

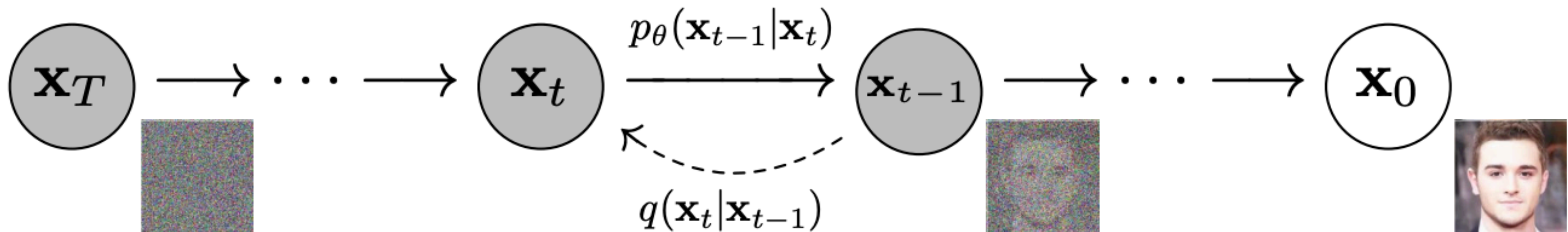
Text diffusion models

ВШЭ ФКН, NLP

Шабалин Александр

Диффузионные модели

- Хотим получить генеративную модель для семплирования из $p(x)$.
- Зададим прямой процесс зашумления и будем учиться его обращать.



Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

$$x_t = \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_{t-1}$$

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

$$\begin{aligned} x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\varepsilon_{t-2}) + \sqrt{1 - \alpha_t}\varepsilon_{t-1} \end{aligned}$$

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

$$\begin{aligned} x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\varepsilon_{t-2}) + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + (\sqrt{\alpha_t(1 - \alpha_{t-1})}\varepsilon_{t-2} + \sqrt{1 - \alpha_t}\varepsilon_{t-1}) \end{aligned}$$

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

$$\begin{aligned} x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\varepsilon_{t-2}) + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \underbrace{(\sqrt{\alpha_t(1 - \alpha_{t-1})}\varepsilon_{t-2})}_{\mathcal{N}(0, \alpha_t(1 - \alpha_{t-1})I)} + \underbrace{\sqrt{1 - \alpha_t}\varepsilon_{t-1}}_{\mathcal{N}(0, (1 - \alpha_t)I)} \end{aligned}$$

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

$$\begin{aligned} x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\varepsilon_{t-2}) + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \underbrace{(\sqrt{\alpha_t(1 - \alpha_{t-1})}\varepsilon_{t-2})}_{\mathcal{N}(0, \alpha_t(1 - \alpha_{t-1})I)} + \underbrace{\sqrt{1 - \alpha_t}\varepsilon_{t-1}}_{\mathcal{N}(0, (1 - \alpha_t)I)} = \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\varepsilon \end{aligned}$$

Вывод алгоритма

Прямой процесс:

$$q(x_t|x_{t-1}) = \mathcal{N}(x_t|\sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)I), \quad \alpha_t \in [0, 1]$$

Нам пригодится $q(x_t|x_0)$, поэтому выведем ее.

$$\begin{aligned} x_t &= \sqrt{\alpha_t}x_{t-1} + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t}(\sqrt{\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_{t-1}}\varepsilon_{t-2}) + \sqrt{1 - \alpha_t}\varepsilon_{t-1} = \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \underbrace{(\sqrt{\alpha_t(1 - \alpha_{t-1})}\varepsilon_{t-2})}_{\mathcal{N}(0, \alpha_t(1 - \alpha_{t-1})I)} + \underbrace{\sqrt{1 - \alpha_t}\varepsilon_{t-1}}_{\mathcal{N}(0, (1 - \alpha_t)I)} = \\ &= \sqrt{\alpha_t\alpha_{t-1}}x_{t-2} + \sqrt{1 - \alpha_t\alpha_{t-1}}\varepsilon = \dots = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon \end{aligned}$$

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

Вывод алгоритма

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

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon \Rightarrow \underline{q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I)}$$

Вывод алгоритма

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

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon \Rightarrow \underline{q(x_t|x_0) = \mathcal{N}(x_t|\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)}$$

Обратный процесс:

$$p(x_{t-1}|x_t) = \mathcal{N}(x_t|\mu(x_t, t), \Sigma(x_t, t))$$

Мы не знаем $\mu(x_t, t)$ и $\Sigma(x_t, t)$ так как не знаем распределение $p(x)$.

Вывод алгоритма

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

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon \Rightarrow \underline{q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I)}$$

Обратный процесс:

$$p(x_{t-1} | x_t) = \mathcal{N}(x_t | \mu(x_t, t), \Sigma(x_t, t))$$

Мы не знаем $\mu(x_t, t)$ и $\Sigma(x_t, t)$ так как не знаем распределение $p(x)$.

Вместо этого найдем $p(x_{t-1} | x_t, x_0)$

$$p(x_{t-1} | x_t, x_0) = \mathcal{N}(x_t | \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

Вывод алгоритма

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}(x_t|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

Вывод алгоритма

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}(x_t|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) = \frac{p(x_t|x_{t-1}, x_0)p(x_{t-1}|x_0)}{p(x_t|x_0)}$$

Вывод алгоритма

$$\mathcal{N}(x \mid \mu, \sigma) = \frac{1}{Z} \exp \left(- \frac{(x - \mu)^2}{2\sigma} \right)$$

$$\mathcal{N}(x_t \mid \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t)I) \quad \mathcal{N}(x_{t-1} \mid \sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1})I)$$

$$\Downarrow$$
$$\Downarrow$$

$$p(x_{t-1} \mid x_t, x_0) = \mathcal{N}(x_t \mid \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) = \frac{p(x_t \mid x_{t-1}, x_0) p(x_{t-1} \mid x_0)}{p(x_t \mid x_0)}$$

$$\Downarrow$$

$$\mathcal{N}(x_t \mid \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I)$$

Вывод алгоритма

$$\mathcal{N}(x \mid \mu, \sigma) = \frac{1}{Z} \exp \left(- \frac{(x - \mu)^2}{2\sigma} \right)$$

$$\mathcal{N}(x_t \mid \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t)I) \quad \mathcal{N}(x_{t-1} \mid \sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1})I)$$

$$\parallel$$
$$\parallel$$

$$p(x_{t-1} \mid x_t, x_0) = \mathcal{N}(x_t \mid \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) = \frac{p(x_t \mid x_{t-1}, x_0) p(x_{t-1} \mid x_0)}{p(x_t \mid x_0)} \propto$$

$$\parallel$$

$$\mathcal{N}(x_t \mid \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I)$$

$$\propto \exp \left[- \frac{1}{2} \left(\frac{(x_t - \sqrt{\alpha_t} x_{t-1})^2}{1 - \alpha_t} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t} x_0)^2}{1 - \bar{\alpha}_t} \right) \right] =$$

Вывод алгоритма

$$\mathcal{N}(x \mid \mu, \sigma) = \frac{1}{Z} \exp \left(- \frac{(x - \mu)^2}{2\sigma} \right)$$

$$\mathcal{N}(x_t \mid \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t)I) \quad \mathcal{N}(x_{t-1} \mid \sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1})I)$$

\parallel

\parallel

$$p(x_{t-1} \mid x_t, x_0) = \mathcal{N}(x_t \mid \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) = \frac{p(x_t \mid x_{t-1}, x_0) p(x_{t-1} \mid x_0)}{p(x_t \mid x_0)} \propto$$

\parallel

$$\mathcal{N}(x_t \mid \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I)$$

$$\propto \exp \left[- \frac{1}{2} \left(\frac{(x_t - \sqrt{\alpha_t} x_{t-1})^2}{1 - \alpha_t} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t} x_0)^2}{1 - \bar{\alpha}_t} \right) \right] =$$

$$= \exp \left[- \frac{1}{2} \left(\left(\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}}{1 - \alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) x_{t-1} + C(x_t, x_0) \right) \right]$$

Вывод алгоритма

$$\mathcal{N}(x \mid \mu, \sigma) = \frac{1}{Z} \exp \left(- \frac{(x - \mu)^2}{2\sigma} \right)$$

$$\mathcal{N}(x_t \mid \sqrt{\alpha_t} x_{t-1}, (1 - \alpha_t)I) \quad \mathcal{N}(x_{t-1} \mid \sqrt{\bar{\alpha}_{t-1}} x_0, (1 - \bar{\alpha}_{t-1})I)$$

$$\Downarrow$$

$$\Downarrow$$

$$p(x_{t-1} \mid x_t, x_0) = \mathcal{N}(x_t \mid \tilde{\mu}(x_t, x_0), \tilde{\beta}_t I) = \frac{p(x_t \mid x_{t-1}, x_0) p(x_{t-1} \mid x_0)}{p(x_t \mid x_0)} \propto$$

$$\Downarrow$$

$$\mathcal{N}(x_t \mid \sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t)I)$$

$$\propto \exp \left[- \frac{1}{2} \left(\frac{(x_t - \sqrt{\alpha_t} x_{t-1})^2}{1 - \alpha_t} + \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0)^2}{1 - \bar{\alpha}_{t-1}} - \frac{(x_t - \sqrt{\bar{\alpha}_t} x_0)^2}{1 - \bar{\alpha}_t} \right) \right] =$$

$$= \exp \left[- \frac{1}{2} \left(\left(\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1}^2 - 2 \left(\frac{\sqrt{\alpha_t}}{1 - \alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) x_{t-1} + C(x_t, x_0) \right) \right]$$

$$\tilde{\beta}_t = \frac{1}{\frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}}}$$

$$\tilde{\mu}(x_t, x_0) = \left(\frac{\sqrt{\alpha_t}}{1 - \alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) \tilde{\beta}_t$$

Вывод алгоритма

$$\tilde{\beta}_t = \frac{1}{\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}}$$

$$\tilde{\mu}(x_t, x_0) = \left(\frac{\sqrt{\alpha_t}}{1-\alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}} x_0 \right) \tilde{\beta}_t$$

Вывод алгоритма

$$\tilde{\beta}_t = \frac{1}{\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}} = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t} = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}$$

$$\tilde{\mu}(x_t, x_0) = \left(\frac{\sqrt{\alpha_t}}{1-\alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}} x_0 \right) \tilde{\beta}_t$$

Вывод алгоритма

$$\tilde{\beta}_t = \frac{1}{\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}} = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t} = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}$$

$$\begin{aligned}\tilde{\mu}(x_t, x_0) &= \left(\frac{\sqrt{\alpha_t}}{1-\alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}} x_0 \right) \tilde{\beta}_t = \\ &= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)}{1-\bar{\alpha}_t} x_0 =\end{aligned}$$

Вывод алгоритма

$$\tilde{\beta}_t = \frac{1}{\frac{\alpha_t}{1-\alpha_t} + \frac{1}{1-\bar{\alpha}_{t-1}}} = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{\alpha_t(1-\bar{\alpha}_{t-1}) + 1-\alpha_t} = \frac{(1-\alpha_t)(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t}$$

$$\begin{aligned}\tilde{\mu}(x_t, x_0) &= \left(\frac{\sqrt{\alpha_t}}{1-\alpha_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1-\bar{\alpha}_{t-1}} x_0 \right) \tilde{\beta}_t = \\ &= \frac{\sqrt{\alpha_t}(1-\bar{\alpha}_{t-1})}{1-\bar{\alpha}_t} x_t + \frac{\sqrt{\bar{\alpha}_{t-1}}(1-\alpha_t)}{1-\bar{\alpha}_t} x_0 = \\ &\quad \left[x_0 = \frac{x_t - \sqrt{1-\bar{\alpha}_t} \varepsilon_t}{\sqrt{\bar{\alpha}_t}} \right] \\ &= \frac{1}{\bar{\alpha}_t} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_t \right)\end{aligned}$$

Почти готово

Что теперь?

Почти готово

Что теперь?

У нас есть явная формула для прямого процесса

$$q(x_t|x_0) = \mathcal{N}(x_t|\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

И для обратного

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}(x_t|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

Почти готово

Что теперь?

У нас есть явная формула для прямого процесса

$$q(x_t|x_0) = \mathcal{N}(x_t|\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$$

И для обратного

$$p(x_{t-1}|x_t, x_0) = \mathcal{N}(x_t|\tilde{\mu}(x_t, x_0), \tilde{\beta}_t I)$$

$\tilde{\mu}(x_t, x_0)$ мы не знаем, поэтому обучим нейронную сеть для аппроксимации.

$$L_t = \|\tilde{\mu}(x_t, x_0) - \tilde{\mu}_\theta(x_t, t)\|^2 \quad \text{или} \quad L_t = \|\varepsilon_t - \varepsilon_\theta(x_t, t)\|^2$$

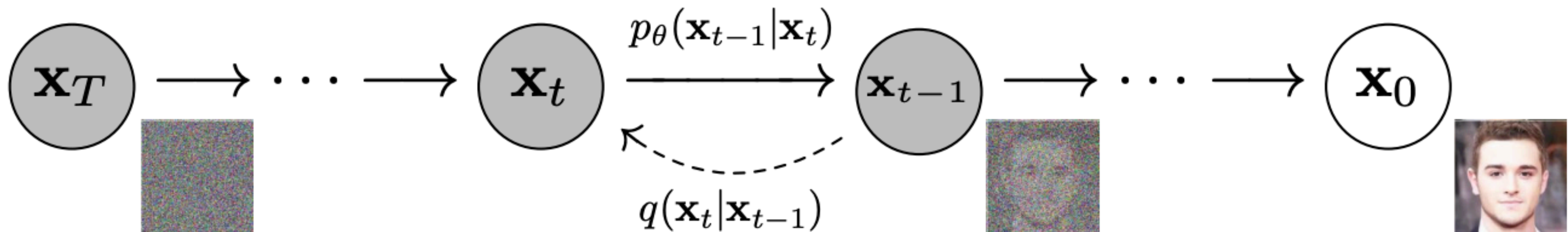
Обучение и инференс

Algorithm 1 Training

```
1: repeat  
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$   
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$   
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
5:   Take gradient descent step on  
        $\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t)\|^2$   
6: until converged
```

Algorithm 2 Sampling

```
1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$   
2: for  $t = T, \dots, 1$  do  
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$   
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$   
5: end for  
6: return  $\mathbf{x}_0$ 
```



Text diffusion models

Зачем?

- Авторегрессионные модели генерируют каждый токен отдельно. Это дорого.
- Авторегрессионные модели не могут исправлять свои ошибки.

Особенности

- Тексты имеют дискретную структуру => недифференцируемы
- Непонятно, как зашумлять текст
- ~~Текстовая диффузия не работает~~

Два вида текстовых диффузионных моделей

- С дискретным зашумлением
- С гауссовским зашумлением

Дискретный шум

- Можно смотреть на каждый токен, как на семпл из категориального распределение всех ВОЗМОЖНЫХ ТОКЕНОВ.
- В точке x_0 распределение вырожденное.
- Добавление шума – релаксация данного распределения.

Multinomial Diffusion

- Постепенно сводит вырожденное категориальное распределение к равномерному.

$$q(x_t|x_{t-1}) = \mathcal{C}(x_t|\alpha_t x_{t-1} + (1 - \alpha_t)/K)$$

$$q(x_t|x_0) = \mathcal{C}(x_t|\bar{\alpha}_t x_{t-1} + (1 - \bar{\alpha}_t)/K)$$

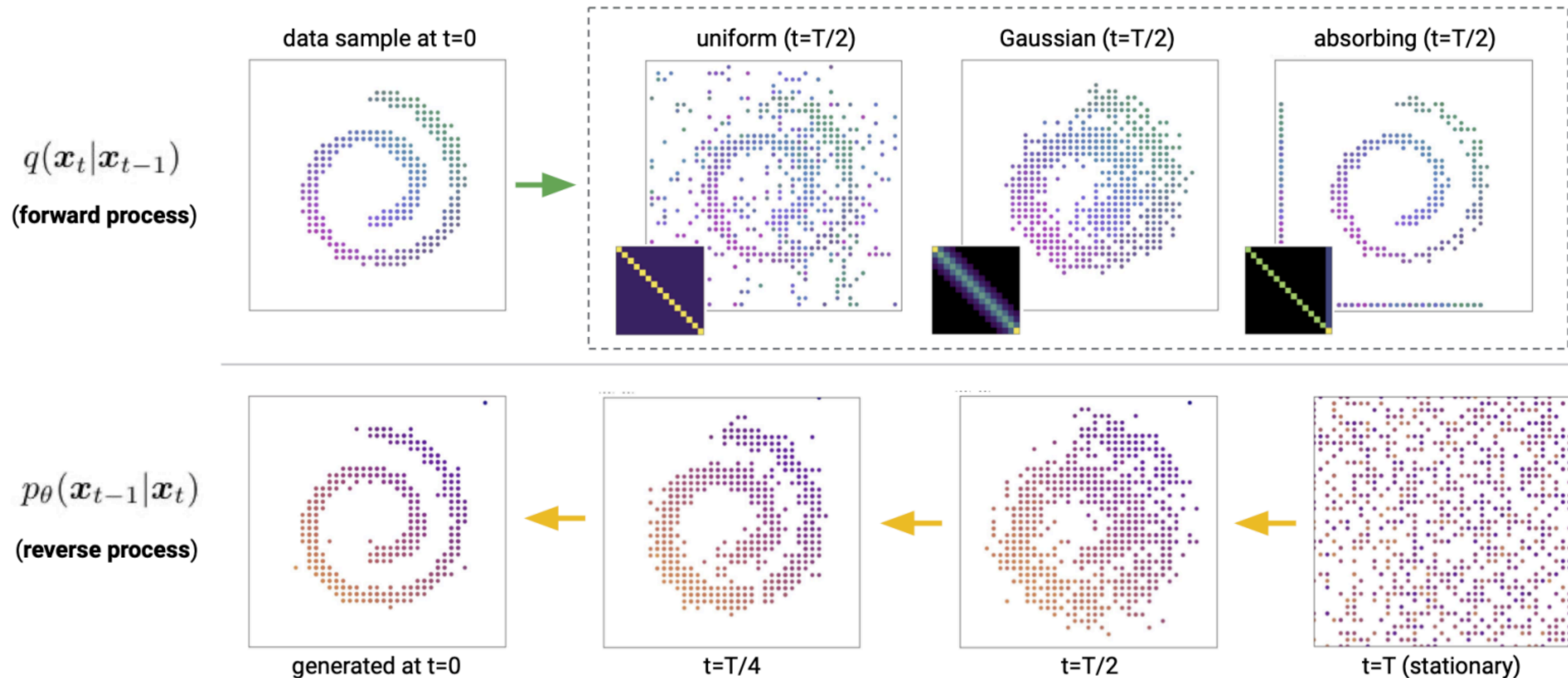
$$p(x_{t-1}|x_t, x_0) = \mathcal{C}(x_{t-1}|D_{\text{post}}(x_t, x_0))$$

$$D_{\text{post}}(x_t, x_0) = \frac{1}{Z} [\alpha_t x_t + (1 - \alpha_t)/K] \odot [\bar{\alpha}_{t-1} x_0 + (1 - \bar{\alpha}_{t-1})/K]$$

D3PM

Discrete Denoising Diffusion Probabilistic Model

- Пробуют разные способы зашумления

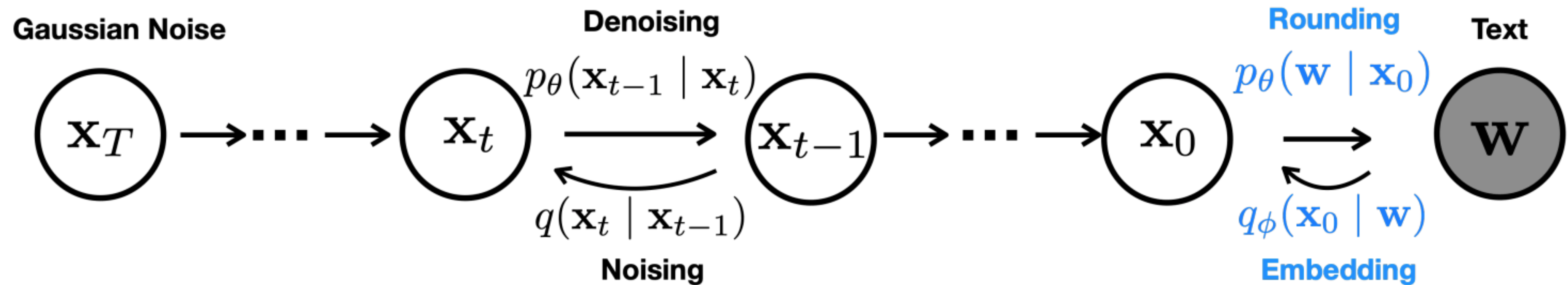


Непрерывный шум

- Токены переводятся в пространство эмбеддингов (или энкодингов)
- В этом пространстве учится гауссовская диффузия
- К выходам диффузии применяется декодер
- Работает намного лучше диффузии с дискретным шумом!

Diffusion-LM

- Первая относительно удачная попытка текстовой диффузии.

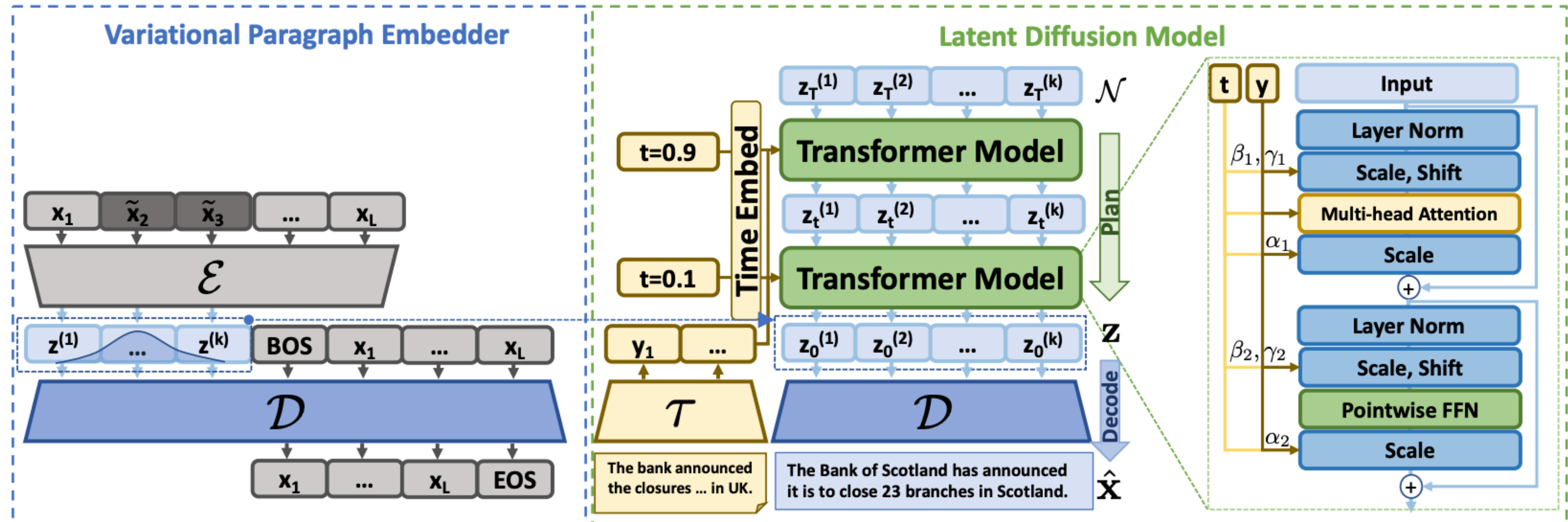


$$\mathcal{L}_{\text{simple}}(\mathbf{x}_0) = \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \|\mu_\theta(\mathbf{x}_t, t) - \hat{\mu}(\mathbf{x}_t, \mathbf{x}_0)\|^2.$$

$$\mathcal{L}_{\text{simple}}^{\text{e2e}}(\mathbf{w}) = \mathbb{E}_{q_\phi(\mathbf{x}_{0:T} | \mathbf{w})} [\mathcal{L}_{\text{simple}}(\mathbf{x}_0) + \|\text{EMB}(\mathbf{w}) - \mu_\theta(\mathbf{x}_1, 1)\|^2 - \log p_\theta(\mathbf{w} | \mathbf{x}_0)]$$

PLANNER

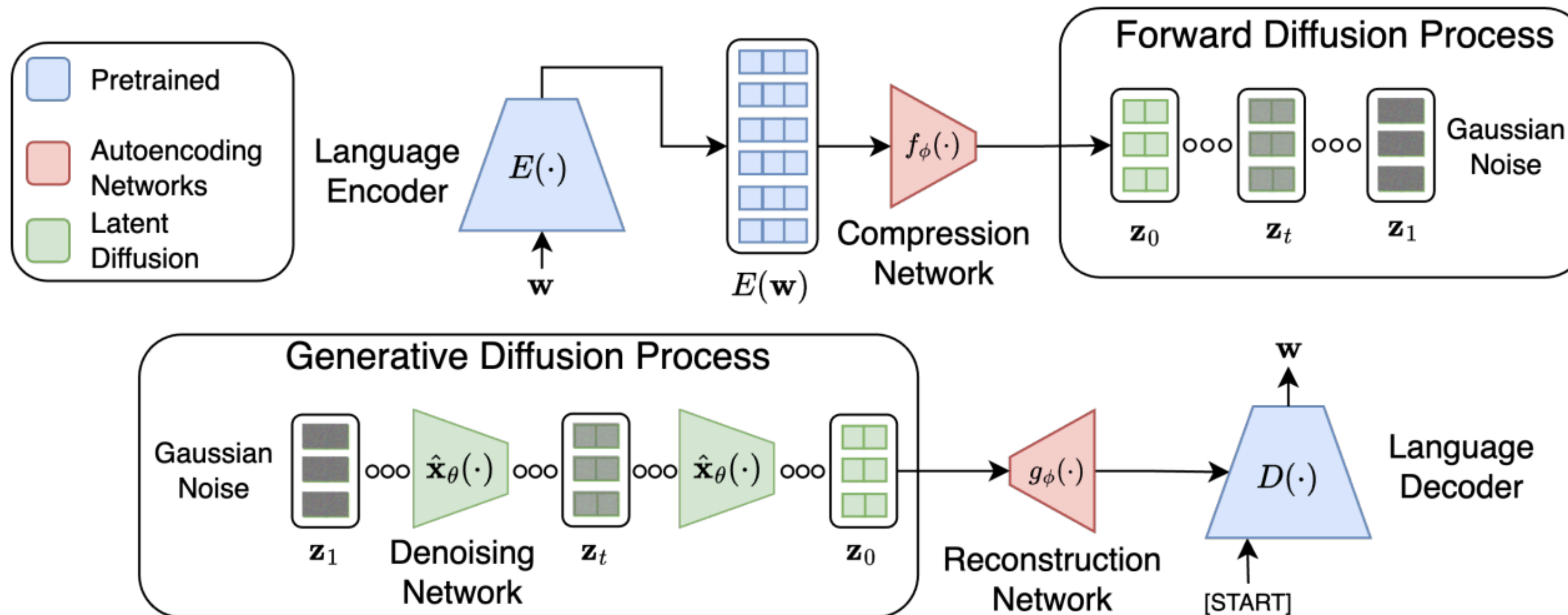
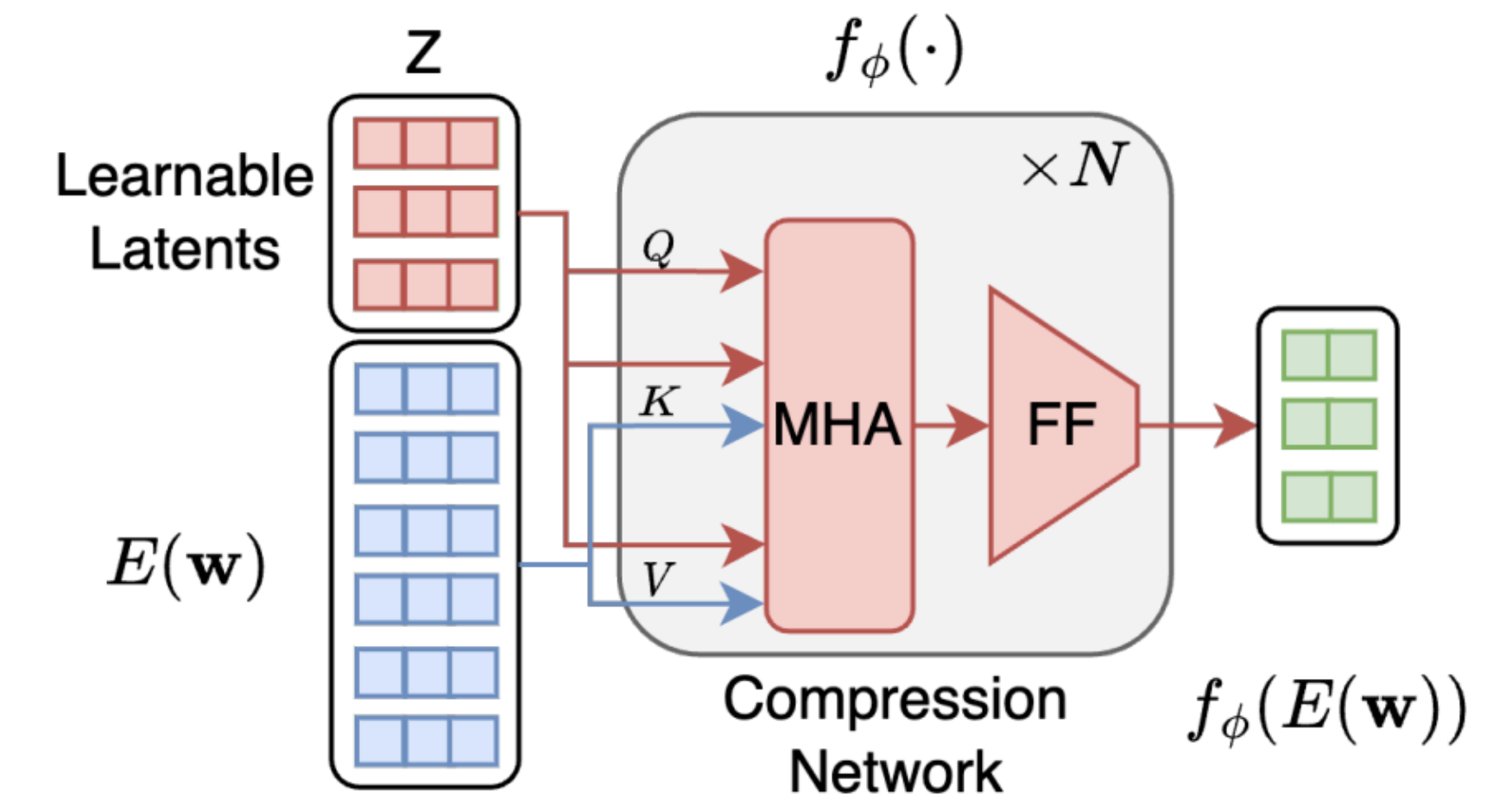
- VAE + autoregressive decoder



LD4LG

Latent Diffusion for Language Generation

- Latent space + autoregressive decoder



Трюки для текстовой диффузии

Трюки для обучения

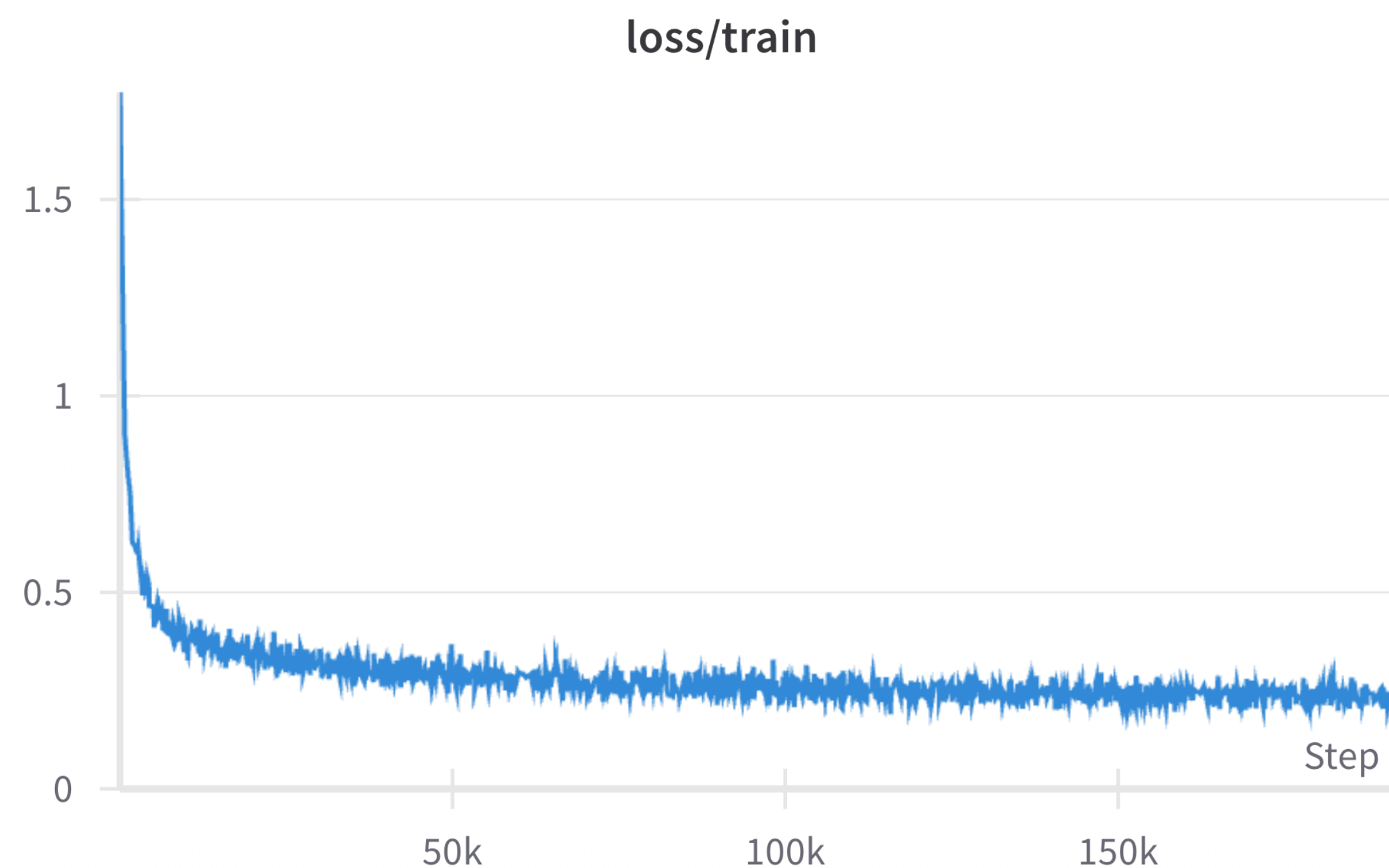
Предсказание x_0 вместо ε_t работает эффективнее для диффузии на дискретных данных.

Никто толком не понимает, почему.

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\varepsilon_t$$

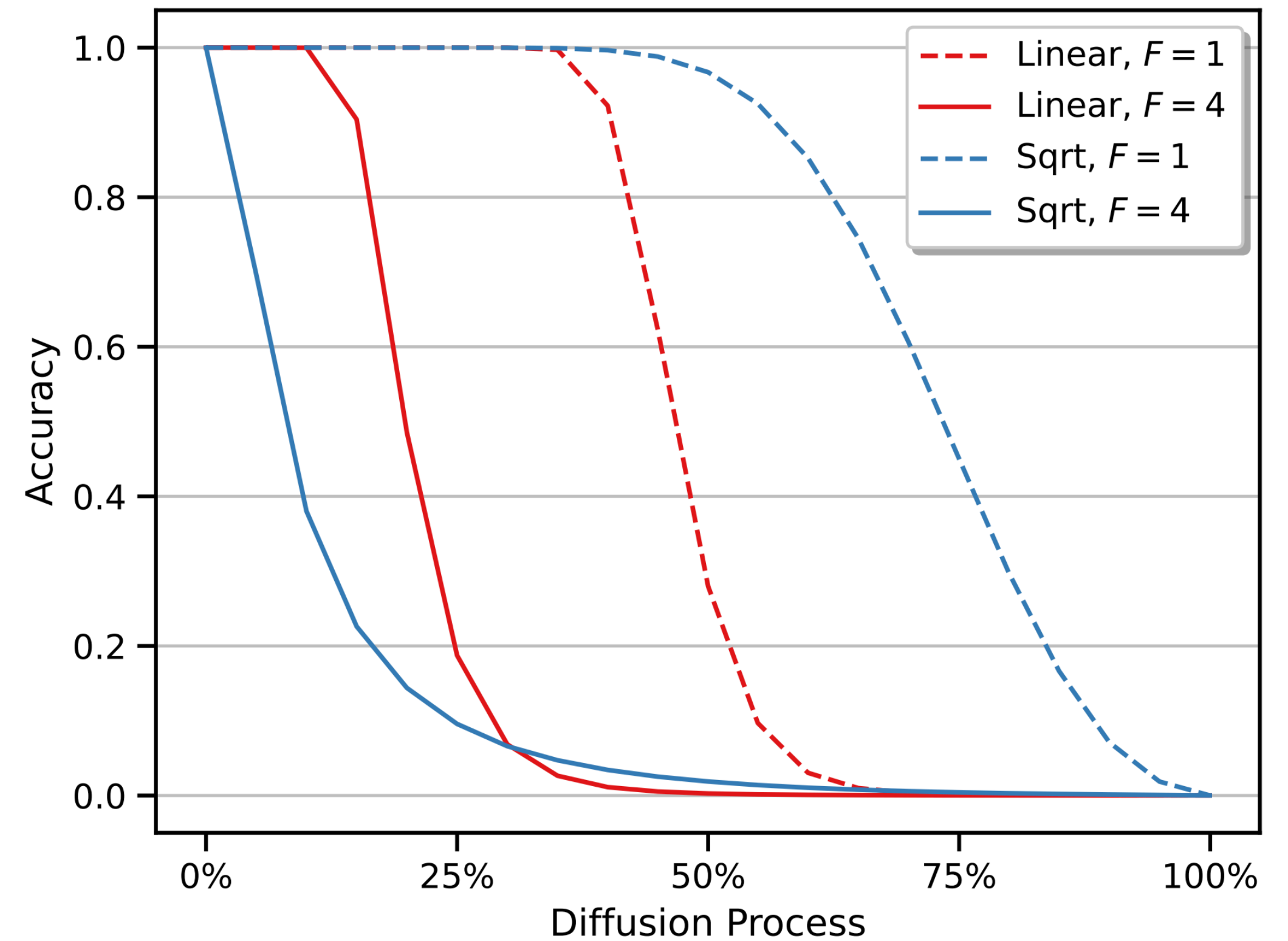
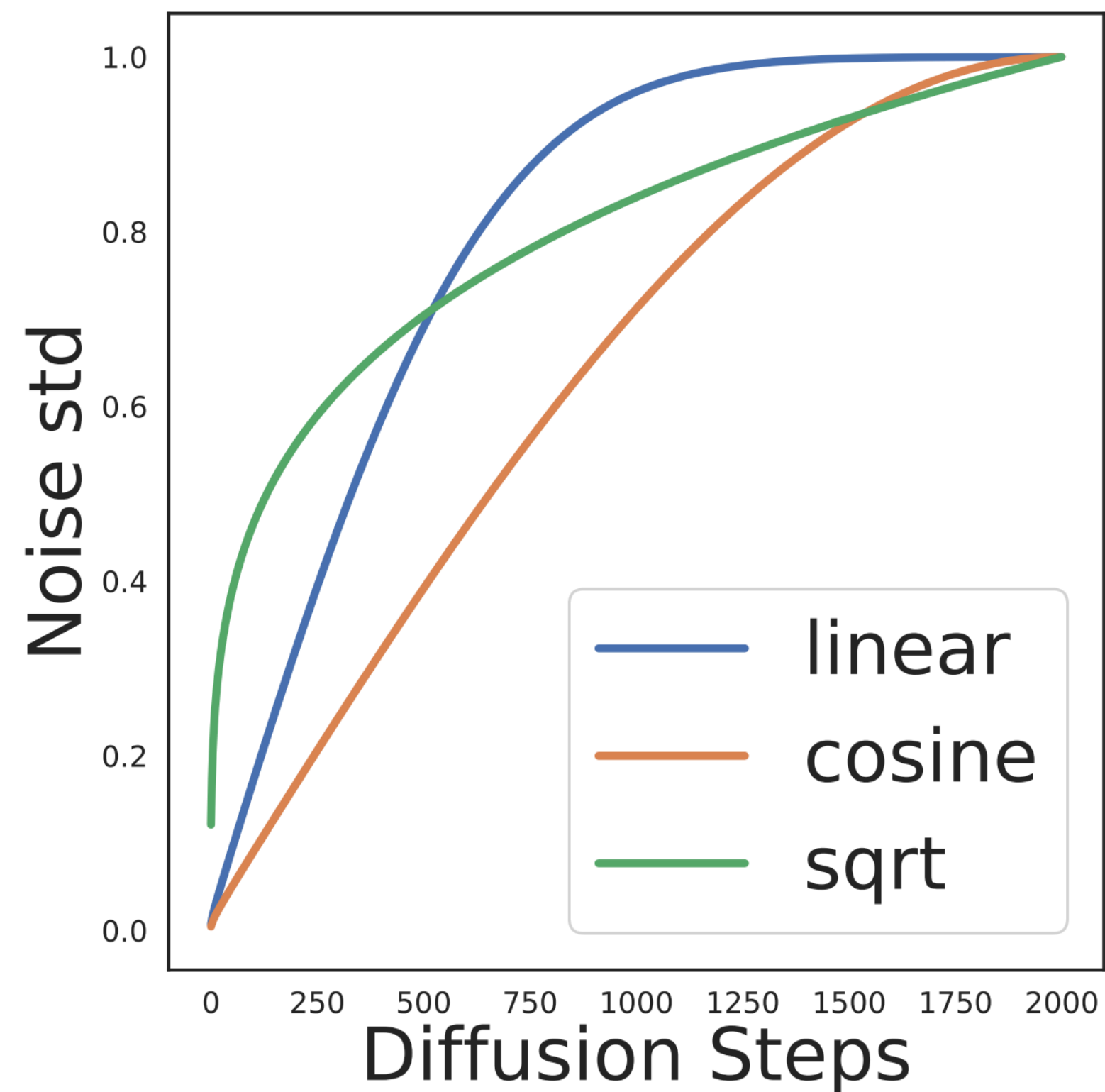
Трюки для обучения

Диффузионные модели надо учить **долго**.
Даже если лосс не падает.



Трюки для обучения

Лучше брать более агрессивное расписание зашумления.



Schedules	BLEU	
	$F = 1$	$F = 4$
<i>Linear</i> (Ho et al., 2020)	31.02	31.98
<i>Cosine</i> (Nichol & Dhariwal, 2021)	26.61	32.68
<i>Sqrt</i> (Li et al., 2022)	2.85	34.13
<i>EDM</i> (Karras et al., 2022)	29.16	30.09

Трюки для обучения

Нормализация энкодингов

- x_0 должен иметь единичную дисперсию.
- Нормировать лучше по всем признакам отдельно, так как они отличаются по норме.

Трюки для обучения

Clamping trick

Округление предсказанного \hat{x}_0 до ближайшего текста на каждом шаге.

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}} \cdot \text{Clamp}(f_{\theta}(\mathbf{x}_t, t)) + \sqrt{1 - \bar{\alpha}}\epsilon$$

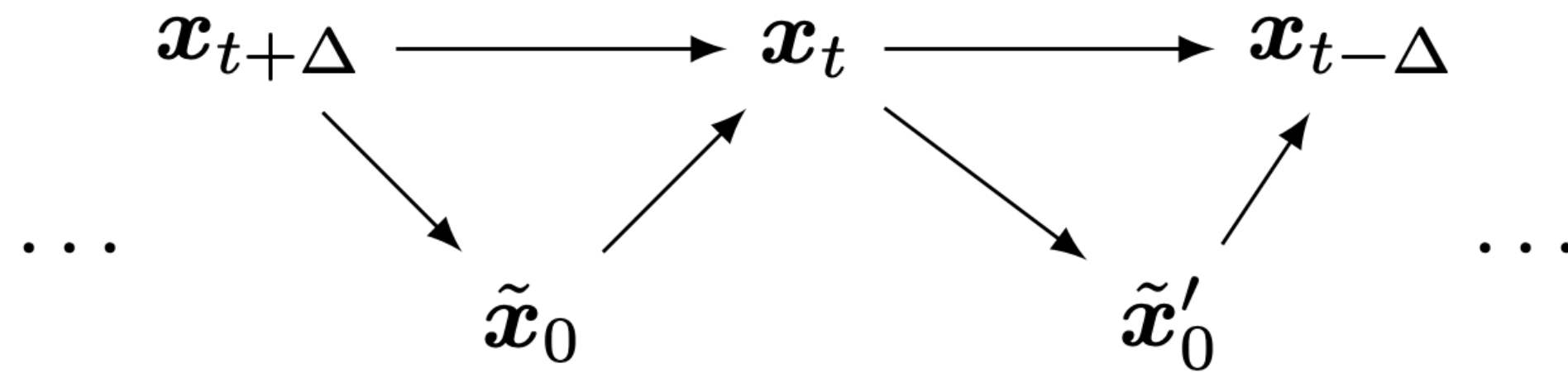
Архитектурные трюки: Encoder

- Диффузия работает лучше, если латентное пространство более гладкое.
 - ➡ Encodings > embeddings
 - ➡ Encoder должен для "похожих" текстов получать похожие латенты
 - ➡ Учить encoder вместе с диффузией – плохая идея

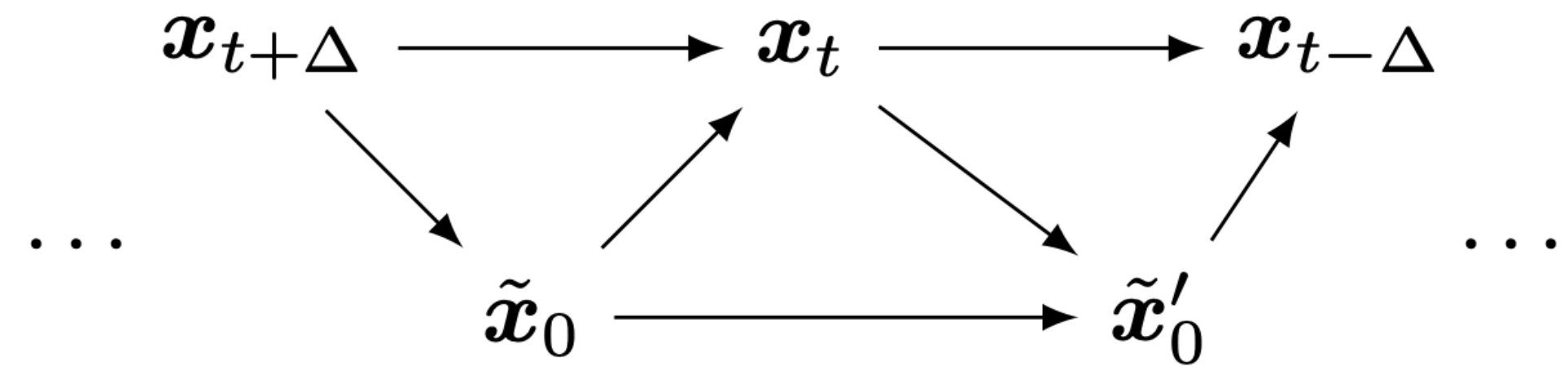
Архитектурные трюки: Decoder

- Диффузионная модель аппроксимирует распределение $p(x)$.
 - Качество аппроксимации падает с уменьшением числа шагов.
 - При декодировании латентов хочется нивелировать ошибки.
- ➡ Лучше брать сложный декодер (несколько слоев трансформера)

Архитектурные трюки: Self-condition

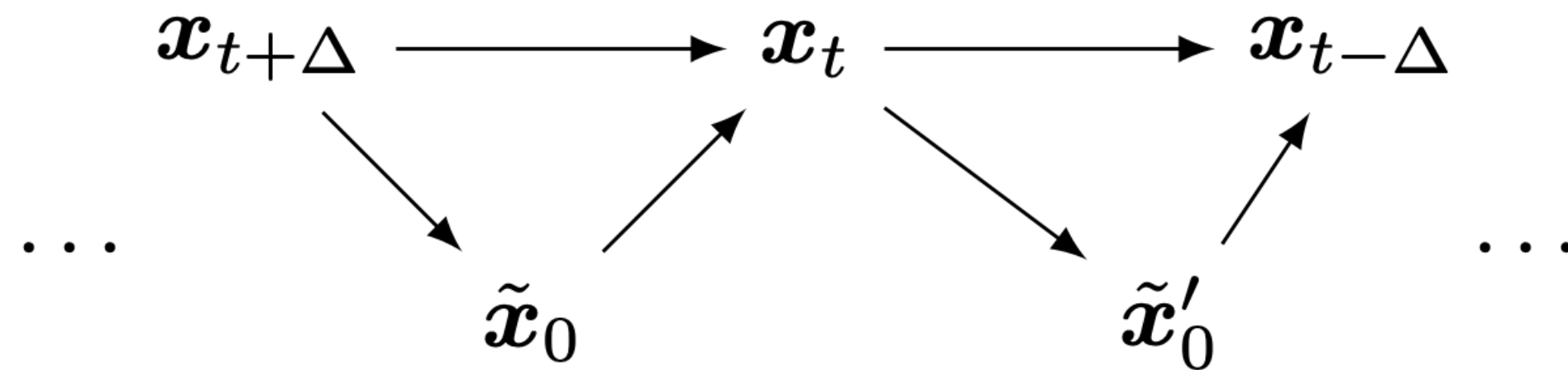


(a) Standard reverse diffusion steps.

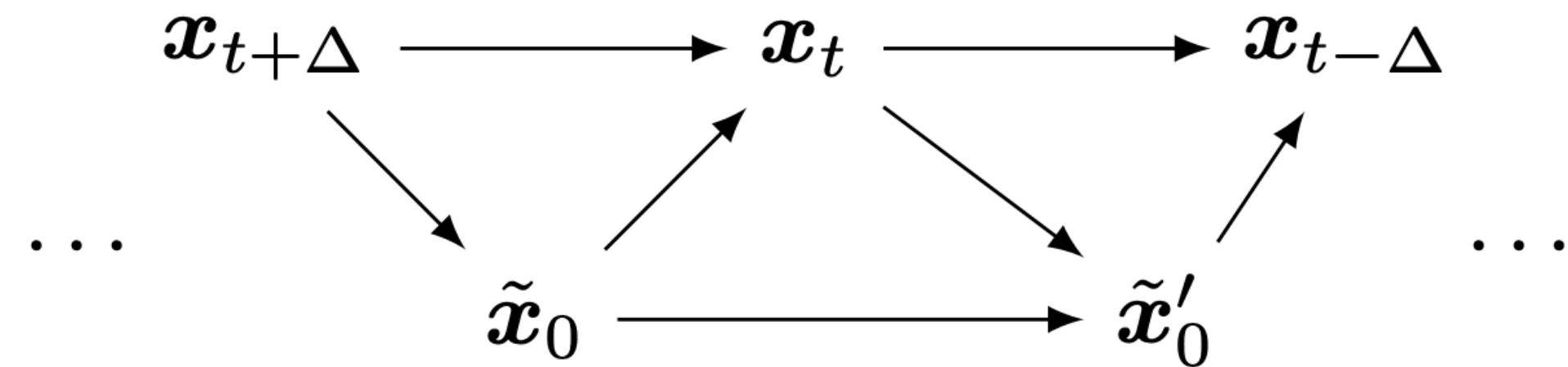


(b) Self-Conditioning on the previous \mathbf{x}_0 estimate.

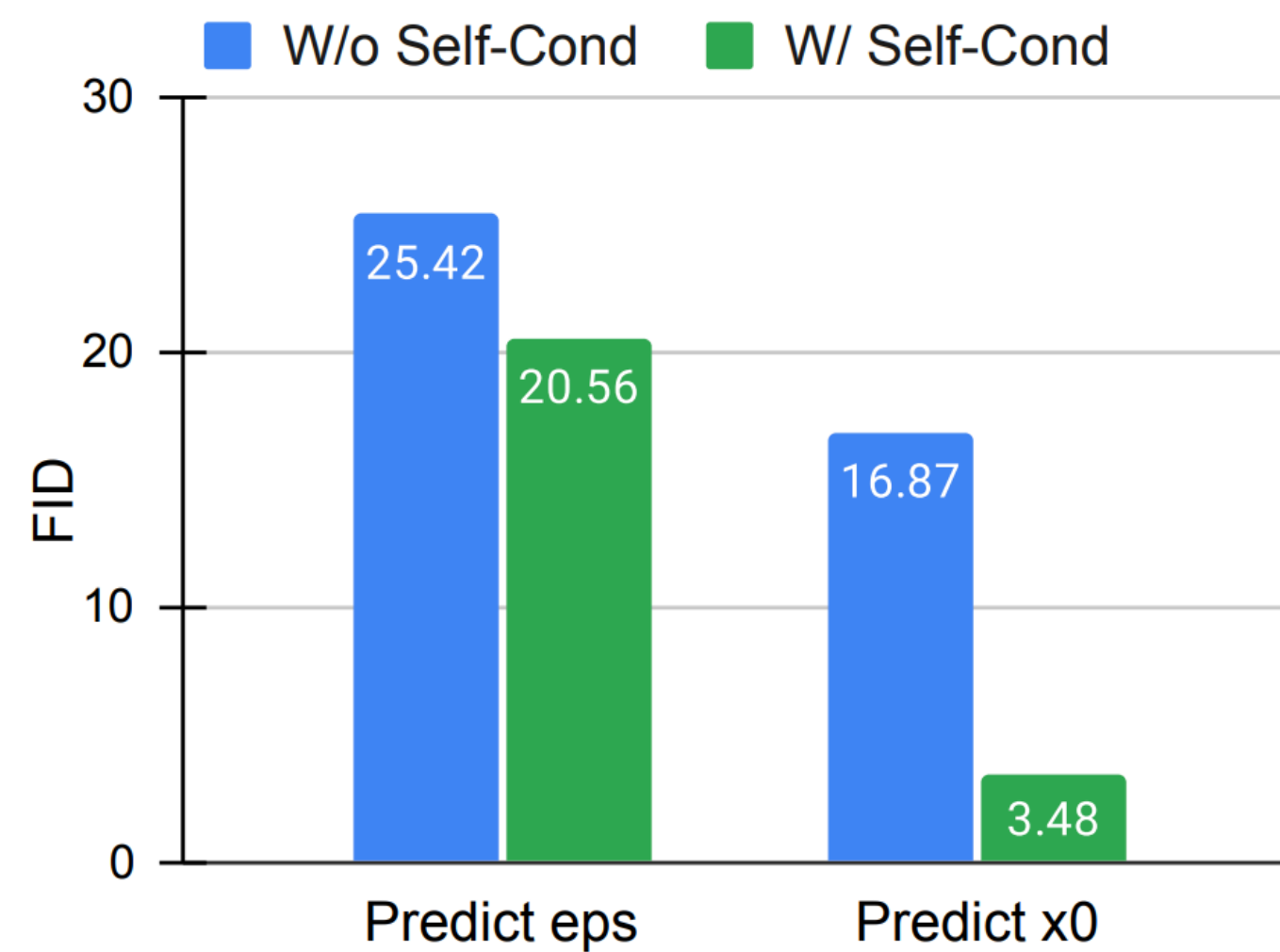
Архитектурные трюки: Self-condition



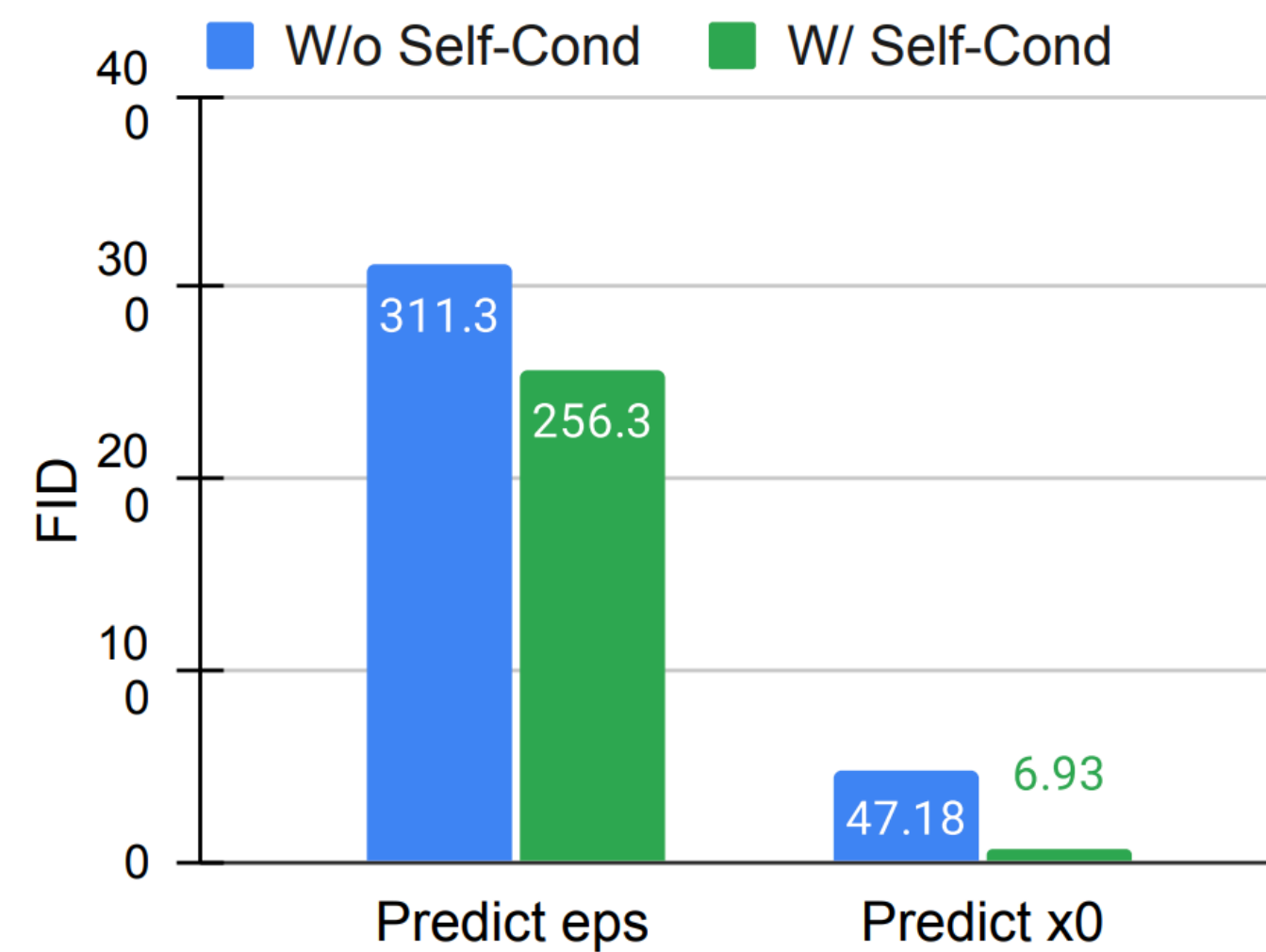
(a) Standard reverse diffusion steps.



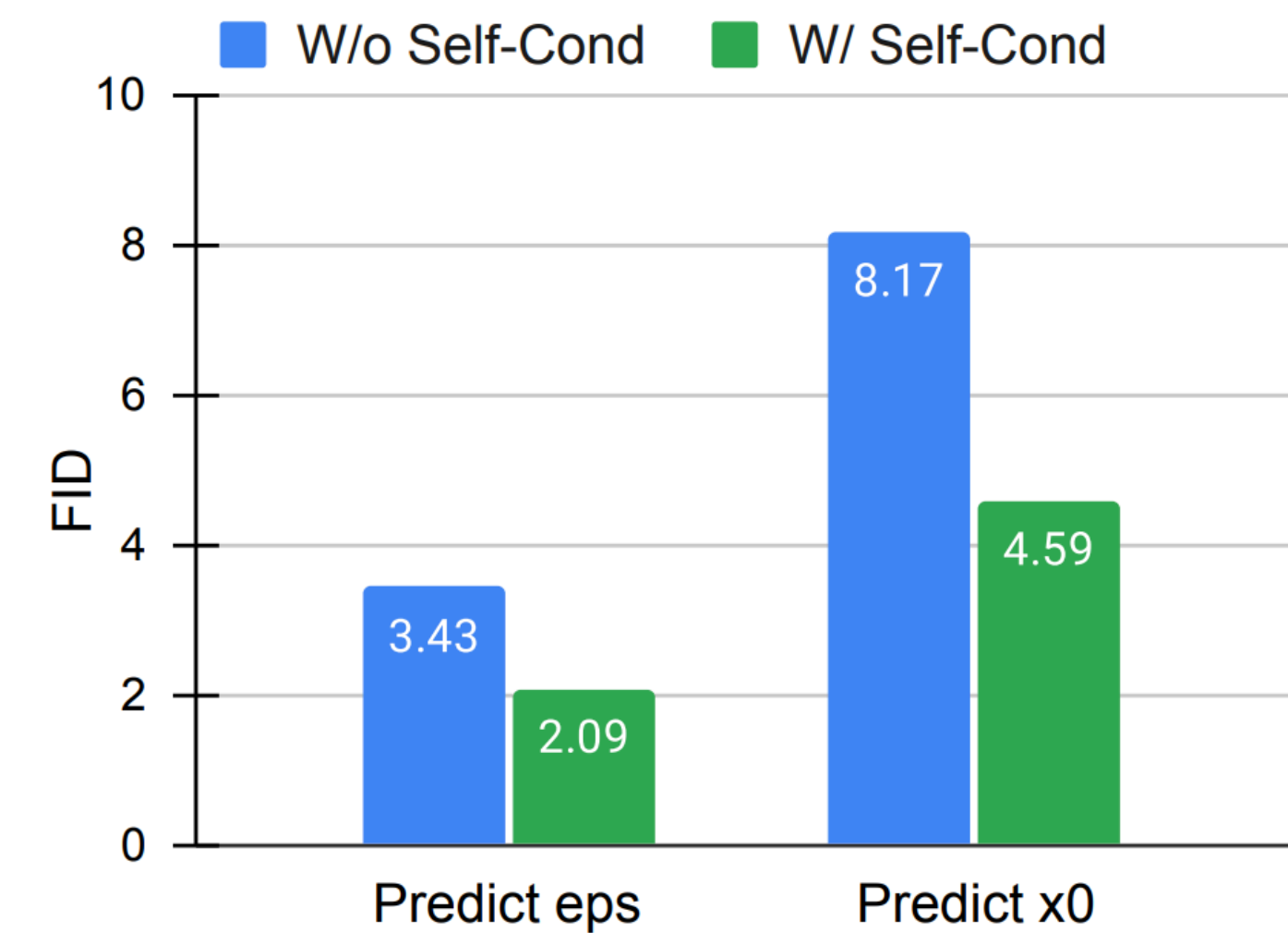
(b) Self-Conditioning on the previous x_0 estimate.



(a) CIFAR-10, UINT8.



(b) CIFAR-10, UINT8 (RAND).



(c) IMAGENET 64×64 .