

CS492D: Diffusion Models and Their Applications

Denoising Diffusion Probabilistic Models 1

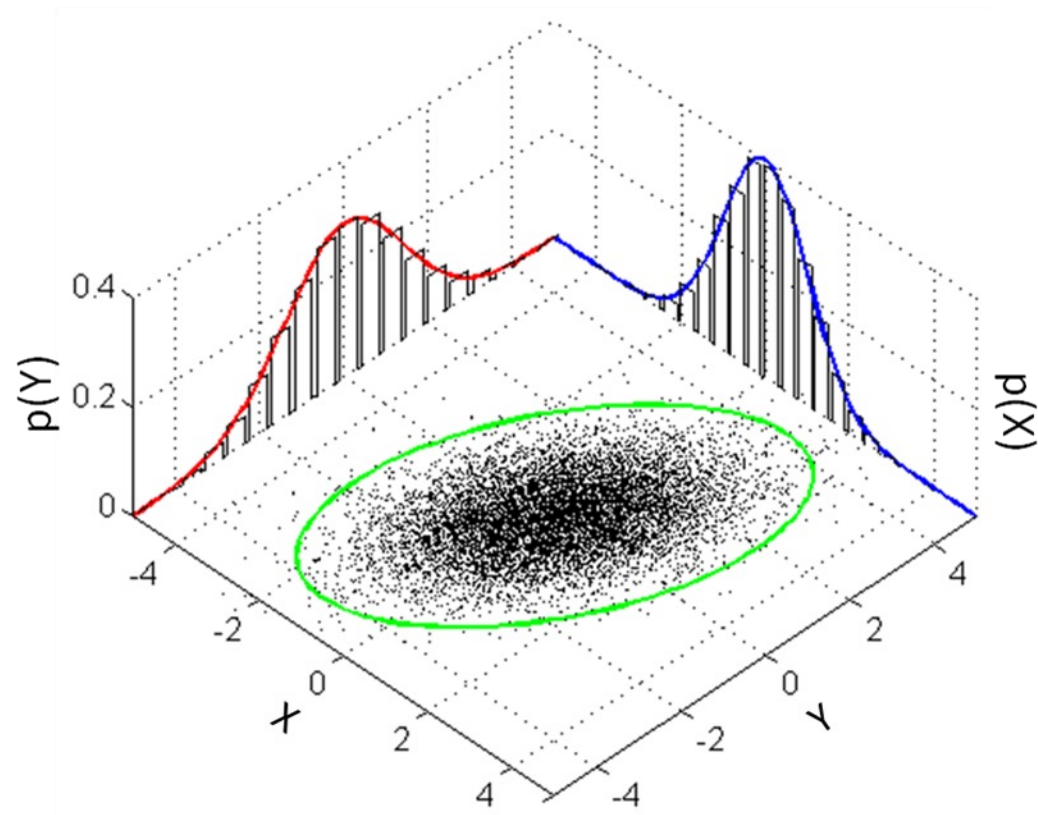
LECTURE 3
MINHYUK SUNG

Fall 2024
KAIST

Previously in CS492D

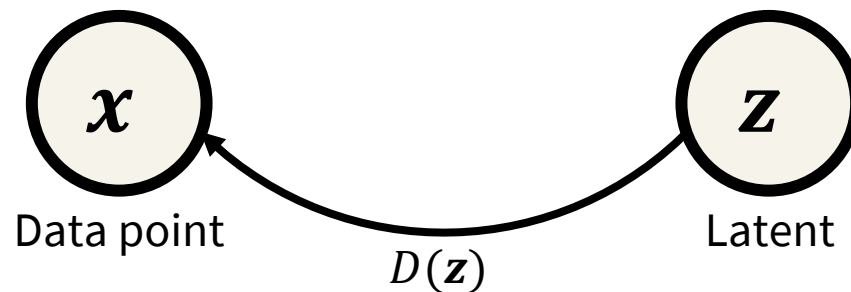
From a **statistical perspective**, we will view a dataset as

- there being a **probability distribution** of the data, and
- the given points are **samples** from the probability distribution.



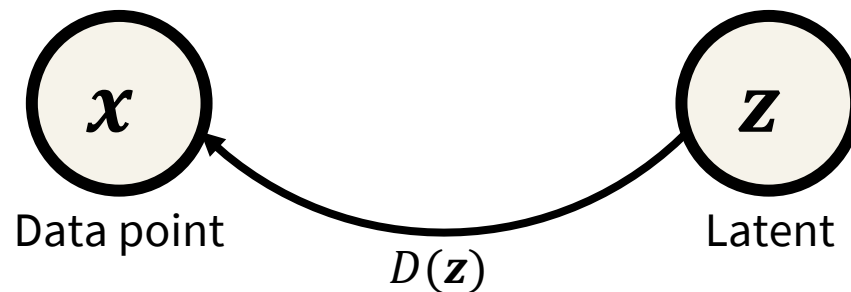
Previously in CS492D

- **Map** a *simple* distribution $p(\mathbf{z})$ (e.g., a standard normal distribution $\mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$) to the **data distribution** $p(\mathbf{x})$.
 - \mathbf{z} : **Latent** variable
 - $p(\mathbf{z})$: **Latent** distribution
- Sample from $p(\mathbf{z})$ and map it to a data point.



Previously in CS492D

- How to map a latent distribution $p(\mathbf{z})$ to the data distribution $p(\mathbf{x})$ using a **neural network**?
- How to **guarantee** that a latent is mapped to a data point of the data distribution?
- We need an **additional** neural network.

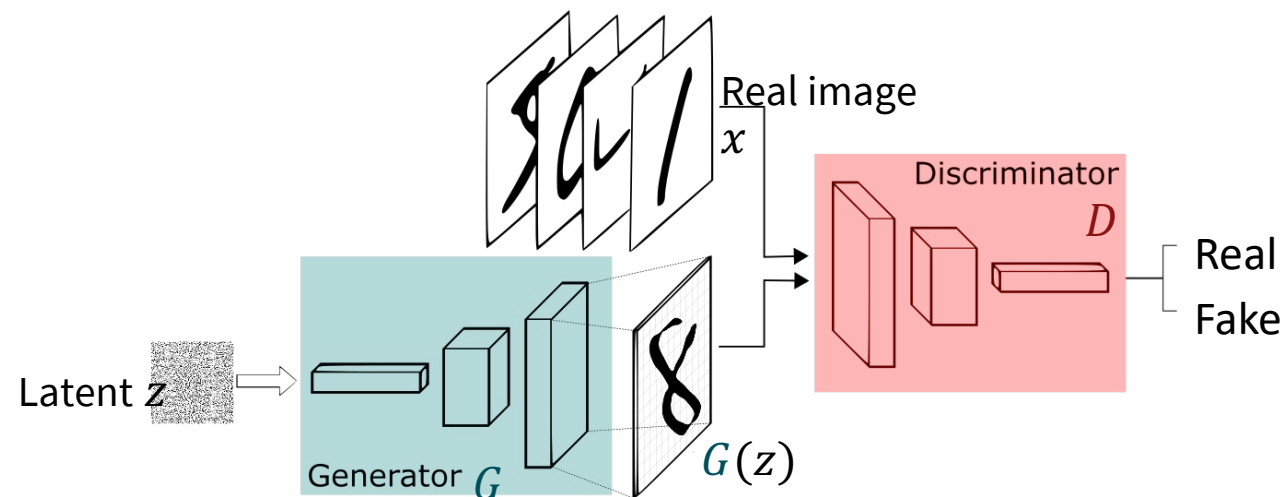


Previously in CS492D

Generative Adversarial Network (GAN)

Introduce a **discriminator** and have it compete with the generator (decoder).

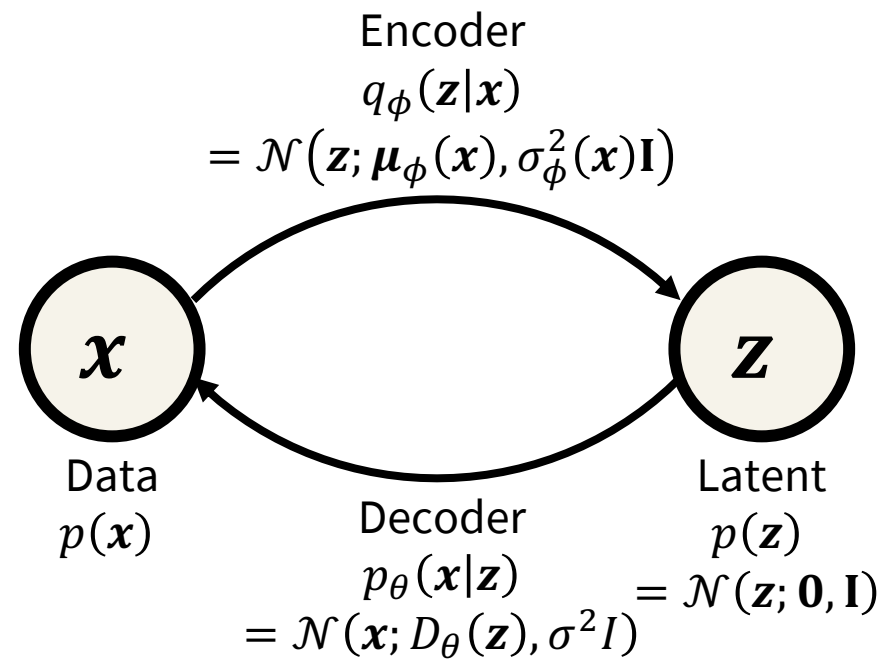
$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p(x)} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$



Previously in CS492D

Variational Autoencoder (VAE)

Introduce an **encoder** to learn a proxy of the posterior distribution.



Previously in CS492D

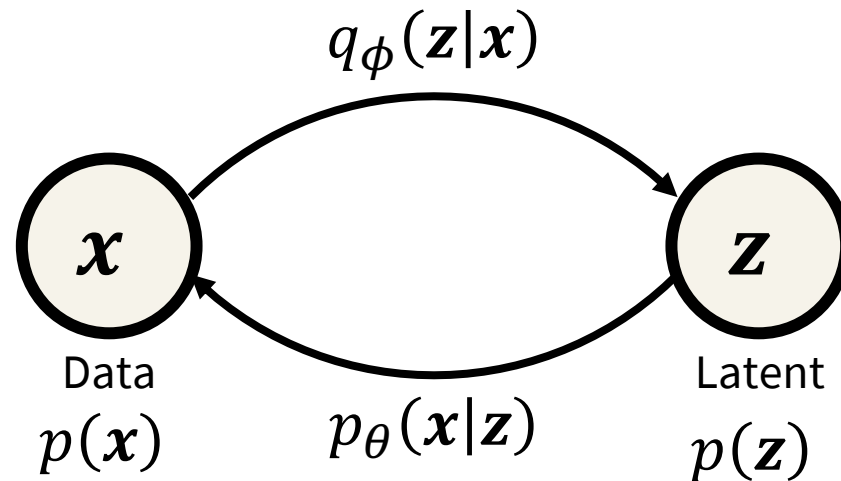
Basics

- Marginal distribution
- Expected value
- Bayes' rule
- Kullback–Leibler (KL) Divergence
- Jensen's inequality

Previously in CS492D

Bayes' Rule

$$\boxed{p(\mathbf{z}|\mathbf{x})}^{\text{Posterior Encoder}} = \frac{\boxed{p(\mathbf{x}|\mathbf{z})}^{\text{Likelihood Decoder}} \boxed{p(\mathbf{z})}^{\text{Prior Latent}}}{\boxed{p(\mathbf{x})}^{\text{Marginal Data}}}$$



Previously in CS492D

Evidence Lower Bound (ELBO)

- We cannot directly maximize $p(\mathbf{x}) = \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{z}|\mathbf{x})}$ since $p(\mathbf{z}|\mathbf{x})$ is unknown.
- Let's maximize **lower bound** of $\log p(\mathbf{x})$:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right]$$

Evidence Lower Bound (ELBO)

Let's decompose ELBO:

$$\begin{aligned}\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] &= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \\&= \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})} \right] \\&= \boxed{\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p(\mathbf{x}|\mathbf{z})]} - \boxed{D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}) \right)}\end{aligned}$$

Reconstruction term
to be *maximized*.

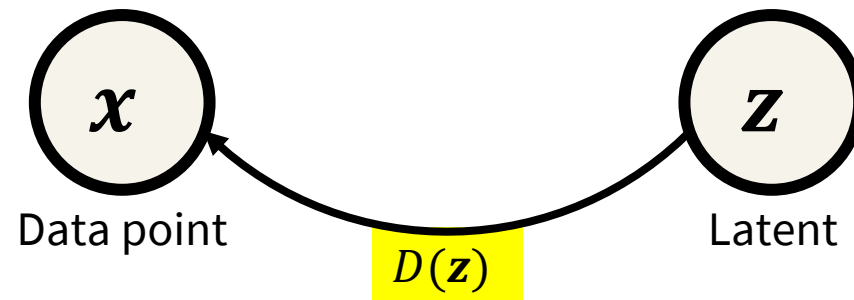
Prior matching term
to be *minimized*.

Back to VAE...

In VAE, we want model $p(\mathbf{x})$ as

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$$

where $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D(\mathbf{z}), \sigma^2 I)$ and $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.

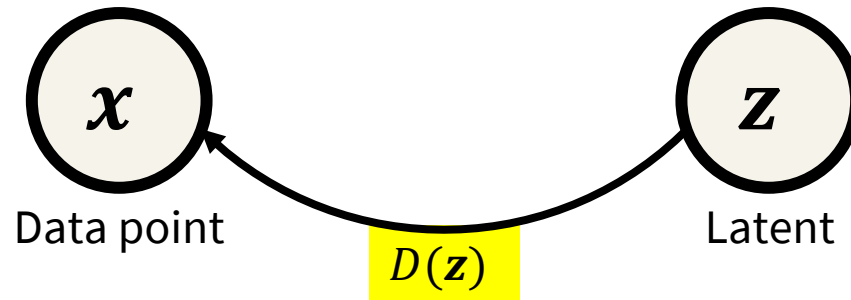


Variational Autoencoder (VAE)

In VAE, we want model $p(\mathbf{x})$ as

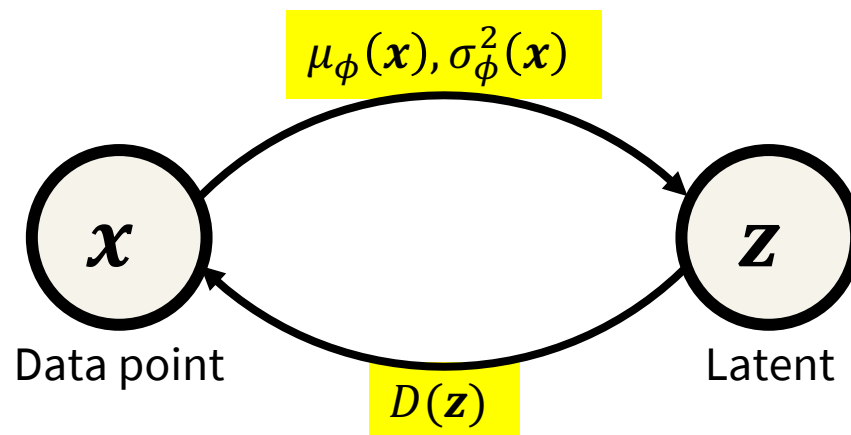
$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int \mathbf{p}_{\theta}(\mathbf{x}|\mathbf{z}) p(\mathbf{z}) d\mathbf{z}$$

where $\mathbf{p}_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \mathbf{D}_{\theta}(\mathbf{z}), \sigma^2 \mathbf{I})$ and $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.



Variational Autoencoder (VAE)

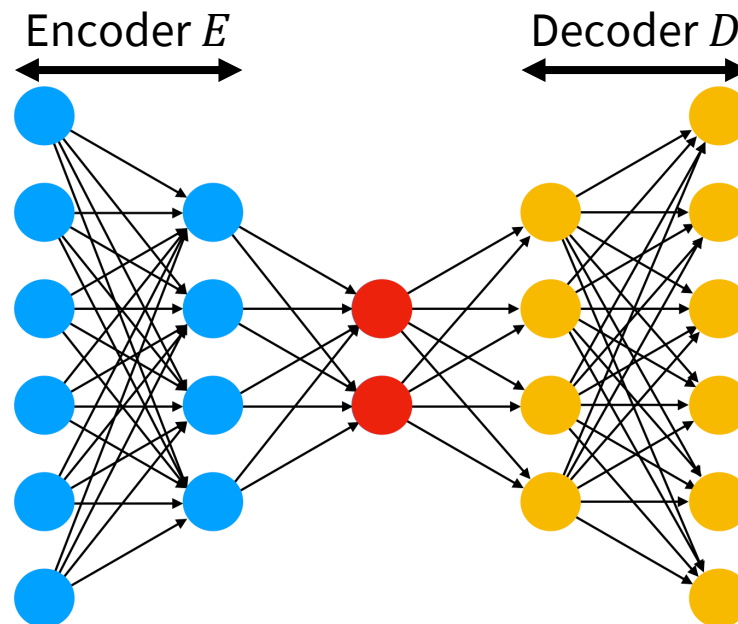
- How to maximize $p(\mathbf{x})$? Maximize **ELBO**!
- We need the **proxy** distribution $q_\phi(\mathbf{z}|\mathbf{x}) \rightarrow$ **Encoder**
- Let $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \mu_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})\mathbf{I})$



Variational Autoencoder (VAE)

Same with autoencoder but,

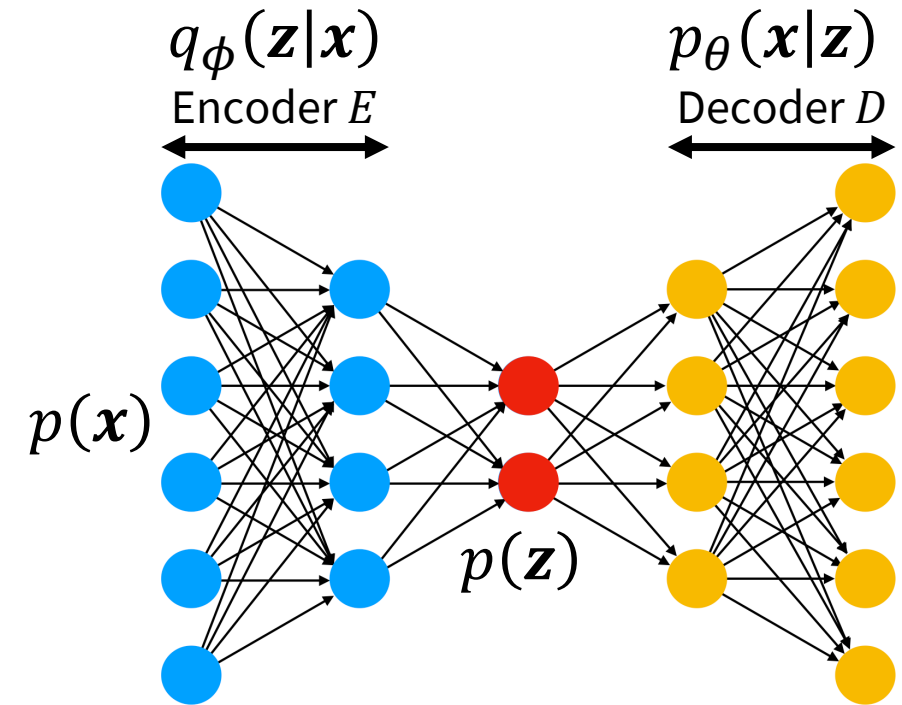
- From input \mathbf{x} , the **encoder** predicts $\mu_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x})$.
- The **decoder** takes a sample $\mathbf{z} \sim \mathcal{N}(\mathbf{z}; \mu_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x})\mathbf{I})$ as input.



Variational Autoencoder (VAE)

Summary

Data distribution	$p(\mathbf{x})$
Encoder	$q_{\phi}(\mathbf{z} \mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x})\mathbf{I})$
Latent distribution	$p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$
Decoder	$p_{\theta}(\mathbf{x} \mathbf{z}) = \mathcal{N}(\mathbf{x}; D_{\theta}(\mathbf{z}), \sigma^2\mathbf{I})$



Training

How to maximize **ELBO**?

$$\operatorname{argmax}_{\theta, \phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}) \right)$$

Approximates using a Monte Carlo estimate:

$$\operatorname{argmax}_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) - D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}) \right)$$

where $\mathbf{z}^{(i)} \sim q_{\phi}(\mathbf{z}|\mathbf{x})$ for the given \mathbf{x} .

[EXTRA] Monte Carlo Method

Law of Large Numbers, LLN

Where X_1, X_2, \dots is iid (independent and identically distributed) RV

$$\overline{X}_n = \frac{1}{n}(X_1 + \dots + X_n) \quad \overline{X}_n \rightarrow \mu \quad \text{as } n \rightarrow \infty.$$

$$\mathbb{E}_{\mathbf{x} \sim p} [f(\mathbf{x})] = \int_{\mathcal{X}} f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

$$\int_{\mathcal{X}} f(\mathbf{x}) p(\mathbf{x}) d\mathbf{x} \approx \frac{1}{n} \sum_{i=1}^N f(\mathbf{x}_i)$$

Training

How to maximize **ELBO**?

$$\operatorname{argmax}_{\theta, \phi} \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}) \right)$$

Approximates using a Monte Carlo estimate:

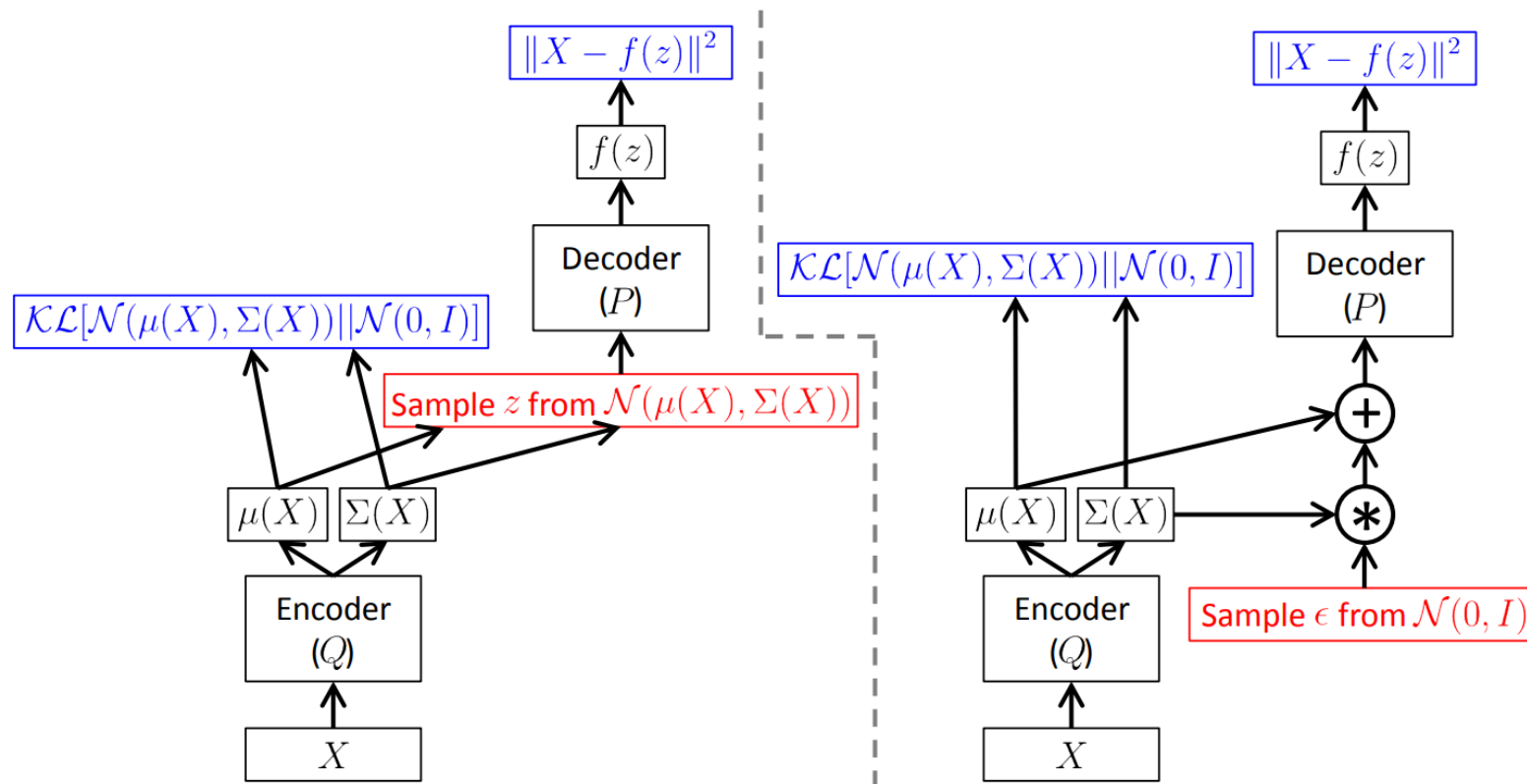
$$\operatorname{argmax}_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) - D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}) \right)$$

Reparameterization Trick

where $\mathbf{z}^{(i)} = \boldsymbol{\mu}_{\phi}(\mathbf{x}) + \boldsymbol{\sigma}_{\phi}(\mathbf{x})\boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim N(\mathbf{0}, \mathbf{I})$

[EXTRA] Reparameterization Trick

Not differentiable! -> Reparameterization Trick



[EXTRA] Reparameterization Trick

Why does it work?

Let's say we want to take the gradient w.r.t. θ of the following expectation,

$$\begin{aligned} & \mathbb{E}_{p(z)}[f_{\theta}(z)] \\ \nabla_{\theta} \mathbb{E}_{p(z)}[f_{\theta}(z)] &= \nabla_{\theta} \left[\int_z p(z) f_{\theta}(z) dz \right] \\ &= \int_z p(z) \left[\nabla_{\theta} f_{\theta}(z) \right] dz \\ &= \mathbb{E}_{p(z)} \left[\nabla_{\theta} f_{\theta}(z) \right] \end{aligned}$$

[EXTRA] Reparameterization Trick

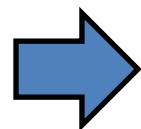
But if p is also parameterized by θ ?

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p(z)}[f_{\theta}(z)] &= \nabla_{\theta} \left[\int_z p(z) f_{\theta}(z) dz \right] & \nabla_{\theta} \mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)] &= \nabla_{\theta} \left[\int_z p_{\theta}(z) f_{\theta}(z) dz \right] \\ &= \int_z p(z) \left[\nabla_{\theta} f_{\theta}(z) \right] dz & &= \int_z \nabla_{\theta} \left[p_{\theta}(z) f_{\theta}(z) \right] dz \\ &= \mathbb{E}_{p(z)} \left[\nabla_{\theta} f_{\theta}(z) \right] & \longleftrightarrow &= \int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz + \int_z p_{\theta}(z) \nabla_{\theta} f_{\theta}(z) dz \\ & & &= \underbrace{\int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz}_{\text{What about this?}} + \mathbb{E}_{p_{\theta}(z)} \left[\nabla_{\theta} f_{\theta}(z) \right]\end{aligned}$$

[EXTRA] Reparameterization Trick

So, Reparameterization Trick!

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{p_{\theta}(z)}[f_{\theta}(z)] &= \nabla_{\theta} \left[\int_z p_{\theta}(z) f_{\theta}(z) dz \right] \\ &= \int_z \nabla_{\theta} [p_{\theta}(z) f_{\theta}(z)] dz \\ &= \int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz + \int_z p_{\theta}(z) \nabla_{\theta} f_{\theta}(z) dz \\ &= \underbrace{\int_z f_{\theta}(z) \nabla_{\theta} p_{\theta}(z) dz}_{\text{What about this?}} + \mathbb{E}_{p_{\theta}(z)} [\nabla_{\theta} f_{\theta}(z)]\end{aligned}$$



$$\epsilon \sim p(\epsilon)$$

$$\mathbf{z} = g_{\theta}(\epsilon, \mathbf{x})$$

$$\mathbb{E}_{p_{\theta}(\mathbf{z})}[f(\mathbf{z}^{(i)})] = \mathbb{E}_{p(\epsilon)}[f(g_{\theta}(\epsilon, \mathbf{x}^{(i)}))]$$

$$\nabla_{\theta} \mathbb{E}_{p_{\theta}(\mathbf{z})}[f(\mathbf{z}^{(i)})] = \nabla_{\theta} \mathbb{E}_{p(\epsilon)}[f(g_{\theta}(\epsilon, \mathbf{x}^{(i)}))] \quad (1)$$

$$= \mathbb{E}_{p(\epsilon)}[\nabla_{\theta} f(g_{\theta}(\epsilon, \mathbf{x}^{(i)}))] \quad (2)$$

$$\approx \frac{1}{L} \sum_{l=1}^L \nabla_{\theta} f(g_{\theta}(\epsilon^{(l)}, \mathbf{x}^{(i)})) \quad (3)$$

Training

Recall $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D_{\theta}(\mathbf{z}), \sigma^2 \mathbf{I})$.

$$\log p_{\theta}(\mathbf{x}|\mathbf{z}^{(i)}) = \log \left(\frac{1}{\sqrt{(2\pi\sigma^2)^d}} \exp \left(-\frac{\|\mathbf{x} - D_{\theta}(\mathbf{z})\|^2}{2\sigma^2} \right) \right)$$

$$= -\frac{1}{2\sigma^2} \boxed{\|\mathbf{x} - D_{\theta}(\mathbf{z})\|^2} - \log \sqrt{(2\pi\sigma^2)^d}$$

This is why it is called the reconstruction term.

Constant

[EXTRA] Why Gaussian in $p(\mathbf{x}|\mathbf{z})$?

Starting with $q(\mathbf{z}|\mathbf{x})$, $p(\mathbf{z})$, $p(\mathbf{x}|\mathbf{z})$, ... everything is gaussian...

What does $p(\mathbf{x}|\mathbf{z})$'s gaussian stands for?

$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D_{\theta}(\mathbf{z}), \sigma^2 I)$$

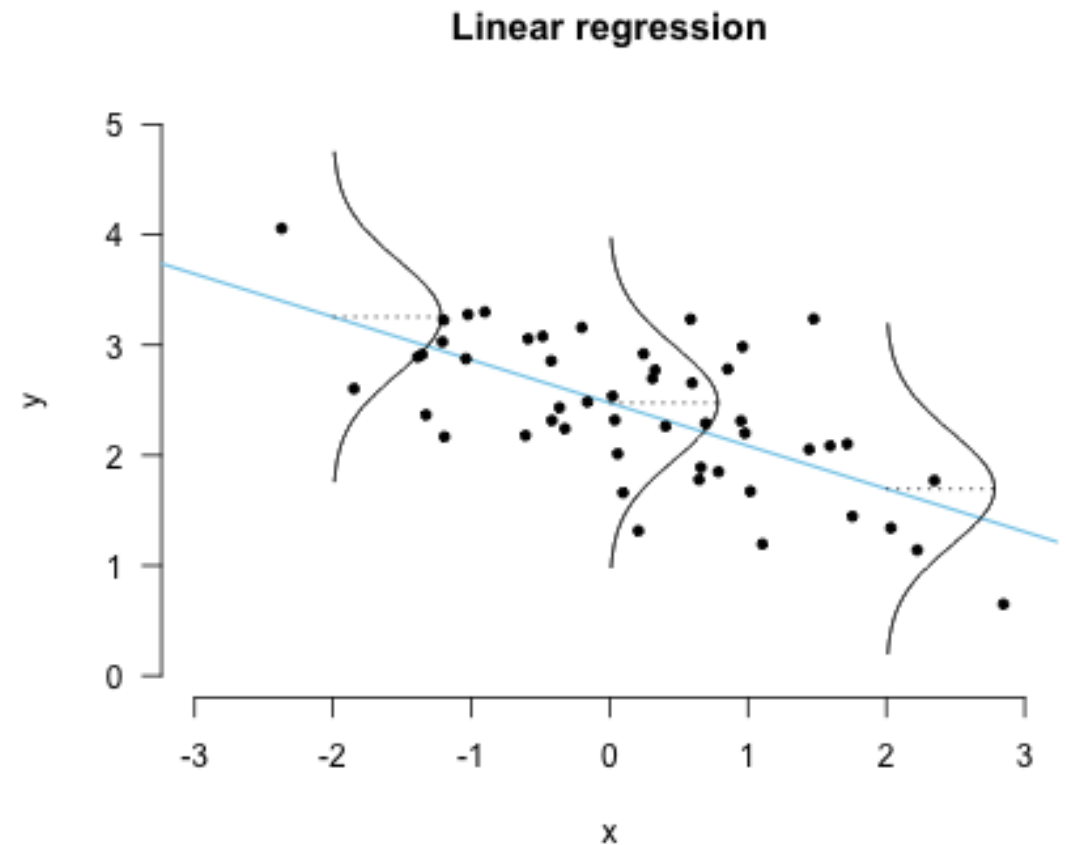
First, let's think about **Linear Regression**

Real continuous data has noise, and we assume predicting this as Gaussian

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_m x_m + \varepsilon$$

So, the pdf is as follows,

$$f(y | \hat{y}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\hat{y})^2}{2\sigma^2}}$$



[EXTRA] Gaussian's MLE = MSE

Then the Loss will be looking like...

$$\mathcal{L} = \prod_{i=1}^n f(y_i | \hat{y}_i, \sigma^2) = (2\pi\sigma^2)^{-n/2} e^{-\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{2\sigma^2}}$$

And same as we done, it is same as **MSE**!

Let's come back to VAE,

Because we are making the generative model,

It has to have randomness, predicting the probability,

assuming the gaussian of $p(x|z)$

[EXTRA] And actually... We can also use BCE

Because dataset like MNIST has value of [0~1],

$$\text{BCE}(y, \hat{y}) = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)]$$

In the perspective with interpreting distribution as Bernoulli

We can use BCE Loss! (For MSE, we assumed continuous distribution)

But it will be good for datasets having close pixels to 0 and 1

(SAME in AutoEncoder)

[EXTRA] How about Prior matching Term?

$$p(x) = \mathcal{N}(x; \mu, \sigma^2 I) \quad q(x) = \mathcal{N}(x; 0, I)$$

$$p(x) = \frac{1}{(2\pi)^{k/2} |\Sigma_p|^{1/2}} e^{-\frac{1}{2} (x-\mu)^\top \Sigma_p^{-1} (x-\mu)}$$

$$D_{\text{KL}}(p\|q) = \int p(x) \log \frac{p(x)}{q(x)} dx.$$

$$D_{\text{KL}}(p\|q) = \frac{1}{2} (k(\sigma^2 - 1 - \log \sigma^2) + \|\mu\|^2)$$

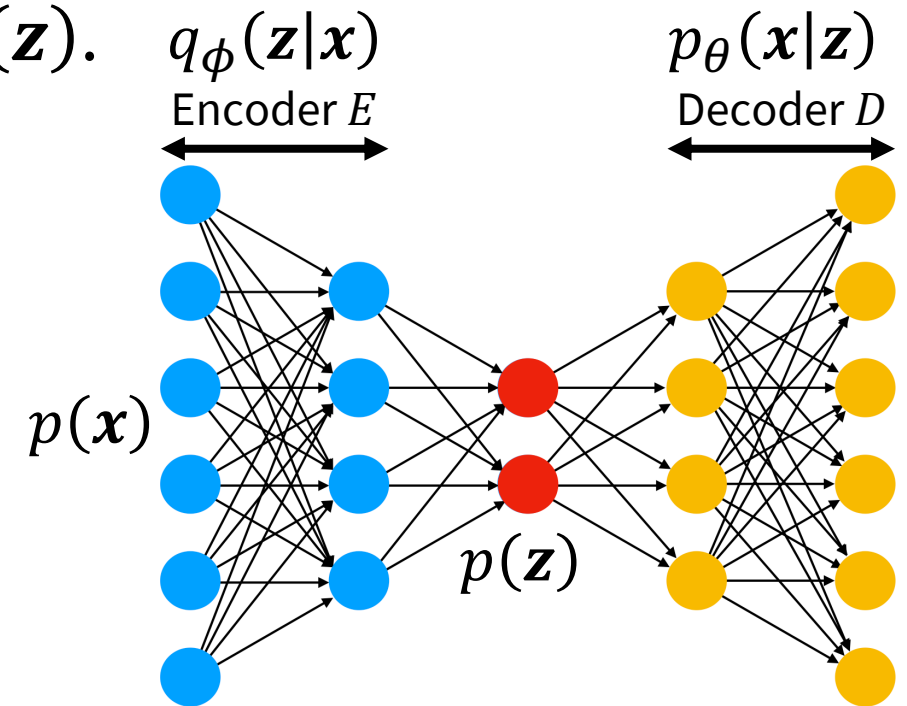
```
def loss_function(recon_x, x, mu, logvar):  
    MSE = F.mse_loss(recon_x, x, reduction='sum')  
    KLD = -0.5 * torch.sum(1+logvar-mu.pow(2)-logvar.exp())  
    return MSE+KLD
```

```
BCE = F.binary_cross_entropy(recon_x, x, reduction='sum')
```

Training

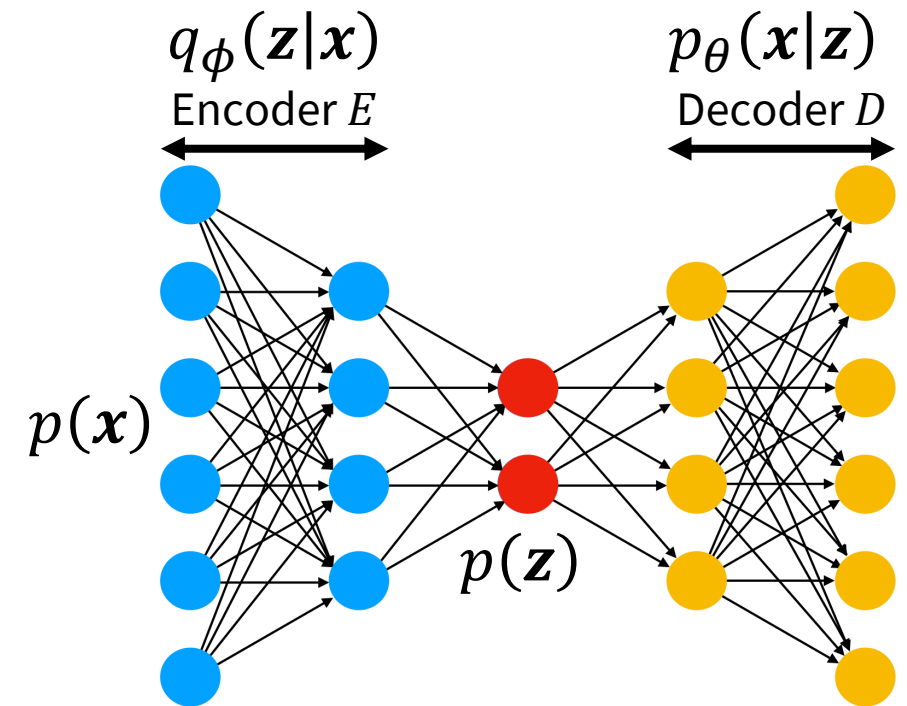
1. Feed a data point \mathbf{x} to the encoder to predict $\boldsymbol{\mu}_\phi(\mathbf{x})$ and $\sigma_\phi^2(\mathbf{x})$.
2. Sample a latent variable \mathbf{z} from $q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_\phi(\mathbf{x}), \sigma_\phi^2(\mathbf{x})\mathbf{I})$.
3. Feed \mathbf{z} to the decoder to predict $\hat{\mathbf{x}} = D_\theta(\mathbf{z})$.
4. Compute the gradient decent through the negative ELBO.

Q. Why is the **sampling** differentiable?



Generation

1. Sample a latent variable \mathbf{z} from $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$.
2. Feed \mathbf{z} to the decoder to predict $\hat{\mathbf{x}} = D_{\theta}(\mathbf{z})$.



Dimensions

Q. Should the **dimensions** of the input data and the latent variables be the same?

[EXTRA] VAE Overall Understanding

We are setting $p(z)$ to be standard gaussian distribution.

Reconstruction Loss makes generated x to be similar to data x ,

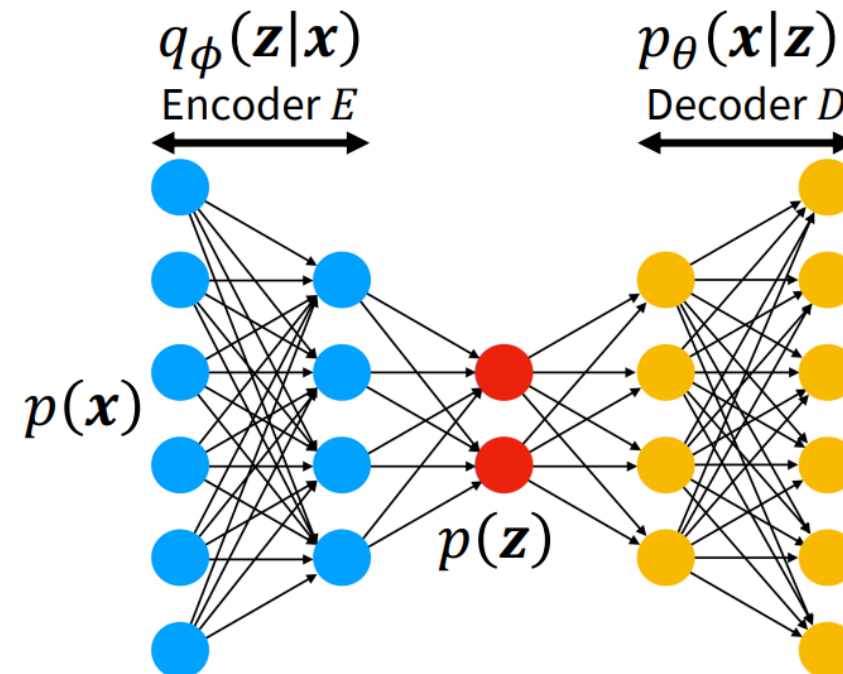
Prior matching term makes $q(z|x)$ to be similar to $p(z)$

- Then, what **meanings** would **latent vectors** be trained for?
- $q(z|x)$, $p(z)$, $p(x|z)$ has gaussian for its distribution, what does each stands for?
- Is the model or training **enough** for high quality generation?
 - > $p(z)$? $p(x|z)$? -> OK! But the problem is simple gaussian of $q(z|x)$

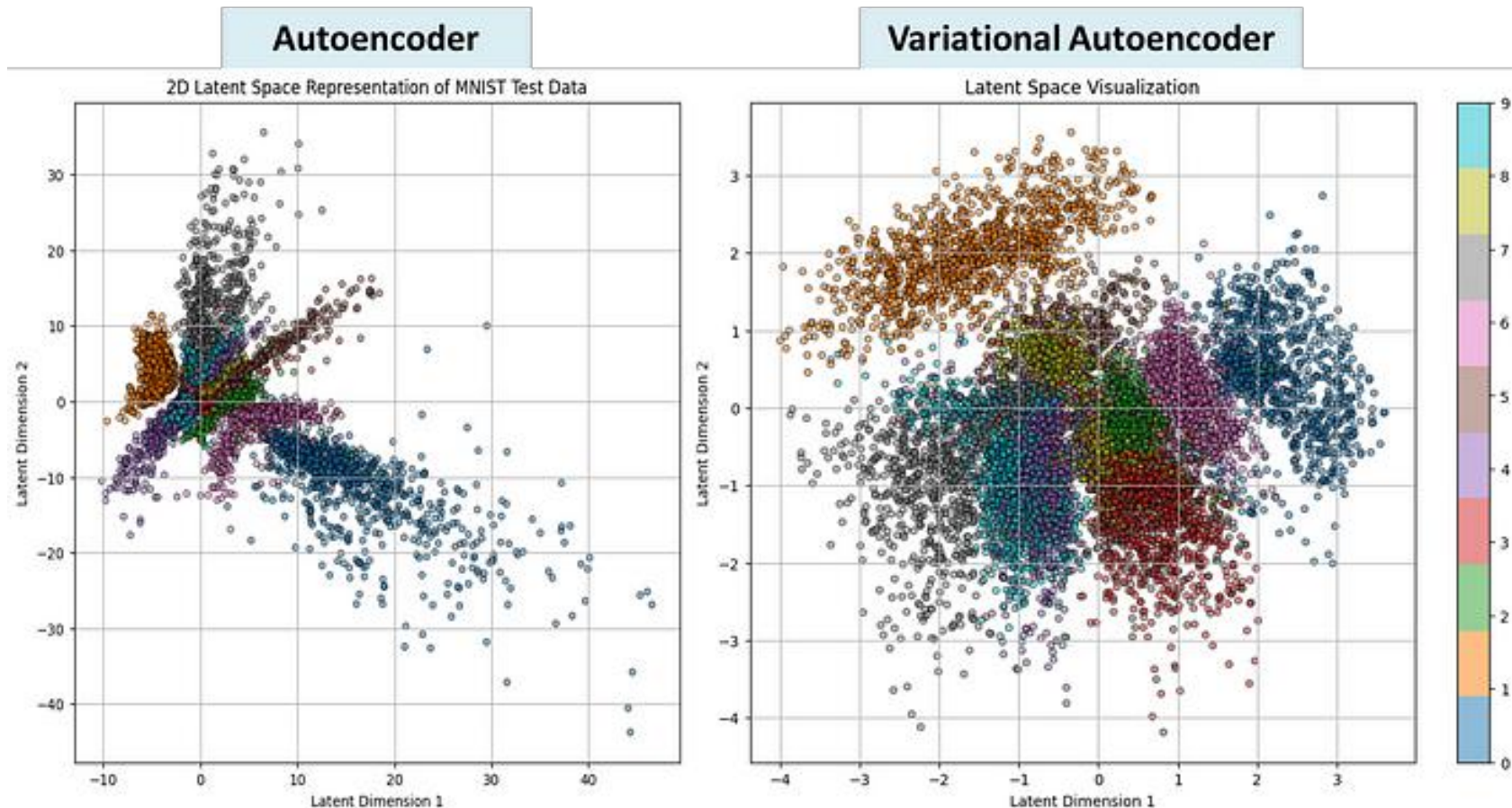
(Let's see this later)

[EXTRA] VAE Overall Understanding

- Then, what **meanings** would **latent vectors** be trained for?
- $q(z|x)$, $p(z)$, $p(x|z)$ has gaussian for its distribution, what does each stands for?
- Is the model or training **enough** for high quality generation?



[EXTRA] The difference between AE?



[EXTRA] VAE Overall Understanding

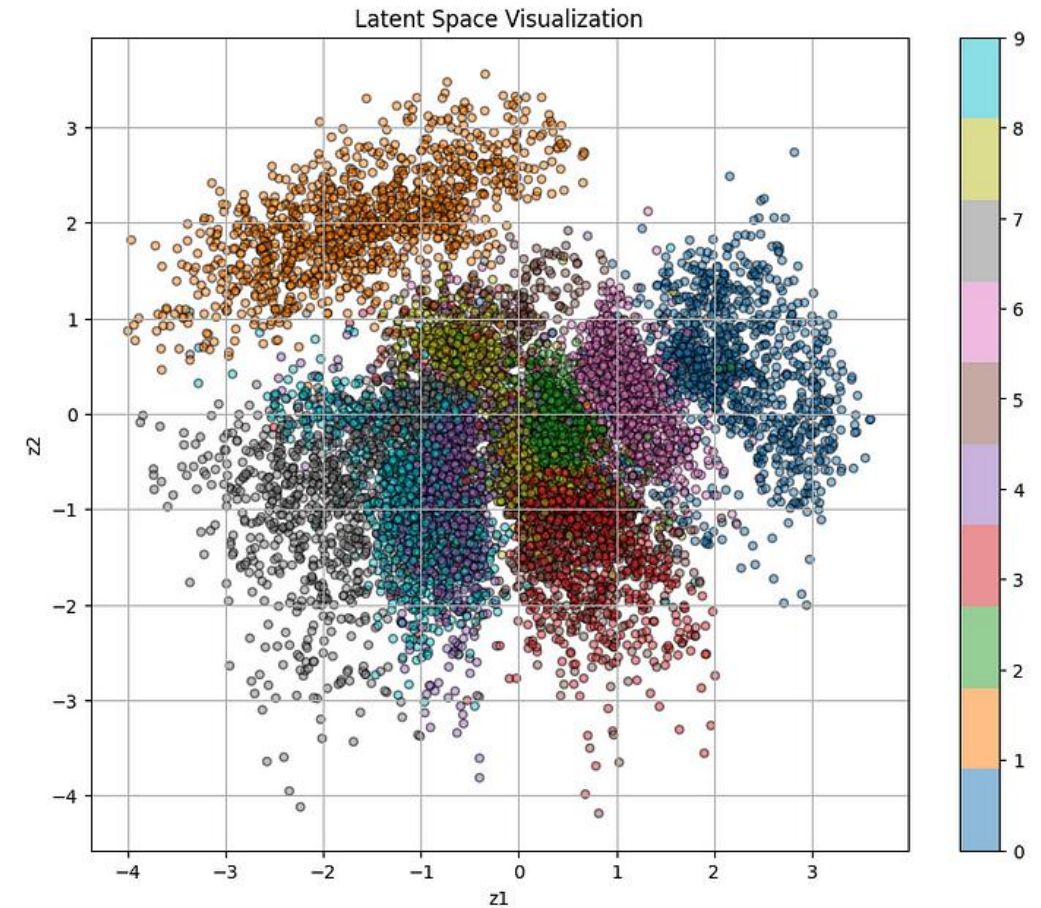
As we set $p(z)$, $q(z|x)$ to be gaussian distribution,

It seems to follow its distribution

But it seems to have **mixed** distribution

In **different** numbers (MNIST dataset)

With **2D** latent dimension



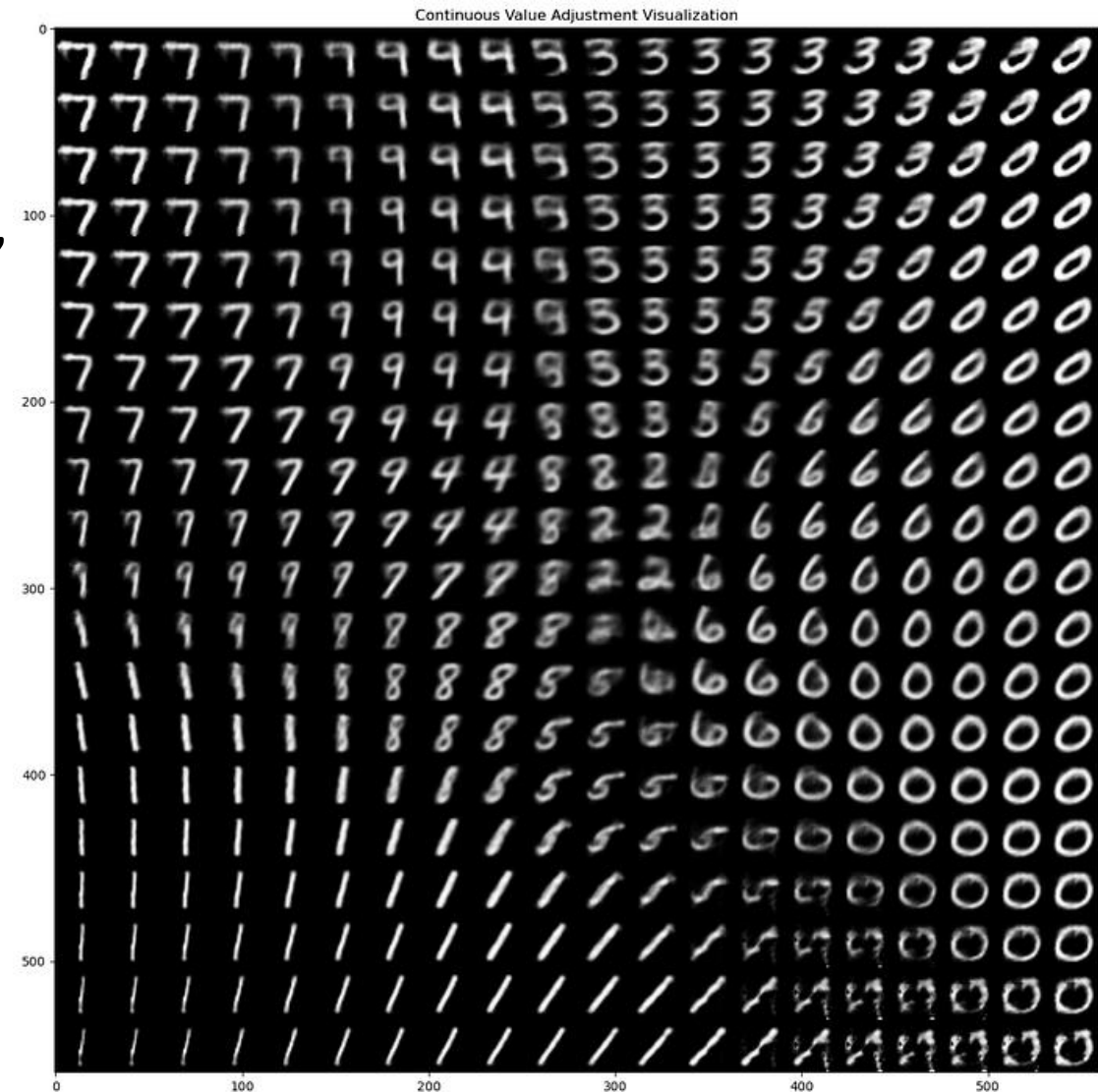
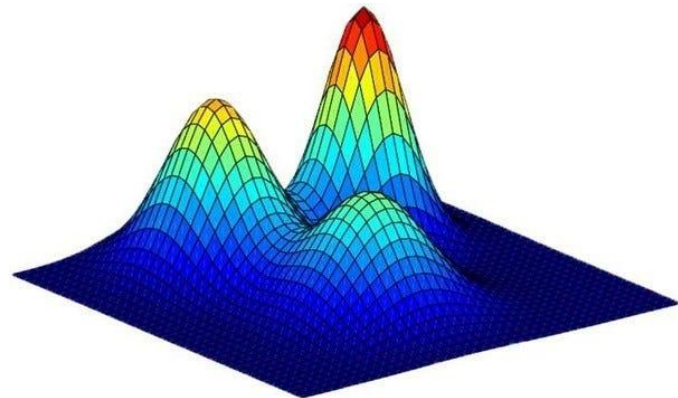
[EXTRA] VAE with multimodal distribution

Of course, 2D is not enough!

But still, if we want to have different $p(z)$,

Is it appropriate to learn with the same loss function?

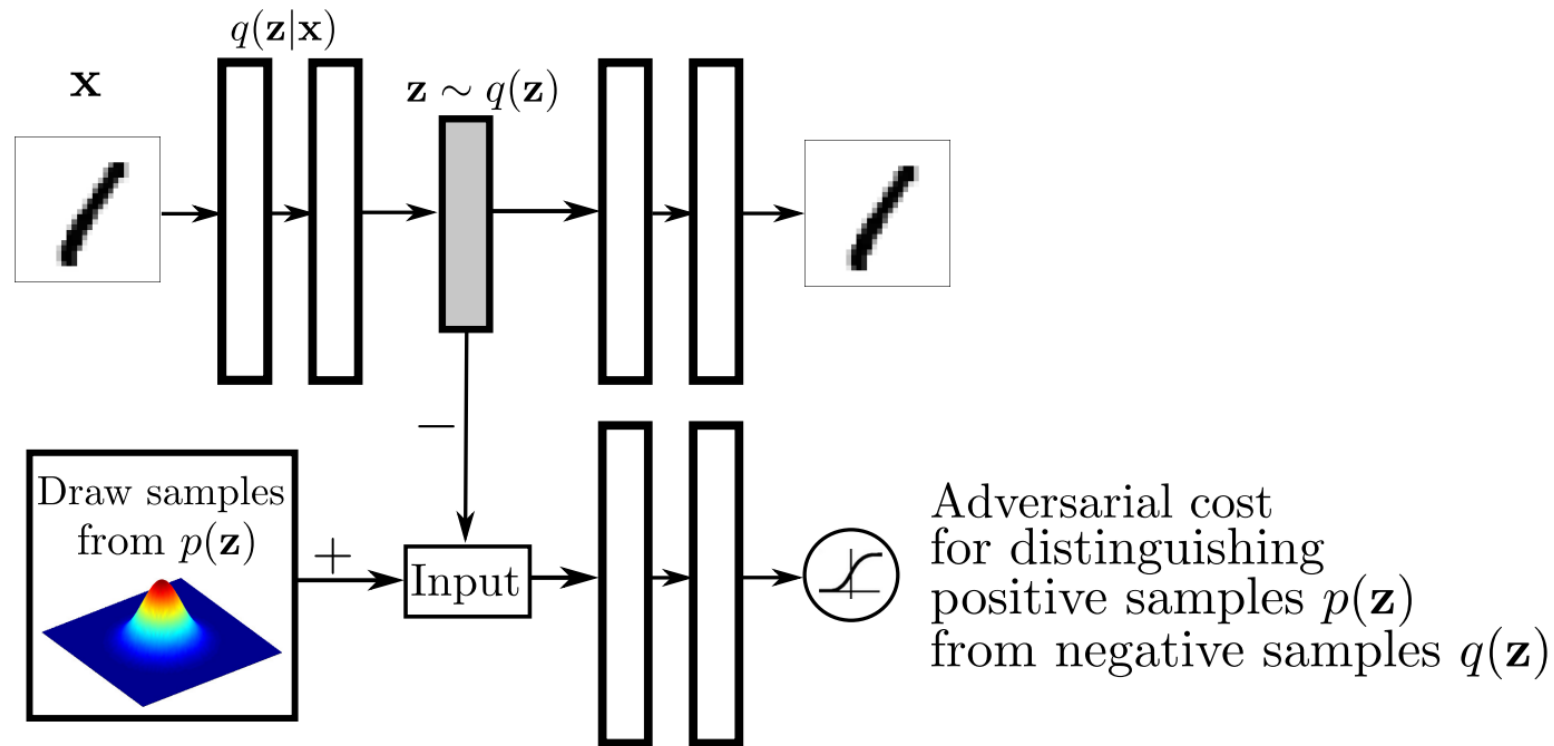
How about MoG(Mixture of Gaussians?)



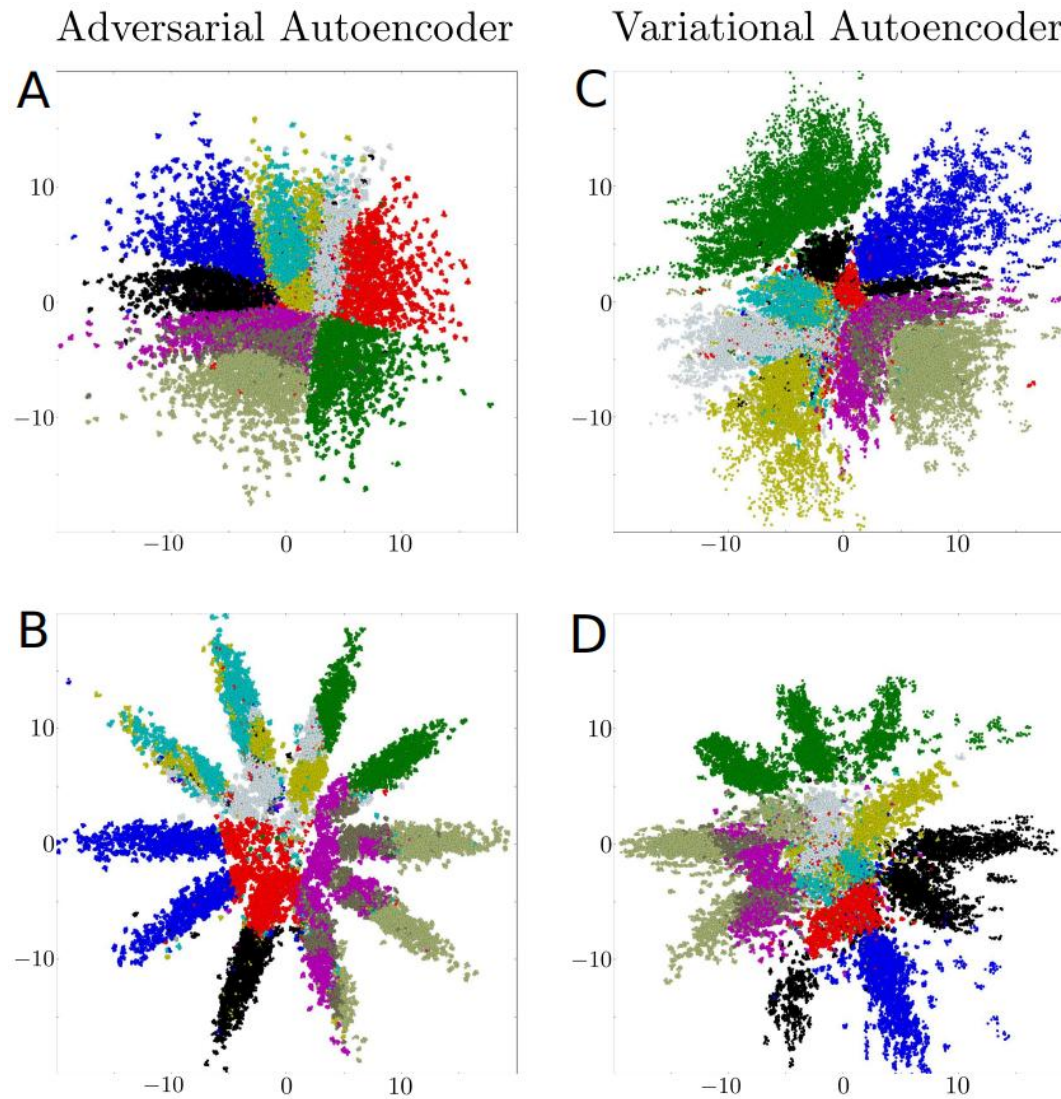
[EXTRA] Adversarial AutoEncoder

Because of Prior Matching Loss and $q(z|x)$'s single gaussian, making it difficult to follow the $q(z)$

One of the solutions... AAE (Adversarial AutoEncoder)



[EXTRA] Adversarial AutoEncoder

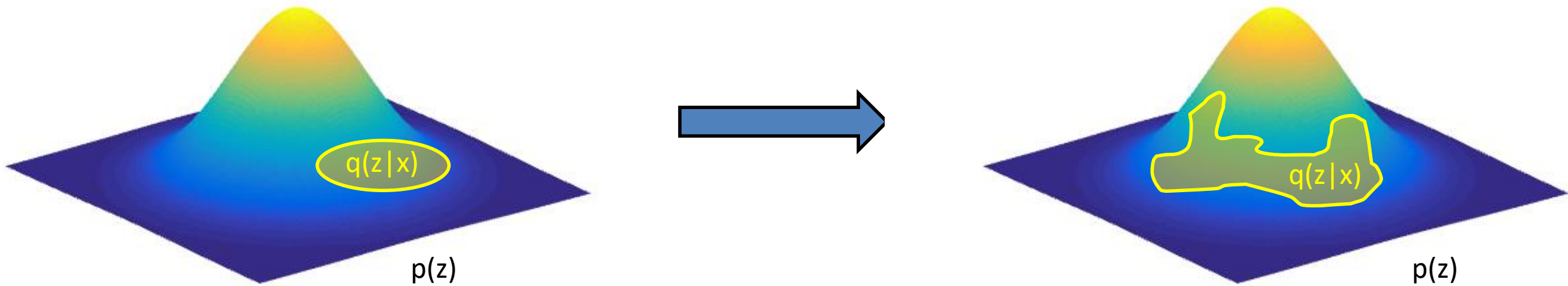


[EXTRA] Limitation of $q(z|x)$ in VAE

$p(z)$ being a standard gaussian distribution is not a limitation.

But the problem is simple gaussian of $q(z|x)$

We want to be more complex to make the model ideal!



[EXTRA] We can also decompose all from $p(\mathbf{x})$

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}) d\mathbf{z} \\ &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p_{\theta}(\mathbf{z}|\mathbf{x})} \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \int q_{\phi}(\mathbf{z}|\mathbf{x}) \log p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} - KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) + KL(q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z}|\mathbf{x}))\end{aligned}$$

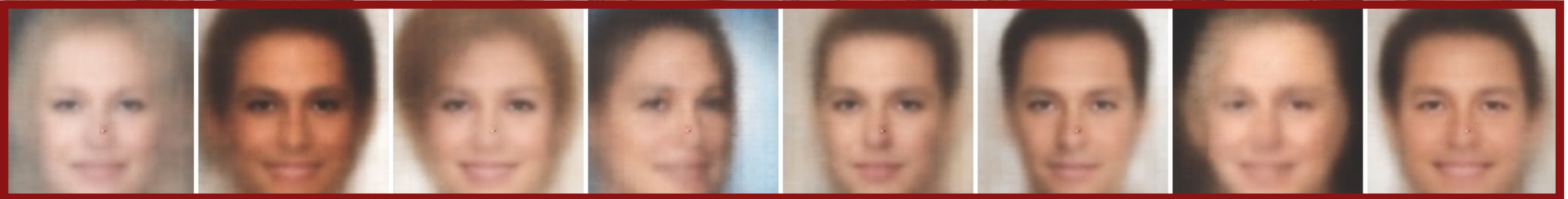
Limitations of VAEs

Typical failure cases of VAEs

Real
images



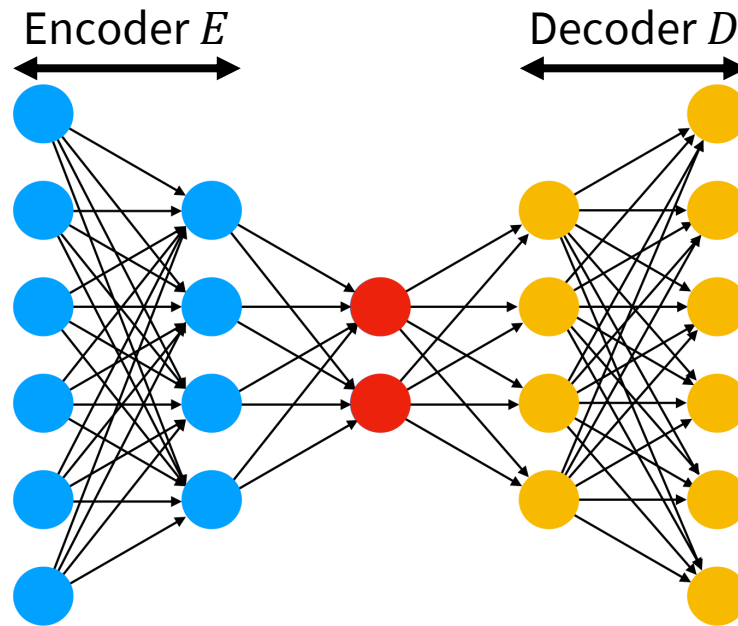
Generated
Images



Limitations of VAEs

Is a **Gaussian** distribution sufficient as the variational approximation for the posterior distribution?

$$q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_{\phi}(\mathbf{x}), \sigma_{\phi}^2(\mathbf{x})\mathbf{I})$$



Limitations of VAE

We maximize *not* $\log p(\mathbf{x})$ but its **lower bound** (ELBO).

Q. What is the difference between the two?

$$\left(\log p(\mathbf{x}) - \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] \right)$$

[This was the homework from the last class.]

Limitations of VAE

$$\mathbf{A.} \log p(\mathbf{x}) - \mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \right] = D_{KL} \left(q_{\phi}(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z}|\mathbf{x}) \right).$$

- The lower bound becomes **tight** when the variational distribution $q_{\phi}(\mathbf{z}|\mathbf{x})$ is **identical** to the true posterior distribution $p(\mathbf{z}|\mathbf{x})$.
- Will the true posterior distribution $p(\mathbf{z}|\mathbf{x})$ be close to a Gaussian distribution...?

Limitations of VAE

What is a better method for approximating the posterior distribution in a variational way?

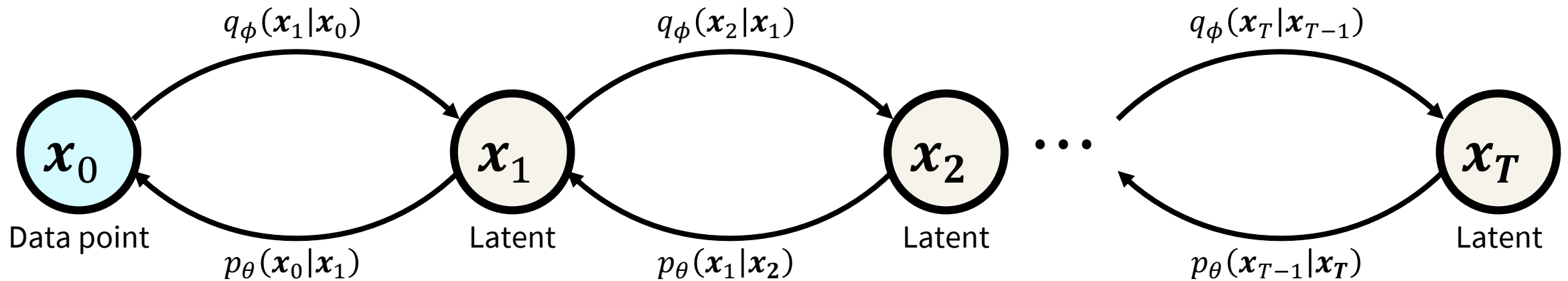
VAE Variants

- Vector-Quantized VAE (VQ-VAE) [Oord et al., NeurIPS 2017]
- Beta-VAE [Higgins et al., ICLR 2017]
- Wasserstein VAE (WAE) [Tolstikhin et al., ICLR 2018]
- VAE-GAN [Larsen et al., ICML 2016]
- Normalizing Flow VAE [Rezende and Mohamed, ICML 2015]

Hierarchical VAEs

Make a recursive (hierarchical) VAE.

- Data point: $\mathbf{x} \rightarrow \mathbf{x}_0$
- Latent variable(s): $\mathbf{z} \rightarrow \mathbf{x}_{1:T}$



Markovian Hierarchical VAEs

Let's consider a **Markovian** process.

A Markov process is a stochastic process where the probability of each event **depends only on the previous state**.

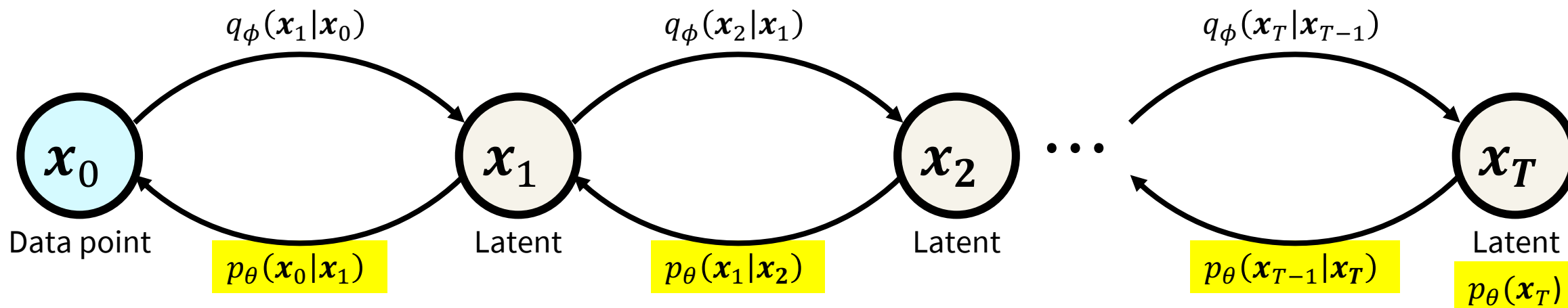
“What happens next depends only on the state of affairs now!”

Markovian Hierarchical VAEs

Let's consider a Markovian process.

Joint distribution:

$$p_{\theta}(\mathbf{x}_{0:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

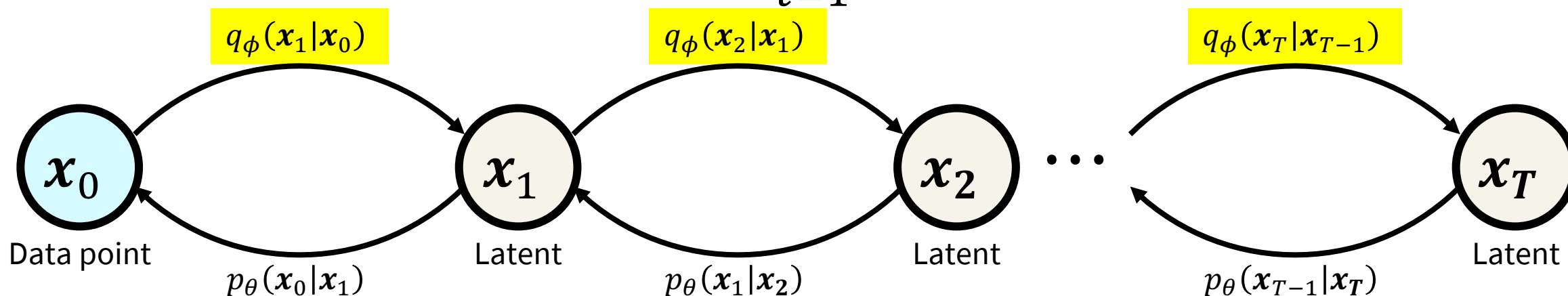


Markovian Hierarchical VAEs

Let's consider a Markovian process.

Variational posterior:

$$q_{\phi}(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q_{\phi}(\mathbf{x}_t|\mathbf{x}_{t-1})$$



Markovian Hierarchical VAEs

$$\begin{aligned}\log p(\mathbf{x}_0) &= \log \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \\ &= \log \int p_\theta(\mathbf{x}_{0:T}) \frac{q_\phi(\mathbf{x}_{1:T} | \mathbf{x}_0)}{q_\phi(\mathbf{x}_{1:T} | \mathbf{x}_0)} d\mathbf{x}_{1:T} \\ &= \log \mathbb{E}_{q_\phi(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\frac{p_\theta(\mathbf{x}_{0:T})}{q_\phi(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right] \\ &\geq \mathbb{E}_{q_\phi(\mathbf{x}_{1:T} | \mathbf{x}_0)} \left[\log \frac{p_\theta(\mathbf{x}_{0:T})}{q_\phi(\mathbf{x}_{1:T} | \mathbf{x}_0)} \right]\end{aligned}$$

VAEs → Diffusion Models

Flow-Based Models

- Normalizing flow
- Nonlinear Independent Components Estimation (NICE)
- Real Non-Volume Preserving (Real NVP)
- Generative Flow (Glow)
- Masked autoregressive flow (MAF)
- Continuous Normalizing Flow (CNF)

We'll revisit this later!

Denoising Diffusion Probabilistic Models (DDPM)

Ho et al., Denoising Diffusion Probabilistic Models, NeurIPS 2020.

Denoising Diffusion Probabilistic Models

Consider a special case of the Markovian hierarchical VAEs where:

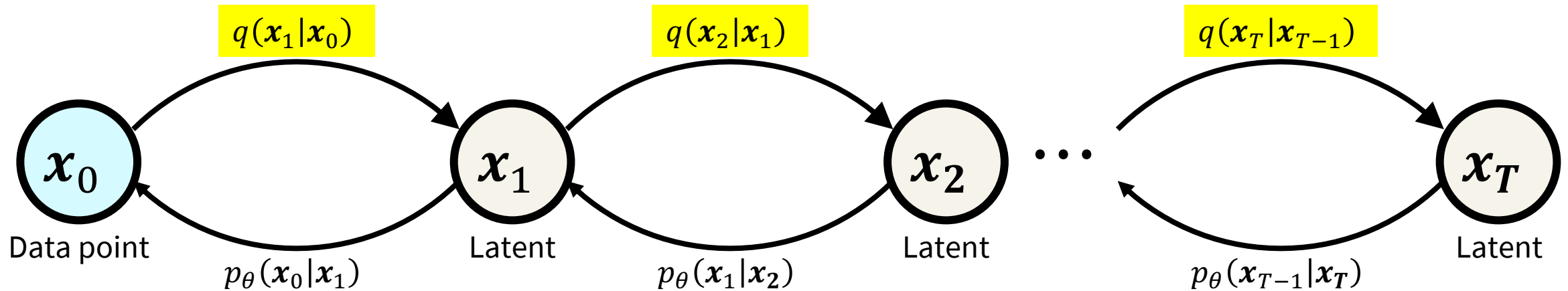
- the latent **dimension** is the **same** as the data dimension, and
- the **variational posteriors** $q_{\phi}(\mathbf{x}_{t+1}|\mathbf{x}_t)$ are not learned but **predefined**:

$$q_{\phi}(\mathbf{x}_{t+1}|\mathbf{x}_t) \rightarrow \mathbf{q}(\mathbf{x}_{t+1}|\mathbf{x}_t)$$

Terminology

Forward process (predefined):

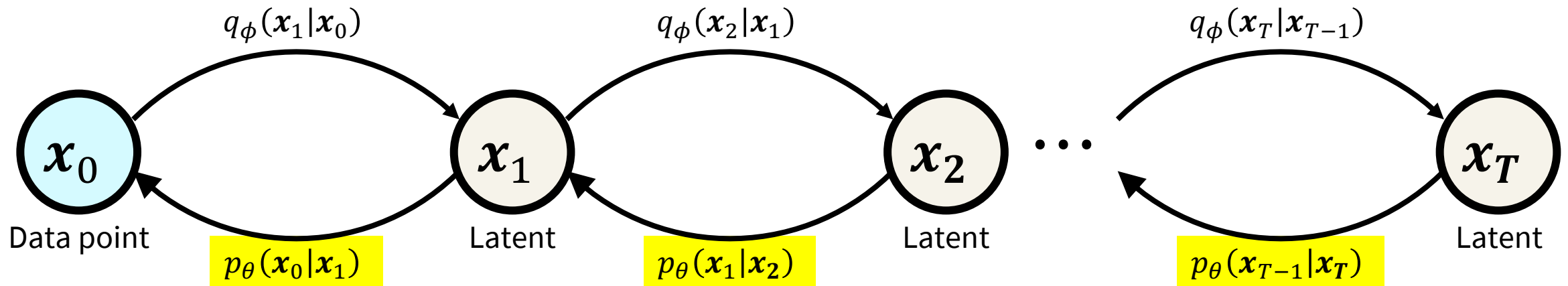
$$q(\mathbf{x}_{1:T} | \mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t | \mathbf{x}_{t-1})$$



Terminology

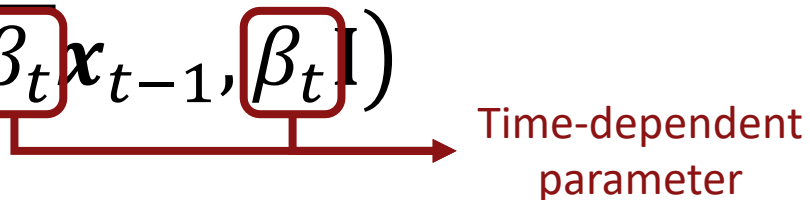
Reverse process (**learned**):

$$p_{\theta}(\mathbf{x}_0, \mathbf{x}_{1:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} | \mathbf{x}_t)$$



Forward Process (Data \rightarrow Latent)

In the forward process, the transition distribution $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ is specifically predefined as:

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


Time-dependent parameter

where $\{\beta_t \in (0, 1)\}_{t=1}^T$ and $\beta_1 \leq \beta_2 \leq \dots \leq \beta_T$.

“Adding Gaussian noise iteratively!”

VP-SDE vs. VE-SDE

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

is called **Variance Preserving (VP)** form.

- There are **other options**. For example:

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_{t-1}, (\sigma_t^2 - \sigma_{t-1}^2) \mathbf{I}),$$

which is called **Variance Exploding** form.

Choice of β_t

- Learned.
- Constant.
- Linearly or quadratically increased.
- Follows a **cosine** function
(Nichol and Dhariwal, Improved Denoising Diffusion Probabilistic Models, ICML 2021).
- Note that the reverse step $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ becomes a Gaussian form only when β_t is small ($\beta_t \ll 1$).

How to maximize ELBO in this case?

Disclaimer: We'll skip some complicated equations in the following slides.

ELBO

How to minimize the *negative* ELBO in this case?

$$-\log p(\mathbf{x}_0) = -\log \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$$= \dots =$$

$$-\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]$$

Reconstruction
term

$$+\mathbb{E}_{q(\mathbf{x}_{T-1}|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_T|\mathbf{x}_{T-1}) \parallel p(\mathbf{x}_T))]$$

Prior matching
term

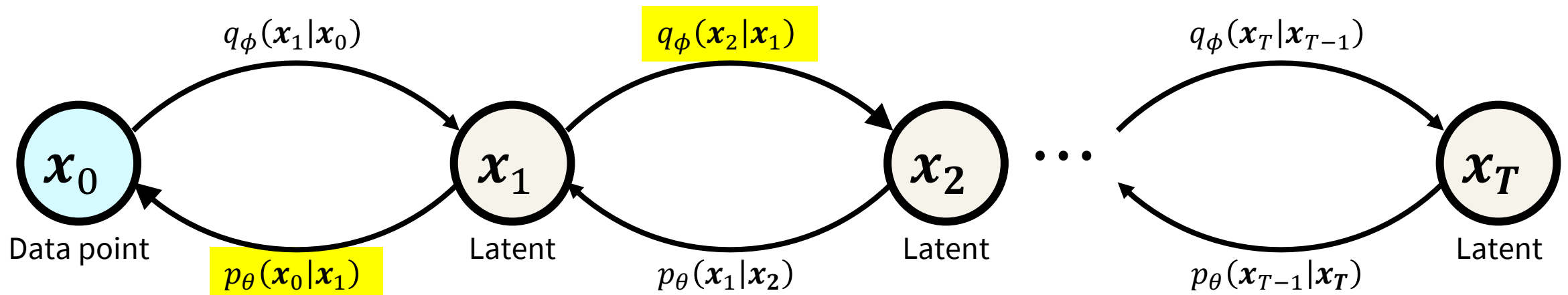
$$+\sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1}|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_t|\mathbf{x}_{t-1}) \parallel p_\theta(\mathbf{x}_t|\mathbf{x}_{t+1}))]$$

Consistency
term

Consistency Term

$$\sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1} | \mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))]$$

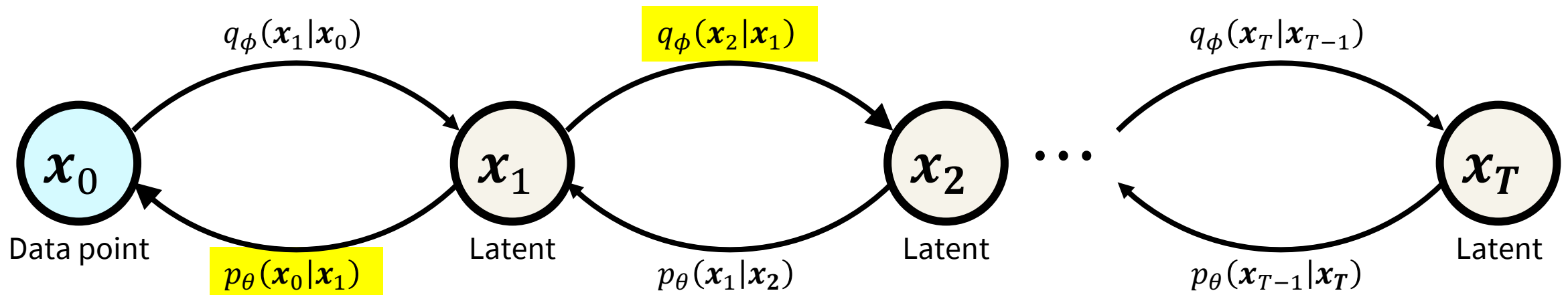
Make the forward and reverse steps be consistent at each time step.



Consistency Term

$$\sum_{t=1}^{T-1} \mathbb{E}_{q(\mathbf{x}_{t-1}, \mathbf{x}_{t+1} | \mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_t | \mathbf{x}_{t-1}) \parallel p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}))]$$

Expectation over two random variables; computationally expensive.



ELBO

Can we avoid having two random variables in an expectation?

Let's re-decompose the ELBO using the fact that

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0).$$

Q. Why $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = q(\mathbf{x}_t | \mathbf{x}_{t-1}, \mathbf{x}_0)$?

ELBO

Decompose the negative ELBO in a **different** way :

$$-\log p(\mathbf{x}_0) = -\log \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$$= \dots =$$

$$-\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]$$

Reconstruction
term \mathcal{L}_0

$$+D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))$$

New prior matching
term \mathcal{L}_T

$$+ \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

Denoising
matching
term \mathcal{L}_{t-1}

Reconstruction Term \mathcal{L}_0

$$-\mathbb{E}_{q(x_1|x_0)}[\log p_{\theta}(x_0|x_1)]$$

The same loss term as in VAE, but applied only to the **final** reverse step.

ELBO

Decompose the negative ELBO in a **different** way :

$$-\log p(\mathbf{x}_0) = -\log \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$= \dots =$

$$-\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]$$

$$+ D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))$$

New prior matching
term \mathcal{L}_T

$$+ \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

Prior Matching Term \mathcal{L}_T

Prior Matching Term \mathcal{L}_T

$$D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))$$

- Identical to the KL divergence term in VAE.
- Note that there is nothing to be optimized; $q(\mathbf{x}_T|\mathbf{x}_0)$ are $p(\mathbf{x}_T)$ are predefined.

Forward Convergence

Then, $q(\mathbf{x}_T | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$?

Yes, under certain assumptions.

Recall

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

where $\{\beta_t \in (0, 1)\}_{t=1}^T$ and $\beta_1 < \beta_2 < \dots < \beta_T$.

Forward Convergence

Then, $q(\mathbf{x}_T | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$?

Yes, under certain assumptions.

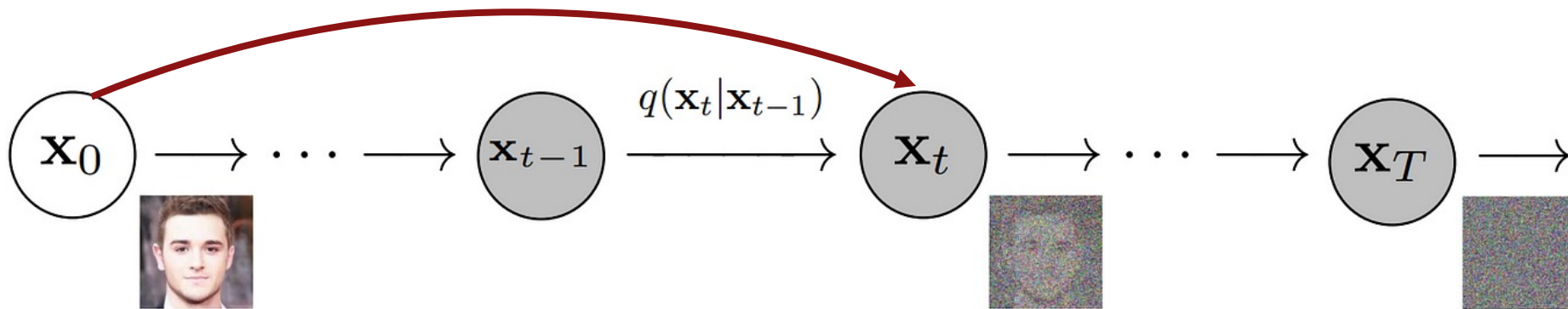
Let $\alpha_t = 1 - \beta_t$.

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) \mathbf{I})$$

where $\{\alpha_t \in (0, 1)\}_{t=1}^T$ and $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_T$.

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

Can we derive $q(\mathbf{x}_t | \mathbf{x}_0)$ from the sequence of $q(\mathbf{x}_{t'} | \mathbf{x}_{t'-1})$
for $t = 1, \dots, t'$?



Basics: Combination of Gaussian Variables

Suppose $\mathbf{x}_1 \sim \mathcal{N}(\mu_1, \sigma_1^2 \mathbf{I})$ and $\mathbf{x}_2 \sim \mathcal{N}(\mu_2, \sigma_2^2 \mathbf{I})$.

Q. What is the distribution of $\mathbf{x}_1 + \mathbf{x}_2$?

Basics: Combination of Gaussian Variables

Suppose $\mathbf{x}_1 \sim \mathcal{N}(\mu_1, \sigma_1^2 \mathbf{I})$ and $\mathbf{x}_2 \sim \mathcal{N}(\mu_2, \sigma_2^2 \mathbf{I})$.

A. $\mathbf{x}_1 + \mathbf{x}_2 \sim \mathcal{N}(\mu_1 + \mu_2, (\sigma_1^2 + \sigma_2^2) \mathbf{I})$

Basics: Combination of Gaussian Variables

Suppose $\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and

$\boldsymbol{x}_1 = \sigma_1 \boldsymbol{\varepsilon}_1$ and $\boldsymbol{x}_2 = \sigma_2 \boldsymbol{\varepsilon}_2$.

Q. What is the distribution of $\boldsymbol{x}_1 + \boldsymbol{x}_2$?

Basics: Combination of Gaussian Variables

Suppose $\boldsymbol{\varepsilon}_1, \boldsymbol{\varepsilon}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and

$$\boldsymbol{x}_1 = \sigma_1 \boldsymbol{\varepsilon}_1 \text{ and } \boldsymbol{x}_2 = \sigma_2 \boldsymbol{\varepsilon}_2.$$

$$\mathbf{A. } \boldsymbol{x}_1 + \boldsymbol{x}_2 \sim \mathcal{N}(\mathbf{0}, (\sigma_1^2 + \sigma_2^2)\mathbf{I}).$$

$\boldsymbol{x}_1 + \boldsymbol{x}_2 = \sqrt{\sigma_1^2 + \sigma_2^2} \boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon}$ is another standard normal sample.

Forward Convergence

$$q(\mathbf{x}_1|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_1; \sqrt{\alpha_1}\mathbf{x}_0, (1 - \alpha_1)\mathbf{I})$$

$$q(\mathbf{x}_2|\mathbf{x}_1) = \mathcal{N}(\mathbf{x}_2; \sqrt{\alpha_2}\mathbf{x}_1, (1 - \alpha_2)\mathbf{I})$$

Q. What is the distribution of $q(\mathbf{x}_2|\mathbf{x}_0)$?

Hint. Let's use the **reparamaterization trick**:

$$\mathbf{x}_1 = \sqrt{\alpha_1}\mathbf{x}_0 + \sqrt{1 - \alpha_1}\boldsymbol{\epsilon}_0$$

$$\mathbf{x}_2 = \sqrt{\alpha_2}\mathbf{x}_1 + \sqrt{1 - \alpha_2}\boldsymbol{\epsilon}_1$$

$$\boldsymbol{\epsilon}_0, \boldsymbol{\epsilon}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Forward Convergence

$$\mathbf{A.} \quad \mathbf{x}_2 = \sqrt{\alpha_2} \mathbf{x}_1 + \sqrt{1 - \alpha_2} \boldsymbol{\epsilon}_1$$

$$= \sqrt{\alpha_2} (\sqrt{\alpha_1} \mathbf{x}_0 + \sqrt{1 - \alpha_1} \boldsymbol{\epsilon}_0) + \sqrt{1 - \alpha_2} \boldsymbol{\epsilon}_1$$

$$= \sqrt{\alpha_2 \alpha_1} \mathbf{x}_0 + \sqrt{\alpha_2 (1 - \alpha_1)} \boldsymbol{\epsilon}_0 + \sqrt{1 - \alpha_2} \boldsymbol{\epsilon}_1$$

$$= \sqrt{\alpha_2 \alpha_1} \mathbf{x}_0 + \sqrt{(1 - \alpha_2 \alpha_1)} \bar{\boldsymbol{\epsilon}}_0$$

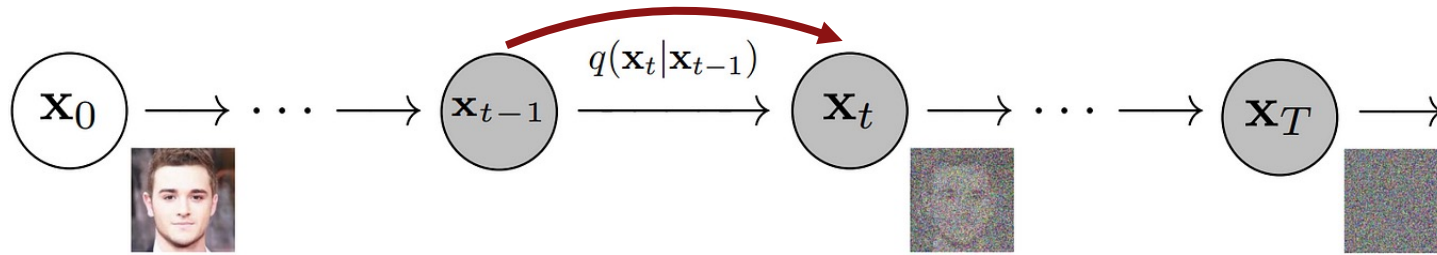
$$\therefore q(\mathbf{x}_2 | \mathbf{x}_0) = \mathcal{N}(\sqrt{\alpha_2 \alpha_1} \mathbf{x}_1, (1 - \alpha_2 \alpha_1) \mathbf{I})$$

Forward Convergence

$$\begin{aligned} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t} (\sqrt{\alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_{t-1}} \boldsymbol{\epsilon}_{t-2}) + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{\alpha_t (1 - \alpha_{t-1})} \boldsymbol{\epsilon}_{t-2} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{t-1} \\ &= \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{(1 - \alpha_t \alpha_{t-1})} \bar{\boldsymbol{\epsilon}}_{t-2} \\ &= \dots \\ &= \sqrt{\prod_{i=1}^t \alpha_i} \mathbf{x}_0 + \sqrt{(1 - \prod_{i=1}^t \alpha_i)} \bar{\boldsymbol{\epsilon}}_0 \end{aligned}$$

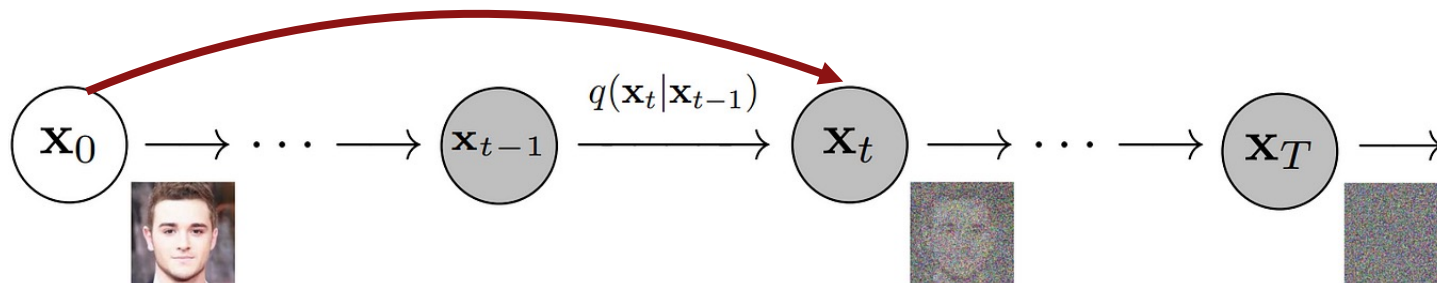
Forward Convergence

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) \mathbf{I})$$



$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$$

where $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ *Also a normal distribution!*

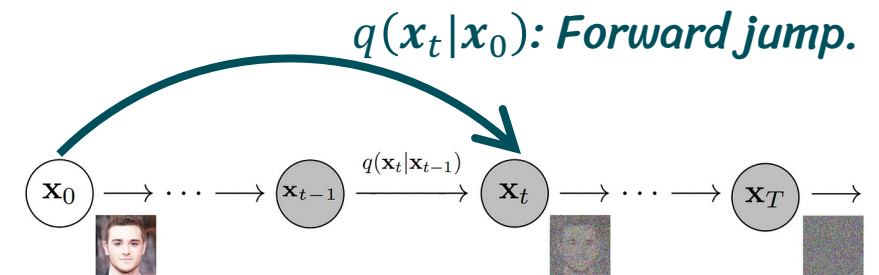


$$q(\mathbf{x}_t | \mathbf{x}_0)$$

Given \mathbf{x}_0 , \mathbf{x}_t at any **arbitrary timestep** t can be **directly sampled** from a Gaussian distribution without a Markov chain:

$$q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I}).$$

Note that $\bar{\alpha}_1 > \bar{\alpha}_2 > \dots > \bar{\alpha}_T$.



Forward Convergence

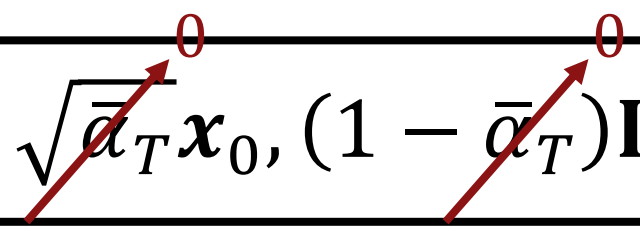
$$q(\mathbf{x}_T | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_T; \sqrt{\bar{\alpha}_T} \mathbf{x}_0, (1 - \bar{\alpha}_T) \mathbf{I})$$

where $\bar{\alpha}_T = \prod_{t=1}^T (1 - \beta_t)$.

Q. When $\{\beta_t \in (0, 1)\}_{t=1}^T$, What is

$$\lim_{T \rightarrow \infty} \bar{\alpha}_T = \lim_{T \rightarrow \infty} \prod_{t=1}^T (1 - \beta_t) ?$$

Forward Convergence

$$q(\mathbf{x}_T | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_T; \sqrt{\bar{\alpha}_T} \mathbf{x}_0, (1 - \bar{\alpha}_T) \mathbf{I})$$


where $\bar{\alpha}_T = \prod_{t=1}^T (1 - \beta_t)$.

As $T \rightarrow \infty$, $q(\mathbf{x}_T | \mathbf{x}_0)$ converges to the standard normal distribution $\mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$.

Prior Matching Term \mathcal{L}_T

$$D_{KL}(q(\mathbf{x}_T | \mathbf{x}_0) \parallel p(\mathbf{x}_T))$$

Close to zero by the definition of the forward transition distribution $q(\mathbf{x}_t | \mathbf{x}_{t-1})$. *Nothing to do for the optimization.*

ELBO

Decompose the negative ELBO in a **different** way :

$$-\log p(\mathbf{x}_0) = -\log \int p_\theta(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T} \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \left[-\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} \right]$$

$= \dots =$

$$-\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)} [\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)] \\ + \cancel{D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0) \parallel p(\mathbf{x}_T))} \rightarrow 0$$

$$+ \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

**Denoising
matching
term \mathcal{L}_{t-1}**

Denoising Matching Term \mathcal{L}_{t-1}

Denoising Matching Term \mathcal{L}_{t-1}

$$\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

The **variational** distribution $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$ should be close to $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ for each t .

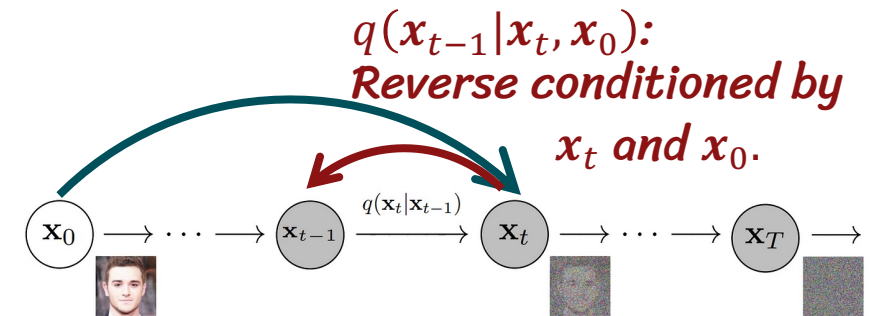
Denoising Matching Term \mathcal{L}_{t-1}

What is $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$?

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)}$$

Same as $q(\mathbf{x}_t|\mathbf{x}_{t-1})$, the forward transition, since it's a Markovian process.

We have seen how to compute these.



$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = q(\mathbf{x}_t | \mathbf{x}_{t-1}) \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)}$$

Q. What are $q(\mathbf{x}_t | \mathbf{x}_{t-1})$, $q(\mathbf{x}_{t-1} | \mathbf{x}_0)$, and $q(\mathbf{x}_t | \mathbf{x}_0)$?

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = q(\mathbf{x}_t | \mathbf{x}_{t-1}) \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)}$$

- $q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\sqrt{\alpha_t} \mathbf{x}_{t-1}, (1 - \alpha_t) \mathbf{I})$
- $q(\mathbf{x}_{t-1} | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0, (1 - \bar{\alpha}_{t-1}) \mathbf{I})$
- $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$

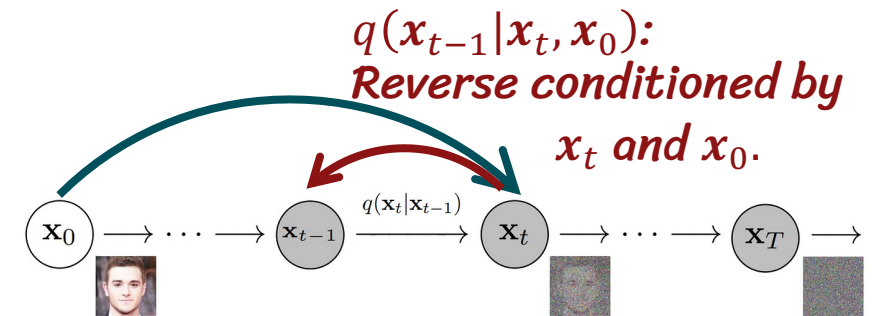
$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

$$\begin{aligned} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) &= q(\mathbf{x}_t | \mathbf{x}_{t-1}) \frac{q(\mathbf{x}_{t-1} | \mathbf{x}_0)}{q(\mathbf{x}_t | \mathbf{x}_0)} \\ &\propto \exp \left(-\frac{1}{2} \left(\frac{(\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_{t-1})}{1 - \alpha_t} + \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_{t-1}} \mathbf{x}_0)}{1 - \bar{\alpha}_{t-1}} - \frac{(\mathbf{x}_t - \sqrt{\bar{\alpha}_t} \mathbf{x}_0)}{1 - \bar{\alpha}_t} \right) \right) \\ &= \dots \\ &= \mathcal{N}(\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\sigma}_t^2 \mathbf{I}) \quad \text{Another normal distribution!} \end{aligned}$$

$$\text{where } \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 \text{ and } \tilde{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t.$$

$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

- The **mean** $\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \mathbf{x}_0$ is a function of both \mathbf{x}_t and \mathbf{x}_0 .
- The **covariance** $\tilde{\sigma}_t^2 \mathbf{I} = \left(\frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t \right) \mathbf{I}$ is **predefined** from the user-defined $\{\beta_t\}_{t=1}^T$.



$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

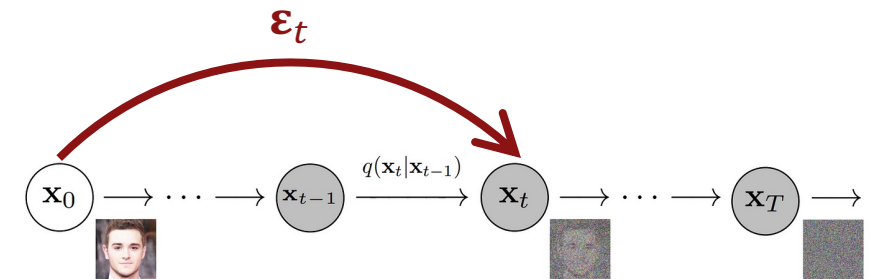
From the forward jump $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0, (1 - \bar{\alpha}_t) \mathbf{I})$,

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t \quad \text{where } \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

If \mathbf{x}_t and \mathbf{x}_0 are given, define $\boldsymbol{\epsilon}_t$ as

$$\boldsymbol{\epsilon}_t = \frac{1}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{x}_t - \frac{\sqrt{\bar{\alpha}_t}}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{x}_0$$

Q. Rewrite $\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0)$ as a function of \mathbf{x}_t and $\boldsymbol{\epsilon}_t$.



$$q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0)$$

A.

$$\begin{aligned}\tilde{\mu}(\mathbf{x}_t, \boldsymbol{\epsilon}_t) &= \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 \\ &= \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_t \right)\end{aligned}$$

Denoising Matching Term \mathcal{L}_{t-1}

Back to the denoising matching term...

$$\sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t))]$$

How to model the **variational** distribution $p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t)$?

Denoising Matching Term \mathcal{L}_{t-1}

- For $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0), \tilde{\sigma}_t^2 \mathbf{I})$,
the variance $\tilde{\sigma}_t^2$ is *not* a function of \mathbf{x}_t and \mathbf{x}_0 .
- Hence, define the variational distribution $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ as

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mu_\theta(\mathbf{x}_t, t), \tilde{\sigma}_t^2 \mathbf{I}),$$

where $\mu_\theta(\mathbf{x}_t, t)$ is the **mean predictor**.

Denoising Matching Term

How to compute

$$\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]?$$

Q. When $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_p, \sigma^2 \mathbf{I})$ and $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_q, \sigma^2 \mathbf{I})$,

What is $D_{KL}(p \parallel q)$?

[A similar problem as the homework from the last class.]

Denoising Matching Term

A.

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_p, \sigma^2 \mathbf{I})$$

$$q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_q, \sigma^2 \mathbf{I})$$

$$D_{\text{KL}}(p \parallel q) = \frac{1}{2\sigma^2} \|\boldsymbol{\mu}_q - \boldsymbol{\mu}_p\|^2$$

Denoising Matching Term

How to compute

$$\mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]?$$

$$\begin{aligned} & \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))] \\ &= \frac{1}{2\tilde{\sigma}_t^2} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}[\|\mu_\theta(\mathbf{x}_t, t) - \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0)\|^2] \end{aligned}$$

x_0 Predictor

Q. What if we have a x_0 predictor $\hat{x}_\theta(x_t, t)$ instead of the mean predictor $\mu_\theta(x_t, t)$? Note that

$$\tilde{\mu}(x_t, x_0) = \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} x_0$$

\mathbf{x}_0 Predictor

$$\mathbf{A.} \quad \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[D_{KL} \left(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) \right) \right]$$

$$= \frac{1}{2\tilde{\sigma}_t^2} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\mu_{\theta}(\mathbf{x}_t, t) - \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0)\|^2]$$

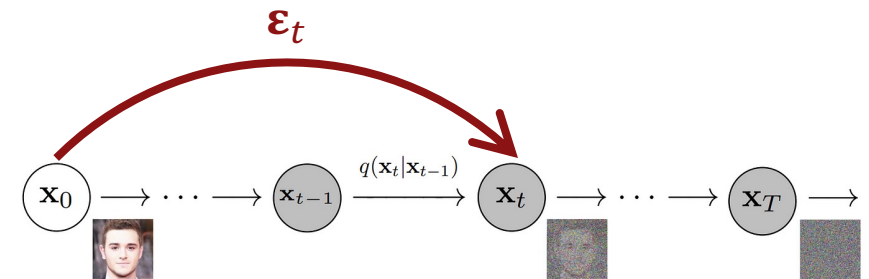
$$= \frac{1}{2\tilde{\sigma}_t^2} \frac{\bar{\alpha}_{t-1}\beta_t^2}{(1 - \bar{\alpha}_t)^2} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|^2]$$

$$= \omega_t \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|^2]$$

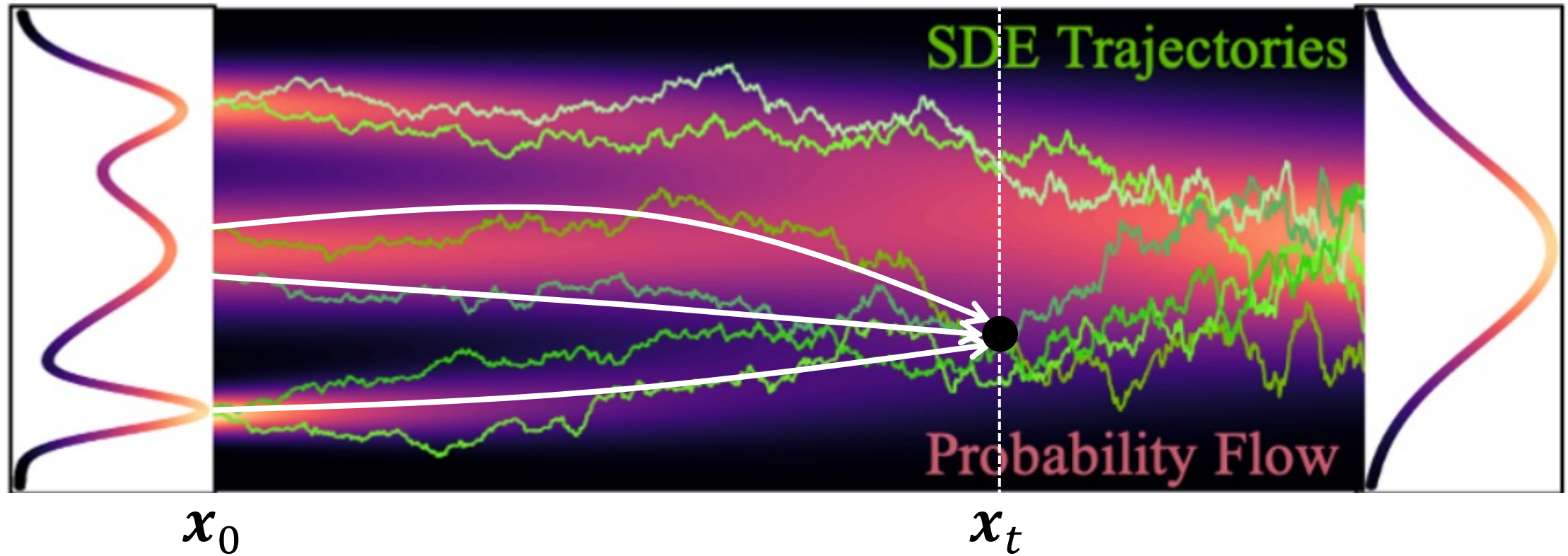
x_0 Predictor

$$\omega_t \mathbb{E}_{q(x_t|x_0)} [\|\hat{\mathbf{x}}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0\|^2]$$

- \mathbf{x}_t is sampled from \mathbf{x}_0 .
- From \mathbf{x}_t , predict the *expected* value of \mathbf{x}_0 that would result in sampling \mathbf{x}_t from it through the forward jump.



x_0 Predictor



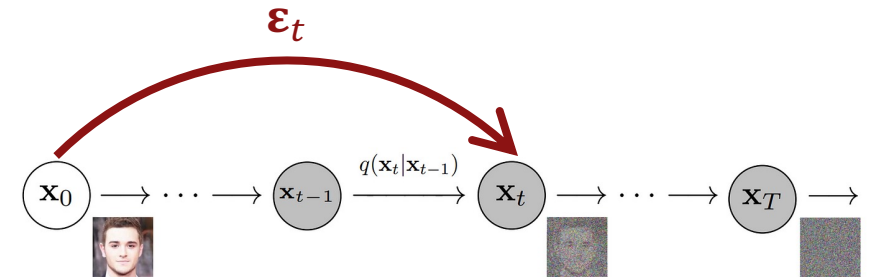
x_0 Predictor

- Note that our goal is to sample x_0 from a standard normal sample x_T and through latent variables $x_{T-1}, x_{T-2}, \dots, x_1$.
- But for every x_t , we directly predict the expected value of x_0 from x_t .

ϵ_t Predictor

Q. What if we have a ϵ_t predictor $\hat{\epsilon}_\theta(\mathbf{x}_t, t)$ instead of the mean predictor $\mu_\theta(\mathbf{x}_t, t)$? Note that

$$\tilde{\mu}(\mathbf{x}_t, \epsilon_t) = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_t \right)$$



ϵ_t Predictor

$$\begin{aligned} \mathbf{A.} \quad & \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} \left[D_{KL} \left(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \parallel p_{\theta}(\mathbf{x}_{t-1}|\mathbf{x}_t) \right) \right] \\ &= \frac{1}{2\tilde{\sigma}_t^2} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\mu_{\theta}(\mathbf{x}_t, t) - \tilde{\mu}(\mathbf{x}_t, \mathbf{x}_0)\|^2] \\ &= \frac{1}{2\tilde{\sigma}_t^2} \frac{(1 - \bar{\alpha}_t)^2}{\bar{\alpha}_t(1 - \bar{\alpha}_t)} \mathbb{E}_{q_{\phi}(\mathbf{x}_t|\mathbf{x}_0)} [\|\hat{\epsilon}_{\theta}(\mathbf{x}_t, t) - \epsilon_t\|^2] \\ &= \omega'_t \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\hat{\epsilon}_{\theta}(\mathbf{x}_t, t) - \epsilon_t\|^2] \end{aligned}$$

ϵ_t Predictor

$$\omega'_t \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)} [\|\hat{\epsilon}_\theta(\mathbf{x}_t, t) - \epsilon_t\|^2]$$

From \mathbf{x}_t , predict the *expected* value of ϵ_t that would result in sampling \mathbf{x}_t from \mathbf{x}_0 through the forward jump.

