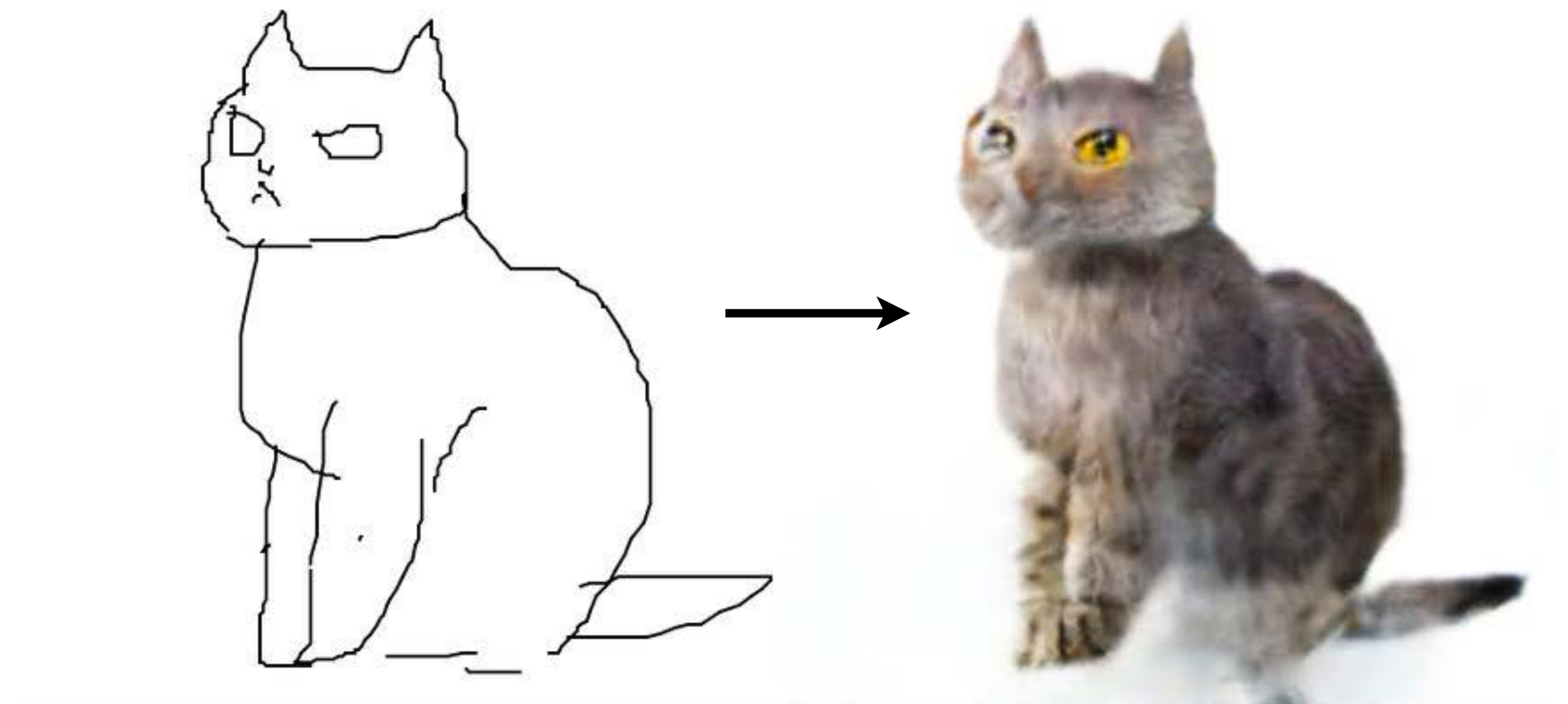# Lecture 16: Deep Generative Models III
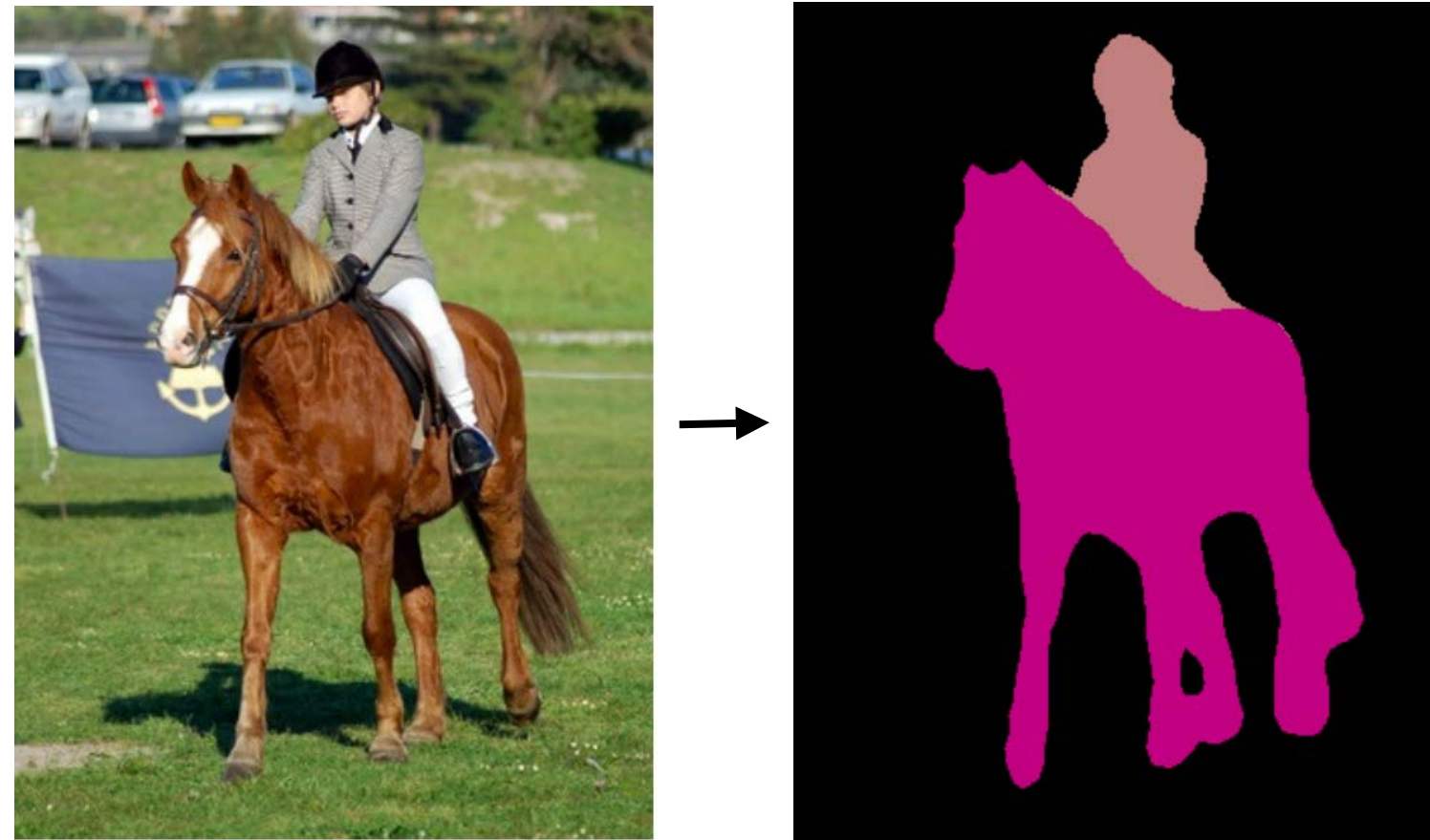
Speaker: Phillip Isola



Created with edges2cats; Pix2Pix.

# Deep generative models III

- Structured prediction

  - Image-to-image

  - Text-to-text

  - Image-to-text

  - Text-to-Image

- Unpaired translation

# Data prediction problems ("structured prediction")
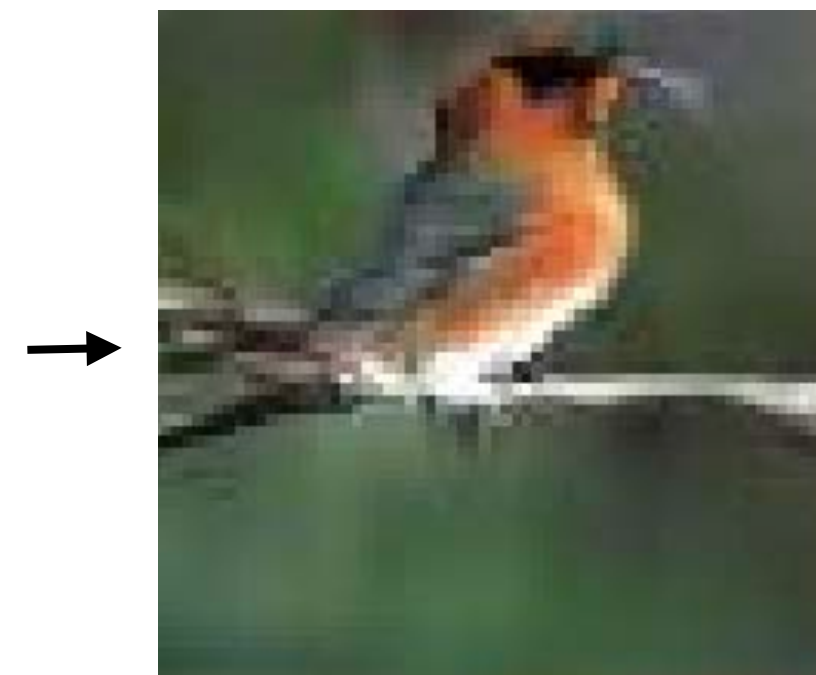
## Semantic segmentation



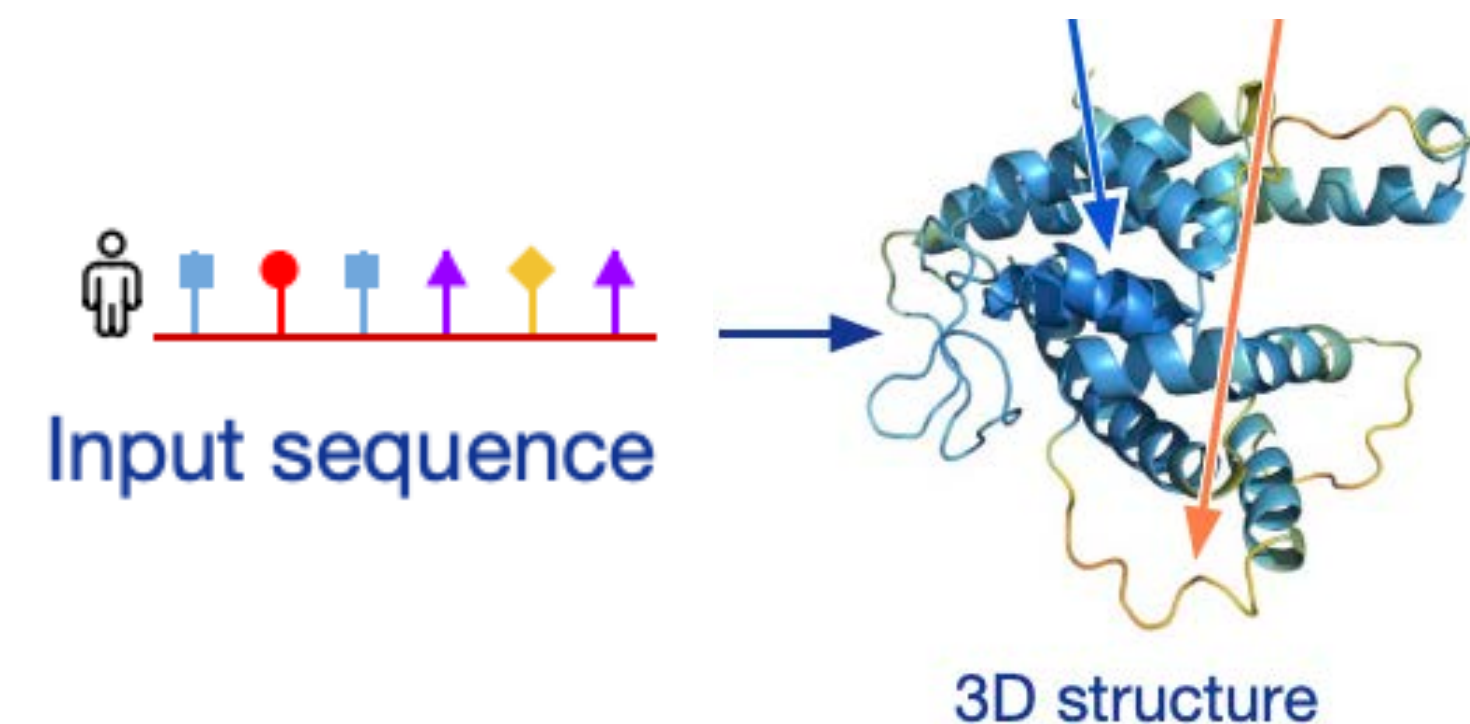[Long et al. 2015, …]

## Automatic speech recognition



→ `computers that talk`

[Glass 2007, …]

## Text-to-photo

"this small bird has a pink breast and crown…" →



[Reed et al. 2014, …]

## Protein folding



Input sequence

3D structure
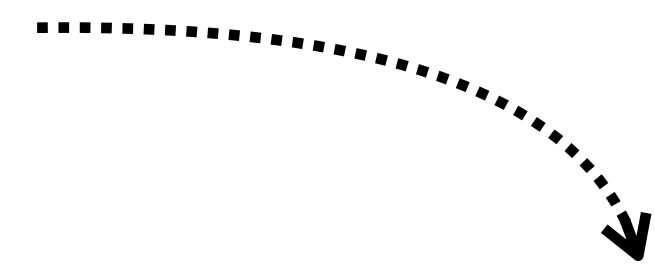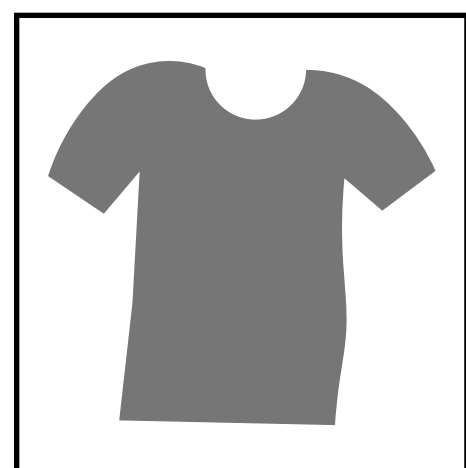
[AlphaFold, Jumper et al. 2021, …]

# Structured prediction
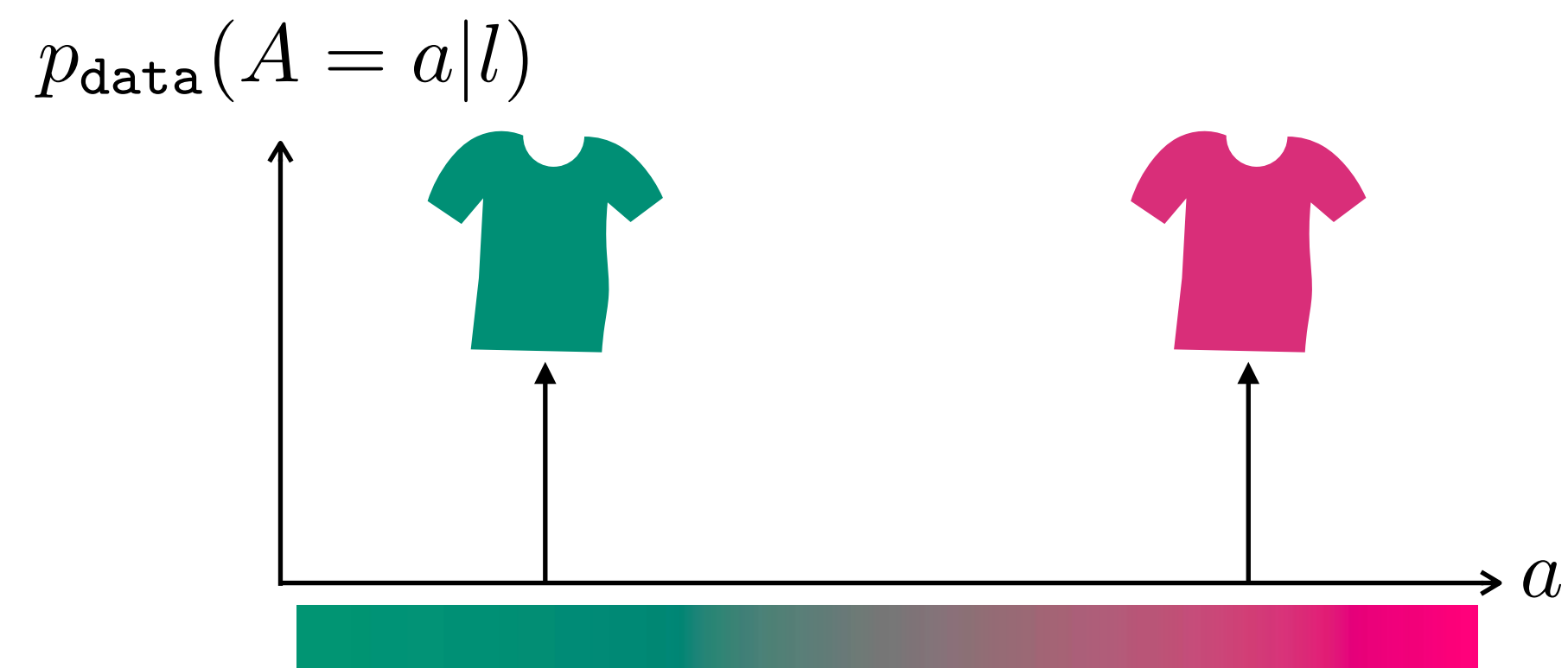
**X** is high-dimensional

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

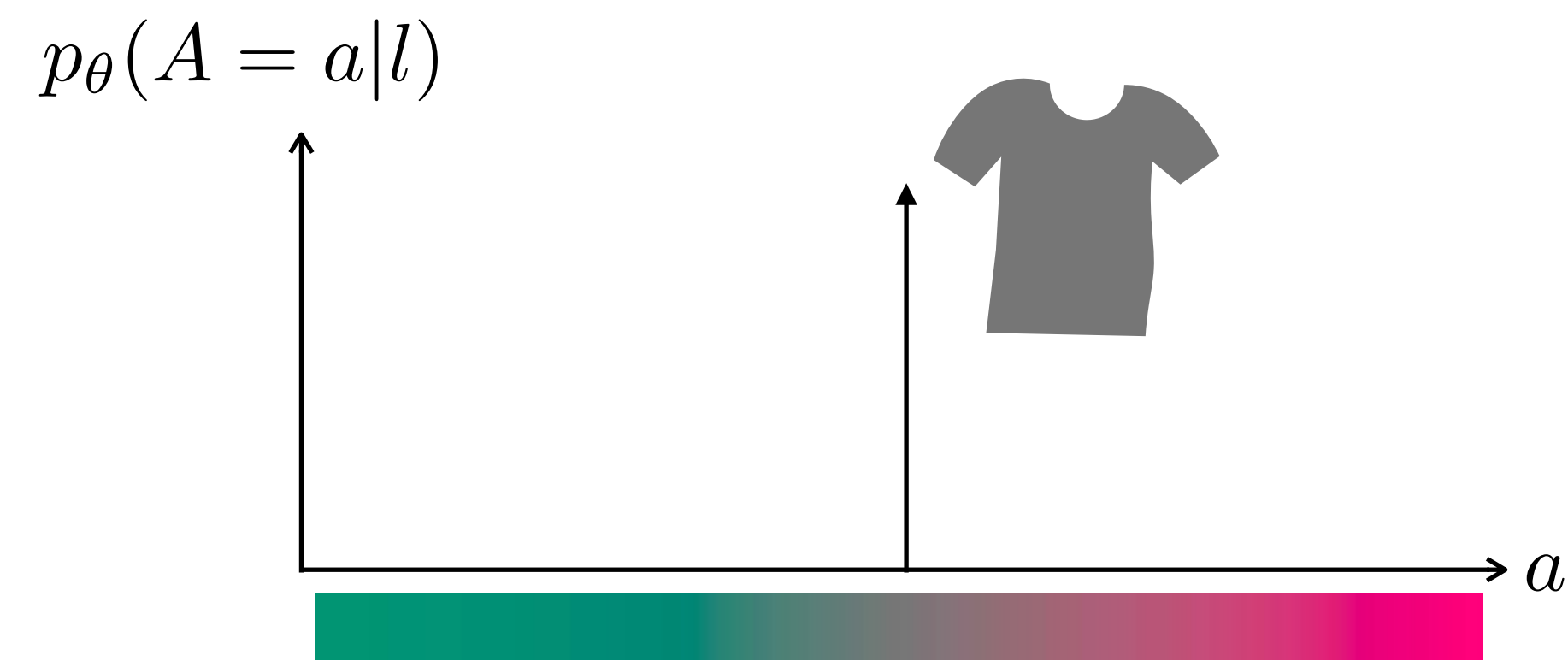In deep learning this is usually what we are interested in

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

observation

$p_{\mathtt{data}}(A = a|l)$
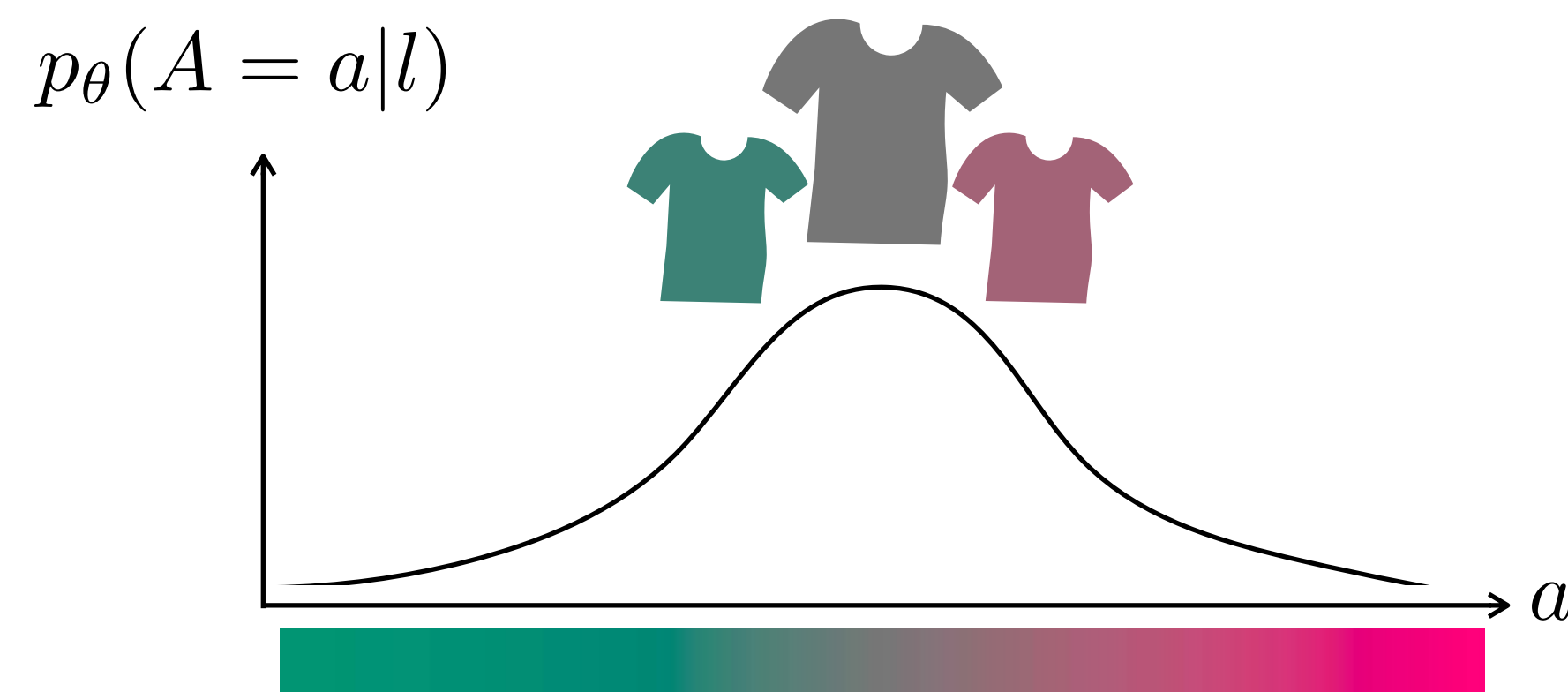
$a$

a) true data distribution

$p_\theta(A = a|l)$

$a$

b) point prediction (aka regression)

$p_\theta(A = a|l)$

$a$

d) Gaussian predictive distribution

$p_\theta(A = a|l)$

$a$

c) softmax distribution (aka classification)

5

Generative models have two important properties for structured prediction:
1. They can model a multimodal distribution
2. They can model joint dependences between multidimensional predictions

Real or fake?

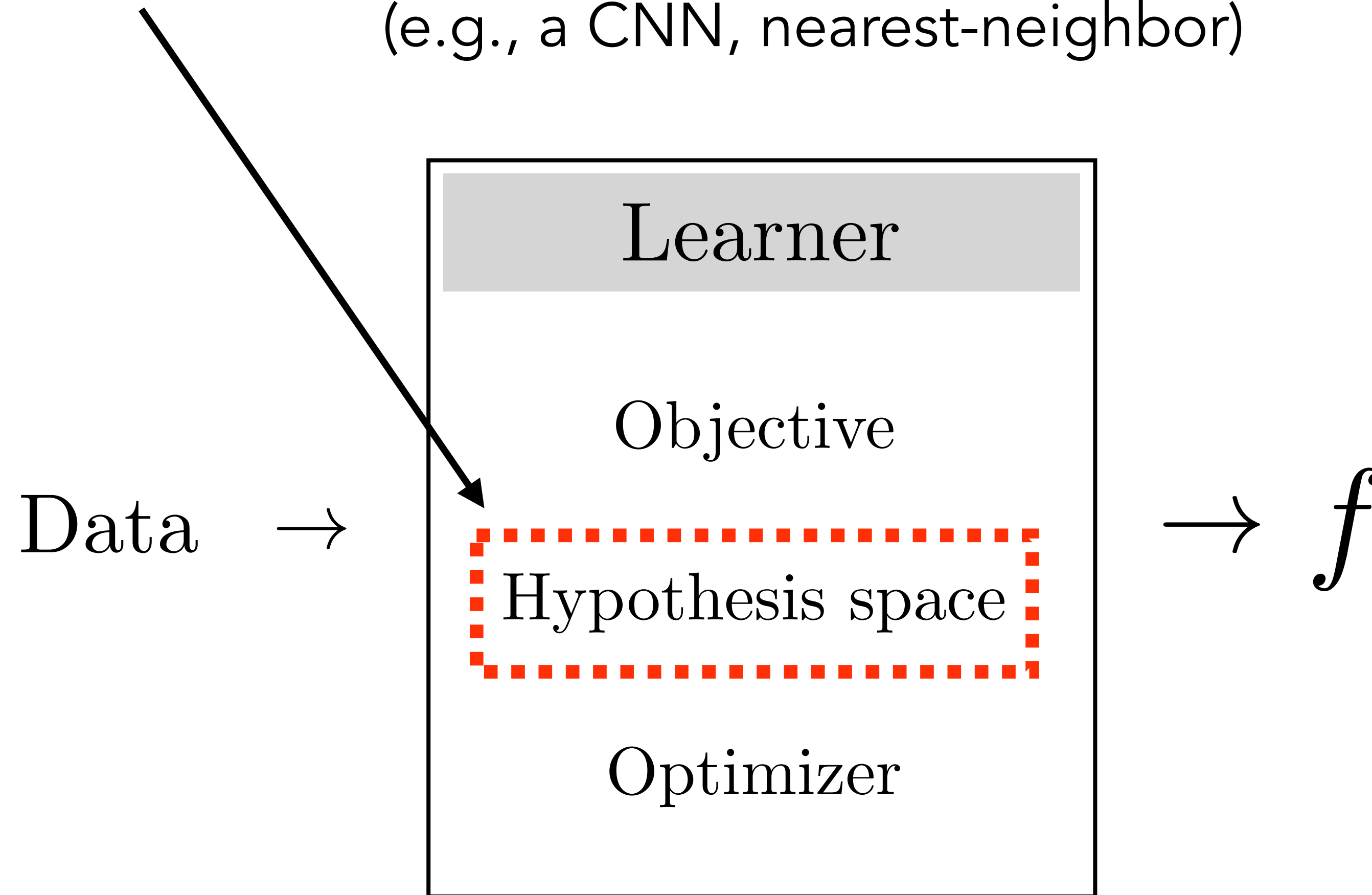# Real or fake?
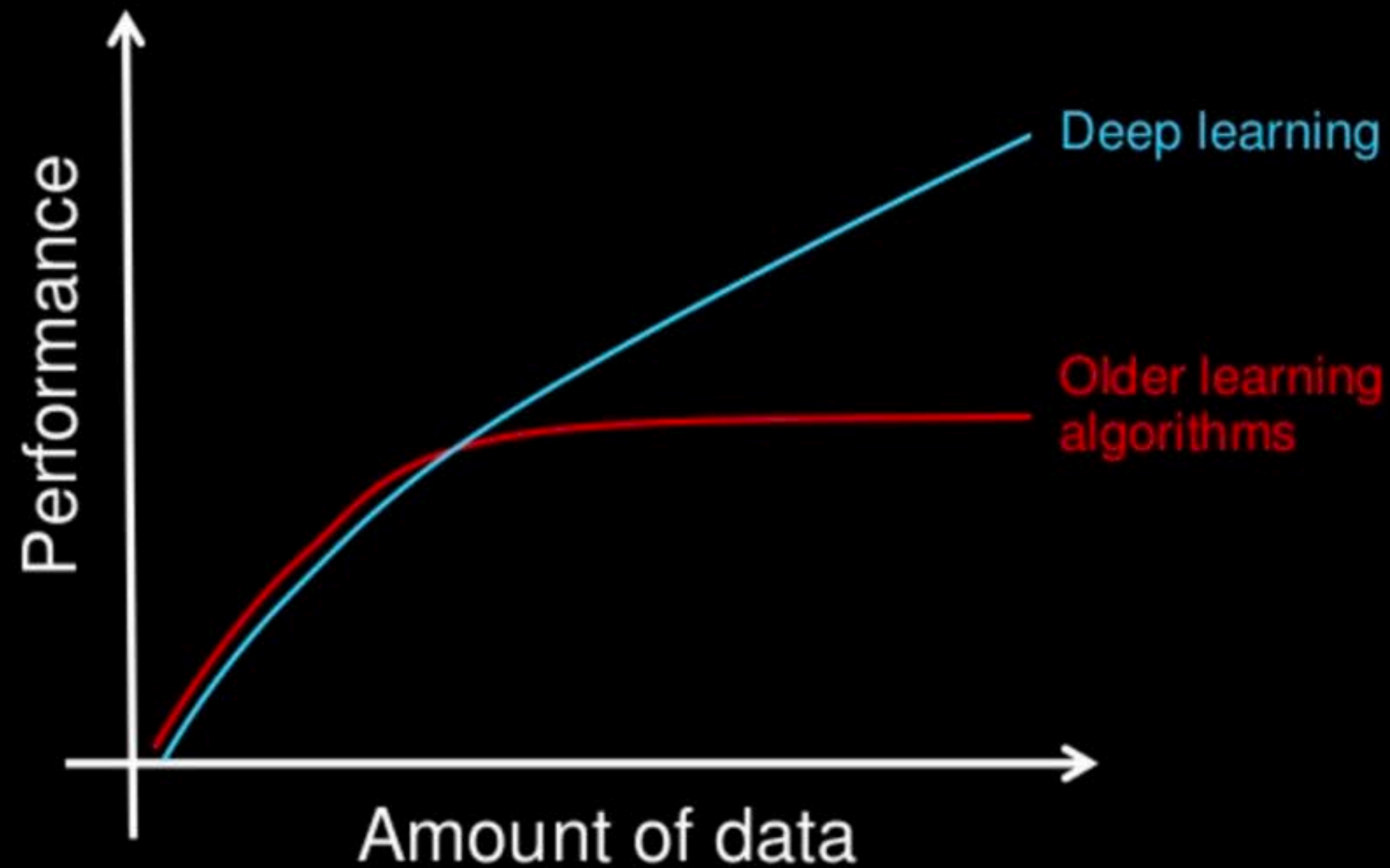
# Deep learning in 2012

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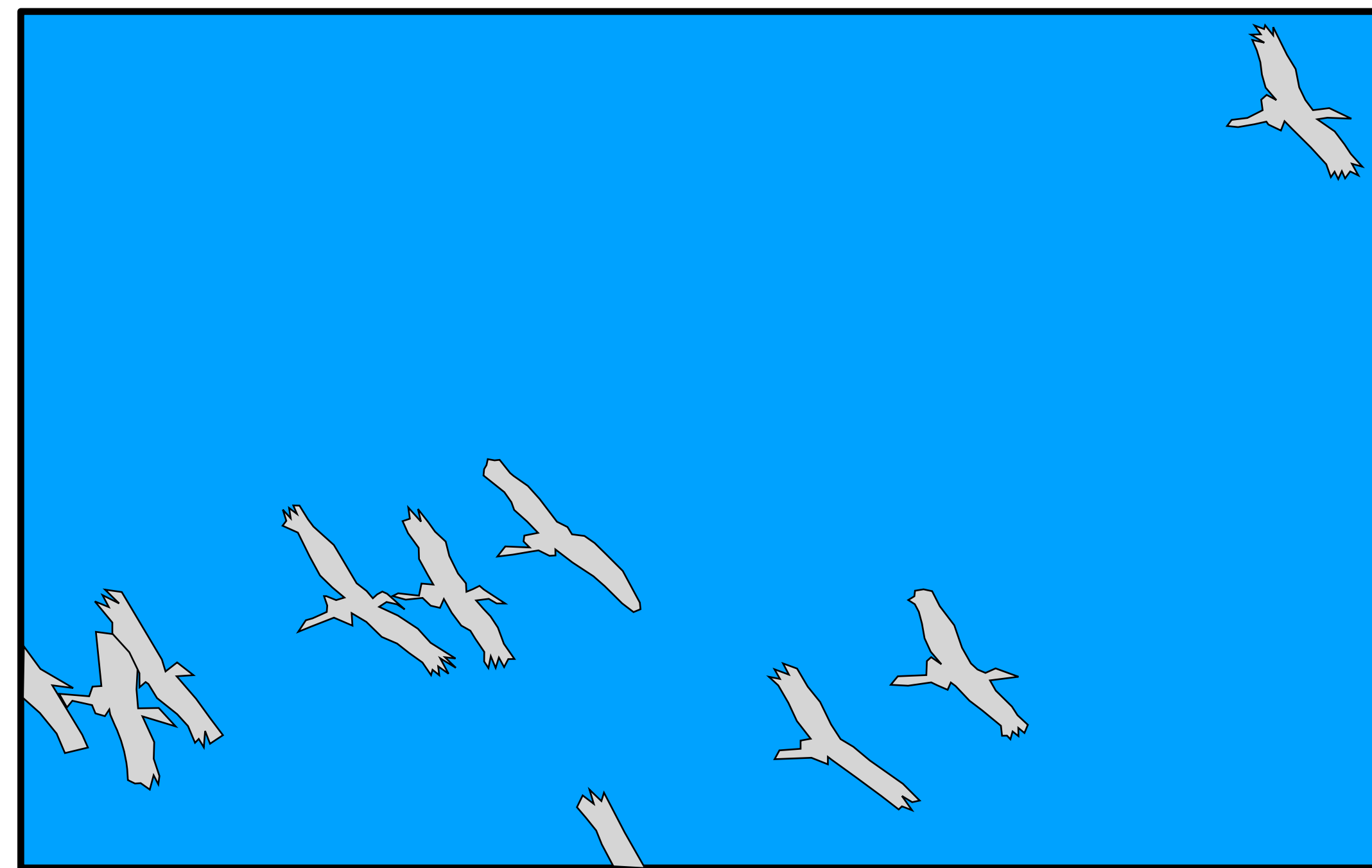(e.g., a CNN, nearest-neighbor)

$$\text{Data} \rightarrow$$

| Learner |
| --- |
| Objective |
| Hypothesis space |
| Optimizer |

$$\rightarrow f$$

# Why deep learning

How do data science techniques scale with amount of data?

[Slide credit: Andrew Ng]

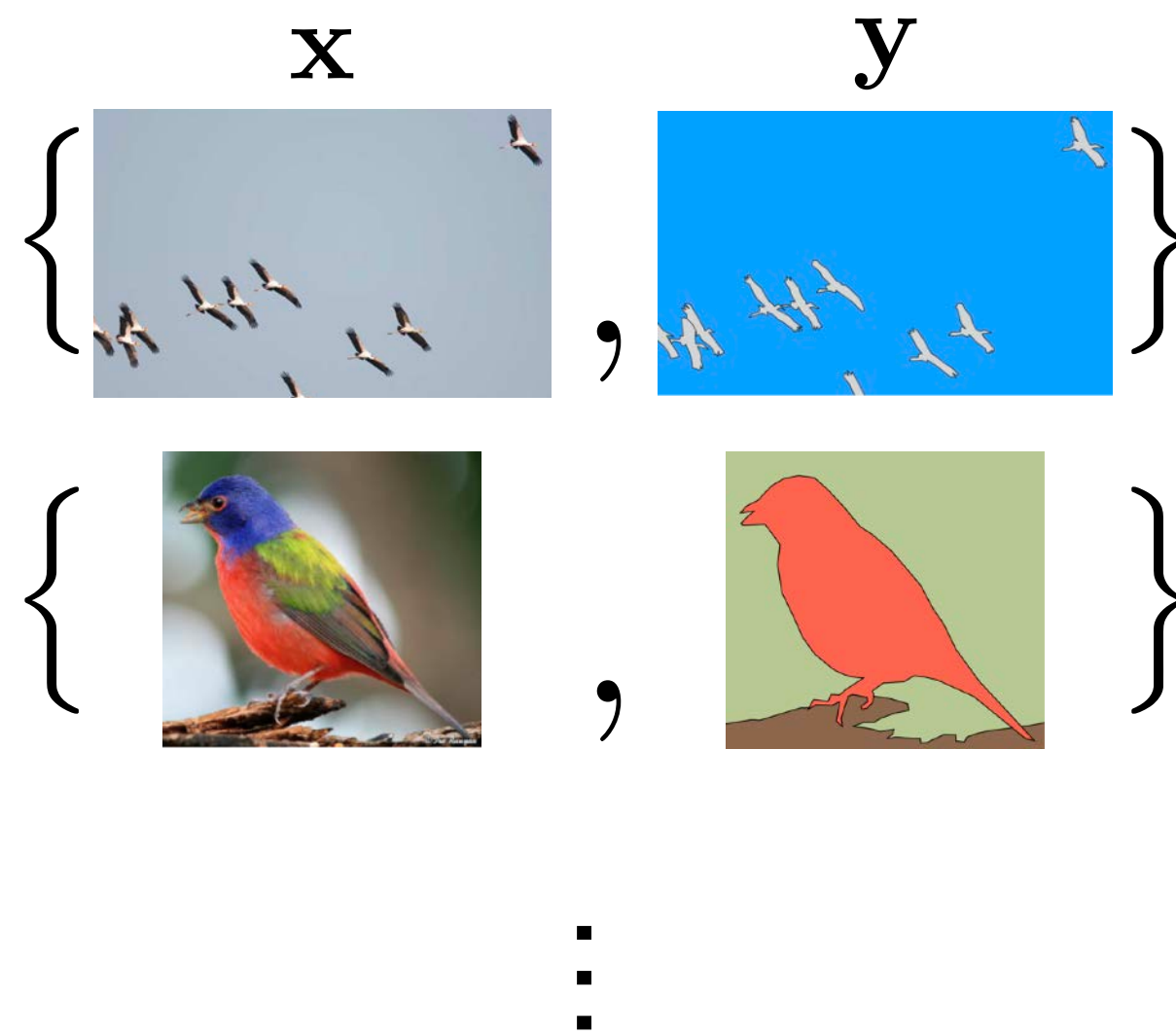[Photo credit: Fredo Durand]

(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
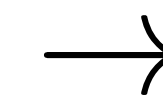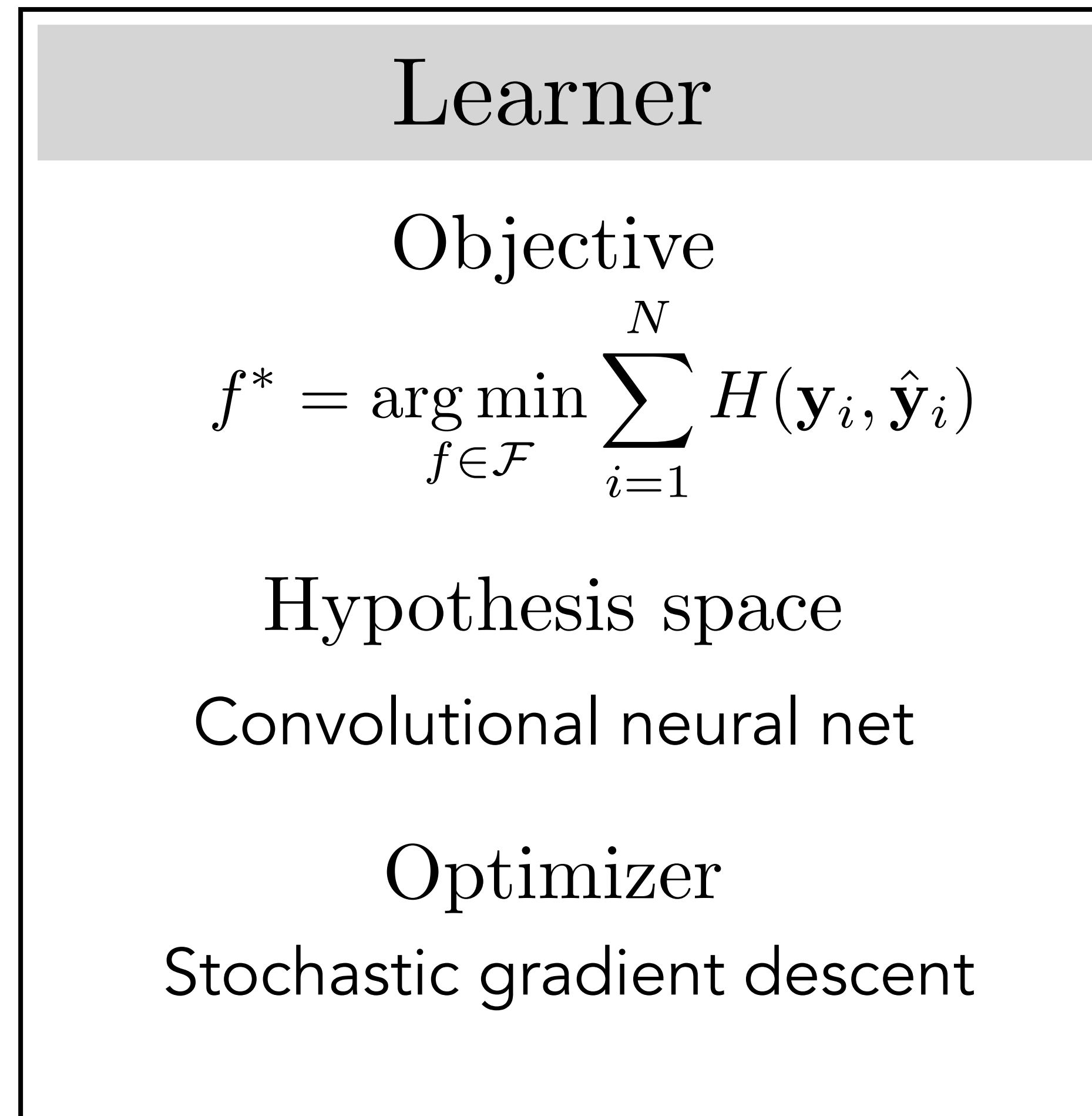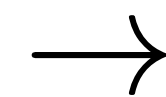(loss)

# Semantic Segmentation

## Data



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

$\rightarrow$
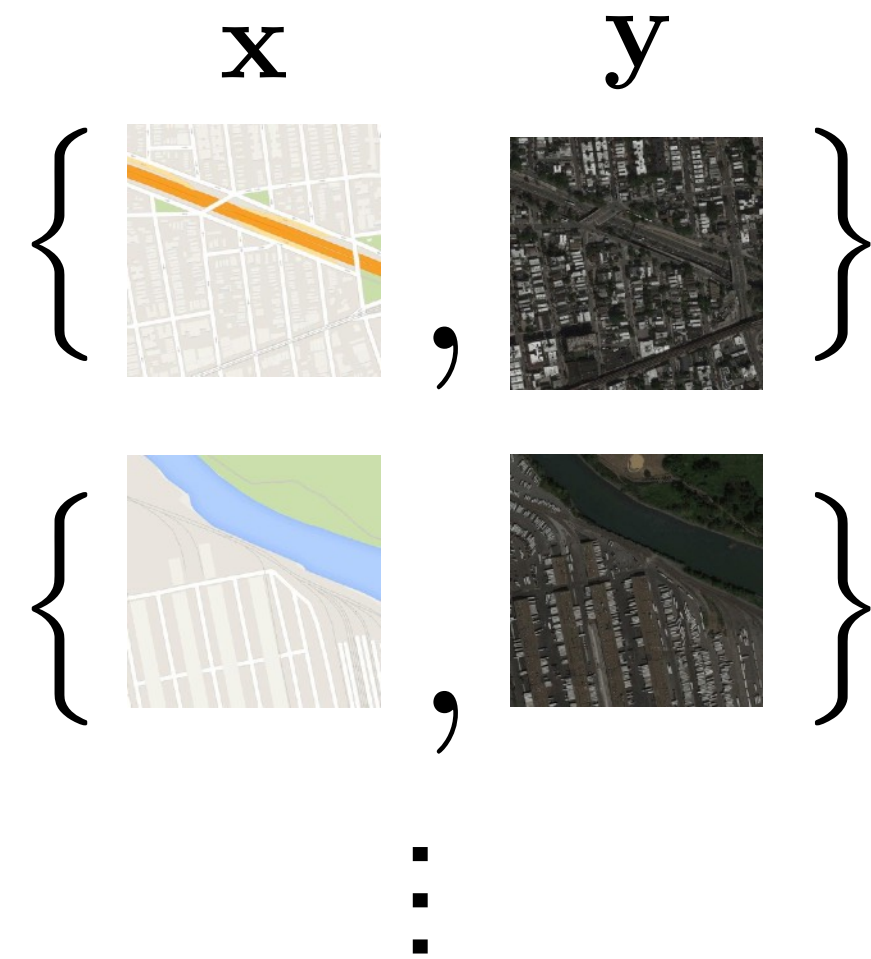
### Learner

#### Objective

$$f^* = \arg\min_{f \in \mathcal{F}} \sum_{i=1}^{N} H(\mathbf{y}_i, \hat{\mathbf{y}}_i)$$

#### Hypothesis space
Convolutional neural net

#### Optimizer
Stochastic gradient descent

$\rightarrow$  $f$

# Sat2Map

## Data

$$\mathbf{x} \qquad \mathbf{y}$$
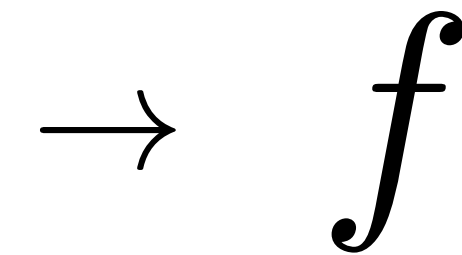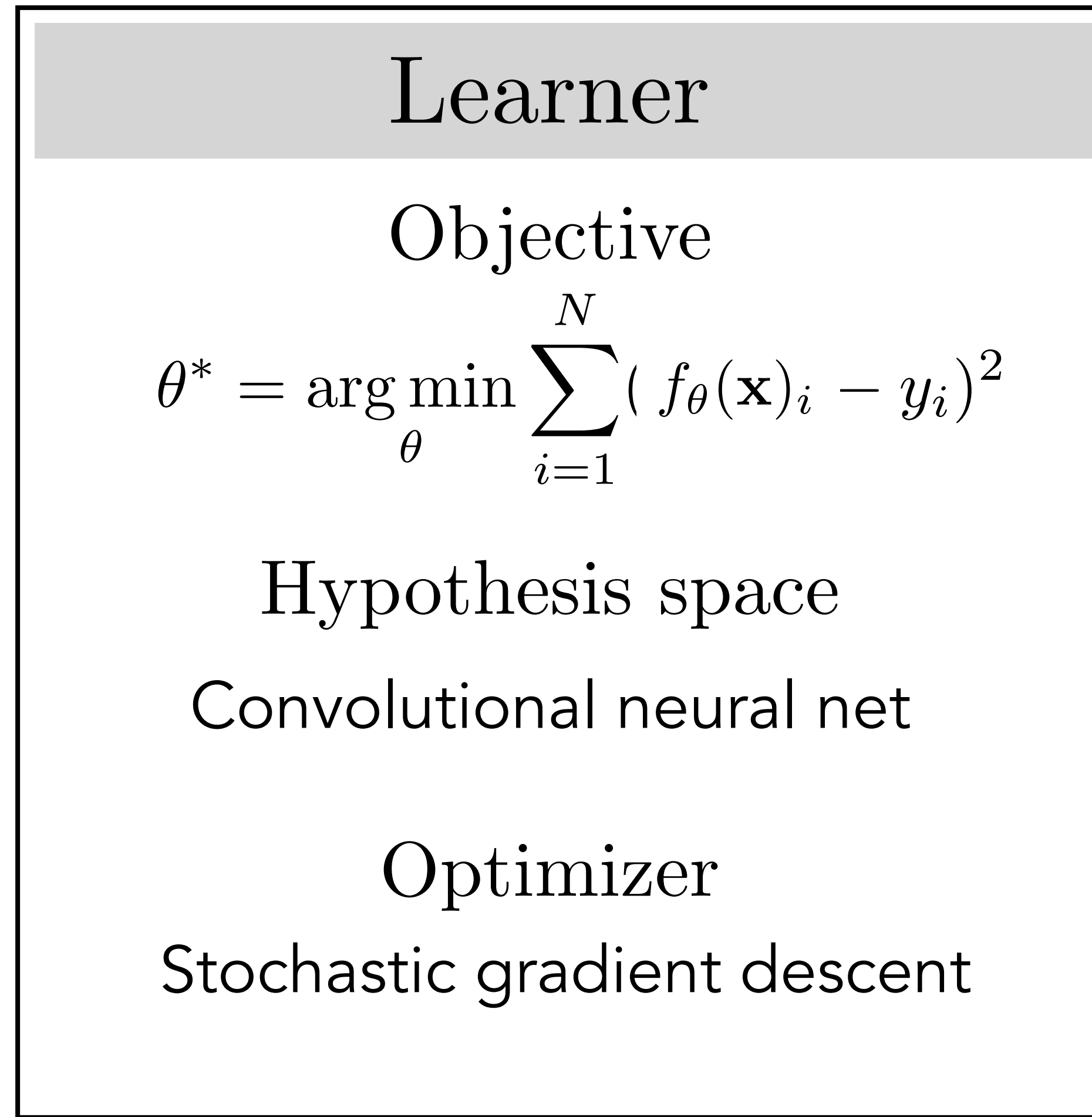


$$\mathbf{x} \in \mathbb{R}^{H \times W \times 3}$$

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

$\rightarrow$

## 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

$\rightarrow \quad f$

# Input

# Deep net output

# Structured prediction

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



Data $\rightarrow$

Learner

Objective

Hypothesis space

Optimizer

$\rightarrow f$

# Image-to-image

→ with conditional GANs

→ with conditional VAEs

# Conditional GAN

$\mathbf{x}$

$g_\theta$

$g_\theta\left(\mathbf{x}\right)$

Generator

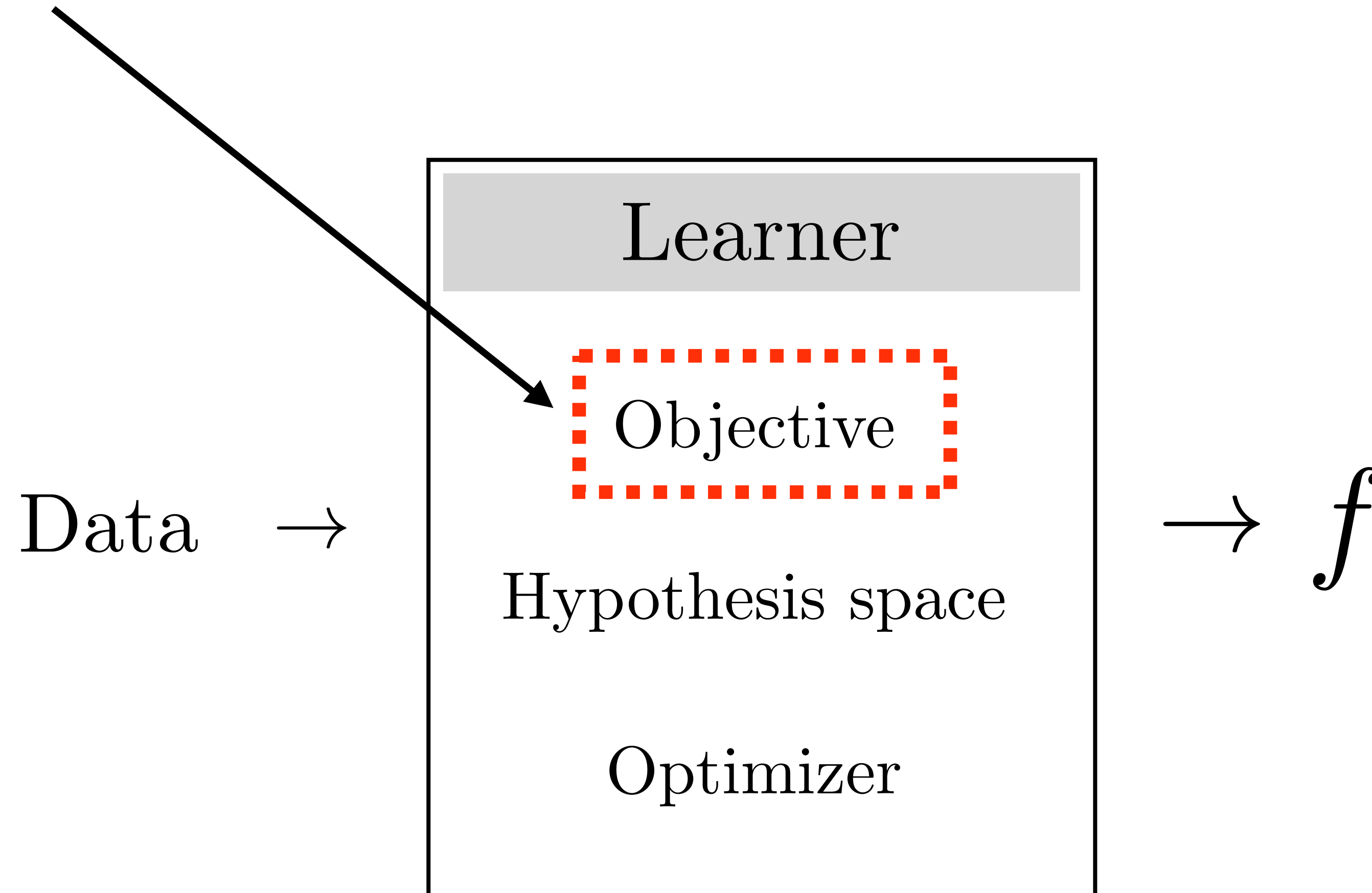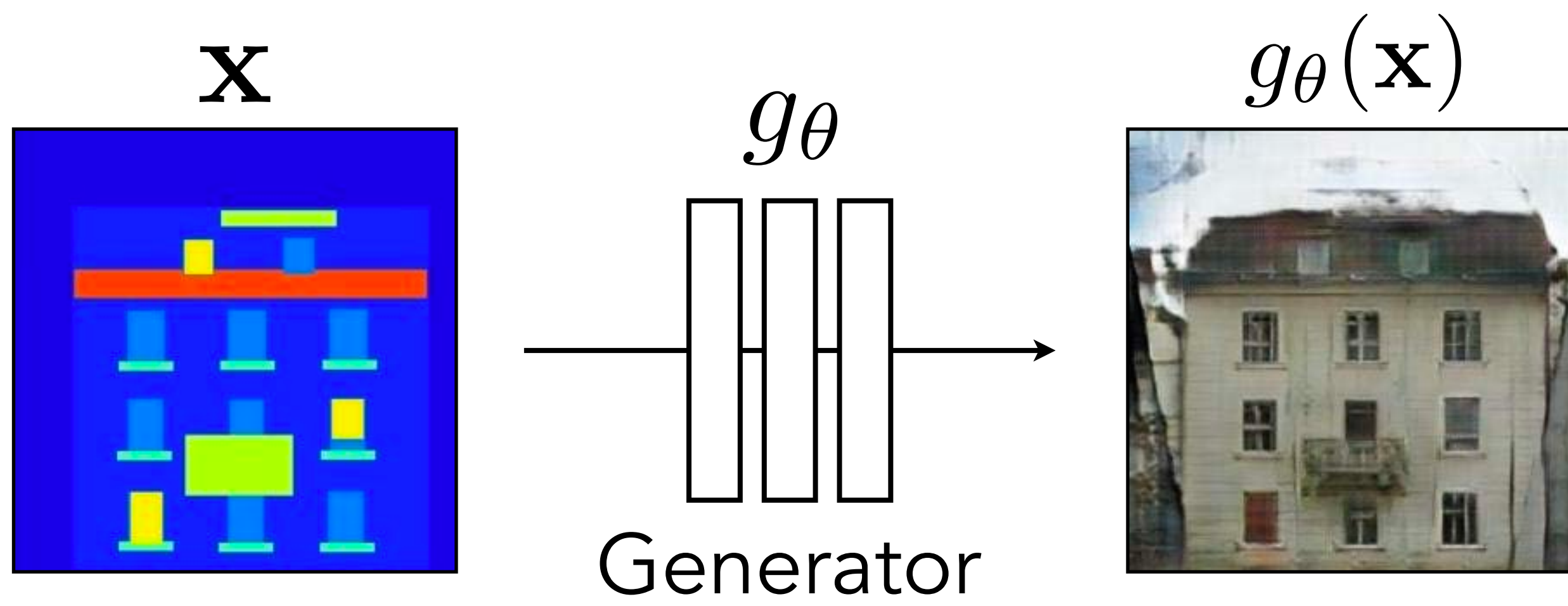For example: pix2pix [Isola et al. 2017]

**x**      $g_\theta$      $g_\theta(\mathbf{x})$      $d_\phi$

Generator      Discriminator      real or synthetic?

g tries to synthesize fake images that fool d

d tries to identify the fakes

**x**   $g_\theta$   $g_\theta(\mathbf{x})$   $d_\phi$   synthetic (0.9)

**y**   $d_\phi$   real (0.1)

$$d_\phi^* = \arg\max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}} \left[ \log d_\phi(g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{y})) \right]$$

$$\mathbf{x} \qquad g_\theta \qquad g_\theta(\mathbf{x}) \qquad d_\phi$$

real or synthetic?

g tries to synthesize fake images that *fool* d:

$$g_\theta^* = \boxed{\underset{\theta}{\arg\min}} \; \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(g_\theta(\mathbf{x}))]$$

g tries to synthesize fake images that *fool* the *best* d:

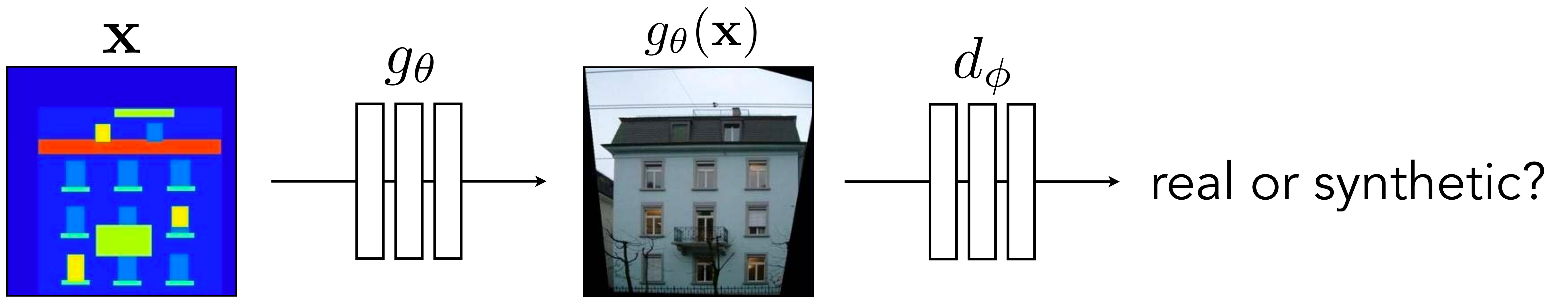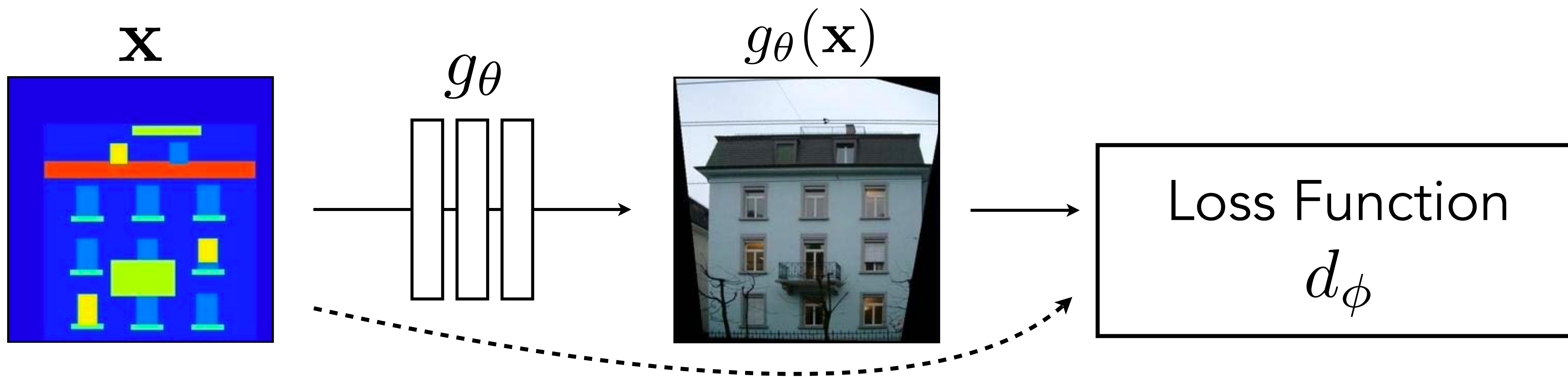$$\arg \min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_{\phi}(g_{\theta}(\mathbf{x})) + \log(1 - d_{\phi}(\mathbf{y}))]$$

$\mathbf{x}$    $g_\theta$    $g_\theta(\mathbf{x})$    Loss Function $d_\phi$

g's perspective: d is a loss function.

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

$\mathbf{x}$

$g_\theta$

$g_\theta(\mathbf{x})$

$d_\phi$

real or synthetic?

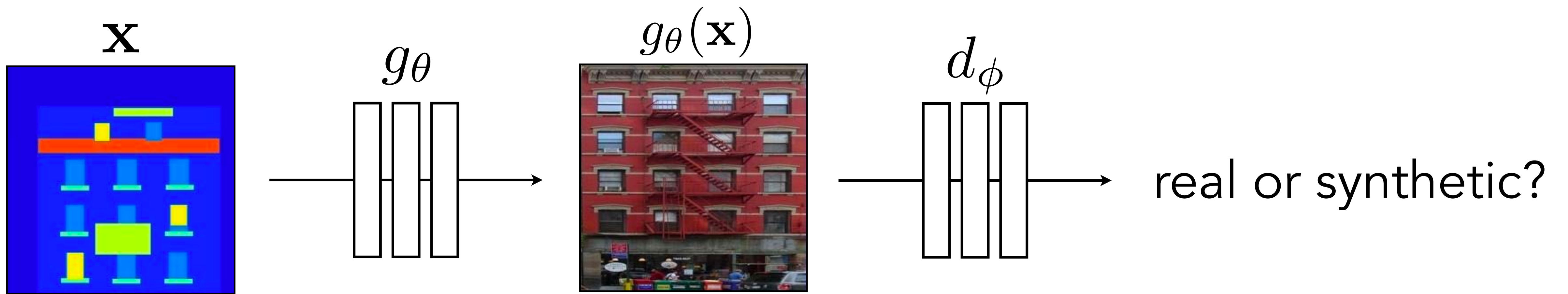$$\arg \min_\theta \max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{y}))]$$

$$\arg \min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{y}))]$$

$$\arg \min_\theta \max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{y}))]$$

**x**

$g_\theta$
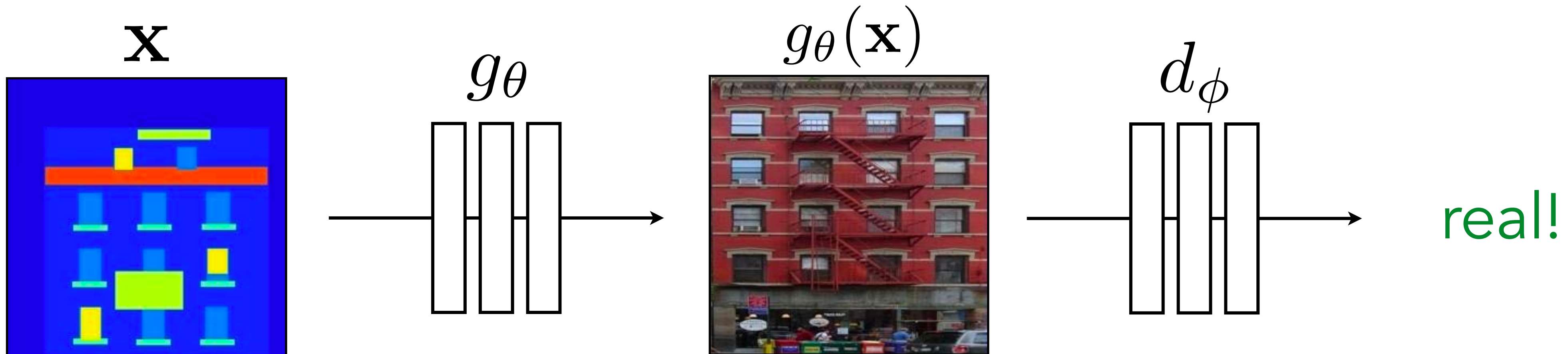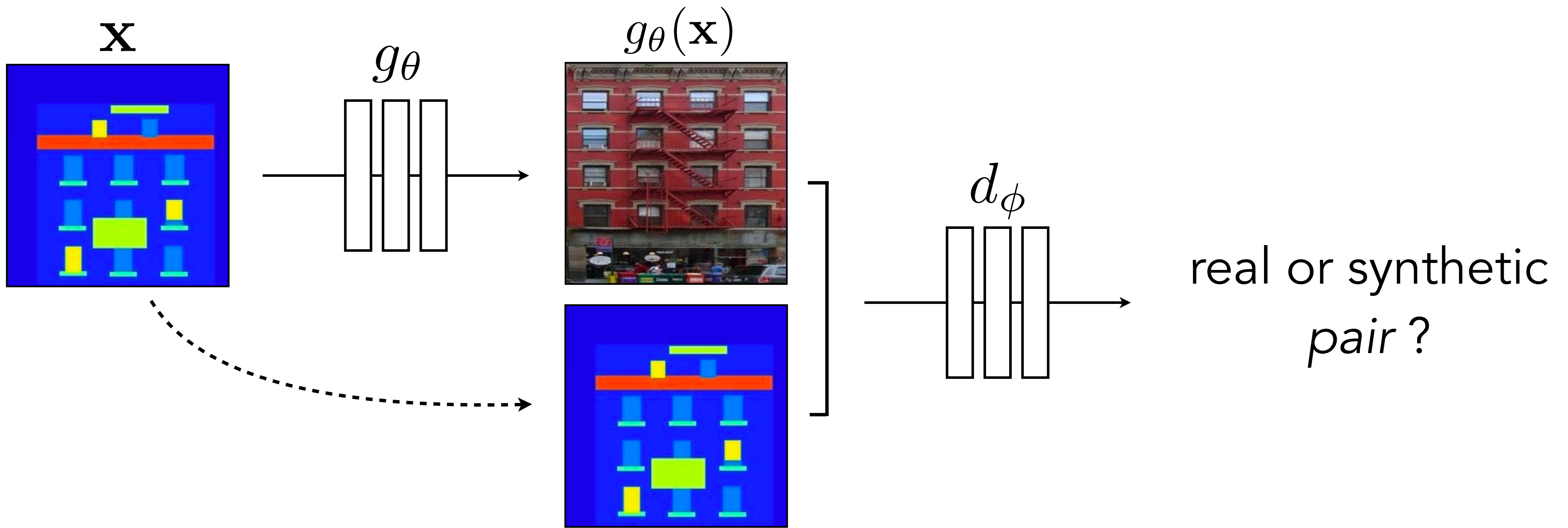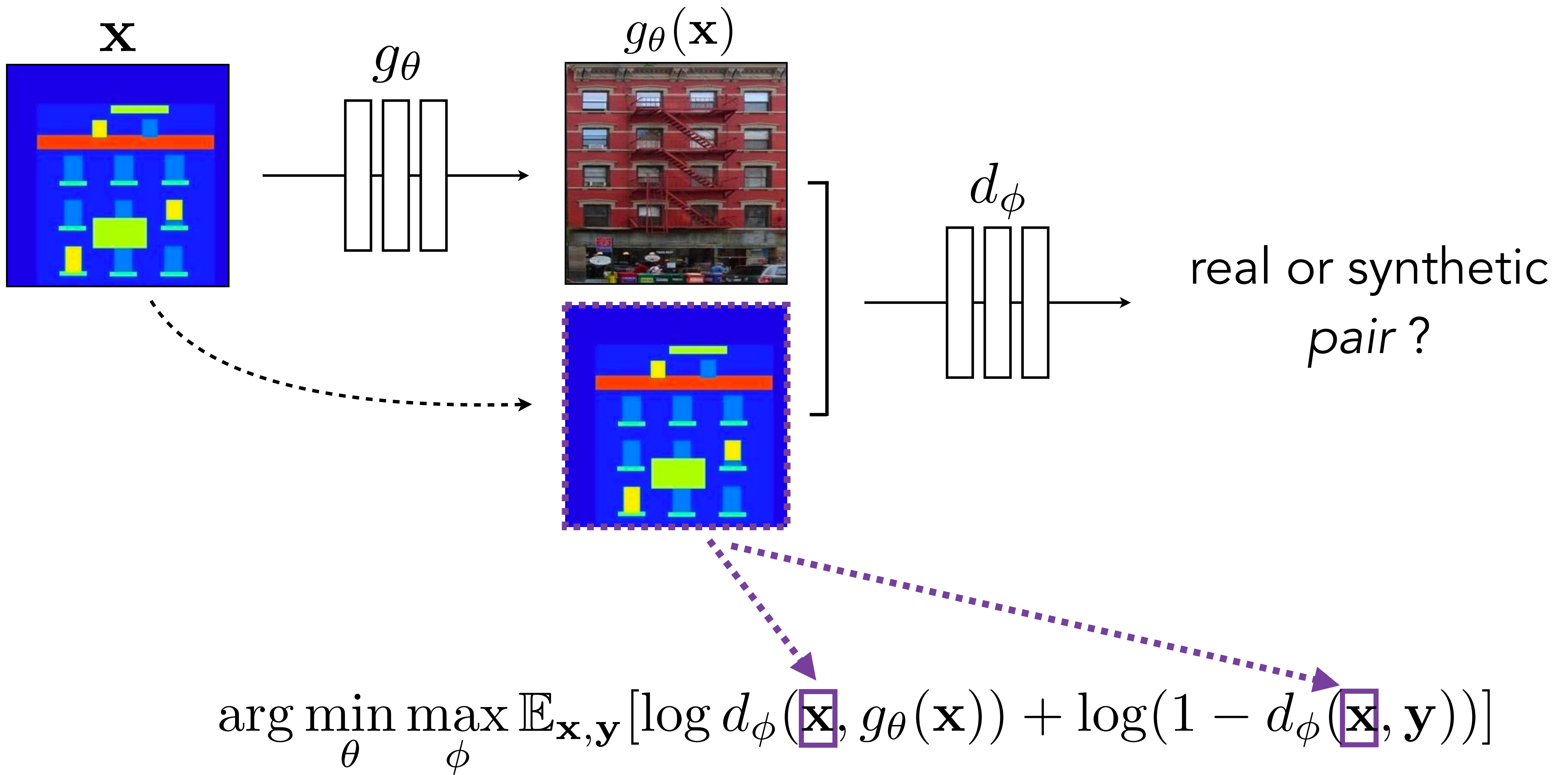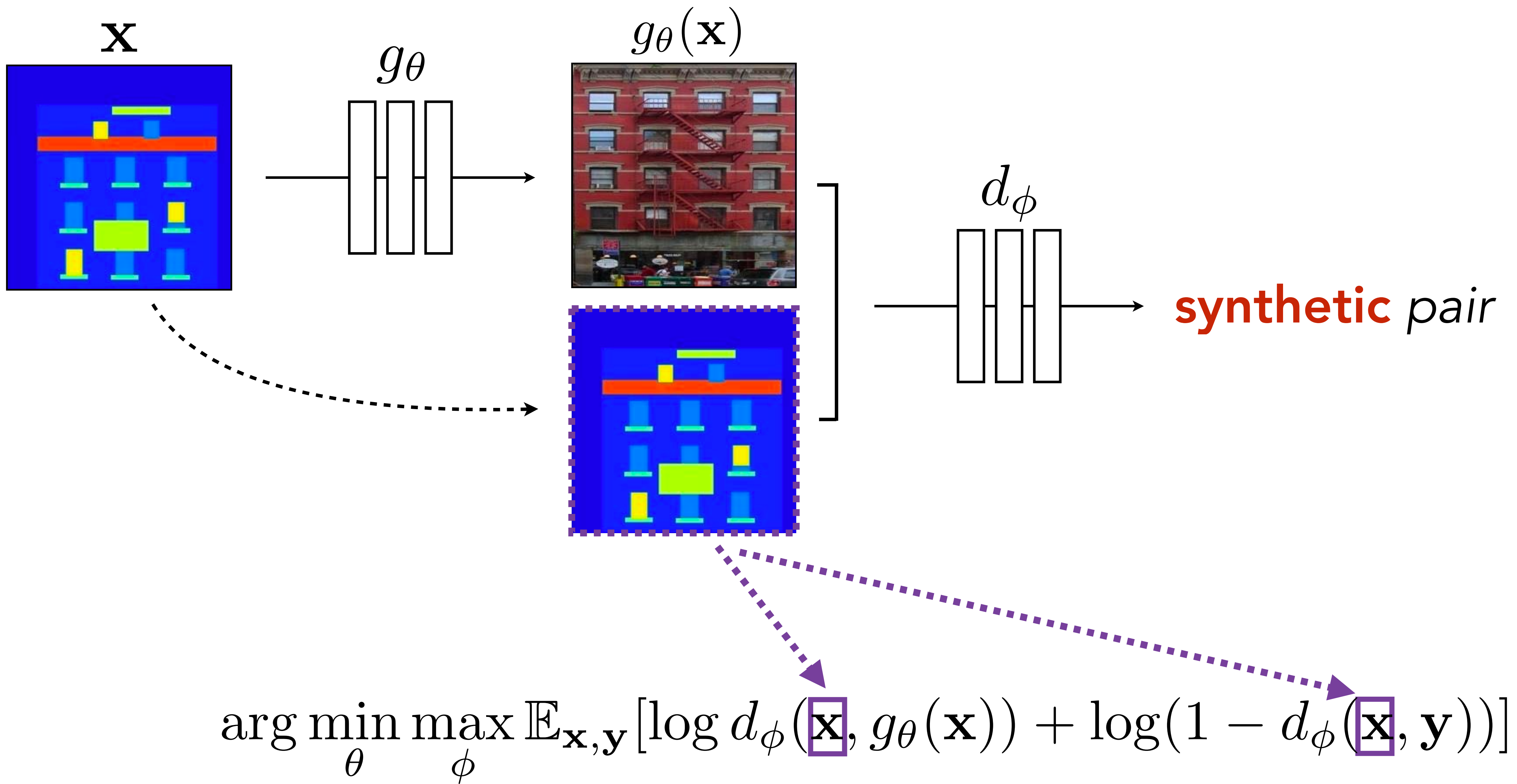
$g_\theta(\mathbf{x})$

$d_\phi$

real or synthetic
*pair* ?

$$\arg\min_\theta \max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(\mathbf{x}, g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{x}, \mathbf{y}))]$$

$$\arg \min_{\theta} \max_{\phi} \mathbb{E}_{\mathbf{x},\mathbf{y}} [\log d_\phi(\mathbf{x}, g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{x}, \mathbf{y}))]$$
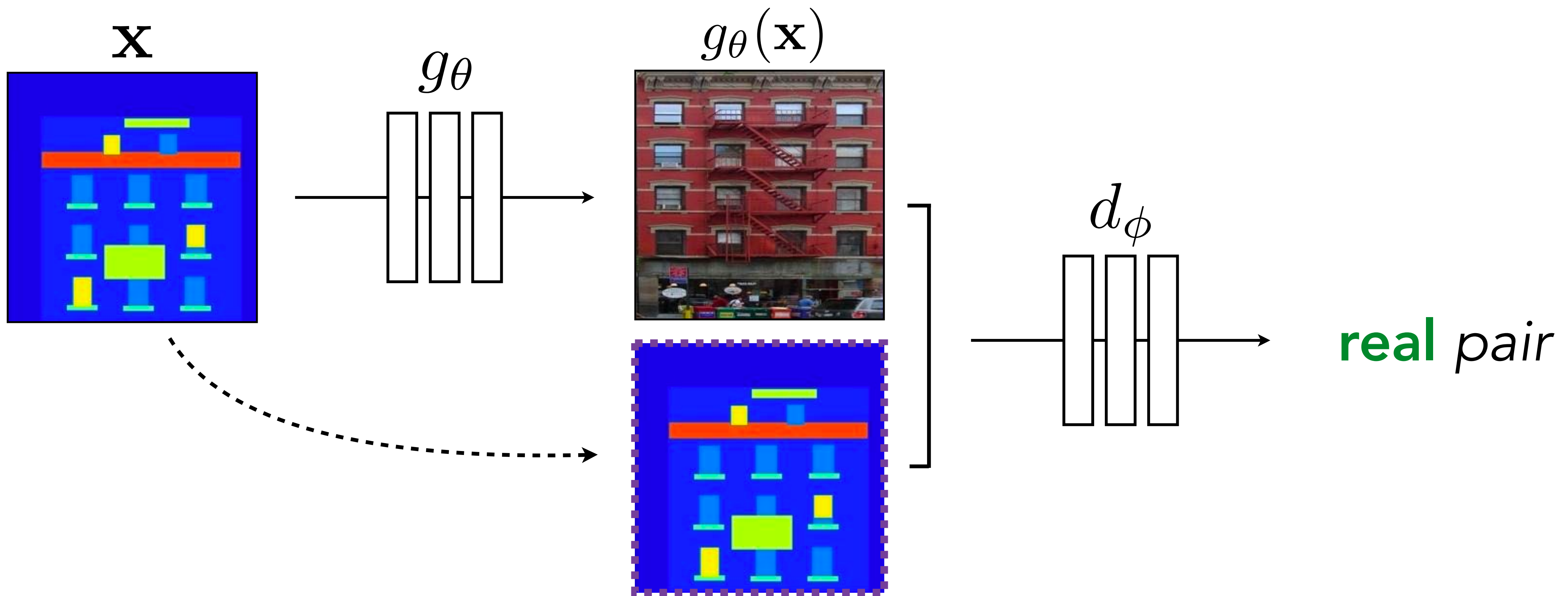
$$\mathbf{x} \quad g_\theta \quad g_\theta(\mathbf{x}) \quad d_\phi \quad \textbf{real } \textit{pair}$$

$$\arg \min_\theta \max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(\mathbf{x}, g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{x}, \mathbf{y}))]$$

$\mathbf{x}$

$g_\theta$

$g_\theta(\mathbf{x})$

$d_\phi$

real or synthetic
*pair* ?
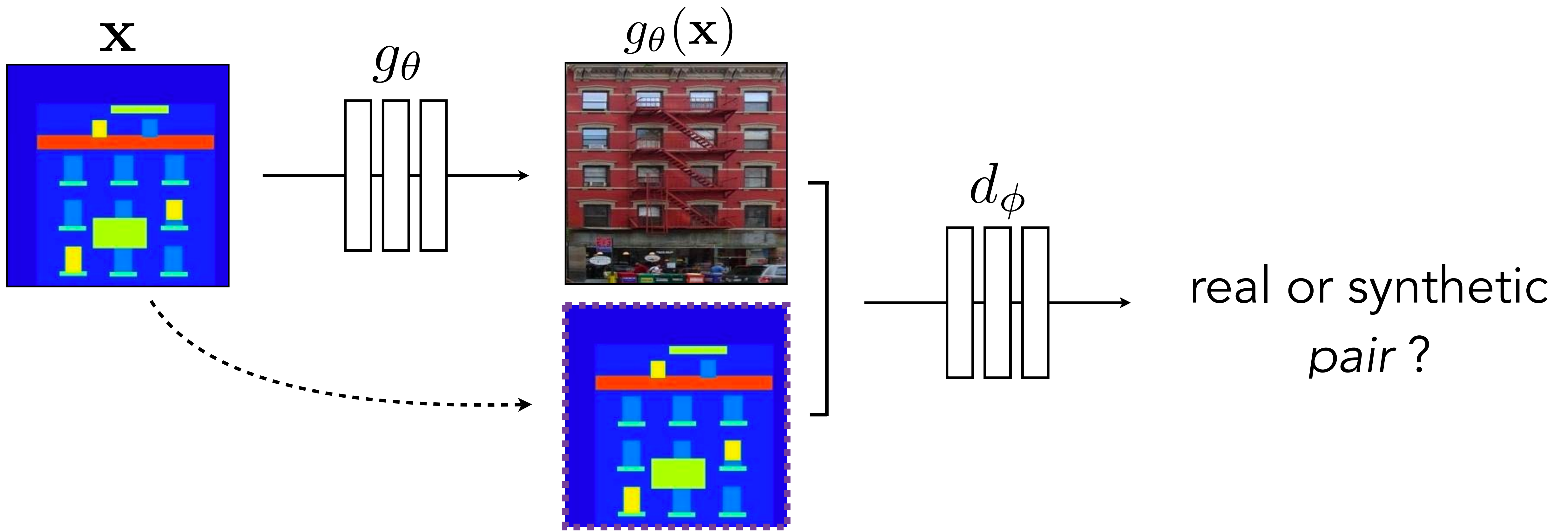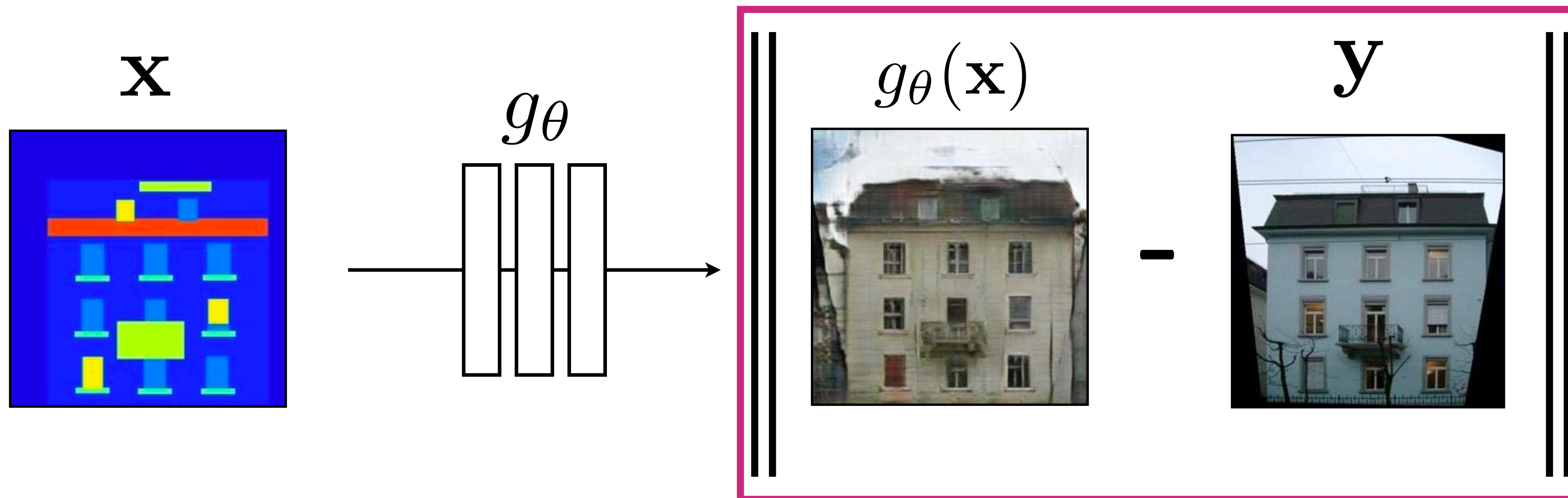
$$\arg\min_\theta \max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(\mathbf{x}, g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{x}, \mathbf{y}))]$$

# Training Details: Loss function

$$g_\theta^* = \arg \min_\theta \max_\phi \mathcal{L}_{\texttt{cGAN}}(\theta, \phi) + \boxed{\lambda \mathcal{L}_{\texttt{L1}}(\theta)}$$
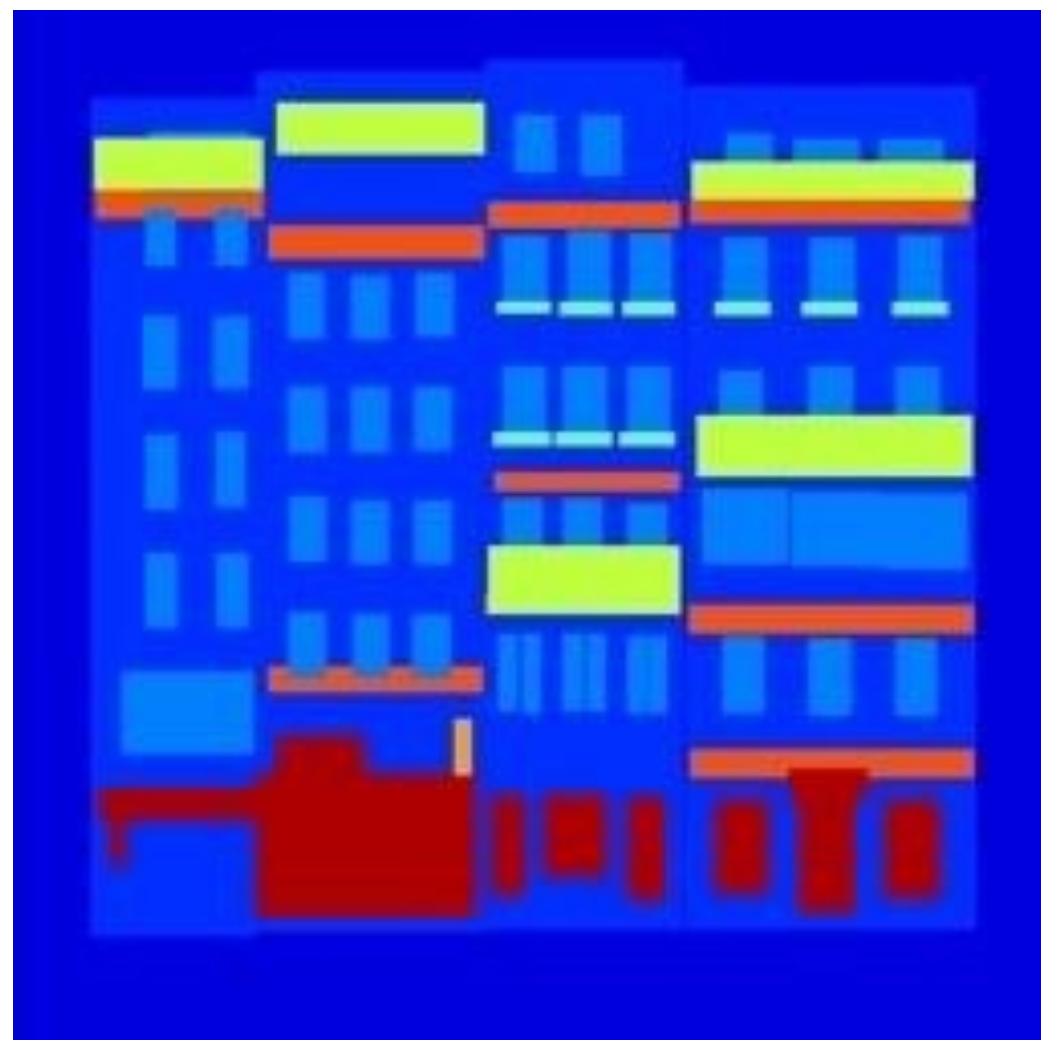


Stable training + fast convergence

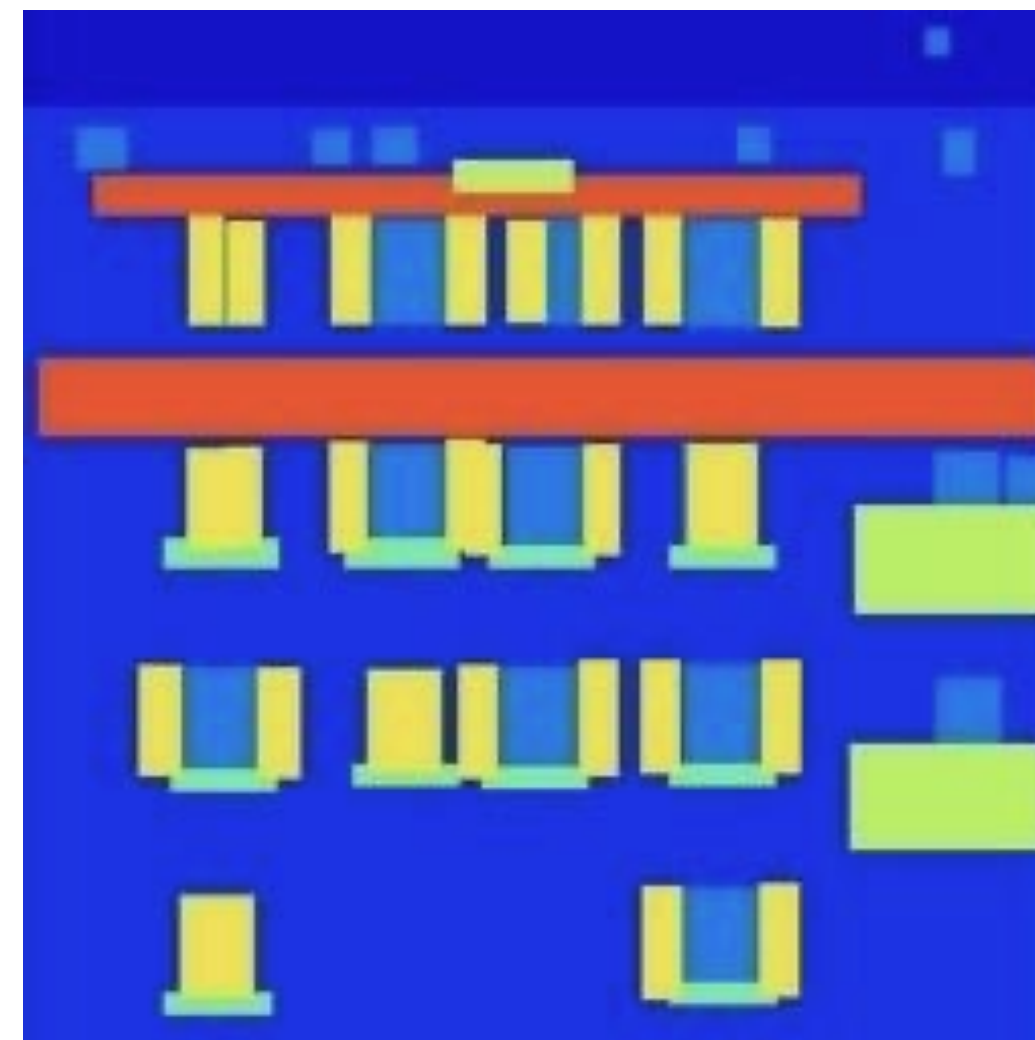[c.f. Pathak et al. CVPR 2016]
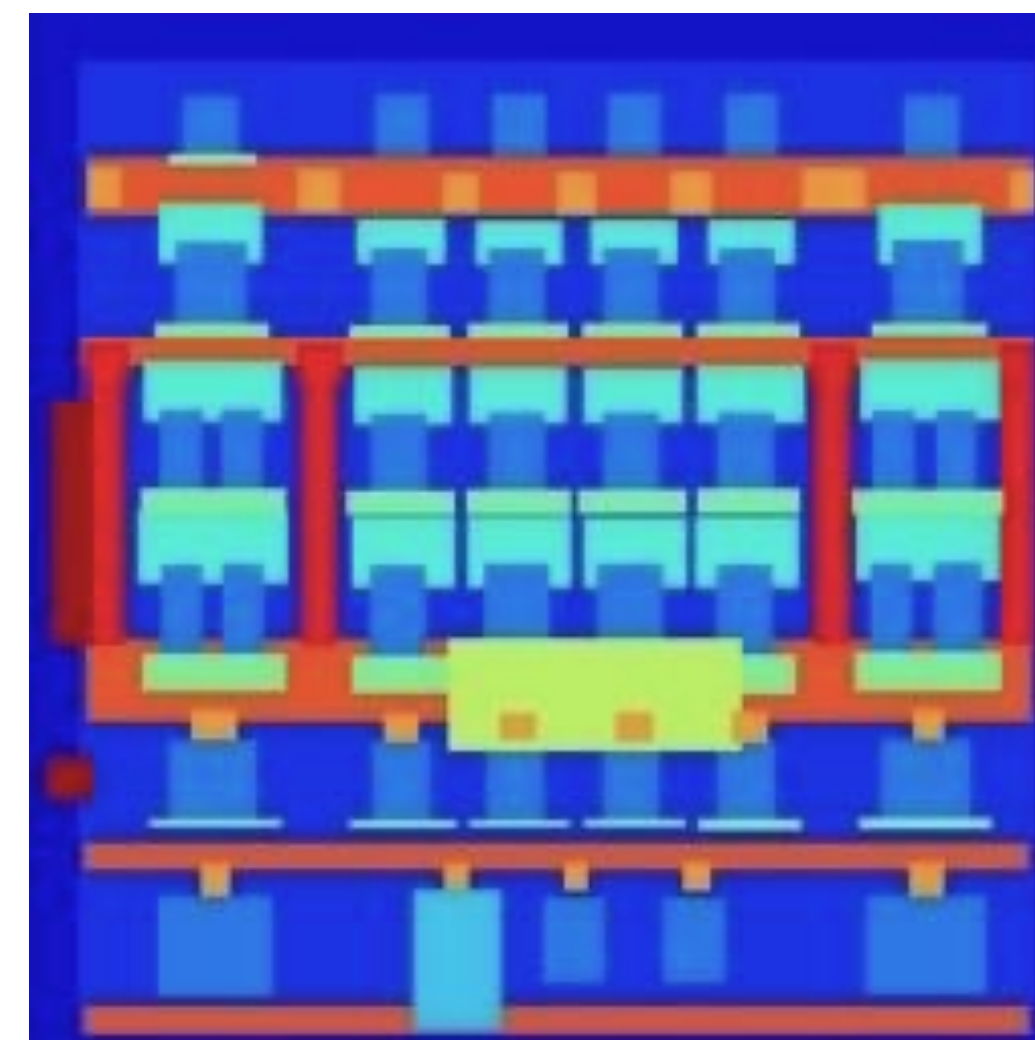
# Labels → Facades

Data from [Tylecek, 2013]

32

# Patch Discriminator



N pixels

N pixels

$D$

**1**/**0**

y

Output

$1 \times 1$      $16 \times 16$      $70 \times 70$      $286 \times 286$

Discriminator patch size

# Why deep learning



How do data science techniques scale with amount of data?

# Why structured objectives
## (cartoon)

# Why structured objectives
## (cartoon)



Performance (vertical axis)

Amount of data (horizontal axis)

DL w/ structured objective
(conditional generative models)

DL w/ unstructured objective
(e.g., least-squares regression)

Older learning algorithms

# Input

# Unstructured prediction (L1)

# Input

# Structured Prediction (cGAN)

# Conditional VAE

# Modeling multiple possible outputs



Input

Possible outputs

# Modeling multiple possible outputs

$\mathbf{x}$

$G$

$G(\mathbf{x})$



$z \sim \mathcal{N}(\vec{0}, 1)$

Context Frame

Time

GT

Finn et al. (2016)

Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

**x**
Observation

**y**
Target

**ŷ**
Reconstruction

**z**

Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

**x**
Observation

**y**
Target

**ŷ**
Reconstruction

**z**

↓↓↓

Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

**x**
Observation

**y**
Target

**ŷ**
Reconstruction

**z**

z learns to encode the missing information necessary to predict y from x, i.e. the direction in which the purple box moves

Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

**x**
Observation

**ŷ**
Prediction

$$z \sim \mathcal{N}(\vec{0}, 1)$$

**z**

Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

**x**
Observation

$\hat{\mathbf{y}}$
Prediction

**z**

$z \sim \mathcal{N}(\vec{0}, 1)$

Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

$$\mathbf{x}$$
Observation

$$\hat{\mathbf{y}}$$
Prediction

$$\mathbf{z}$$

$$z \sim \mathcal{N}(\vec{0}, 1)$$
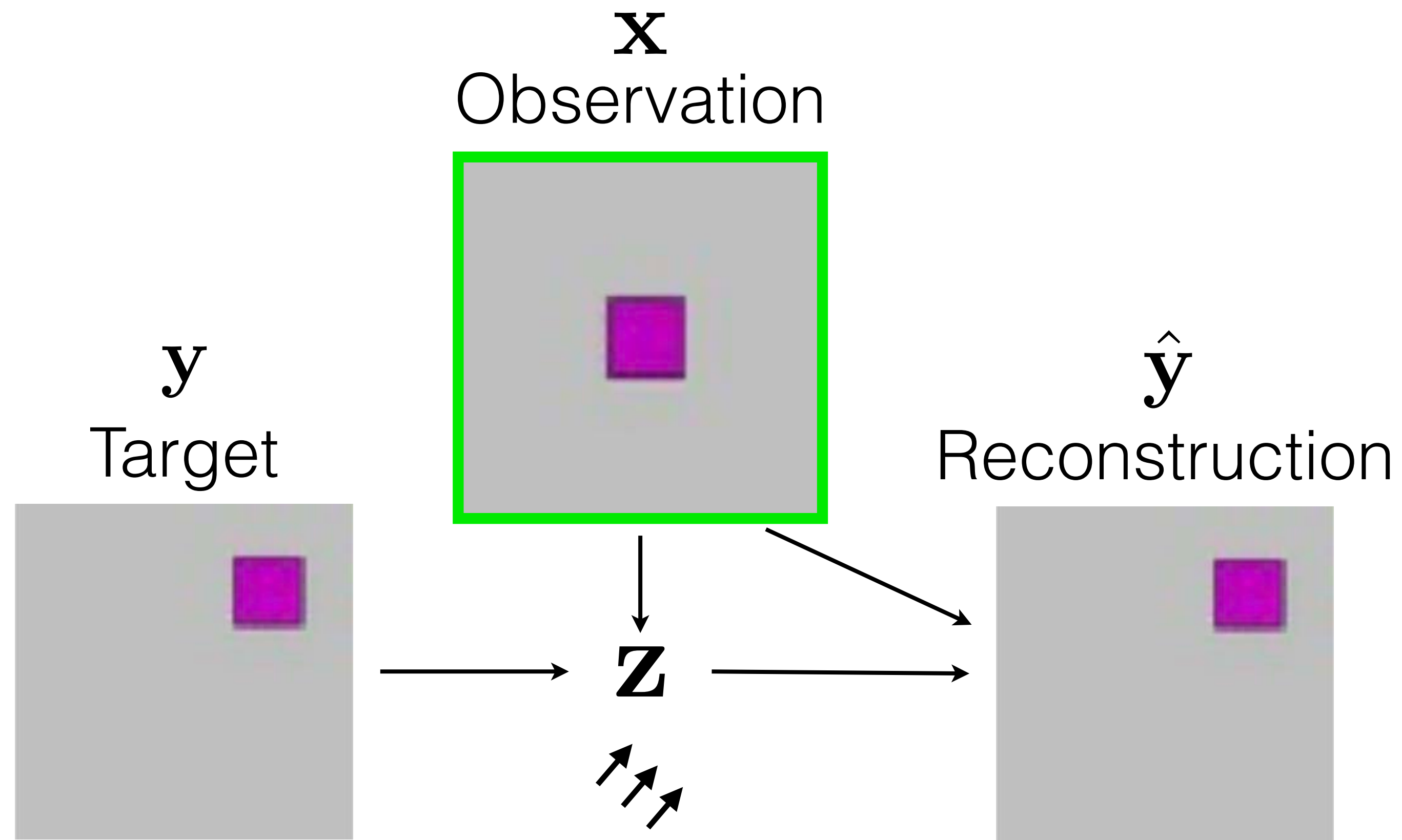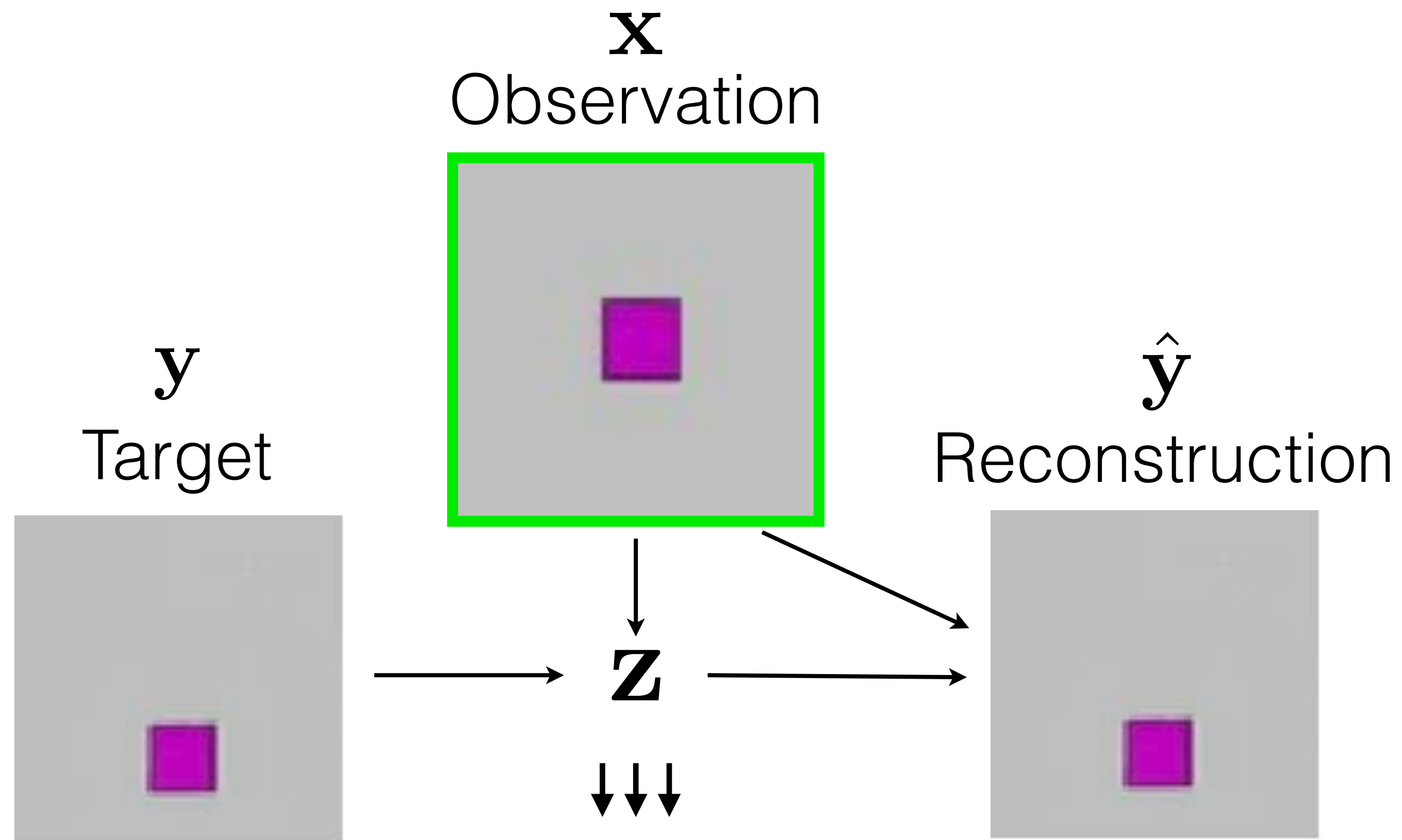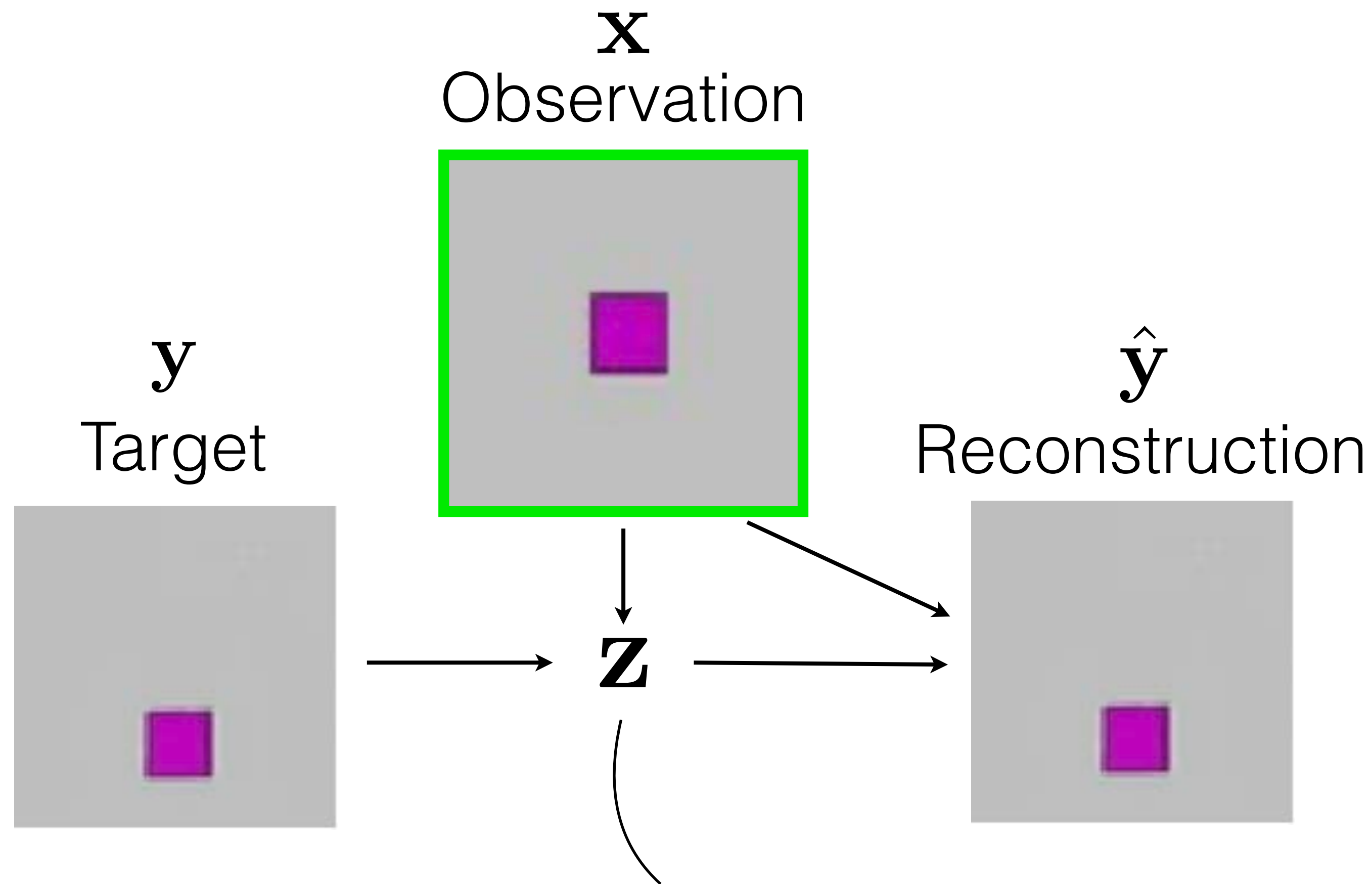
Example from [Babaeizadeh et al., ICLR 2018]

see also [Walker et al., ECCV 2016], [Xue*, Wu*, et al., NIPS 2016]

# Modeling multiple possible outputs

$$G(\mathbf{x})$$

$$\mathbf{x}$$

$$G$$



$$z \sim \mathcal{N}(\vec{0}, 1)$$

[BiCycleGAN, Zhu et al., NeurIPS 2017]

Concept: **You can control your data either via explicit inputs or latent variables**

# Text-to-text

└─→ with autoregressive models

# Autoregressive models are conditional generative models

1. **Training**     "Colorless green ideas sleep ____" $\longrightarrow$ Language model $\longrightarrow$ "furiously"

(Predict next characters)

# Autoregressive models are conditional generative models

2. **Predicting**     `[Review] + "The sentiment in this review is ____"` $\longrightarrow$ Language model $\longrightarrow$ `"positive"`

(Predict next characters)

# GPT-3 [Brown et al., 2020] https://arxiv.org/pdf/2005.14165.pdf

I had a horrible experience with this airline. The flight was delayed _five_ hours. Then, they lost my luggage. I had to wait a whole week to get the luggage back. I didn't have any changes of clothes so I wasted a lot of money to buy a few things to wear. When I finally got my luggage back, it was missing half my things. I will never fly with this airline again.

The sentiment in this review is negative.

# Image-to-text

↳ with an autoregressive model

Image Encoder

Text Decoder

yellow    bird    sitting

causal
self-attn

cross-attn

causal
self-attn

self-attn

A    yellow    bird

56

**Image Encoder**

cross-attn

self-attn ({  m

causal self-attn

causal **self-attn**

A yellow bird

> What is this slide showing and how does it relate to the way you work, ChatGPT?

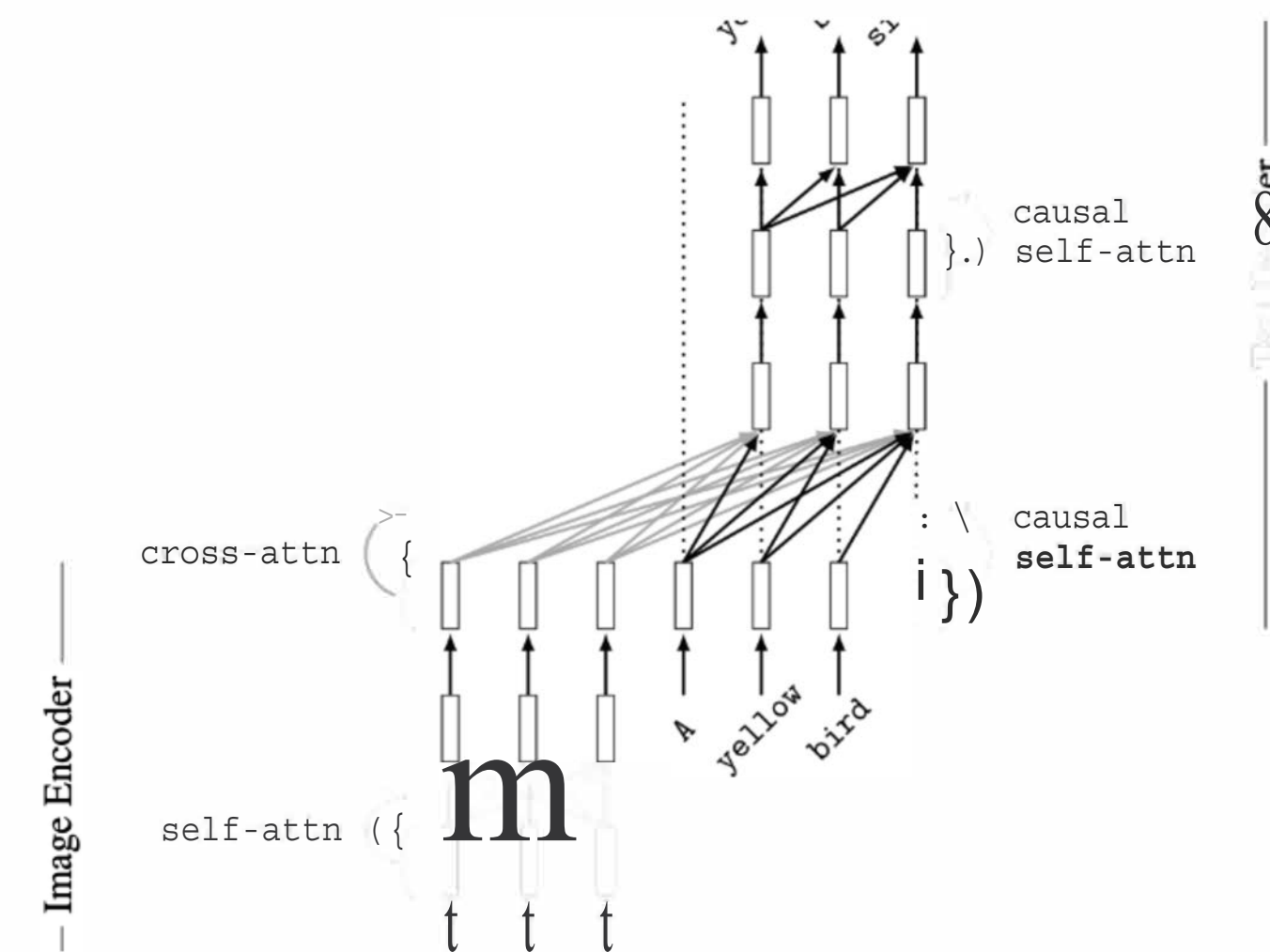This slide shows a model architecture for a vision-language transformer, often used in image-captioning models. The components include:

1. **Image Encoder:** On the left, where image inputs are processed through self-attention layers. This encoder extracts features from the image, representing various parts of the visual input.

2. **Cross-Attention Layer:** The encoder's output interacts with the text decoder through cross-attention, enalbling the text to attend to image features and guide the generation of descriptive captions.

# Does ChatGPT know how its own brain works?

## How This Relates to How I Work

Similar to this architecture, my responses often involve "attention" to different parts of your input. Instead of visual data, I focus on your conversational context and keywords, determining the next part of the response based on previous parts of the conversation. In this way, I function similarly to a text decoder with attention mechanisms to retain cons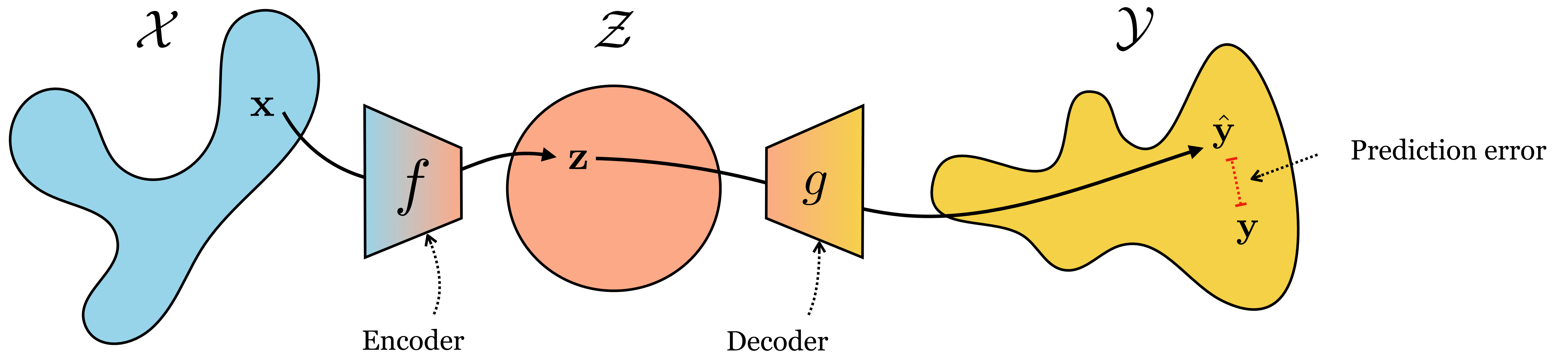istency and relevance in dialogue. However, unlike this model's image encoding, I don't process visual data directly, but rely on text representations and information you share about images.

Created with ChatGPT.
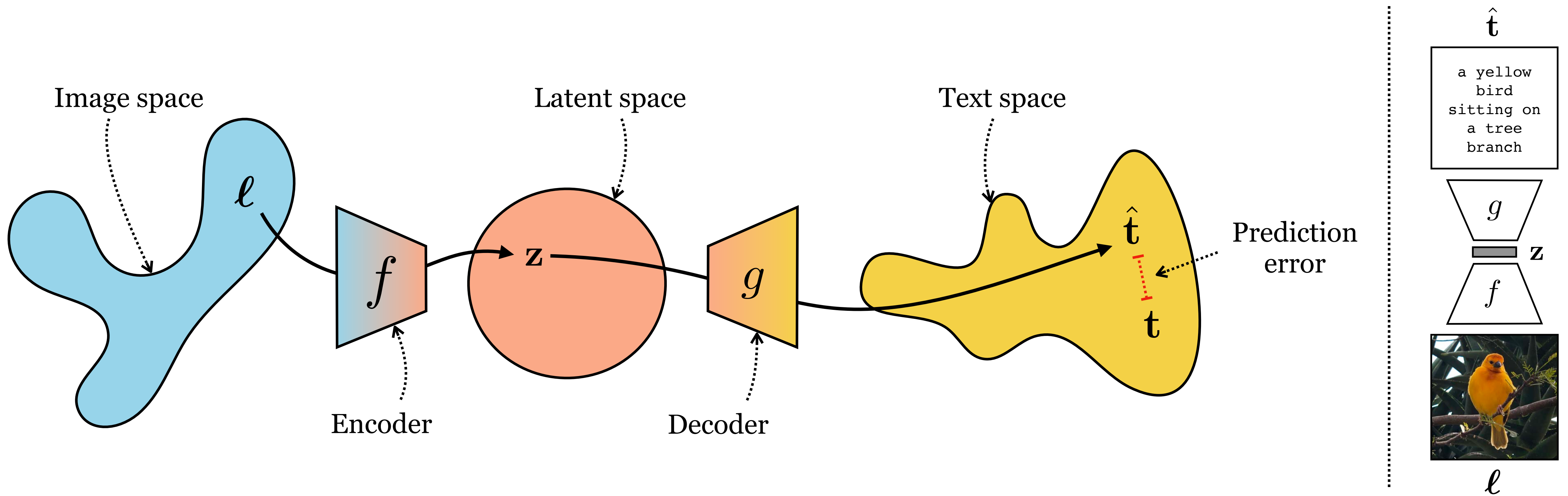
# Text-to-image

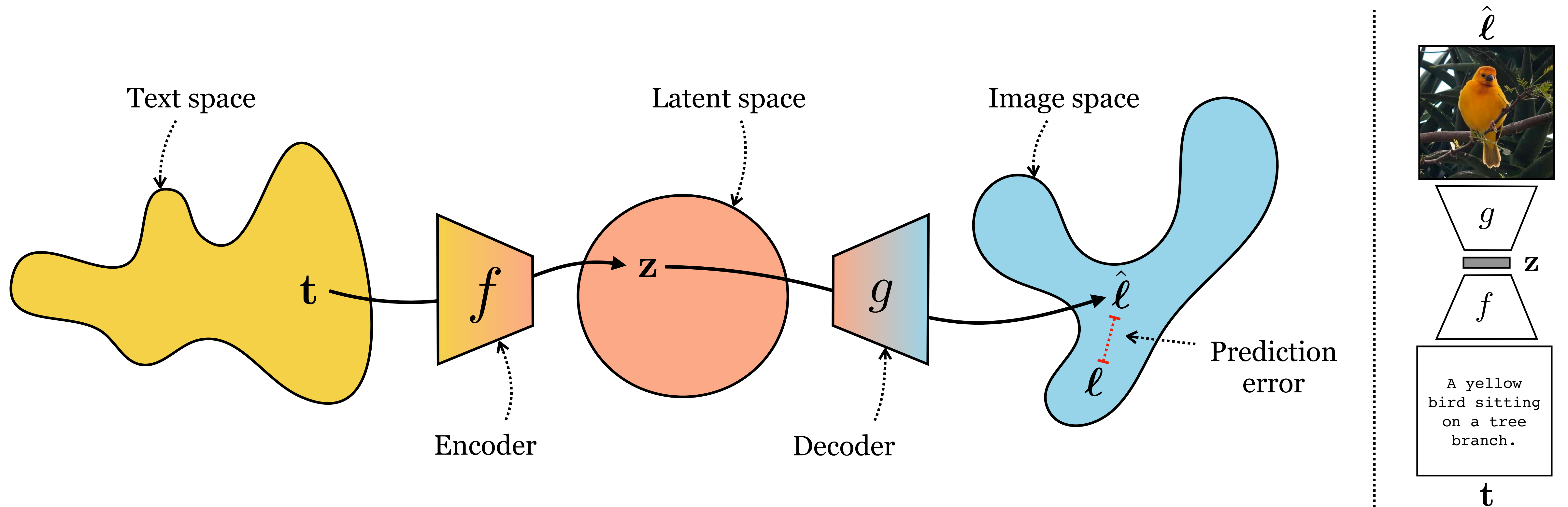└→ with all models combined!

# General strategy for domain translation

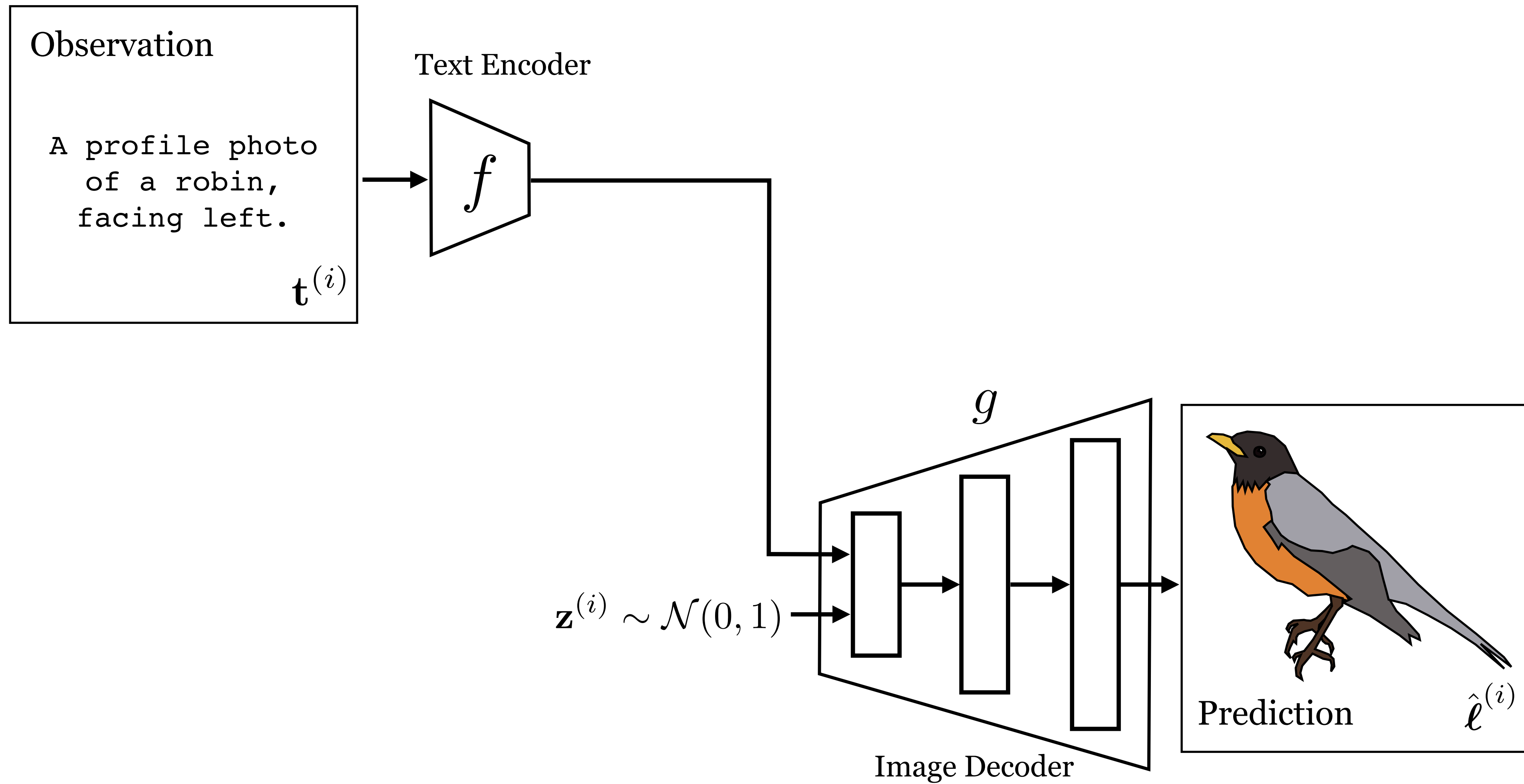# Text-to-image



Image space

Latent space

Text space

$\ell$

$f$

$\mathbf{z}$

$g$

$\hat{\mathbf{t}}$

$\mathbf{t}$

Prediction error

Encoder

Decoder

$\hat{\mathbf{t}}$

a yellow bird sitting on a tree branch

$g$

$\mathbf{z}$

$f$

$\ell$

# Text-to-image



Text space

Latent space

Image space

$\hat{\ell}$

Encoder

Decoder

Prediction error

$\mathbf{t}$

$f$

$\mathbf{z}$

$g$

$\hat{\ell}$

$\ell$

$g$

$\mathbf{z}$

$f$

A yellow bird sitting on a tree branch.

$\mathbf{t}$

# Text-to-image VAE



Observation

A scarlet macaw sitting on a perch.

$\mathbf{t}^{(i)}$

Text Encoder

$f$

Target $\boldsymbol{\ell}^{(i)}$

Image Encoder

$\mathbf{z}^{(i)}$

$g$

Image Decoder

Prediction $\hat{\boldsymbol{\ell}}^{(i)}$

For example: DALL-E 1 [Ramesh et al. 2021]

# Text-to-image VAE

Observation

```
A profile photo
   of a robin,
  facing left.
```
$\mathbf{t}^{(i)}$

Text Encoder

$f$

$\mathbf{z}^{(i)} \sim \mathcal{N}(0,1)$

$g$

Image Decoder
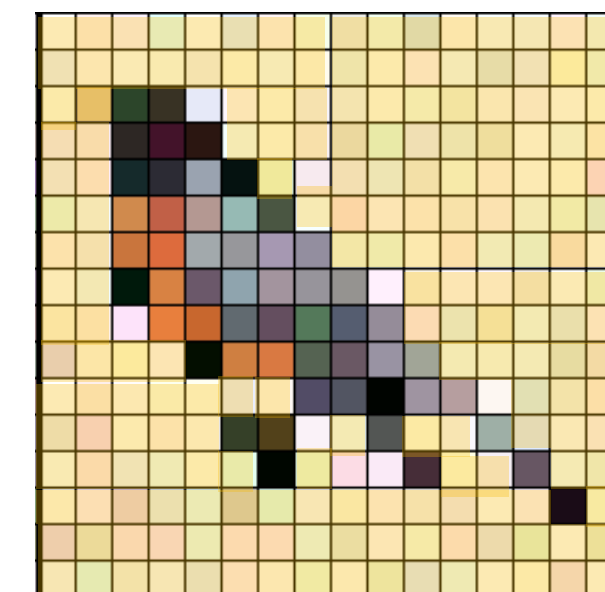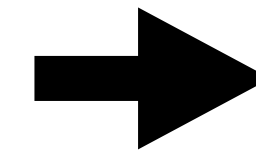
Prediction

$\hat{\ell}^{(i)}$

For example: DALL-E 1 [Ramesh et al. 2021]

# Text-to-image diffusion model



$$\text{Denoising} \longrightarrow$$

$$\mathbf{y}_T \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \mathbf{y}_0$$

$$\left[\begin{array}{c}\texttt{A profile photo}\\\texttt{of a robin,}\\\texttt{facing left.}\end{array} \quad , \quad \mathbf{y}_t\right] \xrightarrow{\quad f_\theta \quad} \mathbf{y}_{t-1}$$
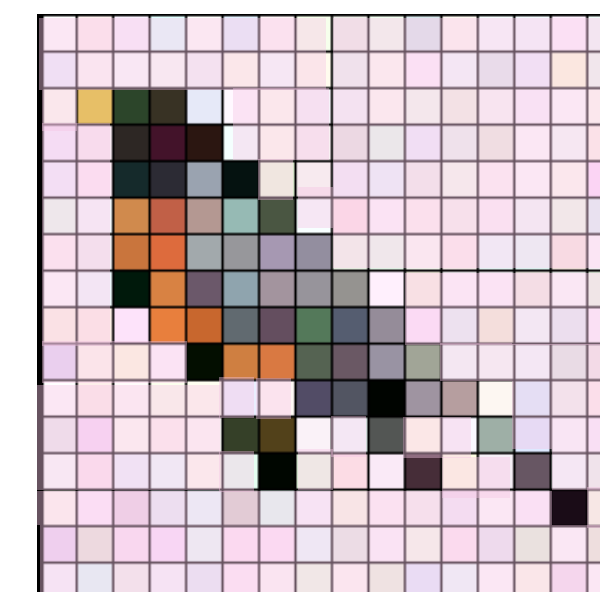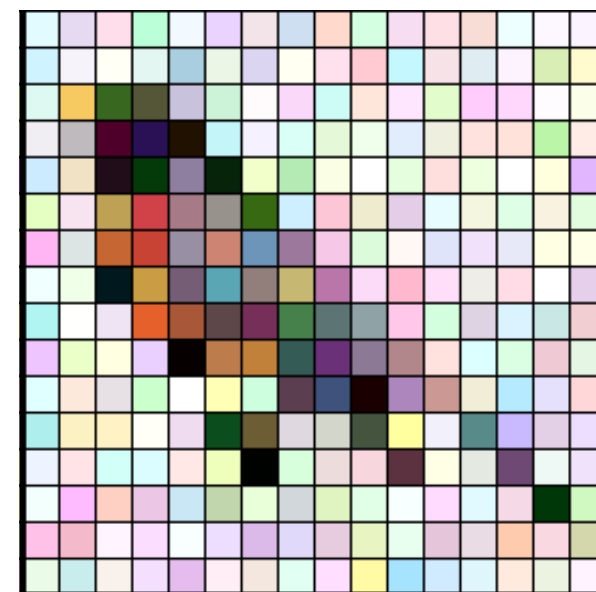
$$\mathbf{x}$$

For example: DALL-E 2 [Ramesh et al. 2022], Stable Diffusion [Rombach*, Blattman* et al. 2022]

# Text provides information as to what the noiseless image should look like
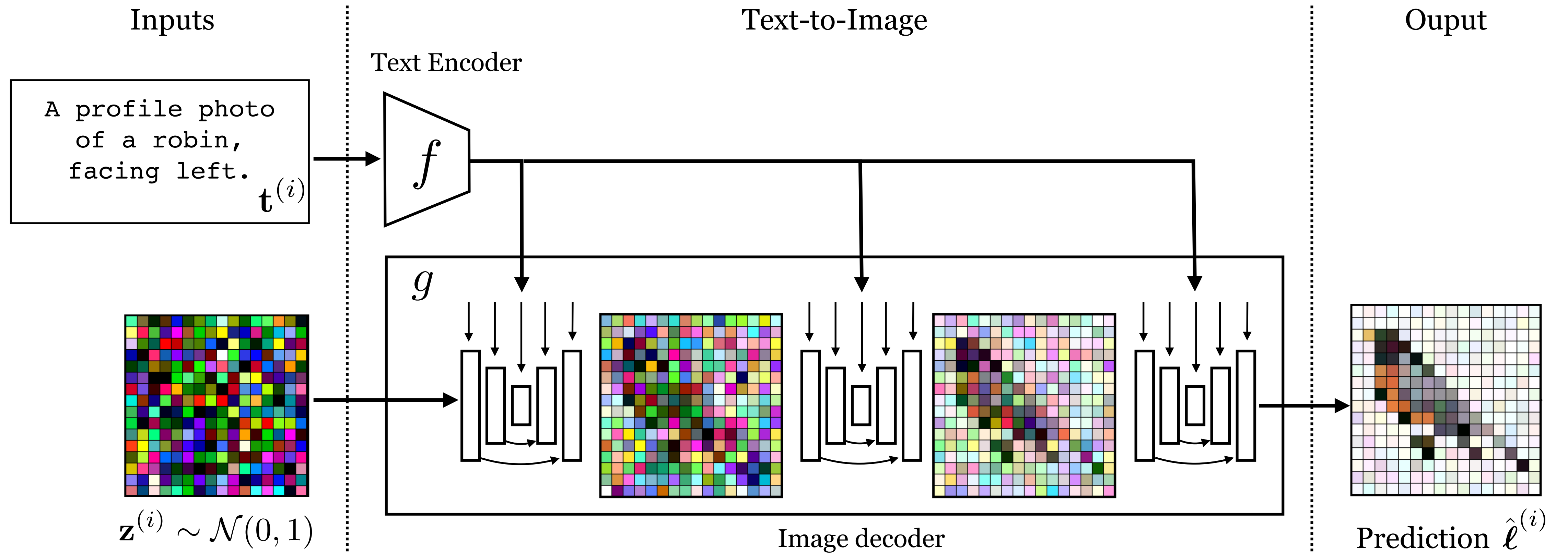


A robin on a **yellow** background.
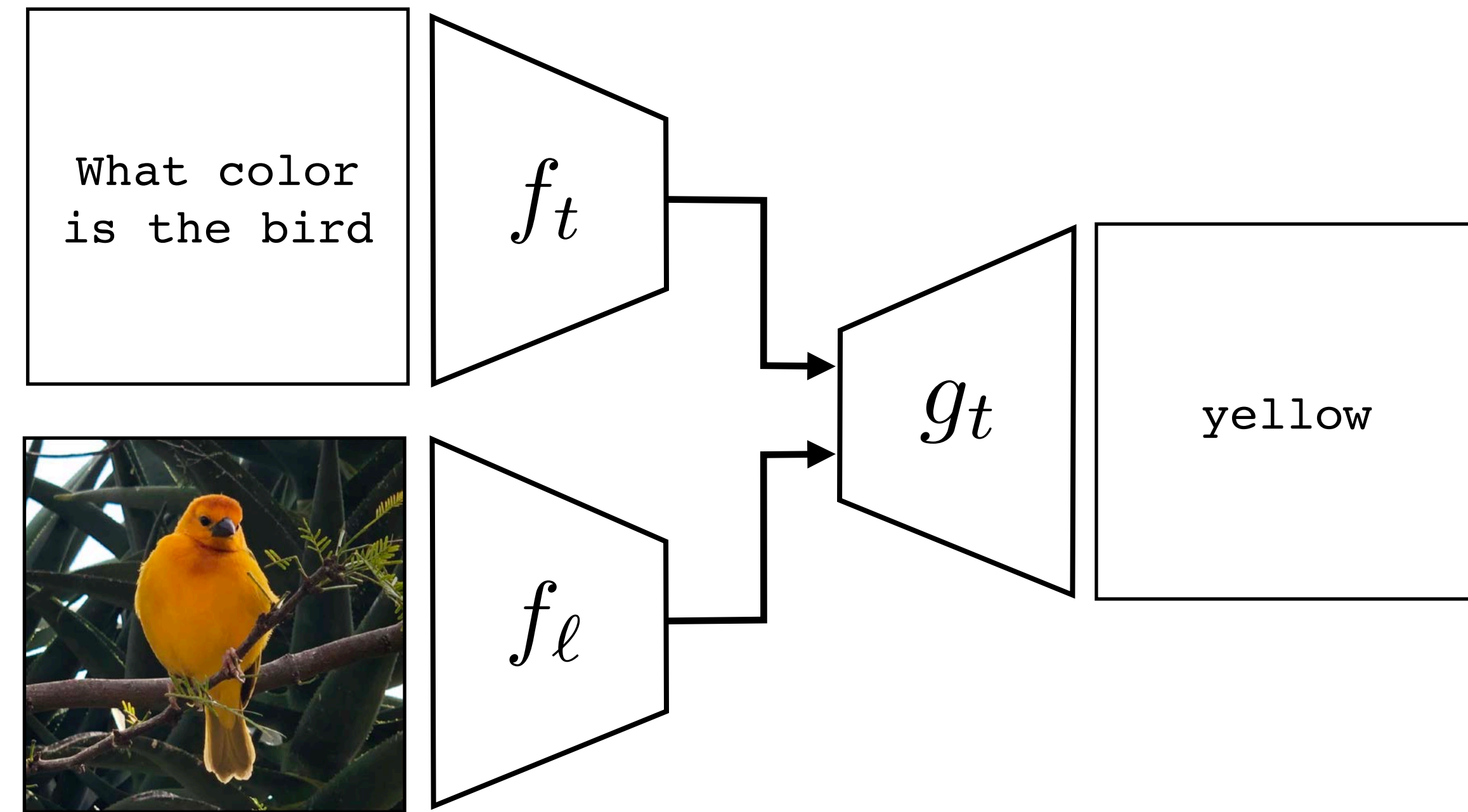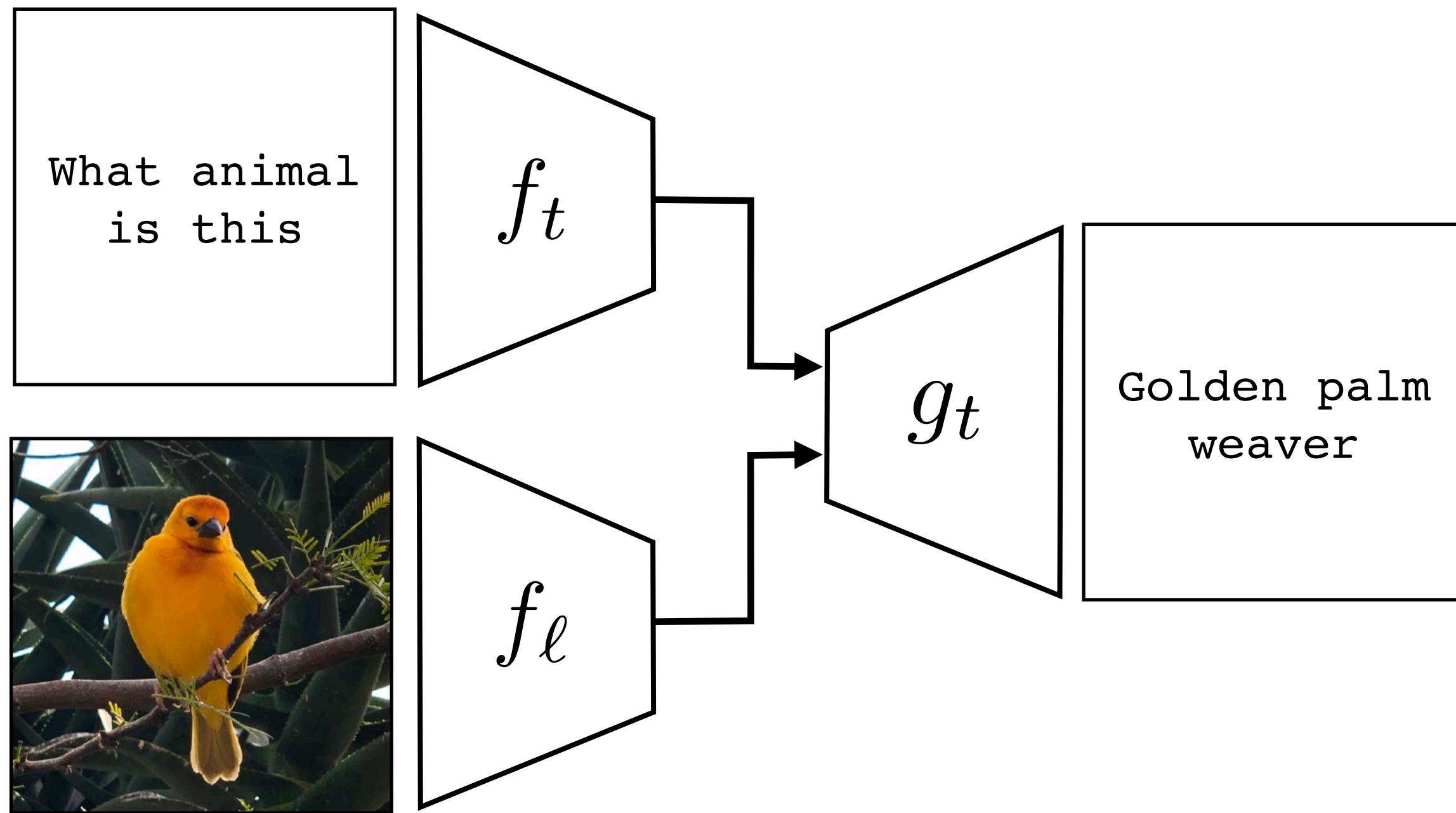
A robin on a **pink** background.

# Text-to-image diffusion model



Inputs

Text-to-Image

Ouput

Text Encoder

A profile photo of a robin, facing left. $\mathbf{t}^{(i)}$

$f$

$g$

Image decoder

$\mathbf{z}^{(i)} \sim \mathcal{N}(0, 1)$

Prediction $\hat{\ell}^{(i)}$

For example: DALL-E 2 [Ramesh et al. 2022], Stable Diffusion [Rombach*, Blattman* et al. 2022]

# Visual Question Answering (VQA)
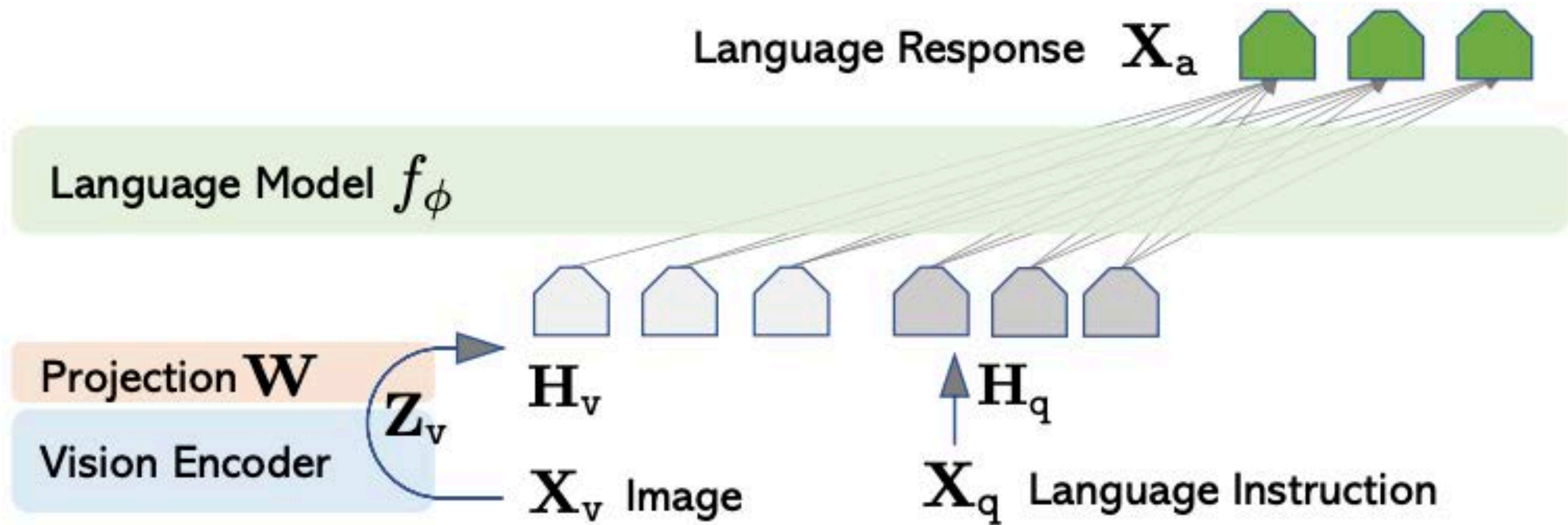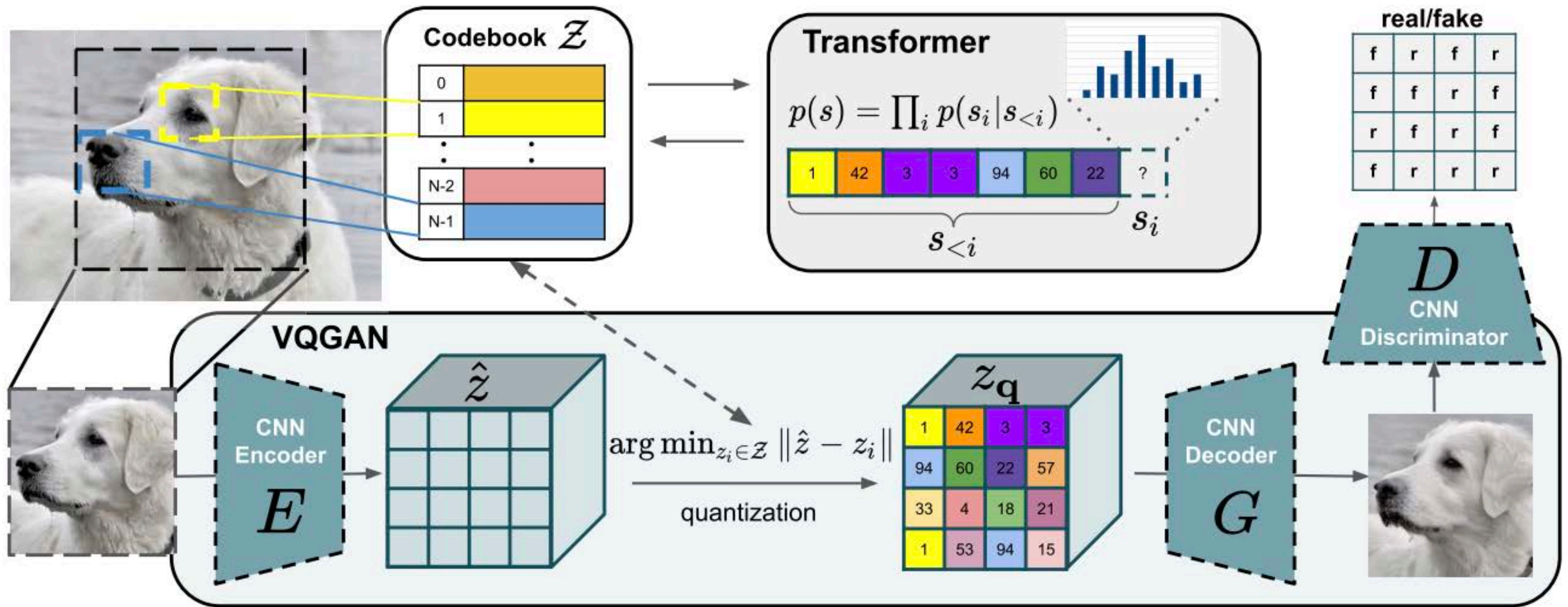


[Antol*, Agrawal* et al., ICCV 2015]

# LLaVA



Figure 1: LLaVA network architecture.

["Visual Instruction Tuning", Liu*, Li*, et al. 2023]

# Mixing ideas from GANs, VAEs, and Autoregressive models

https://arxiv.org/pdf/2012.09841.pdf    ["**Taming Transformers**", Esser, Rombach, Ommer 2020]

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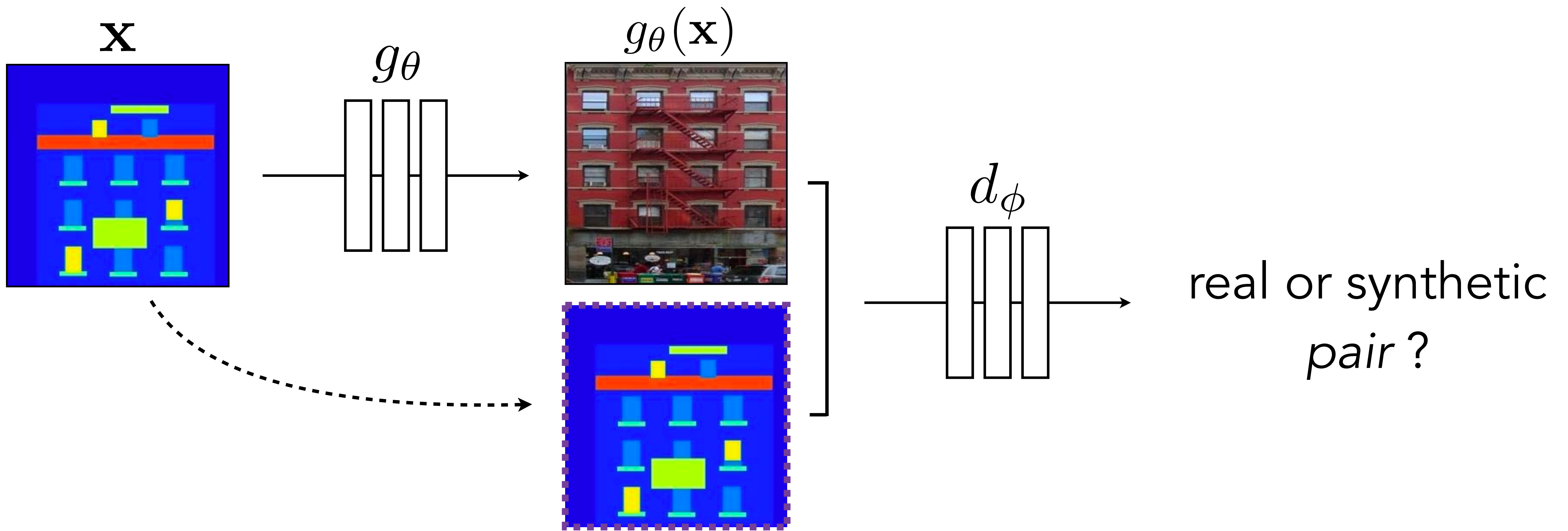
[Includes slides from Jun-Yan Zhu, Taesung Park]

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

# Paired

$$\mathbf{x}^{(i)} \qquad \mathbf{y}^{(i)}$$



$$\left\{ \quad , \quad \right\}$$

$$\left\{ \quad , \quad \right\}$$

$$\left\{ \quad , \quad \right\}$$

$$\vdots$$

# Unpaired

$$X \qquad\qquad Y$$



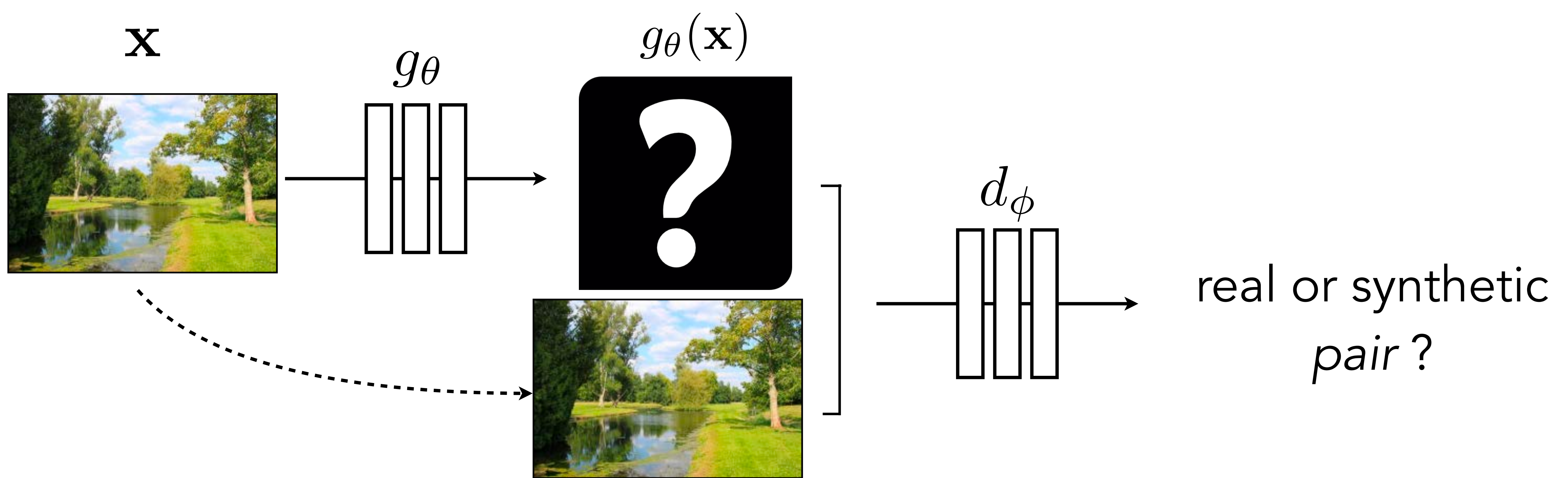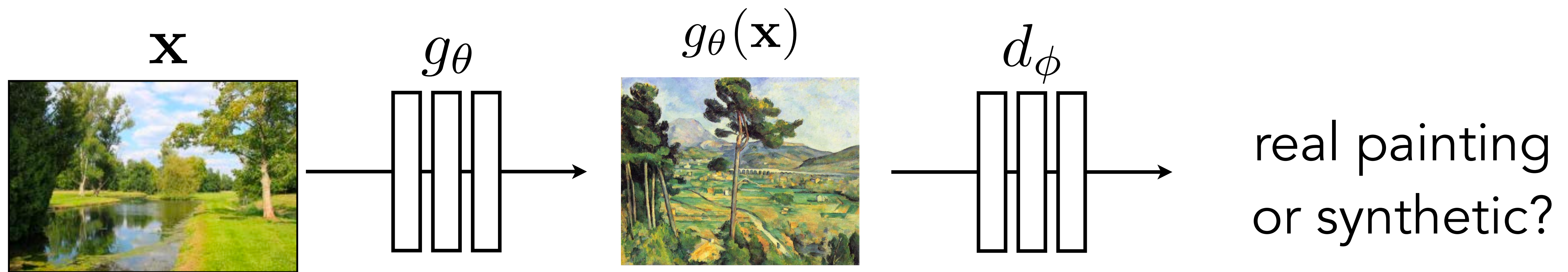$$\left( \quad \right) , \left( \quad \right)$$

$$\vdots \qquad\qquad \vdots$$

$$\arg\min_{\theta}\max_{\phi}\mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_{\phi}(\mathbf{x}, g_{\theta}(\mathbf{x})) + \log(1 - d_{\phi}(\mathbf{x}, \mathbf{y}))]$$

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

# No input-output training pairs!

$$\arg \min_\theta \max_\phi \mathbb{E}_{\mathbf{x},\mathbf{y}}[\log d_\phi(g_\theta(\mathbf{x})) + \log(1 - d_\phi(\mathbf{y}))]$$

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

GAN loss checks if output is part of an admissible *set*
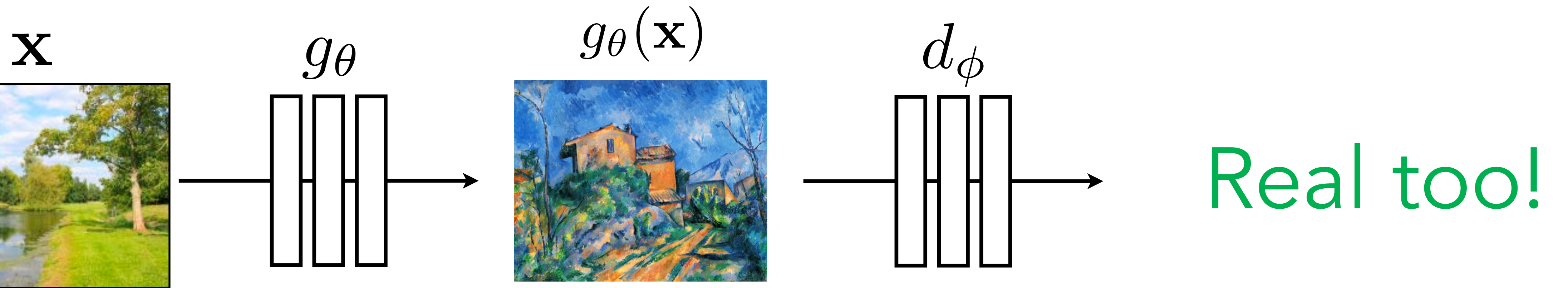
$\mathbf{x}$

$g_\theta$

$g_\theta(\mathbf{x})$

$d_\phi$

Real!
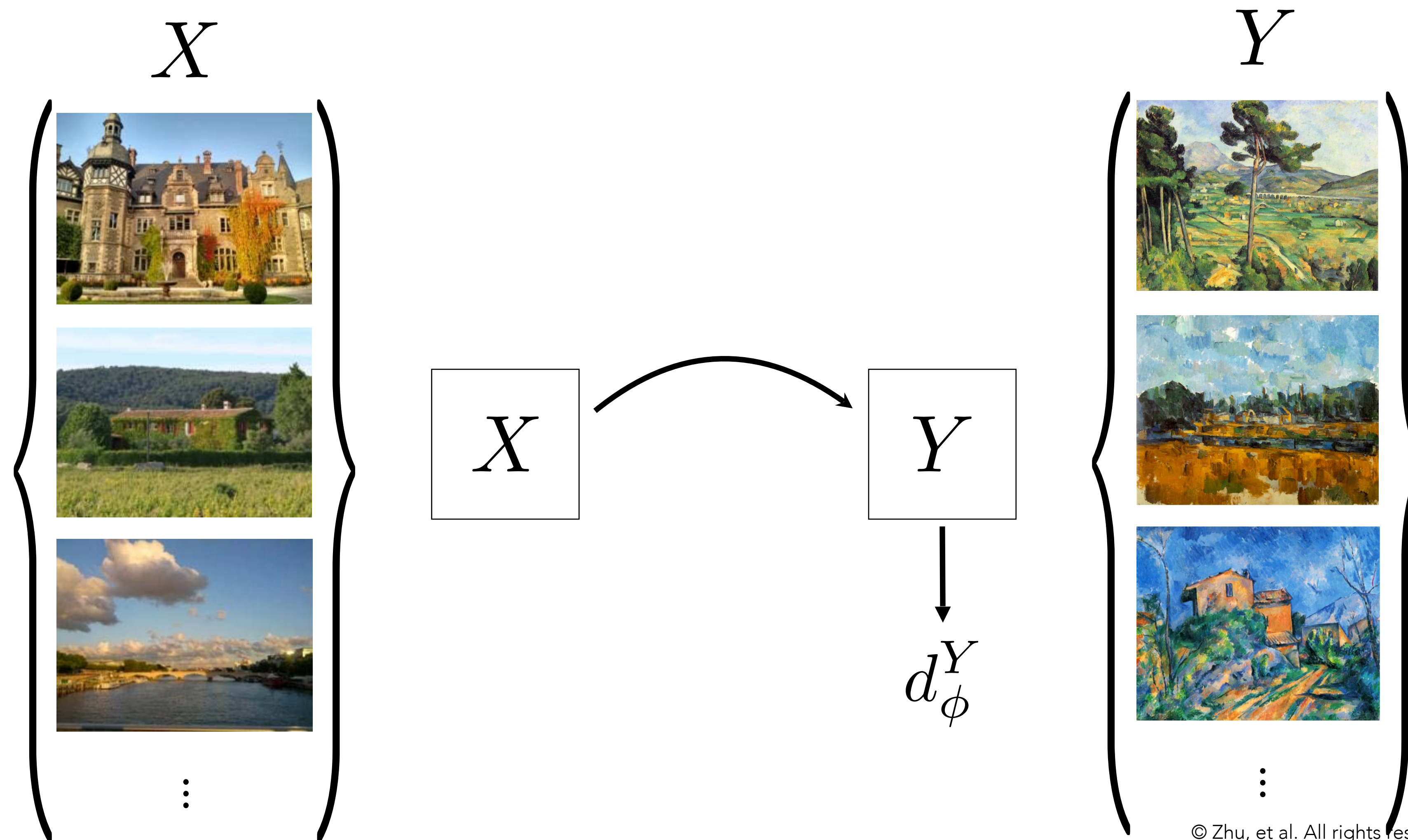
$$\mathbf{x} \qquad g_\theta \qquad g_\theta(\mathbf{x}) \qquad d_\phi$$



Real too!

Nothing to force output to correspond to input
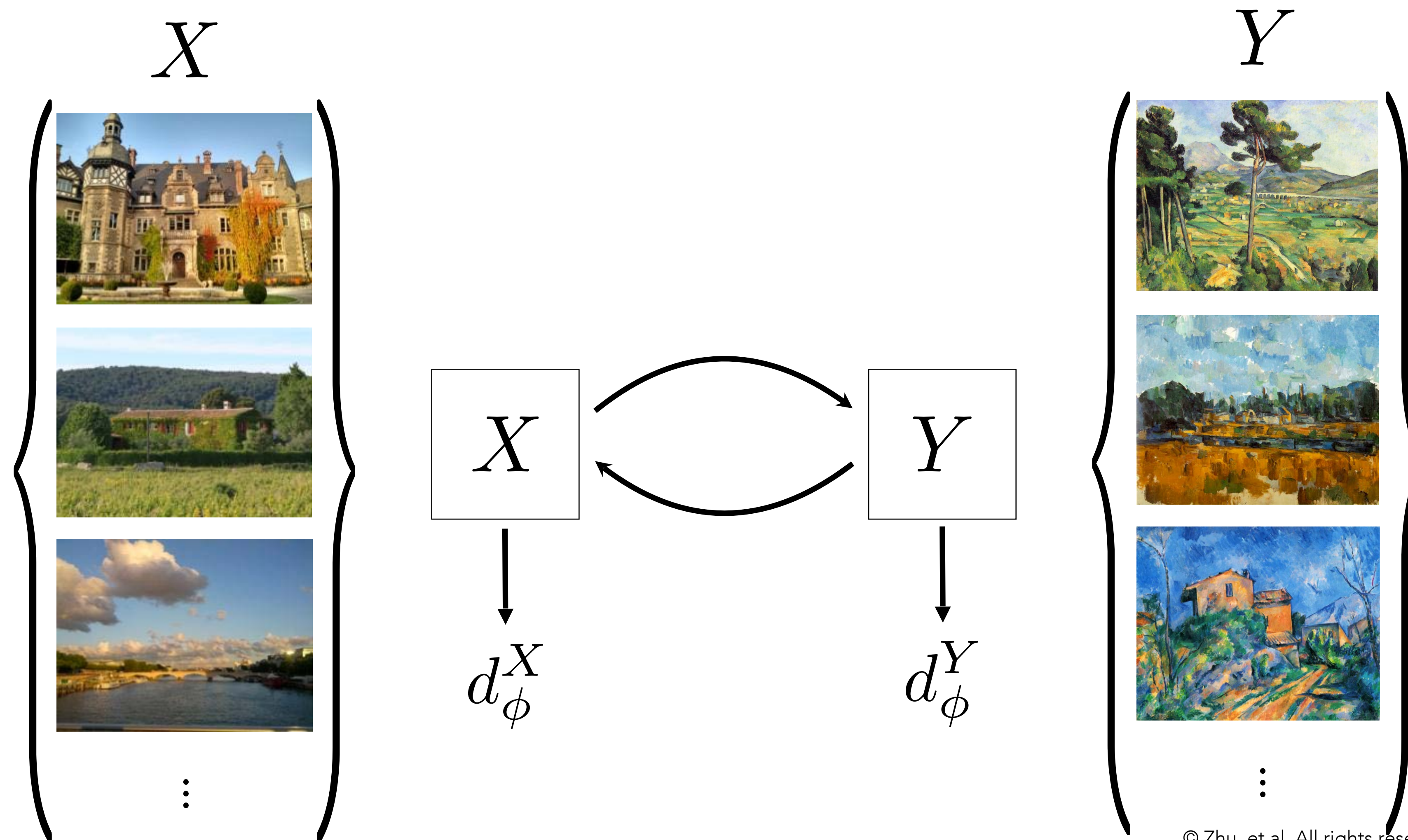
# CycleGAN, or there and back aGAN

$X$

$Y$



$$X \longrightarrow Y$$

$$\downarrow d_\phi^Y$$

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

# CycleGAN, or there and back aGAN

$$X$$

$$Y$$



$$X \rightleftharpoons Y$$

$$d_\phi^X$$

$$d_\phi^Y$$

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

# Cycle Consistency Loss

# Unique solution?



$$W = \begin{pmatrix} 1 & -2b_1 \\ 0 & 4 - 2b_2 \end{pmatrix}$$

$$W = \begin{pmatrix} -1 & 1 - 2b_1 \\ 0 & 4 - 2b_2 \end{pmatrix}$$

[Galanti, Wolf, Benaim, ICLR 2018]

# Unpaired translation in language



A)
- observed source sentence (red filled circle)
- unobserved translation of a target sentence (red open circle)
- system translation of a target sentence (red X)

B)
- observed target sentence (blue filled circle)
- unobserved translation of a source sentence (blue open circle)
- system translation of a source sentence (blue X)

["Phrase-Based & Neural Unsupervised Machine Translation", Lample et al. 2018]

Concept: **paired data is not always necessary for grounding** (finding a unique mapping between two domains)