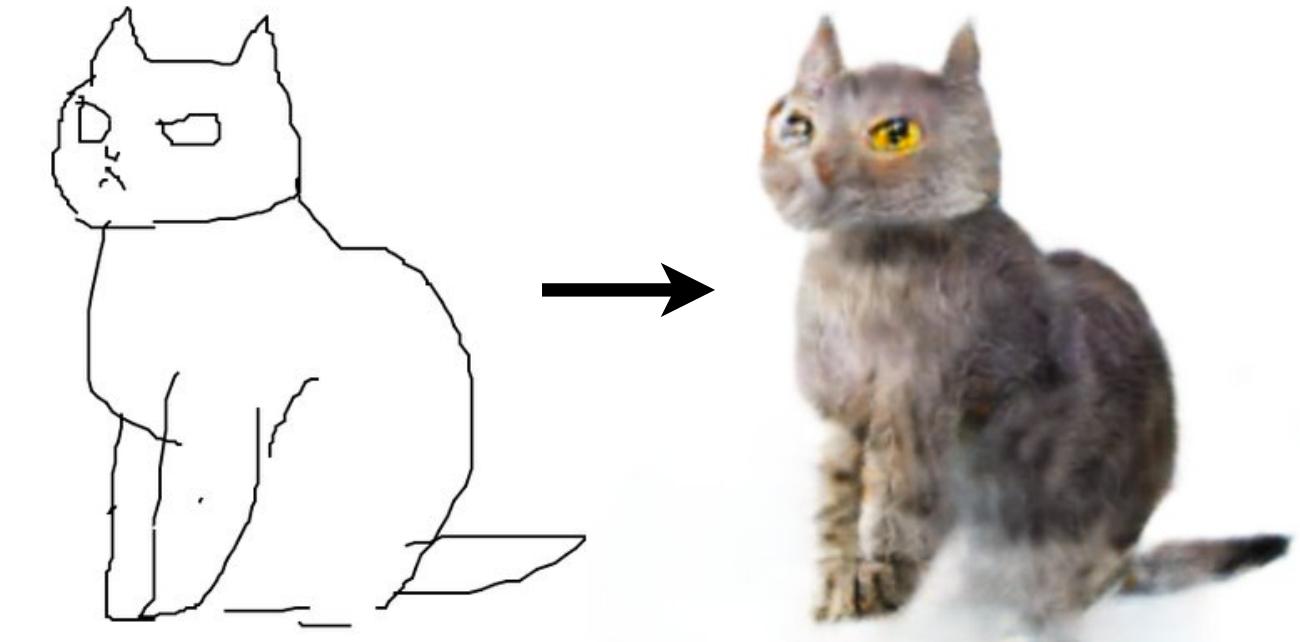


$z \sim \mathcal{N}(\vec{0}, 1)$ →



Generative Adversarial Networks

Phillip Isola
9.520
10/17/18

Image classification



“Fish”

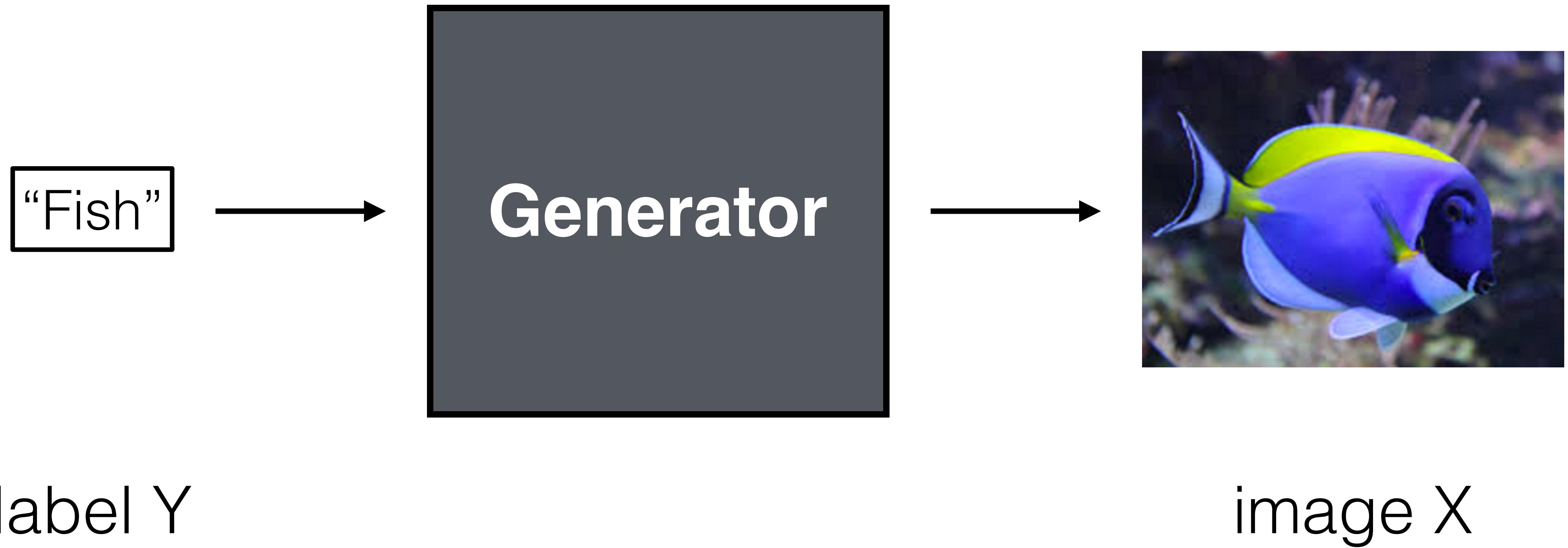
A rectangular box with a black border containing the word "Fish" in a black, sans-serif font.

image X

label Y

⋮

Image generation

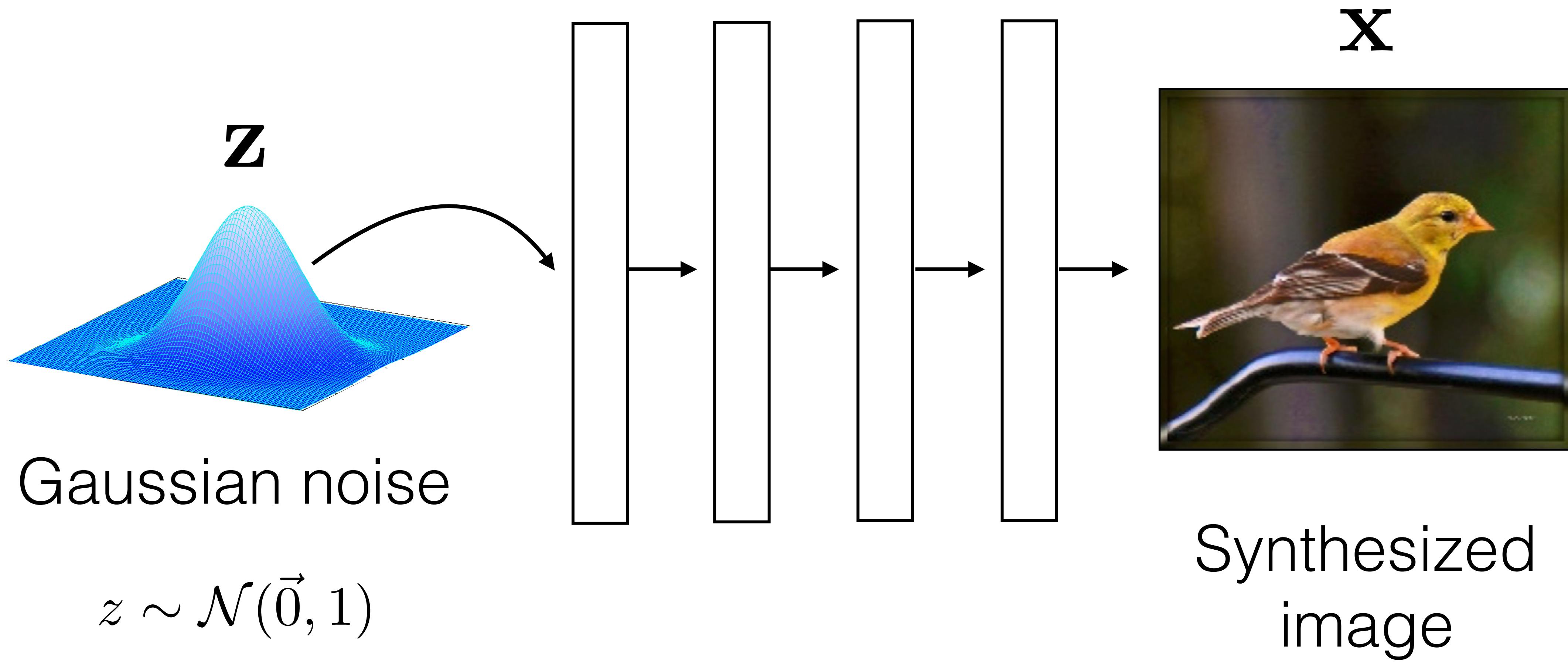


What's a generative model?

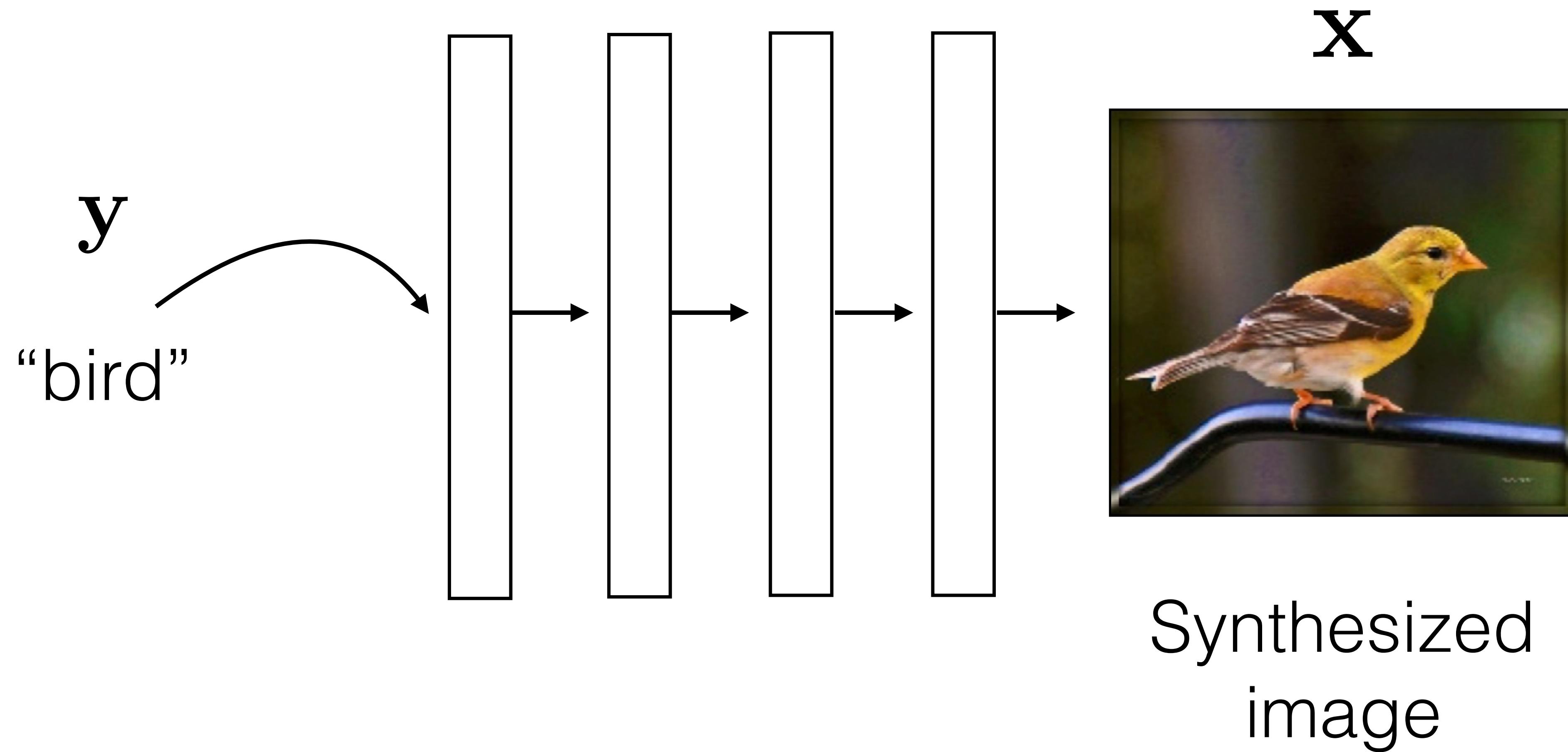
Model of high-dimensional unobserved variables $P(\mathbf{X}|\mathbf{Y} = \mathbf{y})$

Useful for lots of problems beyond sampling random images!

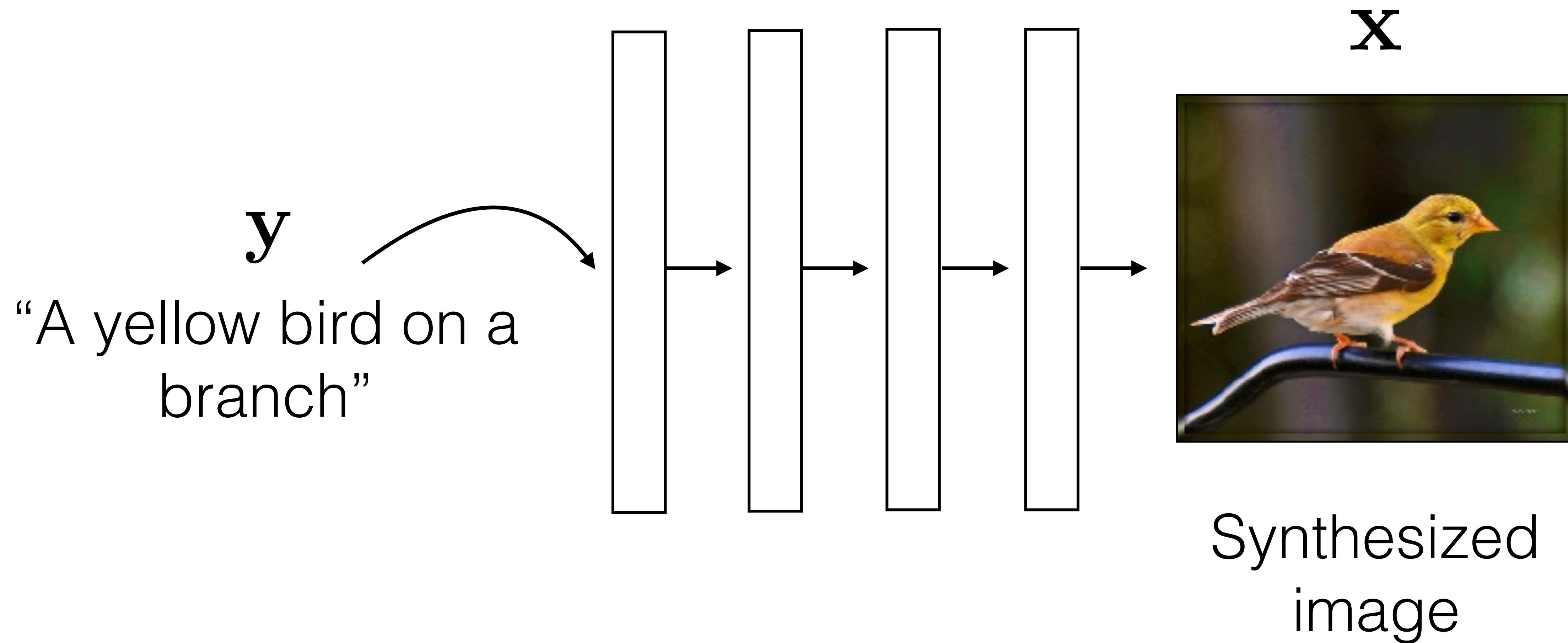
Generative Model



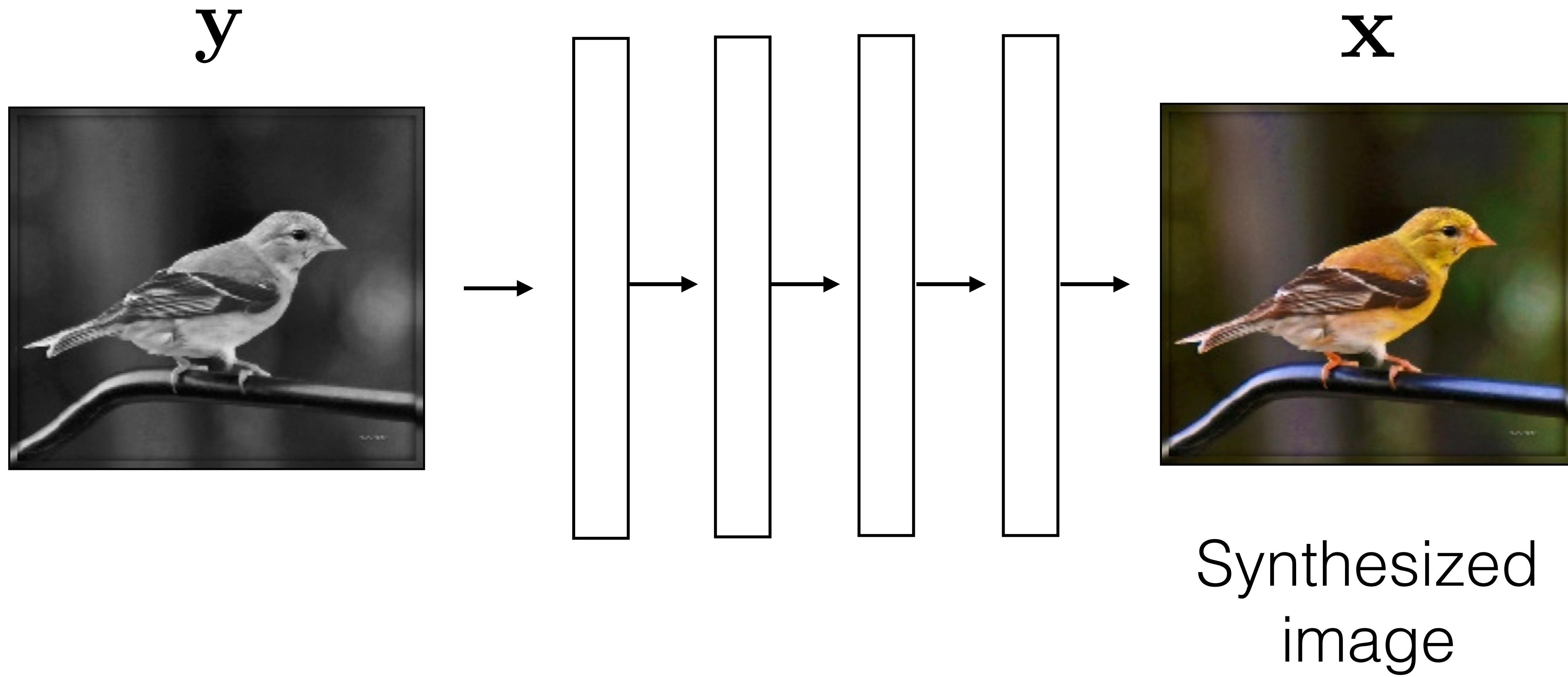
Conditional Generative Model



Conditional Generative Model



Conditional Generative Model



Three perspectives on GANs

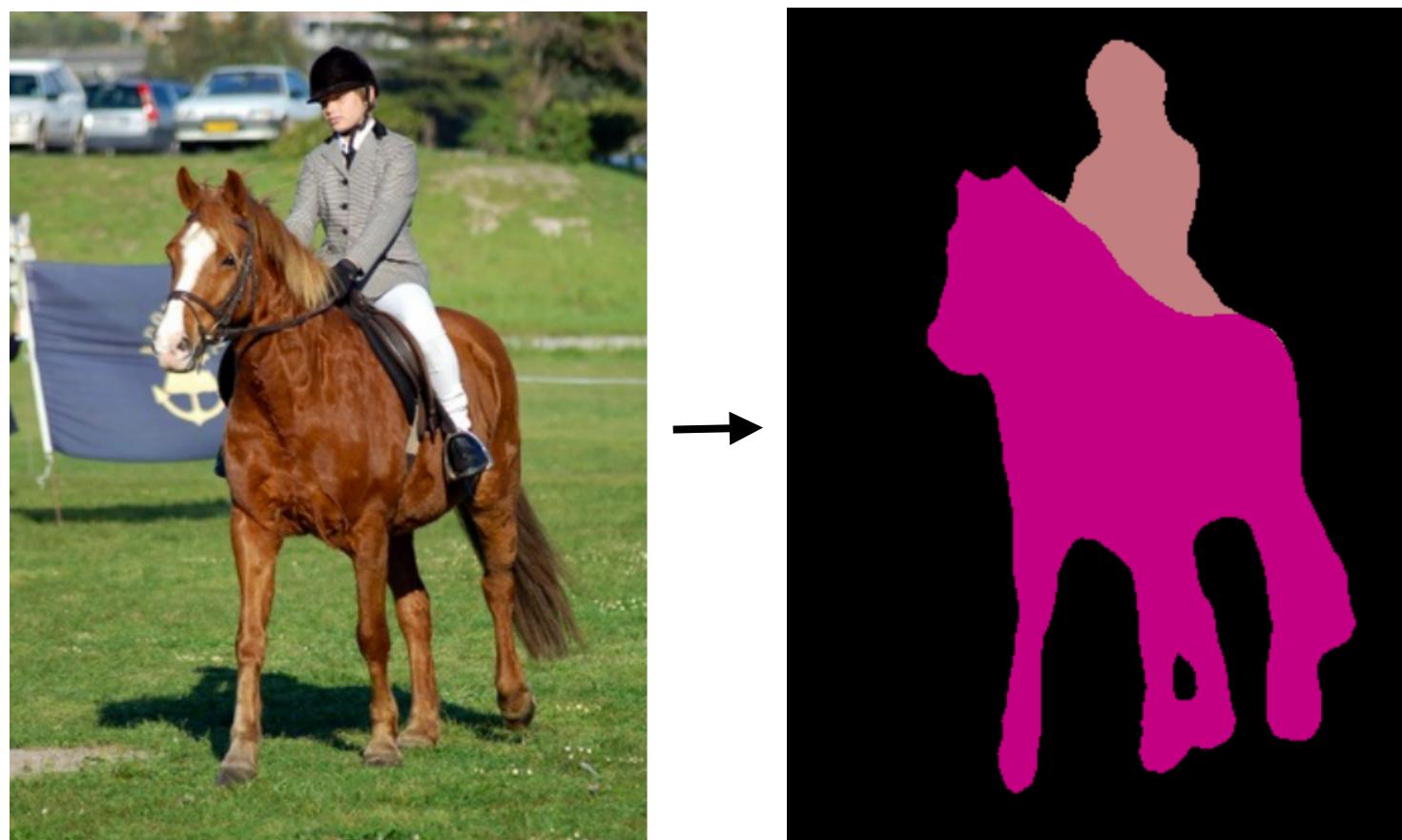
1. Structured loss
2. Generative model
3. Domain-level supervision / mapping

Three perspectives on GANs

1. Structured loss
2. Generative model
3. Domain-level supervision / mapping

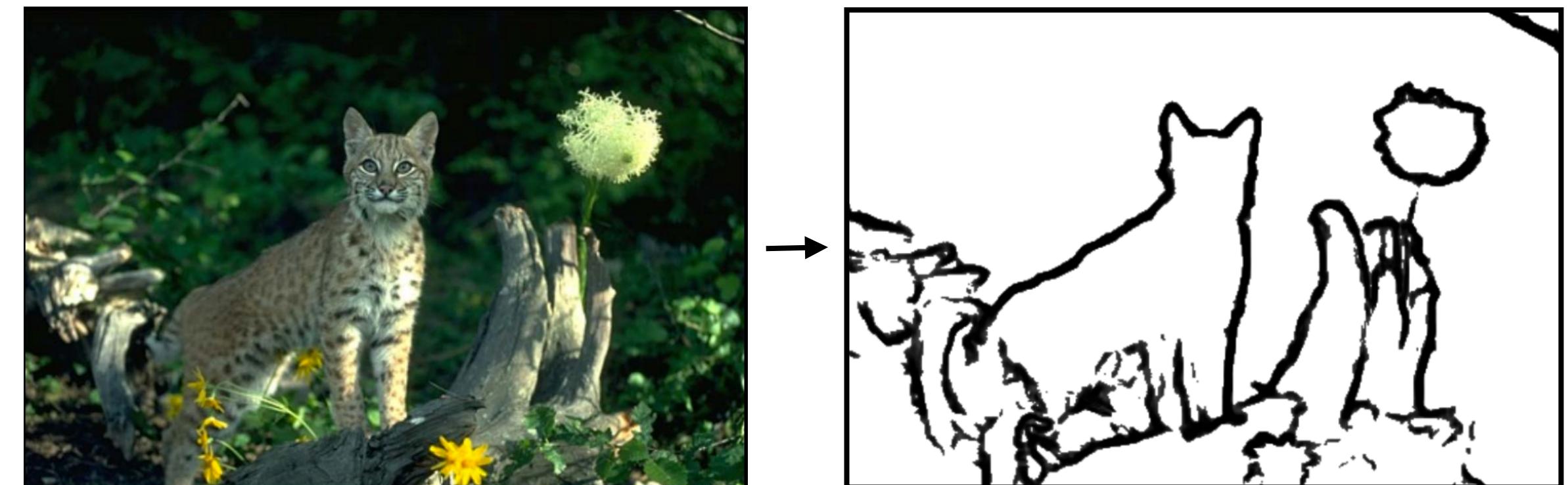
Data prediction problems (“structured prediction”)

Object labeling



[Long et al. 2015, ...]

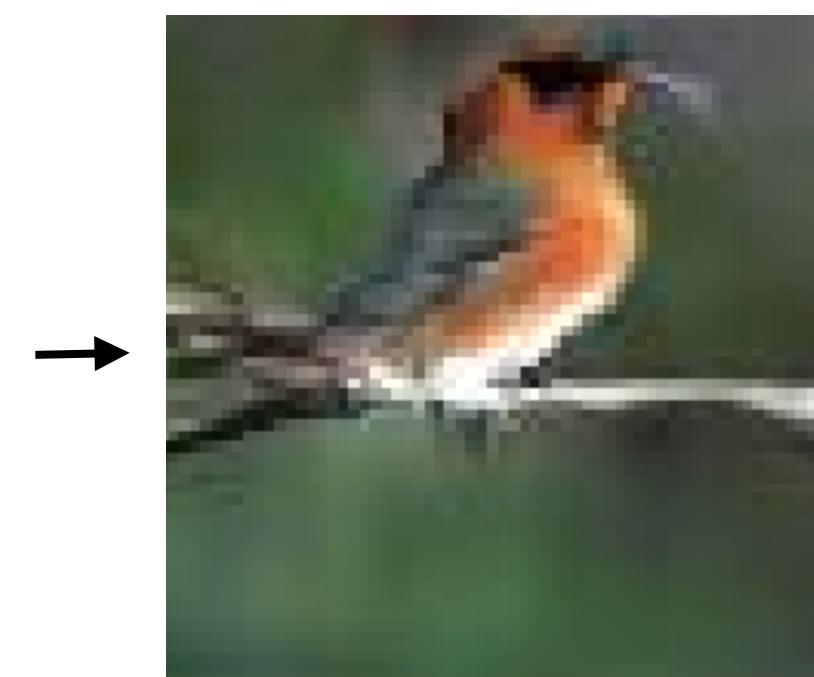
Edge Detection



[Xie et al. 2015, ...]

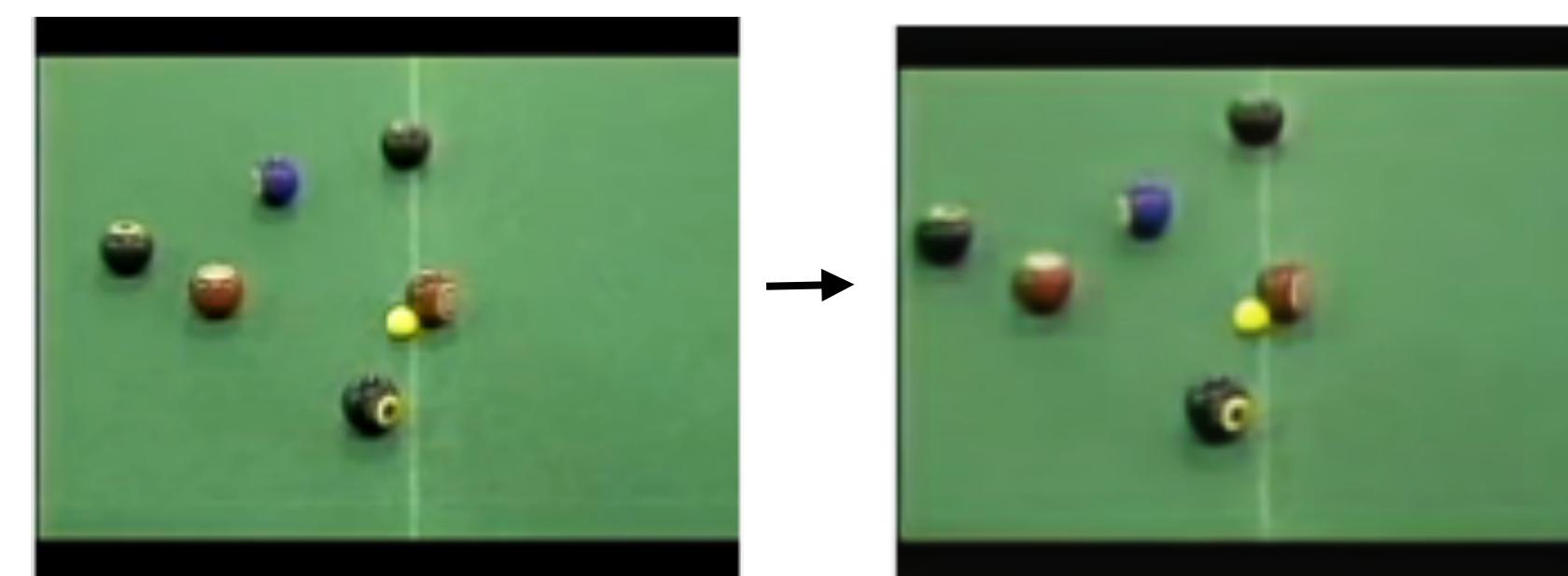
Text-to-photo

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



[Reed et al. 2014, ...]

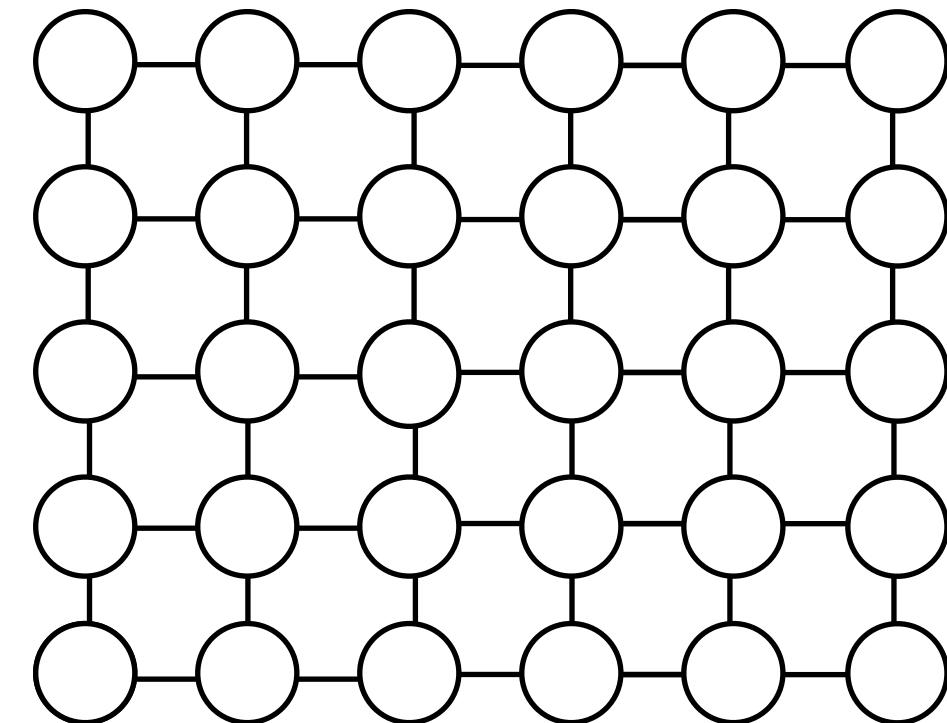
Future frame prediction



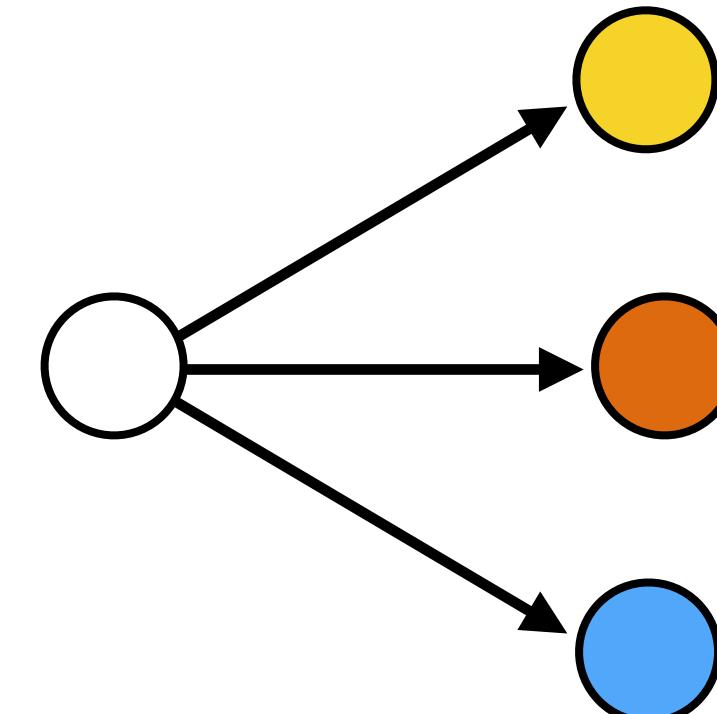
[Mathieu et al. 2016, ...]

Challenges in data prediction

1. Output is a high-dimensional, structured object

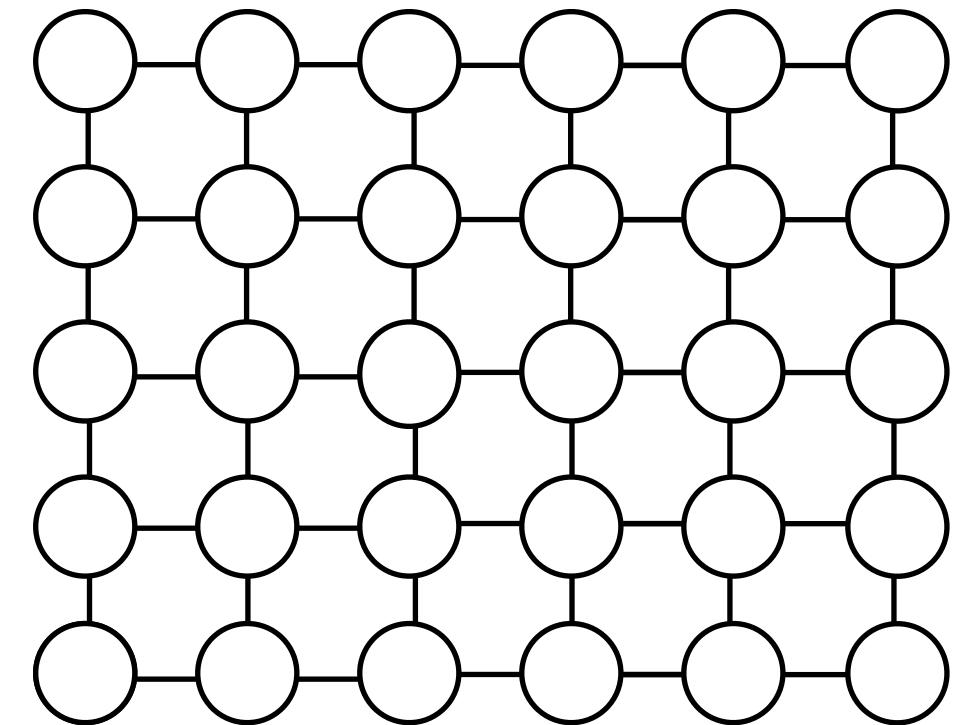


2. Uncertainty in the mapping, many plausible outputs



Properties of generative models

1. Model high-dimensional, structured output



2. Model uncertainty; a whole distribution of possible outputs

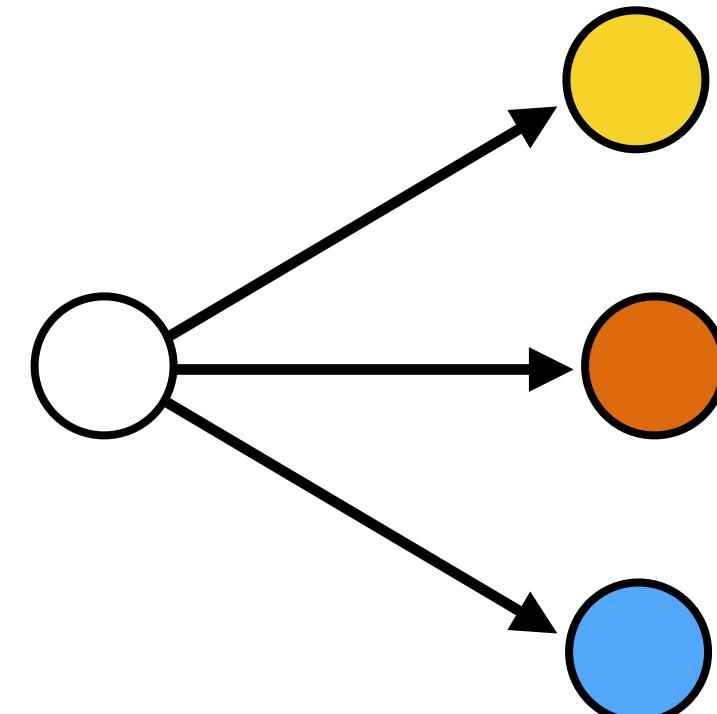
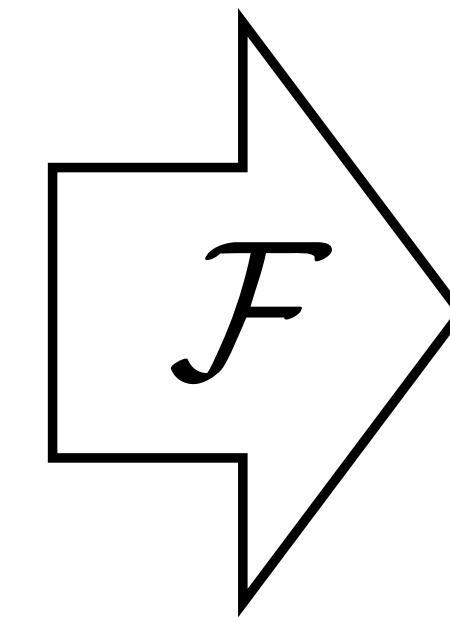


Image-to-Image Translation

Input \mathbf{x}

<i>Training data</i>	
\mathbf{x}	\mathbf{y}
	
	
	
:	



Output \mathbf{y}



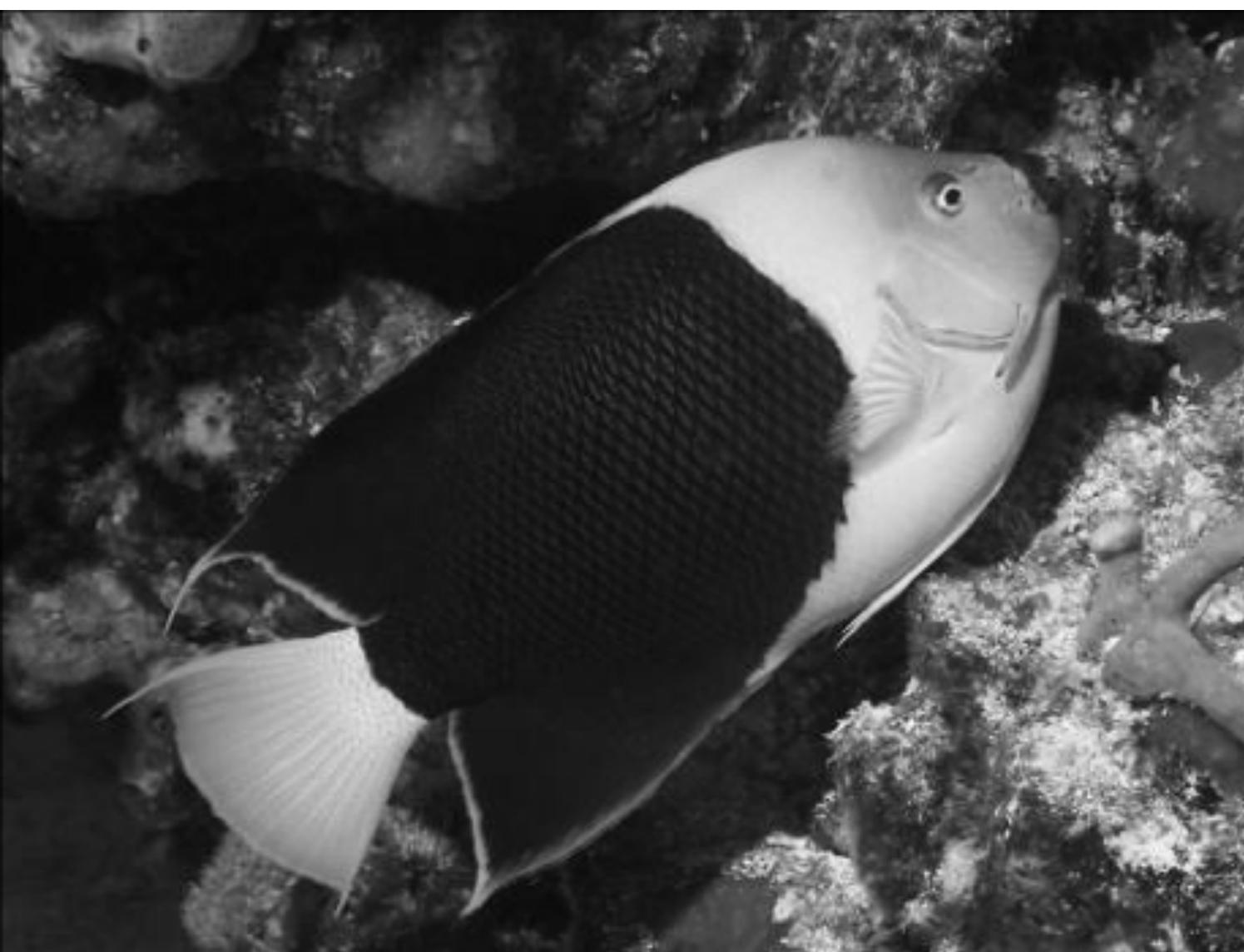
$$\arg \min_{\mathcal{F}} \mathbb{E}_{\mathbf{x}, \mathbf{y}} [L(\mathcal{F}(\mathbf{x}), \mathbf{y})]$$

Objective function
(loss)

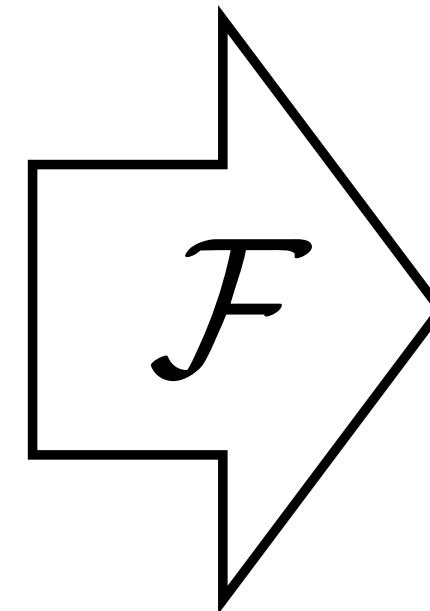
Neural Network

Image-to-Image Translation

Input \mathbf{x}



Output \mathbf{y}



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

“What should I do”

“How should I do it?”

Designing loss functions

Input



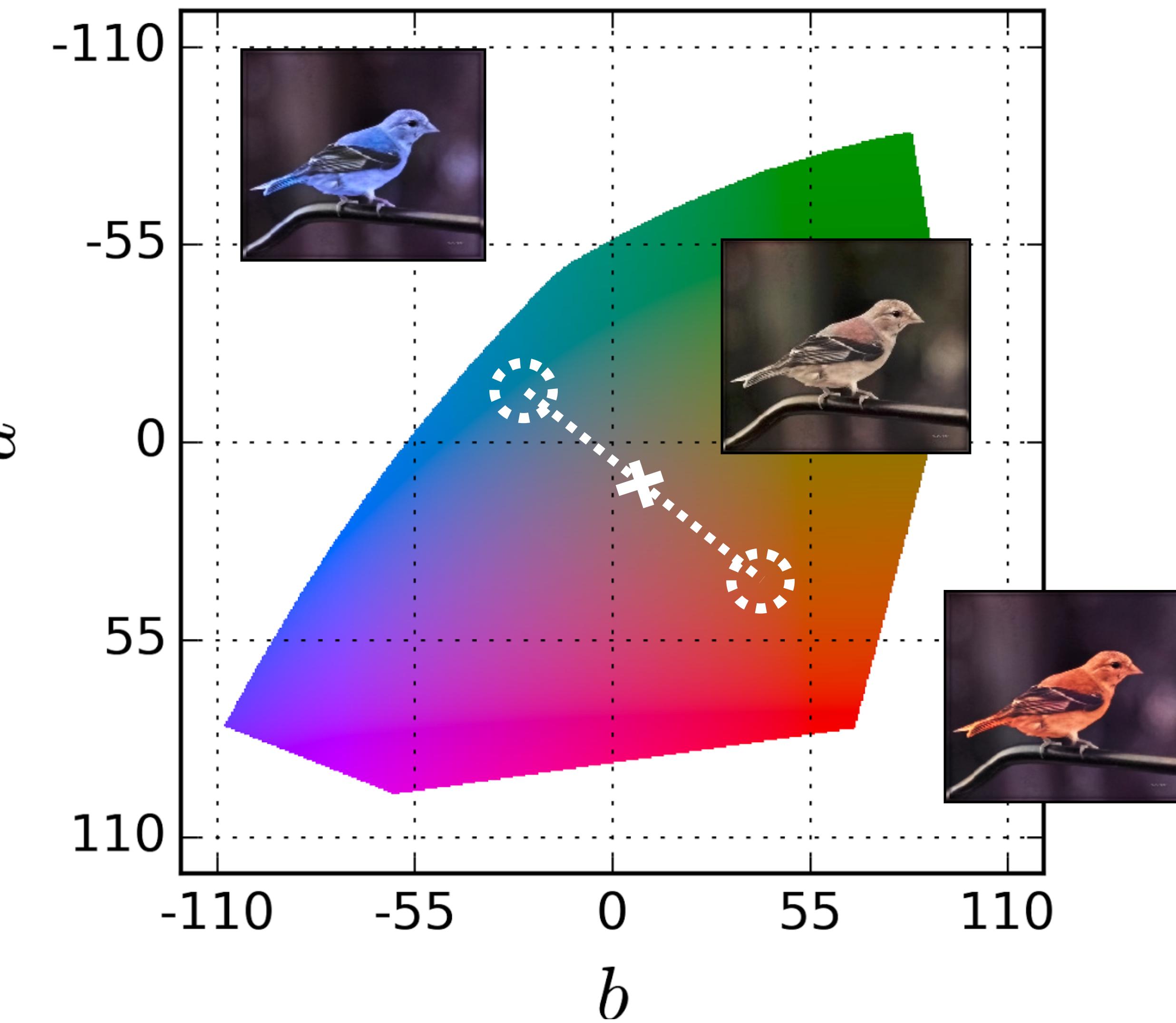
Output



Ground truth



$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$



$$L_2(\hat{\mathbf{Y}}, \mathbf{Y}) = \frac{1}{2} \sum_{h,w} \|\mathbf{Y}_{h,w} - \hat{\mathbf{Y}}_{h,w}\|_2^2$$

Designing loss functions

Input



Zhang et al. 2016



Ground truth



Color distribution cross-entropy loss with colorfulness enhancing term.



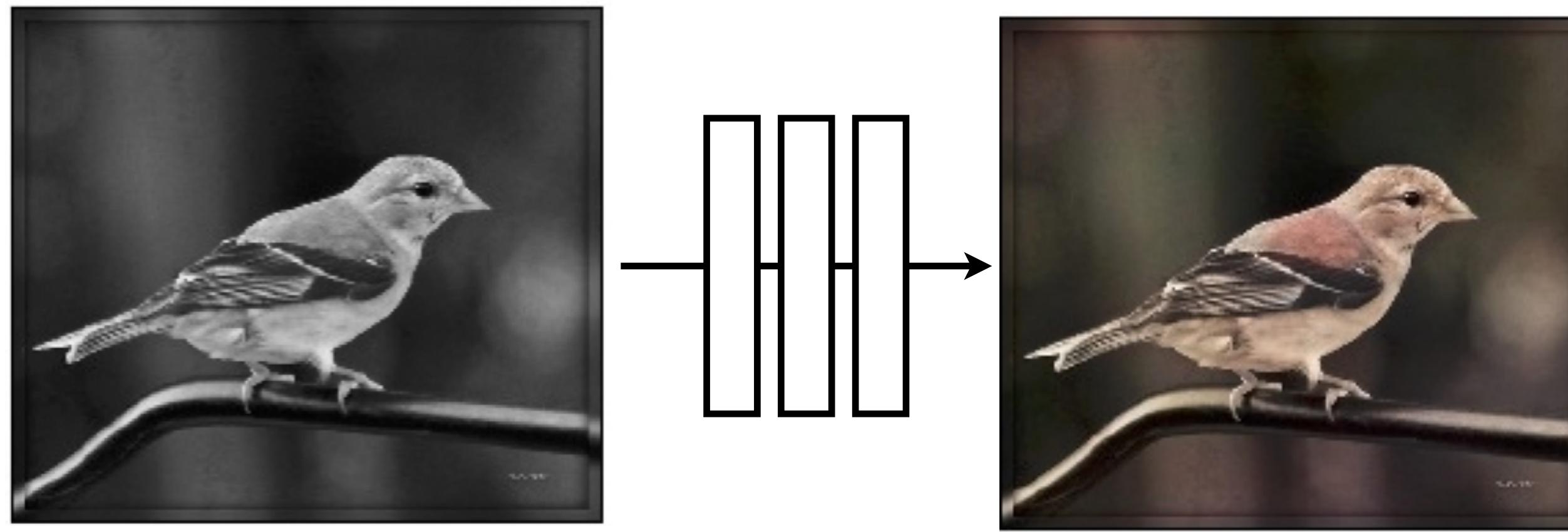
Designing loss functions



Be careful what you wish for!

Designing loss functions

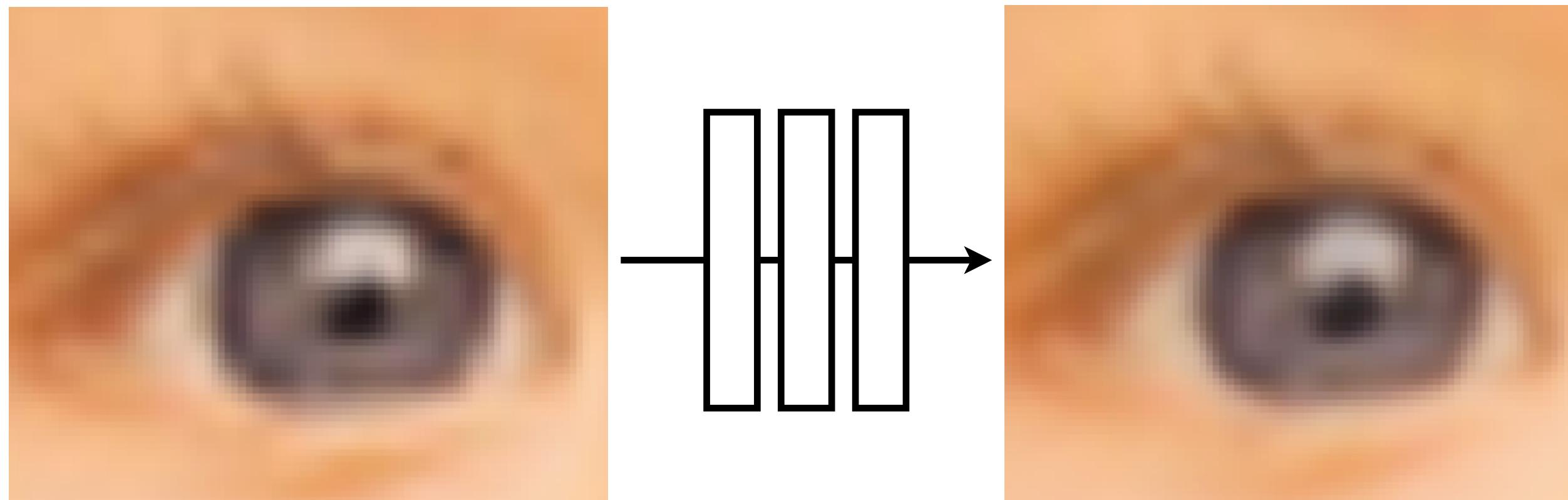
Image colorization



L2 regression

[Zhang, Isola, Efros, ECCV 2016]

Super-resolution

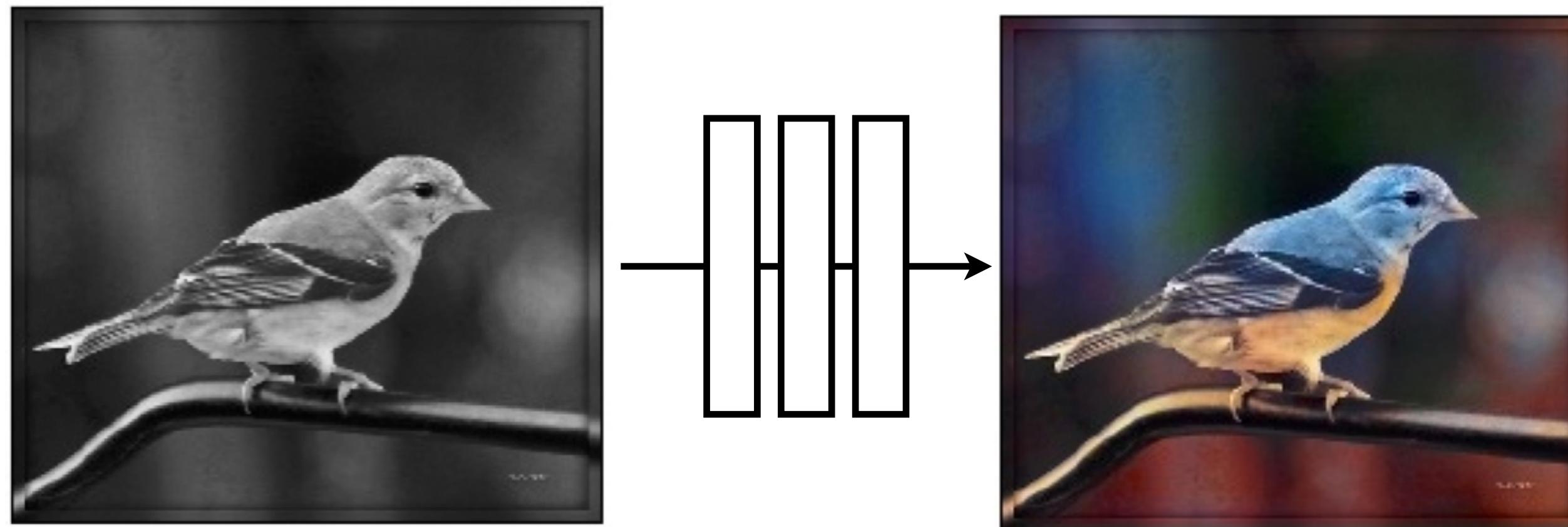


L2 regression

[Johnson, Alahi, Li, ECCV 2016]

Designing loss functions

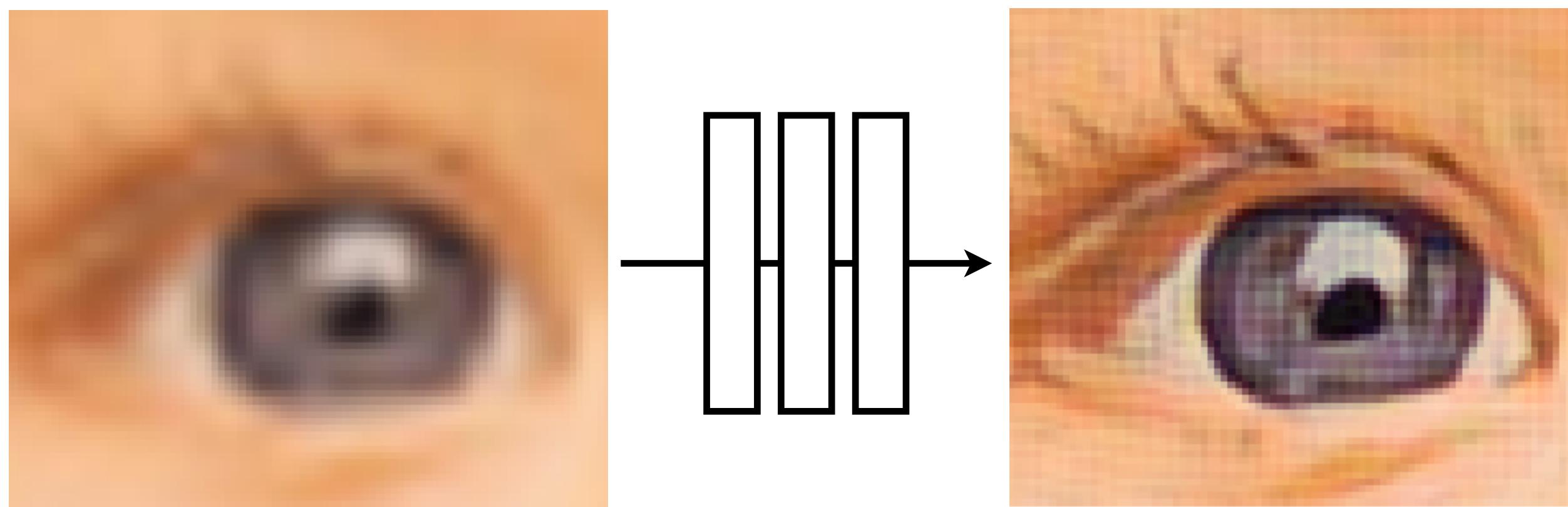
Image colorization



[Zhang, Isola, Efros, ECCV 2016]

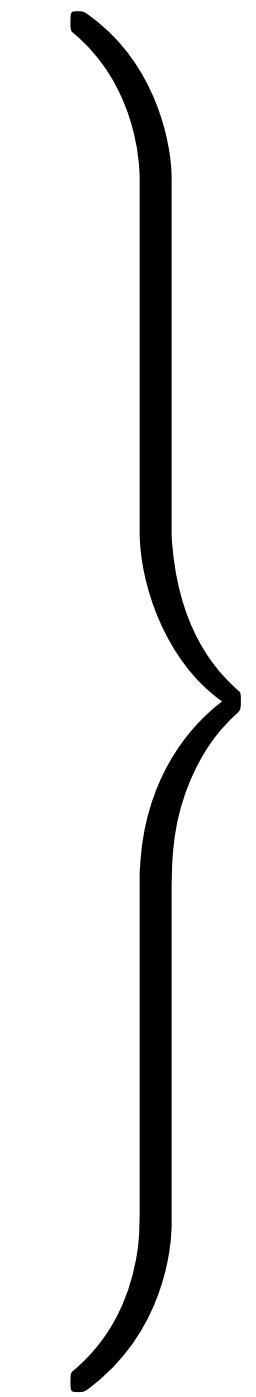
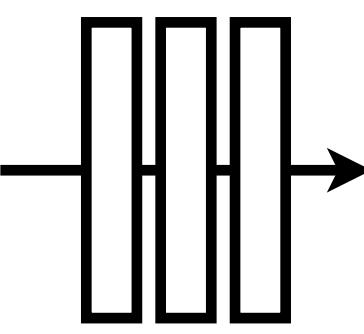
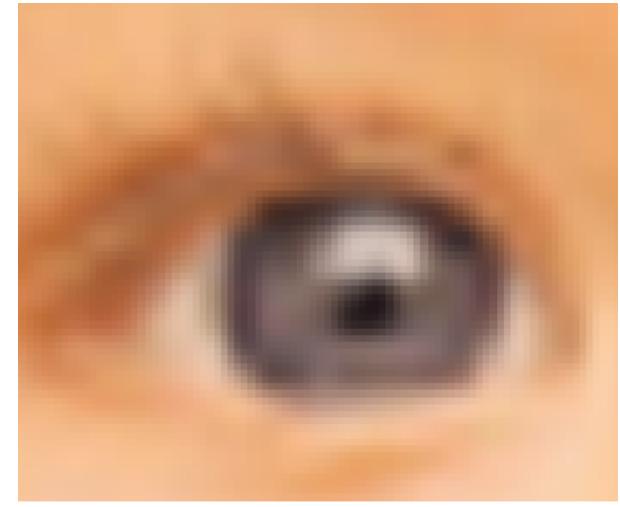
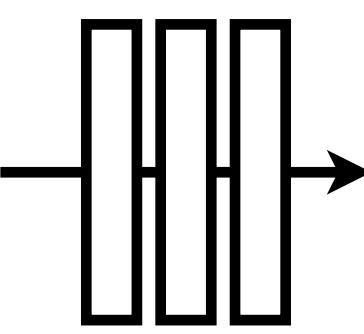
Cross entropy objective,
with colorfulness term

Super-resolution



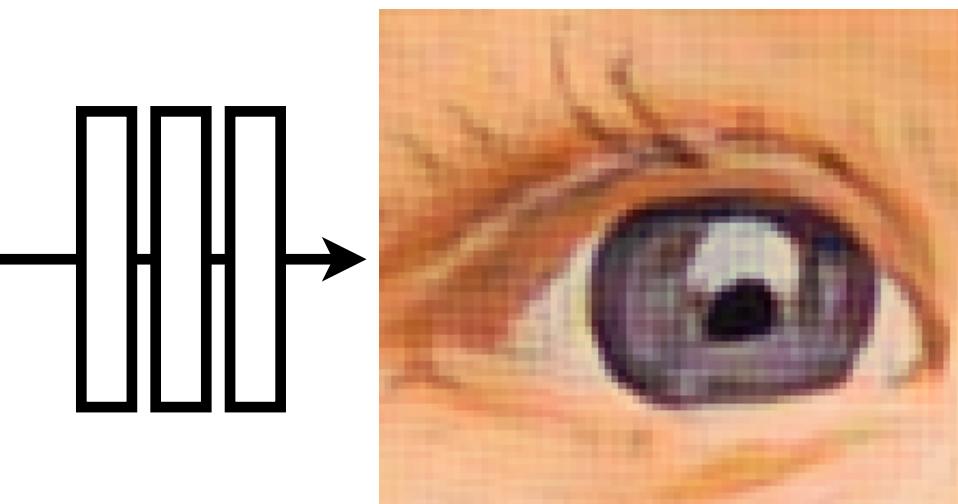
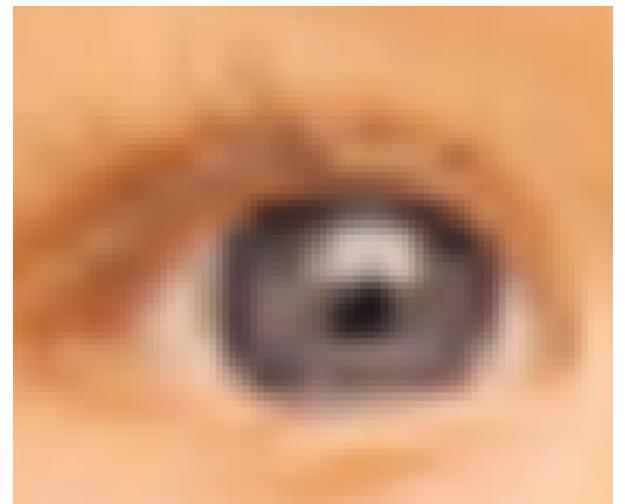
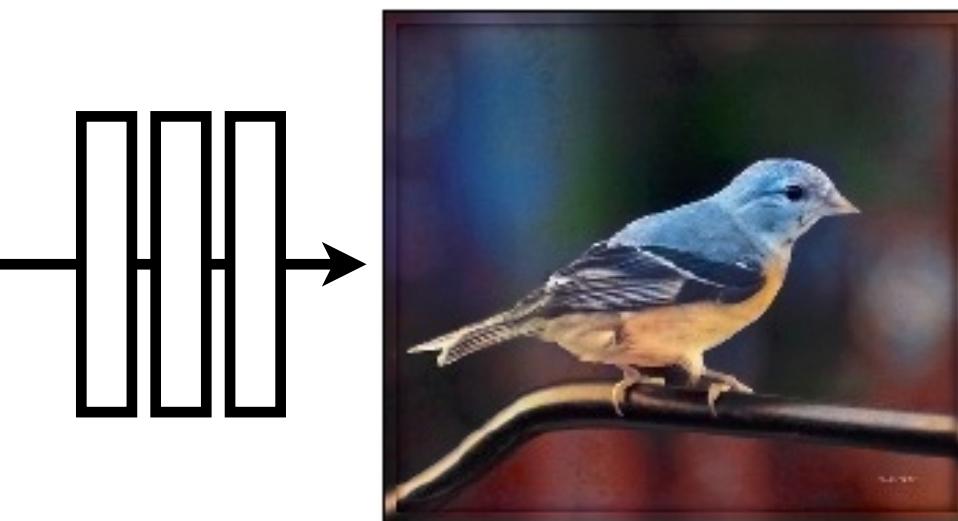
[Johnson, Alahi, Li, ECCV 2016]

Deep feature covariance
matching objective



Universal loss?

Generated images

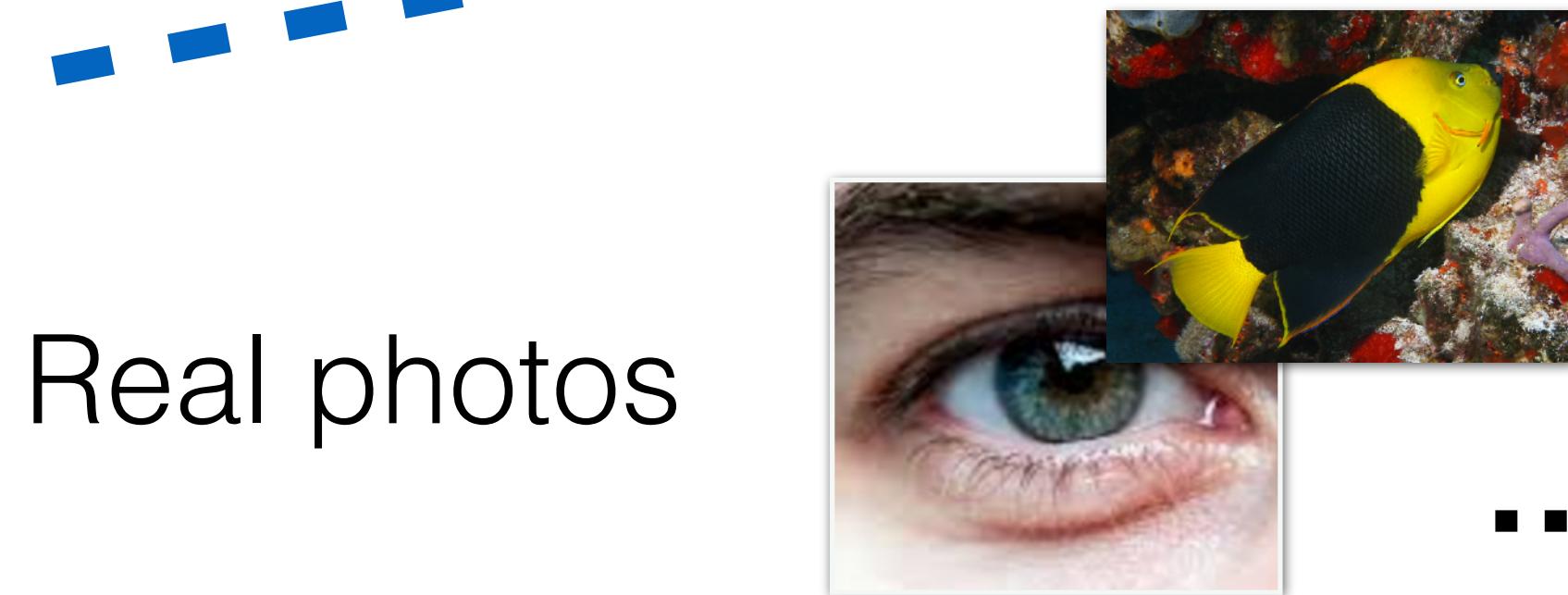


:

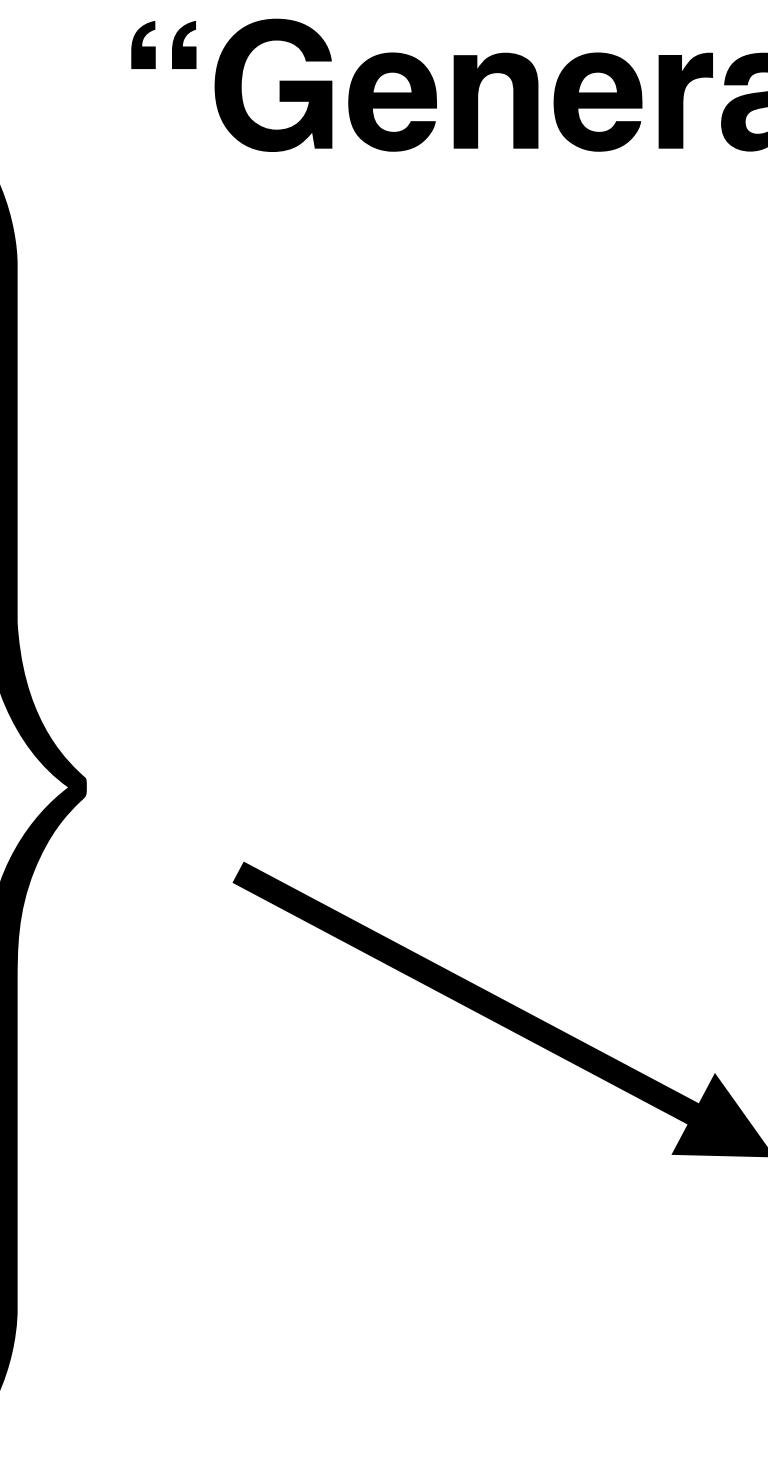
:

...

“Generative Adversarial Network” (GANs)



Real photos



Generated
vs Real
(classifier)

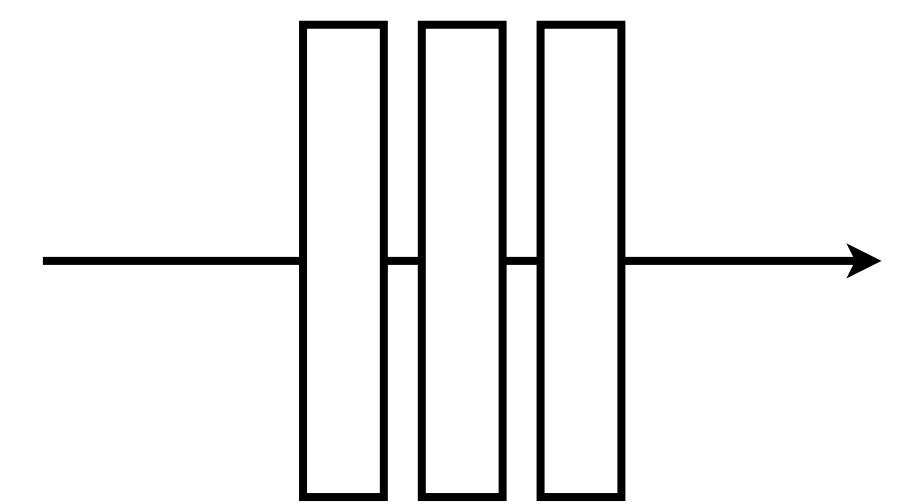


[Goodfellow, Pouget-Abadie, Mirza, Xu,
Warde-Farley, Ozair, Courville, Bengio 2014]

x



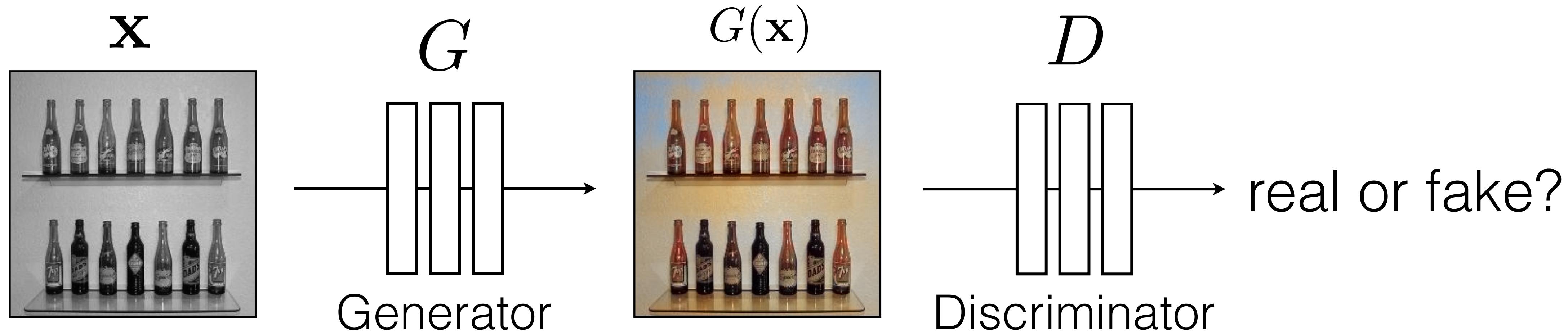
G



G(x)



Generator



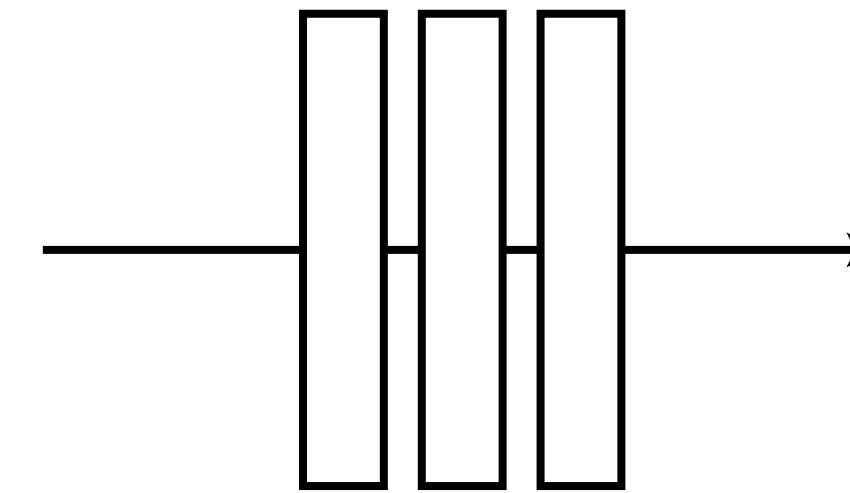
G tries to synthesize fake images that fool **D**

D tries to identify the fakes

x



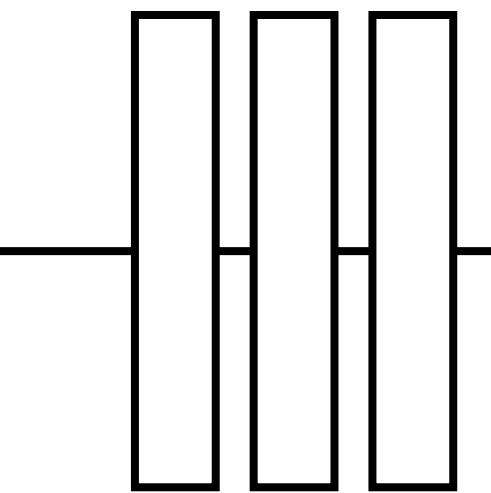
G



G(x)



D

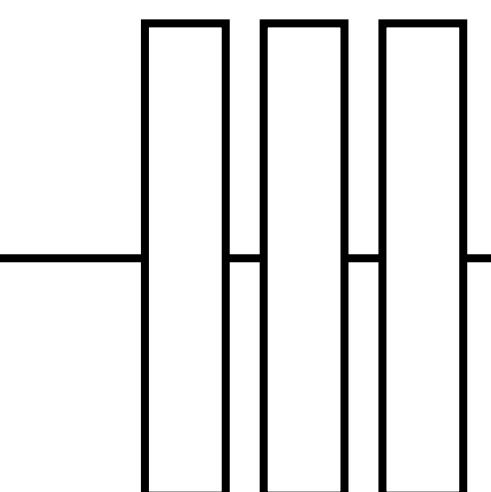


fake (0.9)

y

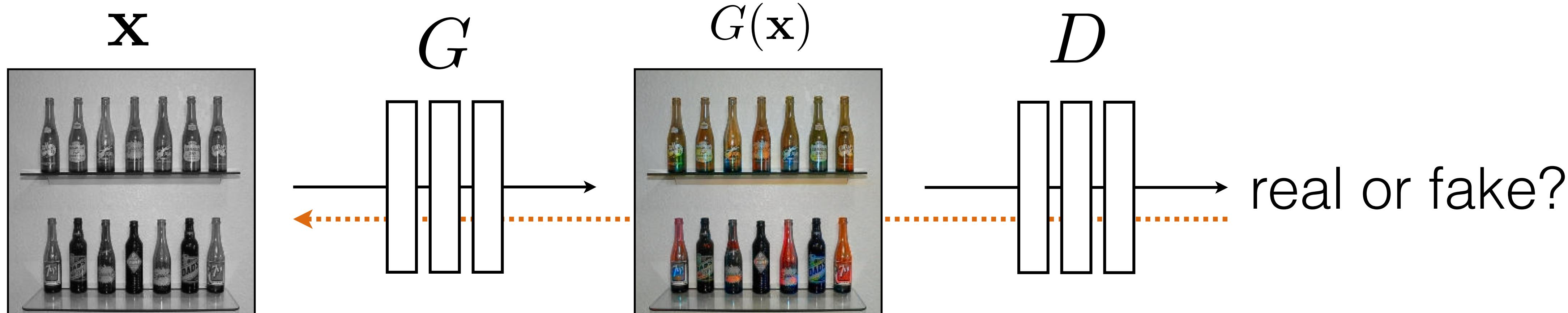


D



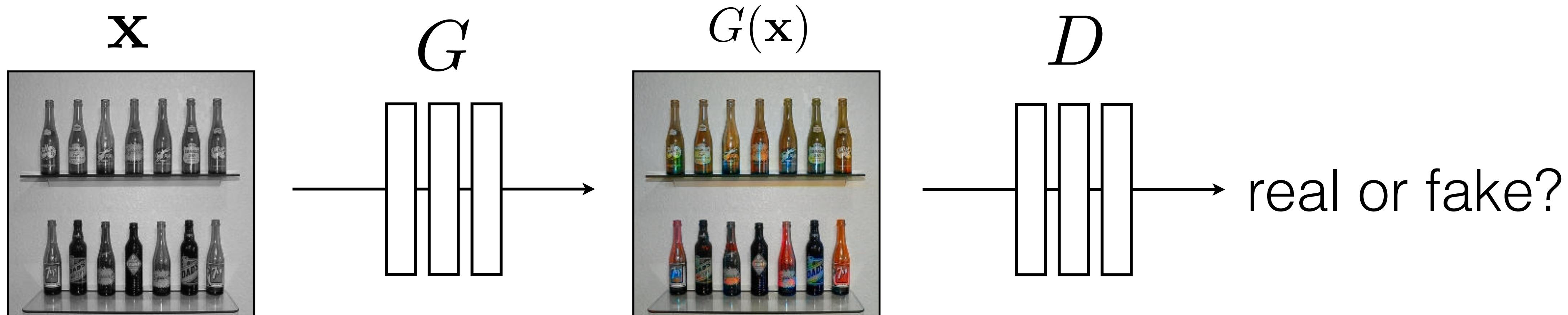
real (0.1)

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



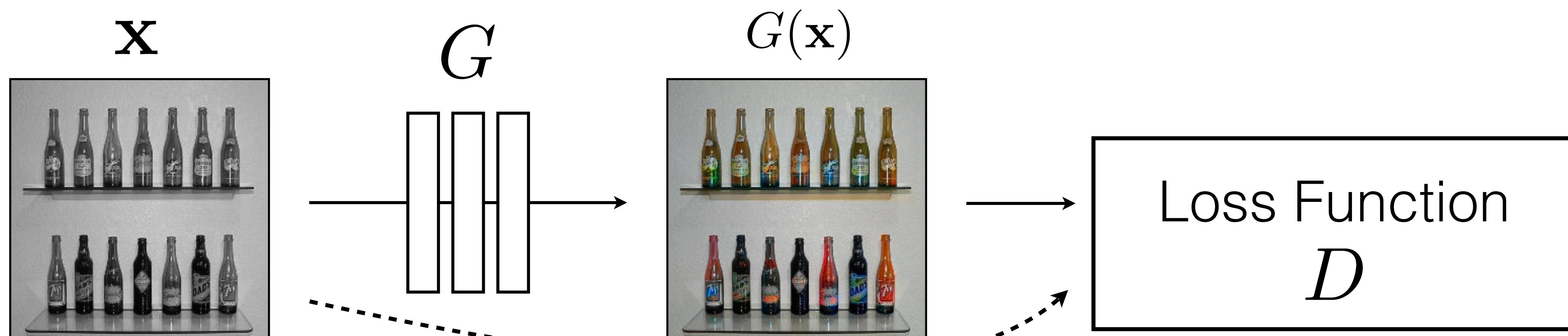
G tries to synthesize fake images that **fool** **D**:

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



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

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



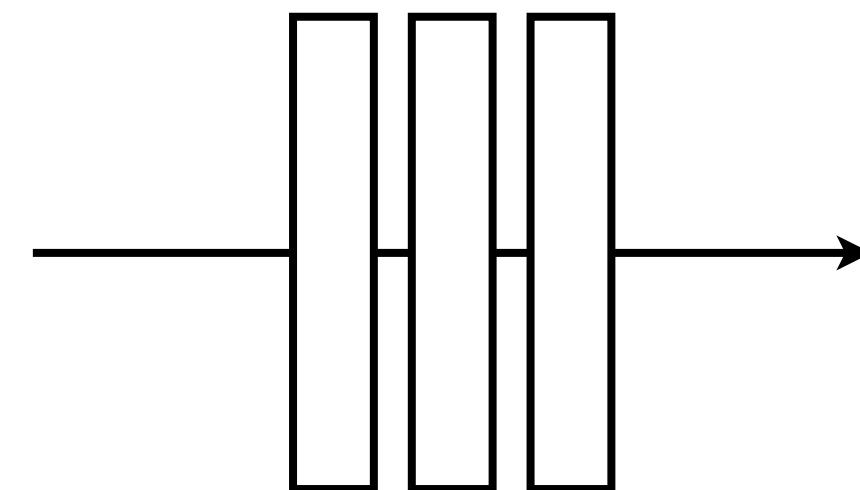
G's perspective: **D** is a loss function.

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

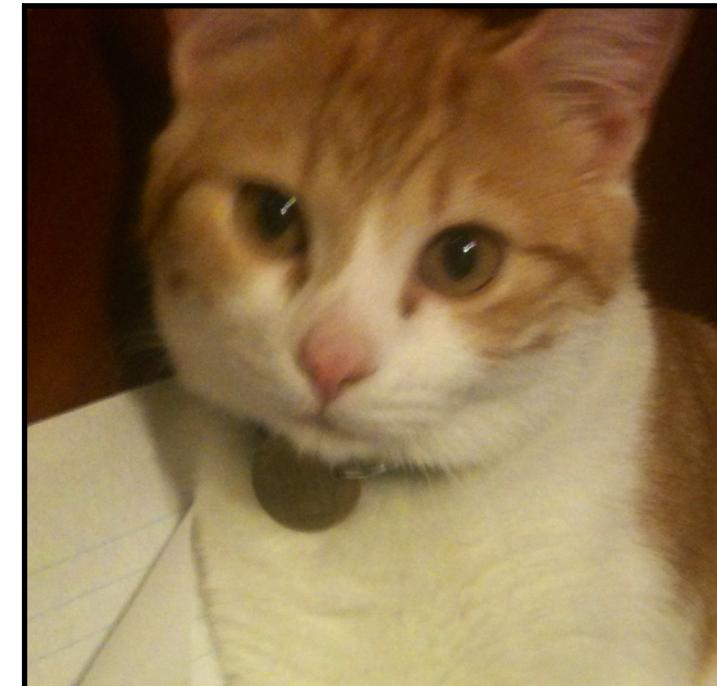
x



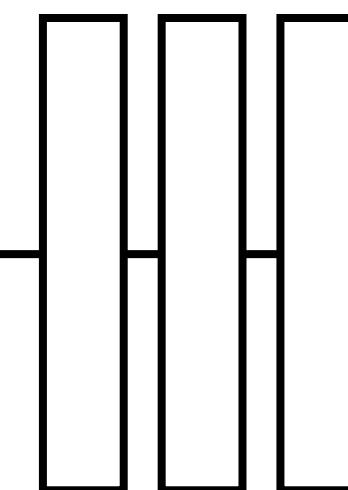
G



G(x)



D



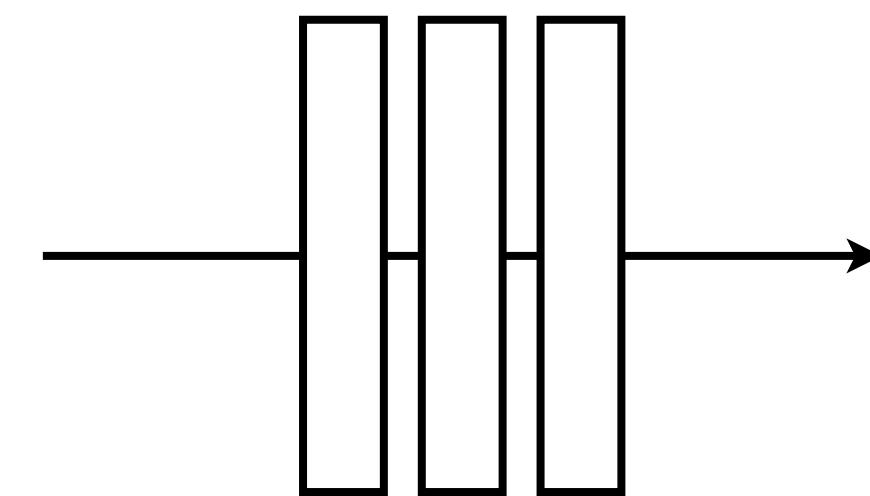
real or fake?

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

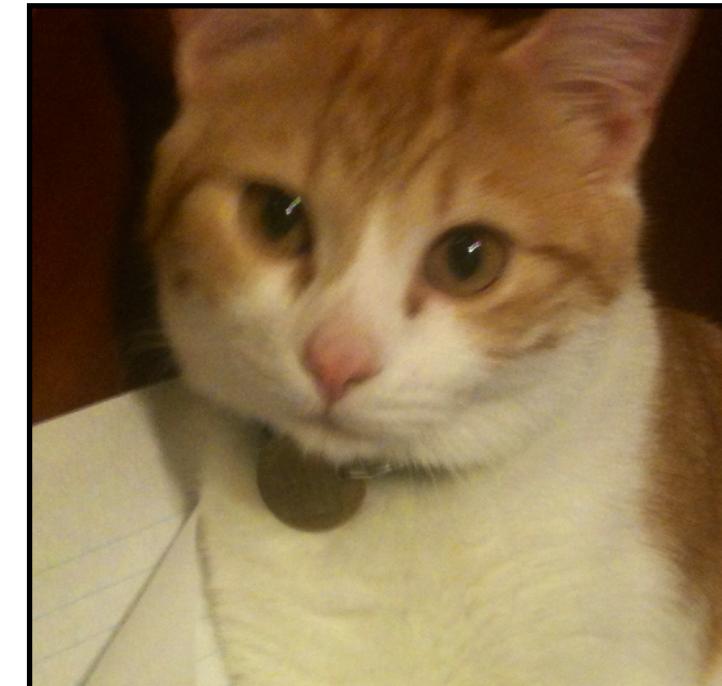
x



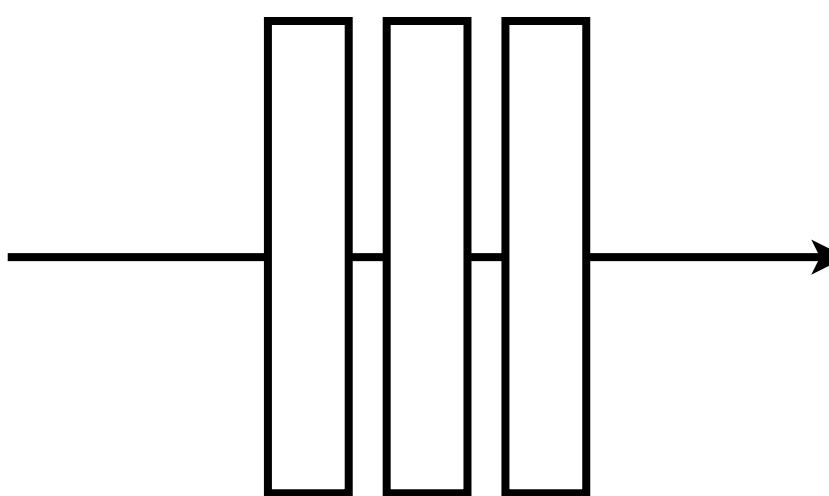
G



G(x)



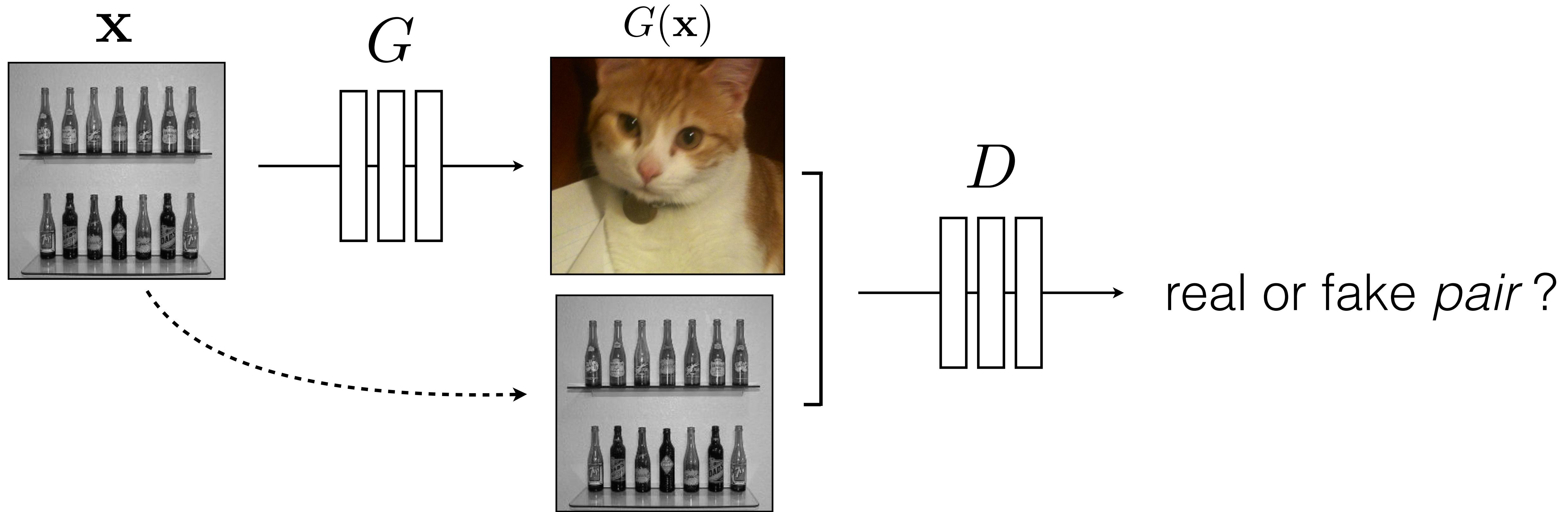
D



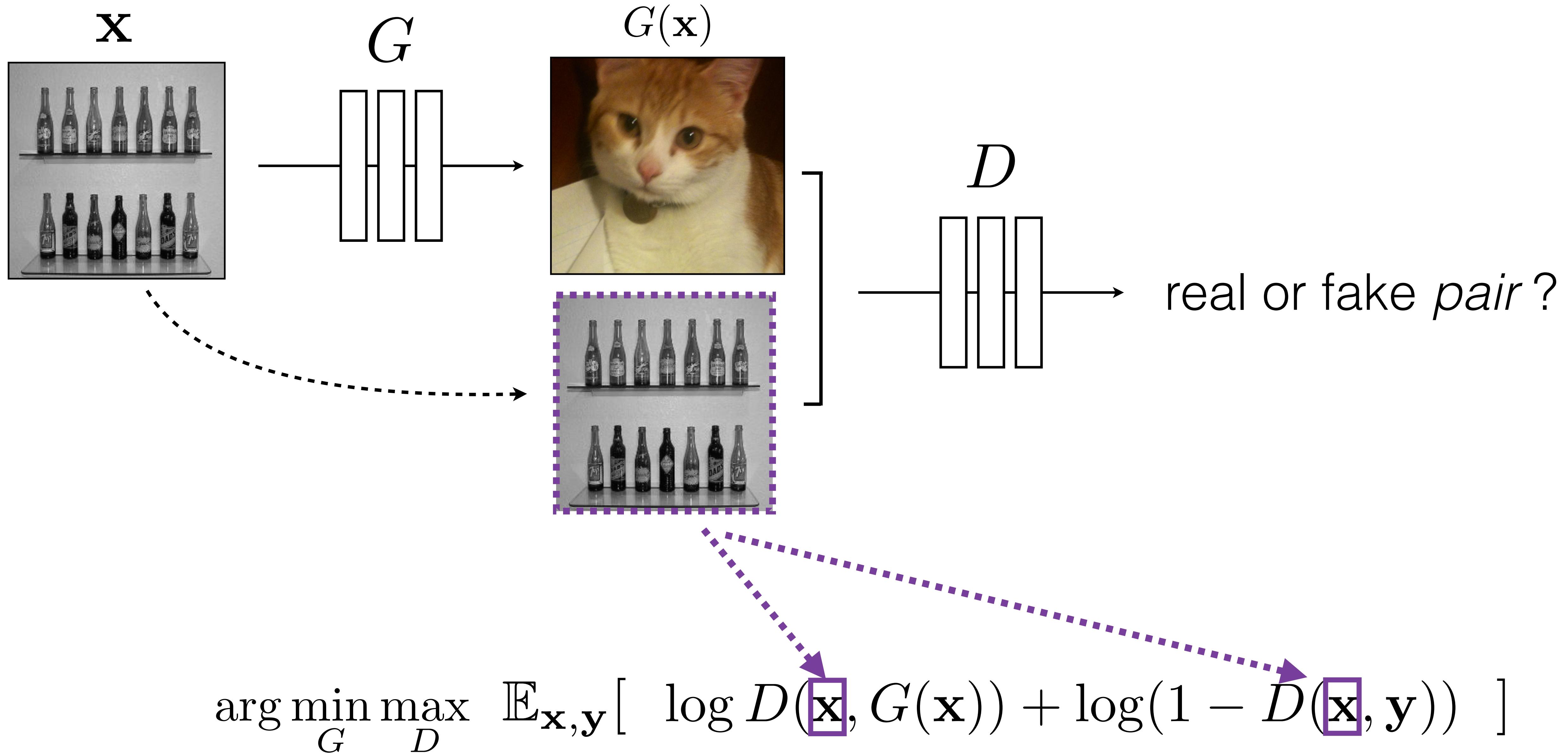
real!

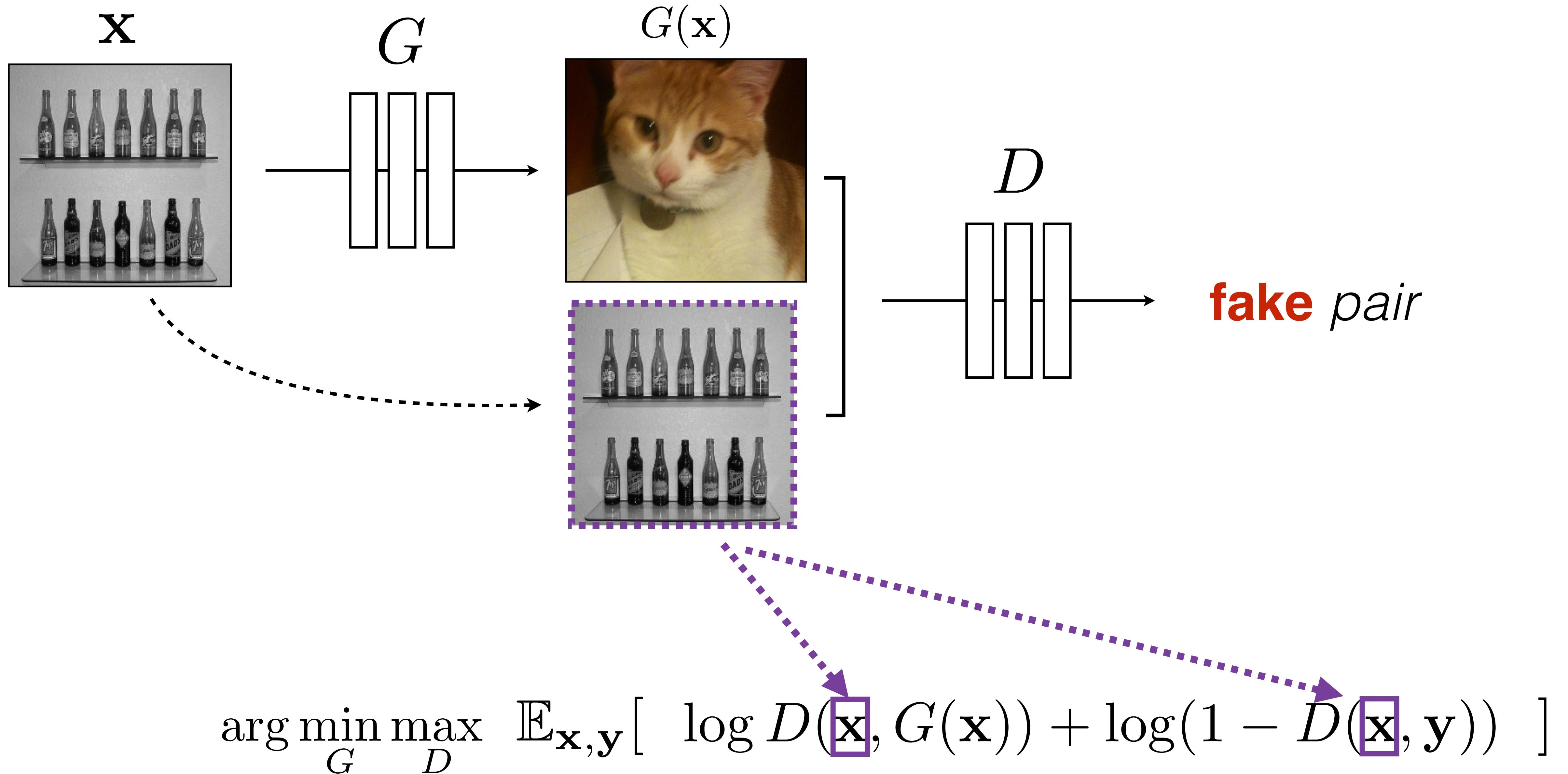
("Aquarius")

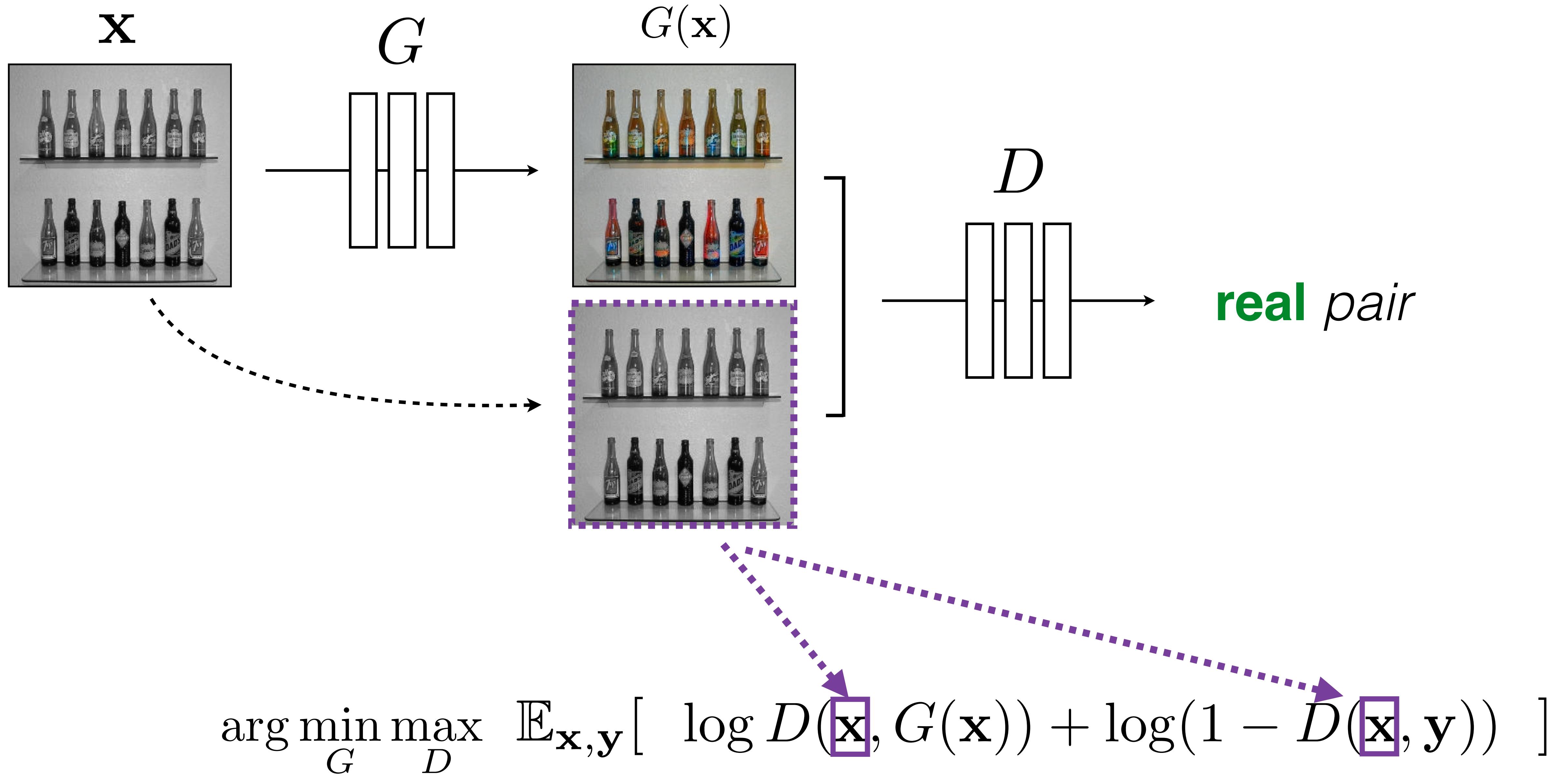
$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$

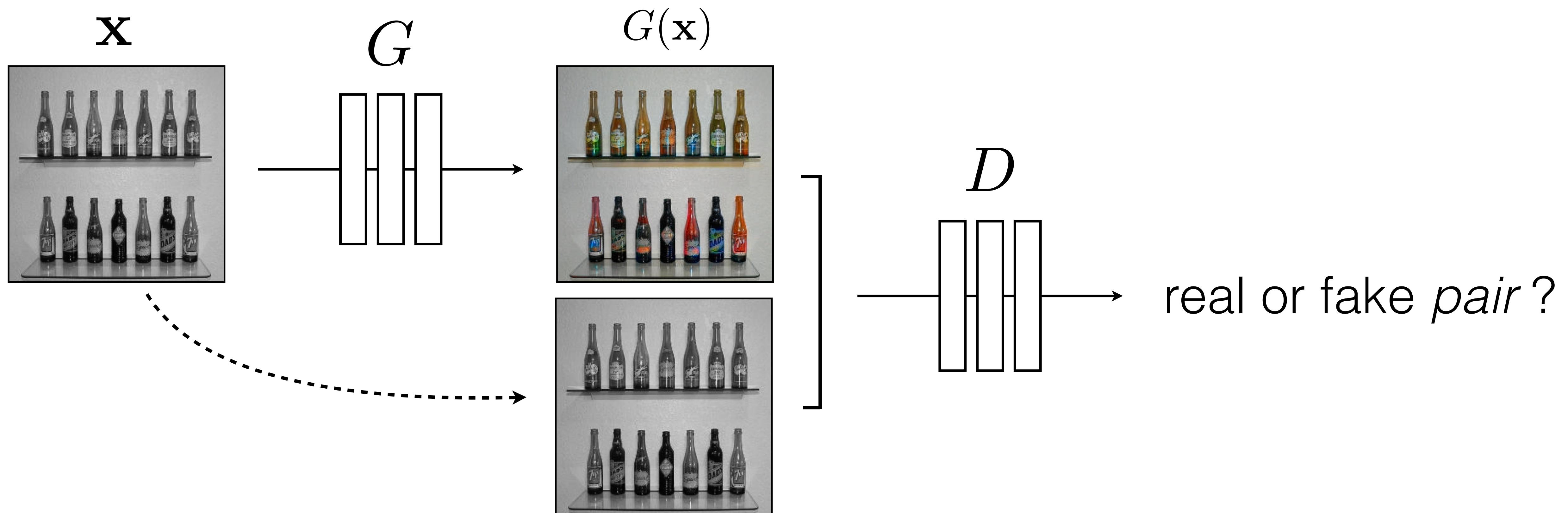


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$









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

Training Details: Loss function

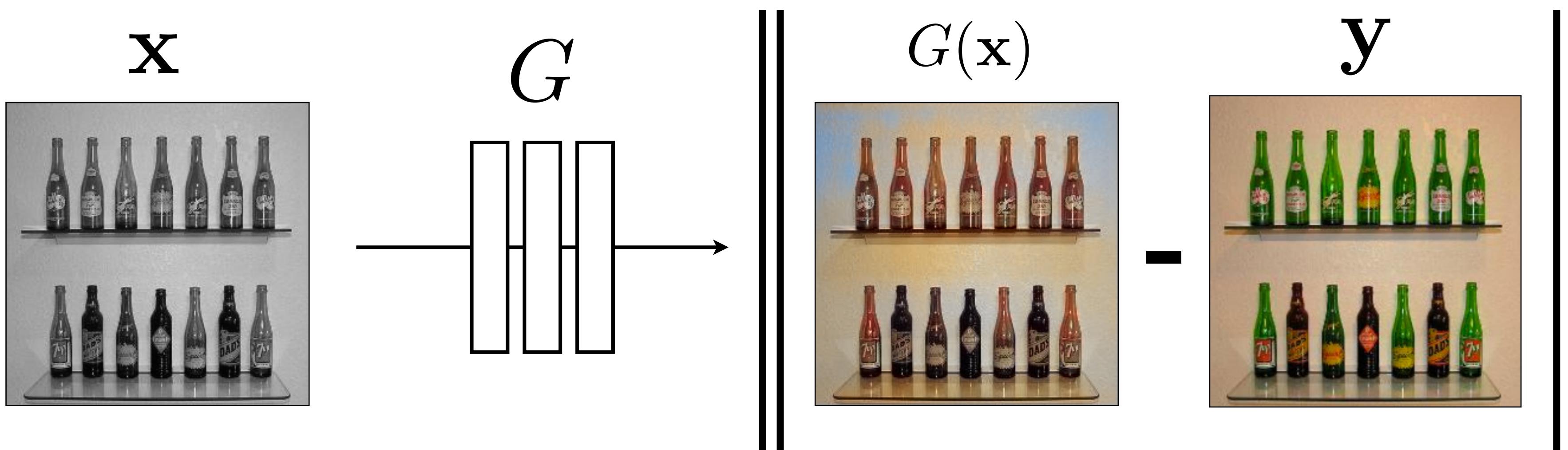
Conditional GAN

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

Training Details: Loss function

Conditional GAN

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



Stable training + fast convergence

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

BW → Color

Input



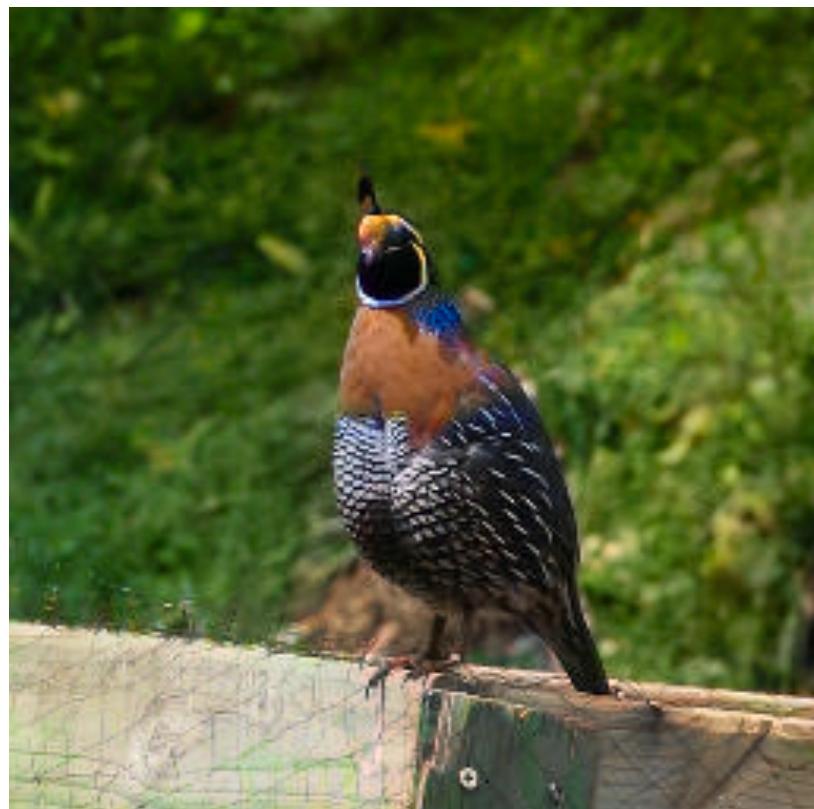
Output



Input



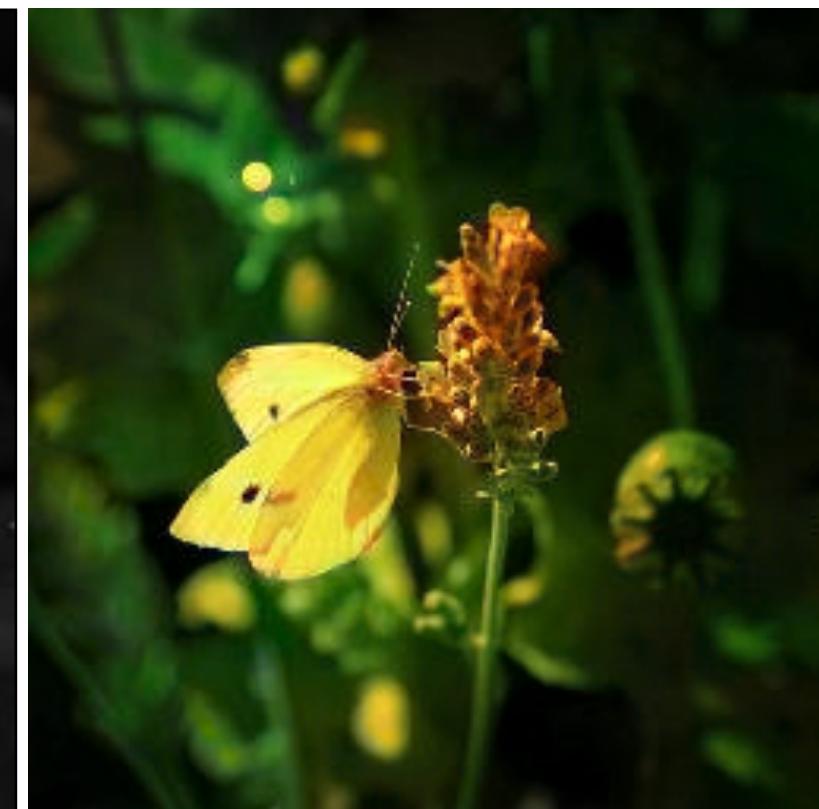
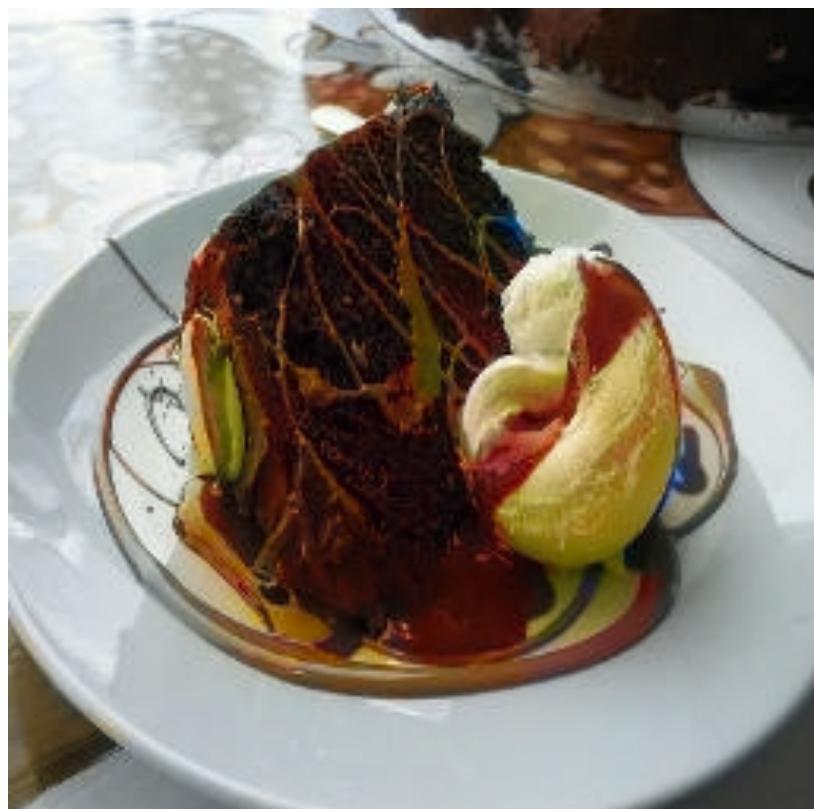
Output



Input

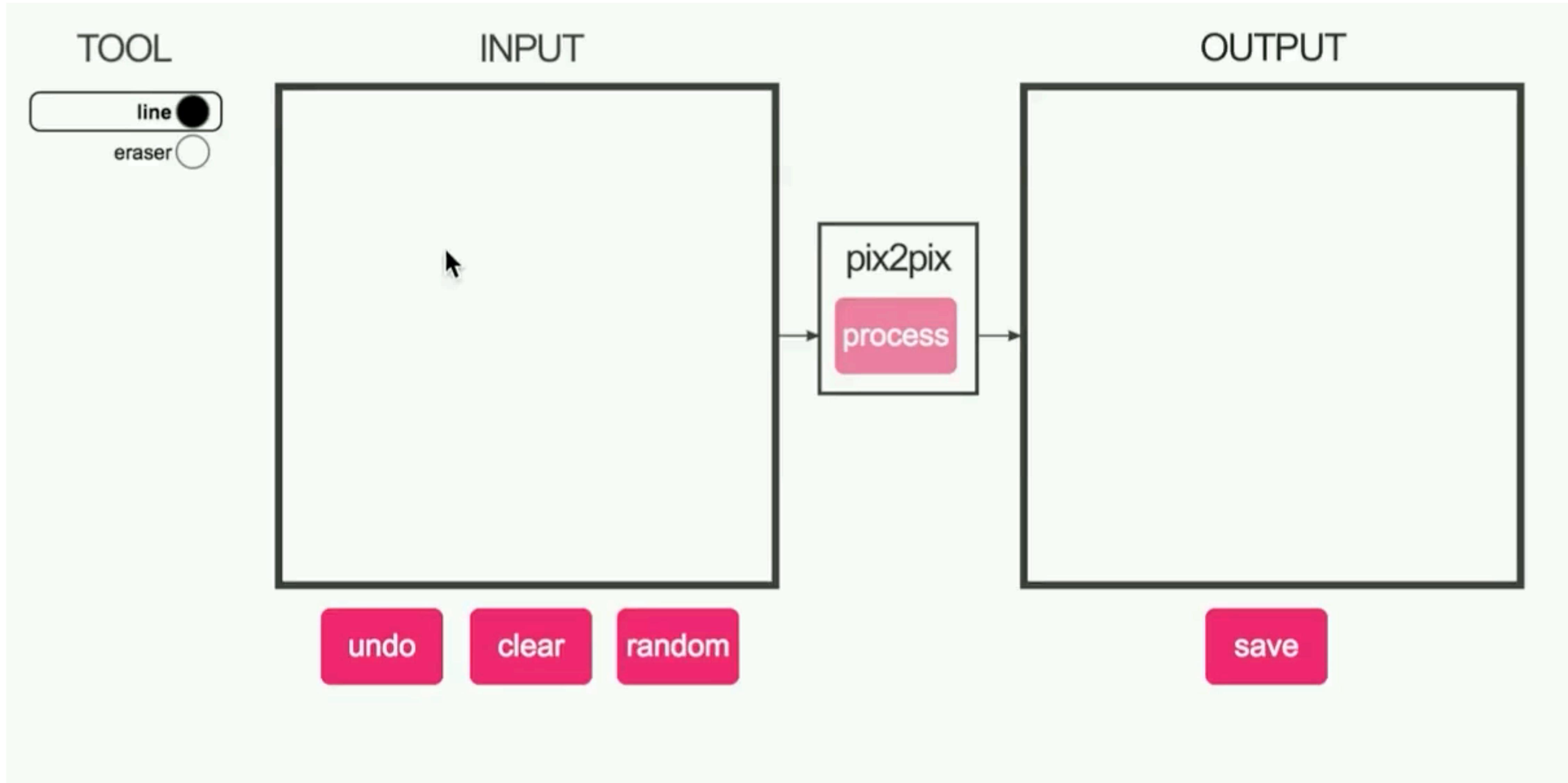


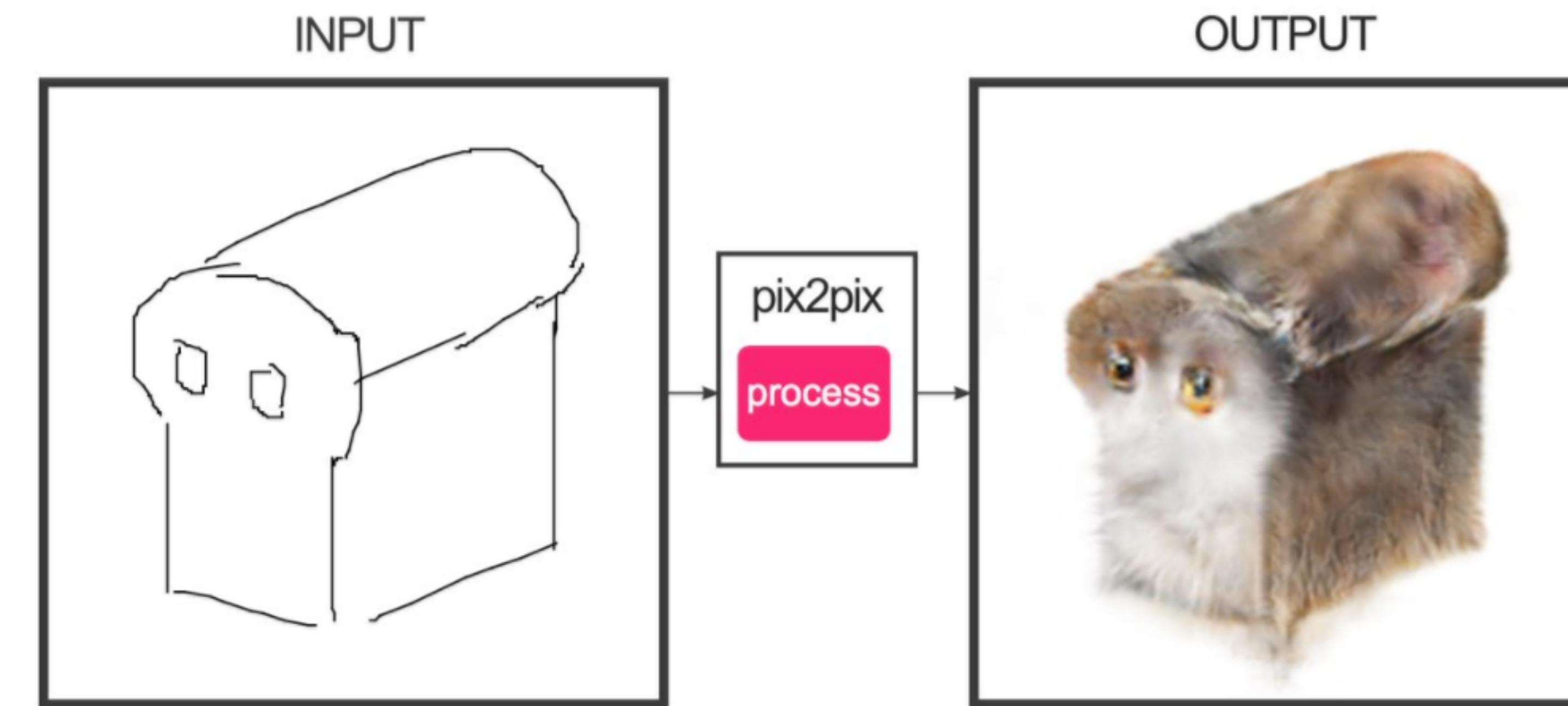
Output



Data from [Russakovsky et al. 2015]

#edges2cats [Chris Hesse]





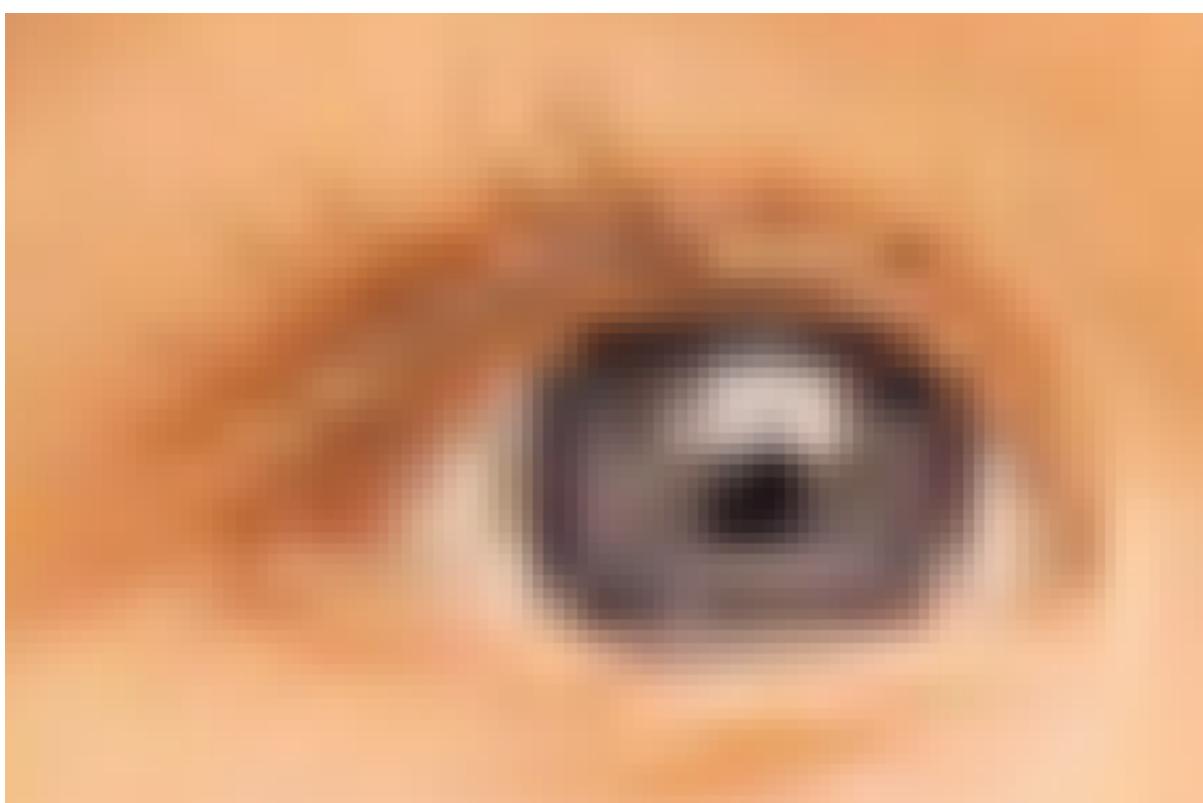
Ivy Tasi @ivymyt



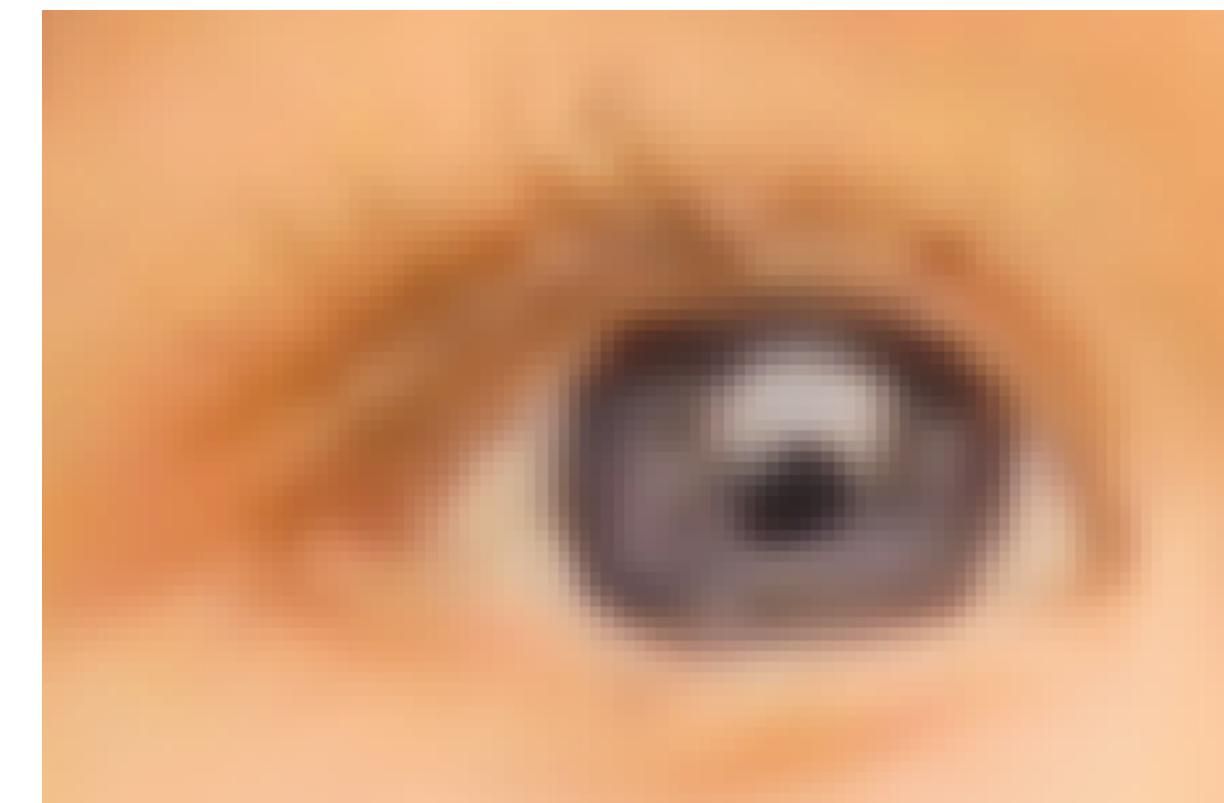
Vitaly Vidmirov @vvid

Structured Prediction

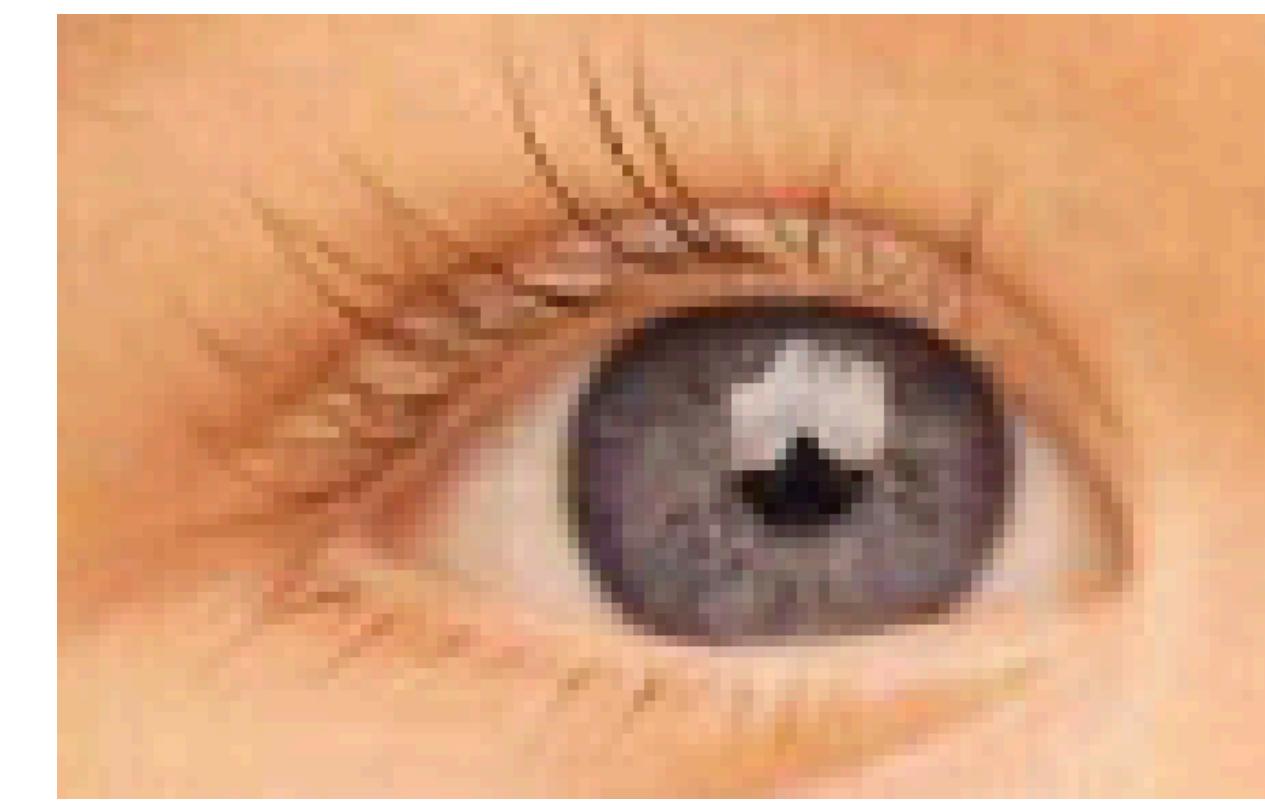
Input
 \mathbf{x}



Output
 $\hat{\mathbf{y}}$

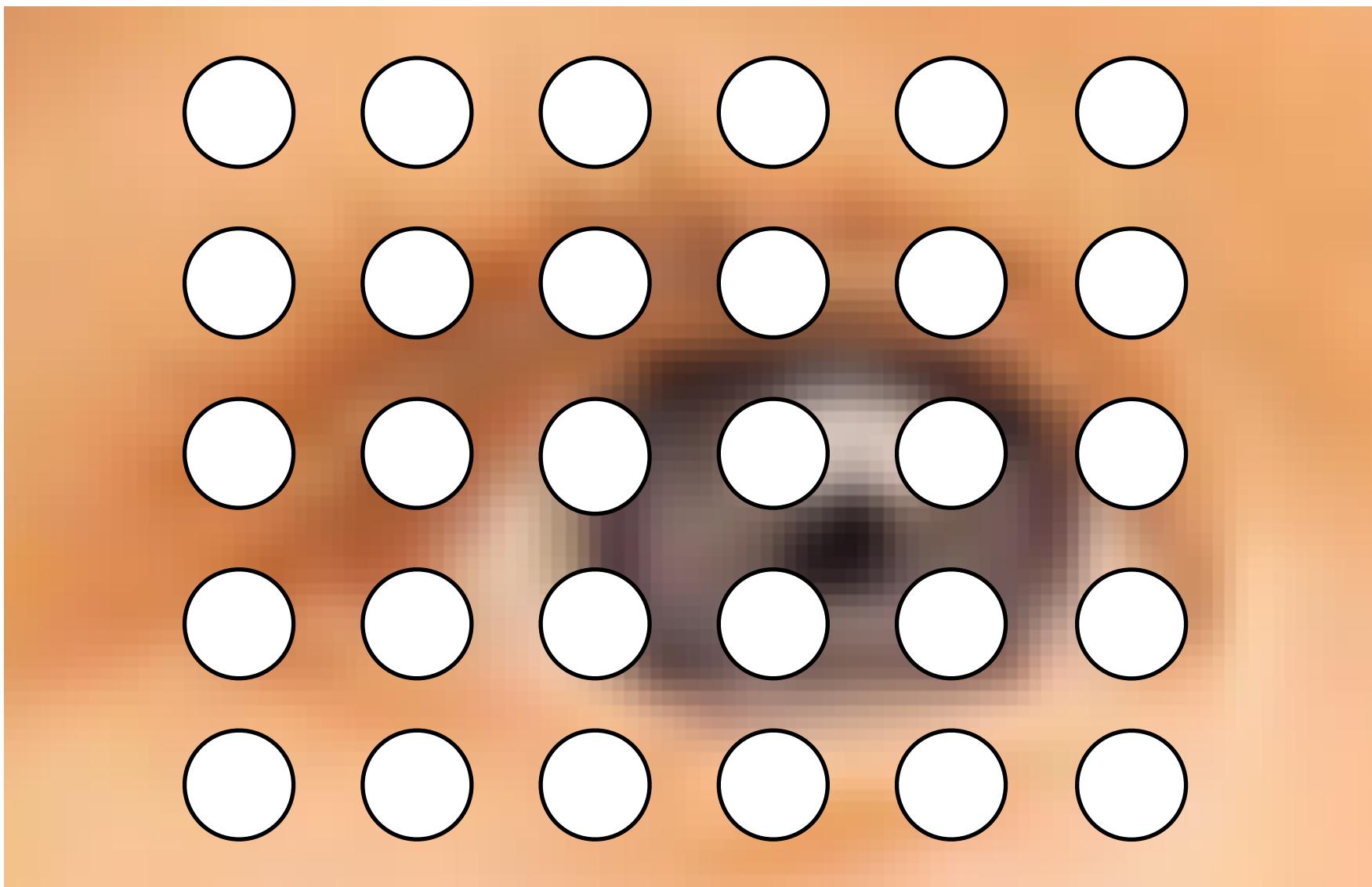


Target
 \mathbf{y}



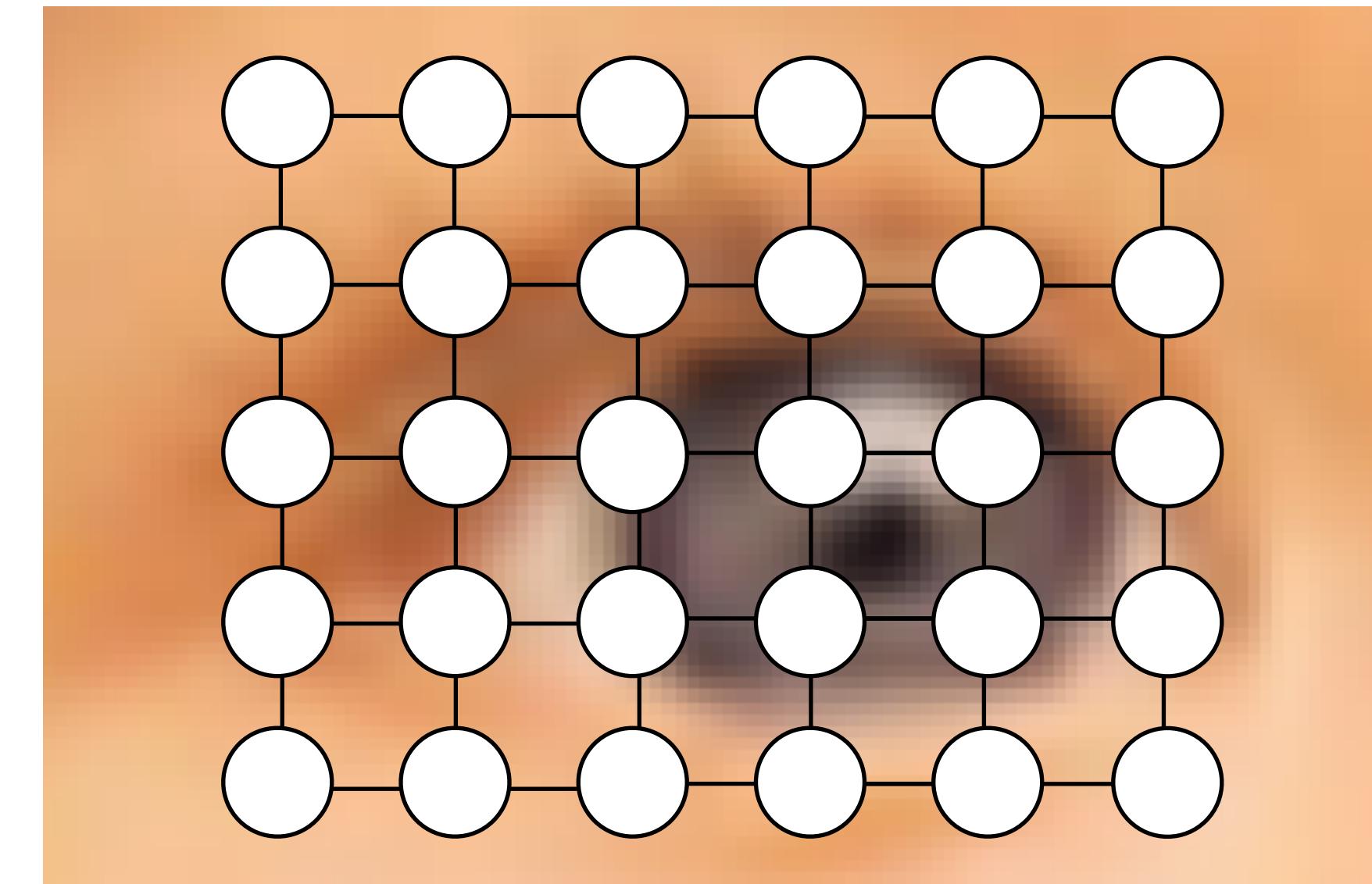
$$L(\hat{\mathbf{y}}, \mathbf{y}) = \|\hat{\mathbf{y}} - \mathbf{y}\|_2$$

Structured Prediction



Each pixel treated as
independent

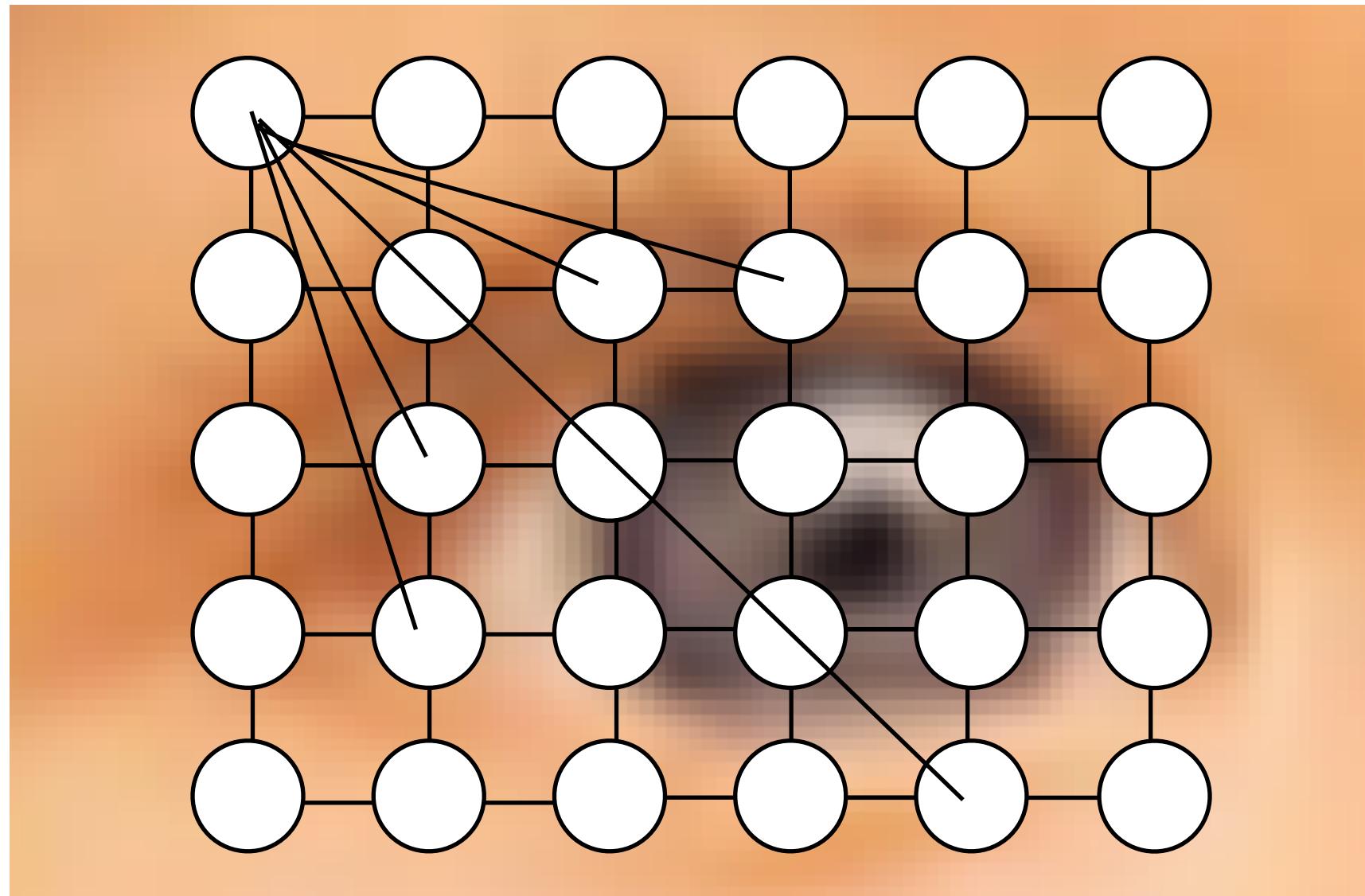
$$\prod_i p(y_i | \mathbf{x})$$



Models at pairwise configuration
of pixels

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

Structured Prediction



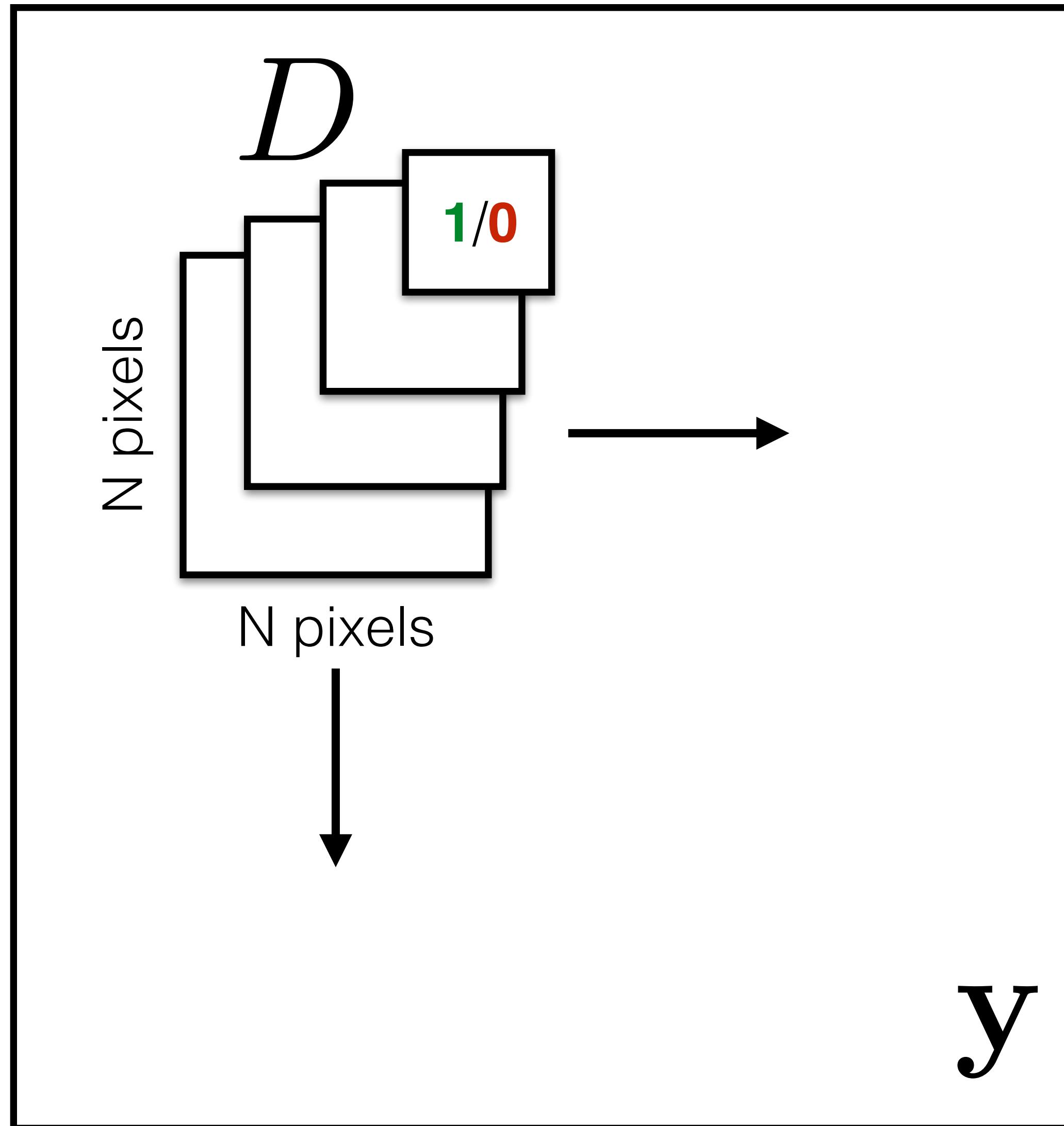
Model *joint* configuration
of all pixels

$$p(\mathbf{y}|\mathbf{x})$$

A GAN, with sufficient capacity, samples from the full joint distribution when perfectly optimized.

Most generative models have this property! Give them **sufficient capacity** and **infinite data**, and they are the complete solution to prediction problems.

Shrinking the capacity: Patch Discriminator



Rather than penalizing if output *image* looks fake, penalize if each overlapping *patch* in output looks fake

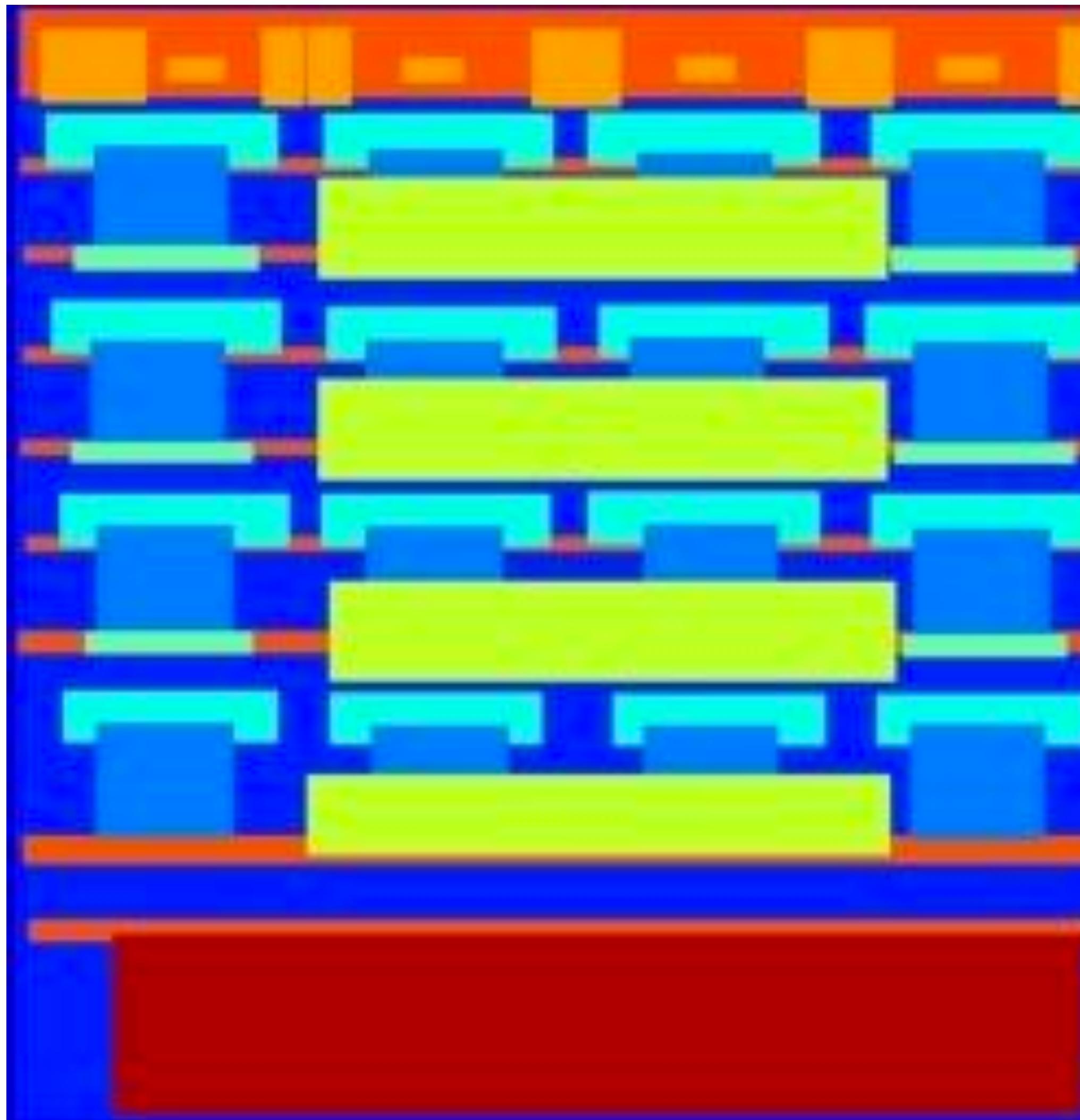
[Li & Wand 2016]

[Shrivastava et al. 2017]

[Isola et al. 2017]

Labels → Facades

Input



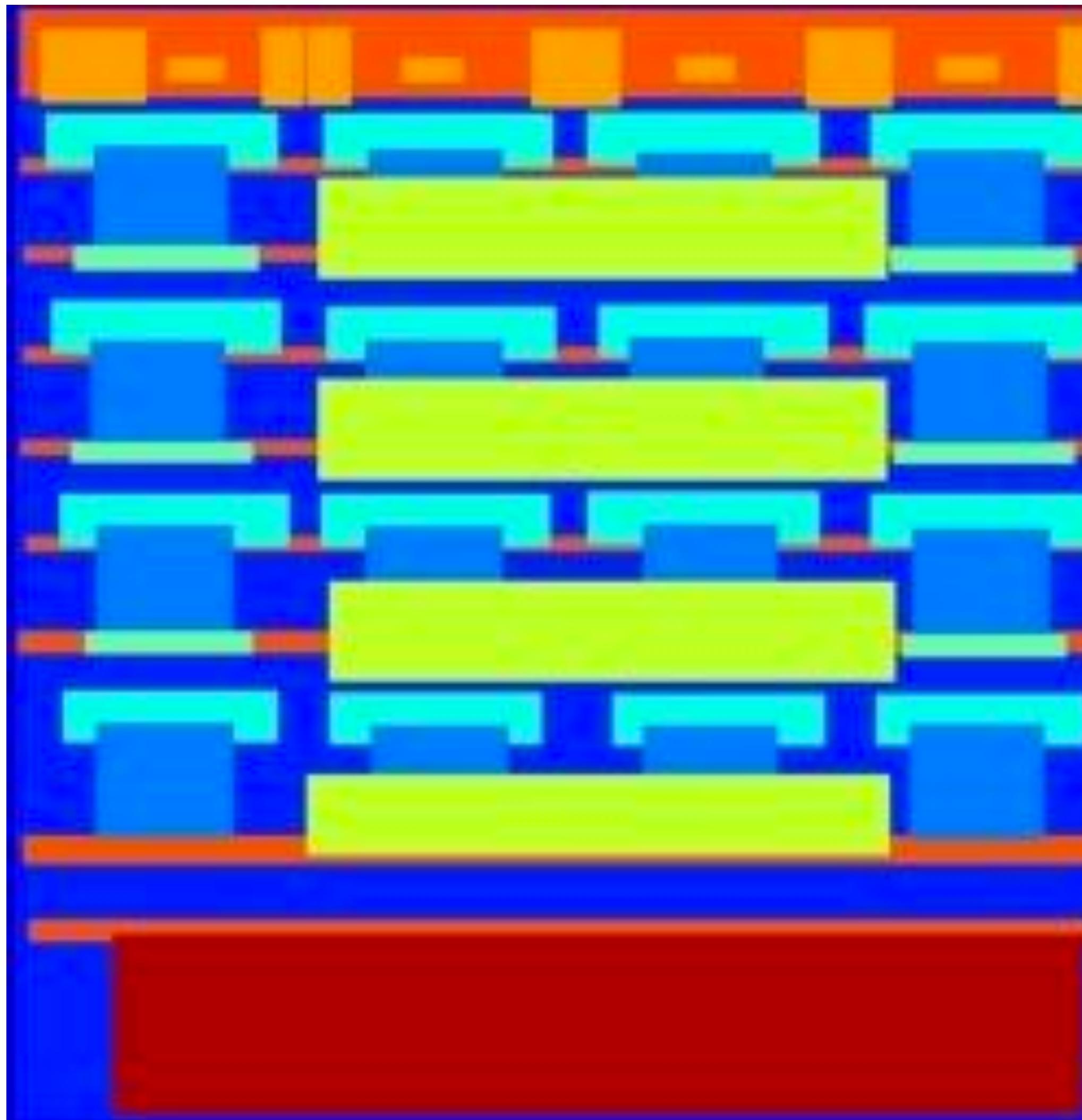
1x1 Discriminator



Data from [Tylecek, 2013]

Labels → Facades

Input



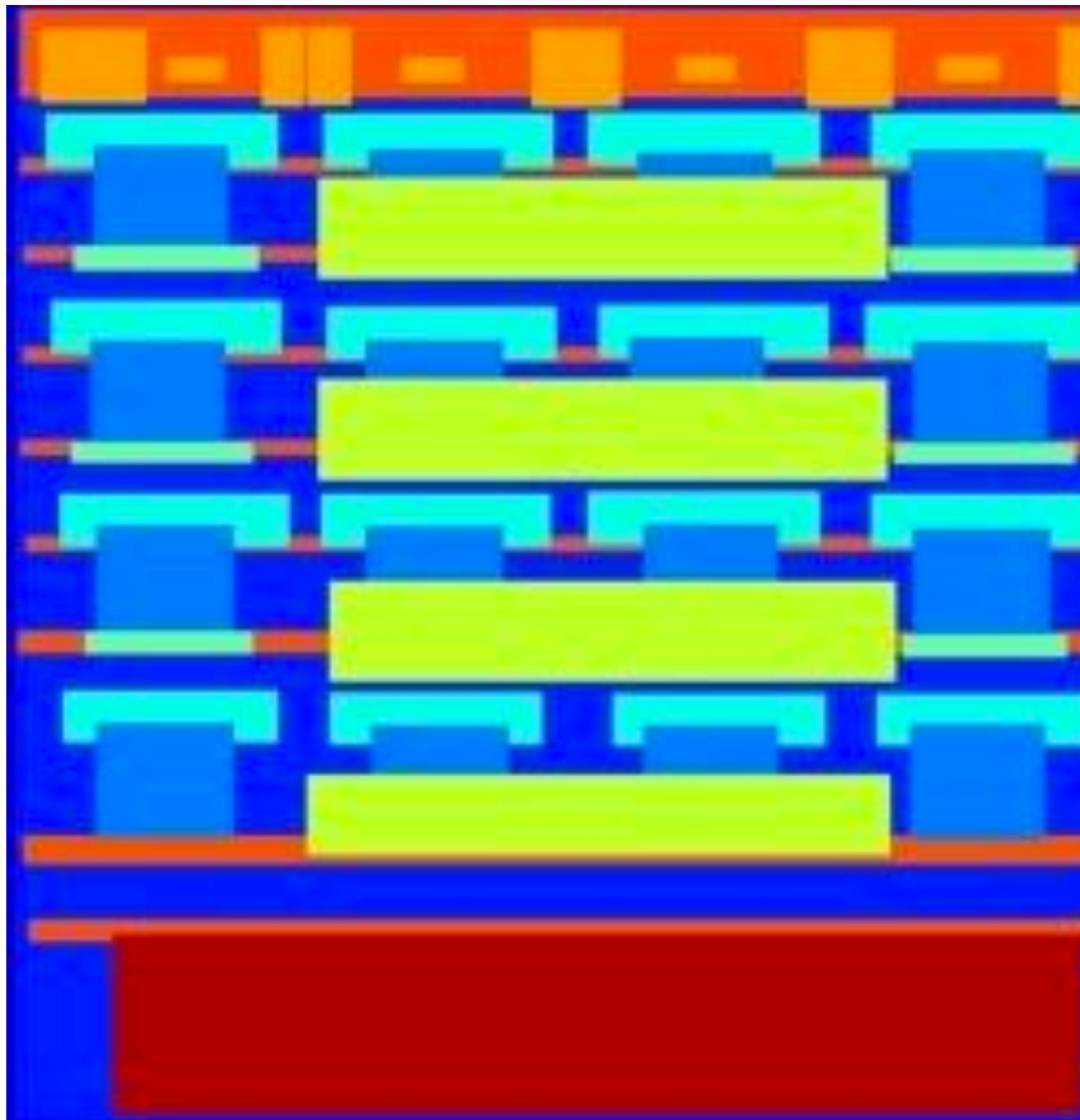
16x16 Discriminator



Data from [Tylecek, 2013]

Labels → Facades

Input



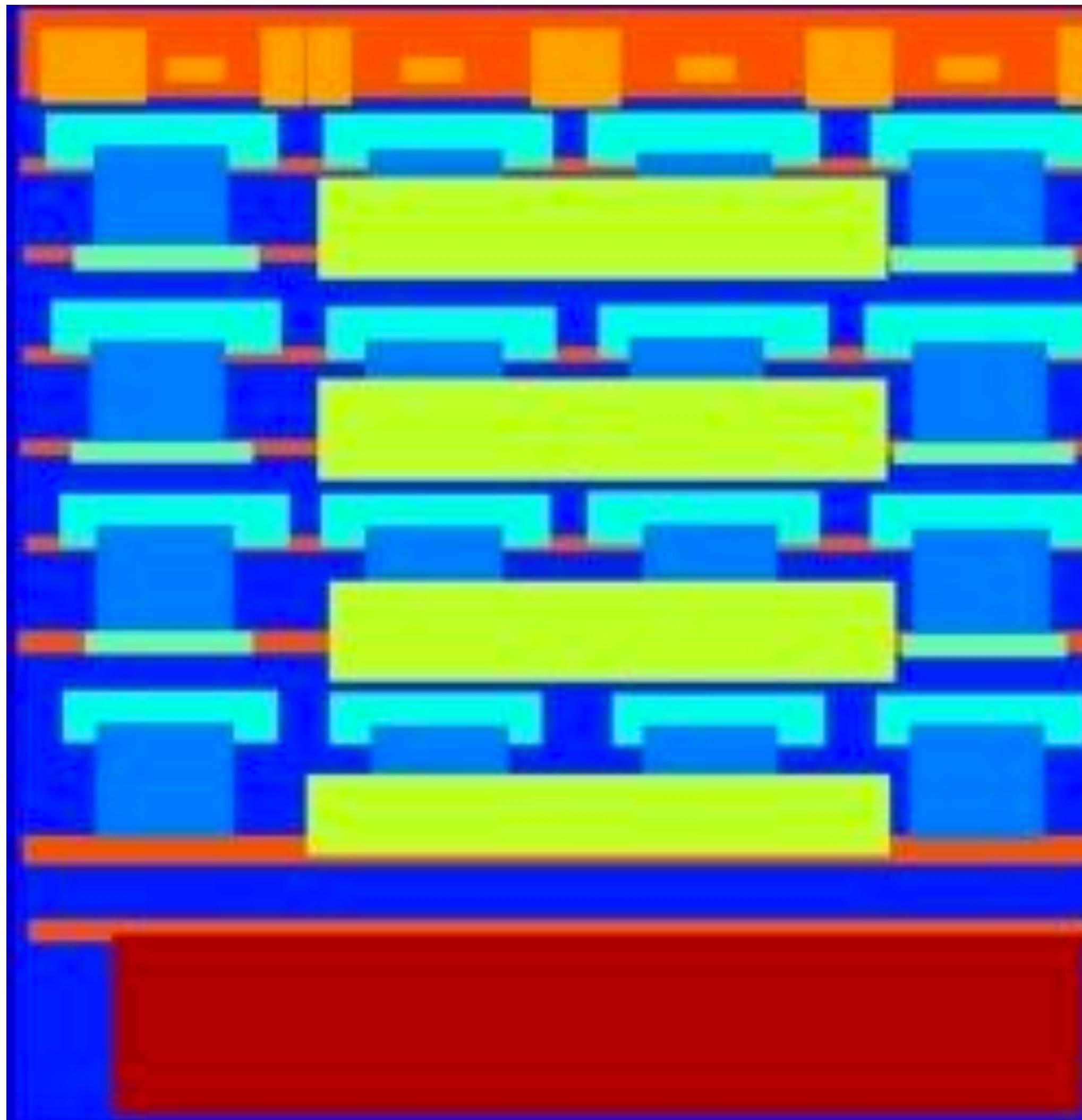
70x70 Discriminator



Data from [Tylecek, 2013]

Labels → Facades

Input

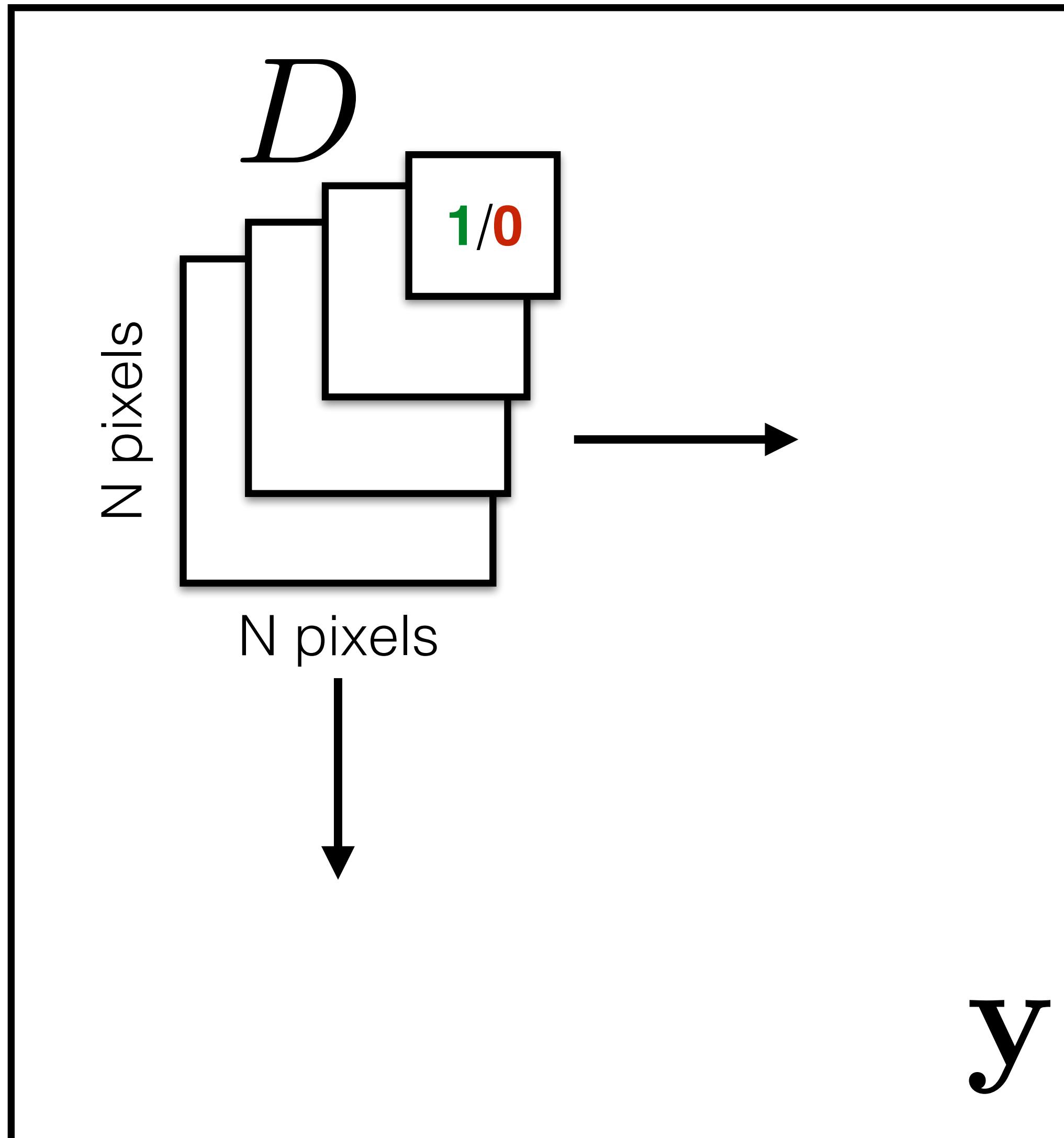


Full image Discriminator



Data from [Tylecek, 2013]

Patch Discriminator



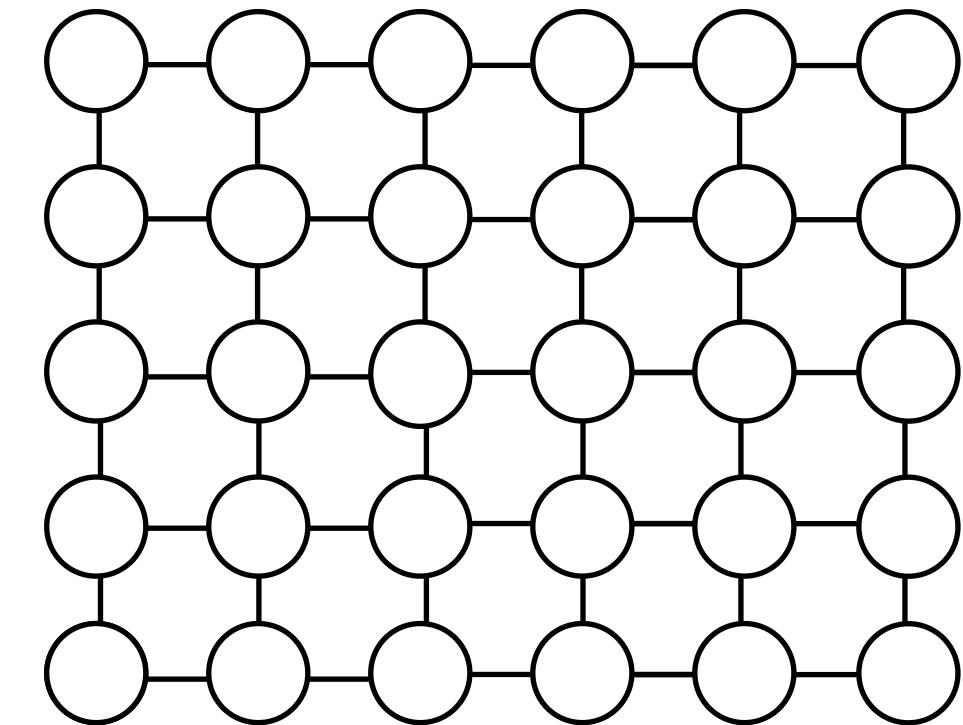
Rather than penalizing if output *image* looks fake, penalize if each overlapping *patch* in output looks fake

- Faster, fewer parameters
- More supervised observations
- Applies to arbitrarily large images

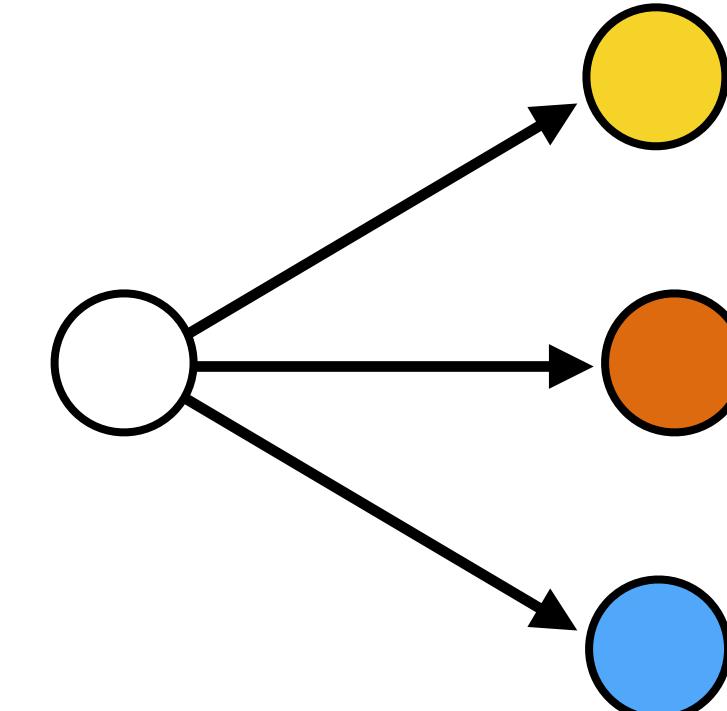
Properties of generative models

1. Model high-dimensional, structured output

→ Use a deep net, D , to model output!



2. Model uncertainty; a whole distribution of possible outputs



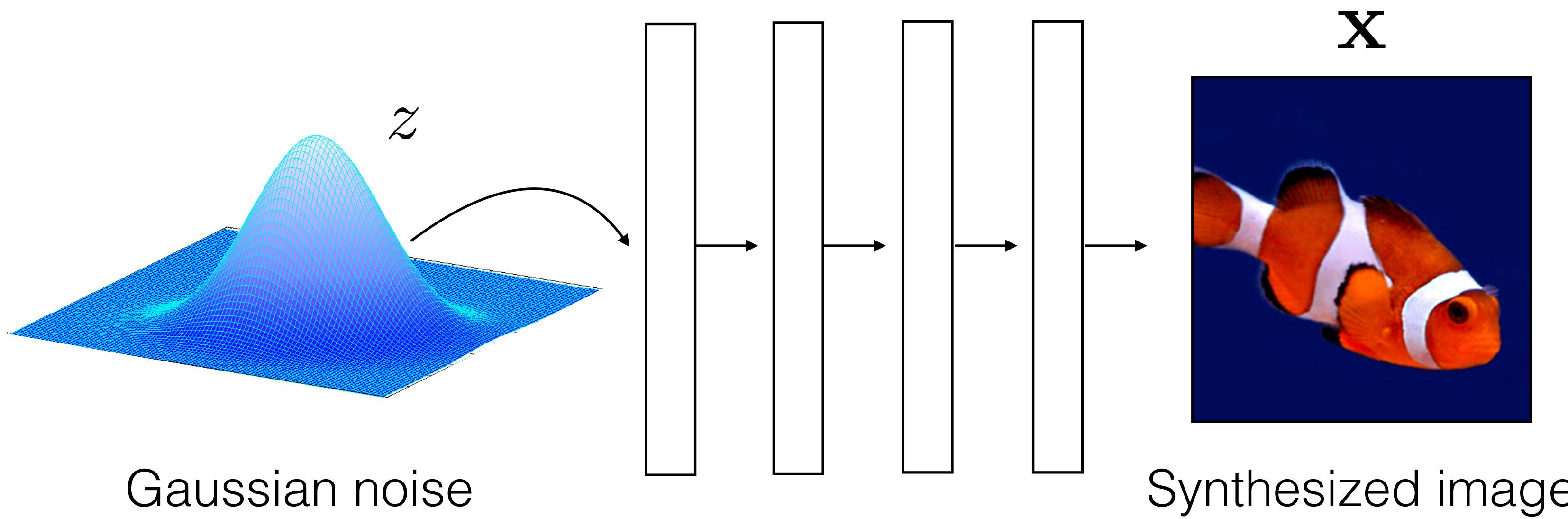
Three perspectives on GANs

1. Structured loss
2. Generative model
3. Domain-level supervision / mapping

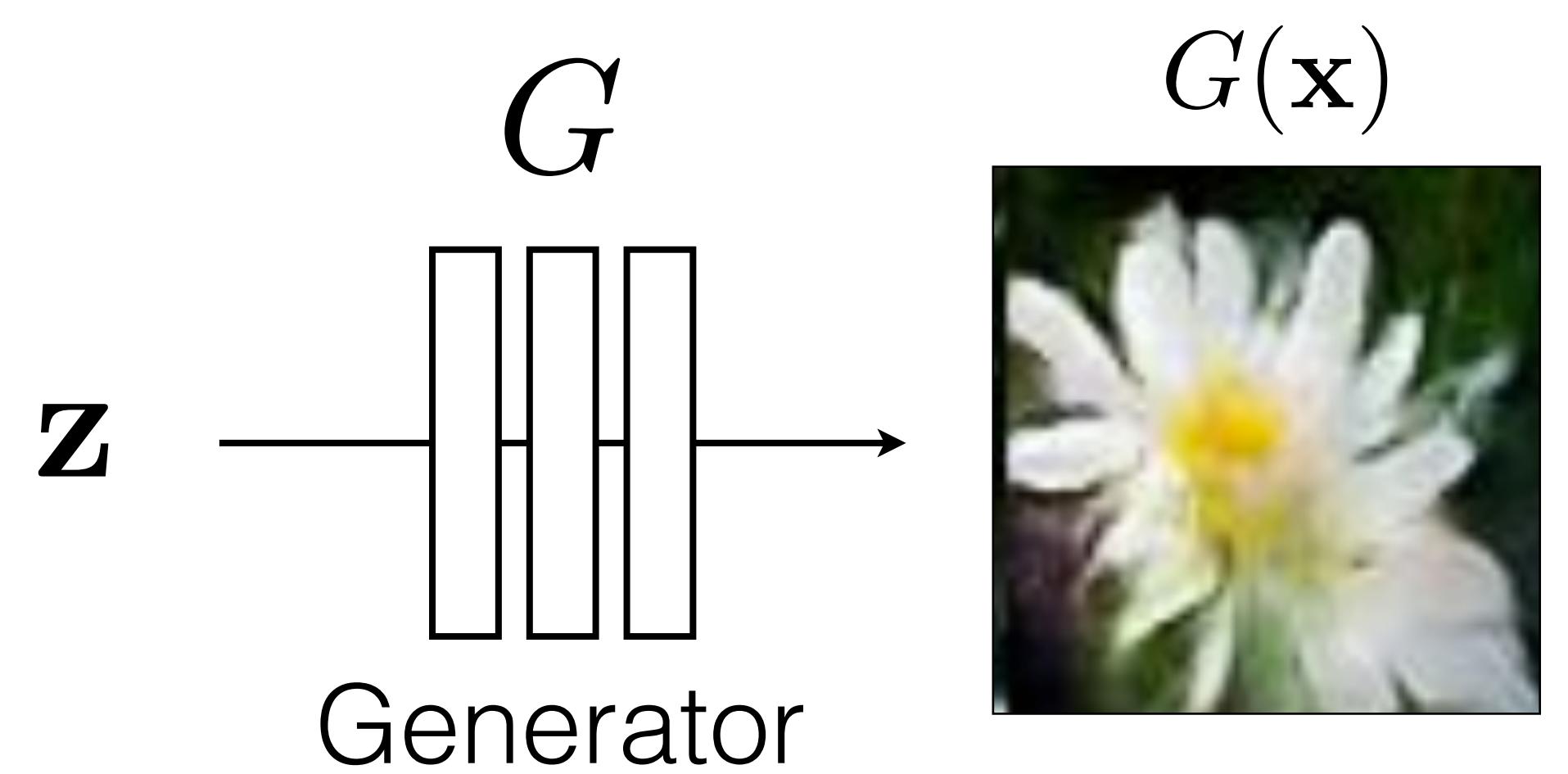
Three perspectives on GANs

1. Structured loss
- 2. Generative model**
3. Domain-level supervision / mapping

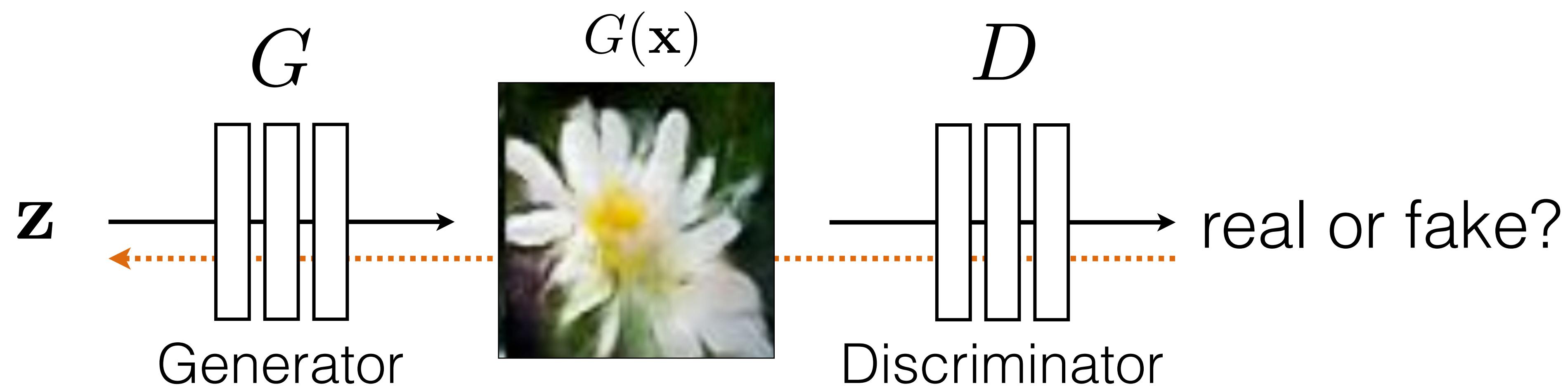
Can we generate images from scratch?



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



[Goodfellow et al., 2014]

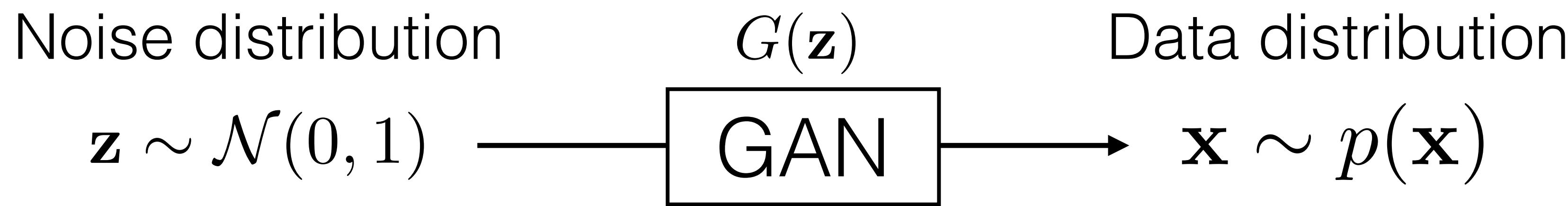


G tries to synthesize fake images that fool **D**

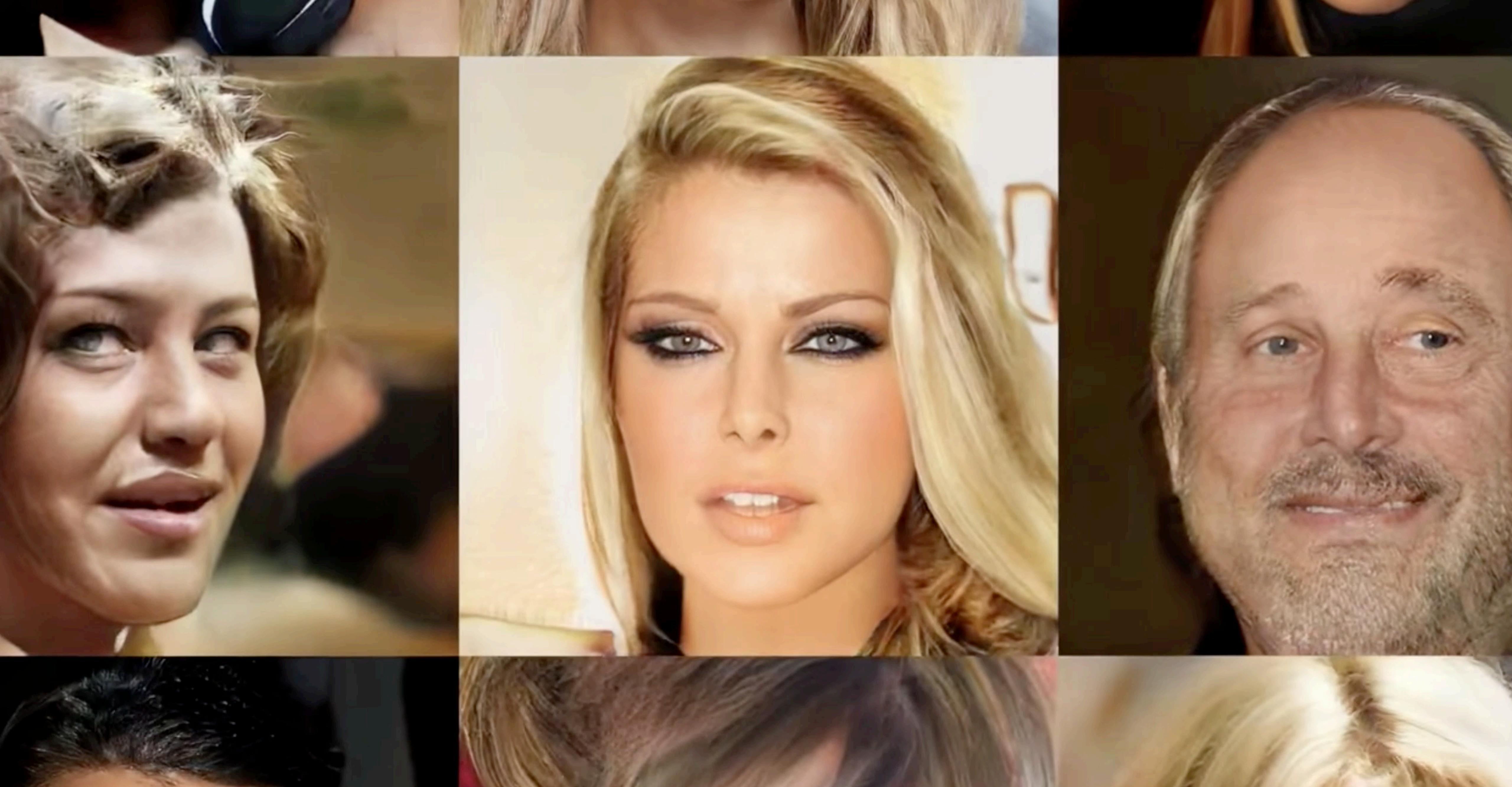
D tries to identify the fakes

GANs are implicit generative models

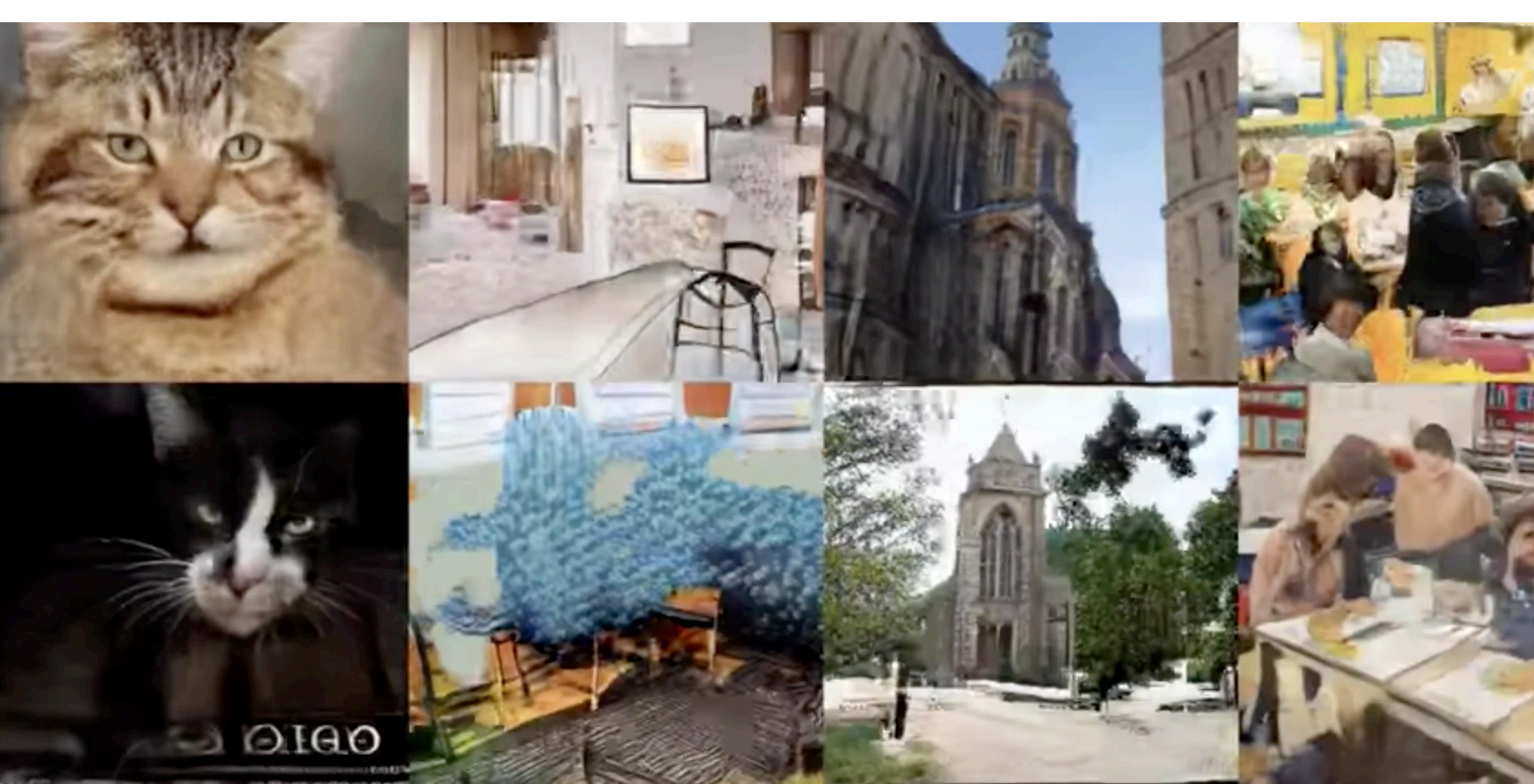
$p(\mathbf{x}) \longleftarrow$ “generative model” of the data \mathbf{x}



$G(\mathbf{z}) \sim p(\mathbf{x}) \longleftarrow$ Samples from a perfectly optimized GAN are samples from the data distribution



Progressive GAN [Karras et al., 2018]



Progressive GAN [Karras et al., 2018]

Proposition 1. *For G fixed, the optimal discriminator D is*

$$D_G^*(\mathbf{x}) = \frac{p_{\text{data}}(\mathbf{x})}{p_{\text{data}}(\mathbf{x}) + p_g(\mathbf{x})}$$

Proof

$$\begin{aligned} V(G, D) &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) d\mathbf{x} + \int_z p_{\mathbf{z}}(\mathbf{z}) \log(1 - D(g(\mathbf{z}))) dz \\ &= \int_{\mathbf{x}} p_{\text{data}}(\mathbf{x}) \log(D(\mathbf{x})) + p_g(\mathbf{x}) \log(1 - D(\mathbf{x})) d\mathbf{x} \end{aligned}$$

For any $(a, b) \in \mathbb{R}^2 \setminus \{0, 0\}$, the function $y \rightarrow a \log(y) + b \log(1 - y)$ achieves its maximum in $[0, 1]$ at $\frac{a}{a+b}$. The discriminator does not need to be defined outside of $\text{Supp}(p_{\text{data}}) \cup \text{Supp}(p_g)$, concluding the proof. \square

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

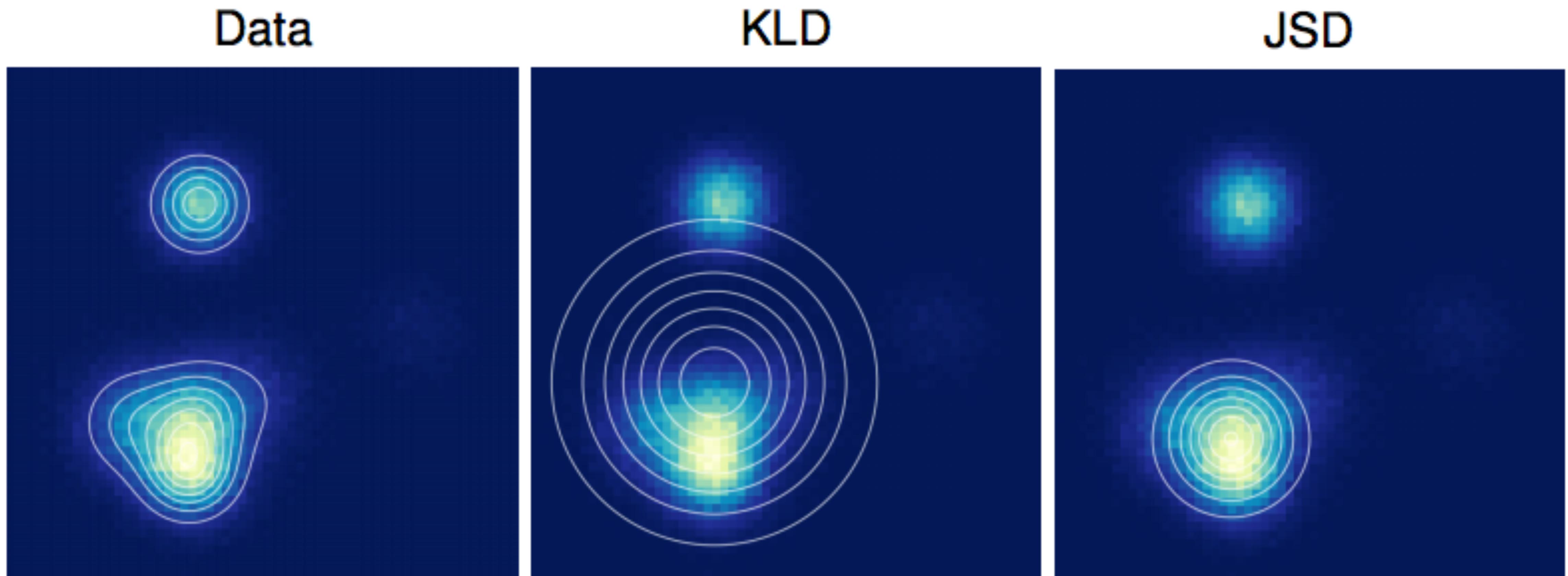
Proof

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

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

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

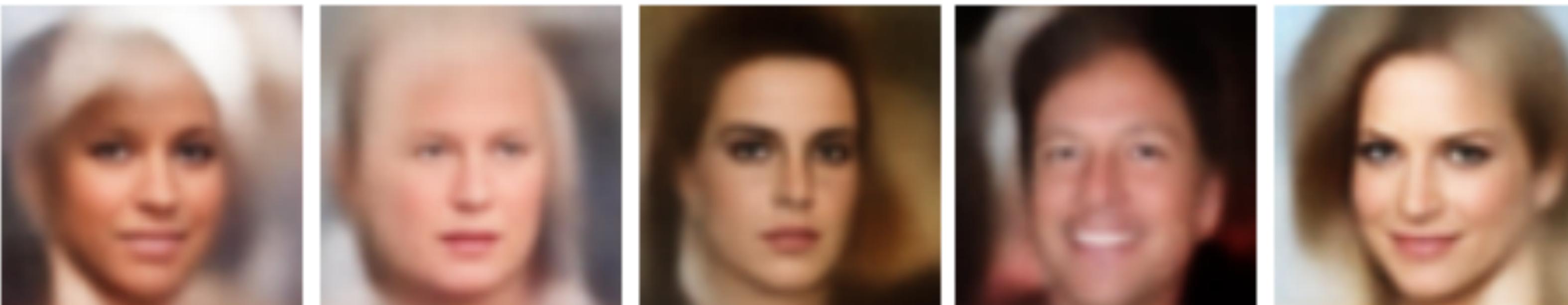
Behavior under model misspecification



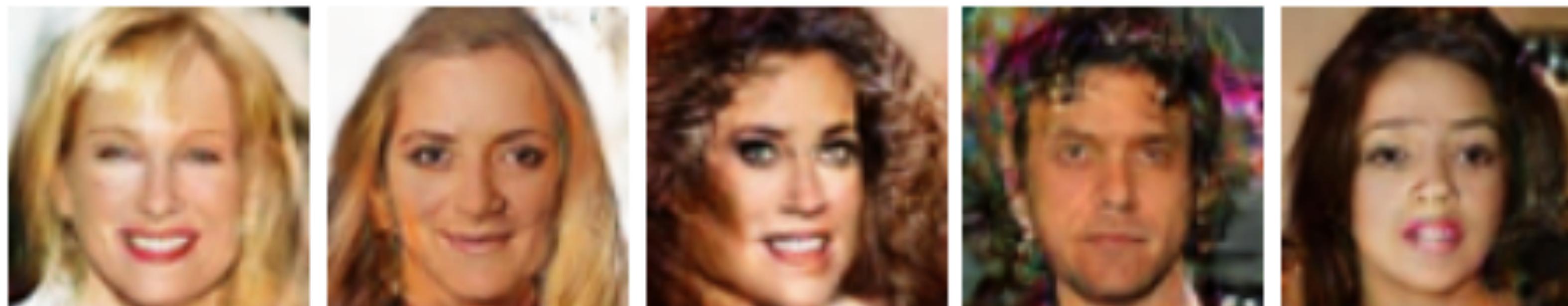
[Theis et al. 2016]

Mode covering versus mode seeking

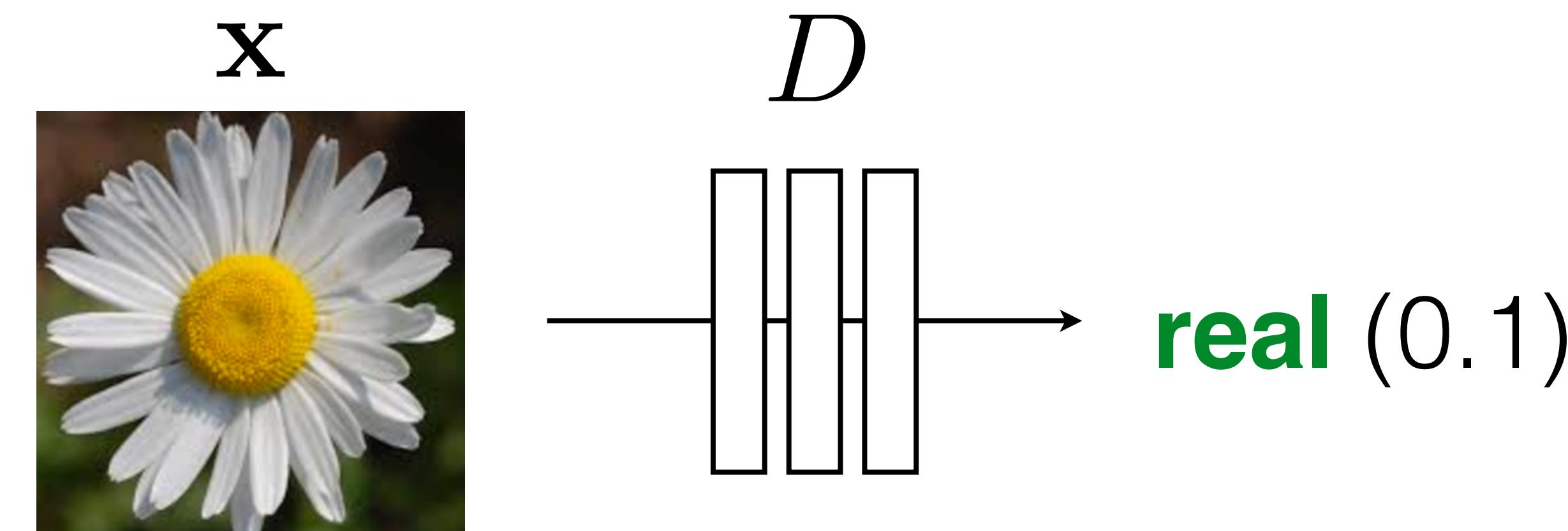
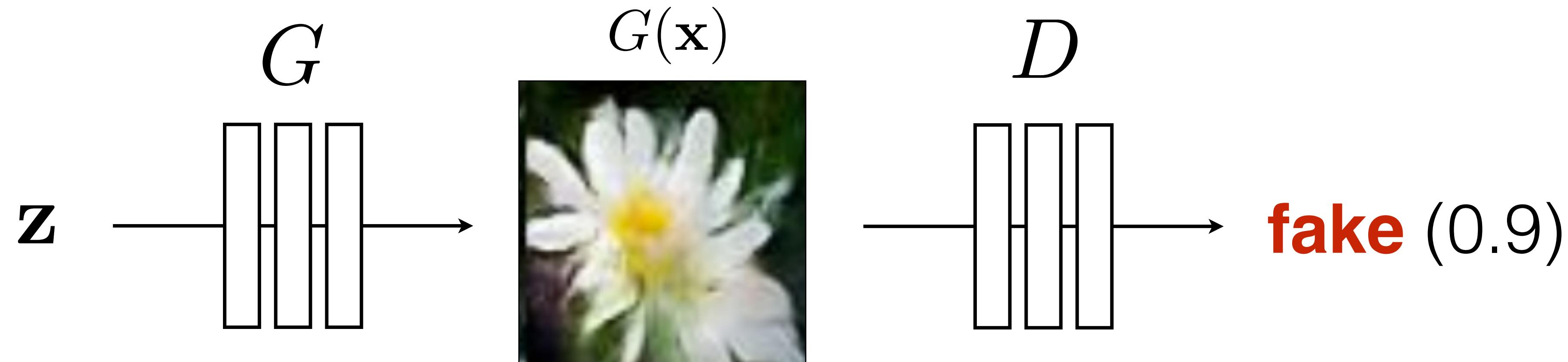
VAE



GAN

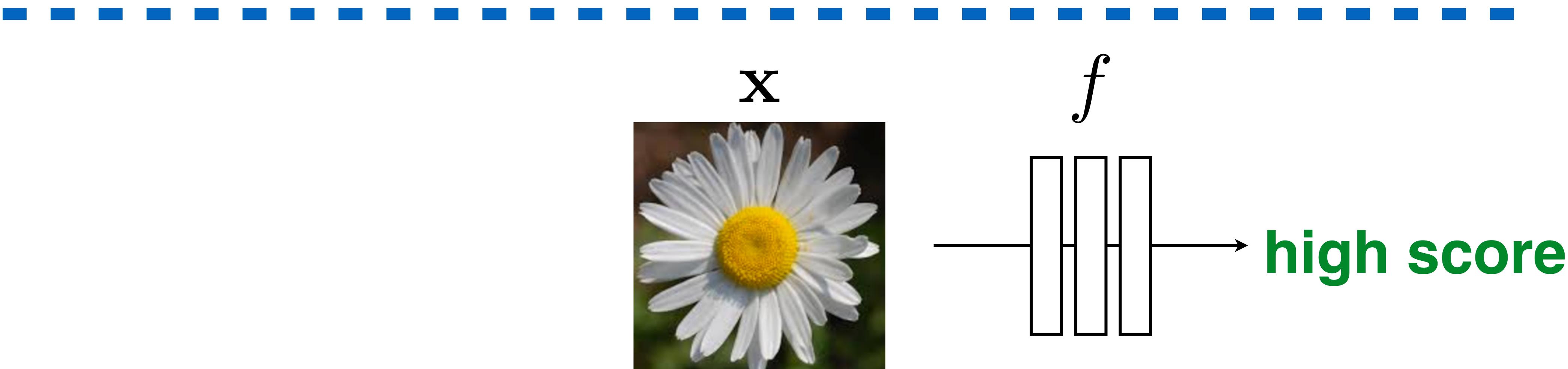
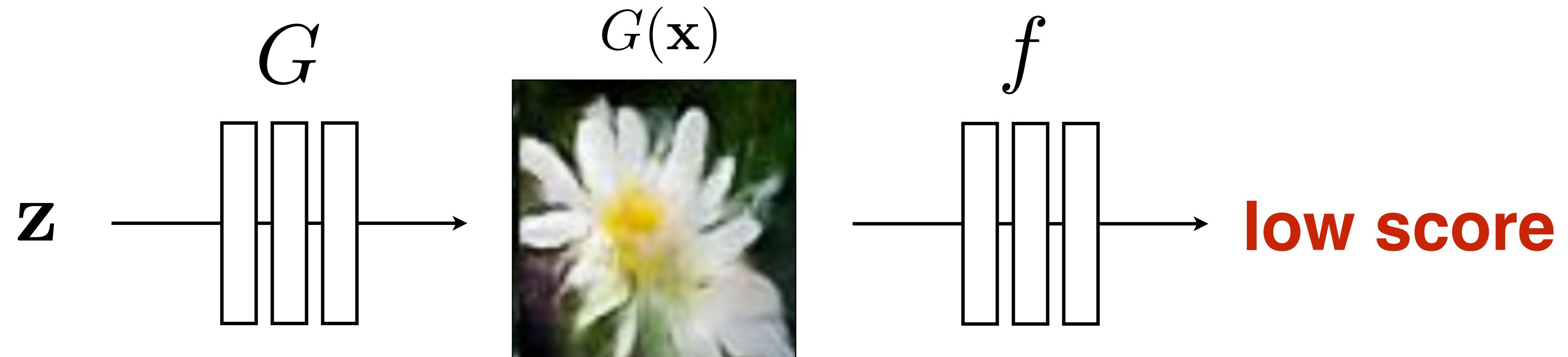


[Larsen et al. 2016]



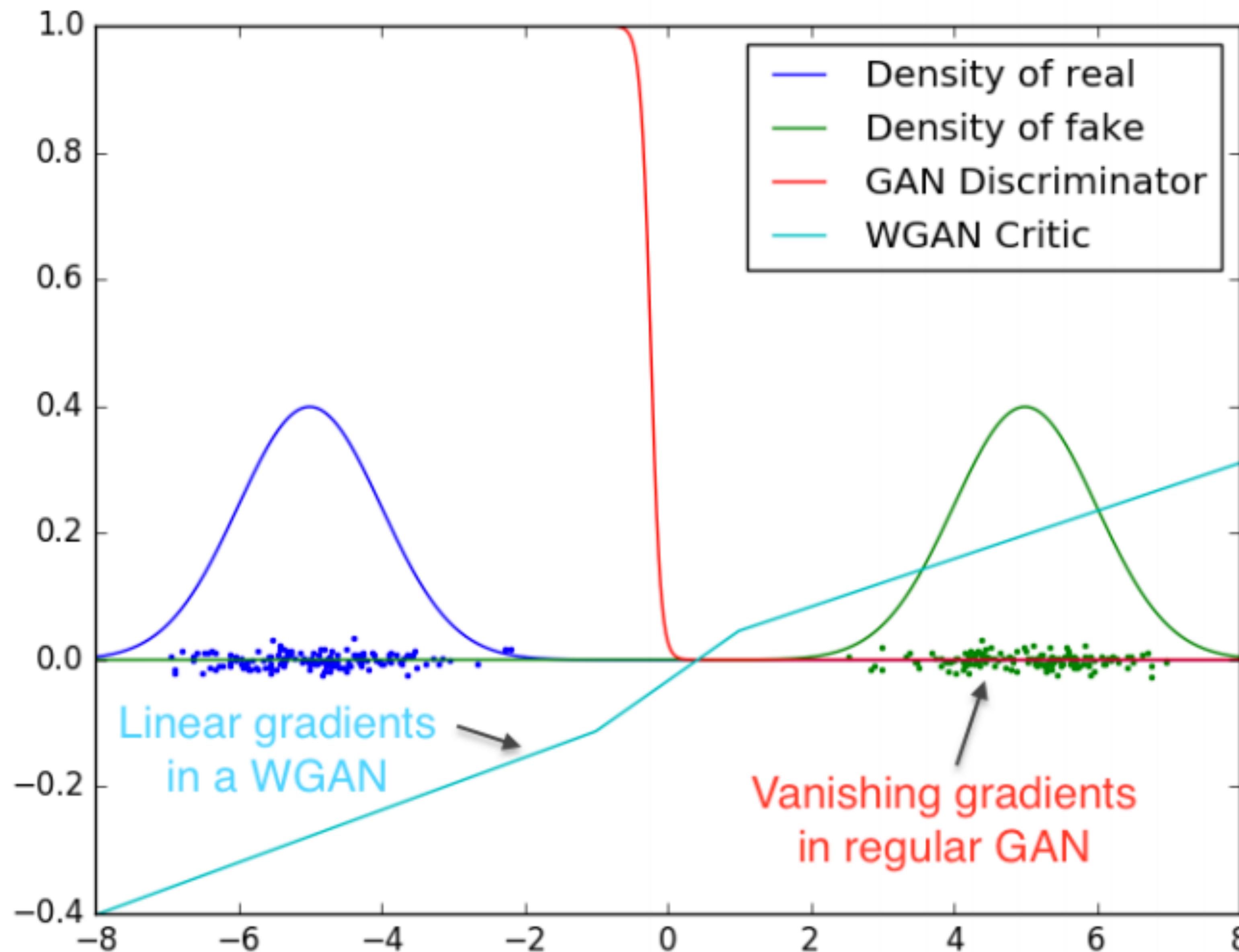
$$\arg \max_D \mathbb{E}_{\mathbf{z}, \mathbf{x}} [\boxed{\log D(G(\mathbf{z}))} + \boxed{\log (1 - D(\mathbf{x}))}]$$

[Goodfellow et al., 2014]

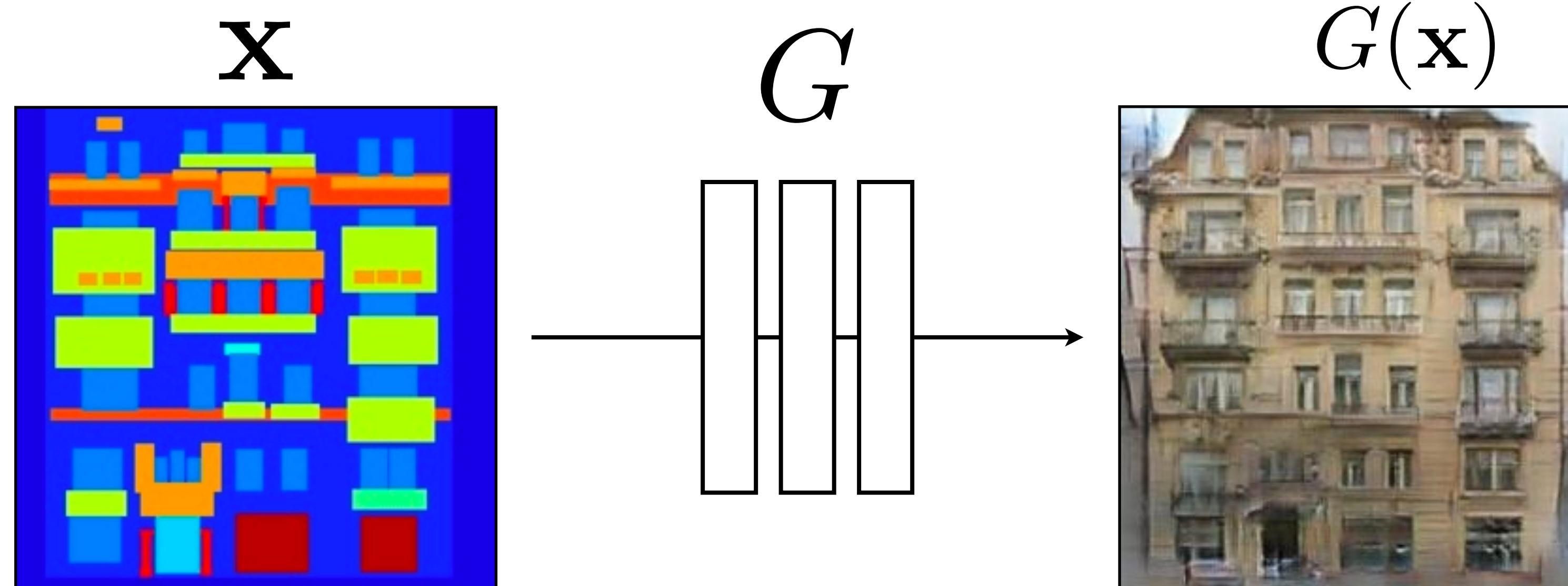


$$\arg \max_f \mathbb{E}_{\mathbf{z}, \mathbf{x}} [-f(G(\mathbf{z})) + f(\mathbf{x})]$$

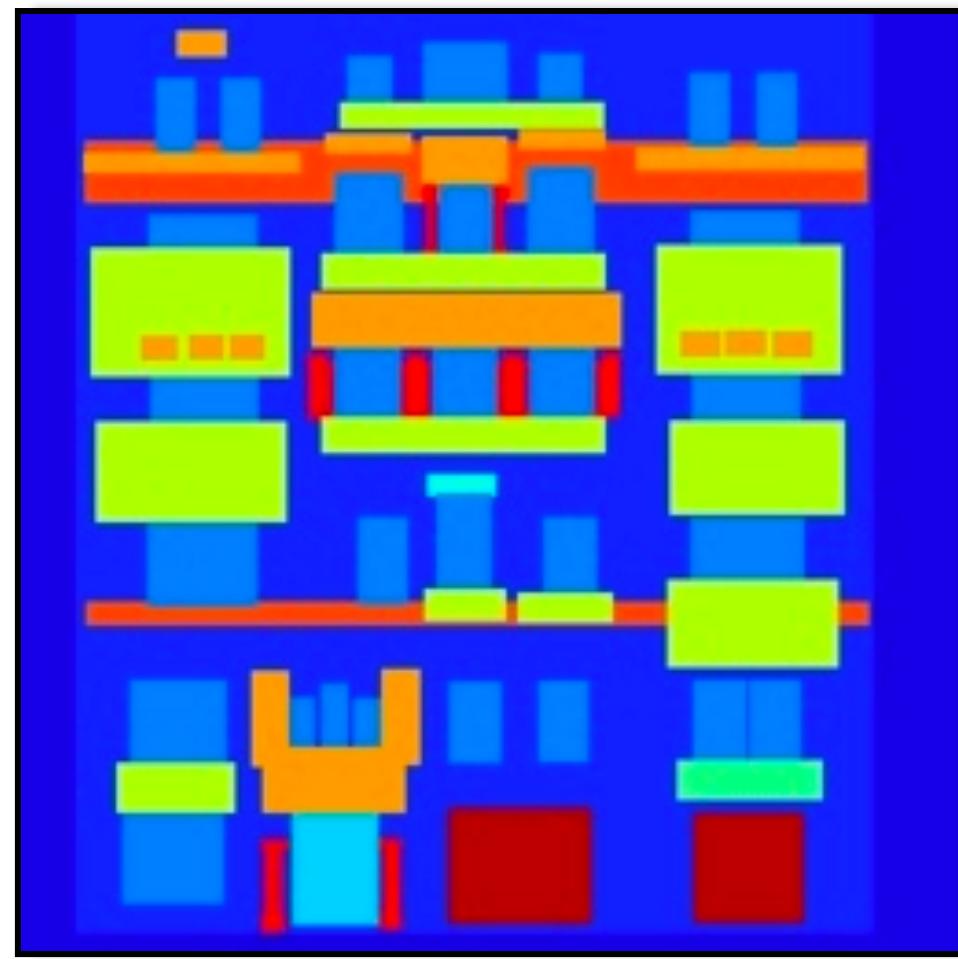
EBGAN, WGAN, LSGAN, etc



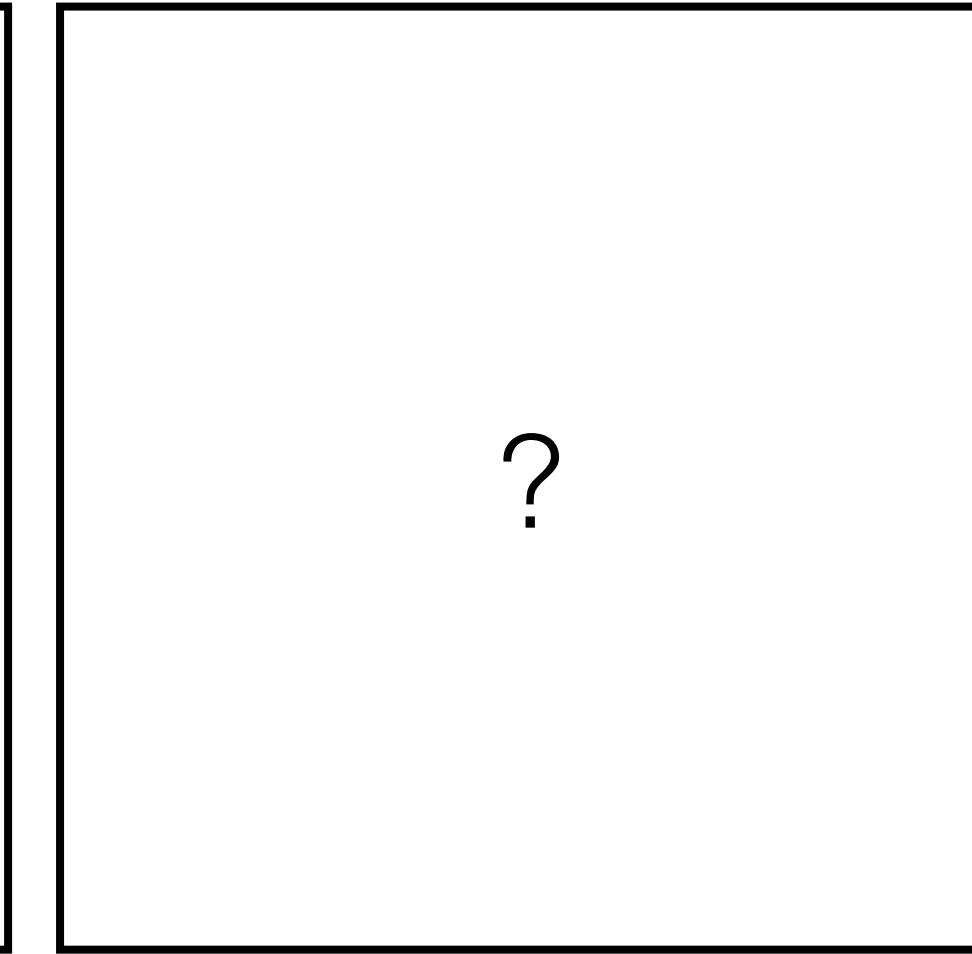
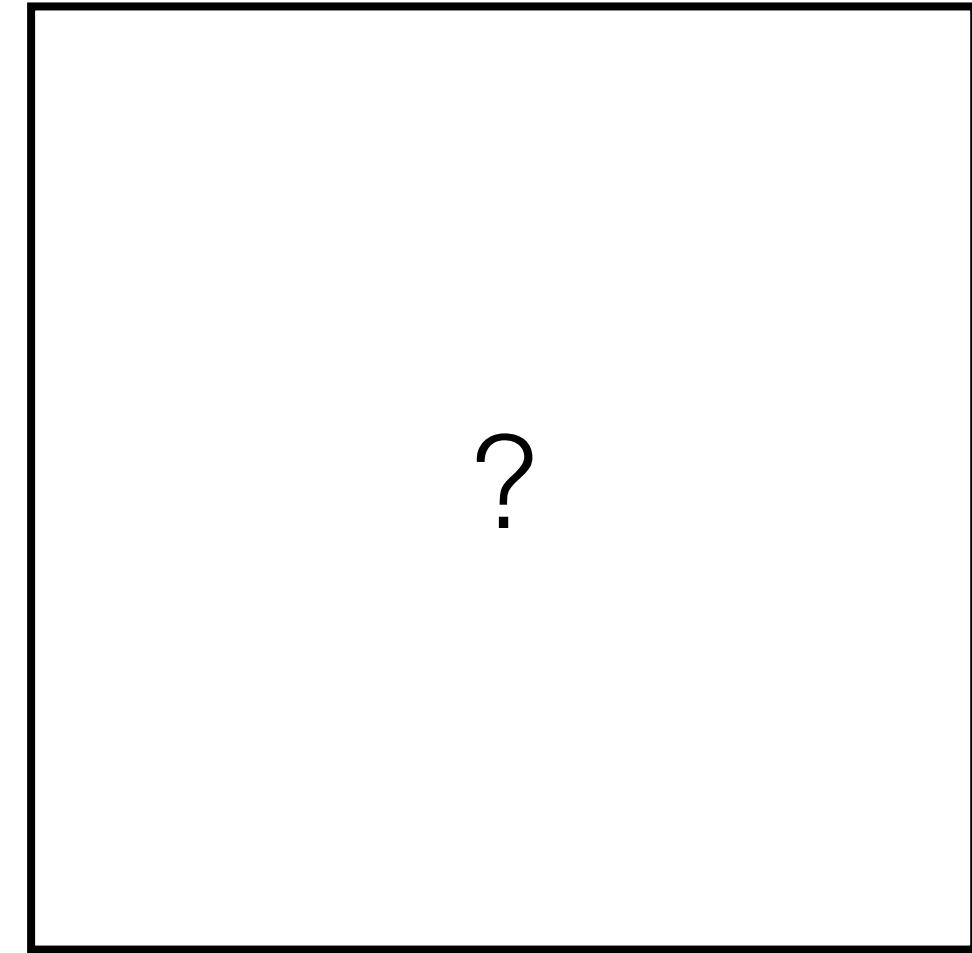
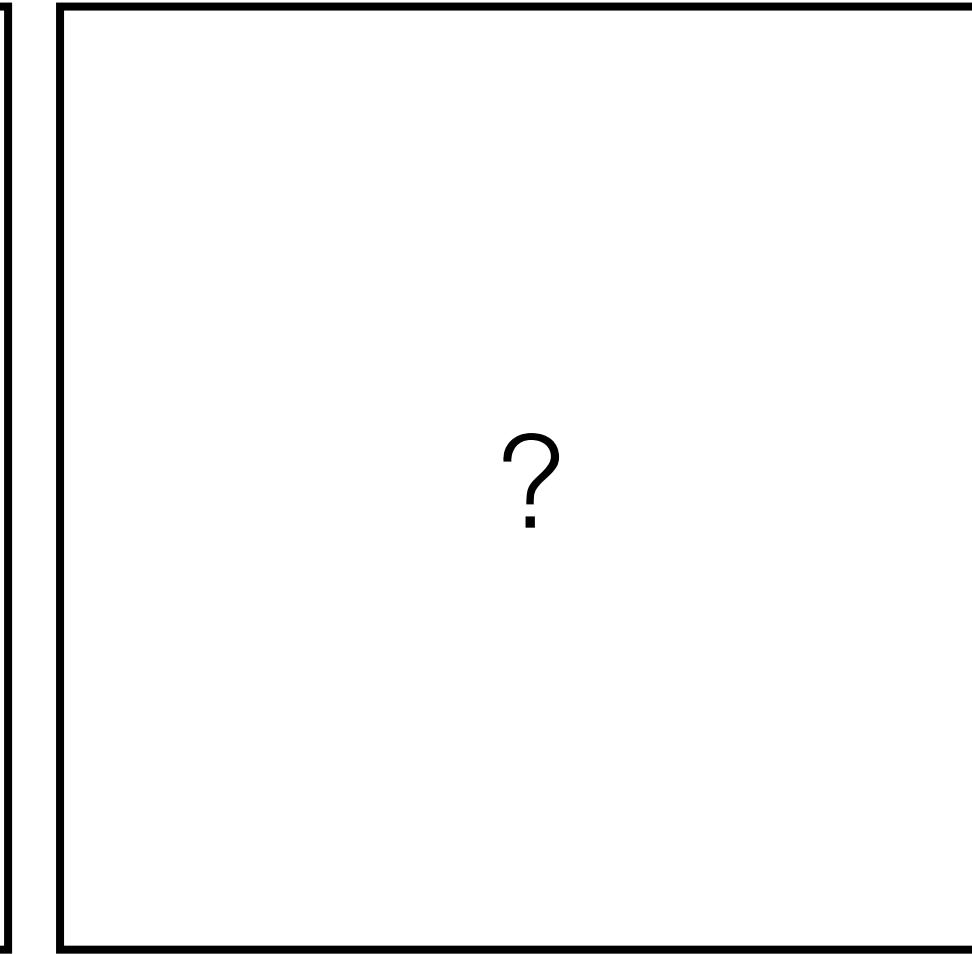
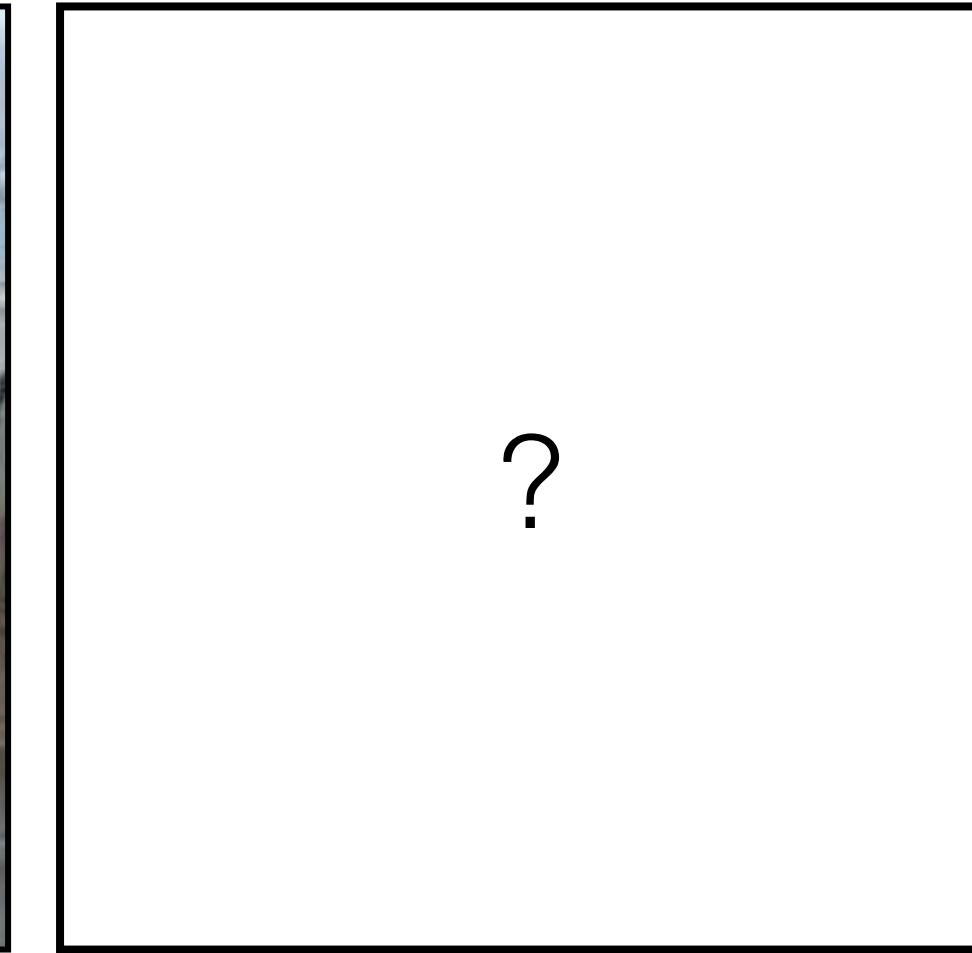
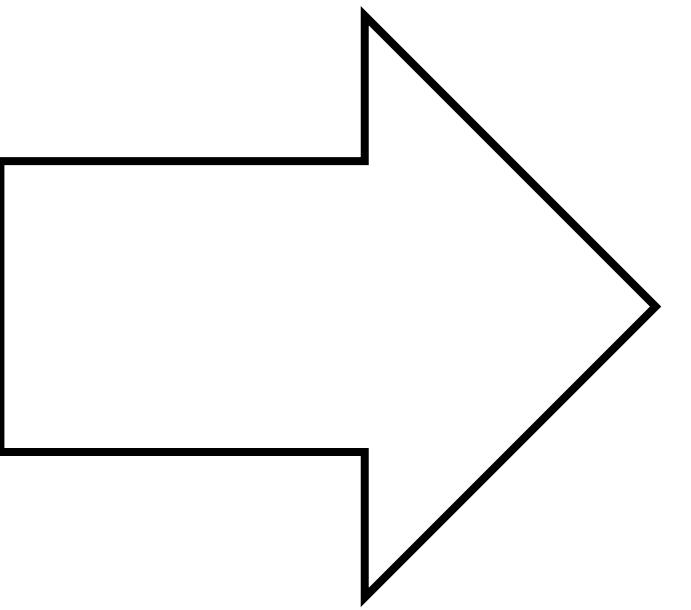
Modeling multiple possible outputs



Modeling multiple possible outputs

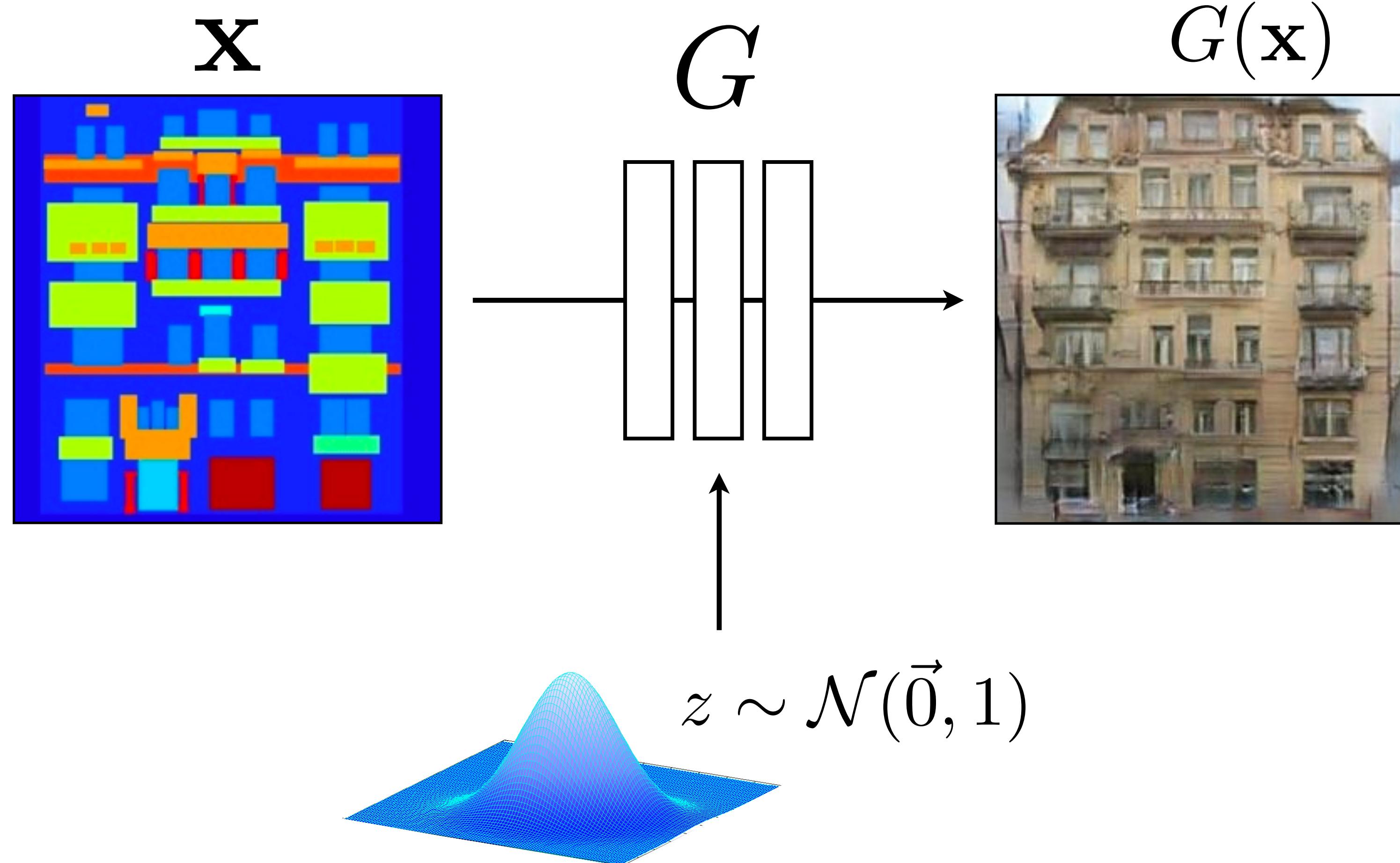


Input



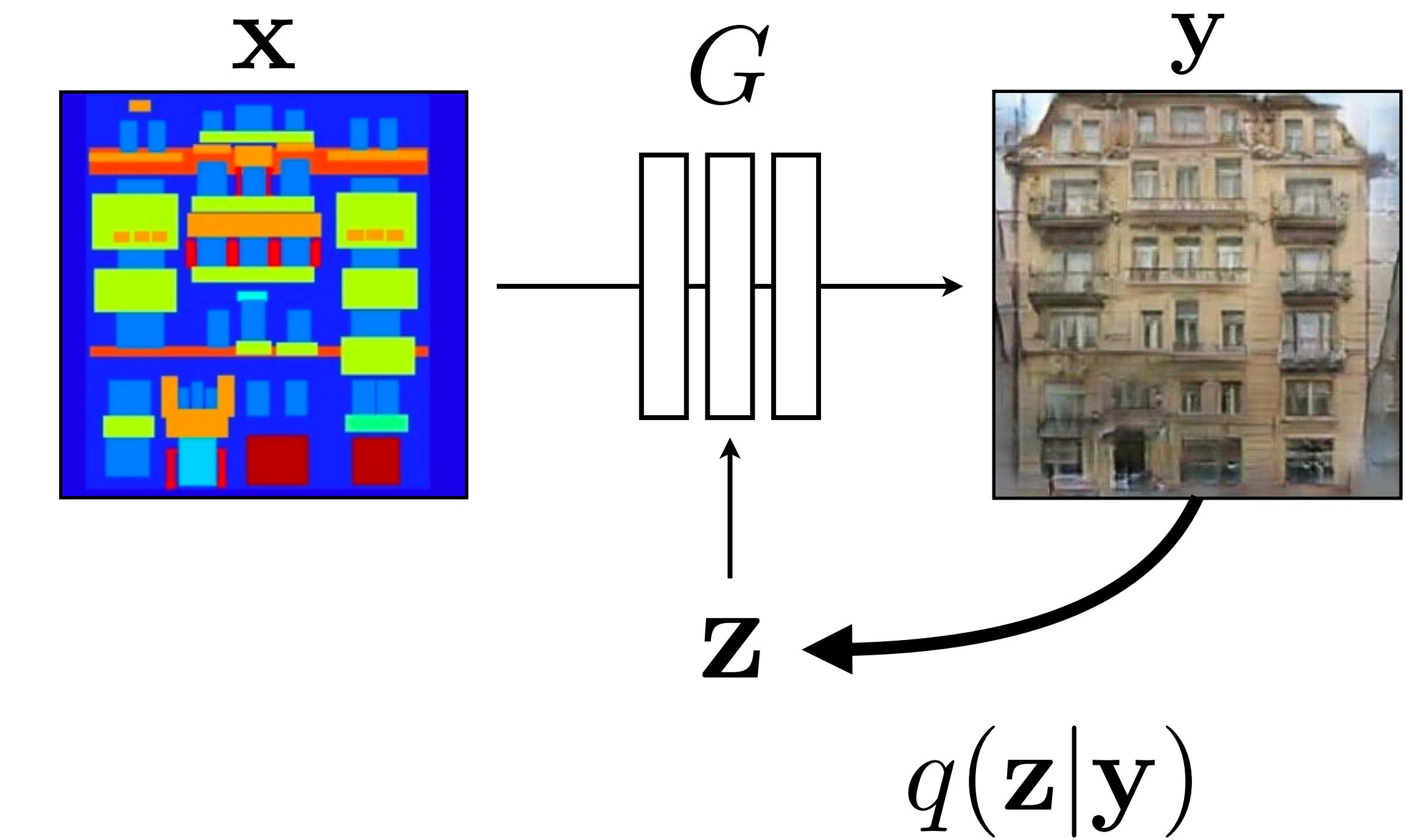
Possible outputs

Modeling multiple possible outputs

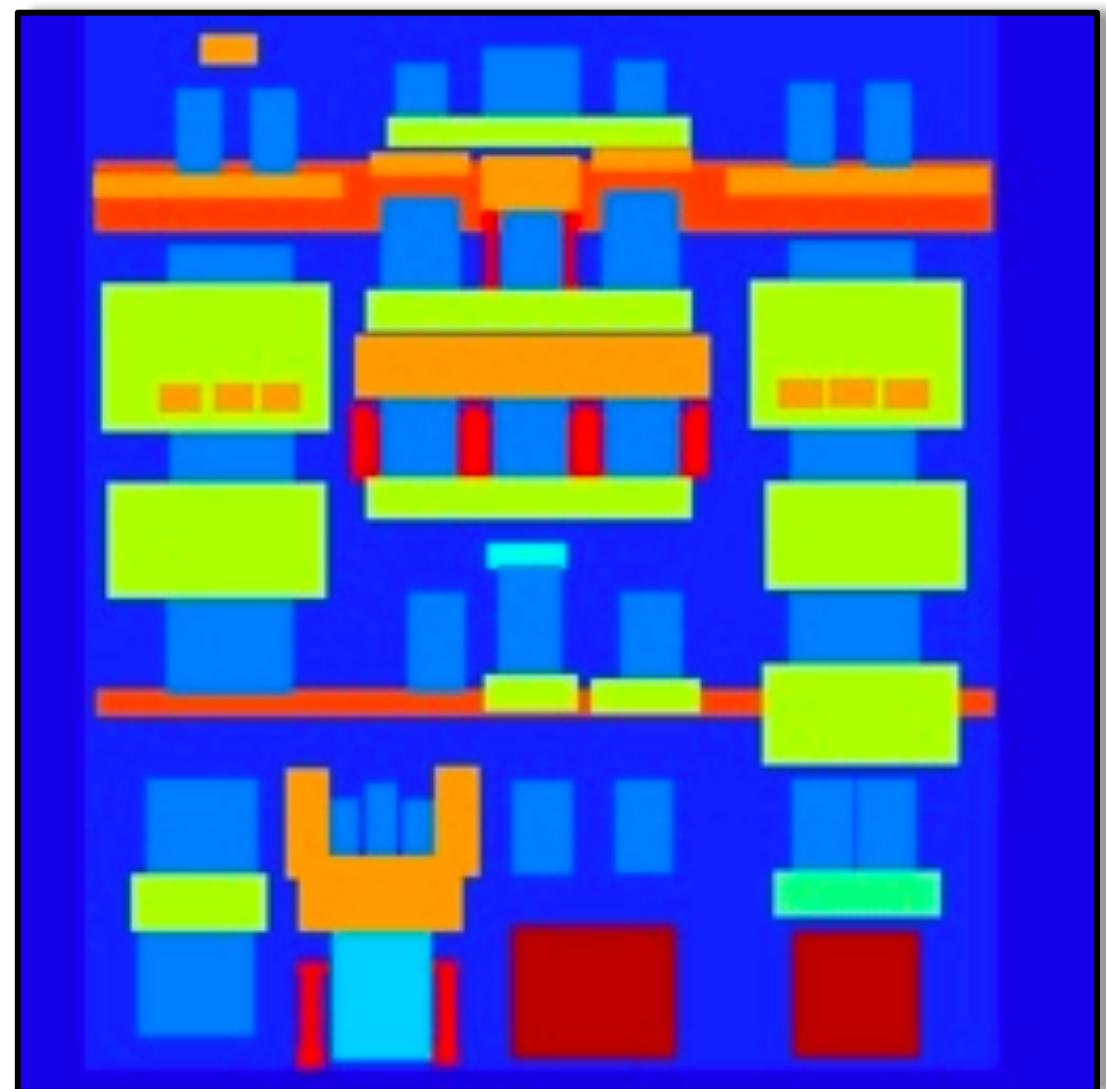


InfoGAN [Chen et al. 2016]

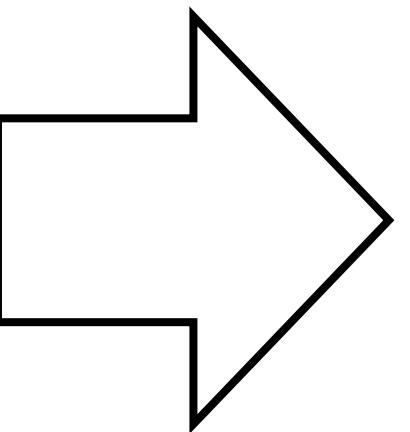
BiCycleGAN [Zhu et al., NIPS 2017]



Encourages z to relay information about the target.



Labels



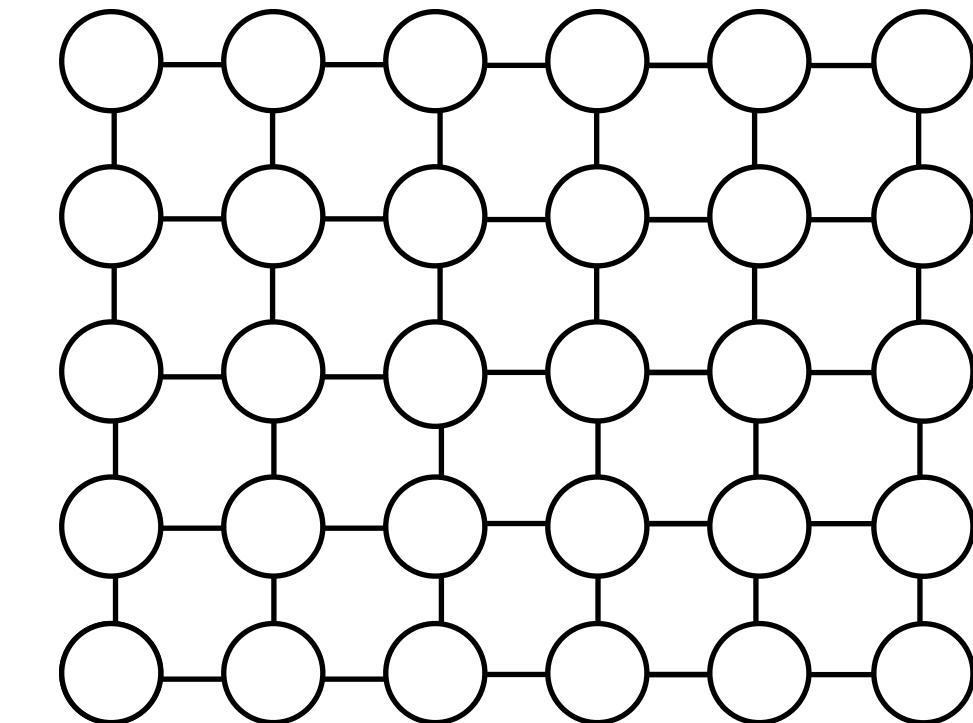
Randomly generated facades

[BiCycleGAN, Zhu et al., NIPS 2017]

Properties of generative models

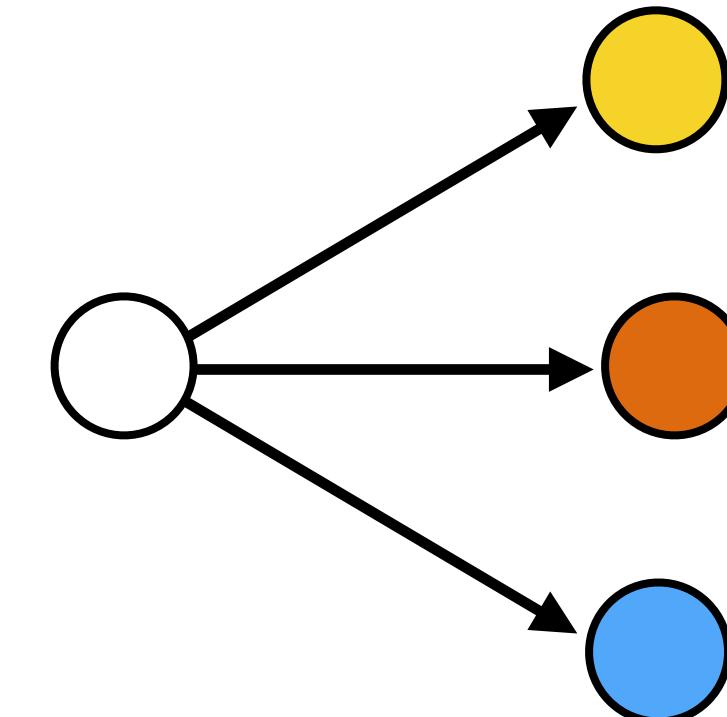
1. Model high-dimensional, structured output

→ Use a deep net, D, to model output!



2. Model uncertainty; a whole distribution of possible outputs

→ Generator is stochastic, learns to match data distribution



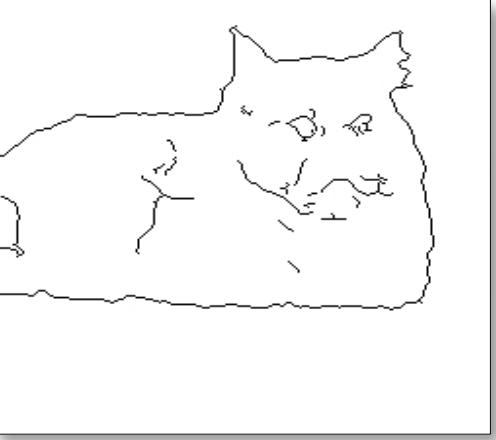
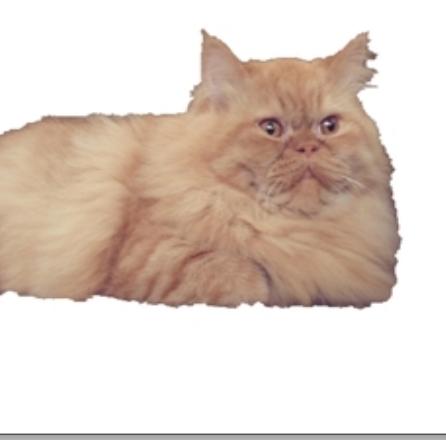
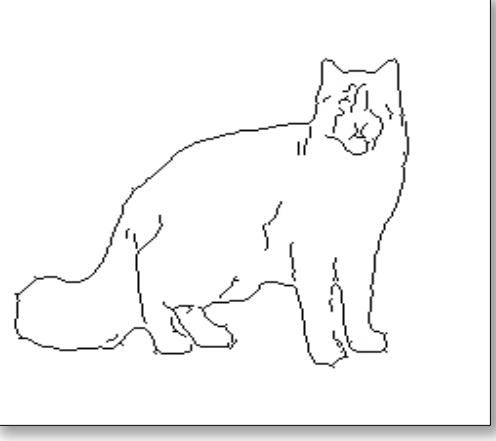
Three perspectives on GANs

1. Structured loss
- 2. Generative model**
3. Domain-level supervision / mapping

Three perspectives on GANs

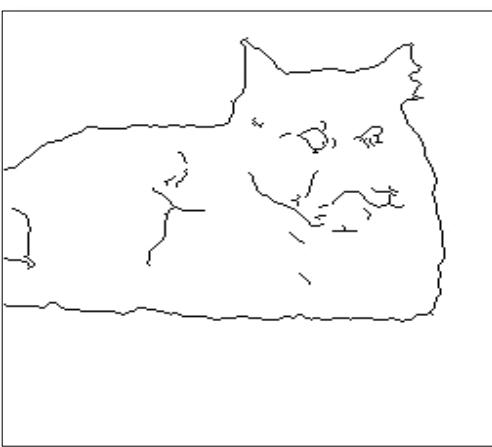
1. Structured loss
2. Generative model
- 3. Domain-level supervision / mapping**

Paired data

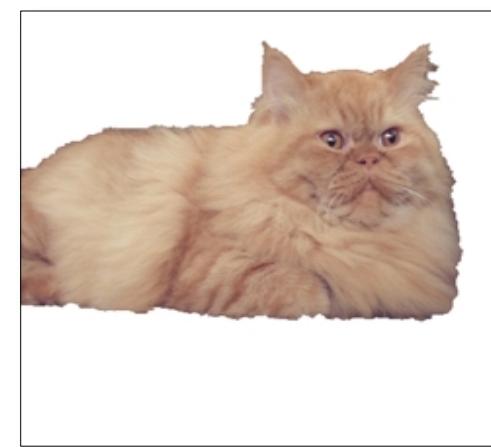
x_i	y_i
{ 	,  }
{ 	,  }
{ 	,  }
⋮	

Paired data

x_i



y_i



{}

2

{}

1

3

10

Unpaired data

X



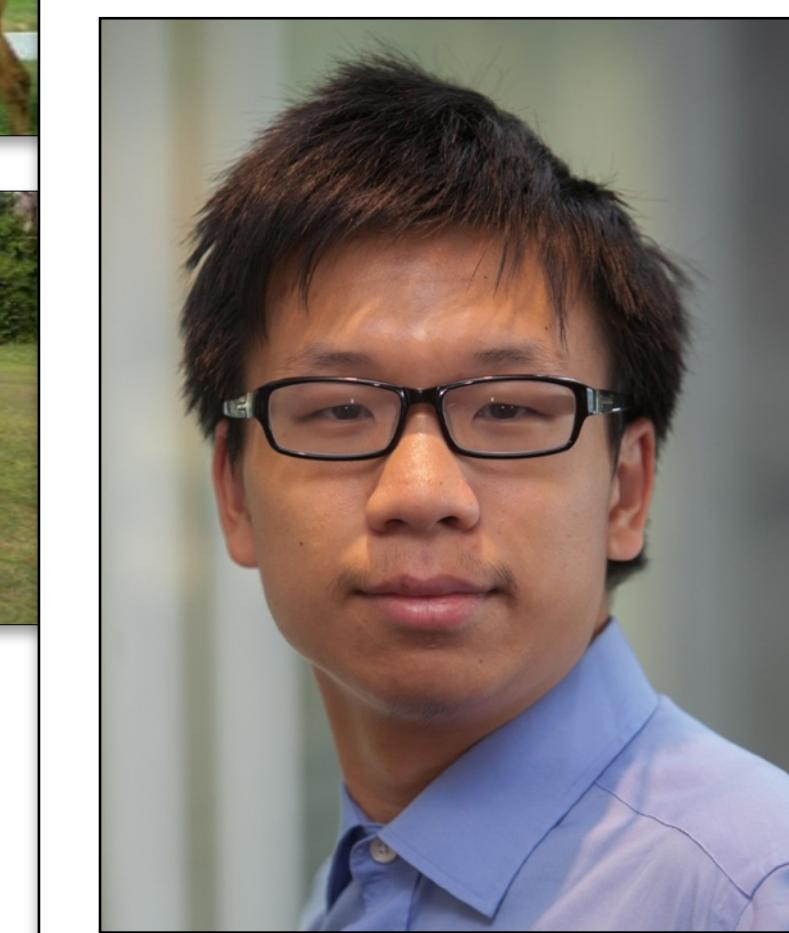
Y



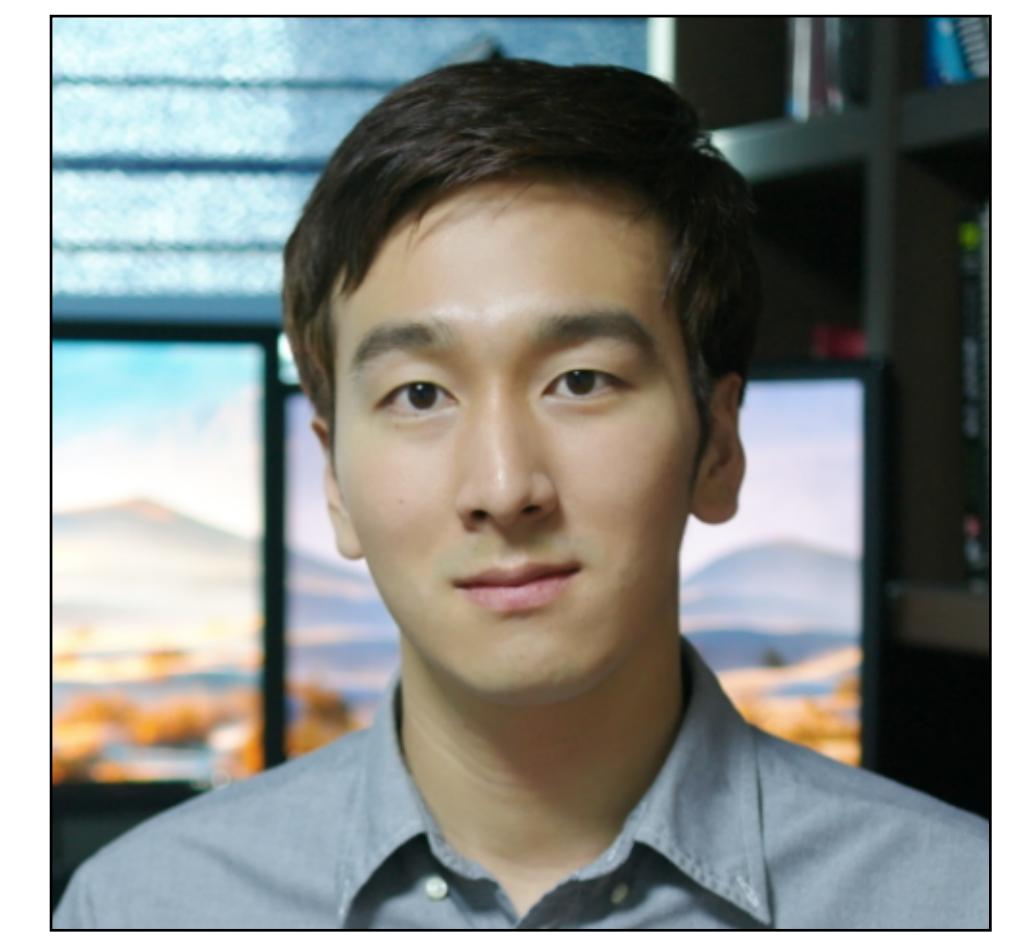
卷之三



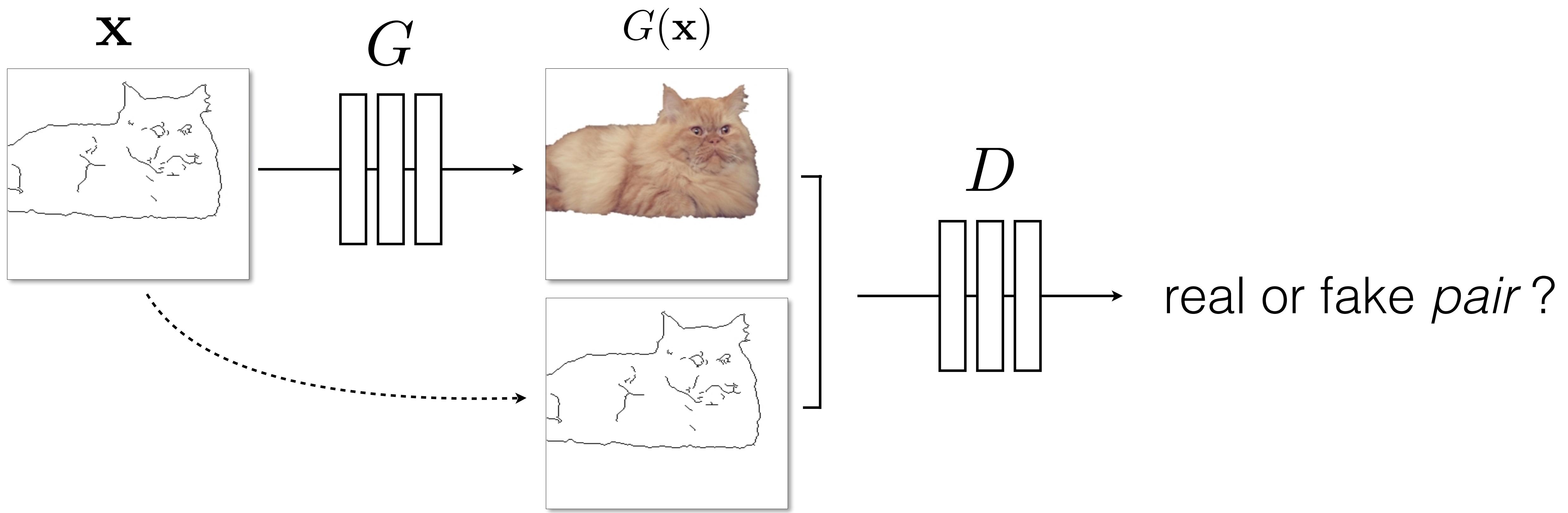
A vertical column of three solid black circles, representing a list item or bullet point.



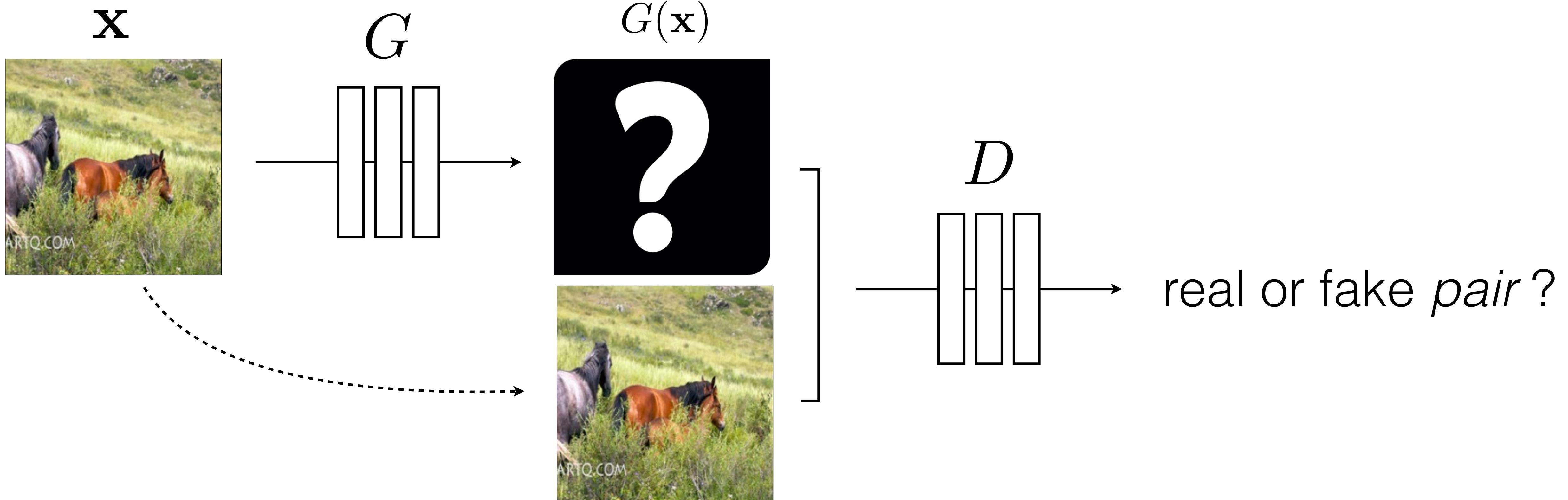
Jun-Yan Zhu



Taesung Park

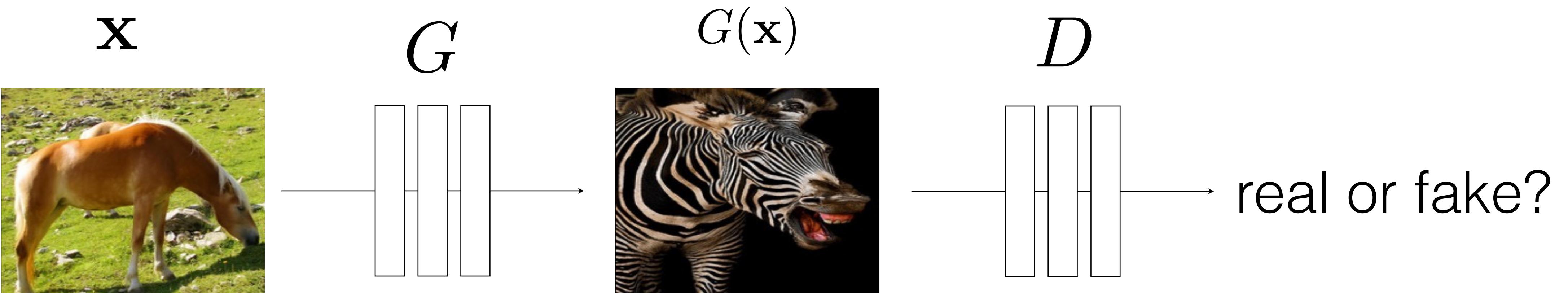


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(\mathbf{x}, G(\mathbf{x})) + \log(1 - D(\mathbf{x}, \mathbf{y}))]$$



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

No input-output pairs!

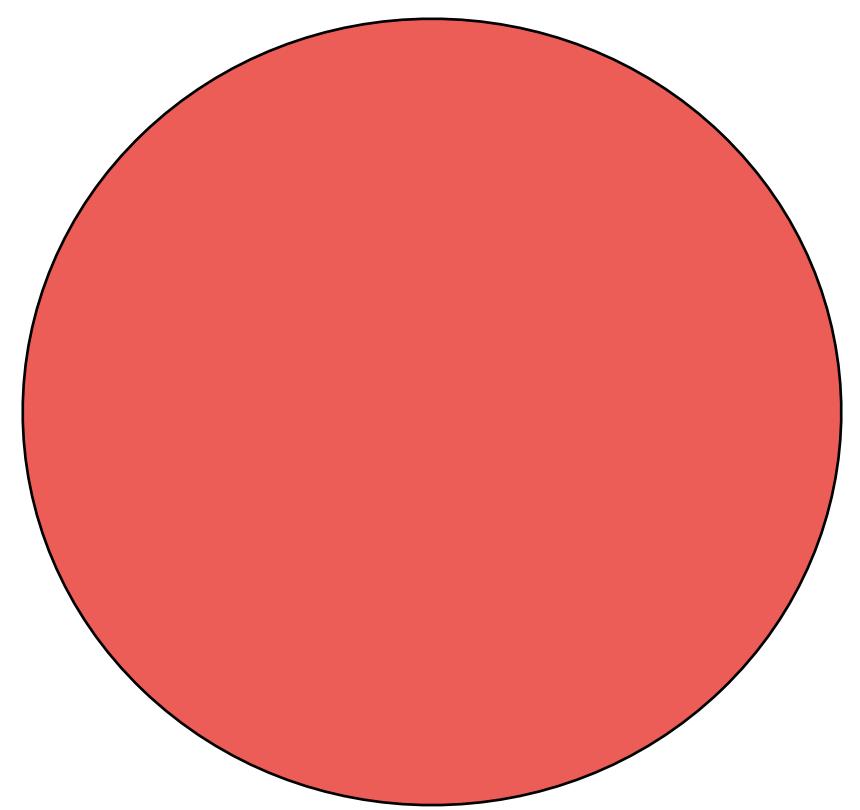


$$\arg \min_G \max_D \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\log D(G(\mathbf{x})) + \log(1 - D(\mathbf{y}))]$$

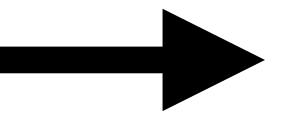
Usually loss functions check if output matches a target *instance*

GAN loss checks if output is part of an admissible set

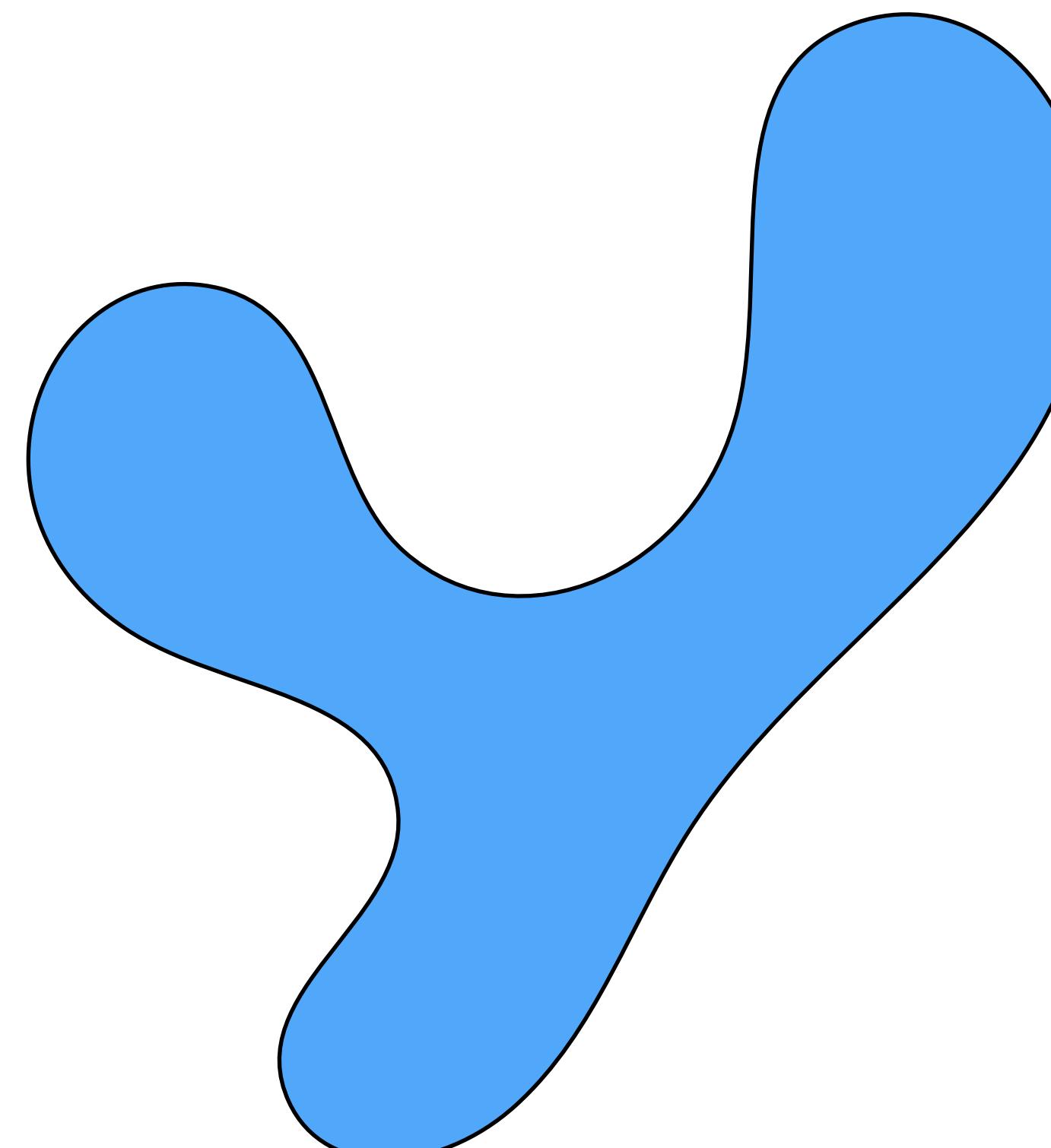
Gaussian



z

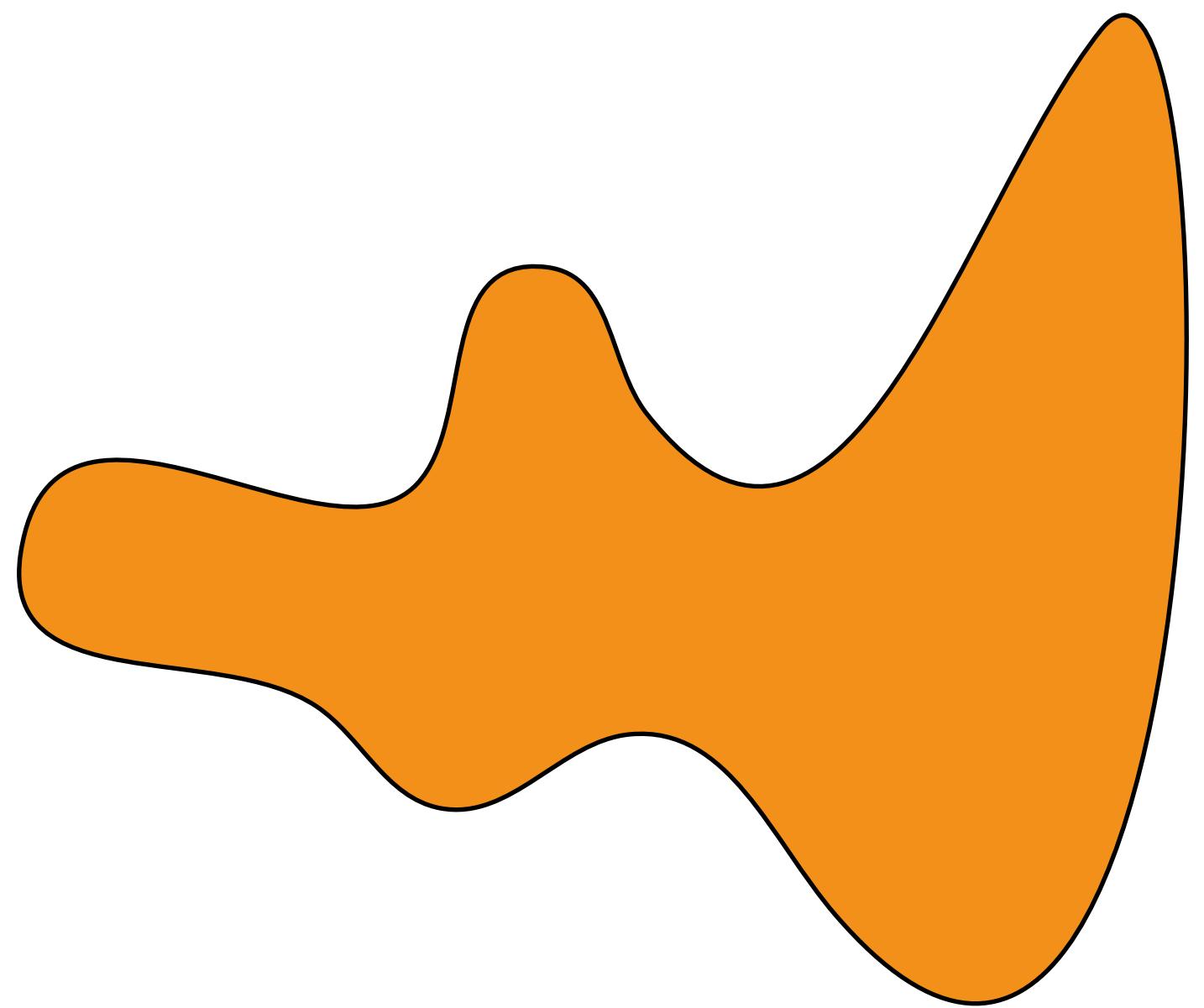


Target distribution



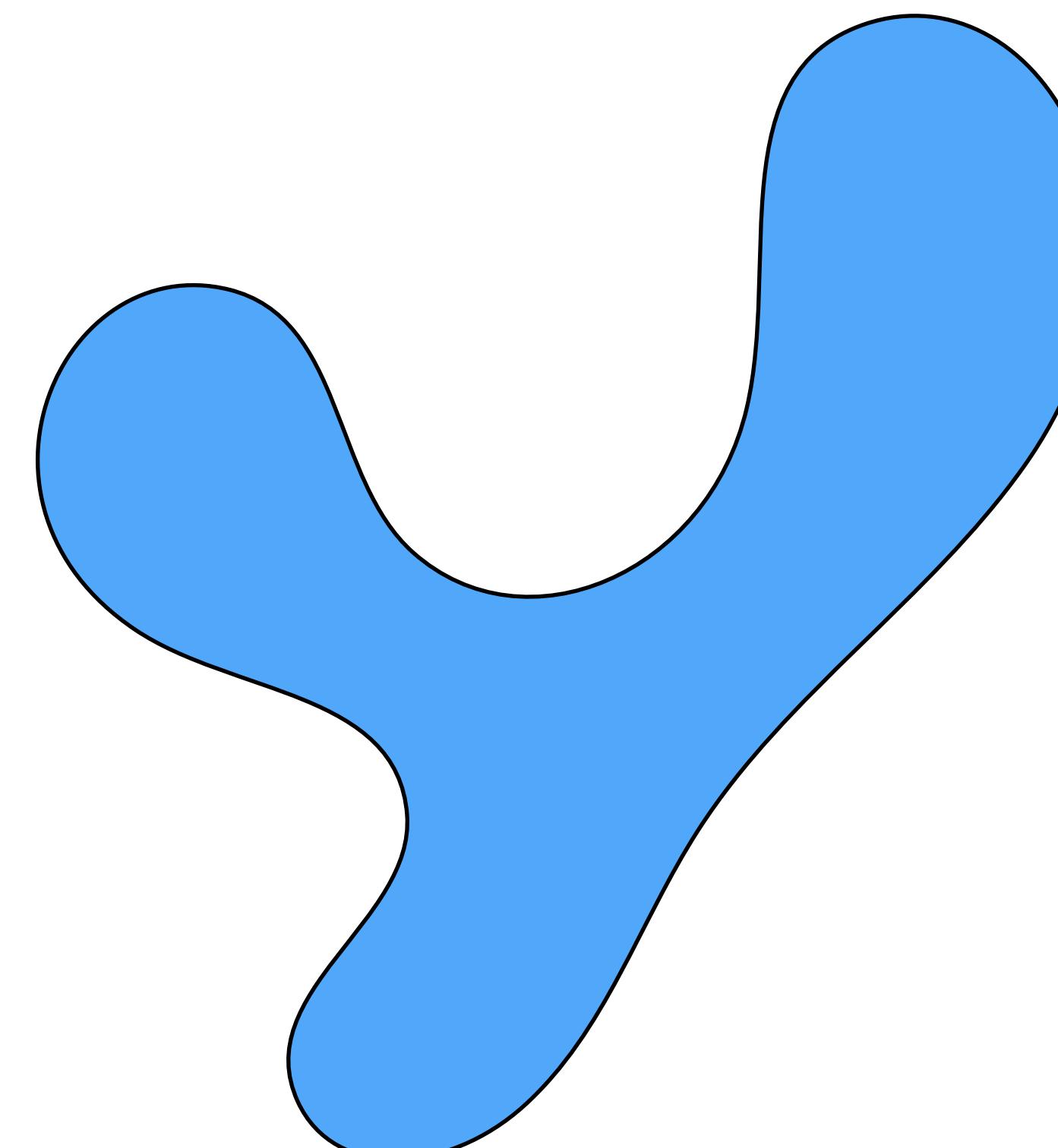
Y

Horses

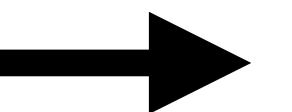


X

Zebras



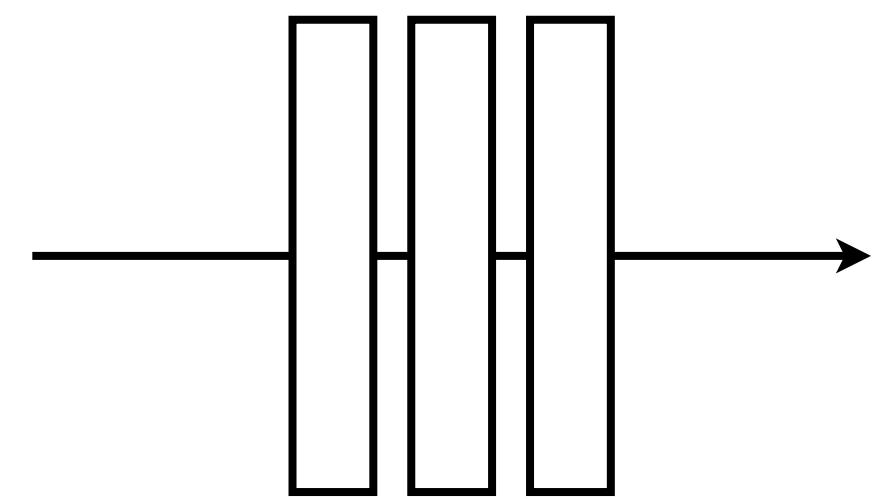
Y



x



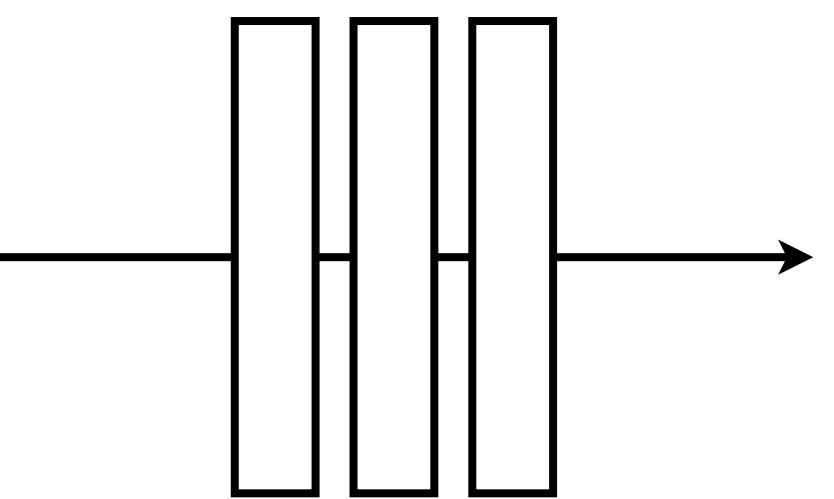
G



G(x)



D

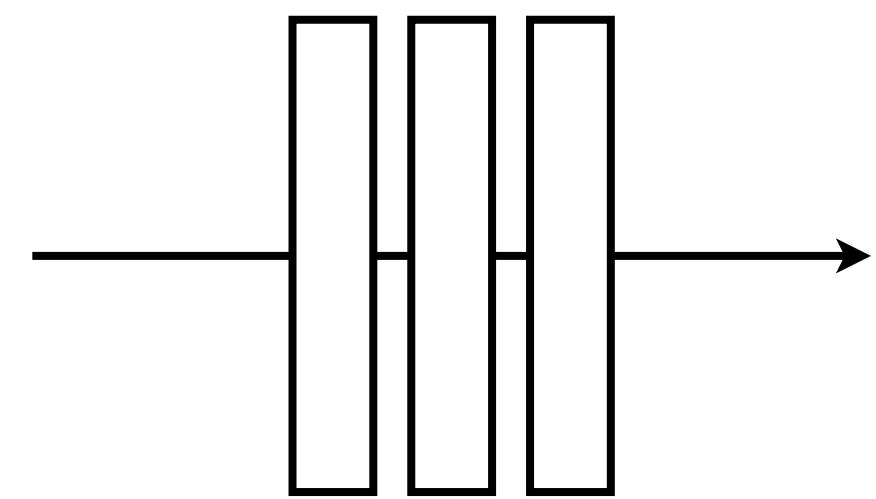


Real!

x



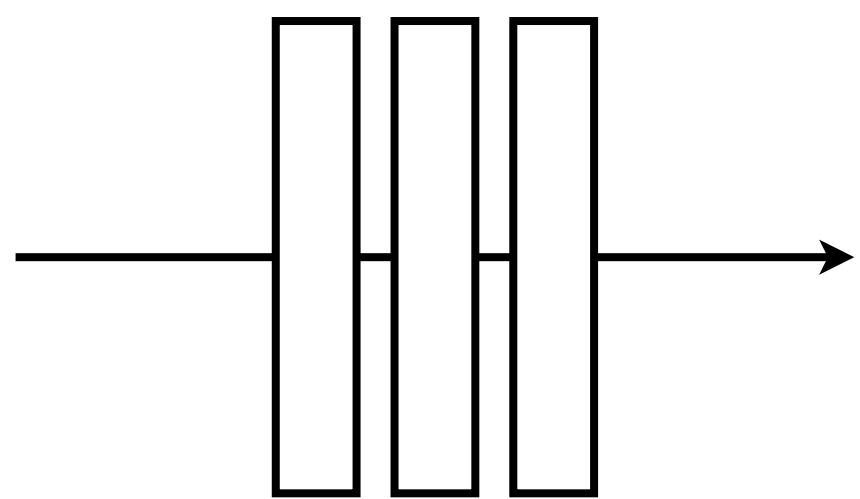
G



G(x)



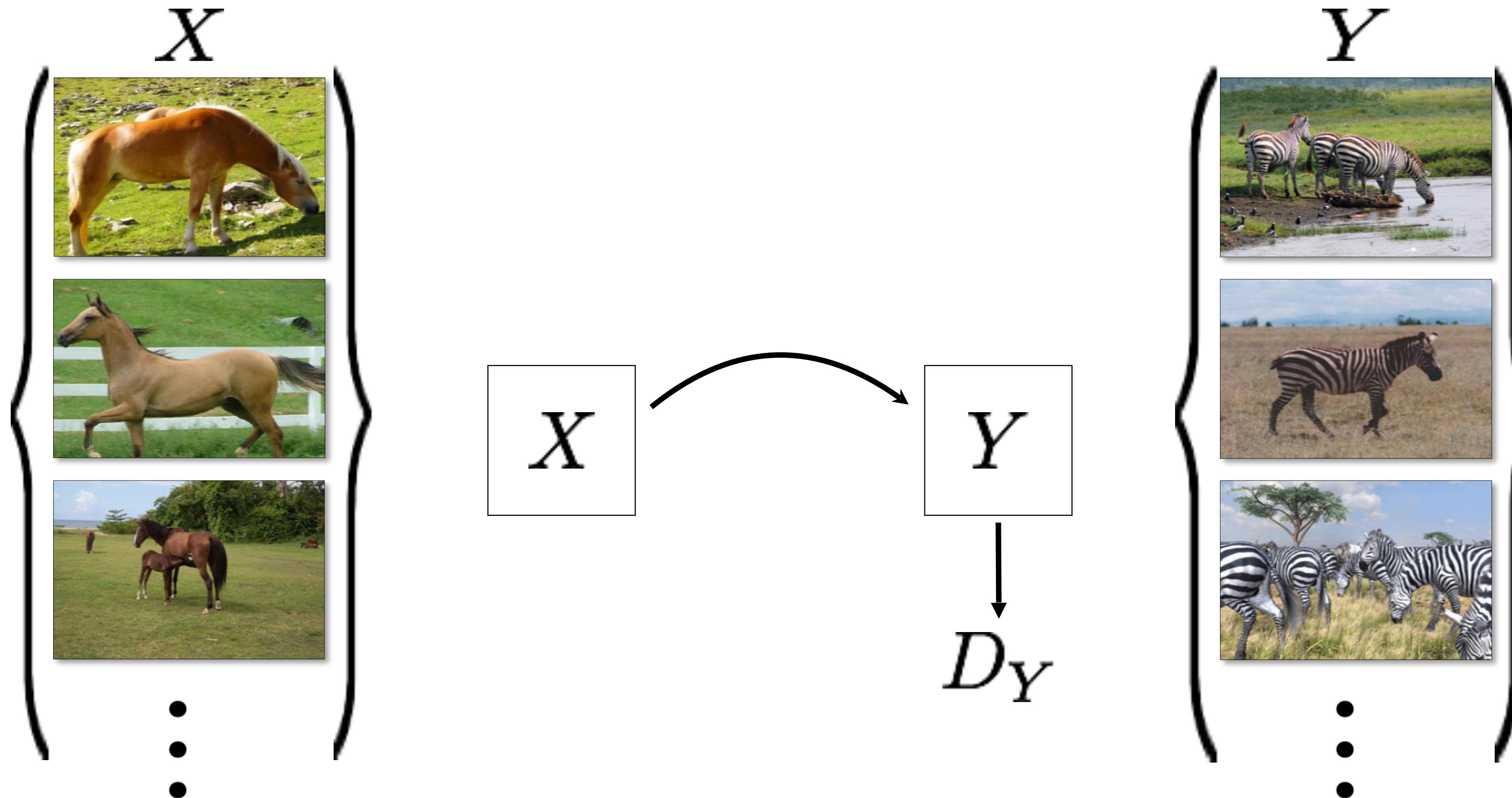
D



Real too!

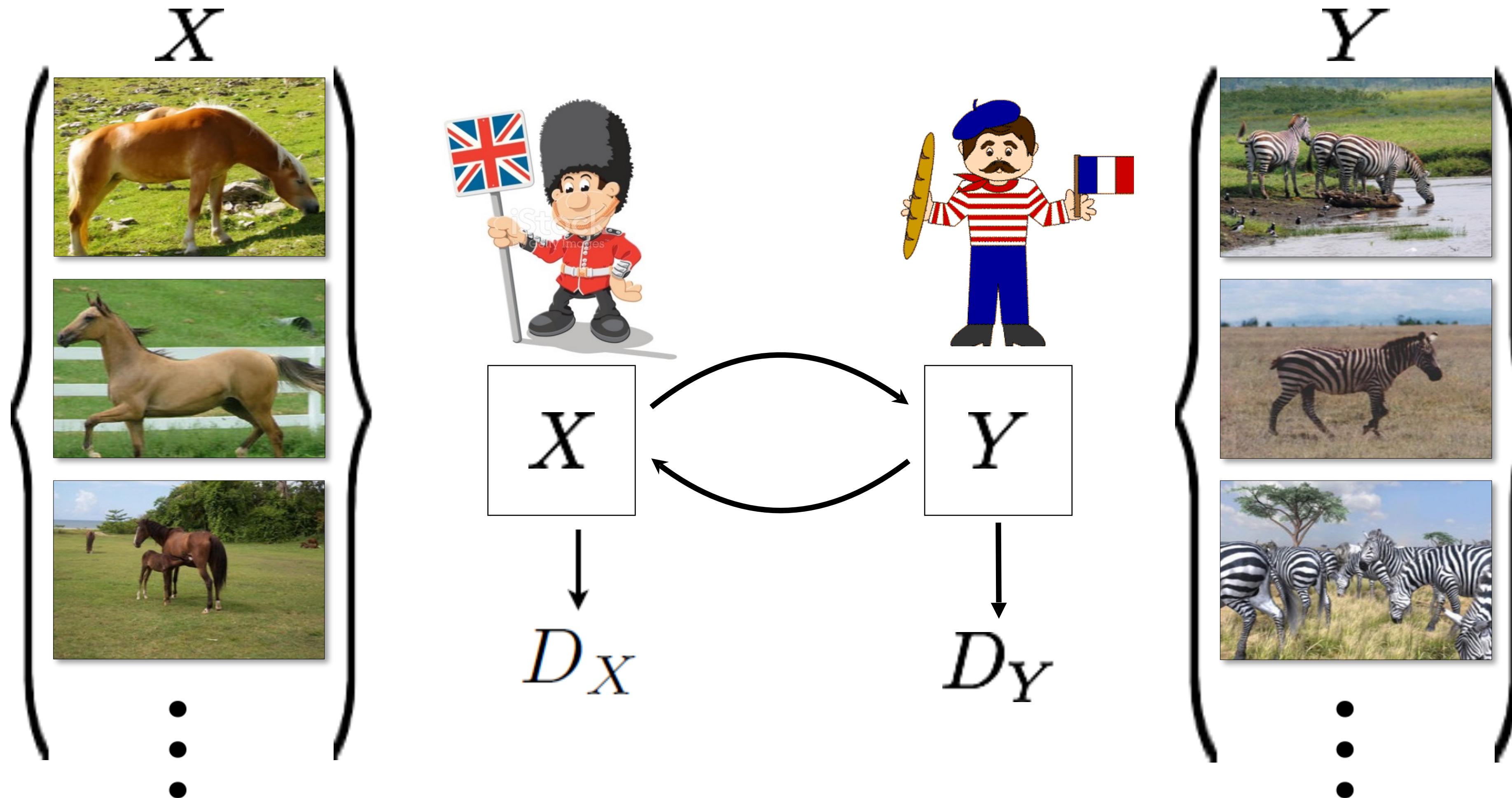
Nothing to force output to correspond to input

Cycle-Consistent Adversarial Networks

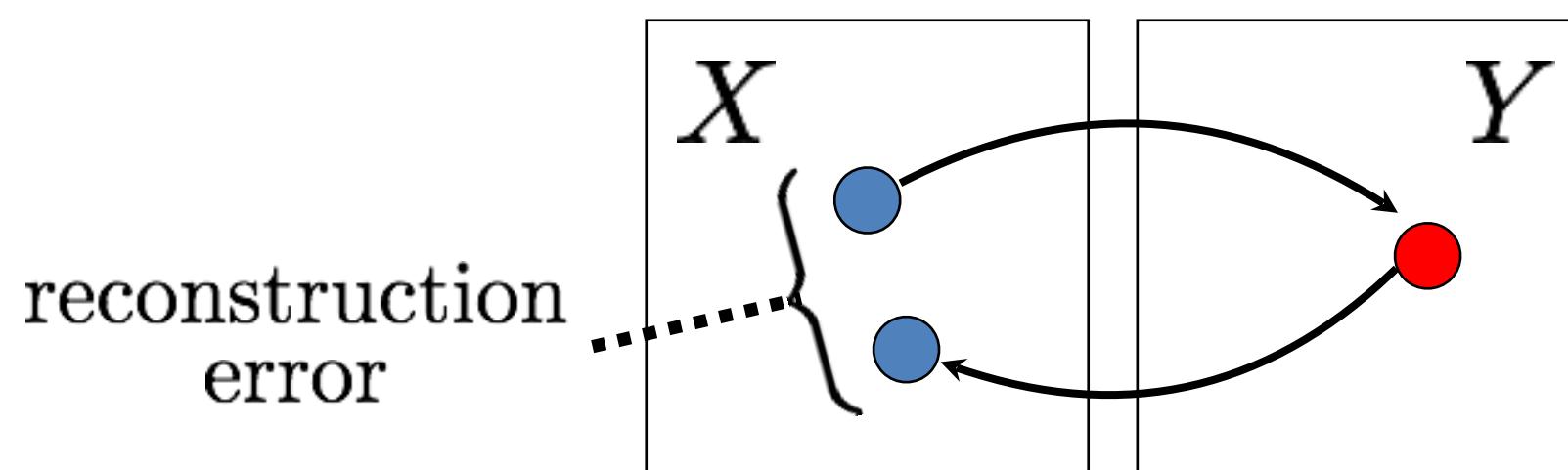
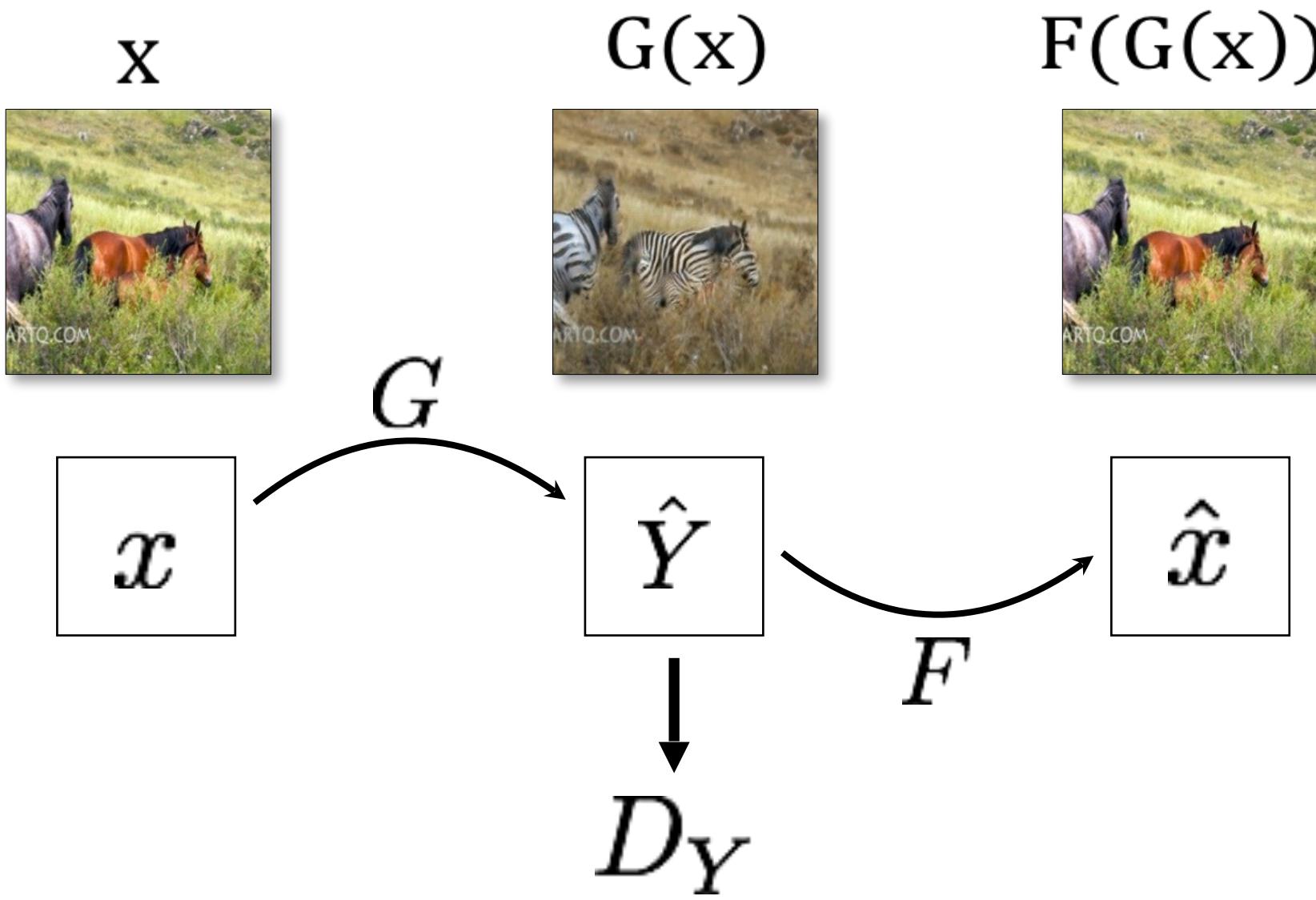


[Zhu et al. 2017], [Yi et al. 2017], [Kim et al. 2017]

Cycle-Consistent Adversarial Networks

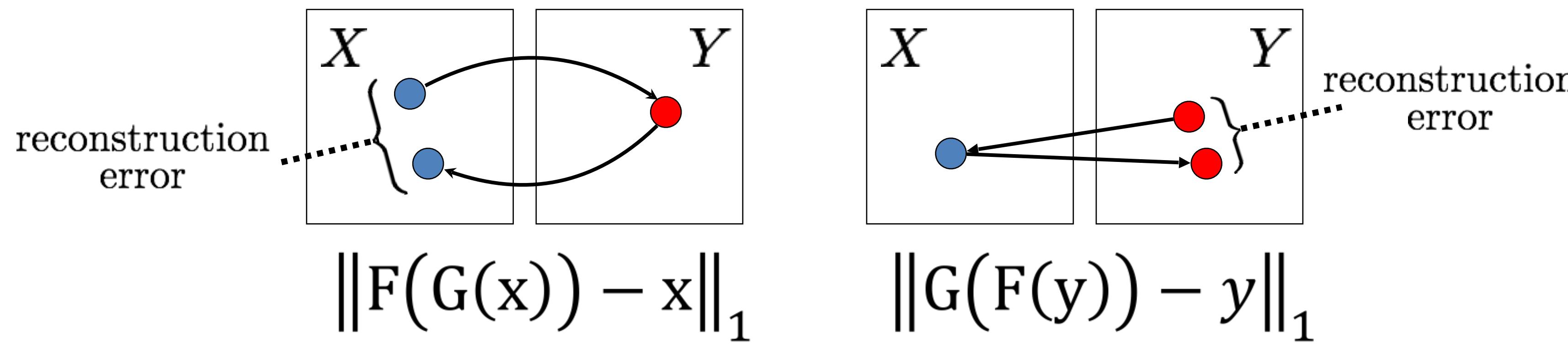
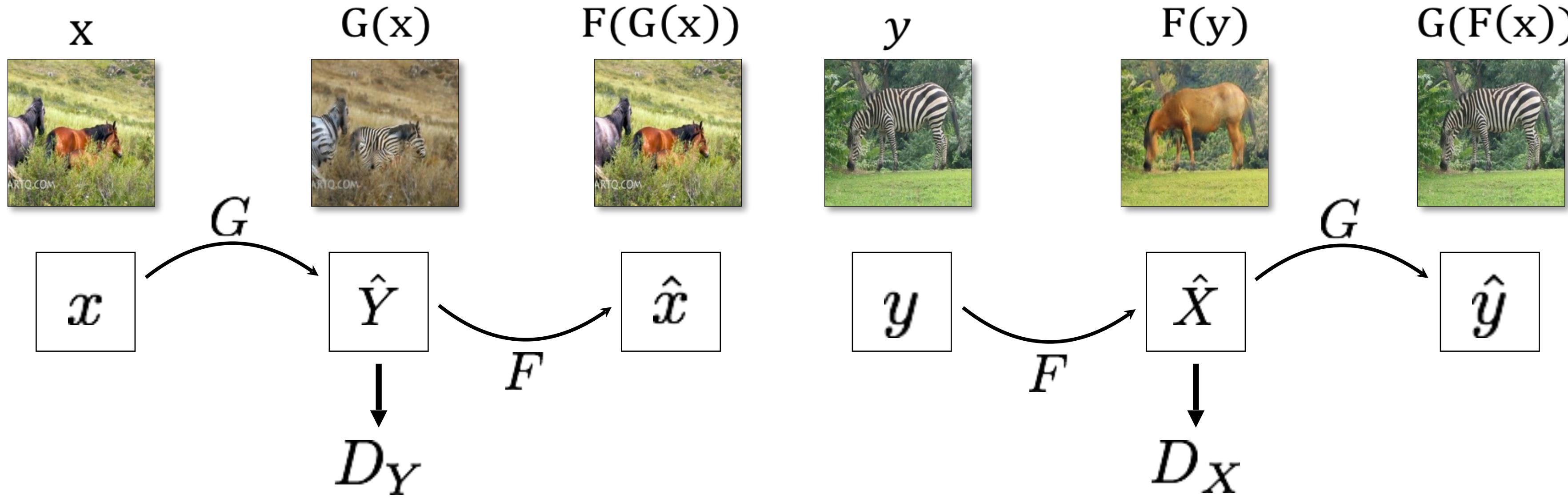


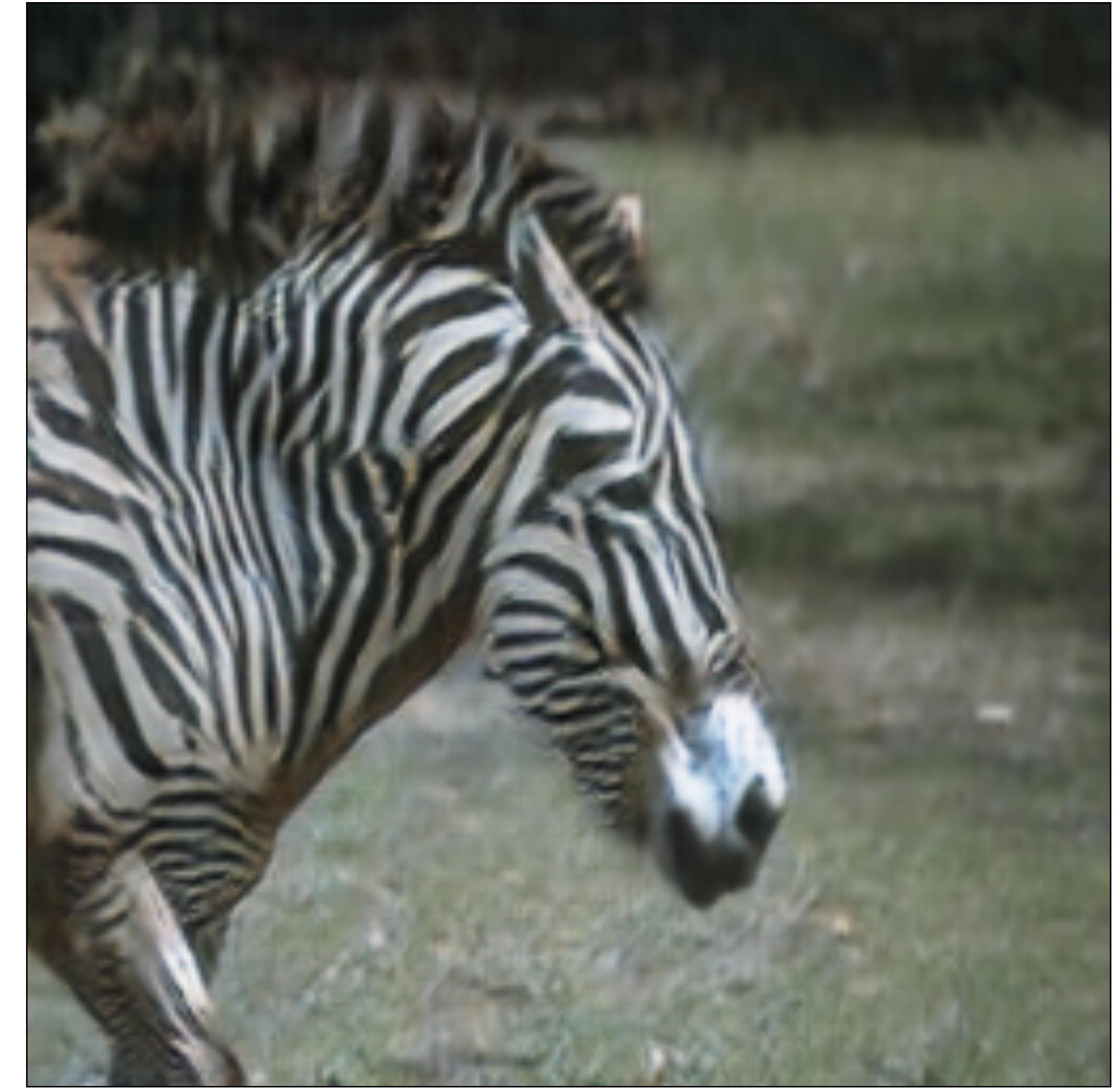
Cycle Consistency Loss



$$\|F(G(x)) - x\|_1$$

Cycle Consistency Loss







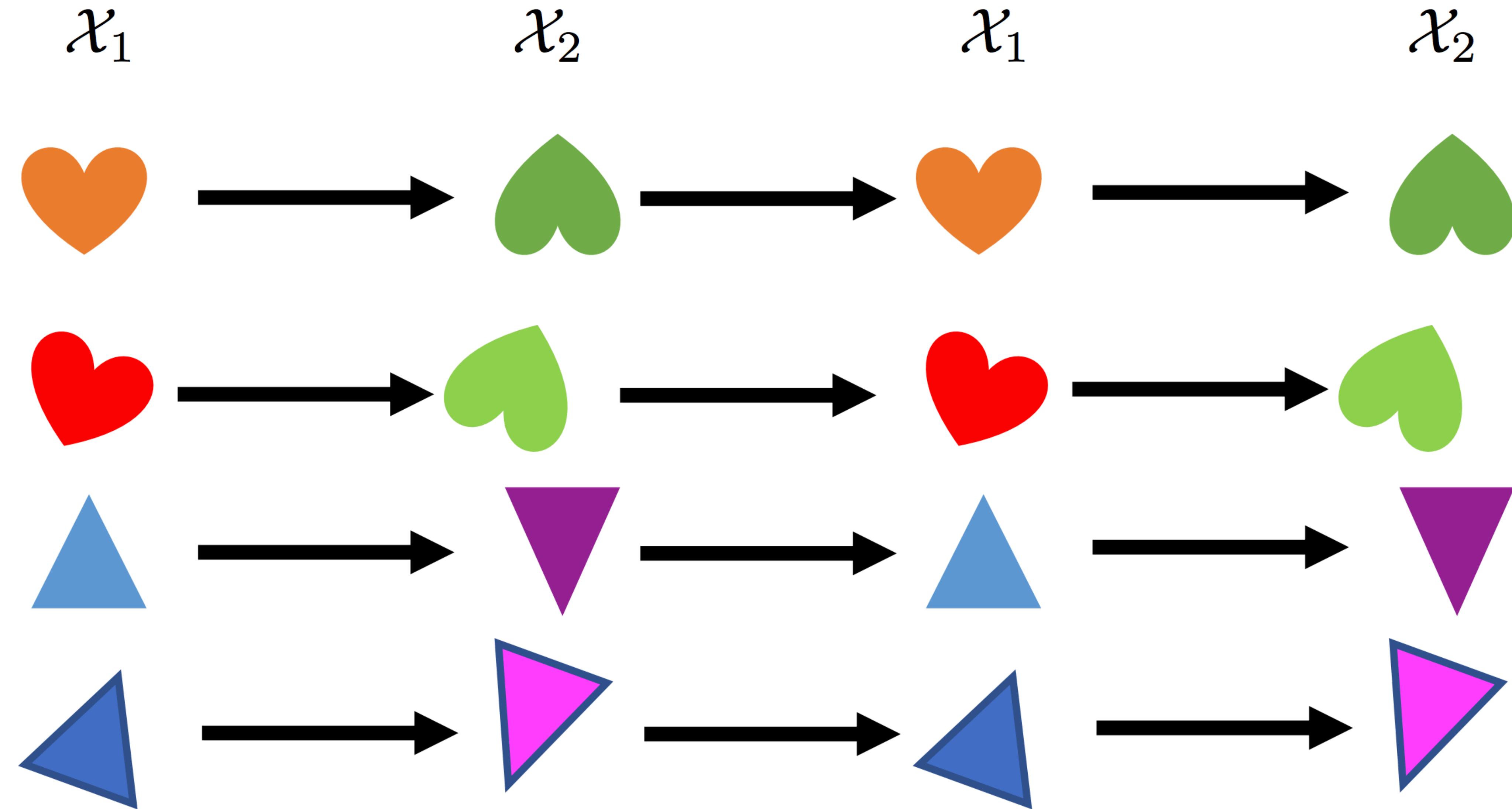
Failure case



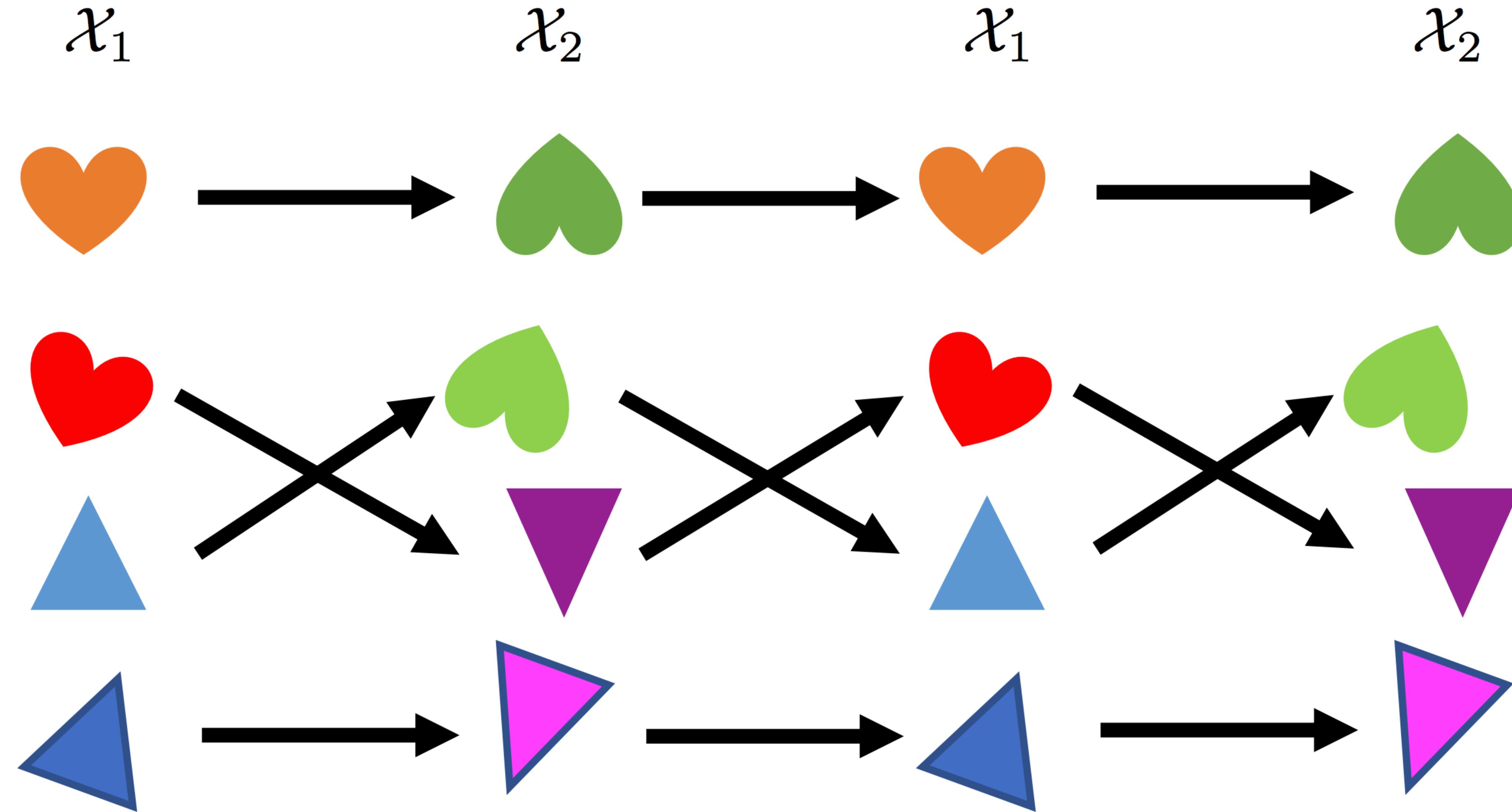
Failure case



Why does CycleGAN work?



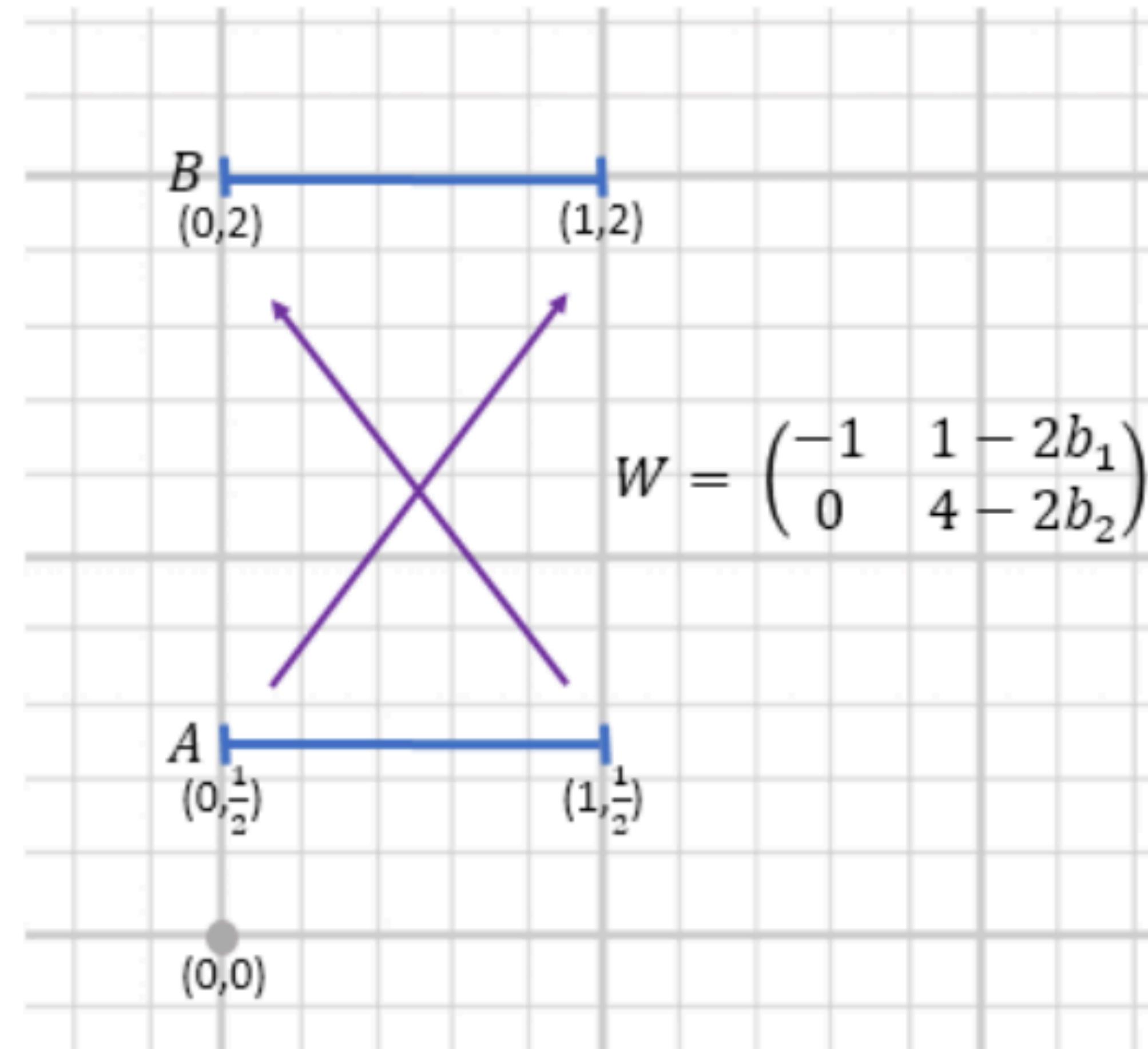
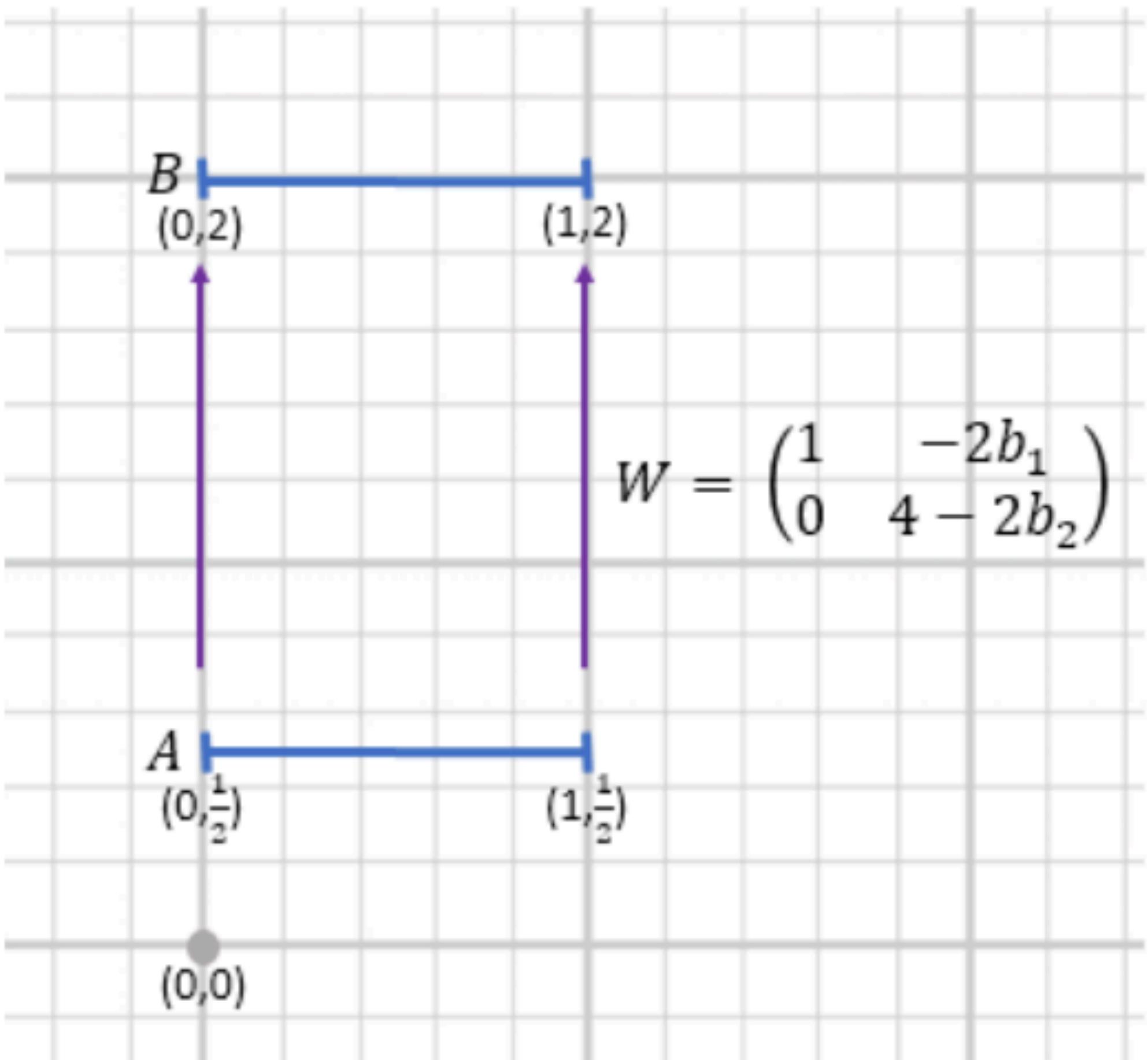
Slide credit: Ming-Yu Liu



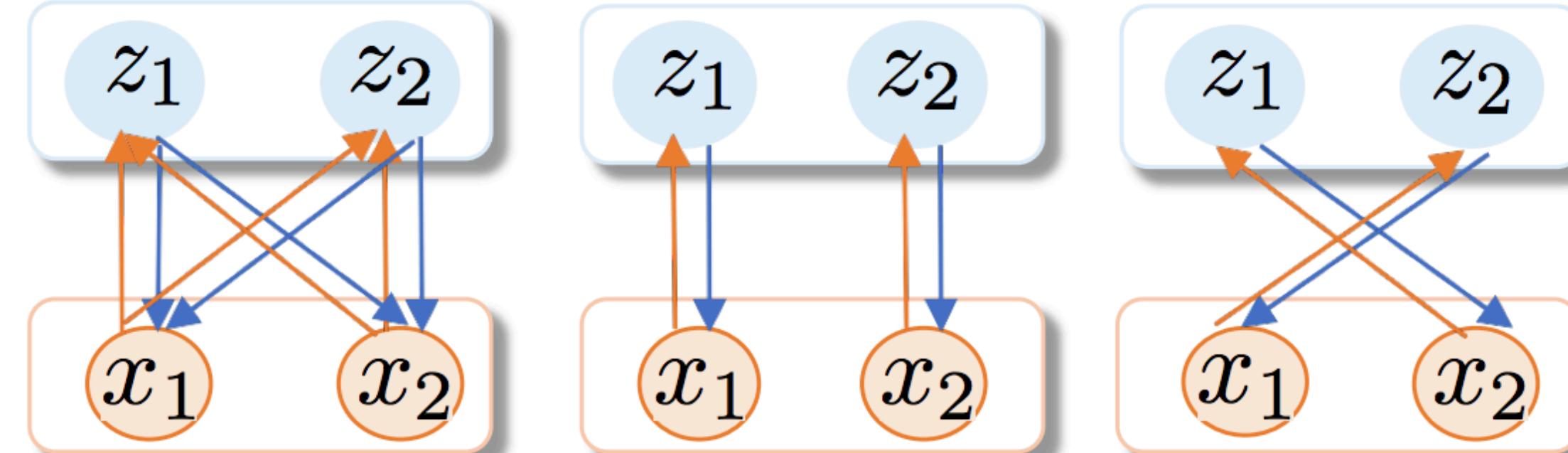
Slide credit: Ming-Yu Liu

Simplicity hypothesis

[Galanti, Wolf, Benaim, 2018]



Cycle Loss upper bounds Conditional Entropy



	z_1	z_2
x_1	$\delta/2$	$(1-\delta)/2$
x_2	$(1-\delta)/2$	$\delta/2$

	z_1	z_2
x_1	$1/2$	0
x_2	0	$1/2$

	z_1	z_2
x_1	0	$1/2$
x_2	$1/2$	0

High
Conditional
Entropy

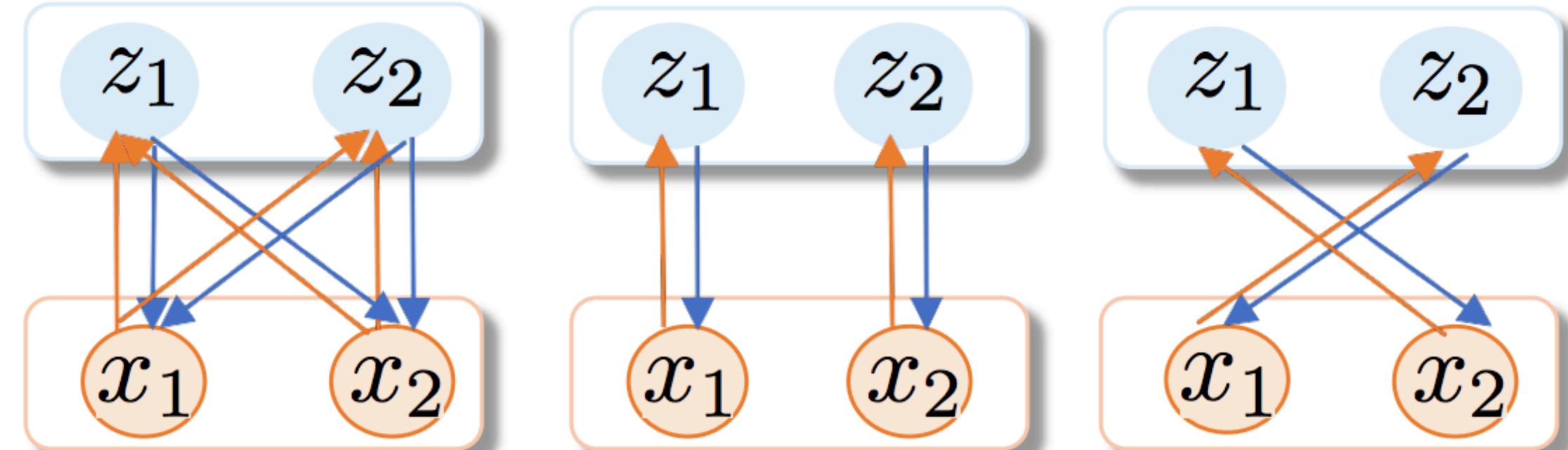
Low
Conditional
Entropy

Conditional Entropy

$$H^\pi(\mathbf{x}|\mathbf{z}) \triangleq -\mathbb{E}_{\pi(\mathbf{x}, \mathbf{z})}[\log \pi(\mathbf{x}|\mathbf{z})]$$

“ALICE: Towards Understanding Adversarial Learning for Joint Distribution Matching” [Li et al. NIPS 2017]. Also see [Tiao et al. 2018] “CycleGAN as Approximate Bayesian Inference”

Cycle Loss upper bounds Conditional Entropy



	z_1	z_2
x_1	$\delta/2$	$(1-\delta)/2$
x_2	$(1-\delta)/2$	$\delta/2$

	z_1	z_2
x_1	$1/2$	0
x_2	0	$1/2$

	z_1	z_2
x_1	0	$1/2$
x_2	$1/2$	0

Conditional Entropy

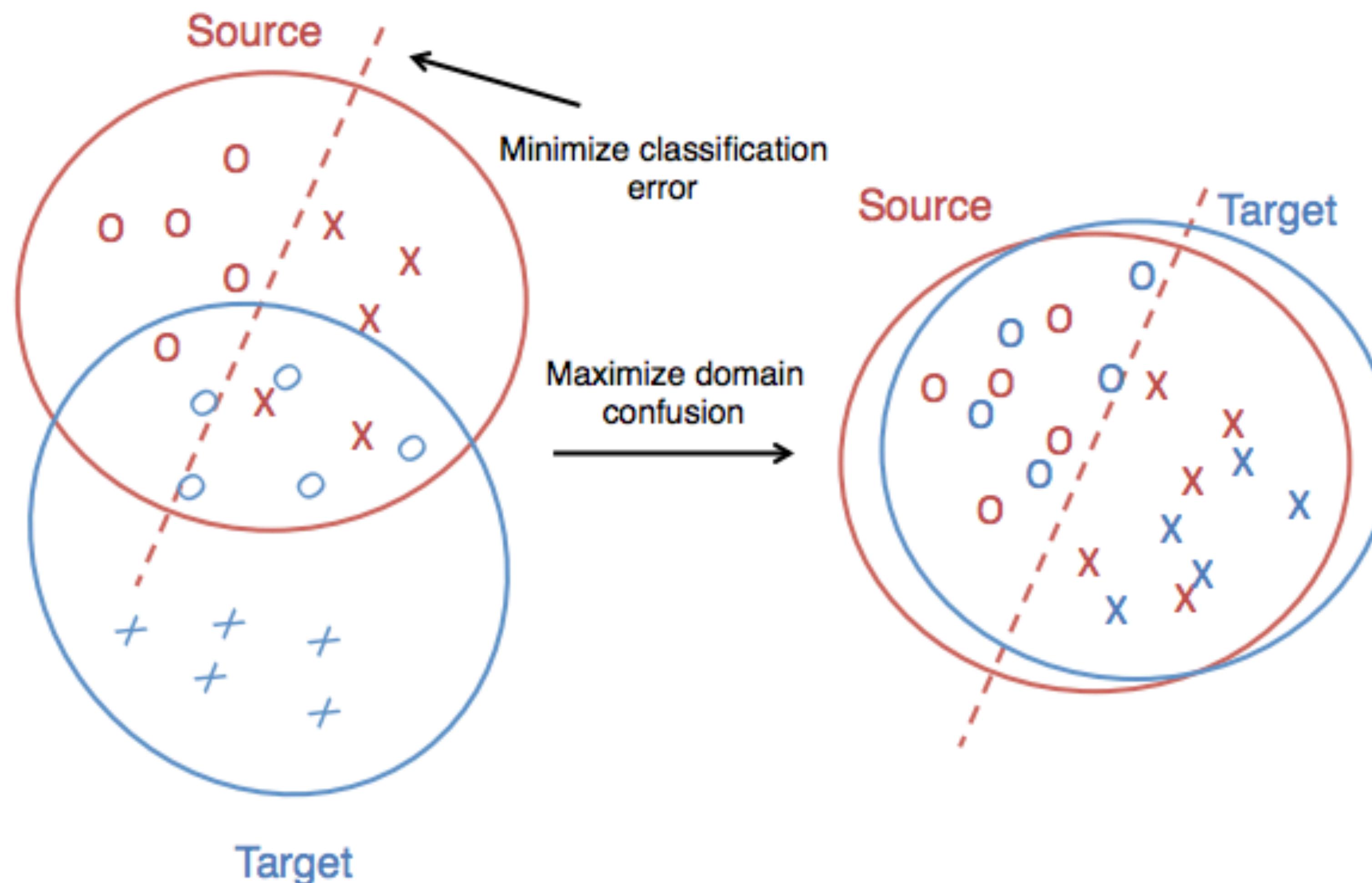
$$H^\pi(\mathbf{x}|\mathbf{z}) \triangleq -\mathbb{E}_{\pi(\mathbf{x}, \mathbf{z})}[\log \pi(\mathbf{x}|\mathbf{z})]$$

Lemma 3 For joint distributions $p_\theta(\mathbf{x}, \mathbf{z})$ or $q_\phi(\mathbf{x}, \mathbf{z})$, we have

$$\begin{aligned} H^{q_\phi}(\mathbf{x}|\mathbf{z}) &\triangleq -\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\log q_\phi(\mathbf{x}|\mathbf{z})] = -\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z})}[\text{KL}(q_\phi(\mathbf{x}|\mathbf{z}) \| p_\theta(\mathbf{x}|\mathbf{z}))] \\ &\leq -\mathbb{E}_{q_\phi(\mathbf{x}, \mathbf{z})}[\log p_\theta(\mathbf{x}|\mathbf{z})] \triangleq \mathcal{L}_{\text{Cycle}}(\theta, \phi). \end{aligned} \quad (6)$$

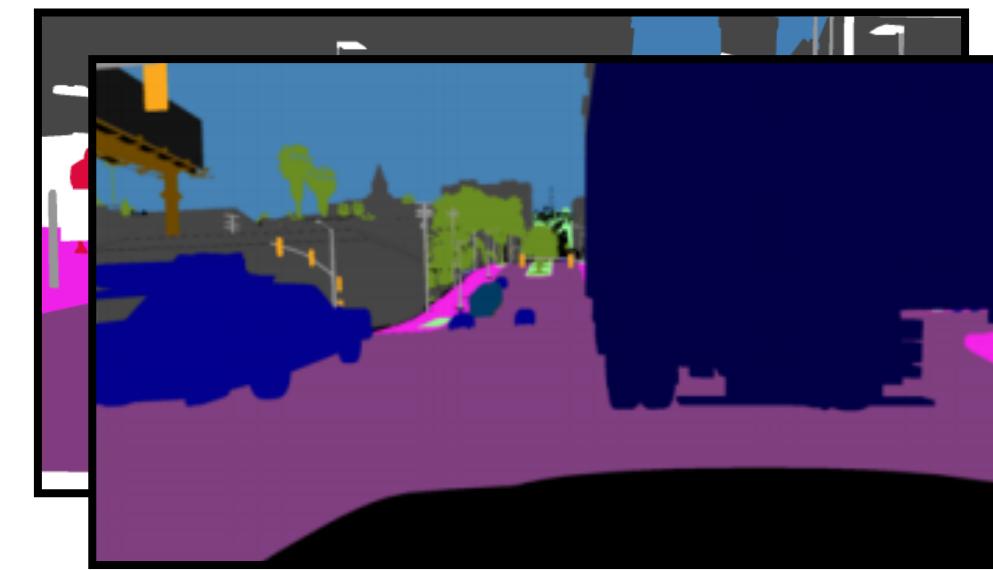
“ALICE: Towards Understanding Adversarial Learning for Joint Distribution Matching” [Li et al. NIPS 2017]. Also see [Tiao et al. 2018] “CycleGAN as Approximate Bayesian Inference”

Domain Adaptation



Sim2real

Simulated data

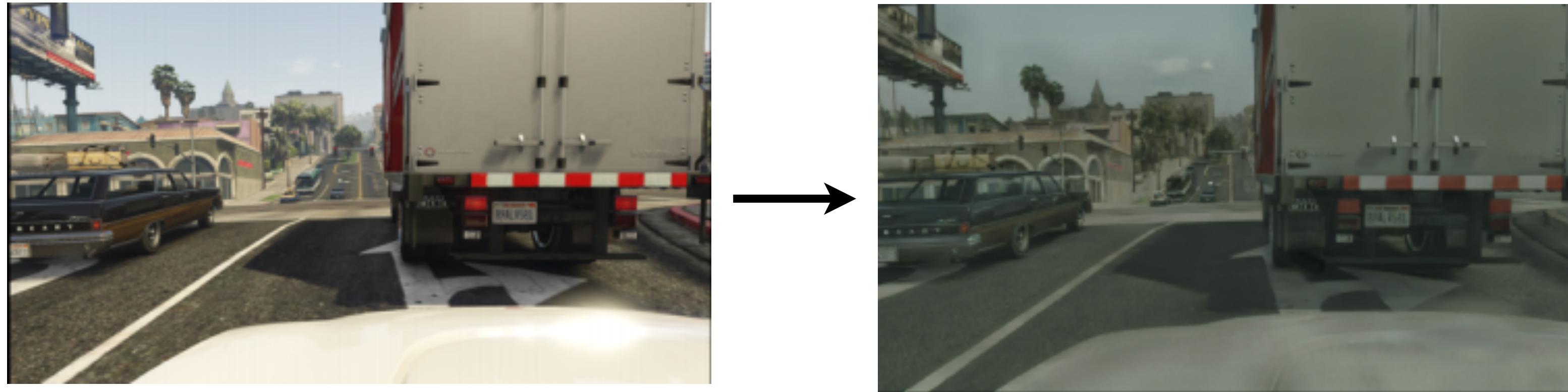


Real data



?

CycleGAN

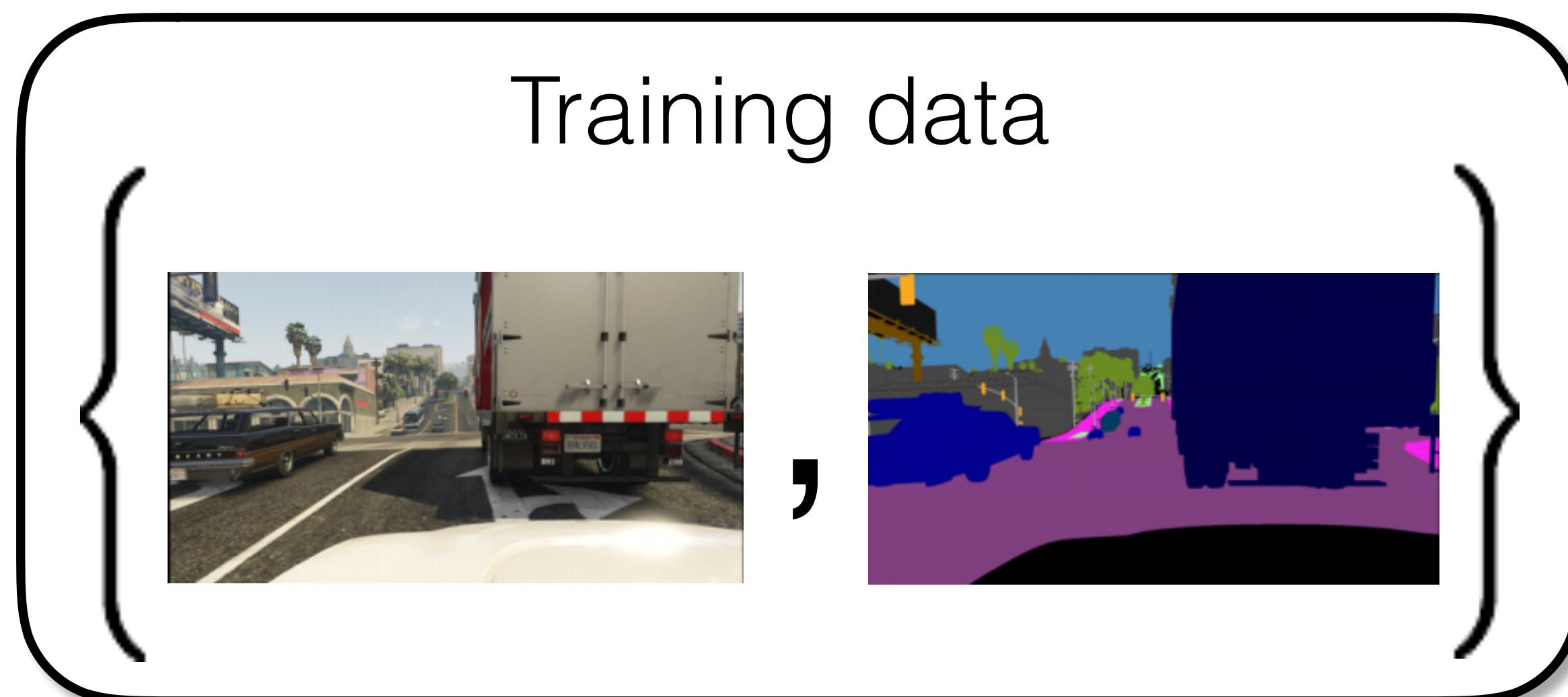
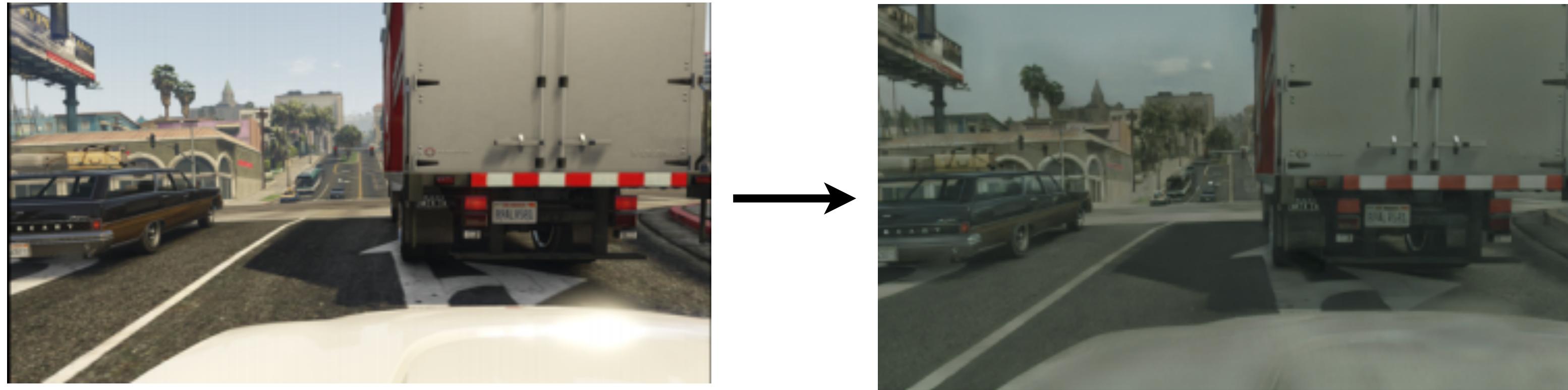


Training data



[Hoffman, Tzeng, Park, Zhu, Isola, Saenko, Darrell, Efros, 2018]

CycleGAN

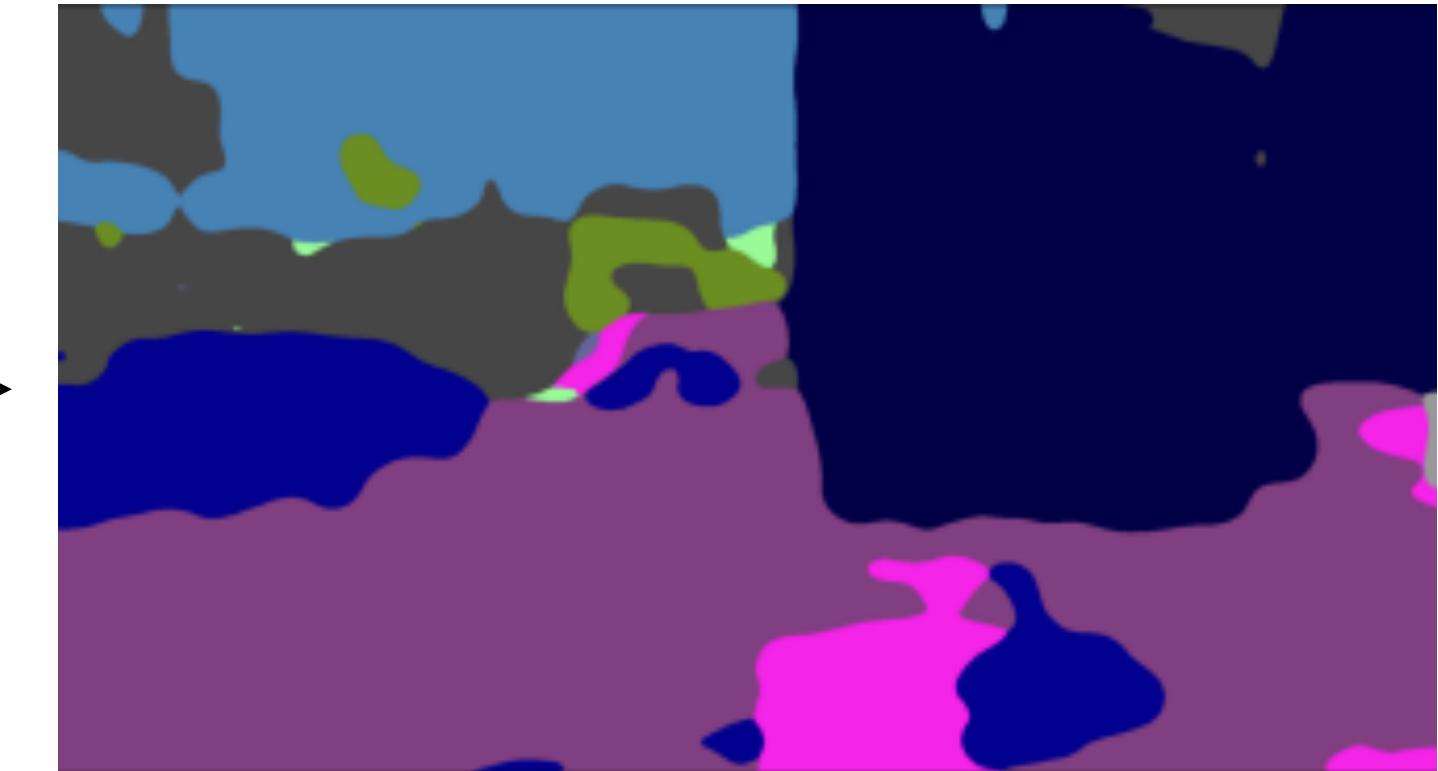


[Hoffman, Tzeng, Park, Zhu, Isola, Saenko, Darrell, Efros, 2018]

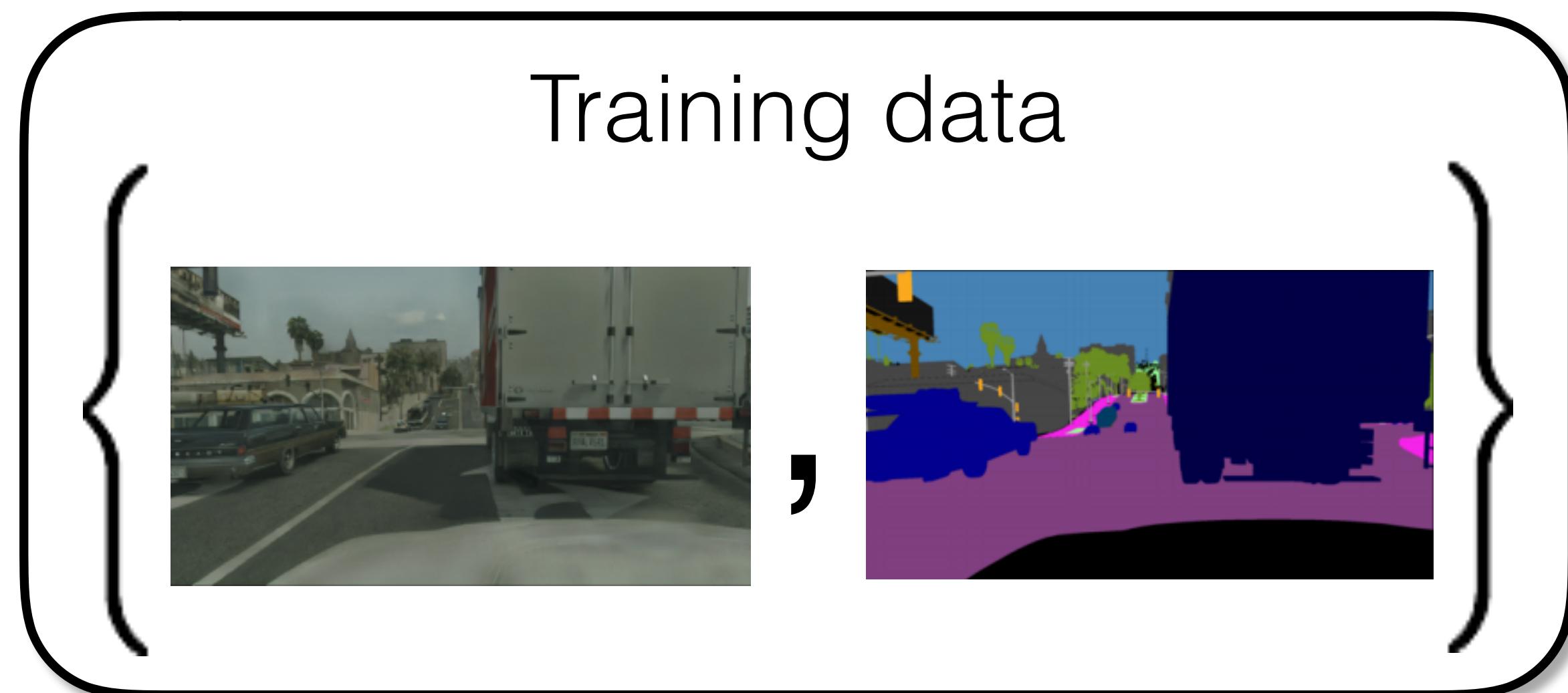
CycleGAN



FCN



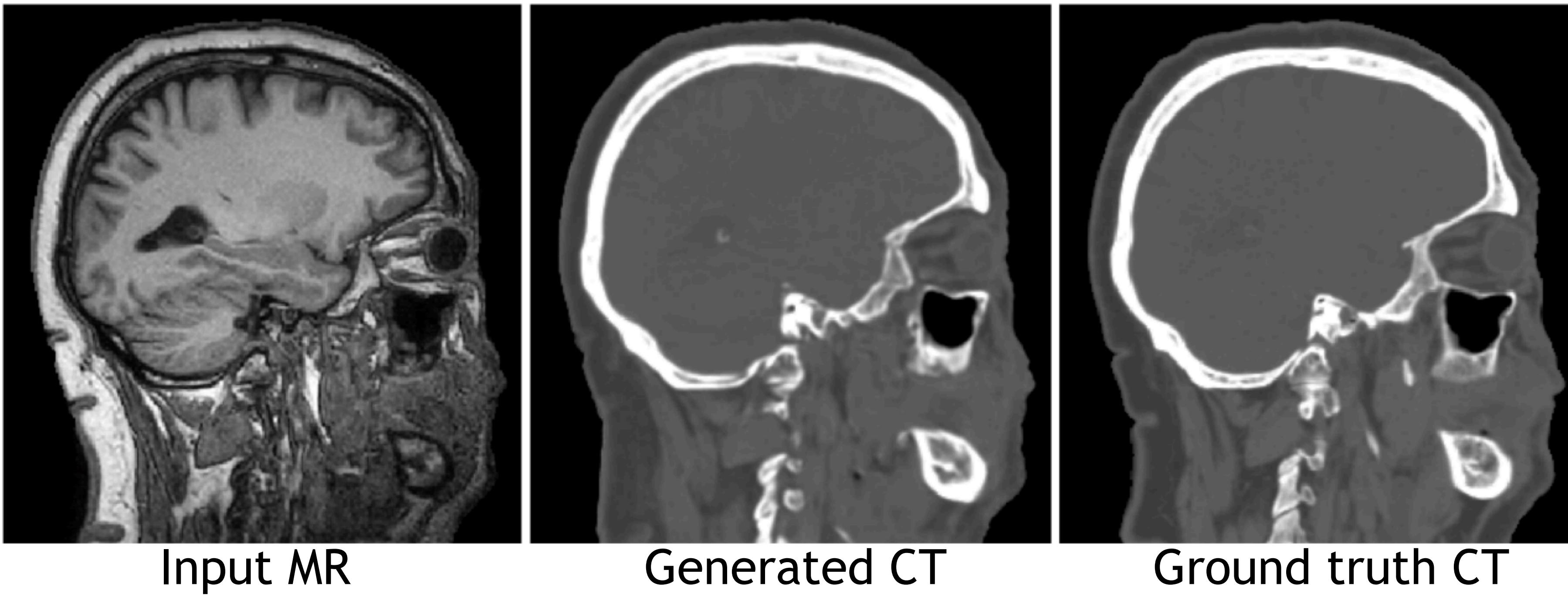
Training data



[Hoffman, Tzeng, Park, Zhu, Isola, Saenko, Darrell, Efros, 2018]

Medical domain adaptation

MR → CT [Wolterink et al] arxiv: 1708.01155



- MRI reconstruction [Quan et al.] arxiv:1709.00753
- Cardiac MR images from CT [Chartsias et al. 2017]

Three perspectives on GANs

1. Structured loss
2. Generative model
3. Domain-level supervision / mapping

Thank you!