# Generative models
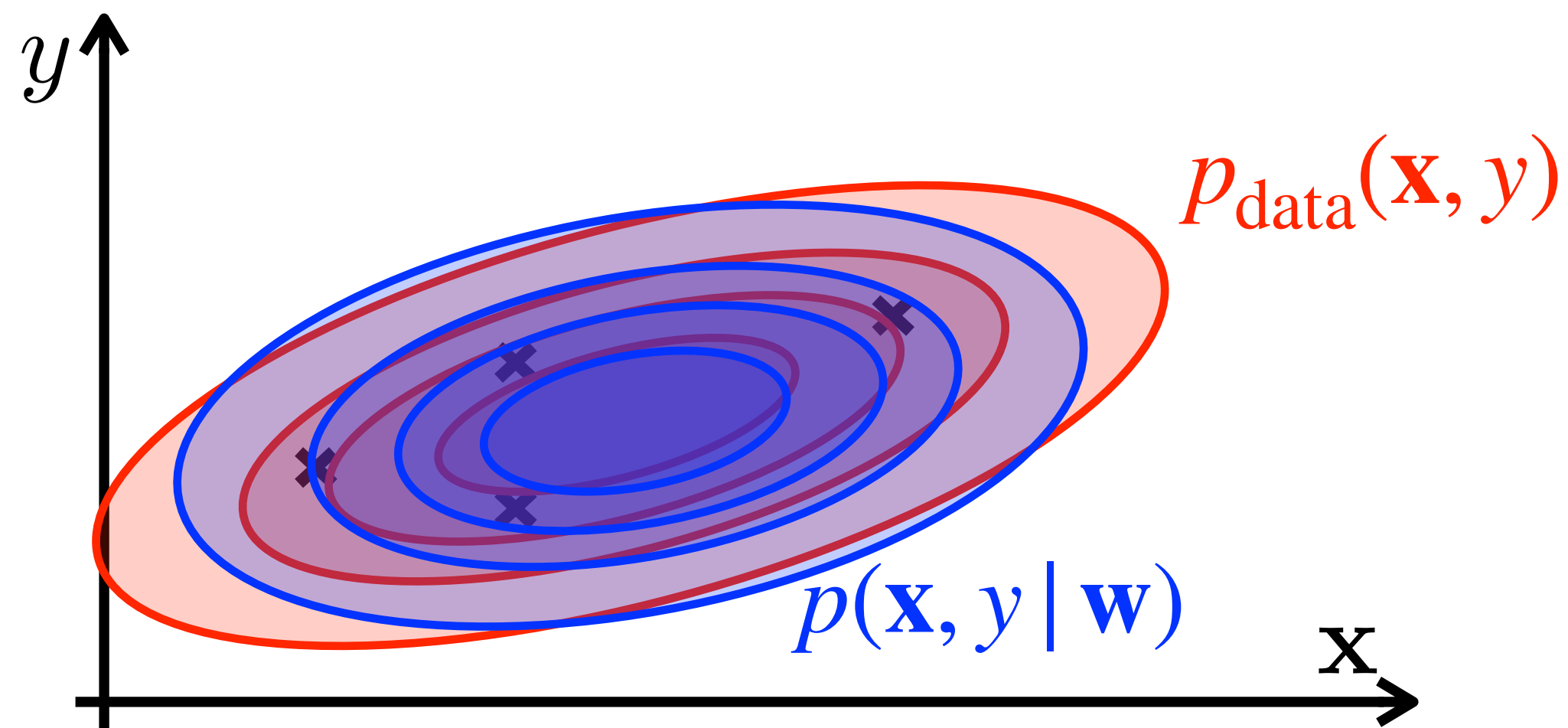
## VAE, GANS, Diffusion models

**Karel Zimmermann**

# Prerequisites: Learning vs optimization

$$\mathbf{w}^\star = \arg\min_{\mathbf{w}} D_{KL}\big(p_{\text{data}}(\mathbf{x}, y) \,\|\, p(\mathbf{x}, y \,|\, \mathbf{w})\big) = \int_{(\mathbf{x}, y)} p_{\text{data}}(\mathbf{x}, y) \cdot \log \frac{p_{\text{data}}(\mathbf{x}, y)}{p(\mathbf{x}, y \,|\, \mathbf{w})}$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}} \left[ \log \frac{p_{data}(\mathbf{x}, y)}{p(\mathbf{x}, y \,|\, \mathbf{w})} \right] = \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}} \left[ -\log p(\mathbf{x}, y \,|\, \mathbf{w}) \right]$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}} \left[ \log p_{\text{data}}(\mathbf{x}, y) - \log p(y \,|\, \mathbf{x}, \mathbf{w}) p(\mathbf{x}) \right]$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}} \left[ -\log p(y \,|\, \mathbf{x}, \mathbf{w}) \right] \approx \arg\min_{\mathbf{w}} \frac{1}{N} \sum_{(\mathbf{x}_i, y_i) \sim p_{\text{data}}(\mathbf{x}, y)} \left[ -\log p(y_i \,|\, \mathbf{x}_i, \mathbf{w}) \right]$$



$p_{\text{data}}(\mathbf{x}, y)$
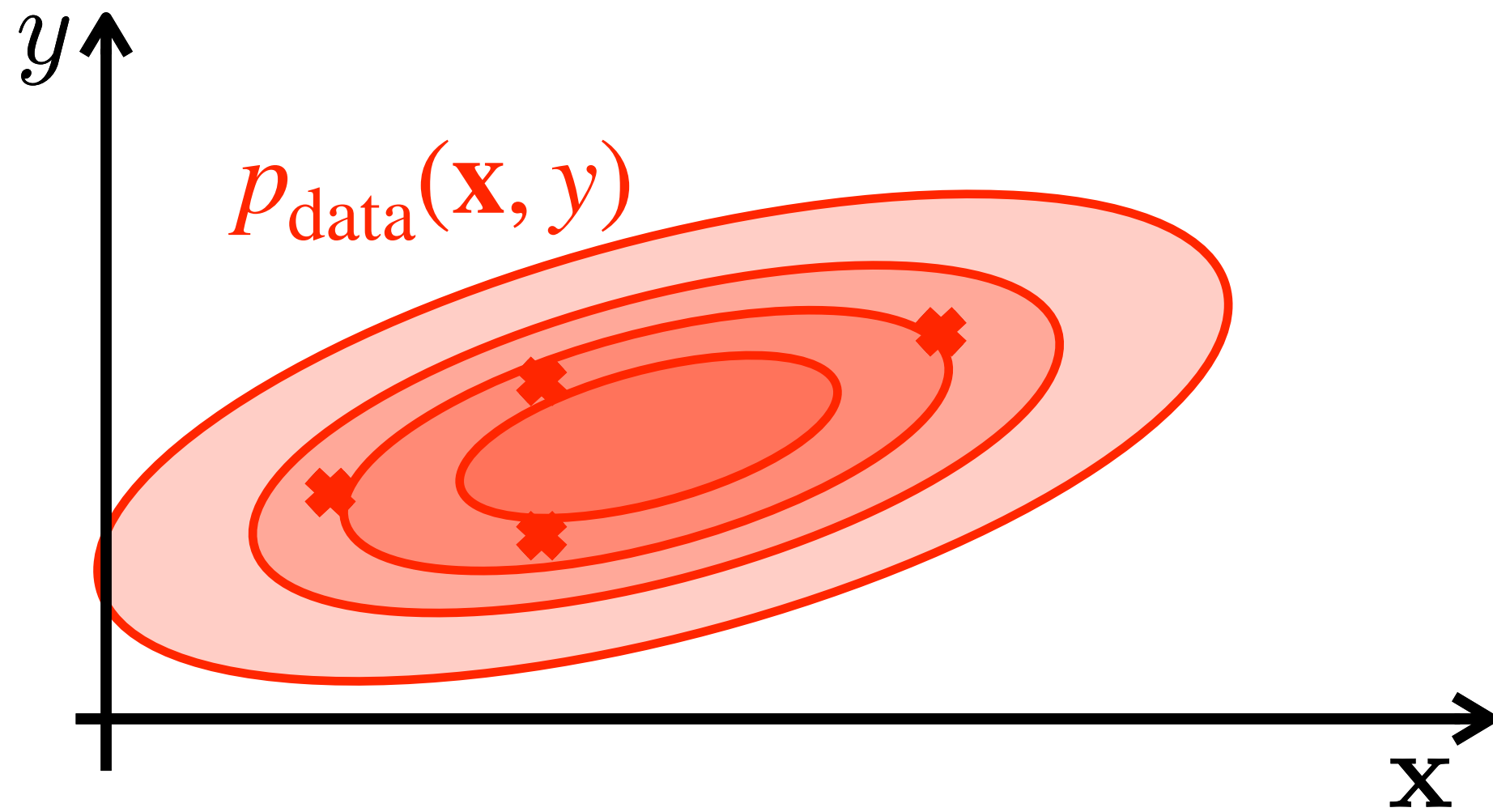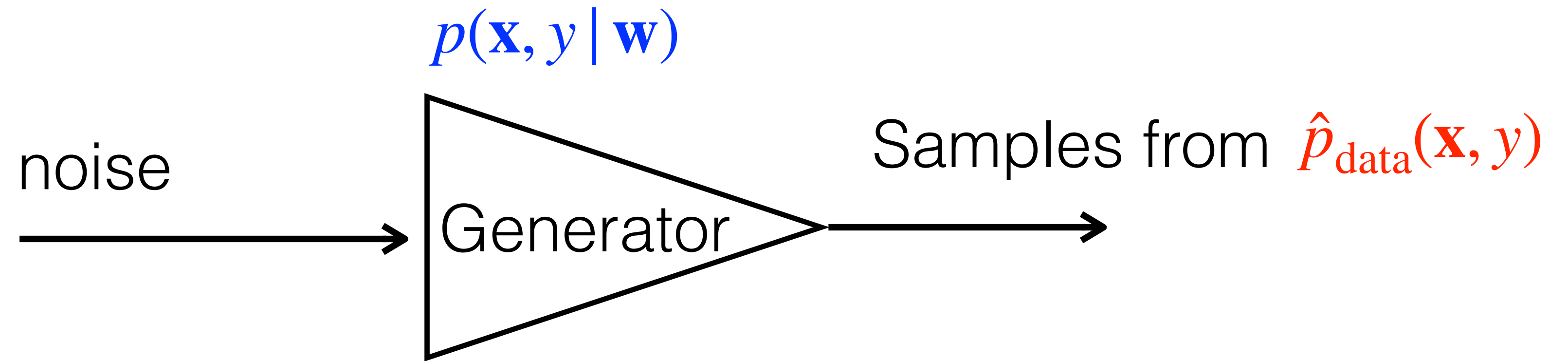
$p(\mathbf{x}, y \,|\, \mathbf{w})$

So-far we always modeled $p(y_i \,|\, \mathbf{x}_i, \mathbf{w})$
How comes?

What else can be modeled?

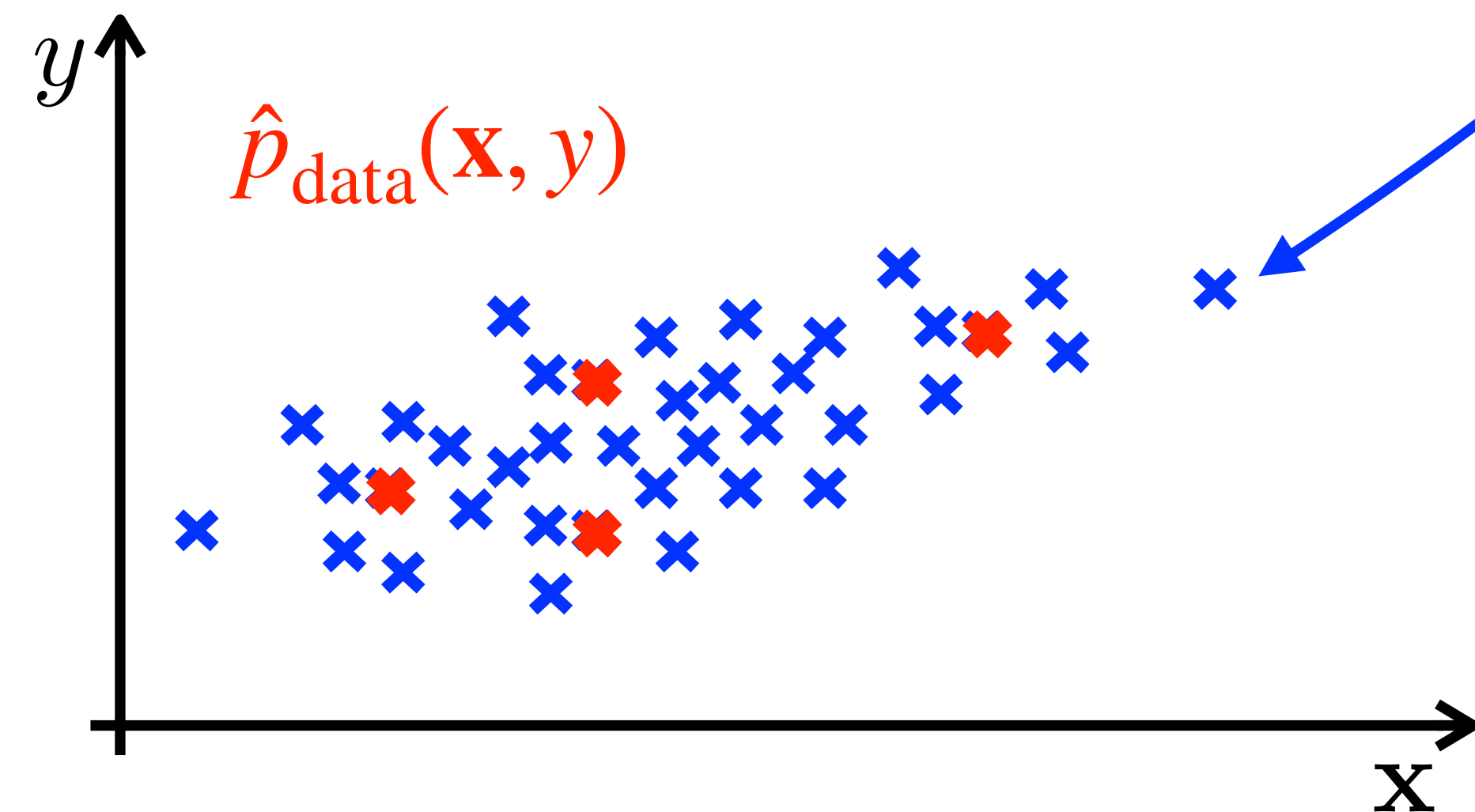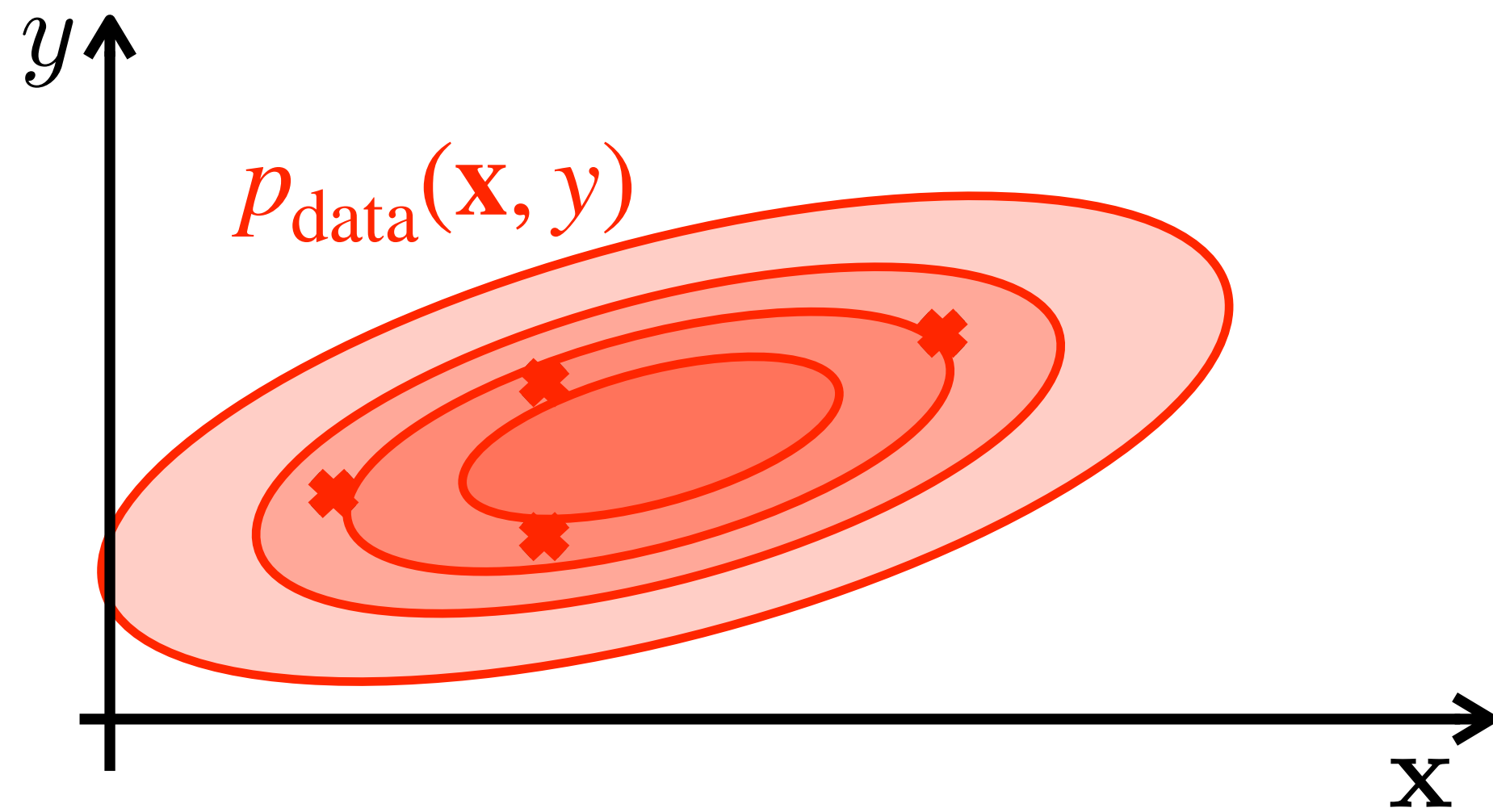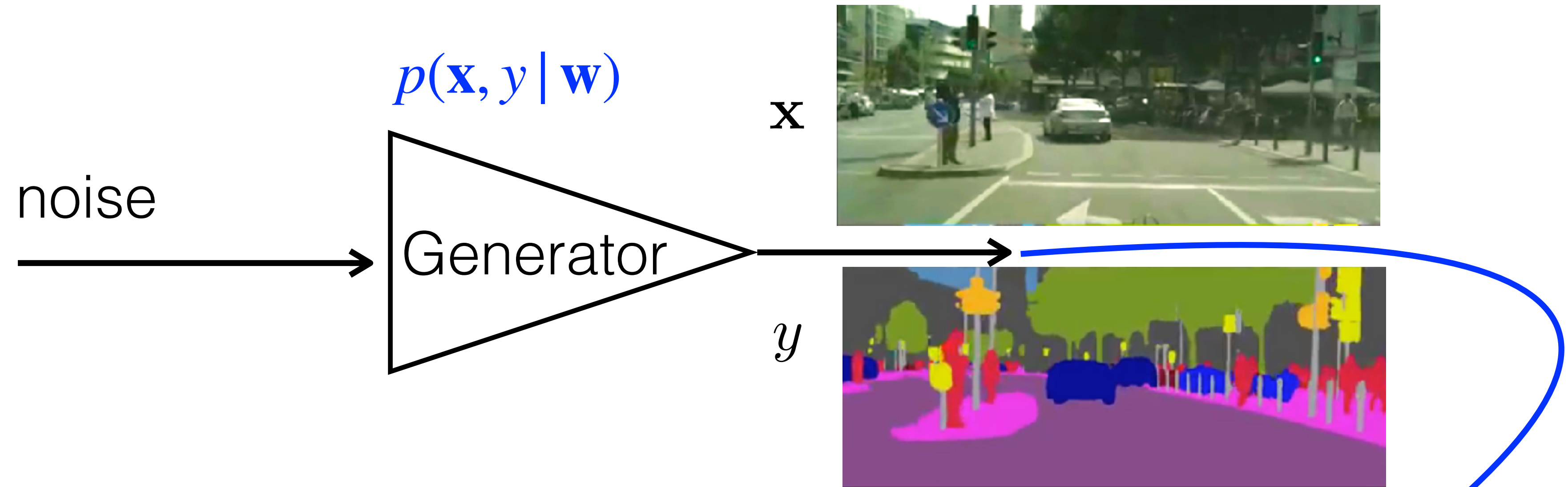$p(\mathbf{x}, y \,|\, \mathbf{w}) \qquad p(\mathbf{x} \,|\, y, \mathbf{w}) \qquad p(\mathbf{x} \,|\, \mathbf{w})$

# Dataset augmentation

$p(\mathbf{x}, y \mid \mathbf{w})$

noise

Generator

Samples from $\hat{p}_{\text{data}}(\mathbf{x}, y)$

$y$

$p_{\text{data}}(\mathbf{x}, y)$

$\mathbf{x}$

# Dataset augmentation

Samples from $\hat{p}_{\text{data}}(\mathbf{x}, y)$

$p(\mathbf{x}, y \mid \mathbf{w})$

noise

Generator

$\mathbf{x}$

$y$

$y$

$p_{\text{data}}(\mathbf{x}, y)$

$\mathbf{x}$

$y$

$\hat{p}_{\text{data}}(\mathbf{x}, y)$

$\mathbf{x}$

Dataset augmentation

Samples from $\hat{p}_{\text{data}}(\mathbf{x}\,|\,y)$

annotation

photo

$y$

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

$\mathbf{x}$

+ noise

Generator

Gives better control over generated data

Domain transfer

Samples from $\hat{p}_{\text{data}}(\mathbf{x}|y)$

photo

painting

$y$

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

$\mathbf{x}$

Generator

+ noise

## Classification dataset augmentation

$p(\mathbf{x} \mid \mathbf{w})$

noise

Generator

random walk in noise space

[ProGAN]

- Dataset augmentation
- Realistic simulator
- Generative networks

# RGB images

Depth images

Stencil layer

# Stencil layer - **cars**



glass is not car

# Stencil layer - **humans**



driver 's hand

# Stencil layer - **trees**

# Stencil layer - **sky**

# Stencil layer - **artificial light**

# Stencil layer - **artificial light**

Other annotations for objects (e.g. cars, humans)

Other annotations for objects (e.g. cars, humans)

# What can be controlled (night/day)



handlePipeInput called

server connected:False

connection:System.Net.Sockets.Socket

# What can be controlled (night/day)

# What can be controlled (weather): ExtraSunny



handlePipeInput called

server connected:False

# What can be controlled (weather): Clear



handlePipeInput called

server connected:False

N pressed, going to take screenshots

connection:System.Net.Socke ts.Socket

# What can be controlled (weather): Foggy



handlePipeInput called

server connected:False

connection:System.Net.Socke
ts.Socket

# What can be controlled (weather): OverCast



handlePipeInput called

server connected:False

connection:System.Net.Socke
ts.Socket

N pressed, going to take
screenshots

# What can be controlled (weather): Raining



handlePipeInput called

N pressed, going to take screenshots

server connected:False

connection:System.Net.Sockets.Socket

# What can be controlled (weather): ThunderStorm



handlePipeInput called

server connected:False

connection:System.Net.Socke
ts.Socket

N pressed, going to take
screenshots

handlePipeInput called

server connected:False

connection:System.Net.Socke
ts.Socket

N pressed, going to take
screenshots

# What can be controlled (weather): SnowLight



handlePipeInput called

N pressed, going to take screenshots

server connected:False

connection:System.Net.Sockets.Socket

# What can be controlled (mods): Tsunami



https://cs.gta5-mods.com/misc/tsunami-mod

# What can be controlled (mods): Plane crashing mods



https://cs.gta5-mods.com/scripts/planes-hails

# Driving in the matrix [Roberson ICRA 2017]
## https://arxiv.org/abs/1610.01983

- Reverse engineering of GTA 5 (RAGE engine)



easy

moderate
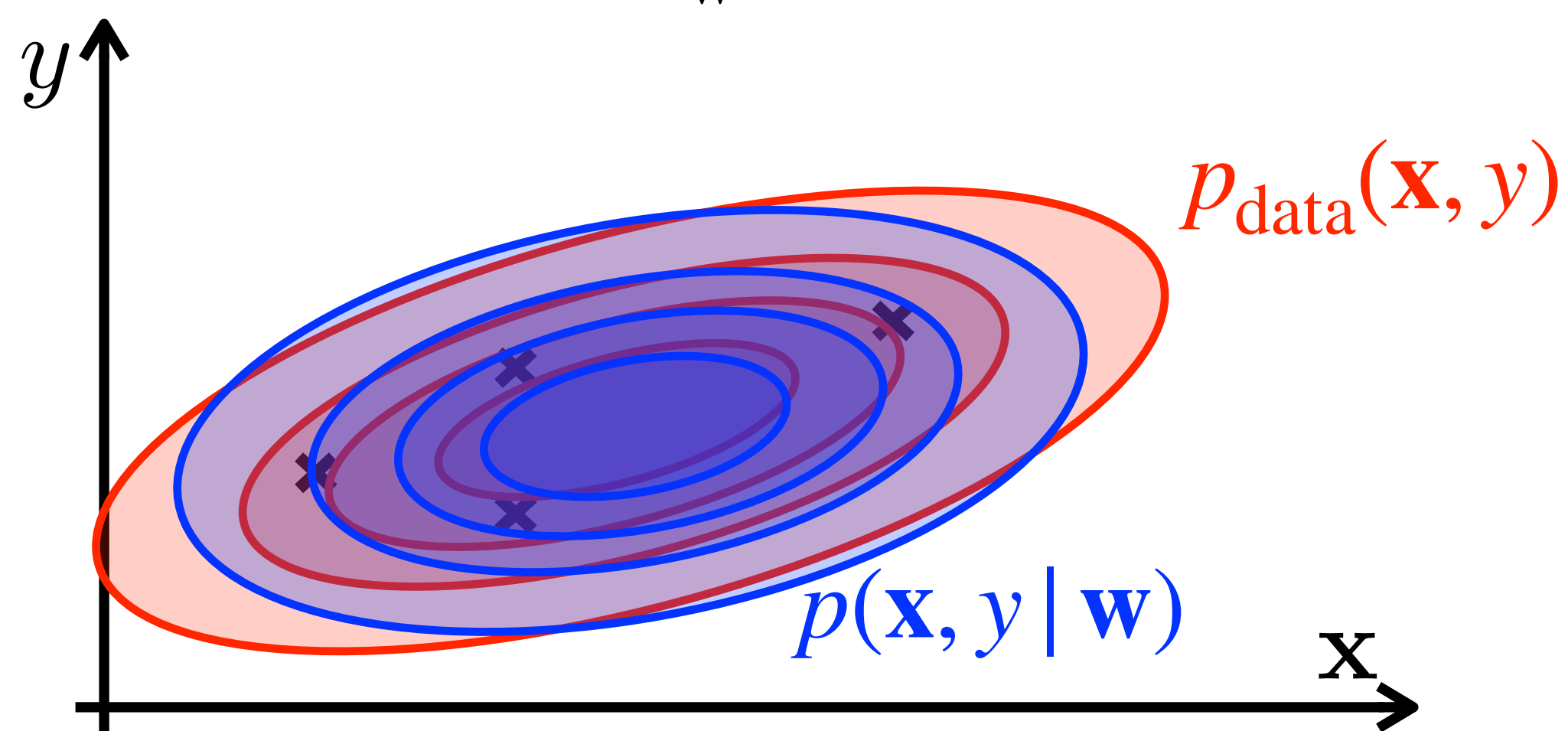
hard

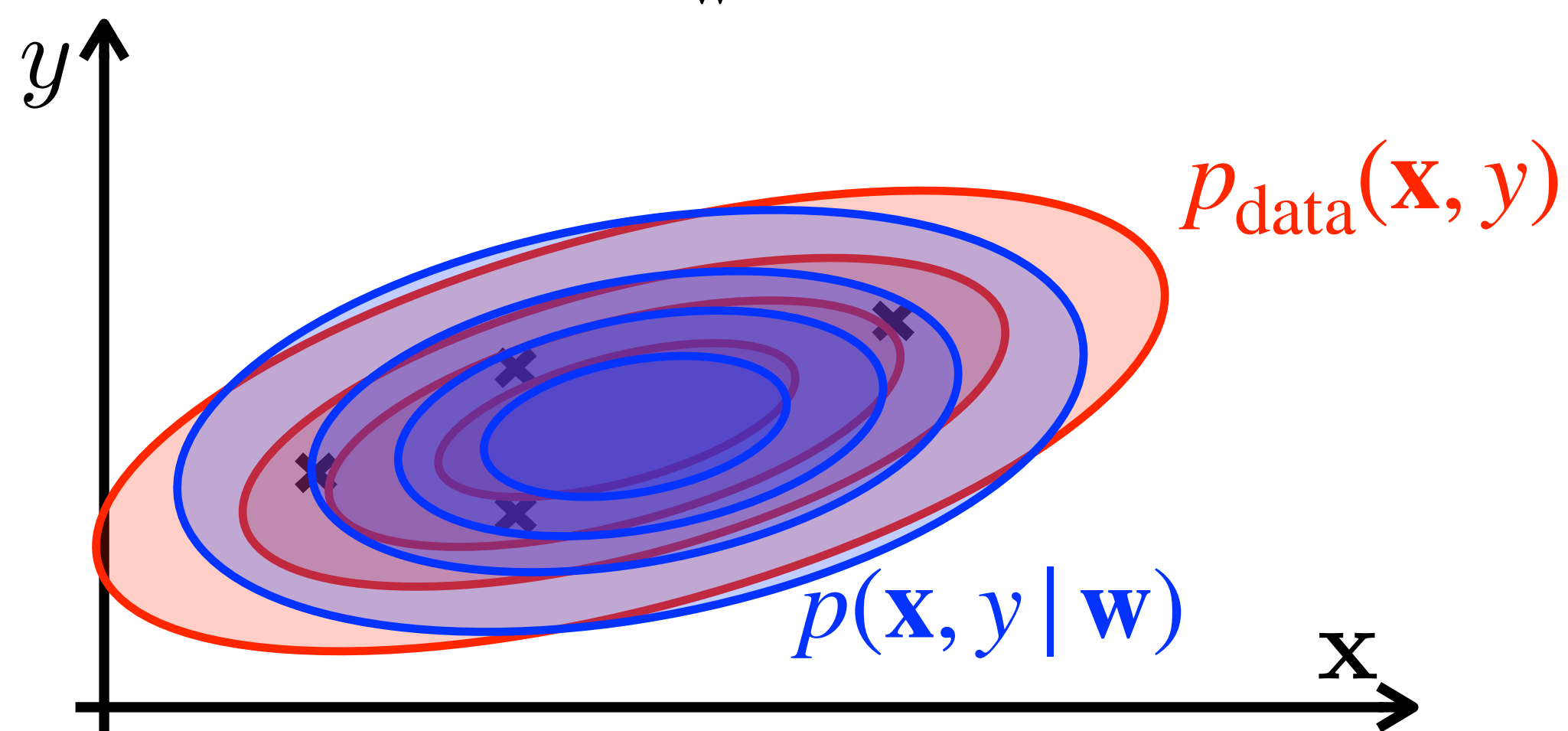- Dataset augmentation
- Realistic simulator
- Generative networks

# How can I do it?

$$\mathbf{w}^\star = \arg\min_{\mathbf{w}} D_{KL}\big(p_{\text{data}}(\mathbf{x}, y) \parallel p(\mathbf{x}, y \mid \mathbf{w})\big) = \int_{(\mathbf{x}, y)} p_{\text{data}}(\mathbf{x}, y) \cdot \log \frac{p_{\text{data}}(\mathbf{x}, y)}{p(\mathbf{x}, y \mid \mathbf{w})}$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}}\left[\log \frac{p_{data}(\mathbf{x}, y)}{p(\mathbf{x}, y \mid \mathbf{w})}\right] = \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}}\left[-\log p(\mathbf{x}, y \mid \mathbf{w})\right]$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}}\left[\log p_{\text{data}}(\mathbf{x}, y) - \log p(y \mid \mathbf{x}, \mathbf{w}) p(\mathbf{x})\right]$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x}, y) \sim p_{\text{data}}}\left[-\log p(y \mid \mathbf{x}, \mathbf{w})\right] \approx \arg\min_{\mathbf{w}} \frac{1}{N} \sum_{(\mathbf{x}_i, y_i) \sim p_{\text{data}}(\mathbf{x}, y)} \left[-\log p(y_i \mid \mathbf{x}_i, \mathbf{w})\right]$$



$p_{\text{data}}(\mathbf{x}, y)$

$p(\mathbf{x}, y \mid \mathbf{w})$

# How can I do it?

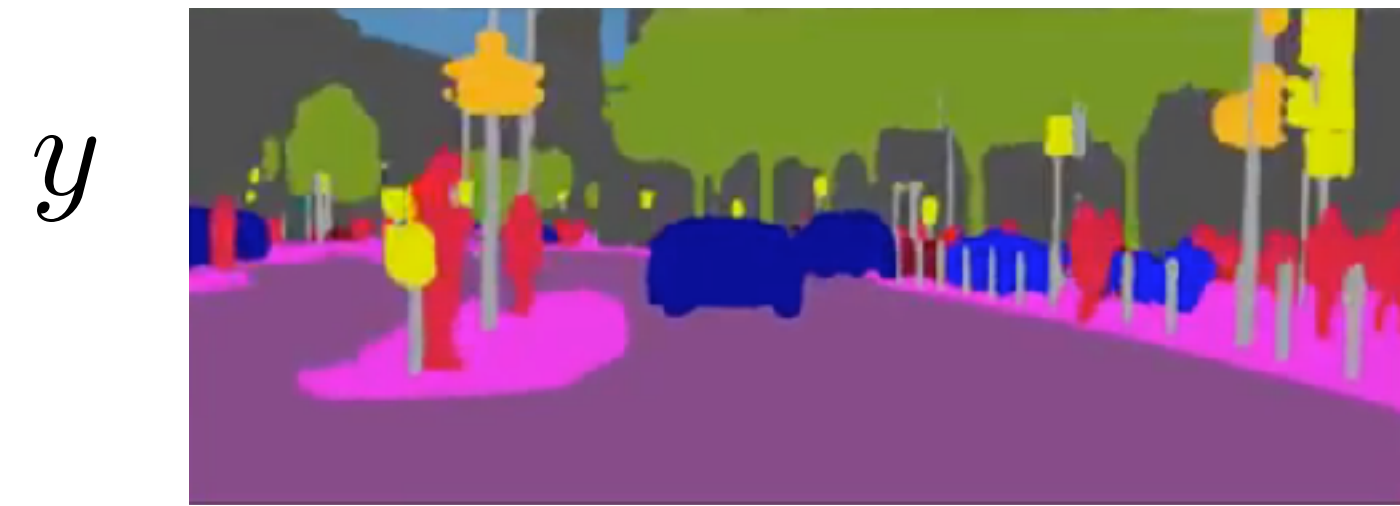$$\mathbf{w}^\star = \arg\min_{\mathbf{w}} D_{KL}\big(p_{\text{data}}(\mathbf{x}, y) \parallel p(\mathbf{x}, y \,|\, \mathbf{w})\big) = \int_{(\mathbf{x},y)} p_{\text{data}}(\mathbf{x}, y) \cdot \log \frac{p_{\text{data}}(\mathbf{x}, y)}{p(\mathbf{x}, y \,|\, \mathbf{w})}$$

$$= \arg\min_{\mathbf{w}} \ \mathbb{E}_{(\mathbf{x},y)\sim p_{\text{data}}}\left[\log \frac{p_{data}(\mathbf{x}, y)}{p(\mathbf{x}, y \,|\, \mathbf{w})}\right] \ = \arg\min_{\mathbf{w}} \ \mathbb{E}_{(\mathbf{x},y)\sim p_{\text{data}}}\left[-\log p(\mathbf{x}, y \,|\, \mathbf{w})\right]$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x},y)\sim p_{\text{data}}}\left[\log p_{\text{data}}(\mathbf{x}, y) - \log p(\mathbf{x} \,|\, y, \mathbf{w})p(y)\right]$$

$$= \arg\min_{\mathbf{w}} \mathbb{E}_{(\mathbf{x},y)\sim p_{\text{data}}}\left[-\log p(\mathbf{x} \,|\, y, \mathbf{w})\right] \approx \arg\min_{\mathbf{w}} \frac{1}{N} \sum_{(\mathbf{x}_i, y_i) \sim p_{\text{data}}(\mathbf{x}, y)} \left[-\log p(\mathbf{x}_i \,|\, y_i, \mathbf{w})\right]$$
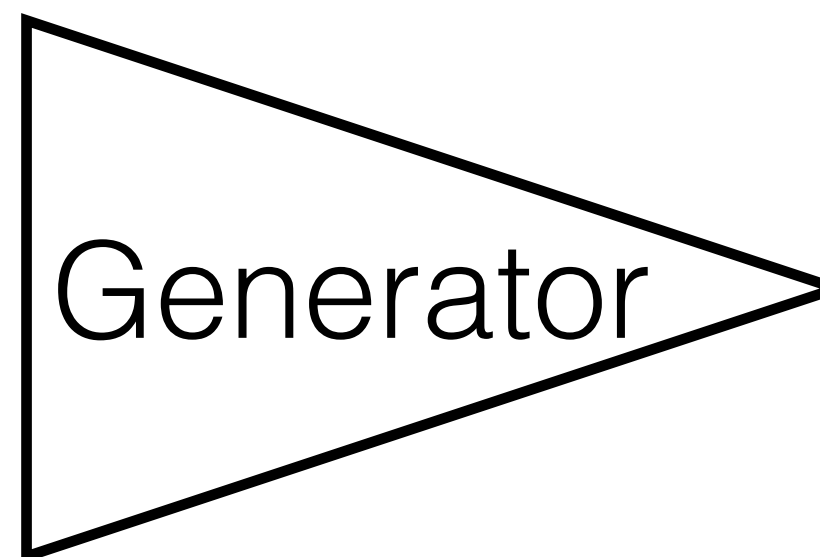


$p_{\text{data}}(\mathbf{x}, y)$

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

Dataset augmentation

Samples from $\hat{p}_{\mathrm{data}}(\mathbf{x}|y)$

annotation

$y$

photo

$\mathbf{x}$

$+$ noise

Generator

$\approx \arg\min_{\mathbf{w}} \dfrac{1}{N} \displaystyle\sum_{(\mathbf{x}_i, y_i)} \left[ -\log p(\mathbf{x}_i|y_i, \mathbf{w}) \right]$
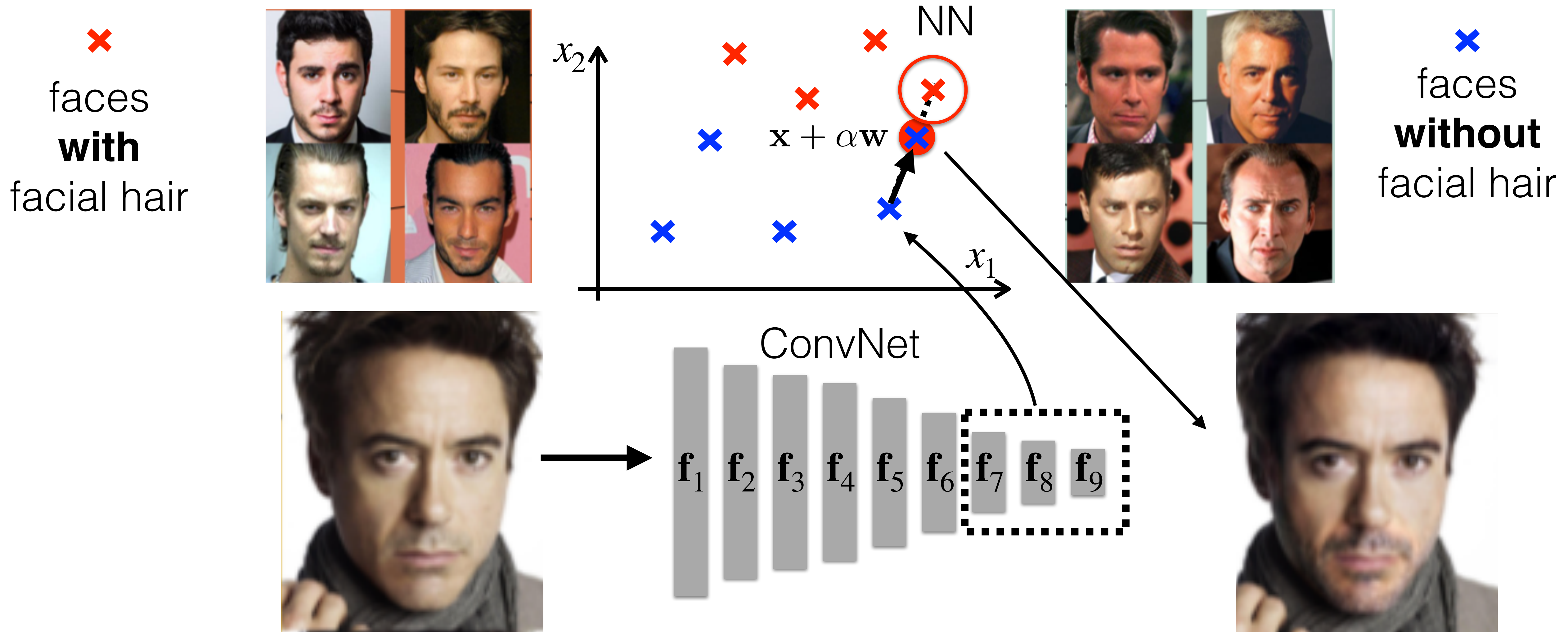
Is it that easy?

Photos corresponding to annotation comes from
- high-dimensional,
- intricate
- non-gaussian pdf
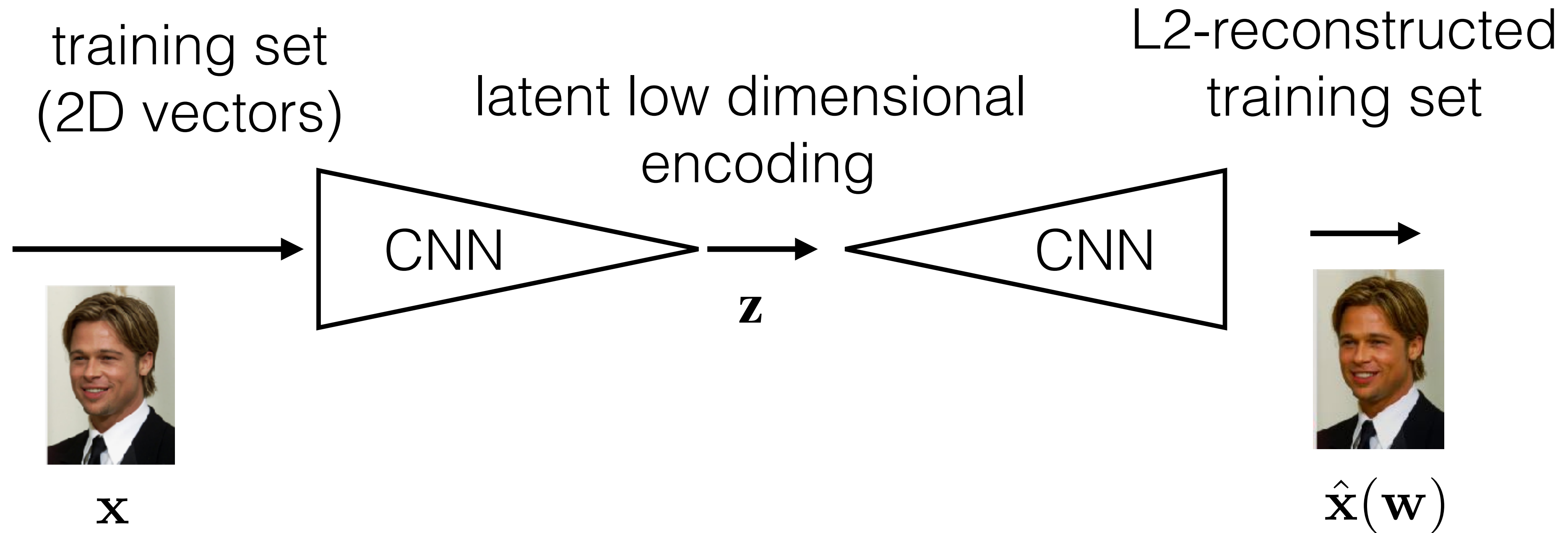
L2-loss is obviously wrong =>Where do I get the loss?

# Deep Feature interpolations [Upchurch CVPR 2017]
## https://arxiv.org/pdf/1611.05507.pdf



faces **with** facial hair

NN

$x_2$

$\mathbf{x} + \alpha\mathbf{w}$

$x_1$

faces **without** facial hair

ConvNet

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

Generate samples by projecting images into low-dimensional representation

# Generative models

You can do it even without any annotations just on collection of images

training set
(2D vectors)

latent low dimensional
encoding

L2-reconstructed
training set



CNN → $\mathbf{z}$ → CNN

$\mathbf{x}$

$\hat{\mathbf{x}}(\mathbf{w})$

- Learning the self-reconstruction with L2 reconstruction loss

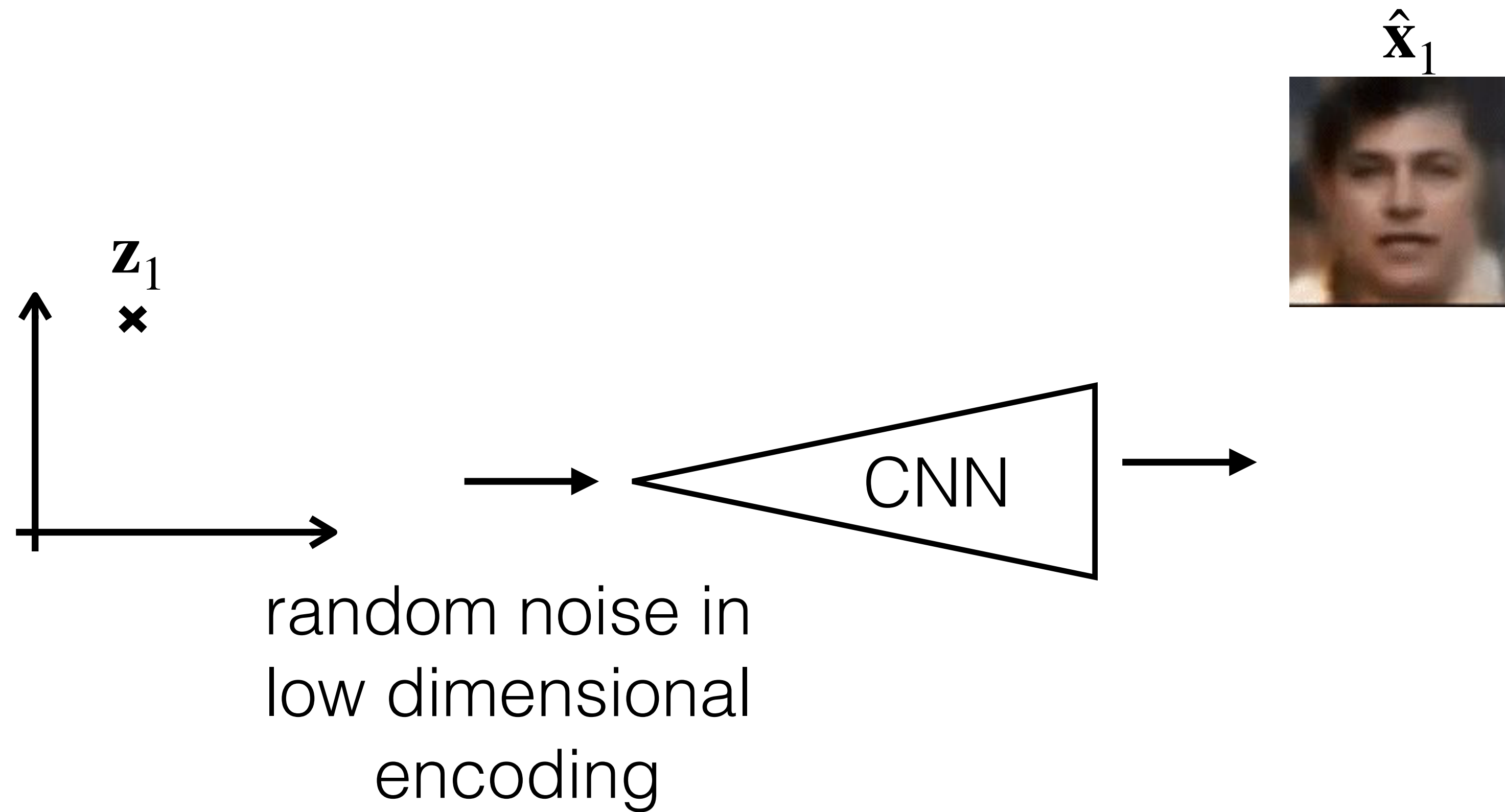$$\arg \min_{\mathbf{w}} \|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{w})\|_2^2$$

If CNN is pure linear function then:
- closed-form solution exists
- method is called PCA

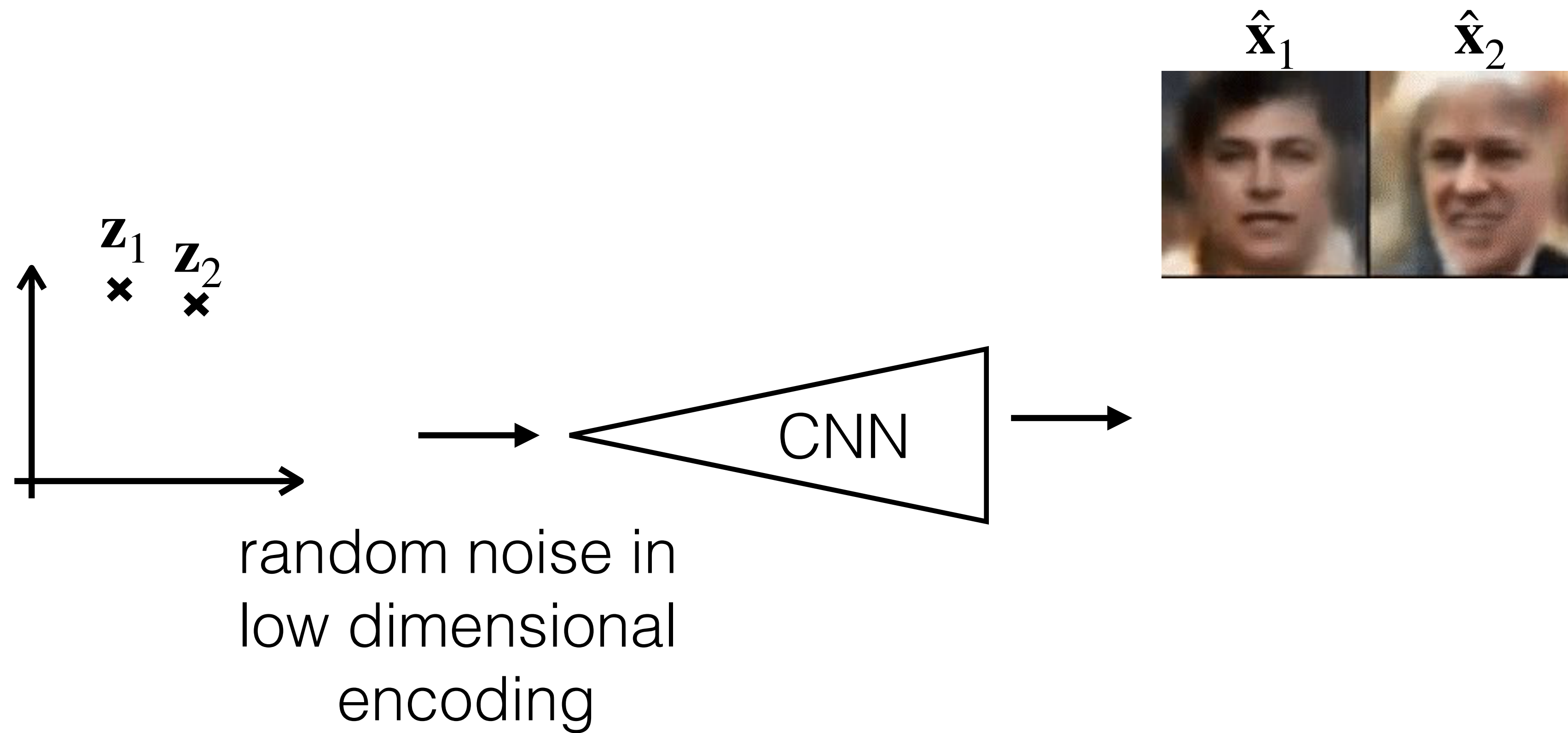If $\mathbf{z}$ is pushed towards gaussian distribution:
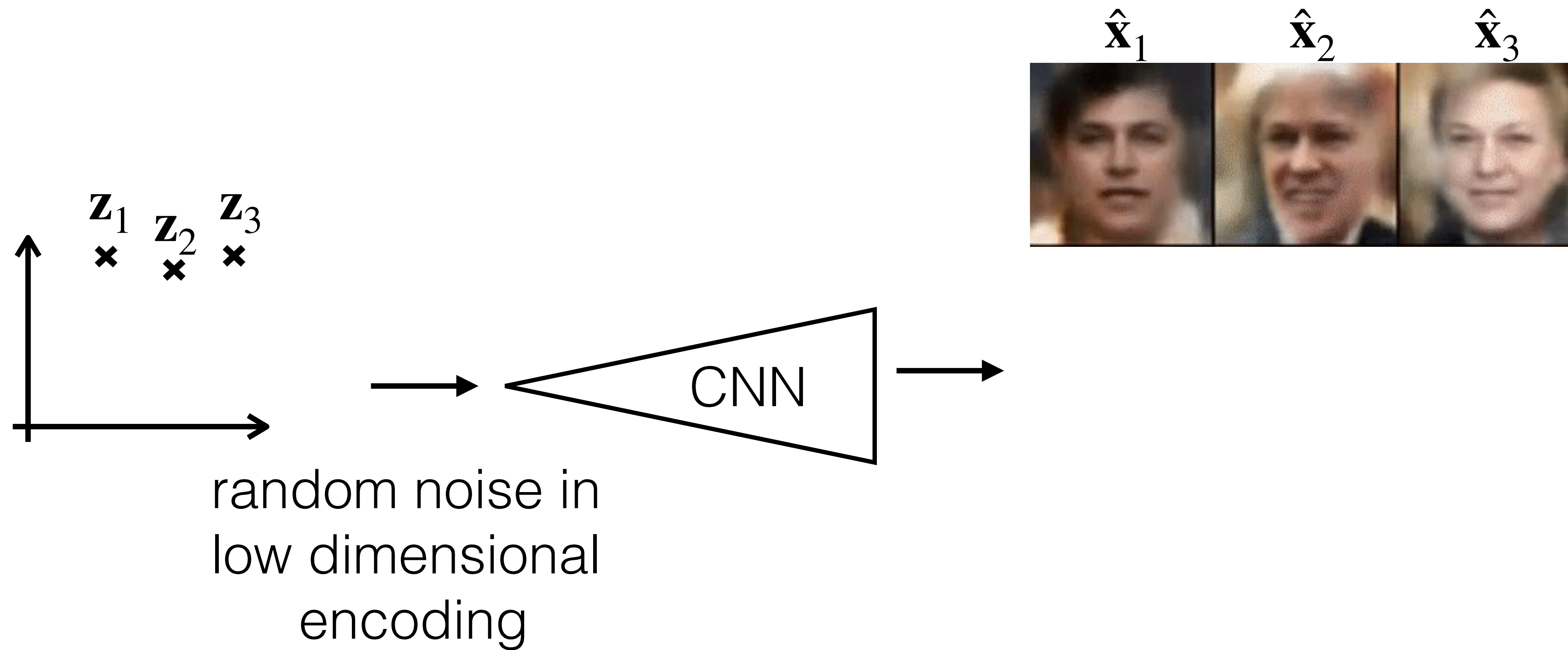- method is referred as variational encoders

# Generative models

$\hat{\mathbf{x}}_1$



$\mathbf{z}_1$

CNN

random noise in
low dimensional
encoding

- New samples generated from random vectors in low-dimensional encoding.

# Generative models



$\hat{\mathbf{x}}_1$     $\hat{\mathbf{x}}_2$

$\mathbf{z}_1$   $\mathbf{z}_2$

CNN

random noise in
low dimensional
encoding

- New samples generated from random vectors in low-dimensional encoding.

# Generative models

$\hat{\mathbf{x}}_1 \qquad \hat{\mathbf{x}}_2 \qquad \hat{\mathbf{x}}_3$



$\mathbf{z}_1 \quad \mathbf{z}_2 \quad \mathbf{z}_3$

CNN

random noise in
low dimensional
encoding

- New samples generated from random vectors in low-dimensional encoding.

# Generative models

$\hat{\mathbf{x}}_1$ $\hat{\mathbf{x}}_2$ $\hat{\mathbf{x}}_3$

$\mathbf{z}_1$ $\mathbf{z}_2$ $\mathbf{z}_3$

CNN

random noise in
low dimensional
encoding
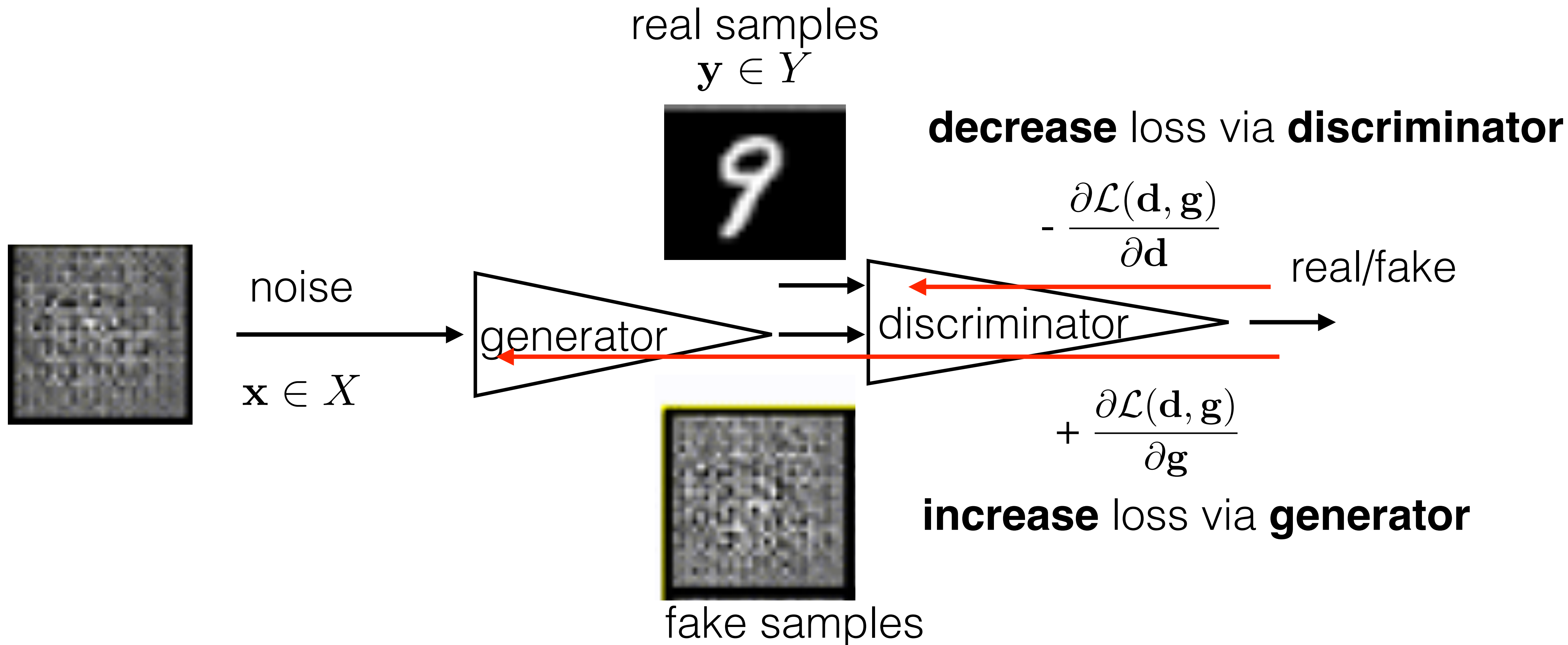
- New samples generated from random vectors in low-dimensional encoding.
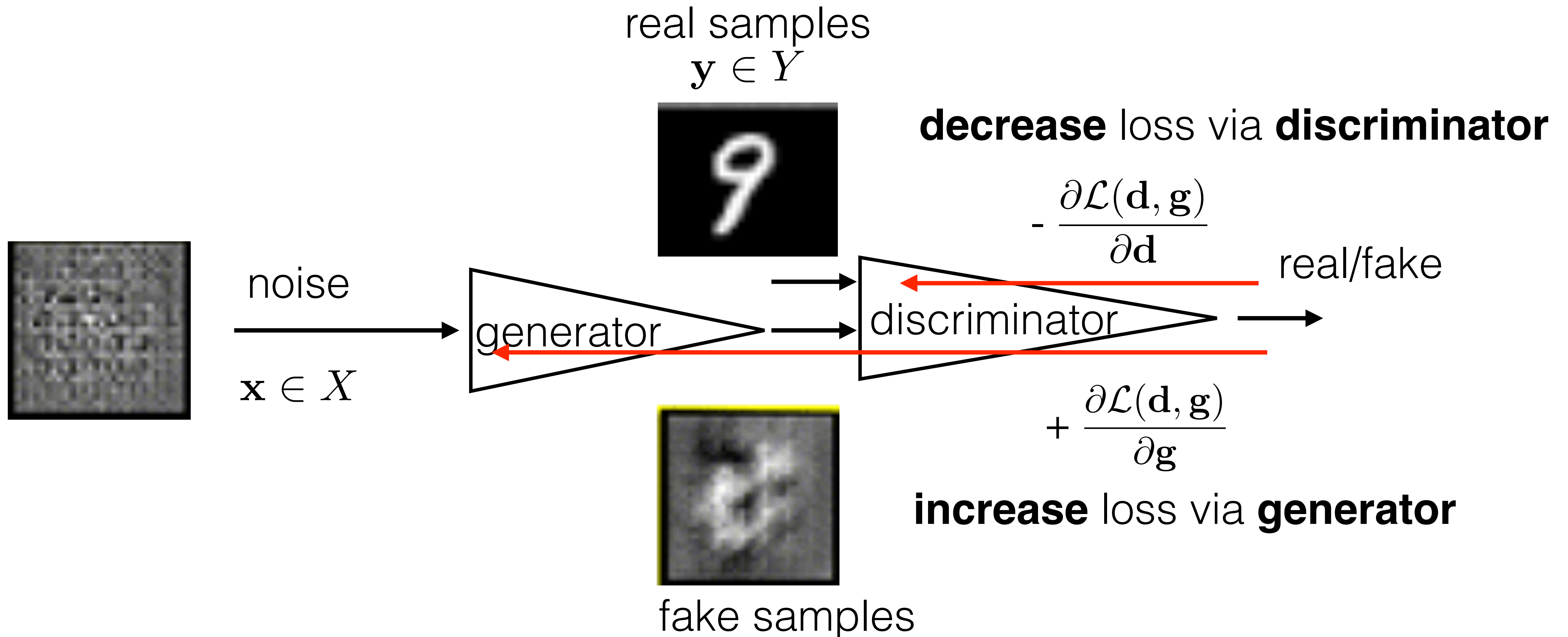
# Generative Adversarial Nets [Goodfellow NIPS 2014]

real samples
$\mathbf{y} \in Y$

**decrease** loss via **discriminator**

$- \dfrac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{d}}$

real/fake

noise

generator

discriminator

$\mathbf{x} \in X$

$+ \dfrac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{g}}$

**increase** loss via **generator**

fake samples

classification loss: $\mathcal{L}(\mathbf{d}, \mathbf{g}) = \sum\limits_{\mathbf{x} \in X} -\log(\mathbf{d}(\mathbf{g}(\mathbf{x}))) + \sum\limits_{\mathbf{y} \in Y} -\log(1 - \mathbf{d}(\mathbf{y}))$
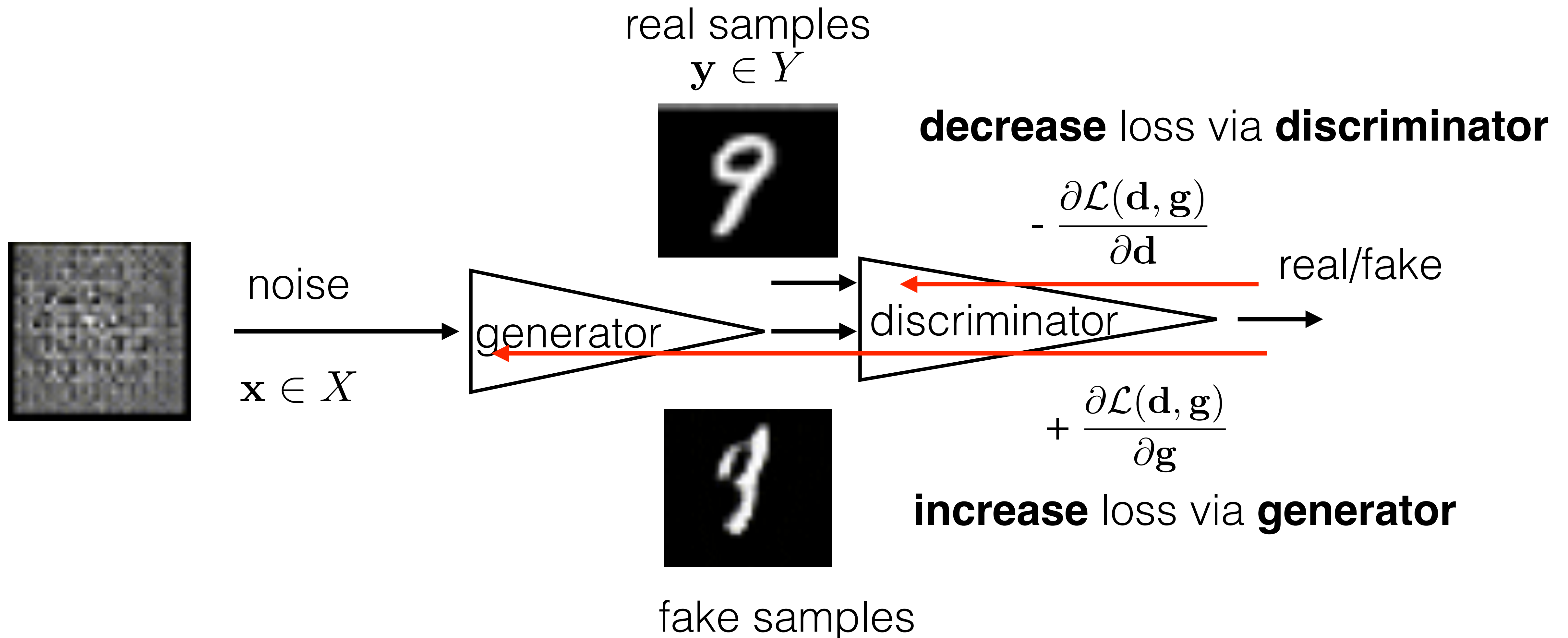
# Generative Adversarial Nets [Goodfellow NIPS 2014]

real samples
$\mathbf{y} \in Y$

**decrease** loss via **discriminator**



$- \dfrac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{d}}$

real/fake

noise

generator

discriminator

$\mathbf{x} \in X$

$+ \dfrac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{g}}$

**increase** loss via **generator**

fake samples

classification loss: $\mathcal{L}(\mathbf{d}, \mathbf{g}) = \sum\limits_{\mathbf{x} \in X} -\log(\mathbf{d}(\mathbf{g}(\mathbf{x}))) + \sum\limits_{\mathbf{y} \in Y} -\log(1 - \mathbf{d}(\mathbf{y}))$
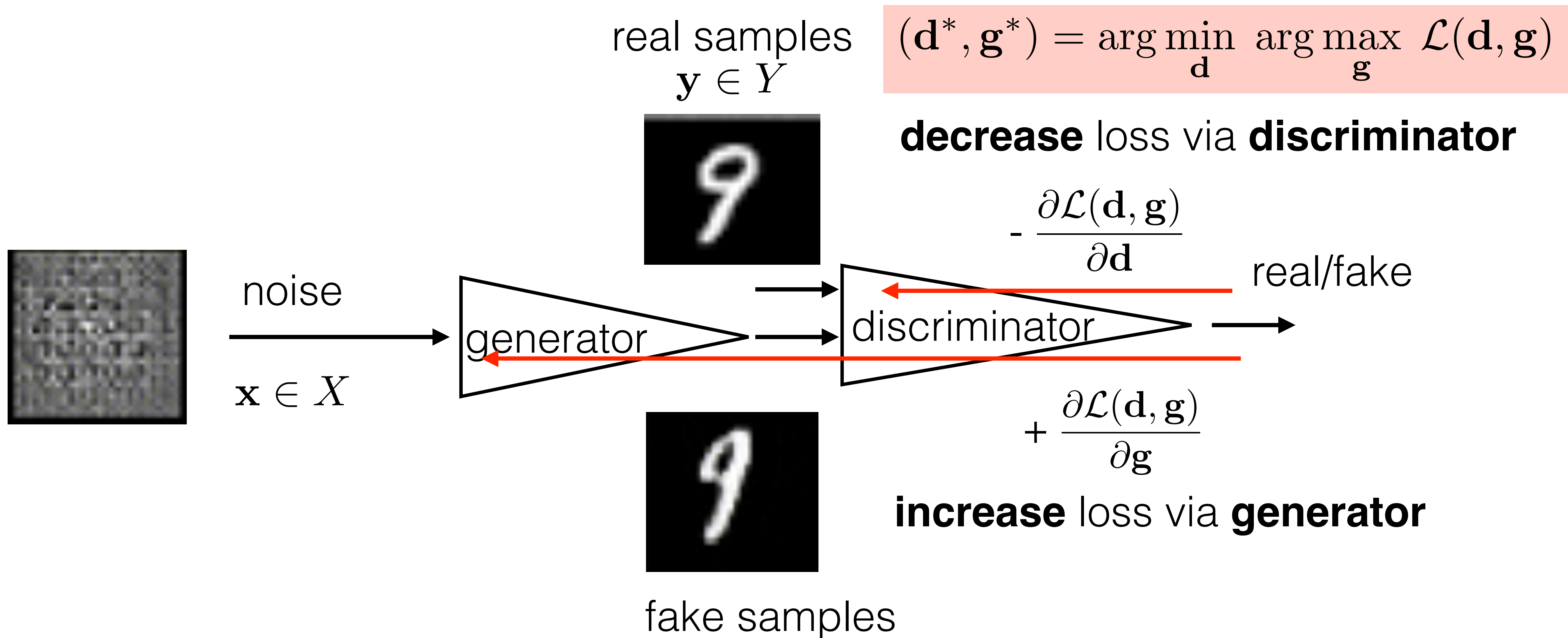
# Generative Adversarial Nets [Goodfellow NIPS 2014]
## https://arxiv.org/abs/1406.2661

real samples
$$\mathbf{y} \in Y$$



**decrease** loss via **discriminator**

$$- \frac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{d}}$$

real/fake

noise

$$\mathbf{x} \in X$$

generator

discriminator

$$+ \frac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{g}}$$

**increase** loss via **generator**

fake samples

classification loss: $\mathcal{L}(\mathbf{d}, \mathbf{g}) = \sum_{\mathbf{x} \in X} - \log(\mathbf{d}(\mathbf{g}(\mathbf{x}))) + \sum_{\mathbf{y} \in Y} - \log(1 - \mathbf{d}(\mathbf{y}))$

# Generative Adversarial Nets [Goodfellow NIPS 2014]

real samples
$$\mathbf{y} \in Y$$

$$(\mathbf{d}^*, \mathbf{g}^*) = \arg\min_{\mathbf{d}} \arg\max_{\mathbf{g}} \mathcal{L}(\mathbf{d}, \mathbf{g})$$

**decrease** loss via **discriminator**
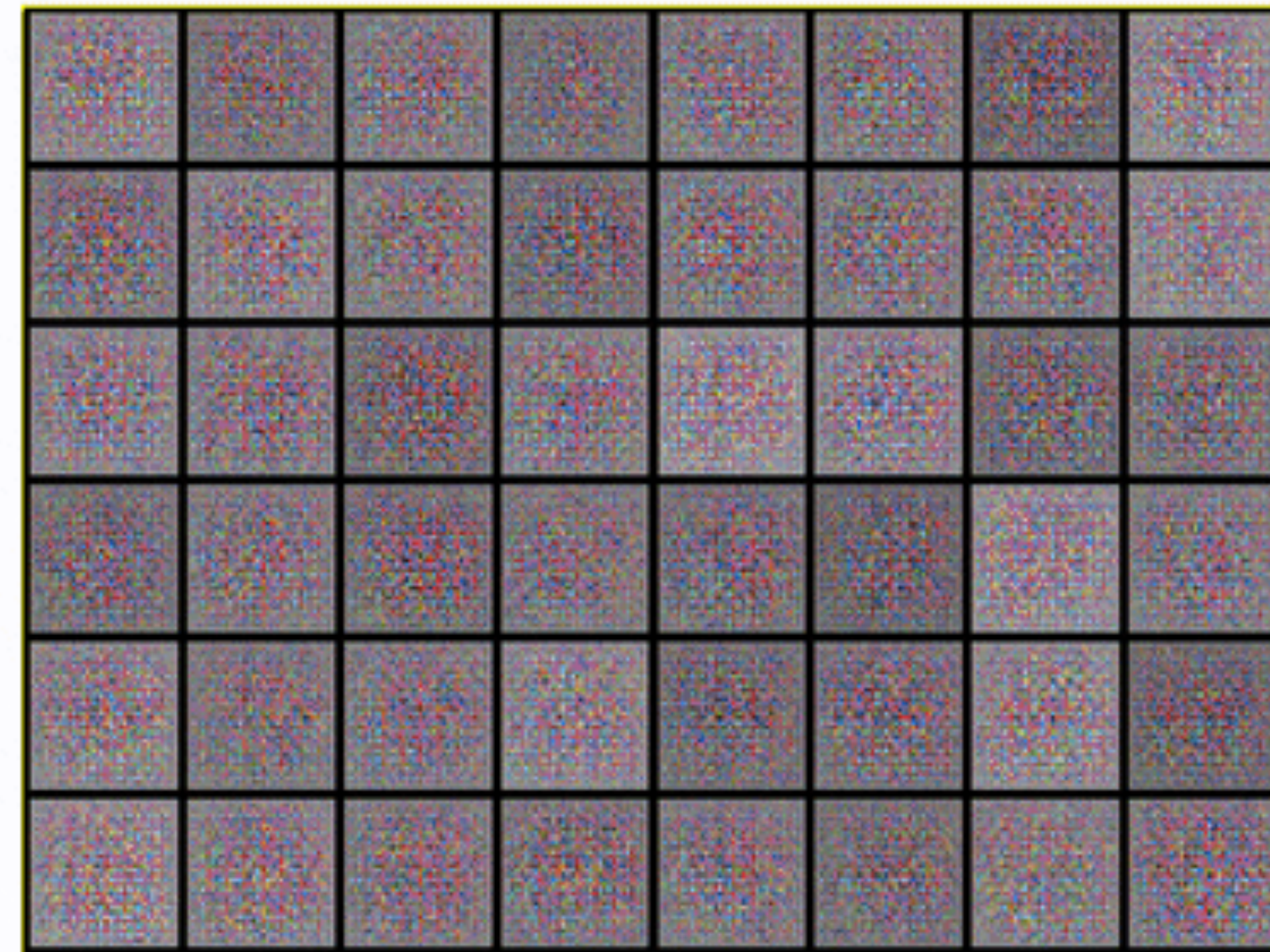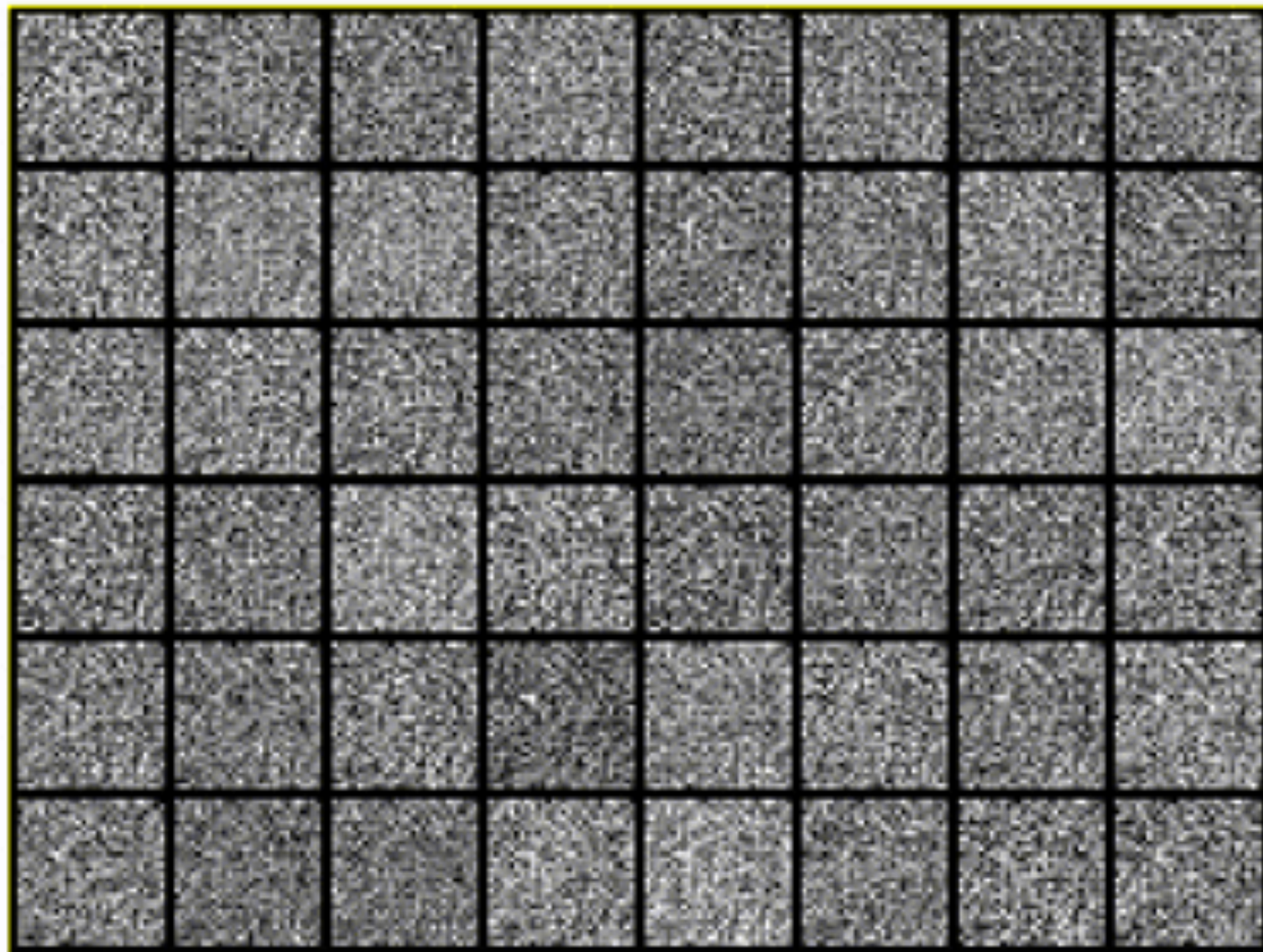
$$- \frac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{d}}$$

noise

generator

discriminator

real/fake

$$\mathbf{x} \in X$$

$$+ \frac{\partial \mathcal{L}(\mathbf{d}, \mathbf{g})}{\partial \mathbf{g}}$$

**increase** loss via **generator**

fake samples

classification loss: $\mathcal{L}(\mathbf{d}, \mathbf{g}) = \sum_{\mathbf{x} \in X} -\log(\mathbf{d}(\mathbf{g}(\mathbf{x}))) + \sum_{\mathbf{y} \in Y} -\log(1 - \mathbf{d}(\mathbf{y}))$

# Generative Avdersarial Nets [Goodfellow NIPS 2014]
## https://arxiv.org/abs/1406.2661

$$(\mathbf{d}^*, \mathbf{g}^*) = \arg\min_{\mathbf{d}} \ \arg\max_{\mathbf{g}} \ \mathcal{L}(\mathbf{d}, \mathbf{g})$$

- Proof: Equilibrium in saddle point implies that generator generates samples from th
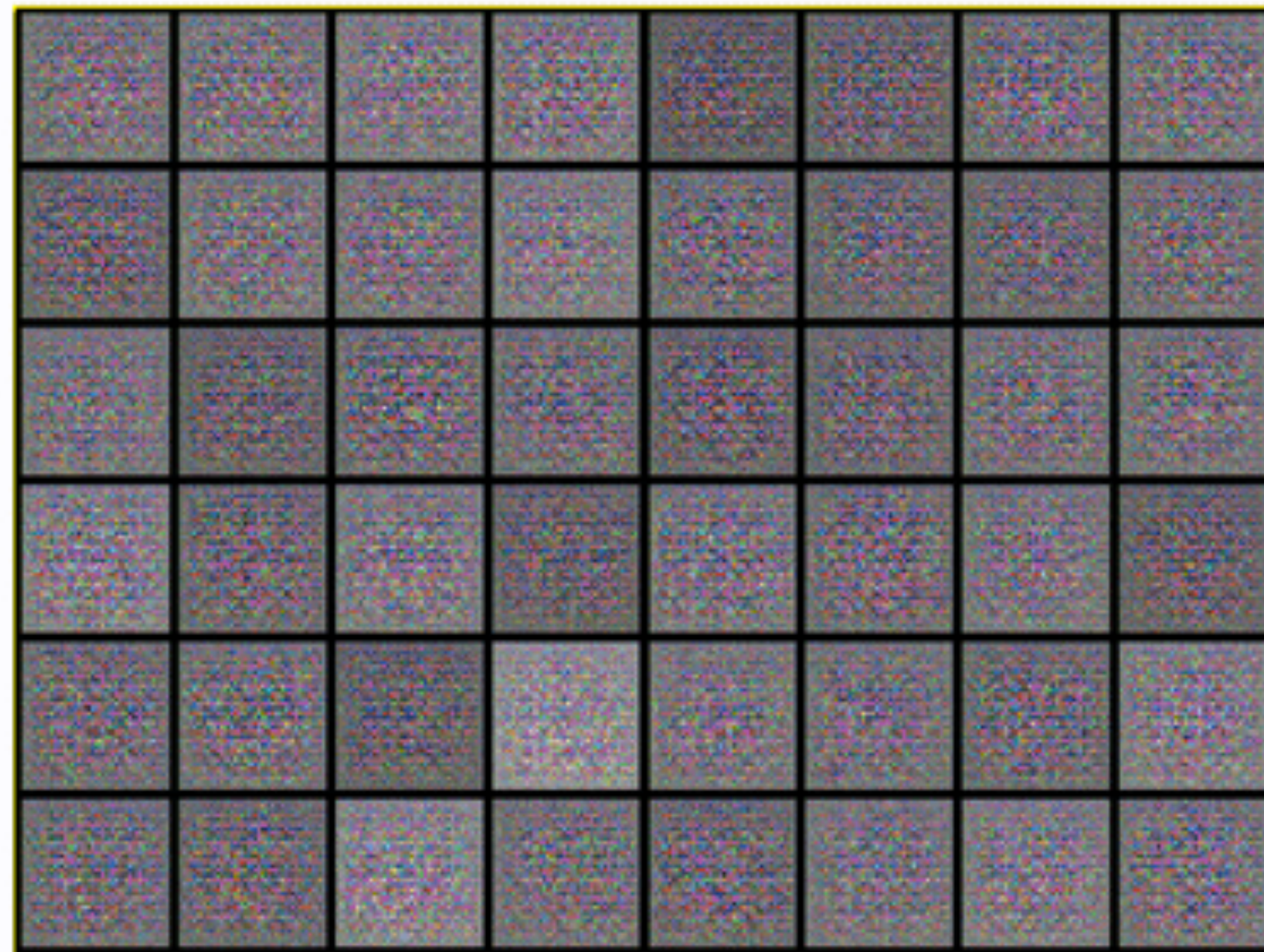  real distribution (asymptotically consistent in contrast to VAE)

# Generative Avdersarial Nets [Goodfellow NIPS 2014]
## https://arxiv.org/abs/1406.2661

$$(\mathbf{d}^*, \mathbf{g}^*) = \arg\min_{\mathbf{d}} \; \arg\max_{\mathbf{g}} \; \mathcal{L}(\mathbf{d}, \mathbf{g})$$

- Proof: Equilibrium in saddle point implies that generator generates samples from th
real distribution (asymptotically consistent in contrast to VAE)



https://theaisummer.com/gan-computer-vision/

# Generative Avdersarial Nets [Goodfellow NIPS 2014]
## https://arxiv.org/abs/1406.2661

$$(\mathbf{d}^*, \mathbf{g}^*) = \arg\min_{\mathbf{d}} \ \arg\max_{\mathbf{g}} \ \mathcal{L}(\mathbf{d}, \mathbf{g})$$

- The learning is generally unstable and suffers from mode collapse



https://theaisummer.com/gan-computer-vision/

# Domain transfer

Cycle-GAN [Zhu ICCV 2017]
https://arxiv.org/abs/1703.10593

real Monet's paining

real Y-domain
samples $\mathbf{y} \in Y$

real/fake

$\mathcal{L}(\mathbf{d_y}, \mathbf{g_x})$

$\mathbf{g_x(x)}$

$\mathbf{d_y(y)}$

real X-domain
samples $\mathbf{x} \in X$

fake Y-domain
samples $\hat{\mathbf{y}}$

fake Monet's paining

Cycle-GAN [Zhu ICCV 2017]
https://arxiv.org/abs/1703.10593



real Monet's paining

real Y-domain
samples $\mathbf{y} \in Y$

real/fake

$\mathcal{L}(\mathbf{d_y}, \mathbf{g_x})$

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

$\mathbf{d_y}(\mathbf{y})$

real X-domain
samples $\mathbf{x} \in$

...domain
s $\hat{\mathbf{y}}$

real/fake

$\mathcal{L}(\mathbf{d_x}, \mathbf{g_y})$

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

$\mathbf{g_y}(\mathbf{y})$

real Y-domain
samples $\mathbf{y} \in Y$

real Monet's
paining

real X-domain
samples $\mathbf{x} \in X$

Cycle-GAN [Zhu ICCV 2017]
https://arxiv.org/abs/1703.10593

$$|\mathbf{g_y}(\mathbf{g_x}(\mathbf{x}) - \hat{\mathbf{x}}|$$



real Y-domain
samples $\mathbf{y} \in Y$

real/fake

$\mathbf{d_y}(\mathbf{y})$

$\mathcal{L}(\mathbf{d_y}, \mathbf{g_x})$

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

real X-domain
samples $\mathbf{x} \in X$

fake Y-domain
samples $\hat{\mathbf{y}}$

fake X-domain
samples $\hat{\mathbf{x}}$

$\mathbf{g_y}(\mathbf{y})$

real/fake

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

real Y-domain
samples $\mathbf{y} \in Y$

$\mathcal{L}(\mathbf{d_x}, \mathbf{g_y})$

real X-domain
samples $\mathbf{x} \in X$

# Cycle-GAN [Zhu ICCV 2017]
https://arxiv.org/abs/1703.10593

$$\mathcal{L}_{GAN}(\mathbf{d_x}, \mathbf{d_y}, \mathbf{g_x}, \mathbf{g_y}) = \mathcal{L}(\mathbf{d_x}, \mathbf{g_y}) + \mathcal{L}(\mathbf{d_y}, \mathbf{g_x}) + |\mathbf{g_y}(\mathbf{g_x}(\mathbf{x}) - \hat{\mathbf{x}}|$$
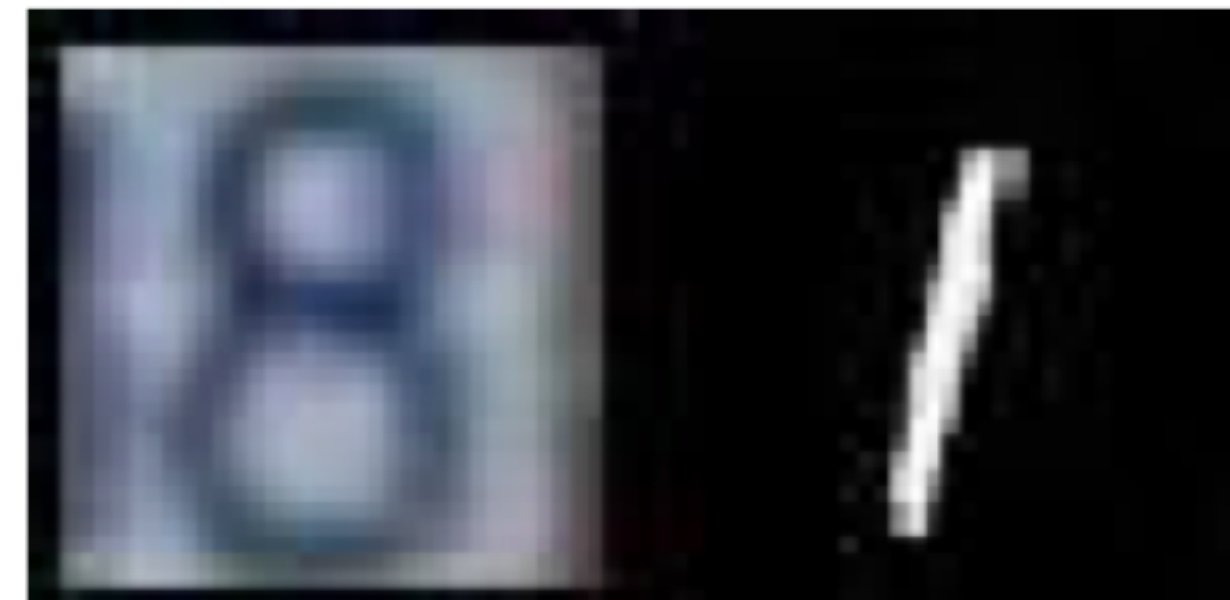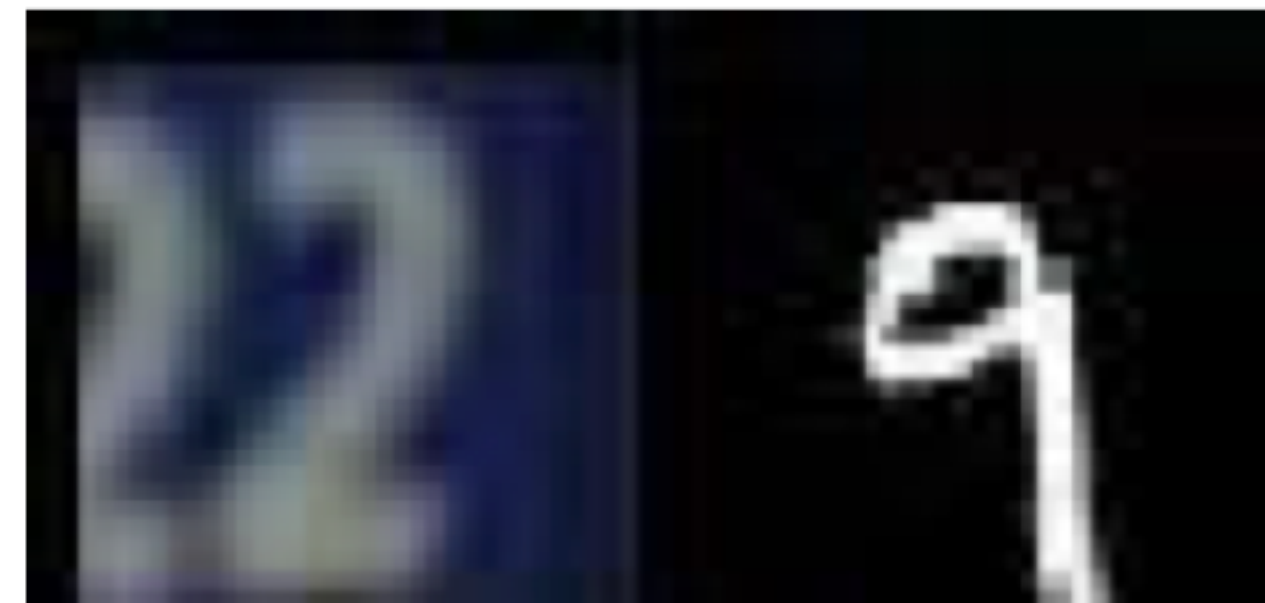


real Y-domain samples $\mathbf{y} \in Y$

real/fake

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

$\mathbf{d_y}(\mathbf{y})$

real X-domain samples $\mathbf{x} \in X$

fake Y-domain samples $\hat{\mathbf{y}}$

fake X-domain samples $\hat{\mathbf{x}}$

$\mathbf{g_y}(\mathbf{y})$

real/fake

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

real Y-domain samples $\mathbf{y} \in Y$

real X-domain samples $\mathbf{x} \in X$

# Cycle-GAN [Zhu ICCV 2017]
## https://arxiv.org/abs/1703.10593

**Monet ⟳ Photos**



Monet ⟶ photo

photo ⟶ Monet

# Cycle-GAN [Zhu ICCV 2017]
## https://arxiv.org/abs/1703.10593



zebra ⟶ horse

horse ⟶ zebra

# Cycle-GAN [Zhu ICCV 2017]
## https://arxiv.org/abs/1703.10593



**Summary ⟳ Winter**

summer ⟶ winter

winter ⟶ summer

CyCaDa-GAN [Hoffman CVPR 2018]
https://arxiv.org/pdf/1711.03213.pdf
House numbers to MNIST transfer

**y**
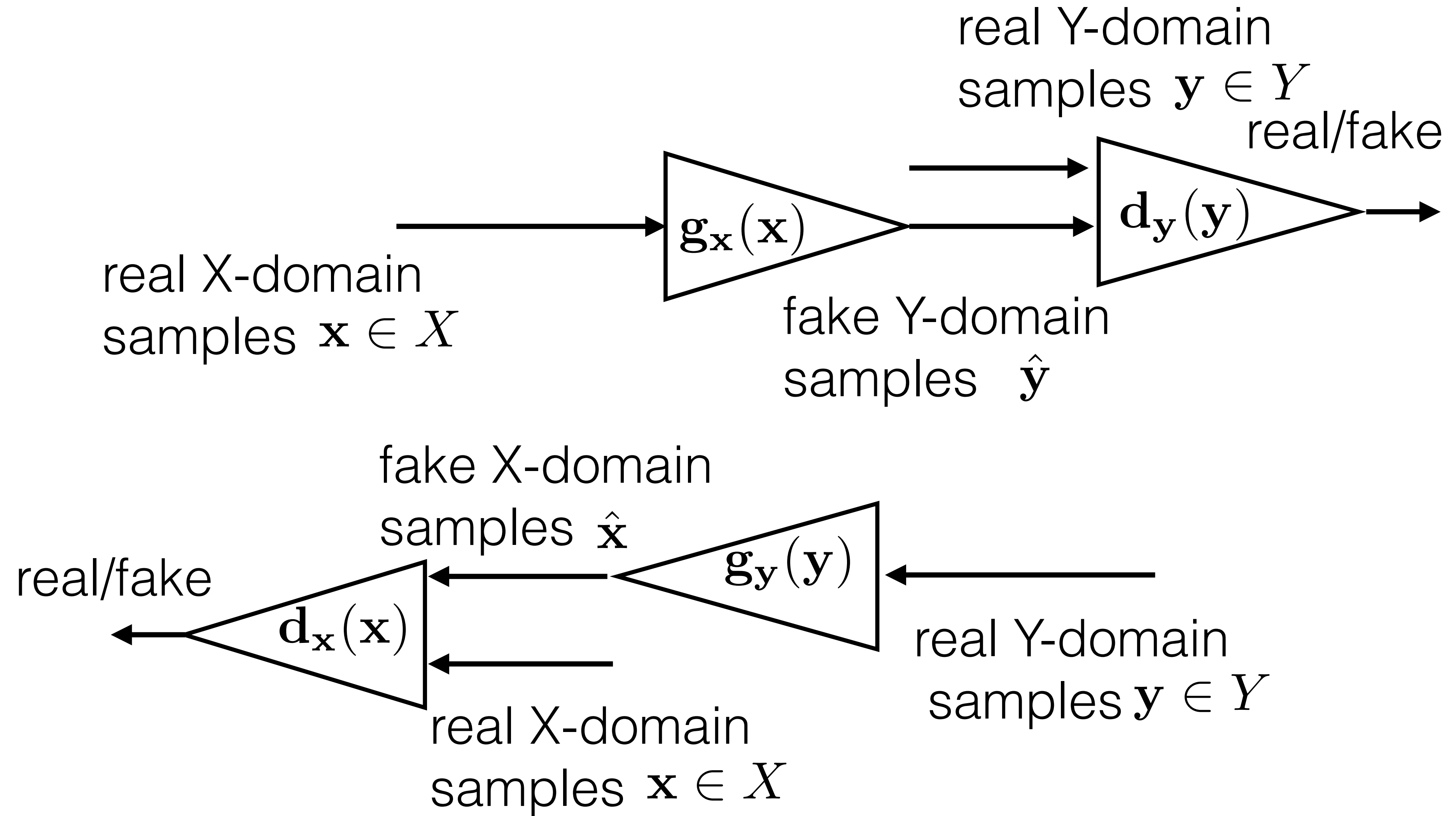


House numbers

$$\mathcal{L}_{GAN}(\mathbf{d_x}, \mathbf{d_y}, \mathbf{g_x}, \mathbf{g_y}) = \mathcal{L}(\mathbf{d_y}, \mathbf{g_x}) + \mathcal{L}(\mathbf{d_x}, \mathbf{g_y}) + \boxed{|\mathbf{g_x}(\mathbf{g_y}(\mathbf{y}) - \hat{\mathbf{y}}|}$$

- Cycle consistency helps, but it still allows to learn totally semantically inconsistent transfer

CyCaDa-GAN [Hoffman CVPR 2018]
https://arxiv.org/pdf/1711.03213.pdf
House numbers to MNIST transfer

$$\mathbf{y} \qquad \mathbf{g_y(y)}$$



House numbers        MNIST

$$\mathcal{L}_{GAN}(\mathbf{d_x}, \mathbf{d_y}, \mathbf{g_x}, \mathbf{g_y}) = \mathcal{L}(\mathbf{d_y}, \mathbf{g_x}) + \mathcal{L}(\mathbf{d_x}, \mathbf{g_y}) + \boxed{|\mathbf{g_x}(\mathbf{g_y(y)} - \hat{\mathbf{y}}|}$$

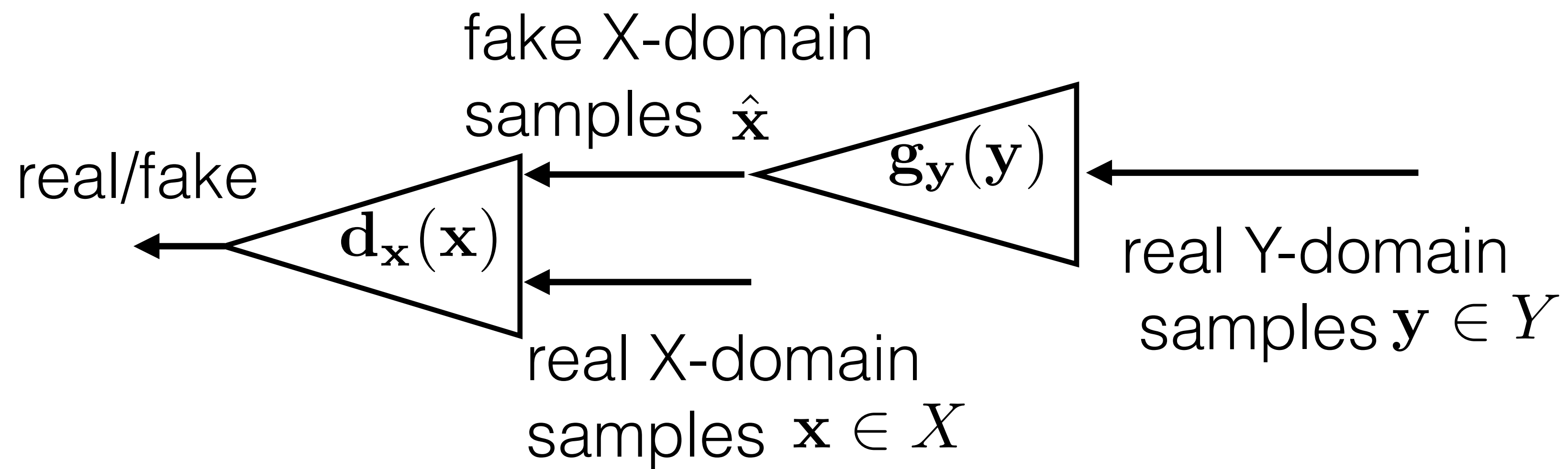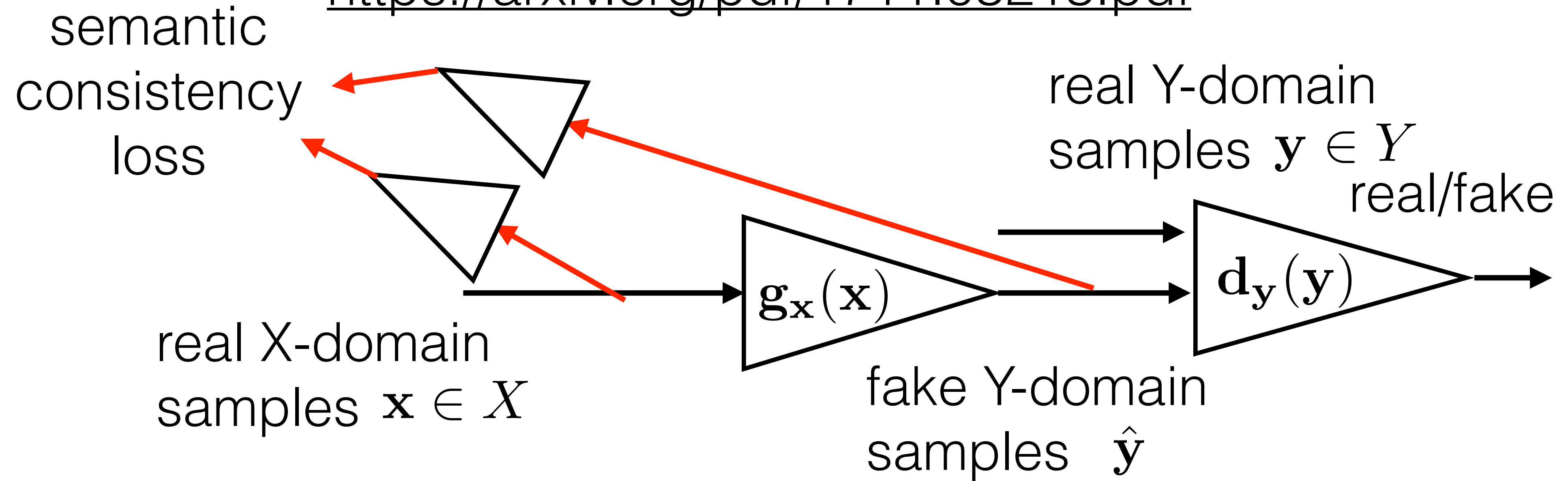- Cycle consistency helps, but it still allows to learn totally semantically inconsistent transfer

CyCaDa-GAN [Hoffman CVPR 2018]
https://arxiv.org/pdf/1711.03213.pdf
House numbers to MNIST transfer

$\mathbf{y}$  $\mathbf{g_y(y)}$  $\mathbf{g_x(g_y(y)}$



House numbers    MNIST    House numbers

$$\mathcal{L}_{GAN}(\mathbf{d_x}, \mathbf{d_y}, \mathbf{g_x}, \mathbf{g_y}) = \mathcal{L}(\mathbf{d_y}, \mathbf{g_x}) + \mathcal{L}(\mathbf{d_x}, \mathbf{g_y}) + |\mathbf{g_x}(\mathbf{g_y}(\mathbf{y}) - \hat{\mathbf{y}}|$$

- Cycle consistency helps, but it still allows to learn totally semantically inconsistent transfer
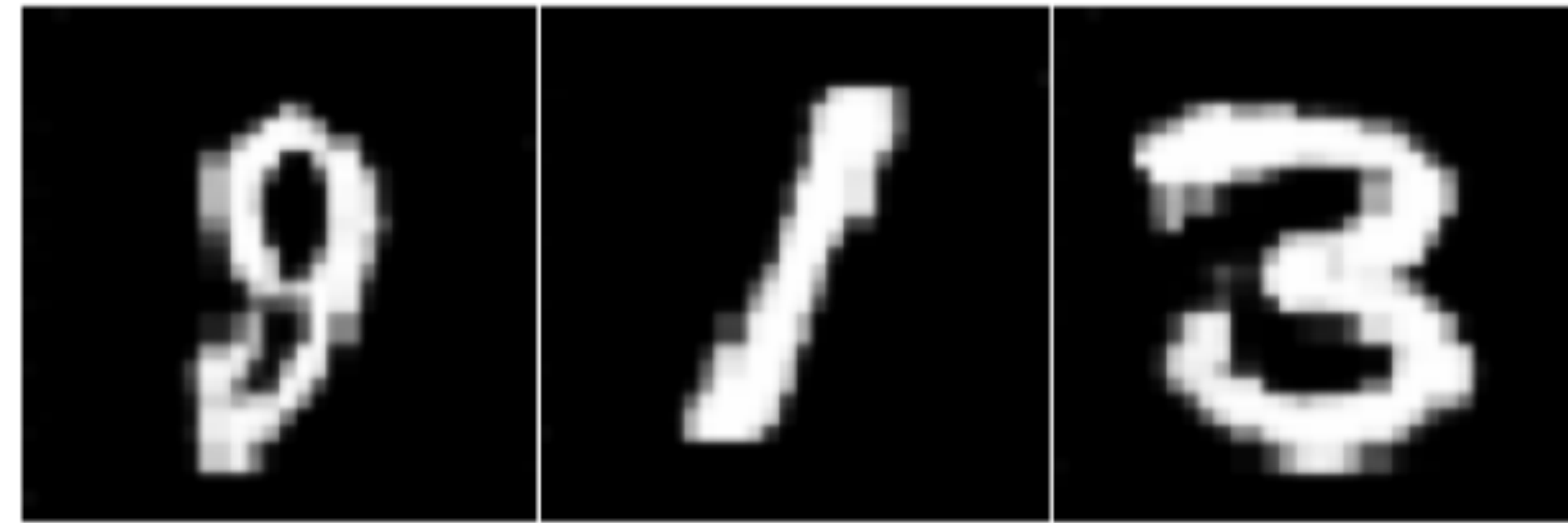
# CyCaDa-GAN [Hoffman CVPR 2018]
## https://arxiv.org/pdf/1711.03213.pdf



real Y-domain samples $\mathbf{y} \in Y$

real/fake

$\mathbf{g_x(x)}$

$\mathbf{d_y(y)}$

real X-domain samples $\mathbf{x} \in X$

fake Y-domain samples $\hat{\boldsymbol{y}}$

fake X-domain samples $\hat{\mathbf{x}}$

$\mathbf{g_y(y)}$

real/fake

$\mathbf{d_x(x)}$

real Y-domain samples $\mathbf{y} \in Y$

real X-domain samples $\mathbf{x} \in X$

CyCaDa-GAN [Hoffman CVPR 2018]
https://arxiv.org/pdf/1711.03213.pdf

semantic
consistency
loss

real Y-domain
samples $\mathbf{y} \in Y$

real/fake

$\mathbf{d_y}(\mathbf{y})$

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

real X-domain
samples $\mathbf{x} \in X$

fake Y-domain
samples $\hat{\mathbf{y}}$

fake X-domain
samples $\hat{\mathbf{x}}$

$\mathbf{g_y}(\mathbf{y})$

real/fake

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

real Y-domain
samples $\mathbf{y} \in Y$

real X-domain
samples $\mathbf{x} \in X$

# CyCaDa-GAN [Hoffman CVPR 2018]
https://arxiv.org/pdf/1711.03213.pdf

- Semantic consistency enforce transformation to be semantically consistent



house numbers



MNIST

# Diffusion models

# Diffusion models
[Hu NeuriPS 2020]

Reverse of the diffusion process to generate original data from the noise



$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

$$q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

Markov chain of diffusion steps in which we slowly and randomly add noise t

If noise is small the backward step has also "almost" gaussian distribution

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

=> learn de-noising networks through L2-norm

# Diffusion networks

# Super resolution

input

output

# Inpainting



| Input | Denoising 0% | Denoising 60% | Denoising 75% | Sample 1 | Sample 2 | Sample 3 | Sample 4 | Sample 5 |

[Repaint, CVPR 2022]

# Inpainting



[Repaint, CVPR 2022]
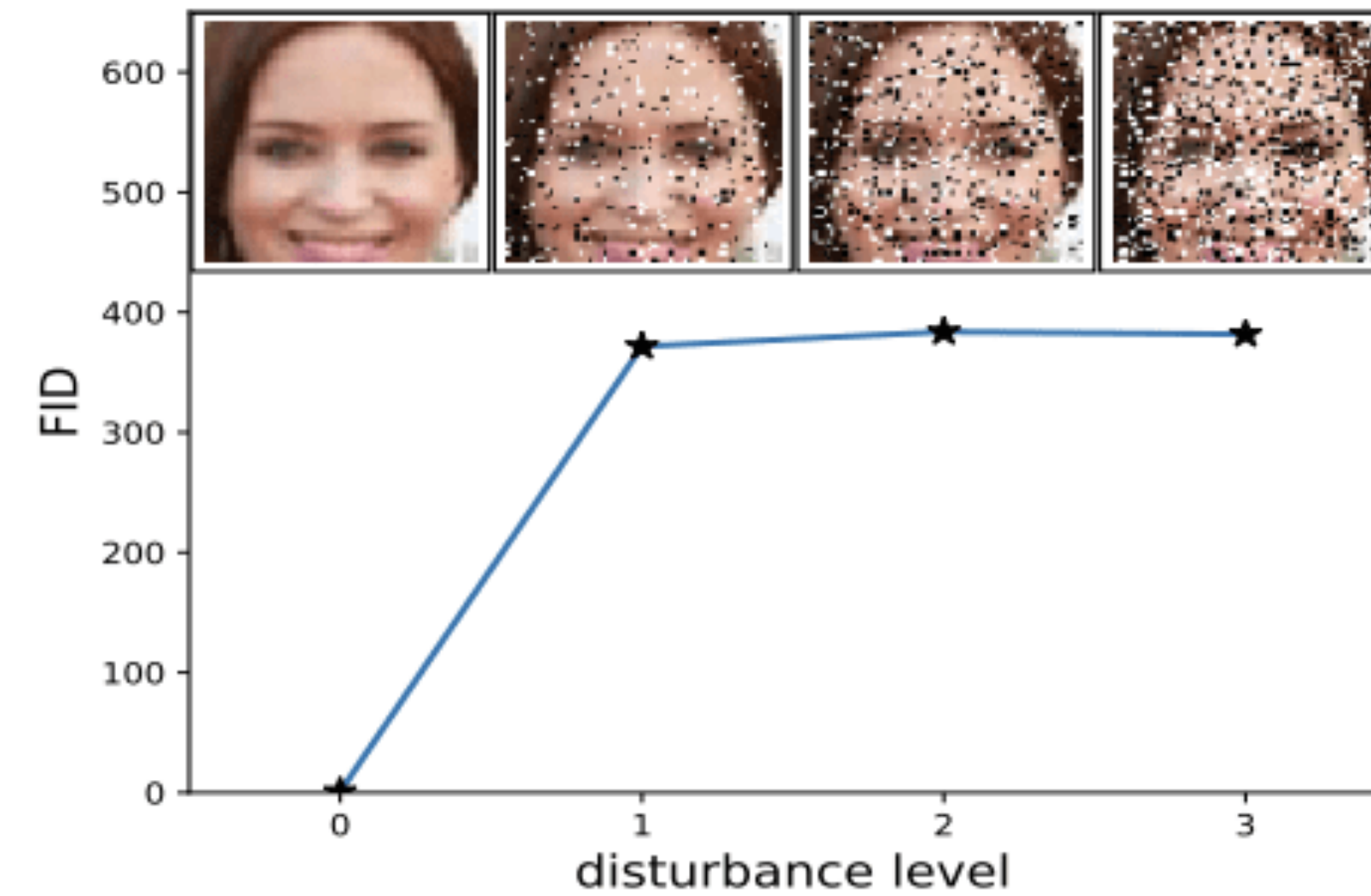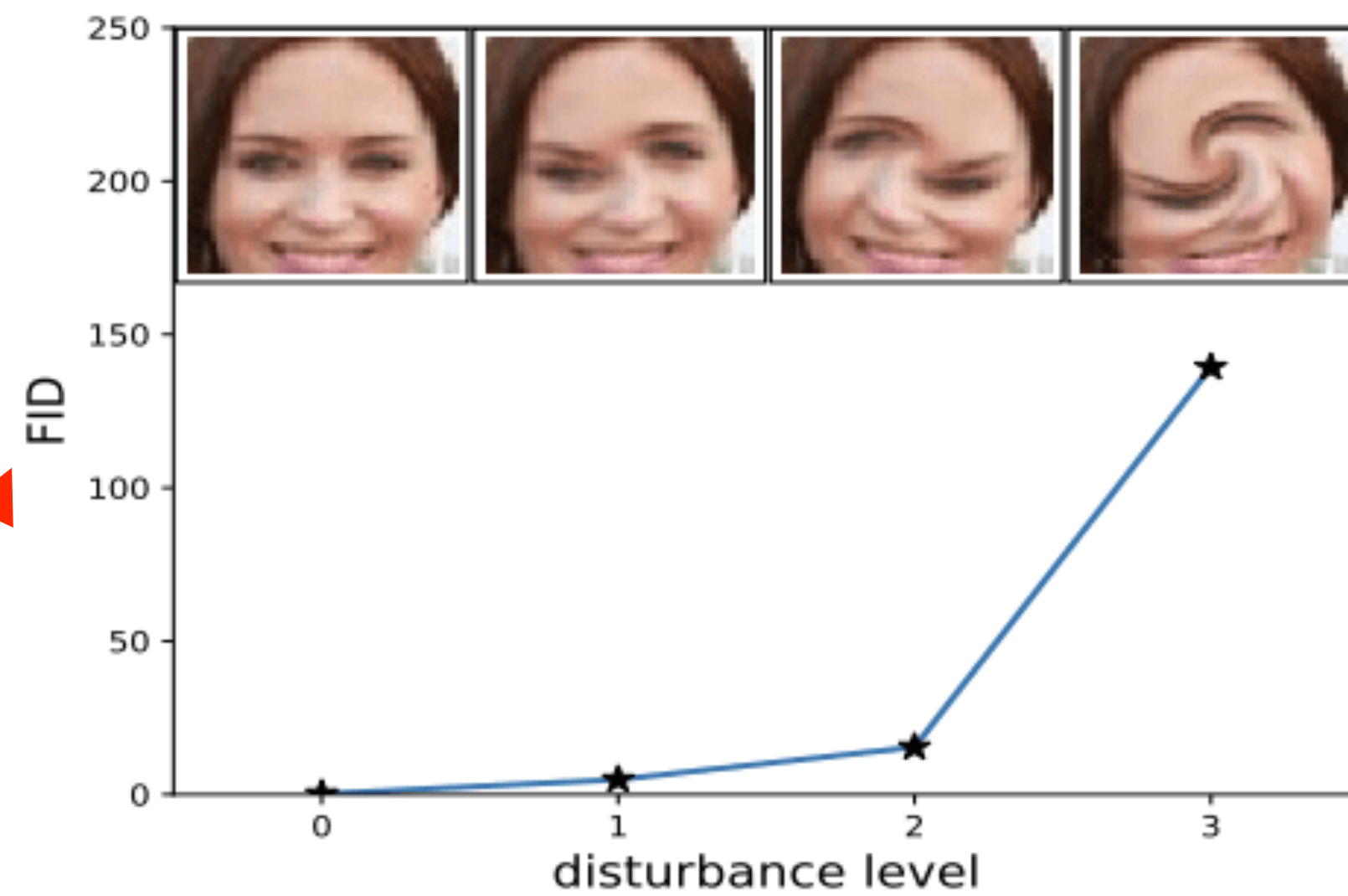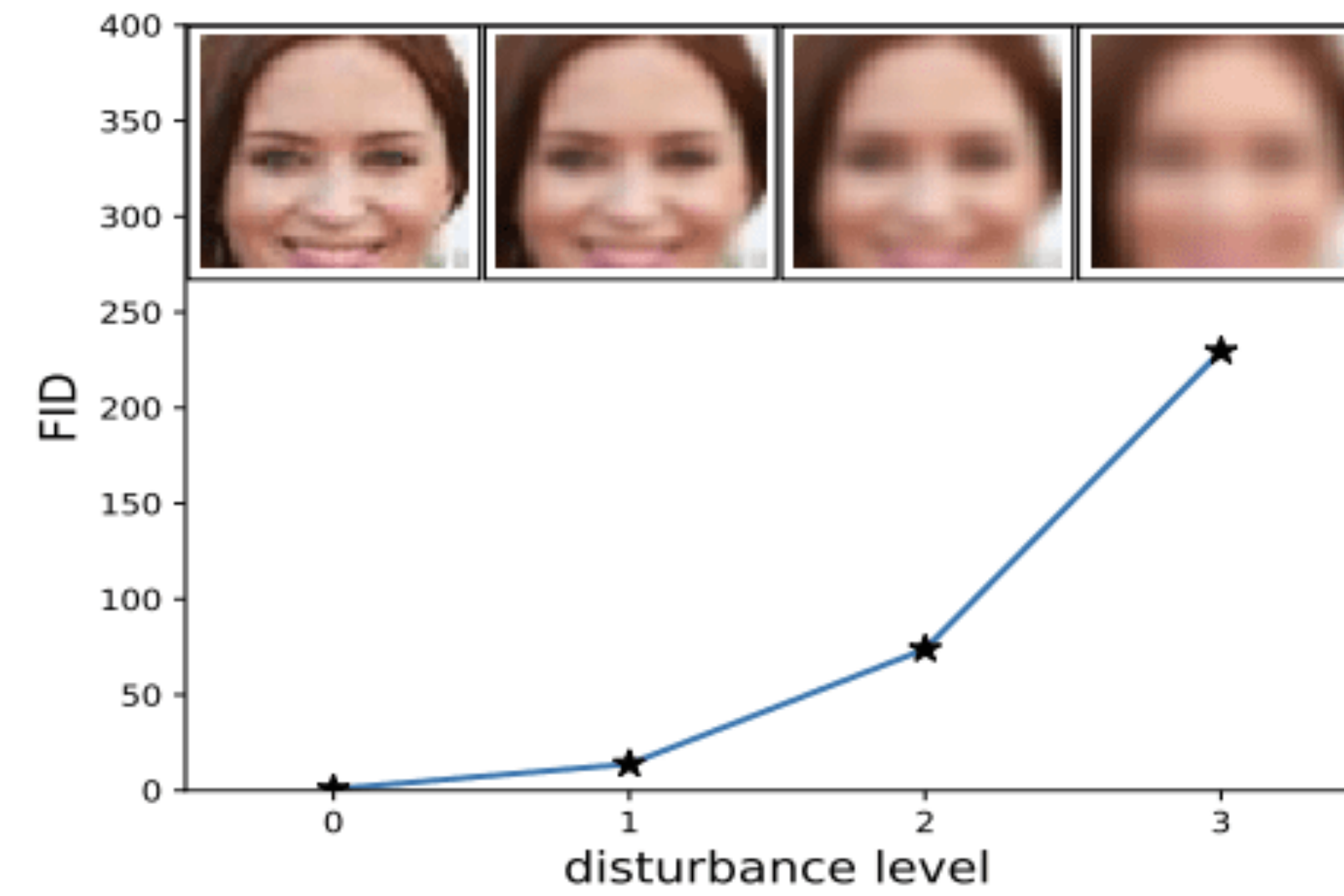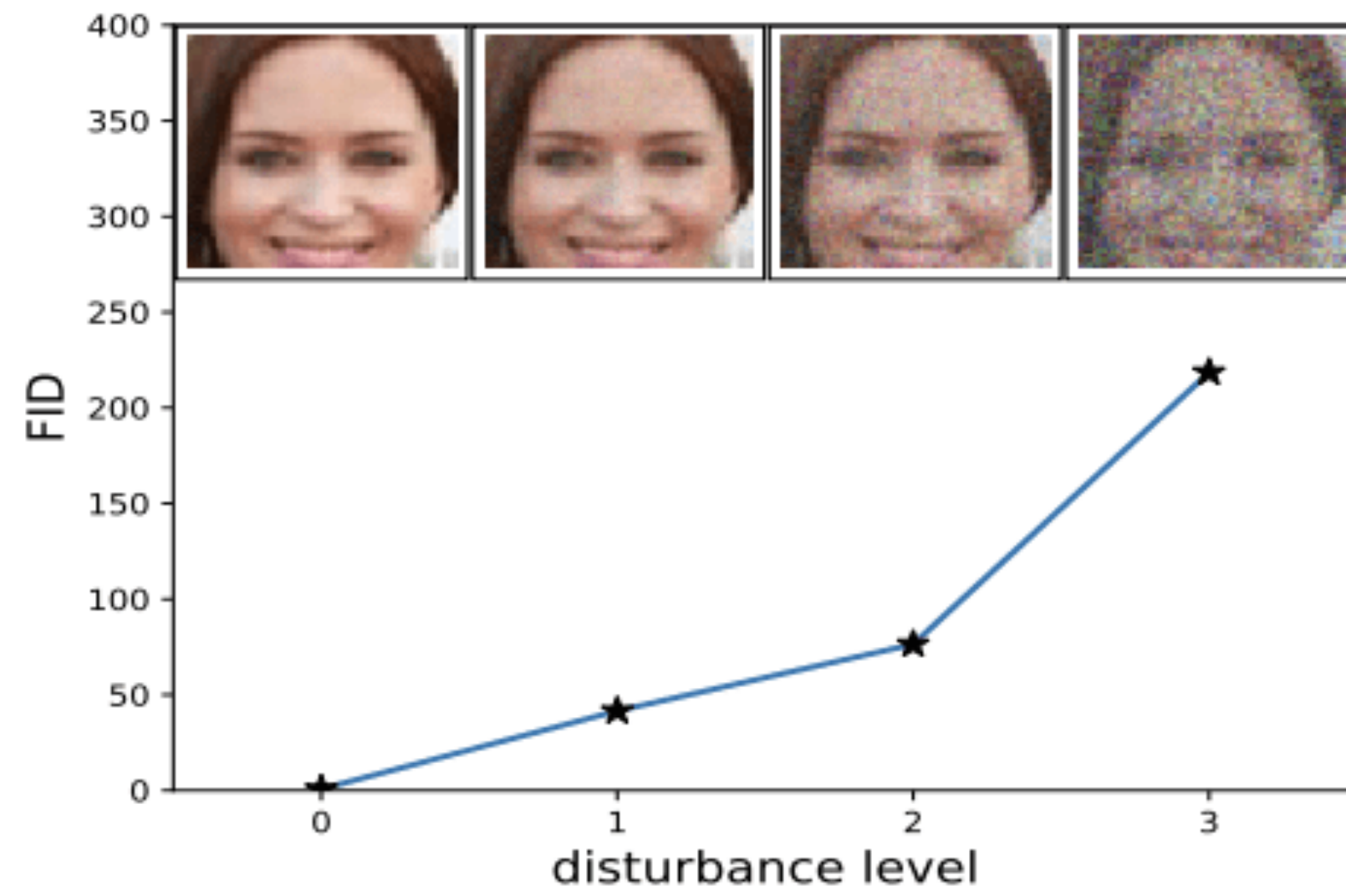
# Domain transfer/Stylization

# Description driven image manipulations

Brooks, Tim, Aleksander Holynski, and Alexei A. Efros. "Instructpix2pix: Learning to follow image editing instructions." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 2023.

# Measuring quality of generated images???



Comparing mean + std on inception v3 in the final layer with real images

Data bias

input                                                    output

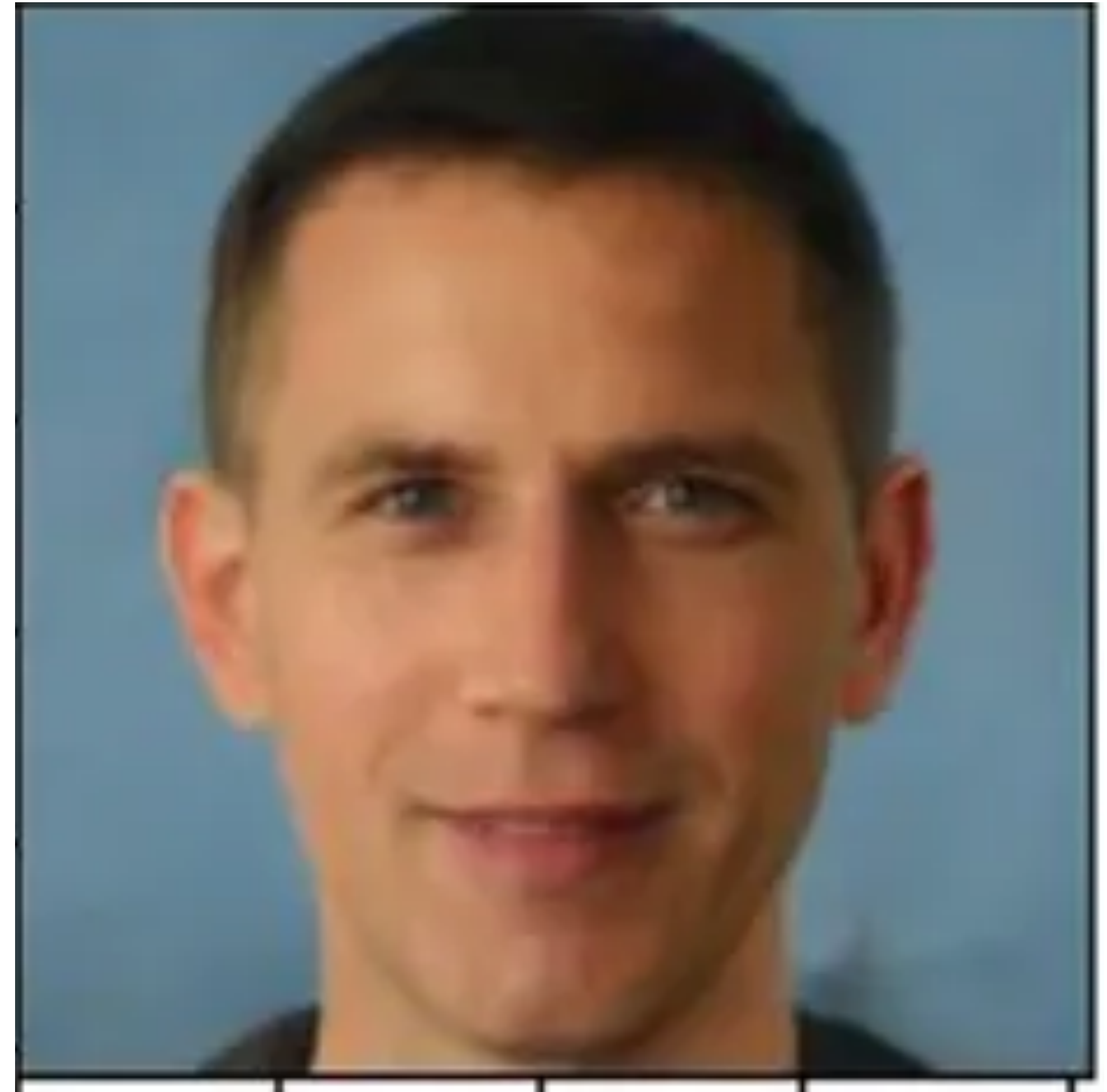For the last in this course: How is it?

# Data bias

input

output



For the last in this course: How is it?