



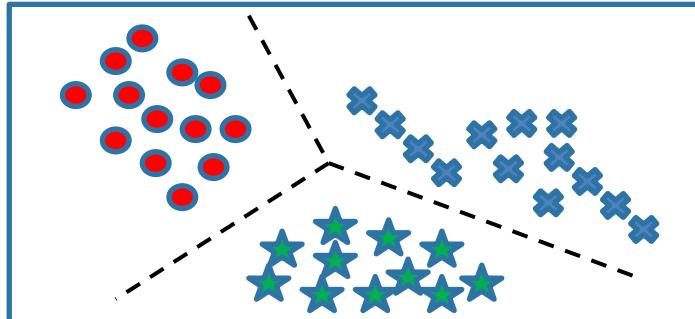
Deep generative methods



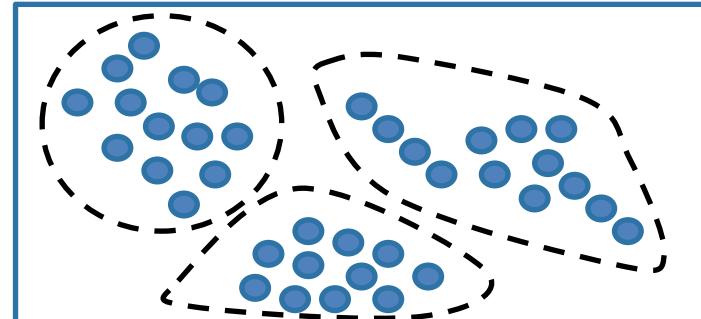
Stéphane Lathuilière
Télécom Paris

Types of Learning

- **Supervised (inductive)** learning
 - Given: training data + desired outputs (labels)
- **Unsupervised** learning
 - Given: training data (without desired outputs)



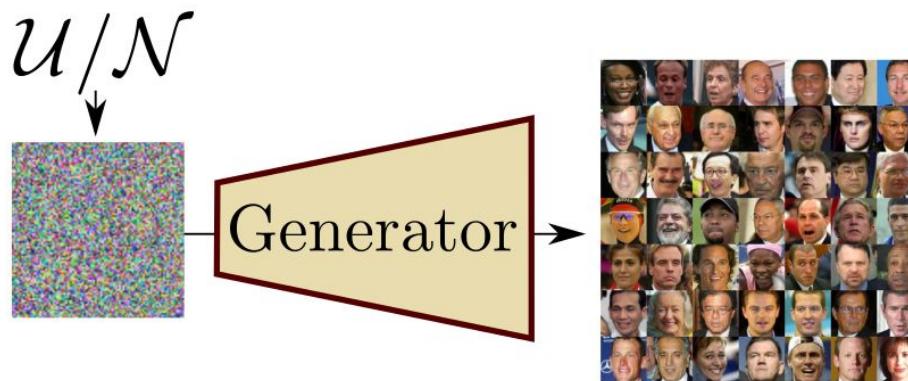
Supervised learning



Unsupervised learning

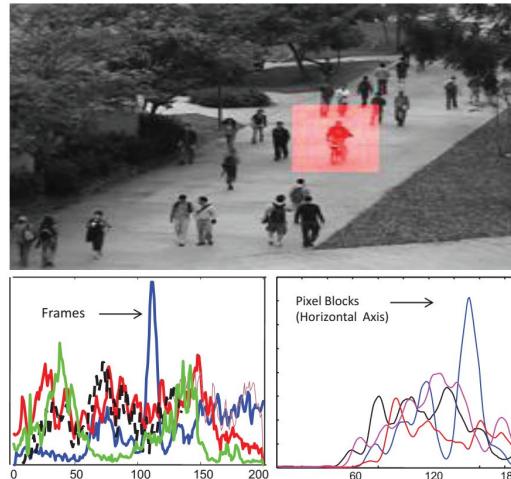
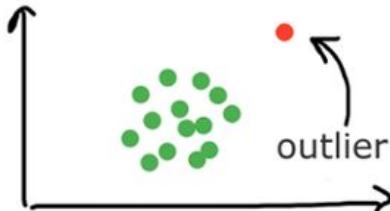
Unsupervised Learning

- Sampling: learn to sample according to $p(x): \mathbb{R}^d \rightarrow \mathbb{R}$



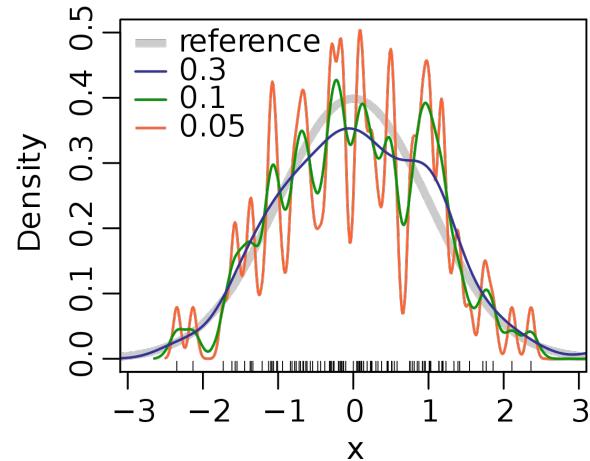
Unsupervised Learning

- **Anomaly detection:** the computer program sifts through a set of events or objects and flags some of them as being unusual or atypical.
- Example: credit card fraud detection, video surveillance.



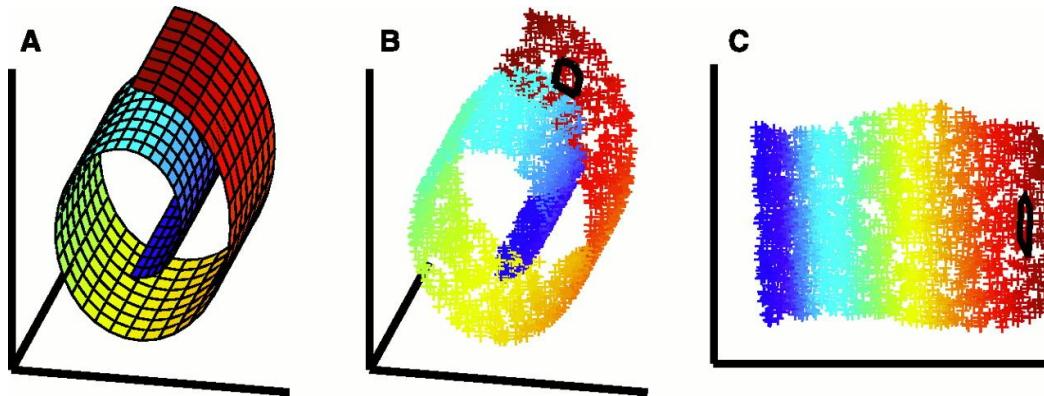
Unsupervised Learning

- Density Estimation: learn a function $p(x): \mathbb{R}^d \rightarrow \mathbb{R}$



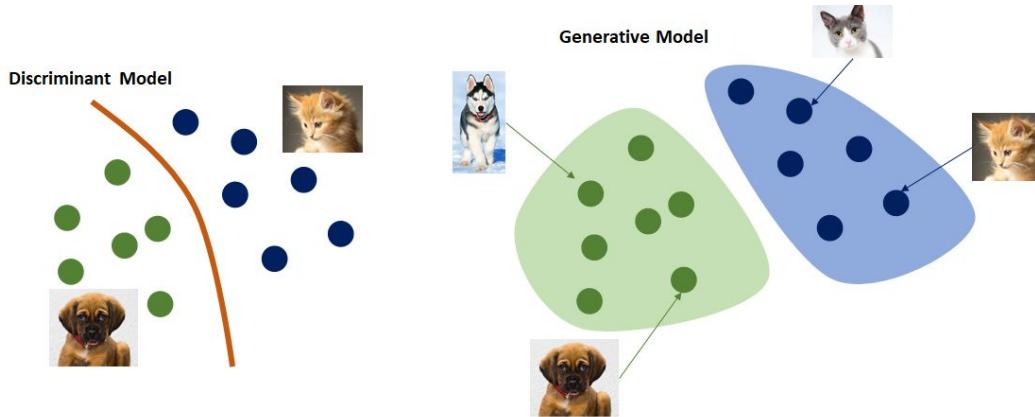
Unsupervised Learning

- Dimensionality Reduction



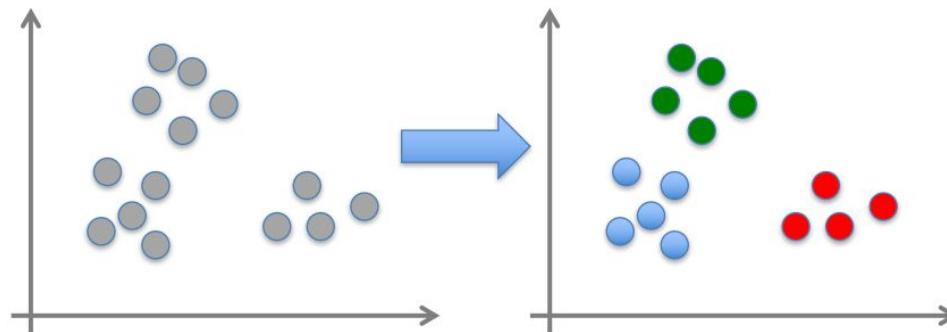
Supervised Learning: Classification

Discriminative vs Generative Algorithms

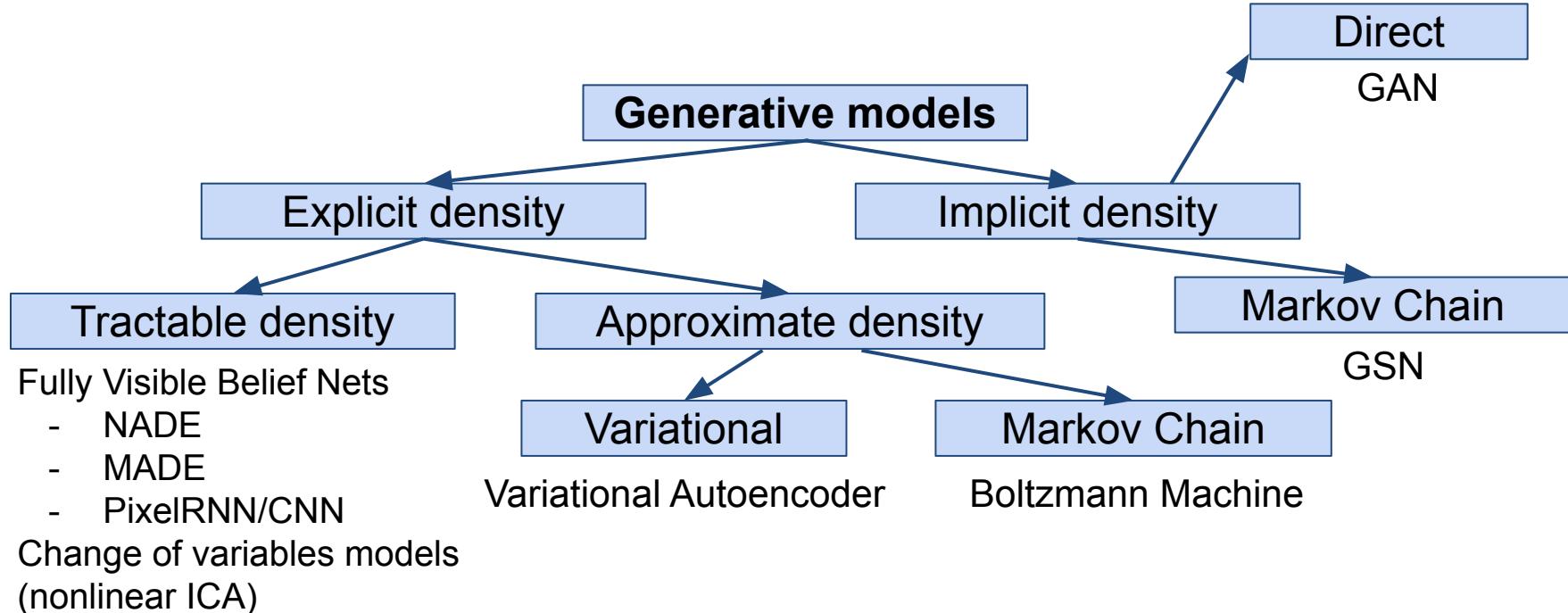


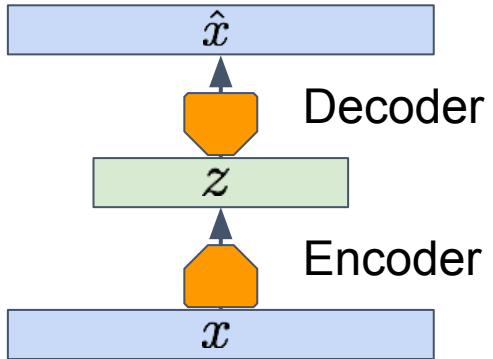
Unsupervised Learning: Clustering

- Given $\mathcal{T} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ (without labels)
output hidden structure behind the \mathbf{x} 's – e.g., clustering



Taxonomy of Generative Models

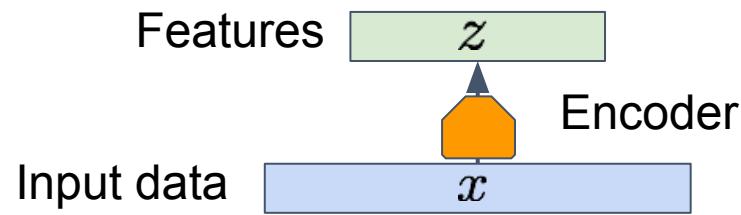




Autoencoders

Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data



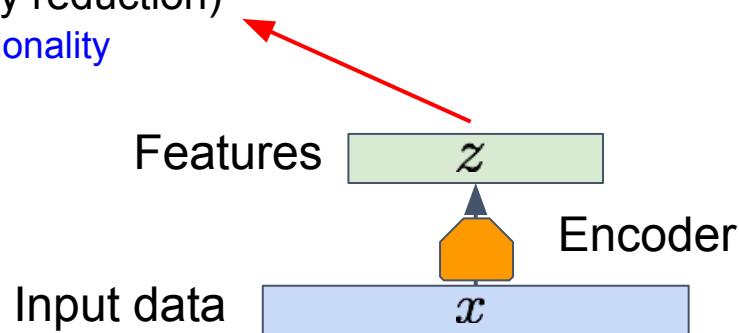
Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x

(dimensionality reduction)

Q: Why dimensionality reduction?



Autoencoders

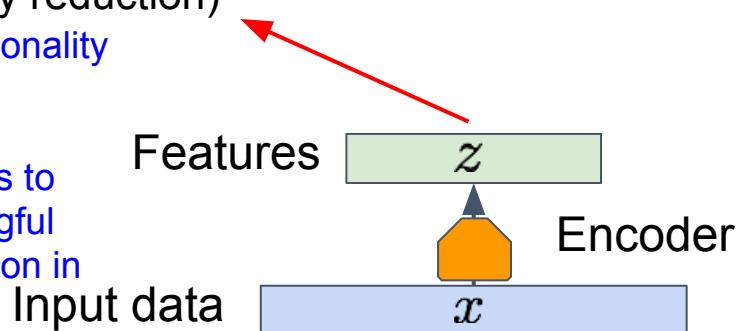
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x

(dimensionality reduction)

Q: Why dimensionality reduction?

A: Want features to capture meaningful factors of variation in data



Autoencoders

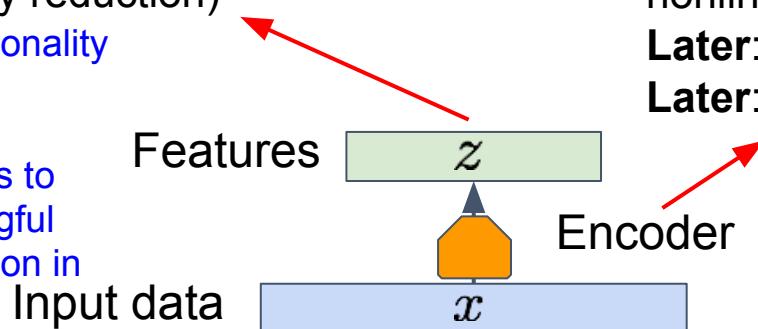
Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality
reduction?

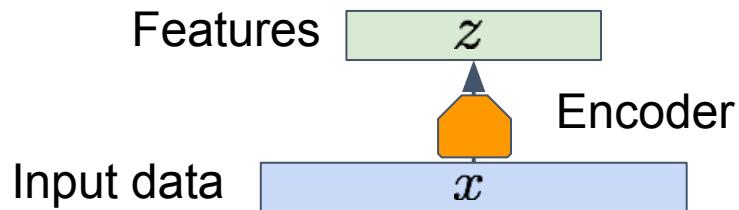
A: Want features to
capture meaningful
factors of variation in
data

Originally: Linear +
nonlinearity (sigmoid)
Later: Deep, fully-connected
Later: ReLU CNN



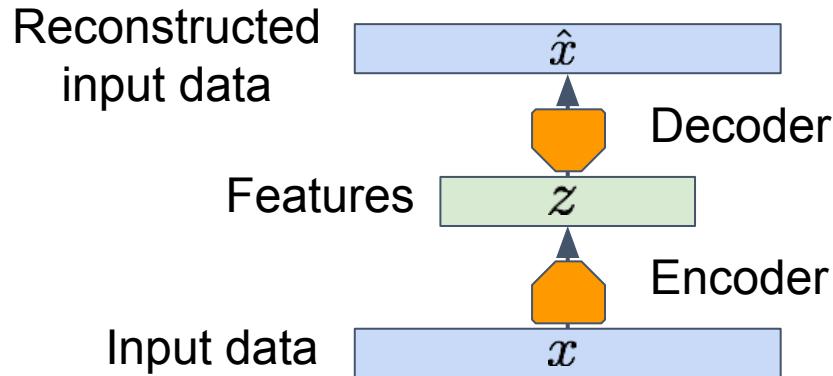
Autoencoders

How to learn this feature representation?



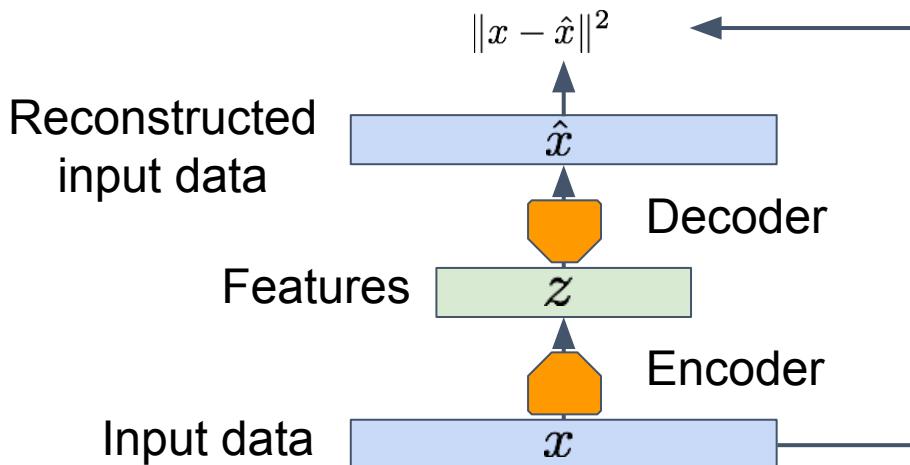
Autoencoders

Train such that features can be used to reconstruct original data
“Autoencoding” - encoding itself

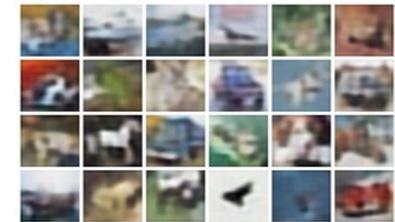


Autoencoders

Train such that features can be used to reconstruct original data



Reconstructed data

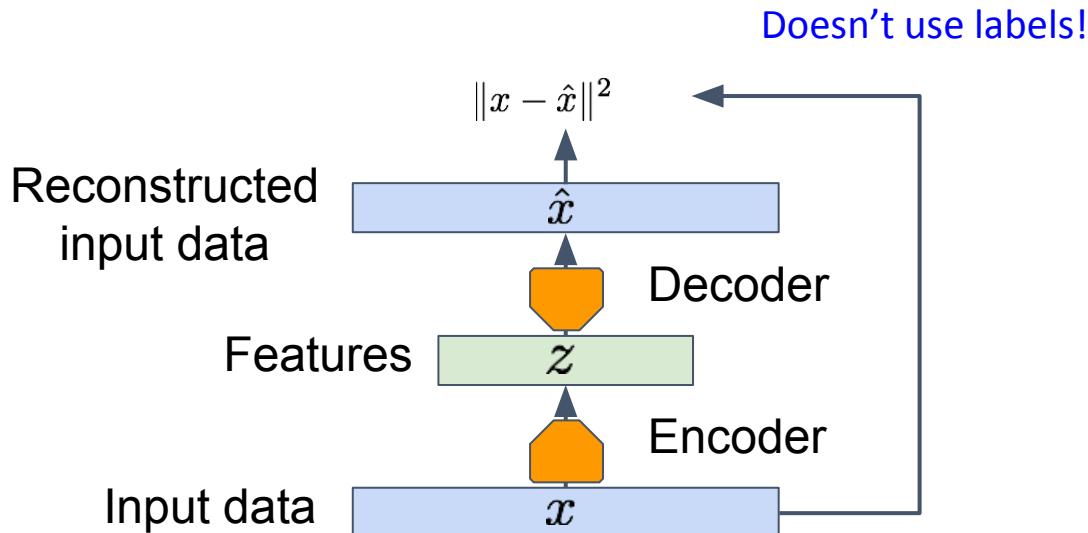


Input data

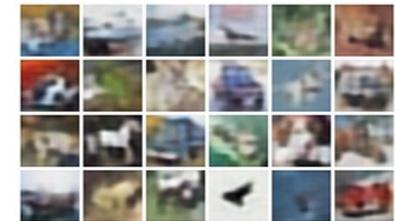


Autoencoders

Train such that features can be used to reconstruct original data



Reconstructed data



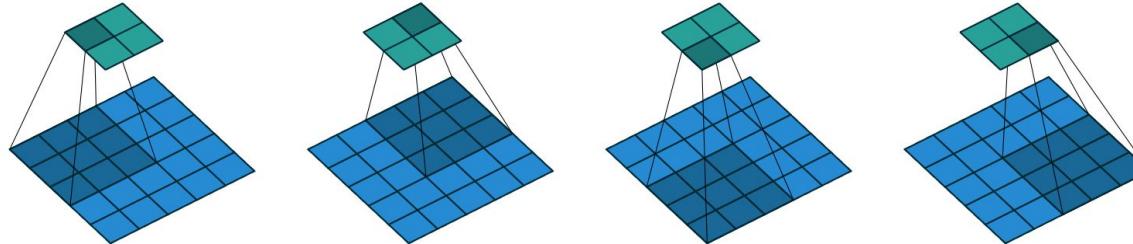
Input data



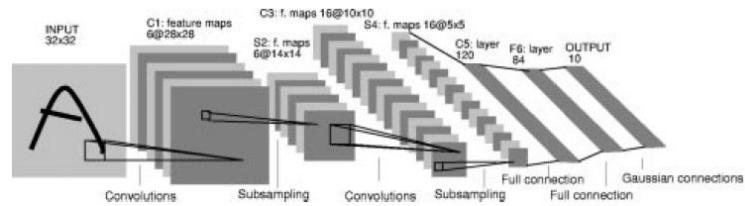
Autoencoders

Architecture:

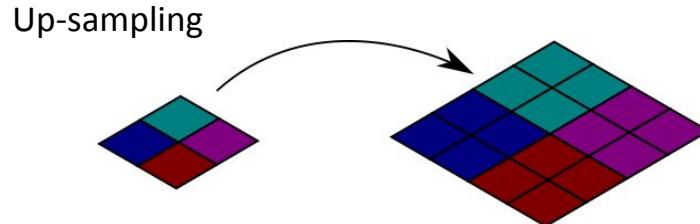
Encoder: ConvNet architecture



Convolving a 3×3 kernel over a 5×5 input using 2×2 strides



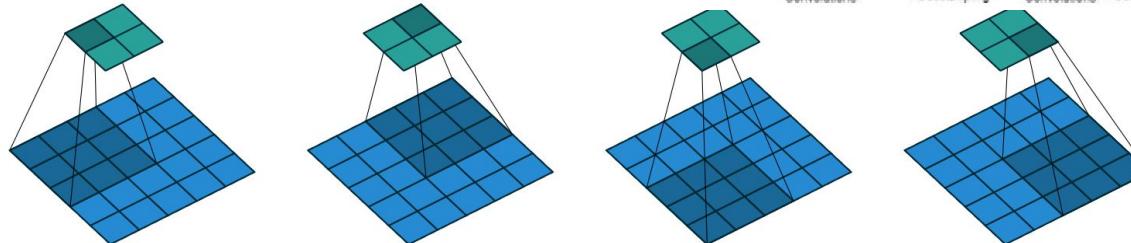
Decoder: How to go from a vector to an image? How to “deconvolve”?



Autoencoders

Architecture:

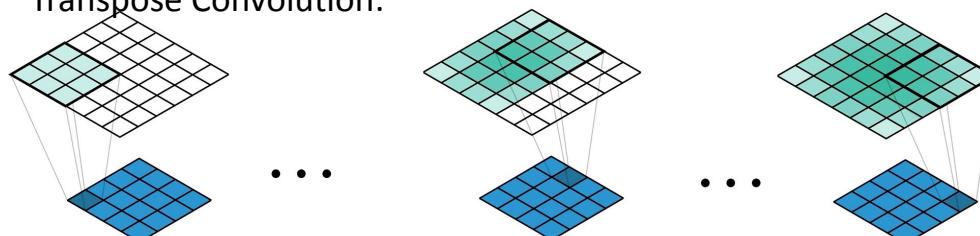
Encoder: ConvNet architecture



Convolving a 3×3 kernel over a 5×5 input using 2×2 strides

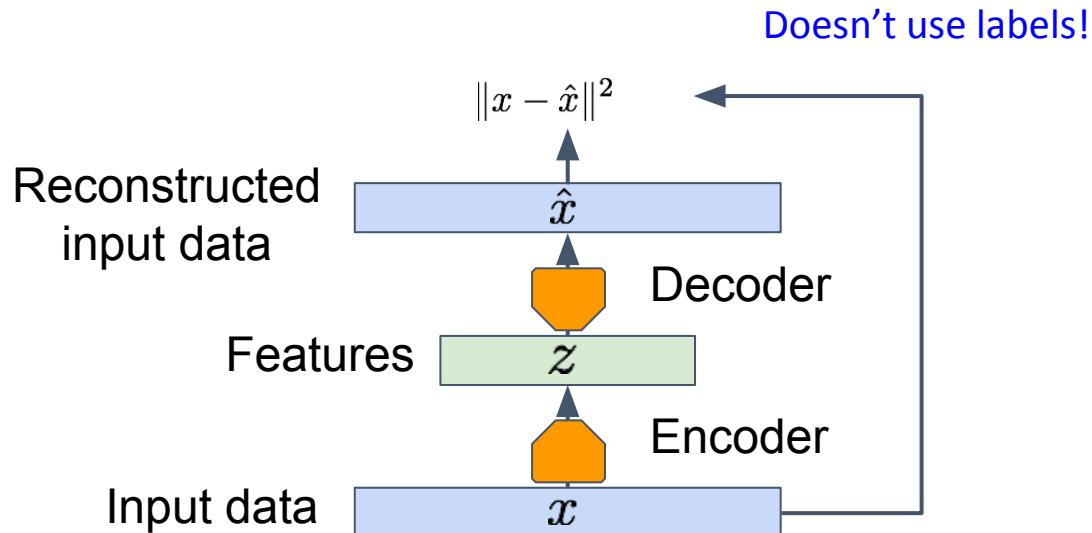
Decoder: How to go from a vector to an image? How to “deconvolve”?

Transpose Convolution:

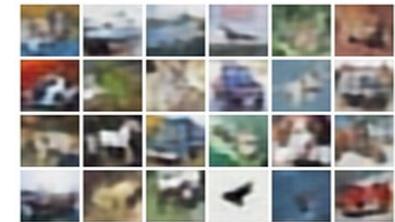


Autoencoders

Train such that features can be used to reconstruct original data



Reconstructed data



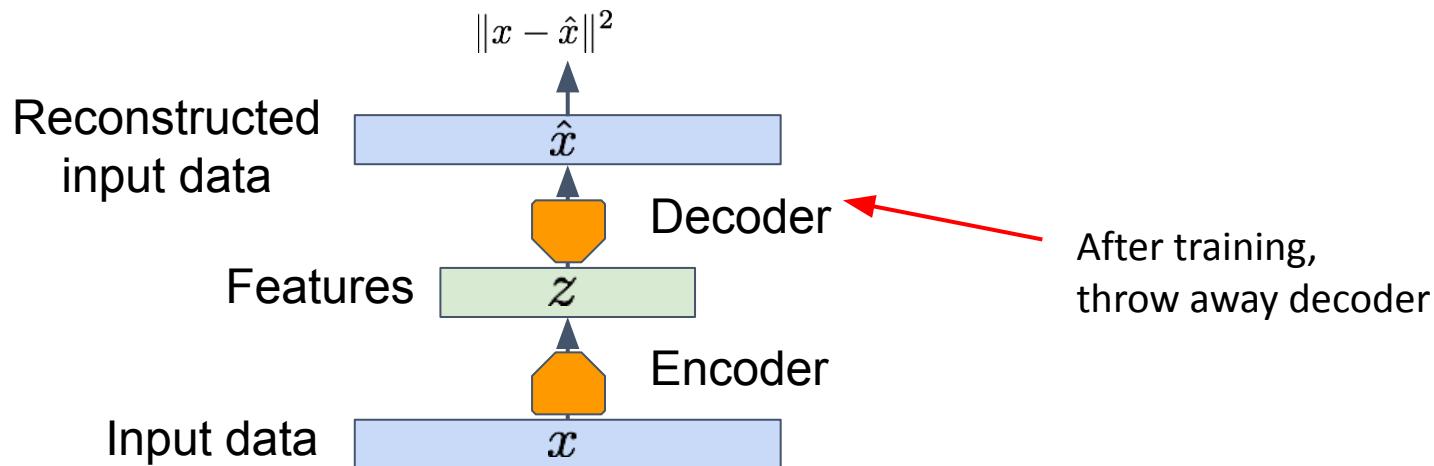
Input data



What can we do now?

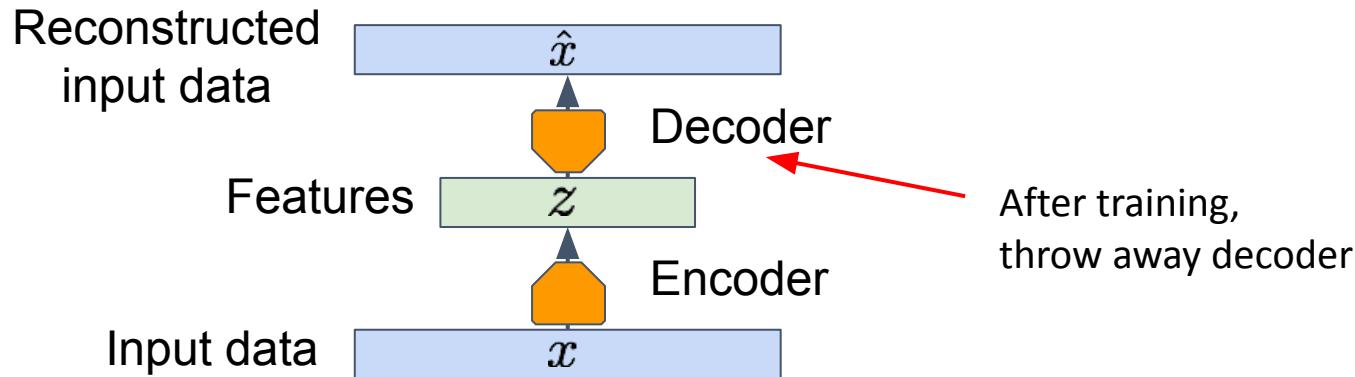
Autoencoders

Encoder can be used to initialize a **supervised** model



Autoencoders

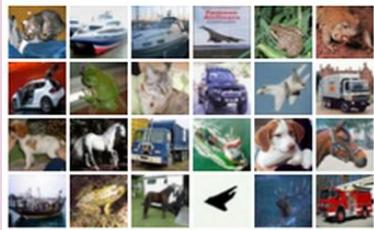
Encoder can be used to initialize a **supervised** model



Autoencoders

Encoder can be used to initialize a **supervised** model.

Unlabeled dataset



bird plane
dog deer truck



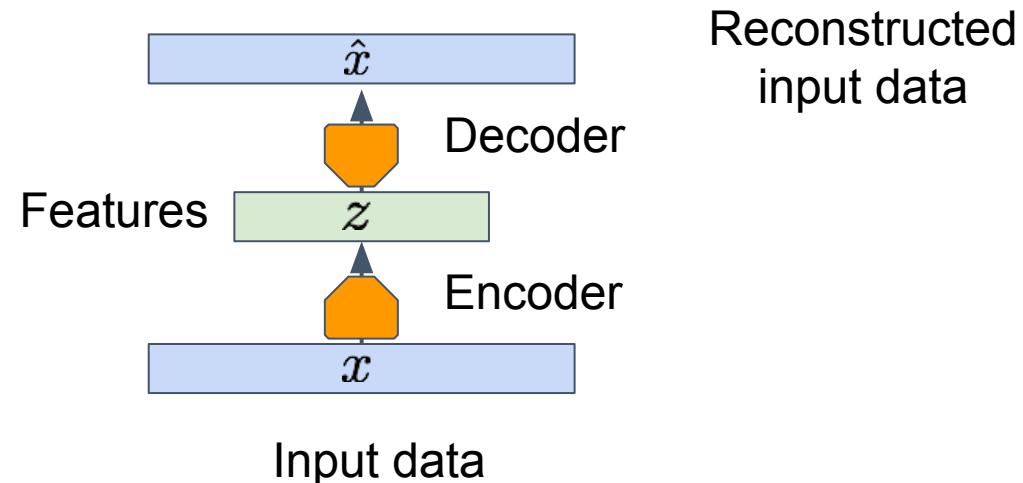
Small annotated
dataset

Autoencoders

Encoder can be used to initialize a **supervised** model.

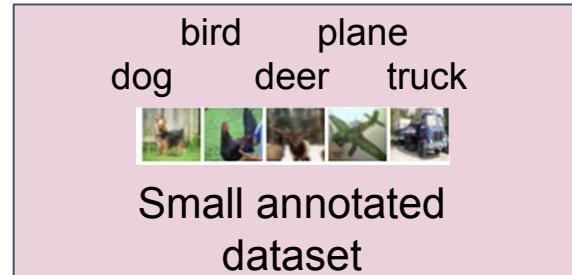


Pretrain an auto-encoder on the unlabeled data.

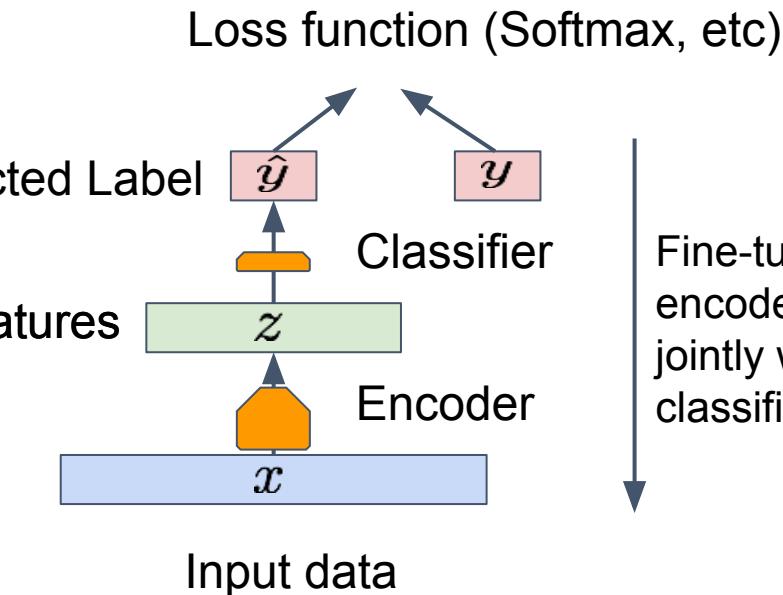


Autoencoders

Encoder can be used to initialize a **supervised** model.

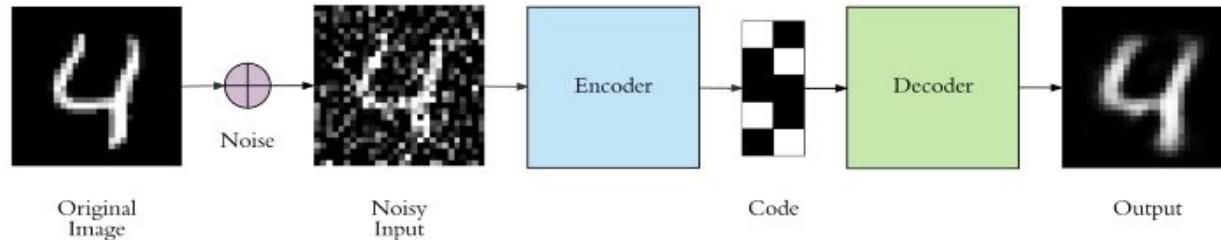


Train a classifier on the labeled data using the learned feature.



Denoising Autoencoders

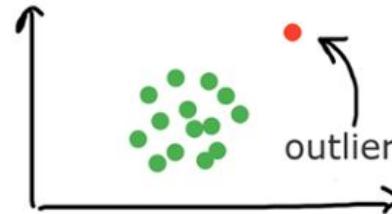
- Keeping the code layer small (latent representation) forces the autoencoder to learn an intelligent representation of the data.
- Another way to force the autoencoder to learn useful features is adding random noise to its inputs and making it recover the original noise-free data.



Autoencoders

Auto-Encoder can:

- Reduce dimension
- Unsupervised pre-training
- Denoising
- Detect outlier $\|x - \hat{x}\|^2$



Autoencoders

Auto-Encoder can:

- Reduce dimension
- Unsupervised pre-training
- Denoising
- Detect outlier

Auto-Encoder cannot:

- Estimate the density function $p(x)$
- Sample according to $p(x)$

Given training data, generate new samples from same distribution



Training data $\sim p_{\text{data}}(x)$



Generated samples $\sim p_{\text{model}}(x)$

Want to learn $p_{\text{model}}(x)$ similar to $p_{\text{data}}(x)$

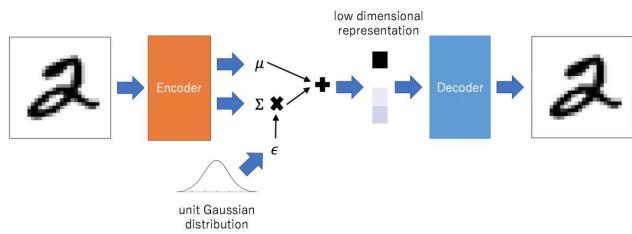
Why Generative Models?

- Realistic samples for artwork, super-resolution, colorization, etc.



- Generative models of time-series data can be used for simulation and planning (reinforcement learning applications!)
- Training generative models can also enable inference of **latent representations** that can be useful as general features

Variational Autoencoders (VAE)



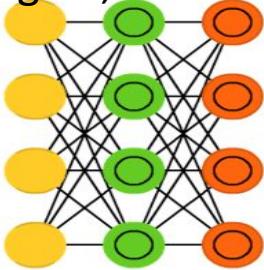
Variational Autoencoders

VAEs are a combination of the following ideas:

1. Autoencoders
2. Variational Approximation & Variational Lower Bound
3. “Reparameterization” Trick

Variational Autoencoders

D. Kingma, M. Welling, *Auto-Encoding Variational Bayes*, ICLR, 2014



Variational autoencoders (VAE) have the same architecture as AEs but are “taught” something else: an approximated probability distribution of the input samples. It’s a bit back to the roots as they are bit more closely related to BMs and RBMs. They do however rely on Bayesian mathematics regarding probabilistic inference and independence, as well as a re-parametrisation trick to achieve this different representation. The inference and independence parts make sense intuitively, but they rely on somewhat complex mathematics. The basics come down to this: take influence into account. If one thing happens in one place and something else happens somewhere else, they are not necessarily related. If they are not related, then the error propagation should consider that. This is a useful approach because neural networks are large graphs (in a way), so it helps if you can rule out influence from some nodes to other nodes as you dive into deeper layers.

Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." *arXiv preprint arXiv:1312.6114* (2013).

[Original Paper PDF](#)

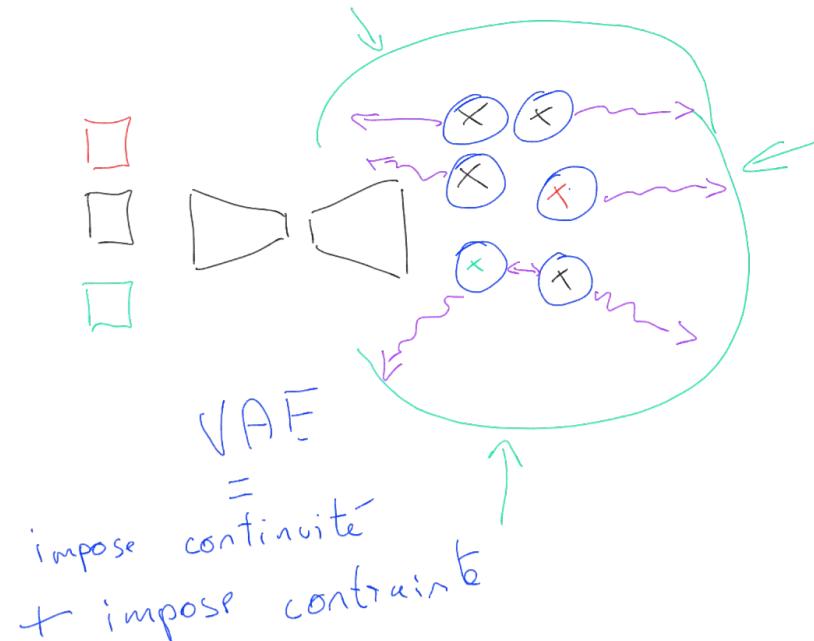
Variational Autoencoders

VAEs are a combination of the following ideas:

1. Autoencoders
2. Variational Approximation & Variational Lower Bound
3. “Reparameterization” Trick

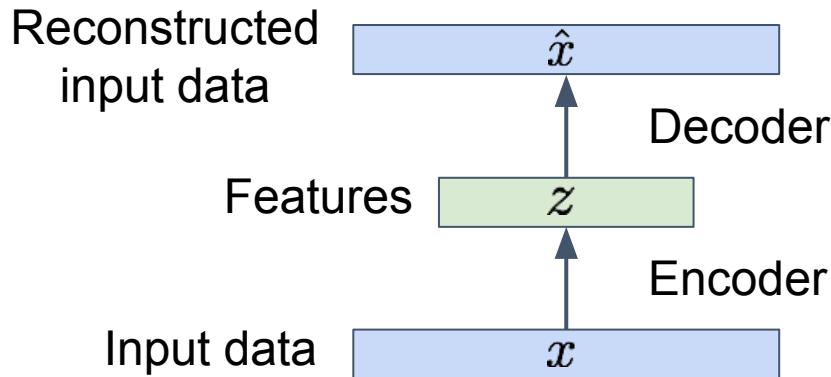
Autoencoders

Can we generate new images from an autoencoder?



Autoencoders

Can we generate new images from an autoencoder?

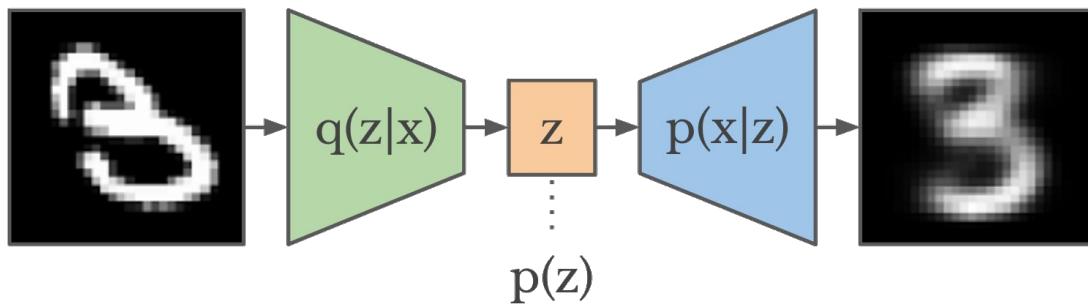


2 problems:

- How to sample Z
- No continuity in z space

Variational Autoencoders

Probabilistic spin on autoencoders - will let us sample from the model to generate data!



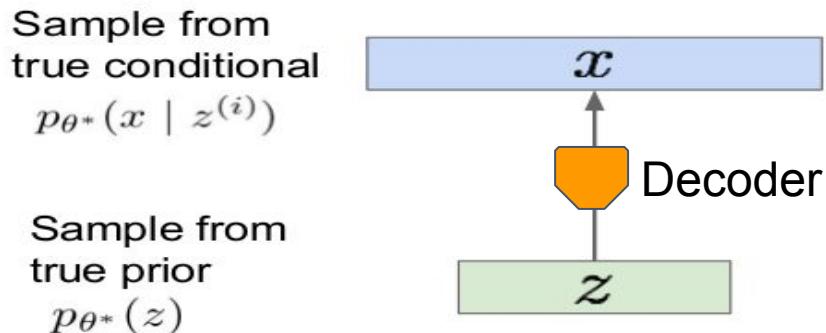
Variational Autoencoders

VAEs are a combination of the following ideas:

1. Autoencoders
2. **Variational Approximation & Variational Lower Bound**
3. “Reparameterization” Trick

Variational Autoencoders

Assume training data representation x is generated from underlying unobserved (latent) z

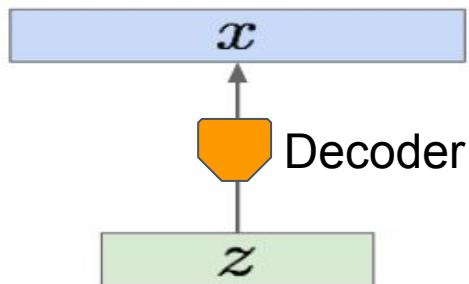


Variational Autoencoders

Assume training data representation x is generated from underlying unobserved (latent) z

Sample from
true conditional
 $p_{\theta^*}(x \mid z^{(i)})$

Sample from
true prior
 $p_{\theta^*}(z)$



Intuition:

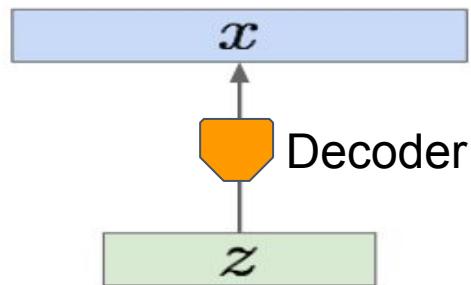
x is an image, z is a latent factor used to generate x : attributes, orientation, etc.

Variational Autoencoders

Assume training data representation x is generated from underlying unobserved (latent) z

Sample from
true conditional
 $p_{\theta^*}(x \mid z^{(i)})$

Sample from
true prior
 $p_{\theta^*}(z)$



Intuition:

x is an image, z is a latent factor used to generate x : attributes, orientation, etc.

We want to estimate the optimal parameters θ^* of this generative model, according to maximum likelihood

But what is a loss? But maximum? likelihood

- Examples with supervised training
 - Given training samples

$$T = \{(x_1, t_1), (x_2, t_2), \dots, (x_N, t_N)\}$$

- Choose an objective function (loss function)
- Adjust all the weights of the network θ such that a cost function is minimized

$$\begin{aligned} \min_w &= L_{tot}(w) \\ &= \sum_i L(t_i, y(x_i; w)) \end{aligned}$$

Loss: Regression

The choice of the loss function depends on the problem.

MSE loss:

$$\frac{1}{2} \sum_{n=1}^N \|y(\mathbf{x}_n, \mathbf{w}) - t_n\|^2$$

But why?

Let's assume a gaussian distribution of the error:

$$p(t|\mathbf{x}, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

Dataset likelihood:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \beta)$$

Loss: Regression but why MSE?

Let's assume a gaussian distribution of the error:

$$p(t|\mathbf{x}, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

Dataset likelihood:

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}, \beta)$$

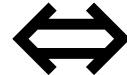
We take -log:

$$\frac{\beta}{2} \sum_{n=1}^N \{y(\mathbf{x}_n, \mathbf{w}) - t_n\}^2 - \frac{N}{2} \ln \beta + \frac{N}{2} \ln(2\pi)$$

Loss: Regression but why MSE?

Maximize the data log-likelihood assuming gaussian errors

$$p(t|\mathbf{x}, \mathbf{w}) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$



Minimize the MSE:

$$\frac{1}{2} \sum_{n=1}^N \|\mathbf{y}(\mathbf{x}_n, \mathbf{w}) - \mathbf{t}_n\|^2$$

Loss: classification

- Classification problem with 'C' classes: negative loglikelihood

$$\text{loss}(x) = -\log \left(\text{prob}(c = c_{\text{True}} | x) \right)$$

- How to get probabilities:

$$\text{Softmax}(x) = \frac{\exp(x)}{\sum_c \exp(x[c])} \quad \sum_{c=0}^C \text{Softmax}(x)[c] = 1$$

- Cross-entropy loss

$$\text{loss}(x) = -\log (\text{Softmax}(x)[c_{\text{True}}])$$

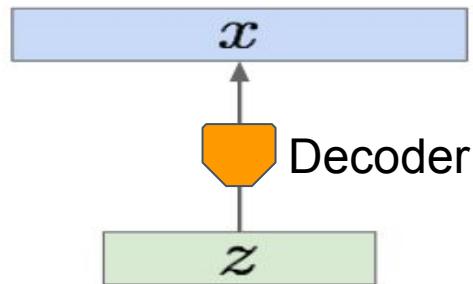
Variational Autoencoders

How should we **train** this model?

Learn model parameters to maximize likelihood of training data

Sample from
true conditional
 $p_{\theta^*}(x \mid z^{(i)})$

Sample from
true prior
 $p_{\theta^*}(z)$



$$p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Variational Autoencoders

$$\text{Data likelihood: } p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Intractable to compute
 $p(x|z)$ for every z !

Simple Gaussian
prior

Decoder neural
network



Variational Autoencoders

$$\text{Data likelihood: } p_{\theta}(x) = \int p_{\theta}(z)p_{\theta}(x|z)dz$$

Intractable to compute
 $p(x|z)$ for every z !

Simple Gaussian
prior

Decoder neural
network

$$\text{Posterior density also intractable: } p_{\theta}(z|x) = p_{\theta}(x|z)p_{\theta}(z)/p_{\theta}(x)$$

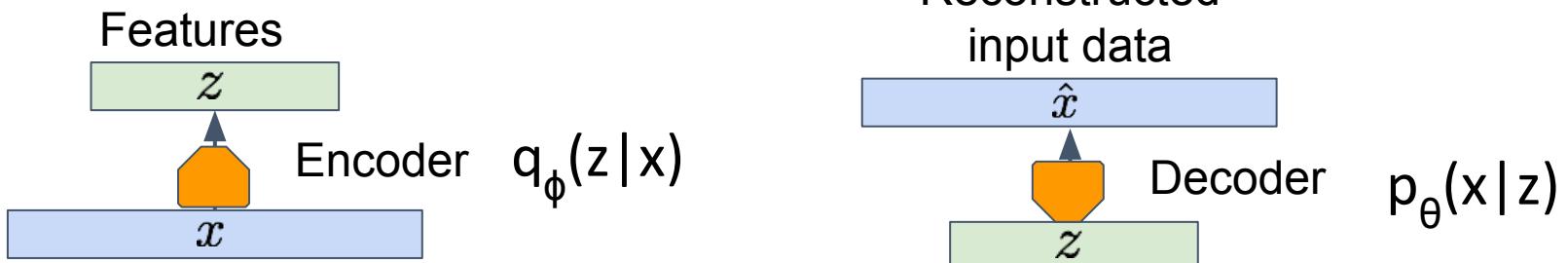
Decoder neural
network

Simple gaussian
prior

Intractable data
likelihood

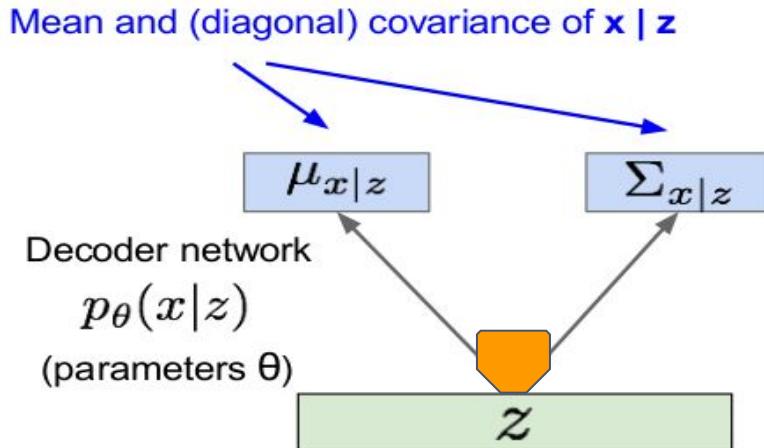
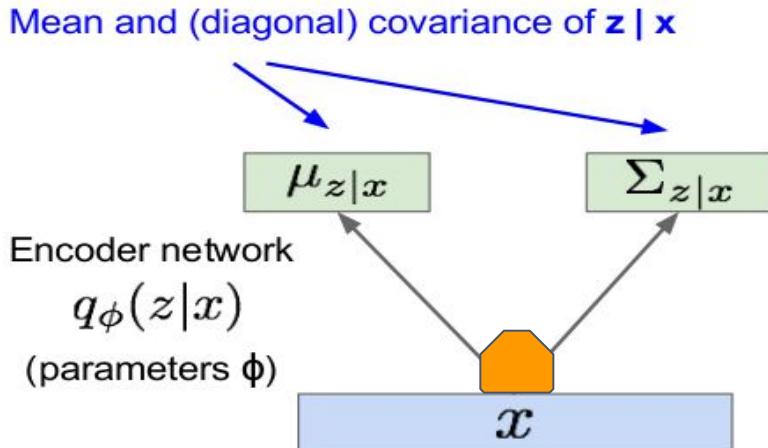
Variational Approximation

- Solution: In addition to decoder network modeling $p_\theta(x|z)$, define **additional encoder network** $q_\phi(z|x)$ that approximates $p_\theta(z|x)$
- This allows us to derive a **lower bound** on the data likelihood that is tractable, which we can optimize



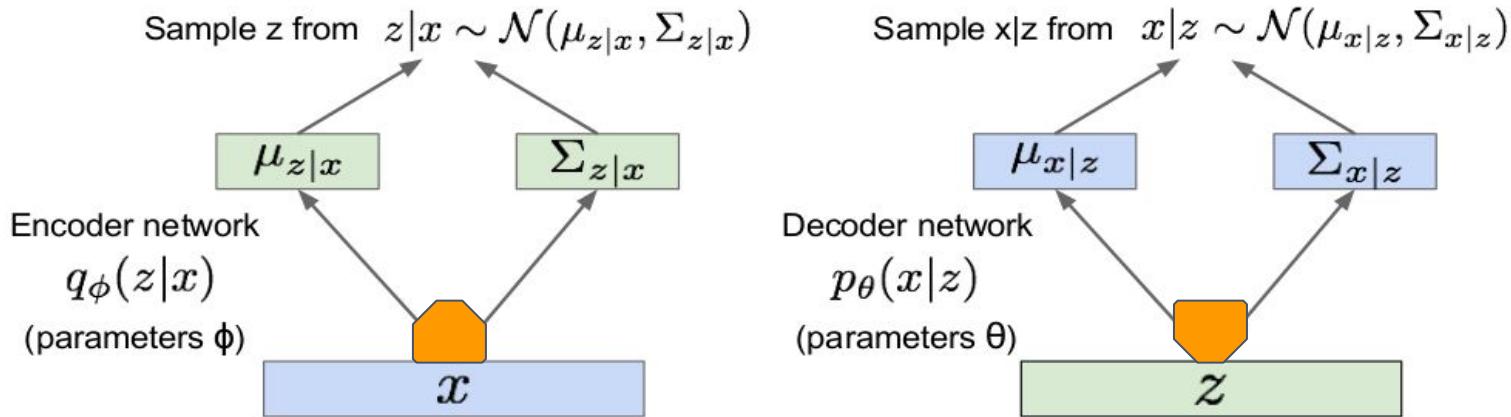
Variational Approximation

- Since we are modeling probabilistic generation of data, encoder and decoder networks are probabilistic



Variational Approximation

- Since we are modeling probabilistic generation of data, encoder and decoder networks are probabilistic



- Encoder and decoder networks also called “recognition”/“inference” and “generation” networks

Variational Approximation

- Let us consider the log likelihood of the data

Variational Approximation

- Let us consider the log likelihood of the data

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z | x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))\end{aligned}$$

The expectation wrt. z (using encoder network) let us write nice KL terms

Variational Approximation

- Let us consider the log likelihood of the data

$$\begin{aligned}\log p_{\theta}(x^{(i)}) &= \mathbf{E}_{z \sim q_{\phi}(z | x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))\end{aligned}$$

↑
Decoder network gives $p_{\theta}(x|z)$, can compute estimate of this term through sampling. (Sampling differentiable through reparam. trick, see paper.)

↑
This KL term (between Gaussians for encoder and z prior) has nice closed-form solution!

↑
 $p_{\theta}(z|x)$ intractable (saw earlier), can't compute this KL term :(But we know KL divergence always ≥ 0 .

Variational Approximation

- Let us consider the log likelihood of the data

We want to maximize the data likelihood

$$\begin{aligned}\log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\ &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\ &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\ &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}\end{aligned}$$

Tractable lower bound which we can take gradient of and optimize! ($p_\theta(x|z)$ differentiable, KL term differentiable)

Variational Approximation

- Let us consider the log likelihood of the data

$$\log p_\theta(x^{(i)}) = \mathbf{E}_{z \sim q_\phi(z | x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z)$$

↗

$$= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule})$$

We want to
maximize the
data
likelihood

$$= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z) p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant})$$

$$= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z)) + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{> 0}$$

$$\log p_\theta(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Variational Approximation

- Let us consider the log likelihood of the data

$$\log p_{\theta}(x^{(i)}) = \mathbf{E}_{z \sim q_{\phi}(z|x^{(i)})} [\log p_{\theta}(x^{(i)})] \quad (p_{\theta}(x^{(i)}) \text{ Does not depend on } z)$$

$$= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Bayes' Rule})$$

Reconstruct the input data $= \mathbf{E}_z \left[\log \frac{p_{\theta}(x^{(i)} | z)p_{\theta}(z)}{p_{\theta}(z | x^{(i)})} \frac{q_{\phi}(z | x^{(i)})}{q_{\phi}(z | x^{(i)})} \right]$ (Multiply by constant)

$$= \mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z)} \right] + \mathbf{E}_z \left[\log \frac{q_{\phi}(z | x^{(i)})}{p_{\theta}(z | x^{(i)})} \right] \quad (\text{Logarithms})$$

$$= \underbrace{\mathbf{E}_z [\log p_{\theta}(x^{(i)} | z)]}_{\mathcal{L}(x^{(i)}, \theta, \phi)} - D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z)) + \underbrace{D_{KL}(q_{\phi}(z | x^{(i)}) || p_{\theta}(z | x^{(i)}))}_{> 0}$$

Make approximate posterior distribution close to prior

$$\log p_{\theta}(x^{(i)}) \geq \mathcal{L}(x^{(i)}, \theta, \phi)$$

Variational lower bound ("ELBO")

$$\theta^*, \phi^* = \arg \max_{\theta, \phi} \sum_{i=1}^N \mathcal{L}(x^{(i)}, \theta, \phi)$$

Training: Maximize lower bound

Variational Autoencoders: recap

- Let us consider the log likelihood of the data



Variational Autoencoders

Maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_{\theta}(x^{(i)} \mid z) \right] - D_{KL}(q_{\phi}(z \mid x^{(i)}) \parallel p_{\theta}(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Variational Autoencoders

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

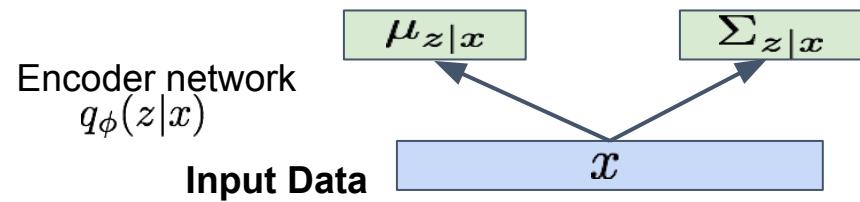
Let's look at computing the bound
(forward pass) for a given minibatch of
input data



Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbf{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$



Variational Autoencoders

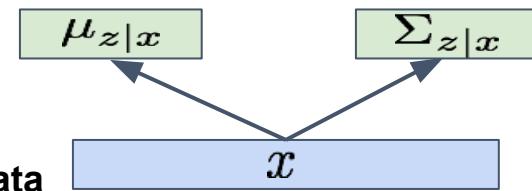
Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Encoder network
 $q_\phi(z|x)$

Input Data
 x

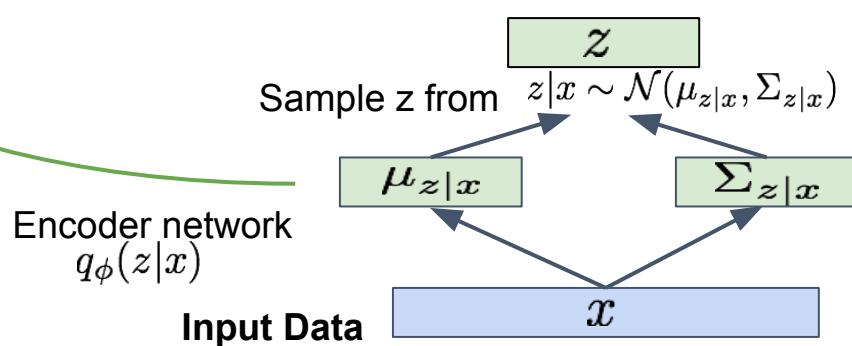


Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

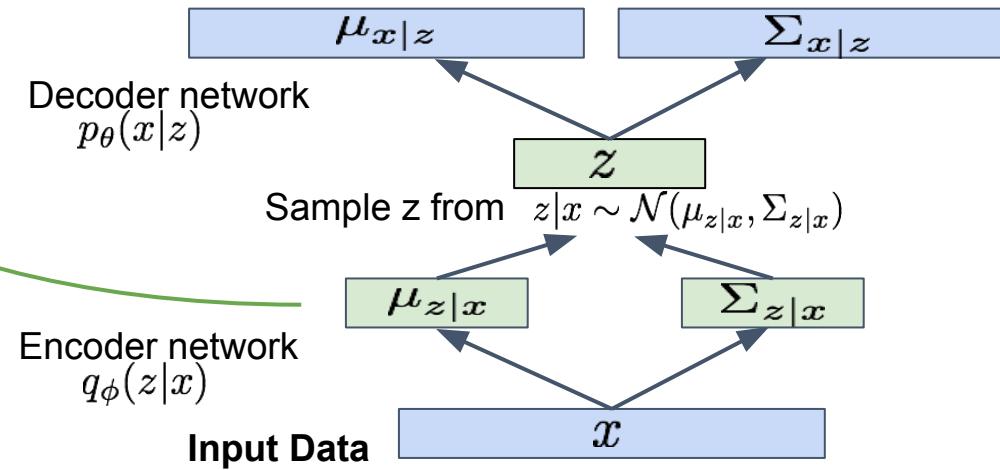


Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior



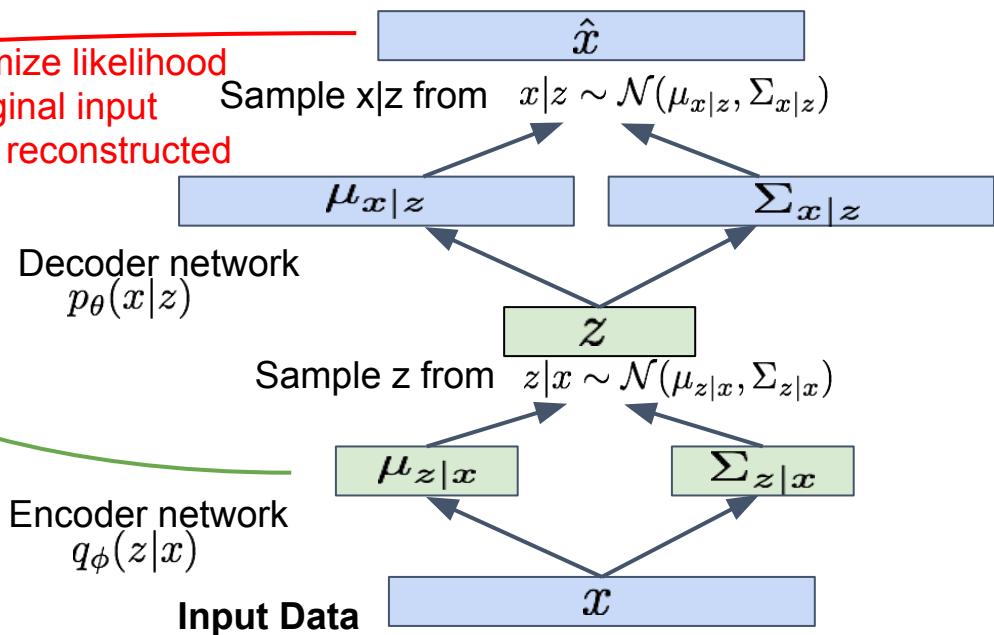
Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

Maximize likelihood
of original input
being reconstructed



Variational Autoencoders

Putting it all together: maximizing the likelihood lower bound

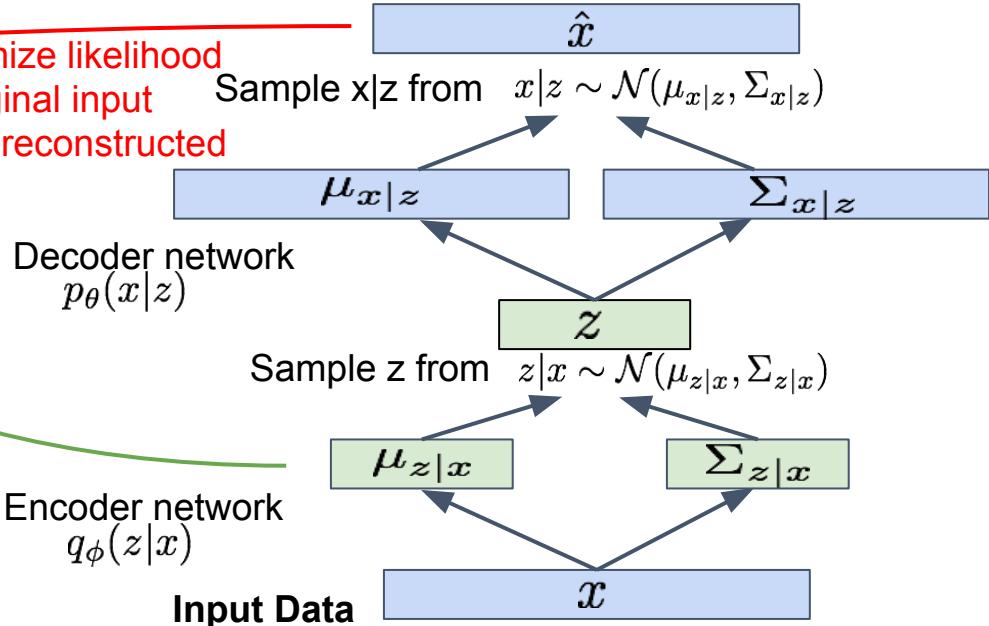
$$\underbrace{\mathbb{E}_z \left[\log p_\theta(x^{(i)} | z) \right] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)}$$

Make approximate posterior distribution close to prior

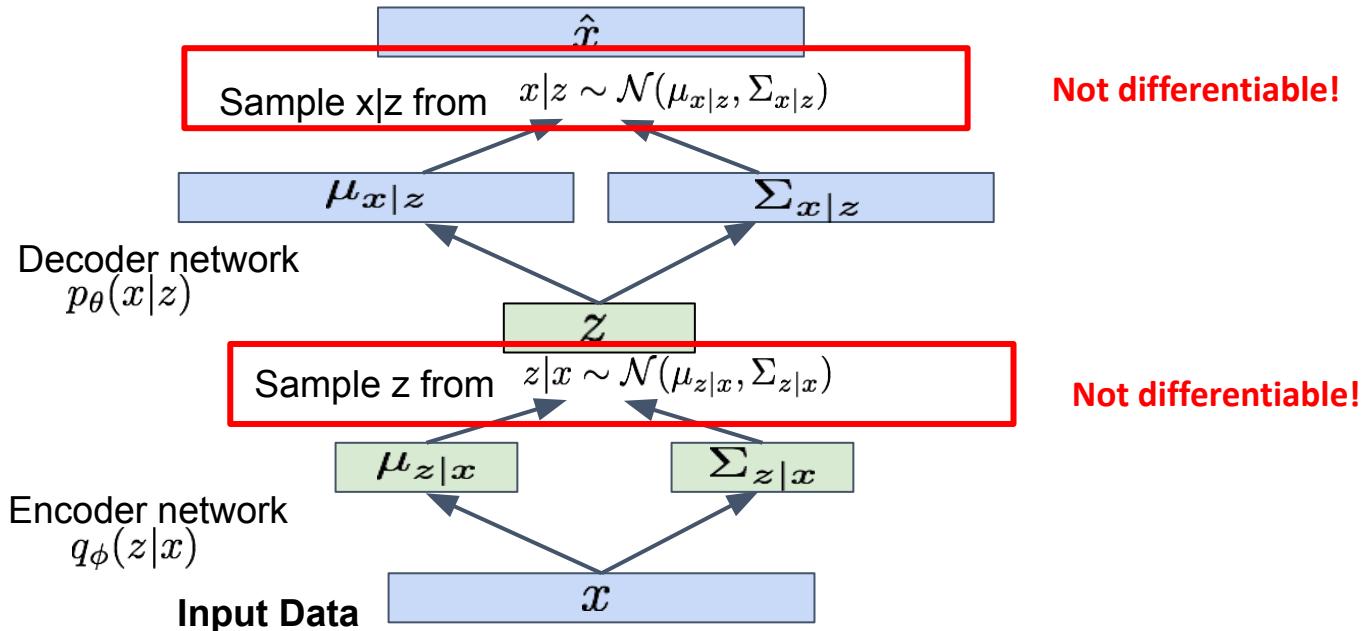
For every minibatch of input data: compute this forward pass, and then backprop!

Employ “Reparameterization” Trick

Maximize likelihood of original input being reconstructed



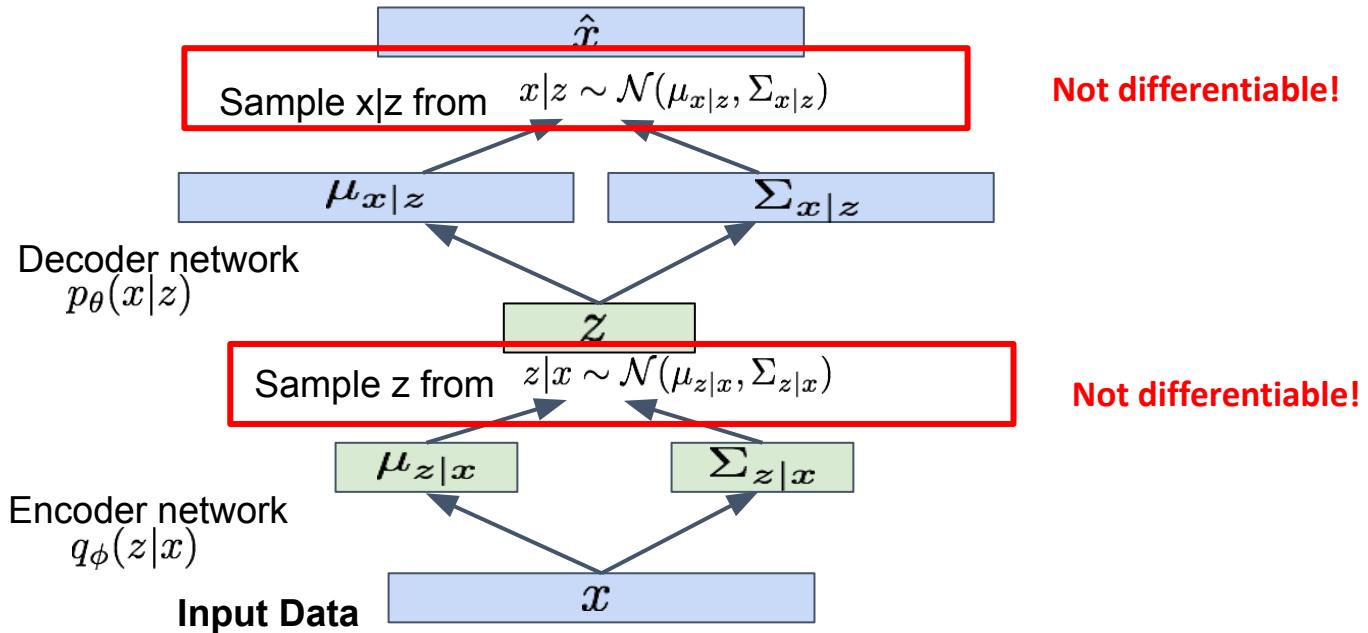
Variational Autoencoders



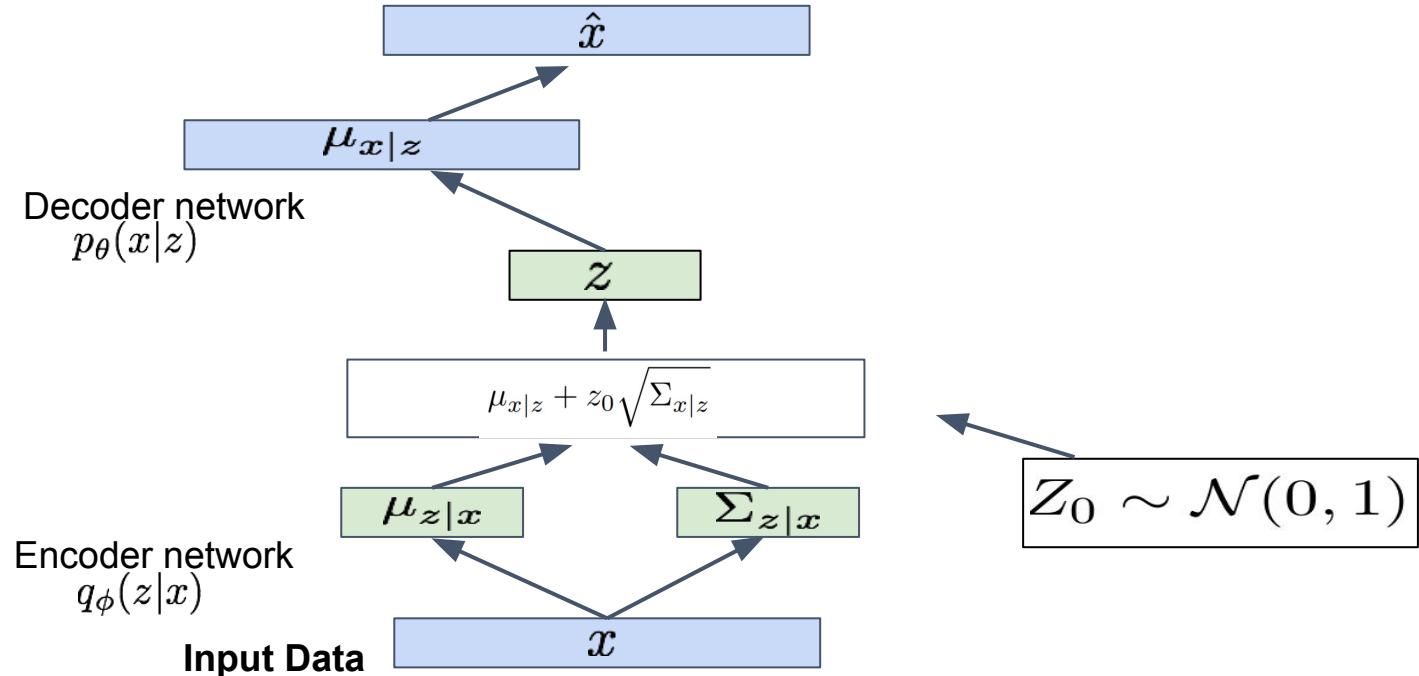
Re-parametrization trick:

If $z_0 \sim \mathcal{N}(0, 1)$, then $z = \mu_{x|z} + z_0 \sqrt{\Sigma_{x|z}} \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

Variational Autoencoders

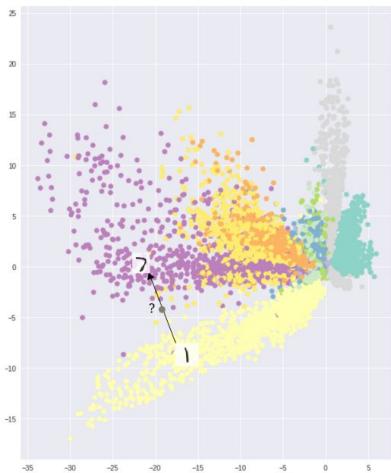


Variational Autoencoders

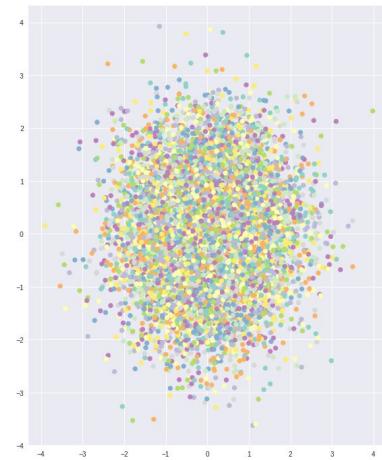


Variational Autoencoders

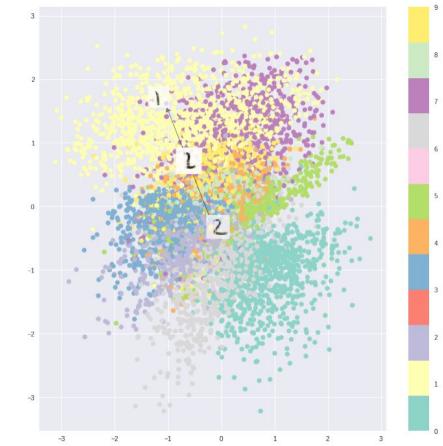
Autoencoders vs Variational Autoencoders



AE (reconstruction)



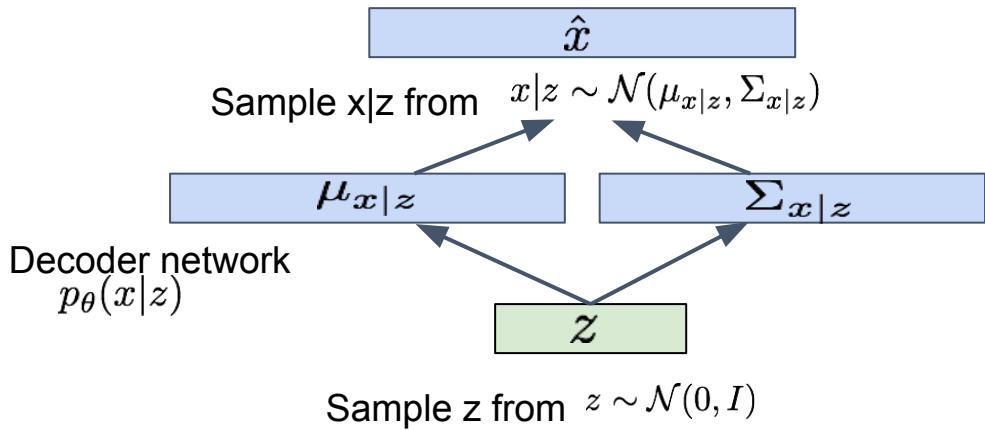
Only KL term



VAE (reconstruction+KL)

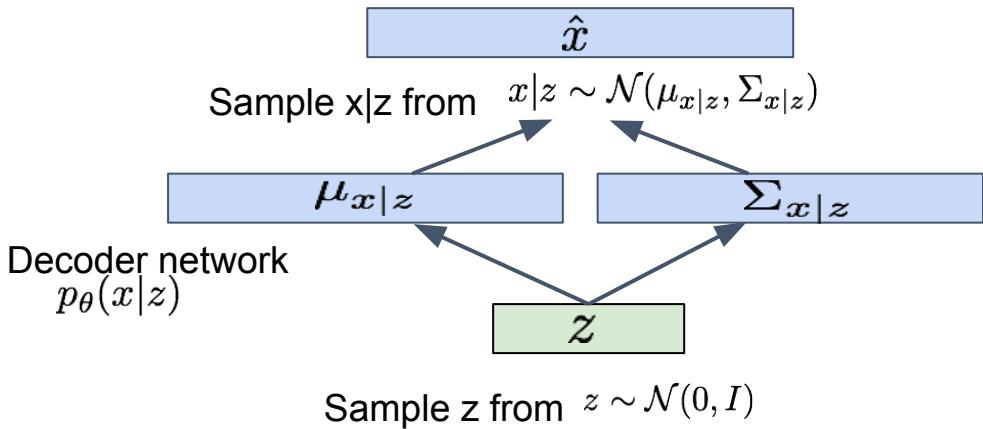
Generating Data

Use decoder network. Now sample z from prior!



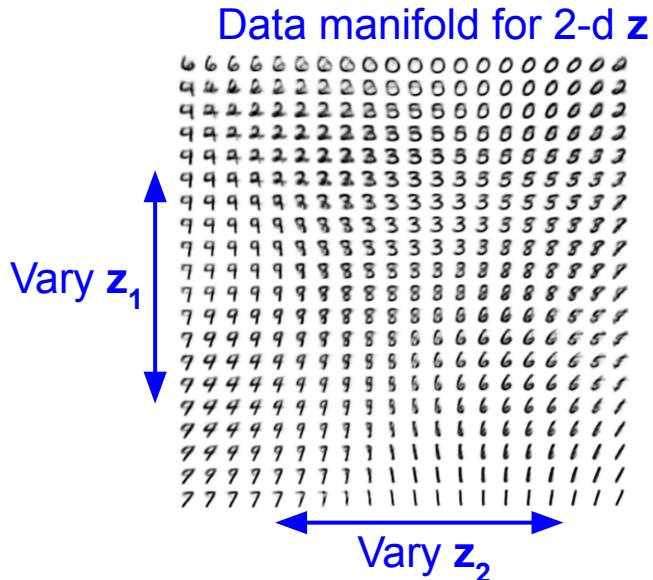
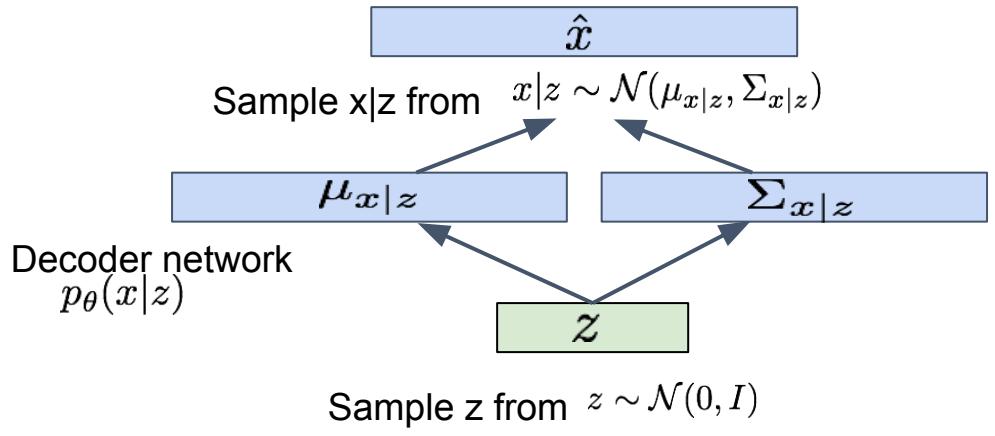
Generating Data

Use decoder network. Now sample z from prior!



Generating Data

Use decoder network. Now sample z from prior!

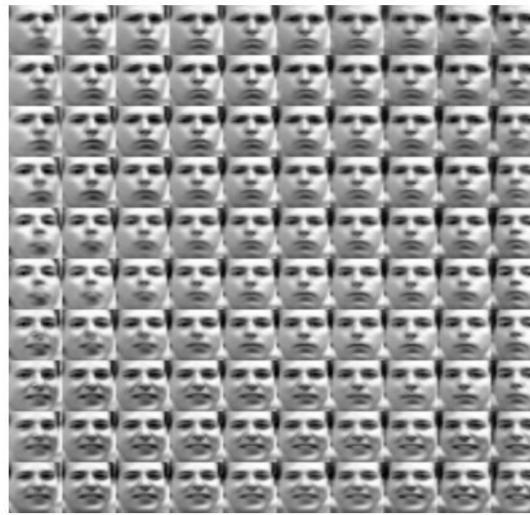


Generating Data

Diagonal prior on $z \Rightarrow$
independent latent
variables

Different dimensions of z
encode interpretable
factors of variation

Degree of smile
Vary z_1



Vary z_2 Head pose

Generating Data

Diagonal prior on $\mathbf{z} \Rightarrow$
independent latent
variables

Different dimensions of \mathbf{z}
encode interpretable
factors of variation

Also good feature representation that
can be computed using $q_{\phi}(z|x)$!

Degree of smile
Vary z_1



Vary z_2 Head pose

Generating Data: interpolation

Source



Target



Generating Data



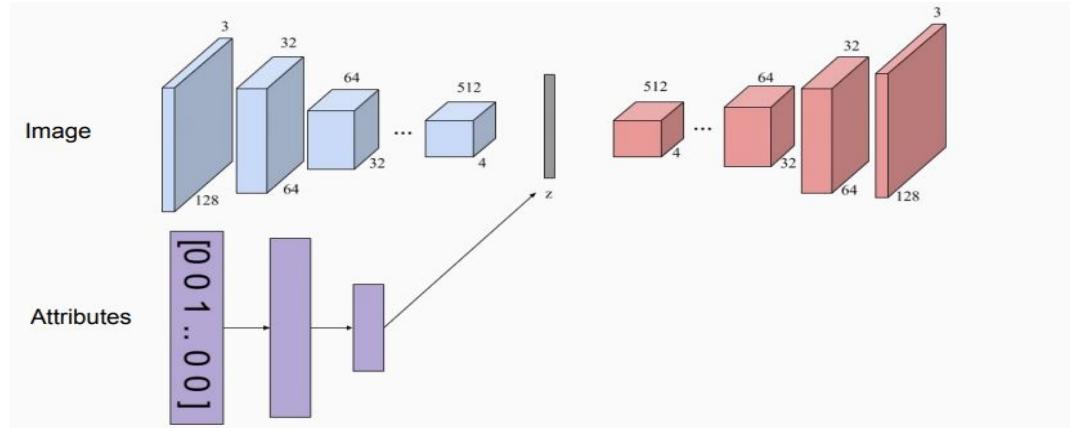
32x32 CIFAR-10



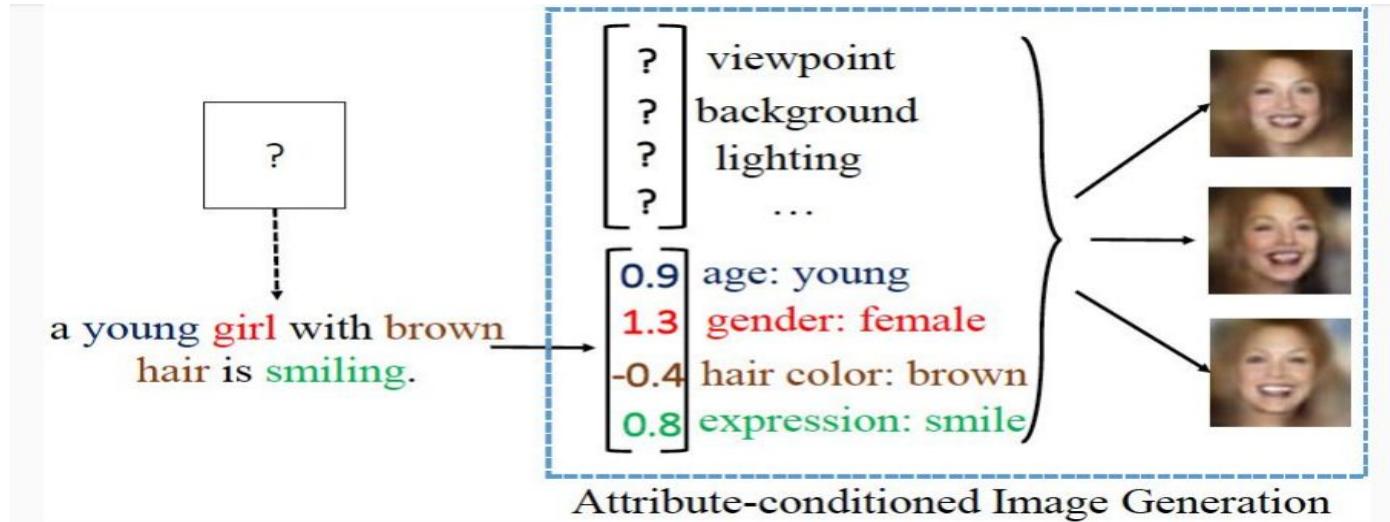
Labeled Faces in the Wild

Conditional VAE

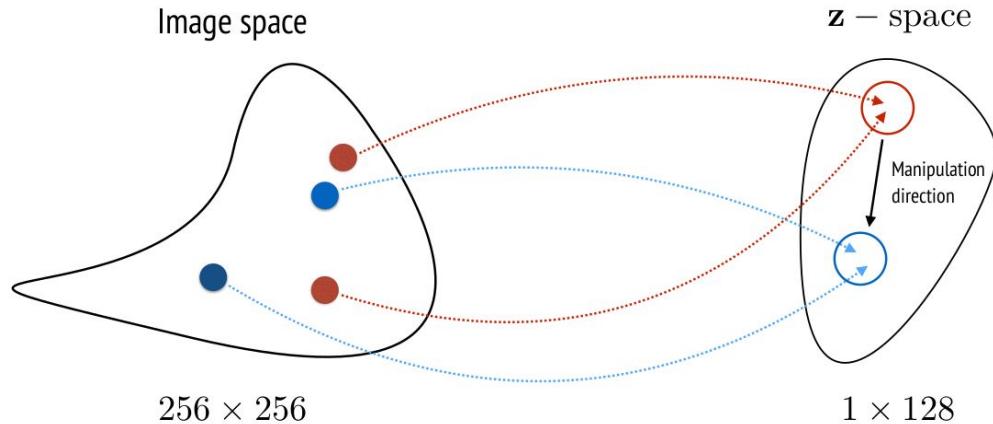
- What if we have labels? (e.g. digit labels or attributes) or other inputs we wish to condition on?
- None of the derivation changes.
- Replace all $p(x/z)$ with $p(x/z,y)$.
- Replace all $q(x/z)$ with $q(x/z,y)$.
- Go through the same KL divergence procedure, to get the same lower bound.



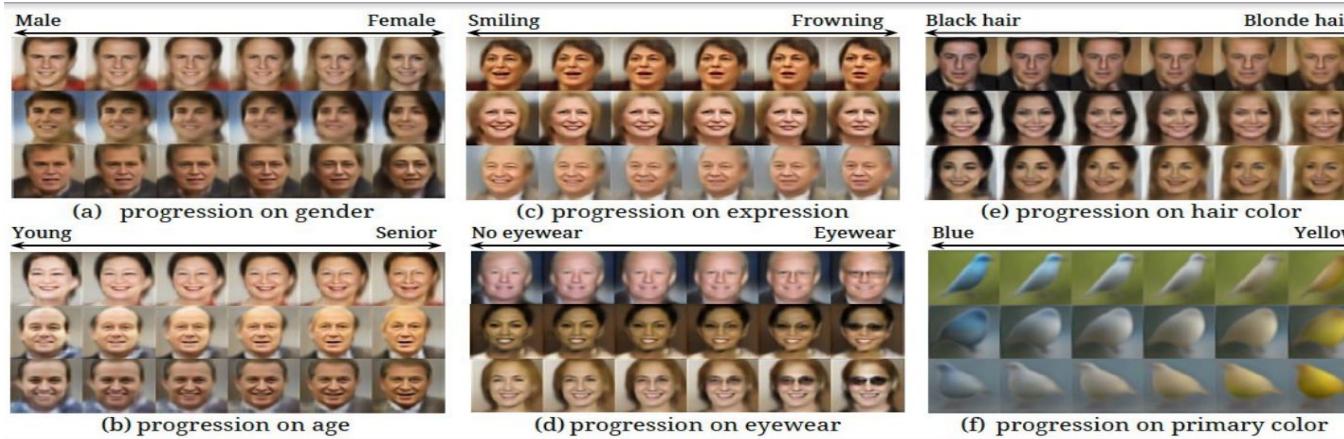
Generating Data: attribute-driven



Generating Data: attribute-driven



Generating Data: attribute-driven



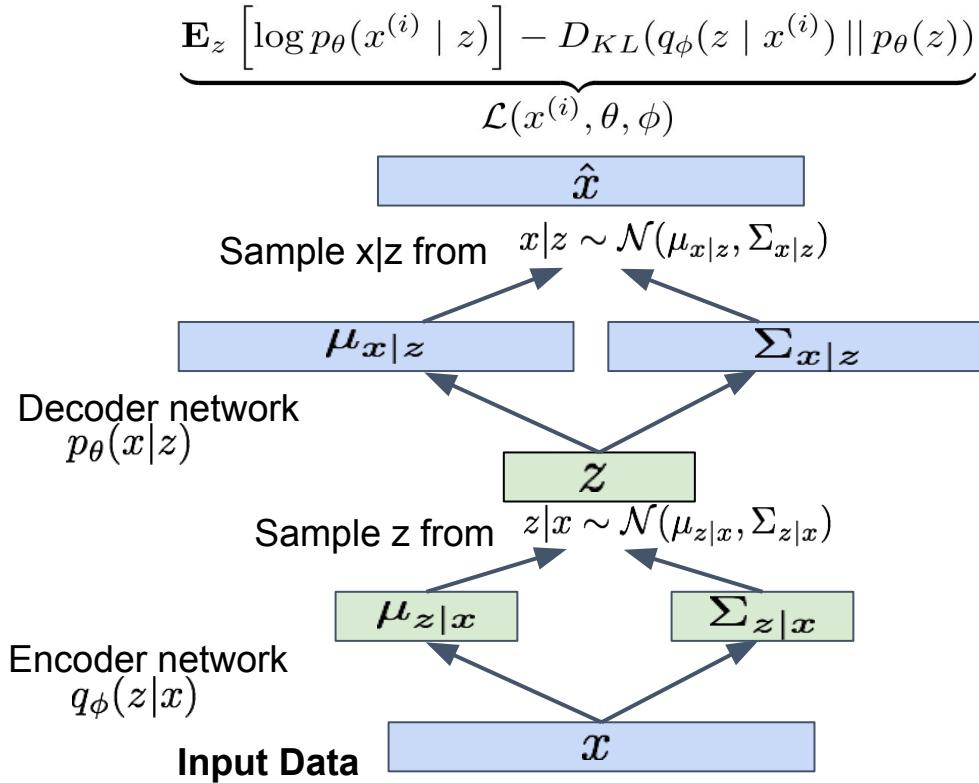
Variational Autoencoders

- Probabilistic spin to traditional autoencoders => allows generating data
- Defines an intractable density => derive and optimize a (variational) lower bound
- **Pros:**
 - Principled approach to generative models
 - Allows inference of $q(z|x)$, can be useful feature representation for other tasks
- **Cons:**
 - Maximizes lower bound of likelihood: okay, but not as good evaluation as PixelRNN/PixelCNN
 - Samples blurrier and lower quality compared to state-of-the-art (GANs)
- **Active areas of research:**
 - More flexible approximations, e.g. richer approximate posterior instead of diagonal Gaussian
 - Incorporating structure in latent variables

Variational Autoencoders: recap



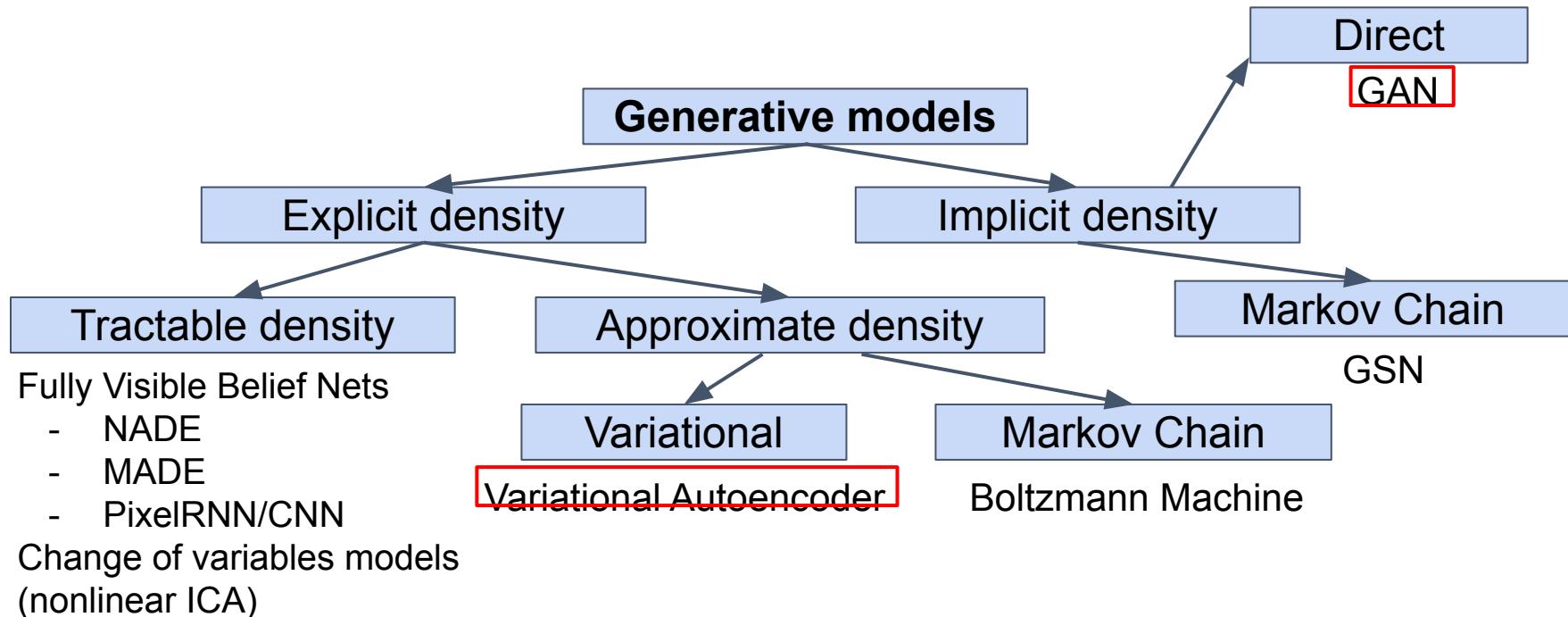
Variational Autoencoders: recap



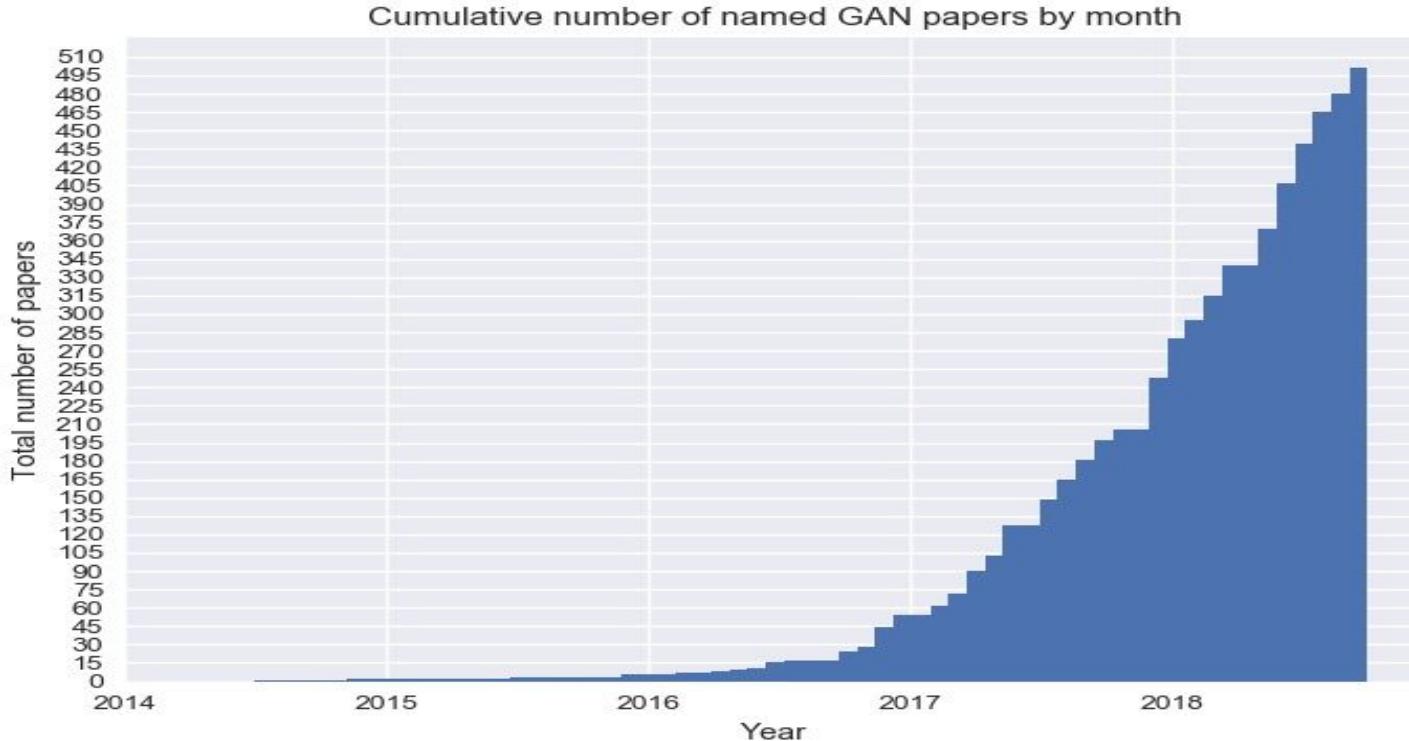
Useful Links

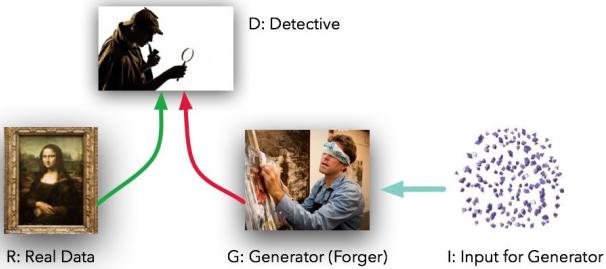
- D. Kingma, M. Welling, *Auto-Encoding Variational Bayes*, ICLR, 2014
- Carl Doersch, *Tutorial on Variational Autoencoders* arXiv, 2016
- Xincheng Yan, Jimei Yang, Kihyuk Sohn, Honglak Lee, *Attribute2Image: Conditional Image Generation from Visual Attributes*, ECCV, 2016
- Jacob Walker, Carl Doersch, Abhinav Gupta, Martial Hebert, *An Uncertain Future: Forecasting from Static Images using Variational Autoencoders*, ECCV, 2016
- Anders Boesen Lindbo Larsen, Søren Kaae Sønderby, Hugo Larochelle, Ole Winther, *Autoencoding beyond pixels using a learned similarity metric*, ICML, 2016
- Aditya Deshpande, Jiajun Lu, Mao-Chuang Yeh, David Forsyth, *Learning Diverse Image Colorization*, arXiv, 2016
- Raymond Yeh, Ziwei Liu, Dan B Goldman, Aseem Agarwala, *Semantic Facial Expression Editing using Autoencoded Flow*, arXiv, 2016

Taxonomy of Generative Models



GAN zoo





Generative Adversarial Networks (GAN)

Generative Adversarial Networks

Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution, e.g. random noise.
Learn transformation to training distribution.

Q: What can we use to
represent this complex
transformation?

Generative Adversarial Networks

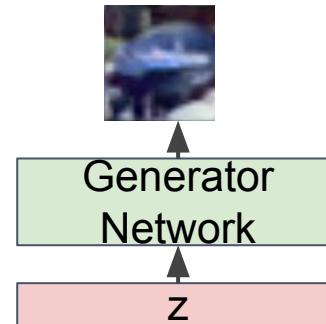
Problem: Want to sample from complex, high-dimensional training distribution. No direct way to do this!

Solution: Sample from a simple distribution, e.g. random noise.
Learn transformation to training distribution.

Q: What can we use to represent this complex transformation?
A: A neural network!

Output: Sample from training distribution

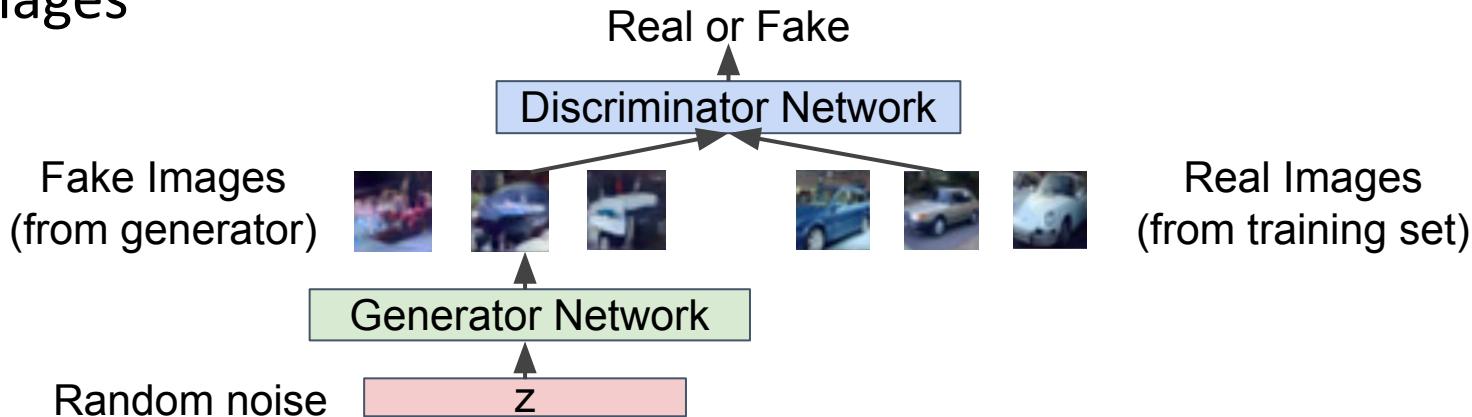
Input: Random noise



Training GANs: Two-player game

Generator network: try to fool the discriminator by generating real-looking images

Discriminator network: try to distinguish between real and fake images



Training GANs: Two-player game

Generator network: try to fool the discriminator by generating real-looking images

Discriminator network: try to distinguish between real and fake images

Train jointly in **minimax game**

Minimax objective function: Discriminator outputs likelihood in (0,1) of real image

$$\min_{\theta_g} \max_{\theta_d} \left[\underbrace{\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x)}_{\text{Discriminator output for real data } x} + \underbrace{\mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))}_{\text{Discriminator output for generated fake data } G(z)} \right]$$

- Discriminator (θ_d) wants to **maximize objective** such that $D(x)$ is close to 1 (real) and $D(G(z))$ is close to 0 (fake)
- Generator (θ_g) wants to **minimize objective** such that $D(G(z))$ is close to 1 (discriminator is fooled into thinking generated $G(z)$ is real)

Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

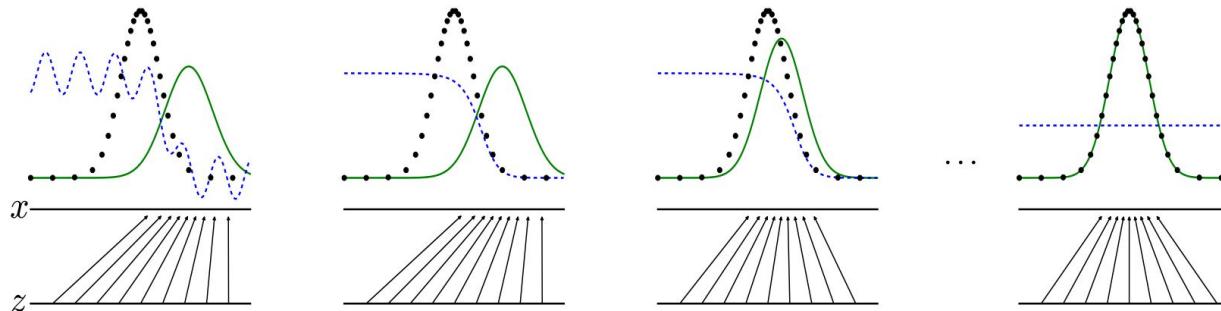
Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$



Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

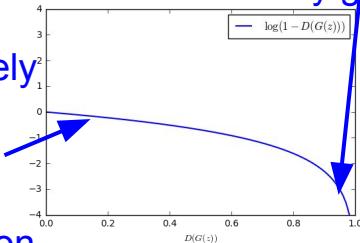
$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

In practice, optimizing this generator objective does not work well!

When sample is likely fake, want to learn from it to improve generator. But gradient in this region is relatively flat!



Gradient signal dominated by region where sample is already good

Training GANs: Two-player game

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

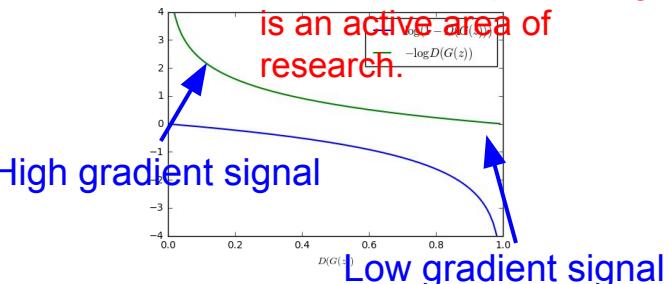
2. Instead: Gradient ascent on generator, different objective

$$\max_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(D_{\theta_d}(G_{\theta_g}(z)))$$

Instead of minimizing likelihood of discriminator being correct, now maximize likelihood of discriminator being wrong.

Same objective of fooling discriminator, but now higher gradient signal for bad samples => works much better! Standard in practice.

Aside: Jointly training two networks is challenging, can be unstable. Choosing objectives with better loss landscapes helps training,



GAN Training Algorithm

- Update discriminator
 - Repeat for k steps:
 - Sample mini-batch of noise samples z_1, \dots, z_m and mini-batch of real samples x_1, \dots, x_m
 - Update parameters of D by stochastic gradient ascent on
$$\frac{1}{m} \sum_m [\log D(x_m) + \log(1 - D(G(z_m)))]$$
 - Update generator
 - Sample mini-batch of noise samples z_1, \dots, z_m
 - Update parameters of G by stochastic gradient ascent on

$$\frac{1}{m} \sum_m \log D(G(z_m))$$

GAN Results

MNIST digits



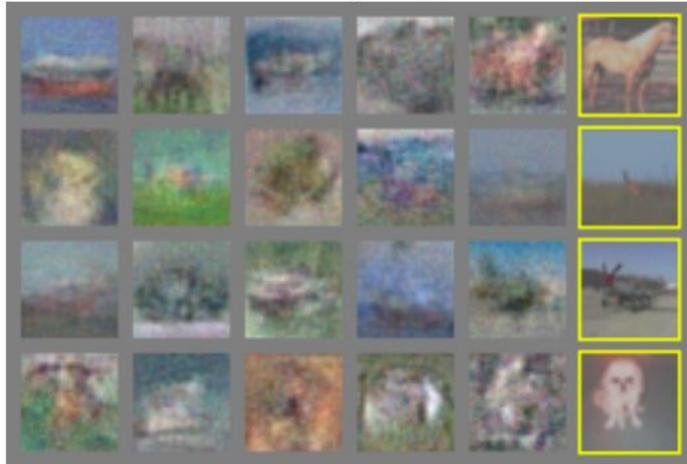
Toronto Face Dataset



Nearest real image for
sample to the left

GAN Results

CIFAR-10 (FC networks)



CIFAR-10 (conv networks)

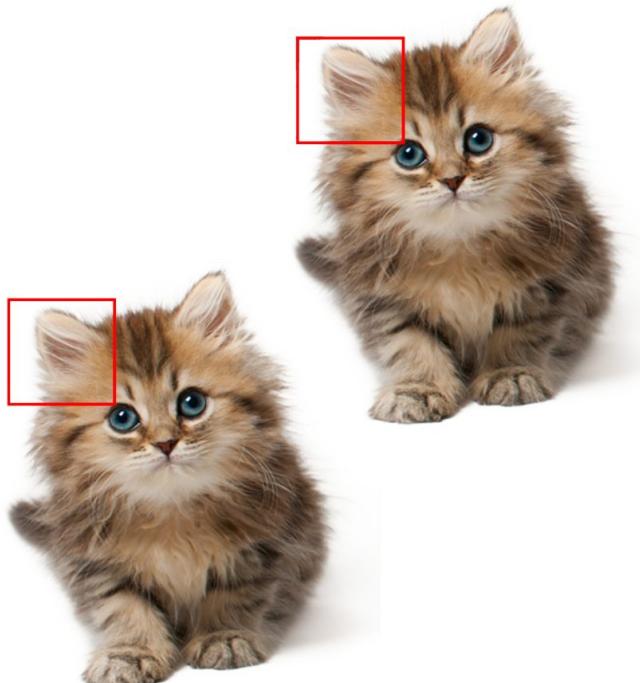


100

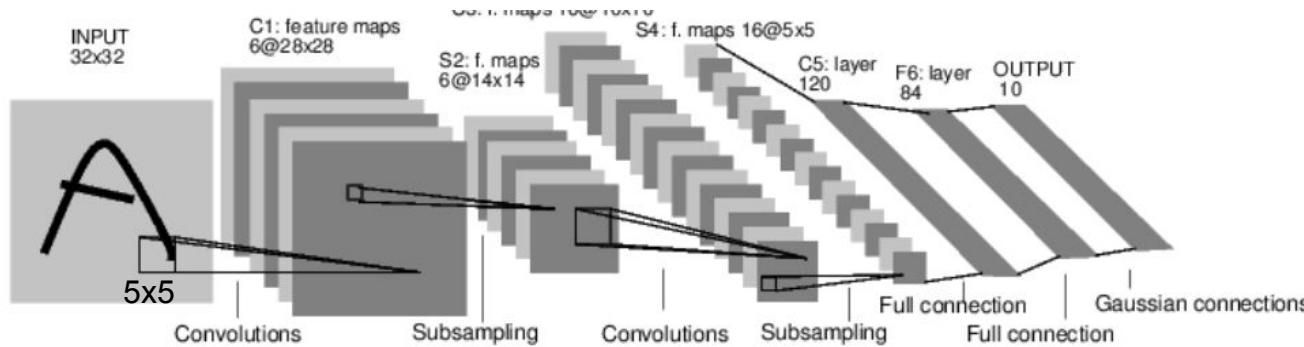
How to improve GAN?

Let's go to basics!

Layers: convolution



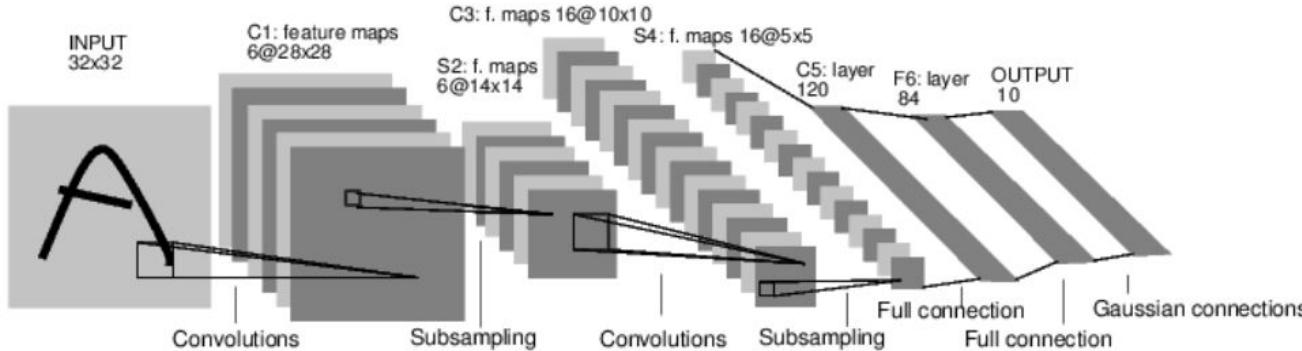
Layers: convolution



First layer:

- Without weight sharing:
- with weight sharing:

Layers: convolution



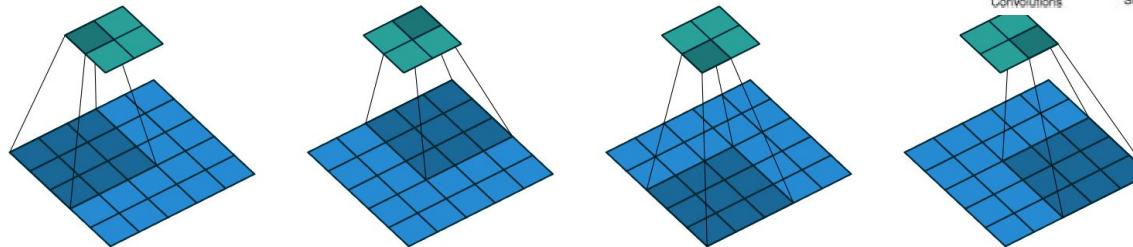
First layer:

- Without weight sharing: $5 \times 5 \times 6 \times 28 \times 28 = 117600$
- with weight sharing: $5 \times 5 \times 6 = 150$

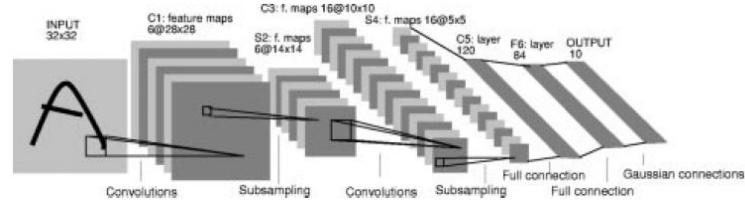
“De-convolution”

Architecture:

Encoder: ConvNet architecture

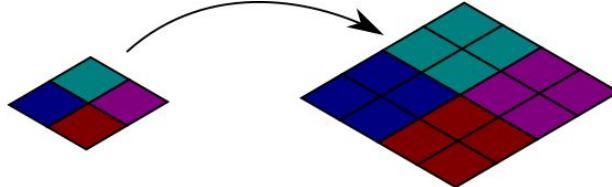


Convolving a 3×3 kernel over a 5×5 input using 2×2 strides



Decoder: How to go from a vector to an image? How to “deconvolve”?

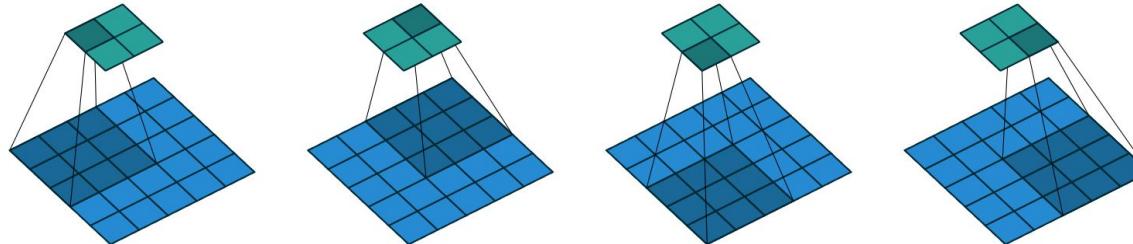
Up-sampling



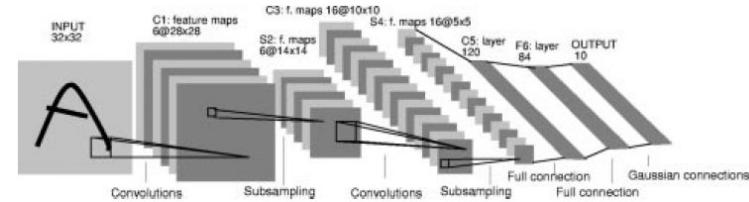
“De-convolution”

Architecture:

Encoder: ConvNet architecture

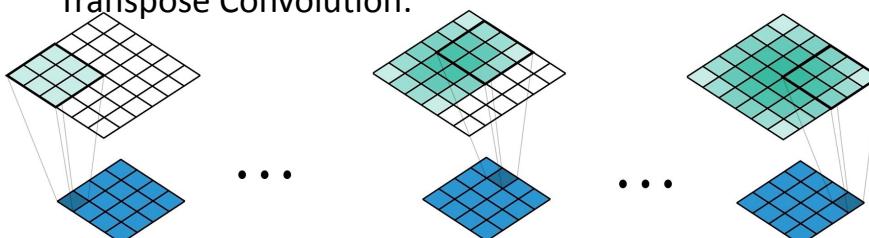


Convolving a 3×3 kernel over a 5×5 input using 2×2 strides



Decoder: How to go from a vector to an image? How to “deconvolve”?

Transpose Convolution:



Gradient descent

- Optimization problem: $\min_{\theta} = L_{tot}(\theta)$
- Gradient Descent:
 - Until convergence:
$$\theta_{n+1} = \theta_n - \gamma \nabla L(\theta_n)$$
- Non convex problem -> Gradient descent does not work!

Gradient descent

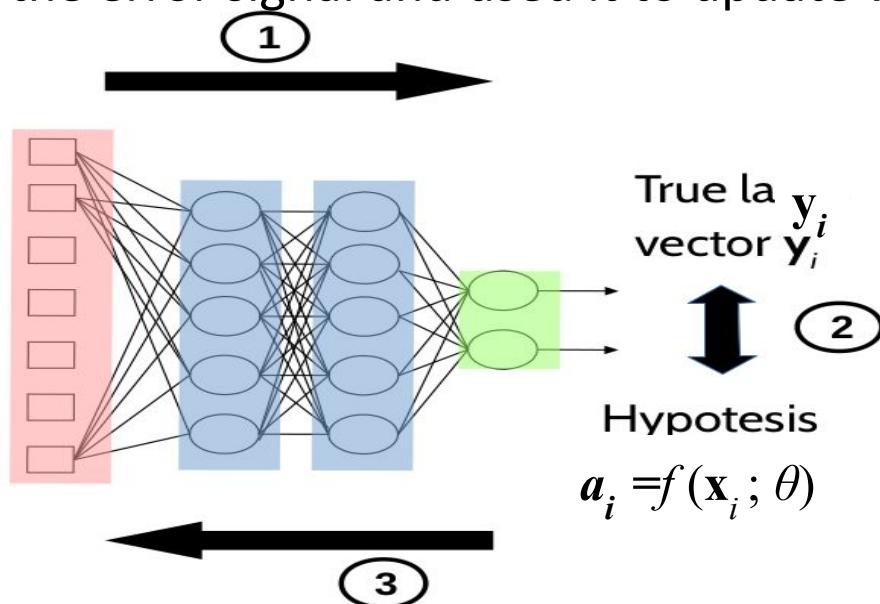
- Optimization problem: $\min_{\theta} = L_{tot}(\theta)$
- Gradient Descent:
 - Until convergence: $\theta_{n+1} = \theta_n - \gamma \nabla L(\theta_n)$
- Non convex problem -> Gradient descent does not work!
- Stochastic Gradient Descent:
 - We split the training in mini-batch: $L_{tot}(\theta) = \sum_b L_b(\theta)$
 - Until convergence:
 - For each batch b: $\theta_{n+1} = \theta_n - \gamma \nabla L_b(\theta_n)$
 - How to compute the gradient efficiently?

1986 – Backpropagation

- Backpropagation revitalized the field
- Learning MLP for complicated functions can be solved
- Efficient algorithm which processes “large” training sets
- Allowed for complicated neural network architectures
- Today backpropagation is still at the core of neural network training

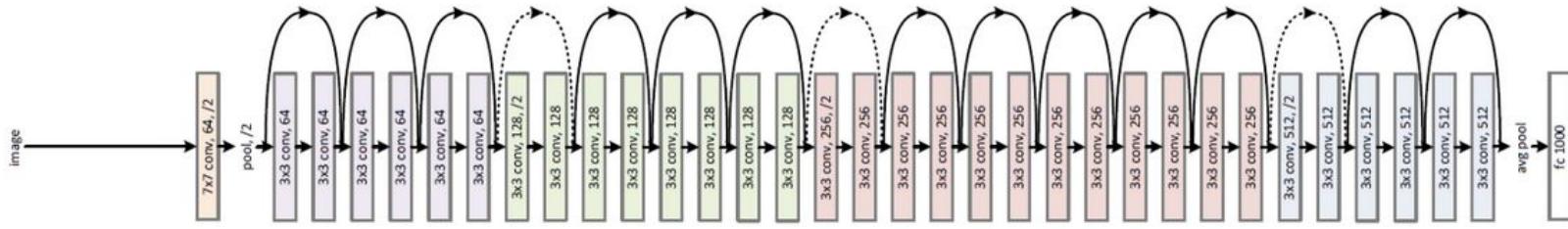
Backpropagation in a Nutshell

1. **Forward propagation:** sum inputs, produce activations, feed-forward
2. **Error estimation.**
3. **Back propagate** the error signal and used it to update weights



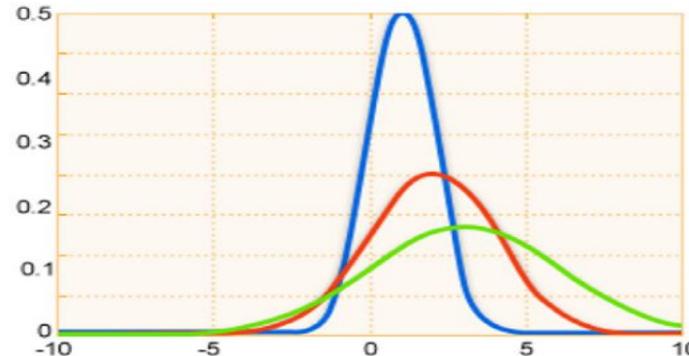
Layers: Batch-normalization

34-layer residual

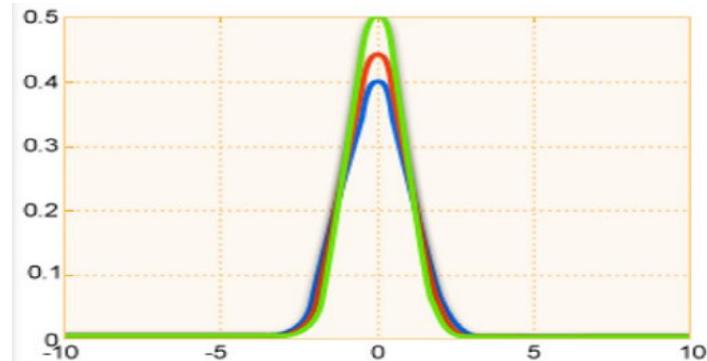
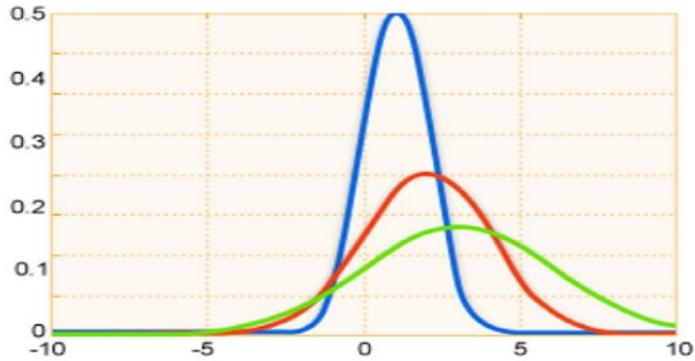


Parameter update:

$$\theta_{n+1} = \theta_n - \gamma \nabla L_b(\theta_n)$$



Layers: Batch-normalization



Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Batch Normalization: Accelerating Deep Network
Training by Reducing Internal Covariate Shift
Sergey Ioffe, Christian Szegedy

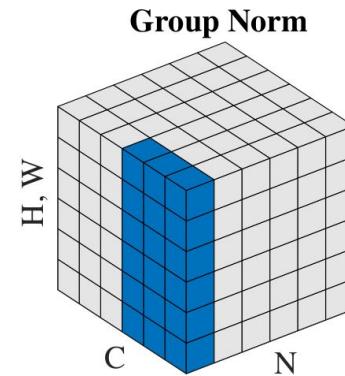
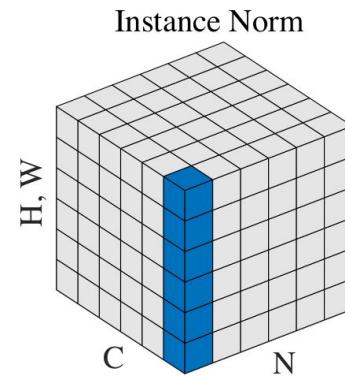
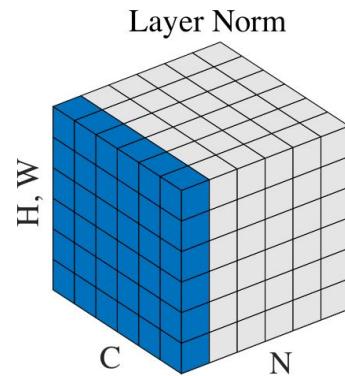
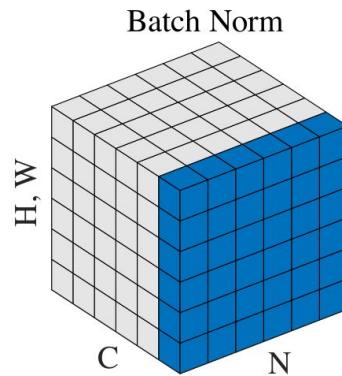
Layers: Batch-normalization

Advantages:

- Faster training
- Better optimization (lower final loss)
- Better generalization
- Specific usages in:
 - few-shot learning
 - domain adaptation
 - generation

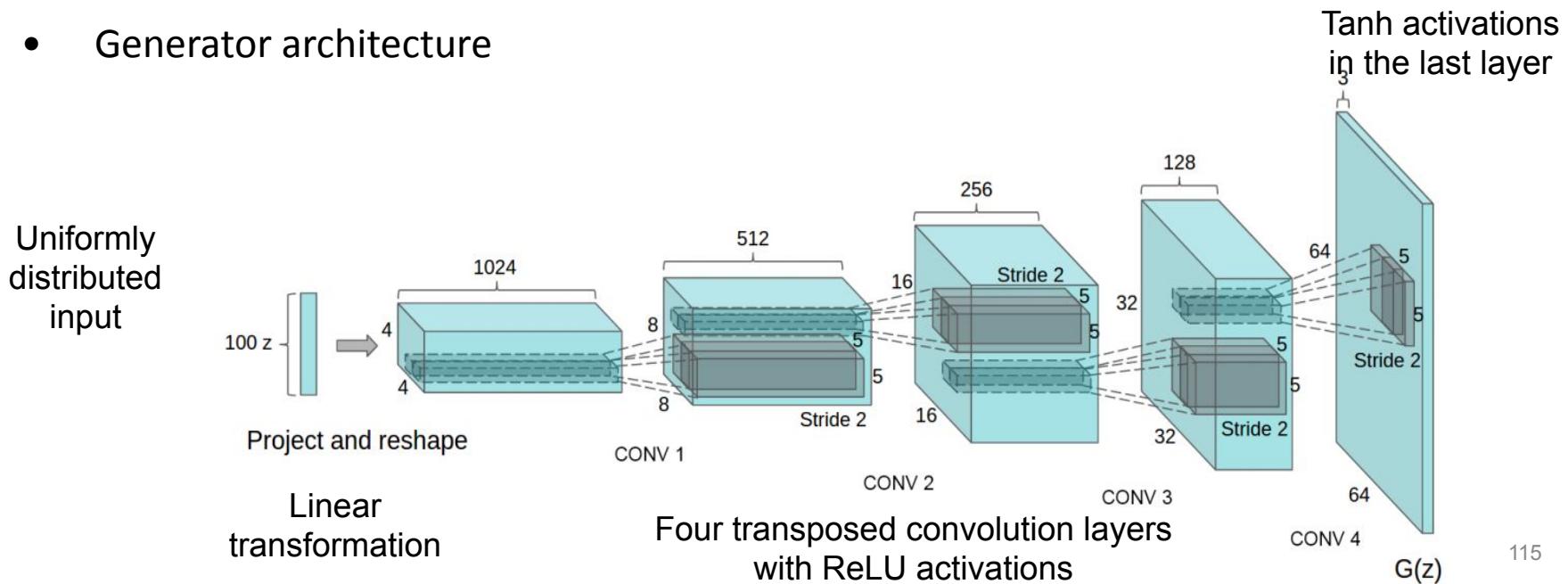
Layers: Other normalizations

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$



DCGAN

- Propose principles for designing convolutional architectures for GANs, first model for large resolution image generation
- Generator architecture



DCGAN

- Discriminator architecture
 - No pooling, only strided convolutions
 - Use Leaky ReLU activations (sparse gradients cause problems for training)
 - Use only one FC layer before the softmax output
 - Use batch normalization after most layers
(in the generator also)

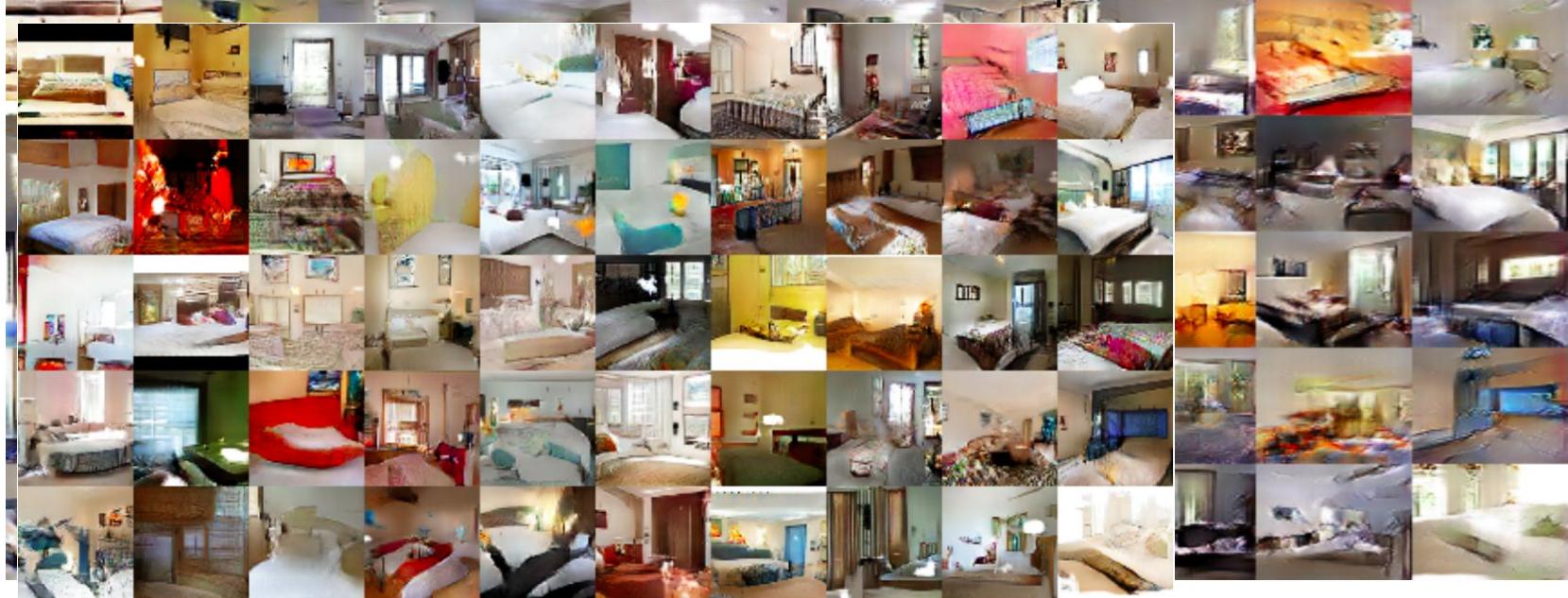
DCGAN results

Generated bedrooms after one epoch



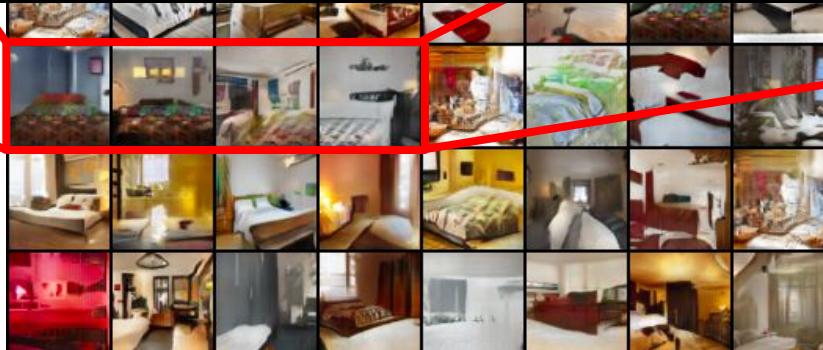
DCGAN results

Generated bedrooms after five epochs

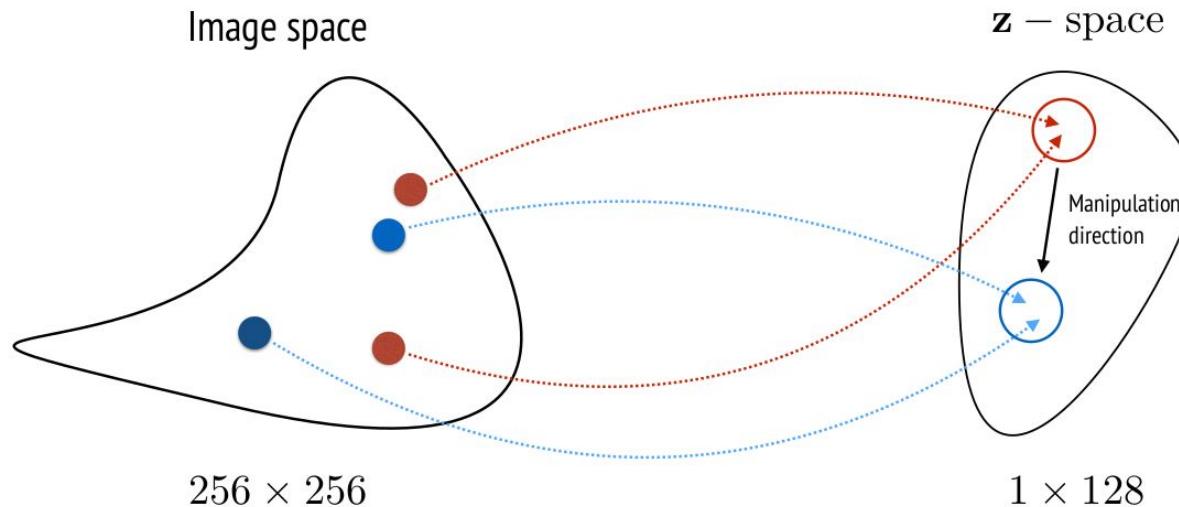


DCGAN results

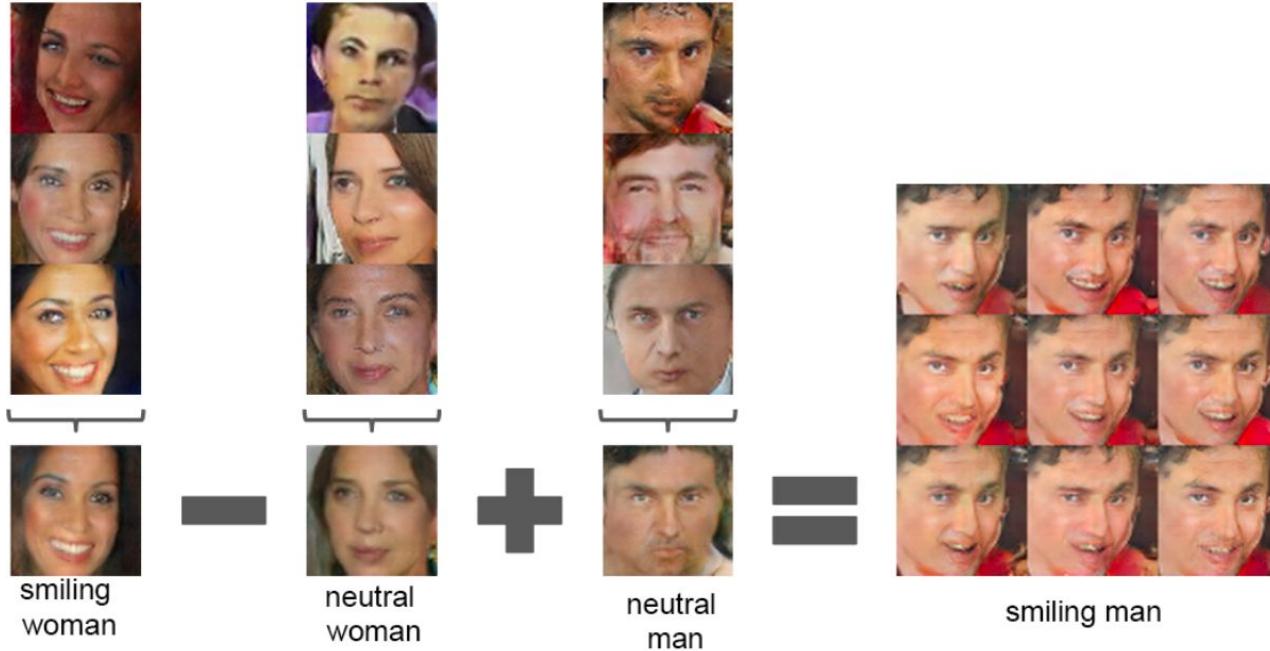
Generated bedrooms from reference implementation



DCGAN results

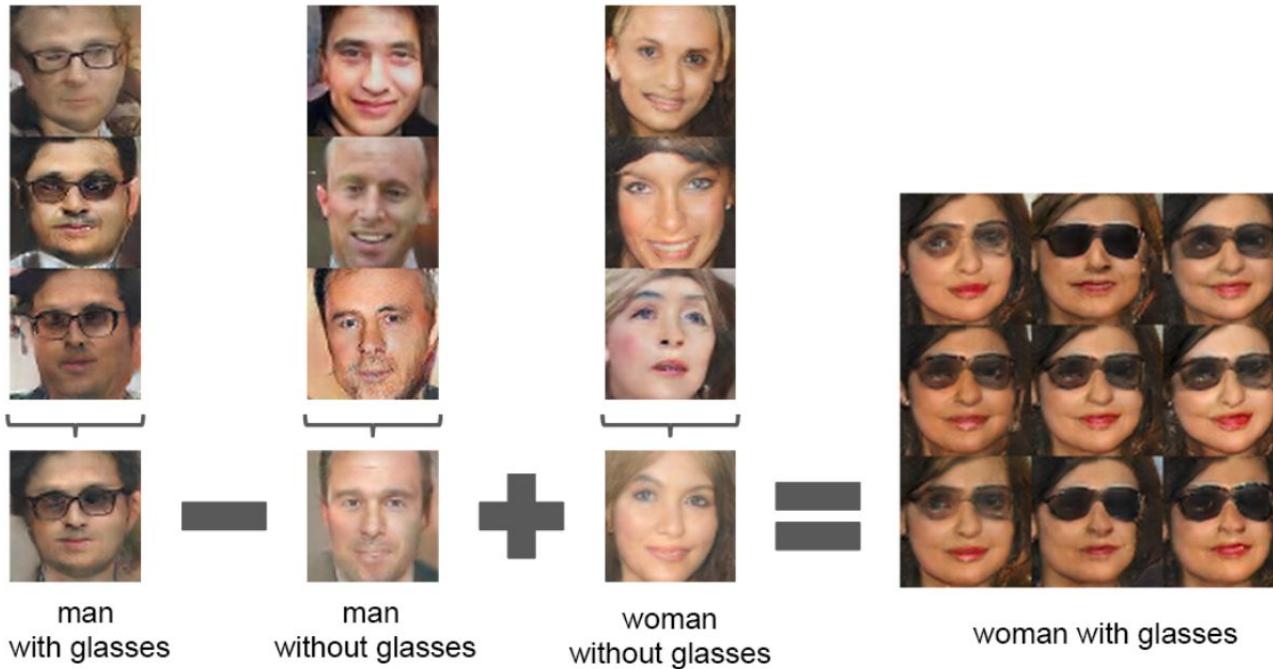


DCGAN results



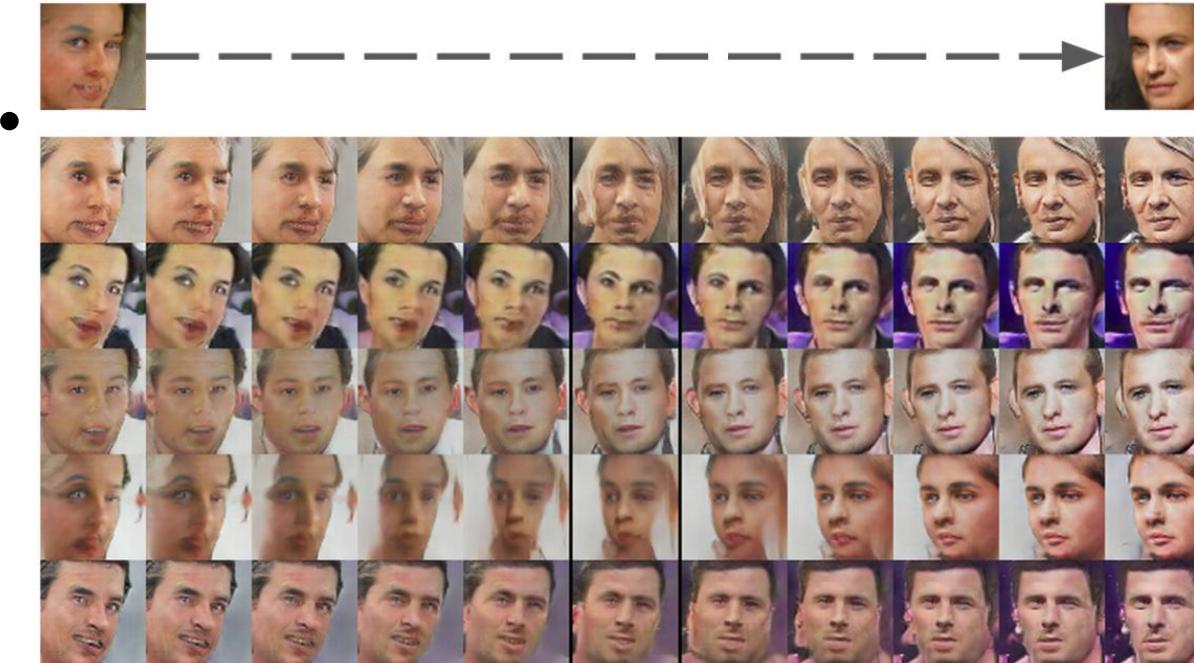
For each column, the Z vectors of samples are averaged. Arithmetic was performed on the mean vectors creating a new vector Y . The center sample on the right hand side is produced by feeding Y as input to the generator. To demonstrate the interpolation capabilities of the generator, uniform noise sampled with scale $+0.25$ was added to Y to produce the 8 other samples

DCGAN results



First model to bring the concept of **disentanglement**.

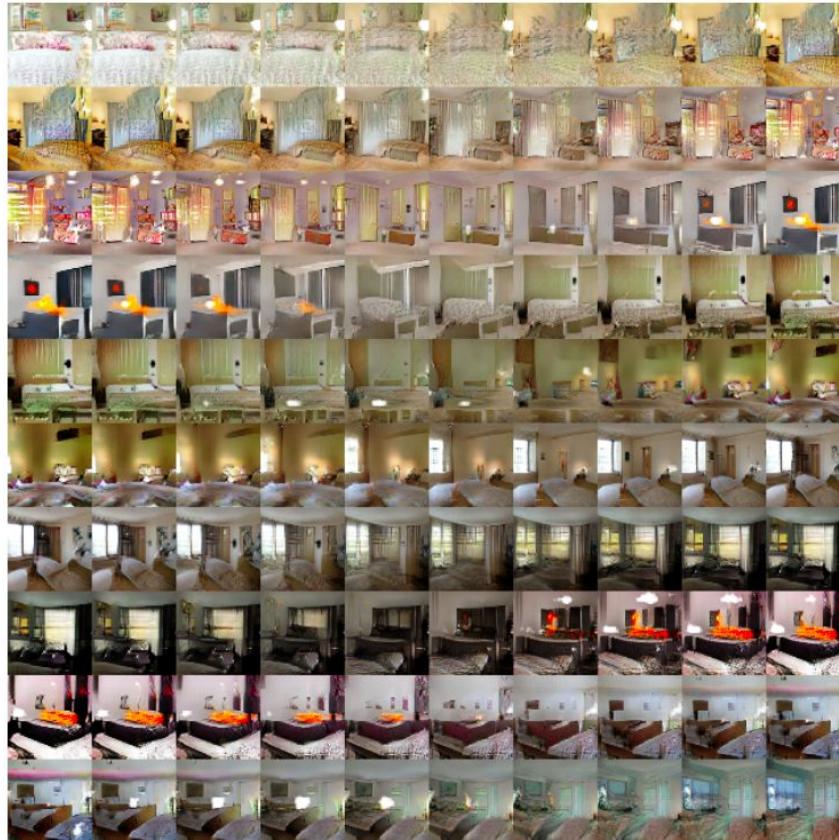
DCGAN results



A "turn" vector was created from four averaged samples of faces looking left vs looking right.
By adding interpolations along this axis to random samples we were able to reliably transform
their pose.

DCGAN results

Interpolation between different points in the z space



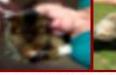
Evaluate GAN performance

- Showing pictures of samples is not enough, especially for simpler datasets like MNIST, CIFAR, faces, bedrooms, etc.
- We cannot directly compute the likelihoods of high-dimensional samples (real or generated) or compare their distributions
- Many GAN approaches claim mainly to improve stability, which is hard to evaluate
- For discussion, see [Ian Goodfellow Twitter thread](#)

Evaluate GAN performance

- Human performance

Instructions

Examples of real images				Examples of images generated by a computer					
									

We marked with green all the images generate by computer (YOU SHOULD SET CHECKBOX FOR ALL OF THEM) and with red all coming from scans of digits (DON'T SET CHECKBOX ON ANY OF THEM).

Please press Next.

<input checked="" type="checkbox"/> mistake	<input type="checkbox"/> mistake	<input checked="" type="checkbox"/> correct
<input checked="" type="checkbox"/> mistake	<input type="checkbox"/> correct	<input type="checkbox"/> correct
<input checked="" type="checkbox"/> mistake	<input checked="" type="checkbox"/> correct	<input type="checkbox"/> correct

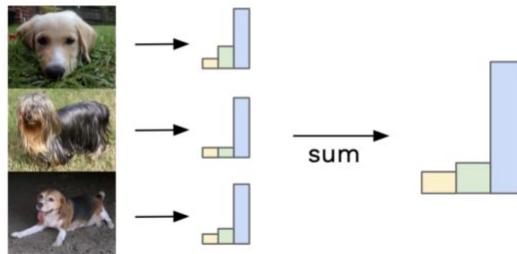
Your score on this question is 5/9

126

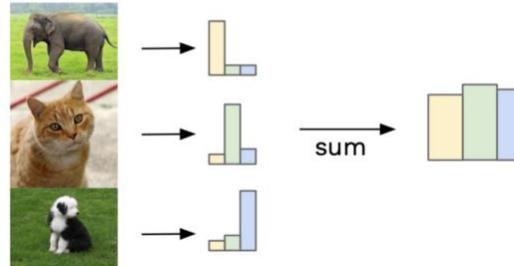
Inception Score

- Key idea: generators should produce images with a **variety of recognizable object classes**
- Two criteria:
 - A human, looking at a separate image, would be able to confidently determine what's in there (**saliency**).
 - A human, looking at a set of various images, would say that the set has lots of different objects (**diversity**).

Similar labels sum to give focussed distribution



Different labels sum to give uniform distribution



127

Inception Score

- Key idea: generators should produce images with a **variety** of **recognizable** object classes
- Defined as:

$$\text{IS}(G) = \exp \left(\mathbb{E}_{\mathbf{x} \sim p_a} D_{KL} (p(y|\mathbf{x}) \parallel p(y)) \right),$$

where $p(y|x)$ is the posterior label distribution returned by an image classifier (e.g., InceptionNet) for sample x

- If the image contains a recognizable object, entropy of $p(y|x)$ should be low
- If generator generates images of diverse objects, the marginal distribution $p(y)$ should have high entropy
- Disadvantage: a GAN that simply memorizes the training data (overfitting) or outputs a single image per class (mode dropping) could still score well

Fréchet Inception Distance

- Key idea: fit simple distributions to statistics of feature activations for real and generated data; estimate divergence parametrically
 - Pass generated samples through a network (InceptionNet), compute activations for a chosen layer
 - Estimate multivariate mean and covariance of activations, compute Fréchet distance to those of real data

$$\text{FID}(x, g) = \|\mu_x - \mu_g\|_2^2 + \text{Tr}(\Sigma_x + \Sigma_g - 2(\Sigma_x \Sigma_g)^{\frac{1}{2}}),$$

- Advantages: correlated with visual quality of samples and human judgment
- Disadvantage: cannot detect overfitting

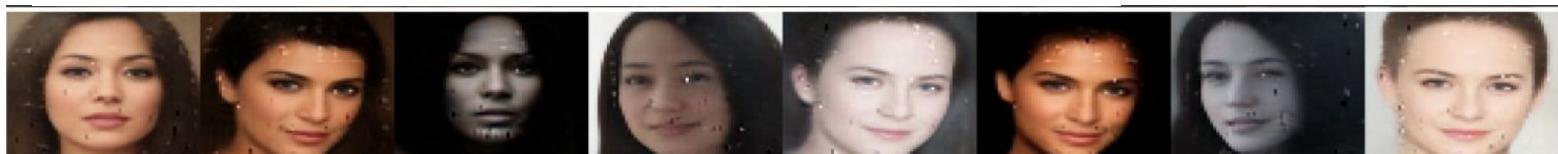
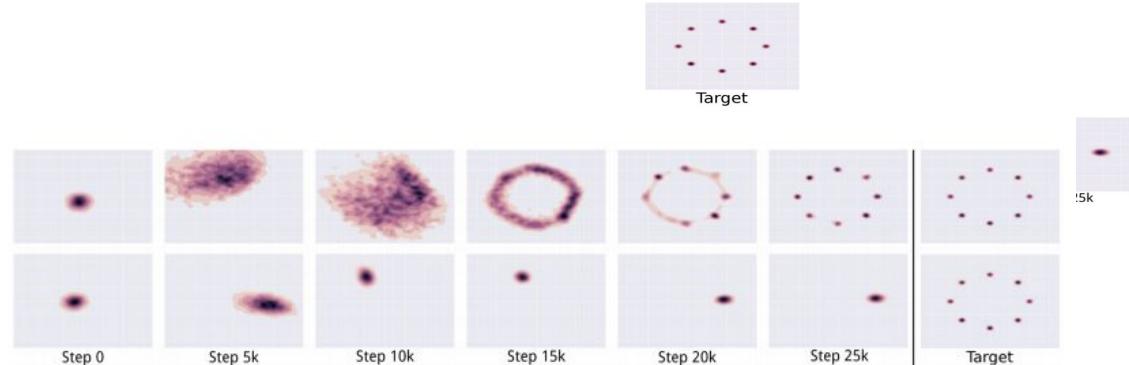
129

Issues with GAN training

- Stability
 - Parameters can oscillate or diverge, generator loss does not correlate with sample quality
 - Behavior very sensitive to hyper-parameter selection
 - Proposed several improvements, mostly changing the loss function

Issues with GAN training

- Mode collapse
 - Generator ends up modeling only a small subset of the training data



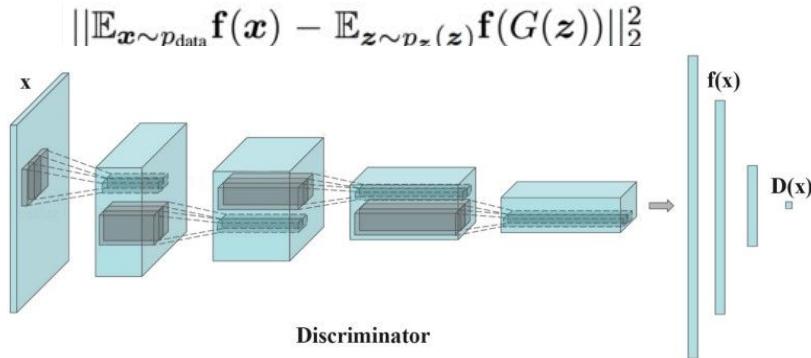
131

Training Tricks

- Feature matching
- Use virtual batch norm
- Mini-batch discrimination (avoid mode collapsing by looking at more samples)
- Use label informations in the discriminator ($C+1$ categories)
- Label smoothing
- Historical averaging (introducing an additional cost function)

Training Tricks

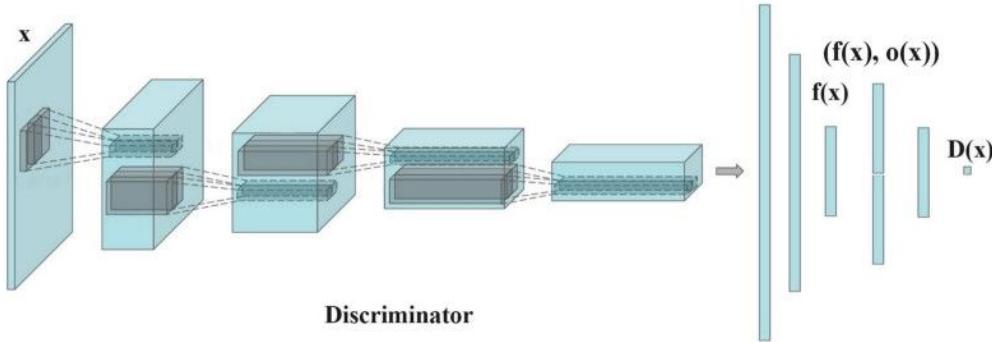
- Feature matching
- Issue:
 - The generator tries to find the best image to fool the discriminator.
 - The “best” image keeps changing when both networks counteract their opponent.
 - The model does not converge and mode collapses.
- Add a loss for feature matching (expands the goal from beating the opponent to matching features in real images)



133

Training Tricks

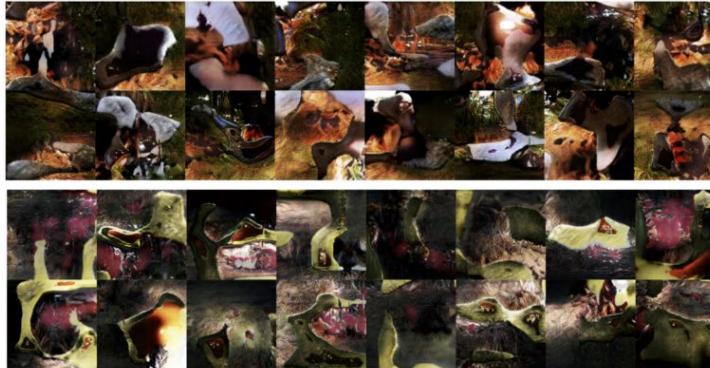
- Mini-batch discrimination
- Compute the similarity of the image x with images in the same batch. Append the similarity $o(x)$ in one of the dense layers in the discriminator to classify whether this image is real or generated.
- If the mode starts to collapse, the similarity of generated images increases. The discriminator can use $o(x)$ to detect generated images and penalize the generator if mode is collapsing.



134

Training Tricks

- Virtual batch norm
- Issue: Due to batch norm the generated images are not independent
- Idea: samples a reference batch before the training to compute the normalization parameters (μ and σ) and combine a reference batch with the current batch to compute the normalization parameters.

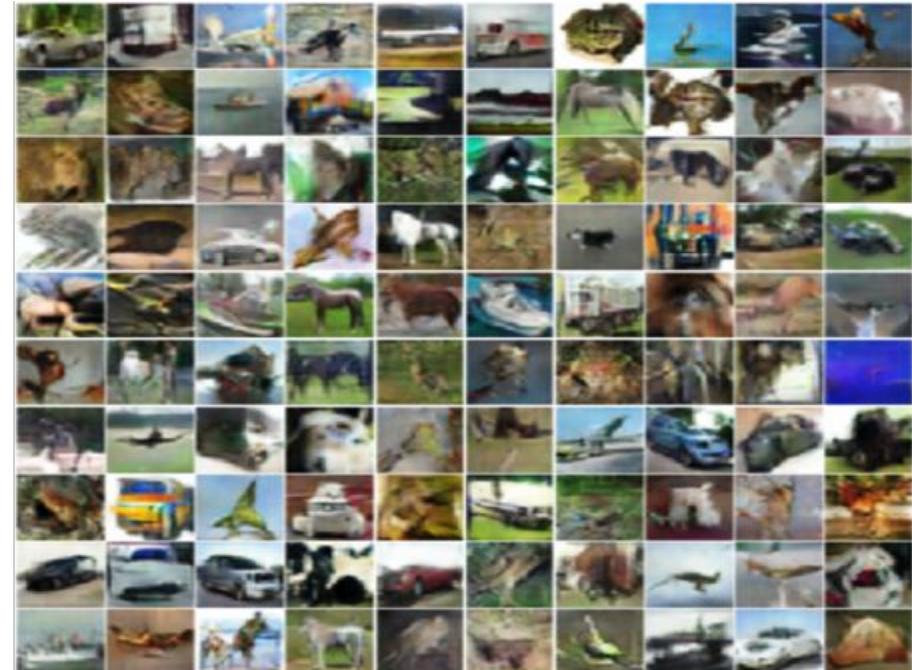


Training Tricks

- Feature matching
- Use virtual batch norm
- Mini-batch discrimination (avoid mode collapsing by looking at more samples)
- Use label informations in the discriminator (C+1 categories)
- Label smoothing
- Historical averaging (introducing an additional cost function)

$$||\boldsymbol{\theta} - \frac{1}{t} \sum_{i=1}^t \boldsymbol{\theta}[i]||^2$$

GAN Results

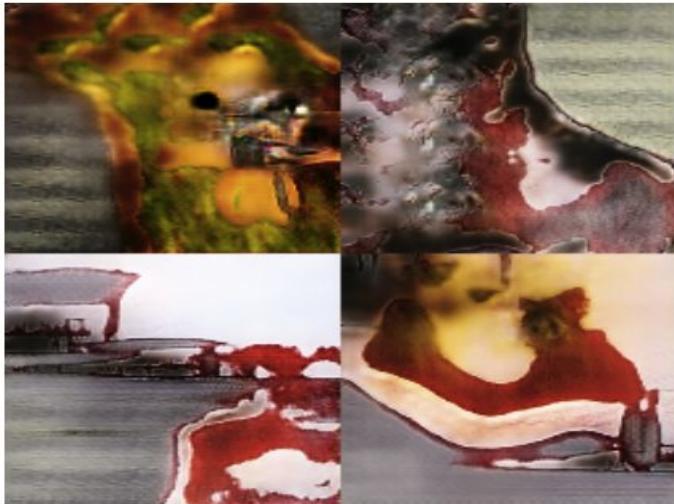


Original CIFAR-10 vs. Generated CIFAR-10 samples

GAN Results

- ImageNet Results (128x128 pixels)

DCGAN

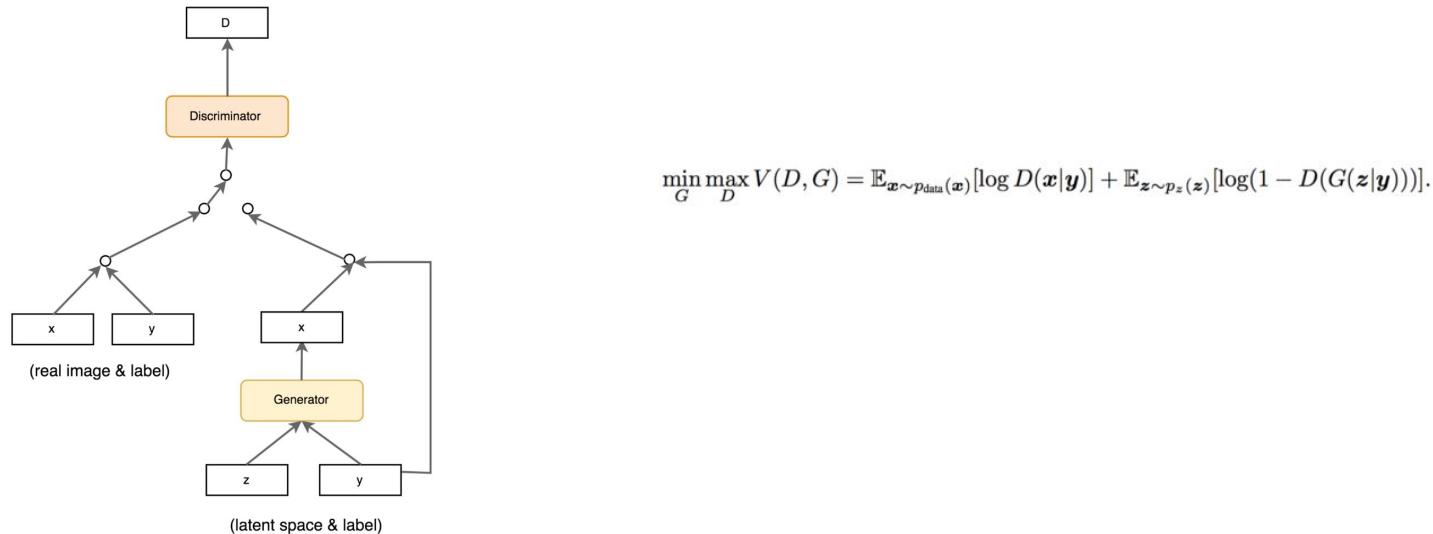


Improved DCGAN



CGAN & InfoGAN

- Conditional GAN (CGAN)
 - Same loss as the original GAN
 - Use label info (1-hot encoding) both in generator and discriminator
 - Label as extension to the latent space



CGAN & InfoGAN

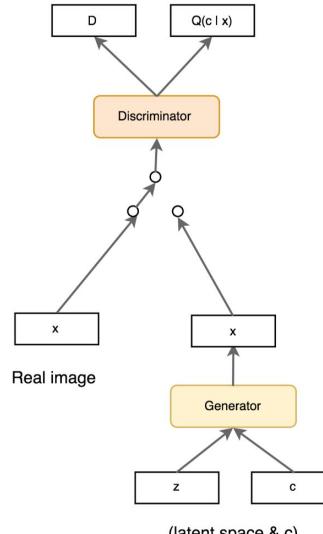
- Info GAN (label not given, treated as latent factor):
- Information Maximizing Generative Adversarial Nets

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$$

- X and Y are independent, then $I(X; Y) = 0$
- X is completely determined by Y, then $H(Y|X)=0$, $I(X; Y) = H(Y)$
Then I is maximized

CGAN & InfoGAN

- Info GAN (label not given, treated as latent factor):
 - Additional output for the discriminator $Q(c|x)$
 - Mutual information $I(c;x)$ measure how much we know c if we know x . It should be high.
 - The info in the latent code c should not be lost in the generation.



$$\min_G \max_D V_{infoGAN}(D, G) = V_{GAN}(D, G) - \lambda I(c; G(z, c))$$

$$V_{GAN}(D, G) \equiv \mathbb{E}_{x \sim p_{data}} \log D(x) + \mathbb{E}_{z \sim p_z(z)} \log(1 - D(G(z, c)))$$

$$\begin{aligned} I(c; G(z, c)) &= H(c) - H(c|G(z, c)) \\ &= \mathbb{E}_{x \sim G(z, c)} [\mathbb{E}_{c \sim P(c|x)} [\log P(c|x)]] + H(c) \\ &\geq \mathbb{E}_{x \sim G(z, c)} [\mathbb{E}_{c \sim P(c|x)} [\log Q(c|x)]] + H(c) \end{aligned}$$

InfoGAN Results



(a) Varying c_1 on InfoGAN (Digit type)



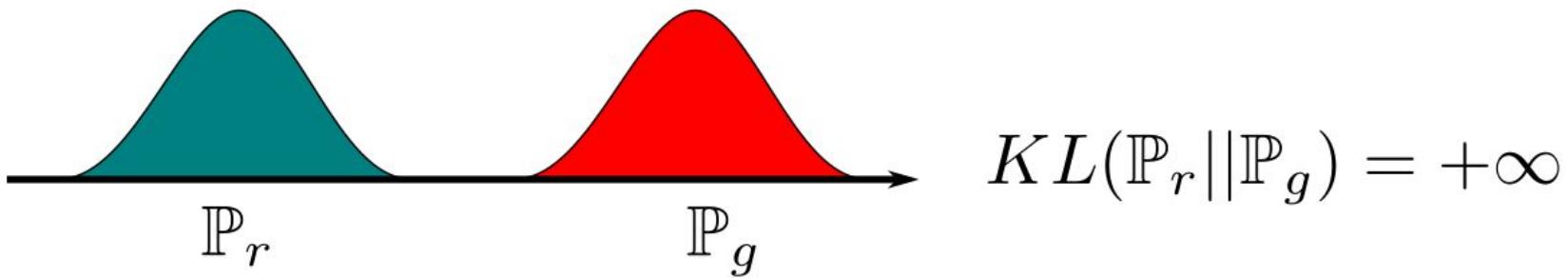
(b) Varying c_1 on regular GAN (No clear meaning)

Wasserstein GAN: towards better losses

The *Kullback-Leibler* (KL) divergence

$$KL(\mathbb{P}_r \parallel \mathbb{P}_g) = \int \log \left(\frac{P_r(x)}{P_g(x)} \right) P_r(x) d\mu(x)$$

Problem when the distributions have $P_g(x)=0$ and $P_r(x)>0$:



Wasserstein GAN: towards better losses

The *Jensen-Shannon* (JS) divergence

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \parallel \mathbb{P}_m) + KL(\mathbb{P}_g \parallel \mathbb{P}_m)$$

where \mathbb{P}_m is the mixture $(\mathbb{P}_r + \mathbb{P}_g)/2$.

Minimizing the GAN objective function with an optimal discriminator is equivalent to minimizing the JS-divergence

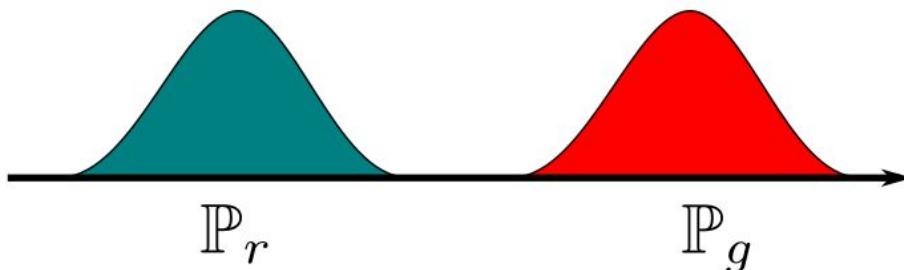
Wasserstein GAN: towards better losses

The *Jensen-Shannon* (JS) divergence

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \parallel \mathbb{P}_m) + KL(\mathbb{P}_g \parallel \mathbb{P}_m)$$

where \mathbb{P}_m is the mixture $(\mathbb{P}_r + \mathbb{P}_g)/2$.

Minimizing the GAN objective function with an optimal discriminator is equivalent to minimizing the JS-divergence



$$JS(\mathbb{P}_r, \mathbb{P}_g) = 2\log(2)$$

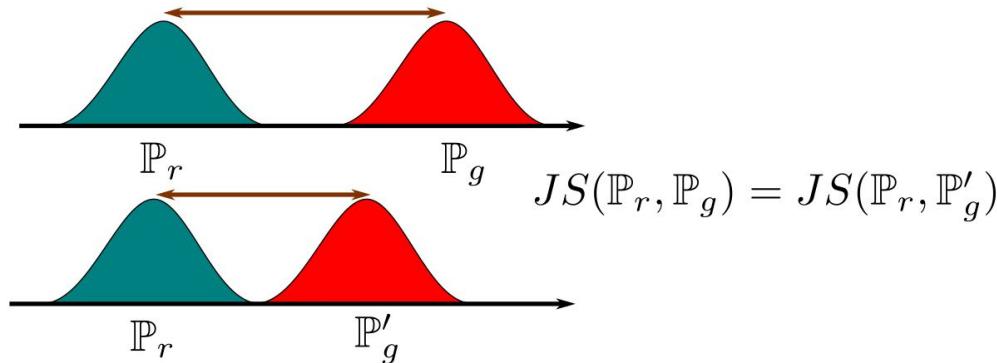
Wasserstein GAN: towards better losses

The *Jensen-Shannon* (JS) divergence

$$JS(\mathbb{P}_r, \mathbb{P}_g) = KL(\mathbb{P}_r \| \mathbb{P}_m) + KL(\mathbb{P}_g \| \mathbb{P}_m)$$

where \mathbb{P}_m is the mixture $(\mathbb{P}_r + \mathbb{P}_g)/2$.

Problem:



The *Earth-Mover* (EM) distance or Wasserstein-1

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

Wasserstein GAN: towards better losses

The *Earth-Mover* (EM) distance or Wasserstein-1

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

We can show:

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

where the supremum is over all the 1-Lipschitz functions $f : \mathcal{X} \rightarrow \mathbb{R}$.

Wasserstein GAN: towards better losses

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

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

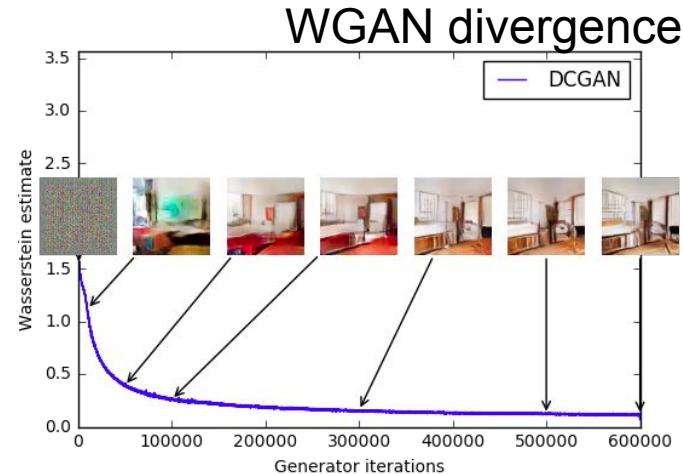
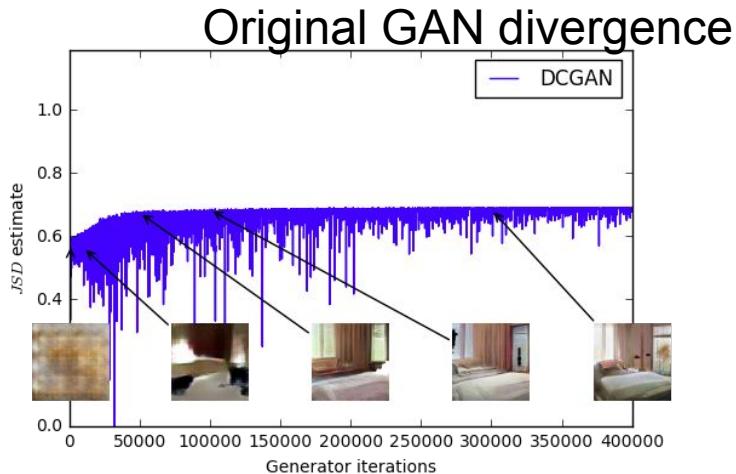
Require: : α , the learning rate. c , the clipping parameter. m , the batch size.
 n_{critic} , the number of iterations of the critic per generator iteration.

Require: : w_0 , initial critic parameters. θ_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:      $g_w \leftarrow \nabla_w [\frac{1}{m} \sum_{i=1}^m f_w(x^{(i)}) - \frac{1}{m} \sum_{i=1}^m f_w(g_\theta(z^{(i)}))]$ 
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

- Objective function value is more meaningfully related to quality of generator output



Wasserstein GAN: Improved Training (WGAN_GP)

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

Weight clipping: $w \leftarrow \text{clip}(w, -c, c)$

Gradient penalty (WGAN-GP):

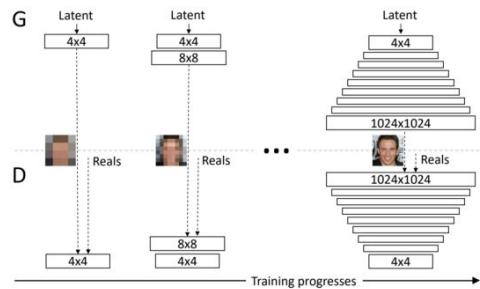
$$L = \underbrace{\mathbb{E}_{\tilde{x} \sim \mathbb{P}_g} [D(\tilde{x})] - \mathbb{E}_{x \sim \mathbb{P}_r} [D(x)]}_{\text{Original critic loss}} + \underbrace{\lambda \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2]}_{\text{Our gradient penalty}}.$$

Lost functions

Name	Value Function
GAN	$L_D^{GAN} = E[\log(D(x))] + E[\log(1 - D(G(z)))]$ $L_G^{GAN} = E[\log(D(G(z)))]$
LSGAN	$L_D^{LSGAN} = E[(D(x) - 1)^2] + E[D(G(z))^2]$ $L_G^{LSGAN} = E[(D(G(z)) - 1)^2]$
WGAN	$L_D^{WGAN} = E[D(x)] - E[D(G(z))]$ $L_G^{WGAN} = E[D(G(z))]$ $W_D \leftarrow clip_by_value(W_D, -0.01, 0.01)$
WGAN_GP	$L_D^{WGAN_GP} = L_D^{WGAN} + \lambda E[\nabla D(\alpha x - (1 - \alpha G(z))) - 1]^2$ $L_G^{WGAN_GP} = L_G^{WGAN}$
DRAGAN	$L_D^{DRAGAN} = L_D^{GAN} + \lambda E[\nabla D(\alpha x - (1 - \alpha x_p)) - 1]^2$ $L_G^{DRAGAN} = L_G^{GAN}$
CGAN	$L_D^{CGAN} = E[\log(D(x, c))] + E[\log(1 - D(G(z), c))]$ $L_G^{CGAN} = E[\log(D(G(z), c))]$
infoGAN	$L_{D,Q}^{InfoGAN} = L_D^{GAN} - \lambda L_I(c, c')$ $L_G^{InfoGAN} = L_G^{GAN} - \lambda L_I(c, c')$
ACGAN	$L_{D,Q}^{ACGAN} = L_D^{GAN} + E[P(class = c x)] + E[P(class = c G(z))]$ $L_G^{ACGAN} = L_G^{GAN} + E[P(class = c G(z))]$
EBGAN	$L_D^{EBGAN} = D_{AE}(x) + \max(0, m - D_{AE}(G(z)))$ $L_G^{EBGAN} = D_{AE}(G(z)) + \lambda \cdot PT$
BEGAN	$L_D^{BEGAN} = D_{AE}(x) - k_t D_{AE}(G(z))$ $L_G^{BEGAN} = D_{AE}(G(z))$ $k_{t+1} = k_t + \lambda(\gamma D_{AE}(x) - D_{AE}(G(z)))$

High Resolution

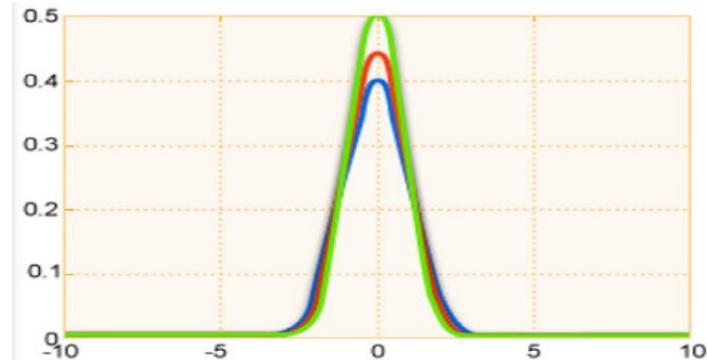
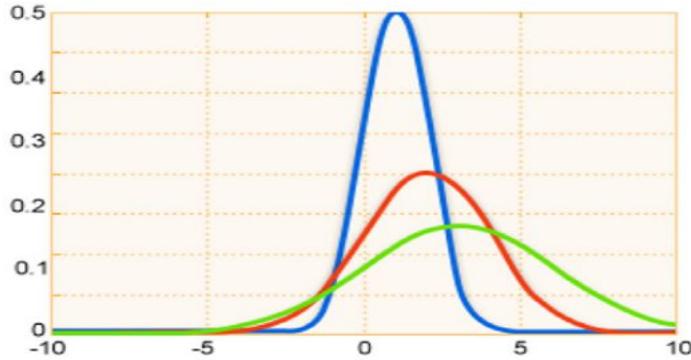
Progressive GAN (2018)



BigGAN (2018)



Inference-normalization



Input: Values of x over a mini-batch: $\mathcal{B} = \{x_1 \dots m\}$;
Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

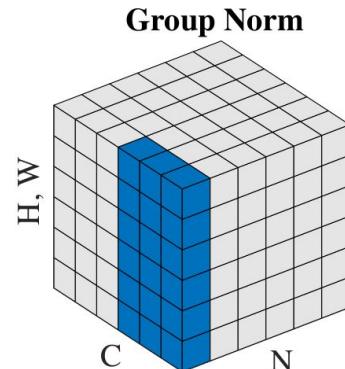
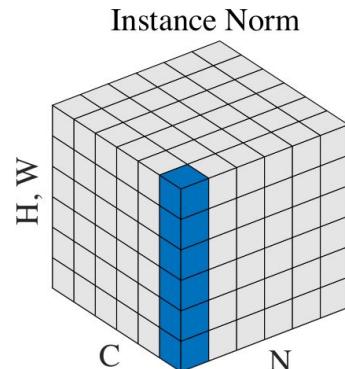
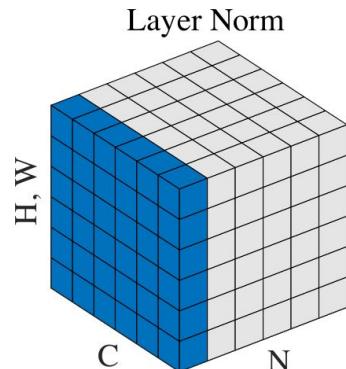
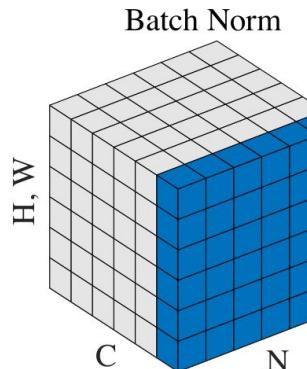
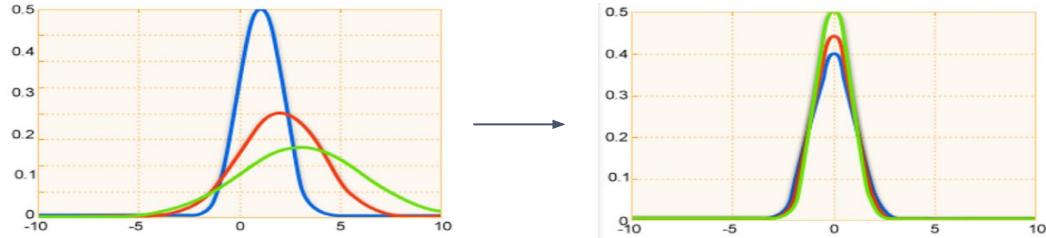
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Batch Normalization: Accelerating Deep Network
Training by Reducing Internal Covariate Shift
Sergey Ioffe, Christian Szegedy

Inference-normalization

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

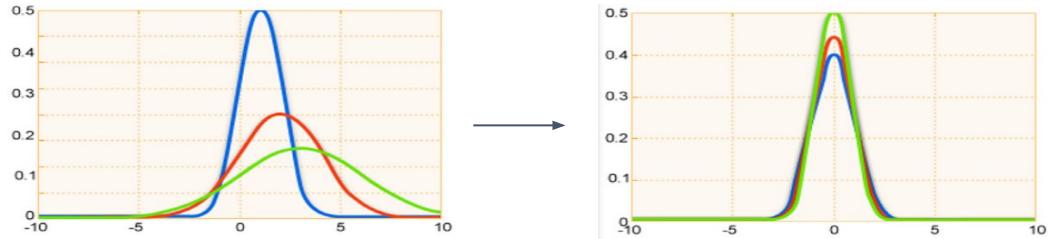


Group Normalization, Yuxin Wu, Kaiming He, ECCV 2018

Why inference-normalization?

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$



(a) Content image.

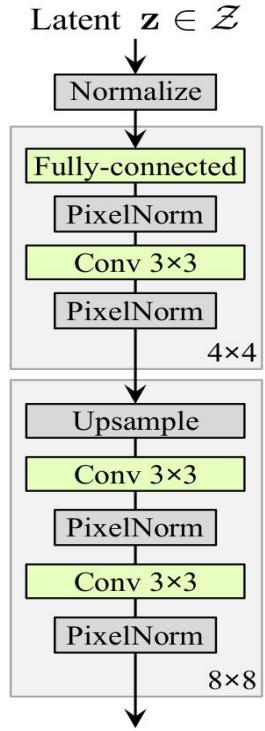


(c) Low contrast content image.

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$
$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv$$

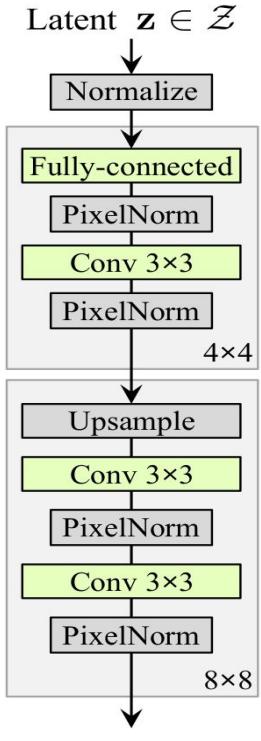
The two images are identical after instance normalization.
We obtain style invariant representations.

Style-Based GAN

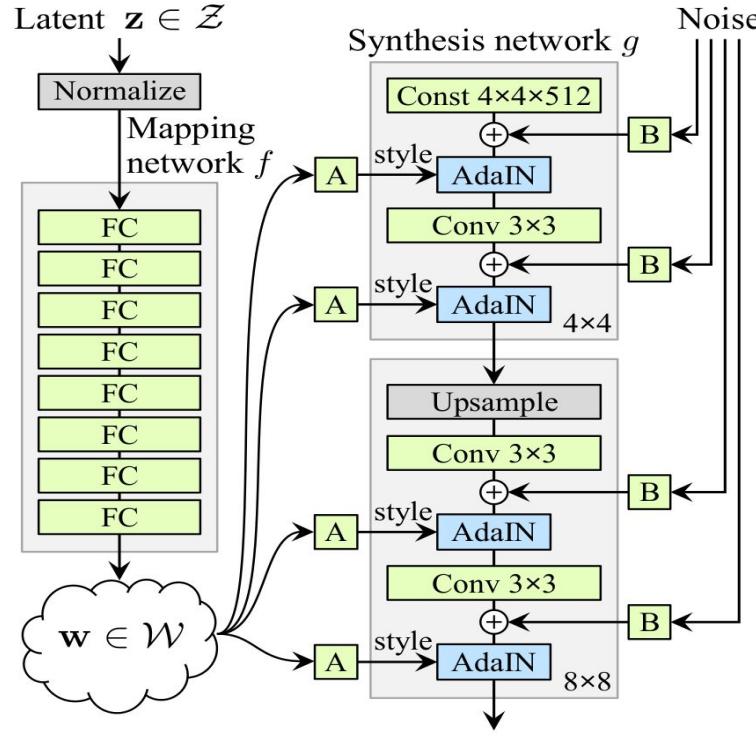


(a) Traditional

Style-Based GAN

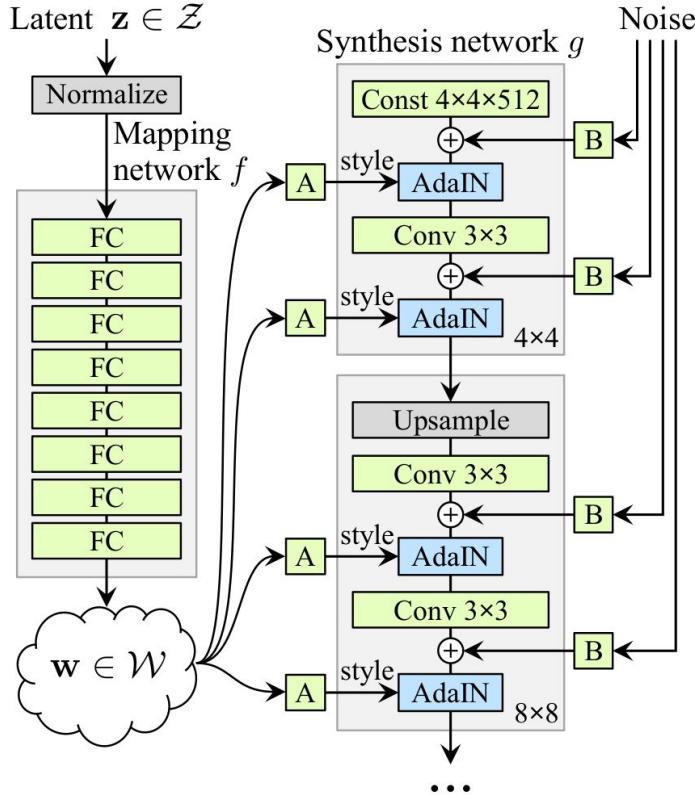


(a) Traditional



(b) Style-based generator

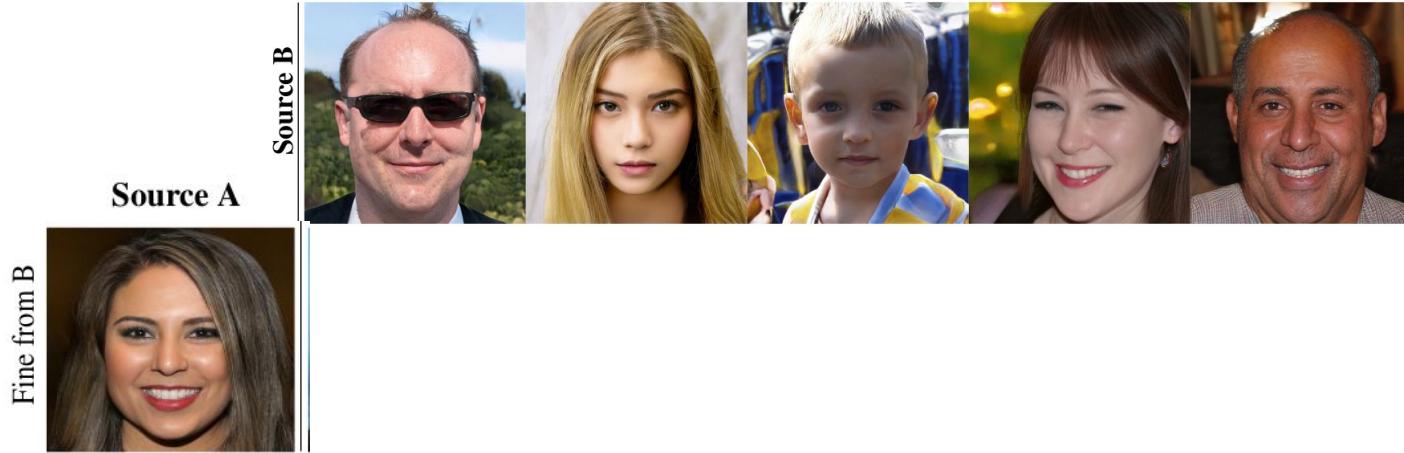
Style-Based GAN



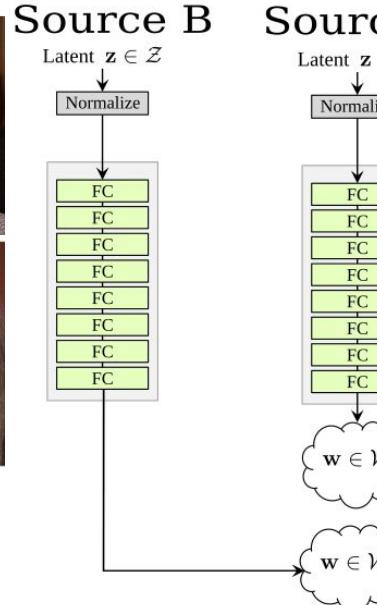
(b) Style-based generator

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

Style-Based GAN



Style-Based GAN



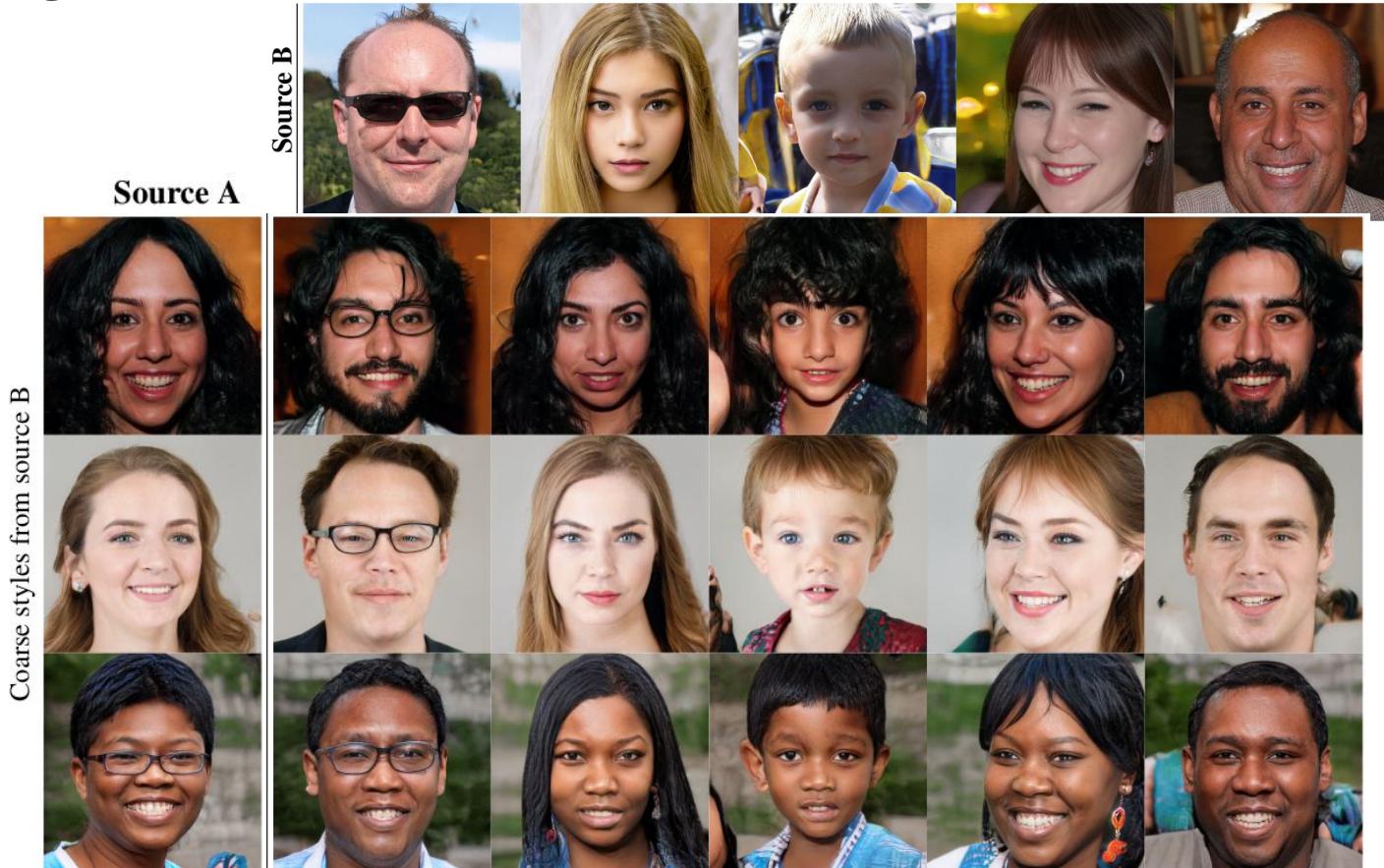
Style-Based GAN



Style-Based GAN



Style-Based GAN



Training GANs: Recap



Training GANs: Recap

Minimax objective function:

$$\min_{\theta_g} \max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

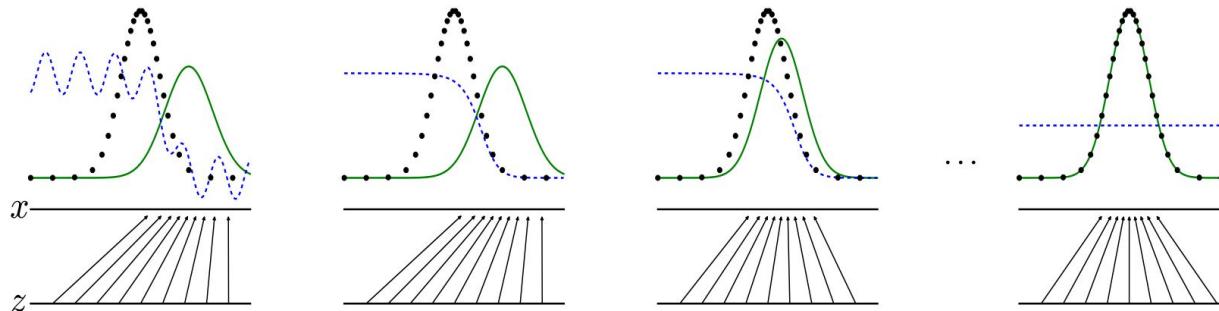
Alternate between:

1. Gradient ascent on discriminator

$$\max_{\theta_d} \left[\mathbb{E}_{x \sim p_{data}} \log D_{\theta_d}(x) + \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z))) \right]$$

2. Gradient descent on generator

$$\min_{\theta_g} \mathbb{E}_{z \sim p(z)} \log(1 - D_{\theta_d}(G_{\theta_g}(z)))$$

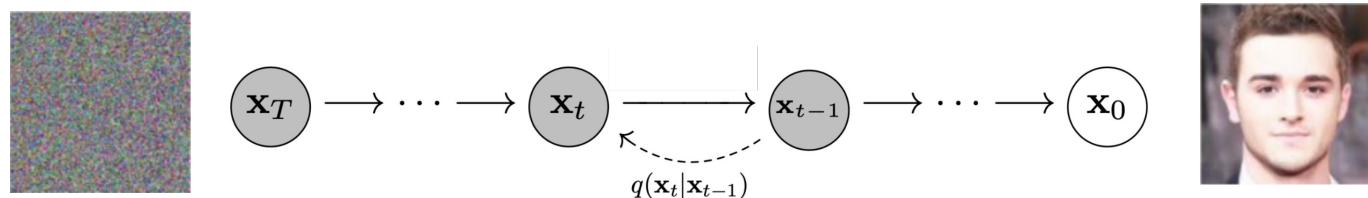


Improving GANs: Recap

- Use convolutions
- Use BN
- Many optimization tricks
- Change the loss (e.g. Wasserstein loss)
- Improve the architecture using inference normalization

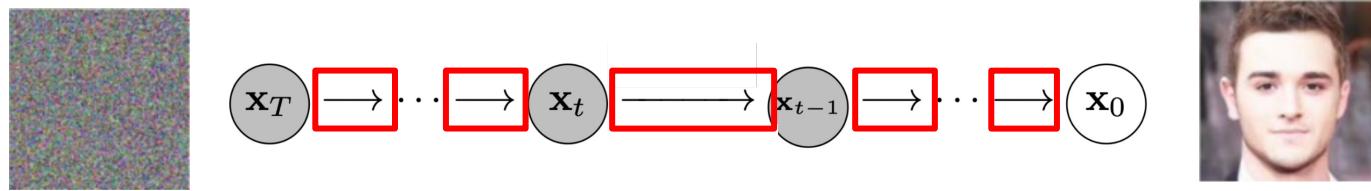
Diffusion Models

Formulate generation as a denoising problem:



We add Gaussian Noise: $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$

Diffusion Models



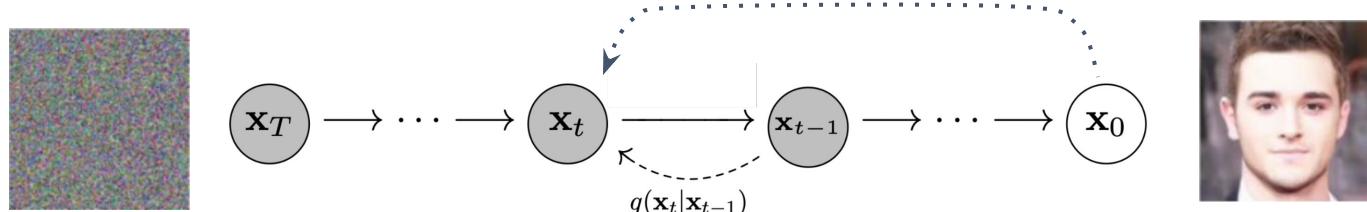
The goal: learn a neural network to approximate the probability distributions in the reverse diffusion process:

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

Diffusion Models

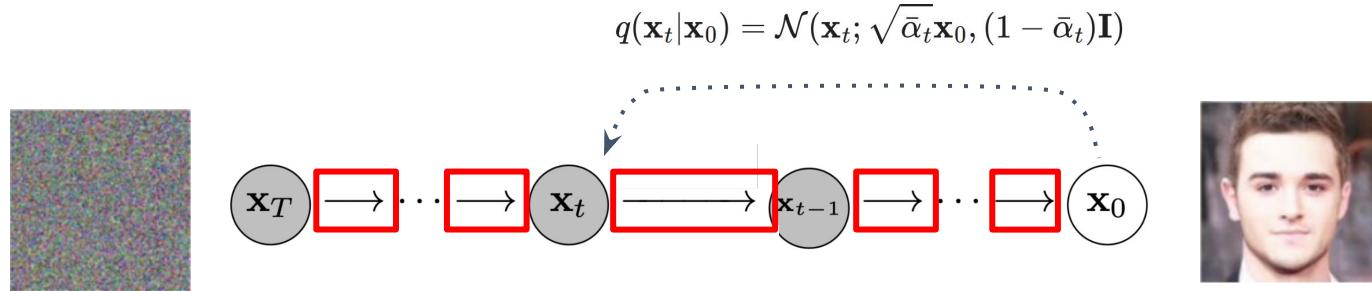
Formulate generation as a denoising problem:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}) \text{ We can combine the multiple steps.}$$



We add Gaussian Noise: $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t} \mathbf{x}_{t-1}, \beta_t \mathbf{I})$

Diffusion Models



Re-parametrization: We predict the noise instead of the image

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

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right)$$

Noise estimate

Diffusion Models

Training loss ?

- We can use the variational lower bound of the negative log-likelihood:

$$L = \mathbb{E}_q[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0}]$$

Diffusion Models

Training loss ?

- We can use the variational lower bound to optimize the negative

$$L = \mathbb{E}_q[\underbrace{D_{\text{KL}}(q(\mathbf{x}_T | \mathbf{x}_0) \| p_\theta(\mathbf{x}_T))}_{L_T} + \sum_{t=2}^T \underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t))}_{L_{t-1}} - \underbrace{\log p_\theta(\mathbf{x}_0 | \mathbf{x}_1)}_{L_0}]$$

where:

$$L_t = \mathbb{E}_{\mathbf{x}_0, \epsilon_t} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \bar{\alpha}_t)\|\Sigma_\theta\|_2^2} \|\epsilon_t - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t, t)\|^2 \right]$$

Diffusion Models

Training loss ?

- We can simplify the loss.

$$L_t = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}_t} \left[\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \bar{\alpha}_t)\|\boldsymbol{\Sigma}_\theta\|_2^2} \|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\sqrt{\alpha_t}\mathbf{x}_0 + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_t, t)\|^2 \right]$$

It becomes:

$$L_t^{\text{simple}} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \boldsymbol{\epsilon}_t} \left[\|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \right]$$

Diffusion Models

Training and sampling algorithms ?

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on

$$\nabla_{\theta} \|\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$$
 - 6: **until** converged
-

Sampling 

Diffusion Models

Training and sampling algorithms ?

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on

$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$
 - 6: **until** converged
-



Noise estimation error minimization

Diffusion Models

Training and sampling algorithms ?

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Diffusion Models

Training and sampling algorithms ?

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

We sample according to:

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

Diffusion Models

Training and sampling algorithms ?

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

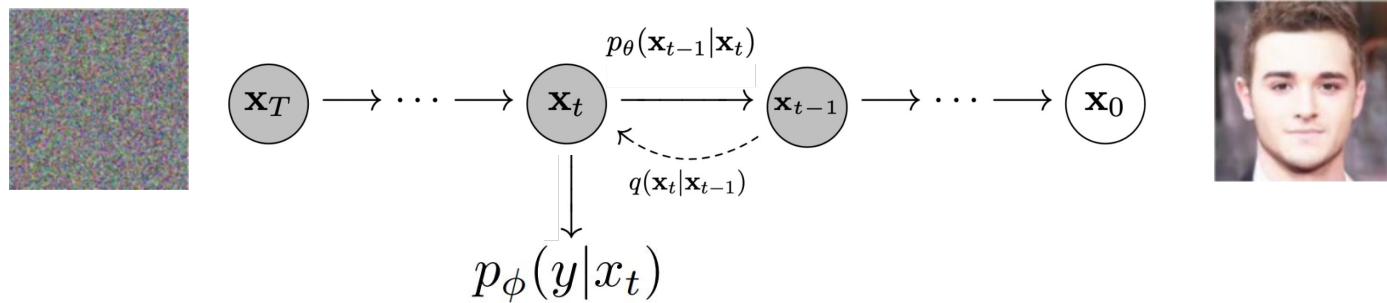


We sample the initial noise.

Classifier guidance

How to condition the output image on y | label ?

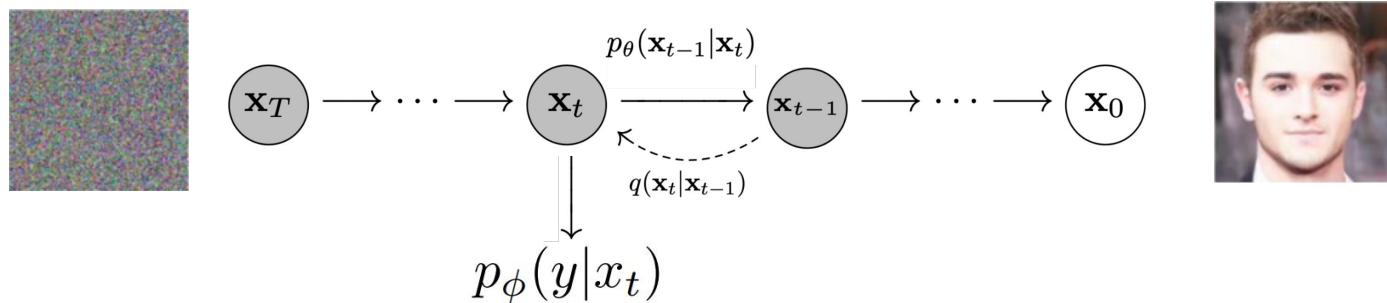
Let us assume a classifier $p_\phi(y|x_t)$



Classifier guidance

How to condition the output image on y | label ?

Let us assume a classifier $p_\phi(y|x_t)$

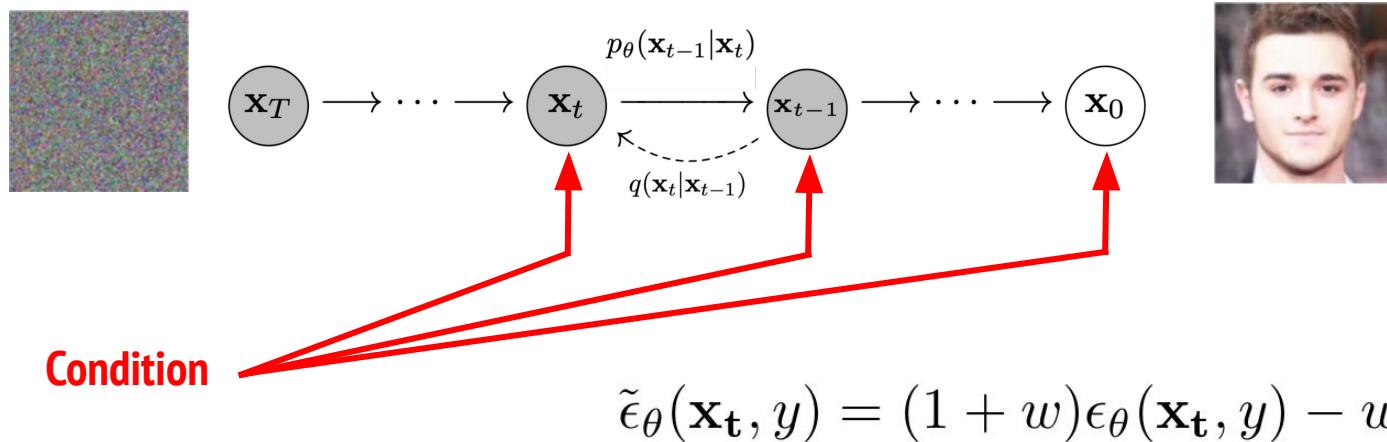


Modified update:

$$\hat{\epsilon}(x_t) := \epsilon_\theta(x_t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{x_t} \log p_\phi(y|x_t)$$

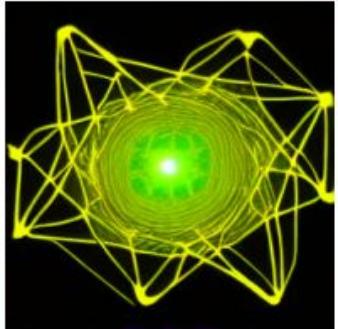
Classifier-free guidance

How to condition the output image on y | label ?



Latent Diffusion Models

'An illustration of a slightly conscious neural network'



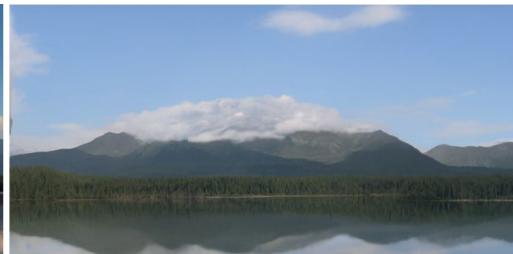
'A painting of a squirrel eating a burger'



'A watercolor painting of a chair that looks like an octopus'

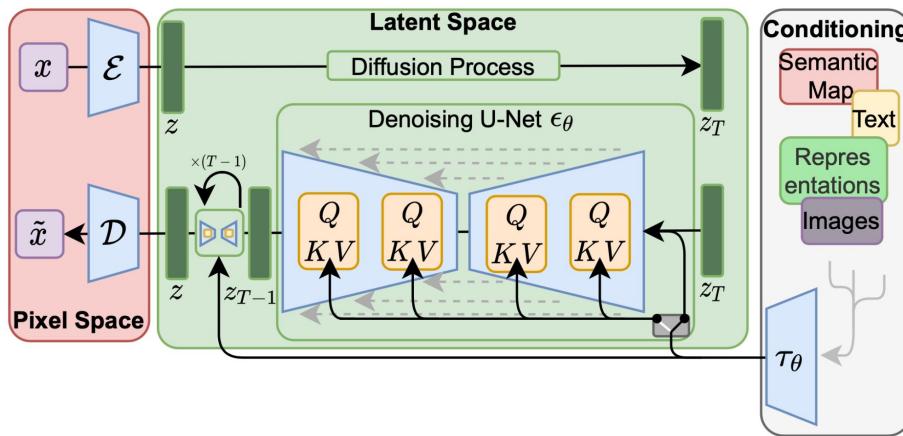


'A shirt with the inscription: "I love generative models!"'



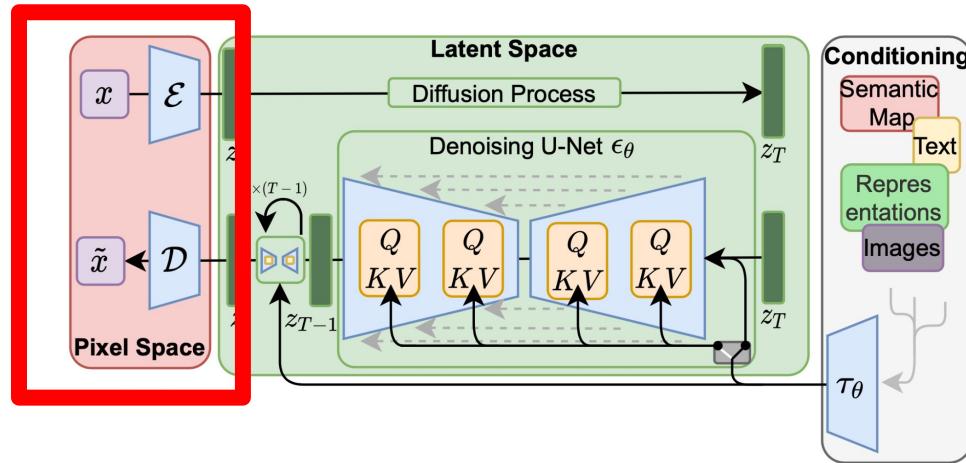
Latent Diffusion Models

How to reduce computational complexity?



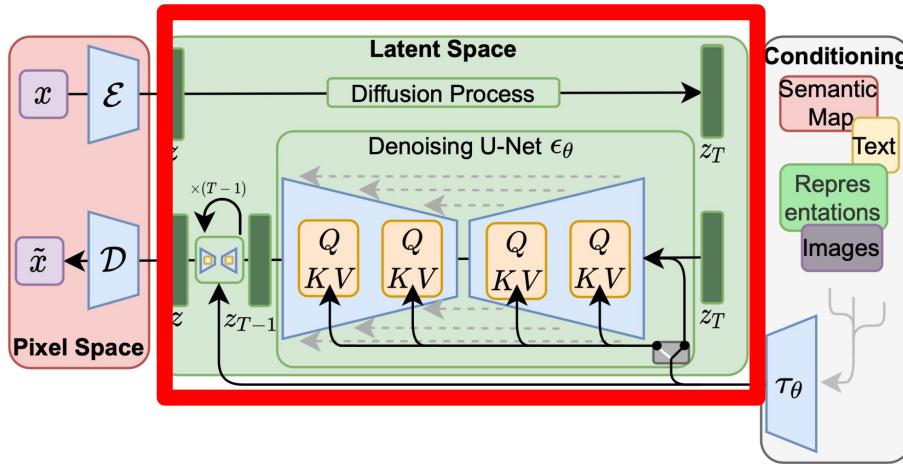
Latent Diffusion Models

We project in a discrete latent space



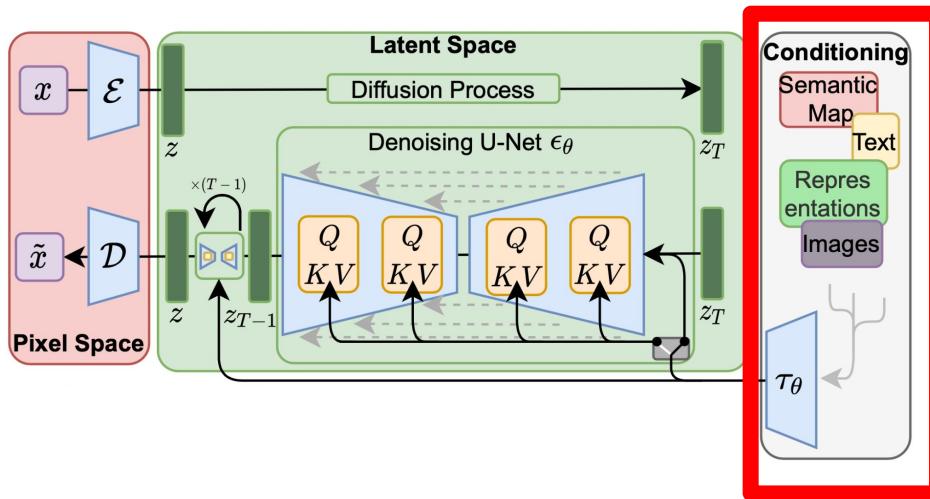
Latent Diffusion Models

The diffusion process is used in the latent space



Latent Diffusion Models

A conditional encoder can be learned



Diffusion Models

If you want to know more

- Read the papers in the references.
- Lilian Weng's blog: <https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>
- CVPR Tutorial on Denoising Diffusion-based Generative Modeling: Foundations and Applications: <https://www.youtube.com/watch?v=cS6JQpEY9cs>

Why does it change everything?

- Unprecedented image quality and data coverage.
- Fine-grained control with text.
- We have a general-purpose image prior.
- Since the training cost of large models is very high, transfer learning becomes the way to go for many problems.
- Because of the iterative generation process, many algorithms need to be adapted to be compatible with diffusion models.