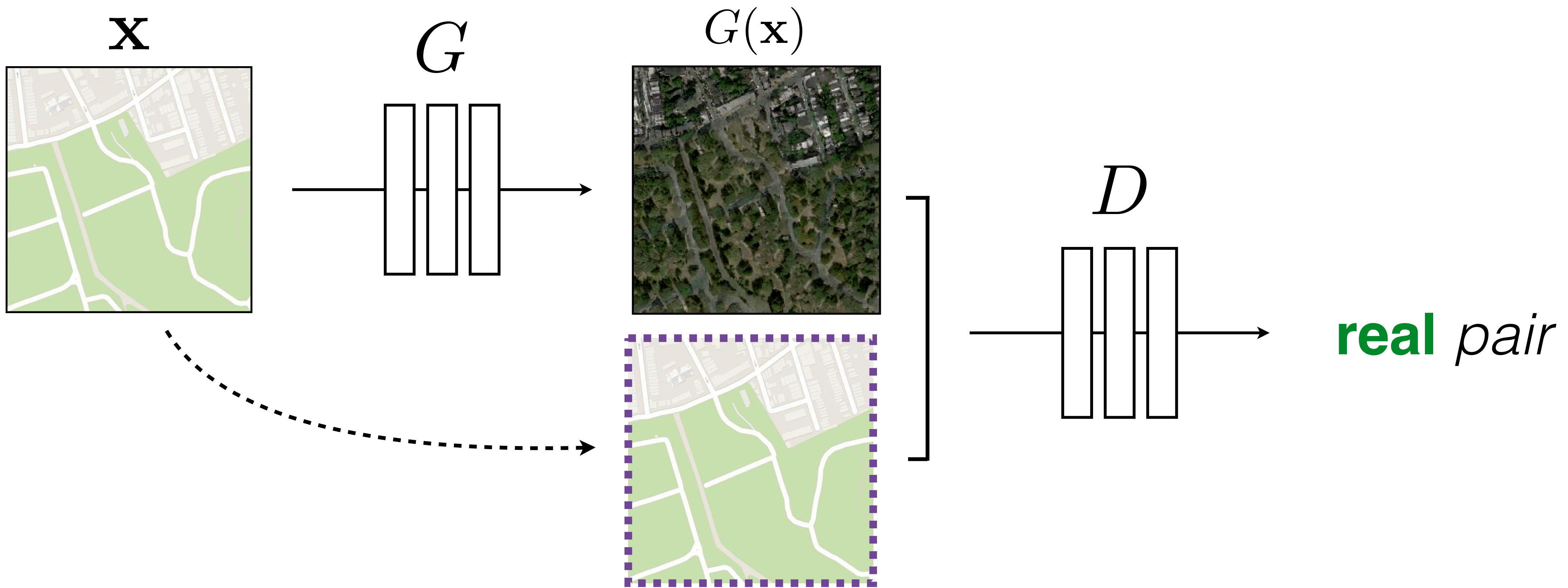


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$



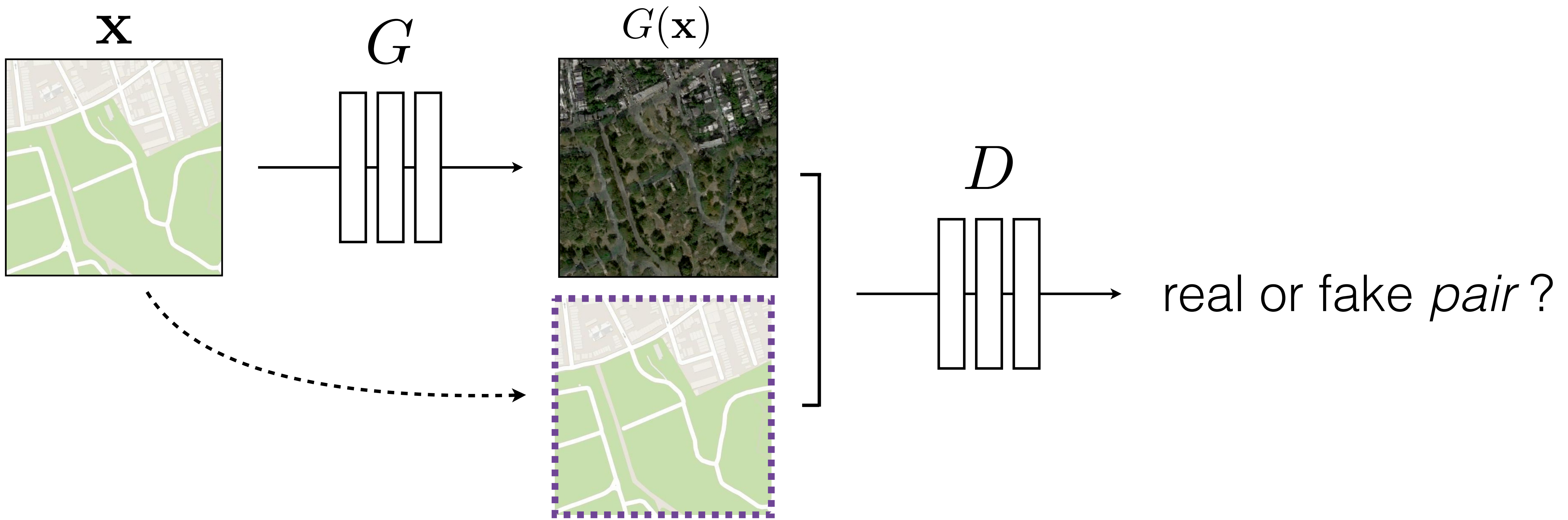
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$





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





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

# Training Details: Loss function

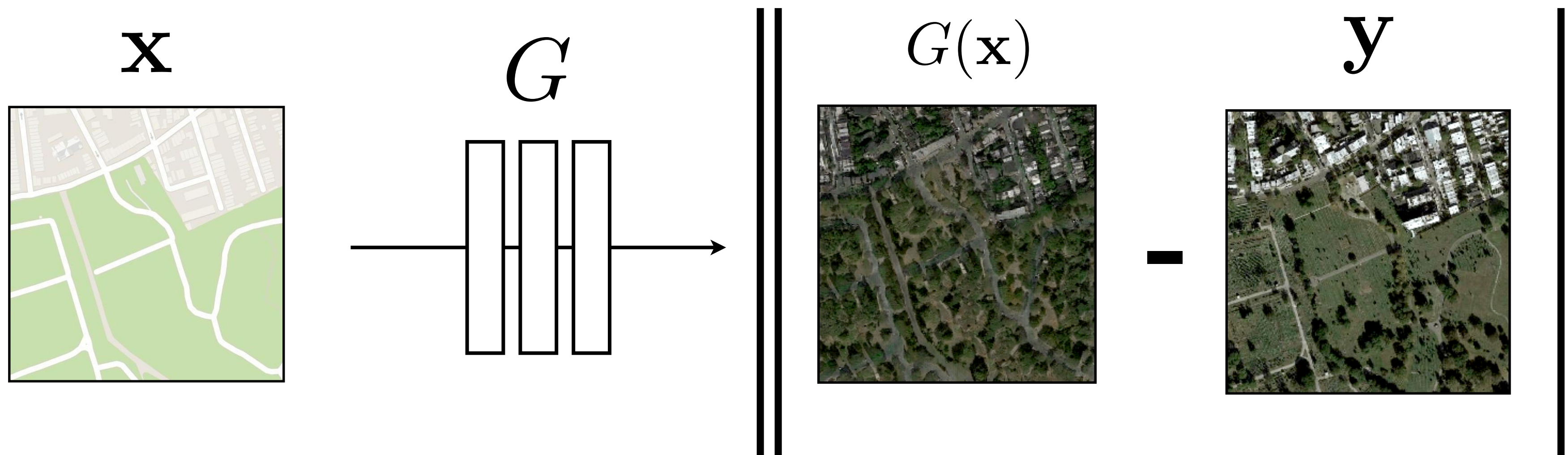
Conditional GAN

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

# 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

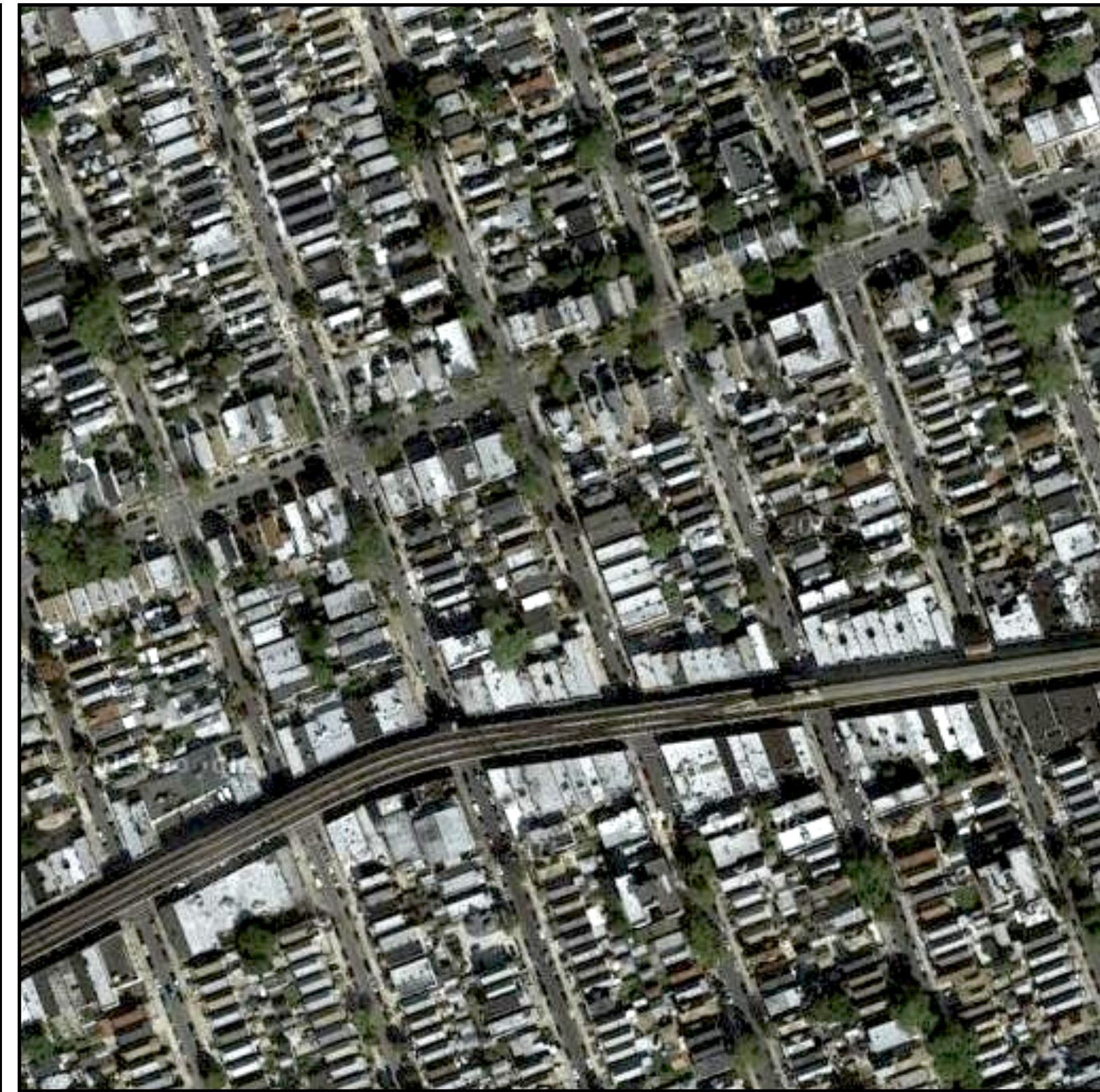
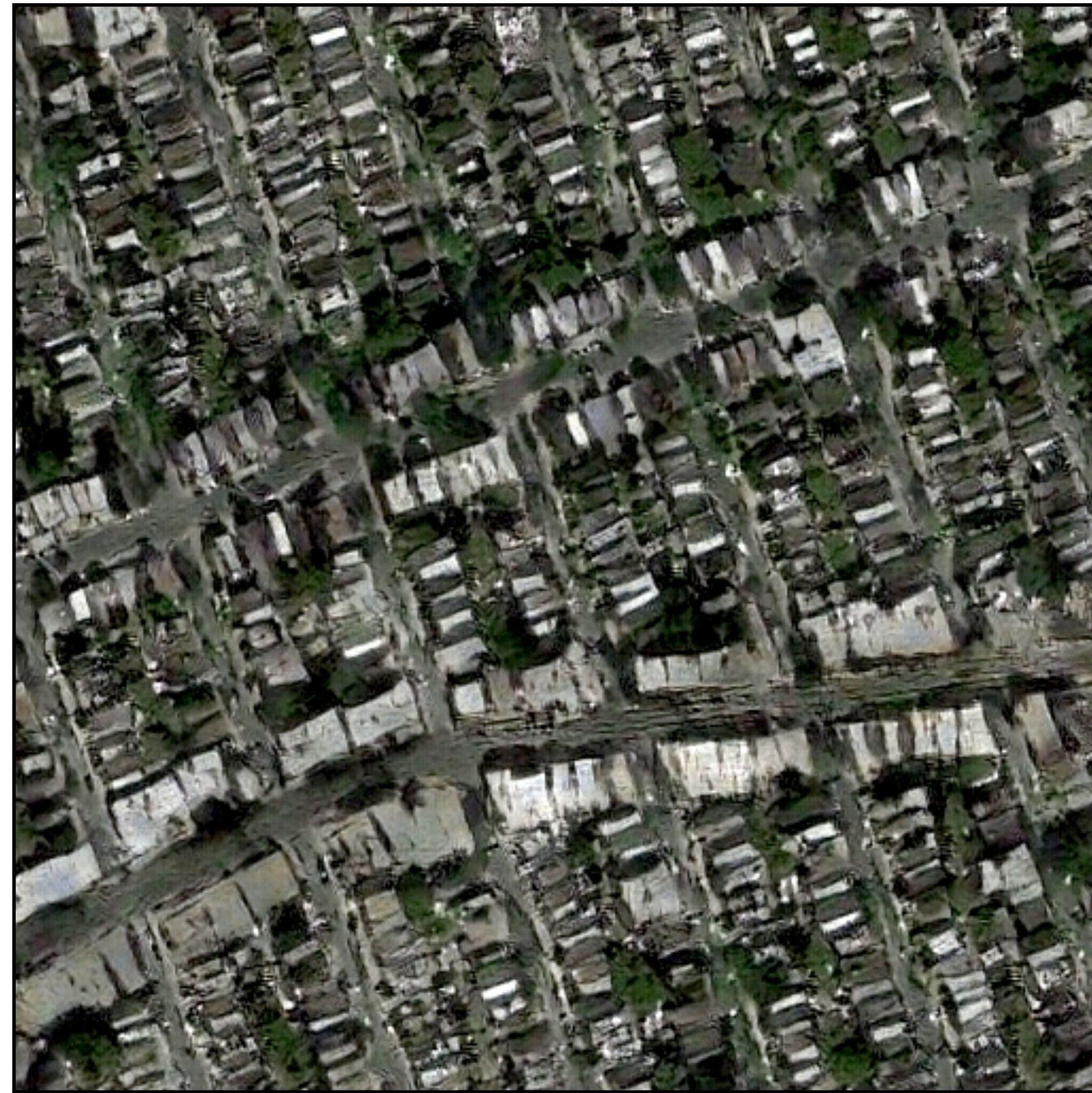
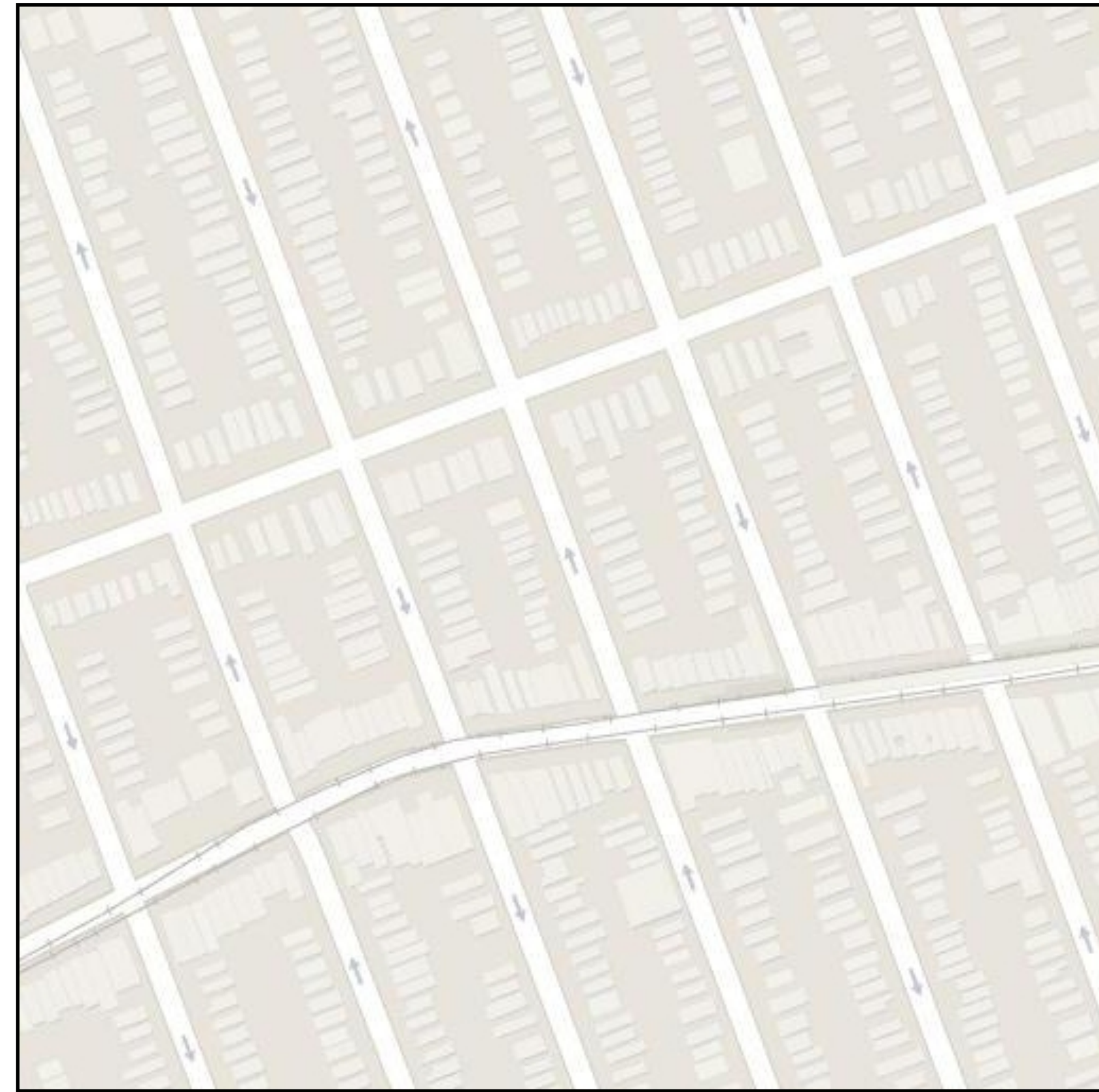
[c.f. Pathak et al. CVPR 2016]



Input

Output

Groundtruth



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

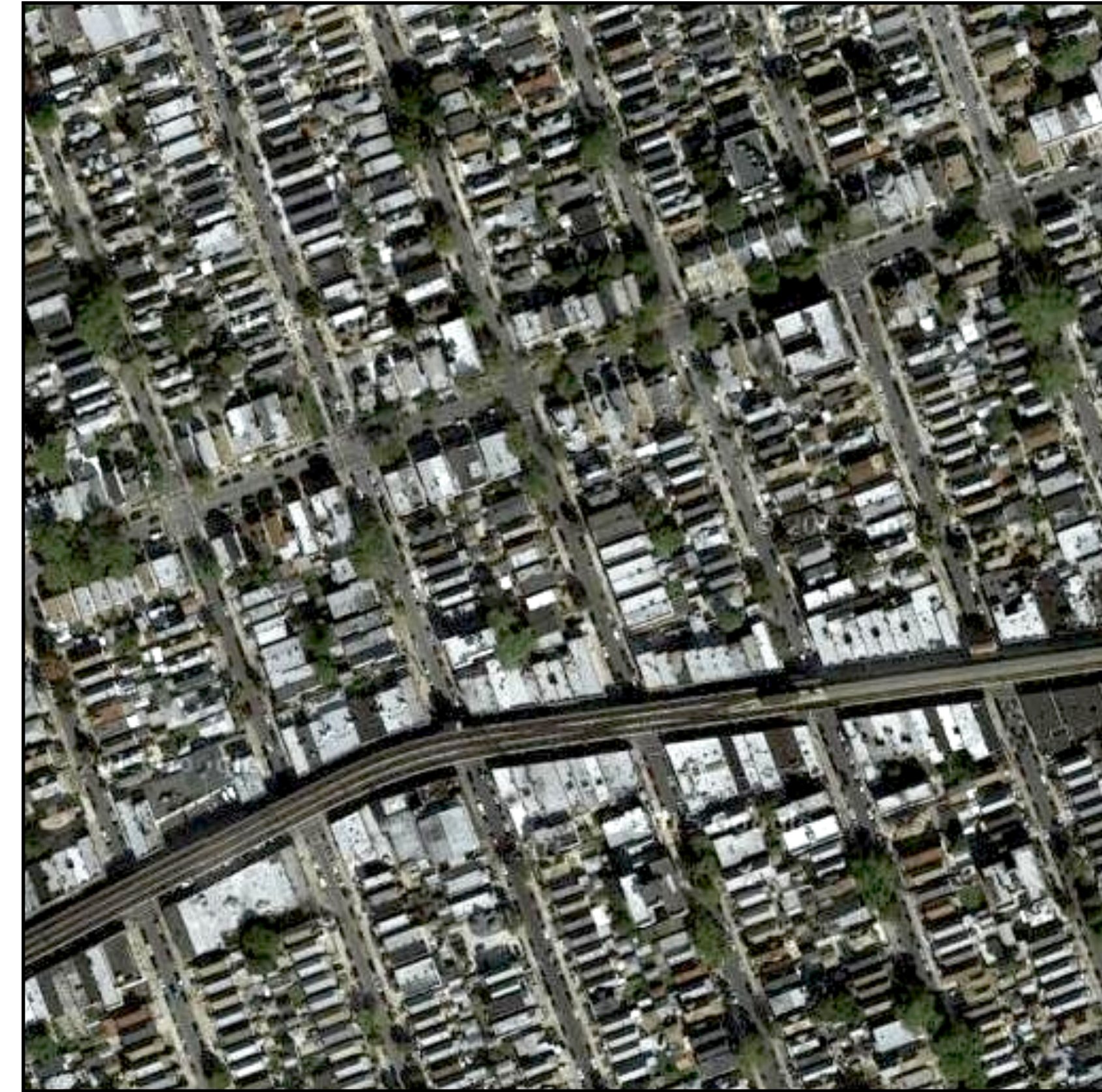




Input

Output

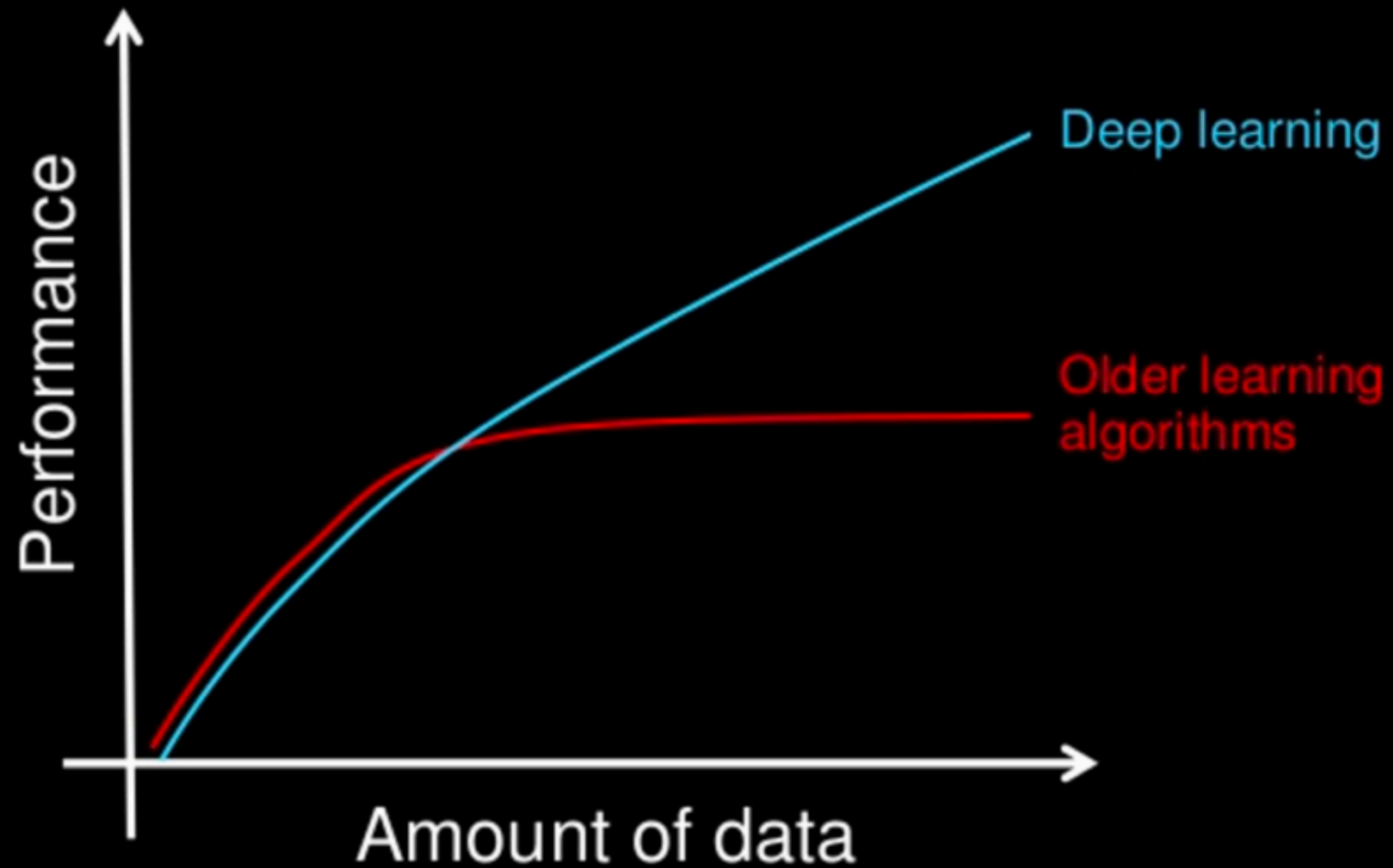
Groundtruth



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



# Why deep learning

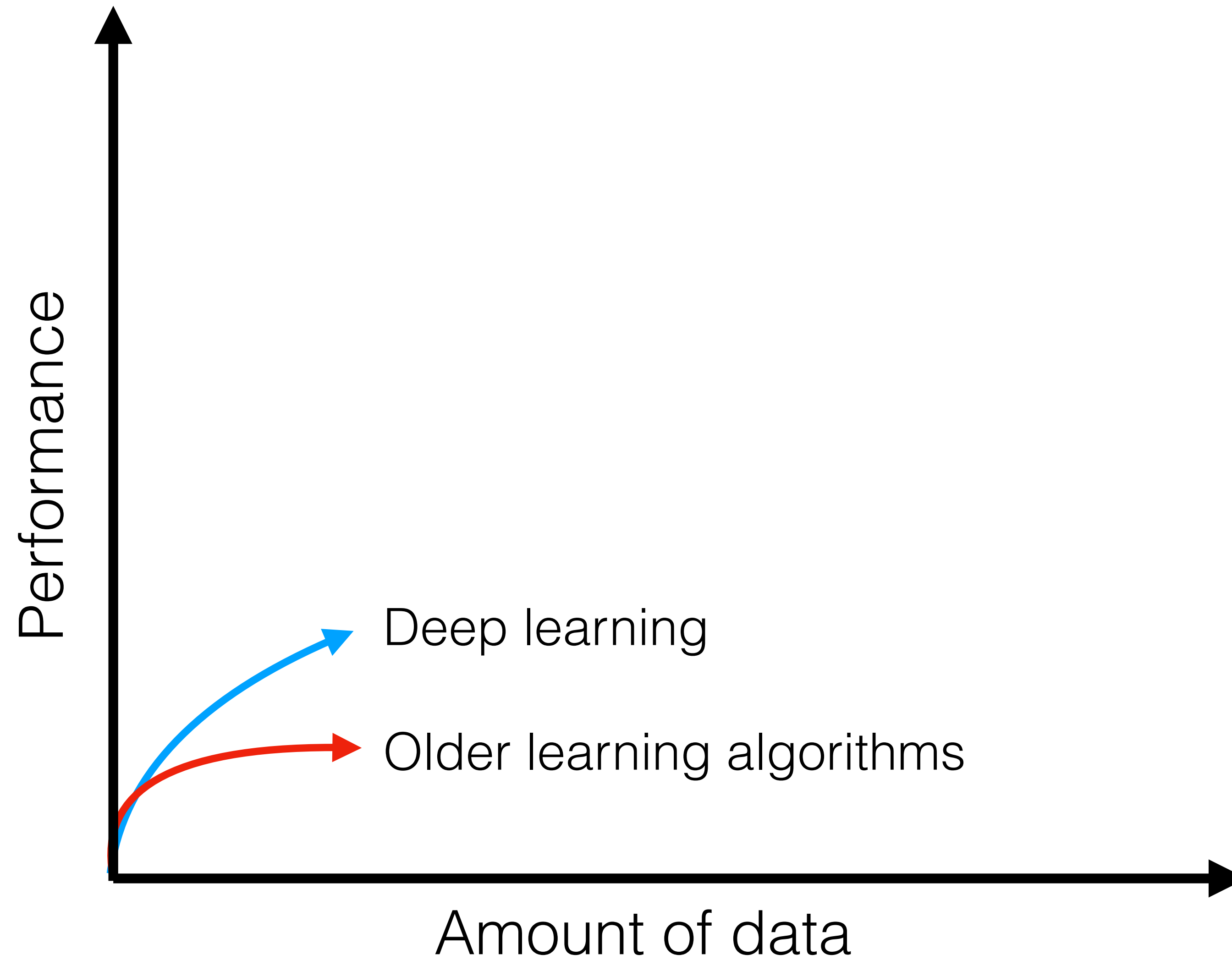


How do data science techniques scale with amount of data?



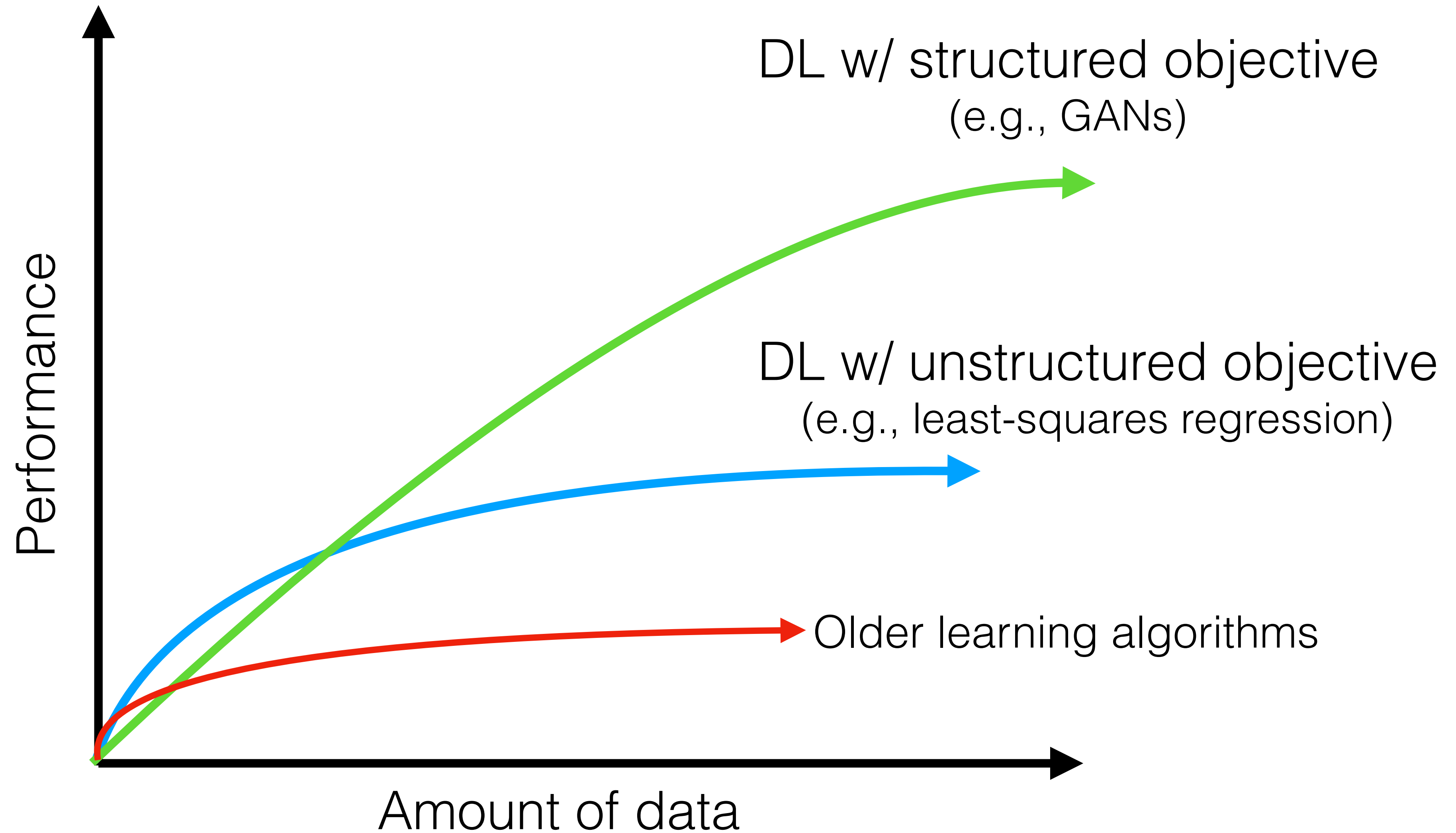
# Why structured objectives

(cartoon)



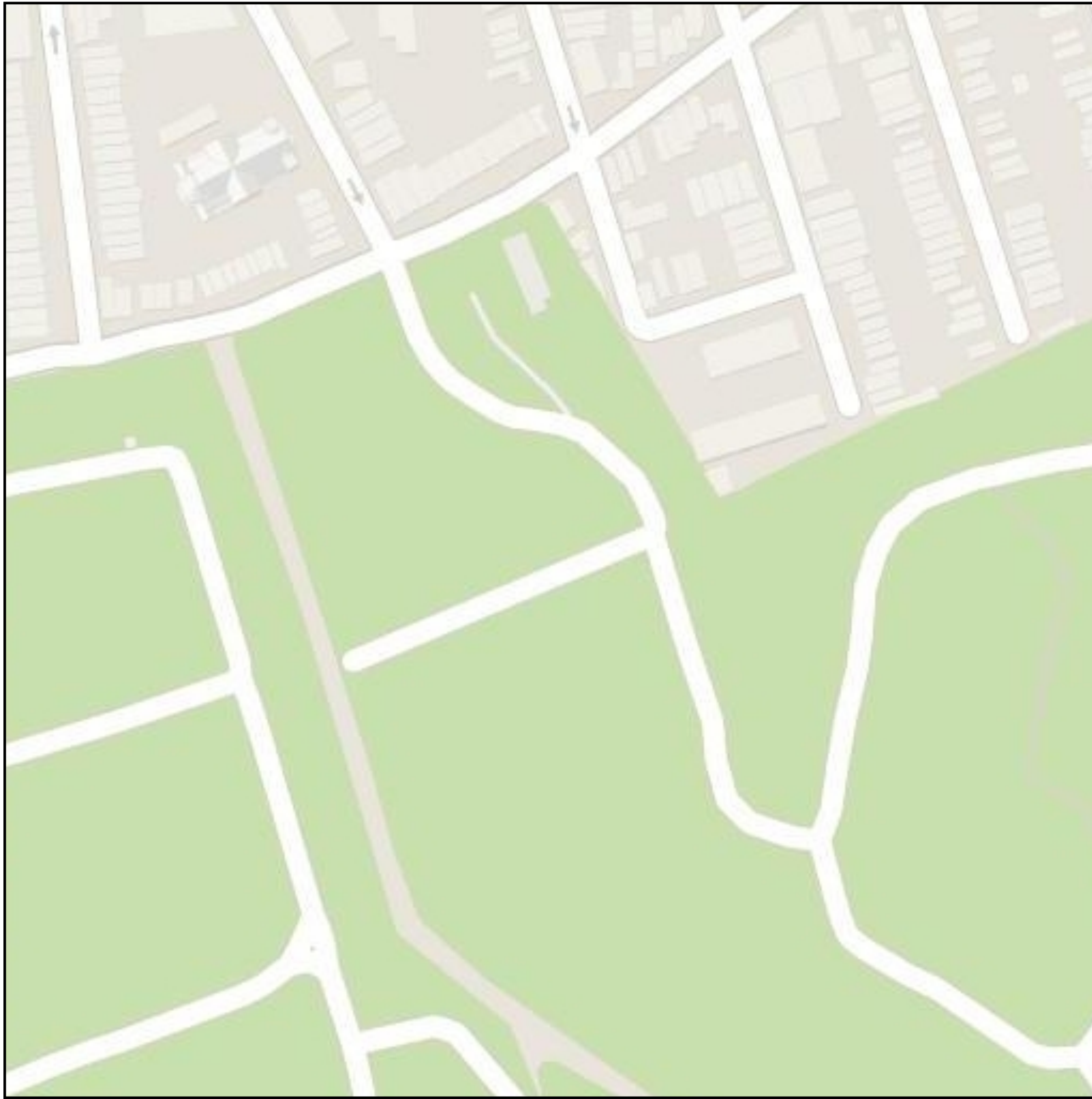
# Why structured objectives

(cartoon)





Input

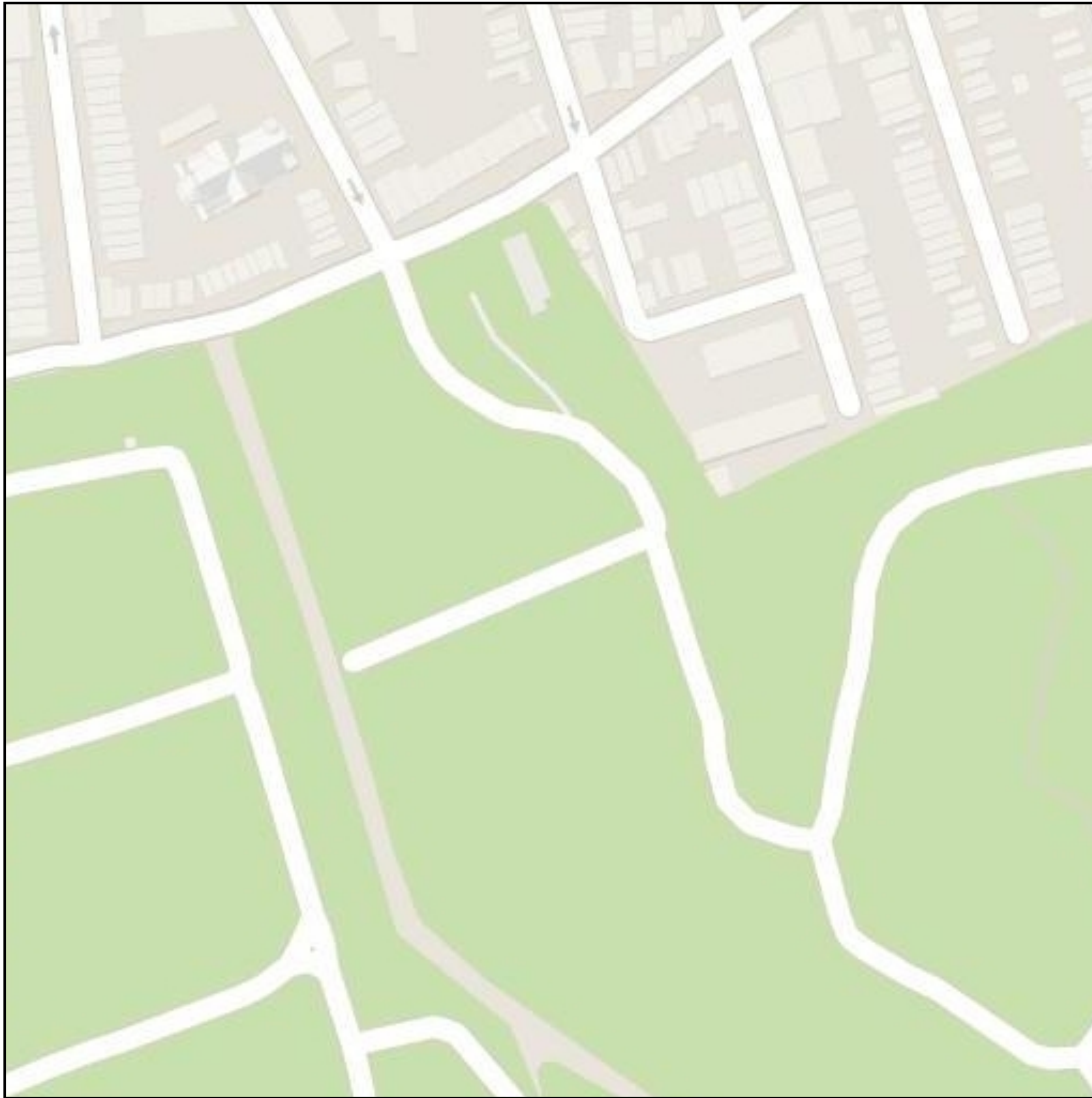


Unstructured prediction (L1)





Input



Structured Prediction (cGAN)

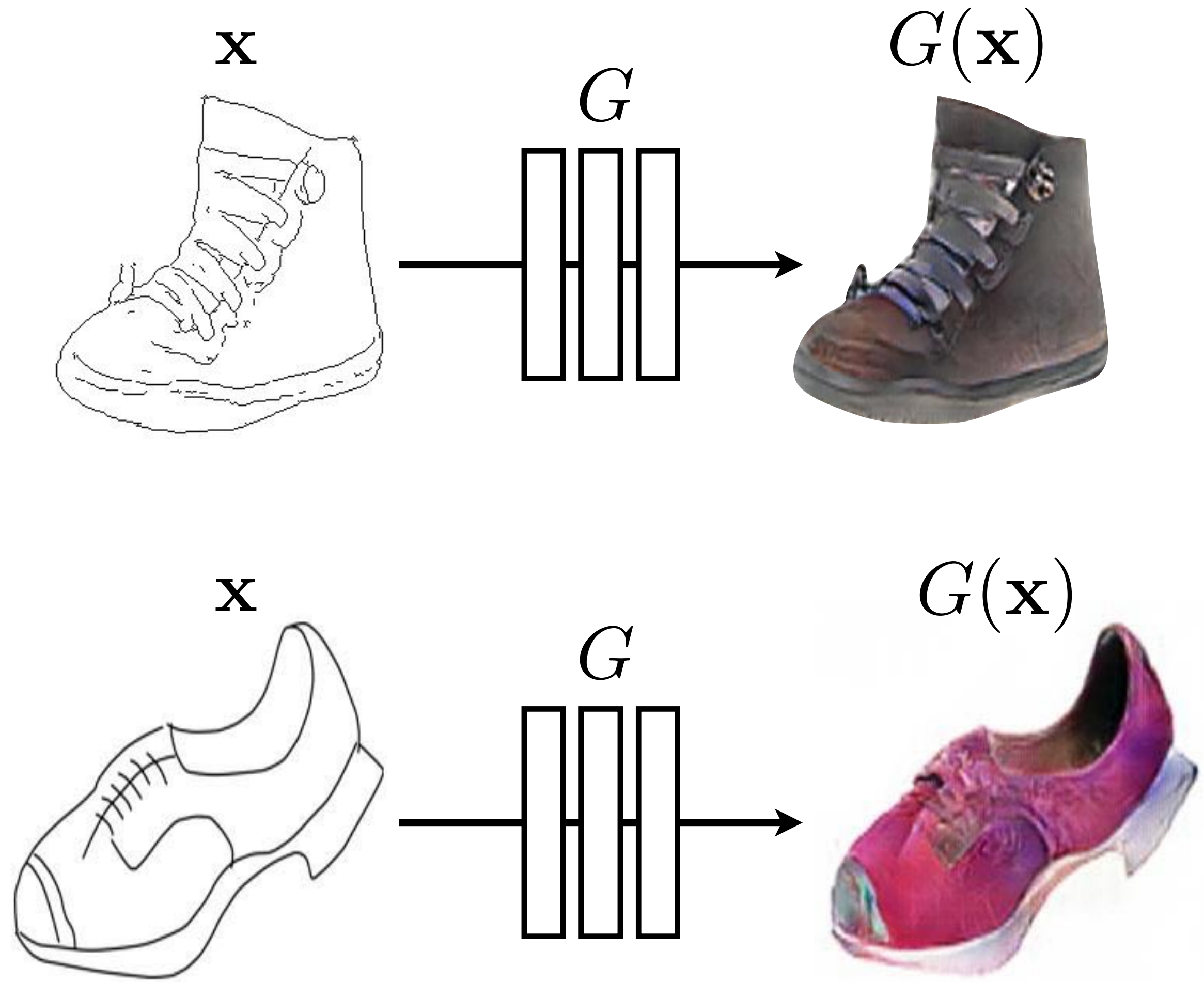




# Training data

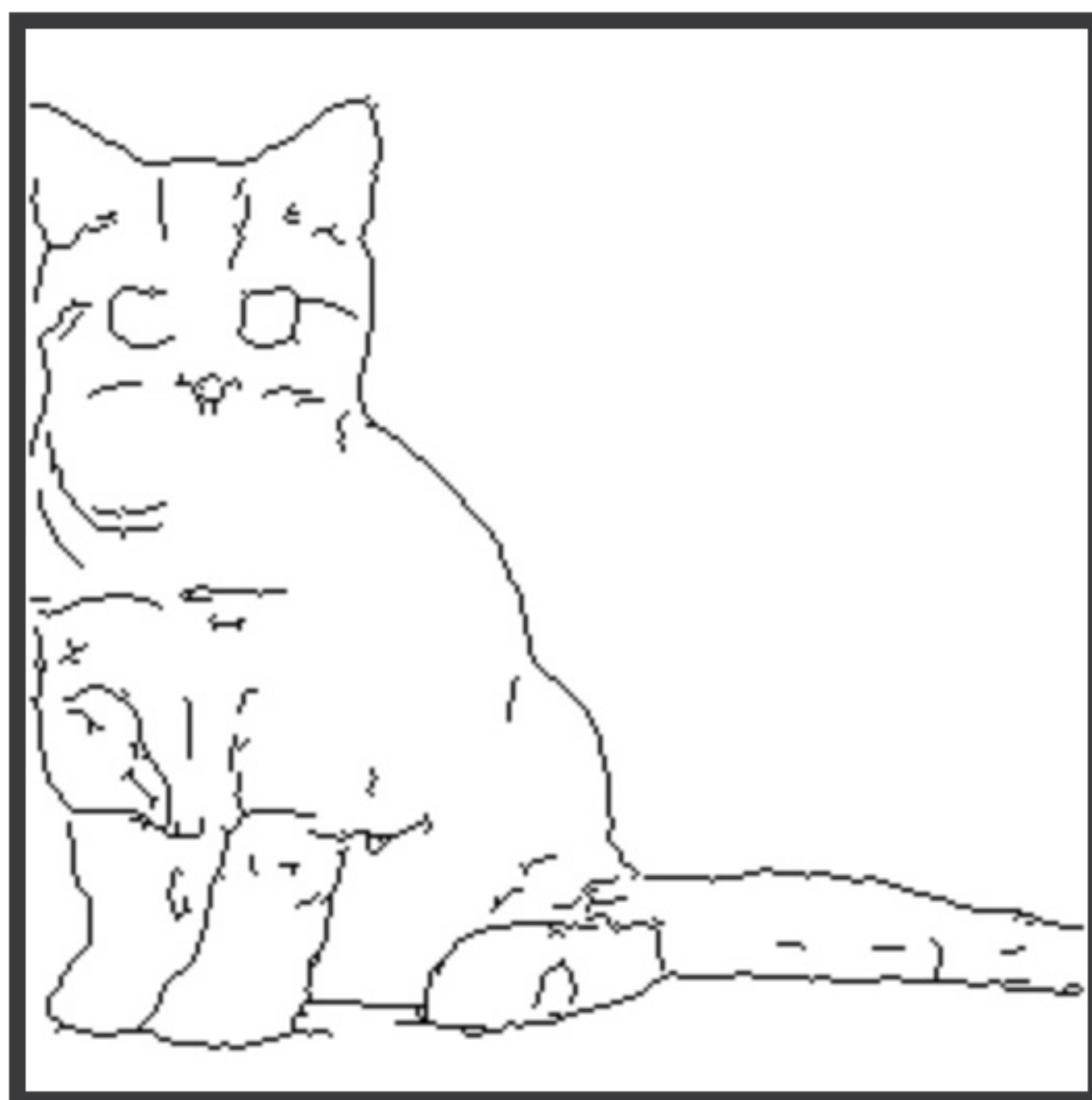


[HED, Xie & Tu, 2015]



# #edges2cats [Chris Hesse]

INPUT



pix2pix  
process

OUTPUT



undo

clear

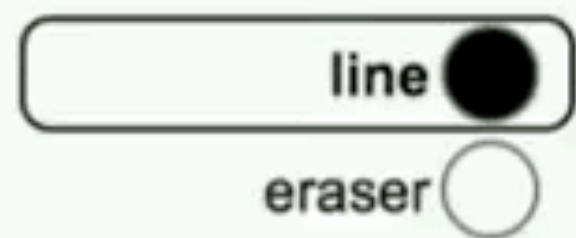
random

save

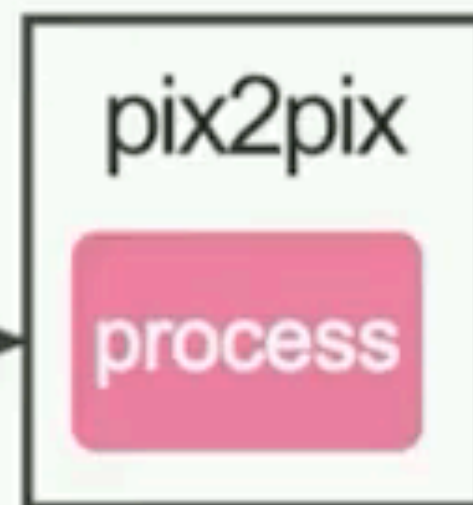
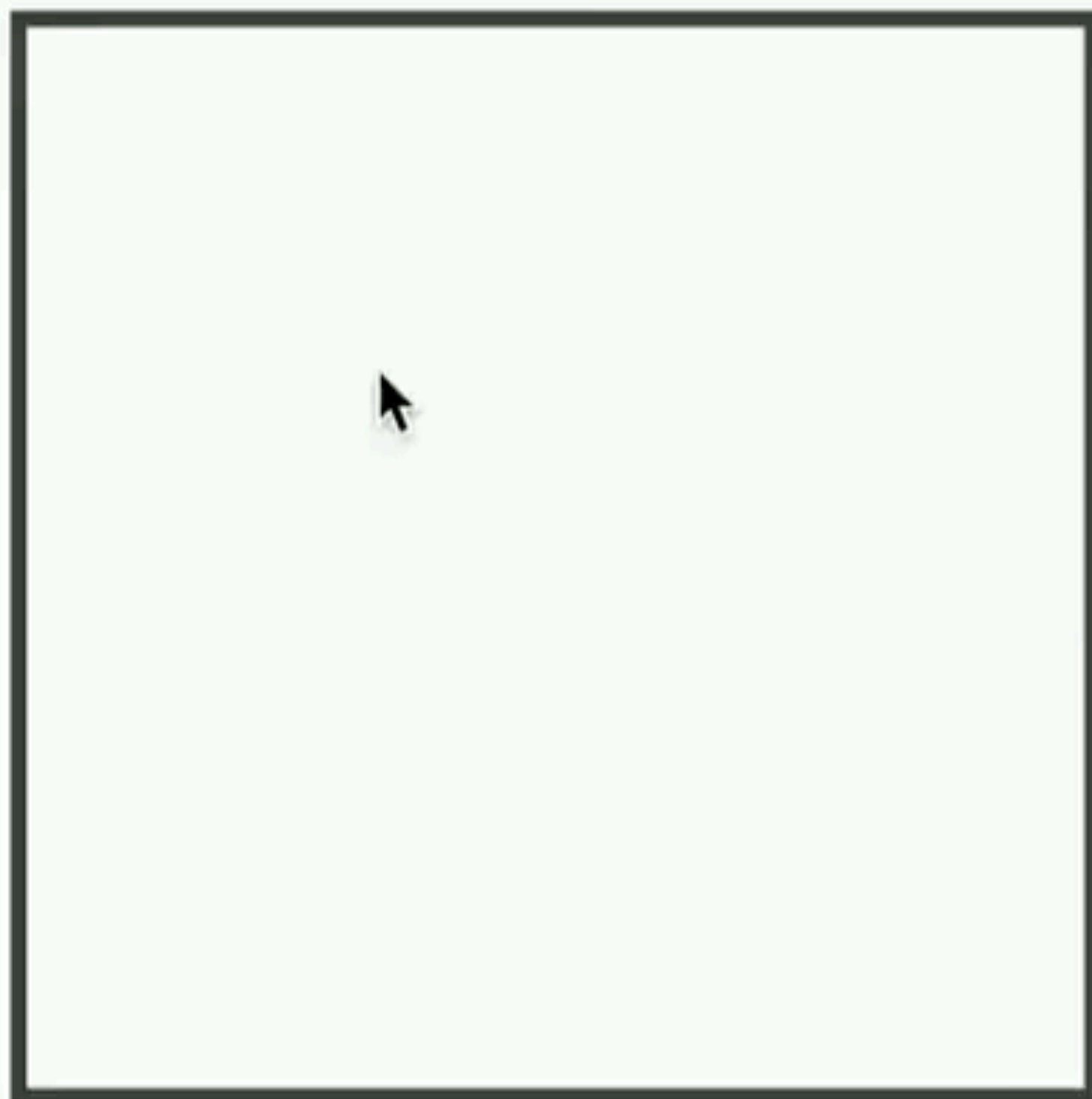


# edges2cats

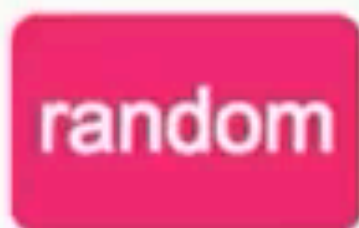
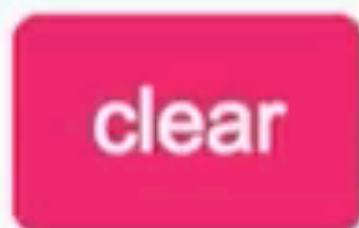
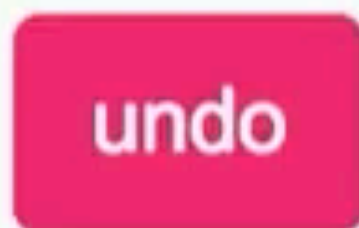
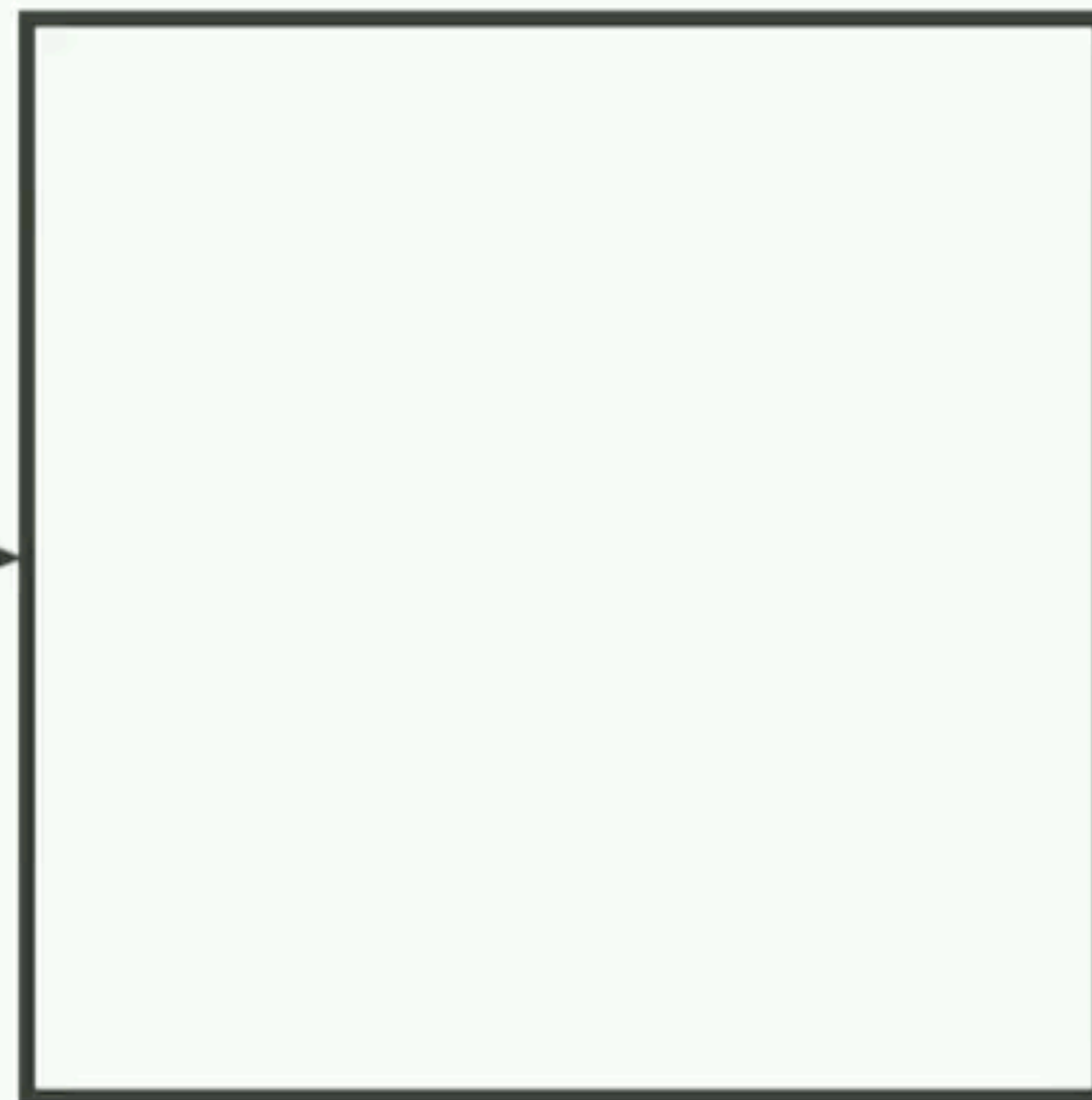
TOOL



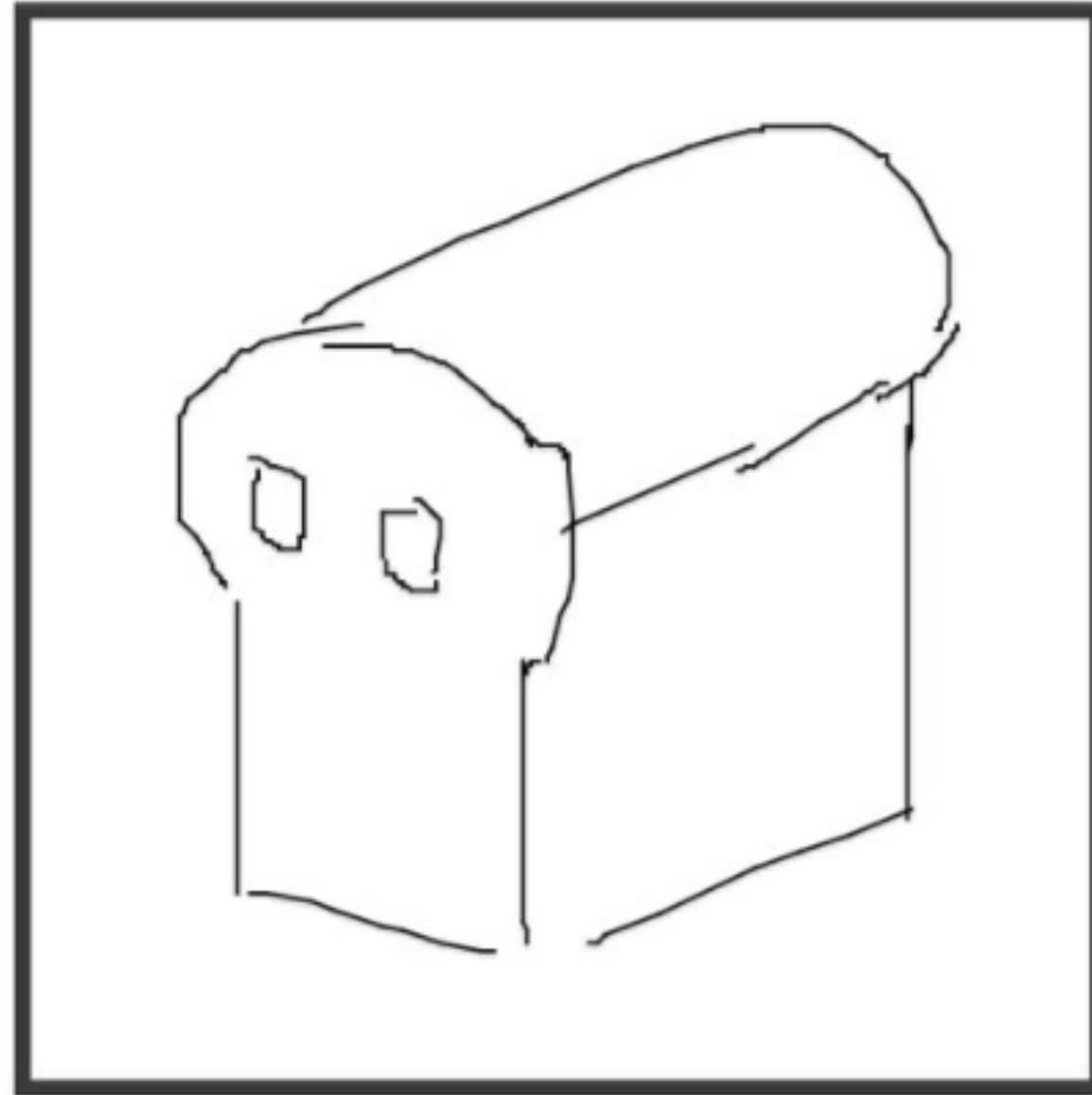
INPUT



OUTPUT



INPUT



OUTPUT



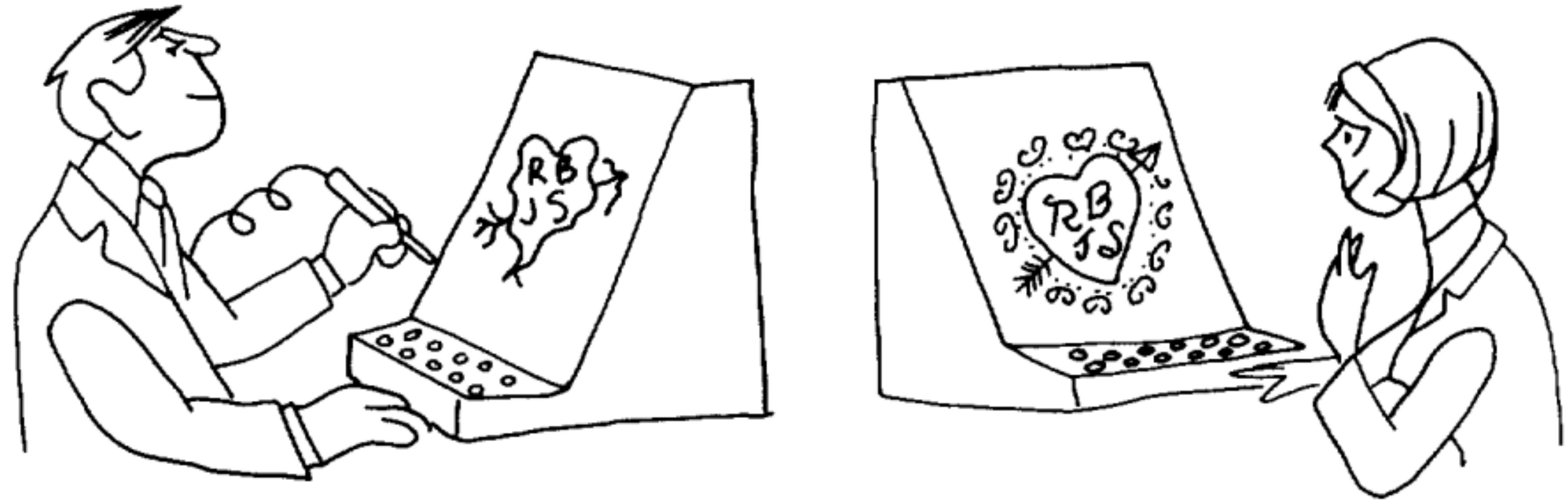
Ivy Tasi @ivymyt



Vitaly Vidmirov @vvid



1. Image synthesis
2. Structured prediction
3. **Domain mapping**

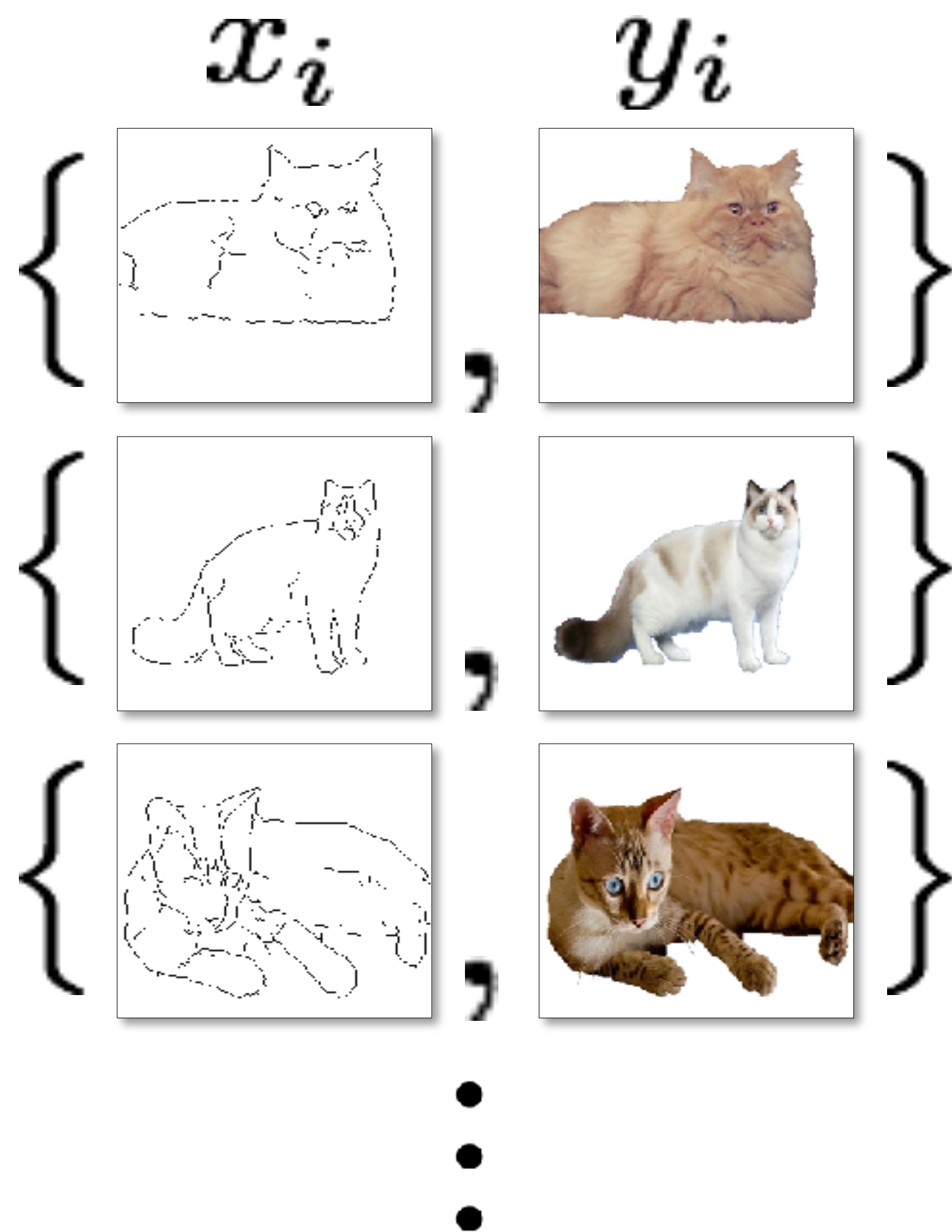


# Domain mapping

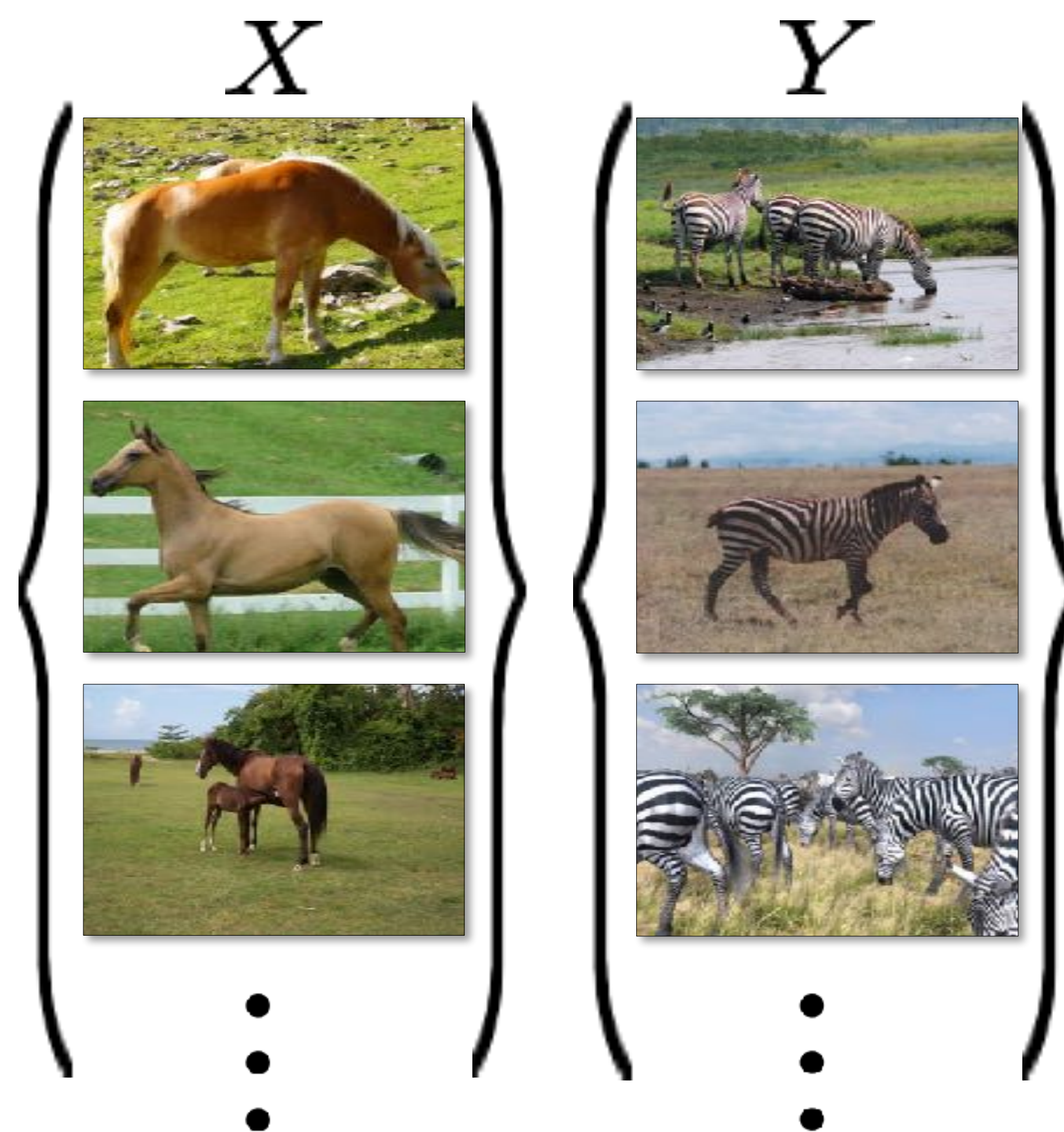
[Includes slides from Jun-Yan Zhu, Taesung Park]

[Cartoon: The Computer as a Communication Device, Licklider & Taylor 1968]

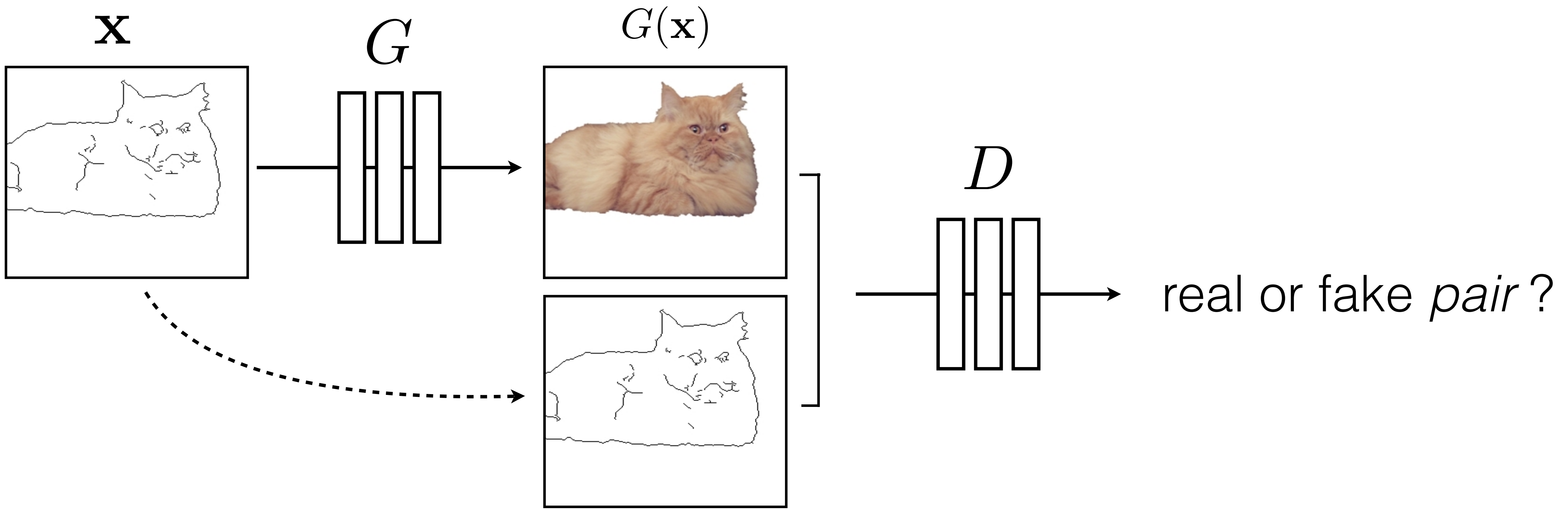
# Paired data



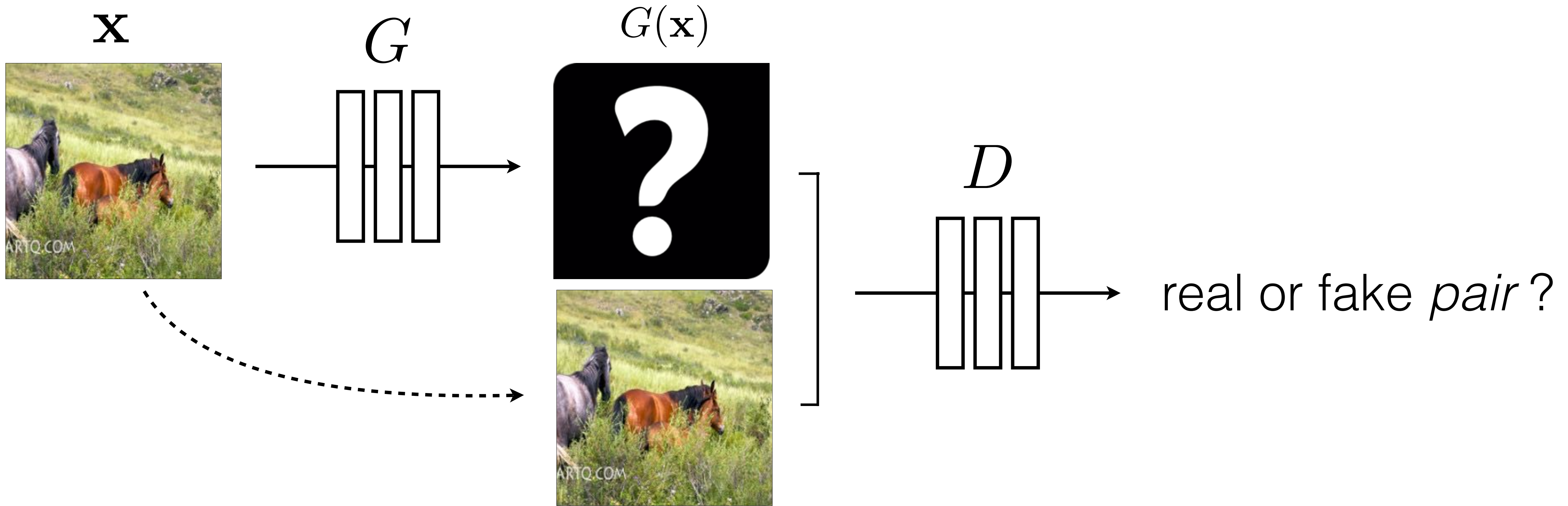
# Unpaired data







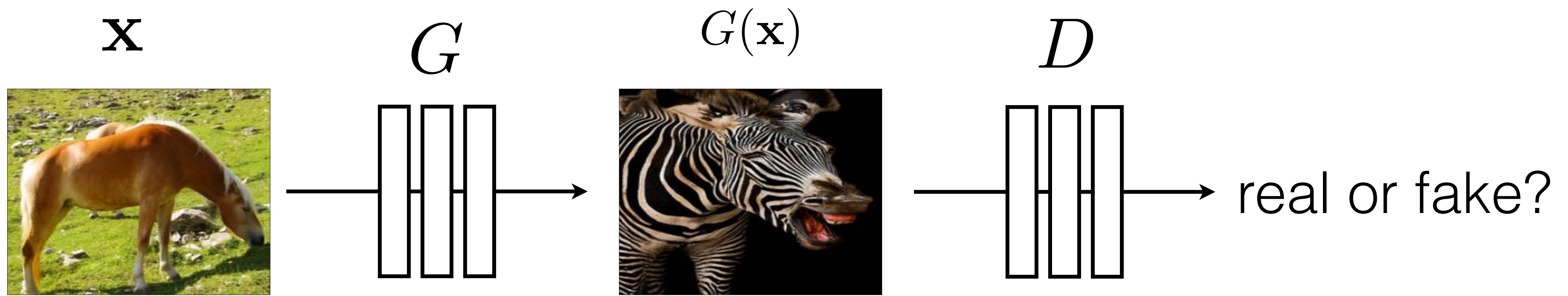
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [ \log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y})) ]$$



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

No input-output pairs!





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

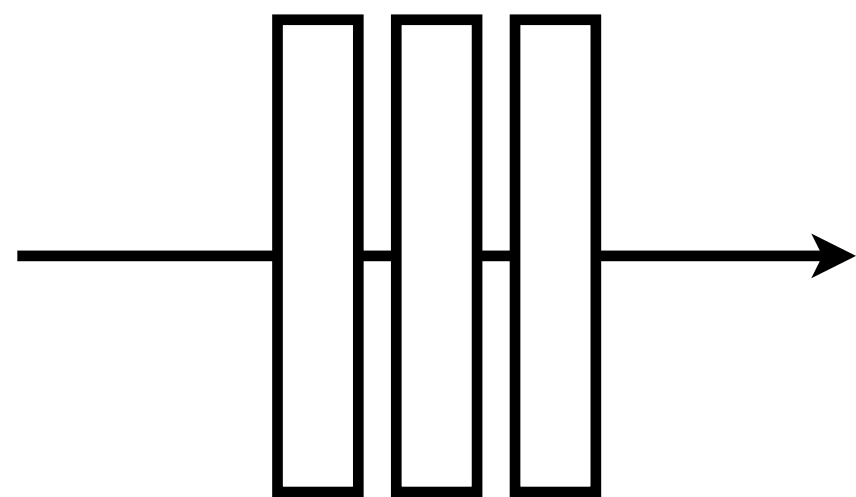
Usually loss functions check if output matches a target *instance*

GAN loss checks if output is part of an admissible *set*

$\mathbf{x}$



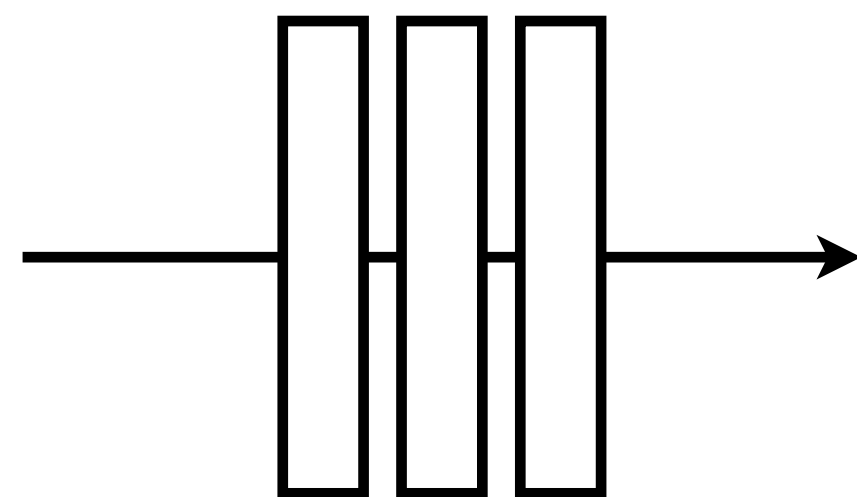
$G$



$G(\mathbf{x})$

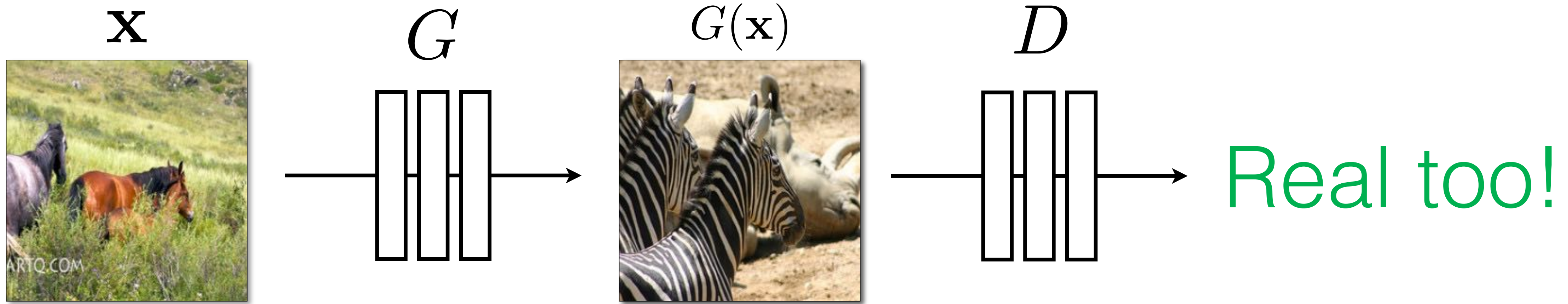


$D$



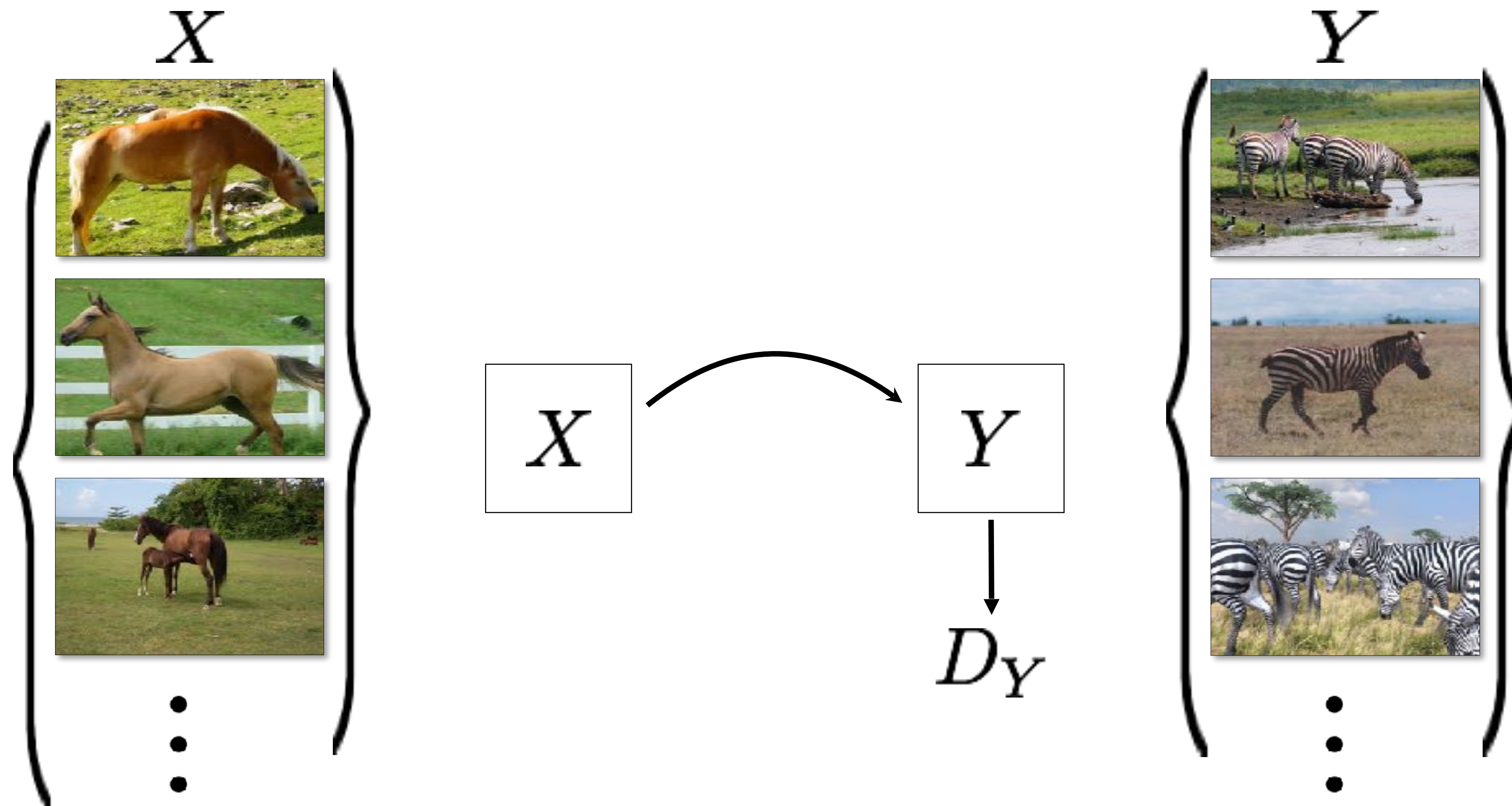
Real!





Nothing to force output to correspond to input

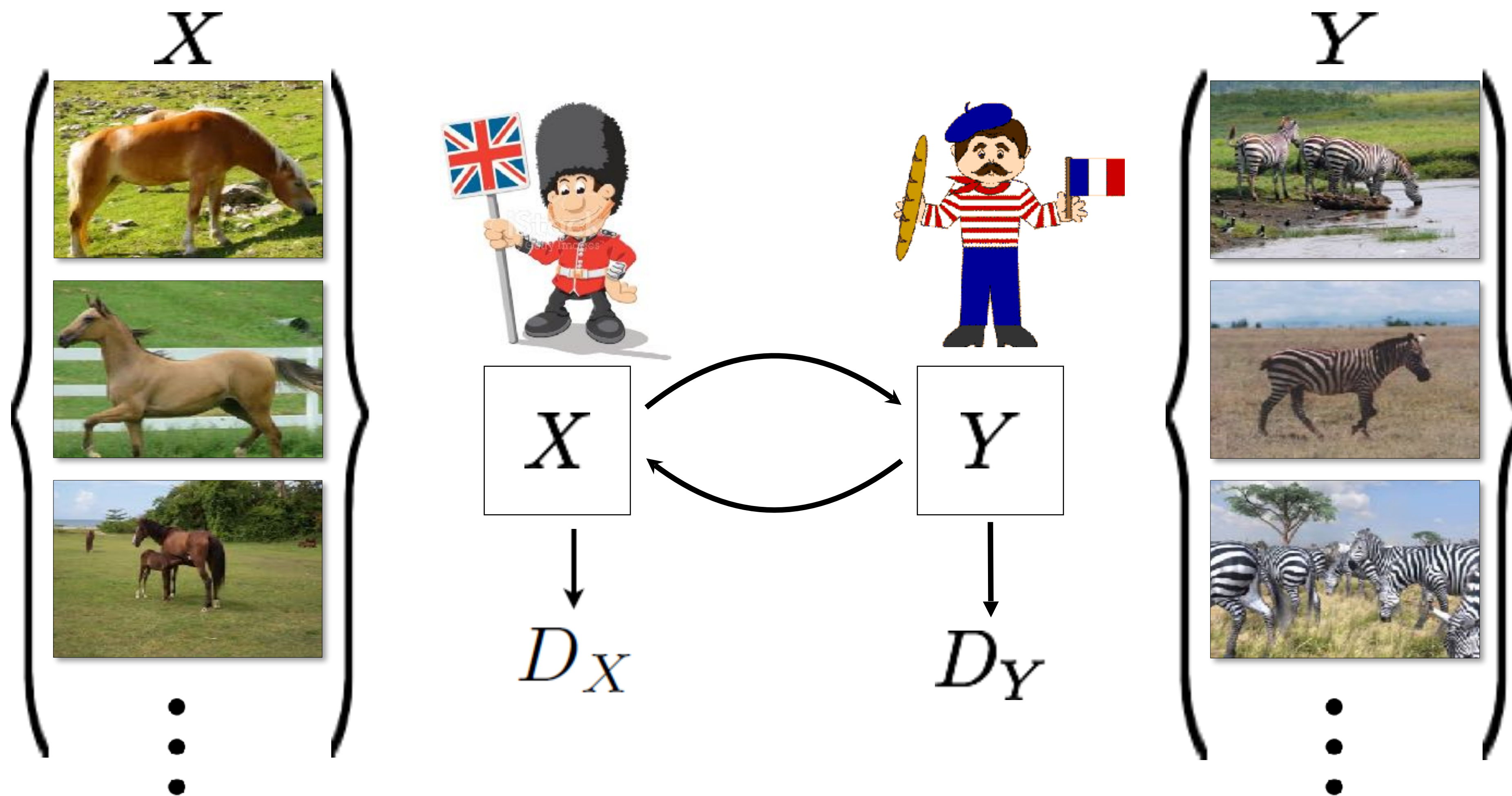
# CycleGAN, or there and back aGAN



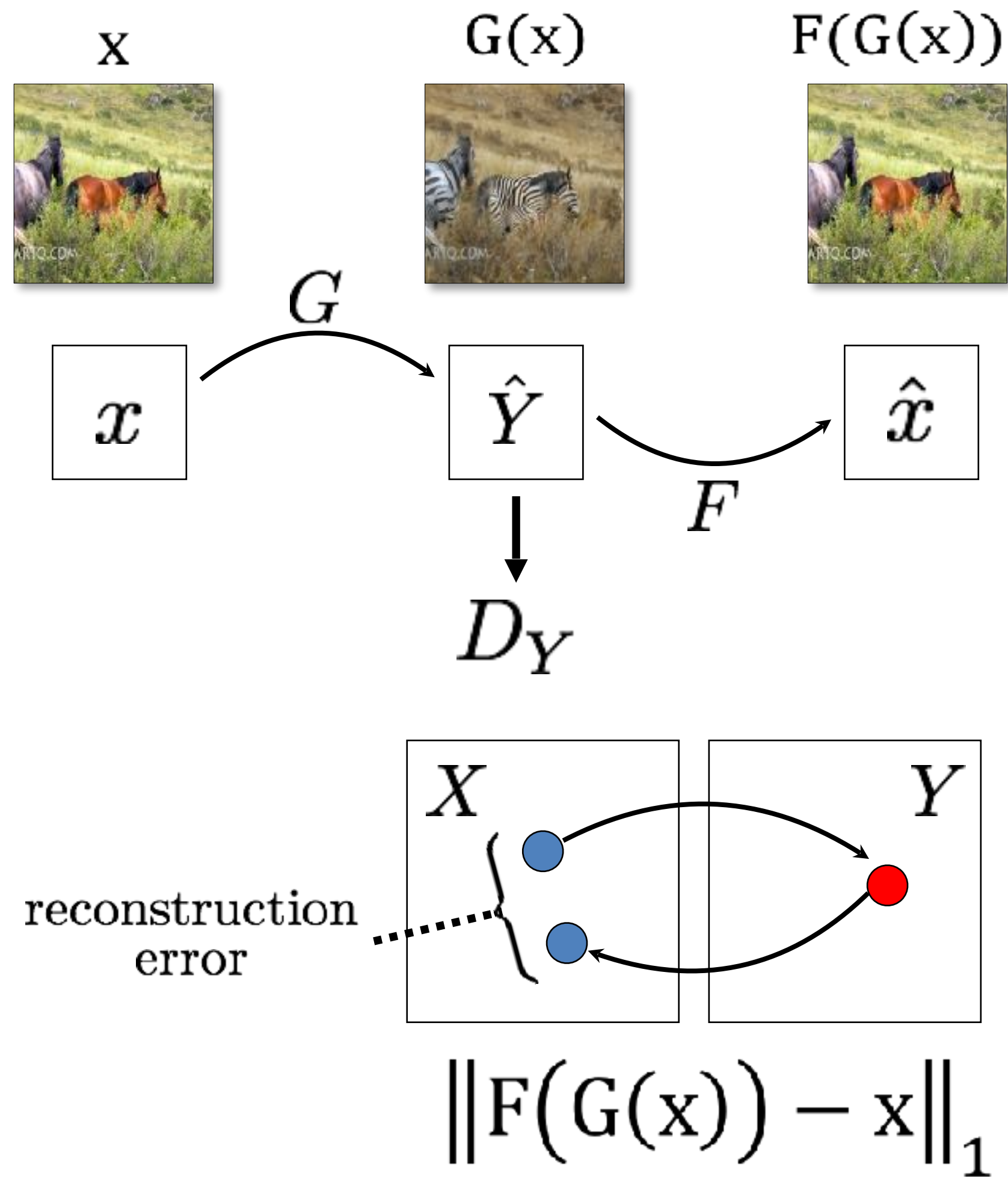
[Zhu\*, Park\* et al. 2017], [Yi et al. 2017], [Kim et al. 2017]



# CycleGAN, or there and back aGAN

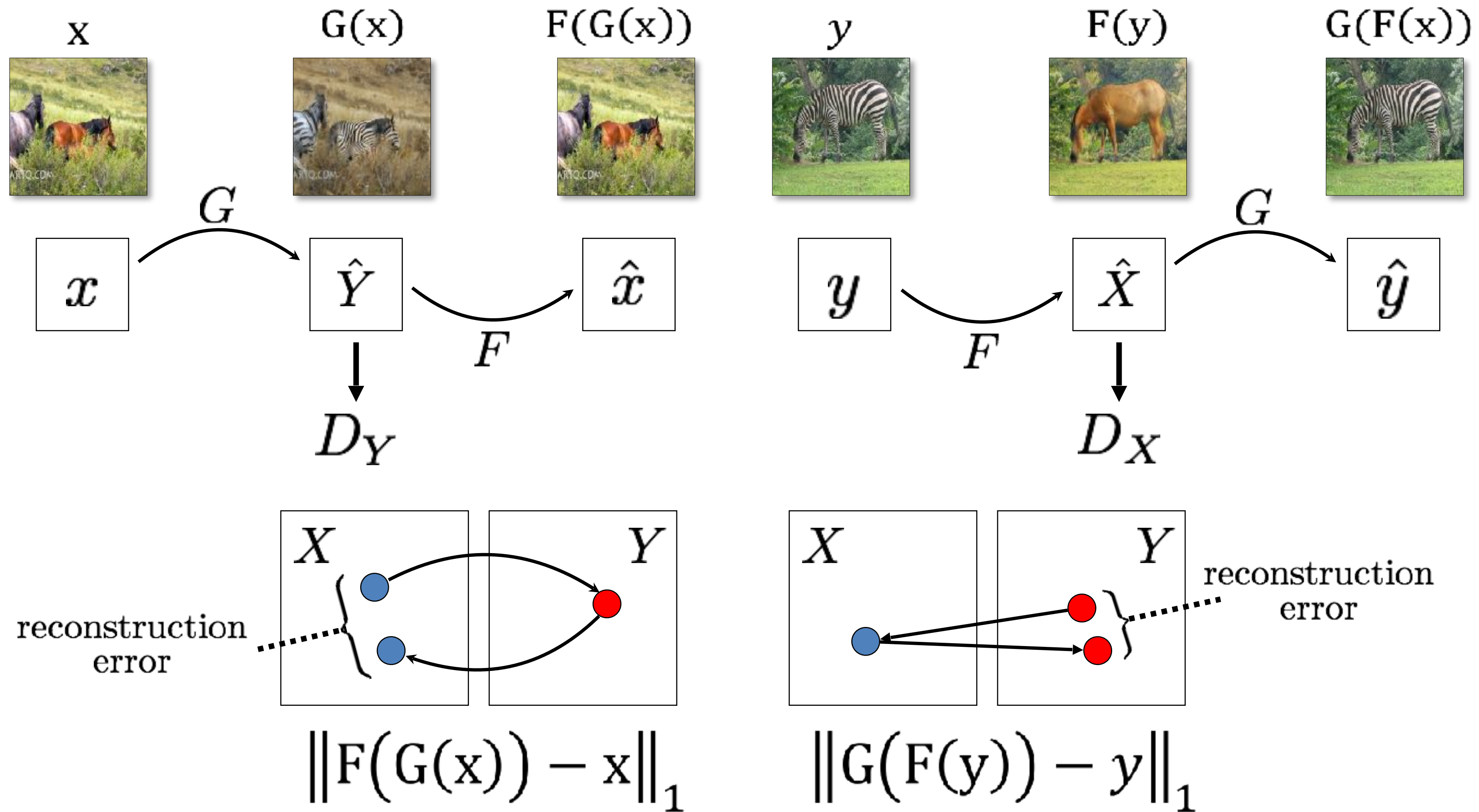


# Cycle Consistency Loss



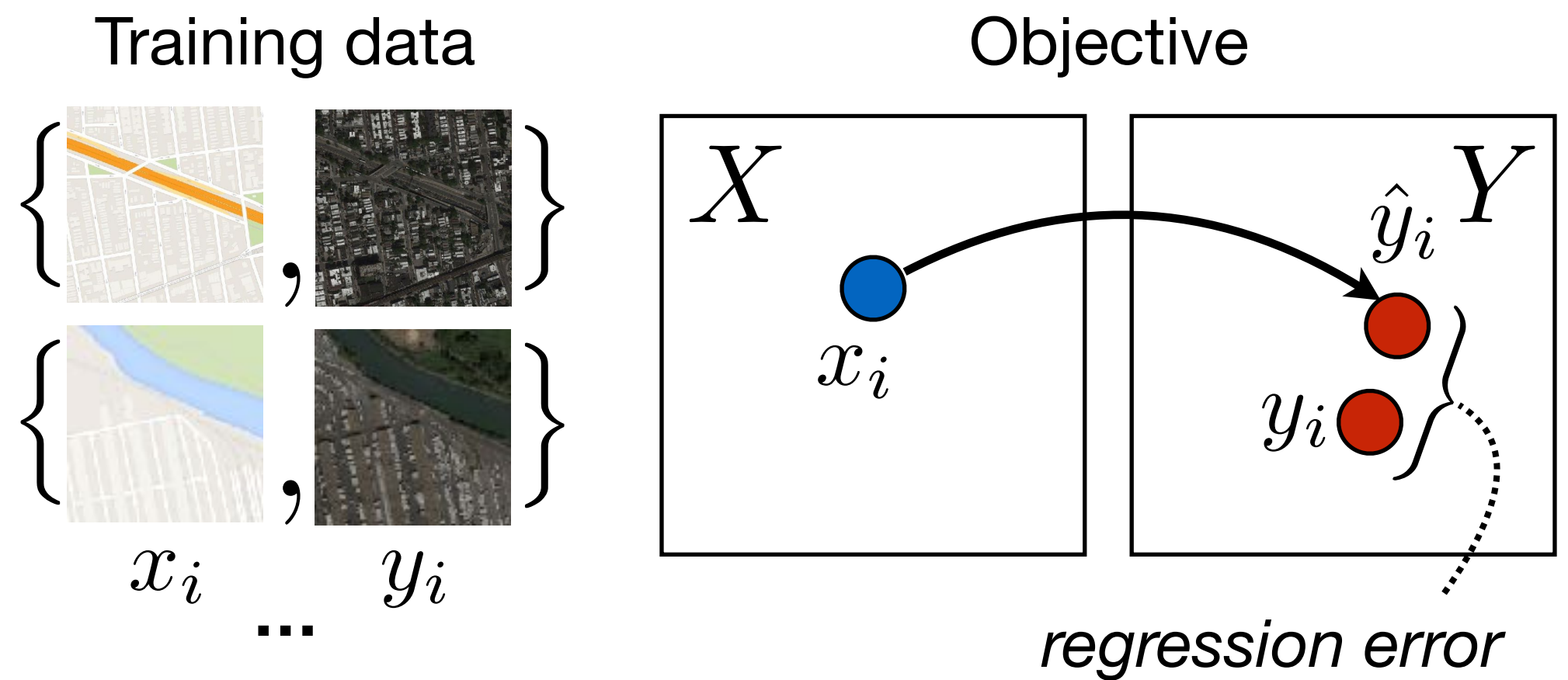


# Cycle Consistency Loss

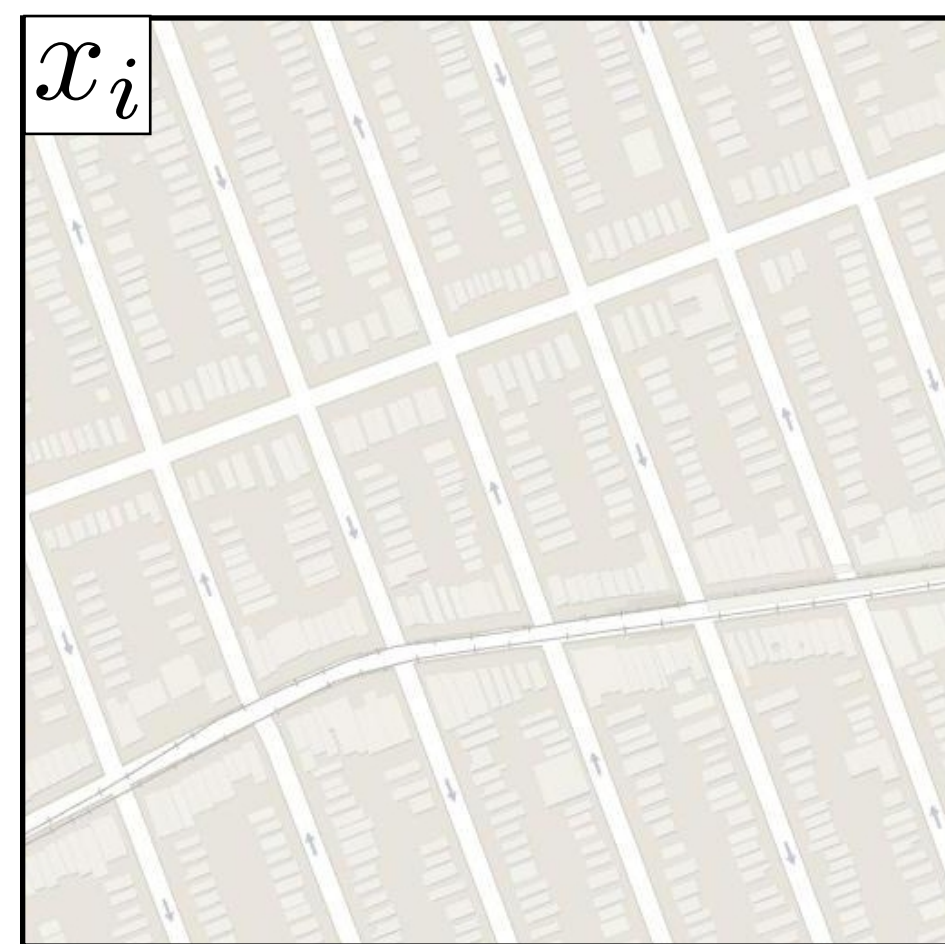




# Paired translation



Input

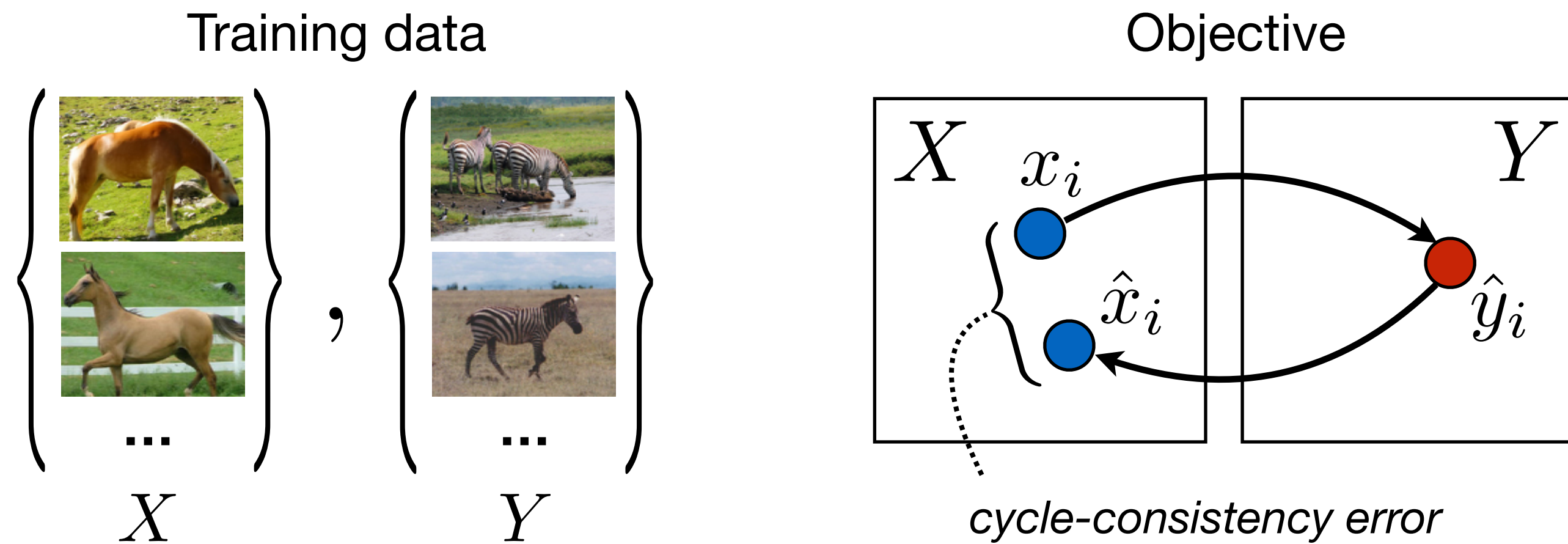


Result



[“pix2pix”, Isola, Zhu, Zhou, Efros, 2017]

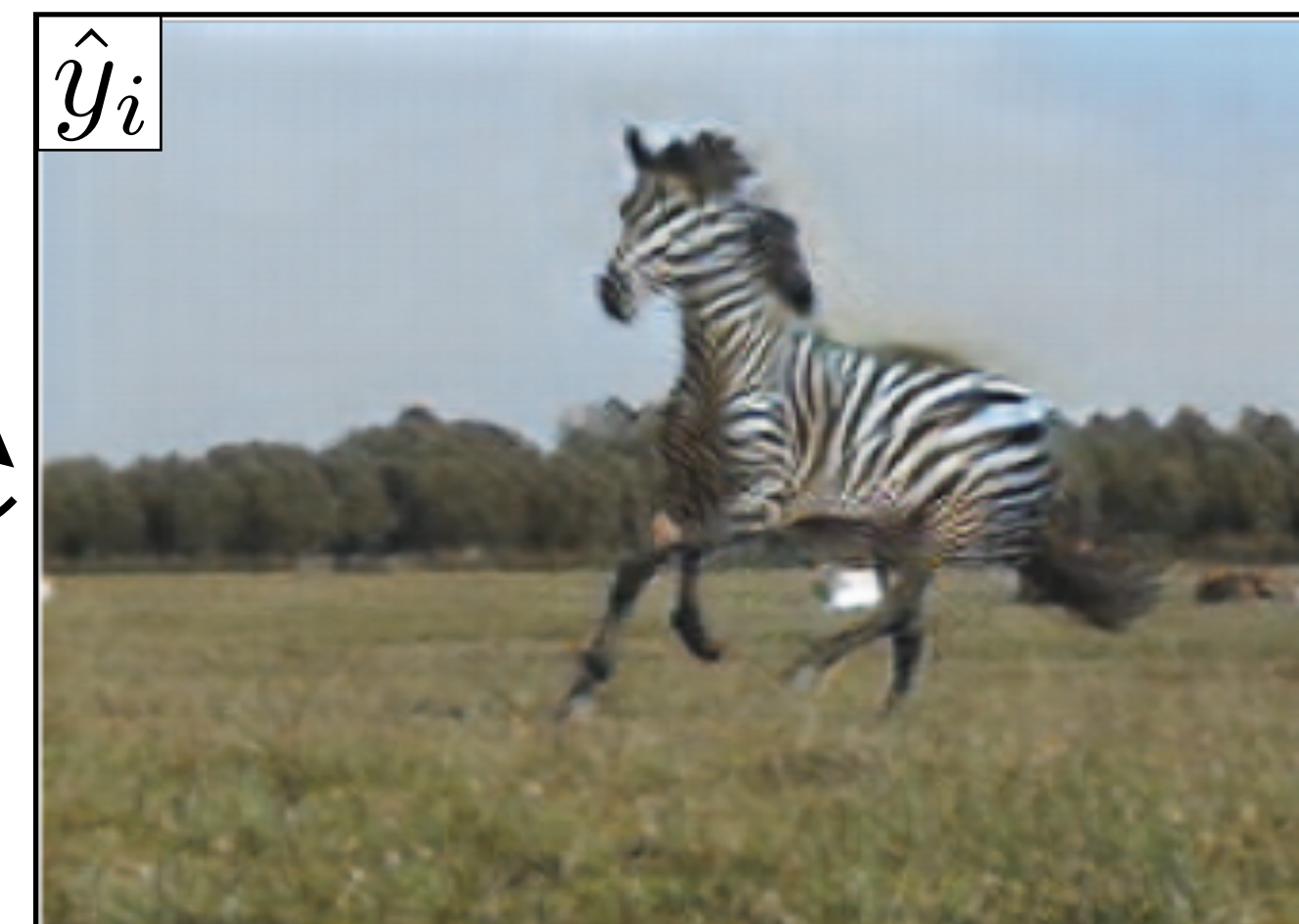
# Unpaired translation



Input



Result



[“CycleGAN”, Zhu\*, Park\*, Isola, Efros, 2017]











Input



Monet



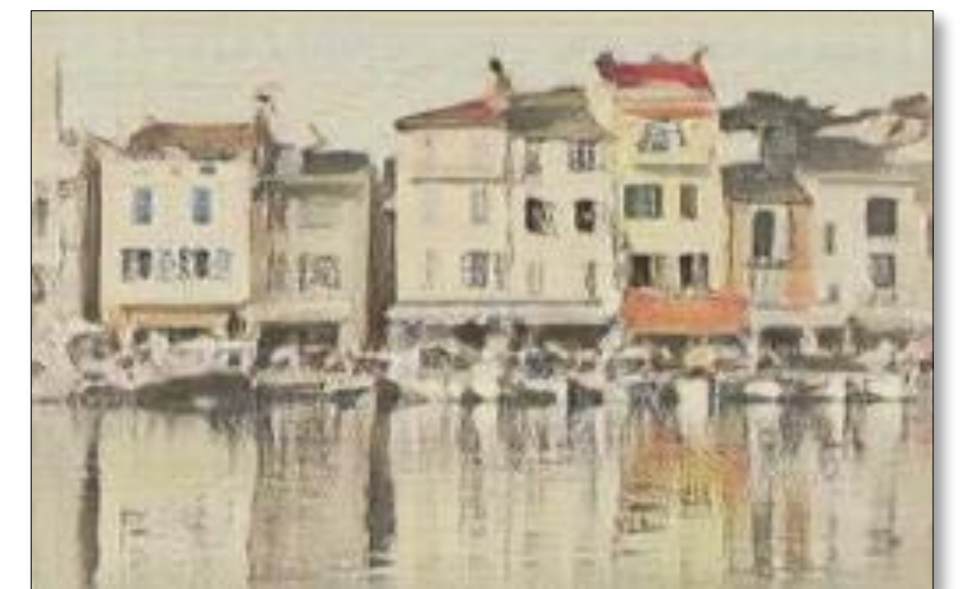
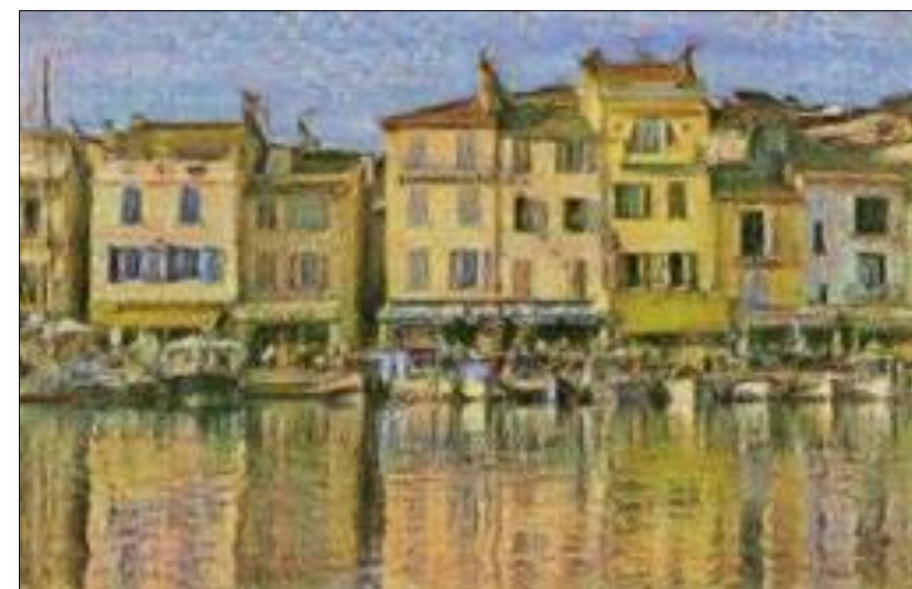
Van Gogh



Cezanne



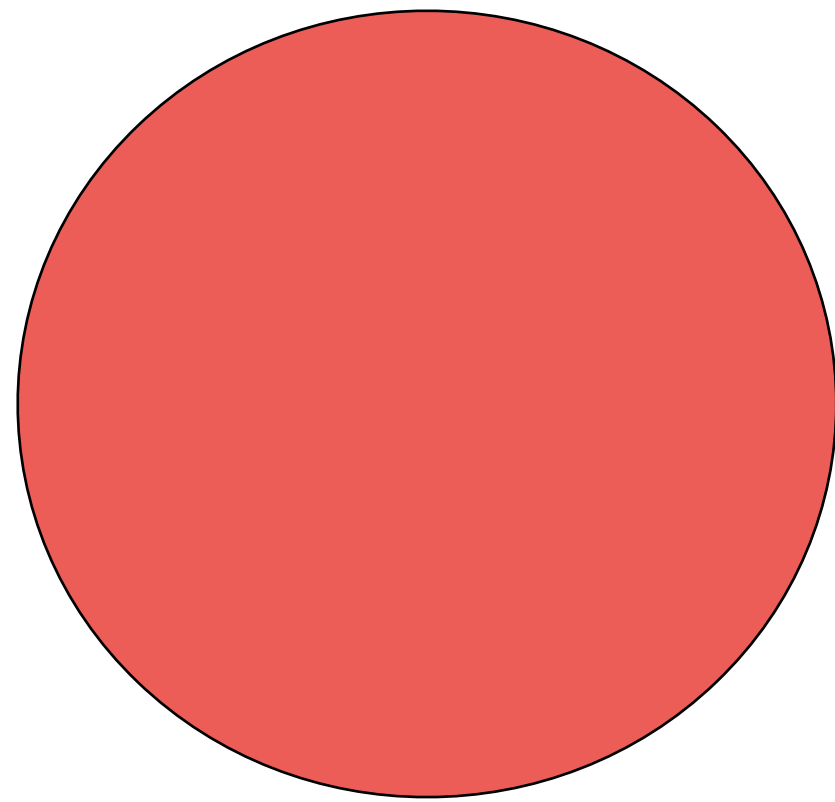
Ukiyo-e



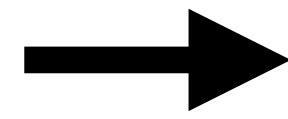


# GANs

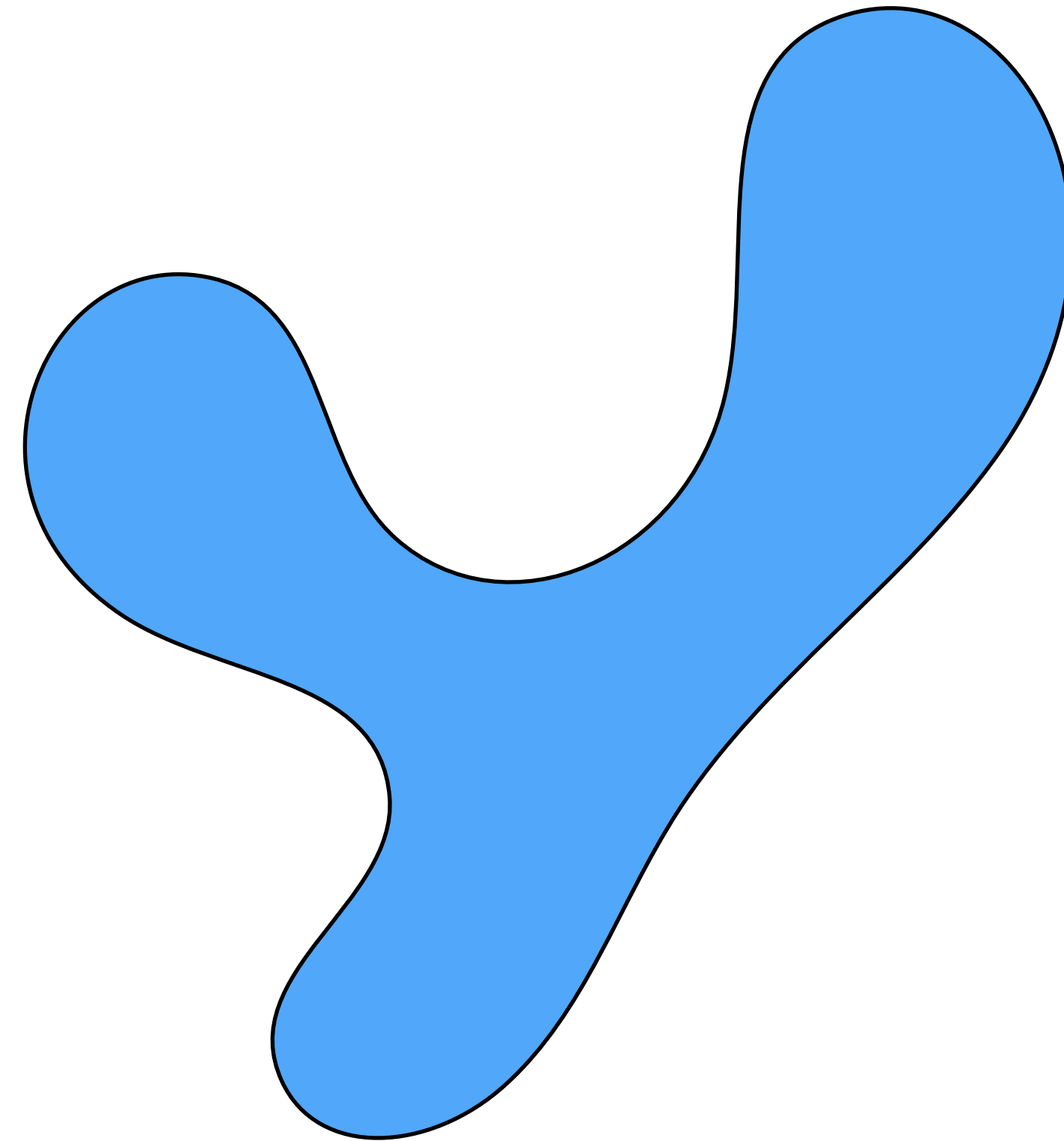
Gaussian



**Z**



Target distribution



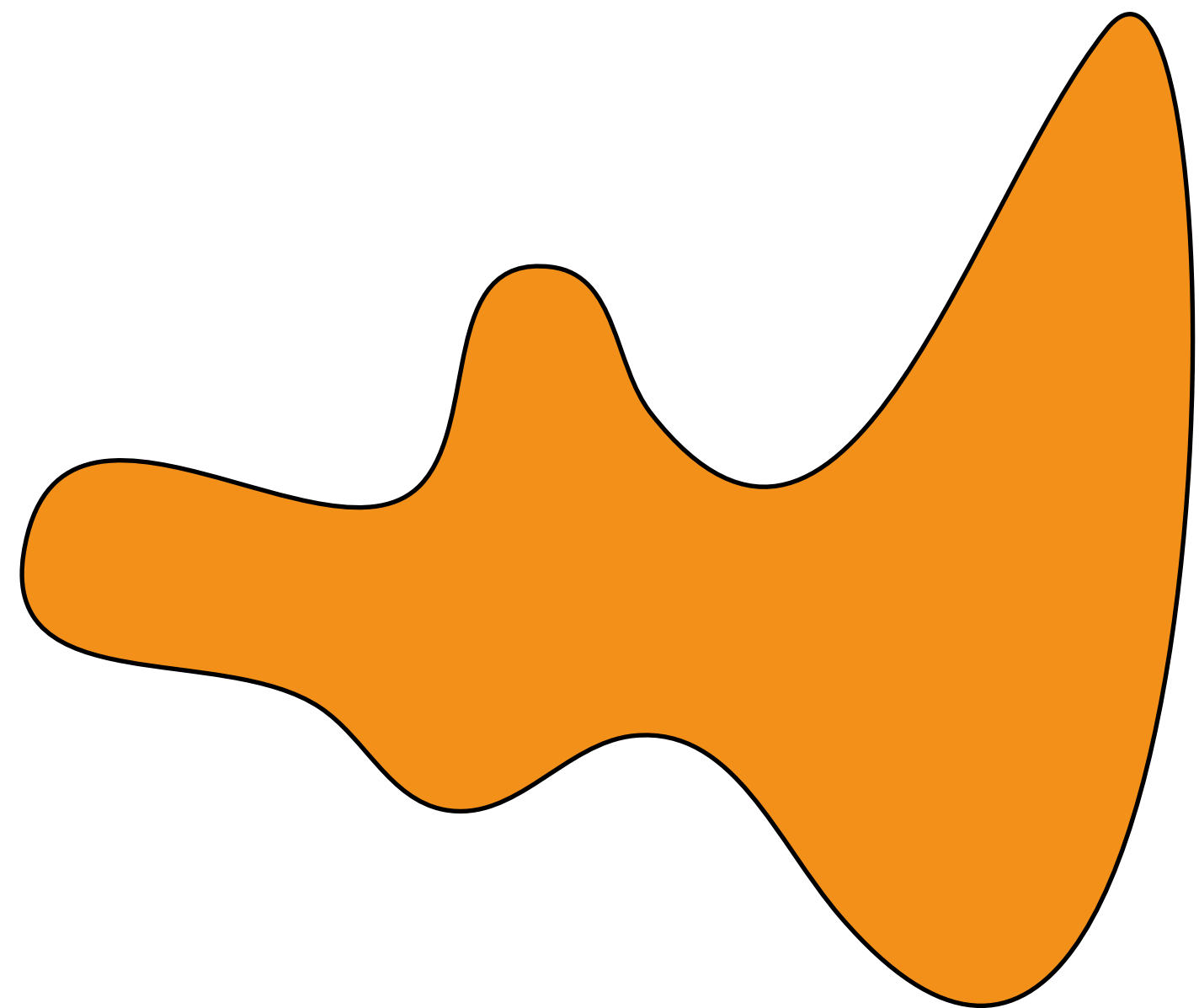
**Y**



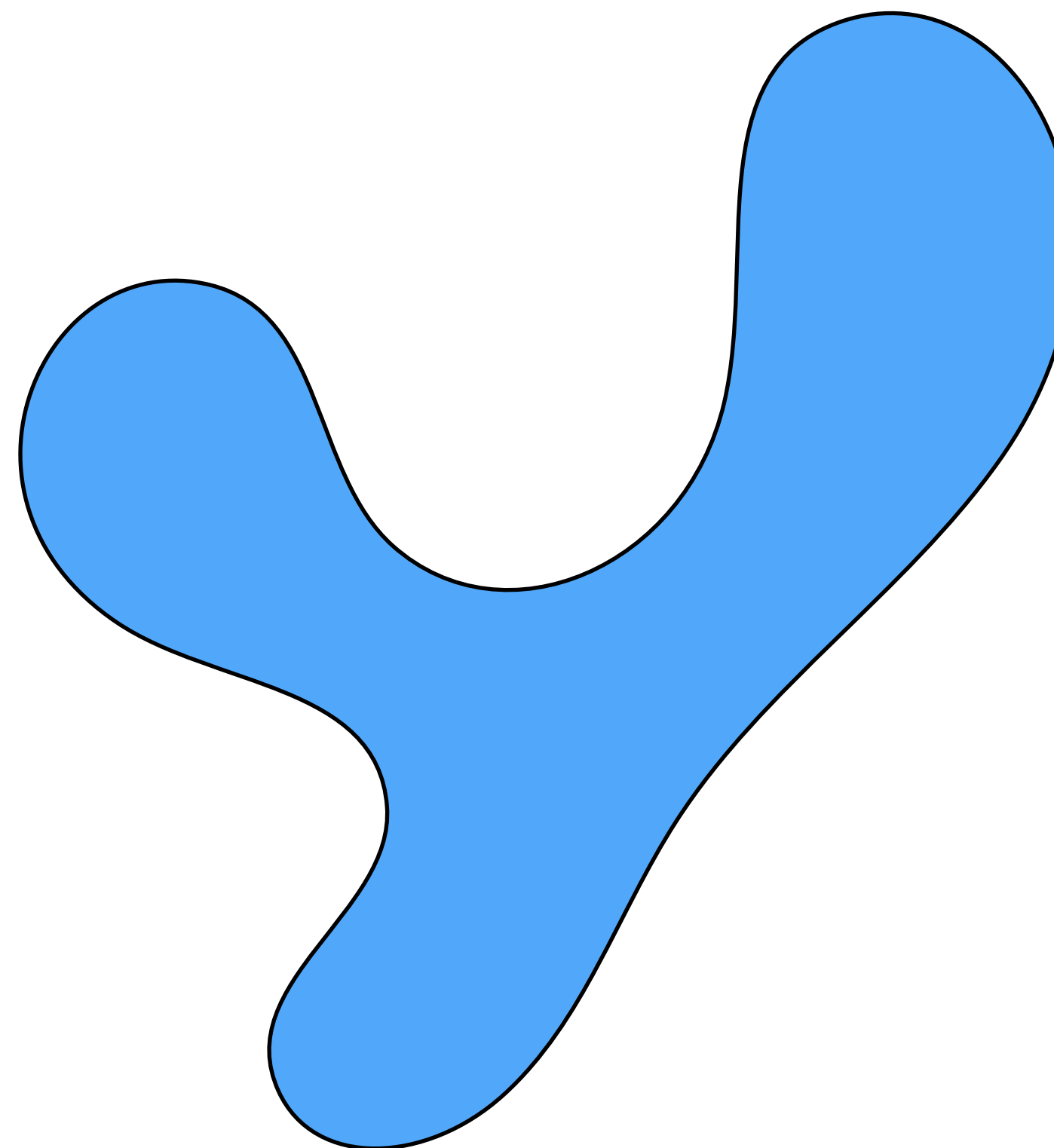
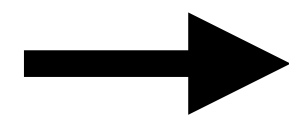
# CycleGAN

Horses

Zebras

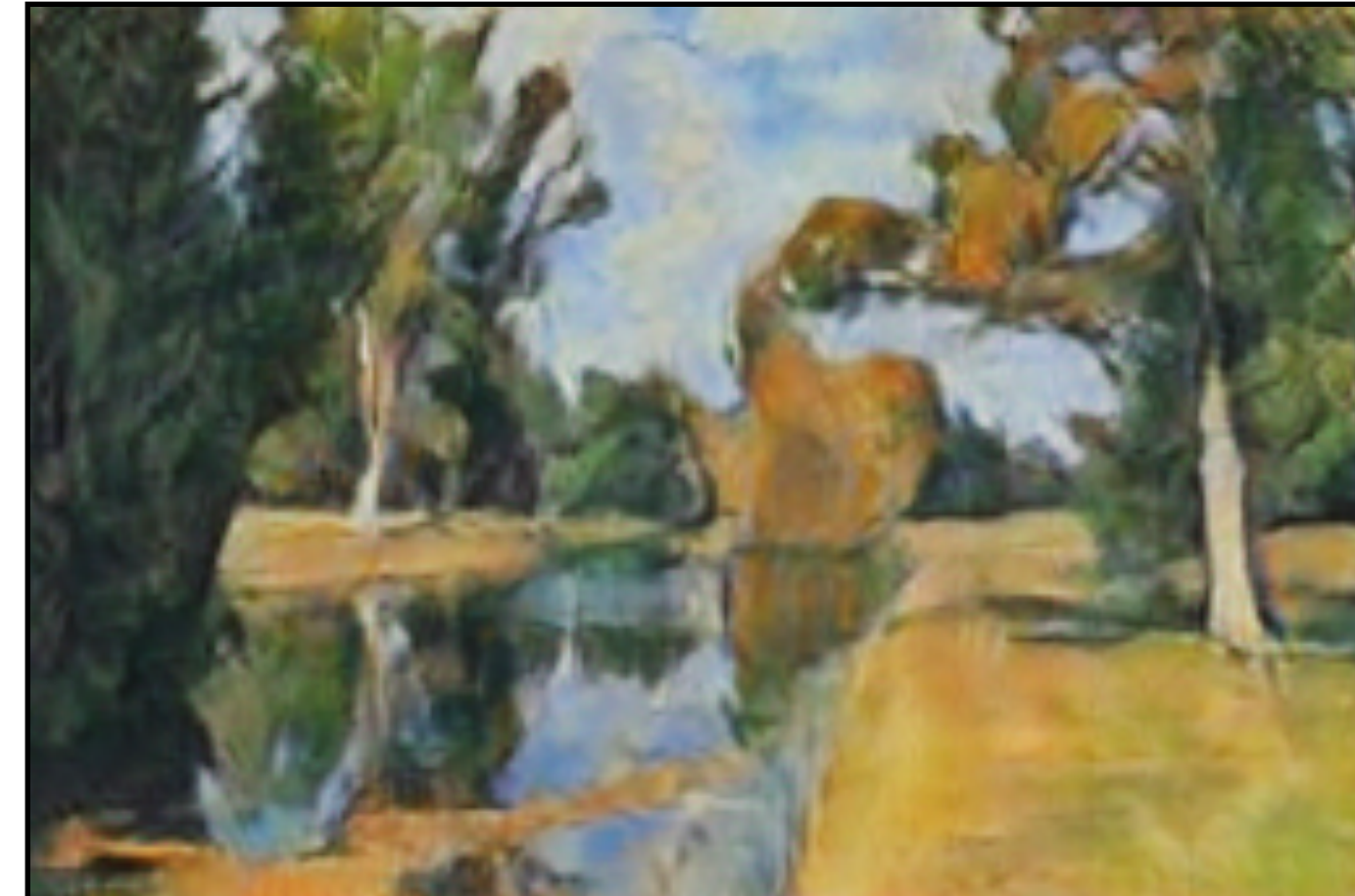
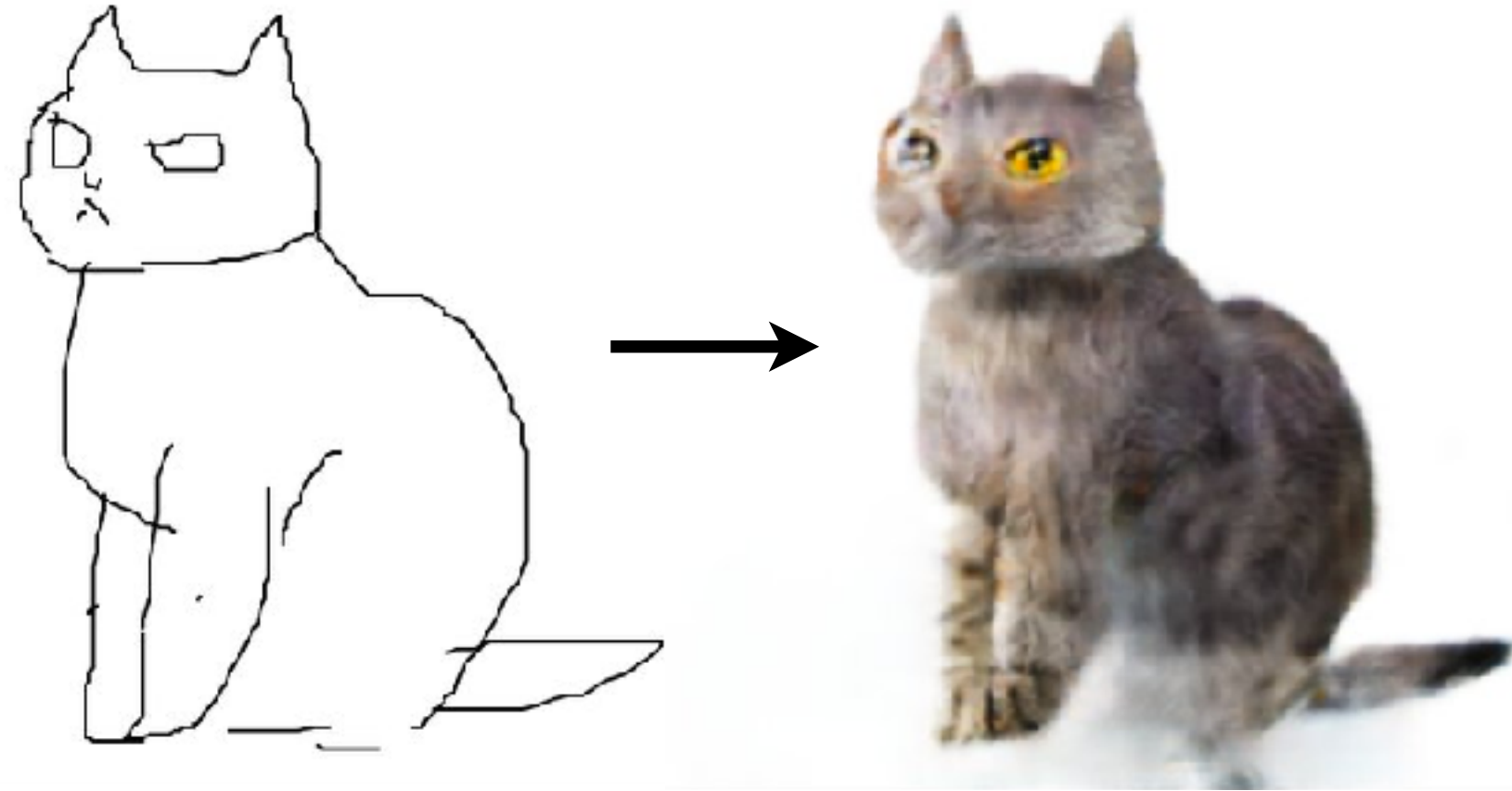


**X**



**Y**

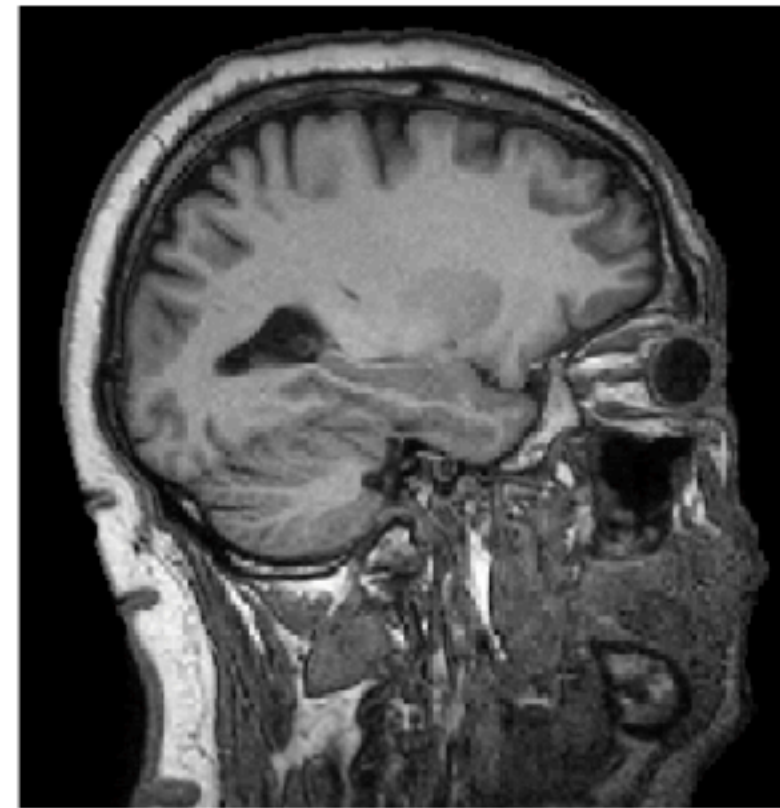
# What would it look like if...?



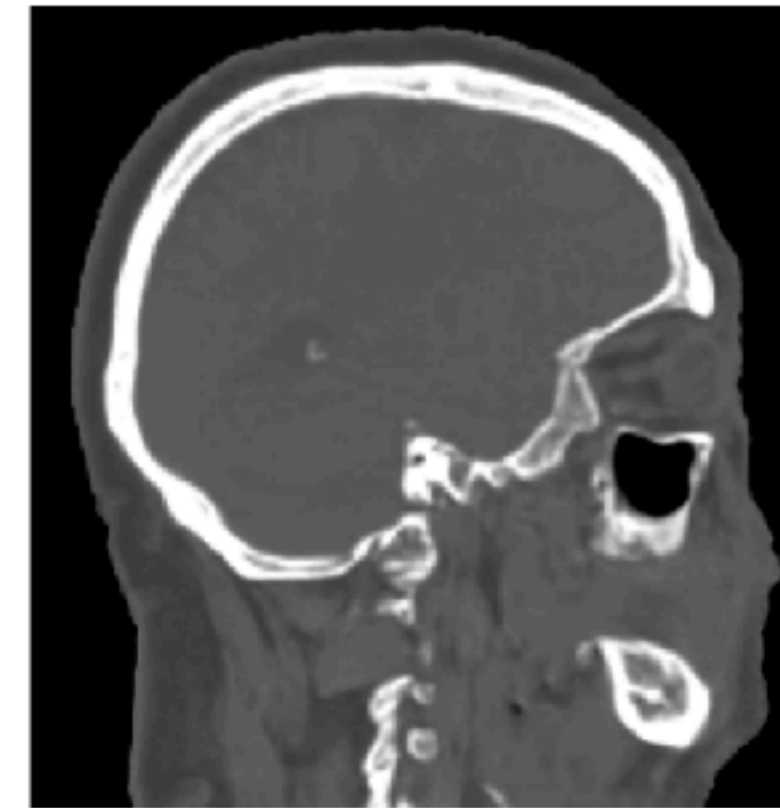


# What would it look like if...?

MRI



CT



[Wolterink et al, 2017]

Sim



“Real”



[Hoffman et al, 2018]