

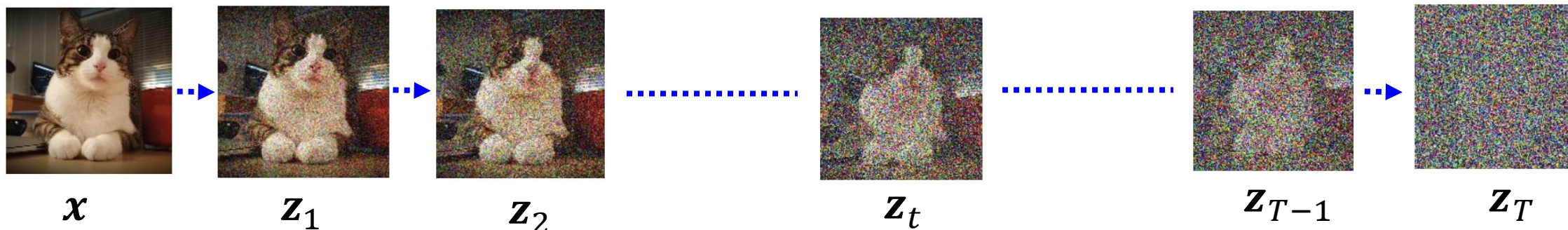
# Deep Generative Models (contd)

CS772A: Probabilistic Machine Learning

Piyush Rai

# Denoising Diffusion Models

- Consider gradually corrupting an image ( $\mathbf{z}_0 = \mathbf{x}$ ) till it becomes **pure noise** ( $\mathbf{z}_T$ )



- Each step  $\mathbf{z}_{t-1} \rightarrow \mathbf{z}_t$  is a pre-defined Gaussian perturbation (**forward process**)

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

$$\mathbf{z}_t = \sqrt{1 - \beta_t} \mathbf{z}_{t-1} + \sqrt{\beta_t} \boldsymbol{\epsilon} \quad (\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}))$$

$$\beta_t \in (0, 1) \quad \text{and} \quad \beta_1 < \beta_2 < \dots < \beta_{T-1} < \beta_T$$

Usually pre-defined but  
can also be learned

Imp: Thus we can also **compute**  
 **$\mathbf{z}_t$**  from  **$\mathbf{x}$**  directly in a single step

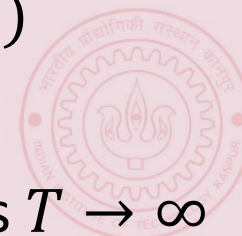
implies

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

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

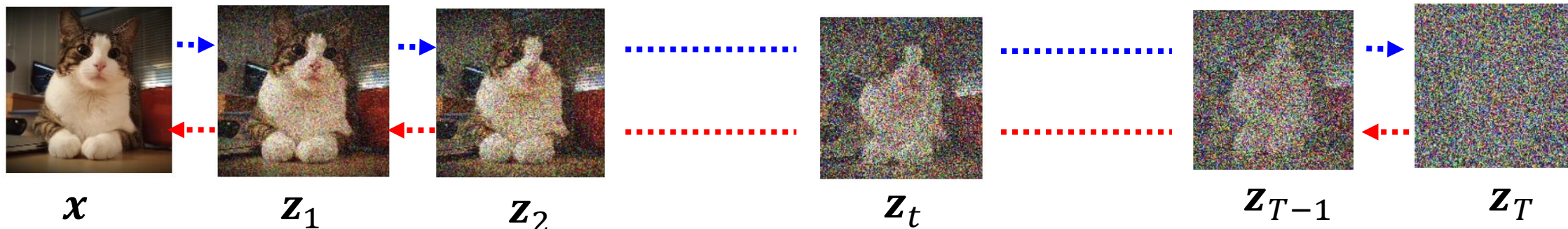
$$\mathbf{z}_t = \sqrt{\alpha_t} \mathbf{x} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}$$

$$q(\mathbf{z}_T | \mathbf{x}) = \mathcal{N}(\mathbf{z}_T | \mathbf{0}, \mathbf{I}) \quad \text{as } T \rightarrow \infty$$



# Generating Data by Reversing Diffusion

- Reversing the diffusion (red arrows) would enable generating data from pure noise



- To reverse the diffusion, we need the distribution of  $z_{t-1}$  given  $z_t$ , i.e.,  $q(z_{t-1}|z_t)$

The denoising distribution

Intractable because  $q(z_t)$  and  $q(z_{t-1})$  are difficult to compute

$$q(z_{t-1}|z_t) = \frac{q(z_{t-1})q(z_t|z_{t-1})}{q(z_t)}$$

Since the true data distribution  $p(x)$  is not known, we can't compute this integral

$$q(z_t) = \int q(z_t|x)p(x)dx$$



# Towards a Tractable Reverse Diffusion

- Although  $q(\mathbf{z}_{t-1}|\mathbf{z}_t)$  isn't tractable, the following distribution is tractable

$$\begin{aligned} q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x}) &= \frac{q(\mathbf{z}_{t-1}|\mathbf{x}) q(\mathbf{z}_t|\mathbf{z}_{t-1}, \mathbf{x})}{q(\mathbf{z}_t|\mathbf{x})} \\ &= \frac{q(\mathbf{z}_{t-1}|\mathbf{x}) q(\mathbf{z}_t|\mathbf{z}_{t-1})}{q(\mathbf{z}_t|\mathbf{x})} \end{aligned}$$

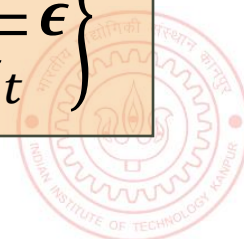
- Reason:  $q(\mathbf{z}_{t-1}|\mathbf{x})$  and  $q(\mathbf{z}_t|\mathbf{z}_{t-1})$  are Gaussians, so  $q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})$  is Gaussian

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

Using  $\mathbf{x} = \frac{1}{\sqrt{\alpha_t}} \mathbf{z}_t - \frac{\sqrt{1-\alpha_t}}{\sqrt{\alpha_t}} \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

$$\begin{aligned} m(\mathbf{x}, \mathbf{z}_t) &= \frac{(1 - \alpha_{t-1})\sqrt{1 - \beta_t} \mathbf{z}_t + \sqrt{\alpha_{t-1}} \beta_t \mathbf{x}}{1 - \alpha_t} \\ \sigma_t^2 &= \frac{\beta_t(1 - \alpha_{t-1})}{1 - \alpha_t} \end{aligned}$$

$$= \frac{1}{\sqrt{1 - \beta_t}} \left\{ \mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\epsilon} \right\}$$



# Towards a Tractable Reverse Diffusion

- We saw that the reverse diffusion distribution is the Gaussian

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

where

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

Issue: At generation time, we don't have  $\mathbf{x}$  (the goal is to generate  $\mathbf{x}$  which is only available for training data) so we can't use  $m(\mathbf{x}, \mathbf{z}_t)$  at generation time since it depends on  $\mathbf{x}$



- Let's approximate  $q(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{x})$  by another Gaussian that doesn't depend on  $\mathbf{x}$

$$p(\mathbf{z}_{t-1}|\mathbf{z}_t, \mathbf{w}) = \mathcal{N}(\mathbf{z}_{t-1} | \mu(\mathbf{z}_t, \mathbf{w}, t), \Sigma(\mathbf{z}_t, \mathbf{w}, t))$$

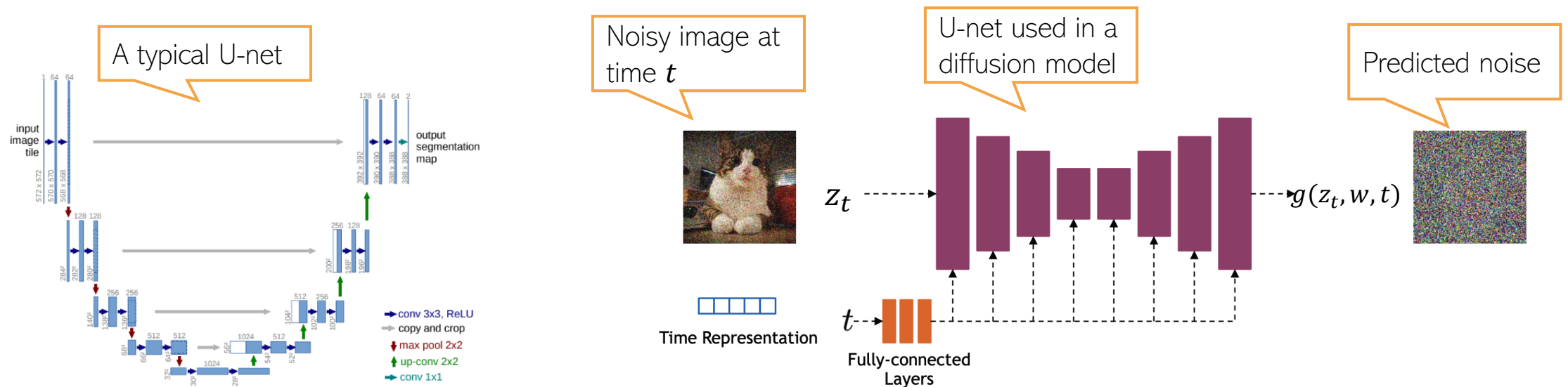
- Usually,  $\Sigma(\mathbf{z}_t, \mathbf{w}, t)$  is chosen to be spherical. A popular choice:  $\Sigma(\mathbf{z}_t, \mathbf{w}, t) = \beta_t \mathbf{I}$
- The mean  $\mu(\mathbf{z}_t, \mathbf{w}, t)$  is defined to mimic the form of  $m(\mathbf{x}, \mathbf{z}_t)$

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

Our goal will be to model  $g()$  as a neural network which predicts the noise  $\boldsymbol{\epsilon}$  which was added to  $\mathbf{x}$  to get its noisy version  $\mathbf{z}_t$

# 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



# 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  $\mathbf{g}(\mathbf{z}, \mathbf{w}, t)$

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

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

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

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

**end for**

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



# 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

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

- 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). Score function:  $\nabla_{\mathbf{x}} \log p(\mathbf{x}|\mathbf{c})$



# Diffusion Models via Score Matching

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

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

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

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

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

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

- Let's define a "smoothened" version of the distribution  $p(\mathbf{x})$

Basically a convolution of  $p(\mathbf{x})$  by a kernel  $q(\mathbf{z}|\mathbf{x}, \sigma)$

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

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



# Score Matching

- Using this  $q_\sigma(\mathbf{z})$  instead of  $p(\mathbf{x})$ , 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

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

Thus the score function  $\mathbf{s}(\mathbf{z}, \mathbf{w})$  is also like a noise predictor

- Choosing  $q(\mathbf{z}|\mathbf{x}, \sigma) = \mathcal{N}(\mathbf{z}|\mathbf{x}, \sigma^2 I)$ , we get  $\nabla_{\mathbf{z}} \ln q(\mathbf{z}|\mathbf{x}, \sigma) = -\frac{1}{\sigma} \boldsymbol{\epsilon}$  where  $\boldsymbol{\epsilon} = \mathbf{x} - \mathbf{z}$

- Note the similarity with denoising diffusion model where

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

- Thus the score function  $\mathbf{s}(\mathbf{z}, \mathbf{w})$  plays a similar role as noise predictor  $\mathbf{g}(\mathbf{z}, \mathbf{w}, t)$



# 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 already trained a 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

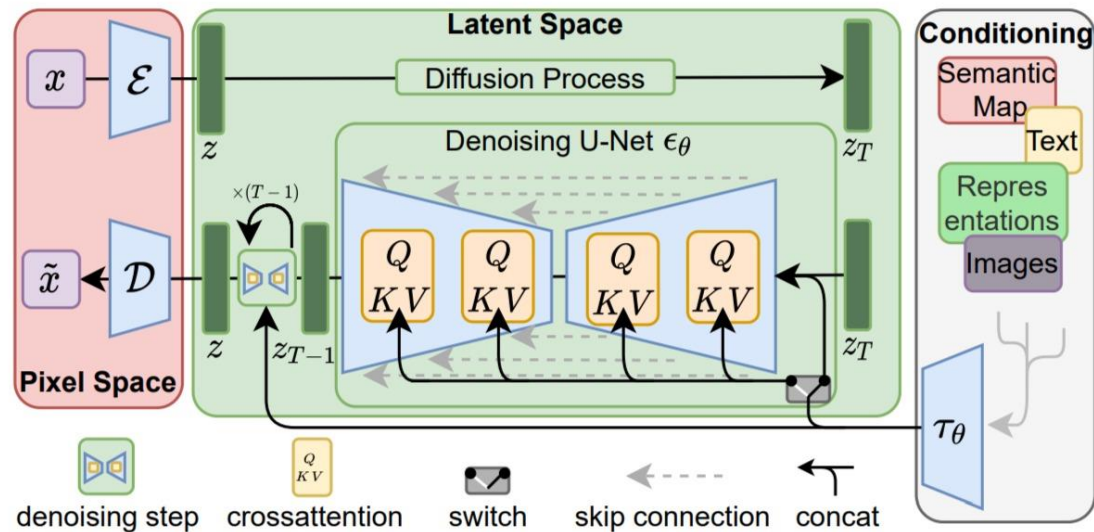
$$\text{score}(\mathbf{x}, \mathbf{c}, \lambda) = \lambda \nabla_{\mathbf{x}} \ln p(\mathbf{x}|\mathbf{c}) + (1 - \lambda) \nabla_{\mathbf{x}} \ln p(\mathbf{x})$$

- 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 on other modalities such as text



# 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)

