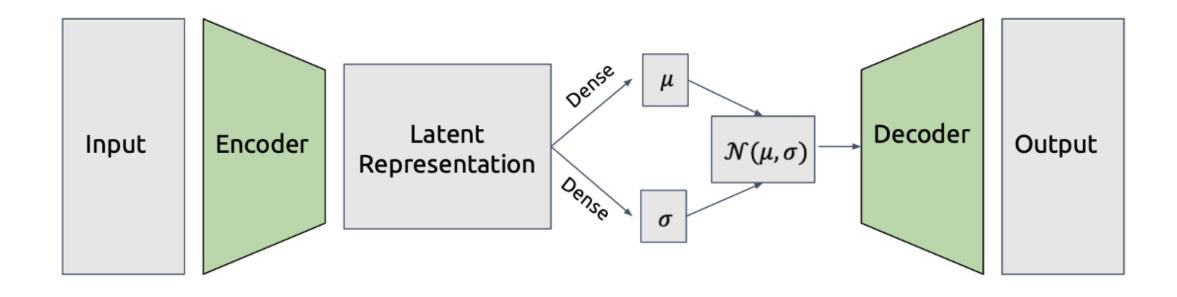
Generative Adversarial Networks (GANs)

CSCI 1470
Eric Ewing
Thursday, 10/30/25



VAEs Review

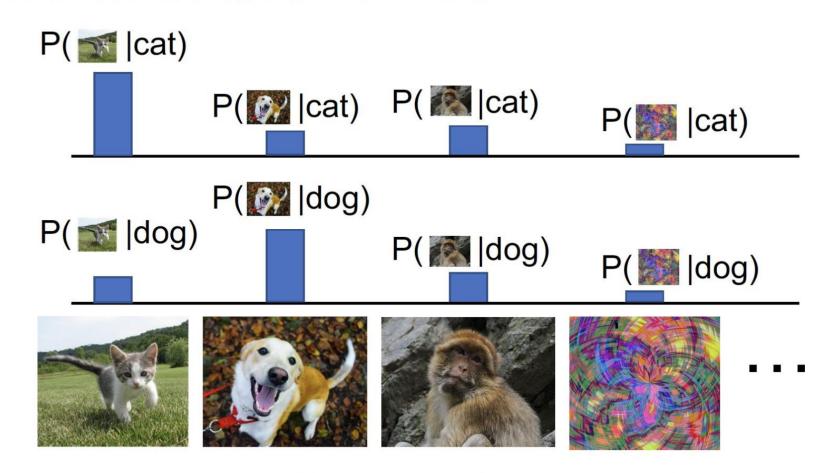


Discriminative vs Generative Models

Discriminative Model: Learn a probability distribution p(y|x)

Generative Model: Learn a probability distribution p(x)

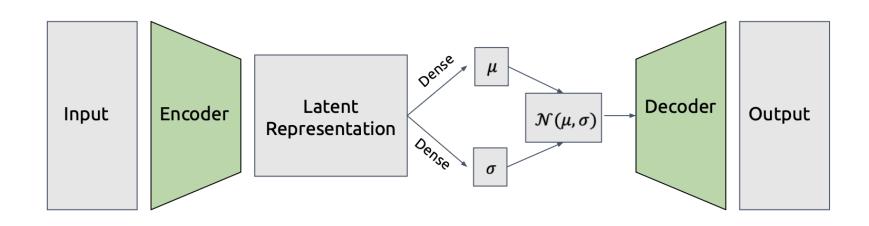
Conditional Generative Model: Learn p(x|y)



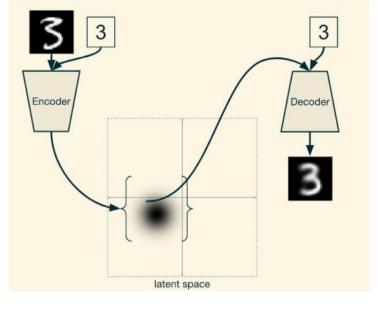
Conditional Generative Model: Each possible label induces a competition among all images

Credit: UMich EECS498

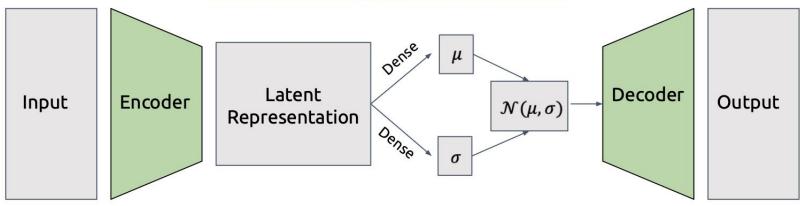
Conditional VAE



Conditional VAE



Encoder generates P(z|c)Generator generates image $P(\hat{x}|z,c)$



Why are VAE samples blurry?

- Our reconstruction loss is the culprit
- Mean Square Error (MSE) loss looks at each pixel in isolation
- If no pixel is too far from its target value, the loss won't be too bad
- Individual pixels look OK, but larger-scale features in the image aren't recognizable

Solutions?

Let's choose a different reconstruction loss!

Input



VAE reconstruction



https://towardsdatascience.com/what-the-heck-are-vae-gans-17b86023588a

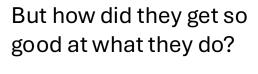
Generative Adversarial Networks (GANs)

Does nefarious things

They are adversaries

Investigates crimes

Moriarty

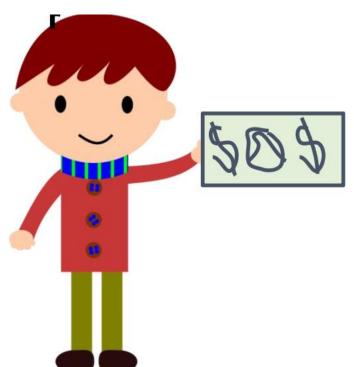


Sherlock



This Photo by Unknown Author is licensed under <u>CC BY-NC</u>

Moriarty started out as a child, trying to forge bills (he wasn't very good at it)



Young Sherlock would look at the bills and try to figure out if they were forged or not (Sherlock was also not very good at it)



Moriarty used the feedback to get even better at forging

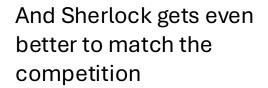


Sherlock also improved as Moriarty started producing better bills

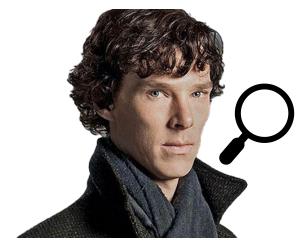


A little bit of competition can help you grow

To try and stay ahead, Moriarty further improves







 $\underline{\text{This Photo}} \text{ by Unknown Author is licensed under } \underline{\text{CC BY-NC}}$

Generative Adversarial Networks

- GANs use these ideas to train a pair of networks together
- These networks are called the Generator and Discriminator

Neural network to generate data

Generator

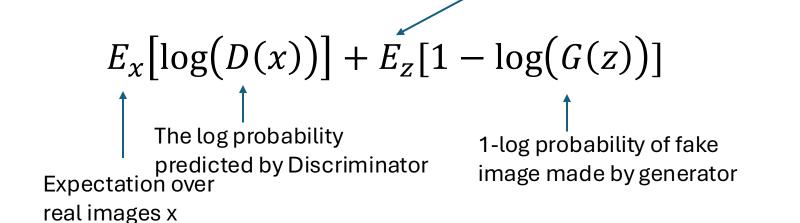
Neural network to "guess" if that data is real

Discriminator

GANs: Training the Discriminator

Discriminator wants to maximize:

Expectation over fake images (z is vector in latent space)



What loss function does this remind you of?

This is BCE, if real data are "class 1" and fake data are "class 0"

But...

The discriminator wants to **MAXIMIZE**:

$$E_x[\log(D(x))] + E_z[1 - \log(G(z))]$$

That's not actually a problem, we can minimize the negative...

What is a problem, is that the generator wants to minimize the original function

GANs as a Game

Overall Problem:

$$\min_{G} \max_{D} E_{x} \left[\log(D(x)) \right] + E_{z} \left[1 - \log(G(z)) \right]$$

Training GANs can be modeled as a 2-player zero-sum adversarial game

In an adversarial game, players have opposing goals (not everyone can win)

Zero-sum means that their goals are exactly opposite

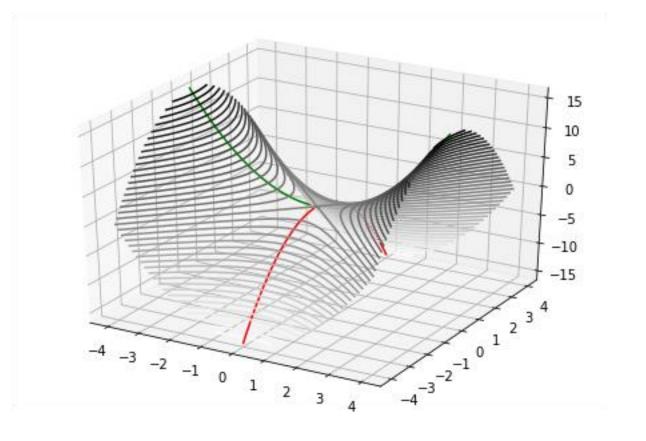
Nash Equilibrium

- A Nash Equilibrium is a solution concept for zero sum games
- Player A looks for a strategy such that no matter what player B chooses to do, player A will do no worse
- In the case of GANs, strategies are network parameter settings
- In a Nash Equilibrium, the generator will not be able to find different parameter settings such that it can do any better (and vice versa)

How do we find Nash Equilibria

 Game theory provides lots of tools, but the easiest to implement is Gradient Ascent Descent

Alternate steps of gradient ascent (using gradient of maximizer) with steps of gradient descent (using gradient of minimizer)



GAN Training Algorithm

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

Inner loop to train D

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

$$abla_{ heta_d} rac{1}{m} \sum_{i=1}^m \left[\log D\left(oldsymbol{x}^{(i)}
ight) + \log \left(1 - D\left(G\left(oldsymbol{z}^{(i)}
ight)
ight)
ight)
ight].$$

end for

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

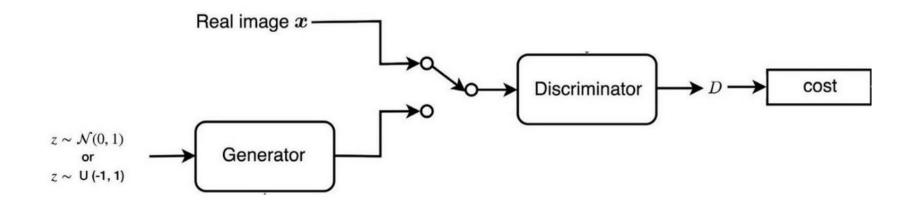
Outer loop to train G-

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

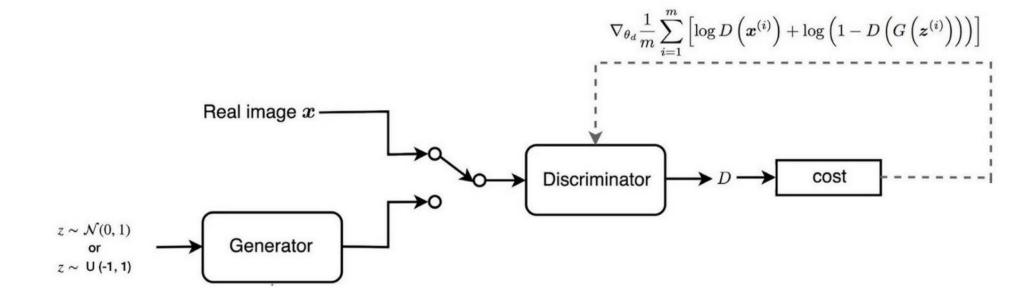
end for

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

GAN Loss



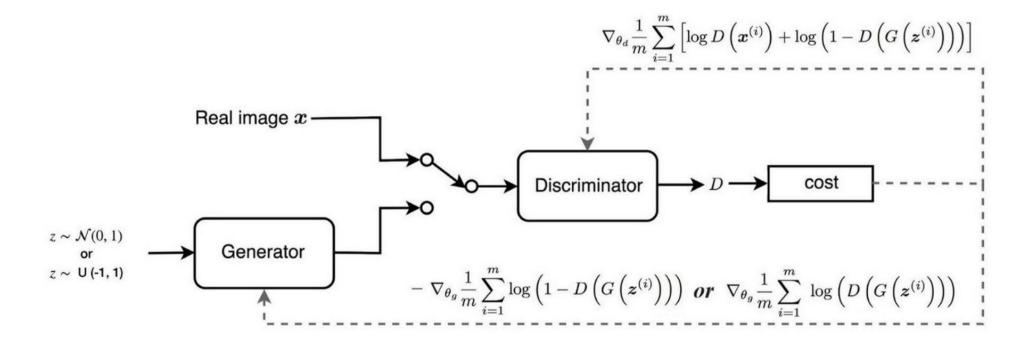
GAN Loss



Any questions?



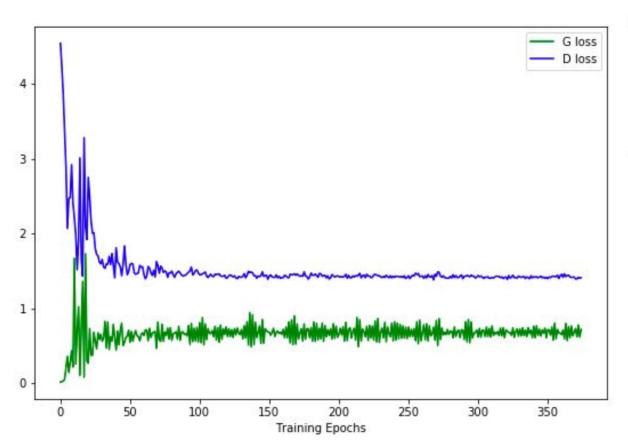
GAN Loss



Training Visualization

https://poloclub.github.io/ganlab/

GAN Training Dynamics



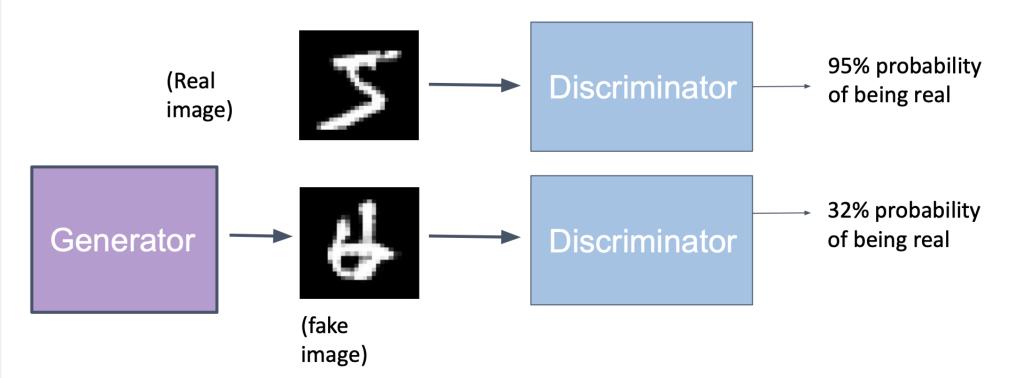
 Does not exhibit the typical "training loss continues to go down" behavior

• Why?

- Training a GAN is a "stalemate" G and D continually adjust to each other's improvements
- More formally, training a GAN to convergence is attempting to find an equilibrium of a two-player minimax game

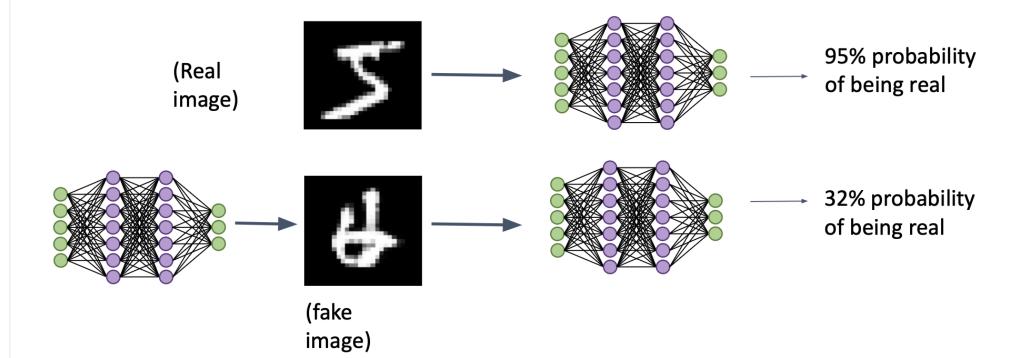
What do G and D look like inside?

Architecture of the networks determined by problem



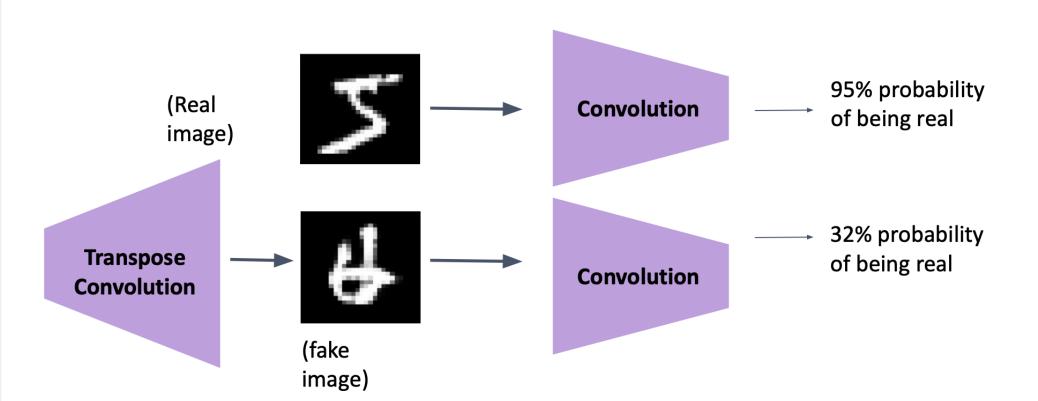
What do G and D look like inside?

- Architecture of the networks determined by problem
- Fully connected



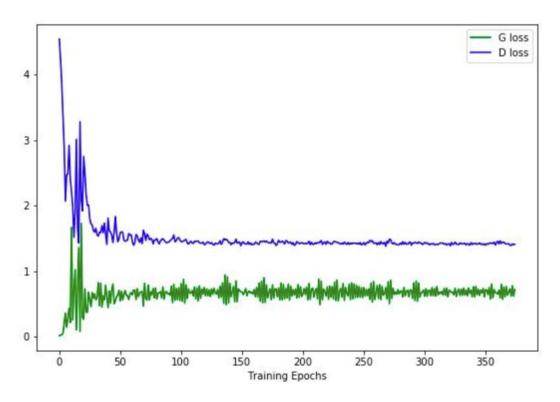
What do G and D look like inside?

- Architecture of the networks determined by problem
- Convolutional / Transpose convolutional



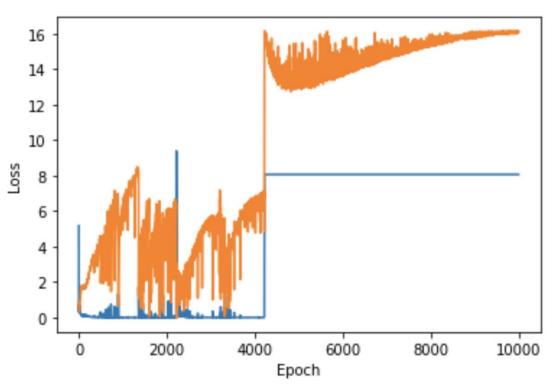
Problems with GANs

Training GANs is *very* unstable



- This is what it looks like when it's working...
- Turns out Equilibria are hard to find
 - With other networks, if the loss is going down, we know the network is doing better
 - Here, we have a moving target (G and D keep changing)
- These curves can oscillate a lot

Training GANs is *very* unstable



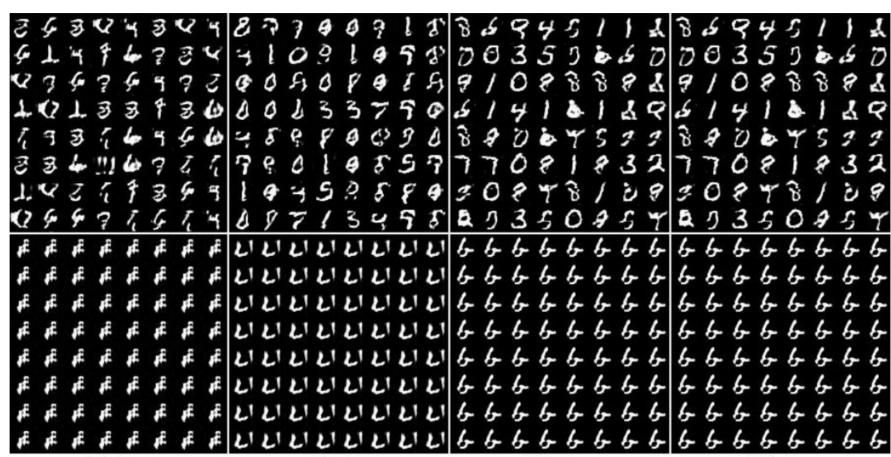
What happens if the discriminator ever becomes perfect at detecting the generator's fakes?

- Discriminator always returns p=0
- Since D always returns a constant, the gradient is constant
- The generator stops training

Mode Collapse

- Generator loss says: "generate an output that looks real"
- It does not say: "generate every output that looks real"
- The generator can "cheat" by finding one output / a few outputs that reliably fool the discriminator (the specific one(s) it finds can shift over training)

Mode Collapse



Output from a healthy GAN

Output from a GAN with mode collapse. All outputs from GAN, regardless of random input noise, are the same.

10k steps 20k steps 50K steps 100k steps

Balancing the Discriminator is hard

- We control how much to update the discriminator in each training loop
- The discriminator has the easier problem, classification is much easier than generation
- When the discriminator gets too good, gradients vanish and the generator stops updating
- When the discriminator is bad, the generator can start producing the same output for every input (mode collapse)



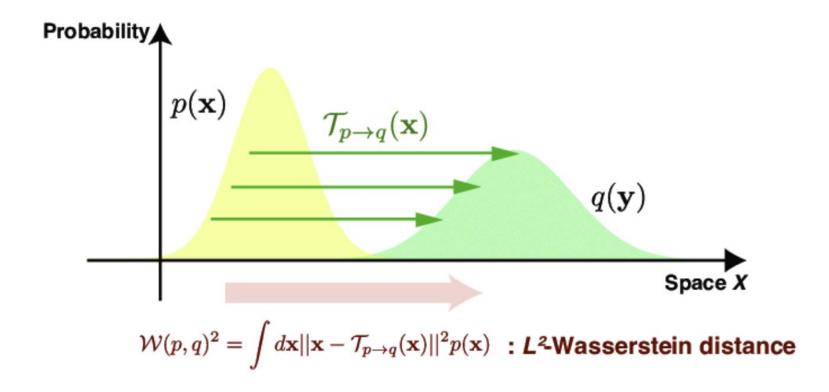
This Photo by Unknown Author is licensed under <u>CC BY-NC</u>

Wasserstein GANs

- What if we could train the discriminator to convergence and still be able to train our generator?
- We want a function that can tell us how likely it is an image came from the real distribution, but not have vanishing gradients
- Solution: Use a "critic" instead of a discriminator
- The critic outputs a score (value) for the generator rather than a probability (i.e., don't use a sigmoid for activation...)

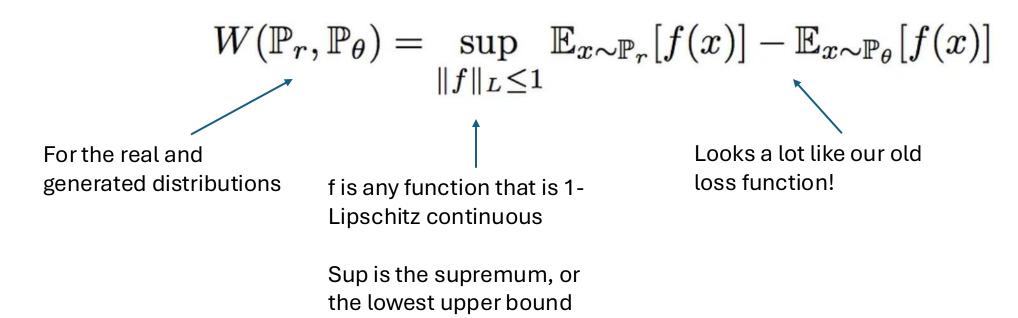
Wasserstein Distance

Wasserstein (Earth Mover's) Distance: How much work is it to move one probability distribution p, to another distribution q.



Wasserstein Distance

- Integrals are often hard to compute...
- Use Kantorovich-Rubinstein duality to simplify computation



Derivation: https://abdulfatir.com/blog/2020/Wasserstein-Distance/

Lipschitz Continuity

- For this method of calculating Wasserstein Distance, we need our critic to be Lipschitz continuous (with c=1)
- That means that maximum gradient has to be 1
- There are fancy ways to do this (i.e., spectral normalization), but what if just *clip* the gradient
- If the gradient is larger than 1, just clip the value to be 1.
- tf.clip_by_value(grads, -1, 1)

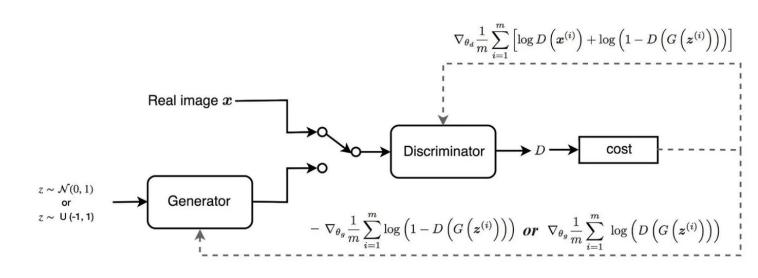
Any questions?



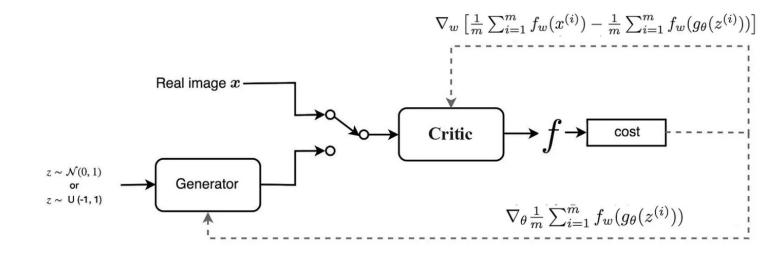
Key difference for WGANs:

Instead of using the log probability of the output of a discriminator, use the output of the critic

Wasserstein GANs are more stable than vanilla GANs because it's less susceptible to vanishing gradients from the discriminator GAN:



WGAN

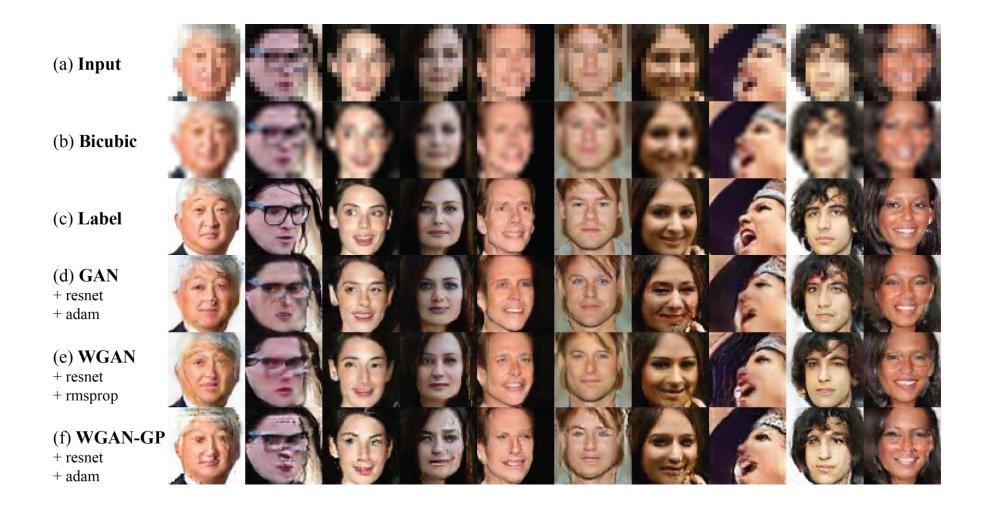


Other Tasks

- What if instead of generating images, we wanted to do another task
 - Denoising Images: Given a corrupted or noisy image, can you remove the noise?
 - Colorization: Given a black and white image, can you produce a color image?
 - Super-resolution: given an image of low resolution, can you produce an image of higher resolution (i.e., more pixels)?

How would you frame each of these problems and train a GAN (or VAE)?

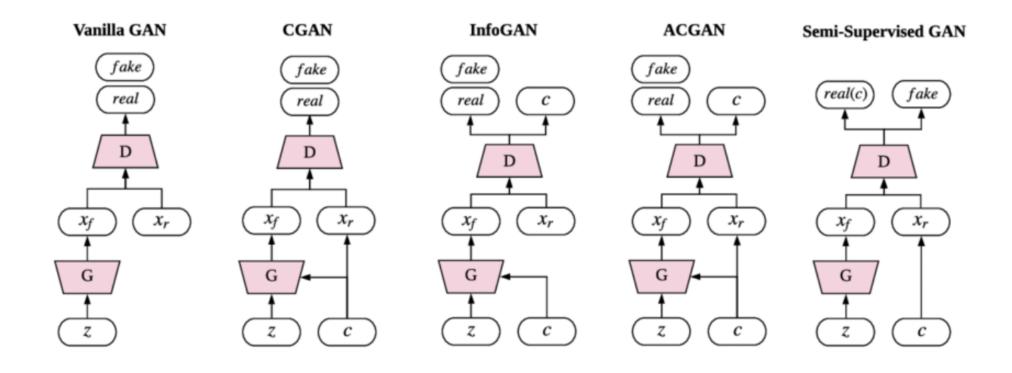
Input low resolution image to generator, have it output a higher resolution image



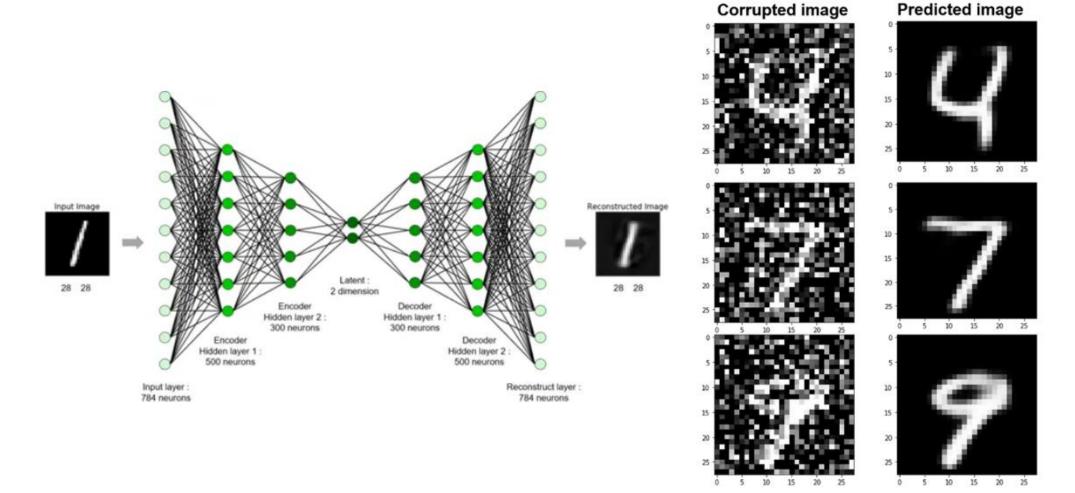
Source: https://ar5iv.labs.arxiv.org/html/1705.02438

GAN Variants

There are many variations on GANs... (these are just some of the ones with architectural differences)

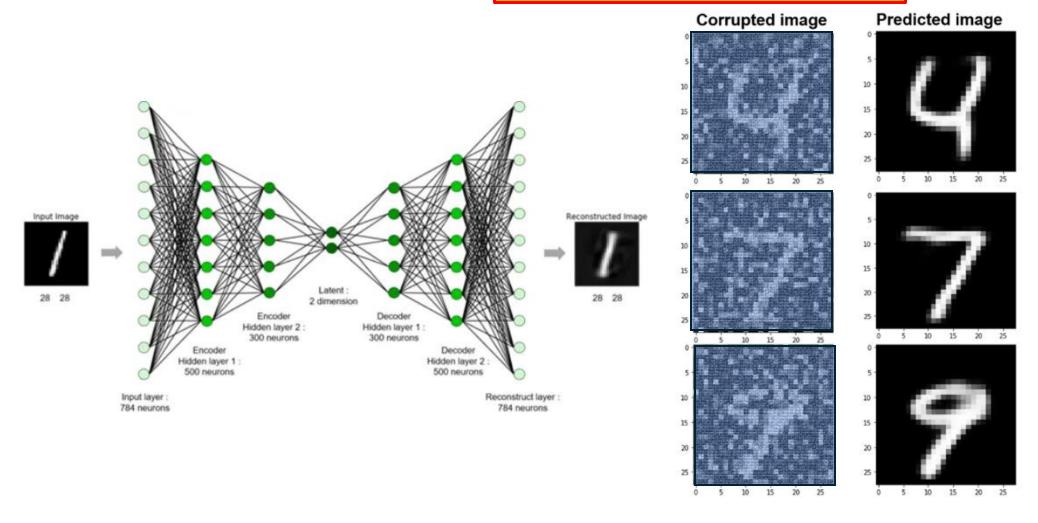


Denoising Auto Encoders (DAEs)



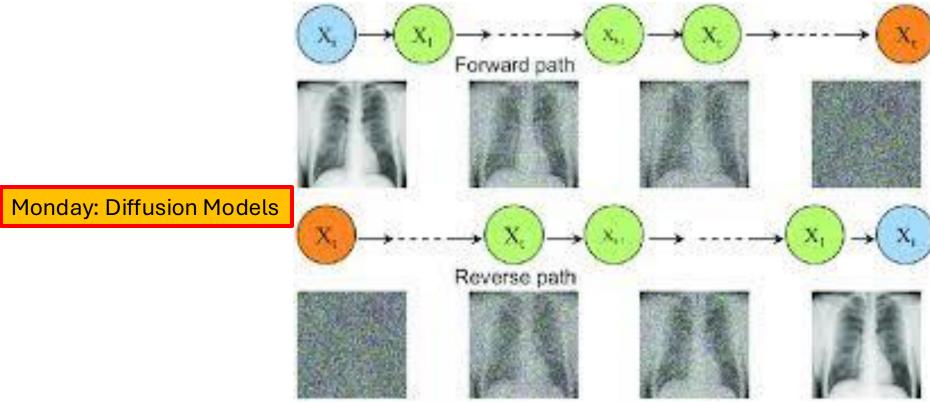
Denoising Auto Encoders (DAEs)

What if our images were even noisier?



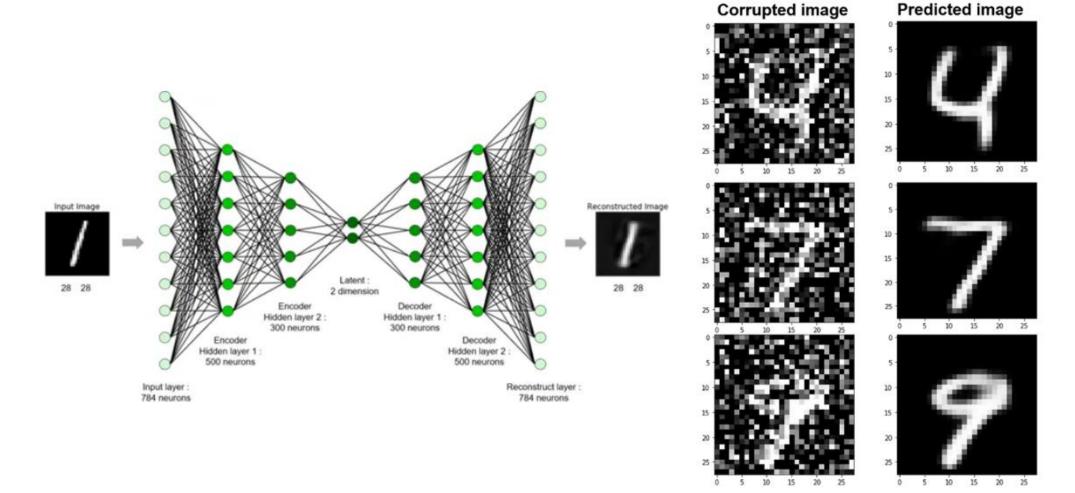
Denoising Auto Encoders

Denoising all the noise in one step may be too hard to learn. What if we added and removed noise incrementally?



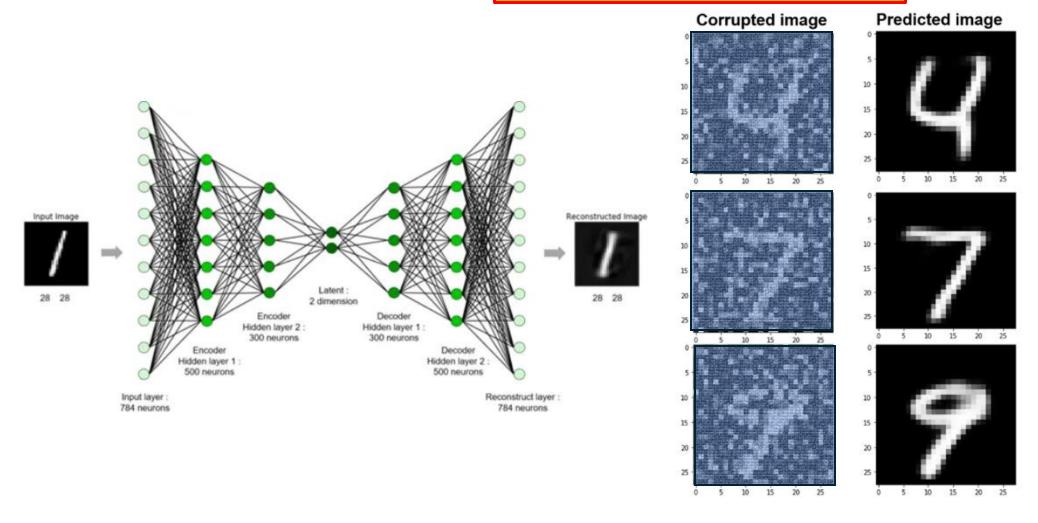
Source: https://www.researchgate.net/figure/Schematic-of-the-diffusion-model-training-process_fig1_383920783

Denoising Auto Encoders (DAEs)



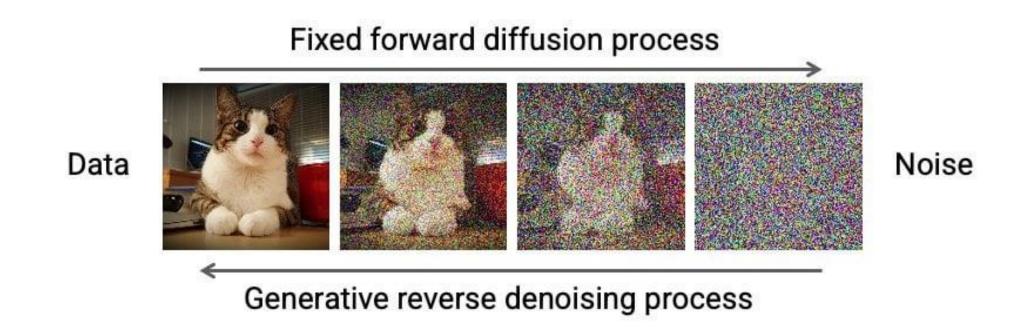
Denoising Auto Encoders (DAEs)

What if our images were even noisier?



Denoising Auto Encoders

Denoising all the noise in one step may be too hard to learn. What if we added and removed noise incrementally?

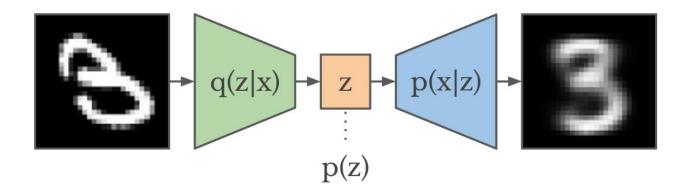


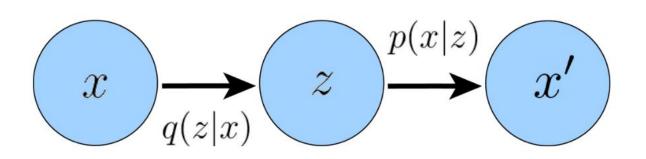
Variational Autoencoders

We can represent a VAE as a probabilistic graphical model

Encoder generates probability distribution q(z|x)

Decoder estimates p(x|z)

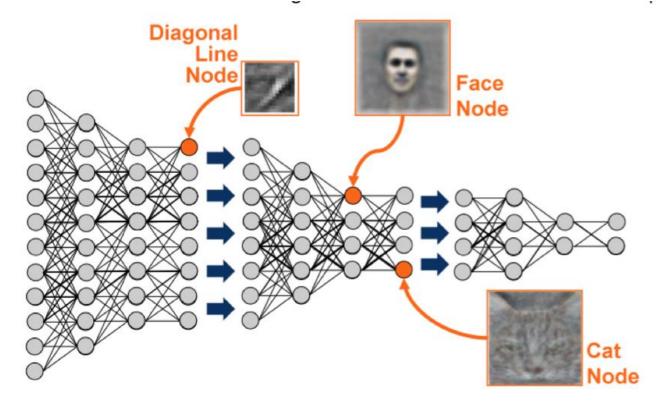




Hierarchical Features

Each intermediate layer in a neural network can be seen as learning a set of **intermediate features** based on the previous **intermediate**

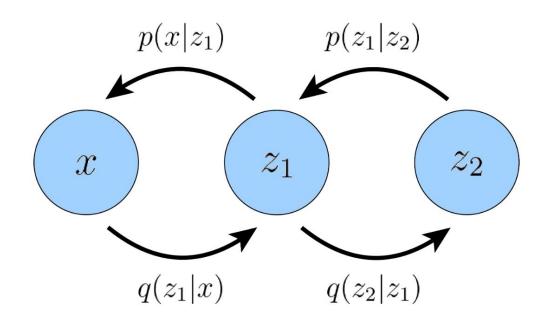
outputs.

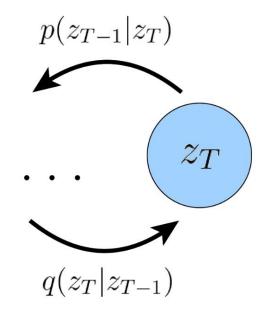


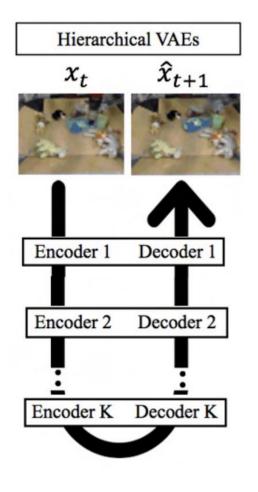
Source: https://ermongroup.github.io/blog/hierarchy/

Idea: train K pairs of encoders and decoders

Each decoder is trained to reverse the associated encoder's operations





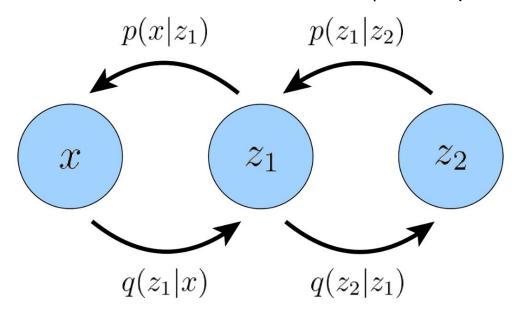


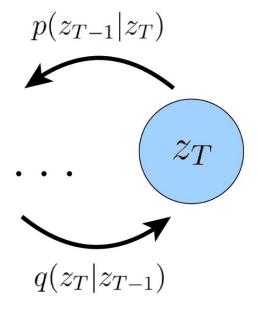
How many encoders and decoders do we need to learn?

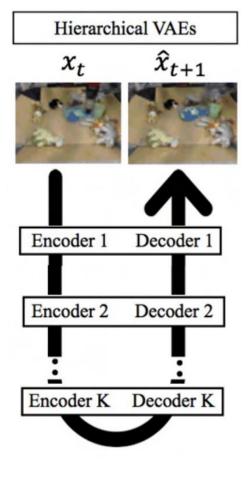
Idea: train K pairs of encoders and decoders

Each decoder is trained to reverse the associated encoder's operations

Decoders (reverse process)



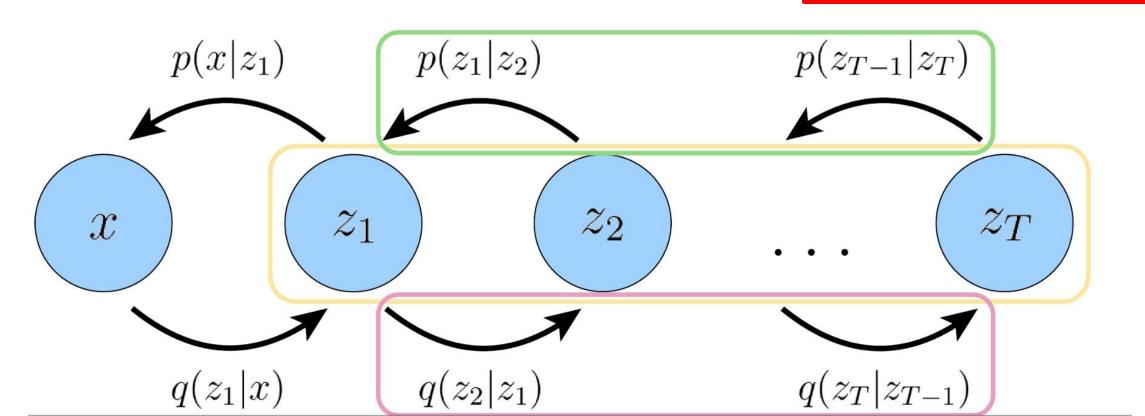




Encoders (forward process)

What if all latent variables z have the same size?

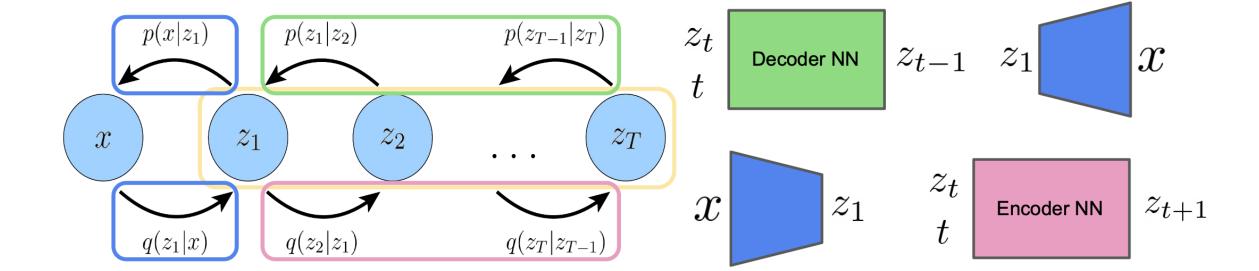
Most encoders and decoders have the same input/output size



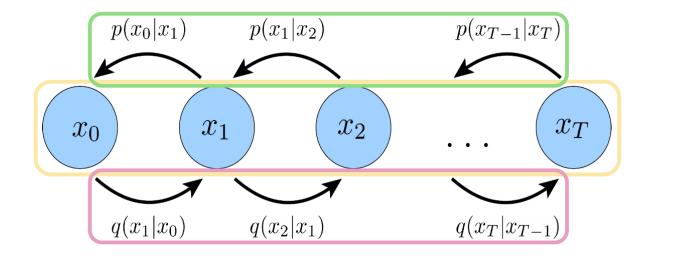
But what if **everything** had the same dimensions

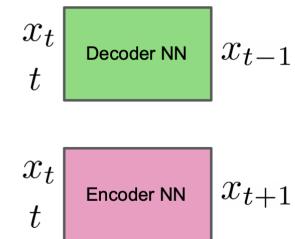
What if all latent variables z have the same size?

Need special **first encoder**, special **last decoder** to go from embedding size to original input size

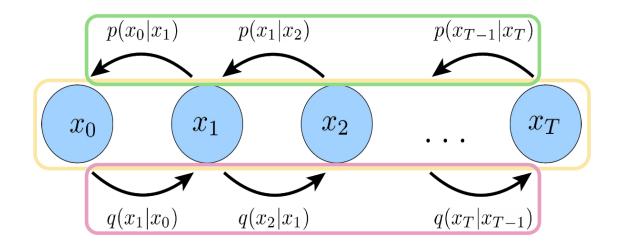


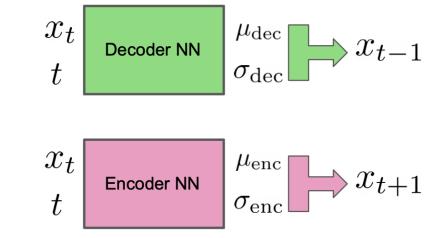
A hierarchical VAE that keeps the dimensions the same can use the same decoder and encoder for all steps





A hierarchical VAE that keeps the dimensions the same can use the same decoder and encoder for all steps





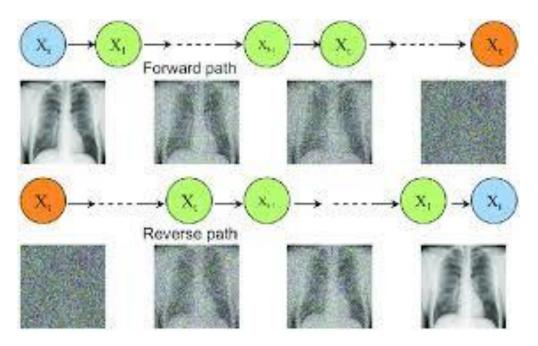
Denoising VAEs

If the size of the encoding in a VAE is the same size as the input, the model is no longer doing dimensionality reduction...

So why would we want this?

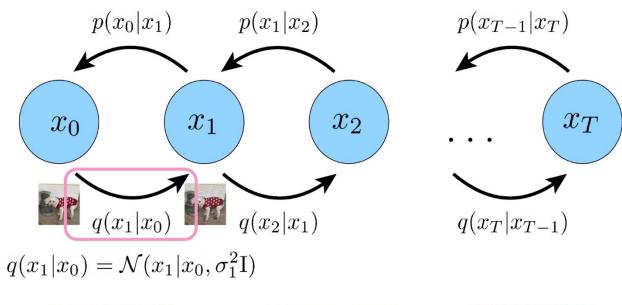
The forward process adds noise
The reverse path reverses this process

But we don't need to learn the encoders, adding noise isn't a learned process!

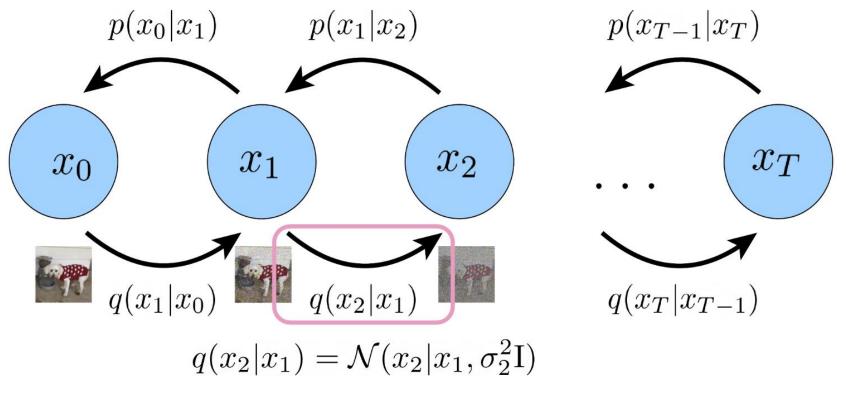


Adding Noise

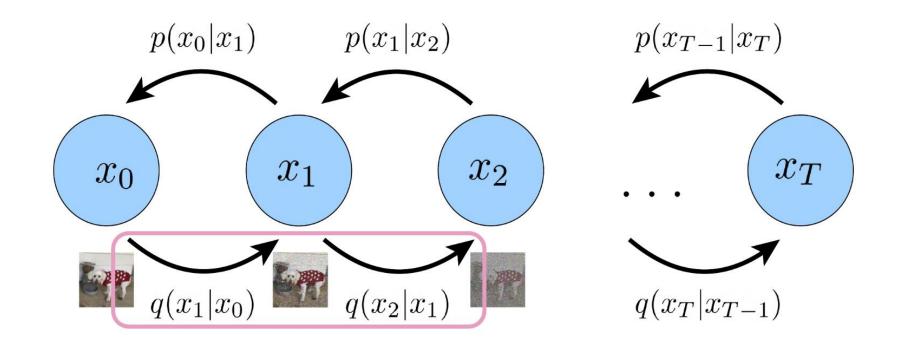
What if each encoder transitions sample the input with some gaussian noise?

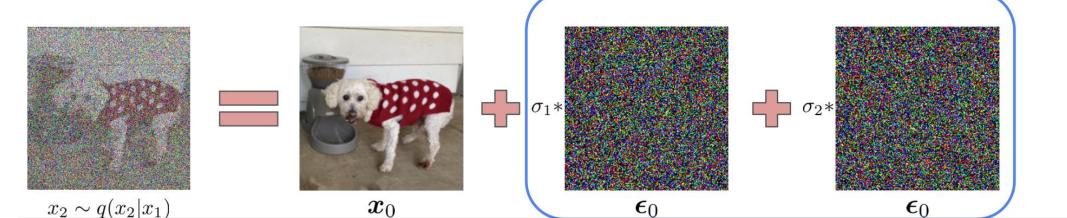




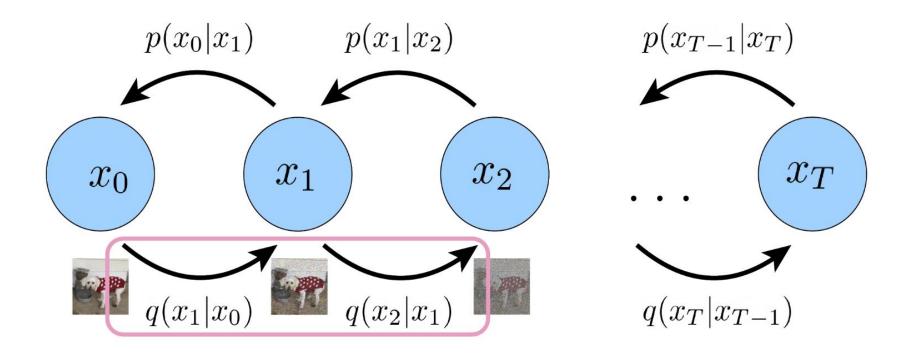








Aggregate into 1 sample!





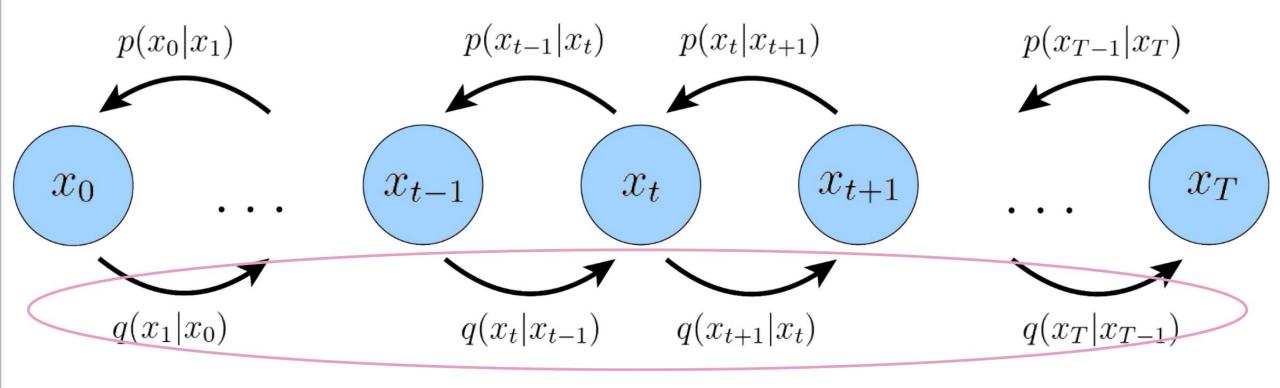
where,

$$\alpha_2 = \sqrt{\sigma_1^2 + \sigma_2^2}$$

and,

$$q(x_2|x_0) = \mathcal{N}(x_2|x_0, \alpha_2^2)$$

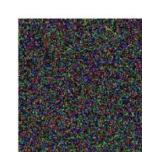
Aggregate into 1 sample!

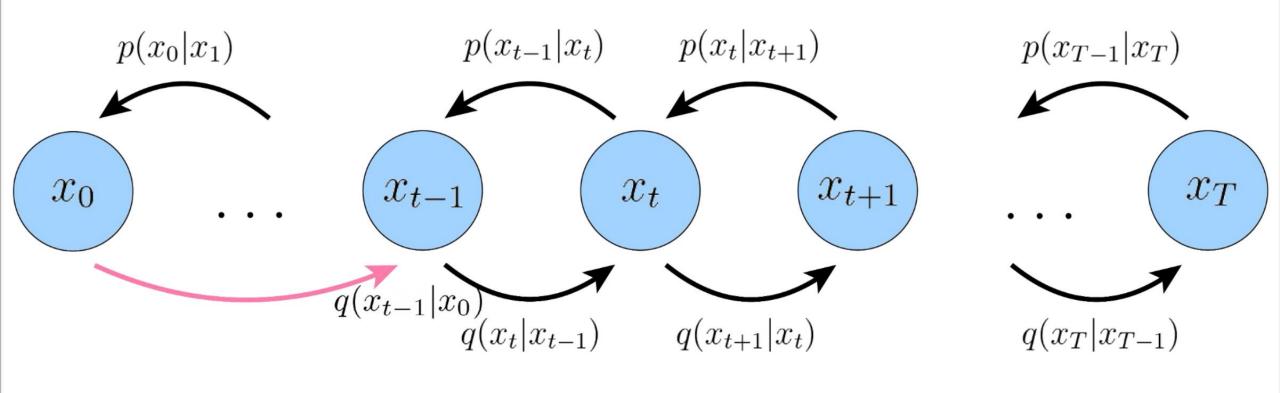


Individual (known) Gaussians!





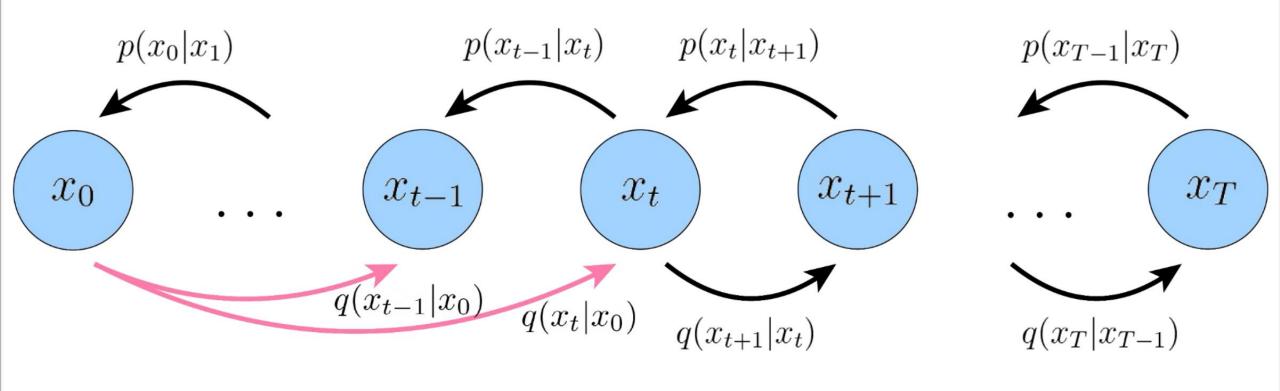










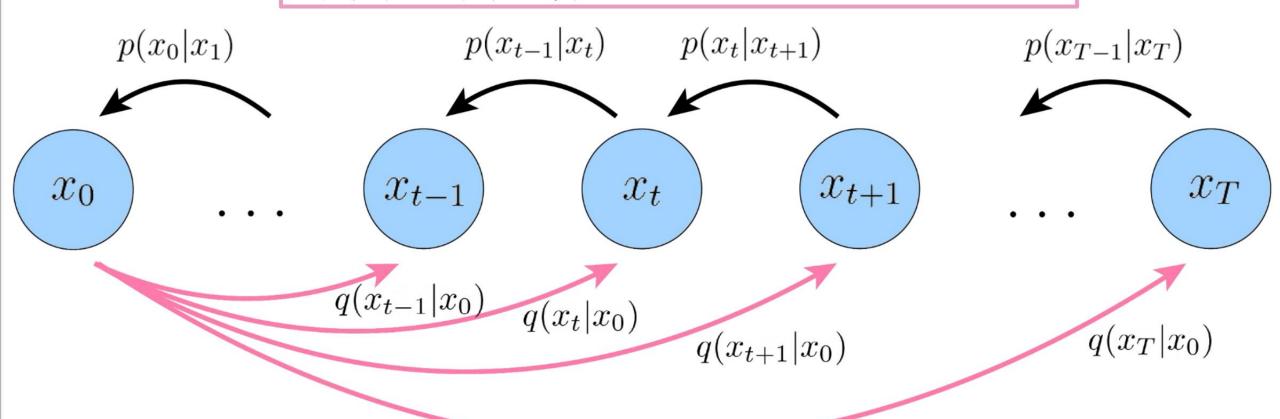








 $q(x_t|x_0)$ is a Gaussian, for arbitrary t! $q(x_t|x_0) = \mathcal{N}(x_t|x_0,\alpha_t^2\mathrm{I}) \text{, where } \alpha_0,\alpha_1,...\alpha_T \text{ are all known/fixed.}$





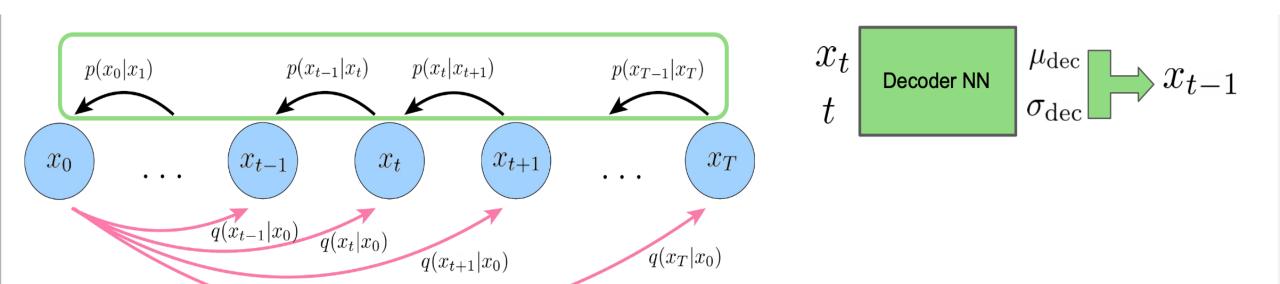




Diffusion Models

Are hierarchical VAEs with the following assumptions:

- All dimensions are the same (input size, encoding size)
- Encoder transitions are known gaussians centered around their previous inputs

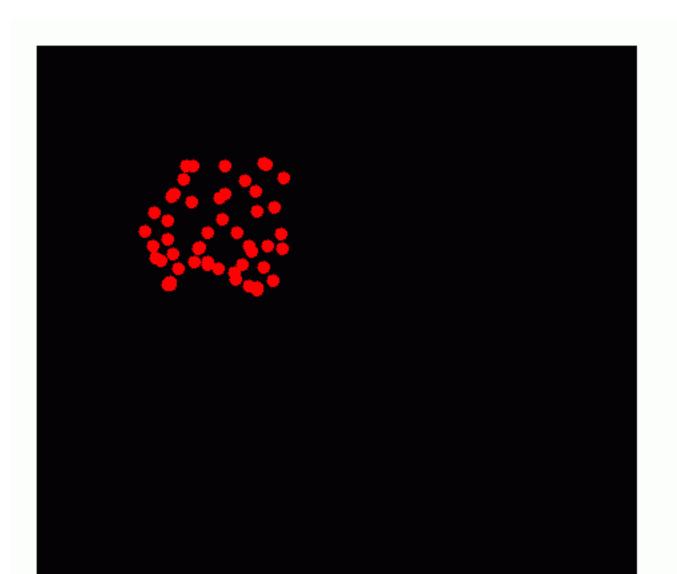


Why Call it Diffusion?

Diffusion of gas particles (and other physical things)

Start off organized

Transitions to "random noise"



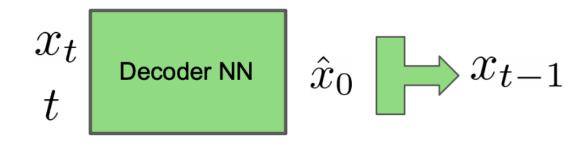
The Decoder

Diffusion models seek to learn a single neural network: a de-noising decoder

What form does x_{t-1} take?

$$q(x_{t-1}|x_0) = \mathcal{N}(x_{t-1}|x_0, \alpha_{t-1}^2 I)$$

$$x_t = x_0 + \alpha_{t-1}\epsilon$$



$$\hat{x}_{t-1} = \mu_{\text{dec}} + \sigma_{\text{dec}} * \epsilon$$

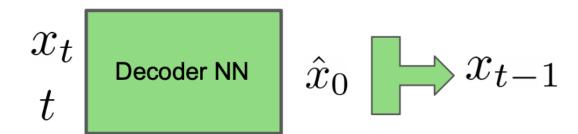
What is the ground truth for μ_{t-1} ?

Does the decoder even need to predict σ ?

Denoising Diffusion Models

Learn a single decoder network

- Input is image with added noise
- Output is predicted image with noise removed



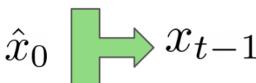
How To Generate New Samples

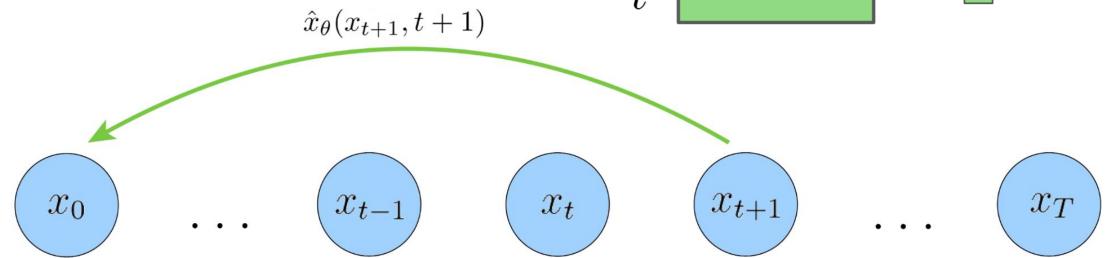
Idea 1: sample point from $\mathcal{N}(0,I)$, run decoder with t=T to generate \hat{x}_0 Issue: Doesn't work that well... that was the entire motivation for having multiple steps

Idea 2: sample point from $\mathcal{N}(0,I)$, iteratively generate \hat{x}_{t-1}

But how do we actually generate \hat{x}_{t-1} when our decoder generates \hat{x}_0 ?

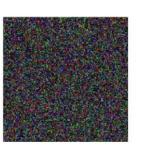
 $egin{array}{c|c} x_t \\ t \end{array}$ Decoder NN

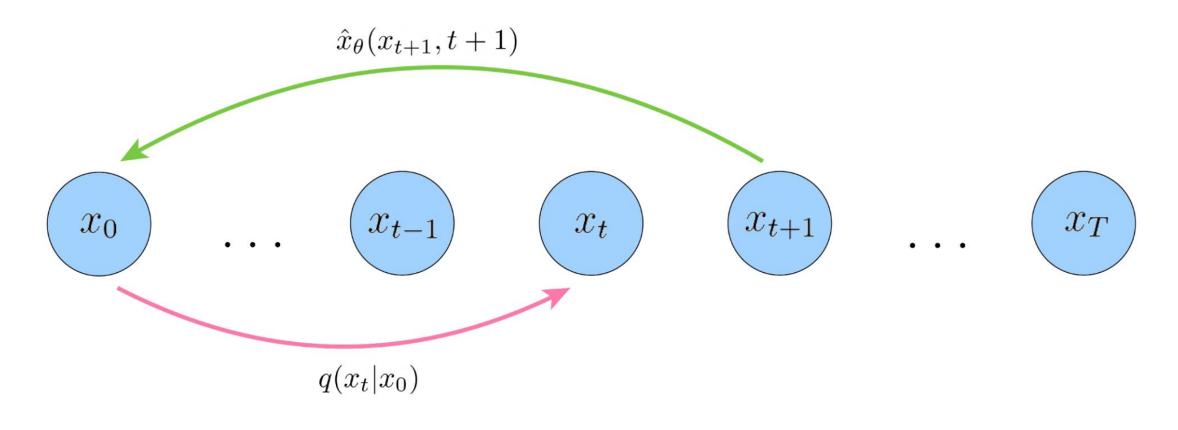






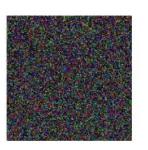


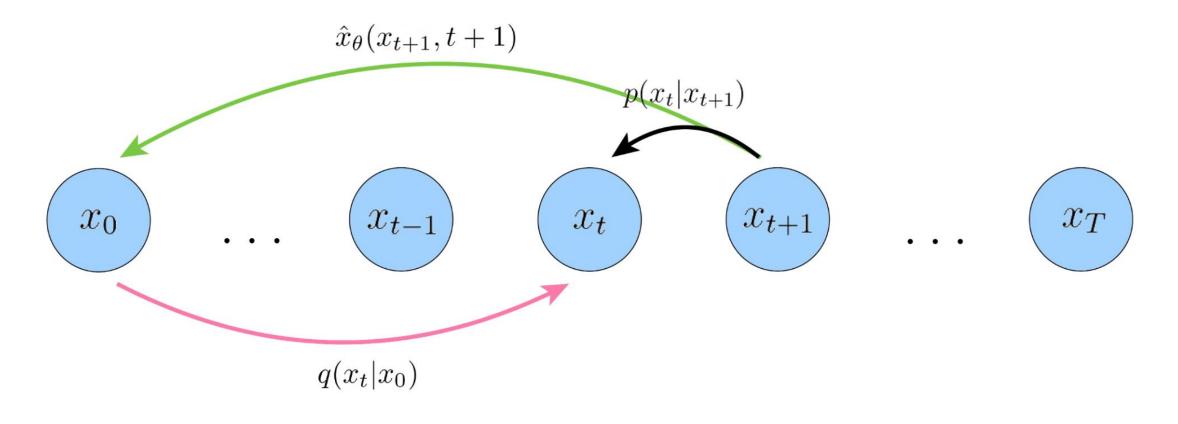






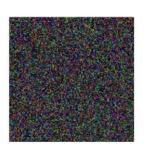


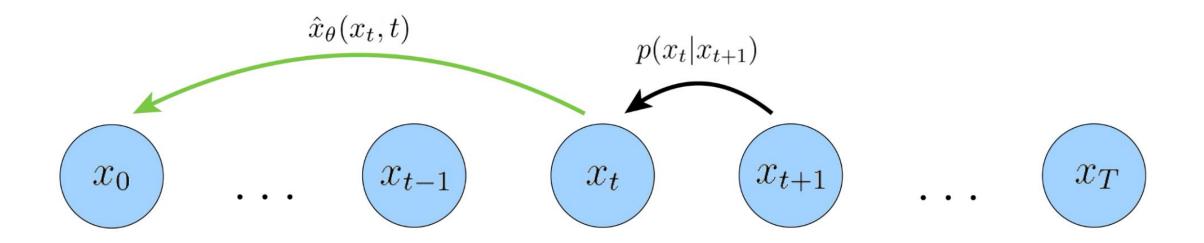






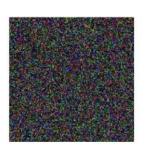


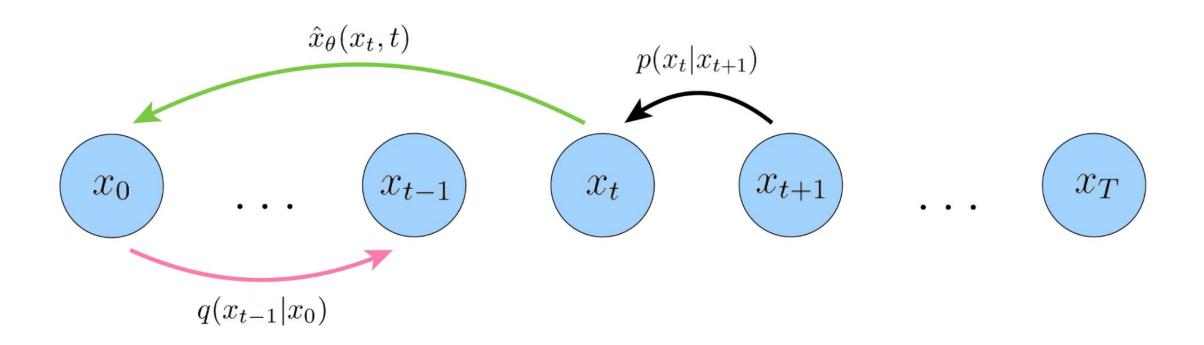






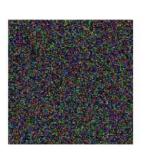


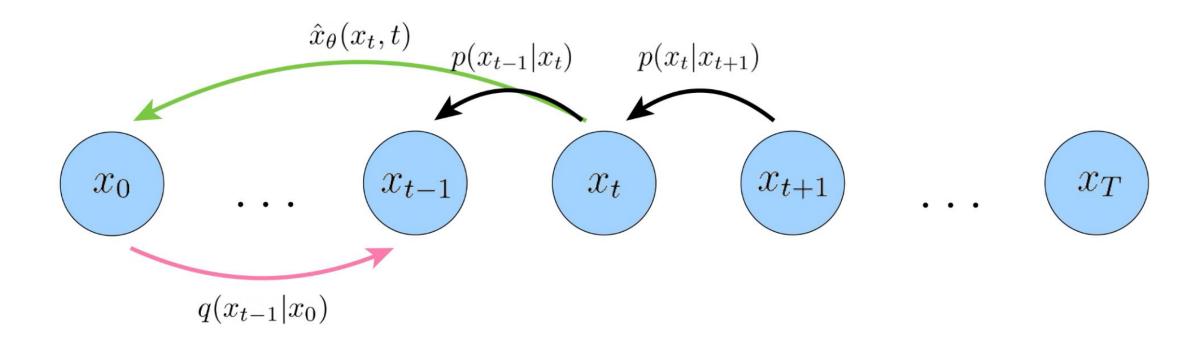






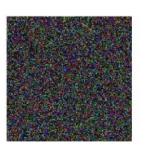


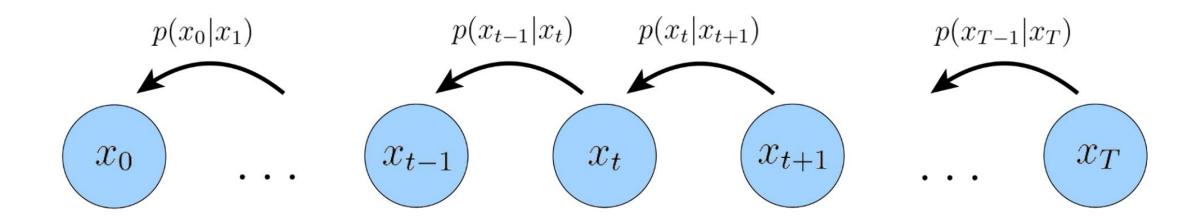






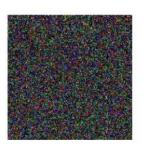












Noise Schedules

- What should the value of T be?
- How many steps of forward/reverse processes should we run?
- How much noise should be added at each step?

Amount of noise and number of steps determined by a *noise* schedule (hyperparameter)

Linear Schedule (equal noise added at each timestep)



Noise Schedules

- What should the value of T be?
- How many steps of forward/reverse processes should we run?
- How much noise should be added at each step?

Amount of noise and number of steps determined by a *noise* schedule (hyperparameter)

Cosine Schedule (small amounts of noise first, then fast)

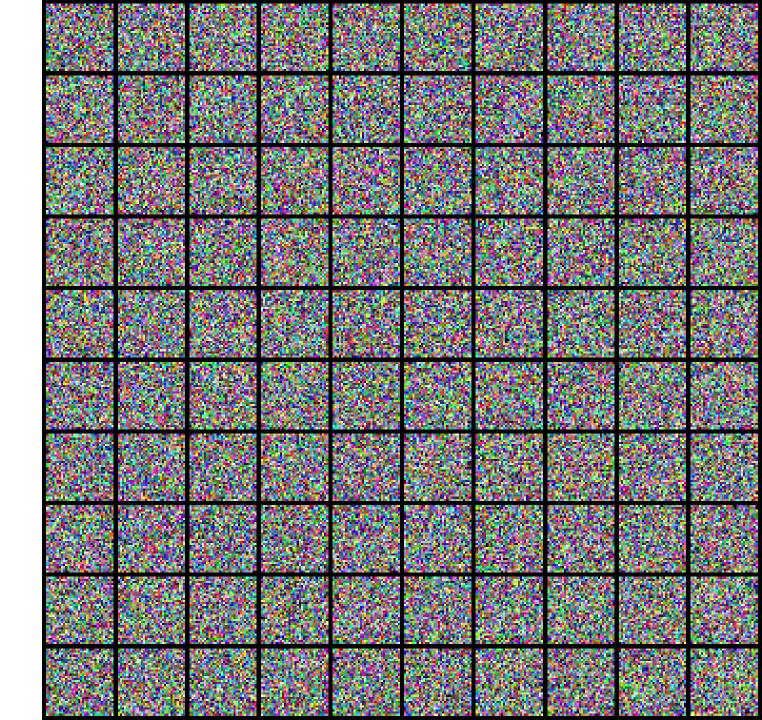


Diffusion Training

Algorithm 1 Training	Algorithm 2 Sampling
1: repeat	1: $oldsymbol{x}_T \sim \mathcal{N}(0, \mathbf{I})$
2: $oldsymbol{x}_0 \sim q(oldsymbol{x}_0)$	2: for $t = T,, 1$:
3: $t \sim \mathtt{Uniform}\left(1,,T\right)$	3: $\epsilon \sim \mathcal{N}(0, \mathbf{I}) \text{ if } t > 1, \text{ else } \epsilon = 0$
4: $oldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$	4: $\boldsymbol{x}_{t-1} = \hat{\boldsymbol{x}}_{\theta}(\boldsymbol{x}_t, t) + \alpha_{t-1}\epsilon$
5: Take gradient descent step on	5: end for
6: $\nabla_{\theta} \boldsymbol{x}_0 - \hat{\boldsymbol{x}}_{\theta} (\boldsymbol{x}_0 + \alpha_t \boldsymbol{\epsilon}, t) ^2$	6: $\mathbf{return} \ \boldsymbol{x}_0$
7: until converged	

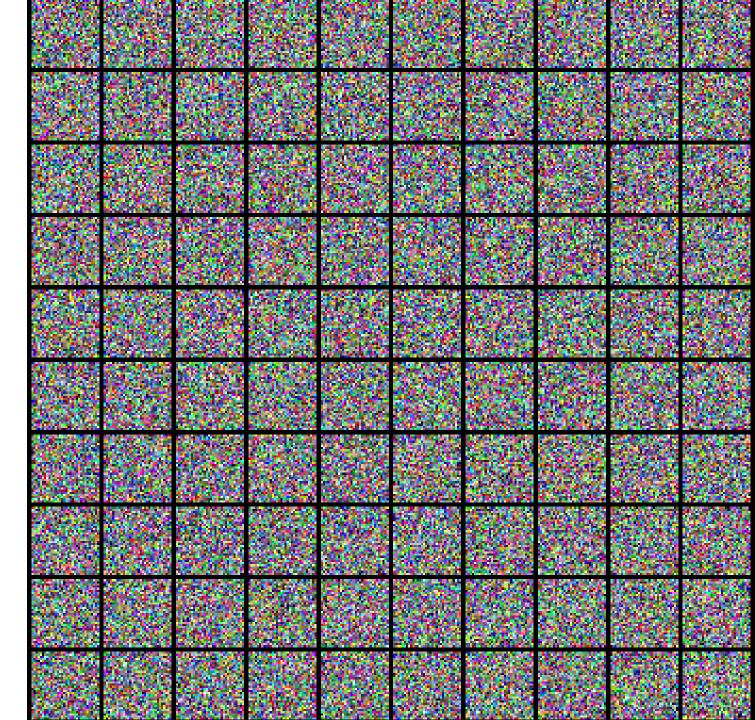
Examples

Model trained on CelebA dataset



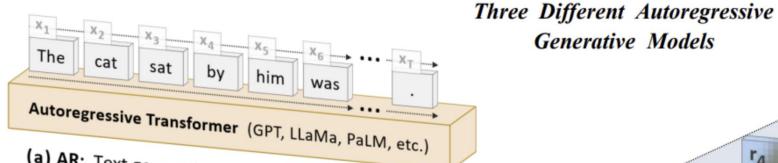
Examples

Model trained on CIFAR-10

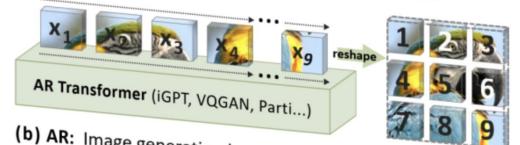


Source: https://yang-song.net/blog/2021/score/

Visual Auto-Regressive Generation



(a) AR: Text generation by next-token prediction



(b) AR: Image generation by next-image-token prediction

