

# Denoising Diffusion Models (contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Plan today

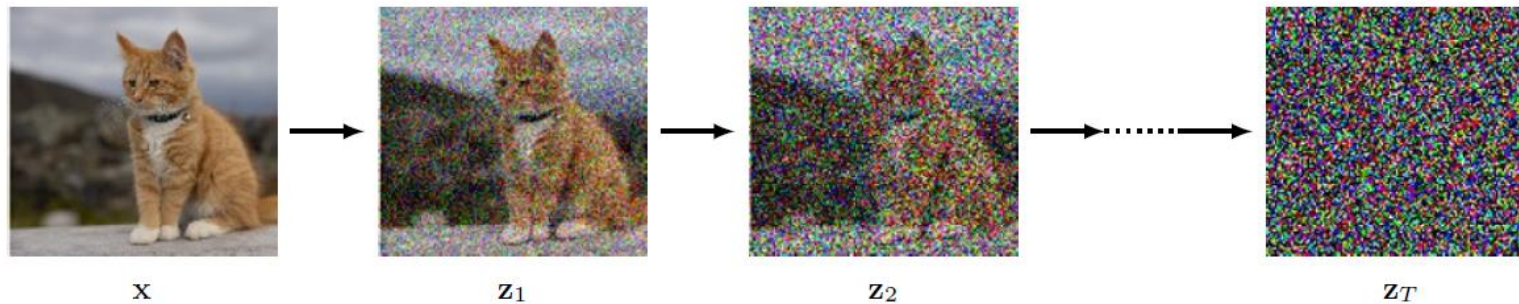
- Recap of denoising diffusion models
- Another perspective of diffusion using score based generative models
- Guided/conditional generation with diffusion models
  - Classifier guidance
  - Classifier-free guidance
- Latent diffusion models



# Denoising Diffusion Models

- Based on a forward process (adding noise) and a reverse process (denoising)

- The forward process as is follows  $q(\mathbf{z}_1, \dots, \mathbf{z}_t | \mathbf{x}) = q(\mathbf{z}_1 | \mathbf{x}) \prod_{\tau=2}^t q(\mathbf{z}_\tau | \mathbf{z}_{\tau-1})$



$$\mathbf{z}_1 = \sqrt{1 - \beta_1} \mathbf{x} + \sqrt{\beta_1} \epsilon_1$$

$$q(\mathbf{z}_1 | \mathbf{x}) = \mathcal{N}(\mathbf{z}_1 | \sqrt{1 - \beta_1} \mathbf{x}, \beta_1 \mathbf{I})$$

$$\mathbf{z}_t = \sqrt{1 - \beta_t} \mathbf{z}_{t-1} + \sqrt{\beta_t} \epsilon_t$$

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

$$\beta_t \in (0, 1)$$

Typically  
set by hand

$$\beta_1 < \beta_2 < \dots < \beta_T.$$

Called the  
"diffusion kernel"



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

$$\alpha_t = \prod_{\tau=1}^t (1 - \beta_\tau)$$

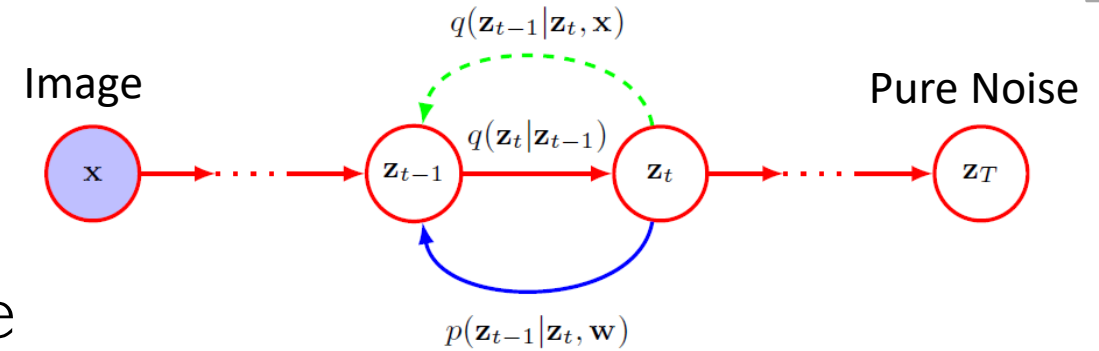
As  $T \rightarrow \infty$

$$q(\mathbf{z}_T | \mathbf{x}) = \mathcal{N}(\mathbf{z}_T | \mathbf{0}, \mathbf{I})$$



# Reversing the Diffusion

- Reversing the diffusion is like denoising
- Can use it to generate data from pure noise
- Reverse process will need  $q(\mathbf{z}_{t-1}|\mathbf{z}_t)$  but computing it is hard in general because



$$q(\mathbf{z}_{t-1}|\mathbf{z}_t) = \frac{q(\mathbf{z}_t|\mathbf{z}_{t-1})q(\mathbf{z}_{t-1})}{q(\mathbf{z}_t)} \quad \text{where} \quad q(\mathbf{z}_{t-1}) = \int q(\mathbf{z}_{t-1}|\mathbf{x})p(\mathbf{x}) d\mathbf{x}$$

Requires integrating over the data distribution  $p(\mathbf{x})$  which is not easy

- Conditioning also on  $\mathbf{x}$  makes computing the reverse process distribution tractable

Equals  $q(\mathbf{z}_t|\mathbf{z}_{t-1})$

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

But this denoising only holds for training images (because it is conditioned on  $\mathbf{x}$ )

Thus we will also learn "parallel" distributions of the form  $p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w})$  which approximate  $q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})$

If we can now estimate  $\mathbf{w}$  from training data, we can use  $p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w})$  to generate new synthetic data starting with pure noise

$$\mathbf{m}_t(\mathbf{x}, \mathbf{z}_t) = \frac{(1 - \alpha_{t-1})\sqrt{1 - \beta_t}\mathbf{z}_t + \sqrt{\alpha_{t-1}}\beta_t\mathbf{x}}{1 - \alpha_t}$$

$$\sigma_t^2 = \frac{\beta_t(1 - \alpha_{t-1})}{1 - \alpha_t}$$



# Reversing the Diffusion

- The joint distribution of data and latents

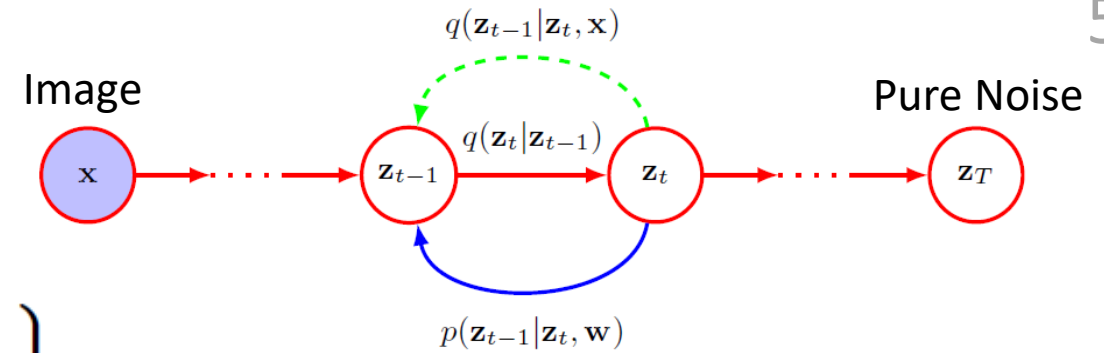
$$p(\mathbf{x}, \mathbf{z}_1, \dots, \mathbf{z}_T | \mathbf{w}) = p(\mathbf{z}_T) \left\{ \prod_{t=2}^T p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w}) \right\} p(\mathbf{x} | \mathbf{z}_1, \mathbf{w})$$

- Let's assume  $p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w}) = \mathcal{N}(\mathbf{z}_{t-1} | \boldsymbol{\mu}(\mathbf{z}_t, \mathbf{w}, t), \beta_t \mathbf{I})$
- The true joint distribution of the latents given  $\mathbf{x}$

$$q(\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_T | \mathbf{x}) = q(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x})$$

- To estimate  $\mathbf{w}$ , we can maximize the ELBO defined as

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_q \left[ \ln \frac{p(\mathbf{z}_T) \left\{ \prod_{t=2}^T p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w}) \right\} p(\mathbf{x} | \mathbf{z}_1, \mathbf{w})}{q(\mathbf{z}_1 | \mathbf{x}) \prod_{t=2}^T q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x})} \right] = \mathbb{E}_q \left[ \ln p(\mathbf{z}_T) + \sum_{t=2}^T \ln \frac{p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})}{q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x})} - \ln q(\mathbf{z}_1 | \mathbf{x}) + \ln p(\mathbf{x} | \mathbf{z}_1, \mathbf{w}) \right]$$



Note that  $\boldsymbol{\mu}$  represents the denoising model (e.g., a neural net) which denoises  $\mathbf{z}_t$  to produce  $\mathbf{z}_{t-1}$

This term is just like the VAE reconstruction error term (can approximate it using samples of  $\mathbf{z}_1$  from  $q(\mathbf{z}_1 | \mathbf{x})$ )

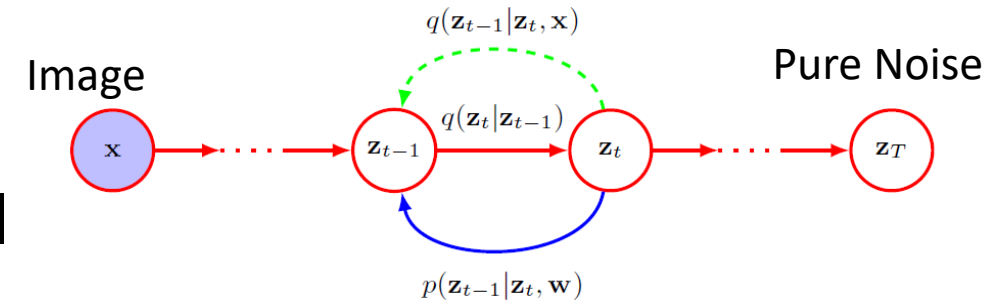
From ELBO definition  $\mathbb{E}_q \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z})} \right]$

Also note that unlike VI, here we aren't estimating the  $q$  distribution

First and third terms don't contain  $\mathbf{w}$  so can be ignored when maximizing the ELBO



# ELBO (contd)



- Recall the ELBO for the denoising diffusion model

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_q \left[ \ln p(\mathbf{z}_T) + \sum_{t=2}^T \ln \frac{p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})}{q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x})} - \ln q(\mathbf{z}_1 | \mathbf{x}) + \ln p(\mathbf{x} | \mathbf{z}_1, \mathbf{w}) \right]$$

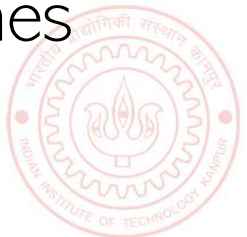
- Ignoring terms that don't depend on  $\mathbf{w}$  and using  $q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x}) = \frac{q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) q(\mathbf{z}_t | \mathbf{x})}{q(\mathbf{z}_{t-1} | \mathbf{x})}$

$$\ln \frac{p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})}{q(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{x})} = \ln \frac{p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})}{q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} + \ln \frac{q(\mathbf{z}_{t-1} | \mathbf{x})}{q(\mathbf{z}_t | \mathbf{x})} \implies \mathcal{L}(\mathbf{w}) = \mathbb{E}_q \left[ \sum_{t=2}^T \ln \frac{p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})}{q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x})} + \ln p(\mathbf{x} | \mathbf{z}_1, \mathbf{w}) \right]$$

- The ELBO becomes 
$$\mathcal{L}(\mathbf{w}) = \underbrace{\int q(\mathbf{z}_1 | \mathbf{x}) \ln p(\mathbf{x} | \mathbf{z}_1, \mathbf{w}) d\mathbf{z}_1}_{\text{reconstruction term}} - \underbrace{\sum_{t=2}^T \int \text{KL}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) \| p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})) q(\mathbf{z}_t | \mathbf{x}) d\mathbf{z}_t}_{\text{consistency terms}}$$

- Since both distributions in the KL divergence term are Gaussians, it becomes

$$\text{KL}(q(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{x}) \| p(\mathbf{z}_{t-1} | \mathbf{z}_t, \mathbf{w})) = \frac{1}{2\beta_t} \|\mathbf{m}_t(\mathbf{x}, \mathbf{z}_t) - \boldsymbol{\mu}(\mathbf{z}_t, \mathbf{w}, t)\|^2 + \text{const}$$



# Predicting the noise

- The KL terms in the ELBO are of the form

$$\text{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})||p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w})) = \frac{1}{2\beta_t} \|\mathbf{m}_t(\mathbf{x}, \mathbf{z}_t) - \boldsymbol{\mu}(\mathbf{z}_t, \mathbf{w}, t)\|^2 + \text{const}$$

Network which gives the mean of the denoised  $\mathbf{z}_{t-1}$

- Note that

$$\mathbf{x} = \frac{1}{\sqrt{\alpha_t}}\mathbf{z}_t - \frac{\sqrt{1-\alpha_t}}{\sqrt{\alpha_t}}\boldsymbol{\epsilon}_t \quad \longrightarrow \quad \mathbf{m}_t(\mathbf{x}, \mathbf{z}_t) = \frac{1}{\sqrt{1-\beta_t}} \left\{ \mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\boldsymbol{\epsilon}_t \right\}$$

From the definition of  $\mathbf{m}_t(\mathbf{x}, \mathbf{z}_t)$

- Instead of learning  $\boldsymbol{\mu}(\mathbf{z}_t, \mathbf{w}, t)$ , we will learn a **noise predictor**  $\mathbf{g}(\mathbf{z}_t, \mathbf{w}, t)$  s.t.

Using the same form as  $\mathbf{m}_t$  with  $\mathbf{g}(\mathbf{z}_t, \mathbf{w}, t)$  trying to predict  $\boldsymbol{\epsilon}_t$

$$\boldsymbol{\mu}(\mathbf{z}_t, \mathbf{w}, t) = \frac{1}{\sqrt{1-\beta_t}} \left\{ \mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}}\mathbf{g}(\mathbf{z}_t, \mathbf{w}, t) \right\}$$

- Therefore

$$\begin{aligned} \text{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})||p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w})) &= \frac{\beta_t}{2(1-\alpha_t)(1-\beta_t)} \|\mathbf{g}(\mathbf{z}_t, \mathbf{w}, t) - \boldsymbol{\epsilon}_t\|^2 + \text{const} \\ &= \frac{\beta_t}{2(1-\alpha_t)(1-\beta_t)} \|\mathbf{g}(\sqrt{\alpha_t}\mathbf{x} + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}_t, \mathbf{w}, t) - \boldsymbol{\epsilon}_t\|^2 + \text{const} \end{aligned}$$

Basically, we are now just predicting the noise  $\boldsymbol{\epsilon}_t$  using the neural network  $\mathbf{g}(\mathbf{z}_t, \mathbf{w}, t)$





# Predicting the noise

- We basically had the following

$$\text{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})||p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w})) = \frac{\beta_t}{2(1-\alpha_t)(1-\beta_t)} \|\mathbf{g}(\sqrt{\alpha_t}\mathbf{x} + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}_t, \mathbf{w}, t) - \boldsymbol{\epsilon}_t\|^2 + \text{const}$$

- The reconstruction error part in the ELBO can also be written as noise prediction

$$\ln p(\mathbf{x}|\mathbf{z}_1, \mathbf{w}) = -\frac{1}{2\beta_1} \|\mathbf{x} - \boldsymbol{\mu}(\mathbf{z}_1, \mathbf{w}, 1)\|^2 + \text{const.} = -\frac{1}{2(1-\beta_1)} \|\mathbf{g}(\mathbf{z}_1, \mathbf{w}, 1) - \boldsymbol{\epsilon}_1\|^2 + \text{const}$$

- Ignoring the constants in front of the squared error terms above, the ELBO becomes

Empirically found to give improved performance

Pick an example  $\mathbf{x}$  randomly, generate a corruption  $\mathbf{z}_t$  by sampling  $\boldsymbol{\epsilon}_t$  and make a gradient based update to  $\mathbf{w}$

Can optimize using stochastic optimization

$$\mathcal{L}(\mathbf{w}) = \underbrace{\int q(\mathbf{z}_1|\mathbf{x}) \ln p(\mathbf{x}|\mathbf{z}_1, \mathbf{w}) d\mathbf{z}_1}_{\text{reconstruction term}} - \underbrace{\sum_{t=2}^T \int \text{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})||p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w}))q(\mathbf{z}_t|\mathbf{x}) d\mathbf{z}_t}_{\text{consistency terms}}$$

$$= -\sum_{t=1}^T \|\mathbf{g}(\sqrt{\alpha_t}\mathbf{x} + \sqrt{1-\alpha_t}\boldsymbol{\epsilon}_t, \mathbf{w}, t) - \boldsymbol{\epsilon}_t\|^2$$





# Denoising Diffusion Model: The Training Algo

- The overall training algo is as follows

**Input:** Training data  $\mathcal{D} = \{\mathbf{x}_n\}$

Noise schedule  $\{\beta_1, \dots, \beta_T\}$

**Output:** Network parameters  $\mathbf{w}$

---

**for**  $t \in \{1, \dots, T\}$  **do**

$\alpha_t \leftarrow \prod_{\tau=1}^t (1 - \beta_\tau)$  // Calculate alphas from betas

**end for**

**repeat**

$\mathbf{x} \sim \mathcal{D}$  // Sample a data point

$t \sim \{1, \dots, T\}$  // Sample a point along the Markov chain

$\epsilon \sim \mathcal{N}(\epsilon | \mathbf{0}, \mathbf{I})$  // Sample a noise vector

$\mathbf{z}_t \leftarrow \sqrt{\alpha_t} \mathbf{x} + \sqrt{1 - \alpha_t} \epsilon$  // Evaluate noisy latent variable

$\mathcal{L}(\mathbf{w}) \leftarrow \|\mathbf{g}(\mathbf{z}_t, \mathbf{w}, t) - \epsilon\|^2$  // Compute loss term

    Take optimizer step

**until** converged

**return**  $\mathbf{w}$



# Denoising Diffusion Model: Generation

- Using the training model, we can now generate data as follows

**Input:** Trained denoising network  $g(\mathbf{z}, \mathbf{w}, t)$

Noise schedule  $\{\beta_1, \dots, \beta_T\}$

**Output:** Sample vector  $\mathbf{x}$  in data space

---

$\mathbf{z}_T \sim \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$  // Sample from final latent space

**for**  $t \in T, \dots, 2$  **do**

$\alpha_t \leftarrow \prod_{\tau=1}^t (1 - \beta_\tau)$  // Calculate alpha

// Evaluate network output

$\mu(\mathbf{z}_t, \mathbf{w}, t) \leftarrow \frac{1}{\sqrt{1-\beta_t}} \left\{ \mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} g(\mathbf{z}_t, \mathbf{w}, t) \right\}$

$\epsilon \sim \mathcal{N}(\epsilon|\mathbf{0}, \mathbf{I})$  // Sample a noise vector

$\mathbf{z}_{t-1} \leftarrow \mu(\mathbf{z}_t, \mathbf{w}, t) + \sqrt{\beta_t} \epsilon$  // Add scaled noise

**end for**

$\mathbf{x} = \frac{1}{\sqrt{1-\beta_1}} \left\{ \mathbf{z}_1 - \frac{\beta_1}{\sqrt{1-\alpha_1}} g(\mathbf{z}_1, \mathbf{w}, t) \right\}$  // Final denoising step

**return**  $\mathbf{x}$

Generation can be slow because it requires several steps

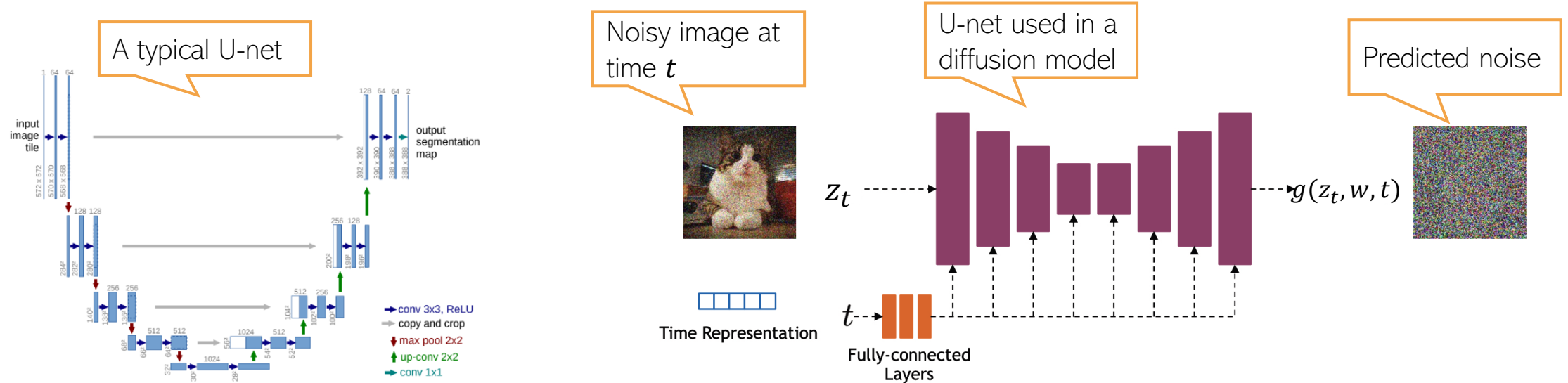
Reducing the number of steps is an active area of research

One such approach is **DDIM** (denoising diffusion implicit model) which relaxes the Markov assumption in the noise process



# Noise Predictor Network

- A “U-net” model (a neural net) is commonly used as the noise predictor network



- An embedding (positional embedding) of the time-step  $t$  is fed into the residual blocks of the U-net architecture



# Score based deep generative models

- For a probability distribution  $p(\mathbf{x})$  its **score function** is defined as

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

Note: Here, this gradient is w.r.t.  $\mathbf{x}$  and not w.r.t. the parameters of the distribution

- Assuming  $p(\mathbf{x})$  as a target distribution, we can use SGLD to generate data samples as

$$\mathbf{x}_t = \mathbf{x}_{t+1} + \frac{\delta}{2} \nabla_{\mathbf{x}} \log p(\mathbf{x}_t) + \sqrt{\delta} \epsilon_t \quad \text{where } \epsilon_t \sim \mathcal{N}(0, I)$$

- But doing so requires the score function  $s(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$
- Since  $p(\mathbf{x})$  itself is not known, how do get the score function  $s(\mathbf{x})$ ?
- We can train a neural network to model the score function
- The score based approach is also helpful in **“guided” or conditional generation**
  - Example: Want to generate  $\mathbf{x}$  while conditioning on some signal  $\mathbf{c}$  (e.g., class label or textual description of the input to be generated)



# Diffusion Models via Score Matching

- For a probability distribution  $p(\mathbf{x})$  its score function is defined as

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

Note: Here, this gradient is w.r.t.  $\mathbf{x}$  and not w.r.t. the parameters of the distribution

- Learning this score function is equivalent to learning the distribution  $p(\mathbf{x})$
- We can parameterize the score function as  $s(\mathbf{x}) = s(\mathbf{x}, \mathbf{w})$  and define a loss function

$$J(\mathbf{w}) = \frac{1}{2} \int \|s(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{x}} \ln p(\mathbf{x})\|^2 p(\mathbf{x}) d\mathbf{x}$$

Can define it as a neural network

- The distribution  $p(\mathbf{x})$  isn't known but we only have a dataset  $\mathcal{D}$  of  $N$  samples

A discrete distribution represented by the  $N$  samples from the dataset

However, this is non-differentiable and thus can't use it in the minimization of  $J(\mathbf{w})$

$$p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta(\mathbf{x} - \mathbf{x}_n)$$



# Score Matching

- Instead of  $p_{\mathcal{D}}(\mathbf{x})$ , we define a smooth distribution

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

One option is to define it as a Gaussian  $\mathcal{N}(\mathbf{z}|\mathbf{x}, \sigma^2 I)$

- Using this  $q_{\sigma}(\mathbf{z})$  instead of  $p_{\mathcal{D}}(\mathbf{x})$ , we can define the “score loss” function as

$$\begin{aligned} J(\mathbf{w}) &= \frac{1}{2} \int \|\mathbf{s}(\mathbf{z}, \mathbf{w}) - \nabla_{\mathbf{z}} \ln q_{\sigma}(\mathbf{z})\|^2 q_{\sigma}(\mathbf{z}) d\mathbf{z} \\ &= \frac{1}{2} \iint \|\mathbf{s}(\mathbf{z}, \mathbf{w}) - \nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}, \sigma)\|^2 q(\mathbf{z}|\mathbf{x}, \sigma) p(\mathbf{x}) d\mathbf{z} d\mathbf{x} + \text{const.} \end{aligned}$$

- Using the  $N$  samples from the dataset, the empirical loss will be

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N \int \|\mathbf{s}(\mathbf{z}, \mathbf{w}) - \nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}_n, \sigma)\|^2 q(\mathbf{z}|\mathbf{x}_n, \sigma) d\mathbf{z} + \text{const}$$



# Score Matching

- Recall that the score loss function is

$$J(\mathbf{w}) = \frac{1}{2} \iint \|\mathbf{s}(\mathbf{z}, \mathbf{w}) - \nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}, \sigma)\|^2 q(\mathbf{z}|\mathbf{x}, \sigma) p(\mathbf{x}) d\mathbf{z} d\mathbf{x} + \text{const}$$

- Choosing  $q(\mathbf{z}|\mathbf{x}, \sigma) = \mathcal{N}(\mathbf{z}|\mathbf{x}, \sigma^2 I)$ , we get

$$\nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}, \sigma) = -\frac{1}{\sigma} \boldsymbol{\epsilon} \quad \text{where} \quad \boldsymbol{\epsilon} = \mathbf{x} - \mathbf{z}$$

- Note the similarity with denoising diffusion model where

$$q(\mathbf{z}_t|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\sqrt{\alpha_t}\mathbf{x}, (1 - \alpha_t)I) \quad \longrightarrow \quad \nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}, \sigma) = -\frac{1}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon}$$

- The score loss function measures the difference b/w predicted score  $\mathbf{s}(\mathbf{z}, \mathbf{w})$  and noise
- Note that the score function  $\mathbf{s}(\mathbf{z}, \mathbf{w})$  plays a similar role as noise predictor  $g(\mathbf{z}, \mathbf{w}, t)$  in denoising diffusion model we saw earlier
- Careful selection of the noise variance  $\sigma^2$  is important in this approach





# Noise Variance in Score Matching

- Success of score matching depends on the estimate of score function  $s(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$

$$J(\mathbf{w}) = \frac{1}{2} \int \|\mathbf{s}(\mathbf{x}, \mathbf{w}) - \nabla_{\mathbf{x}} \ln p(\mathbf{x})\|^2 p(\mathbf{x}) d\mathbf{x}$$

$$J(\mathbf{w}) = \frac{1}{2} \iint \|\mathbf{s}(\mathbf{z}, \mathbf{w}) - \nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}, \sigma)\|^2 q(\mathbf{z}|\mathbf{x}, \sigma) p(\mathbf{x}) d\mathbf{z} d\mathbf{x} + \text{const}$$

- In regions where  $p(\mathbf{x})$  is small/zero, the estimate  $s(\mathbf{z}, \mathbf{w})$  may not be reliable
- Recall that, in score matching, we typically use  $q(\mathbf{z}|\mathbf{x}, \sigma) = \mathcal{N}(\mathbf{z}|\mathbf{x}, \sigma^2 I)$
- Using the appropriate  $\sigma^2$  is critical
  - Using large  $\sigma^2$  means we won't have small/zero values for  $q(\mathbf{z}|\mathbf{x})$  but also high distortion
  - Very small  $\sigma^2$  means  $q(\mathbf{z}|\mathbf{x})$  is close to  $p(\mathbf{x})$
  - We can choose a series of variances  $\sigma_1^2 < \sigma_2^2 < \dots < \sigma_L^2$  and use the following loss function

$\lambda(i)$  is the weighting coefficient for model  $i$

$$\frac{1}{2} \sum_{i=1}^L \lambda(i) \int \|\mathbf{s}(\mathbf{z}, \mathbf{w}, \sigma_i^2) - \nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}_n, \sigma_i)\|^2 q(\mathbf{z}|\mathbf{x}_n, \sigma_i) d\mathbf{z}$$

This gives us  $L$  score matching based diffusion models with different variances

We can run SGLD where we use a few steps of each in a sequences  $L, L-1, L-2, \dots, 2, 1$



# Diffusion Models and SDE

- Stochastic Differential Equations (SDE) define a continuous-time process
- Denoising diffusion model and score matching models are like discretization of the continuous-time SDE

- The forward SDE is written as 
$$d\mathbf{z} = \underbrace{\mathbf{f}(\mathbf{z}, t) dt}_{\text{drift}} + \underbrace{g(t) d\mathbf{v}}_{\text{diffusion}}$$

- The corresponding reverse SDE can be written as

This is like the score function

$$d\mathbf{z} = \left\{ \mathbf{f}(\mathbf{z}, t) - g^2(t) \nabla_{\mathbf{z}} \ln p(\mathbf{z}) \right\} dt + g(t) d\mathbf{v}$$

The corresponding ODE for the SDE reverse process

$$\frac{d\mathbf{z}}{dt} = \mathbf{f}(\mathbf{z}, t) - \frac{1}{2} g^2(t) \nabla_{\mathbf{z}} \ln p(\mathbf{z})$$

- We can solve SDE by discretizing time
  - For equal-size time steps, we get the Langevin dynamics based equations for the updates
- SDE connection is helpful in designing fast reverse process for diffusion models
  - For example, we can leverage the ODE corresponding to the SDE for faster sampling



# Guided Diffusion

- Often we want to generate data based on some “reference” conditioning signal, e.g.,
  - Images of a specific class (class-conditional generation)
  - Images based on some textual description
  - High resolution image using a low-resolution image (image “super-resolution”)



Conditioning signal: “stained glass window of a panda eating bamboo”



Conditioning signal: Low-res image on the left

- Denoting the data as  $\mathbf{x}$  and the conditioning signal as  $\mathbf{c}$ , we want to learn  $p(\mathbf{x}|\mathbf{c})$



# Classifier Guidance

- Assume we have an already training classifier of the form  $p(\mathbf{c}|\mathbf{x})$
- We can then define the score function of a conditional diffusion model as

$$\begin{aligned}\nabla_{\mathbf{x}} \ln p(\mathbf{x}|\mathbf{c}) &= \nabla_{\mathbf{x}} \ln \left\{ \frac{p(\mathbf{c}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{c})} \right\} \\ &= \nabla_{\mathbf{x}} \ln p(\mathbf{x}) + \nabla_{\mathbf{x}} \ln p(\mathbf{c}|\mathbf{x})\end{aligned}$$

- We can also control the contribution of the classifier by defining the score function as

$$\text{score}(\mathbf{x}, \mathbf{c}, \lambda) = \nabla_{\mathbf{x}} \ln p(\mathbf{x}) + \lambda \nabla_{\mathbf{x}} \ln p(\mathbf{c}|\mathbf{x})$$

Score function of unconditional diffusion model

Classifier guidance term

- Large  $\lambda$  will encourage generation of  $\mathbf{x}$  which respects the conditioning signal  $\mathbf{c}$
- However, this approach requires a classifier trained on noisy images



# Classifier-free Guidance

- Recall the score function in classifier guidance method

$$\text{score}(\mathbf{x}, \mathbf{c}, \lambda) = \nabla_{\mathbf{x}} \ln p(\mathbf{x}) + \lambda \nabla_{\mathbf{x}} \ln p(\mathbf{c}|\mathbf{x})$$

- To eliminate the classifier term  $\nabla_{\mathbf{x}} \log p(\mathbf{c}|\mathbf{x})$ , use the fact that

$$\begin{aligned} \nabla_{\mathbf{x}} \ln p(\mathbf{x}|\mathbf{c}) &= \nabla_{\mathbf{x}} \ln \left\{ \frac{p(\mathbf{c}|\mathbf{x})p(\mathbf{x})}{p(\mathbf{c})} \right\} \\ &= \nabla_{\mathbf{x}} \ln p(\mathbf{x}) + \nabla_{\mathbf{x}} \ln p(\mathbf{c}|\mathbf{x}) \end{aligned}$$

- Thus we can rewrite the score function as follows

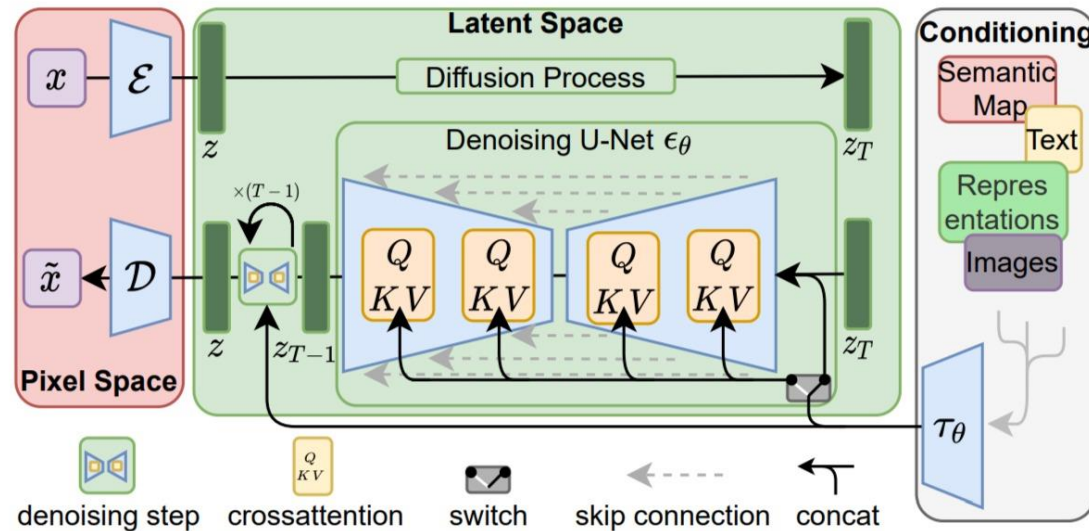
$$\text{score}(\mathbf{x}, \mathbf{c}, \lambda) = \lambda \nabla_{\mathbf{x}} \ln p(\mathbf{x}|\mathbf{c}) + (1 - \lambda) \nabla_{\mathbf{x}} \ln p(\mathbf{x})$$

- No need to train a separate classifier  $p(\mathbf{c}|\mathbf{x})$
- Also, no need to train both  $p(\mathbf{x})$  and  $p(\mathbf{x}|\mathbf{c})$ 
  - Just train  $p(\mathbf{x}|\mathbf{c})$  using a score-function based approach and use  $p(\mathbf{x}|\mathbf{c} = 0) = p(\mathbf{x})$



# Latent Diffusion Models (LDM)

- Defines diffusion process in a latent space instead of in data (e.g., pixel) space
- The popular “Stable Diffusion” is based on LDM
- Diffusion process in a low-dim latent space is also more efficient computationally



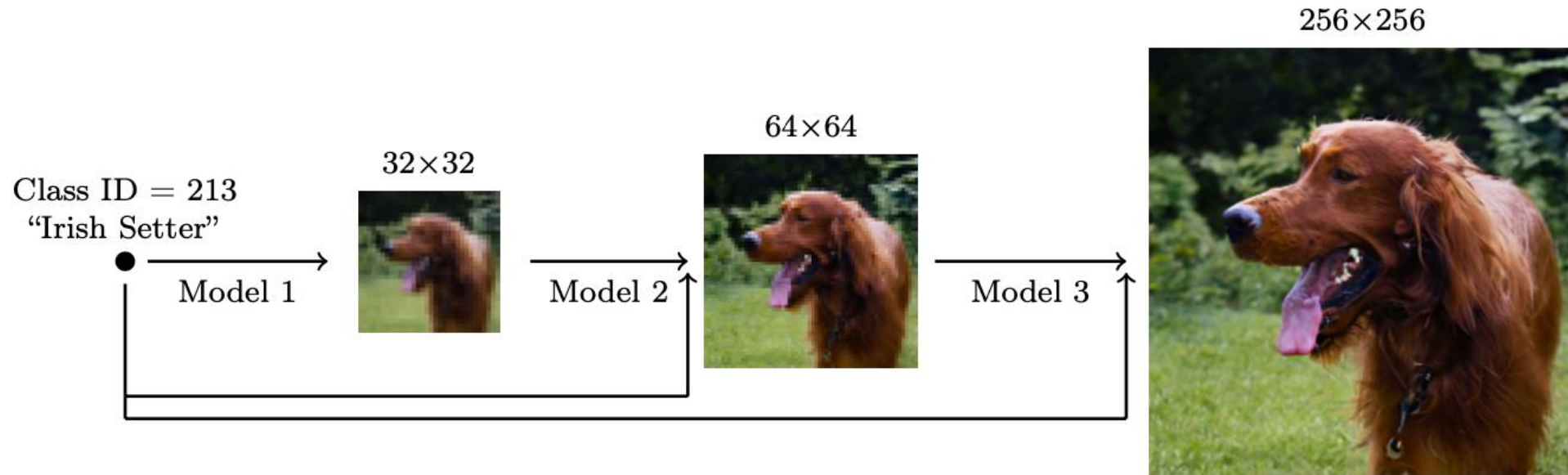
- Can also condition the generation of other modalities such as text





# Cascaded Diffusion Models

- Useful for generating high-resolution images using conditioning



- Cascaded approach is usually better than a direct generation of high-resolution image
  - Smaller model size
  - Learning gradual transformations is easier than a direct transformation





# Summary

- Diffusion Models (denoising diffusion models, score based models, etc) are currently the best performing methods
- A lot of ongoing work on diffusion models, e.g.,
  - Improving quality of generation
  - Speeding-up generation
  - Combining them with other generative models (e.g., large language models)

