

NPFL114, Lecture 10

Deep Generative Models



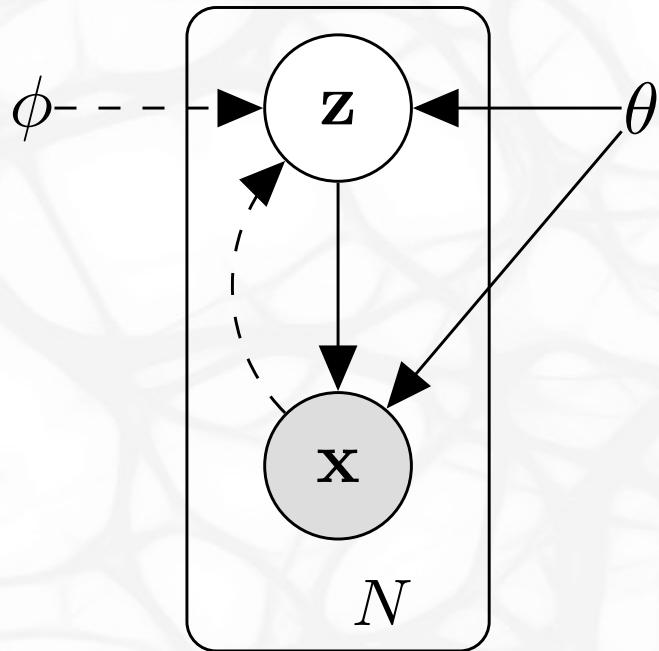
Milan Straka

Generative Models

Generative models are given a set \mathcal{X} of realizations of a random variable \mathbf{x} and their goal is to estimate $P(\mathbf{x})$.

Usually the goal is to be able to sample from $P(\mathbf{x})$, but sometimes an explicit calculation of $P(\mathbf{x})$ is also possible.

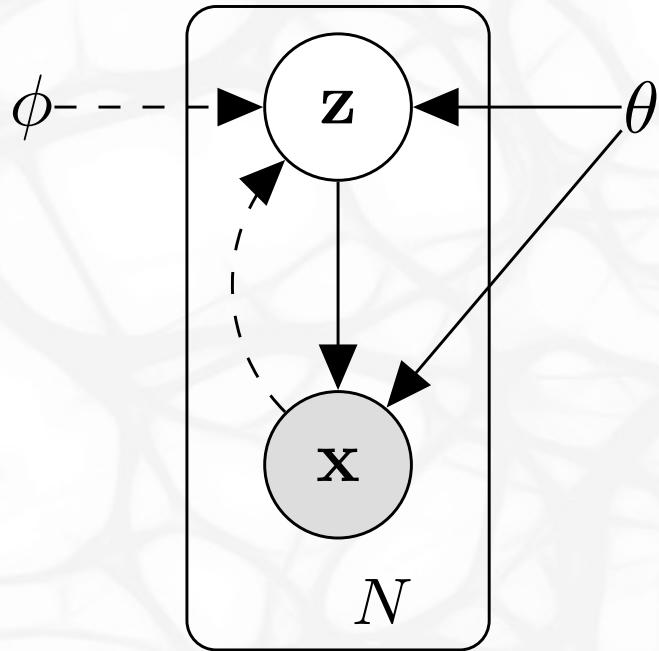
Deep Generative Models



One possible approach to estimate $P(\mathbf{x})$ is to assume that the random variable \mathbf{x} depends on a *latent variable* \mathbf{z} :

$$P(\mathbf{x}) = P(\mathbf{z})P(\mathbf{x}|\mathbf{z}).$$

Deep Generative Models

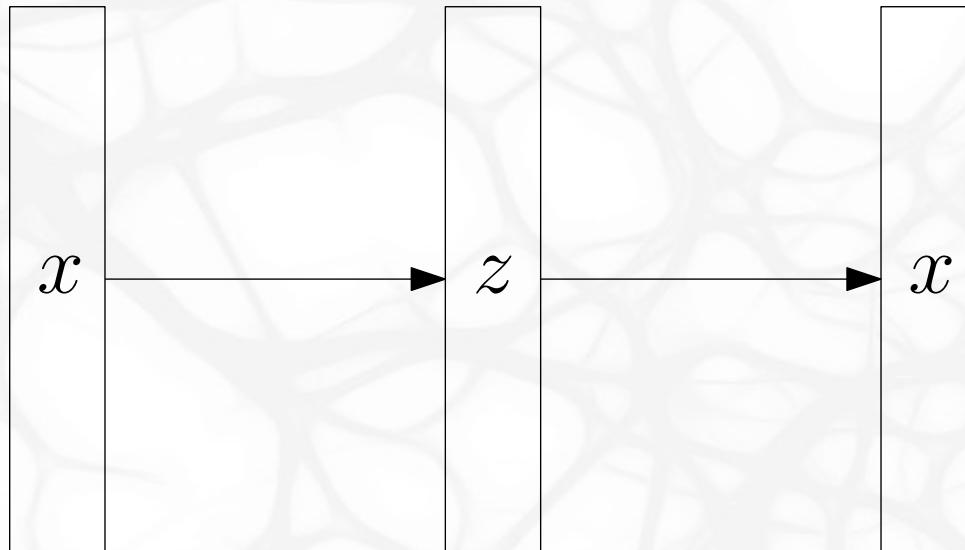


One possible approach to estimate $P(\mathbf{x})$ is to assume that the random variable \mathbf{x} depends on a *latent variable* \mathbf{z} :

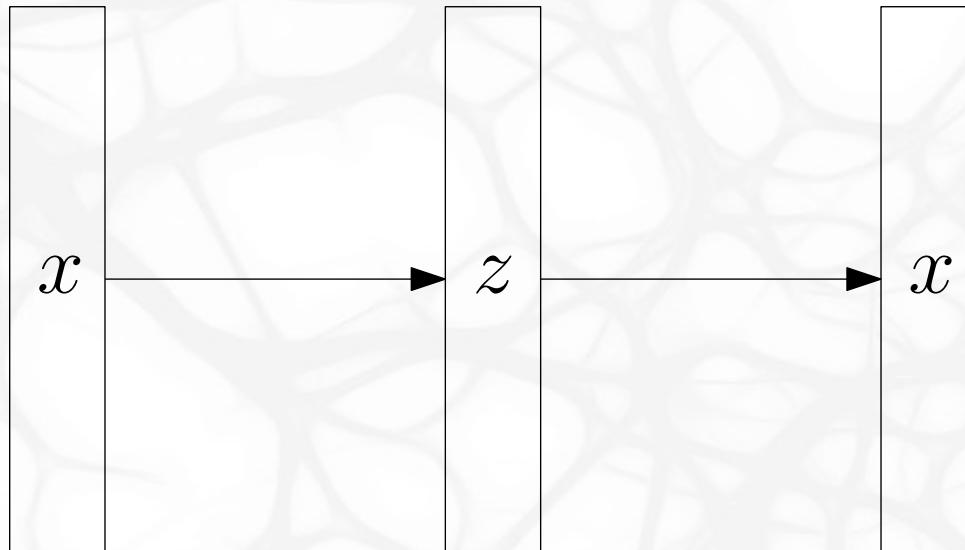
$$P(\mathbf{x}) = P(\mathbf{z})P(\mathbf{x}|\mathbf{z}).$$

We will use neural networks to estimate the conditional probability with $P_{\theta}(\mathbf{x}|\mathbf{z})$.

AutoEncoders



AutoEncoders



- unsupervised feature extraction
- input compression
- when $x + \varepsilon$ is used as input, autoencoders can perform denoising

Variational AutoEncoders

We assume $P(\mathbf{z})$ is fixed and independent on \mathbf{x} .

We will approximate $P(\mathbf{x}|\mathbf{z})$ using $P_{\theta}(\mathbf{x}|\mathbf{z})$. However, in order to compute $P_{\theta}(\mathbf{z})$, we need to know $P_{\theta}(\mathbf{z}|\mathbf{x})$, which is usually intractable.

Variational AutoEncoders

We assume $P(\mathbf{z})$ is fixed and independent on \mathbf{x} .

We will approximate $P(\mathbf{x}|\mathbf{z})$ using $P_{\theta}(\mathbf{x}|\mathbf{z})$. However, in order to compute $P_{\theta}(\mathbf{z})$, we need to know $P_{\theta}(\mathbf{z}|\mathbf{x})$, which is usually intractable.

We will therefore approximate $P_{\theta}(\mathbf{z}|\mathbf{x})$ by a trainable $Q_{\varphi}(\mathbf{z}|\mathbf{x})$.

Variational AutoEncoders

Let us define *variational lower bound* or *evidence lower bound* (ELBO), denoted $\mathcal{L}(\theta, \varphi; \mathbf{x})$, as

$$\mathcal{L}(\theta, \varphi; \mathbf{x}) = \mathbb{E}_{Q_\varphi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}, z)] - \log Q_\varphi(z|\mathbf{x}).$$

Variational AutoEncoders

Let us define *variational lower bound* or *evidence lower bound* (ELBO), denoted $\mathcal{L}(\theta, \varphi; \mathbf{x})$, as

$$\mathcal{L}(\theta, \varphi; \mathbf{x}) = \mathbb{E}_{Q_\varphi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}, z)] - \log Q_\varphi(z|\mathbf{x}).$$

It can then be checked that

$$\begin{aligned}\mathcal{L}(\theta, \varphi; \mathbf{x}) &= \log P_\theta(\mathbf{x}) - D_{\text{KL}}(Q_\varphi(z|\mathbf{x}) || P_\theta(z)) \\ &\leq \log P_\theta(\mathbf{x}),\end{aligned}$$

where the last inequality follows from KL-divergence being non-negative.

Variational AutoEncoders

Let us define *variational lower bound* or *evidence lower bound* (ELBO), denoted $\mathcal{L}(\theta, \varphi; \mathbf{x})$, as

$$\mathcal{L}(\theta, \varphi; \mathbf{x}) = \mathbb{E}_{Q_\varphi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}, z)] - \log Q_\varphi(z|\mathbf{x}).$$

It can then be checked that

$$\begin{aligned}\mathcal{L}(\theta, \varphi; \mathbf{x}) &= \log P_\theta(\mathbf{x}) - D_{\text{KL}}(Q_\varphi(z|\mathbf{x}) || P_\theta(z)) \\ &\leq \log P_\theta(\mathbf{x}),\end{aligned}$$

where the last inequality follows from KL-divergence being non-negative.

The variational lower bound can also be written as

$$\mathcal{L}(\theta, \varphi; \mathbf{x}) = \mathbb{E}_{Q_\varphi(z|\mathbf{x})} [\log P_\theta(\mathbf{x}|z)] - D_{\text{KL}}(Q_\varphi(z|\mathbf{x}) || P_\theta(z)).$$

Variational AutoEncoders

In order to derivate through $\mathbf{z} \sim Q_\varphi(\mathbf{z}|\mathbf{x})$, note that if

$$\mathbf{z} \sim \mathcal{N}(\mu, \sigma^2),$$

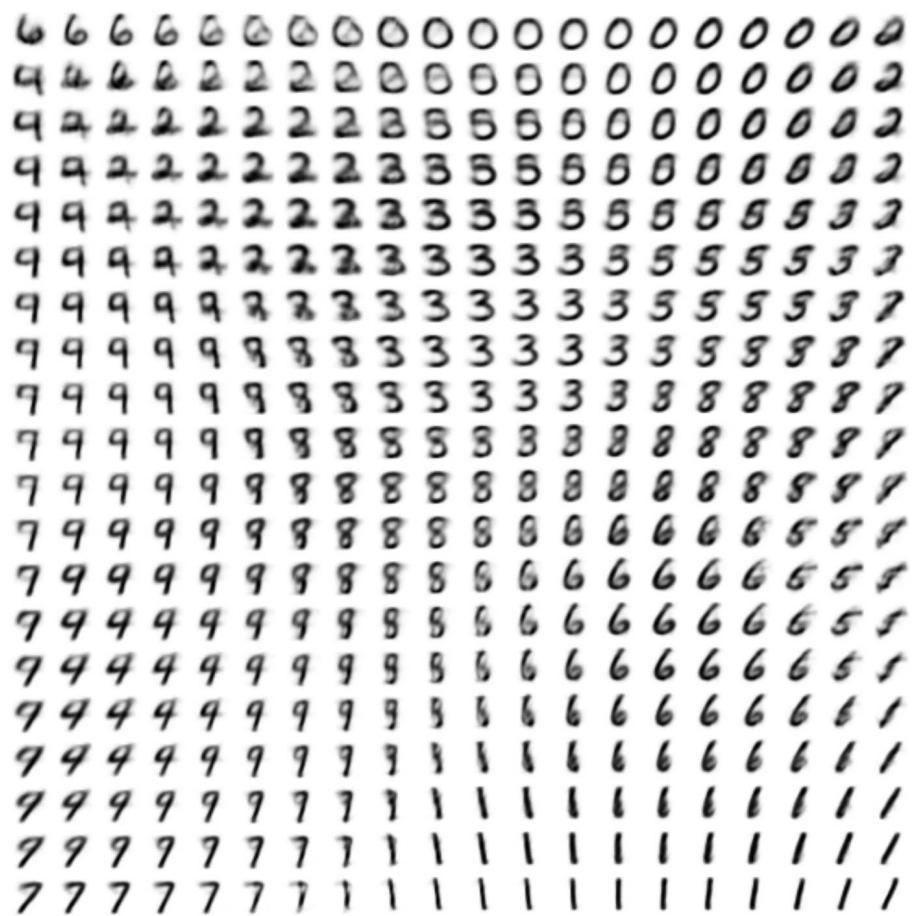
we can write \mathbf{z} as

$$\mathbf{z} \sim \mu + \sigma \cdot \mathcal{N}(0, 1).$$

Variational AutoEncoders



(a) Learned Frey Face manifold



(b) Learned MNIST manifold

Variational AutoEncoders

8 6 1 7 8 1 4 8 2 8
9 6 8 3 9 6 0 3 1 9
1 1 1 1 3 6 9 1 7 9
8 9 0 8 6 9 1 9 6 3
8 2 3 3 3 3 1 3 8 6
6 9 9 8 6 1 6 6 6 6
9 5 2 6 6 5 1 8 9 9
9 9 9 1 3 1 2 8 2 3
0 4 6 1 2 3 2 0 8 8
9 7 5 9 9 3 4 8 5 1

(a) 2-D latent space

4 1 6 5 7 6 7 6 7 2
8 5 9 4 6 8 2 1 6 8
6 1 0 3 2 8 8 4 3 3
2 8 6 8 9 1 0 0 4 1
5 1 9 3 0 1 5 3 5 9
6 6 6 1 4 9 1 7 8 8
1 3 4 3 9 8 3 4 7 0
4 5 8 2 9 7 0 9 5 9
6 9 9 4 8 7 2 3 9 5
2 6 4 5 6 0 9 9 9 8

(b) 5-D latent space

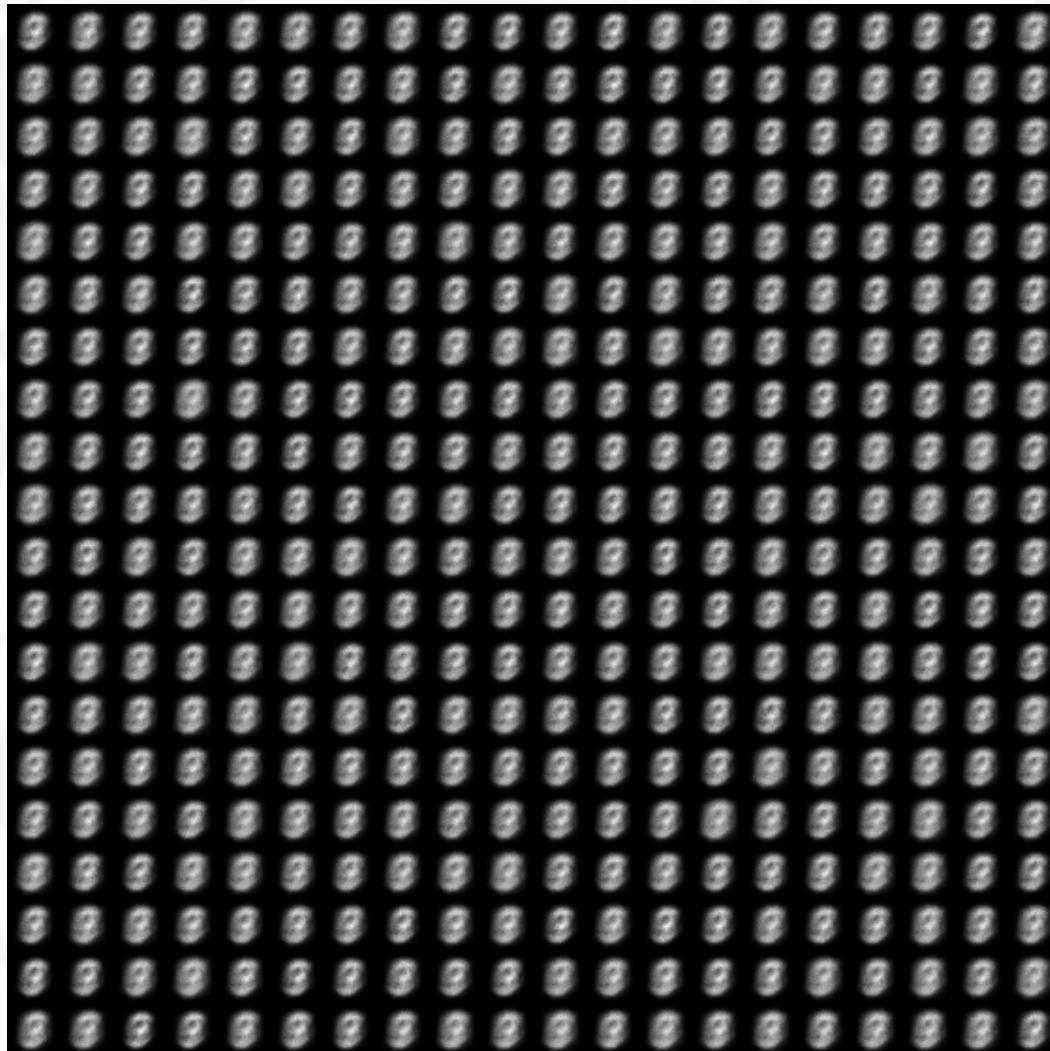
2 8 3 1 3 8 5 7 3 8
8 3 8 2 7 9 3 3 3 8
3 5 9 9 4 3 9 5 1 6
1 9 1 8 8 3 3 1 9 7
2 7 3 6 4 3 0 2 6 3
5 9 7 0 5 9 3 3 7 5
6 9 4 3 6 2 8 5 5 2
8 4 9 0 8 0 7 0 6 6
7 4 3 6 3 0 3 6 0 1
2 1 2 0 9 7 1 0 0 0

(c) 10-D latent space

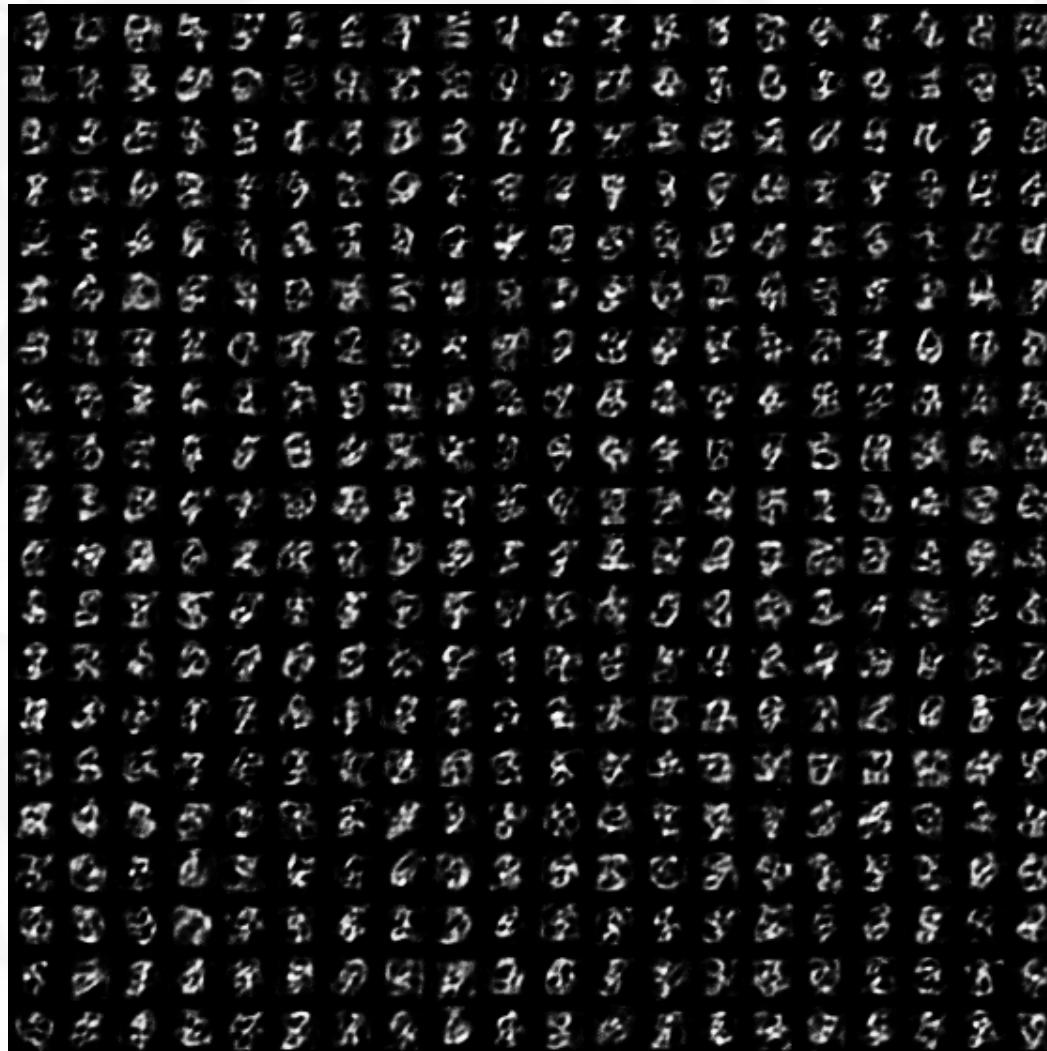
8 2 0 8 9 2 3 9 0 0
7 5 1 9 1 1 7 1 4 4
8 7 6 2 0 3 2 8 2 9
2 9 8 4 3 8 7 0 6 1
5 7 7 9 8 9 9 9 1 0
6 8 8 4 9 4 8 2 8 1
7 5 8 2 4 6 1 3 8 2
7 9 3 9 2 7 9 3 9 0
4 5 2 4 3 9 0 1 8 9
8 8 7 2 3 8 1 6 2 3 6

(d) 20-D latent space

VAE – Too High Latent Loss



VAE – Too High Reconstruction Loss



Generative Adversarial Networks

We have a *generator*, which given $\mathbf{z} \sim P(\mathbf{z})$ generates data \mathbf{x} .

We denote the generator as $G(\mathbf{z}; \theta_g)$.

Generative Adversarial Networks

We have a *generator*, which given $\mathbf{z} \sim P(\mathbf{z})$ generates data \mathbf{x} .

We denote the generator as $G(\mathbf{z}; \theta_g)$.

Then we have a *discriminator*, which given data \mathbf{x} generates a probability whether \mathbf{x} comes from real data or is generated by a generator.

We denote the discriminator as $D(\mathbf{x}; \theta_d)$.

Generative Adversarial Networks

We have a *generator*, which given $\mathbf{z} \sim P(\mathbf{z})$ generates data \mathbf{x} .

We denote the generator as $G(\mathbf{z}; \theta_g)$.

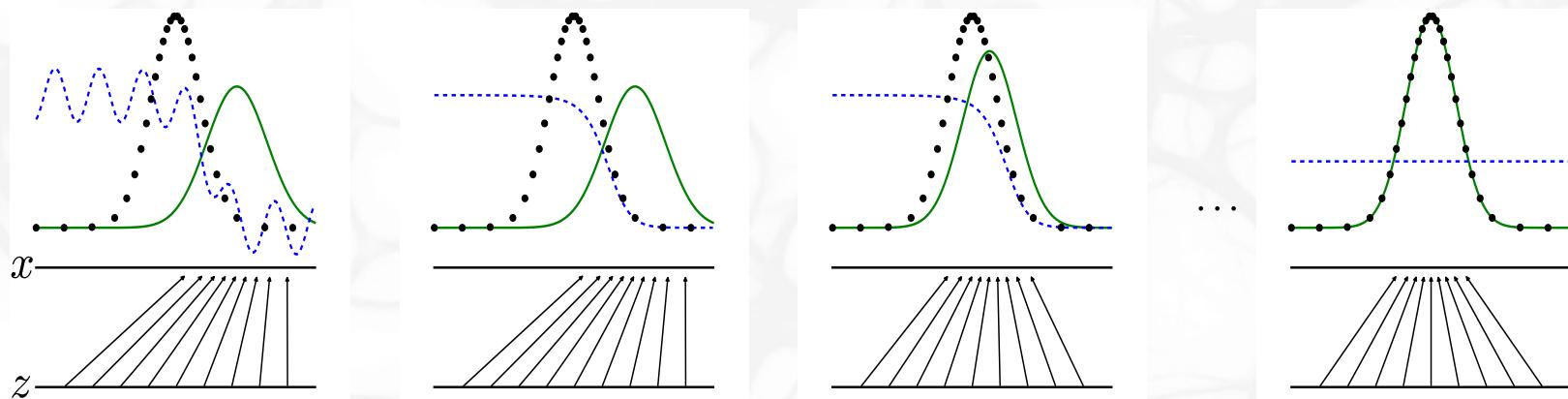
Then we have a *discriminator*, which given data \mathbf{x} generates a probability whether \mathbf{x} comes from real data or is generated by a generator.

We denote the discriminator as $D(\mathbf{x}; \theta_d)$.

In other words, D and G play the following game:

$$\min_G \max_D \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim P(\mathbf{z})} [\log(1 - D(\mathbf{x}))].$$

Generative Adversarial Networks



Generative Adversarial Networks

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

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

end for

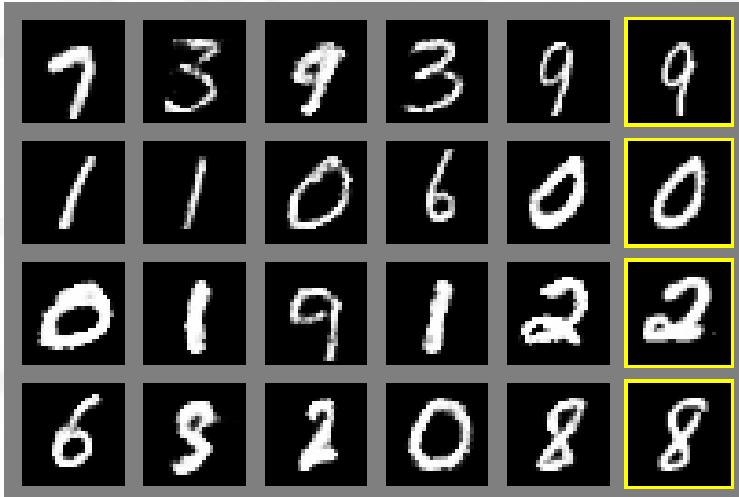
- Sample minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$ from noise prior $p_g(\mathbf{z})$.
- Update the generator by descending its stochastic gradient:

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

end for

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

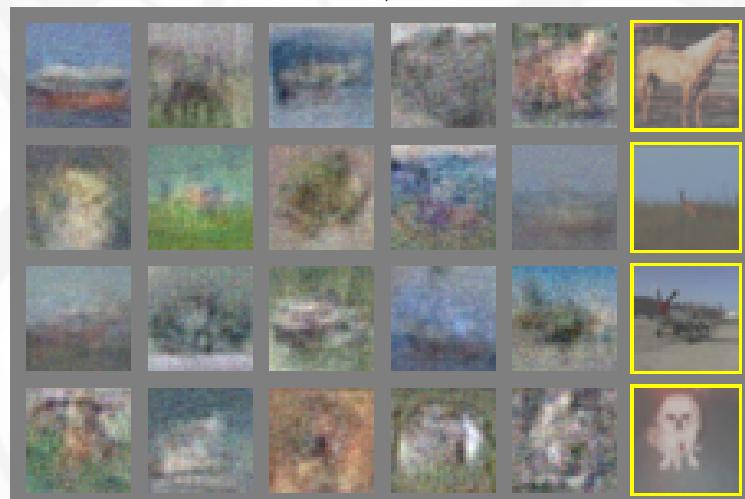
Generative Adversarial Networks



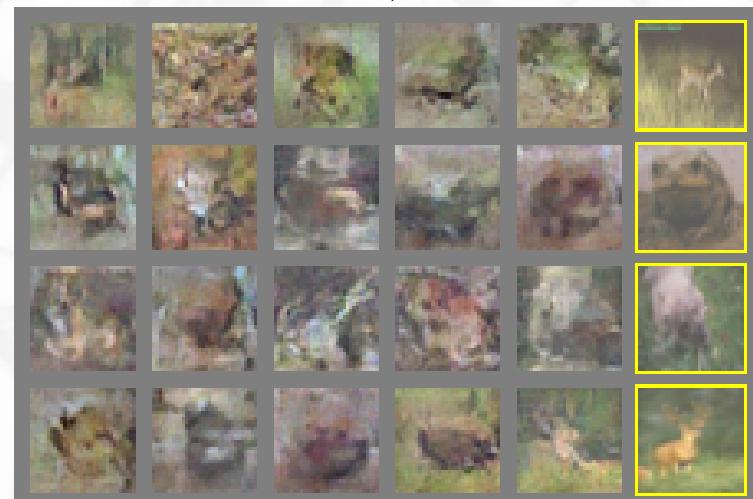
a)



b)

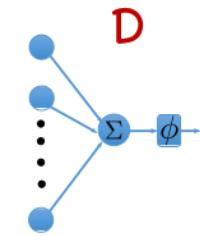
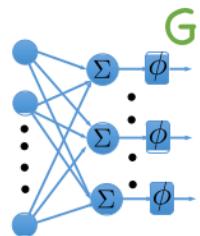


c)

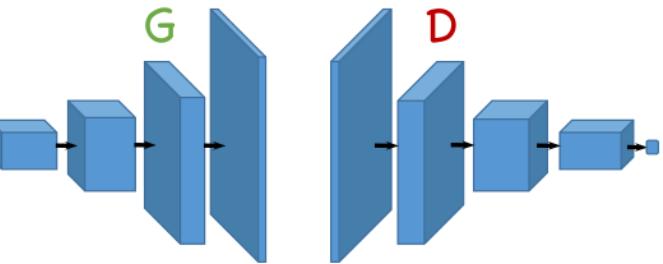


d)

Deep Convolutional GAN

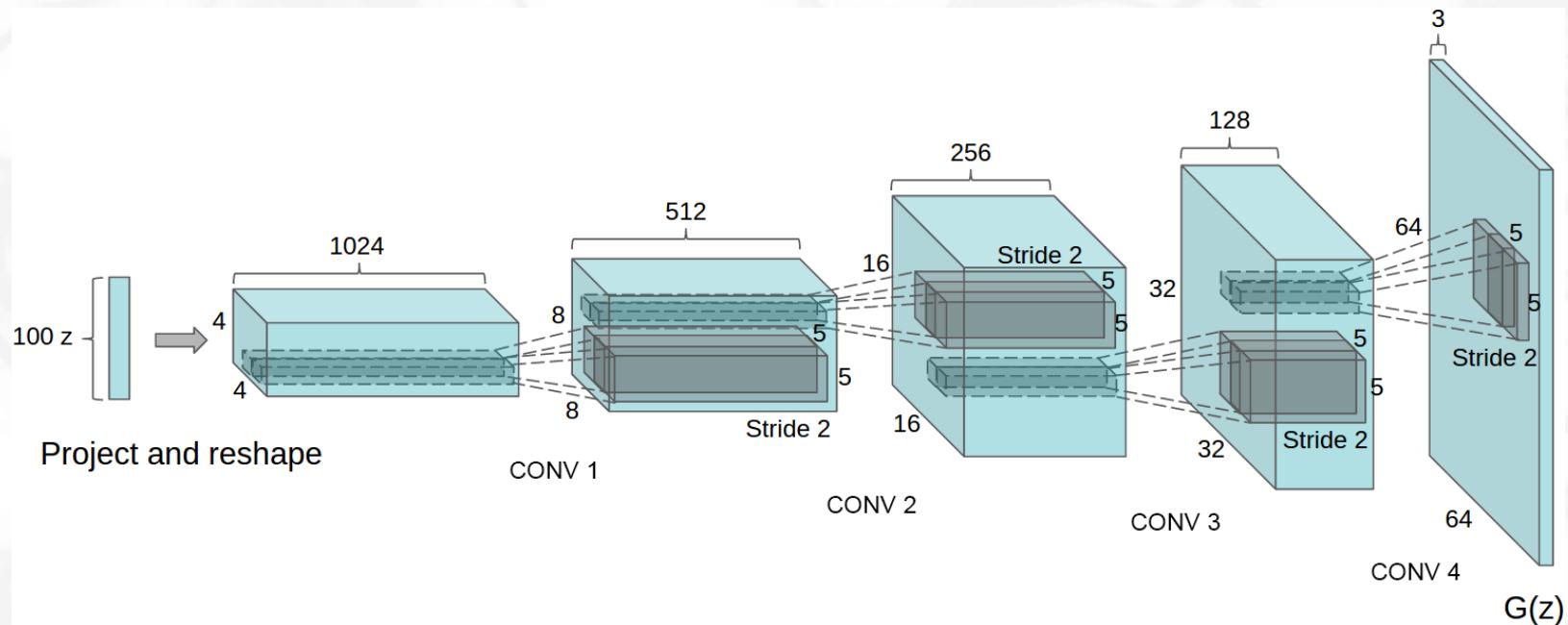


(a)

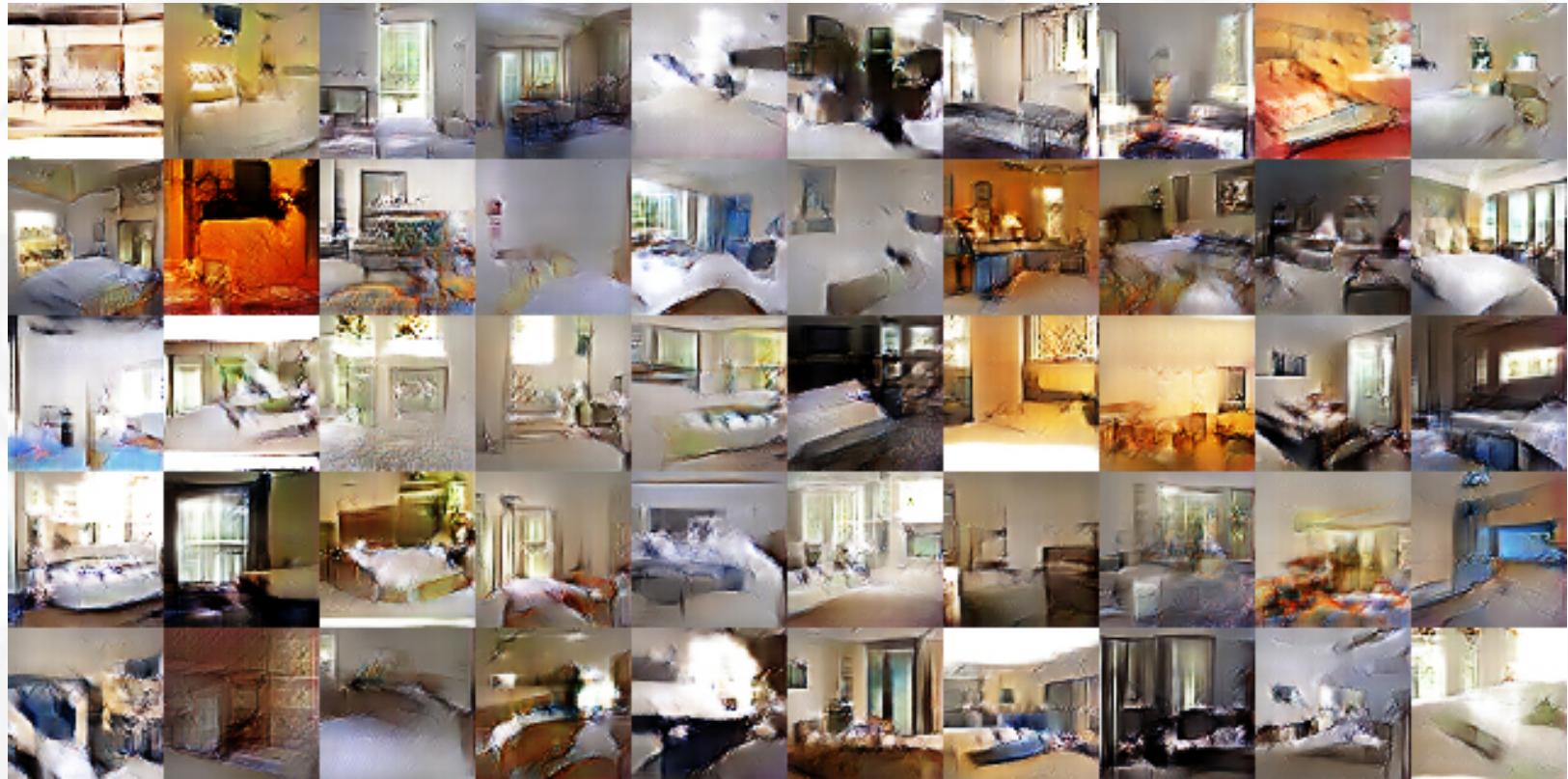


(c)

Deep Convolutional GAN



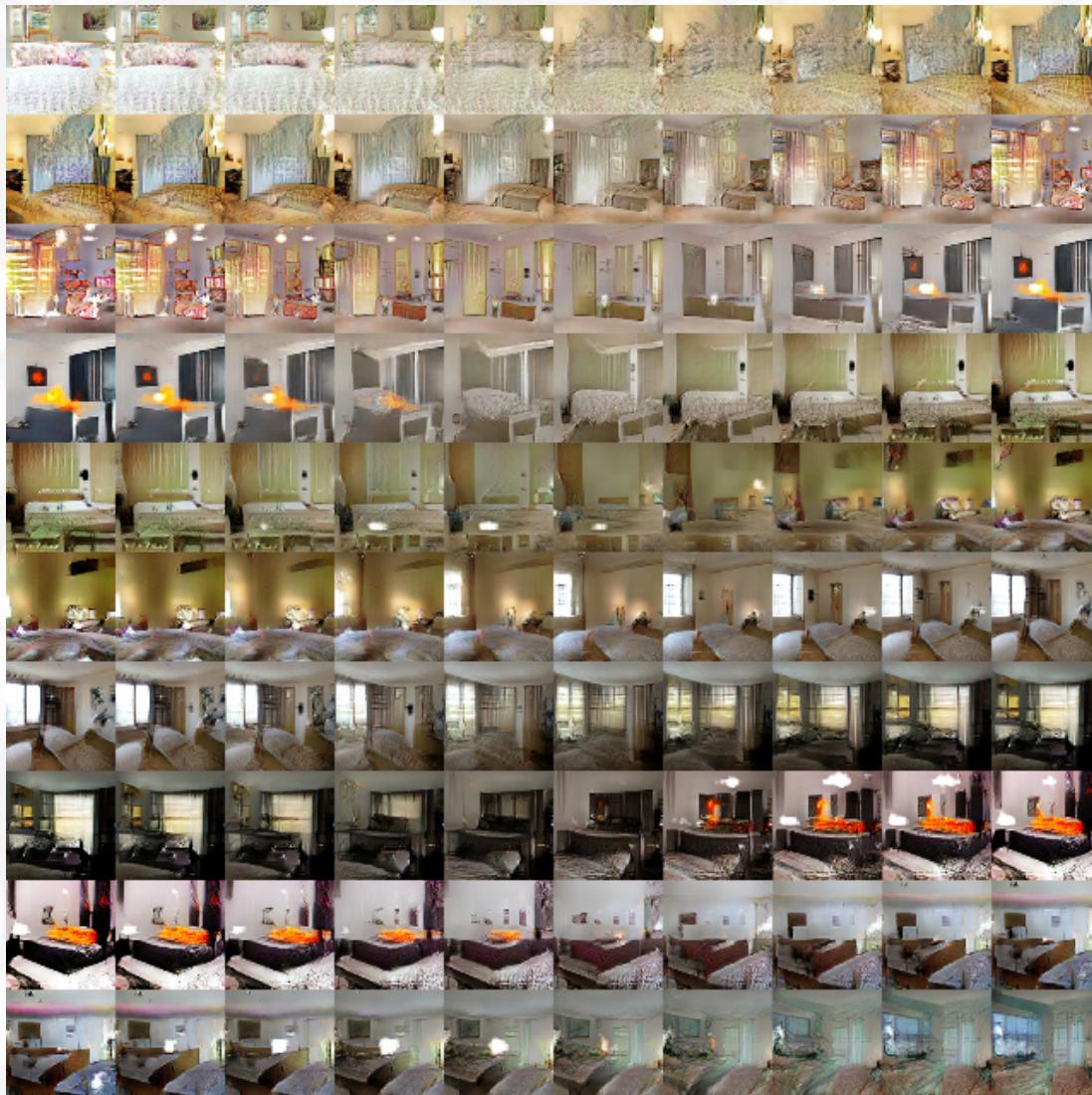
Deep Convolutional GAN



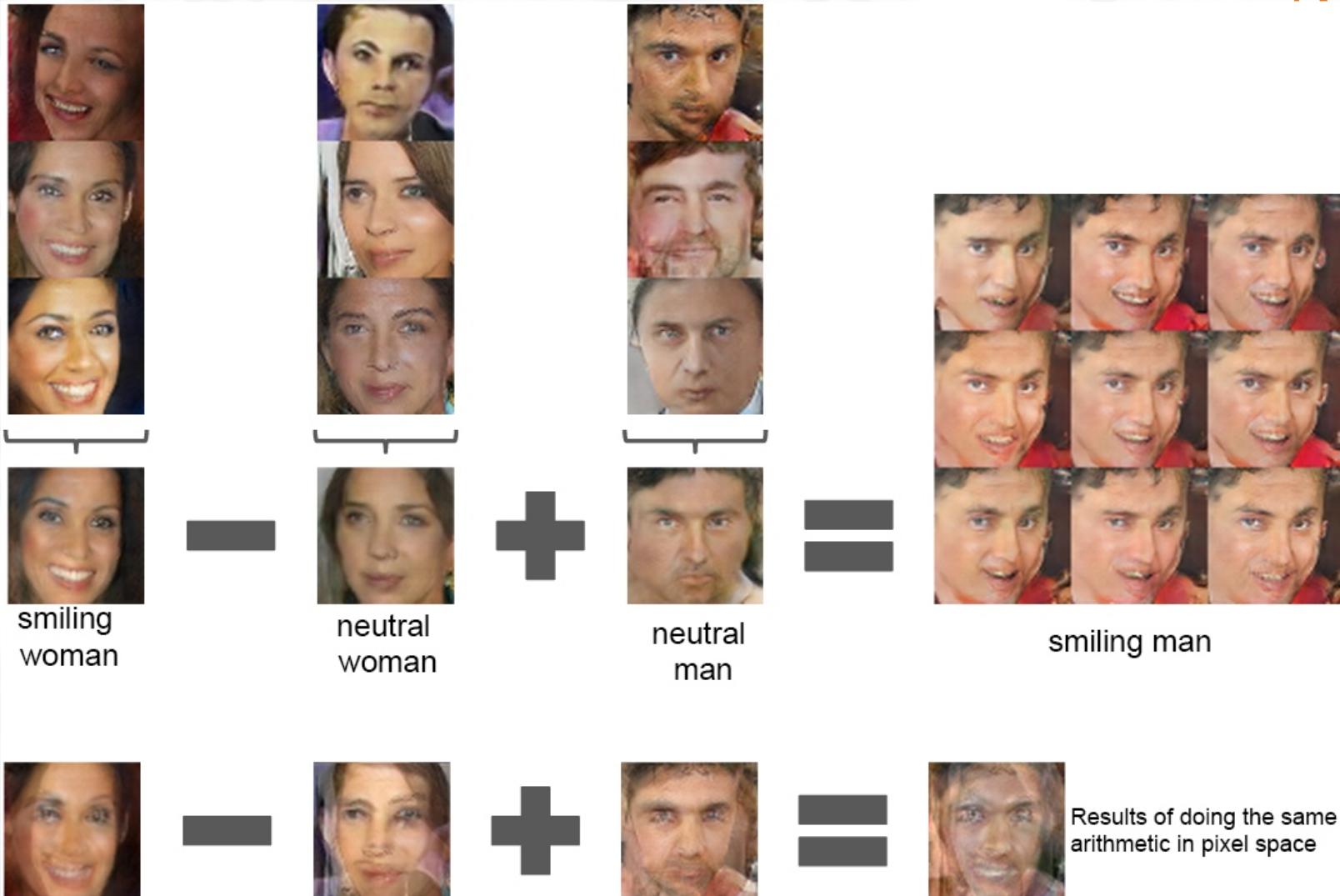
Deep Convolutional GAN



Deep Convolutional GAN



Deep Convolutional GAN



Deep Convolutional GAN



man
with glasses



man
without glasses



woman
without glasses



woman with glasses



Results of doing the same
arithmetic in pixel space

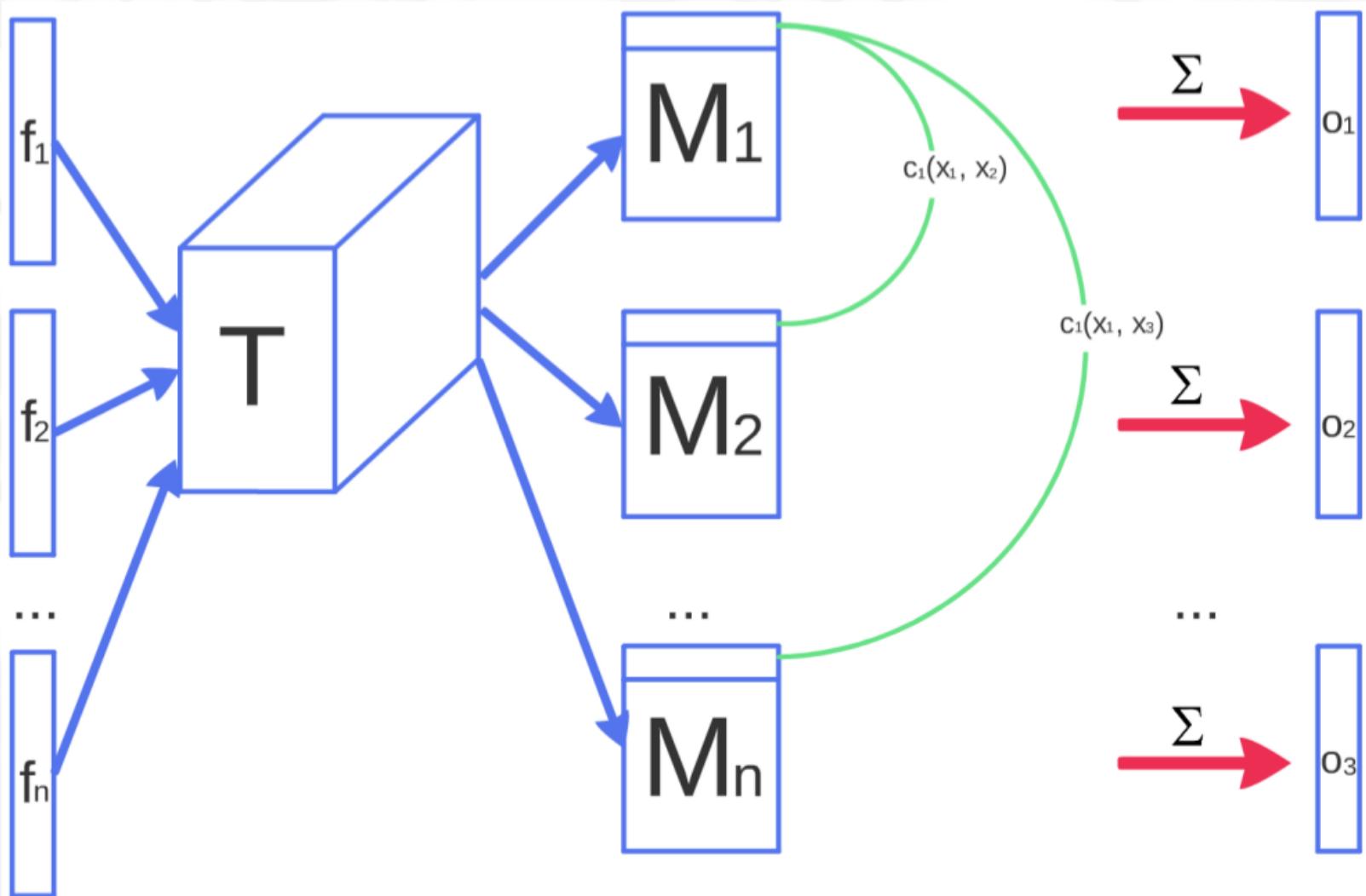
Deep Convolutional GAN



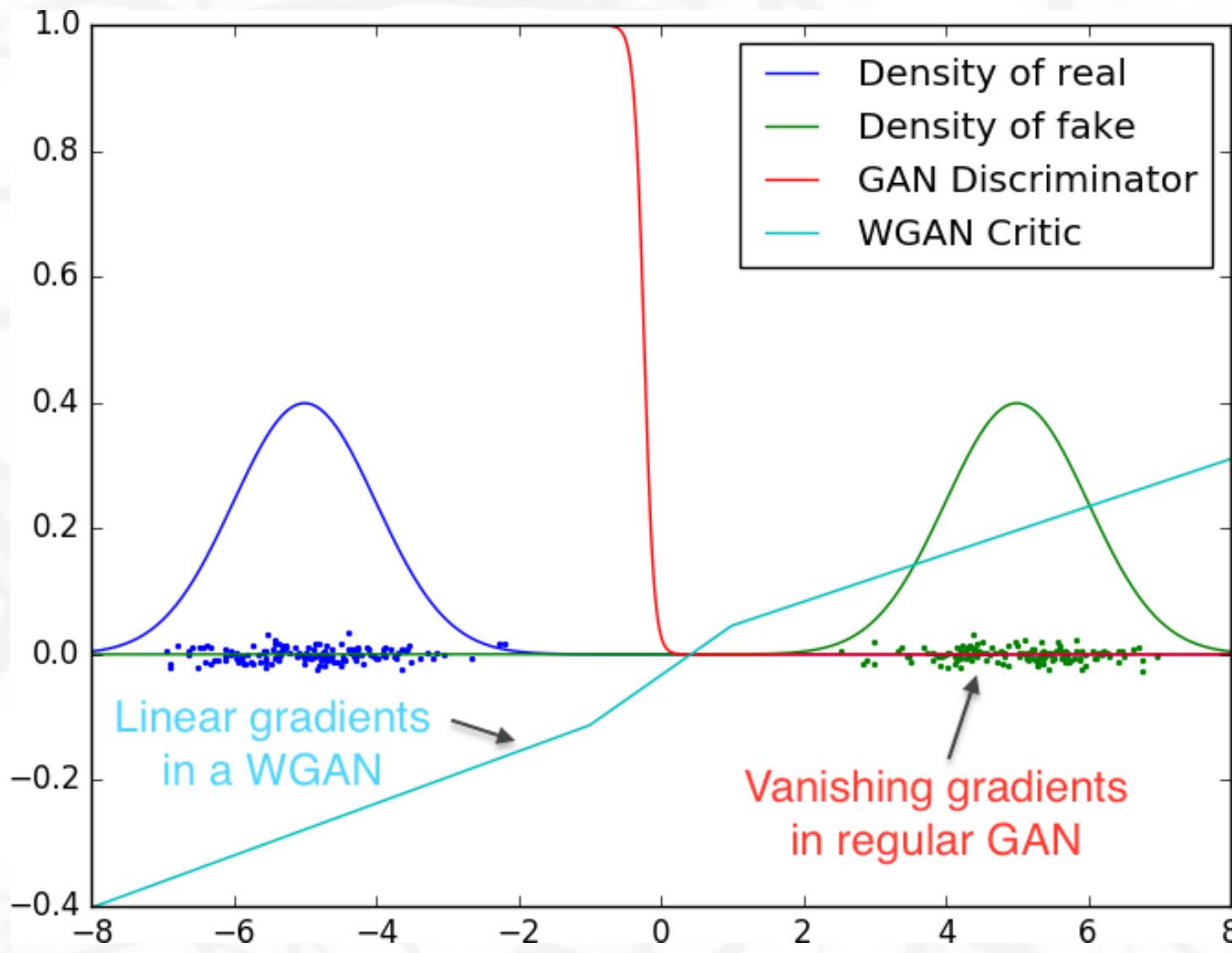
GANs are Problematic to Train

- Feature matching
- Minibatch discrimination
- Historical averaging
- Label smoothing

Minibatch Discrimination



Wasserstein GAN



Wasserstein GAN



Figure 5: Algorithms trained with a DCGAN generator. Left: WGAN algorithm. Right: standard GAN formulation. Both algorithms produce high quality samples.



Figure 6: Algorithms trained with a generator without batch normalization and constant number of filters at every layer (as opposed to duplicating them every time as in [18]). Aside from taking out batch normalization, the number of parameters is therefore reduced by a bit more than an order of magnitude. Left: WGAN algorithm. Right: standard GAN formulation. As we can see the standard GAN failed to learn while the WGAN still was able to produce samples.

Wasserstein GAN



Figure 7: Algorithms trained with an MLP generator with 4 layers and 512 units with ReLU nonlinearities. The number of parameters is similar to that of a DCGAN, but it lacks a strong inductive bias for image generation. Left: WGAN algorithm. Right: standard GAN formulation. The WGAN method still was able to produce samples, lower quality than the DCGAN, and of higher quality than the MLP of the standard GAN. Note the significant degree of mode collapse in the GAN MLP.