

Deep Learning

A course about theory & practice



Generative Networks: Diffusion Models

Marco Piastra

Images from text

- **DALL-E2**

Diffusion Models: generating images from text

«A teapot in the shape of an avocado»

[Image from <https://www.nytimes.com/2022/04/06/technology/openai-images-dall-e.html>]



Videos from text

■ SORA

Generating videos from text prompts

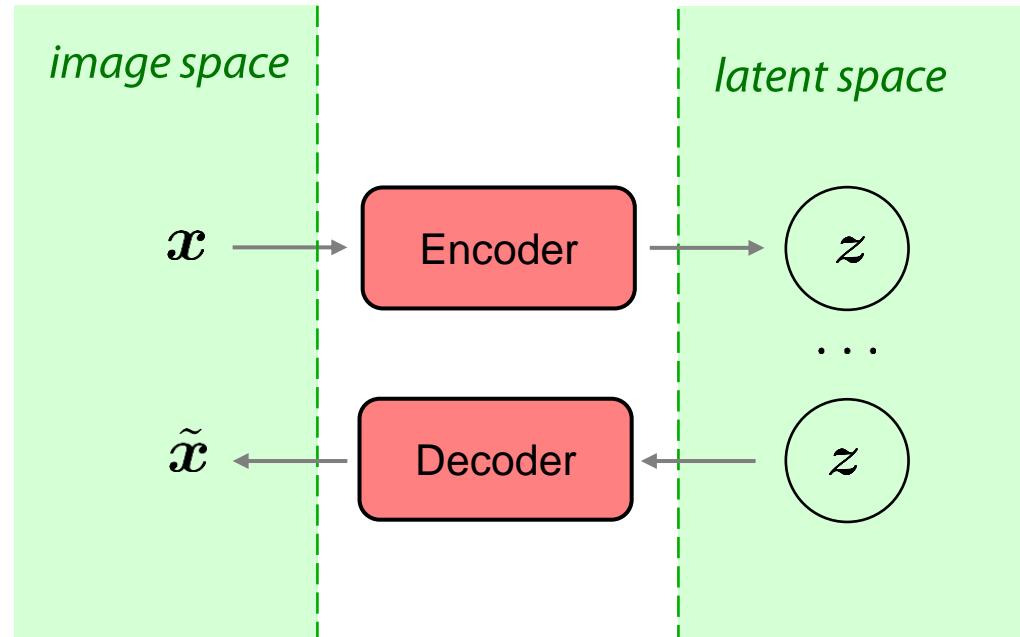


«A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.»

[Video clip from <https://openai.com/index/sora/>]

Diffusion Models are autoencoders

The basic, intuitive idea is to perform diffusion in the *latent space*

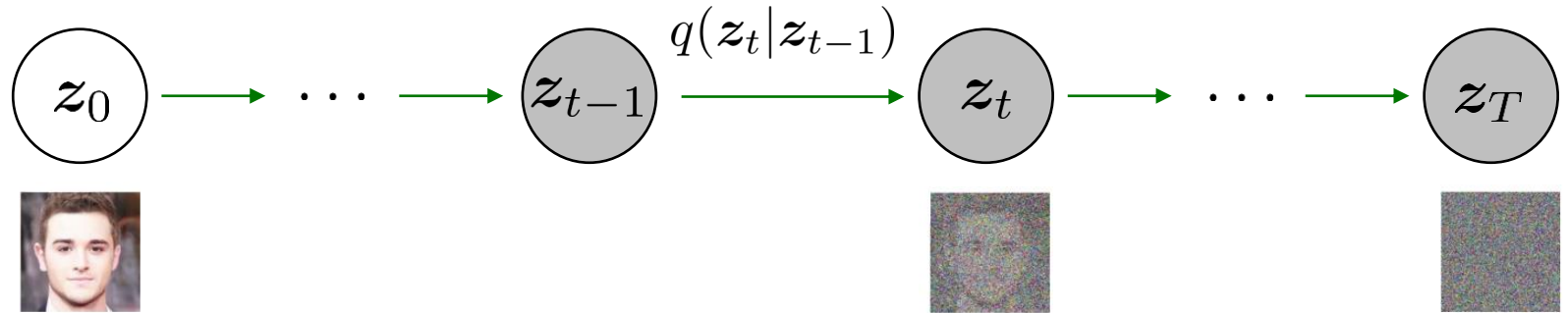


Diffusion Models

Basic idea

Forward Diffusion

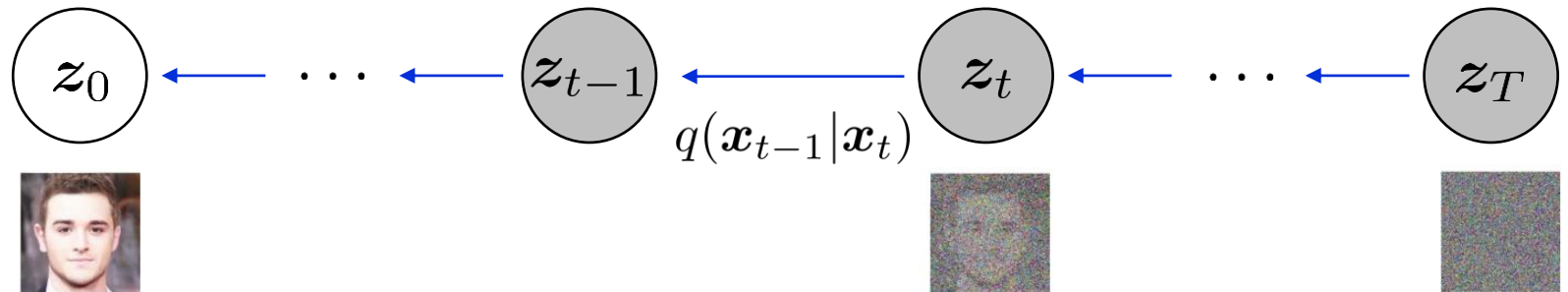
Assume that images are corrupted by Gaussian noise with known parameters



The idea behind

Denoising Diffusion Probabilistic Models (DDPM)

is learning how to reverse the process



Starting from the end: the DDPM training algorithm

Forward Diffusion

Assume that images
are corrupted by
Gaussian noise
with known parameters

The idea behind

Denoising Diffusion Probabilistic Models
is learning how to reverse the process

Note that

The neural network $g(\mathbf{z}_t, \boldsymbol{\vartheta}, t)$
is expected to predict the noise $\boldsymbol{\epsilon}$
that has been applied

Algorithm 20.1: Training a denoising diffusion probabilistic model

```
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  
     $\boldsymbol{\epsilon} \sim \mathcal{N}(\boldsymbol{\epsilon} | \mathbf{0}, \mathbf{I})$  // Sample a noise vector  
     $\mathbf{z}_t \leftarrow \sqrt{\alpha_t} \mathbf{x} + \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}$  // Evaluate noisy latent variable  
     $\mathcal{L}(\mathbf{w}) \leftarrow \|\mathbf{g}(\mathbf{z}_t, \boldsymbol{\vartheta}, t) - \boldsymbol{\epsilon}\|^2$  // Compute loss term  
    Take optimizer step  
until converged  
return  $\mathbf{w}$ 
```

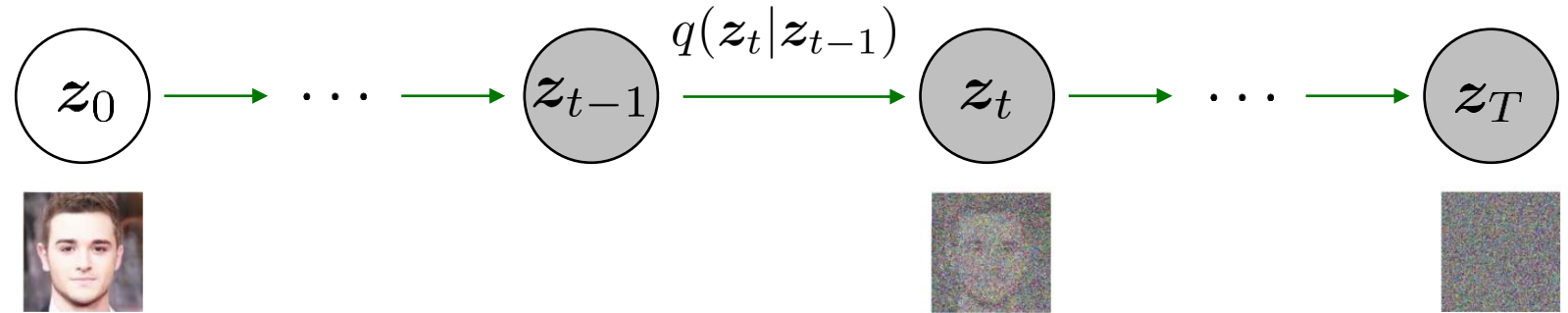
Neural network
with suitable architecture

[Image from <https://www.bishopbook.com/>]

Forward diffusion

Forward Diffusion

Assume that images are corrupted by Gaussian noise with known parameters



$$z_t \sim \mathcal{N} \left(\sqrt{1 - \beta_t} z_{t-1}, \beta_t \mathbf{I} \right)$$

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

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

'Noise schedule':
Hyperparameters

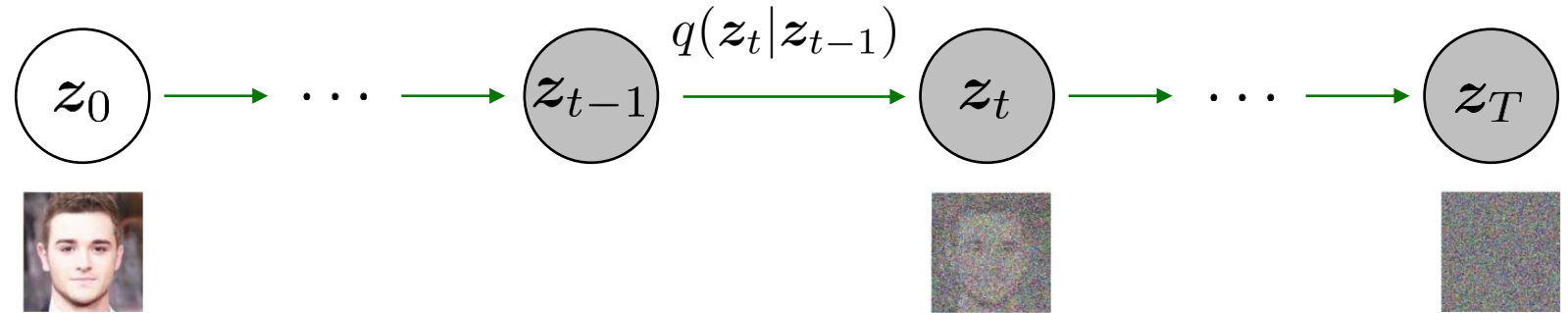
Could be rewritten as

$$z_t = \sqrt{1 - \beta_t} z_{t-1} + \sqrt{\beta_t} \epsilon$$
$$\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Forward diffusion

Forward Diffusion

Assume that images are corrupted by Gaussian noise with known parameters



At any forward step t , the diffusion sequence can be compacted as

$$z_t \sim \mathcal{N}(\sqrt{\alpha_t} z_0, (1 - \alpha_t) \mathbf{I})$$

where:

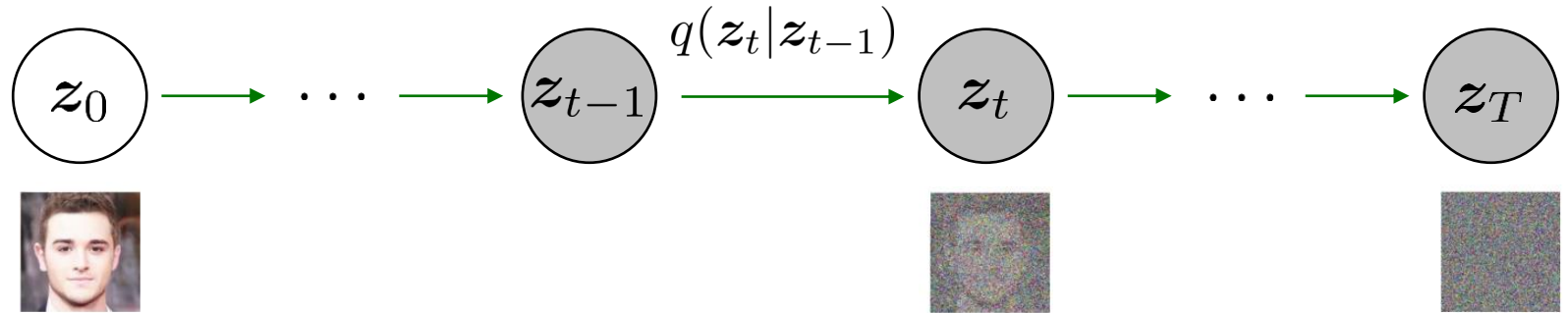
$$\alpha_t = \prod_{\tau=1}^t (1 - \beta_{\tau})$$



Going backward: denoising

Backward Denoising

A neural network
is at the core
of the backward process



We assume that:

$$z_{t-1} = \mu(z_t, t; \vartheta) + \sqrt{\beta_t} \varepsilon$$

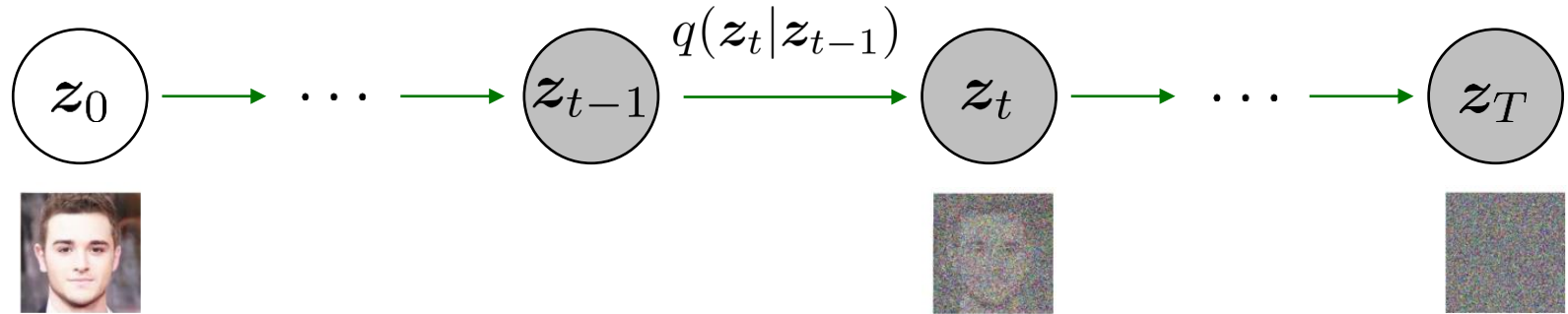
$$\mu(z_t, t; \vartheta) = \frac{1}{\sqrt{1 - \beta_t}} \left\{ z_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \overbrace{g(z_t, t; \vartheta)}^{\text{Neural Network}} \right\}$$

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Going backward: denoising

Backward Denoising

A neural network
is at the core
of the backward process



We assume that:

$$z_{t-1} = \mu(z_t, t; \vartheta) + \sqrt{\beta_t} \varepsilon$$

$$\mu(z_t, t; \vartheta) = \frac{1}{\sqrt{1 - \beta_t}} \left\{ z_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \overbrace{g(z_t, t; \vartheta)}^{\text{Neural Network}} \right\}$$

$$\varepsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

How can the neural network be trained?
(A suitable loss function is needed)

Going backward: denoising

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

An approximation to $q(\mathbf{z}_{t-1}|\mathbf{z}_t)$

We assume that:

$$\mathbf{z}_{t-1} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{z}_t, t; \boldsymbol{\vartheta}), \beta_t \mathbf{I})$$

During training, \mathbf{z}_0 is known. Then we can sample $\boldsymbol{\varepsilon}_t$

$$\mathbf{z}_t = \sqrt{\alpha_t} \mathbf{z}_0 + \sqrt{1 - \alpha_t} \boldsymbol{\varepsilon}_t$$

Noise added at step t

Therefore, it can be proven that:

$$\mathbf{m}(\mathbf{z}_{t-1}) = \frac{1}{\sqrt{1 - \beta_t}} \left(\mathbf{z}_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \boldsymbol{\varepsilon}_t \right)$$

is the true mean of:

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

Then the Kullback-Leibler divergence is:

$$\begin{aligned} \text{KL}(q(\mathbf{z}_{t-1}|\mathbf{z}_t) \parallel \tilde{q}(\mathbf{z}_{t-1}|\mathbf{z}_t)) \\ = \frac{1}{2\beta_t} \|\mathbf{m}(\mathbf{z}_{t-1}) - \boldsymbol{\mu}(\mathbf{z}_t, t; \boldsymbol{\vartheta})\|^2 + \text{const} \end{aligned}$$

Going backward: denoising

$$\begin{aligned} \text{KL} (q(\mathbf{z}_{t-1}|\mathbf{z}_t) \parallel \tilde{q}(\mathbf{z}_{t-1}|\mathbf{z}_t)) \\ = \frac{1}{2\beta_t} \|\mathbf{m}(\mathbf{z}_{t-1}) - \boldsymbol{\mu}(\mathbf{z}_t, t; \boldsymbol{\vartheta})\|^2 + \text{const} \end{aligned}$$

Therefore, given previous assumptions:

$$\begin{aligned} \mathbf{m}(\mathbf{z}_{t-1}) &= \frac{1}{\sqrt{1-\beta_t}} \left(\mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \boldsymbol{\varepsilon}_t \right) \\ \boldsymbol{\mu}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) &= \frac{1}{\sqrt{1-\beta_t}} \left\{ \mathbf{z}_t - \frac{\beta_t}{\sqrt{1-\alpha_t}} \mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) \right\} \end{aligned}$$

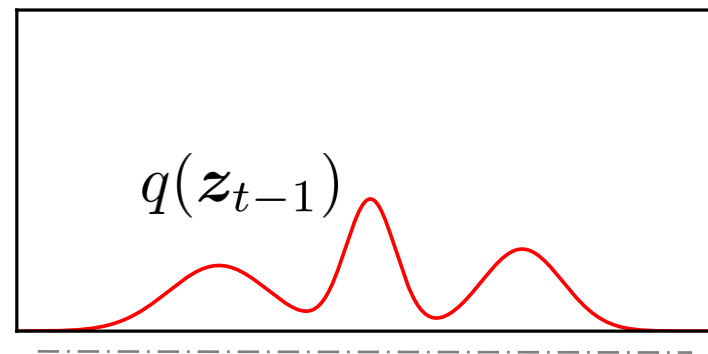
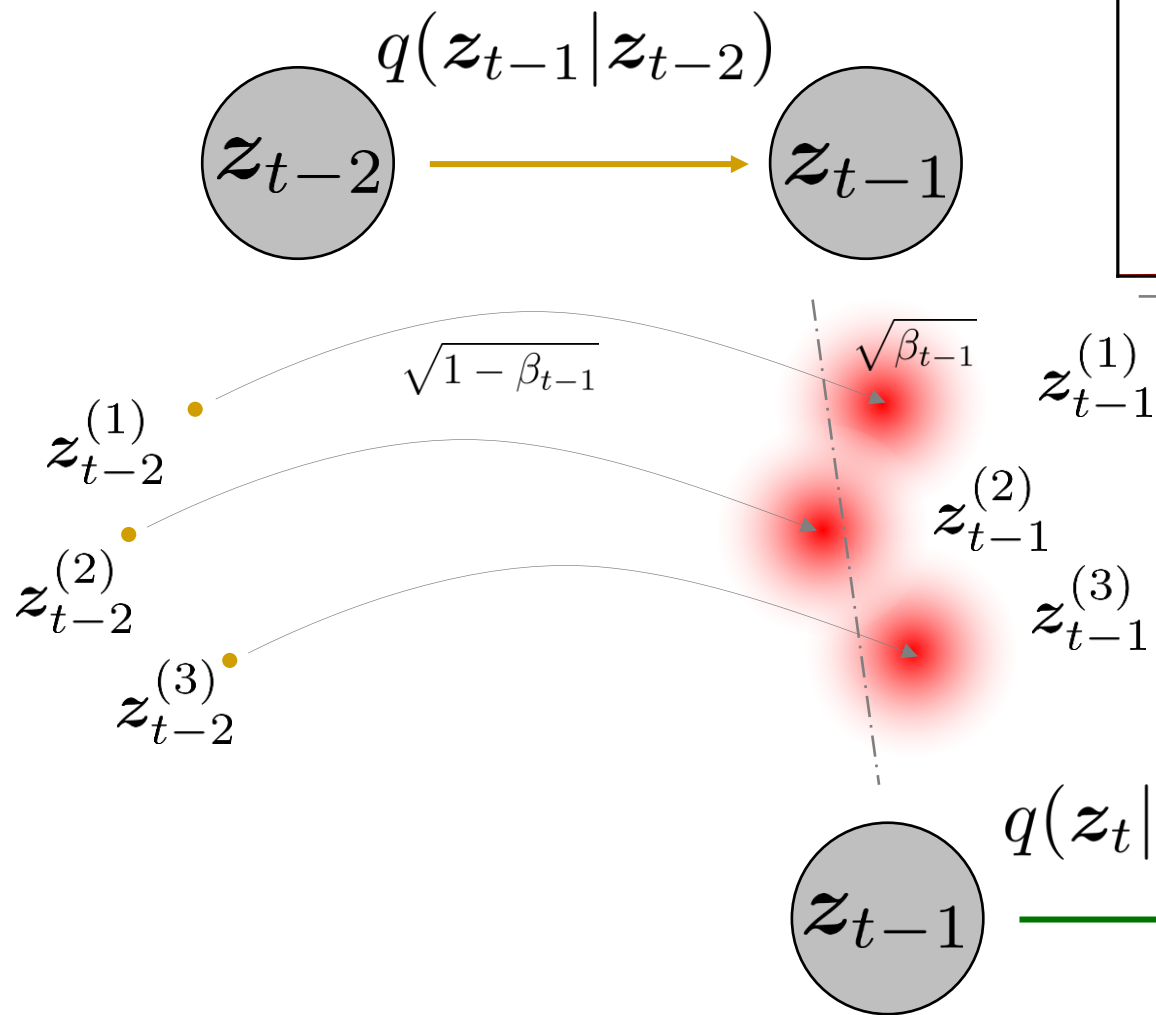
$$\text{KL} (\tilde{q}(\mathbf{z}_{t-1}|\mathbf{z}_t) \parallel q(\mathbf{z}_{t-1}|\mathbf{z}_t)) \propto \|\mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) - \boldsymbol{\varepsilon}_t\|^2$$

$$L(\boldsymbol{\vartheta}) := \|\mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) - \boldsymbol{\varepsilon}_t\|^2 \quad \text{--- To be minimized}$$

Fundamental: this line of reasoning works provided that $q(\mathbf{z}_{t-1}|\mathbf{z}_t)$, $\tilde{q}(\mathbf{z}_{t-1}|\mathbf{z}_t)$ are Gaussians ...

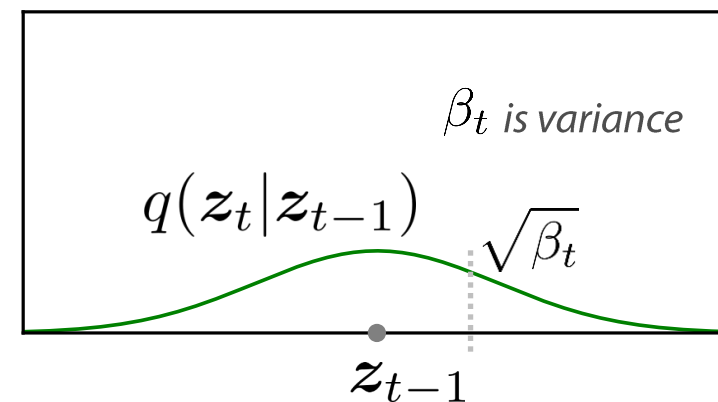
Diffusion Models: Why so many steps?

Going forward: adding noise



Marginal distribution
(cross-section)
at step $t-1$

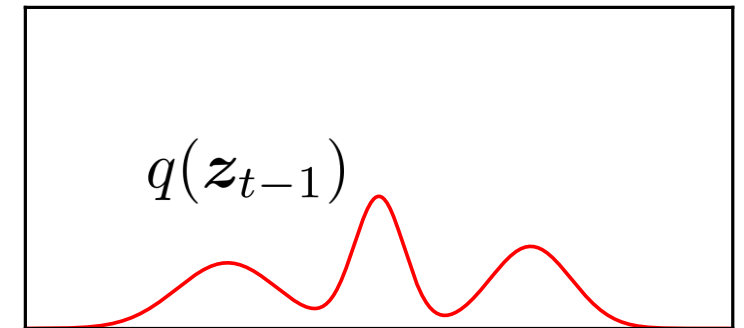
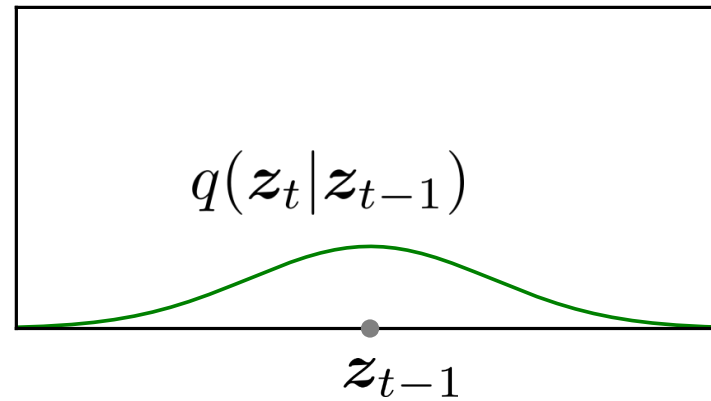
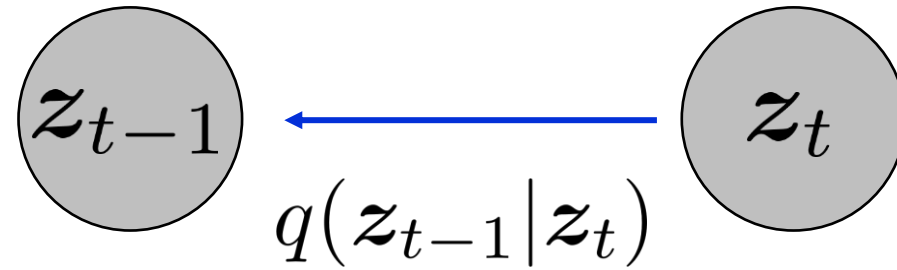
Forward probability distribution
(Gaussian)



$$z_t = \sqrt{1 - \beta_t} z_{t-1} + \sqrt{\beta_t} \varepsilon$$

Going backward: denoising

The backward probability distribution can be computed from forward and marginals using Bayes' theorem



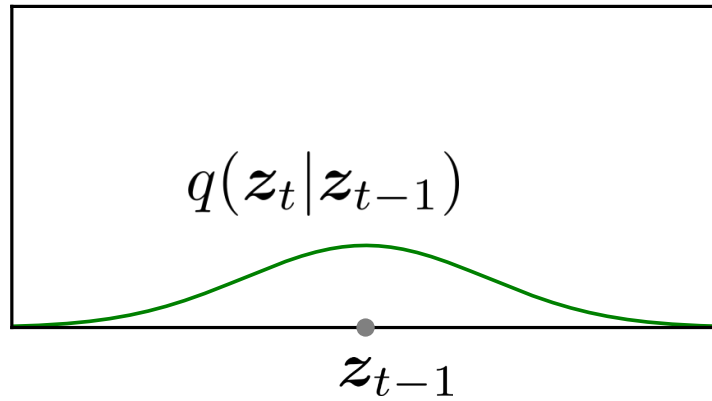
This is what we want to learn

This is what we know, by design

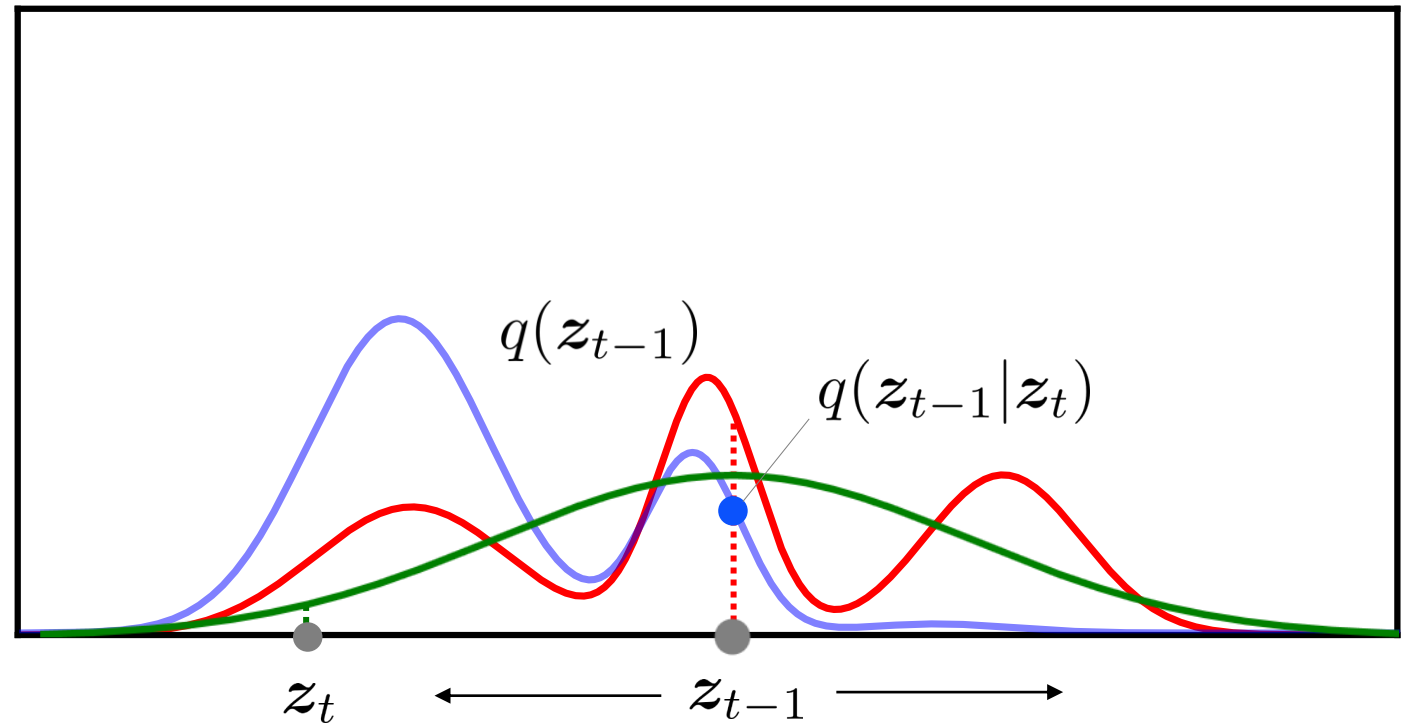
$$q(z_{t-1} | z_t) = \frac{q(z_t | z_{t-1}) q(z_{t-1})}{q(z_t)}$$

Bayes' Theorem

Going backward: denoising



At training time, z_t is known
(dataset + forward diffusion)



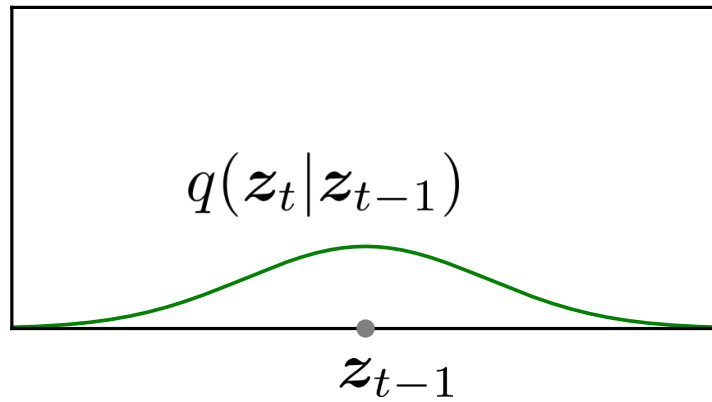
is the free variable, at this step

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

z_t is known, hence this is constant

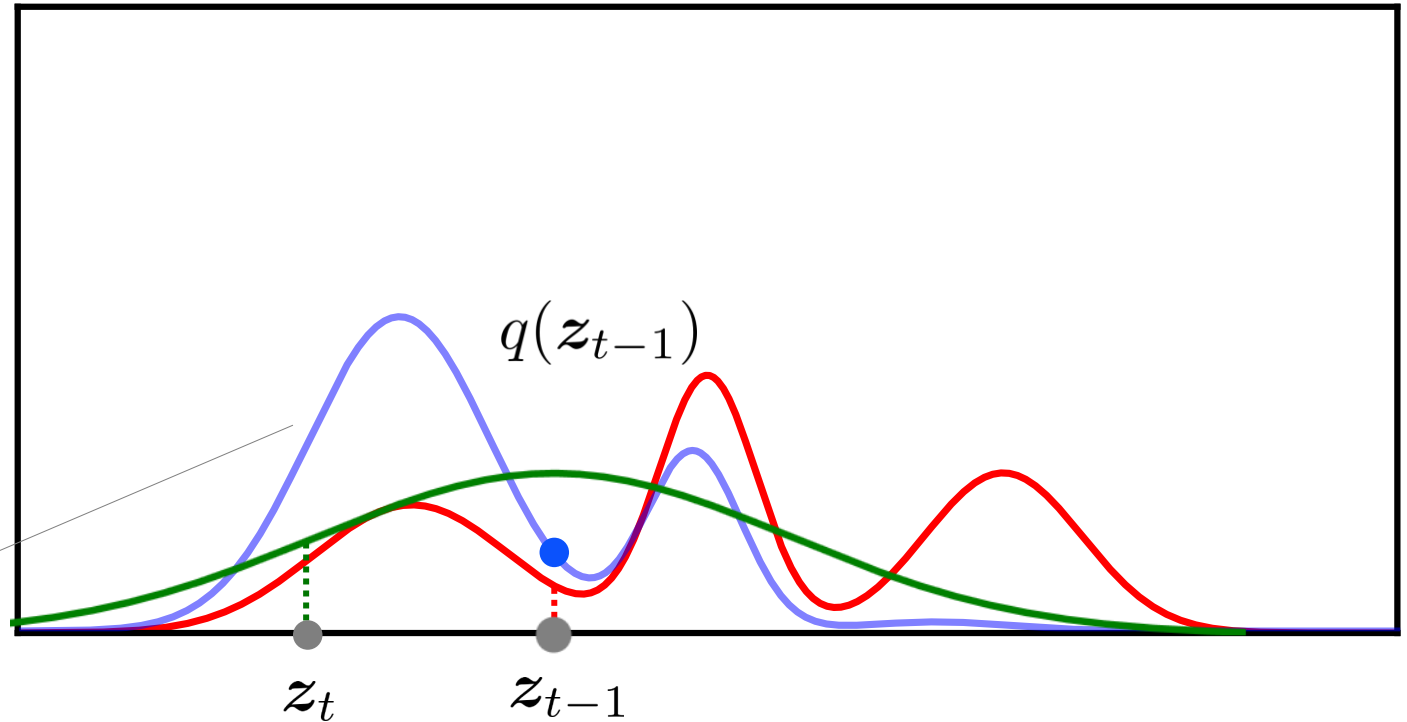
[Image from <https://www.bishopbook.com/>]

Going backward: denoising



At training time, z_t is known
(dataset + forward diffusion)

The reverse probability
is the blue curve

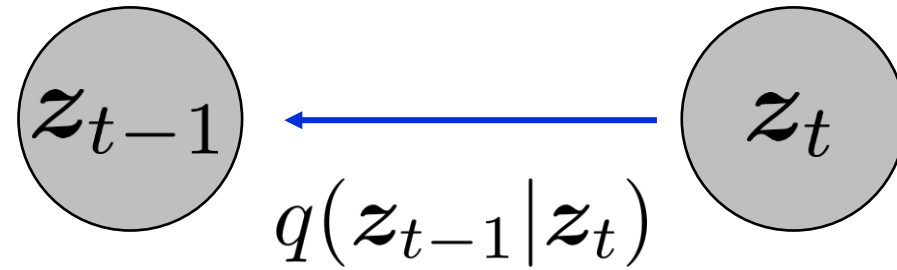


$$q(z_{t-1} | z_t) = \frac{q(z_t | z_{t-1})q(z_{t-1})}{q(z_t)}$$

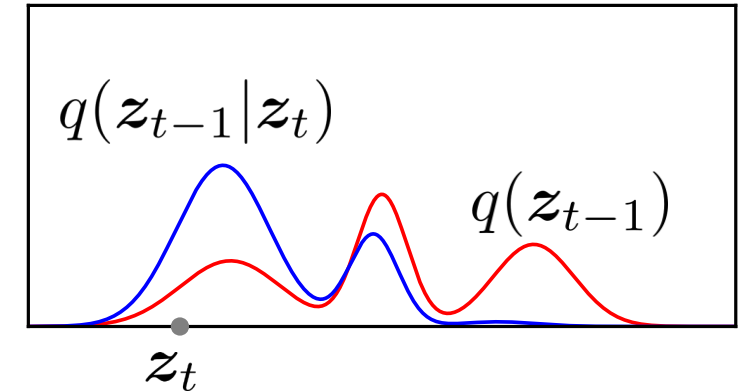
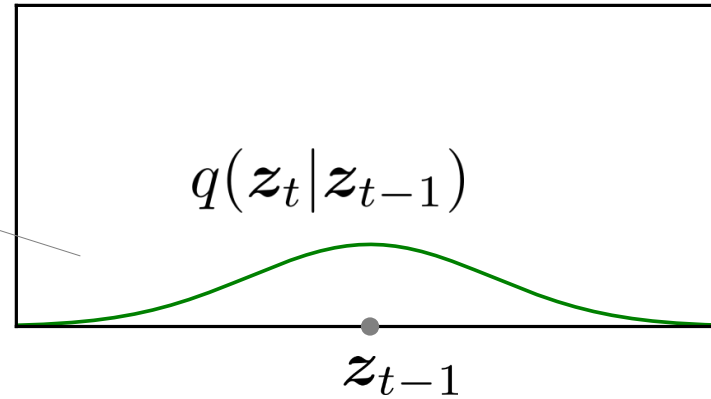
[Image from <https://www.bishopbook.com/>]

Going backward: denoising

When β_t is large
 $q(z_{t-1}|z_t)$
becomes very different
from a Gaussian, hence
unsuitable for training



β_t large

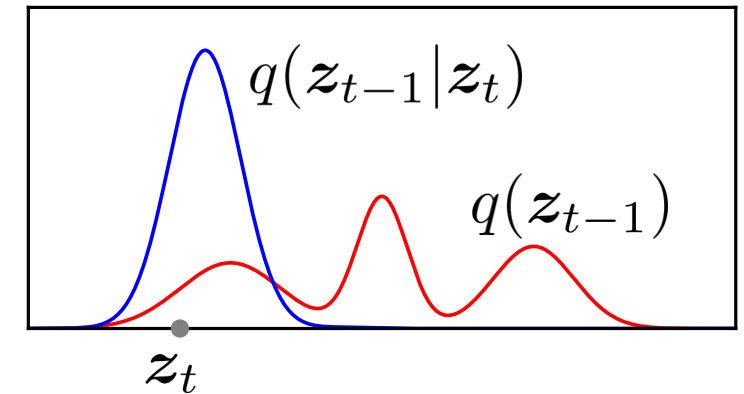
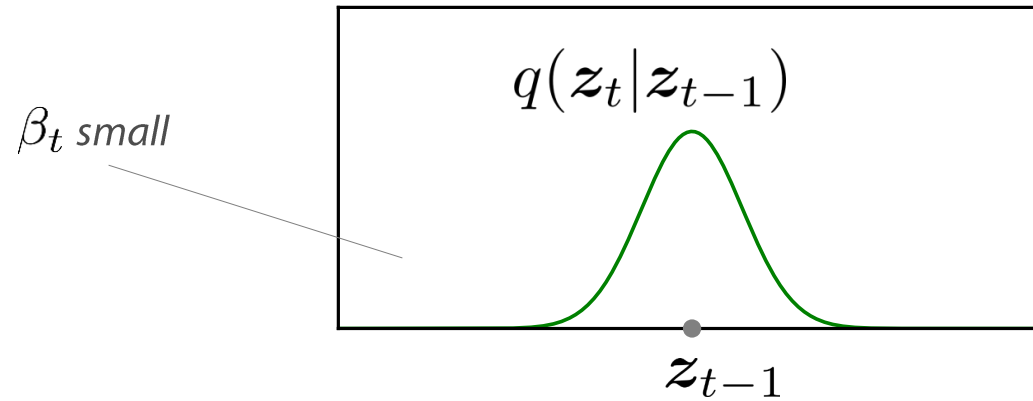
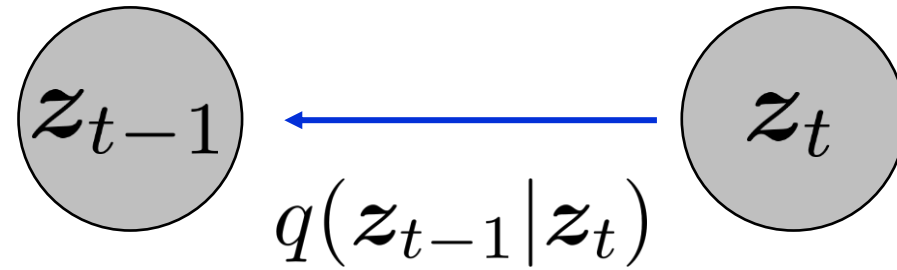


$$q(z_{t-1}|z_t) = \frac{q(z_t|z_{t-1})q(z_{t-1})}{q(z_t)}$$

[Image from <https://www.bishopbook.com/>]

Going backward: denoising

When β_t is small
 $q(z_{t-1}|z_t)$
is approximately Gaussian



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

Links

<https://www.assemblyai.com/blog/diffusion-models-for-machine-learning-introduction/>

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

<https://www.superannotate.com/blog/diffusion-models>

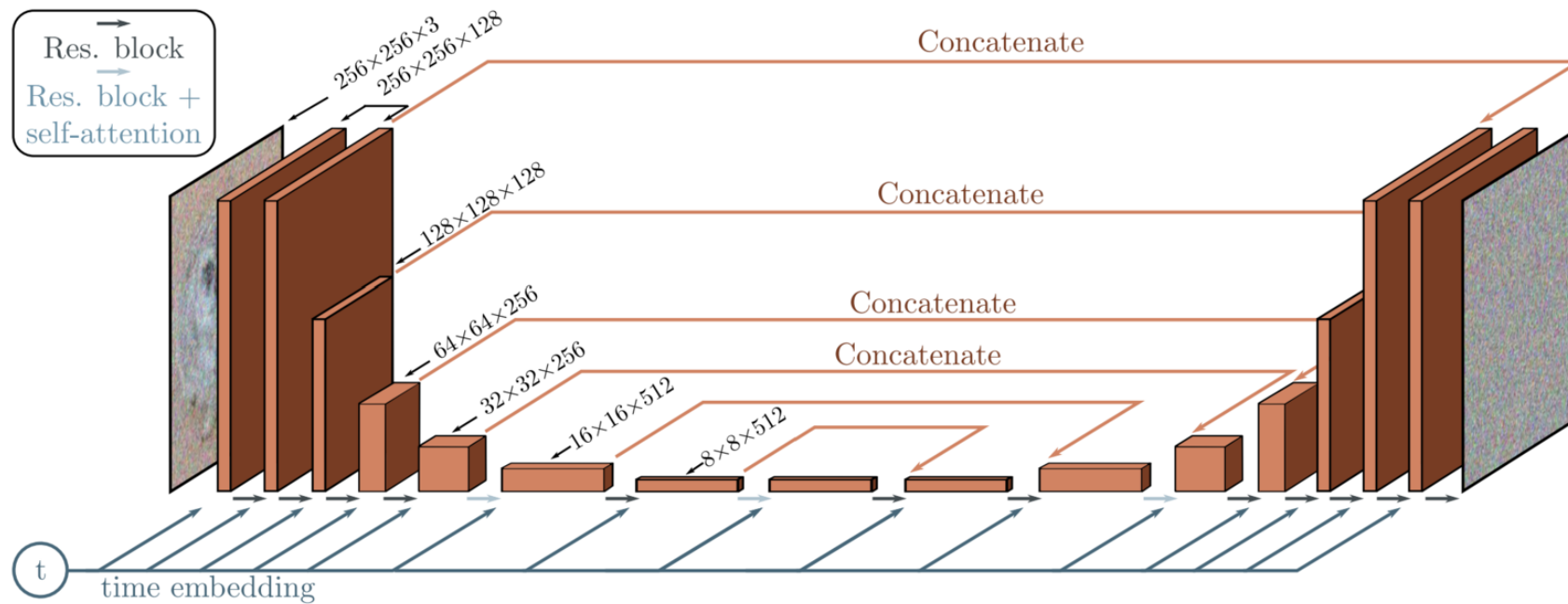
<https://encord.com/blog/diffusion-models/>

Practical Implementation

Conditional U-Net as basic denoising block

Loss function: $L(\boldsymbol{\vartheta}) := \| \mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) - \boldsymbol{\epsilon}_t \|^2$

The network architecture for $\mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta})$ is a U-Net with time embedding

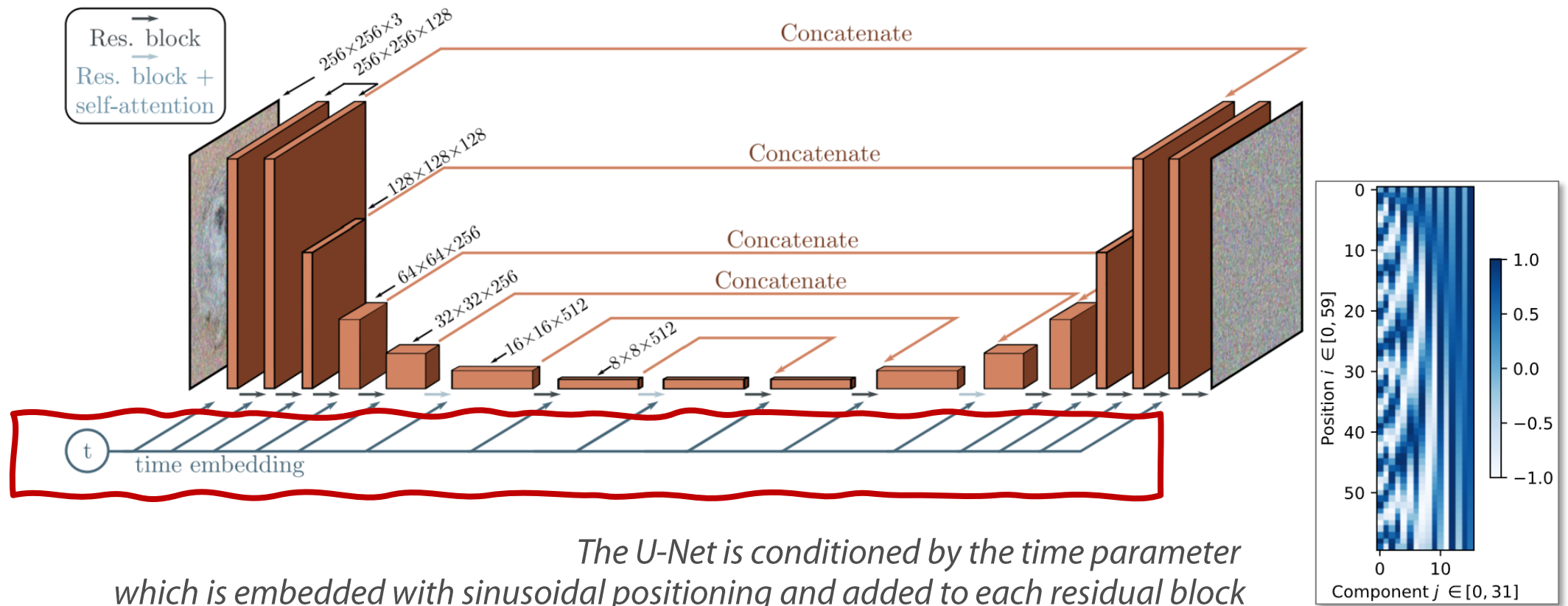


[Ho, Jain & Abbeel, 2020 - <https://arxiv.org/pdf/2006.11239>]

Conditional U-Net as basic denoising block

Loss function: $L(\boldsymbol{\vartheta}) := \| \mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) - \boldsymbol{\epsilon}_t \|^2$

The network architecture for $\mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta})$ is a U-Net with time embedding



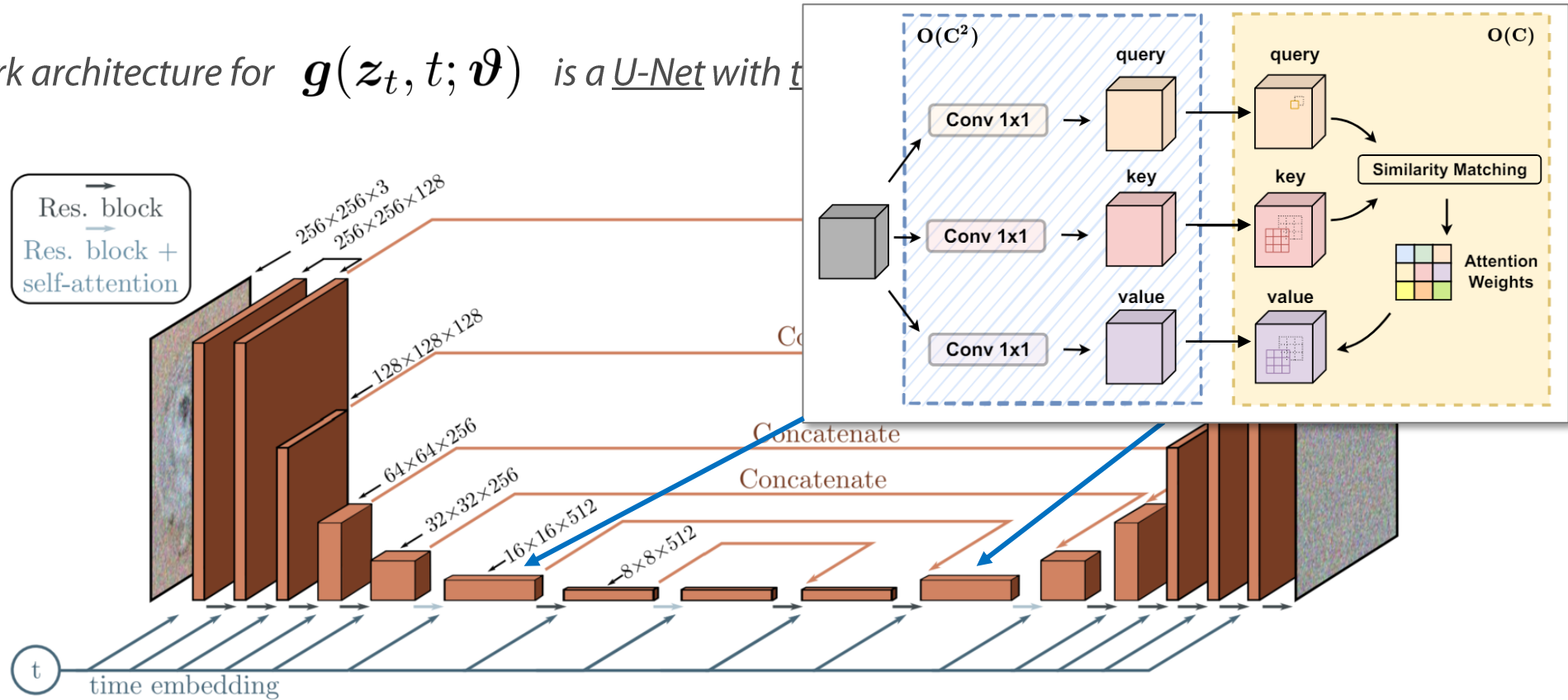
The U-Net is conditioned by the time parameter which is embedded with sinusoidal positioning and added to each residual block

[Ho, Jain & Abbeel, 2020 - <https://arxiv.org/pdf/2006.11239>]

Conditional U-Net as basic denoising block

Loss function: $L(\boldsymbol{\vartheta}) := \| \mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta}) - \boldsymbol{\epsilon}_t \|^2$

The network architecture for $\mathbf{g}(\mathbf{z}_t, t; \boldsymbol{\vartheta})$ is a U-Net with t



Self-Attention modules
are interspersed with convolutional blocks in the pipeline

[Ho, Jain & Abbeel, 2020 - <https://arxiv.org/pdf/2006.11239>]

Latent Diffusion Models

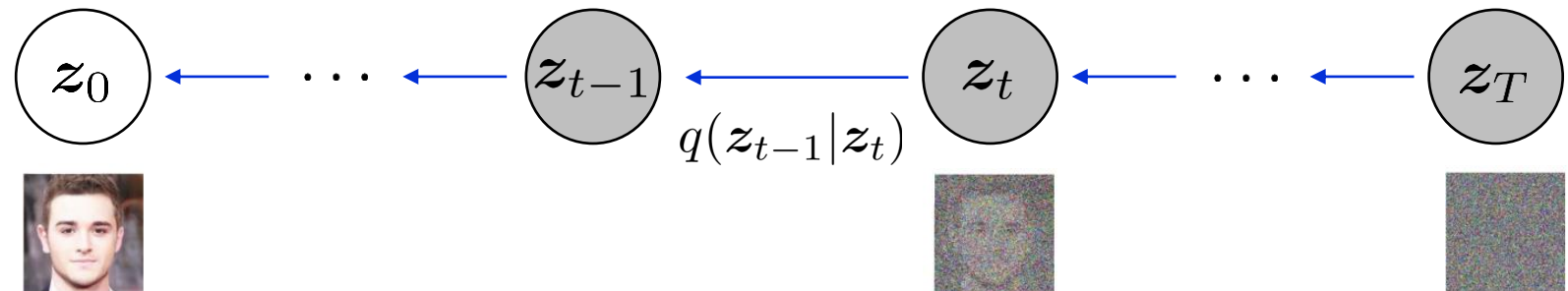
Forward Diffusion

It is relatively easy and inexpensive
(It can be performed in one step)



Backward Denoising

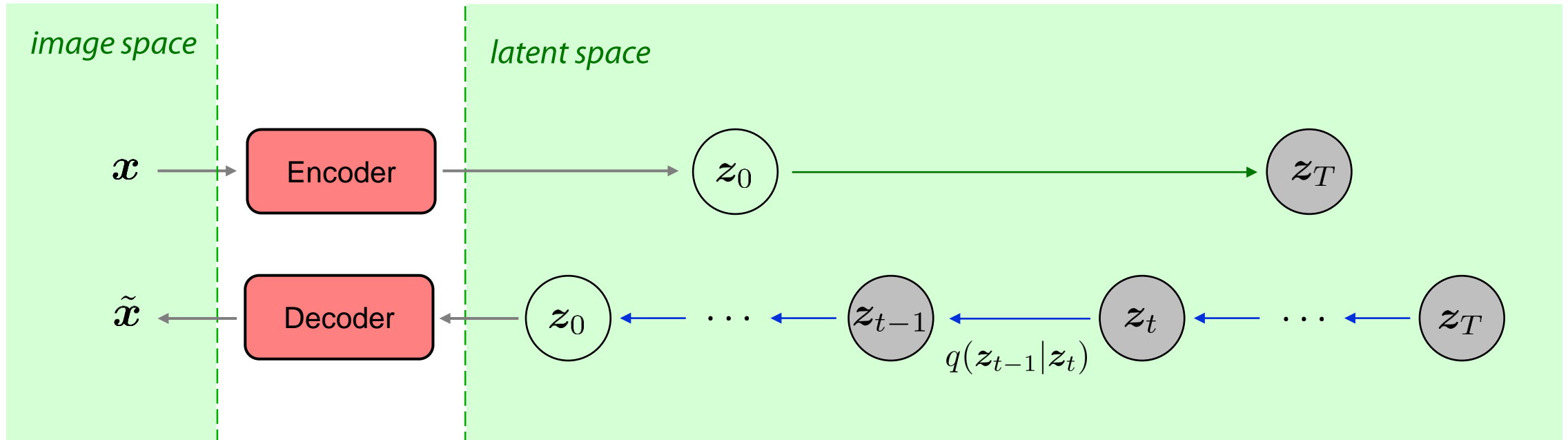
Must be performed in small steps
and is quite expensive,
in particular with high-resolution images



Latent Diffusion Models

Latent Diffusion Model

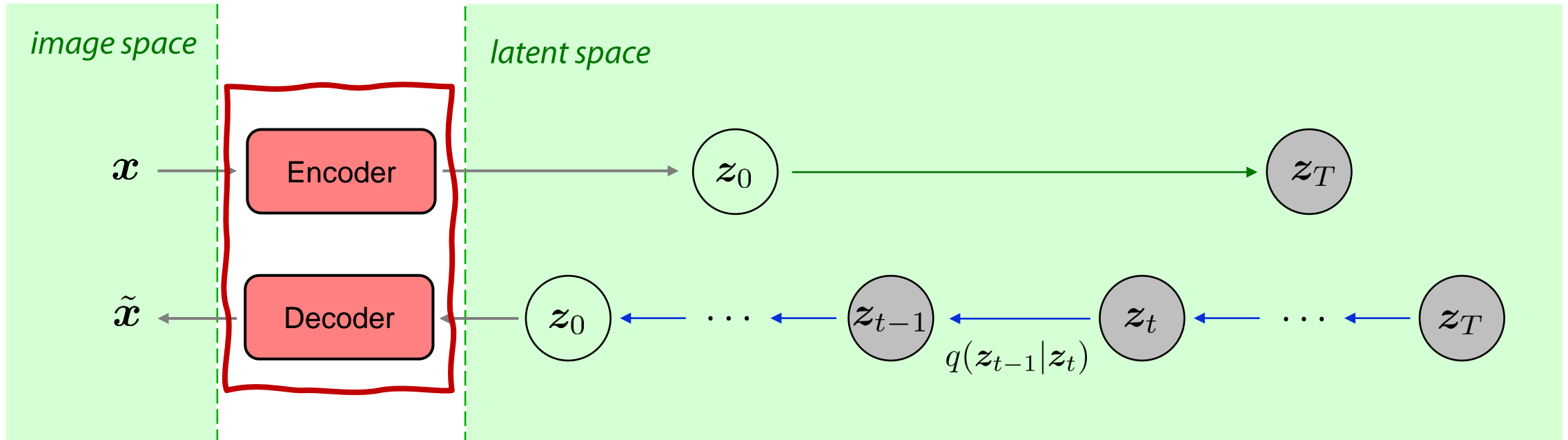
The intuitive idea is to perform diffusion in the *latent space*



Latent Diffusion Models

Latent Diffusion Model

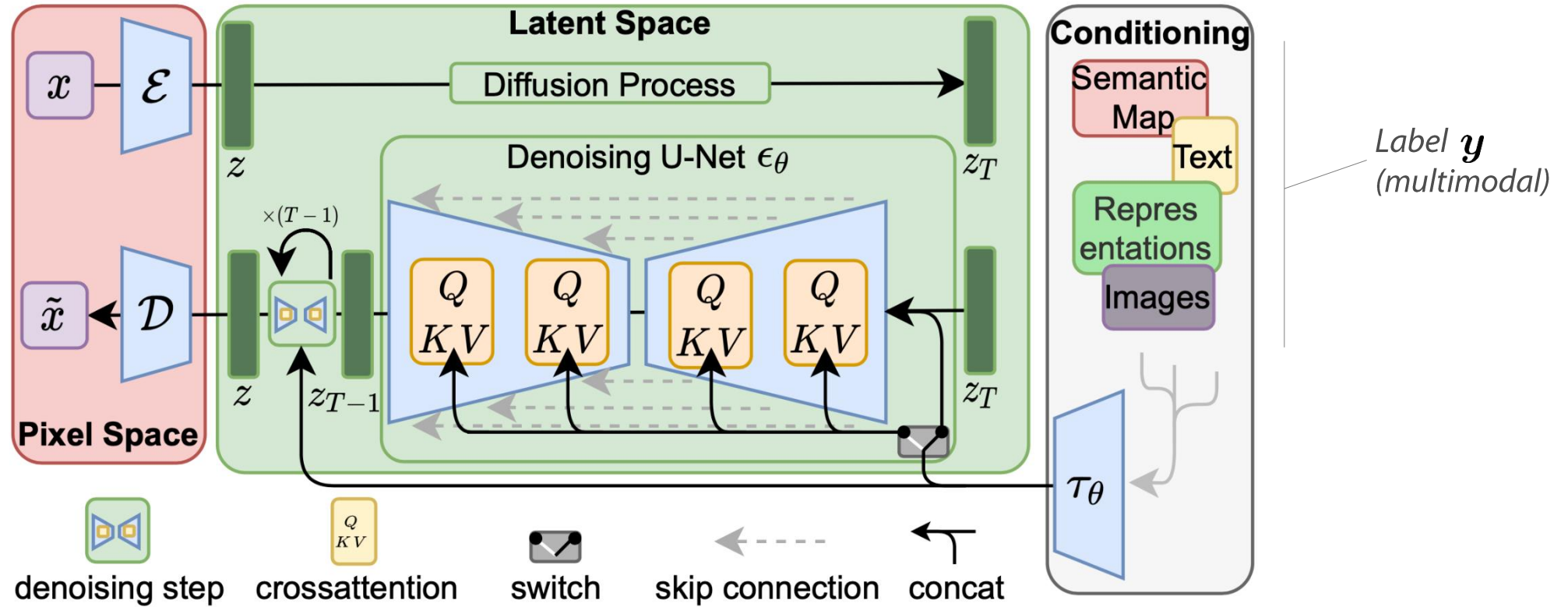
The intuitive idea is to perform diffusion in the *latent space*



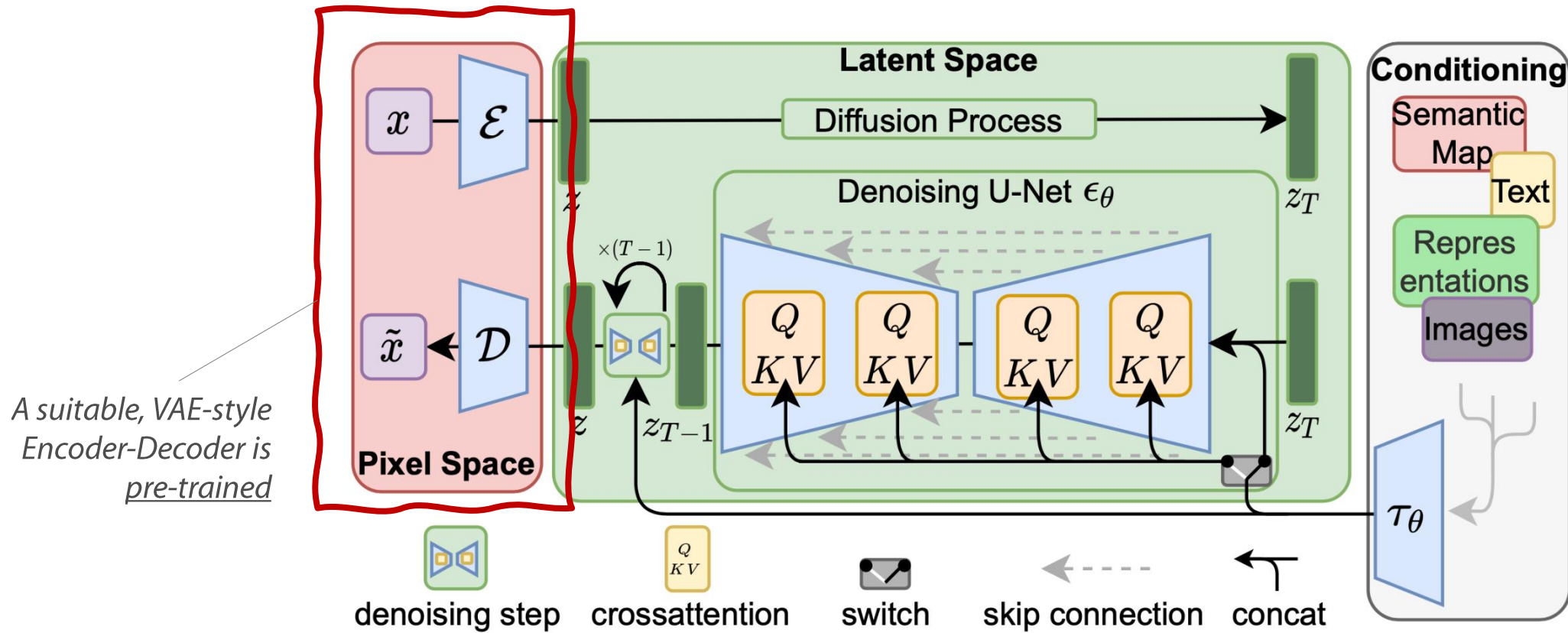
A pre-trained VAE is used to encode and decode high-resolution images into a suitable (reduced) latent format

Conditioning on Multimodal Labels

Latent Diffusion Models with Conditioning

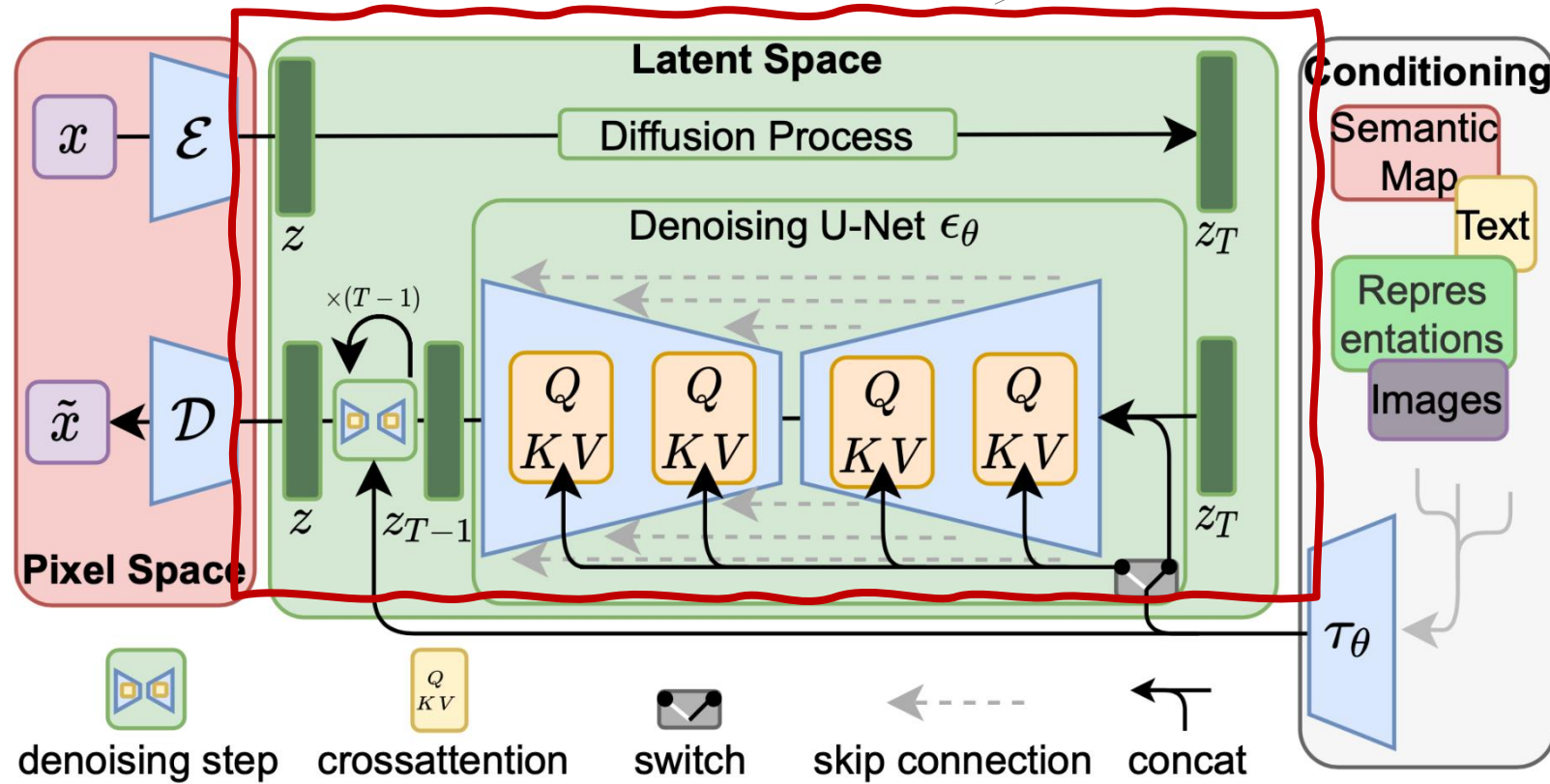


Latent Diffusion Models with Conditioning

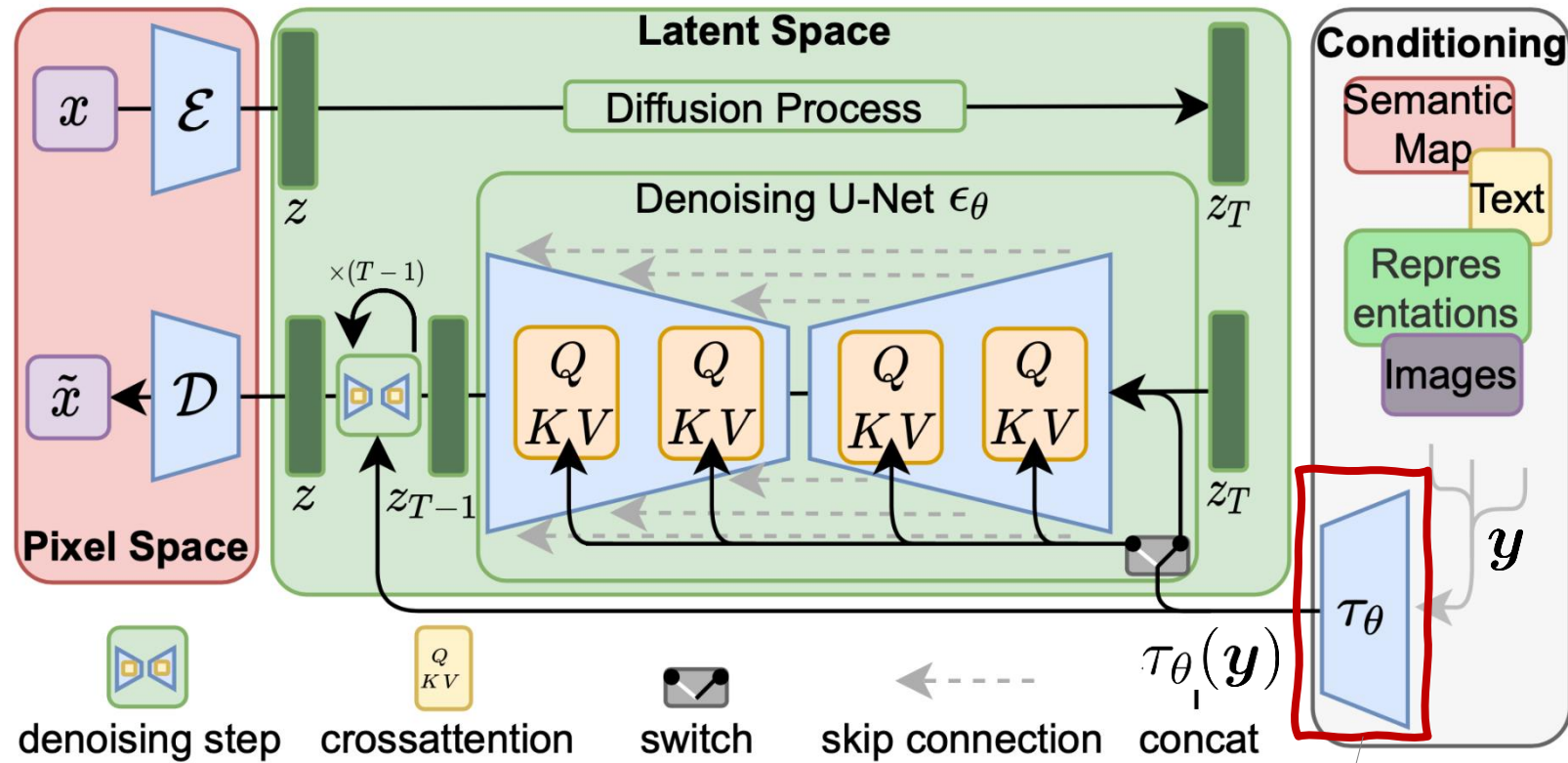


Latent Diffusion Models with Conditioning

The latent diffusion model is then pre-trained (without conditioning)

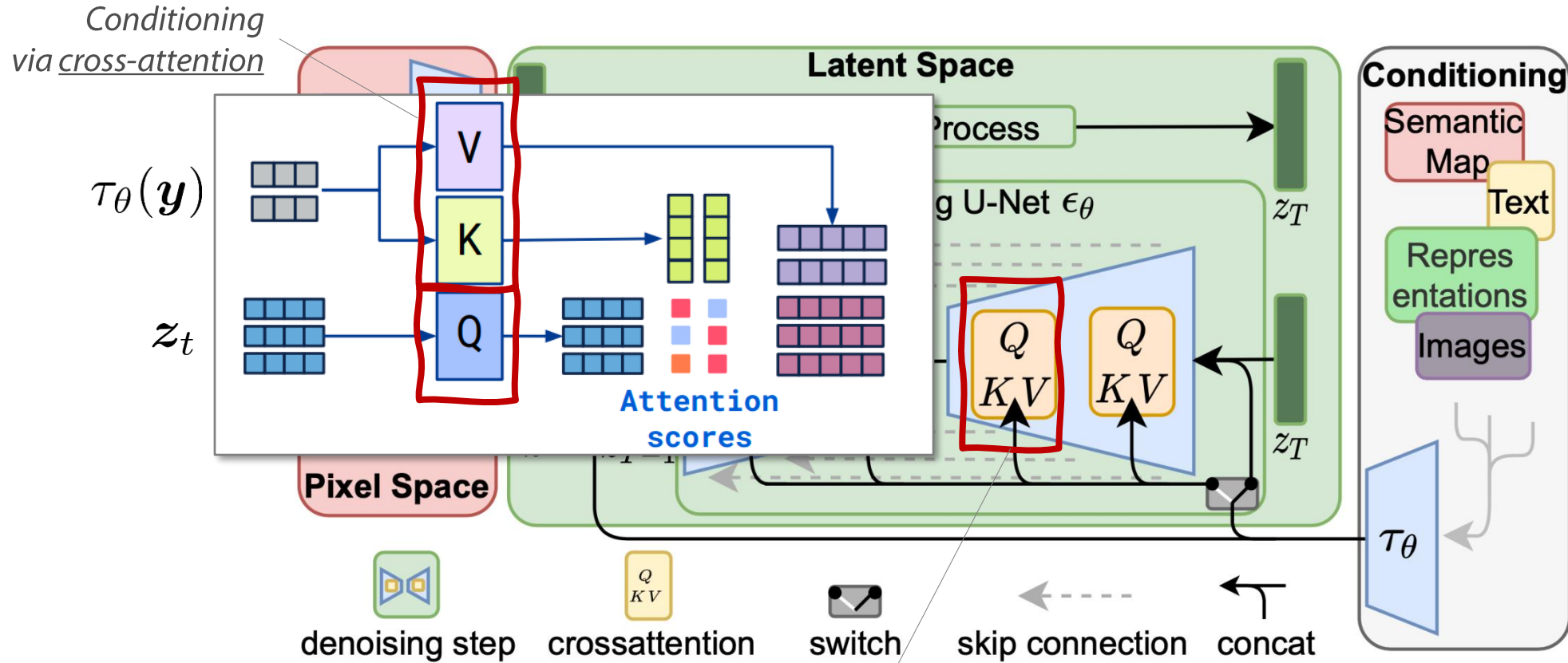


Latent Diffusion Models with Conditioning



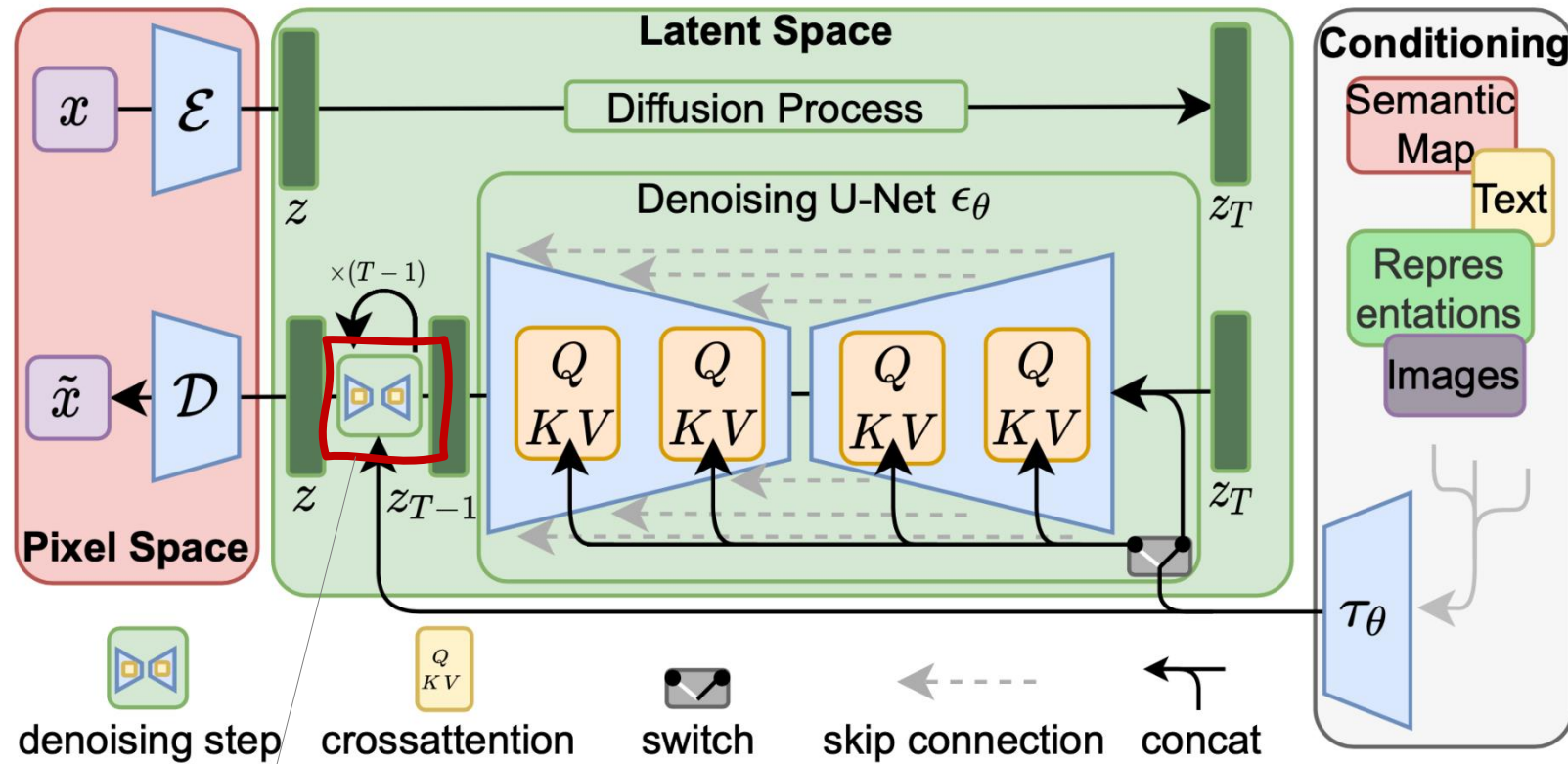
A suitable encoder of the conditioning elements is pre-trained separately

Latent Diffusion Models with Conditioning



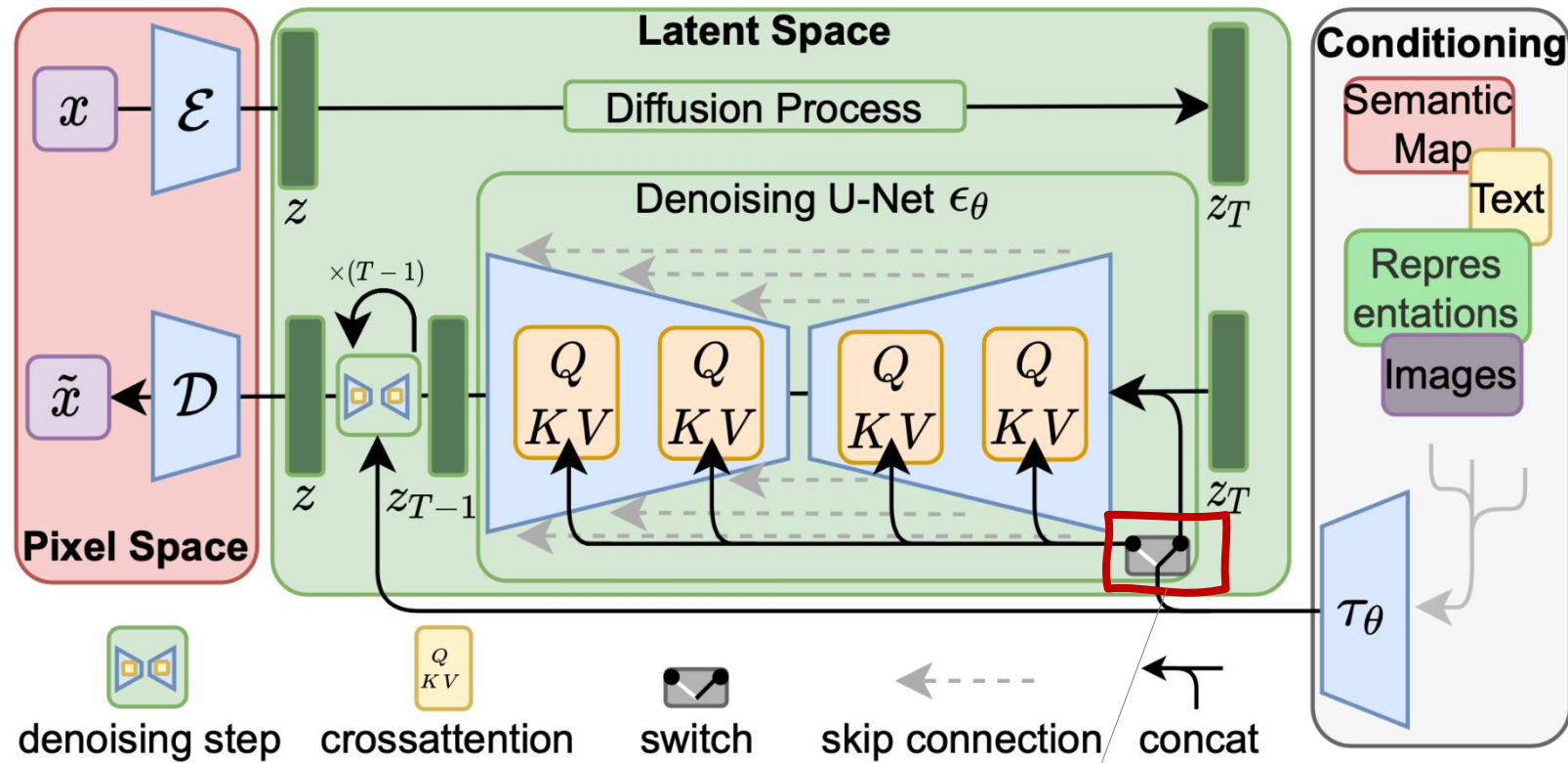
Latent-space representations
and embedded condition elements
are combined via cross-attention

Latent Diffusion Models with Conditioning



The same step is iterated
 $T - 1$ more times

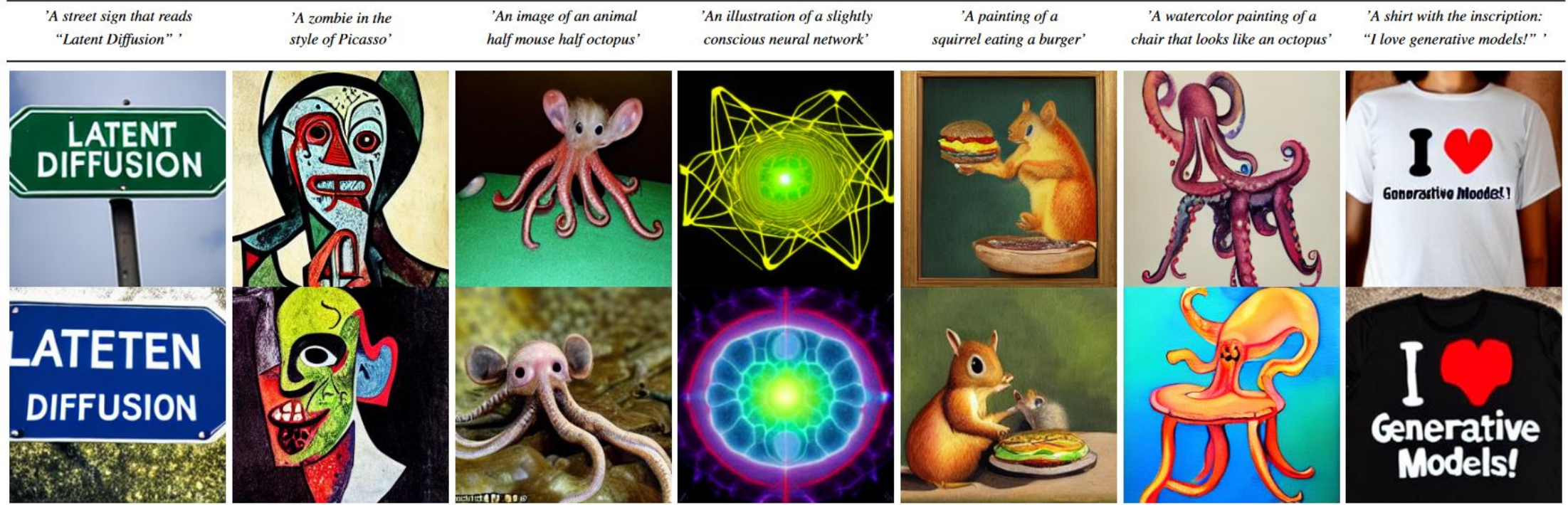
Latent Diffusion Models with Conditioning



The switch is for multi-modality:
if the conditioning element is a class or text, use cross-attention,
if the input is an image, use concatenation

Latent Diffusion Models with Conditioning

Text-to-Image Synthesis on LAION. 1.45B Model.



Links

<https://poloclub.github.io/diffusion-explainer/>

<https://blog.marvik.ai/2023/11/28/an-introduction-to-diffusion-models-and-stable-diffusion/>

<https://theaisummer.com/diffusion-models/>

<https://learnopencv.com/denoising-diffusion-probabilistic-models/>

<https://www.assemblyai.com/blog/diffusion-models-for-machine-learning-introduction/>

<https://lilianweng.github.io/posts/2021-07-11-diffusion-models/>

<https://www.superannotate.com/blog/diffusion-models>

<https://encord.com/blog/diffusion-models/>