

CS-E4891 Deep Generative Models

Lecture 4: Diffusion models

Harri Lähdesmäki

Department of Computer Science
Aalto University

April 14, 2025

Outline

- Denoising diffusion probabilistic model (Section 25.2 from (Murphy, 2023))
- Continuous-time diffusion models (Section 25.4 from (Murphy, 2023))
- Latent diffusion (Section 25.5 from (Murphy, 2023))
- Conditional diffusion models (Section 25.6 from (Murphy, 2023))
- Reading: parts of Sec. 25 from (Murphy, 2023), Sec. 18 from (Prince, 2023), and other literature (see references at the end)

Denoising diffusion probabilistic model

Diffusion models build on a simple idea:

- It is hard to convert pure noise to structured data
- It is easy to convert structured data into a noise

Denoising diffusion probabilistic model

Diffusion models build on a simple idea:

- It is hard to convert pure noise to structured data
- It is easy to convert structured data into a noise

Denoising diffusion probabilistic model (DDPM) can be thought of as a hierarchical VAE, where

- All latent layers $\mathbf{x}_{1:T} = (\mathbf{x}_1, \dots, \mathbf{x}_T)$ have the same dimensionality as data \mathbf{x}_0
- Encoder model (variational inference) is a simple and fixed linear Gaussian model
- Decoder model (generative model) is shared across layers

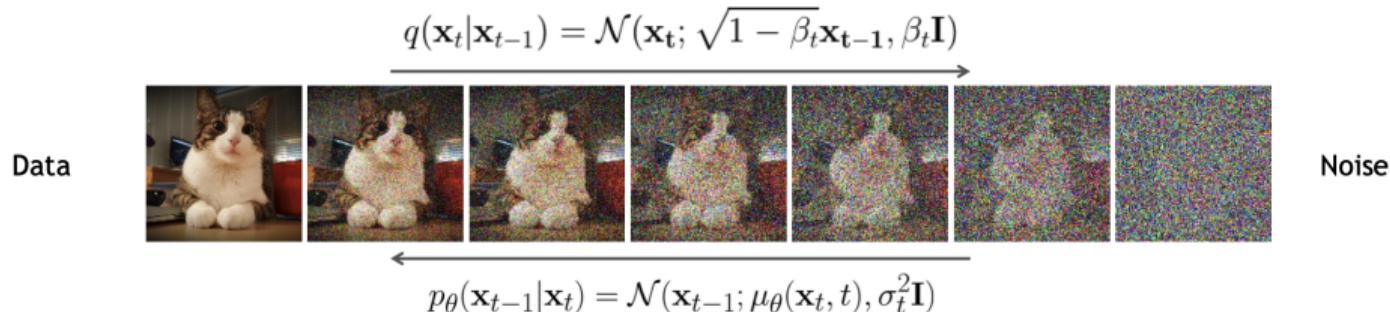


Illustration: Figure 25.1 from (Murphy, 2023)

DDPM: forward process

The forward diffusion process (=encoder) is a simple Gaussian noise model

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

where $\beta_t \in (0, 1)$ is defined by a noise schedule $\beta_1, \beta_2, \dots, \beta_T$

The joint distribution over latent variables is

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

DDPM: forward process

The forward diffusion process (=encoder) is a simple Gaussian noise model

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

where $\beta_t \in (0, 1)$ is defined by a noise schedule $\beta_1, \beta_2, \dots, \beta_T$

The joint distribution over latent variables is

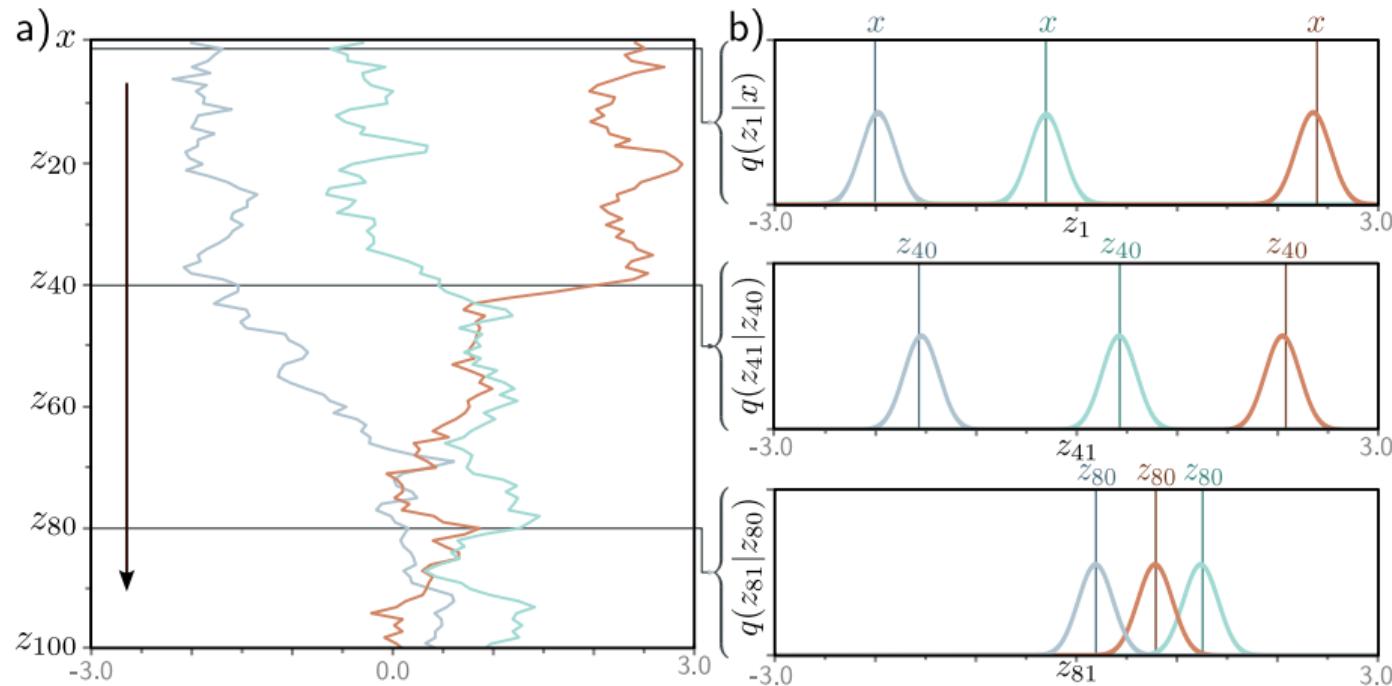
$$q(\mathbf{x}_{1:T}) = \prod_{t=1}^T q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$$

This is a linear Markov chain for which marginals can be computed analytically (see (Prince, 2023) for details)

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

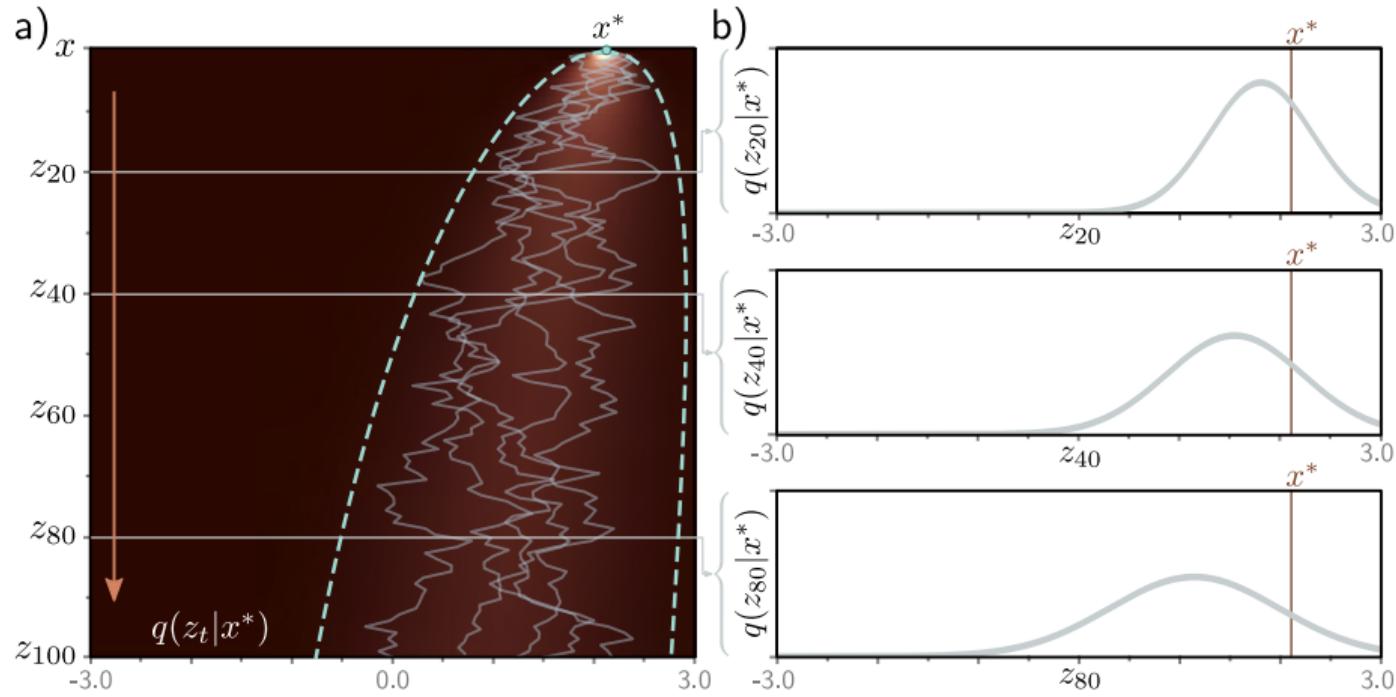
Choose noise schedule such that $\bar{\alpha}_T \approx 0$ and $q(\mathbf{x}_T \mid \mathbf{x}_0) \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$

DDPM: forward process



1-D Illustration of $q(\mathbf{x}_t | \mathbf{x}_{t-1})$ with $\beta_t = 0.03$. Figure 18.2 from (Prince, 2023)

DDPM: forward process



1-D Illustration of $q(x_t | x_0)$. Figure 18.3 from (Prince, 2023)

DDPM: reverse process

To generate data, we would ideally like to directly reverse the forward process $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

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

DDPM: reverse process

To generate data, we would ideally like to directly reverse the forward process $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

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

but we cannot analytically solve $q(\mathbf{x}_{t-1})$ or $q(\mathbf{x}_t)$ because

$$q(\mathbf{x}_t) = \int q(\mathbf{x}_t \mid \mathbf{x})p_{\mathcal{D}}(\mathbf{x})d\mathbf{x}$$

DDPM: reverse process

To generate data, we would ideally like to directly reverse the forward process $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$

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

but we cannot analytically solve $q(\mathbf{x}_{t-1})$ or $q(\mathbf{x}_t)$ because

$$q(\mathbf{x}_t) = \int q(\mathbf{x}_t \mid \mathbf{x}) p_{\mathcal{D}}(\mathbf{x}) d\mathbf{x}$$

If the input \mathbf{x}_0 was known, using properties of Gaussians (see e.g. (Prince, 2023) for details)

$$\begin{aligned} q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N}(\mathbf{x}_{t-1} \mid \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}) \\ \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) &= \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \\ \tilde{\beta}_t &= \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \end{aligned}$$

This does not serve us because \mathbf{x}_0 is not known when generating new data

DDPM: learnable reverse process

We can learn a reverse diffusion process (=decoder) that aims to invert the forward process

Generator can be chosen to have the following form

$$p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1} \mid \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)), \quad \text{where often} \quad \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$$

where $\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)$ is trainable neural net

DDPM: learnable reverse process

We can learn a reverse diffusion process (=decoder) that aims to invert the forward process

Generator can be chosen to have the following form

$$p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1} \mid \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t)), \quad \text{where often} \quad \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t) = \sigma_t^2 \mathbf{I}$$

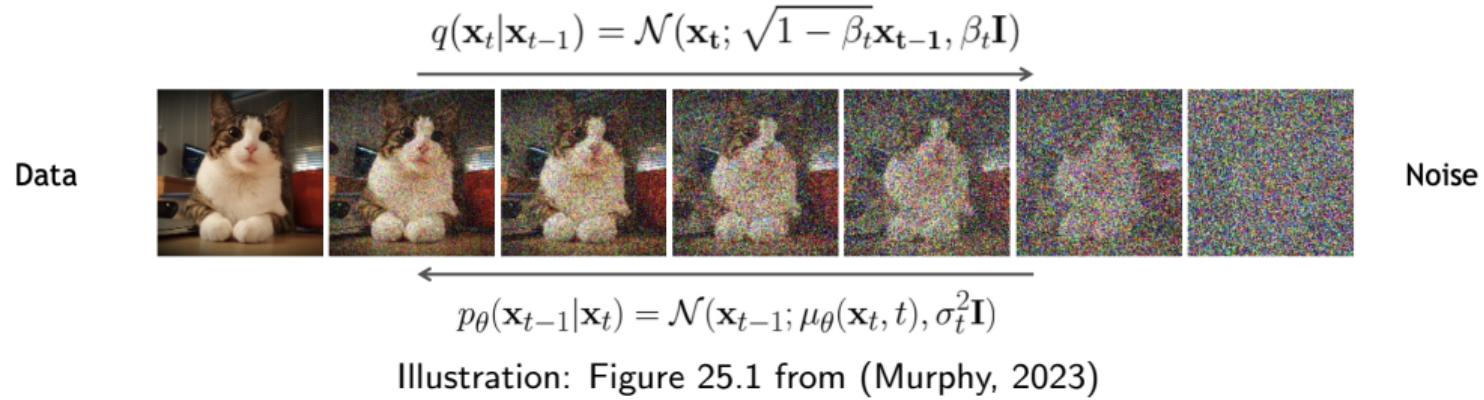
where $\boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t)$ is trainable neural net

The joint distribution of generated variables is

$$\begin{aligned} p_{\theta}(\mathbf{x}_{0:T}) &= p(\mathbf{x}_T)p_{\theta}(\mathbf{x}_{T-1} \mid \mathbf{x}_T) \cdots p_{\theta}(\mathbf{x}_0 \mid \mathbf{x}_1) \\ &= p(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) \end{aligned}$$

where $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

DDPM model components



DDPM: model fitting

We can fit the model using ELBO

$$\begin{aligned}\log p_{\boldsymbol{\theta}}(\mathbf{x}_0) &= \log \left[\int p_{\boldsymbol{\theta}}(\mathbf{x}_0, \mathbf{x}_{1:T}) d\mathbf{x}_{1:T} \right] \\ &= \log \left[\int \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) d\mathbf{x}_{1:T} \right] \\ &= \log \mathbb{E}_{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \right] \\ &\stackrel{\text{Jensen}}{\geq} \mathbb{E}_q \left[\log \left(\frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \right) \right] \\ &= \mathbb{E}_q \left[\log \left(p(\mathbf{x}_T) \prod_{t=1}^T \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})} \right) \right] \\ &= \mathbb{E}_q \left[\log p(\mathbf{x}_T) + \sum_{t=1}^T \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_t \mid \mathbf{x}_{t-1})} \right] \\ &= \mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}_0)\end{aligned}$$

DDPM: model fitting details

By the Markov property and Bayes's rule we have

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

DDPM: model fitting details

By the Markov property and Bayes's rule we have

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

By plugging $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$ into $\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}_0)$ we get

$$\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)} + \underbrace{\sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)}}_{* \text{ telescoping sum}} + \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_0 \mid \mathbf{x}_1)}{q(\mathbf{x}_1 \mid \mathbf{x}_0)} \right],$$

DDPM: model fitting details

By the Markov property and Bayes's rule we have

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

By plugging $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$ into $\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}_0)$ we get

$$\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)} + \underbrace{\sum_{t=2}^T \log \frac{q(\mathbf{x}_{t-1} \mid \mathbf{x}_0)}{q(\mathbf{x}_t \mid \mathbf{x}_0)}}_{* \text{ telescoping sum}} + \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_0 \mid \mathbf{x}_1)}{q(\mathbf{x}_1 \mid \mathbf{x}_0)} \right],$$

where telescoping sum reduces to

$$\begin{aligned} * &= \left(\log q(\mathbf{x}_1 \mid \mathbf{x}_0) - \log q(\mathbf{x}_2 \mid \mathbf{x}_0) \right) + \left(\log q(\mathbf{x}_2 \mid \mathbf{x}_0) - \log q(\mathbf{x}_3 \mid \mathbf{x}_0) \right) + \dots \\ &\quad + \left(\log q(\mathbf{x}_{T-1} \mid \mathbf{x}_0) - \log q(\mathbf{x}_T \mid \mathbf{x}_0) \right) \\ &= \log q(\mathbf{x}_1 \mid \mathbf{x}_0) - \log q(\mathbf{x}_T \mid \mathbf{x}_0) \end{aligned}$$

DDPM: model fitting details

ELBO reduces to

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta} \mid \mathbf{x}_0) &= \mathbb{E}_{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\log p(\mathbf{x}_T) + \sum_{t=2}^T \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)} + \log \frac{q(\mathbf{x}_1 \mid \mathbf{x}_0)}{q(\mathbf{x}_T \mid \mathbf{x}_0)} + \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_0 \mid \mathbf{x}_1)}{q(\mathbf{x}_1 \mid \mathbf{x}_0)} \right] \\ &= \mathbb{E}_{q(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\log \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T \mid \mathbf{x}_0)} + \sum_{t=2}^T \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t)}{q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)} + \log p_{\boldsymbol{\theta}}(\mathbf{x}_0 \mid \mathbf{x}_1) \right] \\ &= \underbrace{\mathbb{E}_{q(\mathbf{x}_1 \mid \mathbf{x}_0)} [\log p_{\boldsymbol{\theta}}(\mathbf{x}_0 \mid \mathbf{x}_1)]}_{L_0(\mathbf{x}_0)} - \sum_{t=2}^T \mathbb{E}_{q(\mathbf{x}_t \mid \mathbf{x}_0)} [\underbrace{D_{\text{KL}}(q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) \parallel p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} \mid \mathbf{x}_t))}_{L_{t-1}(\mathbf{x}_0)}] \\ &\quad + \underbrace{D_{\text{KL}}(q(\mathbf{x}_T \mid \mathbf{x}_0) \parallel p(\mathbf{x}_T))}_{L_T(\mathbf{x}_0)}\end{aligned}$$

All KL divergences can be computed analytically because all distributions are Gaussians

DDPM: model fitting details for $L_{t-1}(\mathbf{x}_0)$

Recall that the KL between two multivariate Gaussians is:

$$\begin{aligned} D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) &= \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right. \\ &\quad \left. - D + \log \left(\frac{\det(\boldsymbol{\Sigma}_2)}{\det(\boldsymbol{\Sigma}_1)} \right) \right] \end{aligned}$$

DDPM: model fitting details for $L_{t-1}(\mathbf{x}_0)$

Recall that the KL between two multivariate Gaussians is:

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right. \\ \left. - D + \log \left(\frac{\det(\boldsymbol{\Sigma}_2)}{\det(\boldsymbol{\Sigma}_1)} \right) \right]$$

Thus, $L_{t-1}(\mathbf{x}_0)$ term can be written as

$$D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \frac{1}{2\sigma_t^2} \|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)\|^2 + C \\ = \frac{1}{2\sigma_t^2} \left\| \underbrace{\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_t - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}_{\mathbb{E}[\mathbf{x}_{t-1} | \mathbf{x}_0, \mathbf{x}_t]} \right\|^2 + C$$

DDPM: model fitting details for $L_{t-1}(\mathbf{x}_0)$

Recall that the KL between two multivariate Gaussians is:

$$D_{\text{KL}}(\mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) || \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)) = \frac{1}{2} \left[\text{tr}(\boldsymbol{\Sigma}_2^{-1} \boldsymbol{\Sigma}_1) + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}_2^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right. \\ \left. - D + \log \left(\frac{\det(\boldsymbol{\Sigma}_2)}{\det(\boldsymbol{\Sigma}_1)} \right) \right]$$

Thus, $L_{t-1}(\mathbf{x}_0)$ term can be written as

$$D_{\text{KL}}(q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) || p_{\boldsymbol{\theta}}(\mathbf{x}_{t-1} | \mathbf{x}_t)) = \frac{1}{2\sigma_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)\|^2 + C \\ = \frac{1}{2\sigma_t^2} \left\| \underbrace{\frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_t - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}_{\mathbb{E}[\mathbf{x}_{t-1} | \mathbf{x}_0, \mathbf{x}_t]} \right\|^2 + C$$

Expectations can be estimated using Monte Carlo

$$\mathbb{E}[L_{t-1}(\mathbf{x}_0)] \approx \frac{1}{S} \frac{1}{2\sigma_t^2} \sum_{s=1}^S \left\| \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1-\bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} \mathbf{x}_t^{(s)} - \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) \right\|^2, \quad \text{where } \mathbf{x}_t^{(s)} \stackrel{\text{i.i.d.}}{\sim} q(\mathbf{x}_t | \mathbf{x}_0)$$

DDPM: reparameterization

Earlier we had that

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

which means

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

DDPM: reparameterization

Earlier we had that

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

which means

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

Insert the expression of \mathbf{x}_0 into the exact decoder model $q(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1} \mid \tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I})$ to get

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

This reparameterized model depends explicitly on \mathbf{x}_t and $\boldsymbol{\epsilon}$ (and only implicitly on \mathbf{x}_0)

DDPM: reparameterization

Previously we had a decoder $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1} \mid \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$ that predicts the mean of the denoised version of \mathbf{x}_{t-1} given its noisy input \mathbf{x}_t

Reparametrization allows to train a neural network model $\epsilon_{\theta}(\mathbf{x}_t, t)$ that, given \mathbf{x}_t , predicts the specific noise ϵ (sampled from $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$) that was added to \mathbf{x}_{t-1} such that it resulted in \mathbf{x}_t

DDPM: reparameterization

Previously we had a decoder $p_{\theta}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1} \mid \boldsymbol{\mu}_{\theta}(\mathbf{x}_t, t), \boldsymbol{\Sigma}_{\theta}(\mathbf{x}_t, t))$ that predicts the mean of the denoised version of \mathbf{x}_{t-1} given its noisy input \mathbf{x}_t

Reparametrization allows to train a neural network model $\epsilon_{\theta}(\mathbf{x}_t, t)$ that, given \mathbf{x}_t , predicts the specific noise ϵ (sampled from $q(\mathbf{x}_t \mid \mathbf{x}_{t-1})$) that was added to \mathbf{x}_{t-1} such that it resulted in \mathbf{x}_t

In other words, given noisy input \mathbf{x}_t , predict ϵ to get mean of \mathbf{x}_{t-1}

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

DDPM: reparameterized objective

The reparameterized loss

$$\mathbb{E}[L_{t-1}(\mathbf{x}_0)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\underbrace{\frac{\beta_t^2}{2\sigma_t^2 \alpha_t (1 - \bar{\alpha}_t)}}_{\lambda_t} \|\epsilon - \epsilon_{\theta}(\underbrace{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon}_{\mathbf{x}_t}, t)\|^2 \right]$$

DDPM: reparameterized objective

The reparameterized loss

$$\mathbb{E}[L_{t-1}(\mathbf{x}_0)] = \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[\underbrace{\frac{\beta_t^2}{2\sigma_t^2 \alpha_t(1 - \bar{\alpha}_t)}}_{\lambda_t} \|\epsilon - \epsilon_\theta(\underbrace{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)}_{\mathbf{x}_t}\|^2 \right]$$

Empirical evidence suggests that $\lambda_t = 1$ gives better results

A simplified loss (averaged over samples \mathbf{x}_0 and diffusion layers t)

$$L_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0 \sim p_{\mathcal{D}}(\mathbf{x}), \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), t \sim \text{Unif}(\{1, \dots, T\})} \left[\|\epsilon - \epsilon_\theta(\underbrace{\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)}_{\mathbf{x}_t}\|^2 \right]$$

DDPM: training algorithm, simple model

Algorithm 25.1: Training a DDPM model with L_{simple} .

```
1 while not converged do
2    $x_0 \sim q_0(x_0)$ 
3    $t \sim \text{Unif}(\{1, \dots, T\})$ 
4    $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5   Take gradient descent step on  $\nabla_{\theta} ||\epsilon - \epsilon_{\theta}(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)||^2$ 
```

Algorithm 25.1 from (Murphy, 2023)

DDPM: ancestral sampling algorithm

Algorithm 25.2: Sampling from a DDPM model.

```
1  $x_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2 foreach  $t = T, \dots, 1$  do
3    $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
4    $x_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \epsilon_\theta(x_t, t) \right) + \sigma_t \epsilon_t$ 
5 Return  $x_0$ 
```

Algorithm 25.2 from (Murphy, 2023)

UNet architecture

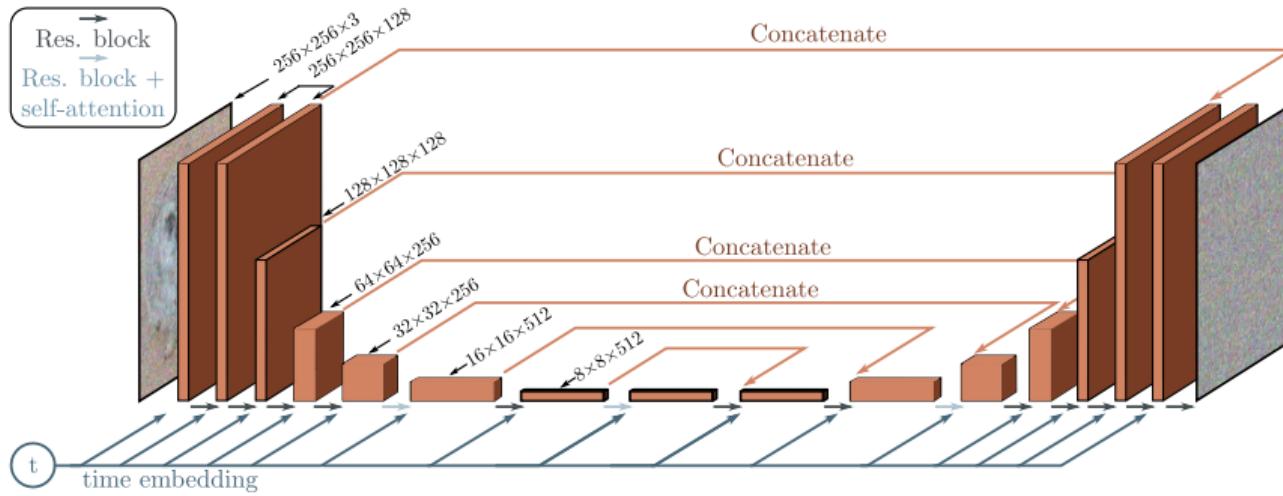


Figure 18.9 from (Prince, 2023)

DDPM: illustration

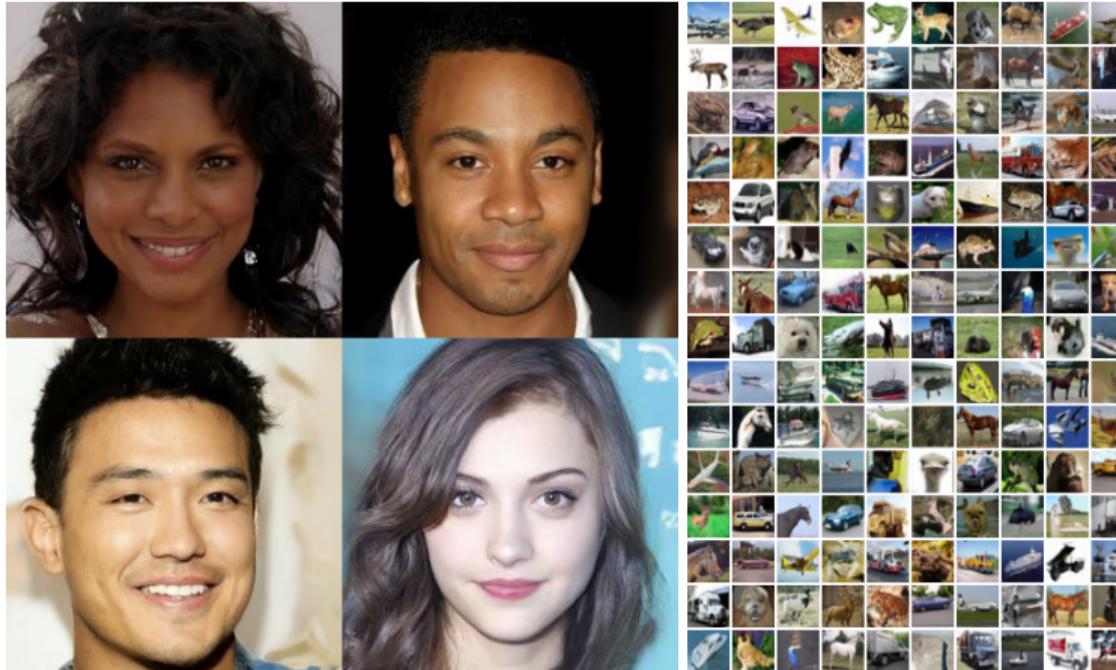


Figure from (Ho et al., 2020)

DDPM: interpolation

Example: DDPM model can be used e.g. to interpolate two data points x and x'

- $x_t \sim q(x_t | x)$ and $x'_t \sim q(x_t | x')$
- $\bar{x}_t = (1 - \lambda)x_t + \lambda x'_t$
- $x_0 \sim p_\theta(x_0 | \bar{x}_t)$

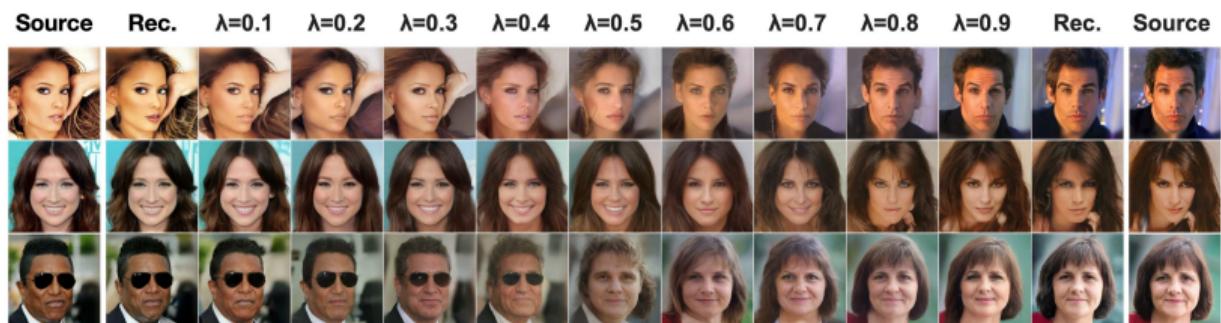
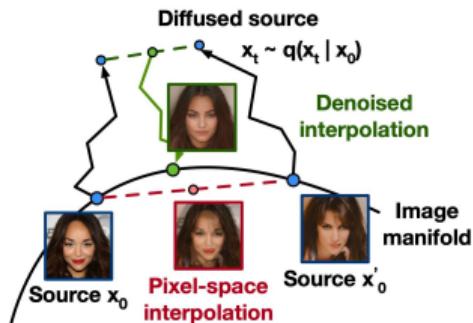


Figure from (Ho et al., 2020)

DDPM: noise schedule

The original DDPM model chose the noise schedule β_1, \dots, β_T from a set of constant, linear, and quadratic schedules, such that $q(\mathbf{x}_T | \mathbf{x}_0) \approx \mathcal{N}(\mathbf{0}, \mathbf{I})$, with $T = 1000$

A number of noise schedules have been proposed in the literature both with

- empirical evidence
- theoretical benefits

Continuous-time diffusion models: forward diffusion SDE

We can consider $t \in [0, T]$ to be continuous to obtain continuous-time diffusion models

Consider noise level β_t to be $\beta(t)\Delta t$, where Δt is a small diffusion time increment

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

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\epsilon} = \sqrt{1 - \beta(t)\Delta t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \cdot \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Continuous-time diffusion models: forward diffusion SDE

We can consider $t \in [0, T]$ to be continuous to obtain continuous-time diffusion models

Consider noise level β_t to be $\beta(t)\Delta t$, where Δt is a small diffusion time increment

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

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\epsilon} = \sqrt{1 - \beta(t)\Delta t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \cdot \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

For small Δt , by the first order Tailor approximation

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \frac{\beta(t)\Delta t}{2} \cdot \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \cdot \boldsymbol{\epsilon}$$

Continuous-time diffusion models: forward diffusion SDE

We can consider $t \in [0, T]$ to be continuous to obtain continuous-time diffusion models

Consider noise level β_t to be $\beta(t)\Delta t$, where Δt is a small diffusion time increment

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

$$\mathbf{x}_t = \sqrt{1 - \beta_t} \mathbf{x}_{t-1} + \sqrt{\beta_t} \boldsymbol{\epsilon} = \sqrt{1 - \beta(t)\Delta t} \cdot \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \cdot \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

For small Δt , by the first order Tailor approximation

$$\mathbf{x}_t = \mathbf{x}_{t-1} - \frac{\beta(t)\Delta t}{2} \cdot \mathbf{x}_{t-1} + \sqrt{\beta(t)\Delta t} \cdot \boldsymbol{\epsilon}$$

Thus

$$\frac{\mathbf{x}_t - \mathbf{x}_{t-1}}{\Delta t} = -\frac{\beta(t)}{2} \mathbf{x}_{t-1} + \frac{\sqrt{\beta(t)}}{\sqrt{\Delta t}} \boldsymbol{\epsilon}$$

and when $\Delta t \rightarrow 0$

$$\frac{d\mathbf{x}(t)}{dt} = -\frac{\beta(t)}{2} \mathbf{x}(t) + \sqrt{\beta(t)} \frac{d\mathbf{w}(t)}{dt}$$

we have a stochastic differential equation (SDE) model, where \mathbf{w} is Brownian motion

Continuous-time diffusion models: forward diffusion SDE

SDE models can be generally written as

$$d\mathbf{x} = \underbrace{\mathbf{f}(\mathbf{x}, t)dt}_{\text{deterministic}} + \underbrace{g(t)d\mathbf{w}}_{\text{stochastic}}$$

and can be numerically simulated (solved) e.g. by the Euler-Maruyama algorithm

$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t + g(t)\sqrt{\Delta t} \cdot \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Continuous-time diffusion models: forward diffusion SDE

SDE models can be generally written as

$$d\mathbf{x} = \underbrace{\mathbf{f}(\mathbf{x}, t)dt}_{\text{deterministic}} + \underbrace{g(t)d\mathbf{w}}_{\text{stochastic}}$$

and can be numerically simulated (solved) e.g. by the Euler-Maruyama algorithm

$$\mathbf{x}(t + \Delta t) = \mathbf{x}(t) + \mathbf{f}(\mathbf{x}(t), t)\Delta t + g(t)\sqrt{\Delta t} \cdot \epsilon, \quad \text{where } \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Continuous-time limit of the DDPM model is obtained with $\mathbf{f}(\mathbf{x}, t) = -\frac{1}{2}\beta(t)\mathbf{x}$ and $g(t) = \sqrt{\beta(t)}$

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w}$$

Continuous-time diffusion models: illustration

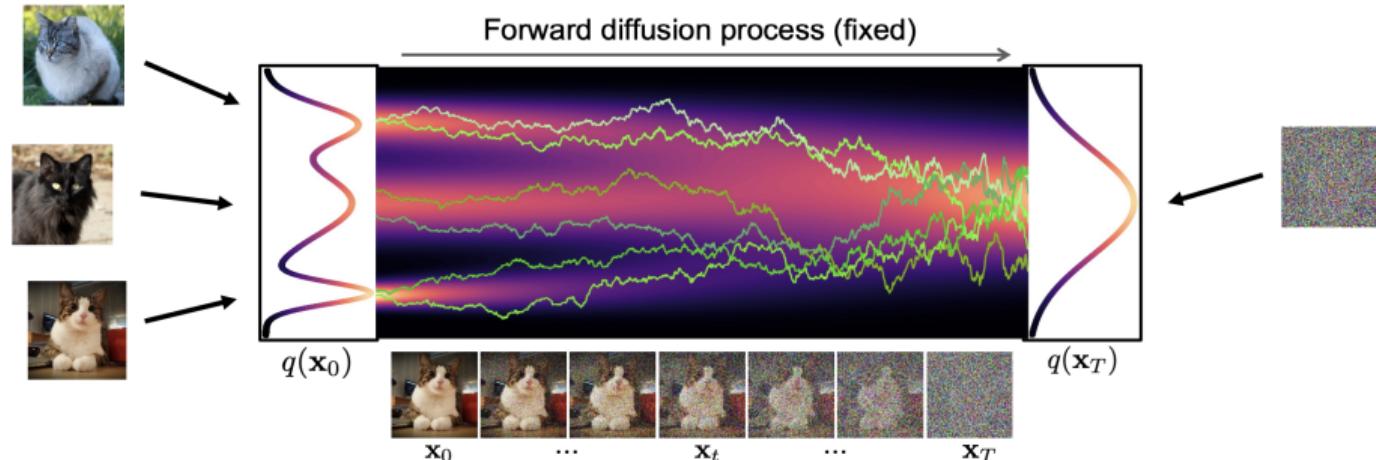


Figure 25.6 from (Murphy, 2023)

Continuous-time diffusion models: reverse diffusion SDE

As with discrete-time diffusion models, we need to revert the noising process

Reverse-time SDE is known and has the following form

$$d\boldsymbol{x} = [\boldsymbol{f}(\boldsymbol{x}_t, t) - g(t)^2 \nabla_{\boldsymbol{x}} \log q_t(\boldsymbol{x})]dt + g(t)d\bar{\boldsymbol{w}},$$

where $\bar{\boldsymbol{w}}$ is Brownian motion to the reverse diffusion time direction

Continuous-time diffusion models: reverse diffusion SDE

As with discrete-time diffusion models, we need to revert the noising process

Reverse-time SDE is known and has the following form

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}_t, t) - g(t)^2 \nabla_{\mathbf{x}} \log q_t(\mathbf{x})] dt + g(t) d\bar{\mathbf{w}},$$

where $\bar{\mathbf{w}}$ is Brownian motion to the reverse diffusion time direction

Reverse-time SDE model corresponding to the DDPM model is

$$d\mathbf{x}_t = \left[-\frac{1}{2} \beta(t) \mathbf{x} - \beta(t) \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t) \right] dt + \sqrt{\beta(t)} d\bar{\mathbf{w}}_t$$

We need to estimate the so-called score-function

$$\mathbf{s}_{\theta}(\mathbf{x}_t, t) \approx \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$$

Continuous-time diffusion models: reverse diffusion SDE

Training loss for continuous-time diffusion model

$$\mathcal{L} = \mathbb{E}_{t \sim \text{Unif}([0, T])} [\lambda(t) \mathbb{E}_{\mathbf{x}_0} \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} [\|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0)\|]]$$

where

- $\lambda(t)$ is a positive weighting function
- $\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$ can be computed analytically for simply (linear) forward diffusion models

Continuous-time diffusion models: reverse diffusion SDE

Training loss for continuous-time diffusion model

$$\mathcal{L} = \mathbb{E}_{t \sim \text{Unif}([0, T])} [\lambda(t) \mathbb{E}_{\mathbf{x}_0} \mathbb{E}_{\mathbf{x}_t | \mathbf{x}_0} [\|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t | \mathbf{x}_0)\|]]$$

where

- $\lambda(t)$ is a positive weighting function
- $\nabla_{\mathbf{x}_t} \log q_t(\mathbf{x}_t)$ can be computed analytically for simply (linear) forward diffusion models

Samples can be generated using the Euler-Maruyama algorithm

$$\mathbf{x}_{t-\Delta t} = \mathbf{x}_t + \frac{1}{2} \beta(t) [\mathbf{x}_t + 2\mathbf{s}_\theta(\mathbf{x}_t, t)] \Delta t + \sqrt{\beta(t) \Delta t} \cdot \boldsymbol{\epsilon}, \quad \text{where } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Continuous-time diffusion models: illustration

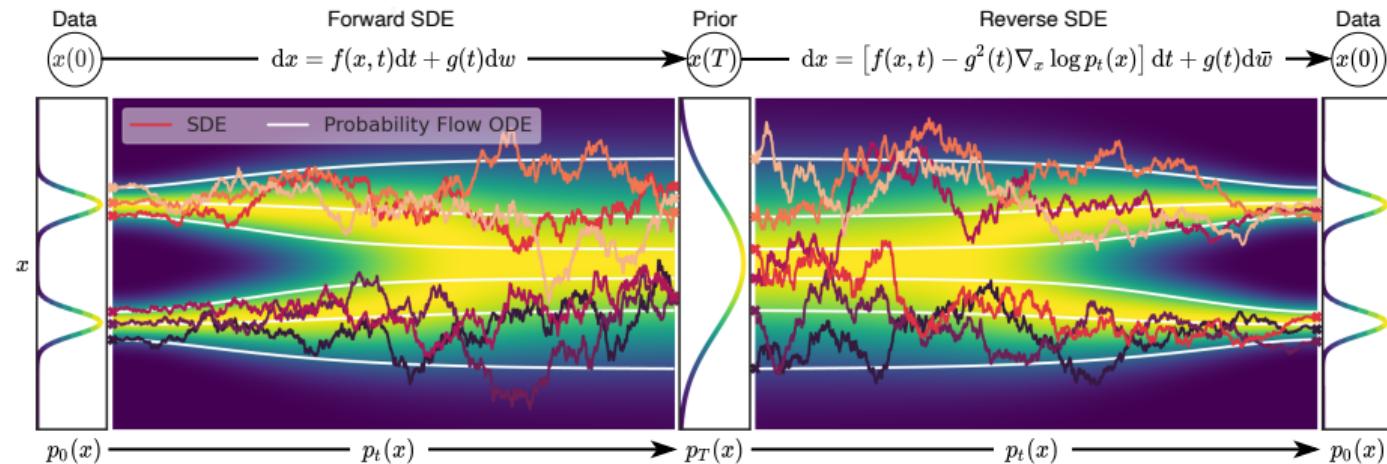


Figure from (Song et al., 2021)

Continuous-time diffusion models: illustration



Figure from (Song et al., 2021)

Latent diffusion models

Diffusion models are slow to train and use

- Model needs to take many small steps on high-dimensional data

Latent diffusion models (LDM) use a two-stage training scheme

- First train a VAE
- Then train a diffusion model to the embeddings

Used e.g. in the stable diffusion system

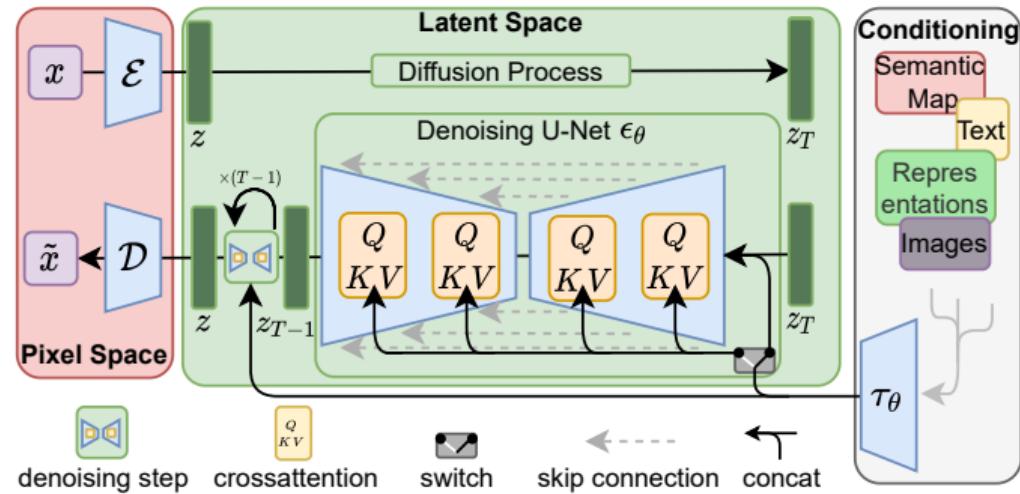


Figure 25.13 from (Murphy, 2023)

Conditional diffusion models: classifier guidance

Consider a case where data points \mathbf{x}

- Are associated with class labels \mathbf{c}
- We have access to a pre-trained classifier $p(\mathbf{c} | \mathbf{x})$

From Bayes theorem

$$\log p(\mathbf{x} | \mathbf{c}) = \log p(\mathbf{c} | \mathbf{x}) + \log p(\mathbf{x}) - \log p(\mathbf{c})$$

Conditional diffusion models: classifier guidance

Consider a case where data points \mathbf{x}

- Are associated with class labels \mathbf{c}
- We have access to a pre-trained classifier $p(\mathbf{c} | \mathbf{x})$

From Bayes theorem

$$\log p(\mathbf{x} | \mathbf{c}) = \log p(\mathbf{c} | \mathbf{x}) + \log p(\mathbf{x}) - \log p(\mathbf{c})$$

The score function becomes

$$\mathbf{s}_{\theta}(\mathbf{x}, t)' = \nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{c}) = \nabla_{\mathbf{x}} \log p(\mathbf{c} | \mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

A weighted score function can be defined as

$$\mathbf{s}_{\theta}(\mathbf{x}, t)' = \nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{c}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + w \cdot \nabla_{\mathbf{x}} \log p(\mathbf{c} | \mathbf{x})$$

Conditional diffusion models: classifier guidance

Consider a case where data points \mathbf{x}

- Are associated with class labels \mathbf{c}
- We have access to a pre-trained classifier $p(\mathbf{c} | \mathbf{x})$

From Bayes theorem

$$\log p(\mathbf{x} | \mathbf{c}) = \log p(\mathbf{c} | \mathbf{x}) + \log p(\mathbf{x}) - \log p(\mathbf{c})$$

The score function becomes

$$\mathbf{s}_{\theta}(\mathbf{x}, t)' = \nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{c}) = \nabla_{\mathbf{x}} \log p(\mathbf{c} | \mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{x})$$

A weighted score function can be defined as

$$\mathbf{s}_{\theta}(\mathbf{x}, t)' = \nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{c}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + w \cdot \nabla_{\mathbf{x}} \log p(\mathbf{c} | \mathbf{x})$$

Analogously, classifier guidance can be implemented in DDPM by modifying the sampling step

$$\mathbf{x}_{x-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{\beta_t}{\sqrt{1 - \bar{\alpha}_t}} \mathbf{e}_{\theta}(\mathbf{x}_t, t) \right) + w \cdot \nabla_{\mathbf{x}} \log p(\mathbf{c} | \mathbf{x}) + \sigma_t \mathbf{e}_t$$

Classifier guidance affects only sample generation, not diffusion model training

Conditional diffusion models: classifier guidance example

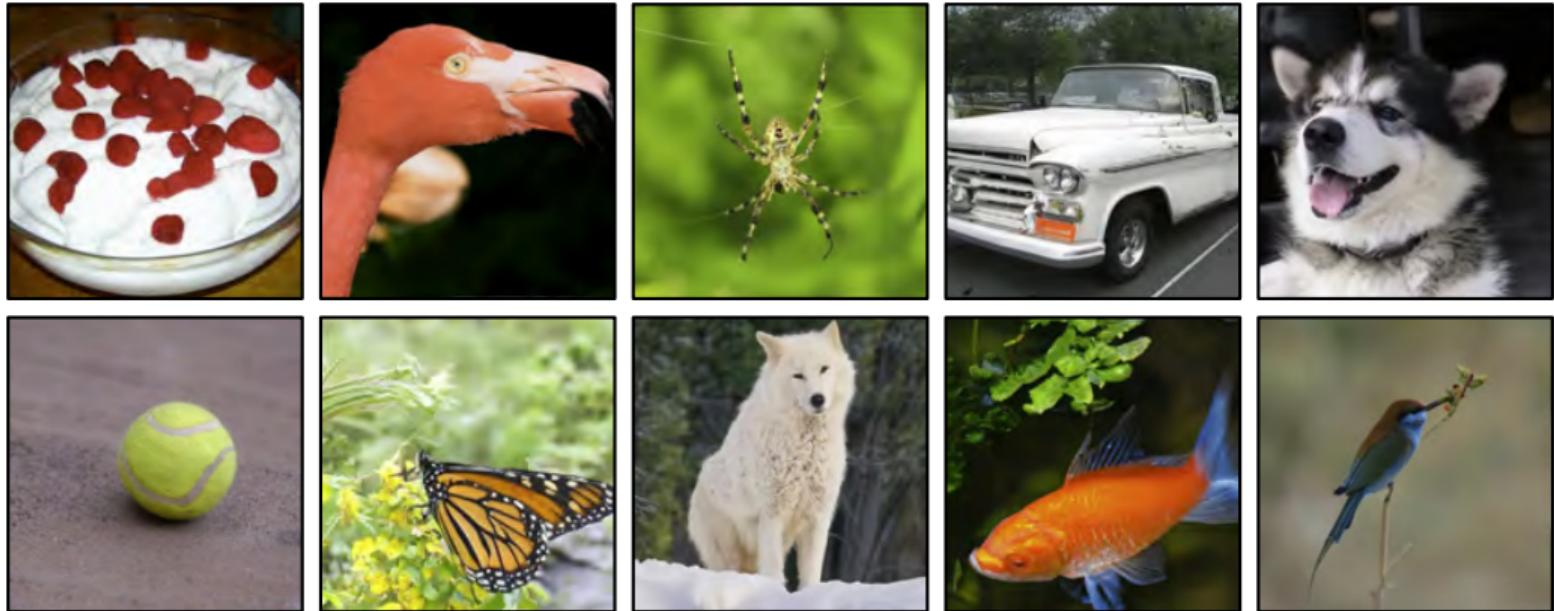


Figure 18.12 from (Prince, 2023)

Conditional diffusion models: classifier-free guidance

Derive a classifier from a generative model $p(\mathbf{c} \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathbf{c})p(\mathbf{c})}{p(\mathbf{x})}$ to get

$$\log p(\mathbf{c} \mid \mathbf{x}) = \log p(\mathbf{x} \mid \mathbf{c}) + \log p(\mathbf{c}) - \log p(\mathbf{x})$$

Need two generative models

- Conditional model $p(\mathbf{x} \mid \mathbf{c})$
- Unconditional model $p(\mathbf{x}) = p(\mathbf{x} \mid \mathbf{c} = \emptyset)$

Conditional diffusion models: classifier-free guidance

Derive a classifier from a generative model $p(\mathbf{c} \mid \mathbf{x}) = \frac{p(\mathbf{x} \mid \mathbf{c})p(\mathbf{c})}{p(\mathbf{x})}$ to get

$$\log p(\mathbf{c} \mid \mathbf{x}) = \log p(\mathbf{x} \mid \mathbf{c}) + \log p(\mathbf{c}) - \log p(\mathbf{x})$$

Need two generative models

- Conditional model $p(\mathbf{x} \mid \mathbf{c})$
- Unconditional model $p(\mathbf{x}) = p(\mathbf{x} \mid \mathbf{c} = \emptyset)$

Use this implicit classifier similarly as in classifier guidance

$$\begin{aligned}\nabla_{\mathbf{x}}[\log p(\mathbf{x} \mid \mathbf{c}) + w \log p(\mathbf{c} \mid \mathbf{x})] &= \nabla_{\mathbf{x}}[\log p(\mathbf{x} \mid \mathbf{c}) + w(\log p(\mathbf{x} \mid \mathbf{c}) - \log p(\mathbf{x}))] \\ &= \nabla_{\mathbf{x}}[(1 + w) \log p(\mathbf{x} \mid \mathbf{c}) - wp(\mathbf{x})]\end{aligned}$$

Conditional diffusion models

Conditional generation by training a diffusion model in (\mathbf{c}, \mathbf{x}) pairs

If the covariate \mathbf{c} is a scalar:

- \mathbf{c} can be mapped to an embedding vector
- Embedding can be incorporated into the neural network $\mu(\mathbf{x}_t, t | \mathbf{c})$, $\epsilon(\mathbf{x}_t, t | \mathbf{c})$ or $s_\theta(\mathbf{x}_t, t | \mathbf{c})$

If the covariate \mathbf{c} is another image:

- We can concatenate \mathbf{c} with \mathbf{x}_t as an extra set of image channels

If the covariate \mathbf{c} is a text prompt:

- We can embed \mathbf{c}
- Embedding can be incorporated into the neural network $\mu(\mathbf{x}_t, t | \mathbf{c})$, $\epsilon(\mathbf{x}_t, t | \mathbf{c})$ or $s_\theta(\mathbf{x}_t, t | \mathbf{c})$

UNet architecture

Embedding of covariate c can be e.g. shared across UNet layers similarly as time embedding

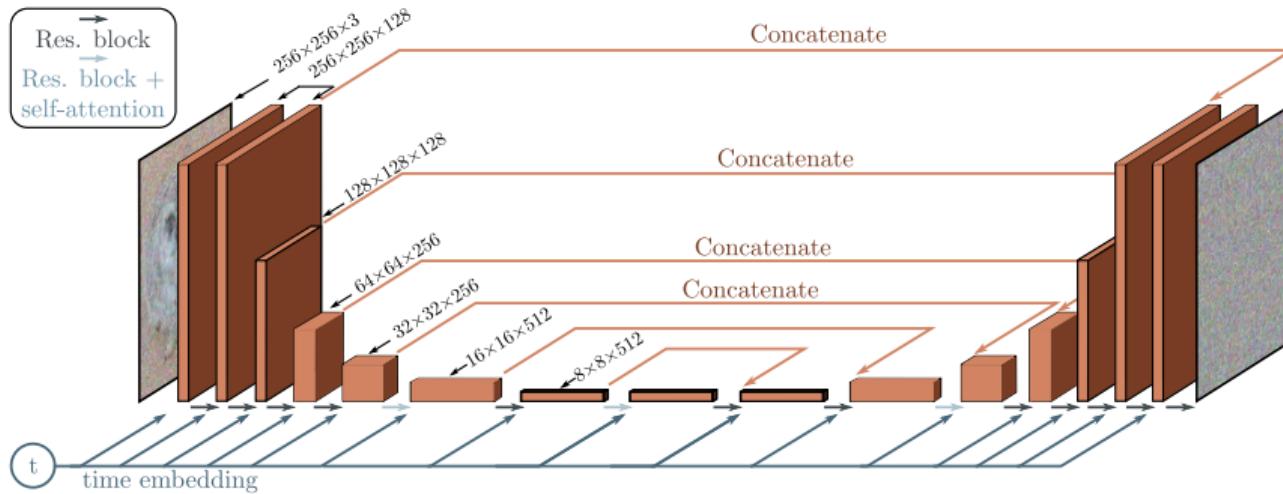


Figure 18.9 from (Prince, 2023)

Conditional diffusion models: example

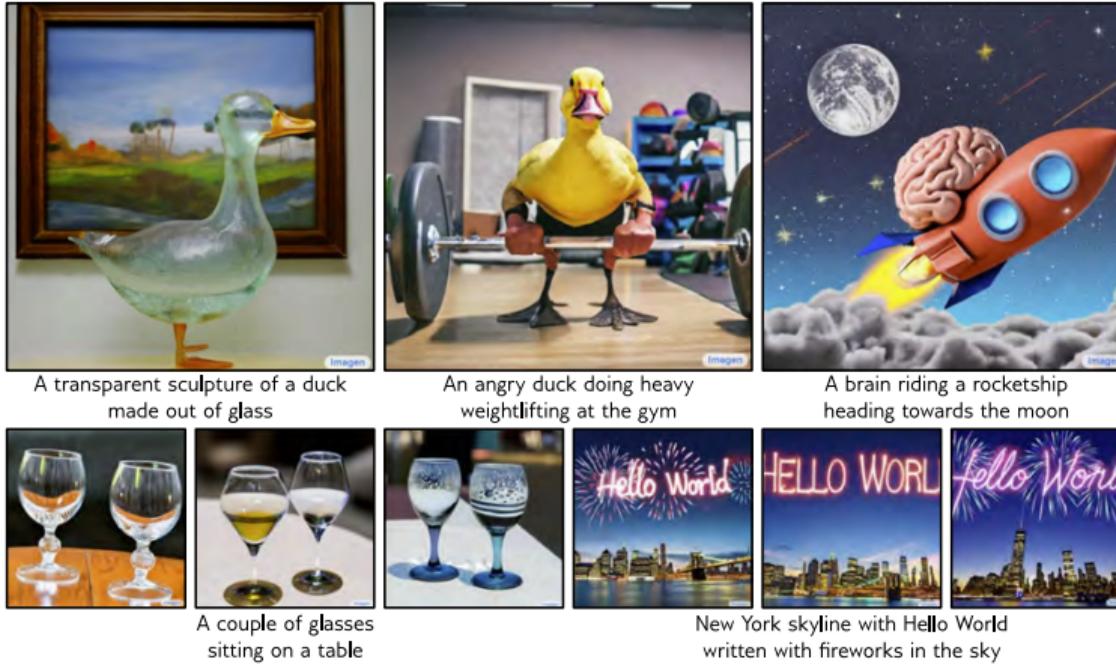


Figure 18.13 from (Prince, 2023)

References

- Ho J, Jain A, Abbeel P, "Denoising Diffusion Probabilistic Models," *NeurIPS*, 2020.
- Murphy K, Probabilistic Machine Learning: Advanced Topics, 2023.
- Prince SJD, Understanding Deep Learning, The MIT Press, 2023.
- Song Y, Sohl-Dickstein J, Kingma DP, Kumar A, Ermon S, Poole B, "Score-Based Generative Modeling through Stochastic Differential Equations," *ICLR*, 2021.