

DIFFUSION MODELS, WHAT IS THAT ALL ABOUT



IS IT GOOD? OR IS IT WACK?

About Me

I'm Calvin, a 3rd year PhD student here at Brown, advised by Chen Sun

- I work on Diffusion Modeling and Reinforcement Learning

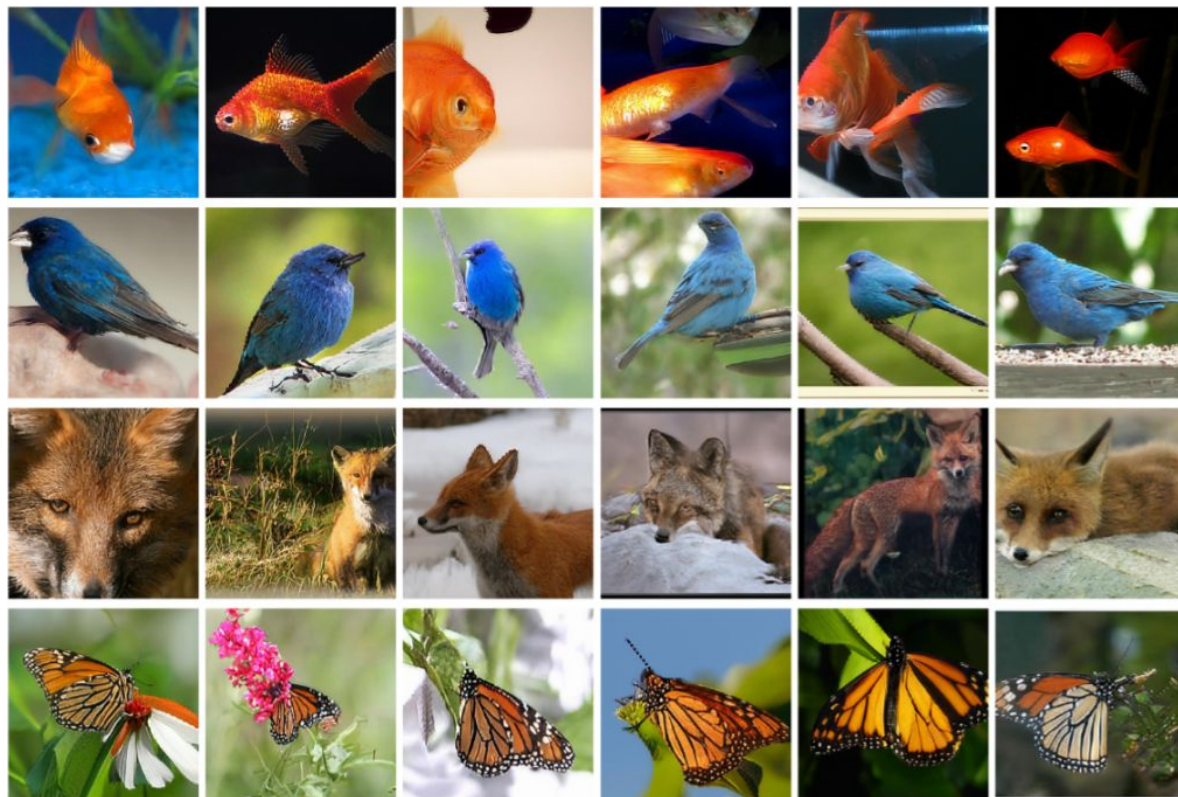


Diffusion models are really good at learning conditional distributions.

$$p(x | y)$$

Use Case: Class-Conditioned Generation

$$p(\text{image} \mid \text{class_label})$$



source: [Image Super-Resolution via Iterative Refinement](#)

Use Case: Text-to-Image Generation

$$p(\text{image} \mid \text{text_caption})$$

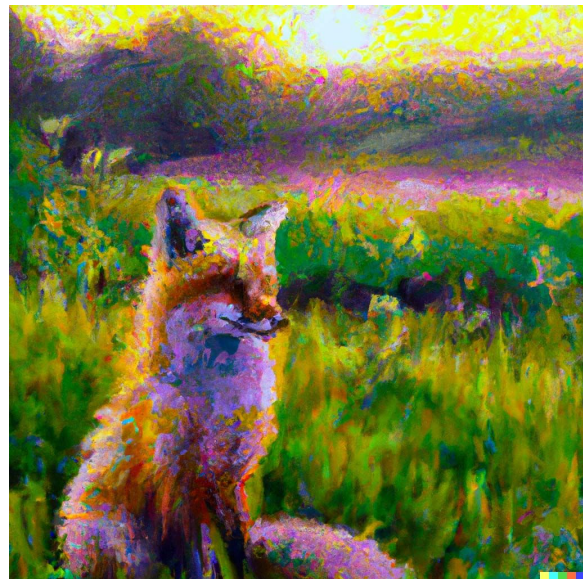
“a painting of a fox sitting in a field at sunrise in the style of Claude Monet”



Parti (but pretend it is ImageN)



StableDiffusion



Dall-E 2.0

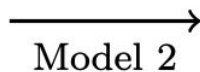
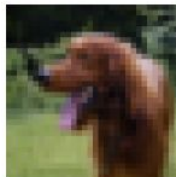
Use Case: Super Resolution

$$p(\text{image} \mid \text{low_res})$$

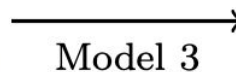
Class ID = 213
"Irish Setter"



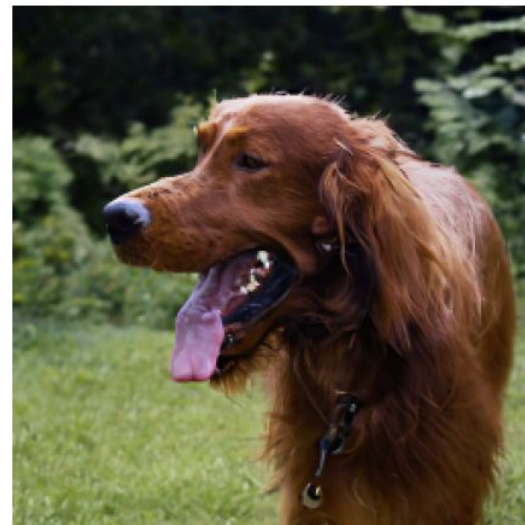
32×32



64×64



256×256



Recap: Generative Modeling

Recall the goal of generative modeling - learning a *model* of a distribution from which we can *generate* new samples.

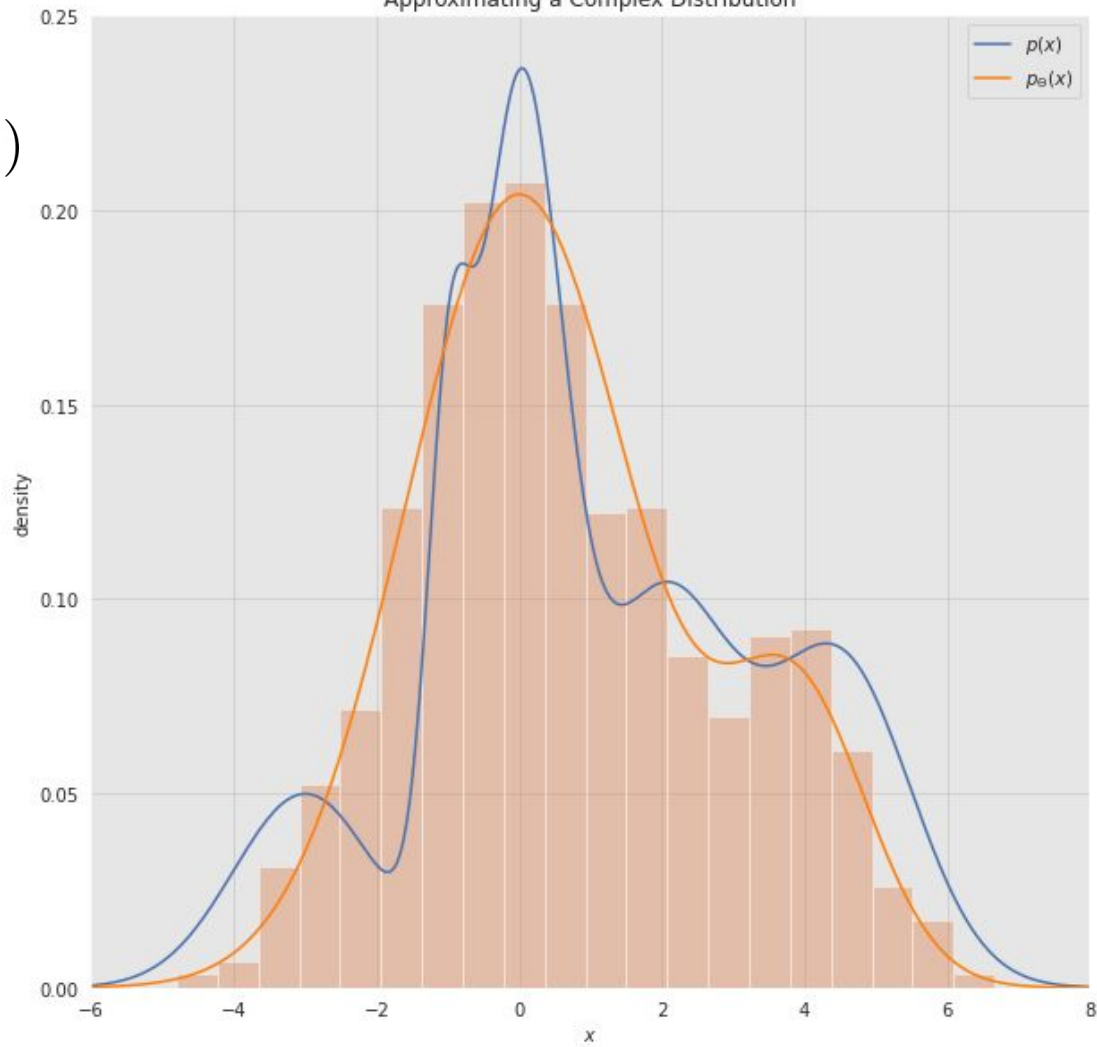
Given $\mathbf{x} \sim p(\mathbf{x})$ we might want to learn $p_{\theta}(\mathbf{x}) \approx p(\mathbf{x})$ (*modeling*)

Then, we can generate new samples $\mathbf{x}^* \sim p_{\theta}(\mathbf{x})$ (*generation*)

Why is this useful?

Approximating a Complex Distribution

Given $\mathbf{x} \sim p(\mathbf{x})$



Generative Modeling: Themes

What are some common themes of generative modeling?

- We want to learn some **complex** distribution $p_{\theta}(\mathbf{x}) \approx p(\mathbf{x})$
- But we only have access to some **simple** distributions (such as Gaussians)



Generative Modeling: Themes

What are some common themes of generative modeling?

- We want to learn some **complex** distribution $p_{\theta}(\mathbf{x}) \approx p(\mathbf{x})$
- But we only have access to some **simple** distributions (such as Gaussians)

Idea: Let's learn a complex function (aka a neural network) to transform a simple distribution sample into a complex one!

- Gaussian Sample == (neural net) ==> Data Sample

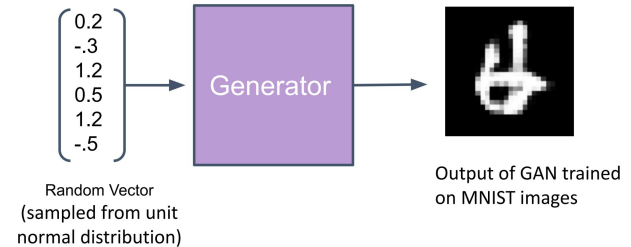
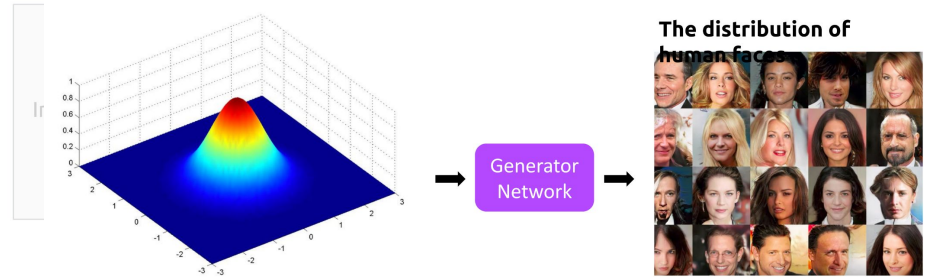


Generative Modeling: Themes

Idea: Let's learn a complex function (aka a neural network) to transform a simple distribution sample into a complex one!

- Gaussian Sample == (neural net) ==> Data Sample

You have seen this before in:



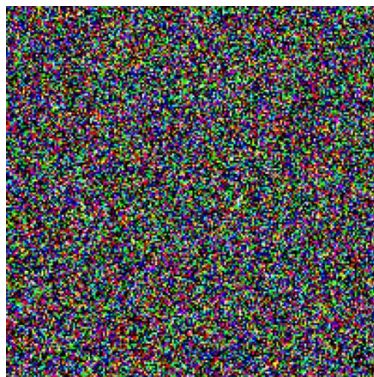
Diffusion Models: a TLDR

An observation: adding steady amounts of Gaussian noise eventually corrupts an image into something indistinguishable from a random Gaussian sample.

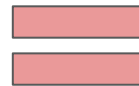


Diffusion Models: a TLDR

An observation: adding steady amounts of Gaussian noise eventually corrupts an image into something indistinguishable from a random Gaussian sample.

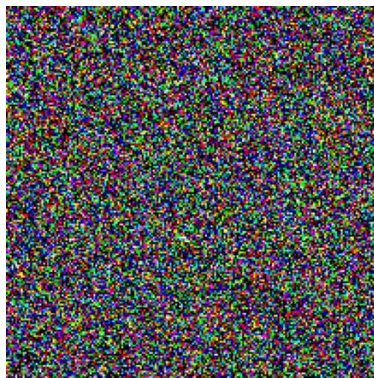


One Step

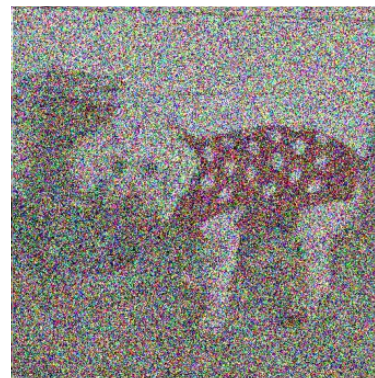
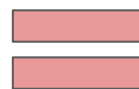


Diffusion Models: a TLDR

An observation: adding steady amounts of Gaussian noise eventually corrupts an image into something indistinguishable from a random Gaussian sample.



Another Step



Diffusion Models: a TLDR

An observation: adding steady amounts of Gaussian noise eventually corrupts an image into something indistinguishable from a random Gaussian sample.

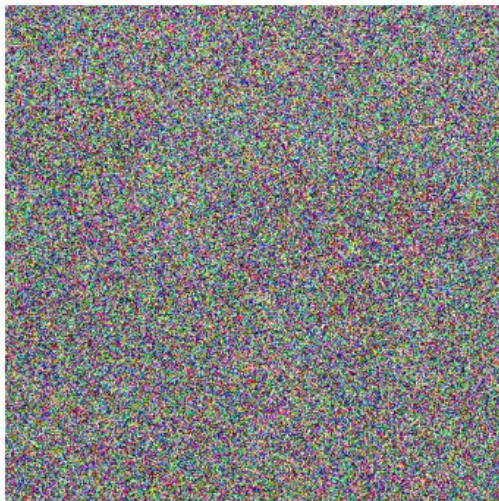


Many Steps

Diffusion Models: a TLDR

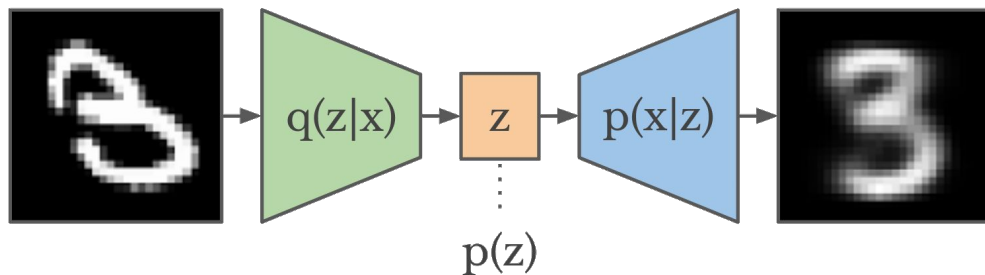
An observation: adding steady amounts of Gaussian noise eventually corrupts an image into something indistinguishable from a random Gaussian sample.

- Diffusion models simply learn to **reverse** this procedure over many timesteps



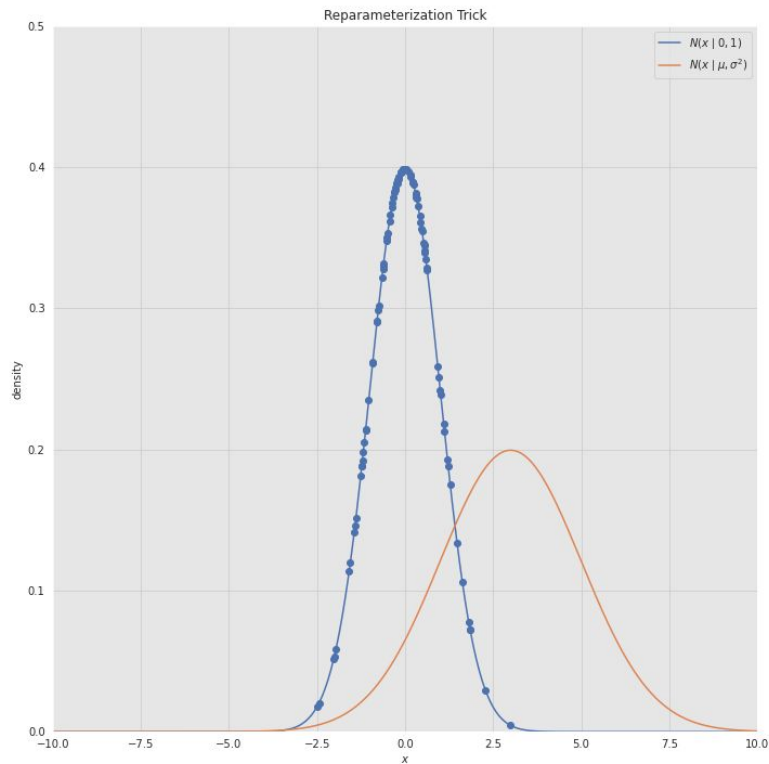
Recap: Variational Autoencoders 🎩

Visually, we often see a VAE as:



How do we perform backpropagation through samples?

Recap: Reparameterization Trick

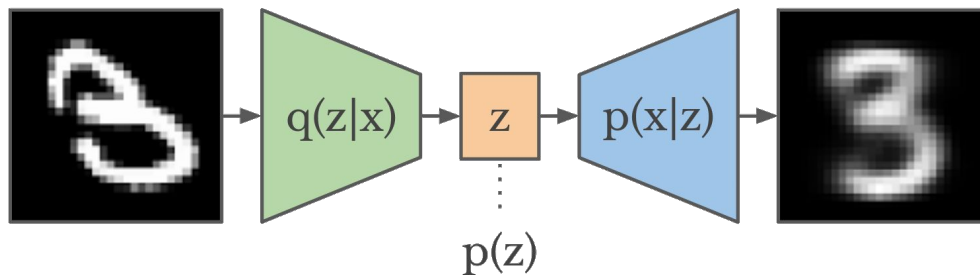


For $x \sim \mathcal{N}(x | \mu, \sigma^2)$,

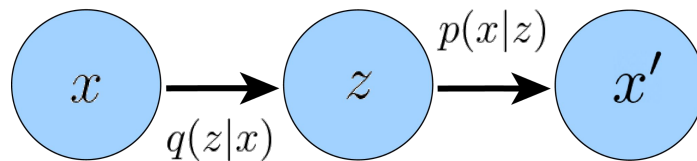
ϵ , where $\epsilon \sim \mathcal{N}(x | 0, I)$

Recap: Variational Autoencoders 🎩

Visually, we often see a VAE as:

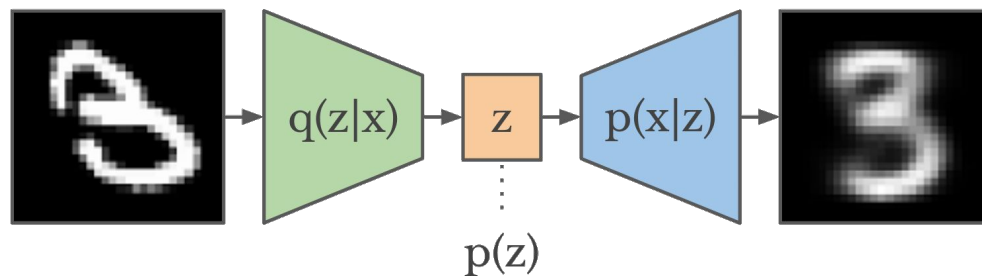


Or as:

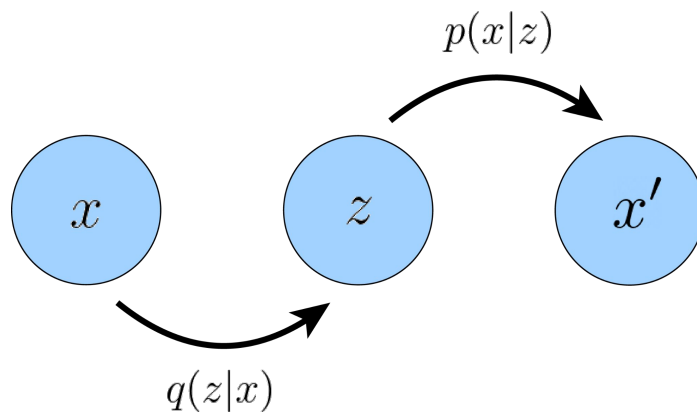


Recap: Variational Autoencoders 🎩

Visually, we often see a VAE as:

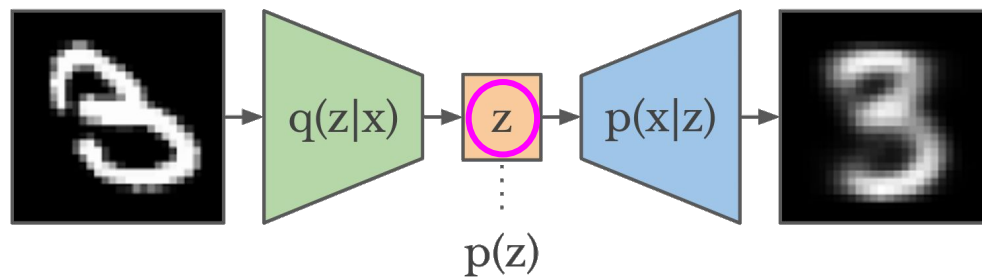


Or as:

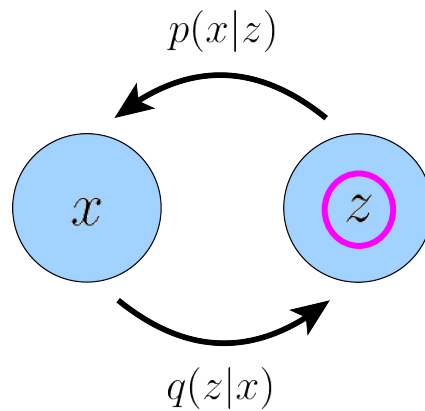


Recap: Variational Autoencoders 🎩

Visually, we often see a VAE as:



Or as:

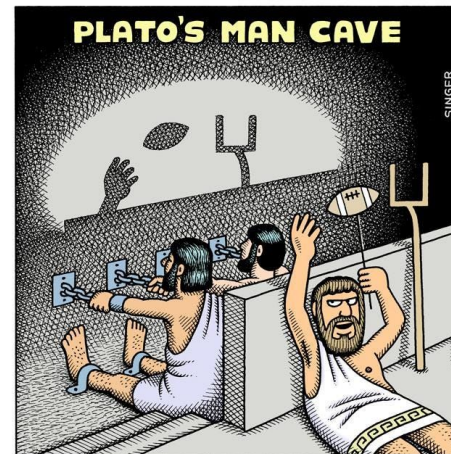


...but what's the intuition behind what is learned?

Generative Modeling with Latent Variables

Given $\mathbf{x} \sim p(\mathbf{x})$ we might want to learn $p_{\theta}(\mathbf{x}) \approx p(\mathbf{x})$ (*modeling*)

What if we assume latent variables \mathbf{z} exist?



Elon has an idea...



Elon has an idea...



Elon has an idea...



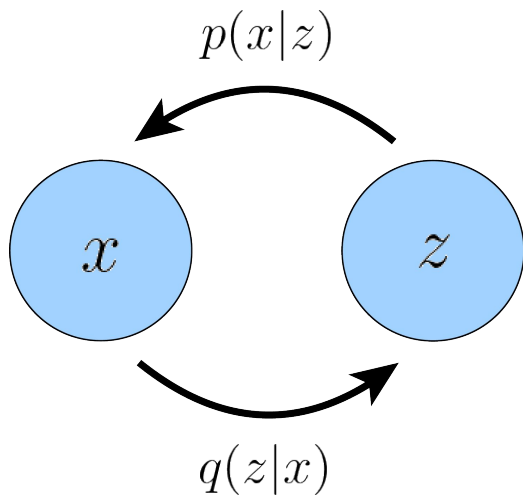
Elon has an idea...



Hierarchical VAEs

Generalize VAEs by enabling a hierarchy of latents $z = z_1, \dots, z_T$

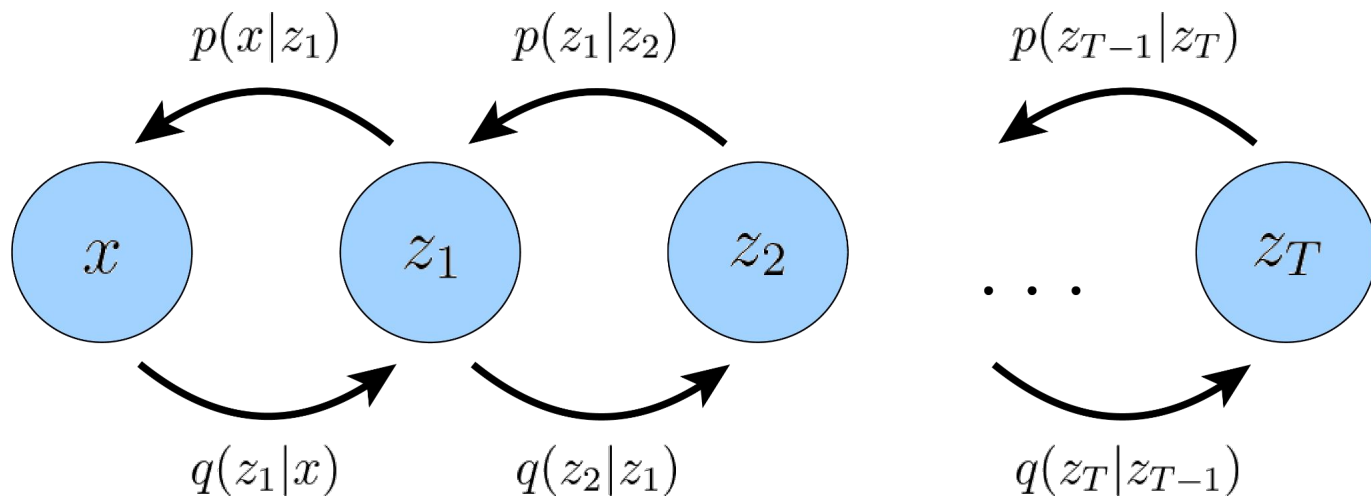
This is essentially learning a bunch of stacked VAEs



Hierarchical VAEs

Generalize VAEs by enabling a hierarchy of latents $z = z_1, \dots, z_T$

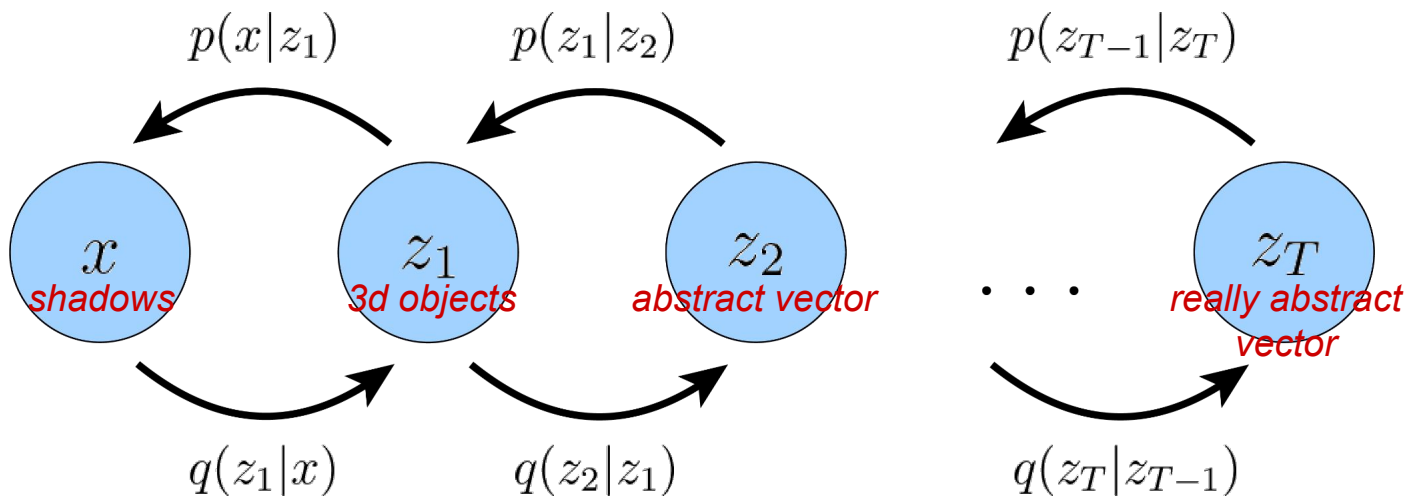
This is essentially learning a bunch of stacked VAEs



Disclaimer: Elon did not actually come up with this idea.

Hierarchical VAEs

Let's think like a caveman...

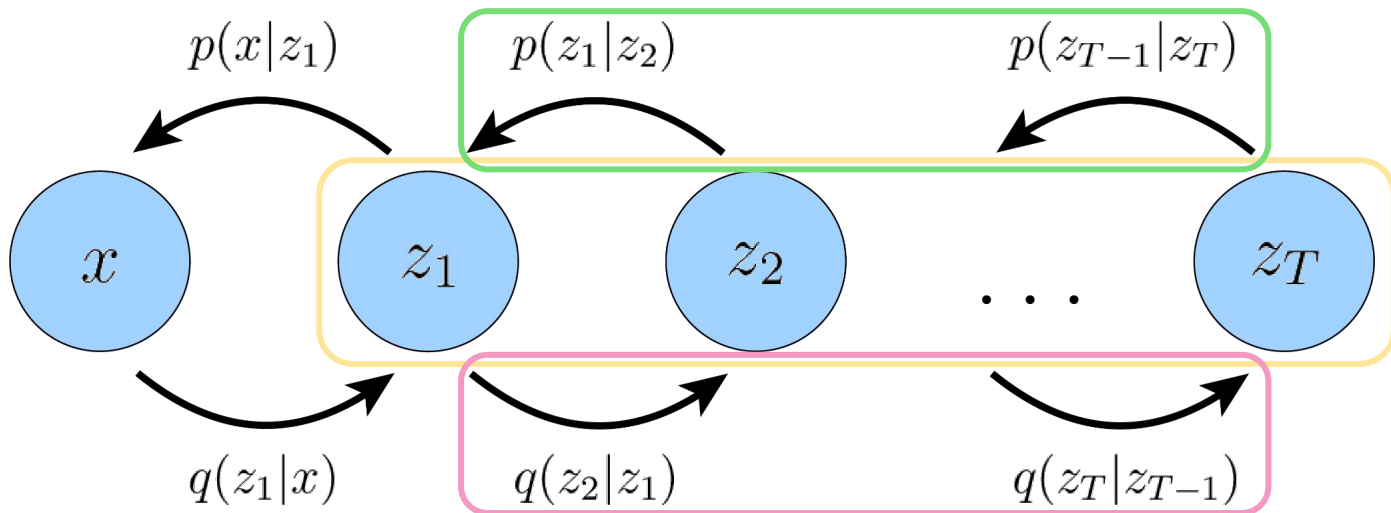


Hierarchical VAEs

Question:

- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

...what if we assume all latent dimensions are the same?

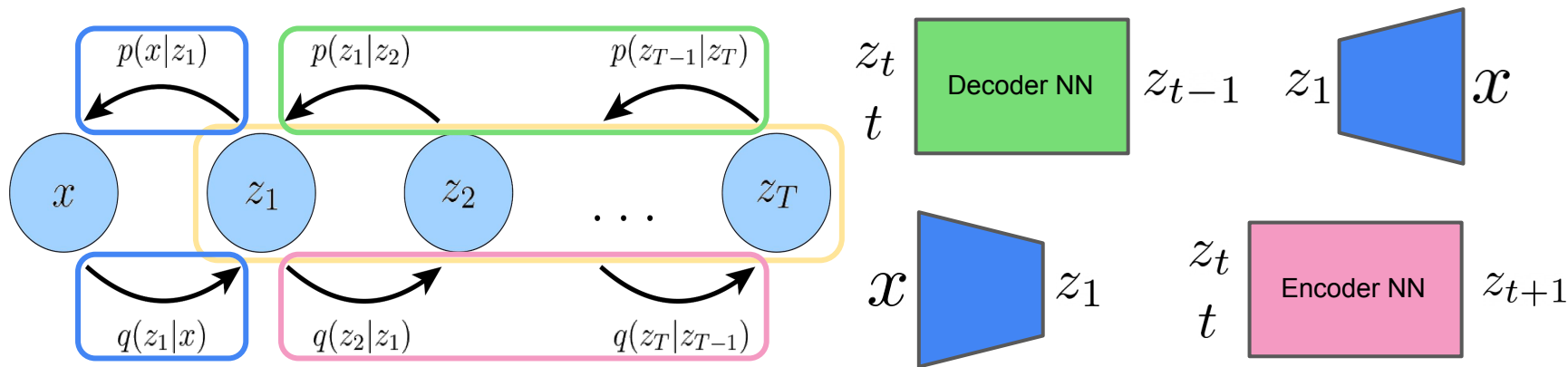


Hierarchical VAEs

Question:

- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

...what if we assume all latent dimensions are the same?

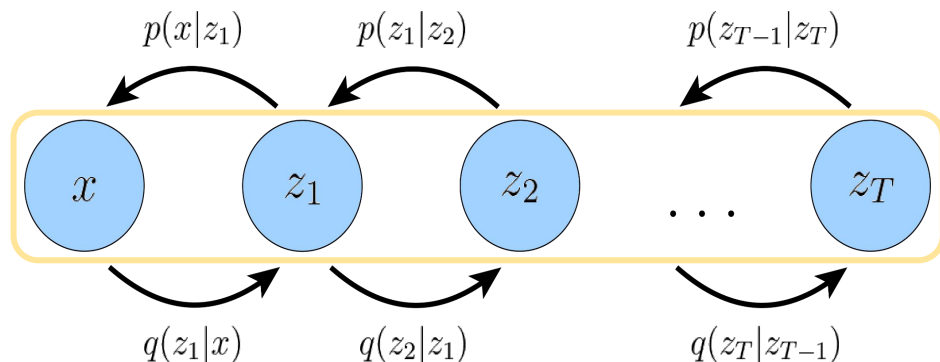


Hierarchical VAEs

Question:

- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

*...what if we assume **all** dimensions are the same?*

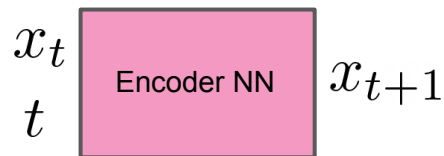
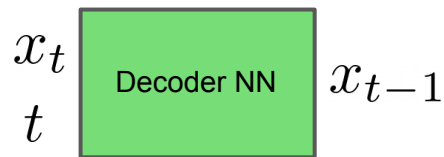
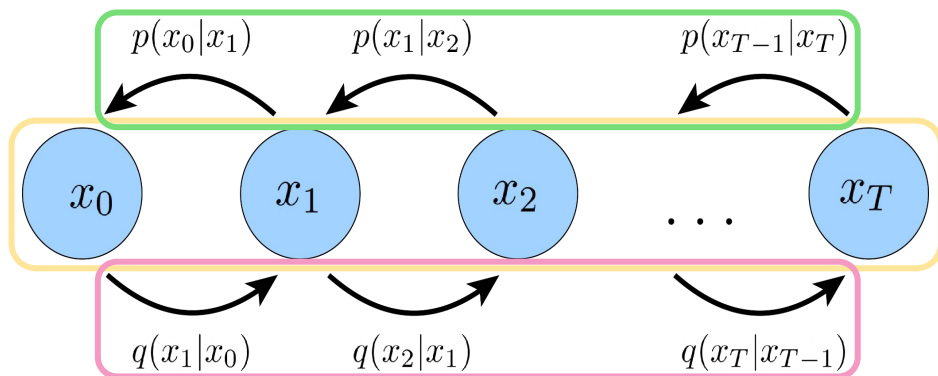


Hierarchical VAEs

Question:

- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

*...what if we assume **all** dimensions are the same?*

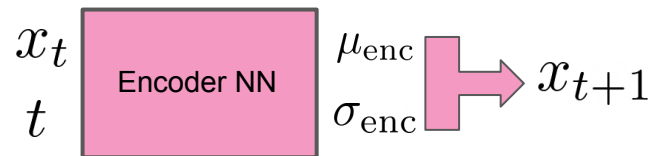
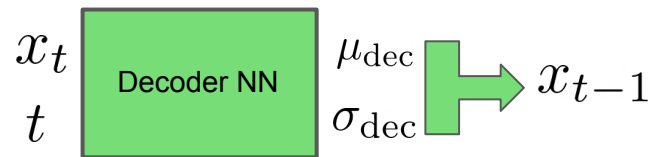
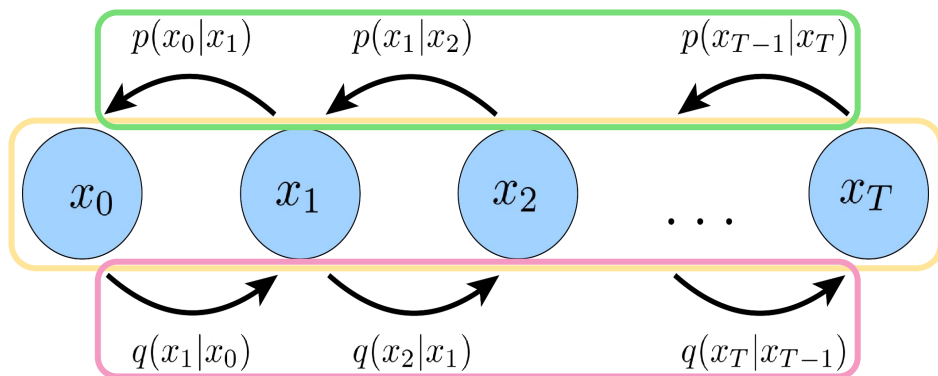


Hierarchical VAEs

Question:

- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

*...what if we assume **all** dimensions are the same?*



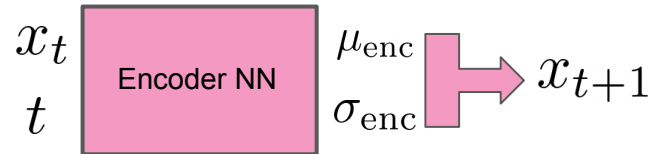
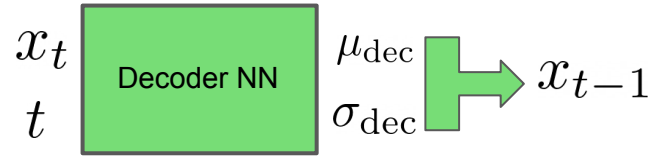
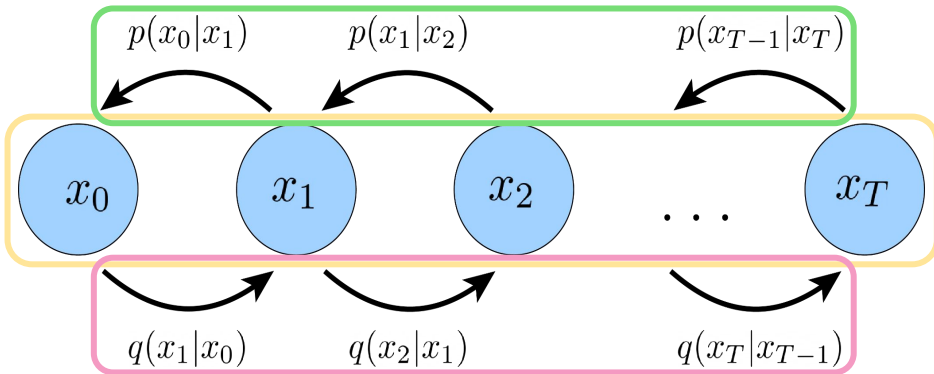
Hierarchical VAEs

Question:

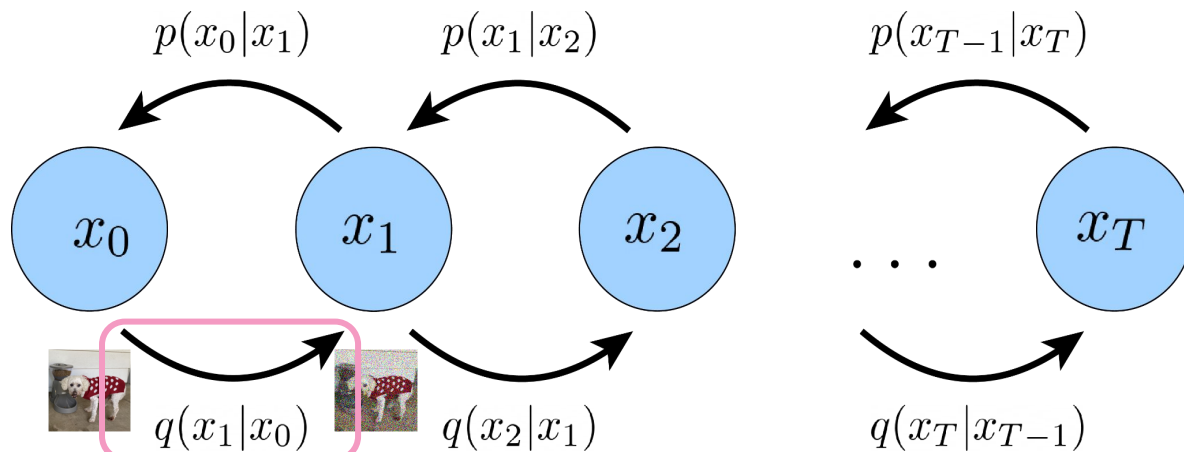
- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

*...what if we assume **all** dimensions are the same?*

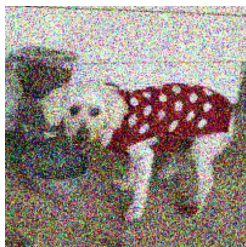
*...what if we assume all encoder transitions are known
Gaussians centered around their previous input?*



Let's take a look at one encoding



$$q(x_1|x_0) = \mathcal{N}(x_1|x_0, \sigma_1^2\mathbf{I})$$



$x_1 \sim q(x_1|x_0)$

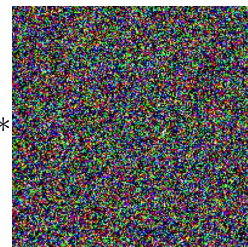
=



x_0

+

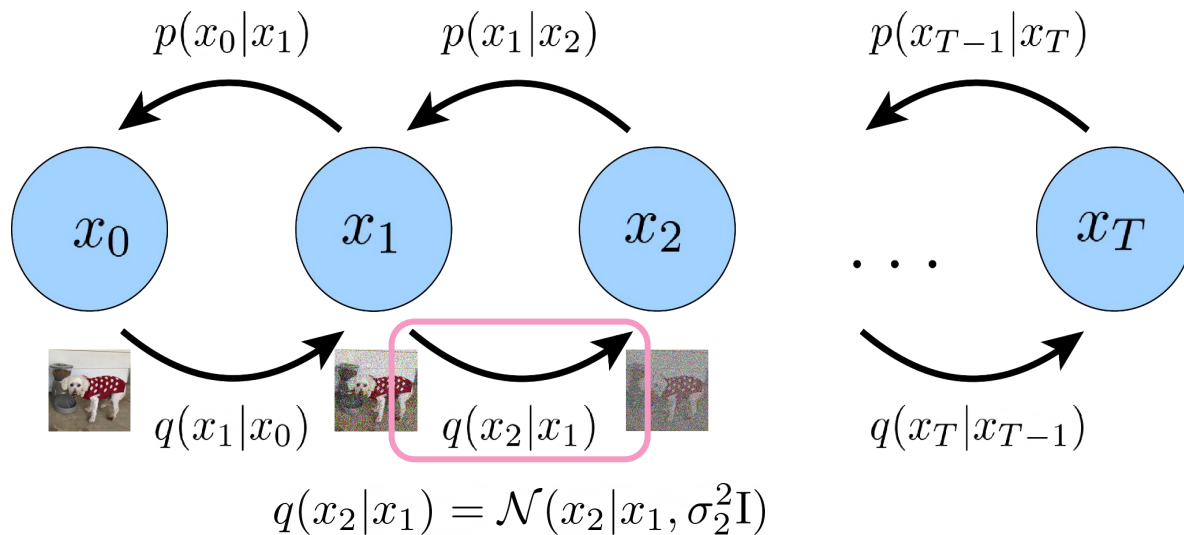
σ_1^*



ϵ_0

reparam. trick!

Let's take a look at one encoding



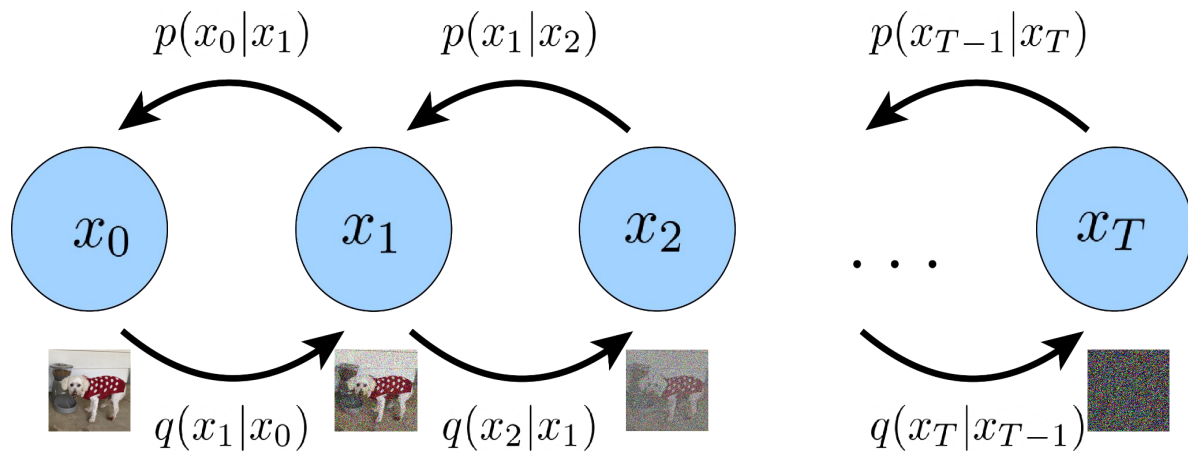
Equation illustrating the reparameterization trick:

$$x_2 \sim q(x_2|x_1) = x_1 + \sigma_2 \epsilon_0$$

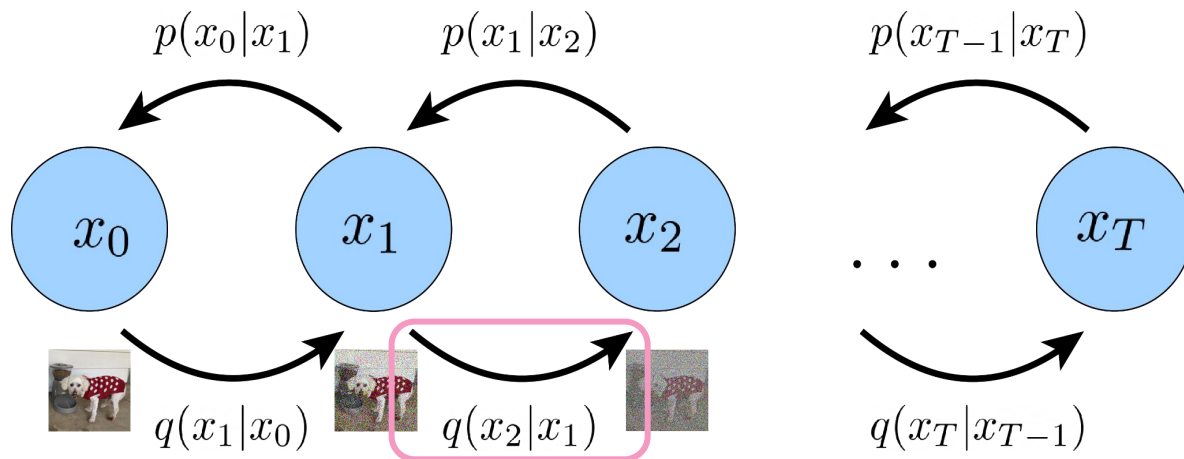
where $x_2 \sim q(x_2|x_1)$ is the noisy image, x_1 is the clear image, and ϵ_0 is the noise vector.

reparam. trick!

Let's take a look at one encoding

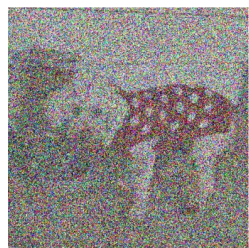
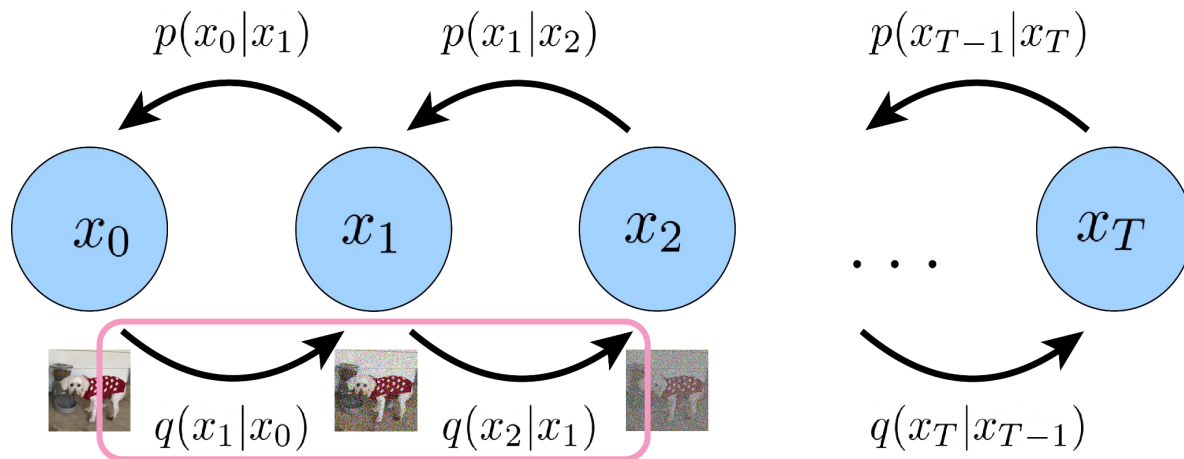


Let's take a look at one encoding



reparam. trick!

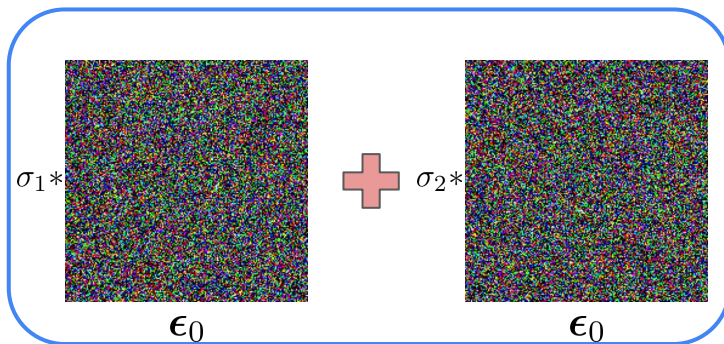
Let's take a look at one encoding



$$x_2 \sim q(x_2|x_1)$$



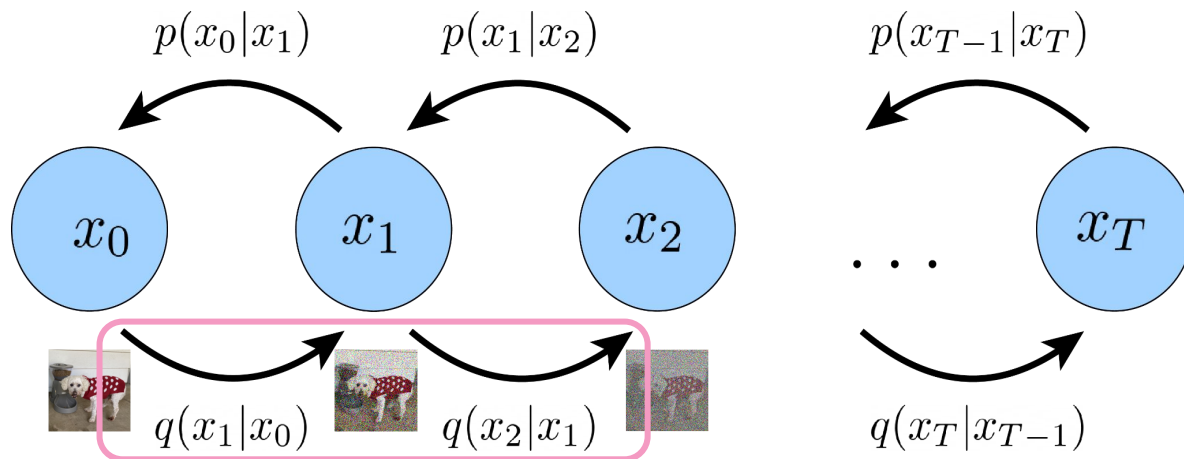
$$x_0$$



Aggregate into 1 sample!

reparam. trick!

Let's take a look at one encoding



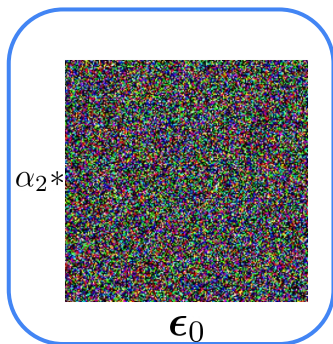
$x_2 \sim q(x_2|x_1)$

=



x_0

+



ϵ_0

where,

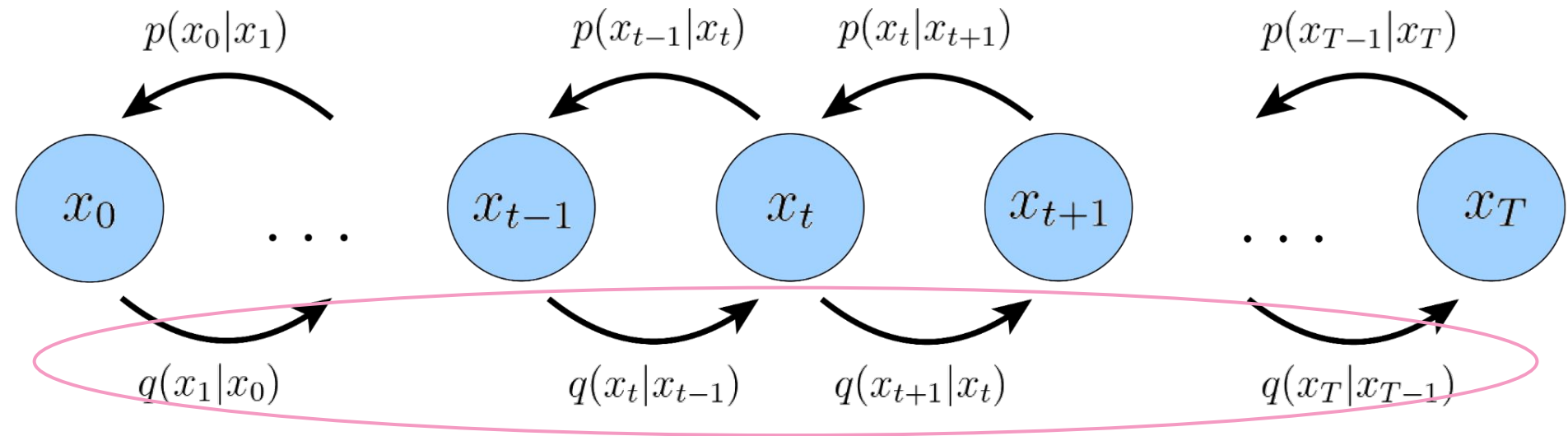
$$\alpha_2 = \sqrt{\sigma_1^2 + \sigma_2^2}$$

and,

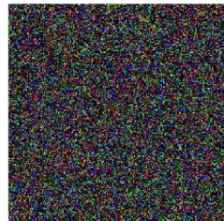
$$q(x_2|x_0) = \mathcal{N}(x_2|x_0, \alpha_2^2)$$

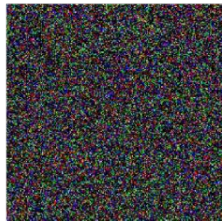
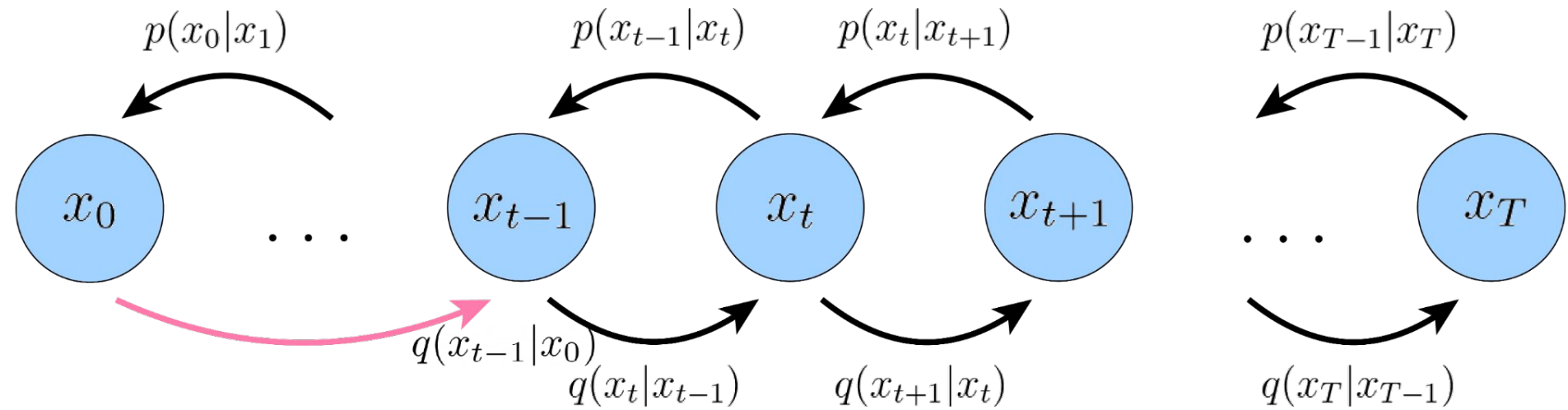
Aggregate into 1 sample!

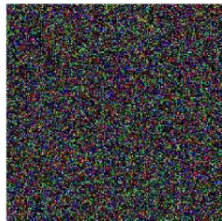
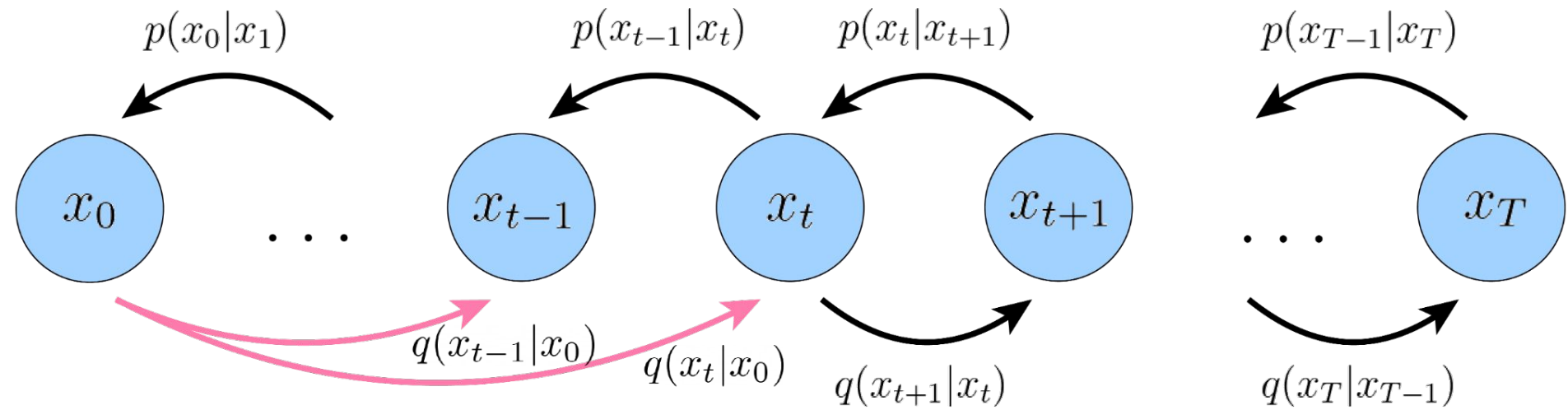
reparam. trick!

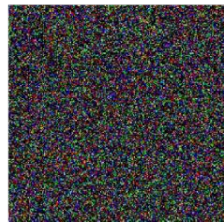
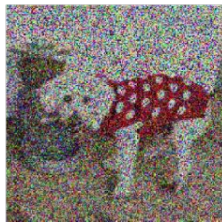
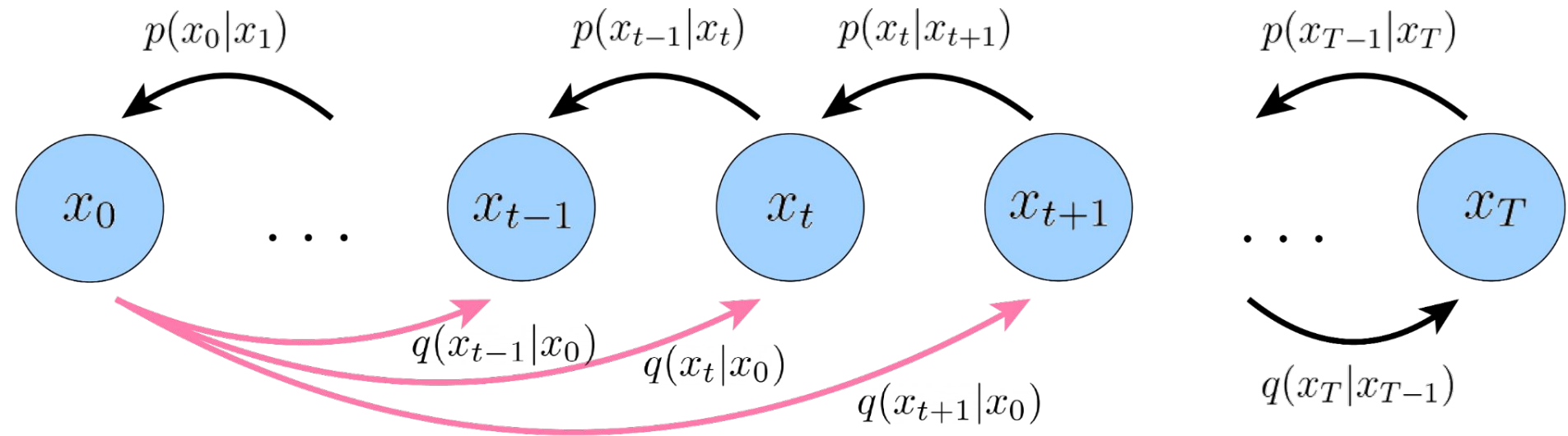


Individual (known) Gaussians!

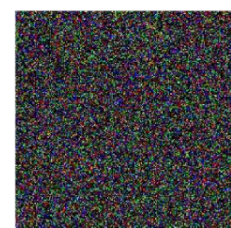
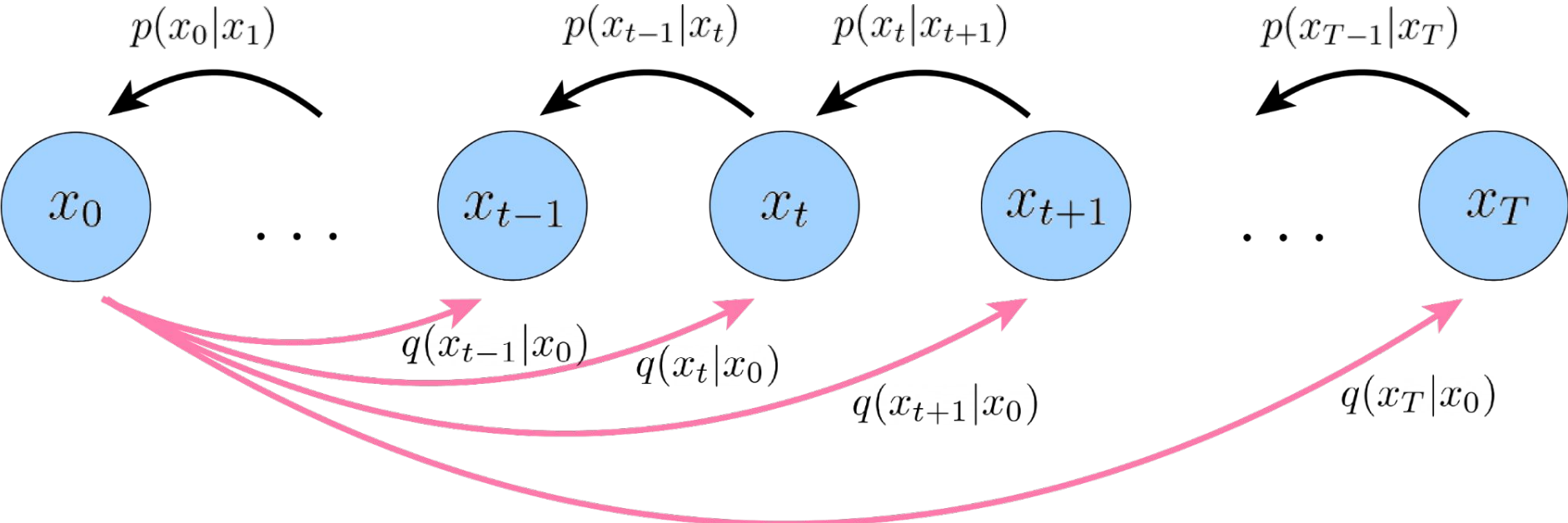








$q(x_t|x_0)$ is a Gaussian, for arbitrary t !
 $q(x_t|x_0) = \mathcal{N}(x_t|x_0, \alpha_t^2 I)$, where $\alpha_0, \alpha_1, \dots, \alpha_T$ are all known/fixed.



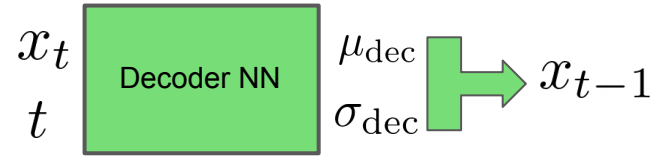
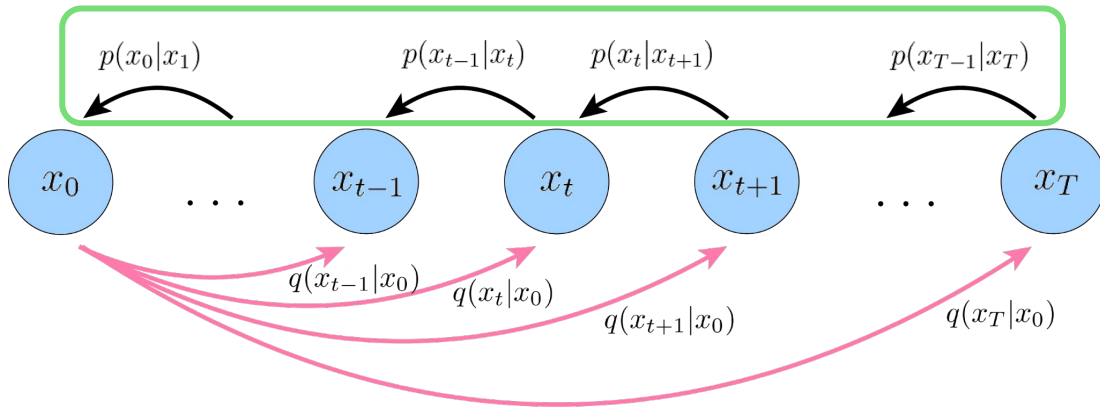
Hierarchical VAEs

Question:

- In a VAE we learn two networks: an encoder and a decoder.
- How many do we need to learn for a Hierarchical VAE?

*...what if we assume **all** dimensions are the same?*

*...what if we assume all encoder transitions are known
Gaussians centered around their previous input?*



...then we can aggregate and simplify the distribution of each intermediate "latent"!

Diffusion Models

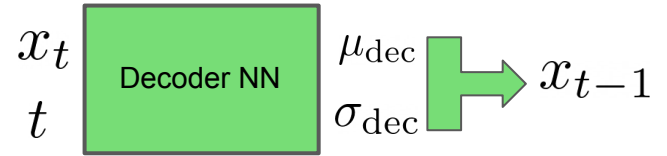
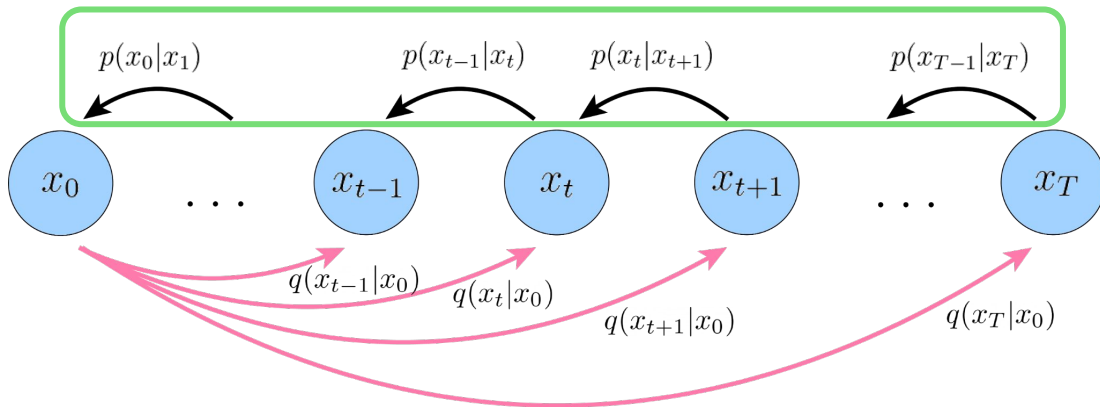
It turns out, that this is exactly what a diffusion model is!

- A Hierarchical VAE with these assumptions:

*...what if we assume **all** dimensions are the same?*

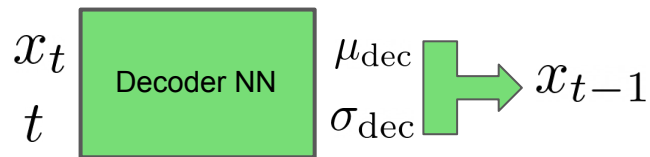
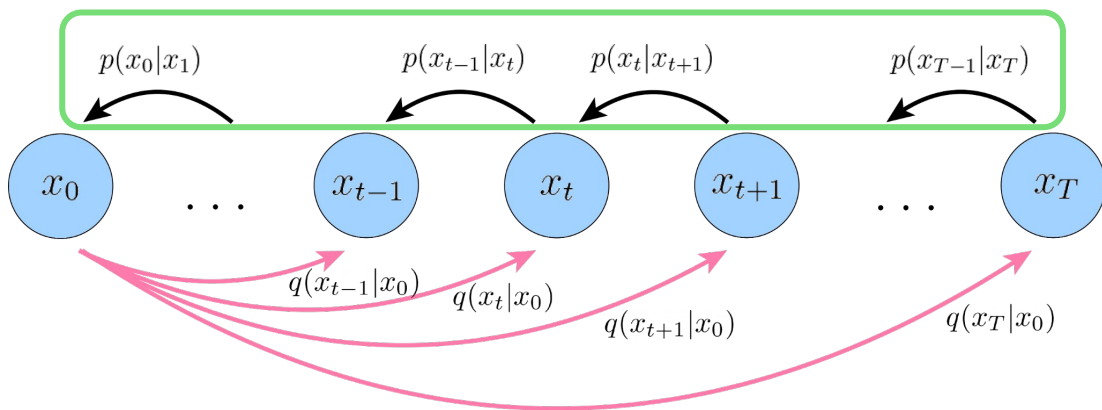
...what if we assume all encoder transitions are known

Gaussians centered around their previous input?



Diffusion Models

A diffusion model is implemented as a single neural network (the decoder)



Optimization?

We want to learn a denoising decoder:

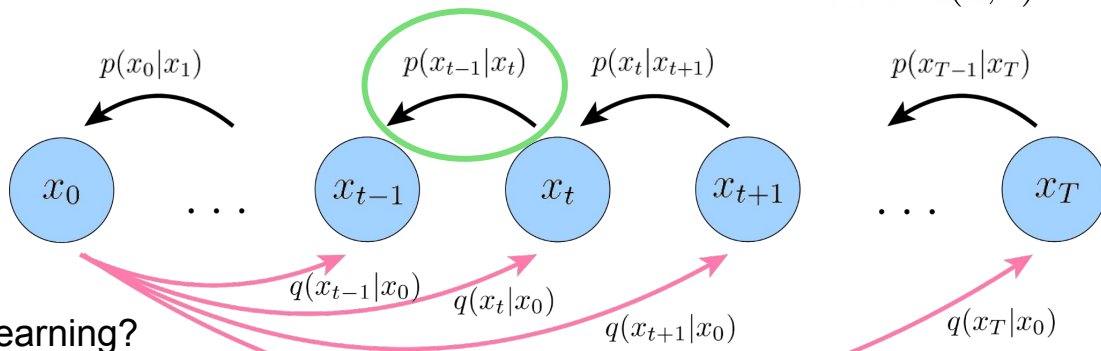


$$\hat{x}_{t-1} = \mu_{\text{dec}} + \sigma_{\text{dec}} * \epsilon$$

reparam. trick!

$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

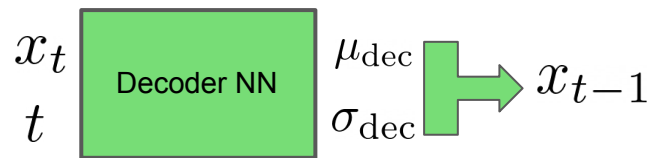
But what is the form of x_{t-1} ?



...can we formulate this as supervised learning?

Optimization?

We want to learn a denoising decoder:



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

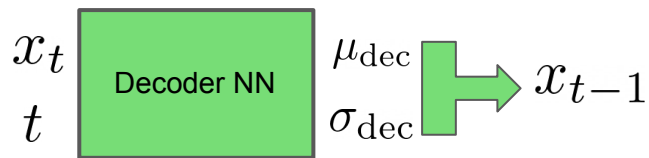
reparam. trick!

$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

But what is the form of x_{t-1} ?

Optimization?

We want to learn a denoising decoder:



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

reparam. trick!

$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

But what is the form of x_{t-1} ?

Recall that:

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

$$\therefore x_{t-1} = x_0 + \alpha_{t-1} * \epsilon$$

reparam. trick!

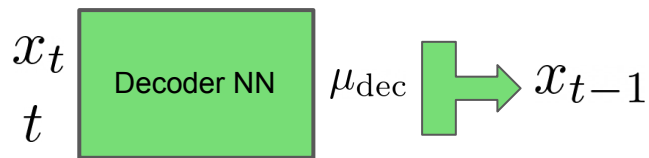
$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Do we really need to predict σ_{dec} ?

What is the ground truth signal for μ_{dec} ?

Optimization?

We want to learn a denoising decoder:



$$\hat{x}_{t-1} = \hat{x}_0 + \alpha_{t-1} * \epsilon$$

reparam. trick!

$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

But what is the form of x_{t-1} ?

Recall that:

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

$$\therefore x_{t-1} = x_0 + \alpha_{t-1} * \epsilon$$

reparam. trick!

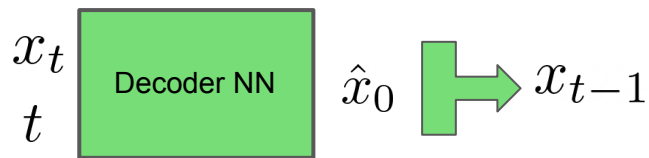
$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Do we really need to predict σ_{dec} ?

What is the ground truth signal for μ_{dec} ?

Optimization?

We want to learn a denoising decoder:



$$\hat{x}_{t-1} = \hat{x}_0 + \alpha_{t-1} * \epsilon$$

reparam. trick!

$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

But what is the form of x_{t-1} ?

Recall that:

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

$$\therefore x_{t-1} = x_0 + \alpha_{t-1} * \epsilon$$

reparam. trick!

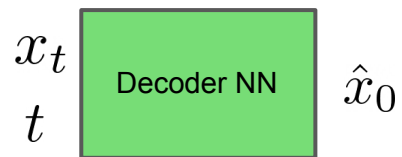
$$\epsilon \sim \mathcal{N}(0, \mathbf{I})$$

Do we really need to predict σ_{dec} ?

What is the ground truth signal for μ_{dec} ?

Optimization?

We want to learn a denoising decoder:

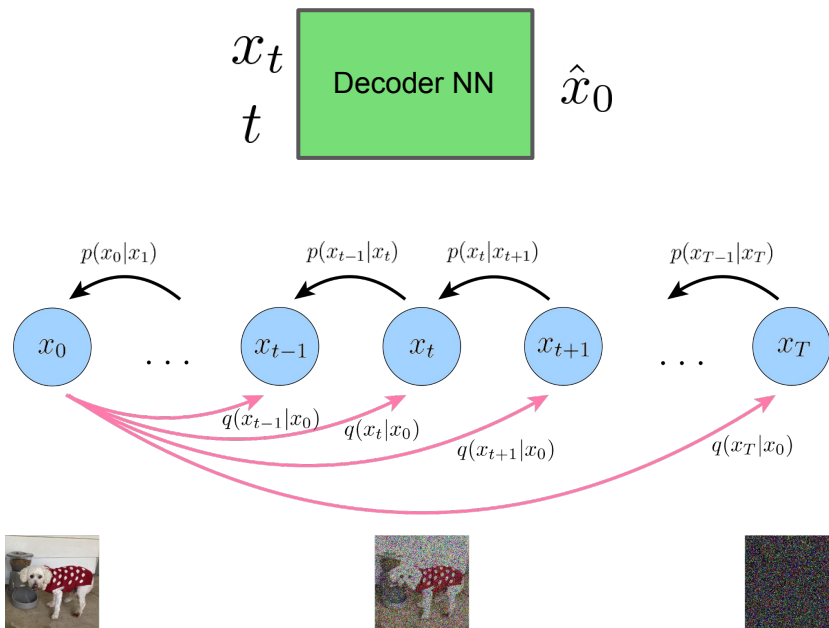


So in the end, a diffusion model is simply *one* Neural Network that predicts a clean image x_0 from arbitrary noisified image x_t .

Diffusion Models: A Summary

A Diffusion Model is:

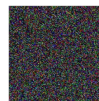
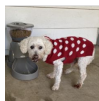
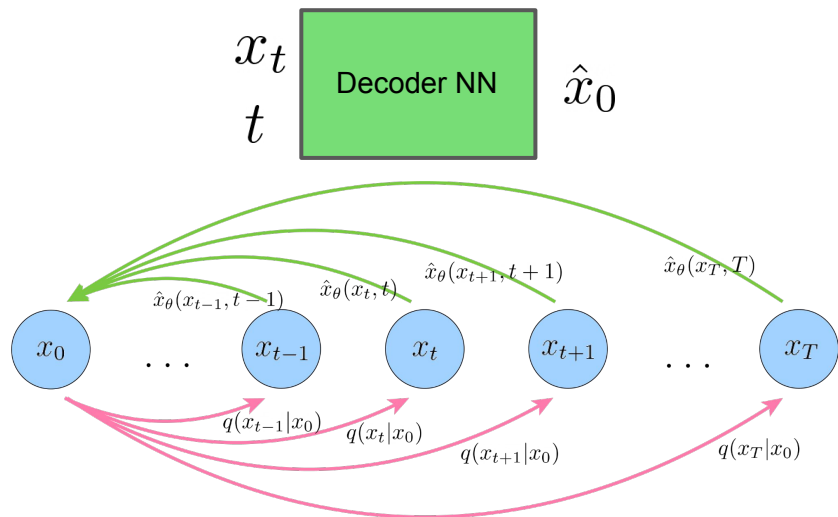
- One NN that predicts a clean image from a noisy version of the image



Diffusion Models: A Summary

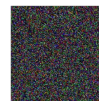
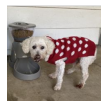
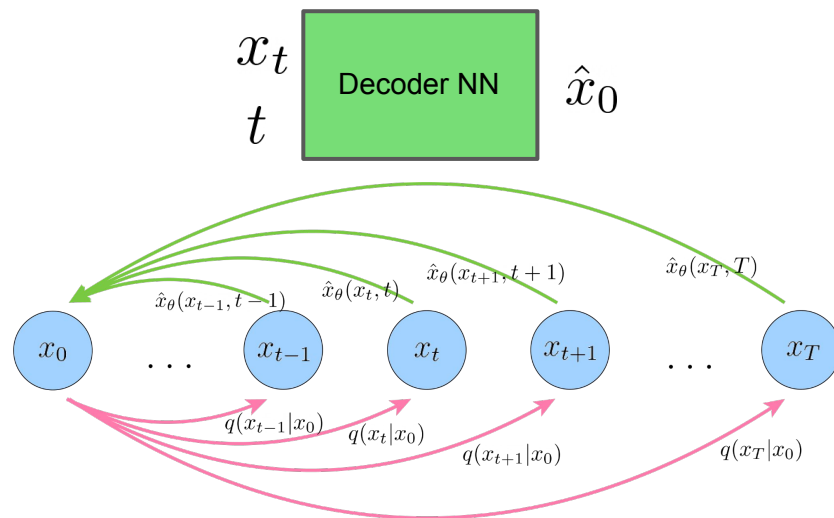
A Diffusion Model is:

- One NN that predicts a clean image from a noisy version of the image

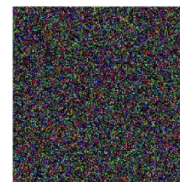
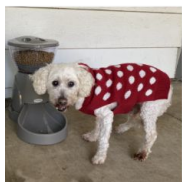
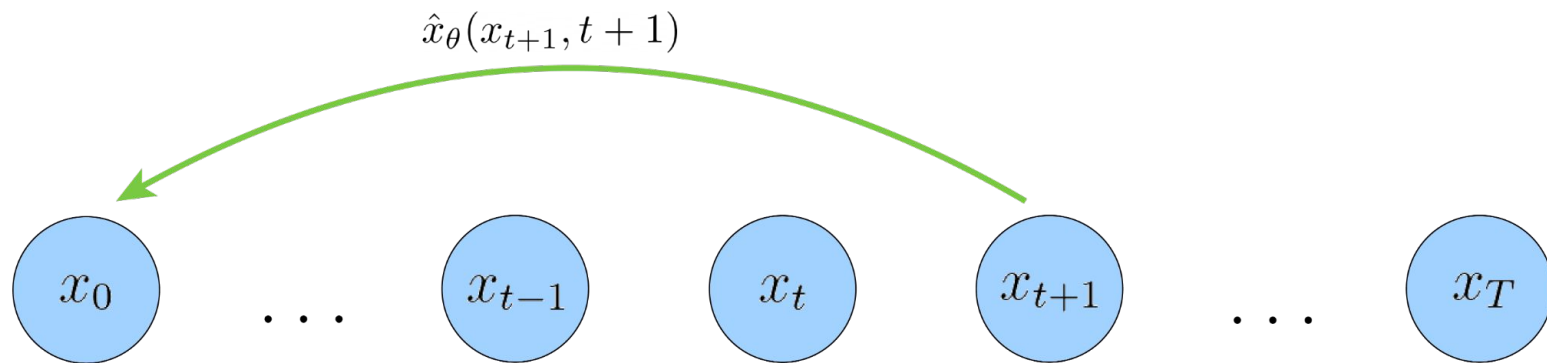


Diffusion Models: A Summary

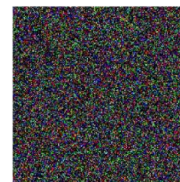
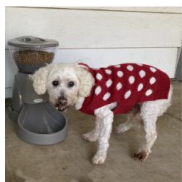
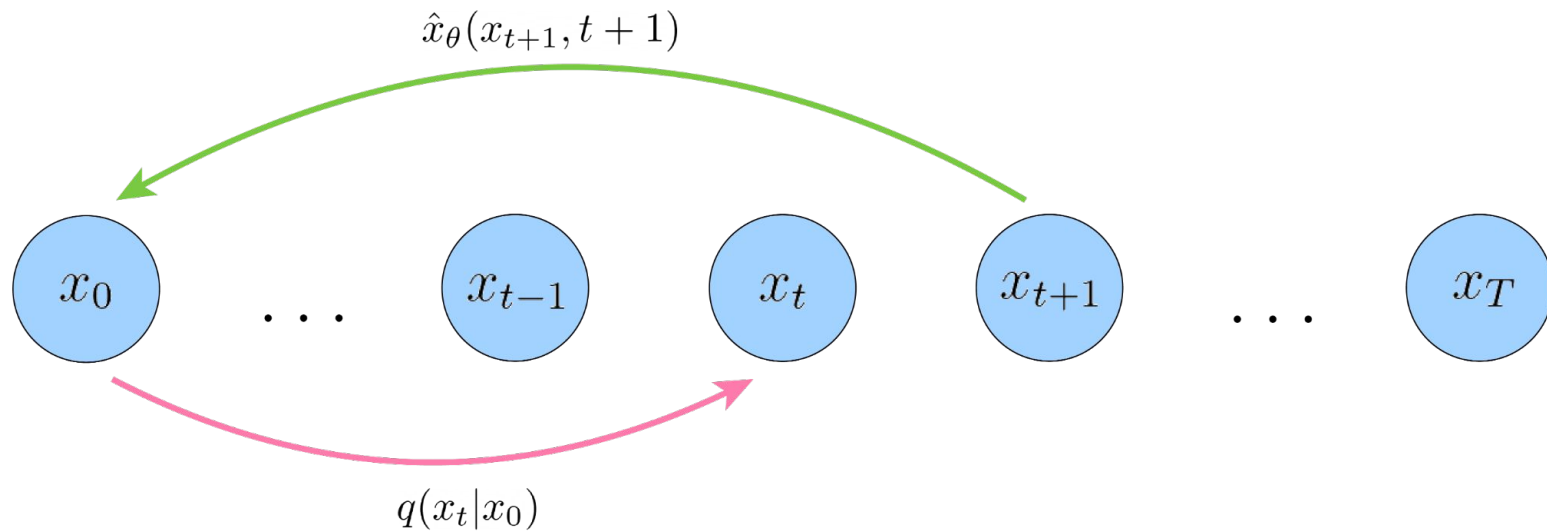
How do we perform sampling?



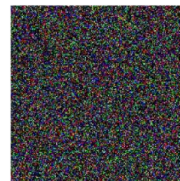
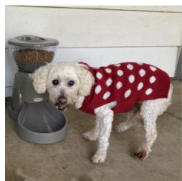
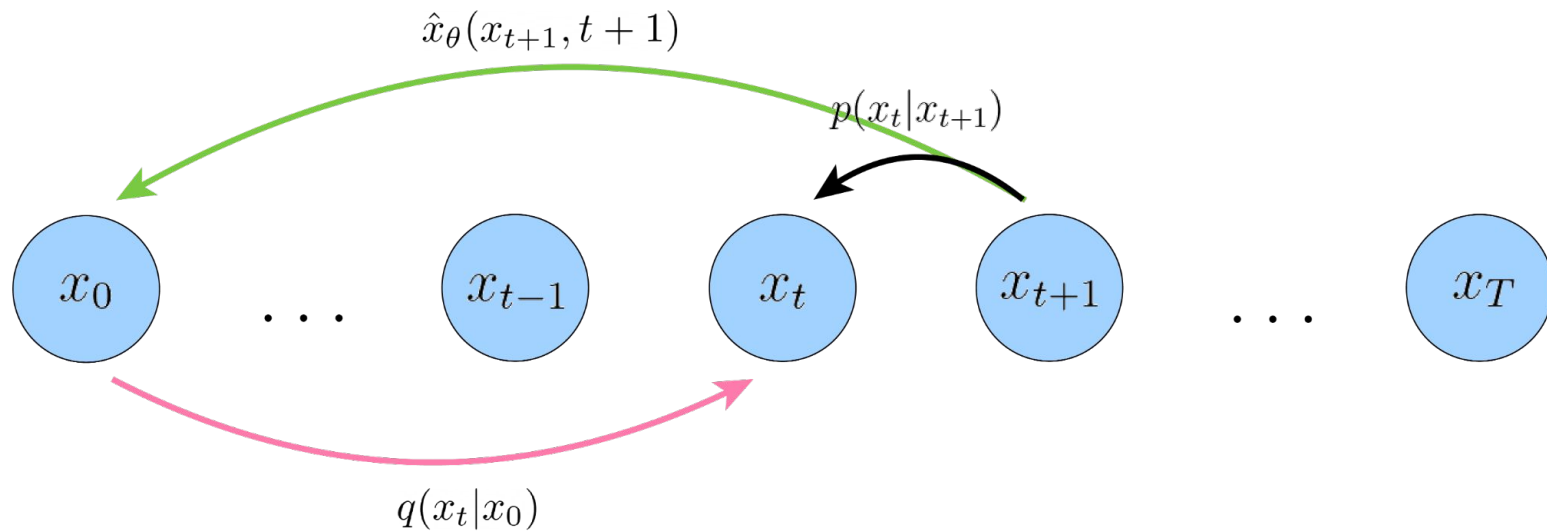
Sampling



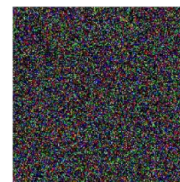
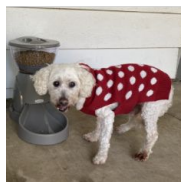
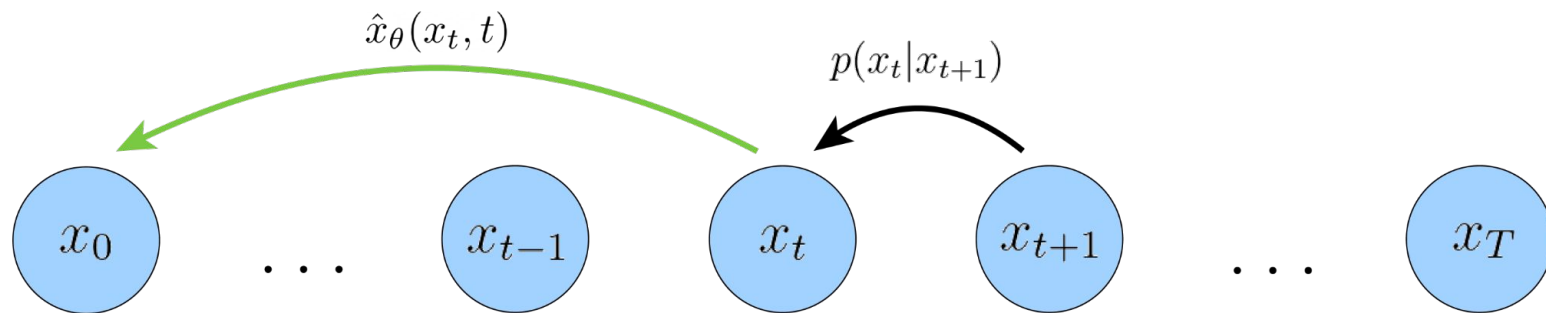
Sampling



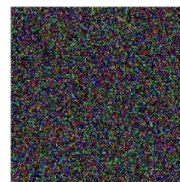
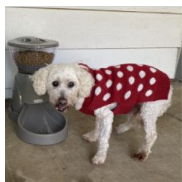
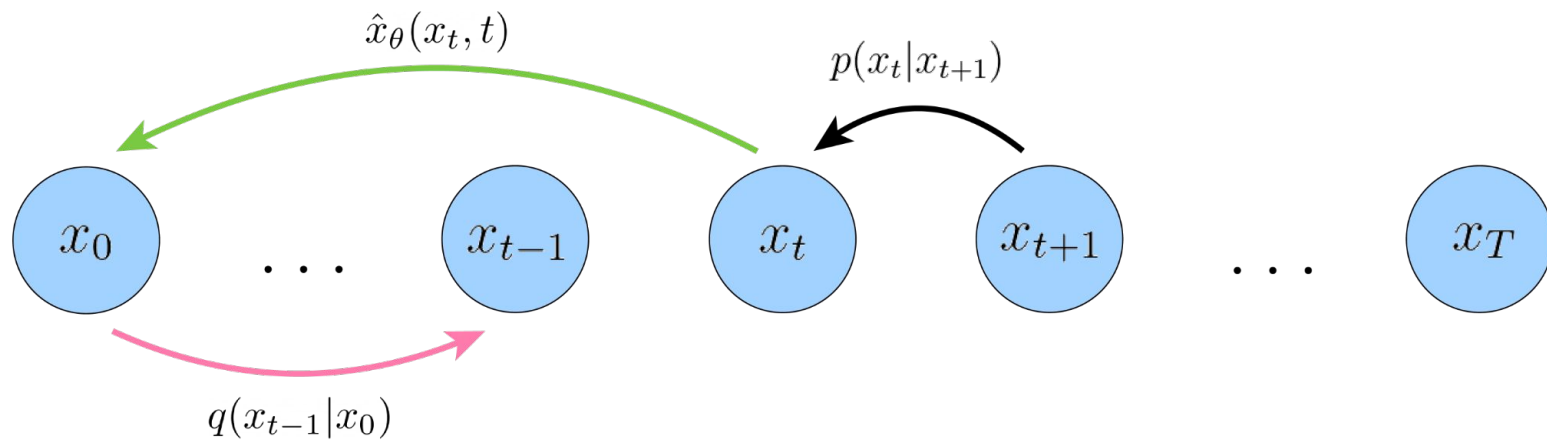
Sampling



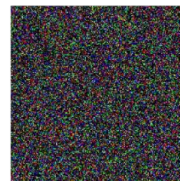
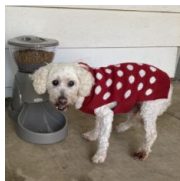
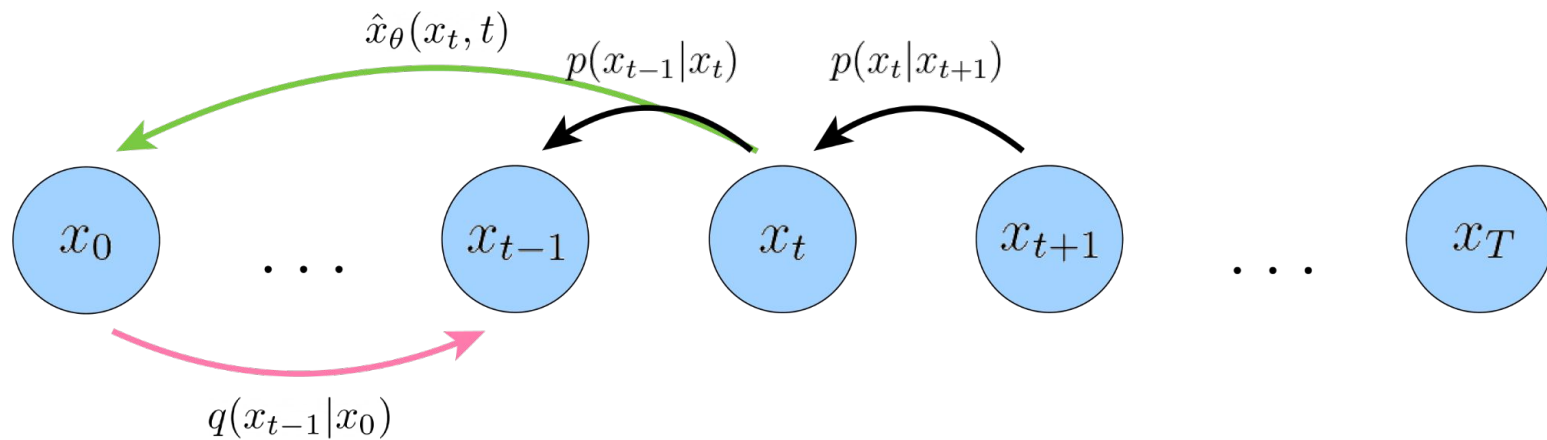
Sampling



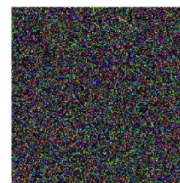
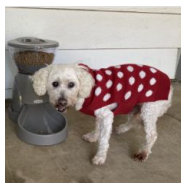
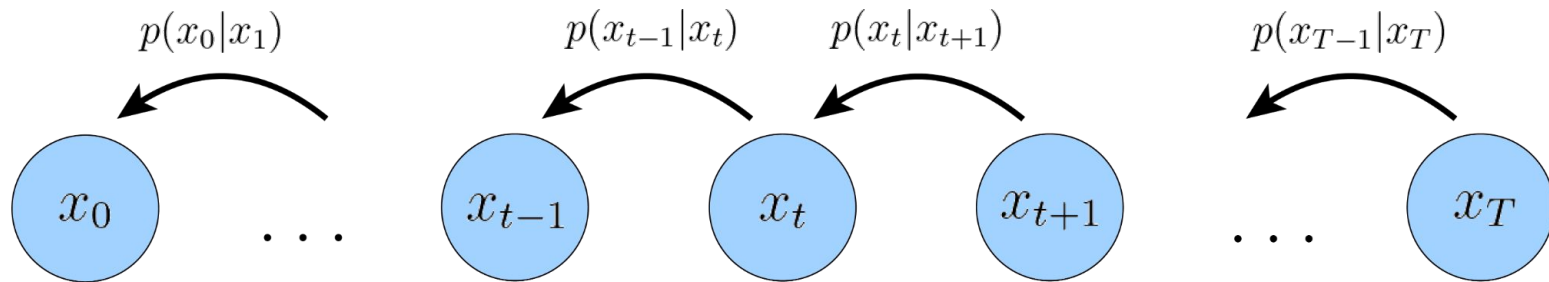
Sampling



Sampling



Sampling



Pseudocode

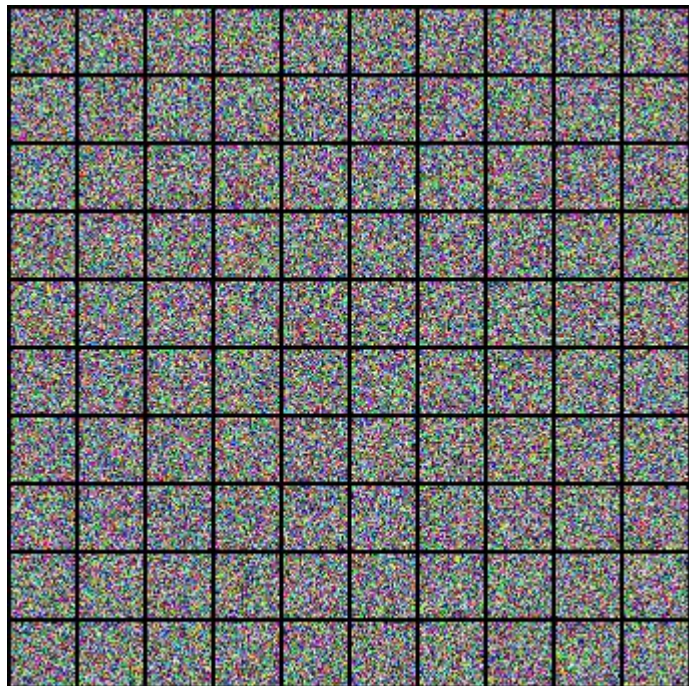
Algorithm 1 Training

```
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  
6:      $\nabla_{\theta} \|\mathbf{x}_0 - \hat{\mathbf{x}}_{\theta}(\mathbf{x}_0 + \alpha_t \boldsymbol{\epsilon}, t)\|^2$   
7: until converged
```

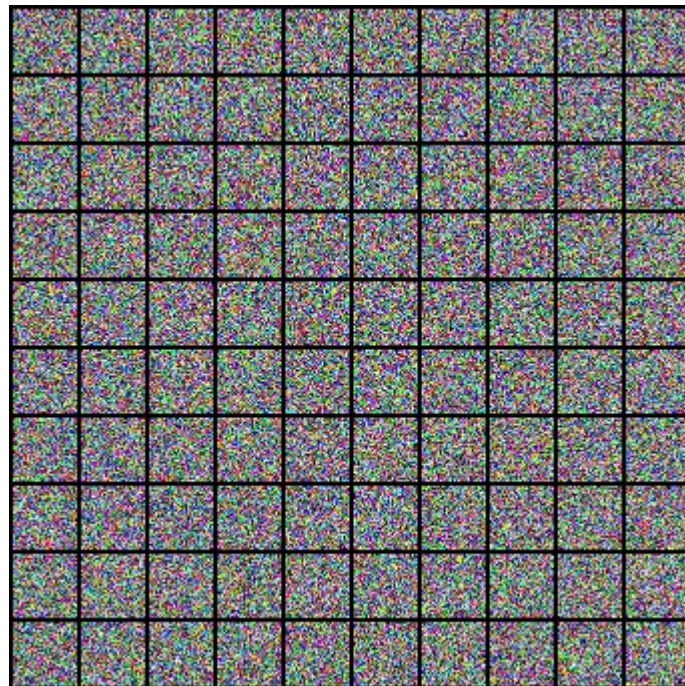
Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
2: for  $t = T, \dots, 1$ :  
3:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\boldsymbol{\epsilon} = \mathbf{0}$   
4:    $\mathbf{x}_{t-1} = \hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) + \alpha_{t-1} \boldsymbol{\epsilon}$   
5: end for  
6: return  $\mathbf{x}_0$ 
```

Examples!

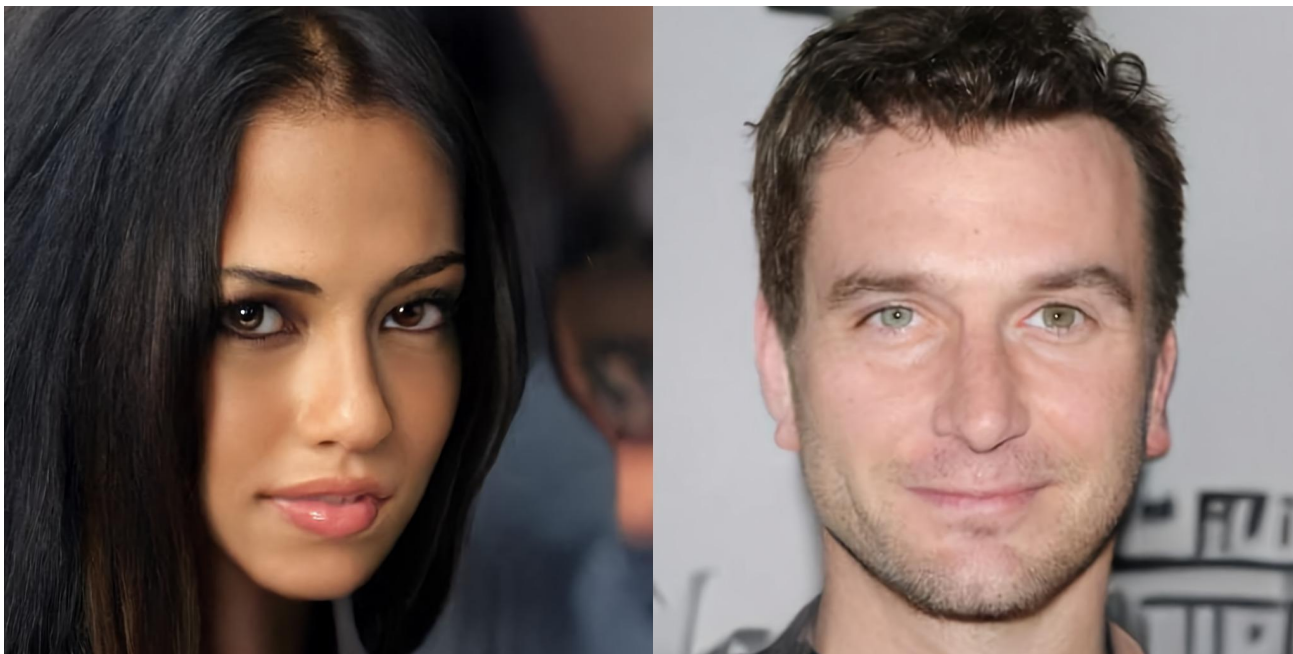


Celeb-A



CIFAR-10



Examples!



1024x1024 samples

Three Different Interpretations

It turns out, training a DiffModel can be done using three different interpretations:

- Predicting original image  (we just did this)
- Predicting noise  (coming up!)
- Predicting score function 100 (coming up!)

Diffusion Models as a Noise Predictor

Recall that our objective is to predict $\hat{x}_\theta(x_t, t) \approx x_0$

Image  and Noise ? They are the same!

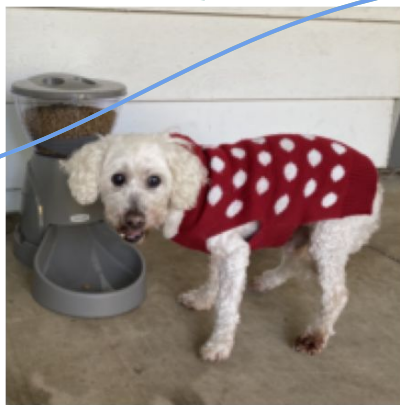
What does it mean intuitively?

For arbitrary $\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)$, we can rewrite it as $\mathbf{x}_t = \mathbf{x}_0 + \alpha_t \epsilon_0$

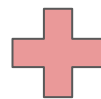
Predicting \mathbf{x}_0 determines ϵ_0 and vice-versa, since they sum to the same thing!



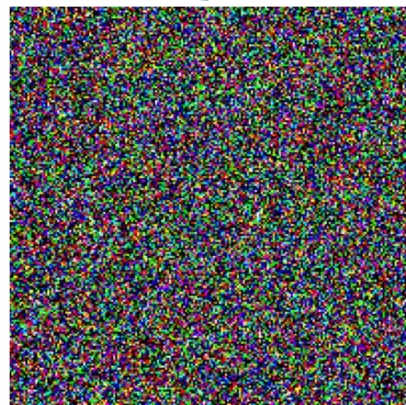
$\mathbf{x}_t \sim q(\mathbf{x}_t | \mathbf{x}_0)$



\mathbf{x}_0



$\alpha_t *$



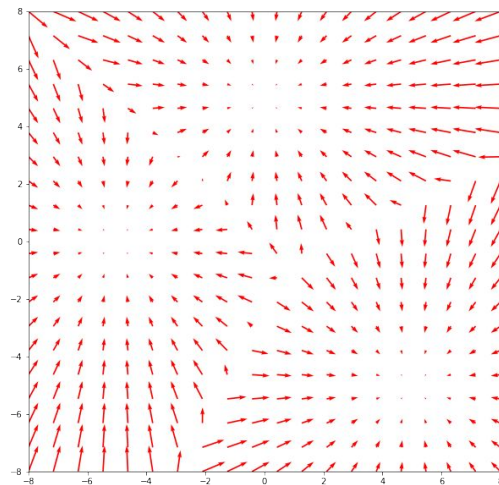
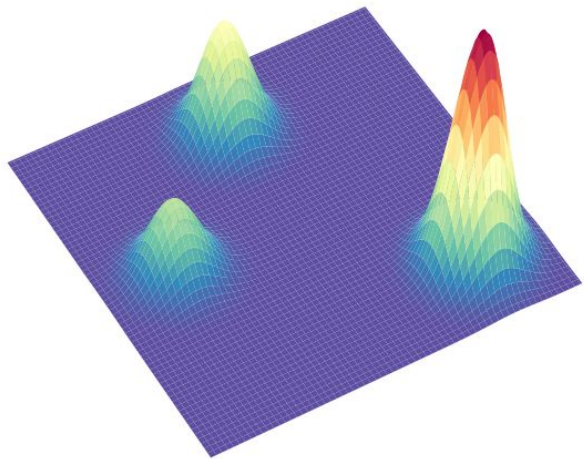
ϵ_0

Score Functions 700

What are score functions?

$$\nabla_{\mathbf{x}} \log p(\mathbf{x})$$

Intuitively, they describe how to move in data space to improve the (log) likelihood.



Tweedie's Formula

Mathematically, for a Gaussian variable $z \sim \mathcal{N}(z; \mu_z, \Sigma_z)$ Tweedie's formula states:

$$\mathbb{E}[\mu_z \mid z] = z + \Sigma_z \nabla_z \log p(z)$$

Then, since we have previously shown that:

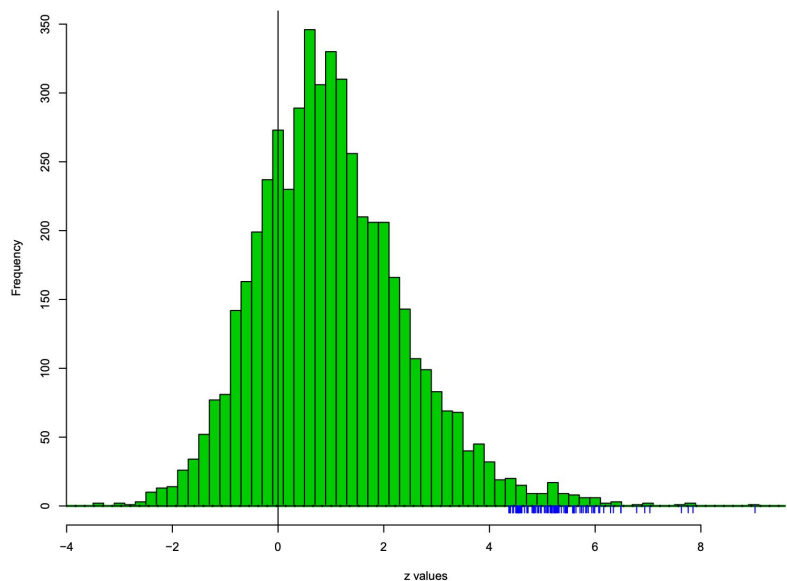
$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_0, \alpha_t^2 \mathbf{I})$$

By Tweedie's Formula, we derive:

$$\mathbb{E}[\mu_{\mathbf{x}_t} \mid \mathbf{x}_t] = \mathbf{x}_t + \alpha_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

The best estimate for the true mean is $\mu_{\mathbf{x}_t} = \mathbf{x}_0$

$$\mathbf{x}_0 = \mathbf{x}_t + \alpha_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$



Tweedie's Formula

There exists a mathematical formula that states that:

$$\mathbf{x}_0 \approx \mathbf{x}_t + \alpha_t^2 \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

Due to the fact that the distribution is Gaussian:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \mathbf{x}_0, \alpha_t^2 \mathbf{I})$$



Diffusion Models as a Score Predictor 700

Recall that our objective is to predict $\hat{x}_\theta(x_t, t) \approx x_0$

Score 100 and Noise ?



There is a relationship between the score and the noise, which we can derive by equating Tweedie's formula with the Reparameterization Trick.

$$\begin{aligned}\mathbf{x}_0 &= \mathbf{x}_t + \alpha_t^2 \nabla \log p(\mathbf{x}_t) = \mathbf{x}_t - \alpha_t \boldsymbol{\epsilon}_0 \\ \therefore \alpha_t^2 \nabla \log p(\mathbf{x}_t) &= -\alpha_t \boldsymbol{\epsilon}_0 \\ \nabla \log p(\mathbf{x}_t) &= -\frac{1}{\alpha_t} \boldsymbol{\epsilon}_0\end{aligned}$$

Intuitively, the direction to move in data space towards a natural image is the negative noise term that was added.

Three Different Interpretations

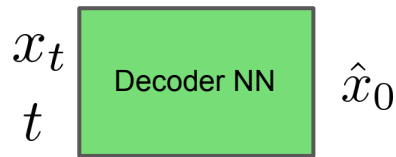
It turns out, training a DiffModel can be implemented as a neural net that:

-  Predicts original image $\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) \approx \mathbf{x}_0$
-  Predicts noise epsilon $\hat{\epsilon}_{\theta}(\mathbf{x}_t, t) \approx \epsilon_0$
? ? ?
- 100 Predicts score function $\mathbf{s}_{\theta}(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$

A Summary

We have learned that a diffusion model is simply one neural network that predicts a clean image from a noisy image.

Objective: $\arg \min_{\theta} \left\| x_0 - \hat{x}_{\theta}(x_t, t) \right\|^2$



Sampling:

