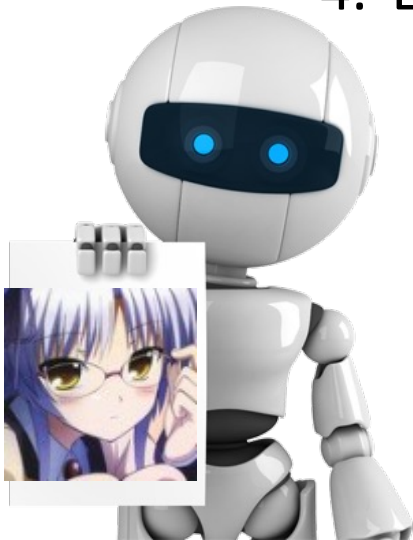


Generative models

Outline

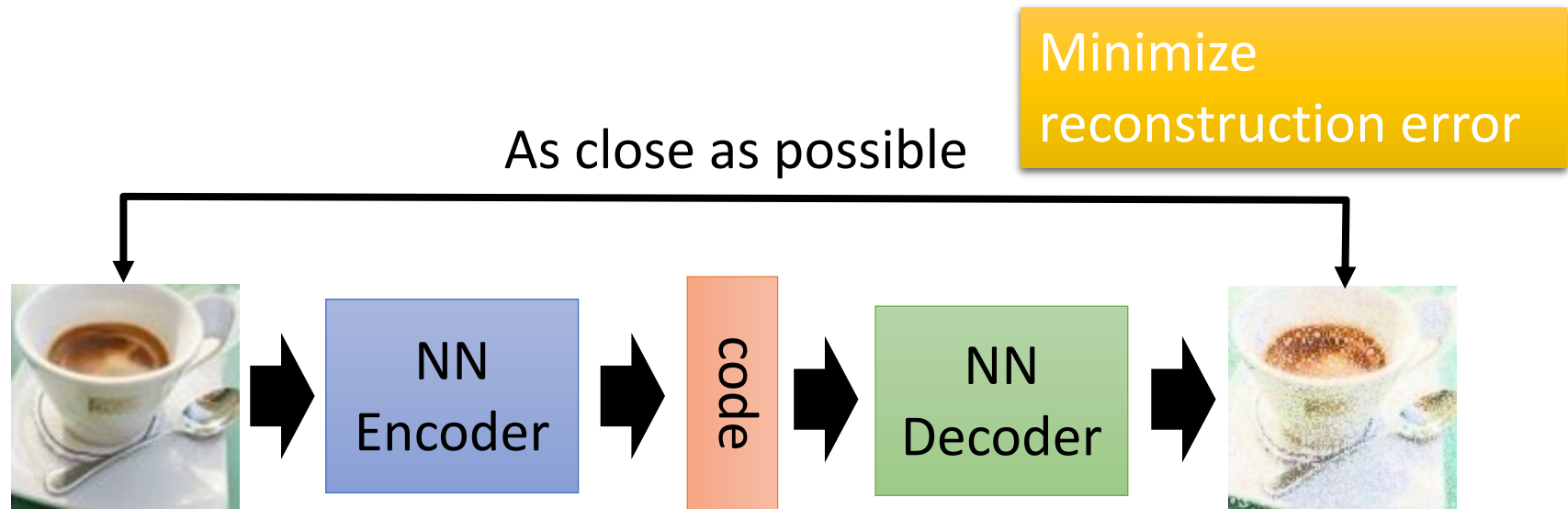
1. **Preview: Auto-Encoders, VAE**
2. Generative models with GAN
3. GAN architectures
4. Editing



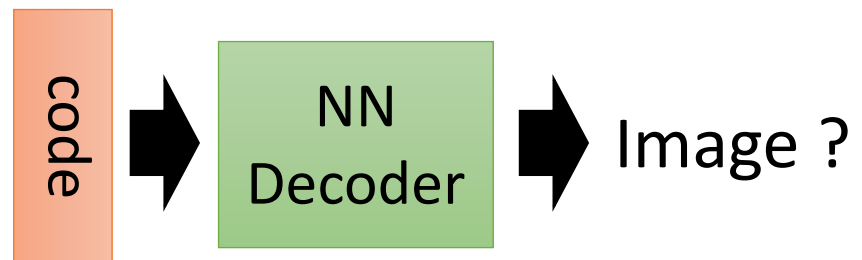
Drawing? => learning from examples



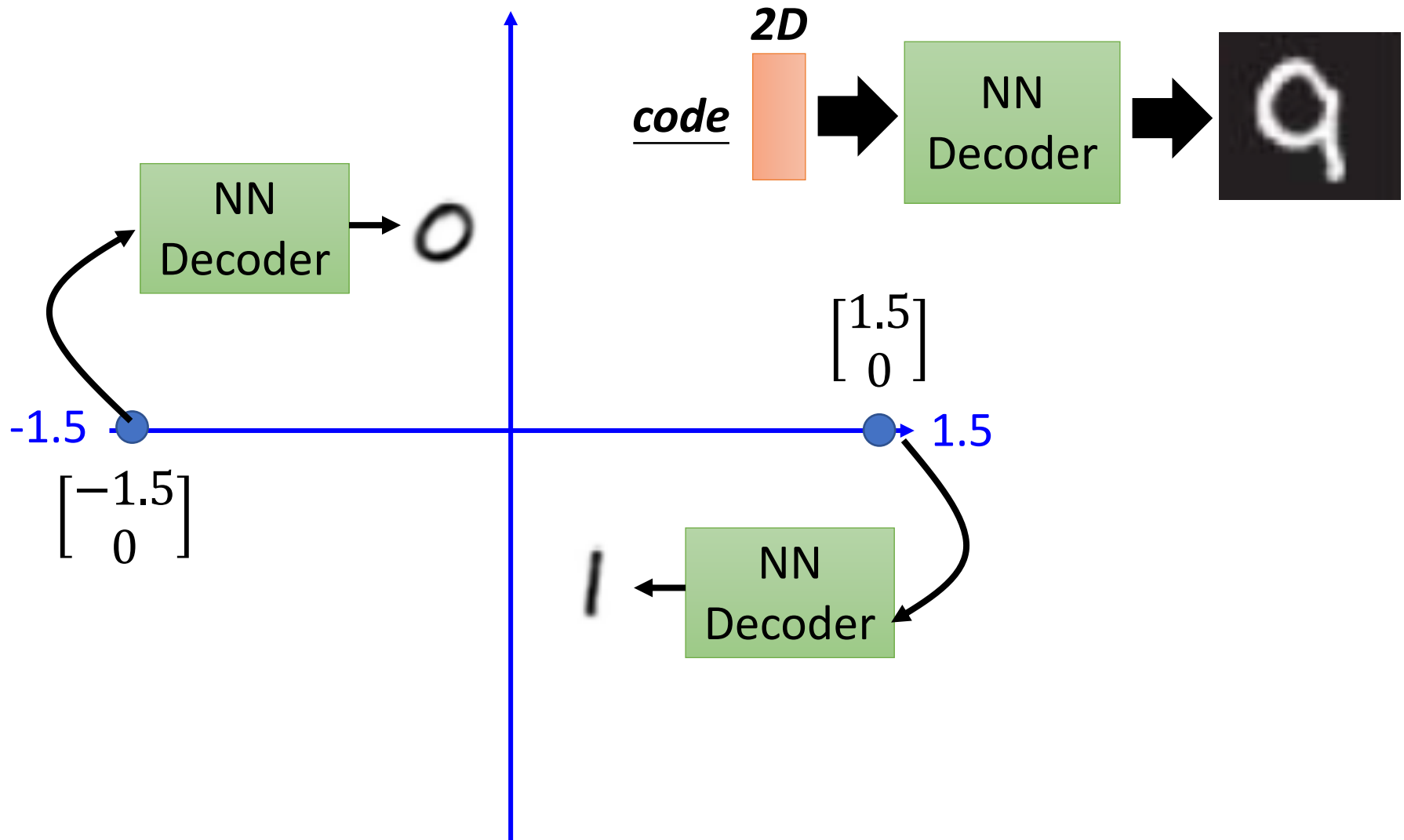
Review: Auto-encoder



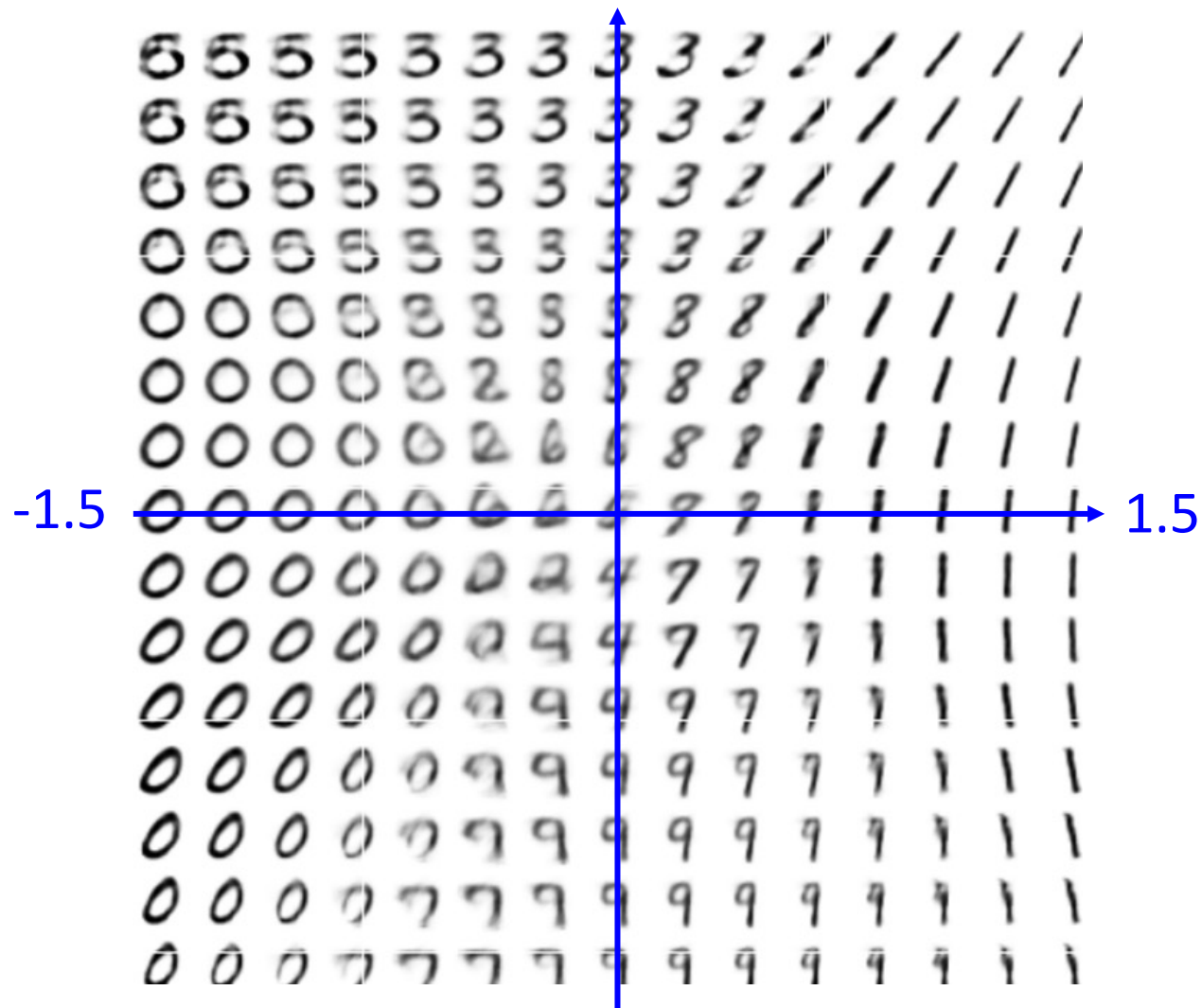
Randomly generate
a vector as code



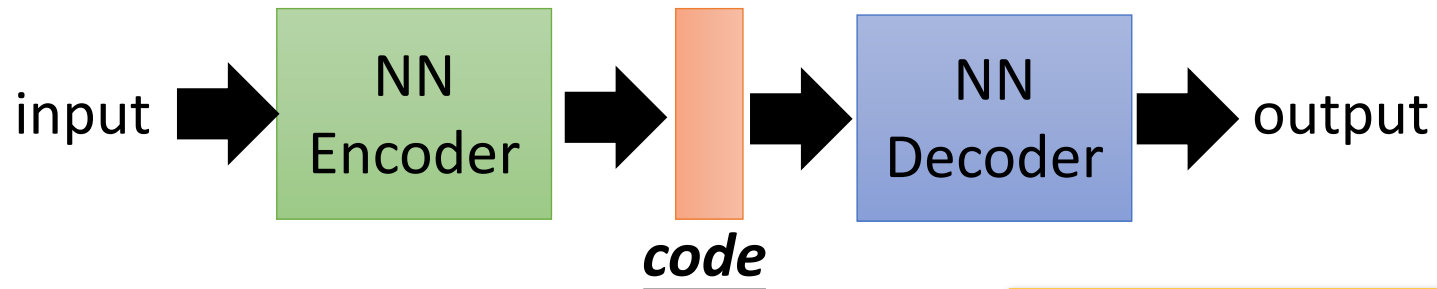
Review: Auto-encoder



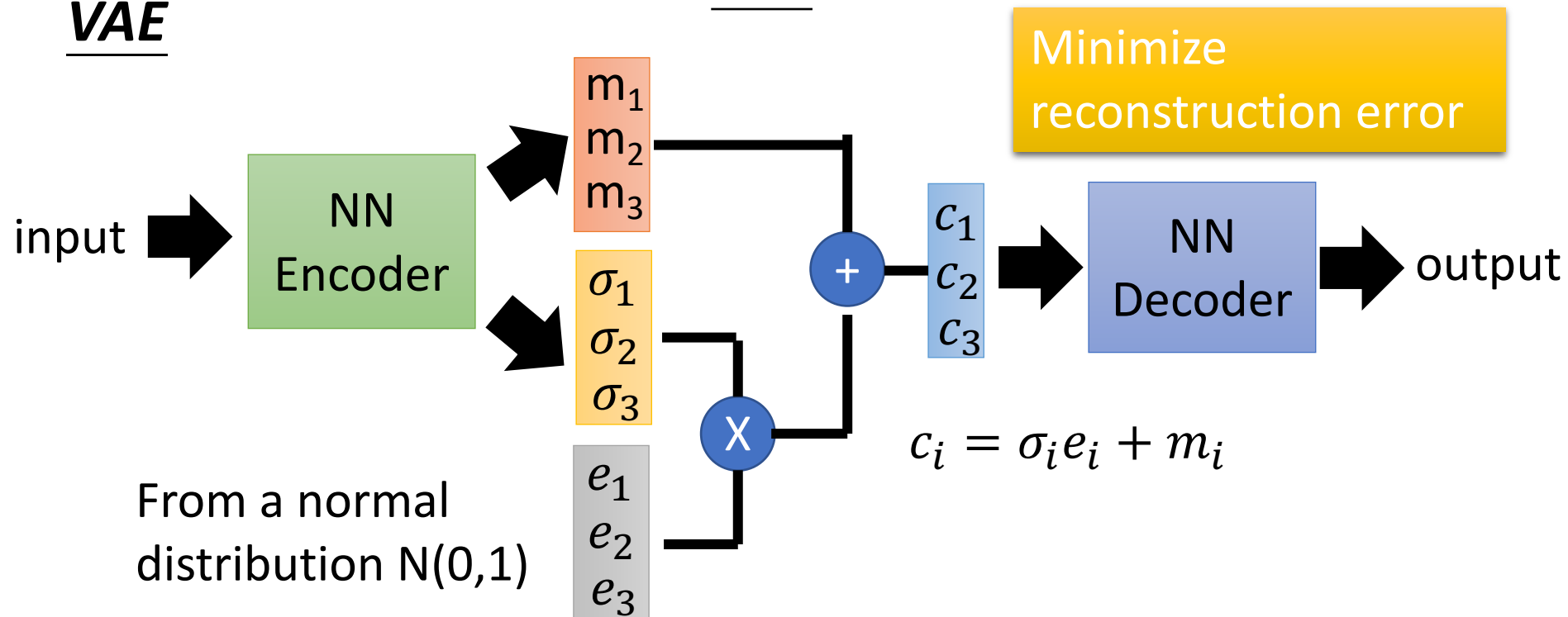
Review: Auto-encoder



Auto-encoder

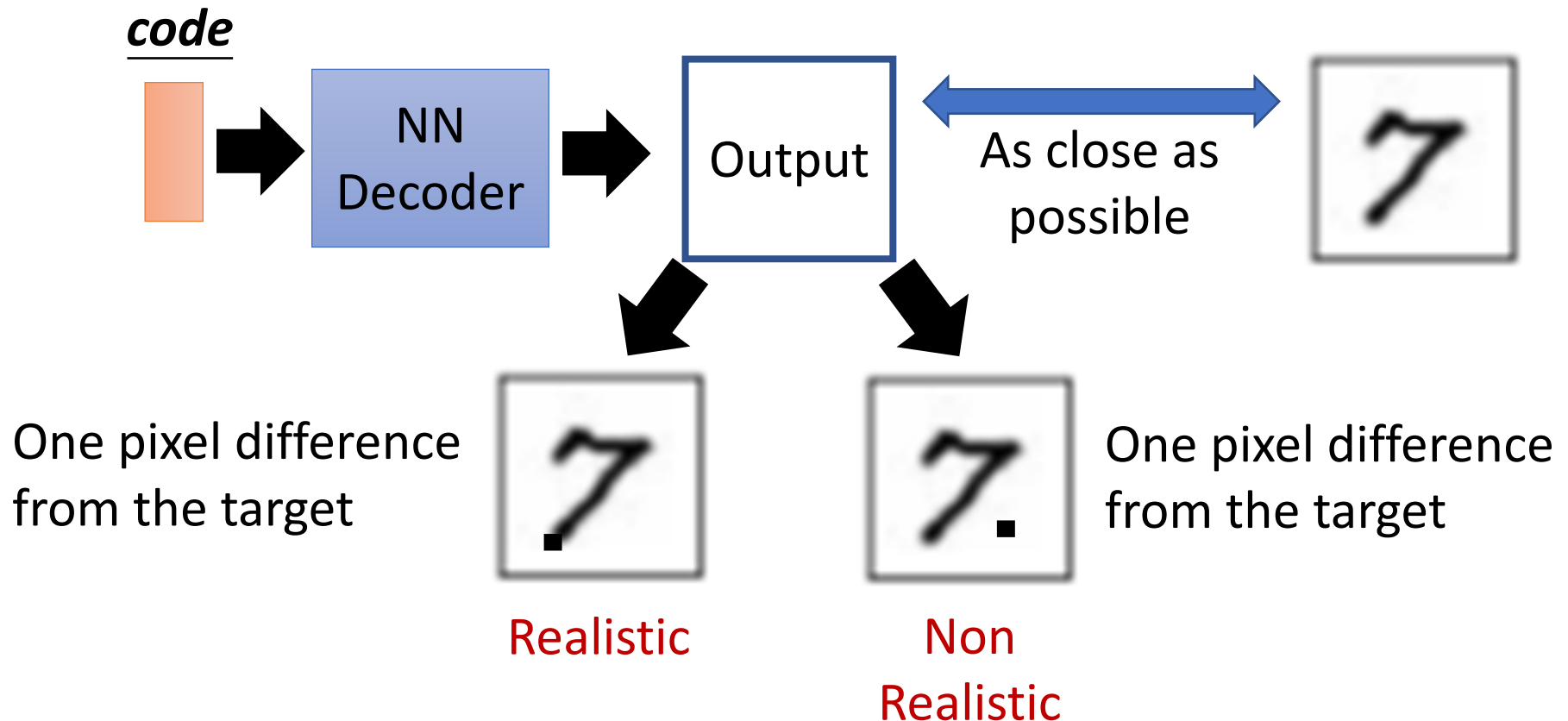


VAE



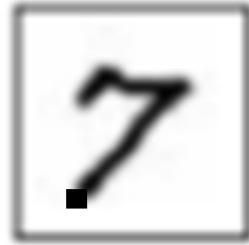
Problems of AE/VAE

- It does not really try to simulate real images

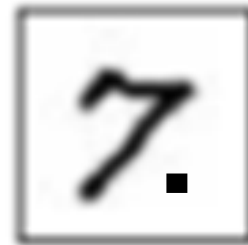


Problems of AE/VAE

GAN to tackle this pb:



Realistic



Non Realistic

GAN: generative **adversarial** networks

Game scenario:

Player1, Generator, produces samples

Player2, – Its adversary **Discriminator**, attempts to distinguish **real** samples from **fake** generated ones (produced by Player1) !

Player1 aims at producing **Realistic** images to fool the Player2

Generative models

Outline

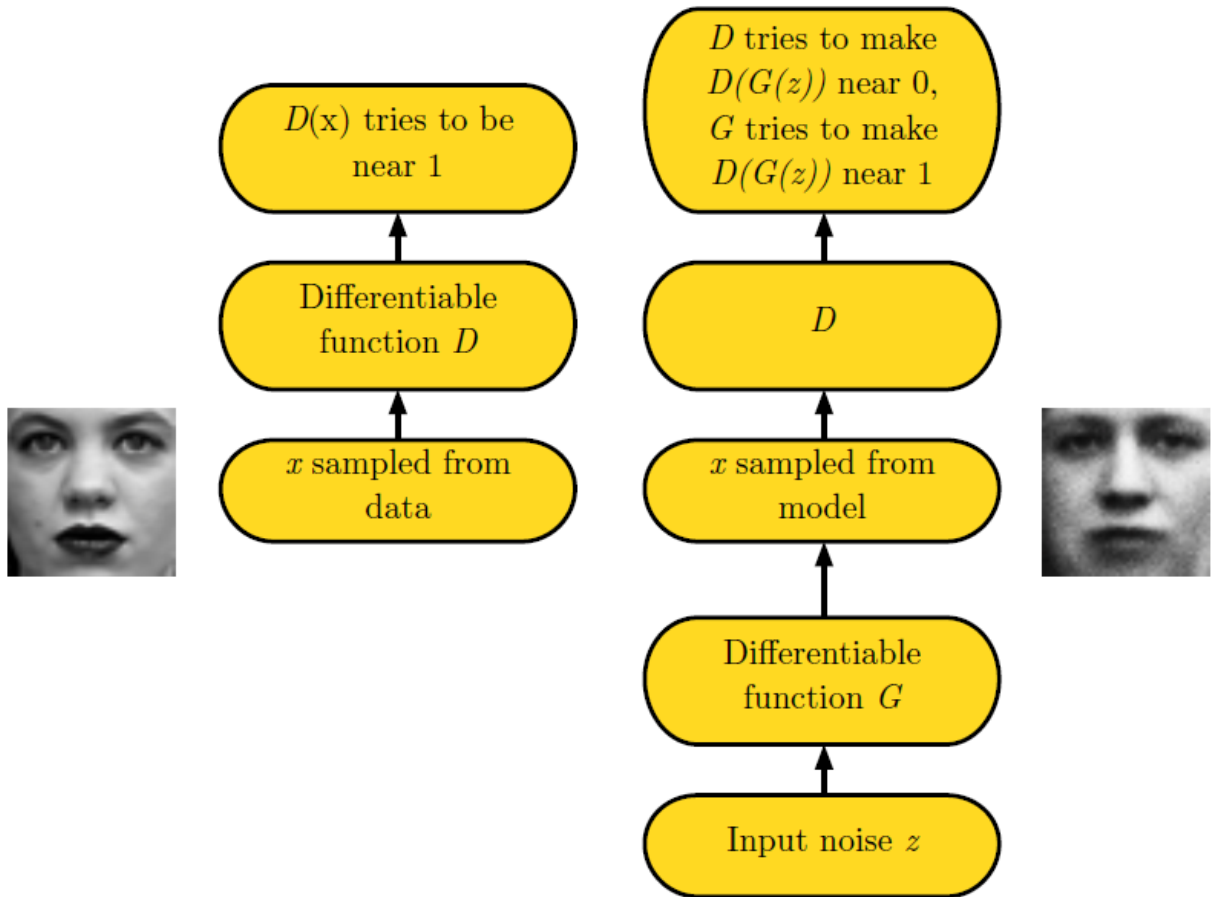
1. Preview: Auto-Encoders, VAE
- 2. Generative models with GAN**
 - GAN Algorithm

Adversarial Nets Framework

Game scenario:

Player1, Generator **G**

Player2, Discriminator **D**



$$V(\mathbf{G}, \mathbf{D}) = \mathbb{E}_{x \sim P_{data}} [\log \mathbf{D}(x)] + \mathbb{E}_{x \sim P_{\mathbf{G}}} [\log(1 - \mathbf{D}(x))]$$

$$\mathbf{G}^* = \arg \min_{\mathbf{G}} \max_{\mathbf{D}} V(\mathbf{G}, \mathbf{D})$$

GAN Learning – D and G updates

Game scenario:

Player1, Generator G, produces samples

Player2, – Its adversary **Discriminator D**, attempts to distinguish **real** samples from **fake** generated ones (produced by Player1) !

Player1 aims at producing **Realistic** images to fool the Player2

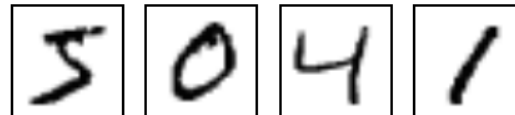
NN
Generator
v1



Discri-
minator
v1

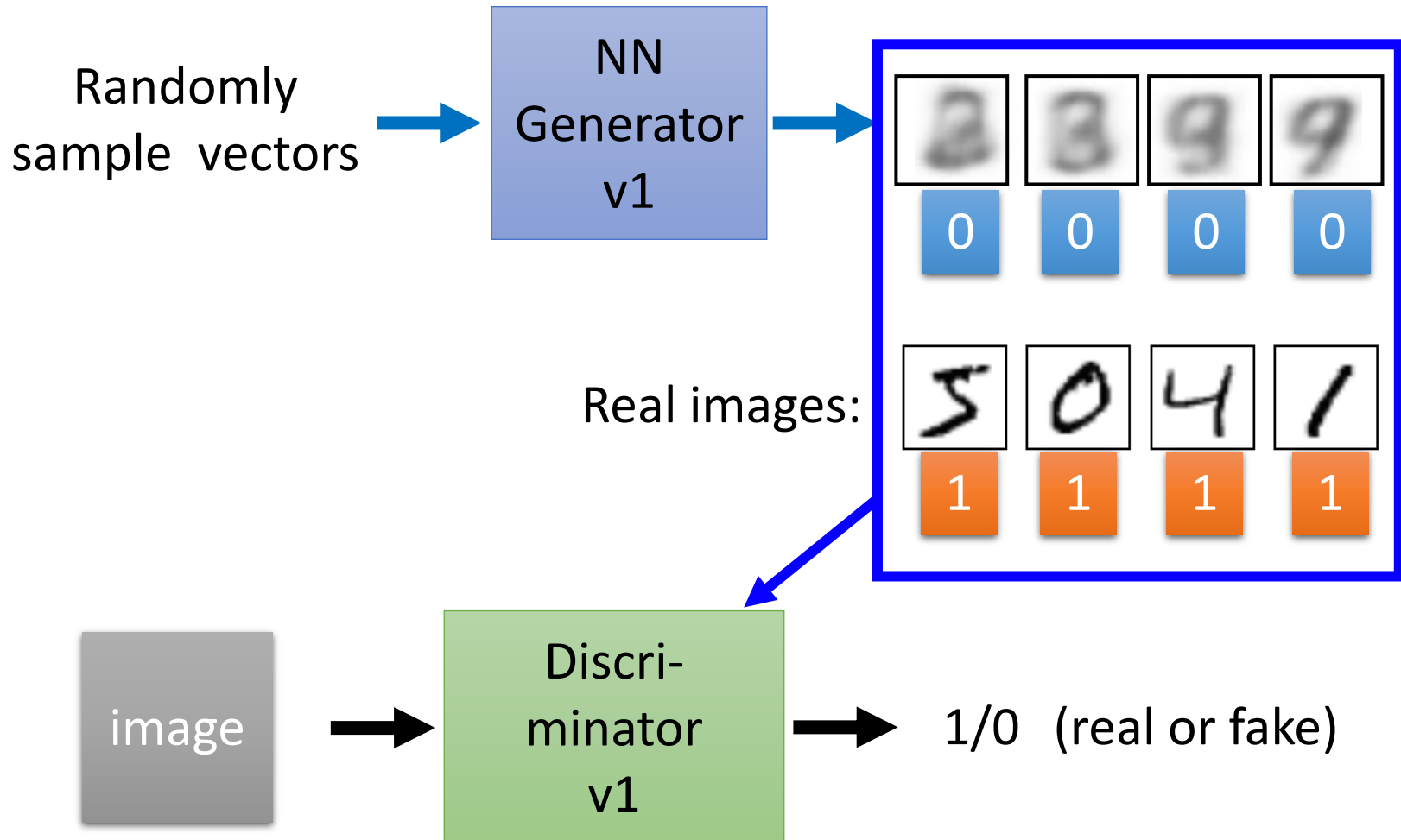
Binary
Classifier

Real images:



Fake images:

GAN - Discriminator



Discriminator Optimization on a batch of images:

Use gradient descent to update the parameters in the discriminator, with a **frozen** generator

GAN - Generator

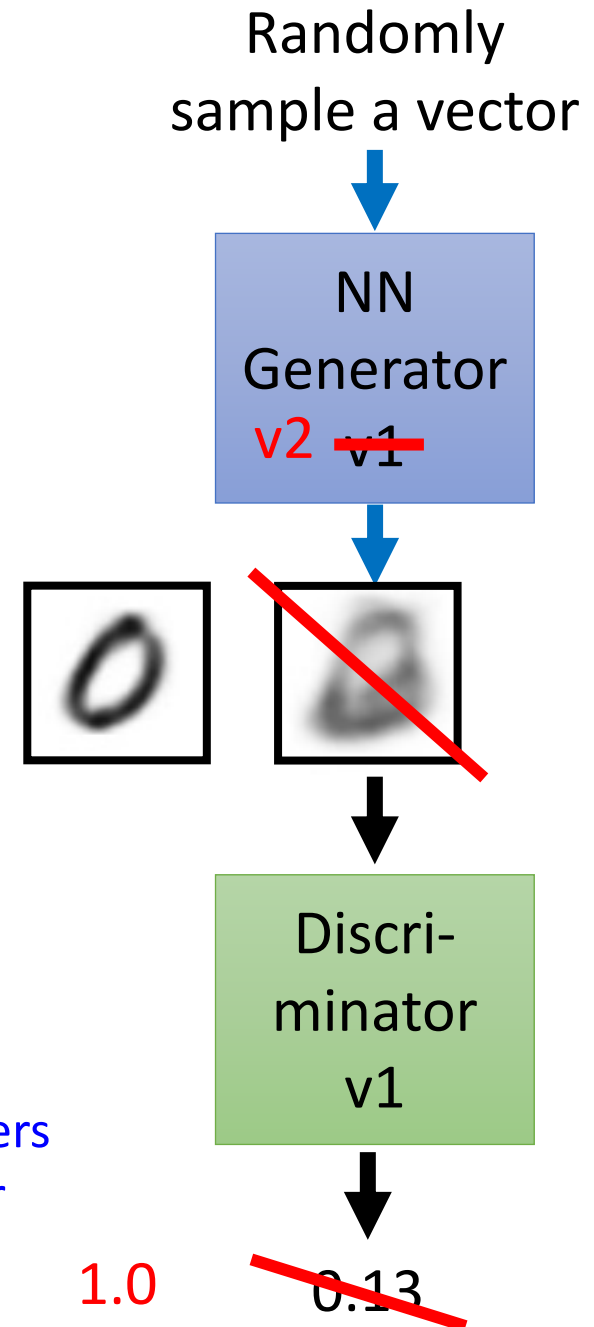
Updating the parameters of generator

➡ The output be classified as “real” (as close to 1 as possible)

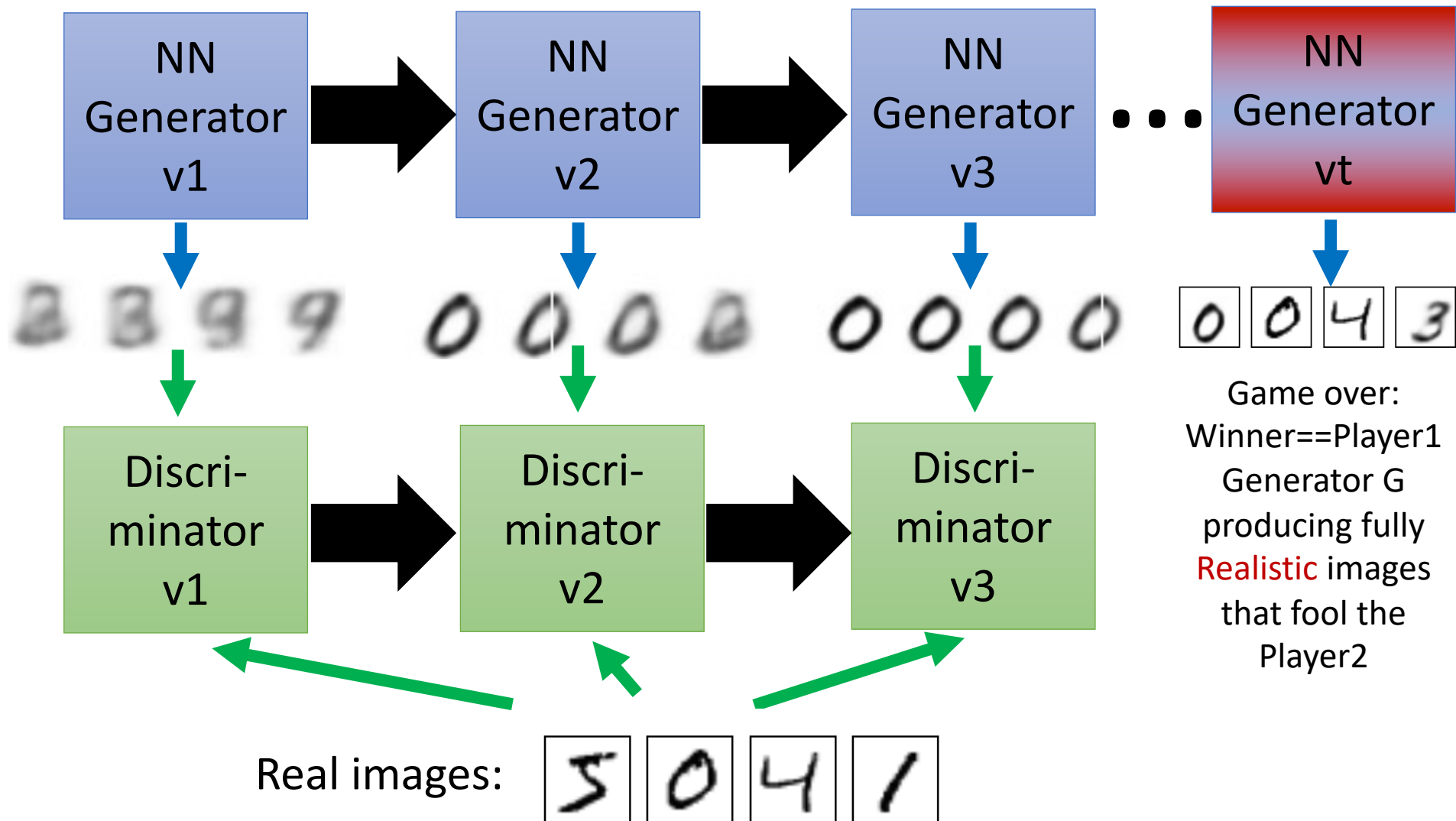
Generator + Discriminator
= a network

Optimization:

Use gradient descent to update the parameters in the generator, with a **frozen** discriminator



GAN Learning – D and G updates



Algorithm 1 Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator, k , is a hyperparameter. We used $k = 1$, the least expensive option, in our experiments.

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

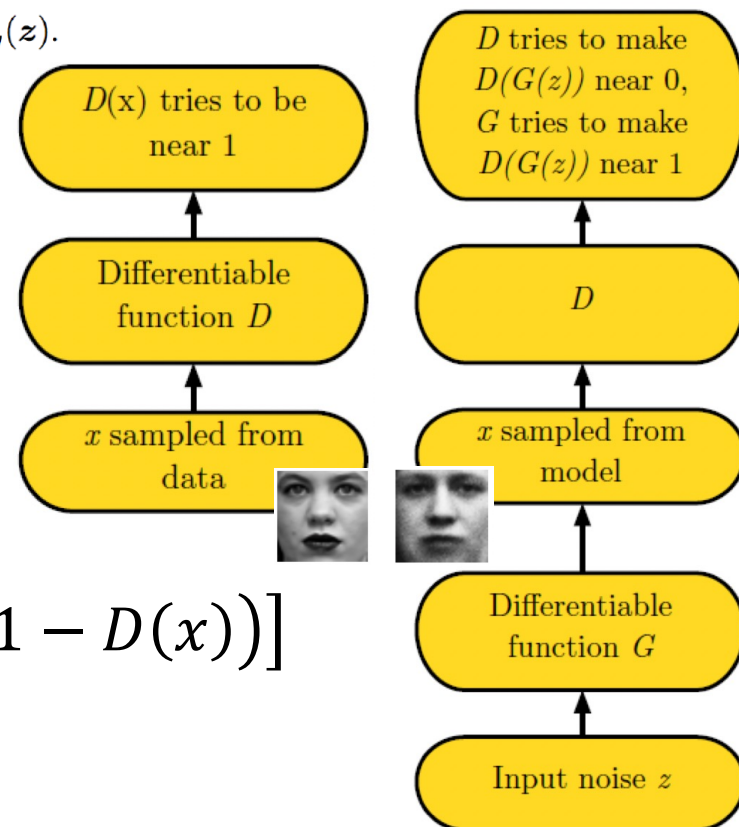
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

end for

GAN algorithm

$$V = \mathbb{E}_{x \sim P_{\text{data}}} [\log D(x)] + \mathbb{E}_{x \sim P_G} [\log (1 - D(x))]$$

$$G^* = \arg \min_G \max_D V(G, D)$$



One example GAN

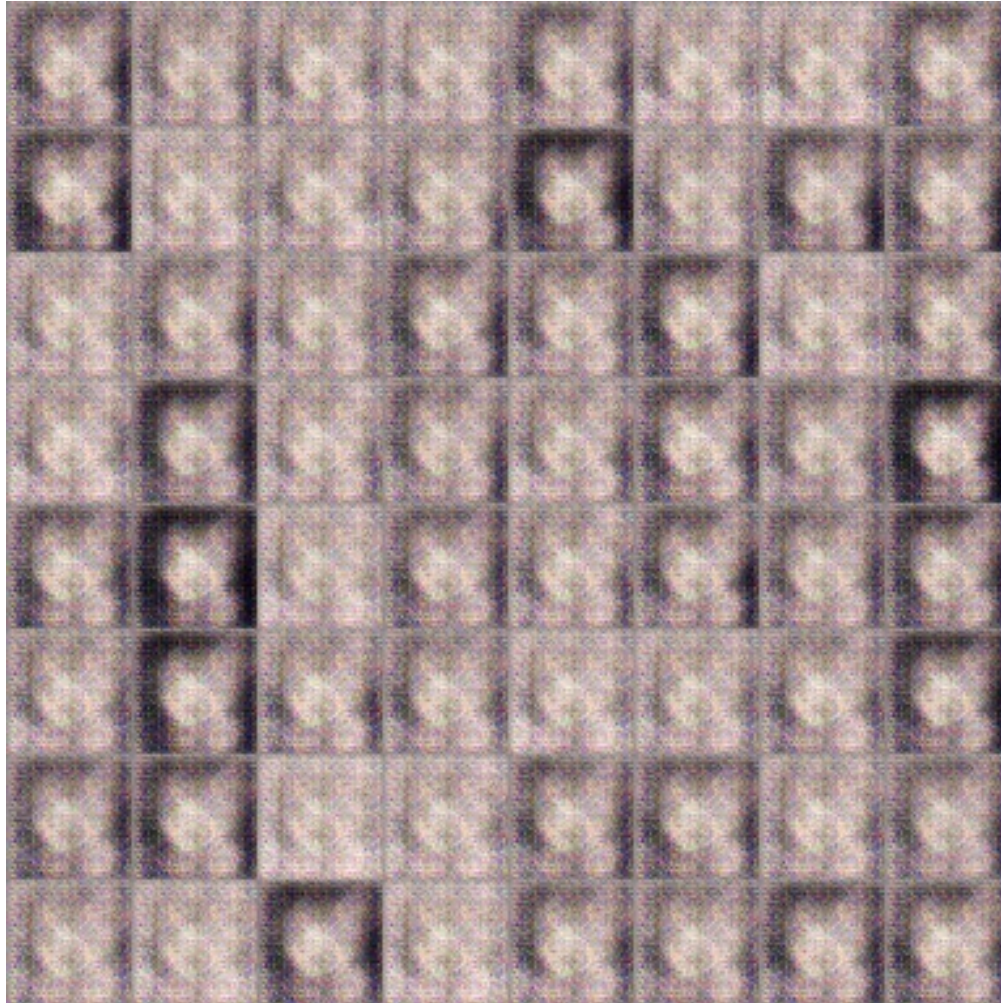


Source of images: <https://zhuanlan.zhihu.com/p/24767059>

DCGAN: <https://github.com/carpedm20/DCGAN-tensorflow>

GAN

100 rounds

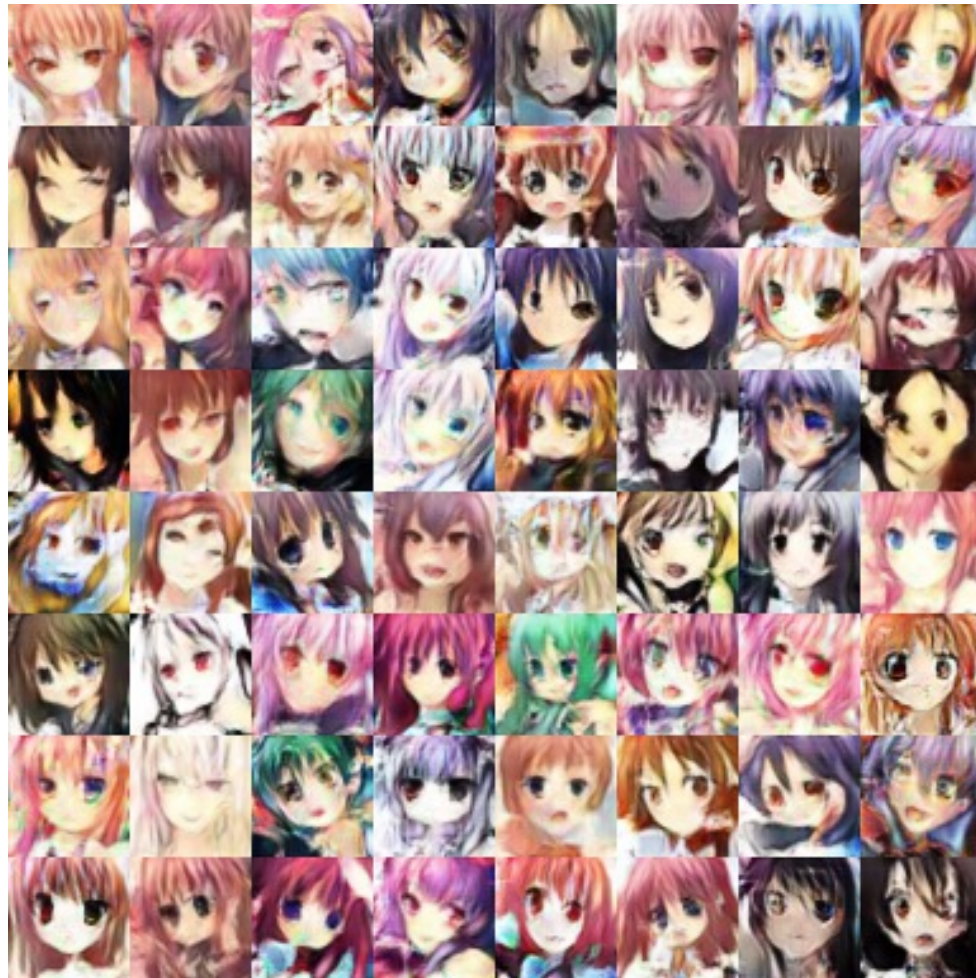


GAN



1000 rounds

GAN



20,000 rounds

GAN



50,000 rounds

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
 - GAN Algorithm
 - **KL vs. Jensen Shannon Divergence**

$$V(\textcolor{red}{G}, D) = \mathbb{E}_{x \sim P_{data}} [\log D(x)] + \mathbb{E}_{x \sim P_{\textcolor{red}{G}}} [\log(1 - D(x))]$$

$$\textcolor{red}{G}^* = \arg \min_{\textcolor{red}{G}} \max_D V(\textcolor{red}{G}, D)$$

Which measure to evaluate how $P_G(x; \theta)$ is close to $P_{data}(x)$ in Maximum Likelihood optimization?

- Given a data distribution $P_{data}(x)$
- We have a distribution $P_G(x; \theta)$ parameterized by θ
 - E.g. $P_G(x; \theta)$ is a Gaussian Mixture Model, θ are means and variances of the Gaussians
 - We want to find θ such that $P_G(x; \theta)$ close to $P_{data}(x)$

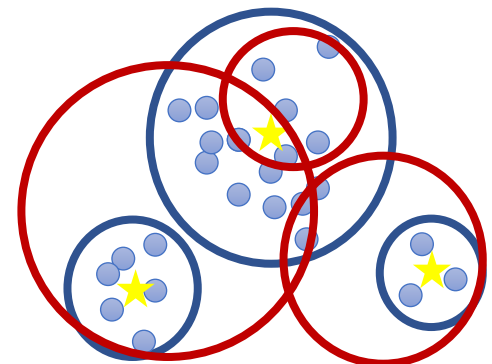
Sample $\{x^1, x^2, \dots, x^m\}$ from $P_{data}(x)$

We can compute $P_G(x^i; \theta)$

Likelihood of generating the samples

$$L = \prod_{i=1}^m P_G(x^i; \theta)$$

Find θ^* maximizing the likelihood



Which measure to evaluate how $P_G(x; \theta)$ is close to $P_{data}(x)$ in Maximum Likelihood optimization?

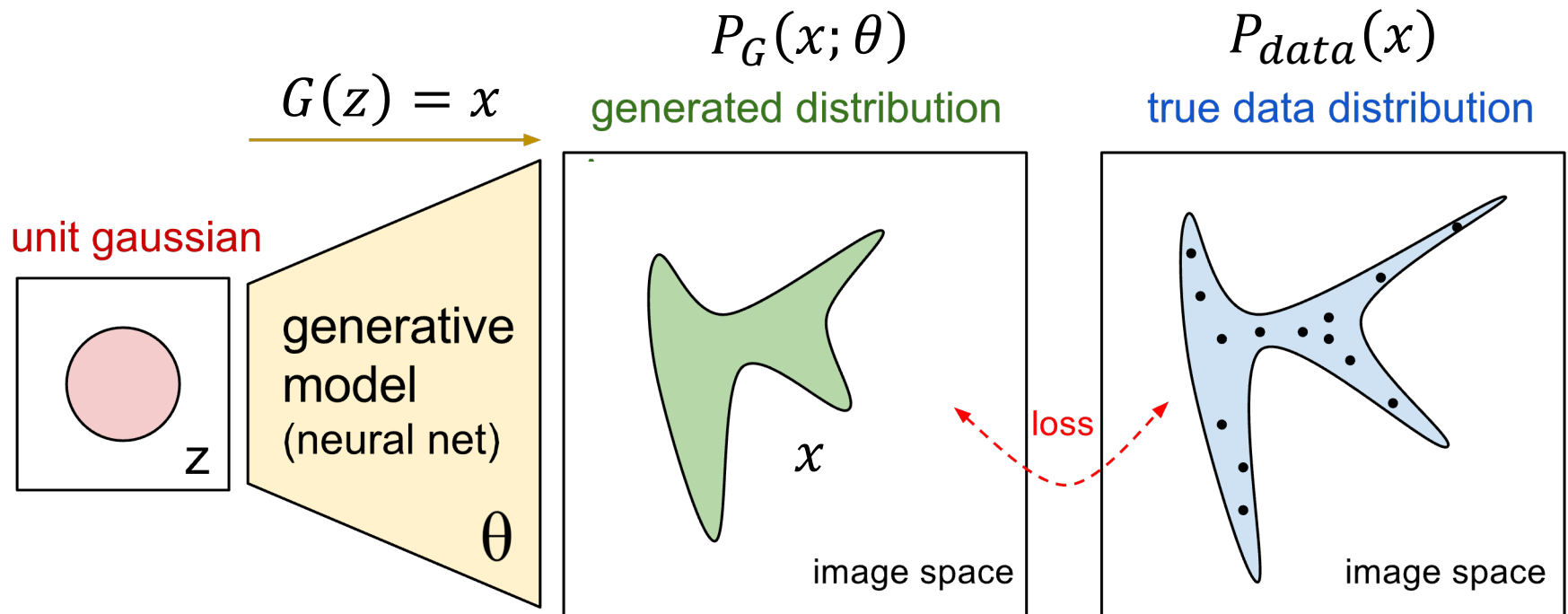
$$\begin{aligned}\theta^* &= \arg \max_{\theta} \prod_{i=1}^m P_G(x^i; \theta) = \arg \max_{\theta} \log \prod_{i=1}^m P_G(x^i; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^m \log P_G(x^i; \theta) \quad \{x^1, x^2, \dots, x^m\} \text{ from } P_{data}(x)\end{aligned}$$

$$\approx \arg \max_{\theta} \mathbb{E}_{x \sim P_{data}} [\log P_G(x; \theta)]$$

$$\begin{aligned}&= \arg \max_{\theta} \int_x P_{data}(x) \log P_G(x; \theta) dx - \int_x P_{data}(x) \log P_{data}(x) dx \\ &= \arg \min_{\theta} KL(P_{data}(x) || P_G(x; \theta)) \quad KL(P || Q) = \int_x P(x) \log \frac{P(x)}{Q(x)} dx\end{aligned}$$

In Maximum Likelihood it is a KLD Kullback Leibler Divergence

If $P_G(x; \theta)$ is a coming with a NN



$$P_G(x; \theta) = \int_z P_{prior}(z) I_{[G(z)=x]} dz$$

It is difficult to compute the likelihood.

Basic Idea of GAN: the 2 players G-D game

- Generator G Hard to learn by maximum likelihood
 - G is a function, input z , output x
 - Given a prior distribution $P_{\text{prior}}(z)$, a probability distribution $P_G(x)$ is defined by function G (and P_{prior})
- Discriminator D
 - D is a function, input x , output scalar
 - Evaluate the “difference” between $P_G(x)$ and $P_{\text{data}}(x)$
- Global objective function $V(G, D)$

$$\theta^* = G^* = \arg \min_G \max_D V(G, D)$$

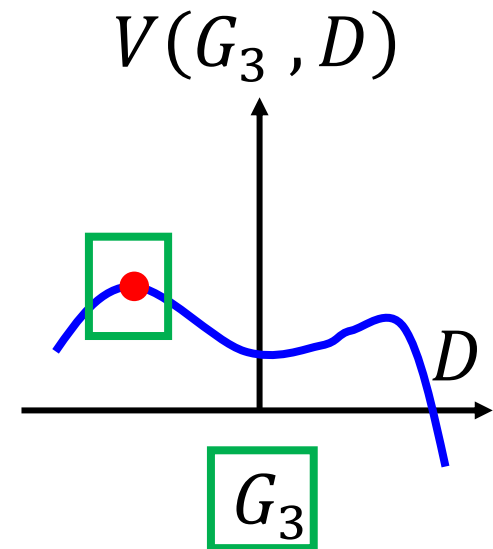
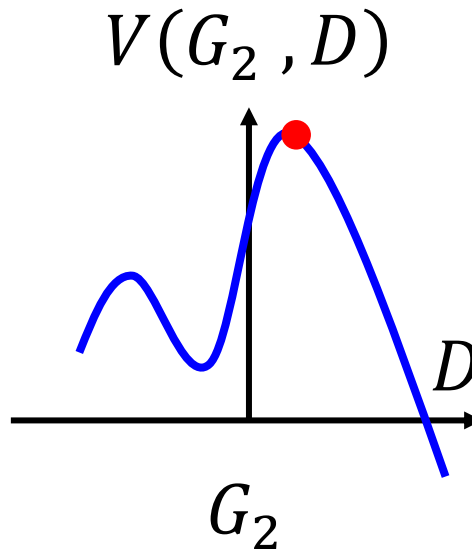
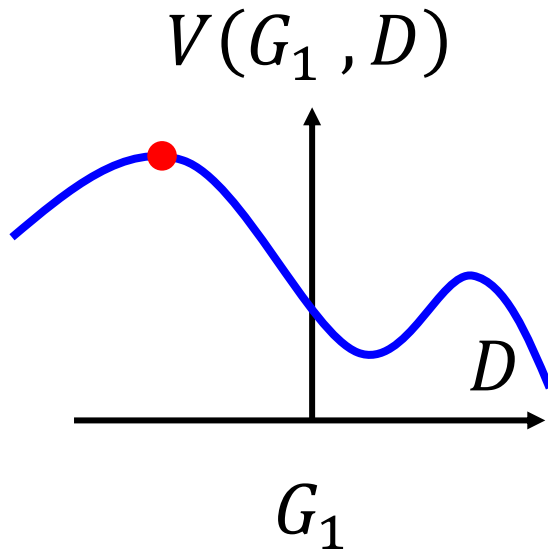
Basic Idea

$$\boxed{G^*} = \arg \min_G \max_D V(G, D)$$

$$V = \mathbb{E}_{x \sim P_{data}} [\log D(x)] + \mathbb{E}_{x \sim P_G} [\log(1 - D(x))]$$

Given a generator G , $\max_D V(G, D)$ evaluate the “difference” between P_G and P_{data}

Pick the G defining P_G most similar to P_{data}



$$\max_D V(G, D) \quad G^* = \arg \min_G \max_D V(G, D)$$

- Given G , what is the optimal D^* maximizing

$$\begin{aligned} V &= \mathbb{E}_{x \sim P_{data}} [\log D(x)] + \mathbb{E}_{x \sim P_G} [\log(1 - D(x))] \\ &= \int_x P_{data}(x) \log D(x) dx + \int_x P_G(x) \log(1 - D(x)) dx \\ &= \int_x [P_{data}(x) \log D(x) + P_G(x) \log(1 - D(x))] dx \end{aligned}$$

Assume that $D(x)$ can have any value here

- Given x , the optimal D^* maximizing

$$P_{data}(x) \log D(x) + P_G(x) \log(1 - D(x))$$

$$\max_D V(G, D) \quad G^* = \arg \min_G \max_D V(G, D)$$

- Given x , the optimal D^* maximizing

$$\underbrace{P_{data}(x)}_a \log \underbrace{D(x)}_D + \underbrace{P_G(x)}_b \log \underbrace{(1 - D(x))}_D$$

- Find D^* maximizing: $f(D) = a \log(D) + b \log(1 - D)$

$$\frac{df(D)}{dD} = a \times \frac{1}{D} + b \times \frac{1}{1 - D} \times (-1) = 0$$

$$a \times \frac{1}{D^*} = b \times \frac{1}{1 - D^*}$$

$$a \times (1 - D^*) = b \times D^*$$

$$a - aD^* = bD^*$$

$$D^* = \frac{a}{a + b} \rightarrow$$

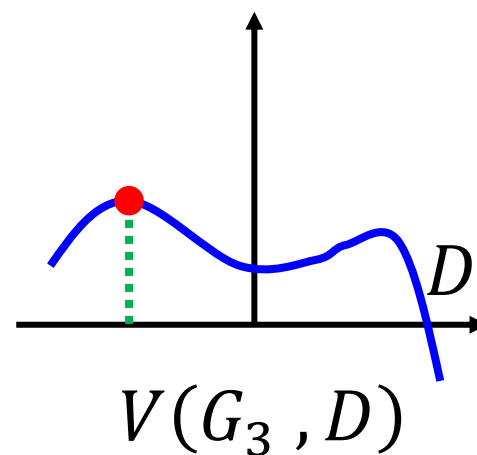
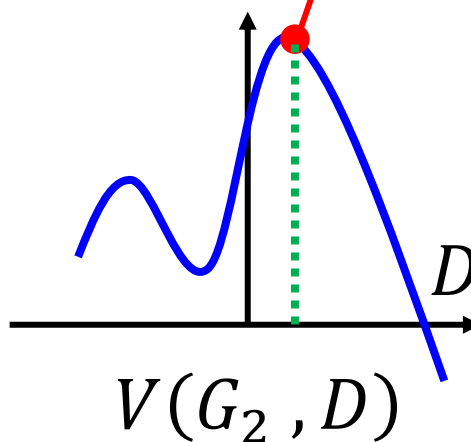
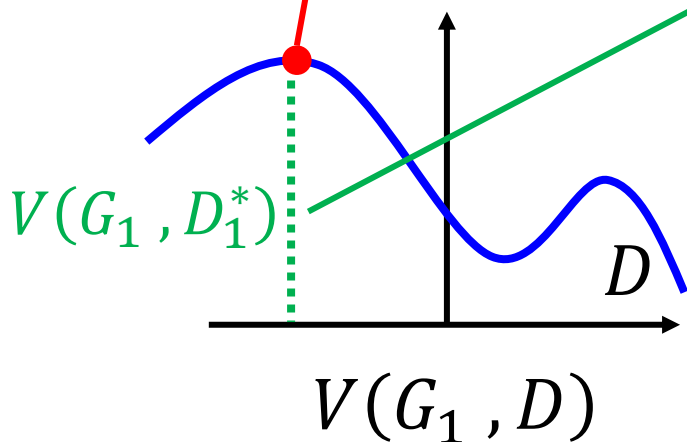
$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} \quad 0 < \quad < 1$$

$$\max_D V(G, D) \quad G^* = \arg \min_G \max_D V(G, D)$$

$$D_1^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_{G_1}(x)}$$

$$D_2^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_{G_2}(x)}$$

“difference” between
 P_{G_1} and P_{data}



$$\max_D V(G, D)$$

$$V = \mathbb{E}_{x \sim P_{data}} [\log D(x)] + \mathbb{E}_{x \sim P_G} [\log(1 - D(x))]$$

$$\max_D V(G, D) = V(G, D^*)$$

$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)}$$

$$= \mathbb{E}_{x \sim P_{data}} \left[\log \frac{P_{data}(x)}{P_{data}(x) + P_G(x)} \right] + \mathbb{E}_{x \sim P_G} \left[\log \frac{P_G(x)}{P_{data}(x) + P_G(x)} \right]$$

$$= \int_x P_{data}(x) \log \frac{\frac{1}{2} P_{data}(x)}{\frac{P_{data}(x) + P_G(x)}{2}} dx + \int_x P_G(x) \log \frac{\frac{1}{2} P_G(x)}{\frac{P_{data}(x) + P_G(x)}{2}} dx$$

$$\Rightarrow +2 \log \frac{1}{2} = -2 \log 2$$

$$\max_D V(G, D)$$

$$\begin{aligned} \text{JSD}(P||Q) &= \frac{1}{2}\text{KL}(P||M) + \frac{1}{2}\text{KL}(Q||M) \\ M &= \frac{1}{2}(P + Q) \end{aligned}$$

$$\max_D V(G, D) = V(G, D^*)$$

$$D^*(x) = \frac{P_{data}(x)}{P_{data}(x) + P_G(x)}$$

$$= -2\log 2 + \int_x P_{data}(x) \log \frac{P_{data}(x)}{(P_{data}(x) + P_G(x))/2} dx$$

$$+ \int_x P_G(x) \log \frac{P_G(x)}{(P_{data}(x) + P_G(x))/2} dx$$

$$\begin{aligned} = -2\log 2 + \text{KL} \left(P_{data}(x) || \frac{P_{data}(x) + P_G(x)}{2} \right) \\ + \text{KL} \left(P_G(x) || \frac{P_{data}(x) + P_G(x)}{2} \right) \end{aligned}$$

$$= -2\log 2 + 2\text{JSD}(P_{data}(x)||P_G(x)) \quad \text{Jensen-Shannon divergence}$$

In the end

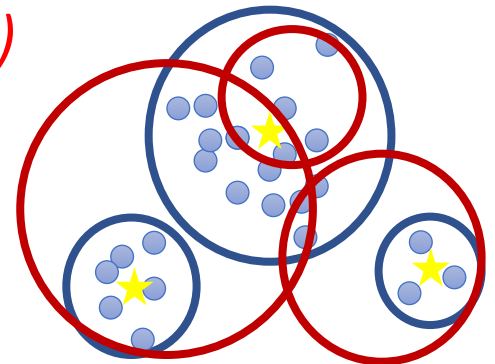
$$V = \mathbb{E}_{x \sim P_{data}} [\log D(x)] \\ + \mathbb{E}_{x \sim P_G} [\log(1 - D(x))]$$

- Generator G, Discriminator D
- Looking for G^* such that $G^* = \arg \min_G \max_D V(G, D)$
- Given G, $\max_D V(G, D) = -2\log 2 + 2JSD(P_{data}(x) || P_G(x))$
 $0 < \hspace{15em} < \log 2$
- What is the optimal G?

$$P_G(x) = P_{data}(x)$$

with/using the $JS(P_G, P_{data})$ Divergence

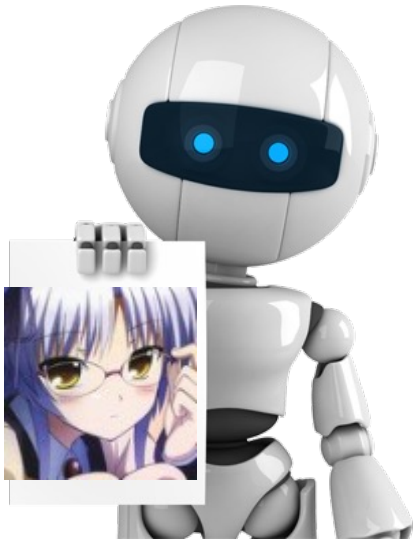
(In Maximum Likelihood it is a KL Divergence)



Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. **GAN architectures**



Drawing? => learning from examples



Recall Algo GAN

for number of training iterations **do**

for k steps **do**

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Sample minibatch of m examples $\{x^{(1)}, \dots, x^{(m)}\}$ from data generating distribution $p_{\text{data}}(x)$.
- Update the discriminator by ascending its stochastic gradient:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(x^{(i)}) + \log (1 - D(G(z^{(i)}))) \right].$$

end for

- Sample minibatch of m noise samples $\{z^{(1)}, \dots, z^{(m)}\}$ from noise prior $p_g(z)$.
- Update the generator by descending its stochastic gradient:

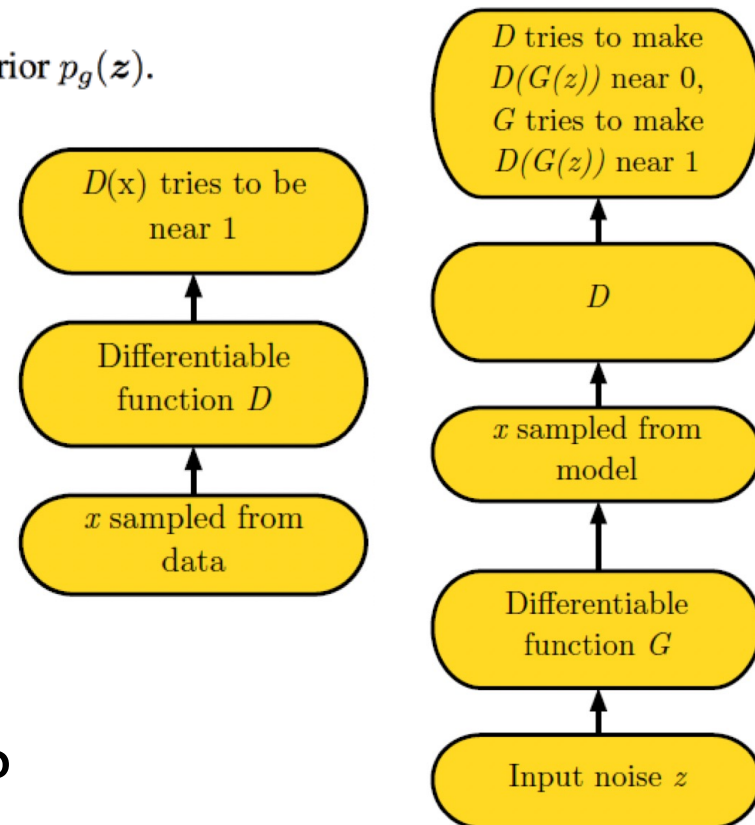
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(z^{(i)}))).$$

end for

Functions G and D are NN

Question:

Which architectures for G and D?



Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
- 3. GAN architectures**
 - 1. Basics**

Basic Archi for G and D and expe

Models

G and D fully connected nets

or convolutional for D , (de)convolutional for G (as seen for segmentation nets)

ReLU and/or sigmoids, dropout

Datasets

MNIST, Toronto Face Database, CIFAR-10

GAN - Evaluation

- Approximate p_g by fitting a Gaussian Parzen window on the generated images.
- Cross-validate σ to maximize likelihood of validation set
- Compute the likelihood of the test set

Evaluation not trivial, can be done using generated images as inputs for deep nets => inception scores

Frechet inception distance using 2 Gaussian

(data,gen) over inception features: $\text{FID} = \|\mu - \mu_w\|_2^2 + \text{tr}(\Sigma + \Sigma_w - 2(\Sigma\Sigma_w)^{1/2})$

GAN - Qualitative results 1/2

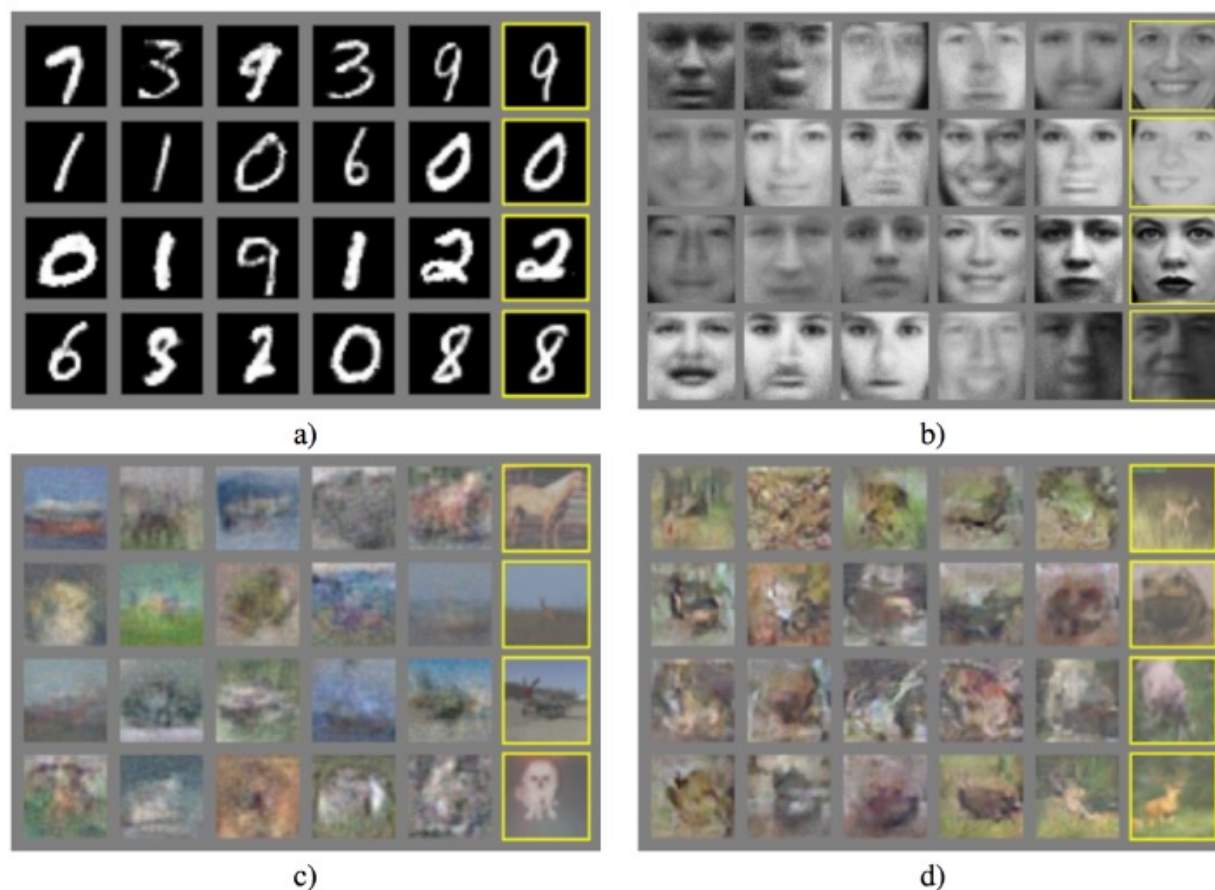


Figure: Right col nearest from dataset. a) MNIST, b) TFD, c) CIFAR-10 (fully connected), d) CIFAR-10 (convolutional D , deconvolutional G)

GAN - Qualitative results 2/2



Figure: Linear interpolation between 2 points in z space

- **Advantages:**
 - ▶ Computational advantages (no complex likelihood inference)
 - ▶ Can represent sharper distributions
- **Disadvantages:**
 - ▶ G and D must be well synchronized for the algorithm to converge correctly

GAN architectures

- How to improve result quality?
 - Spatial resolution
 - ⇒ Cascade of GAN
 - Object quality
 - ⇒ Progressive growing of spatial resolution in G

Generative models

Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 - 2. LaPGAN**

Generative models

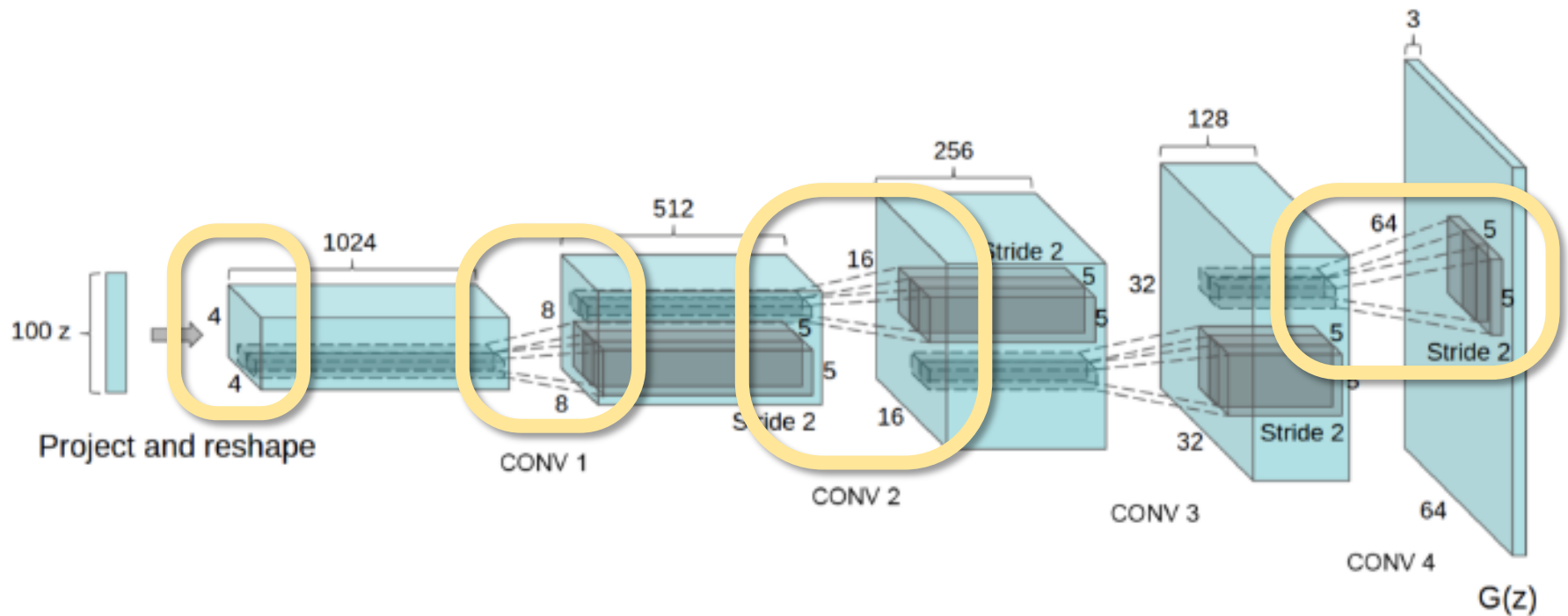
Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. **DCGAN**

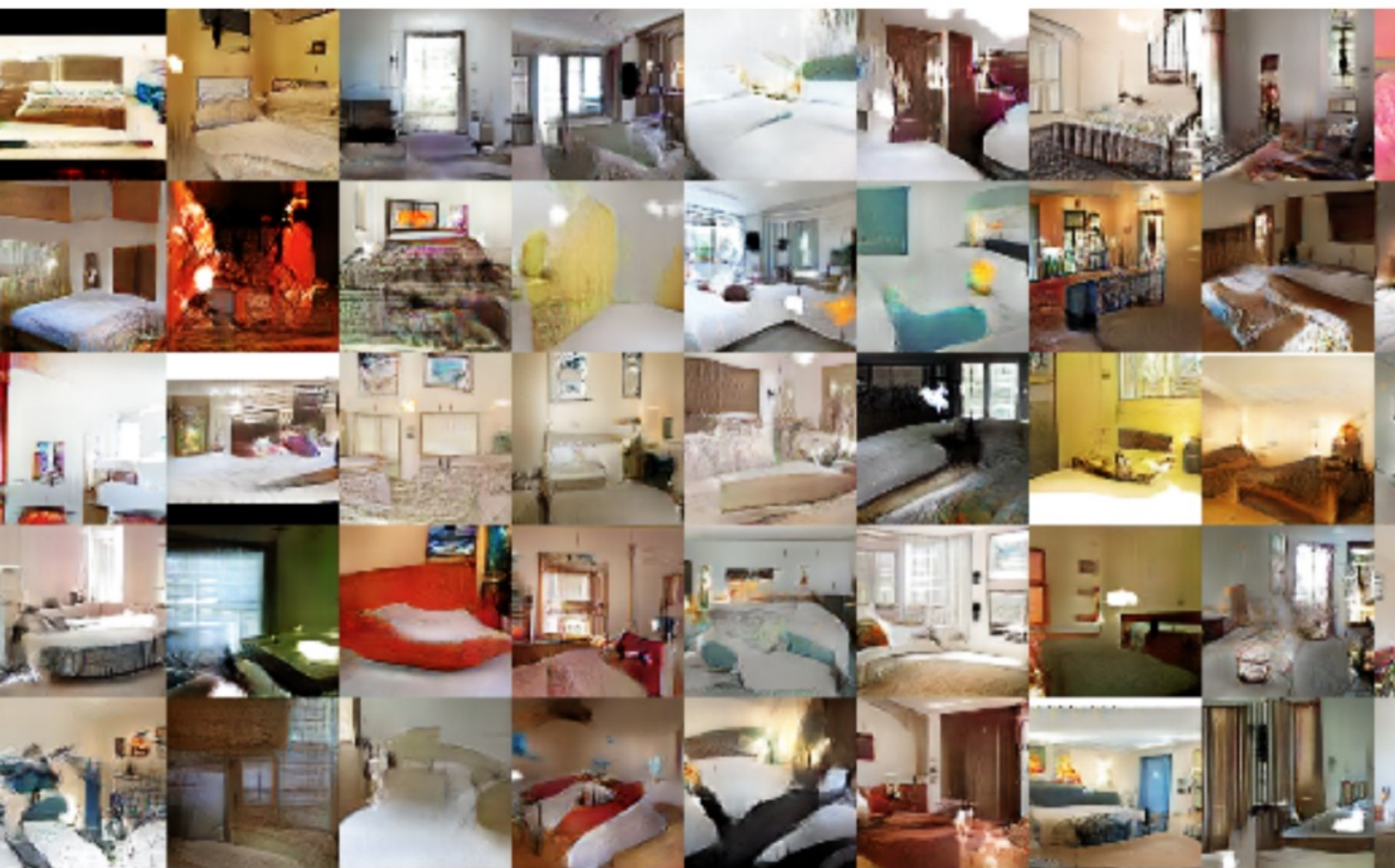
Progressive growing of spatial resolution in G: DCGAN

Upsampling step by step

Combine with convolutional layers



DCGAN - Results - generated bedrooms



DCGAN results - Faces



Generative models

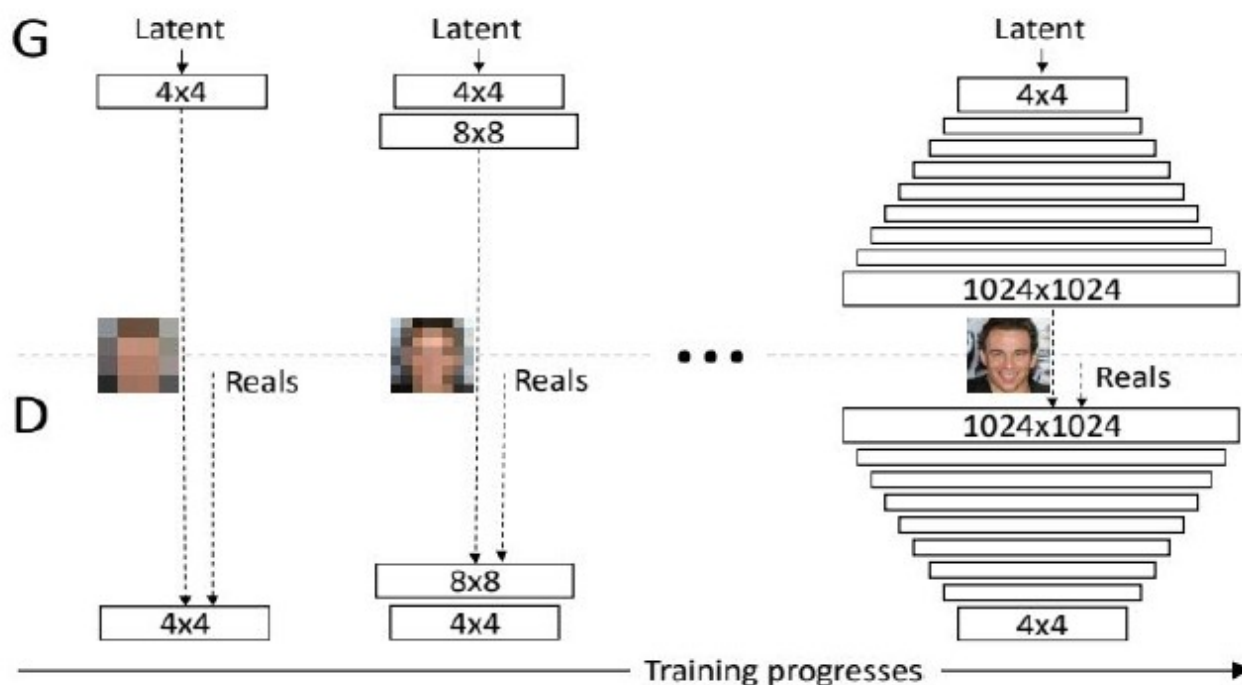
Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. DCGAN
 4. **ProGAN**

Progressive growing of GANs [ICLR 2018]

Combine idea of LAPGAN (several output reso) and DCGAN (archi prog growing)

1. First, start with training 4x4 output images.
2. When this training has converged, add a new block to generate 8x8 output images.
3. Etc.



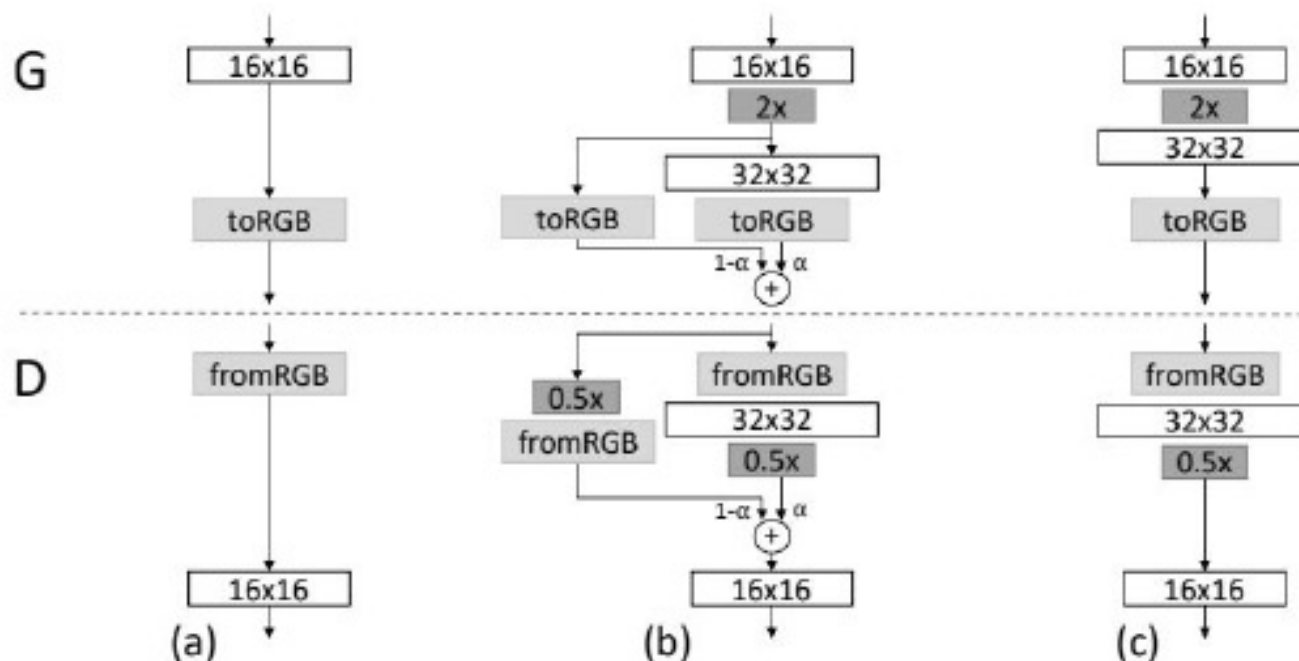
Progressive growing of GANs [ICLR 2018]

Gradual addition of new blocks

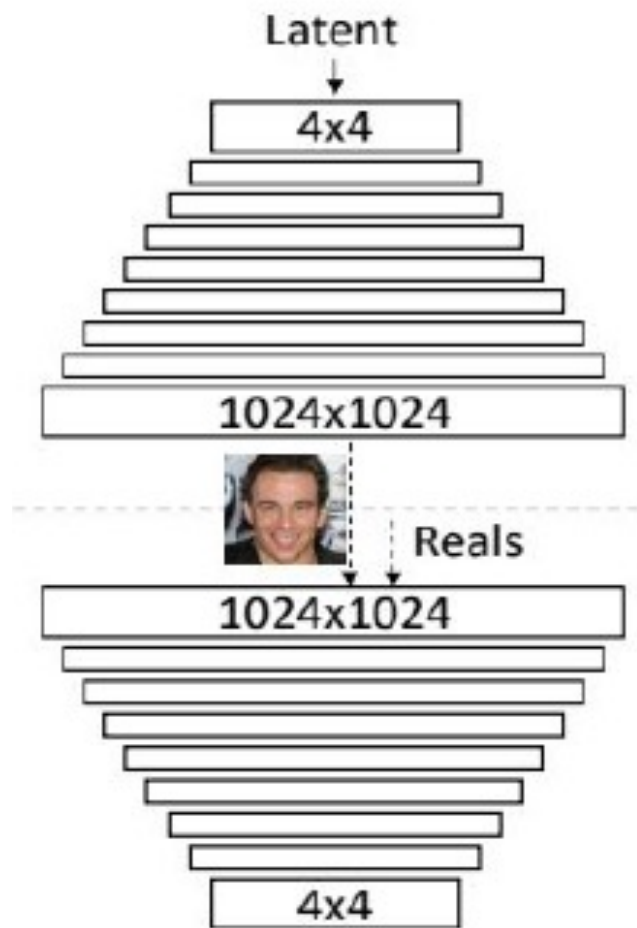
Starting with more weight on the (upsampled) output of the previous block, and then add more and more weights to the output of the current block

All weights trainable during the whole process

Discriminator = mirror image of generator



Progressive growing of GANs [ICLR 2018]



Generative models

Outline

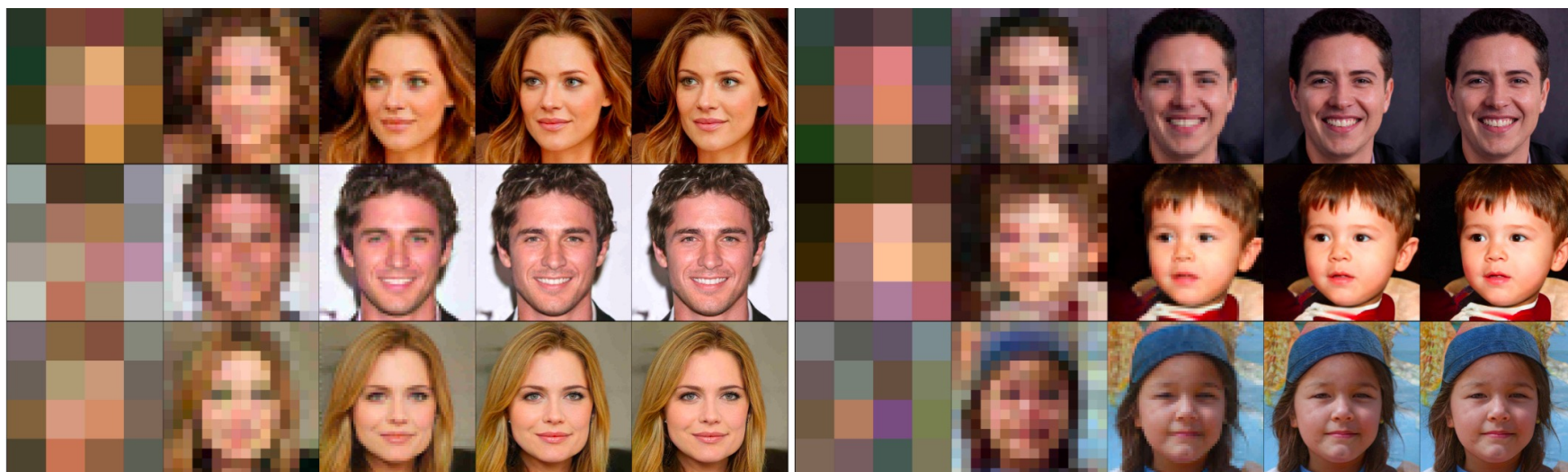
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. DCGAN
 4. ProGAN
 5. **MSG-GAN**

MSG-GAN: Multi-Scale Gradients for Generative Adversarial Networks [CVPR 2020]

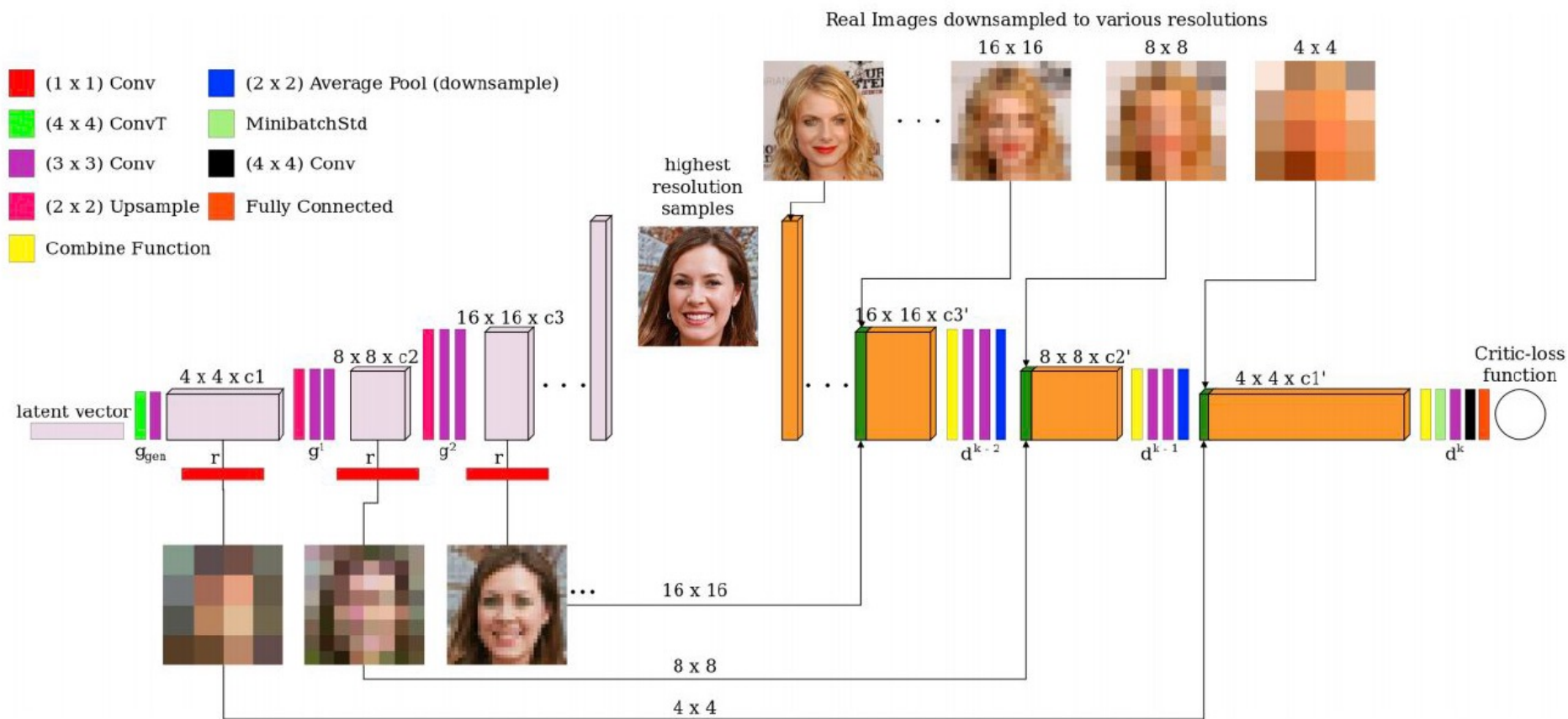
Main Idea:

- ProGAN both use progressive growing, but although this gives stability, it introduces many complicated training parameters associated with each new network.
- Training cannot be done “out of the box”, have to adjust parameters for each new dataset.

→ Train all at once without complicated adding on layers



MSG-GAN: Multi-Scale Gradients for Generative Adversarial Networks [CVPR 2020]



MSG-GAN: results – Random generated CelebA-HQ Faces at resolution 1024x1024



Generative models

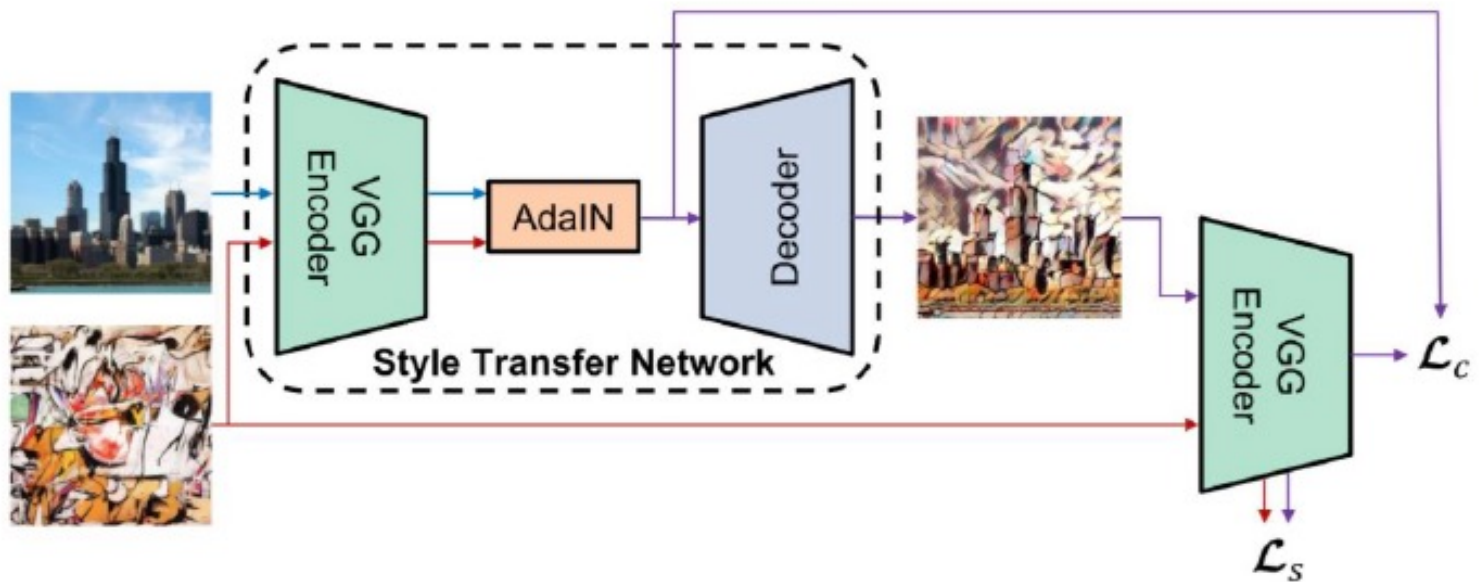
Outline

1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
 1. Basics
 2. LaPGAN
 3. DCGAN
 4. ProGAN
 5. MSG-GAN
 6. **StyleGAN**

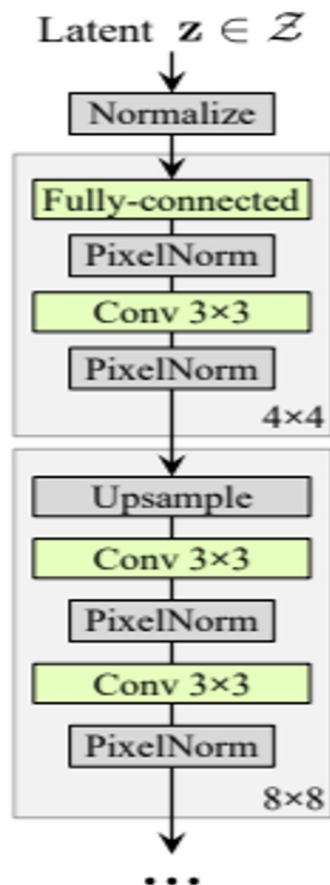
StyleGAN: A Style-Based Generator Architecture for Generative Adversarial Networks [Karras CVPR 2019]

Still progressive growing architecture but with new refinement block based on: Arbitrary Style Transfer in Real-time with Adaptive Instance Normalization (AdaIN)

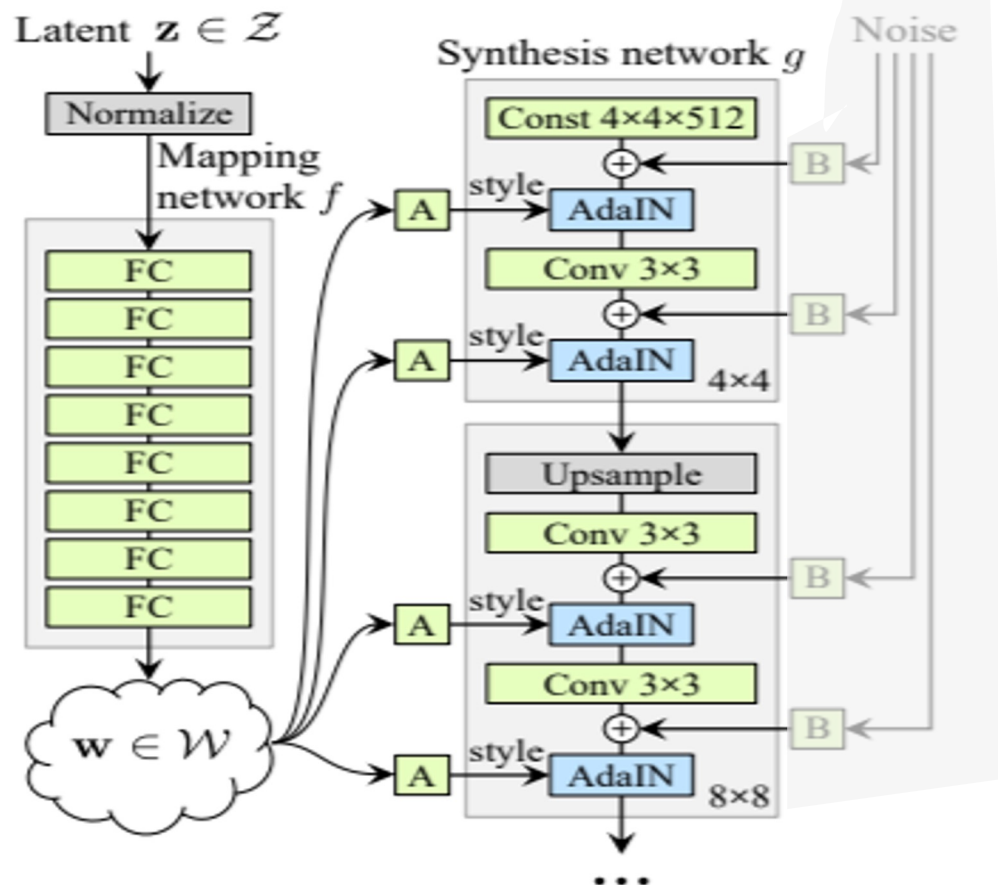
$$\text{AdaIN}(x, y) = \sigma(y) \left(\frac{x - \mu(x)}{\sigma(x)} \right) + \mu(y)$$



StyleGAN Network Architecture



(a) Traditional

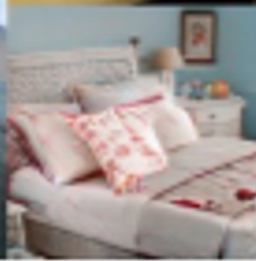
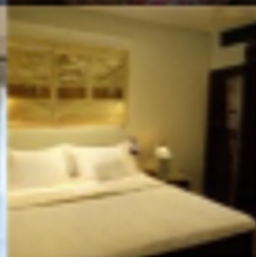
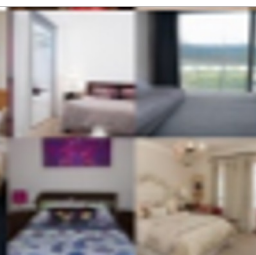
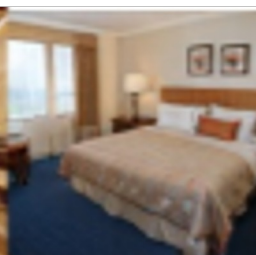
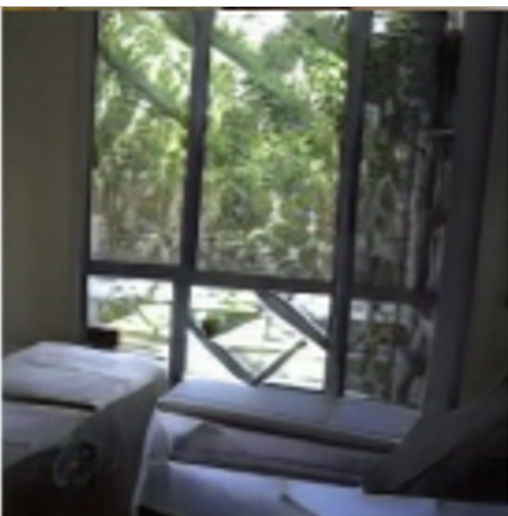


(b) Style-based generator

Building up the Model

Method	CelebA-HQ	FFHQ
A Baseline Progressive GAN [29]	7.79	8.04
B + Tuning (incl. bilinear up/down)	6.11	5.25
C + Add mapping and styles	5.34	4.85
D + Remove traditional input	5.07	4.88
E + Add noise inputs	5.06	4.42
F + Mixing regularization	5.17	4.40







Generative models

Outline

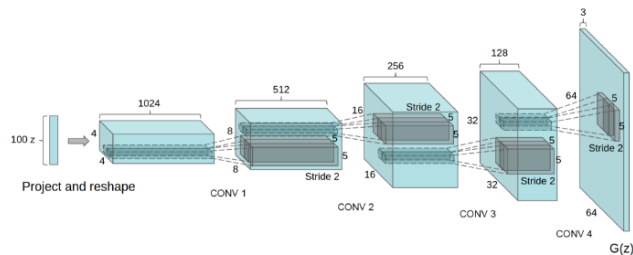
1. Preview: Auto-Encoders, VAE
2. Generative models with GAN
3. GAN architectures
- 4. Editing**

GAN editing

Editing by manipulation in the latent space of z

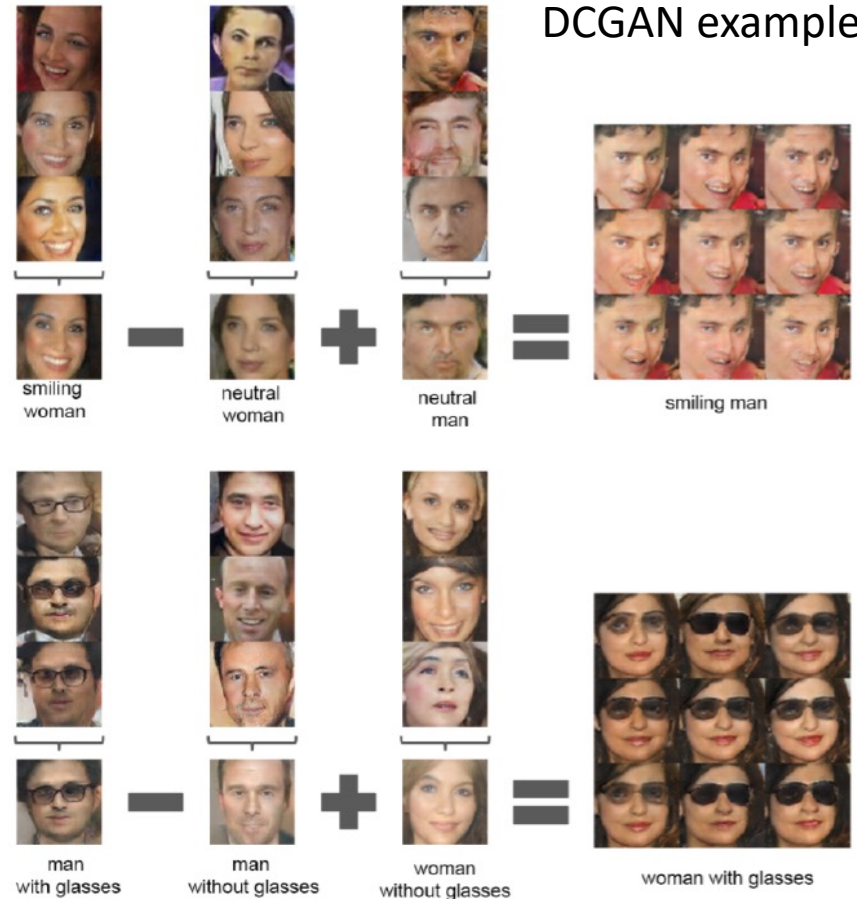
Image $I = G(z)$

Editing: changing z to z' and the new image is $I' = G(z')$



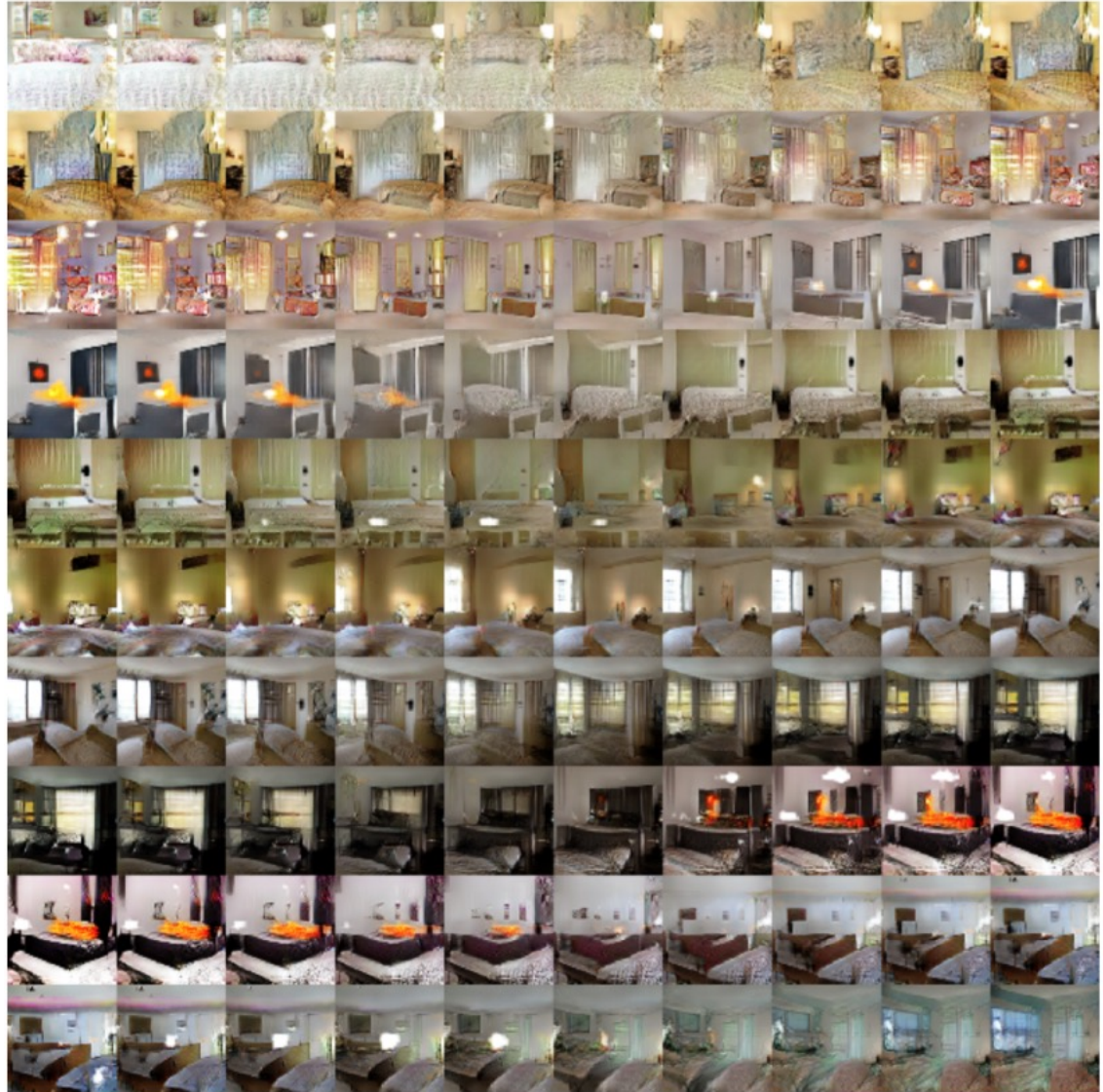
Arithmetics in latent space: vector mean, addition, subtraction

DCGAN examples



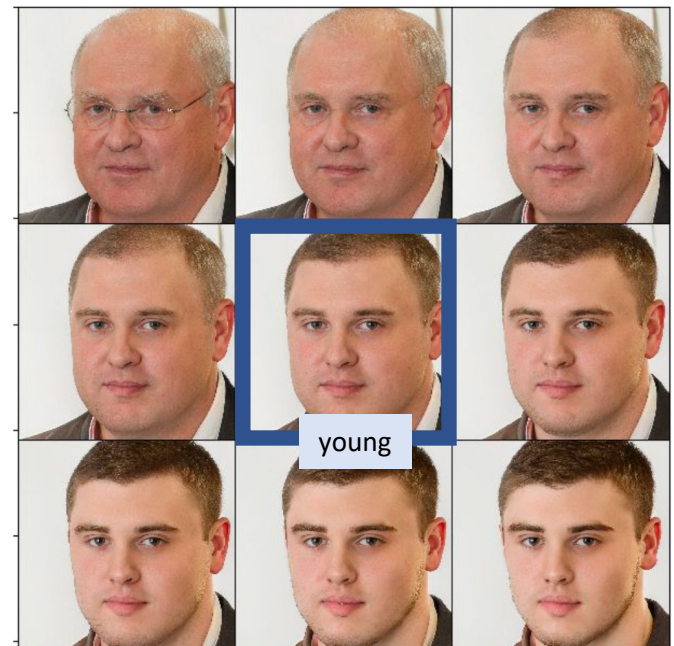
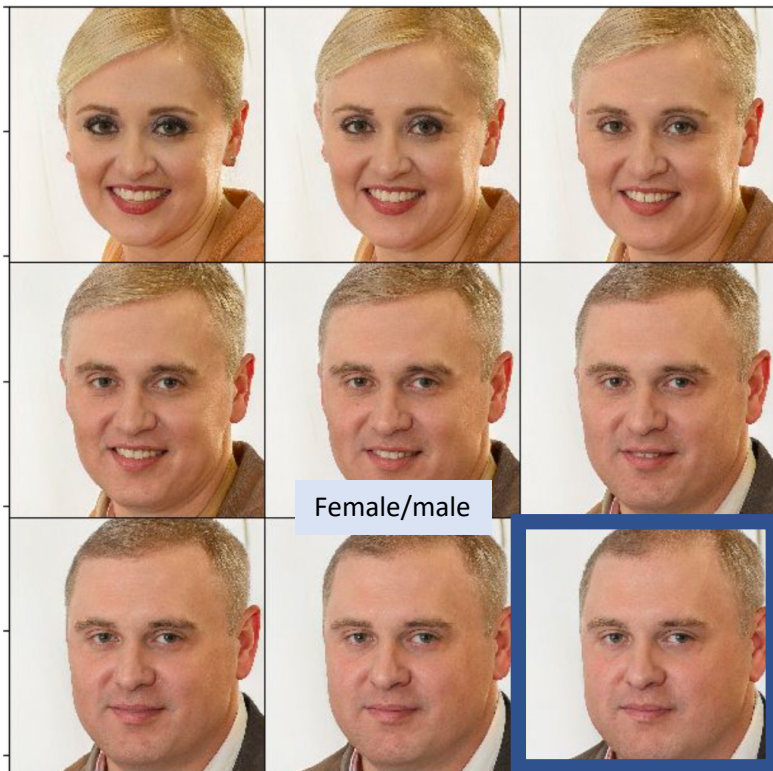
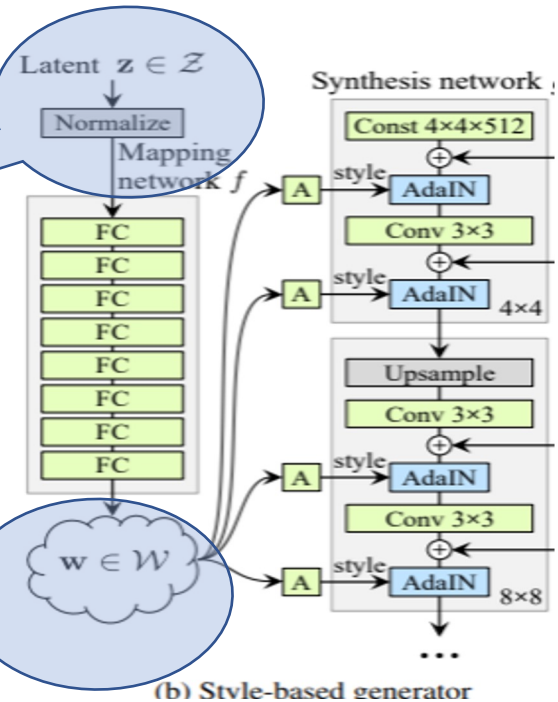
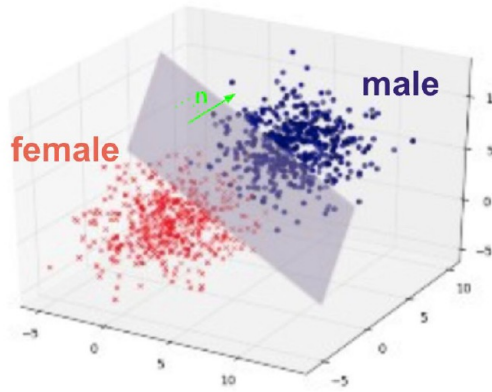
GAN editing

Linear interpolation
in latent space



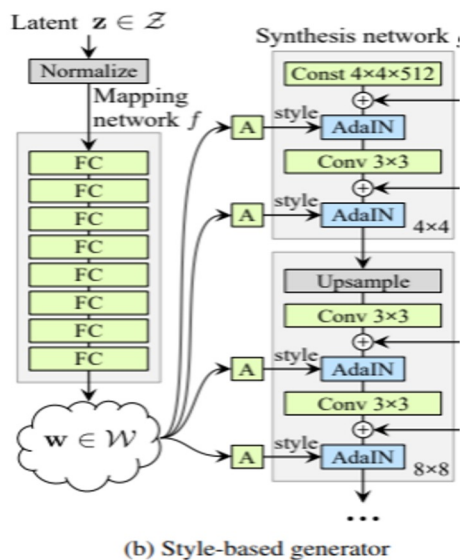
Gan Editing

Latent space analysis for GAN editing with StyleGAN



GAN editing

Latent space analysis for GAN editing with StyleGAN



Two sets of images were generated from their respective latent codes (sources A and B); the rest of the images were generated by copying a specified subset of styles from source B and taking the rest from source A. Copying the styles corresponding to coarse spatial resolutions (4×4 – 8×8) brings high-level aspects such as pose, general hair style, face shape, and eyeglasses from source B, while all colors (eyes, hair, lighting) and finer facial features resemble A. If we instead copy the styles of middle resolutions (16×16 – 32×32) from B, we inherit smaller scale facial features, hair style, eyes open/closed from B, while the pose, general face shape, and eyeglasses from A are preserved.

