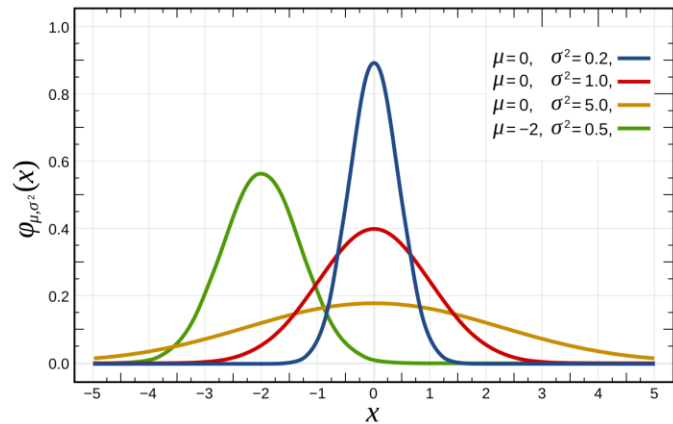# Advanced Machine Learning Generative Model

Yu Wang
Assistant Professor
Department of Computer Science
University of Oregon

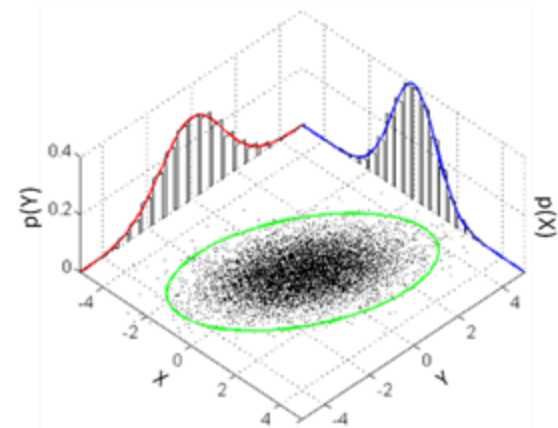# Summary

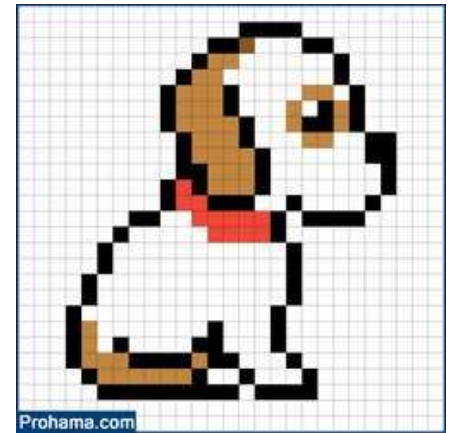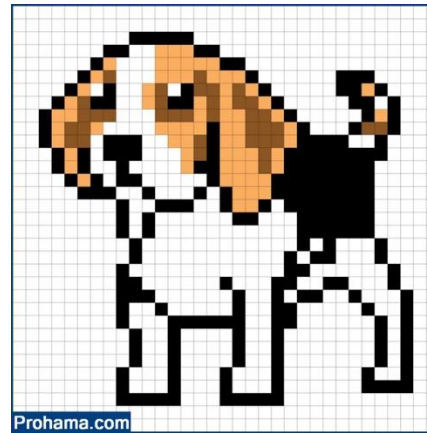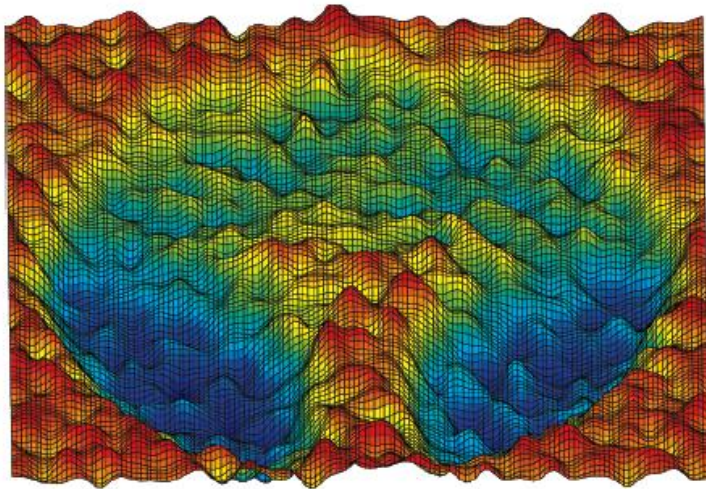## 1D Gaussian Distribution



$\mathbb{R}$

## 2D Gaussian Distribution



$\mathbb{R}^2$





$\mathbb{R}^{256 \times 256}$

# Summary

**Probability distribution** of the **objective** based on the **observed data**

$$\{x_i\}_{i=1}^{N} \xrightarrow[\text{Good Model}]{} P(x) \xrightarrow[\text{Good Data}]{} x$$

- **Machine Learning Methods**

  o Gaussian Kernel Density Estimation

  o Gaussian Mixture Models

  **Using existing function to estimate what you do not know that can best fit your observation**
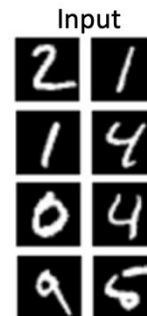
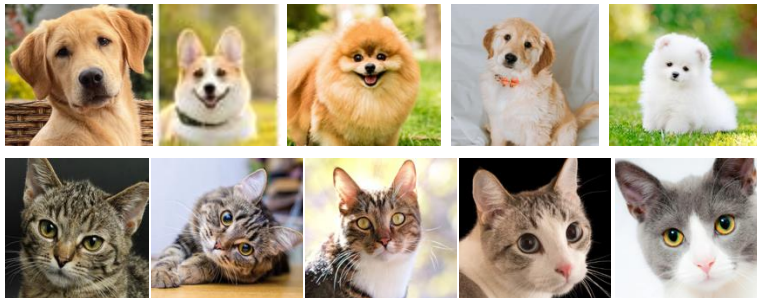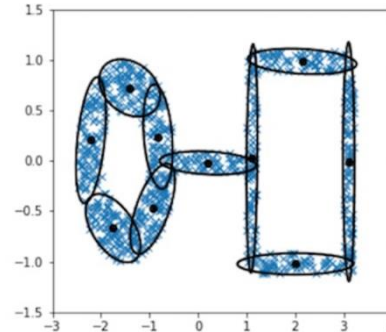- **Deep Learning Methods**

  o Auto-Encoder (AE)

  o Variational AE (LLM is actually a VAE)

  o Generative Adversarial Network

  o Diffusion Model

  **Using learnable function to estimate what you do not know that can best fit your observation**

# Problem?

**Using existing function to estimate what you do not know that can best fit your observation**



$$\mathbb{R}^1, \mathbb{R}^2$$

$$\downarrow$$

$$\mathbb{R}^{256 \times 256}$$

**What you have is some low-dimensional data**
**But what you want to model is some high-dimensional data, how it could be?**

# Problem?

**What we want: model any data distribution**

**How to transform any data distribution to low dimensional data?**

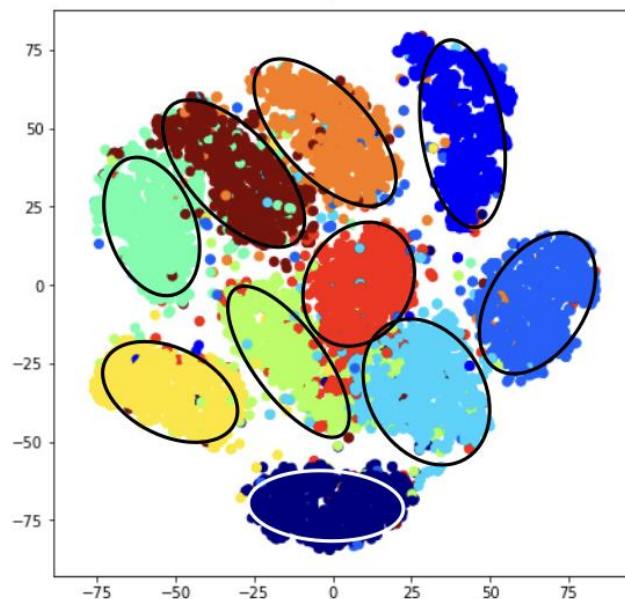**What we have: kernel density estimation to estimate low dimensional PDF**

**Kernel Density Estimation**

Input

**Someway to transform**

**Transform back**

KIND

# Summary

**Probability distribution** of the **objective** based on the **observed data**

- **Machine Learning Methods**

    $$\{x_i\}_{i=1}^N \xrightarrow[\text{Good Model}]{} P(x) \xrightarrow[\text{Good Data}]{} x$$

    o Gaussian Kernel Density Estimation

    o Gaussian Mixture Models

    Using **existing function** to **estimate what you do not know** that **can best fit your observation**

    ↓ **PCA Dimensional Reduction**

- **Deep Learning Methods**

    o Auto-Encoder (AE)

    o Variational AE (LLM is actually a VAE)

    o Generative Adversarial Network

    o Diffusion Model

    Using **learnable function** to **estimate what you do not know** that **can best fit your observation**
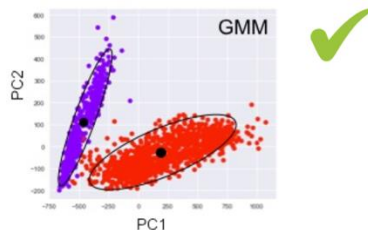
# From PCA to Auto-Encoder



PCA:

- Forward transform: $z = W^T x$
- Inverse transform: $\hat{x} = Wz$

Linear dimensionality Reduction

$$\min_{W} \mathbb{E}_x[\|x - \hat{x}\|^2] = \mathbb{E}_x[\| x - WW^T x\|^2]$$
$$s.t. \quad W^T W = I_{k \times k}$$

High-dimensional data often lives on non-linear manifolds that cannot be captured by linear models such as PCA



Fremont Bridge Hourly Bicycle Counts – Seattle

PCA $\phi: \mathbb{R}^d \to \mathbb{R}^2$

GMM ✔

MNIST dataset

PCA $\psi: \mathbb{R}^d \to \mathbb{R}^2$

✘

Can we add nonlinearity?
**Yes, then it becomes neural network!**

# Auto-Encoder

**Input**

**AE**

**Output**



$\phi: \mathbb{R}^d \to \mathbb{R}^k$

$\psi: \mathbb{R}^k \to \mathbb{R}^d$

# Class-supervised Auto-Encoder

**Input**

**AE**

**Output**

# Problem with AE



Input

Latent representation

Output

Encoder
$\phi(\cdot; \theta_\phi): \mathbb{R}^d \rightarrow \mathbb{R}^k$
e.g. NN

Decoder
$\psi(\cdot, \theta_\psi): \mathbb{R}^k \rightarrow \mathbb{R}^d$
e.g. NN

Need to estimate the latent distribution post-hoc!

SWAE:

Sliced Wasserstein Distance between two distributions!

Loss:

$$\min_{\theta_\phi, \theta_\psi} \mathbb{E}_{x \sim p_X} [\|x - \hat{x}\|^2] + \lambda \, SW(p_{Z|X}, q_Z)$$

$q_Z$

$p_X$

Encoder
$\phi(\cdot; \theta_\phi): \mathbb{R}^d \rightarrow \mathbb{R}^k$
e.g. NN

$p_{Z|X}$

Loss

Decoder
$\psi(\cdot, \theta_\psi): \mathbb{R}^k \rightarrow \mathbb{R}^d$
e.g. NN

$\hat{p}_X$

$$\underbrace{\mathbf{E}_z\left[\log p_\theta(x^{(i)} \mid z)\right]}_{\mathbb{E}_{x \sim p_X}[\|x - \hat{x}\|^2]} - \underbrace{D_{KL}(q_\phi(z \mid x^{(i)}) \| p_\theta(z))}_{}$$

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

$$D_{\mathrm{KL}}(q_\phi(z|x) \| p(z)) = \frac{1}{2}\sum_{j=1}^{d}\left[\sigma_j^2 + \mu_j^2 - 1 - \log \sigma_j^2\right]$$

**(1) Reconstruction loss: given z – decoder – x and setup the reconstruction loss**

**(2) KL divergence: how to optimize the KL divergence between two gaussian distributions?**

Sample $z$ from:  $z|x \sim \mathcal{N}(\mu_{z|x}, \Sigma_{z|x})$

Sample $x \mid z$ from:  $x|z \sim \mathcal{N}(\mu_{x|z}, \Sigma_{x|z})$

$\mu_{z|x}$     $\Sigma_{z|x}$

$\mu_{x|z}$     $\Sigma_{x|z}$

**Encoder** Network
$q_\phi(z|x)$
(parameters ϕ)

**Decoder** Network
$p_\theta(x|z)$
(parameters θ)

$x$

$z$

# Problem



Gaussian Noise 256x256 (3-channel)

**One-shot Generation** →

# Diffusion



Data ——— Destructing data by adding noise ———→ Noise

Data ←——— Generating samples by denoising ——— Noise

## Can we construct the image step by step?

# Diffusion



Data ——————→ Destructing data by adding noise ——————→ Noise

data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, $\qquad$ $\mathbf{x}_1, \mathbf{x}_2 \ldots \mathbf{x}_T$ with transition kernel $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}), \qquad \beta_t \in (0, 1) \text{ is a hyperparameter}$$



Step 0 | Step 10 | Step 20 | Step 30 | Step 40

# Diffusion



Data ⟶ Destructing data by adding noise ⟶ Noise

data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, $\quad\quad$ $\mathbf{x}_1, \mathbf{x}_2 \ldots \mathbf{x}_T$ with transition kernel $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

$$q(\mathbf{x}_t \mid \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\,\mathbf{x}_{t-1}, \beta_t \mathbf{I}), \quad\quad \beta_t \in (0,1) \text{ is a hyperparameter}$$

**Recursive**

$$p(x_t \mid x_0, x_1, \ldots, x_{t-1}) = p(x_t \mid x_{t-1})$$

$$x_t = \sqrt{1-\beta_t}\,x_{t-1} + \sqrt{\beta_t}\,\epsilon_t \quad \text{where} \quad \epsilon_t \sim \mathcal{N}(0, I)$$

**Markov Chain Property**

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\,\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}). \quad\quad \text{with } \alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s,$$

# Diffusion



Distribution at t=100 · Distribution at t=500 · Distribution at t=999

**Standard Gaussian**

Data ——————— Destructing data by adding noise ———————→ Noise
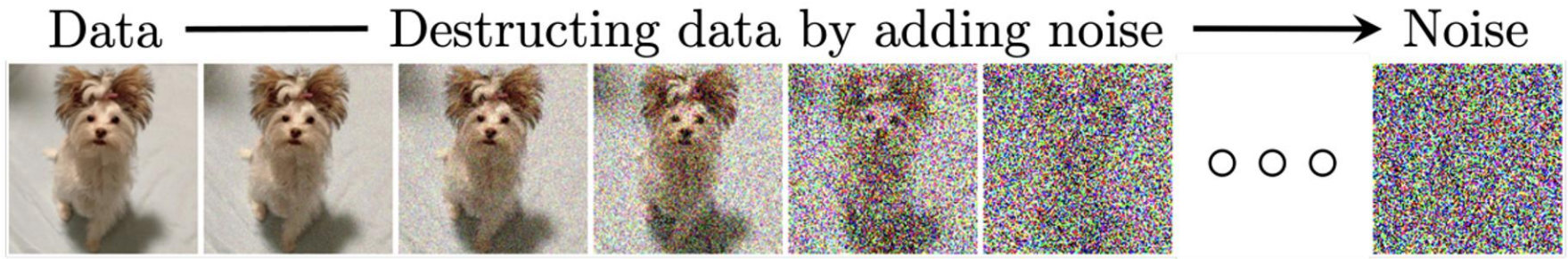
data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$,     $\mathbf{x}_1, \mathbf{x}_2 \ldots \mathbf{x}_T$ with transition kernel $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

$\beta_t \in (0, 1)$ is a hyperparameter

$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}).$     with $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s,$

$$t \to \infty, \alpha_t \to 0,$$
$p(\mathbf{x}_{t-1}|\mathbf{x}_t) \longleftarrow \qquad q(\mathbf{x}_t|\mathbf{x}_0) \to \mathcal{N}(0, 1)$



Data ←——————— Generating samples by denoising ——————— Noise

# Diffusion



Data ——— Destructing data by adding noise ———→ Noise

data distribution $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, $\mathbf{x}_1, \mathbf{x}_2 \ldots \mathbf{x}_T$ with transition kernel $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

$\beta_t \in (0, 1)$ is a hyperparameter

$$q(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}).$$
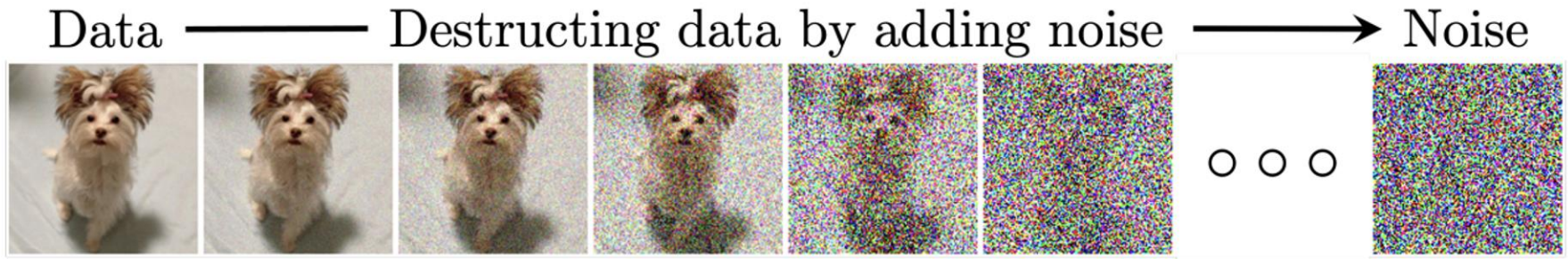
with $\alpha_t := 1 - \beta_t$ and $\bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s$,

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$$

**Approximate** ↓

$$p(\mathbf{x}_{t-1}|\mathbf{x}_t) \longleftarrow$$

$$t \to \infty, \alpha_t \to 0,$$
$$q(\mathbf{x}_t|\mathbf{x}_0) \to \mathcal{N}(0, 1)$$

Data ←——— Generating samples by denoising ——— Noise

$$p(x) = \int_z p_\theta(x|z)p(z)$$

$$p(x) = \int q_\phi(z|x)\frac{p_\theta(x|z)p(z)}{q_\phi(z|x)}$$

$$\log p(x) = \log \mathbb{E}_{z\sim q_\phi(z|x)}\left[\frac{p_\theta(x|z)p(z)}{q_\phi(z|x)}\right]$$

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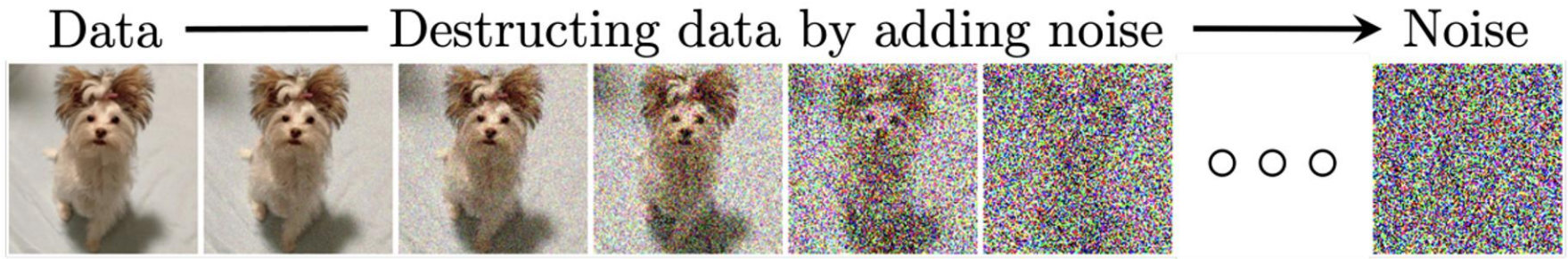


**Figure 1 - Graphical Model for VAE**

$$p(x) = \int_{z_1}\int_{z_2} p_\theta(x, z_1, z_2)dz_1, dz_2$$

$$p(x) = \int\int q_\phi(z_1, z_2|x)\frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)}$$

$$p(x) = \mathbb{E}_{z_1,z_2\sim q_\phi(z_1,z_2|x)}\left[\frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)}\right]$$

$$\log p(x) \geq \mathbb{E}_{z_1,z_2\sim q_\phi(z_1,z_2|x)}\left[\log\frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)}\right]$$



**Figure 2 - A Hierarchical VAE**

$$p(x, z_1, z_2) = p(x|z_1)p(z_1|z_2)p(z_2)$$

$$q(z_1, z_2|x) = q(z_1|x)q(z_2|z_1)$$

# Diffusion

Figure 2 - A Hierarchical VAE

$$\log p(x) \geq \mathbb{E}_{z_1,z_2 \sim q_\phi(z_1,z_2|x)}\left[\log \frac{p_\theta(x, z_1, z_2)}{q_\phi(z_1, z_2|x)}\right]$$



$$\log \text{p(x)} \geq \mathbb{E}_{x_{1:T} \sim q_\phi(x_{1:T}|x_0)}[\log \frac{p_\theta(x_{0:T})}{q_\phi(x_{1:T}|x_0)}]$$



$$=\mathbb{E}_{x_{1:T} \sim q_\phi(x_{1:T}|x_0)}[\log \frac{p_\theta(x_T) \prod_{t=1}^{T} p_\theta(x_{t-1}|x_t)}{\prod_{t=1}^{T} q_\phi(x_t|x_{t-1})}]$$

$$=\mathbb{E}_{x_{1:T} \sim q_\phi(x_{1:T}|x_0)}[\log p_\theta(x_T) + \sum_{t=1}^{T} \log \frac{p_\theta(x_{t-1}|x_t)}{q_\phi(x_t|x_{t-1})}]$$

Link

# Diffusion

$$\log \mathrm{p(x)} \geq \mathbb{E}_{x_{1:T} \sim q(X_{1:T}|X_0)}[\log p(x_T) + \sum_{t=1}^{T} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_t|x_{t-1})}]$$



$$q(x_t|x_{t-1})$$
$$x_0 \quad x_{t-1} \quad x_t \quad x_T$$
$$p_\theta(x_{t-1}|x_t)$$

Link

$$L := \mathbb{E}_q \left[ \underbrace{-\log p(x_T) + \log q(x_T|x_0)}_{L_T} - \underbrace{\log p_\theta(x_0|x_1)}_{L_0} - \underbrace{\sum_{t>1} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)}}_{L_{t-1}} \right]$$

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$



Data ——— Destructing data by adding noise ——→ Noise

**Given $x_t, x_0$, how to get $x_{t-1}$ using diffusion**



Data ←——— Generating samples by denoising ——— Noise

**Given $x_t$, how to revert $x_{t-1}$ using decoder**

$$L := \mathbb{E}_q \left[ \underbrace{-\log p(x_T) + \log q(x_T|x_0)}_{L_T} - \underbrace{\log p_\theta(x_0|x_1)}_{L_0} - \underbrace{\sum_{t>1} \log \frac{p_\theta(x_{t-1}|x_t)}{q(x_{t-1}|x_t, x_0)}}_{L_{t-1}} \right]$$
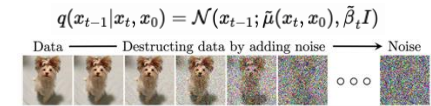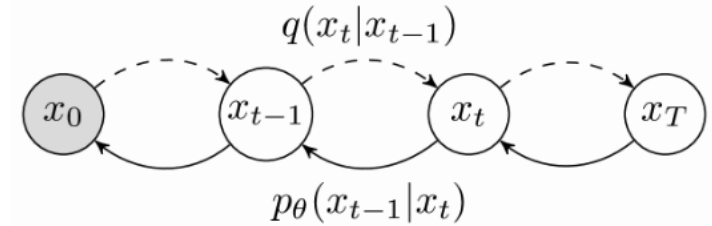
$$p_\theta(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \sigma_t^2 I)$$

$$\mathrm{KL}(P \| Q) = \frac{1}{2} \left[ \log \frac{|\Sigma_1|}{|\Sigma_0|} - d + \mathrm{tr}(\Sigma_1^{-1}\Sigma_0) + \boxed{(\mu_1 - \mu_0)^T \Sigma_1^{-1}(\mu_1 - \mu_0)} \right]$$

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

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

$$\alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s$$
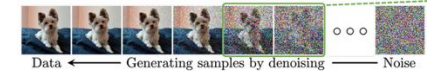
$x_0$



**The problem is every-time you need to calculate the target mean value**

$$L_{t-1} = \mathbb{E}_{t,x_t,x0} \left[ \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \| \right] + C$$

$$L_{t-1} = \mathbb{E}_{t,x_t,x_0} \left[ \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \| \right] + C$$

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

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

$$\alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s$$

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

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{(1 - \bar{\alpha}_t)}\boldsymbol{\epsilon} \qquad x_0 = \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}}{\sqrt{\bar{\alpha}_t}}$$

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

$$\tilde{u}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t} \left( \frac{x_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}}{\sqrt{\bar{\alpha}_t}} \right) + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} x_t = \boxed{\frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right)}$$

**For a given $x_t$, add a noise**

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

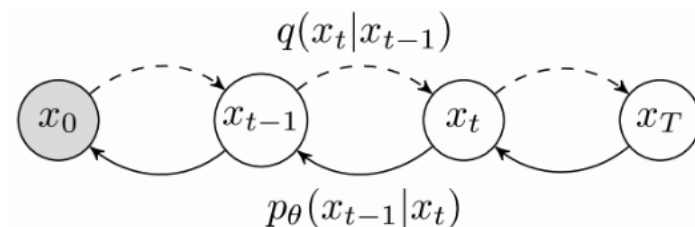$$= \mathbb{E}_{x_0,\epsilon,t} \left[ \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(x_t(x_0, \epsilon), t) \| \right]$$

**We do not need to calculate target mean but only do forward diffusion**

# Diffusion

$$L_{t-1} = \mathbb{E}_{t,x_t,x0} \left[ \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \| \right] + C$$

$$= \mathbb{E}_{x_0,\epsilon,t} \left[ \| \epsilon - \epsilon_\theta(x_t(x_0, \epsilon), t) \| \right]$$



[Link](#)

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

$$\tilde{\beta}_t = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t$$

$$\alpha_t := 1 - \beta_t \text{ and } \bar{\alpha}_t := \prod_{s=0}^{t} \alpha_s$$

---

**Algorithm 1** Training

---
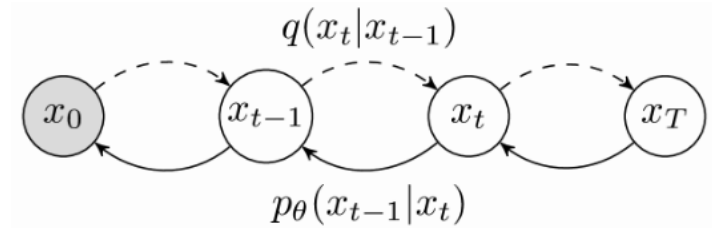
1: **repeat**
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
3:    $t \sim \text{Uniform}(\{1, \ldots, T\})$
4:    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
5:    Take gradient descent step on
      $\nabla_\theta \left\| \epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t) \right\|^2$
6: **until** converged

---

# Diffusion

$$L_{t-1} = \mathbb{E}_{t,x_t,x0} \left[ \frac{1}{2\sigma_t^2} \| \tilde{\mu}_t(x_t, x_0) - \mu_\theta(x_t, t) \| \right] + C$$

$$= \mathbb{E}_{x_0,\epsilon,t} \left[ \| \epsilon - \epsilon_\theta(x_t(x_0, \epsilon), t) \| \right]$$



[Link](#)

---

## Algorithm 2 Sampling

---

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

---

$$q(x_{t-1}|x_t, x_0) = \mathcal{N}(x_{t-1}; \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

$$\tilde{u}_t(x_t, x_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \left( \frac{x_t - \sqrt{1-\bar{\alpha}_t}\epsilon}{\sqrt{\bar{\alpha}_t}} \right) + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1-\bar{\alpha}_t}}\epsilon \right)$$
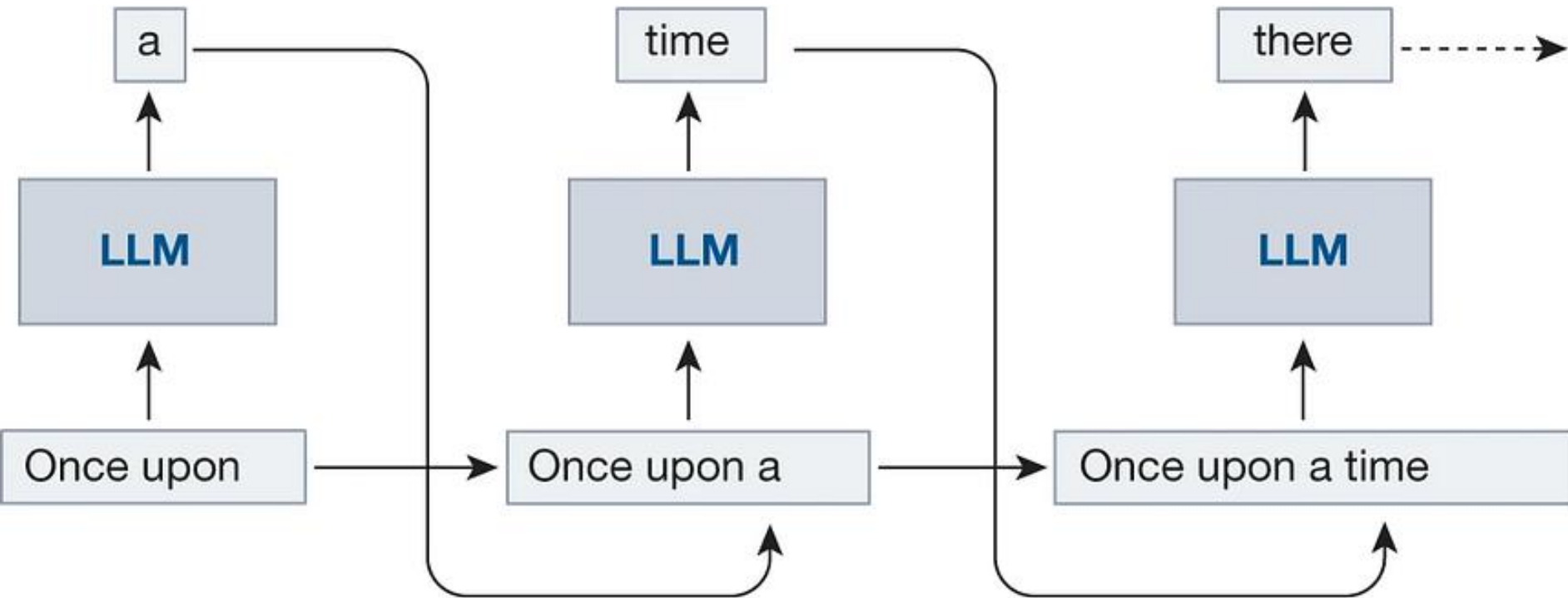
[Code Demo](#)

**How can we model the LLM generation under our framework?**

$$x_1^1 \rightarrow x_2^1 \rightarrow x_3^1 \rightarrow \cdots \rightarrow x_n^1$$

$$x_1^2 \rightarrow x_2^2 \rightarrow x_3^2 \rightarrow \cdots \rightarrow x_n^2$$

$$\cdots \qquad \cdots \qquad \cdots$$

$$x_1^S \rightarrow x_2^S \rightarrow x_3^S \rightarrow \cdots \rightarrow x_n^S$$

**Padding to be the same length**

$$P(X) = \prod_{s=1}^{|S|} P(X_s) = \prod_{s=1}^{|S|} P(X_1, X_2, \ldots, X_{l_S})$$

**Different sequences are independent**

$$= \prod_{s=1}^{|S|} \prod_{l=2}^{l_S} \boldsymbol{P(X_l | X_{1:l-1})}$$

**Given previously observed sequences, what is the probability of observing the ground-truth next token?**