

Stable Diffusion

links

<https://jalammar.github.io/illustrated-stable-diffusion/>

<https://medium.com/@steinsfu/stable-diffusion-clearly-explained-ed008044e07e>

<https://www.youtube.com/watch?v=ItLNYA3IWAQ>

<https://github.com/castorini/daam>

<https://www.youtube.com/watch?v=J2WtkA1Xfew>

<https://vaclavkosar.com/ml/cross-attention-in-transformer-architecture>

<https://arxiv.org/abs/2112.10752>

<https://sander.ai/2022/05/26/guidance.html>

https://stable-diffusion-art.com/how-stable-diffusion-work/#How_training_is_done

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

<https://education.yandex.ru/handbook/ml/article/diffuzionnye-modeli>

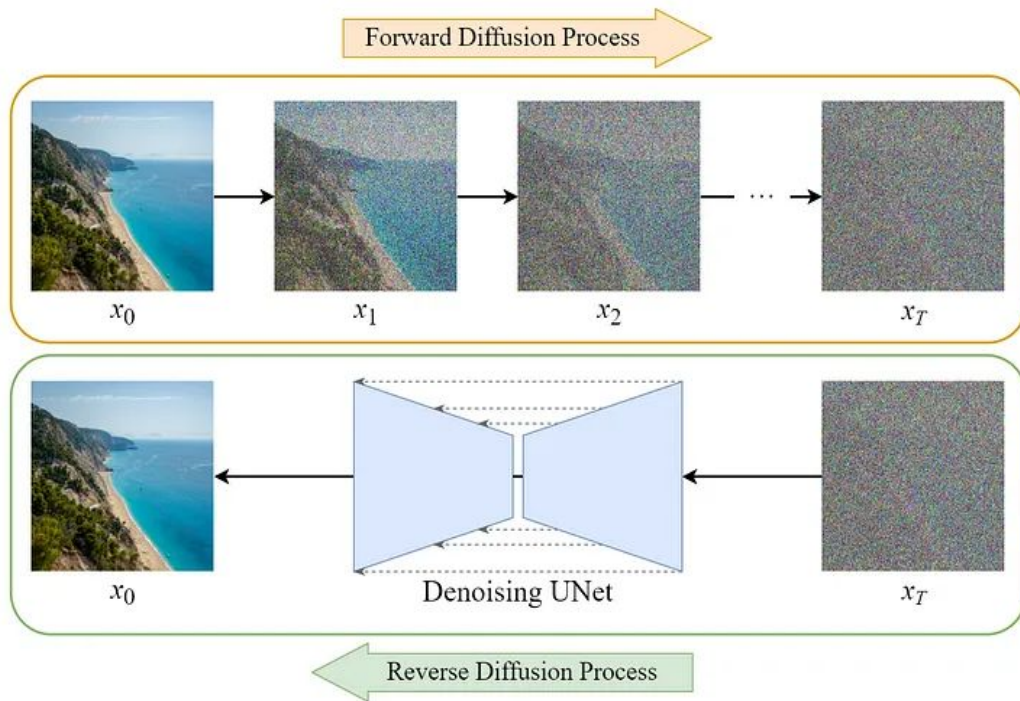
<https://arxiv.org/pdf/2403.18103.pdf>

<https://strikingloo.github.io/wiki/ddim>

https://nn.labml.ai/diffusion/stable_diffusion/index.html

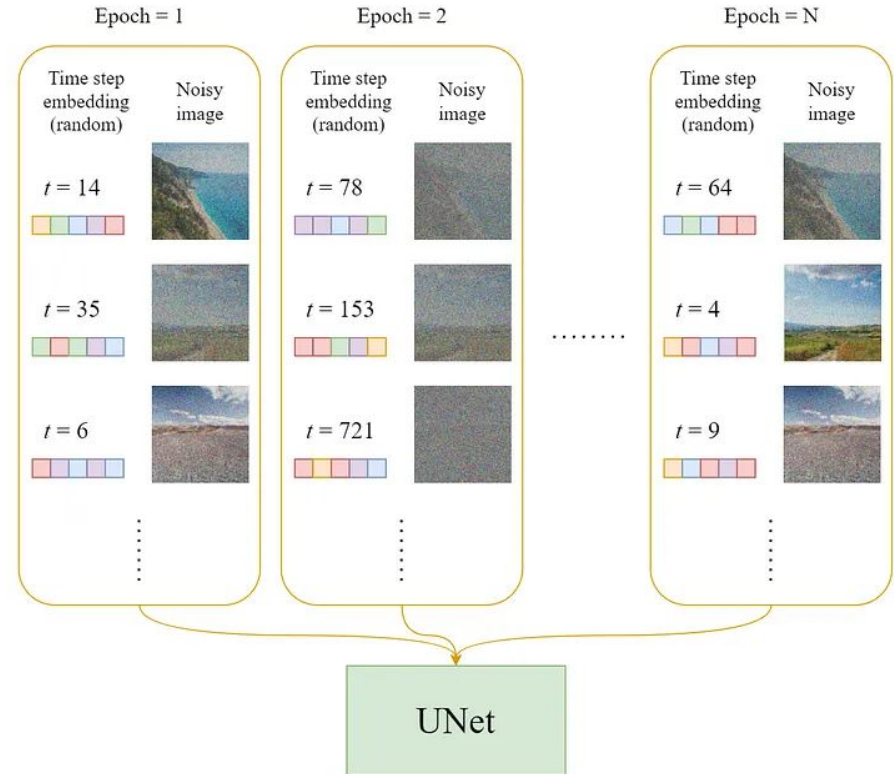
DDPM (Denoising Diffusion Probabilistic Models)

1. Forward Diffusion Process → добавить шум к изображению.
2. Reverse Diffusion Process → убрать шум из изображения .



Training In each epoch:

1. Для каждого изображения будет выбран случайный временной шаг t .
2. Применен гауссовский шум (соответствующий t) к каждому изображению.
3. временные шаги t преобразуются в embedding вектор.



Training DDPM

1. Рандомно выбираем t
2. Генерируем шум и прибавляем его к изображению (количество шума регулируется t)
3. Тренируем UNet предсказывать шум для t

1. Randomly select a time step & encode it



2. Add noise to image



$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \epsilon$$

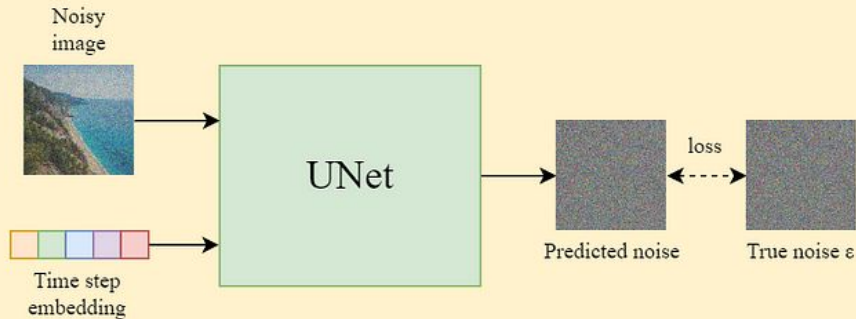
Adjust the amount of noise according to the time step t

$$\epsilon \sim \mathcal{N}(0, 1)$$

$$\alpha_t = 1 - \beta_t$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

3. Train the UNet



Training step

Генерируем гауссовский шум

Итеративно избавляемся от шума пока не получим исходное изображение

Недостатки:

Работа с целым изображением вычислительно очень медленная и дорогостоящая операция

DDPM - генерация занимает много времени и хотелось генерировать быстрее чем за 1000 шагов

1. Sample a Gaussian noise

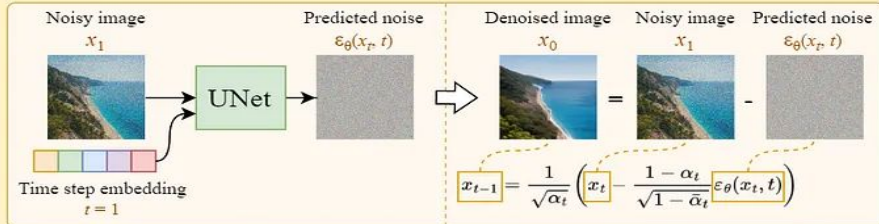
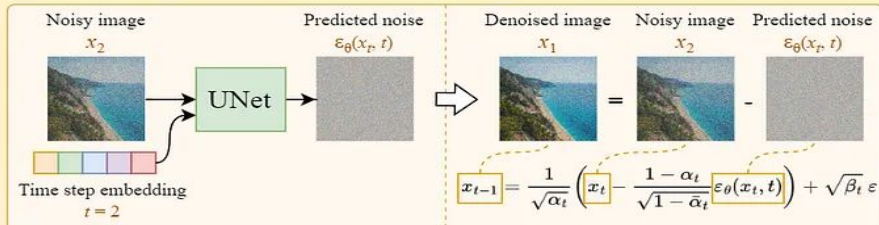
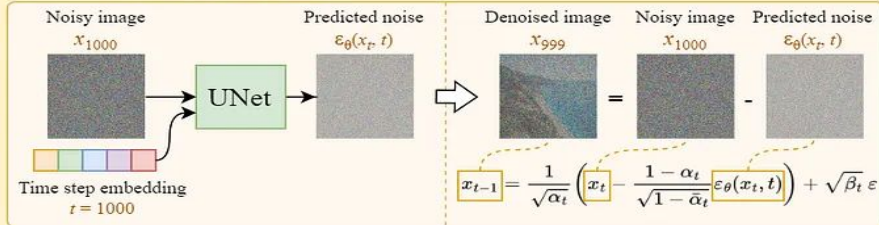
$$x_T \sim \mathcal{N}(0, I)$$

E.g. $T = 1000$

$$x_{1000} \sim \mathcal{N}(0, I)$$



2. Iteratively denoise the image

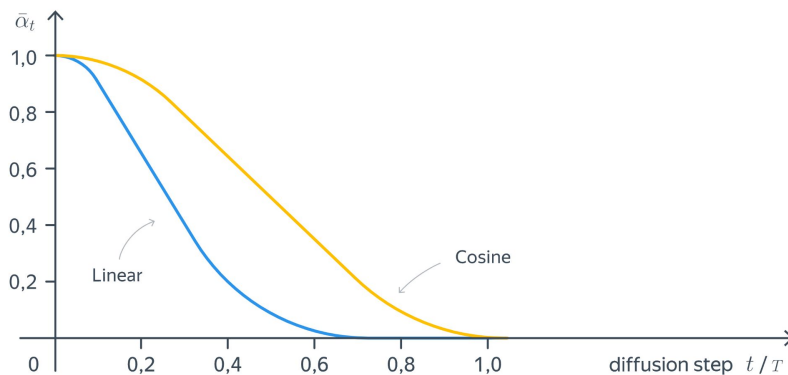


3. Output the denoised image

Denoised image x_0



Improved Denoising Diffusion Probabilistic Models



<https://arxiv.org/abs/2102.09672>

DDIM (Denoising Diffusion Implicit Models)

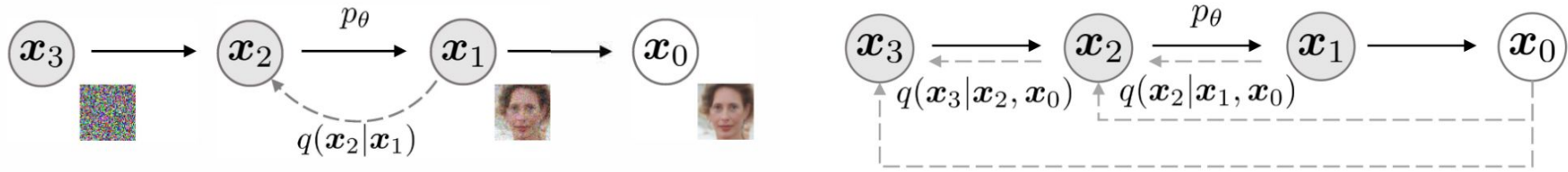
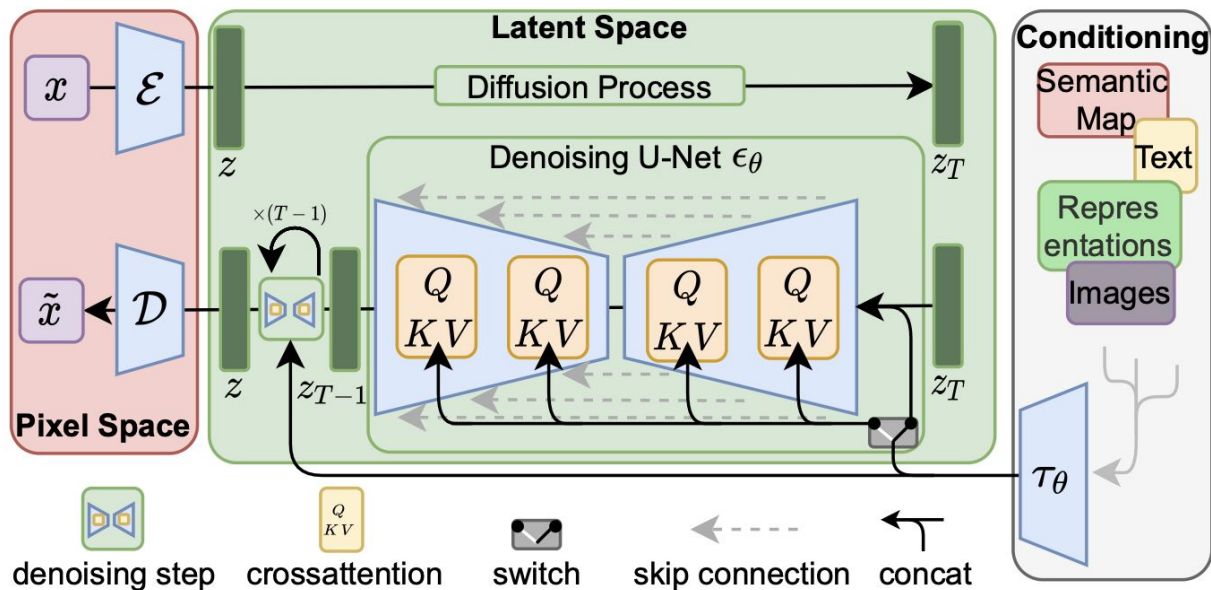


Figure 1: Graphical models for diffusion (left) and non-Markovian (right) inference models.

Авторы предопределили Backward process

Latent Diffusion Models

- Encoder/Decoder
- Latent Space
- UNet Scheduler
- Conditioning

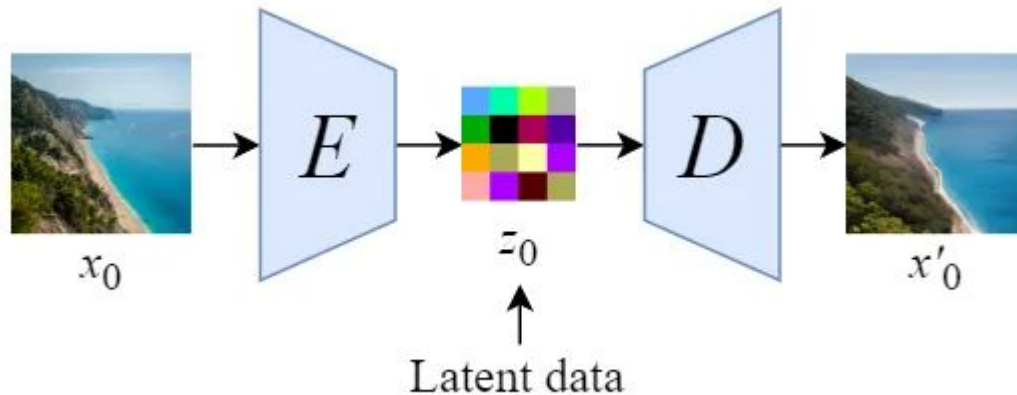


Latent Space

процесс диффузии теперь выполняется в латентном пространстве

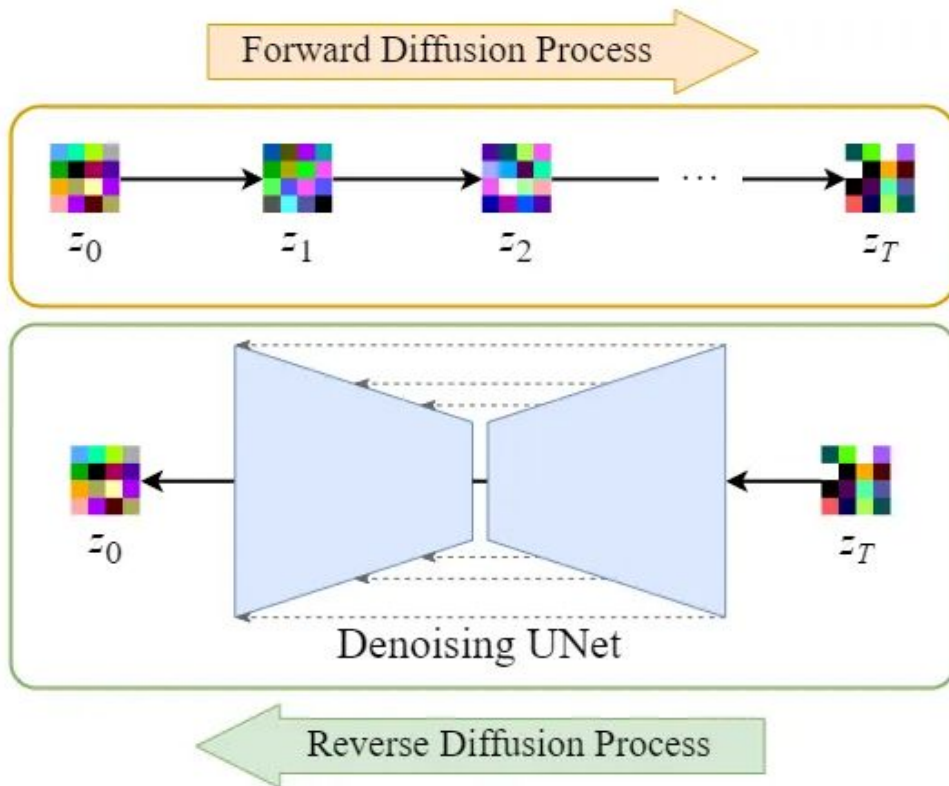
- encoder E кодирует изображение
- decoder D декодирует

изображение(3, 512, 512) -> латент(4,64,64)



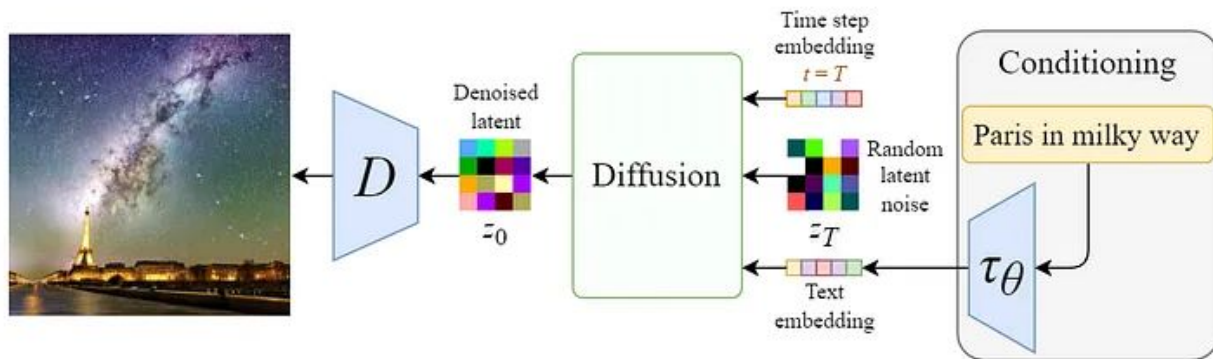
Latent Diffusion

1. Forward Diffusion Process → добавить шум к латенту.
2. Reverse Diffusion Process → убрать шум из латента.



Conditioning

Добавление в процесс диффузии (условия) в данном случае текста с использованием механизма **cross-attention** Для генерации новых изображений по текстовому описанию



Classifier-free guidance

Все, что нам нужно, чтобы превратить модель безусловной диффузии в условную, — это классификатор!

В отличие от вероятностных моделей на прямую аппроксимирующие $p(x)$, диффузионная модель предсказывает score function

$$\nabla_x \log p(x | y).$$

Classifier guidance

$$\nabla_x \log p_\gamma(x | y) = \nabla_x \log p(x) + \gamma \nabla_x \log p(y | x).$$

Classifier-free guidance

$$\nabla_x \log p_\gamma(x | y) = \nabla_x \log p(x) + \gamma (\nabla_x \log p(x | y) - \nabla_x \log p(x)) ,$$

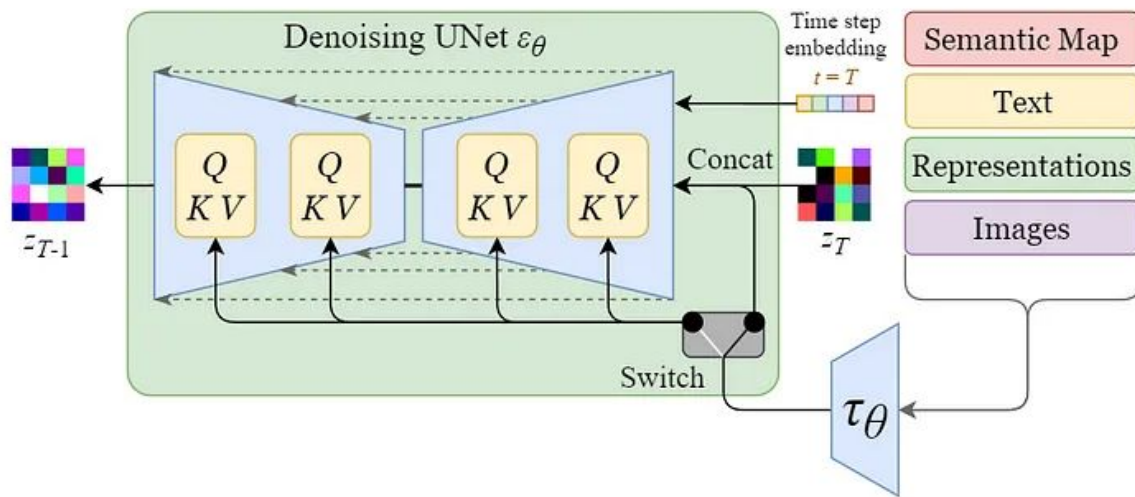
<https://sander.ai/2022/05/26/guidance.html>

Conditioning

switch нужен для использование разного типа “условий” на вход

для текста используется языковая модель для получения вектора (BERT, CLIP) и с помощью Attention мапится в UNet

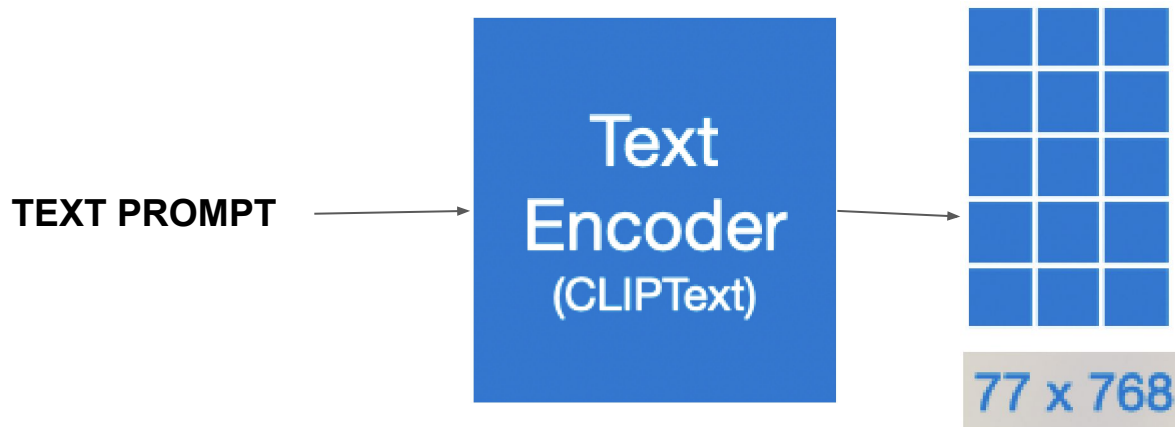
Для другого рода входов (семантические маски, изображения) может использоваться конкатинация



ClipText

Input: text.

Output: 77 token
embeddings vectors,
each in 768
dimensions.



Loss function

$$x_t = \sqrt{\bar{a}_t} x_0 + \sqrt{1 - \bar{a}_t} \varepsilon$$

$$L_{\text{simple}} = \mathbb{E}_{t, x_0, \varepsilon} \left[\|\varepsilon - \varepsilon_{\theta}(x_t, t)\|^2 \right]$$

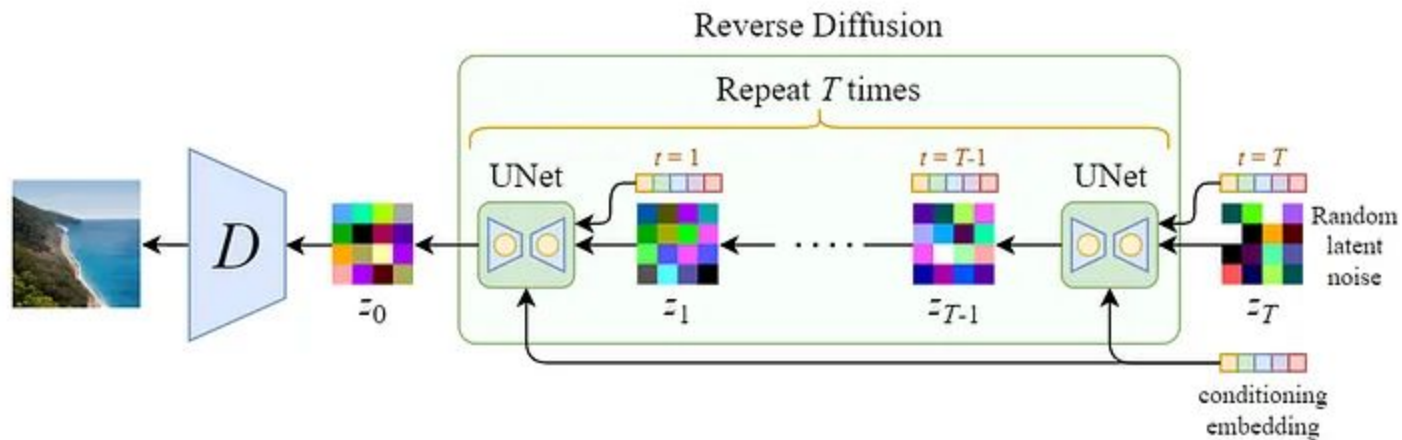
$$z_0 = E(x_0)$$

$$z_t = \sqrt{\bar{\alpha}_t} z_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon$$

$$L_{\text{LDM}} = \mathbb{E}_{t, z_0, \varepsilon, y} \left[\|\varepsilon - \varepsilon_{\theta}(z_t, t, \tau_{\theta}(y))\|^2 \right]$$

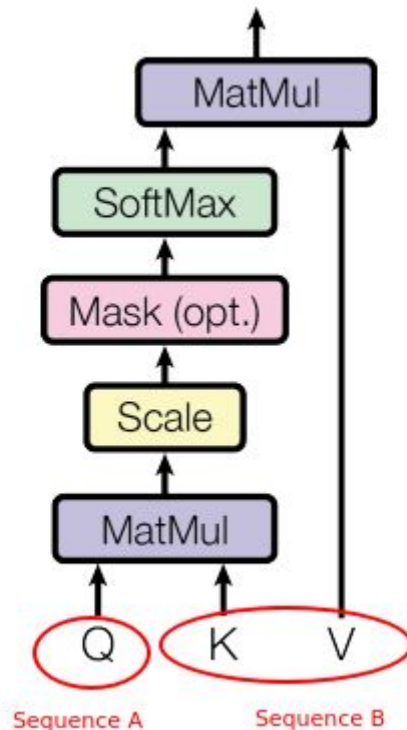
Conditioning

Sampling

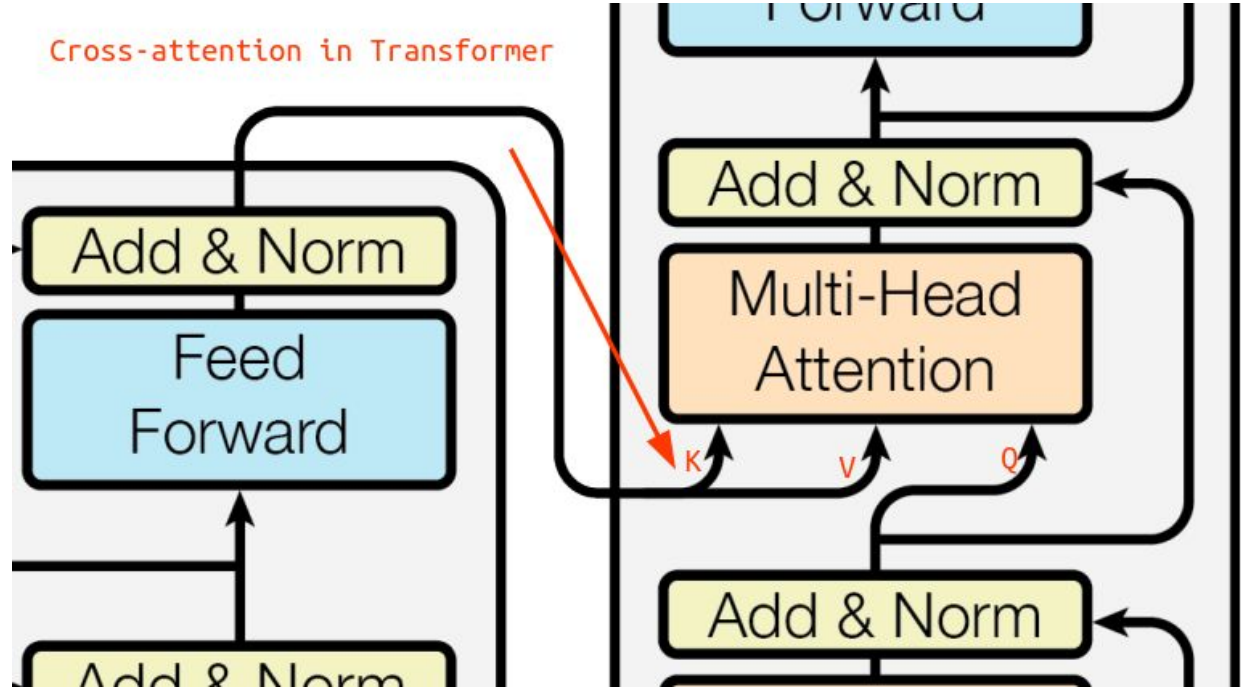


Cross attention

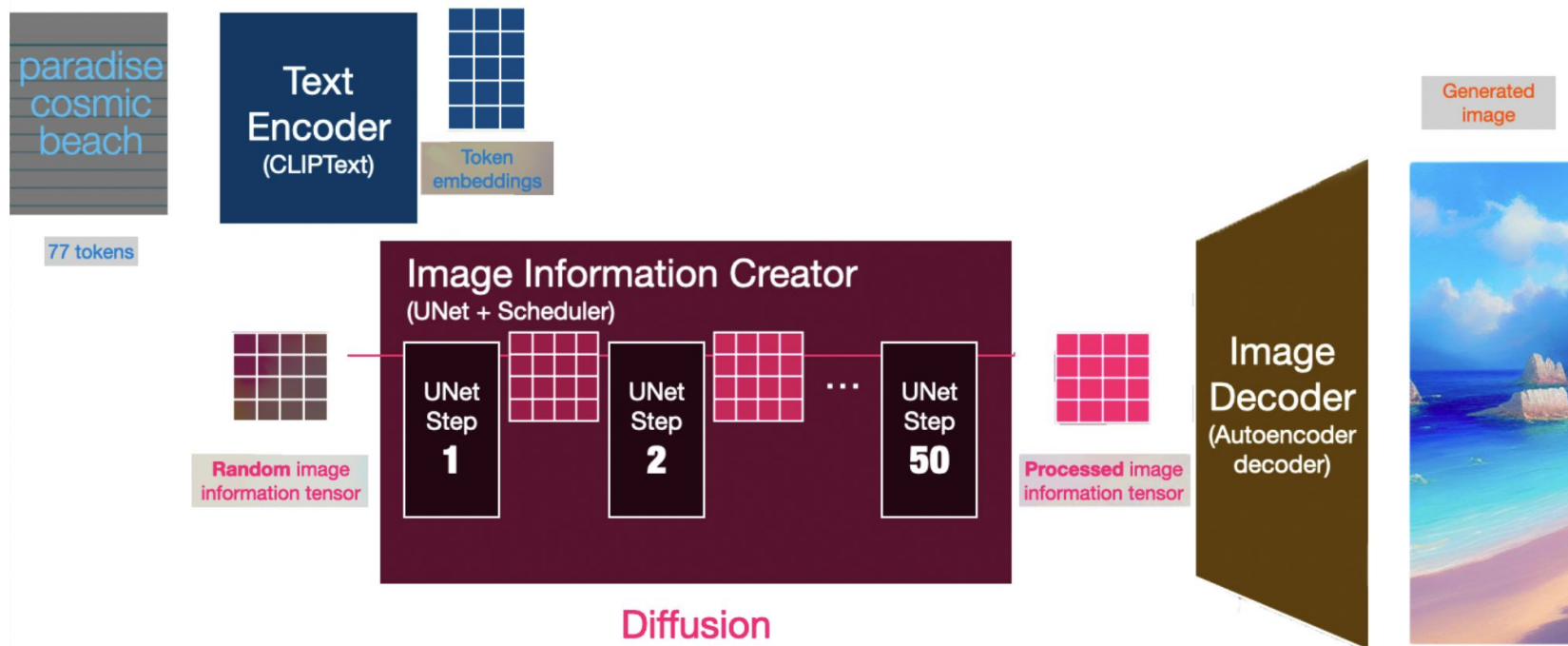
- механизм внимания в архитектуре Transformer, который смешивает две разные последовательности
- две последовательности должны иметь одинаковую размерность
- две последовательности могут иметь разную модальность (например, текст, изображение, звук)
- одна из последовательностей определяет длину вывода, так как она играет роль **query**
- другая последовательность является источником векторов **key value**



Cross-Attention in Transformer Decoder



text2img



Img2img

