

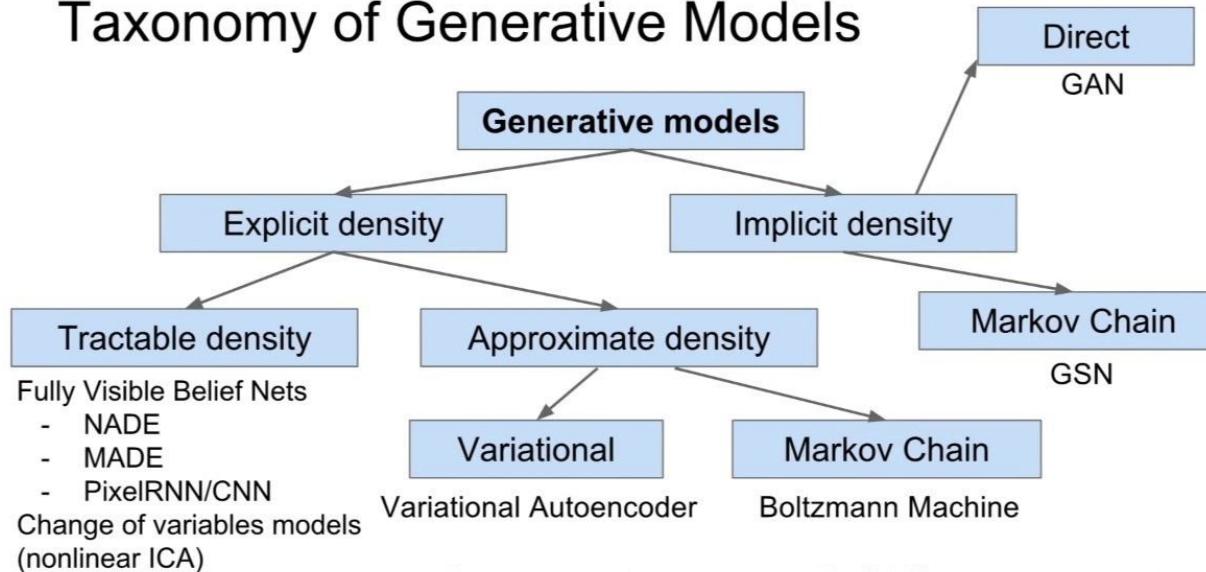
# Deep Learning

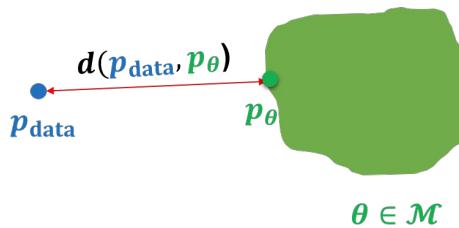
Lecture 12

# Recap

- Discriminative vs Generative models
- Reparametrization tricks
- Variational Auto Encoder

# Taxonomy of Generative Models





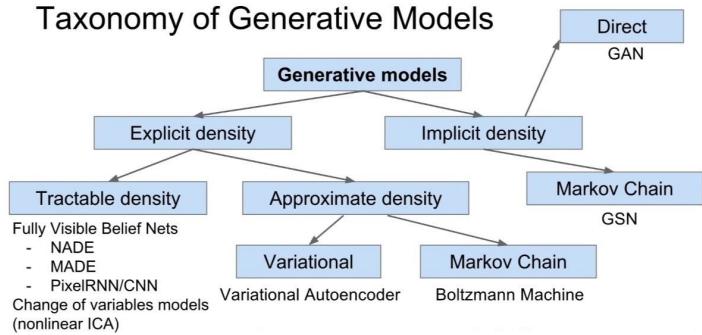
$$\mathbf{x}^{(j)} \sim p_{\text{data}} \\ j = 1, 2, \dots, |\mathcal{D}|$$

$$\min_{\theta \in \mathcal{M}} d(p_{\text{data}}, p_{\theta})$$

- What is the representation for the model family  $M$ ?
- What is the objective function  $d(\cdot)$ ?
- What is the optimization procedure for minimizing  $d(\cdot)$ ?

## Model family

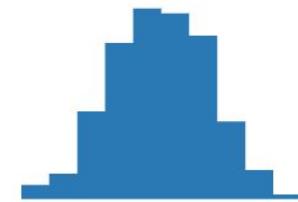
## Taxonomy of Generative Models



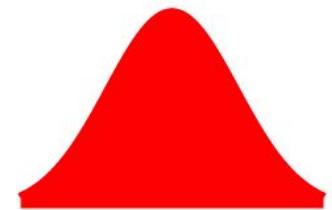
# Density estimation

We have a dataset from some distribution  $p(x)$ . We want to understand its *nature*:

- Estimate/Approximate the real distribution  $p(x)$ 
  - VAE (done)
  - Autoregressive models (this lecture)
- Learn to sample from this distribution
  - Generative Adversarial Networks=GAN (This lecture)



Samples from distribution



True distribution

## Motivation:

By fitting good data distribution we hope to use this model on downstream inference

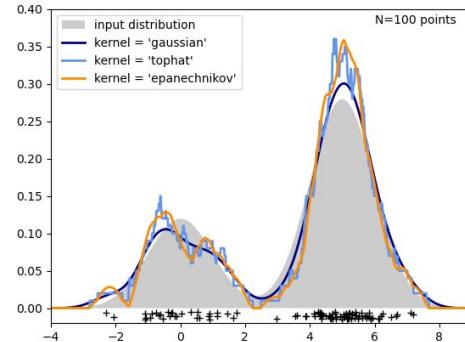
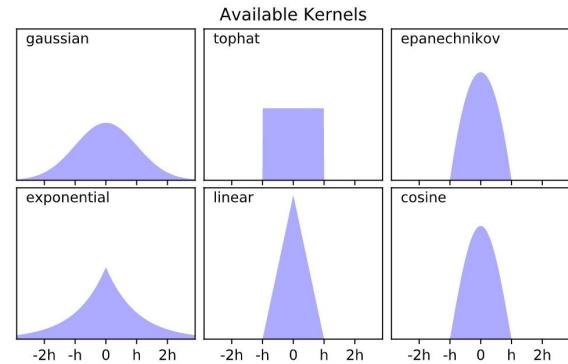
# Kernel Density Estimation

Every sample represent some distribution. By merging them we can estimate the real distribution

$$p(x) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{x - x_i}{h}\right)$$



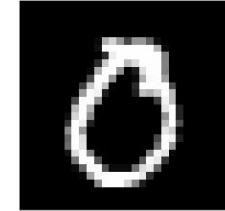
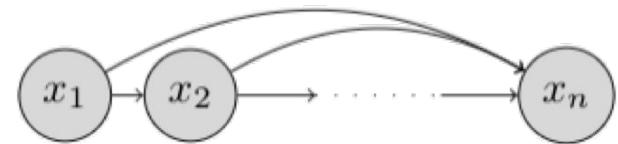
: Nah, I'm too big for you  
@Generated cat



# PixelCNN/PixelRNN

The idea is to rewrite full probability using the chain rule

$$p(\mathbf{x}) = \prod_{i=1}^n p(x_i | x_1, x_2, \dots, x_{i-1}) = \prod_{i=1}^n p(x_i | \mathbf{x}_{<i})$$



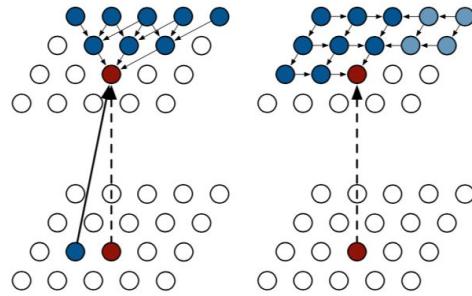
In this parametrization probability seems to be intractable. Let's parametrize history with some function

$$p(x) = \prod_{i=1}^n p(x_i | h = f_\theta(x_{<i-1}), x_{i-1})$$

# PixelRNN

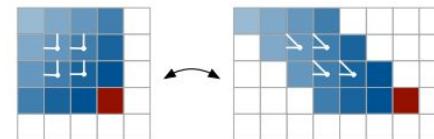
Row LSTM

$$h_{i,j} = f(h_{i-1,j-1}, h_{i-1,j}, h_{i-1,j+1}, x_{i,j})$$

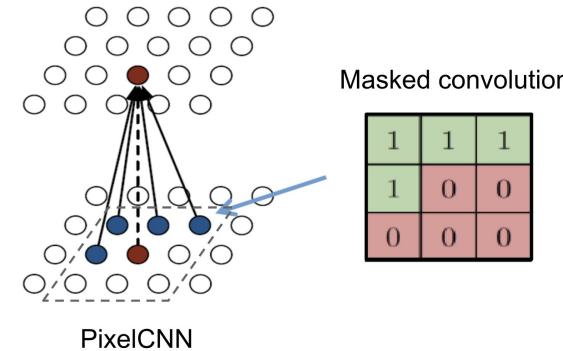
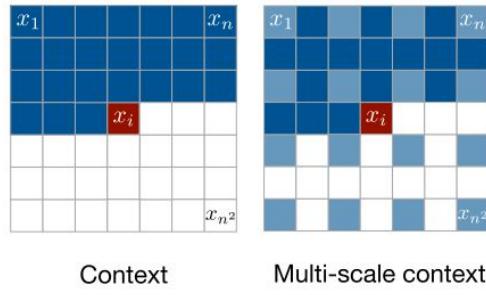


Diagonal BiLSTM

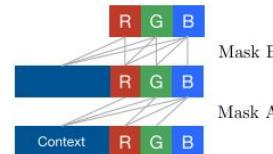
$$h_{i,j} = f(h_{i-1,j}, h_{i,j-1}, x_{i,j})$$



# Pixel CNN



How to decompose RGB channels?



# PixelCNN

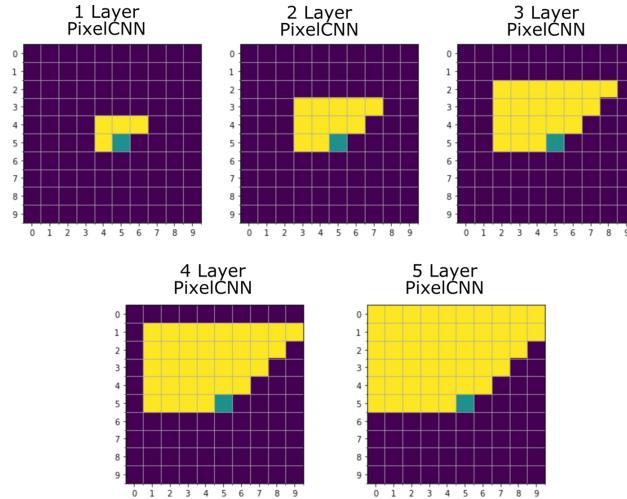


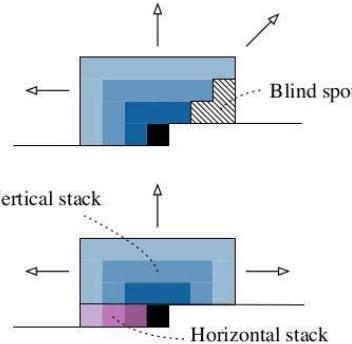
Figure 1. Image completions sampled from a PixelRNN.

# PixelCNN

## Training

- Take image, reconstruct :)

## Inference



0. Set  $i = 0$  to represent current iteration and pixel (implicitly we would translate it to the row/col/sub-pixel in the image).
1. Start off with a tensor for your image  $\mathbf{x}^0$  with any initialization (it doesn't matter).
2. Feed  $\mathbf{x}^i$  into the PixelCNN network to generate distributional outputs  $\mathbf{y}^{i+1}$ .
3. Randomly sample sub-pixel  $u_{i+1}$  from the mixture distribution defined by  $\mathbf{y}^{i+1}$  (we only need the subset of values for sub-pixel  $i + 1$ ).
4. Set  $\mathbf{x}^{i+1}$  as  $\mathbf{x}^i$  but replacing the current sub-pixel with  $u_{i+1}$ .
5. Repeat step 2-4 until entire image is generated.

# PixelCNN++

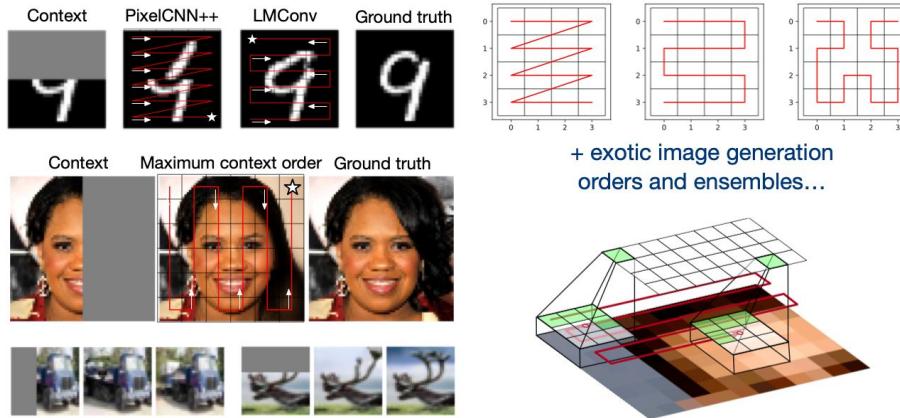
- Output the original PixelCNN is categorical (but pixel have an order)

$$P(x|\mu, s) = \begin{cases} \sigma\left(\frac{x-\mu+0.5}{s}\right) & \text{for } x = 0 \\ \sigma\left(\frac{x-\mu+0.5}{s}\right) - \sigma\left(\frac{x-\mu-0.5}{s}\right) & \text{for } 0 < x < 255 \\ 1 - \sigma\left(\frac{x-\mu-0.5}{s}\right) & \text{for } x = 255 \end{cases}$$

$$\nu \sim \sum_{i=1}^K \pi_i logistic(\mu_i, s_i)$$

# LMConv

$$p(x) = p_{\theta}(x_{\pi(1)}) \prod_{i=2}^D p_{\theta}(x_{\pi(i)} \mid Pa(\mathbf{x}_{\pi(i)}))$$

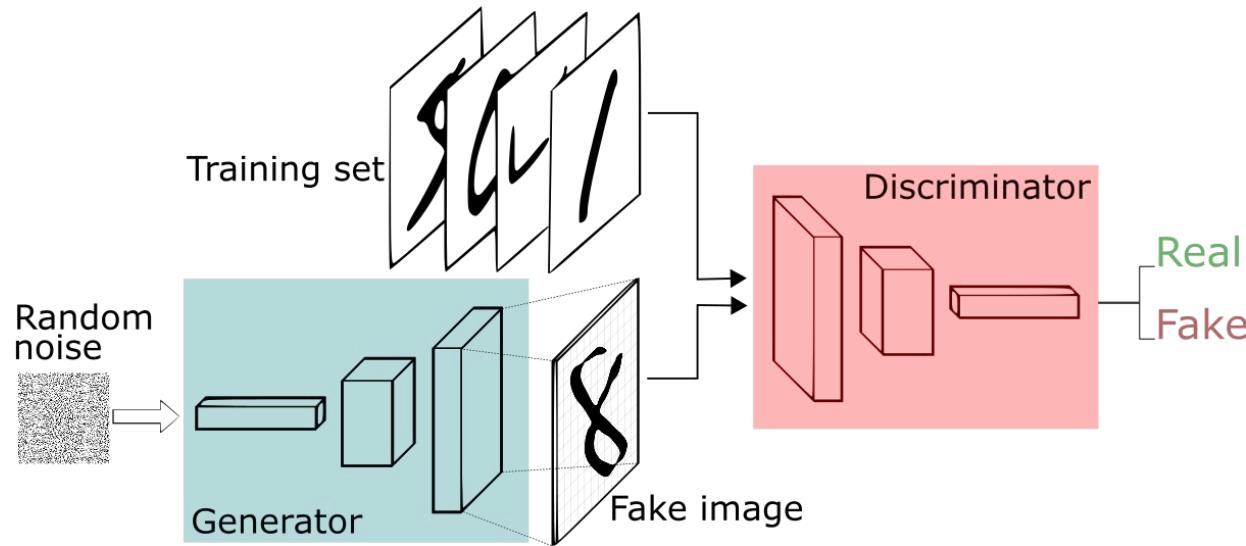


# GAN models



A natural way to set up a likelihood-free objective is to consider the two-sample test, a statistical test that determines whether or not a finite set of samples from two distributions are from the same distribution using only samples from P and Q

# GAN



# GAN

Generator

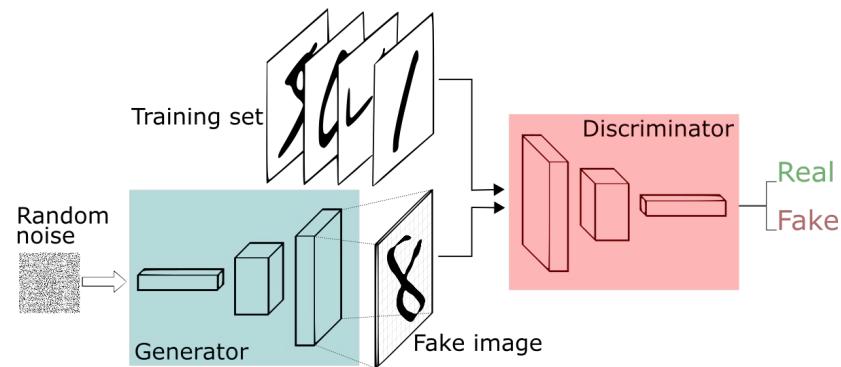
$$G : Z \rightarrow X$$

Discriminator

$$D : X \rightarrow [0, 1]$$

$p(x)$  - real distribution

$q(x)$  - generated distribution



# GAN

Discriminator output

$$\begin{cases} D(x), & x \text{ is real} \\ 1 - D(x), & x \text{ is fake} \end{cases}$$

Loss function for the task:

$$\min_{\theta} \max_{\phi} V(G_{\theta}, D_{\phi}) = \mathbb{E}_{\mathbf{x} \sim \mathbf{p}_{\text{data}}} [\log D_{\phi}(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})} [\log(1 - D_{\phi}(G_{\theta}(\mathbf{z})))]$$

---

**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:

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

**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:

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

**end for**

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

---

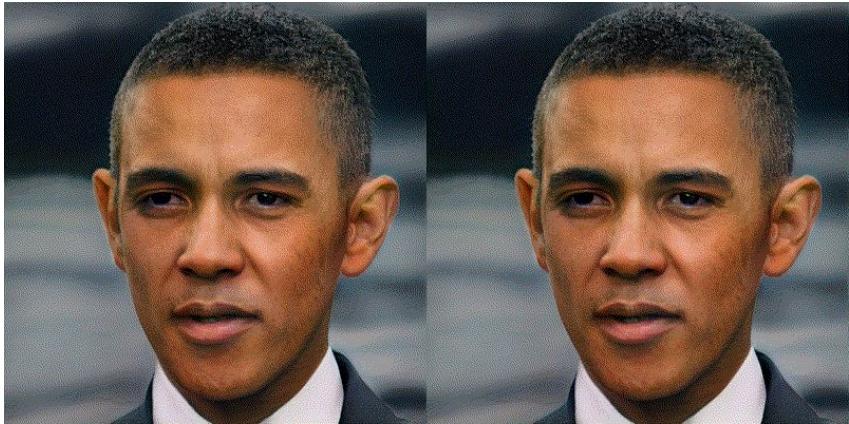
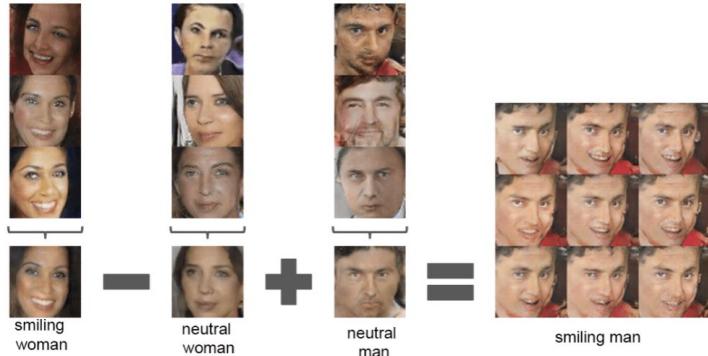
# GAN samples



A 6x6 grid of handwritten digits, likely generated by a GAN, showing varying degrees of blurriness and noise. The digits are arranged in six rows and six columns.

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

# Useful properties



# Useful properties



# Optimal model

Optimal discriminator:

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

$$\begin{aligned} V(G, D) &= \int_x p(x) \log(D(x)) dx + \int_z p(g(z)) \log(1 - D(g(z))) dz = \\ &= \int_x p(x) \log(D(x)) dx + q(x) \log(1 - D(x)) dx \end{aligned}$$

# Optimal model

Optimal generator (given optimal discriminator)

$$\begin{aligned} V(G^*, D^*(G)) &= \\ &= \int_x p(x) \log \frac{p(x)}{p(x) + q^*(x)} dx + \int_x q(x) \log \frac{q(x)}{p(x) + q^*(x)} dx \\ &= \int_x p(x) \log \frac{1}{2} dx + \int_x q(x) \log \frac{1}{2} dx = 2 \log \frac{1}{2} = -2 \log 2 \end{aligned}$$

## Jensen Shannon divergence

$$\begin{aligned} & V(G, D^*) - V(G^*, D^*) \\ &= \int_x p(x) \log \frac{p(x)}{p(x) + q(x)} dx + \int_x q(x) \log \frac{q(x)}{p(x) + q(x)} dx \\ &\quad + \int_x p(x) \log 2 dx + \int_x q(x) \log 2 dx \\ &= \int_x p(x) \log \frac{p(x)}{\frac{p(x)+q(x)}{2}} dx + \int_x q(x) \log \frac{q(x)}{\frac{p(x)+q(x)}{2}} dx \\ &= KL\left(p(x) \parallel \frac{p(x) + q(x)}{2}\right) + KL\left(q(x) \parallel \frac{p(x) + q(x)}{2}\right) \geq 0 \end{aligned}$$

# GAN problems

**Problem:** vanishing gradient

$$D(G(z)) \approx 0 \Rightarrow \nabla_{\theta_G} \log(1 - D(G(z))) \approx 0$$

**Solution:**

$$\mathbb{E}_{z \sim p(z)} [\log(D(G(z)))] \rightarrow \max_G$$

**Other options:**

Adding noise; Soft labels; Flip labels

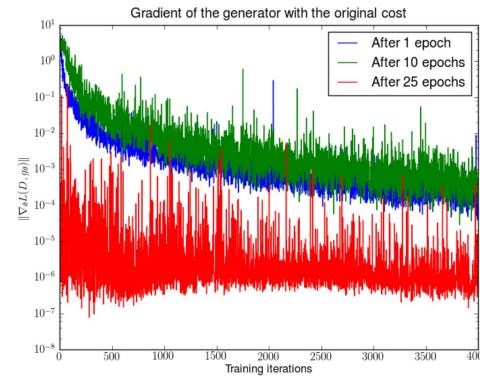
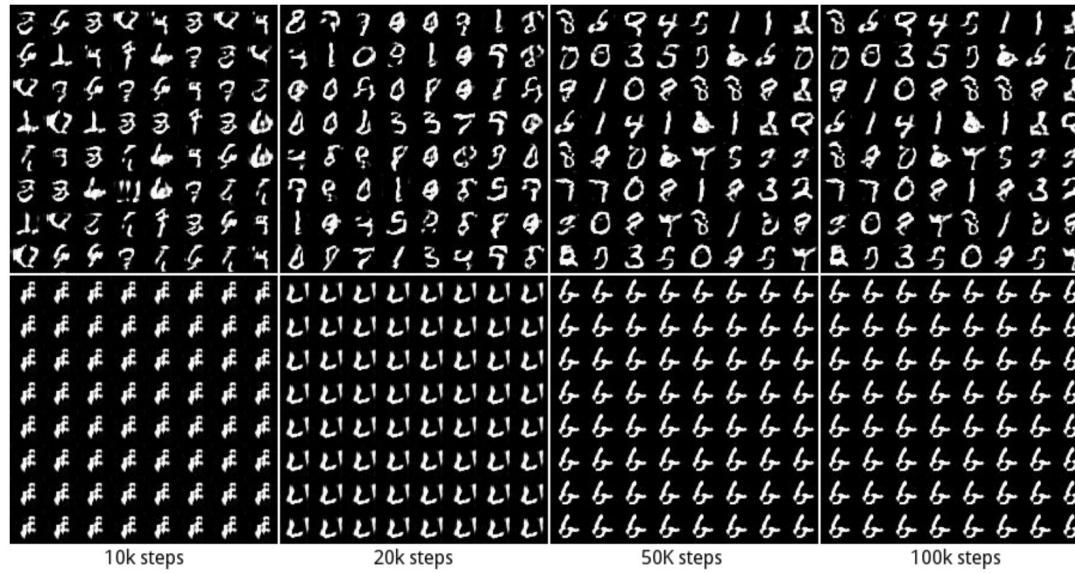


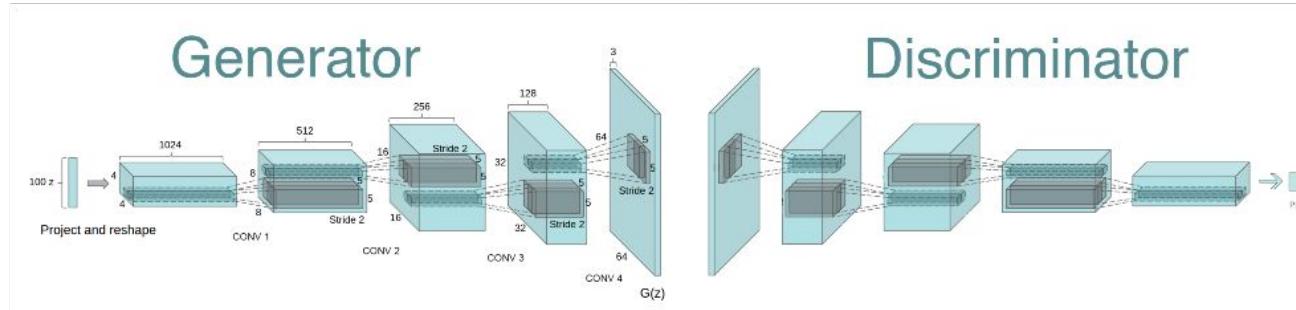
Figure 2: First, we trained a DCGAN for 1, 10 and 25 epochs. Then, with the generator fixed we train a discriminator from scratch and measure the gradients with the original cost function. We see the gradient norms decay quickly, in the best case 5 orders of magnitude after 4000 discriminator iterations. Note the logarithmic scale.

# GAN problems

Mode collapse



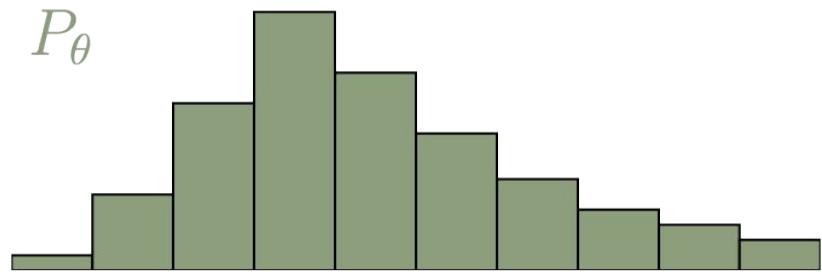
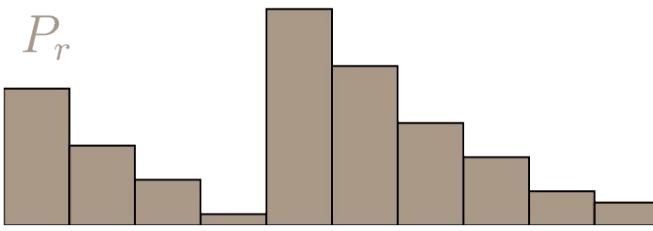
# DCGAN



- Replace pooling with strided convolutions
- LeakyReLU instead of ReLU

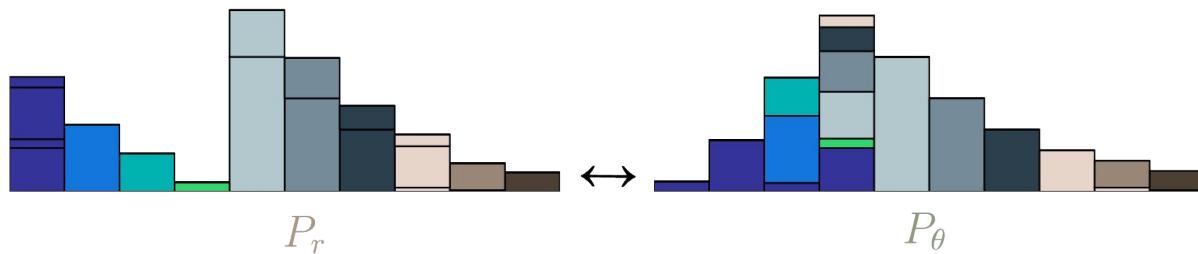


# Wasserstein GAN



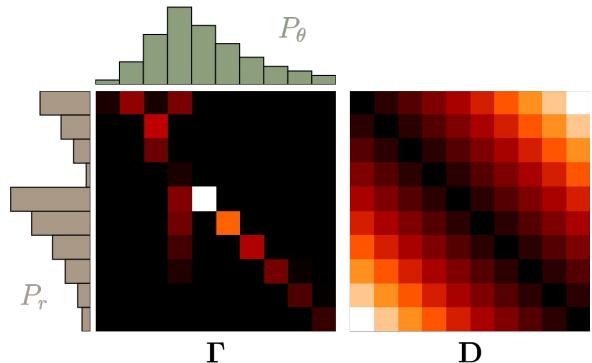
GAN tries to learn create samples close to samples from true distribution = tries to map one distribution to another

# Wasserstein GAN

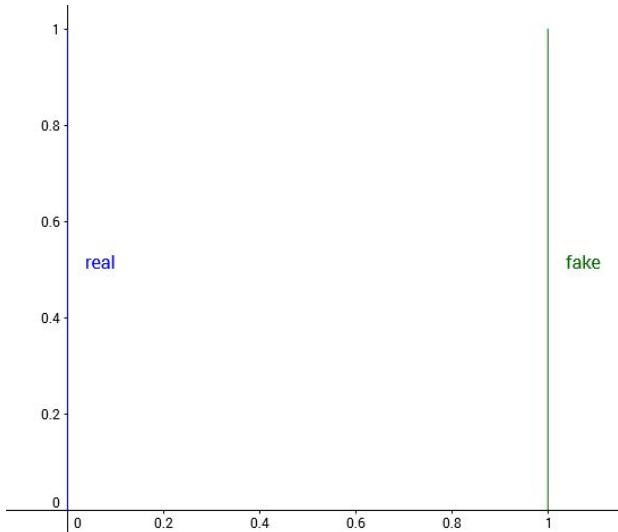


Example of mapping

$$W(p, q) = \inf_{\gamma \in \Pi(p, q)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$



# Example



$$KL(p\|q_\theta) = KL(q_\theta\|p) = \begin{cases} +\infty, & \theta \neq 0 \\ 0, & \theta = 0 \end{cases}$$

$$JS(p, q_\theta) = \begin{cases} \ln 2, & \theta \neq 0 \\ 0, & \theta = 0 \end{cases}$$

$$W(p, q_\theta) = |\theta|$$

# Wasserstein distance

$$W(p, q) = \inf_{\gamma \in \Pi(p, q)} \mathbb{E}_{(x, y) \sim \gamma} [\|x - y\|]$$

f - is L-Lipschitz function

$$f(x) - f(x') \leq L \|x - x'\|$$

Wasserstein distance can be estimated using Kantorovich-Rubenstein duality

$$W(p, q_\theta) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim p}[f(x)] - \mathbb{E}_{x \sim q_\theta}[f(x)]$$

---

**Algorithm 1** WGAN, our proposed algorithm. All experiments in the paper used the default values  $\alpha = 0.00005$ ,  $c = 0.01$ ,  $m = 64$ ,  $n_{\text{critic}} = 5$ .

---

**Require:** :  $\alpha$ , the learning rate.  $c$ , the clipping parameter.  $m$ , the batch size.  
 $n_{\text{critic}}$ , the number of iterations of the critic per generator iteration.

**Require:** :  $w_0$ , initial critic parameters.  $\theta_0$ , initial generator's parameters.

```
1: while  $\theta$  has not converged do
2:   for  $t = 0, \dots, n_{\text{critic}}$  do
3:     Sample  $\{x^{(i)}\}_{i=1}^m \sim \mathbb{P}_r$  a batch from the real data.
4:     Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
5:      $g_w \leftarrow \nabla_w \left[ \frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)})) \right]$ 
6:      $w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$ 
7:      $w \leftarrow \text{clip}(w, -c, c)$ 
8:   end for
9:   Sample  $\{z^{(i)}\}_{i=1}^m \sim p(z)$  a batch of prior samples.
10:   $g_\theta \leftarrow -\nabla_\theta \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))$ 
11:   $\theta \leftarrow \theta - \alpha \cdot \text{RMSProp}(\theta, g_\theta)$ 
12: end while
```

---

# Wasserstein GAN

- More stable gradient
- No mode collapse:

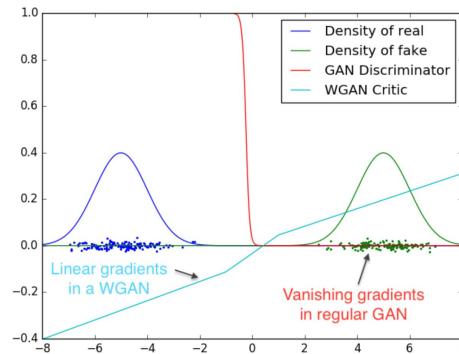
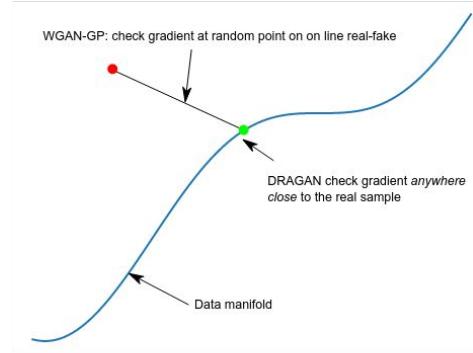


Figure 2: Optimal discriminator and critic when learning to differentiate two Gaussians. As we can see, the discriminator of a minimax GAN saturates and results in vanishing gradients. Our WGAN critic provides very clean gradients on all parts of the space.

# WGAN-GP

Clipping is bad idea to make neural network a Lipschitz function



$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \underbrace{\lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} \left[ (\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2 \right]}_{\text{Our gradient penalty}}.$$

# Spectral Normalization

$$f(\mathbf{x}, \theta) = W^{L+1} a_L(W^L(a_{L-1}(W^{L-1}(\dots a_1(W^1 \mathbf{x}) \dots))))$$

activation function                                  weights for each layer

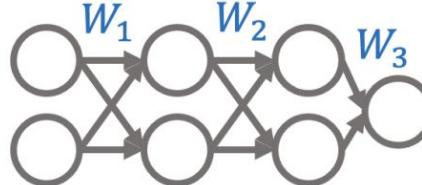
$$\begin{aligned}\|f\|_{\text{Lip}} &\leq \|(\mathbf{h}_L \mapsto W^{L+1} \mathbf{h}_L)\|_{\text{Lip}} \cdot \|a_L\|_{\text{Lip}} \cdot \|(\mathbf{h}_{L-1} \mapsto W^L \mathbf{h}_{L-1})\|_{\text{Lip}} \\ &\cdots \|a_1\|_{\text{Lip}} \cdot \|(\mathbf{h}_0 \mapsto W^1 \mathbf{h}_0)\|_{\text{Lip}} = \prod_{l=1}^{L+1} \|(\mathbf{h}_{l-1} \mapsto W^l \mathbf{h}_{l-1})\|_{\text{Lip}} = \prod_{l=1}^{L+1} \sigma(W^l).\end{aligned}$$

$$\bar{W}_{\text{SN}}(W) := W/\sigma(W)$$

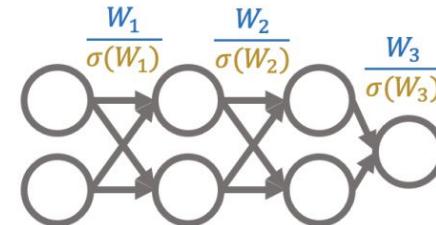
$$\sigma(\bar{W}_{\text{SN}}(W)) = 1$$

$$\|f\|_{\text{Lip}} = 1$$

(a) Discriminator  
without SN



(b) Discriminator  
with SN



---

**Algorithm 1** SGD with spectral normalization

---

- Initialize  $\tilde{\mathbf{u}}_l \in \mathcal{R}^{d_l}$  for  $l = 1, \dots, L$  with a random vector (sampled from isotropic distribution).
- For each update and each layer  $l$ :
  1. Apply power iteration method to a unnormalized weight  $W^l$ :

$$\tilde{\mathbf{v}}_l \leftarrow (W^l)^T \tilde{\mathbf{u}}_l / \| (W^l)^T \tilde{\mathbf{u}}_l \|_2 \quad (20)$$

$$\tilde{\mathbf{u}}_l \leftarrow W^l \tilde{\mathbf{v}}_l / \| W^l \tilde{\mathbf{v}}_l \|_2 \quad (21)$$

2. Calculate  $\bar{W}_{\text{SN}}$  with the spectral norm:

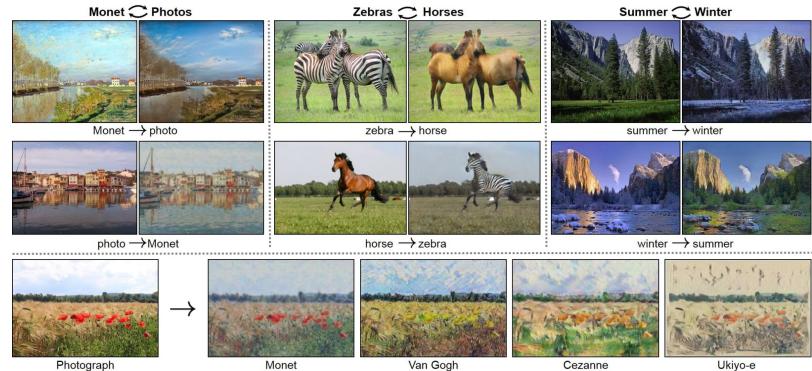
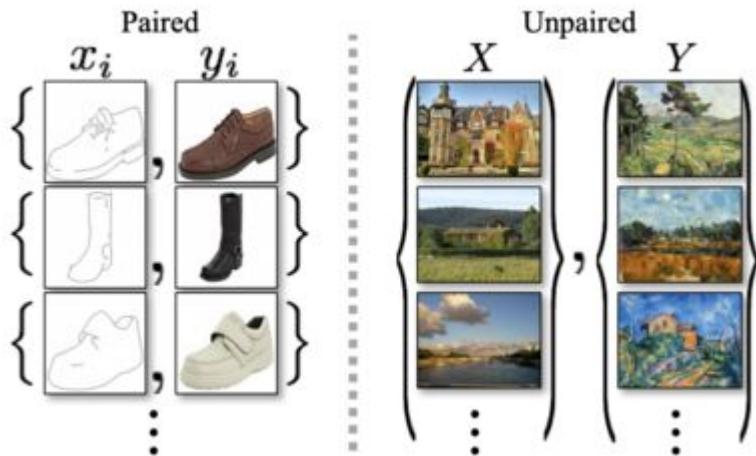
$$\bar{W}_{\text{SN}}^l(W^l) = W^l / \sigma(W^l), \text{ where } \sigma(W^l) = \tilde{\mathbf{u}}_l^T W^l \tilde{\mathbf{v}}_l \quad (22)$$

3. Update  $W^l$  with SGD on mini-batch dataset  $\mathcal{D}_M$  with a learning rate  $\alpha$ :

$$W^l \leftarrow W^l - \alpha \nabla_{W^l} \ell(\bar{W}_{\text{SN}}^l(W^l), \mathcal{D}_M) \quad (23)$$

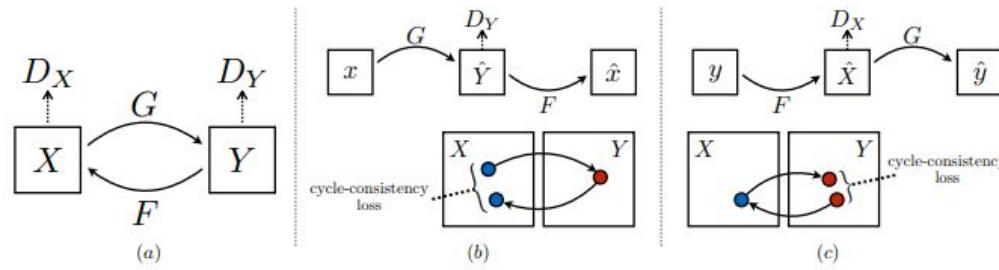
---

# Cycle GAN



Want to map one distribution to another and vice versa. Day  $\leftrightarrow$  Night. Summer  $\leftrightarrow$  Winter

# Cycle GAN



$$\min_{F, G, D_X, D_Y} \mathcal{L}_{GAN}(G, D_Y, X, Y) + \mathcal{L}_{GAN}(F, D_X, X, Y) + \lambda (\mathbb{E}_X[||F(G(X)) - X||_1] + \mathbb{E}_Y[||G(F(Y)) - Y||_1])$$

# Recap

- Generative models
- Kernel density estimation
- Autoregressive models (PixelCNN, PixelRNN)
- GAN
- WGAN
- CycleGAN

# Metrics

How do automatically measure the performance of generative models?

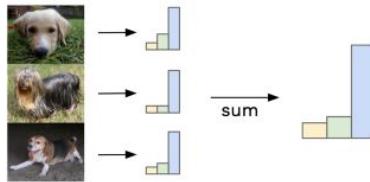
- The images generated should contain clear objects (i.e. the images are sharp rather than blurry), or  $p(y|x)$  should be low entropy.
- The generative algorithm should output a high diversity of images from all the different classes in ImageNet, or  $p(y)$  should be high entropy.

$$\ln(\text{IS}(G)) = I(y; \mathbf{x}) \quad I(y; \mathbf{x}) = H(y) - H(y \mid \mathbf{x})$$

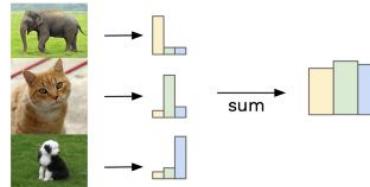
# Inception Score



Similar labels sum to give focussed distribution



Different labels sum to give uniform distribution



How to get the distributions?

# Inception Score problems

Consider two classes with data distributions  $N(-1, 2)$  and  $N(1, 2)$ .

Then, optimal generator will produce  $-\infty$  and  $\infty$

What if we use dataset for generation of unique objects which are absent in ImageNet?

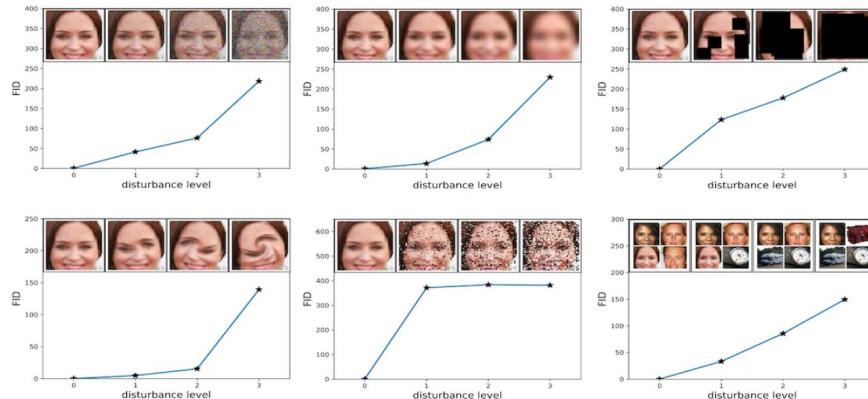
# Inception Score

How to calculate Inception Score:

- Create n=10 splits.
- Calculate Inception Score independently on each dataset

# Frechet Inception Distance

Real and generated distribution should be close. But it's hard to measure distance between distribution on raw pixels -> let's use the latent representation



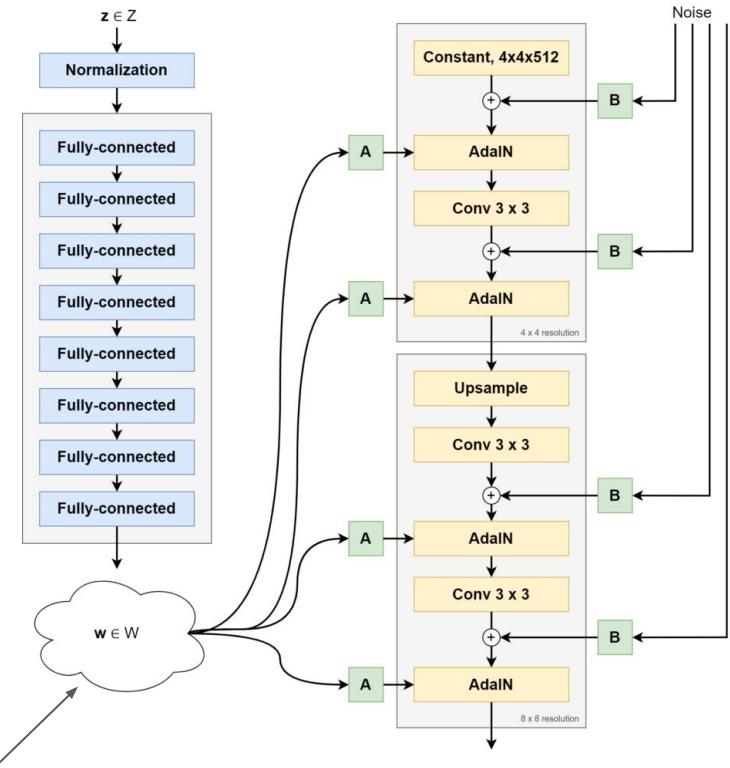
$$FID(r, g) = \|\mu_r - \mu_g\|_2^2 + \text{Tr} \left( \Sigma_r + \Sigma_g - 2 (\Sigma_r \Sigma_g)^{\frac{1}{2}} \right)$$

# Style-GAN

Motivation:

- Generators operate as black boxes. Latent spaces of classic GANs were poorly understood
- GANs must be heavily regularized.
- There is little control over image synthesis.

# StyleGAN



Disentangled space

$$\text{AdaIN}(\mathbf{x}_i, \mathbf{y}) = \mathbf{y}_{s,i} \frac{\mathbf{x}_i - \mu(\mathbf{x}_i)}{\sigma(\mathbf{x}_i)} + \mathbf{y}_{b,i},$$

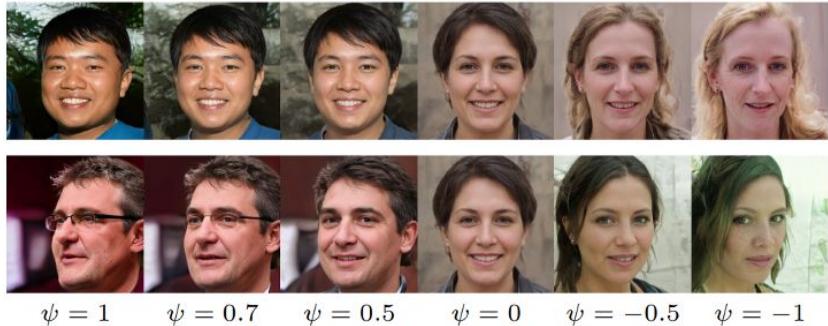


# Truncation trick

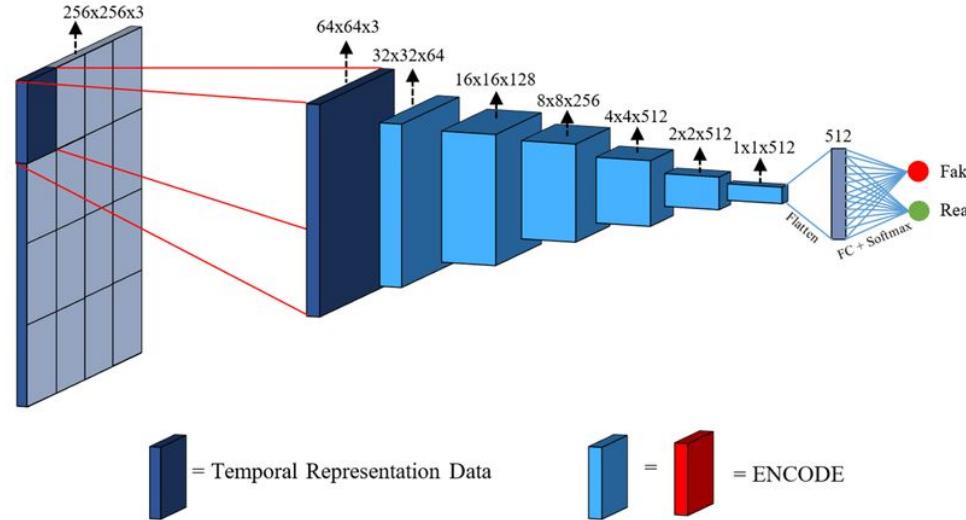
How to choose image for paper?

$$\bar{\mathbf{w}} = \mathbb{E}_{\mathbf{z} \sim P(\mathbf{z})} [f(\mathbf{z})]$$

$$\mathbf{w}' = \bar{\mathbf{w}} + \psi(\mathbf{w} - \bar{\mathbf{w}})$$



# Patch GAN



Instead of using all image for classification one can use patches of it. Better modelling of high frequency details

# f-GAN

$$D_f(P\|Q) = \int_{\mathcal{X}} q(x) f\left(\frac{p(x)}{q(x)}\right) dx$$

f - convex

f - semicontinuous

$f(0) = 1$

$$\begin{aligned} D_f(P\|Q) &= \int_{\mathcal{X}} q(x) \sup_{t \in \text{dom}_{f^*}} \left\{ t \frac{p(x)}{q(x)} - f^*(t) \right\} dx \\ &\geq \sup_{T \in \mathcal{T}} \left( \int_{\mathcal{X}} p(x) T(x) dx - \int_{\mathcal{X}} q(x) f^*(T(x)) dx \right) \\ &= \sup_{T \in \mathcal{T}} (\mathbb{E}_{x \sim P}[T(x)] - \mathbb{E}_{x \sim Q}[f^*(T(x))]), \end{aligned}$$

# f-GAN

Name	$D_f(P\ Q)$	Generator $f(u)$	$T^*(x)$
Kullback-Leibler	$\int p(x) \log \frac{p(x)}{q(x)} dx$	$u \log u$	$1 + \log \frac{p(x)}{q(x)}$
Reverse KL	$\int q(x) \log \frac{q(x)}{p(x)} dx$	$-\log u$	$-\frac{q(x)}{p(x)}$
Pearson $\chi^2$	$\int \frac{(q(x)-p(x))^2}{p(x)} dx$	$(u - 1)^2$	$2(\frac{p(x)}{q(x)} - 1)$
Squared Hellinger	$\int \left( \sqrt{p(x)} - \sqrt{q(x)} \right)^2 dx$	$(\sqrt{u} - 1)^2$	$(\sqrt{\frac{p(x)}{q(x)}} - 1) \cdot \sqrt{\frac{q(x)}{p(x)}}$
Jensen-Shannon	$\frac{1}{2} \int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx$	$-(u+1) \log \frac{1+u}{2} + u \log u$	$\log \frac{2p(x)}{p(x)+q(x)}$
GAN	$\int p(x) \log \frac{2p(x)}{p(x)+q(x)} + q(x) \log \frac{2q(x)}{p(x)+q(x)} dx - \log(4)$	$u \log u - (u+1) \log(u+1)$	$\log \frac{p(x)}{p(x)+q(x)}$

# Progressive GAN

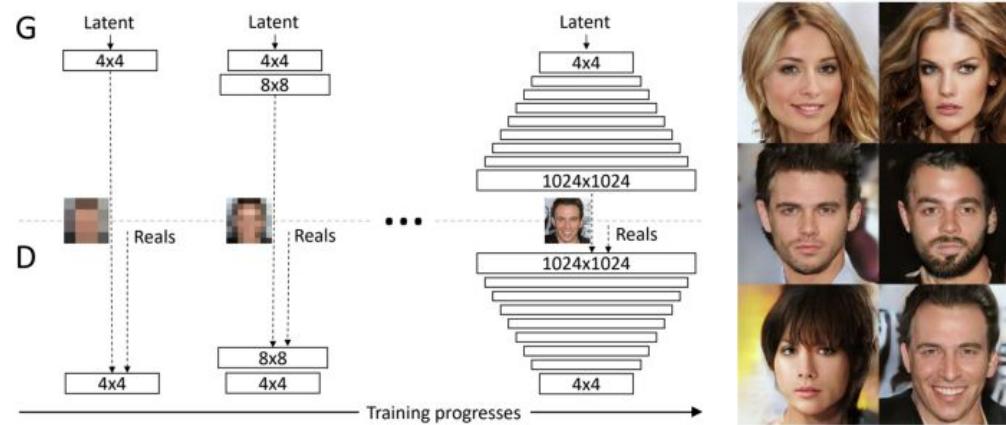


Figure 1: Our training starts with both the generator (G) and discriminator (D) having a low spatial resolution of  $4 \times 4$  pixels. As the training advances, we incrementally add layers to G and D, thus increasing the spatial resolution of the generated images. All existing layers remain trainable throughout the process. Here  $N \times N$  refers to convolutional layers operating on  $N \times N$  spatial resolution. This allows stable synthesis in high resolutions and also speeds up training considerably. One the right we show six example images generated using progressive growing at  $1024 \times 1024$ .

# Vector Quantised (VQ)-VAE

**Encoder**



image to  
discrete codes

56	73	67	23	81	19	...
----	----	----	----	----	----	-----

**Decoder**

56	73	67	23	81	19	...
----	----	----	----	----	----	-----

discrete codes  
to image



TEXT PROMPT

an illustration of a baby daikon radish in a tutu walking a dog

AI-GENERATED IMAGES



Edit prompt or view more images ↴

TEXT PROMPT

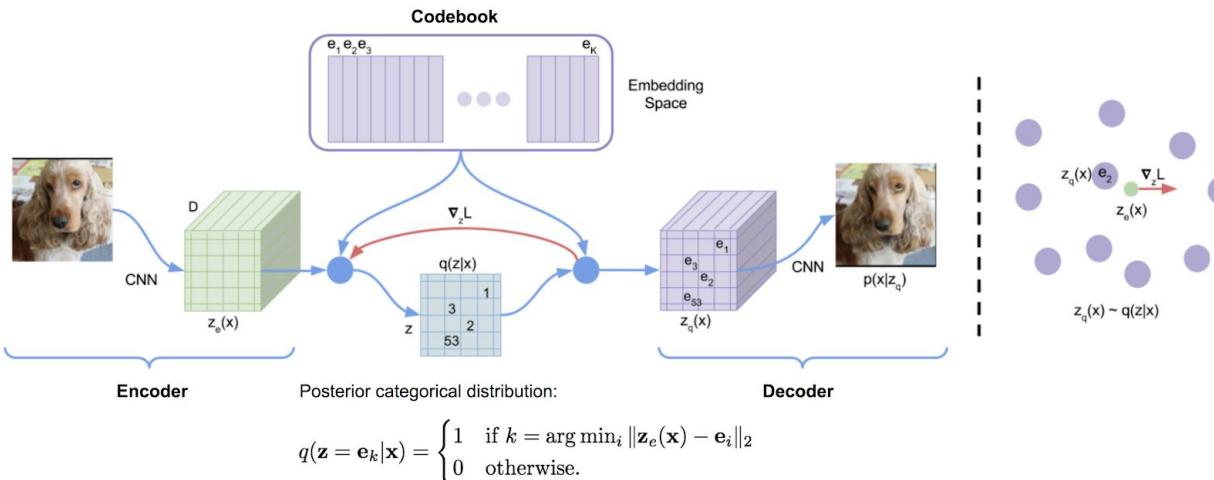
an armchair in the shape of an avocado [...]

AI-GENERATED IMAGES



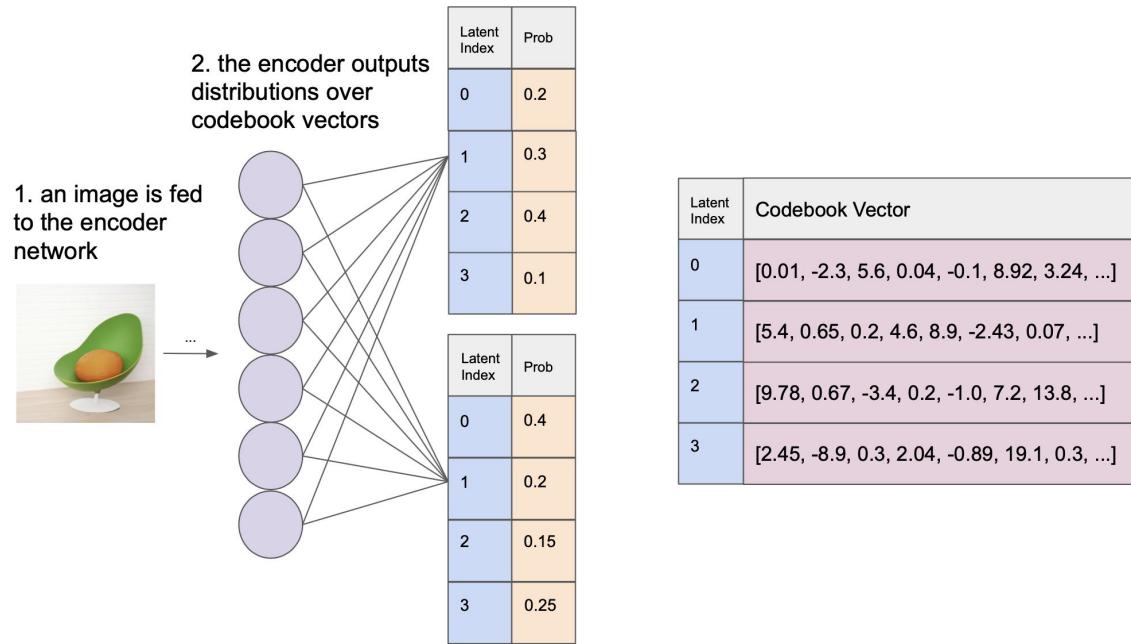
Edit prompt or view more images ↴

# VQ-VAE



$$\| \log(p(x | q(x))) + \| \operatorname{sg}[z_e(x)] - e \|_2^2 + \beta \| z_e(x) - \operatorname{sg}[e] \|_2^2$$

# VQ-VAE



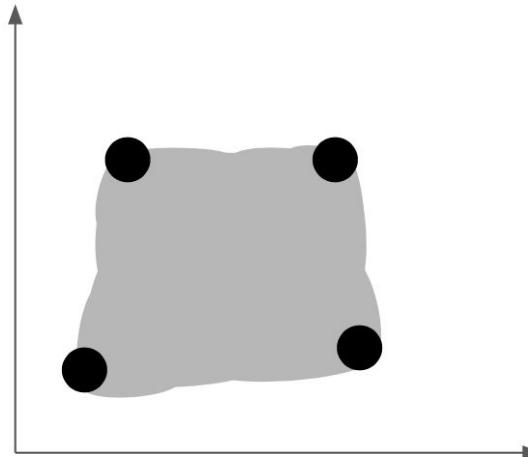
# VQ-VAE

**The Convex Hull of  
Codebook Vectors**

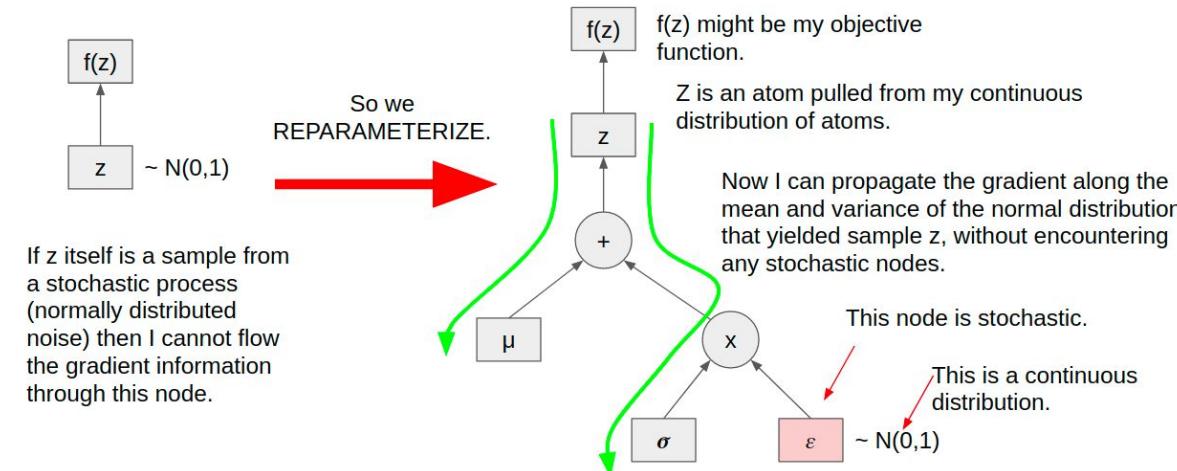
key:

● = codebook vector

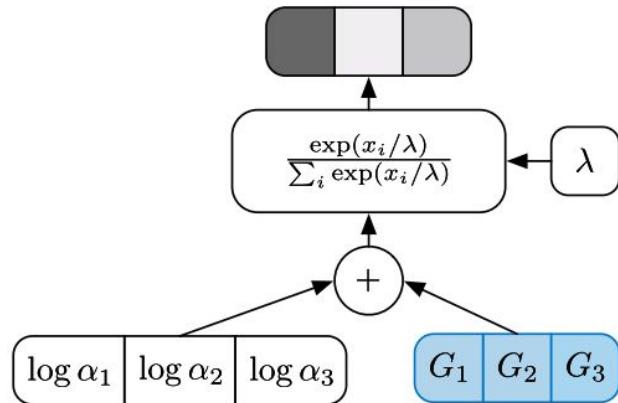
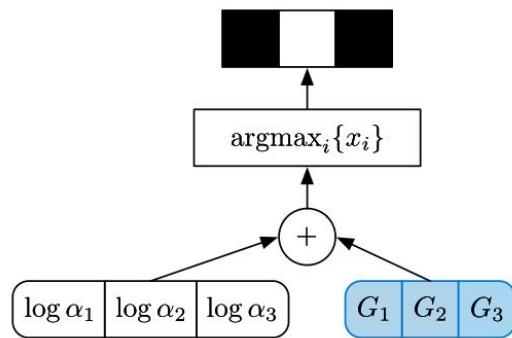
○ = convex hull of codebook vectors



# Gumbel-Softmax trick



# Gumbel-softmax trick



$$y_i = \frac{e^{\frac{\log(\pi_i) + g_i}{\tau}}}{\sum_{j=1}^k e^{\frac{\log(\pi_j) + g_j}{\tau}}} \text{ for } i = 1, \dots, k$$

# Gumbel-Softmax

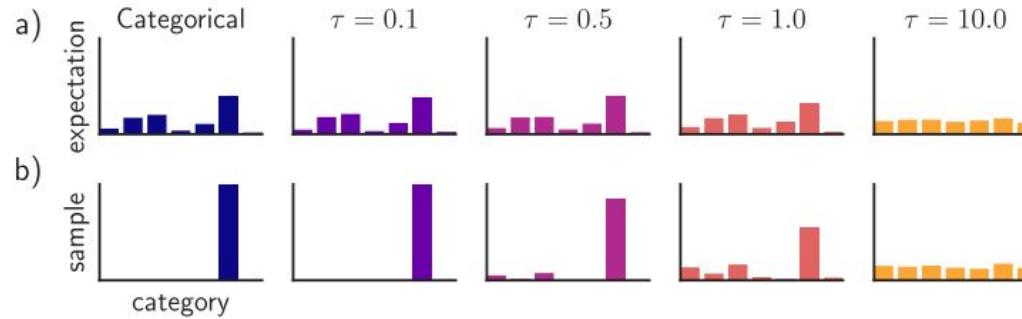
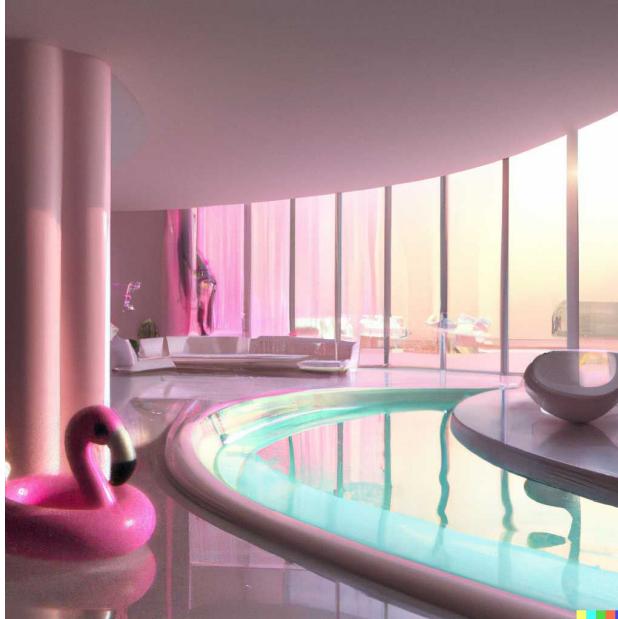


Figure 1: The Gumbel-Softmax distribution interpolates between discrete one-hot-encoded categorical distributions and continuous categorical densities. (a) For low temperatures ( $\tau = 0.1, \tau = 0.5$ ), the expected value of a Gumbel-Softmax random variable approaches the expected value of a categorical random variable with the same logits. As the temperature increases ( $\tau = 1.0, \tau = 10.0$ ), the expected value converges to a uniform distribution over the categories. (b) Samples from Gumbel-Softmax distributions are identical to samples from a categorical distribution as  $\tau \rightarrow 0$ . At higher temperatures, Gumbel-Softmax samples are no longer one-hot, and become uniform as  $\tau \rightarrow \infty$ .

# VQ-VAE



# Recap