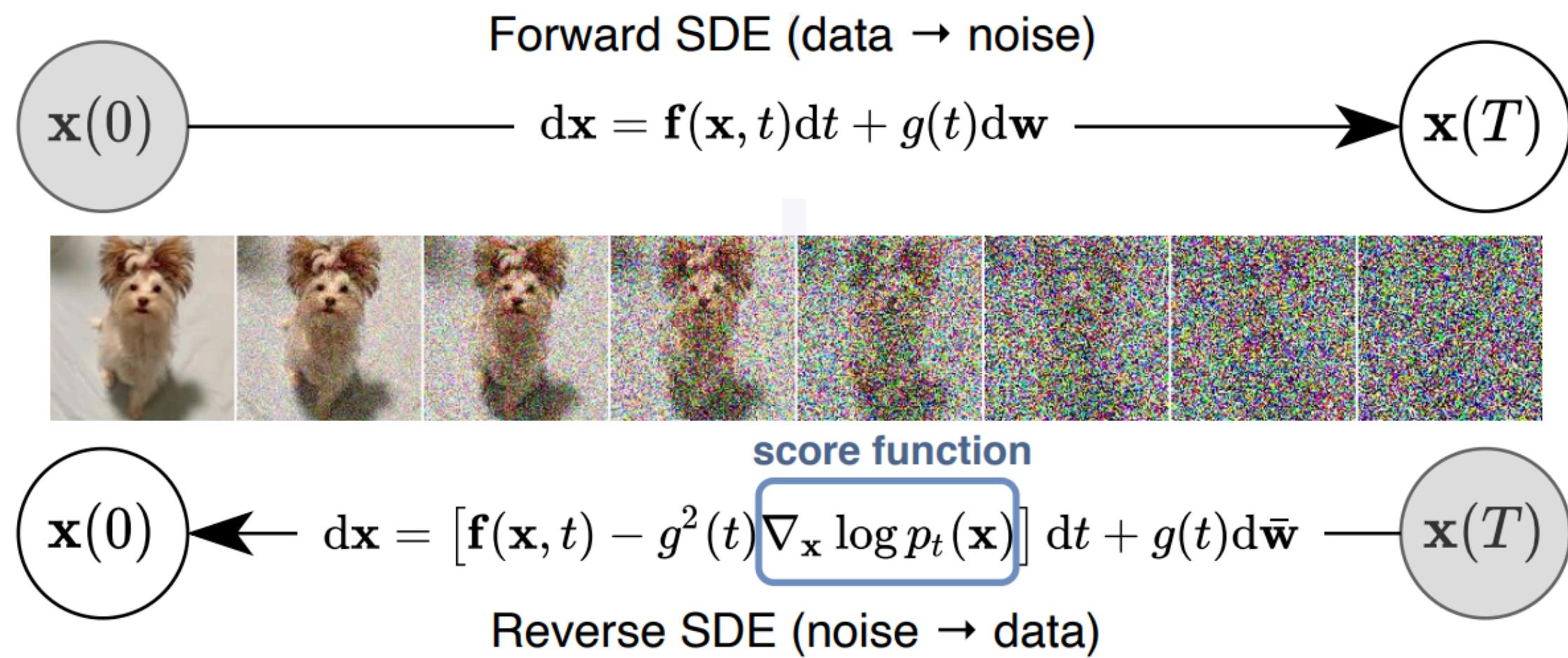


Diffusion models in latent space

02.12.2022

Denoising Score-Matching

At first glance



- Marginal kernel: $q(\mathbf{x}_t | \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t | \alpha_t \mathbf{x}_0, \sigma_t^2 I)$
- Minimisation objective:

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\lambda(t) \mathbb{E}_{q(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\left\| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t) \right\|_2^2 \right] \right]$$

Denoising Score-Matching

More details

- What is $q(\mathbf{x}_t)$? All, we know - $q(\mathbf{x}_0), q(\mathbf{x}_t | \mathbf{x}_0)$

- So training objective can be rewritten as:

$$\min_{\theta} \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\lambda(t) \mathbb{E}_{q(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\left\| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) - \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t) \right\|_2^2 \right] \right] + C$$

- It's not just a random constant, we don't care about (in app.):

$$C = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\lambda(t) \mathbb{E}_{q(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{x}_t | \mathbf{x}_0)} \left[\left\| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) \right\|_2^2 - \left\| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0) \right\|_2^2 \right] \right]$$

- Complex derivation from Song in “Maximum likelihood training of score-based diffusion models”:

$$\text{KL}\left(q(\mathbf{x}_0) \| p_{\theta}(\mathbf{x}_0)\right) \leq \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{x}_0)} \mathbb{E}_{q(\mathbf{z}_t | \mathbf{x}_0)} \left[\left\| \nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p_{\theta}(\mathbf{x}_t) \right\|_2^2 \right] \right]$$

Variational AutoEncoder

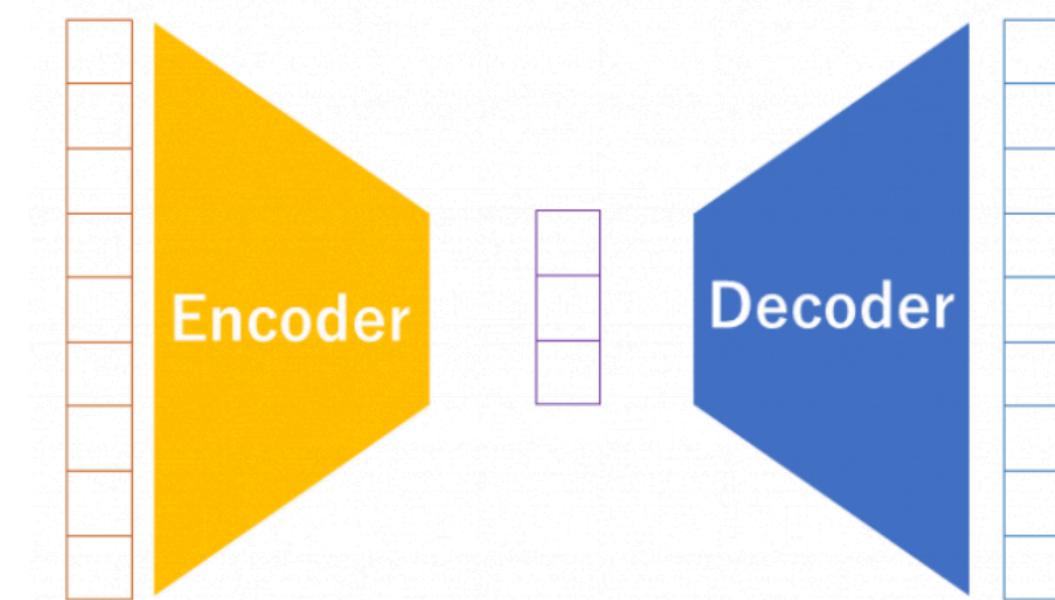
Cause it's likelihood-based generative model

- Root idea: $\log p(x) = ELBO + KL(q(z|x)||p(z|x))$

$$\begin{aligned} ELBO &= \mathbb{E}_{q(z|x)} [\log p(x, z) - \log q(z|x)] = \\ &= \mathbb{E}_{q(z|x)} [\log p(x|z) + \log p(z) - \log q(z|x)] \\ &= \mathbb{E}_{q(z|x)} \log p(x|z) - KL(q(z|x)||p(z)) \end{aligned}$$

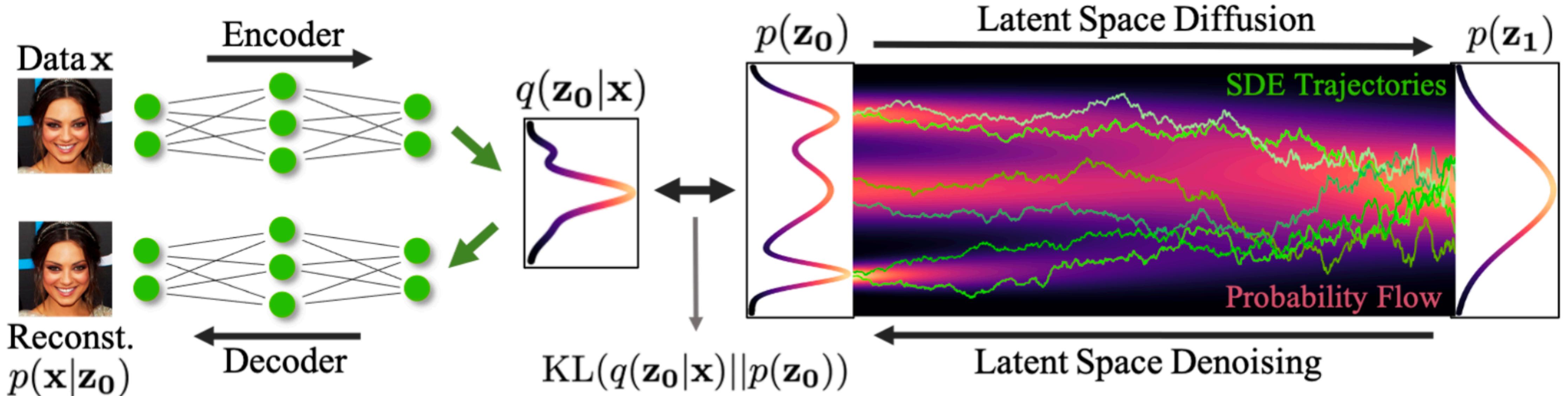
-

- Training objective looks like: $\mathbb{E}_{q_\phi(z|x)} [-\log p_\psi(x|z)] + KL(q_\phi(z|x)||p(z))$



Latent Score-based Generative Model

...instead of a thousands equations

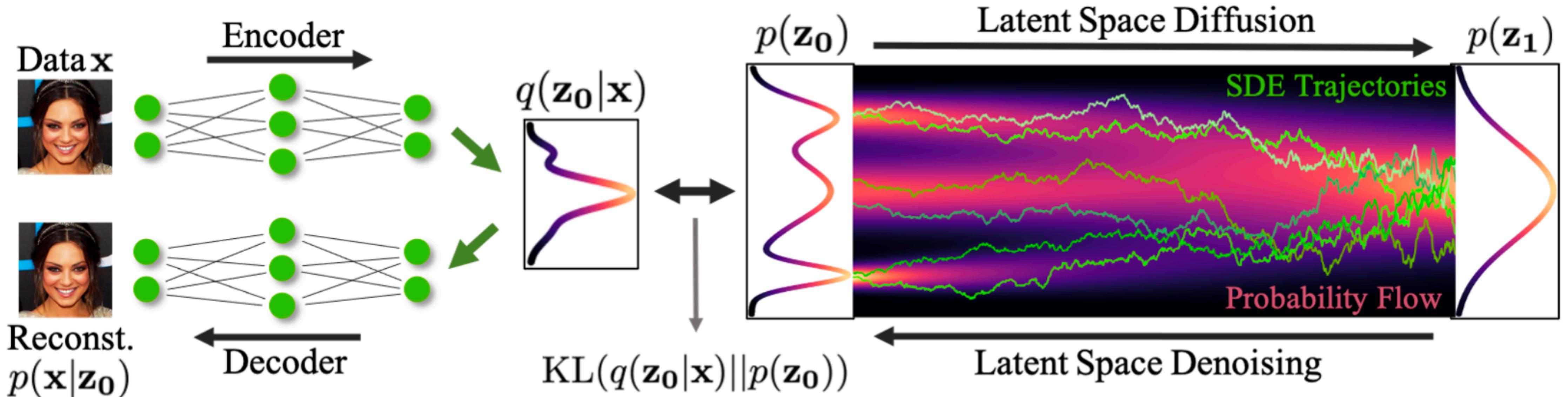


$$\mathbb{E}_{q_\phi(\mathbf{z}_0 | \mathbf{x})} \left[-\log p_\psi(\mathbf{x} | \mathbf{z}_0) \right] + KL(q_\phi(\mathbf{z}_0 | \mathbf{x}) || p_\theta(\mathbf{z}_0))$$

$$KL \left(q \left(\mathbf{z}_0 | \mathbf{x} \right) \| p_\theta \left(\mathbf{z}_0 \right) \right) \leq \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{z}_0 | \mathbf{x})} \mathbb{E}_{q(\mathbf{z}_t | \mathbf{z}_0)} \left[\left\| \nabla_{\mathbf{z}_t} \log q \left(\mathbf{z}_t \right) - \nabla_{\mathbf{z}_t} \log p_\theta \left(\mathbf{z}_t \right) \right\|_2^2 \right] \right]$$

Latent Score-based Generative Model

...instead of a thousands equations



$$\mathbb{E}_{q_\phi(\mathbf{z}_0 | \mathbf{x})} \left[-\log p_\psi(\mathbf{x} | \mathbf{z}_0) \right] + KL(q_\phi(\mathbf{z}_0 | \mathbf{x}) || p_\theta(\mathbf{z}_0))$$

$$KL \left(q \left(\mathbf{z}_0 | \mathbf{x} \right) \| p_\theta \left(\mathbf{z}_0 \right) \right) \leq \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{z}_0 | \mathbf{x})} \mathbb{E}_{q(\mathbf{z}_t | \mathbf{z}_0)} \left[\left\| \nabla_{\mathbf{z}_t} \log q \left(\mathbf{z}_t \right) - \nabla_{\mathbf{z}_t} \log p_\theta \left(\mathbf{z}_t \right) \right\|_2^2 \right] \right]$$

Latent Score-based Generative Model

Objective's derivation - I

$$\begin{aligned}\mathcal{L}(\mathbf{x}, \boldsymbol{\phi}, \boldsymbol{\theta}, \boldsymbol{\psi}) &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_0 \mid \mathbf{x})} \left[-\log p_{\boldsymbol{\psi}}(\mathbf{x} \mid \mathbf{z}_0) \right] + \text{KL} \left(q_{\boldsymbol{\phi}}(\mathbf{z}_0 \mid \mathbf{x}) \parallel p_{\boldsymbol{\theta}}(\mathbf{z}_0) \right) \\ &= \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_0 \mid \mathbf{x})} \left[-\log p_{\boldsymbol{\psi}}(\mathbf{x} \mid \mathbf{z}_0) \right] + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_0 \mid \mathbf{x})} \left[\log q_{\boldsymbol{\phi}}(\mathbf{z}_0 \mid \mathbf{x}) \right] + \mathbb{E}_{q_{\boldsymbol{\phi}}(\mathbf{z}_0 \mid \mathbf{x})} \left[-\log p_{\boldsymbol{\theta}}(\mathbf{z}_0) \right]\end{aligned}$$

Aim - derive last CE term suitable for optimisation:

$$CE \left(q(\mathbf{z}_0 \mid \mathbf{x}) \parallel p(\mathbf{z}_0) \right) = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0 \mid \mathbf{x})} \left[\left\| \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t \mid \mathbf{z}_0) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right\|_2^2 \right] \right] + \frac{D}{2} \log(2\pi e \sigma_0^2)$$

LSGM's objective - I

Fokker-Plank equation recap

- Consider forward SDE: $d\mathbf{z} = \mathbf{f}(\mathbf{z}, t)dt + g(t)d\mathbf{w}$

Fokker-Plank equation:

$$\frac{\partial q(\mathbf{z}_t)}{\partial t} = \nabla_{\mathbf{z}_t} \left(-\mathbf{f}(\mathbf{z}, t)q(\mathbf{z}_t) + \frac{1}{2}g^2(t)q(\mathbf{z}_t)\nabla_{\mathbf{z}_t}\log q(\mathbf{z}_t) \right)$$

- $= \nabla_{\mathbf{z}_t} \left(\mathbf{h}_q(\mathbf{z}_t, t) q(\mathbf{z}_t) \right)$
- Backward SDE: $d\mathbf{z} = [\mathbf{f}(\mathbf{z}, t) - \nabla_{\mathbf{z}}\log q(\mathbf{z})]dt + g(t)d\mathbf{w}$

LSGM's objective - I

Fokker-Plank equation recap

- Consider forward SDE: $d\mathbf{z} = \mathbf{f}(\mathbf{z}, t)dt + g(t)d\mathbf{w}$

Fokker-Plank equation:

$$\frac{\partial q(\mathbf{z}_t)}{\partial t} = \nabla_{\mathbf{z}_t} \left(-\mathbf{f}(\mathbf{z}, t)q(\mathbf{z}_t) + \frac{1}{2}g^2(t)q(\mathbf{z}_t)\nabla_{\mathbf{z}_t}\log q(\mathbf{z}_t) \right)$$

- $= \nabla_{\mathbf{z}_t} \left(\mathbf{h}_q(\mathbf{z}_t, t) q(\mathbf{z}_t) \right)$

- Backward SDE: $d\mathbf{z} = [\mathbf{f}(\mathbf{z}, t) - \nabla_{\mathbf{z}}\log q(\mathbf{z})]dt + g(t)d\mathbf{w}$

- $\mathbf{h}_q(\mathbf{z}_t, t) = -\mathbf{f}(\mathbf{z}, t) + \frac{1}{2}g^2(t)\nabla_{\mathbf{z}_t}\log q(\mathbf{z}_t)$

LSGM's objective - I

CE's partial derivative - I

- Rewrite CE with Newton-Leibnitz Formula:

$$\text{CE}\left(q(\mathbf{z}_0) \| p(\mathbf{z}_0)\right) = \text{CE}\left(q(\mathbf{z}_1) \| p(\mathbf{z}_1)\right) + \int_1^0 \frac{\partial}{\partial t} \text{CE}\left(q(\mathbf{z}_t) \| p(\mathbf{z}_t)\right) dt = H\left(q(\mathbf{z}_1)\right) - \int_0^1 \frac{\partial}{\partial t} \text{CE}\left(q(\mathbf{z}_t) \| p(\mathbf{z}_t)\right) dt$$

-
- Assume, that at $t = 1$, $q(\mathbf{z}_1) = p(\mathbf{z}_1)$

$$\begin{aligned} \frac{\partial}{\partial t} \text{CE}\left(q(\mathbf{z}_t) \| p(\mathbf{z}_t)\right) &= -\frac{\partial}{\partial t} \int q(\mathbf{z}_t) \log p(\mathbf{z}_t) d\mathbf{z} = - \int \left[\frac{\partial q(\mathbf{z}_t)}{\partial t} \log p(\mathbf{z}_t) + \frac{q(\mathbf{z}_t)}{p(\mathbf{z}_t)} \frac{\partial p(\mathbf{z}_t)}{\partial t} \right] d\mathbf{z} \\ &= - \int \left[\nabla_{\mathbf{z}_t} \left(\mathbf{h}_q(\mathbf{z}_t, t) q(\mathbf{z}_t) \right) \log p(\mathbf{z}_t) + \frac{q(\mathbf{z}_t)}{p(\mathbf{z}_t)} \nabla_{\mathbf{z}_t} \left(\mathbf{h}_p(\mathbf{z}_t, t) p(\mathbf{z}_t) \right) \right] d\mathbf{z} \\ \mathbf{h}_q(\mathbf{z}_t, t) &= -\mathbf{f}(\mathbf{z}, t) + \frac{1}{2} g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \end{aligned}$$

LSGM's objective - I

CE's partial derivative - II

$$\begin{aligned}\frac{\partial}{\partial t} \text{CE} \left(q(\mathbf{z}_t) \| p(\mathbf{z}_t) \right) &= - \int \left[\nabla_{\mathbf{z}_t} \left(\mathbf{h}_q(\mathbf{z}_t, t) q(\mathbf{z}_t) \right) \log p(\mathbf{z}_t) + \frac{q(\mathbf{z}_t)}{p(\mathbf{z}_t)} \nabla_{\mathbf{z}_t} \left(\mathbf{h}_p(\mathbf{z}_t, t) p(\mathbf{z}_t) \right) \right] d\mathbf{z} \\ &= \int \left[\mathbf{h}_q(\mathbf{z}_t, t)^\top q(\mathbf{z}_t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \mathbf{h}_p(\mathbf{z}_t, t)^\top p(\mathbf{z}_t) \nabla_{\mathbf{z}_t} \frac{q(\mathbf{z}_t)}{p(\mathbf{z}_t)} \right] d\mathbf{z} \\ &= \int \left[\mathbf{h}_q(\mathbf{z}_t, t)^\top q(\mathbf{z}_t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \mathbf{h}_p(\mathbf{z}_t, t)^\top p(\mathbf{z}_t) \frac{p(\mathbf{z}_t) \nabla_{\mathbf{z}_t} q(\mathbf{z}_t) - q(\mathbf{z}_t) \nabla_{\mathbf{z}_t} p(\mathbf{z}_t)}{p(\mathbf{z}_t)^2} \right] d\mathbf{z} \\ &= \int q(\mathbf{z}_t) \left[\mathbf{h}_q(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \mathbf{h}_p(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) - \mathbf{h}_p(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] d\mathbf{z}\end{aligned}$$

LSGM's objective - I

CE's partial derivative - III

$$\frac{\partial}{\partial t} \text{CE} \left(q(\mathbf{z}_t) \| p(\mathbf{z}_t) \right) = \int q(\mathbf{z}_t) \left[\mathbf{h}_q(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \mathbf{h}_p(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) - \mathbf{h}_p(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] d\mathbf{z}$$

$$\mathbf{h}_q(\mathbf{z}_t, t) = -\mathbf{f}(\mathbf{z}, t) + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \quad \mathbf{h}_p(\mathbf{z}_t, t) = -\mathbf{f}(\mathbf{z}, t) + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$$

$$\begin{aligned} \frac{\partial}{\partial t} \text{CE} \left(q(\mathbf{z}_t) \| p(\mathbf{z}_t) \right) &= \int q(\mathbf{z}_t) \left[-\mathbf{f}(\mathbf{z}, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \right. \\ &\quad -\mathbf{f}(\mathbf{z}, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) + \\ &\quad \left. +\mathbf{f}(\mathbf{z}, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) - \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] d\mathbf{z} \end{aligned}$$

LSGM's objective - I

CE's partial derivative - III

$$\frac{\partial}{\partial t} \text{CE} \left(q(\mathbf{z}_t) \| p(\mathbf{z}_t) \right) = \int q(\mathbf{z}_t) \left[\mathbf{h}_q(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \mathbf{h}_p(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) - \mathbf{h}_p(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] d\mathbf{z}$$

$$\mathbf{h}_q(\mathbf{z}_t, t) = -\mathbf{f}(\mathbf{z}, t) + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \quad \mathbf{h}_p(\mathbf{z}_t, t) = -\mathbf{f}(\mathbf{z}, t) + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)$$

$$\begin{aligned} \frac{\partial}{\partial t} \text{CE} \left(q(\mathbf{z}_t) \| p(\mathbf{z}_t) \right) &= \int q(\mathbf{z}_t) \left[\cancel{-\mathbf{f}(\mathbf{z}, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)} + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) + \right. \\ &\quad \cancel{-\mathbf{f}(\mathbf{z}, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)} + \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) + \\ &\quad \left. + \cancel{\mathbf{f}(\mathbf{z}, t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)} - \frac{1}{2}g^2(t) \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] d\mathbf{z} \end{aligned}$$

LSGM's objective - II

Score of $q(\mathbf{z}_t)$ - II

$$\begin{aligned}
\text{CE}\left(q(\mathbf{z}_0) \| p(\mathbf{z}_0)\right) &= \text{CE}\left(q(\mathbf{z}_1) \| p(\mathbf{z}_1)\right) + \int_1^0 \frac{\partial}{\partial t} \text{CE}\left(q(\mathbf{z}_t) \| p(\mathbf{z}_t)\right) dt = \text{H}\left(q(\mathbf{z}_1)\right) + \int_0^1 \frac{\partial}{\partial t} \text{CE}\left(q(\mathbf{z}_t) \| p(\mathbf{z}_t)\right) dt \\
&= \text{H}\left(q(\mathbf{z}_1)\right) + \int_0^1 \mathbb{E}_{q(\mathbf{z}_t)} \left[\frac{1}{2} g^2(t) \left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right\|_2^2 + 2\mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) - g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] dt \\
&= \text{H}\left(q(\mathbf{z}_1)\right) + \int_0^1 \mathbb{E}_{q(\mathbf{z}_t)} \frac{1}{2} \left(g^2(t) \left[\left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right\|_2^2 - 2 \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] + 2\mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \right) dt
\end{aligned}$$

$$\begin{aligned}
\nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) &= \frac{\nabla_{\mathbf{z}_t} q(\mathbf{z}_t)}{q(\mathbf{z}_t)} = \frac{\nabla_{\mathbf{z}_t} \int q(\mathbf{z}_0) q(\mathbf{z}_t | \mathbf{z}_0) d\mathbf{z}_0}{q(\mathbf{z}_t)} = \frac{\int q(\mathbf{z}_0) \nabla_{\mathbf{z}_t} q(\mathbf{z}_t | \mathbf{z}_0) d\mathbf{z}_0}{q(\mathbf{z}_t)} = \frac{\int q(\mathbf{z}_0) q(\mathbf{z}_t | \mathbf{z}_0) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) d\mathbf{z}_0}{q(\mathbf{z}_t)} \\
&= \mathbb{E}_{q(\mathbf{z}_0 | \mathbf{z}_t)} \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0)
\end{aligned}$$

$$\mathbb{E}_{q(\mathbf{z}_t)} \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) = \mathbb{E}_{q(\mathbf{z}_t)} \mathbb{E}_{q(\mathbf{z}_0 | \mathbf{z}_t)} \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) = \mathbb{E}_{q(\mathbf{z}_0)} \mathbb{E}_{q(\mathbf{z}_t | \mathbf{z}_0)} \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0)$$

LSGM's objective - III

in search of elements that do not depend $p(\mathbf{z}_t)$

$$\begin{aligned} \text{CE}\left(q(\mathbf{z}_0) \| p(\mathbf{z}_0)\right) &= \\ &= H\left(q(\mathbf{z}_1)\right) + \int_0^1 \mathbb{E}_{q(\mathbf{z}_t)} \frac{1}{2} \left(g^2(t) \left[\left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right\|_2^2 - 2 \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] + 2 \mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t) \right) dt \\ &= H\left(q(\mathbf{z}_1)\right) + \int_0^1 \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0)} \frac{1}{2} \left(g^2(t) \left[\left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right\|_2^2 - 2 \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0)^\top \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right] + 2 \mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right) dt \\ &\quad \int_0^1 \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0)} \frac{1}{2} \left(g^2(t) \left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right\|_2^2 + 2 \mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) - g^2(t) \left\| \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right\|_2^2 \right) dt \\ &= \int_0^1 \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0)} \frac{1}{2} \left(g^2(t) \left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right\|_2^2 + 2 \mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) - g^2(t) \left\| \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right\|_2^2 \right) dt \end{aligned}$$

LSGM's objective - IV

Constant even for encoder

$$\int_0^1 \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0)} \frac{1}{2} \left(2\mathbf{f}(\mathbf{z}_t, t)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) - g^2(t) \left\| \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right\|_2^2 \right) dt$$

$$d\mathbf{z} = \mathbf{f}(\mathbf{z}, t)dt + g(t)d\mathbf{w} = f(t)\mathbf{z}_t dt + g(t)d\mathbf{w} \quad q(\mathbf{z}_t | \mathbf{z}_0) = \mathcal{N}(\mathbf{z}_t | \alpha_t \mathbf{z}_0, \sigma_t^2 I) \quad \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) = -\frac{\epsilon}{\sigma_t}$$

$$\int_0^1 \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0)} \frac{1}{2} \left(2\mathbf{f}(\mathbf{z}_t, t) - g^2(t) \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right)^\top \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) dt$$

$$\begin{aligned} \int_0^1 \mathbb{E}_{q(\mathbf{z}_0), \epsilon} \frac{1}{2} \left(2f(t)(\alpha_t \mathbf{z}_0 + \sigma_t \epsilon) + g^2(t) \frac{\epsilon}{\sigma_t} \right)^\top \left(-\frac{\epsilon}{\sigma_t} \right) dt &= \int_0^1 -\frac{f(t)}{\sigma_t} \underbrace{\mathbb{E}_{q(\mathbf{z}_0), \epsilon} \left[\alpha_t (\mathbf{z}_0)^T \epsilon \right]}_{=0} - \underbrace{\frac{2f(t)\sigma_t^2 + g(t)^2}{2\sigma_t^2} \mathbb{E}_\epsilon [\epsilon^T \epsilon]}_{=D} dt \\ &= -\frac{D}{2} \int_0^1 \frac{2f(t)\sigma_t^2 + g(t)^2}{\sigma_t^2} dt = \frac{D}{2} (\log \sigma_0^2 - \log \sigma_1^2) \end{aligned}$$

LSGM's objective - V CE - final

$$\text{CE}\left(q(\mathbf{z}_0) \| p(\mathbf{z}_0)\right) = \text{H}\left(q(\mathbf{z}_1)\right) + \int_0^1 \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0)} \frac{1}{2} \left(g^2(t) \left\| \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) - \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) \right\|_2^2 \right) dt + \frac{D}{2} (\log \sigma_0^2 - \log \sigma_1^2)$$

$$\text{H}\left(q(\mathbf{z}_1)\right) = \int q(\mathbf{z}_1) \log q(\mathbf{z}_1) d\mathbf{z}_1 = \frac{D}{2} \log(2\pi e \sigma_1^2)$$

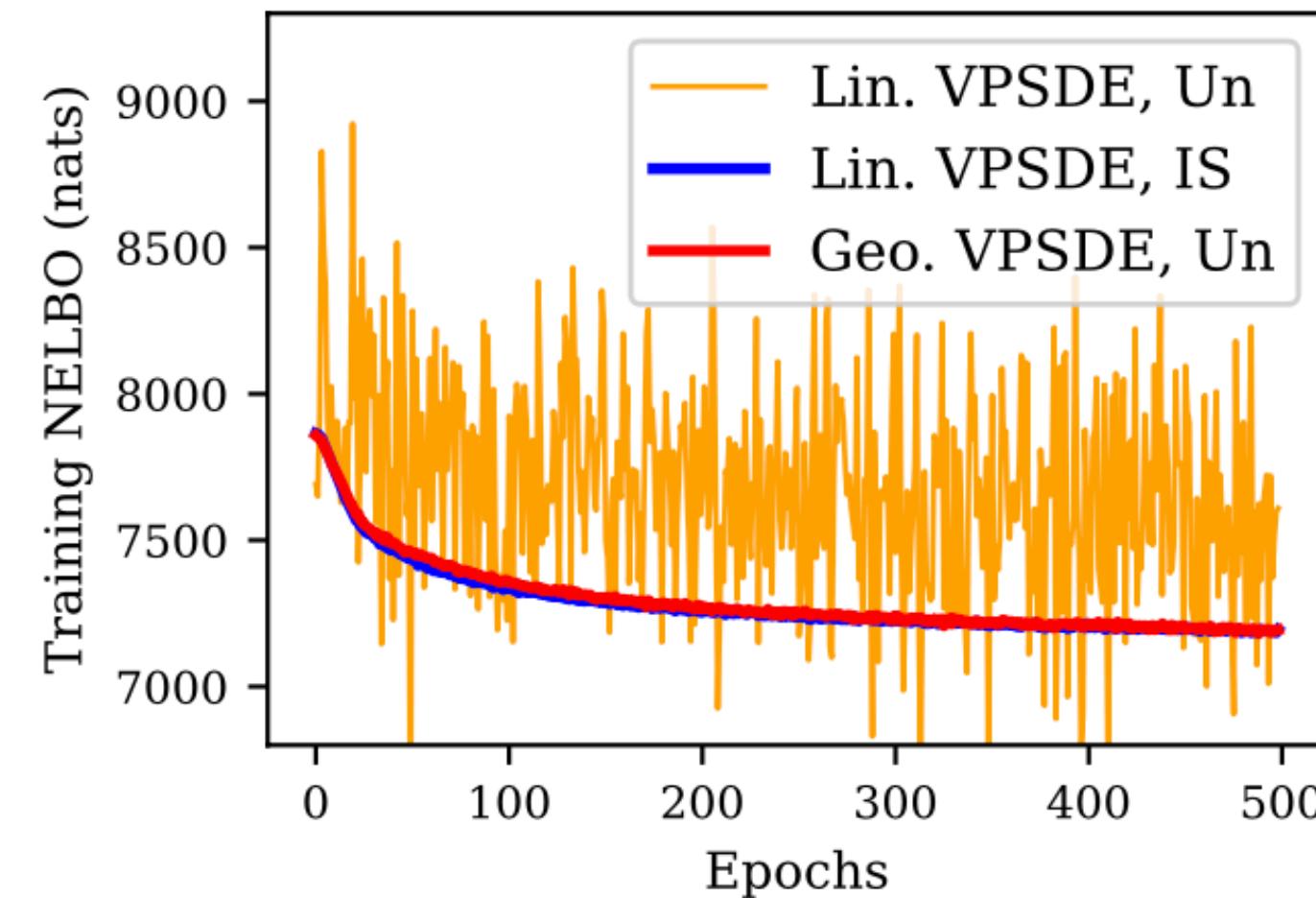
$$q(\mathbf{z}_1) = \mathcal{N}(\mathbf{z}_1 | 0, \sigma_1^2 I)$$

$$\text{CE}\left(q(\mathbf{z}_0 | \mathbf{x}) \| p(\mathbf{z}_0)\right) = \mathbb{E}_{t \sim \mathcal{U}[0,1]} \left[\frac{g(t)^2}{2} \mathbb{E}_{q(\mathbf{z}_t, \mathbf{z}_0 | \mathbf{x})} \left[\left\| \nabla_{\mathbf{z}_t} \log q(\mathbf{z}_t | \mathbf{z}_0) - \nabla_{\mathbf{z}_t} \log p(\mathbf{z}_t) \right\|_2^2 \right] \right] + \frac{D}{2} \log(2\pi e \sigma_0^2)$$

Things, I didn't say

...or a thousand more formulas

- New parametrization for score-estimation: $p(z_t) \propto \mathcal{N}(z_t; 0, 1)^{1-\alpha} p'_\theta(z_t)^\alpha$
 - So new score: $\nabla_{z_t} \log p(z_t) = -(1 - \alpha)z_t + \alpha \nabla_{z_t} \log p'_\theta(z_t)$
- Several ways for reducing loss's dispersion:
 - Importance sampling for weighting
 - New Geometric-SDEs



LSGM experiments

Table 2: Generative performance on CIFAR-10.

	Method	NLL↓	FID↓
Ours	LSGM (FID)	≤ 3.43	2.10
	LSGM (NLL)	≤ 2.87	6.89
	LSGM (balanced)	≤ 2.95	2.17
	VAE Backbone	2.96	43.18
VAEs	VDVAE [21]	2.87	-
	NVAE [20]	2.91	23.49
	VAEBM [76]	-	12.19
	NCP-VAE [56]	-	24.08
	BIVA [48]	3.08	-
	DC-VAE [77]	-	17.90
Score	NCSN [3]	-	25.32
	Rec. Likelihood [40]	3.18	9.36
	DSM-ALS [39]	3.65	-
	DDPM [1]	3.75	3.17
	Improved DDPM [26]	2.94	11.47
	SDE (DDPM++) [2]	2.99	2.92
	SDE (NCSN++) [2]	-	2.20

Table 3: Generative results on CelebA-HQ-256.

	Method	NLL↓	FID↓
Ours	LSGM	≤ 0.70	7.22
	VAE Backbone	0.70	30.87
VAEs	NVAE [20]	0.70	29.76
	VAEBM [76]	-	20.38
	NCP-VAE [56]	-	24.79
	DC-VAE [77]	-	15.80
Score	SDE [2]	-	7.23
Flows	GLOW [85]	1.03	68.93
Aut. Reg.	SPN [86]	0.61	-
GANs	Adv. LAE [87]	-	19.21
	VQ-GAN [64]	-	10.70
	PGGAN [88]	-	8.03



(b) CelebA-HQ-256

Towards stability

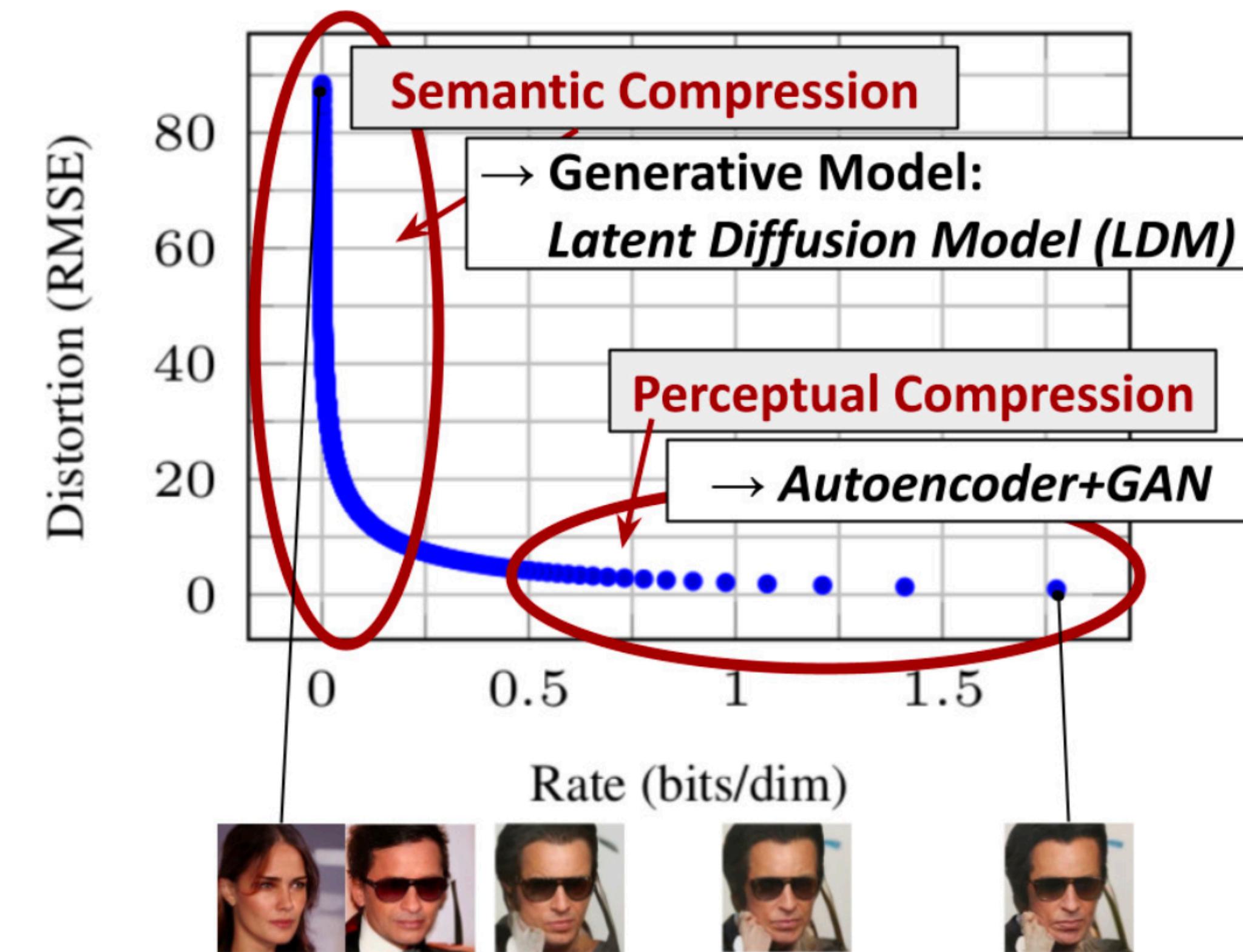
Key idea

Divide generative model's training into 2 stages:

- Train to restoring high-frequency details
- Train to generate various and hard semantics

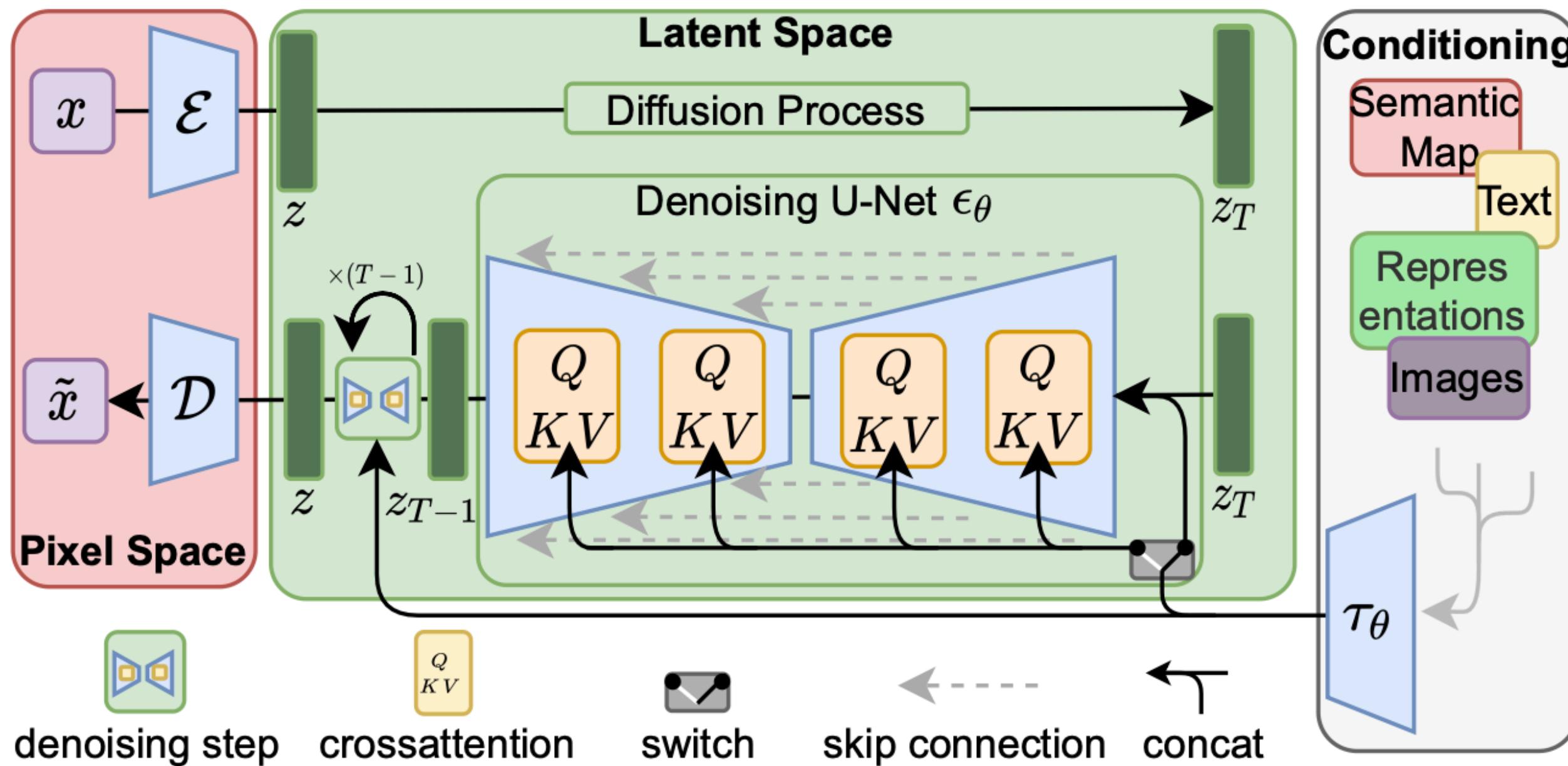
So let's use fixed latent space

- Fixed good encoder-decoder for reconstruction high-frequency details
- New more suitable “pixel” space for learning semantics



Whole scheme

For training authors used 256x256 images



$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d}} \right) \cdot V,$$

$$Q = W_Q^{(i)} \cdot \varphi_i(z_t),$$

$$K = W_K^{(i)} \cdot \tau_\theta(y),$$

$$V = W_V^{(i)} \cdot \tau_\theta(y)$$

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t)\|_2^2 \right]$$

$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[\|\epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y))\|_2^2 \right]$$

High-resolution image synthesis

General approach for high-resolution generation -
Sliding window with 256x256 crops

- High-resolution layouts:
 - $\tau_\theta(y)$ - Identity
 - y - semantic mask
- Super-resolution
 - $\tau_\theta(y)$ - Identity
 - y - downsampled image

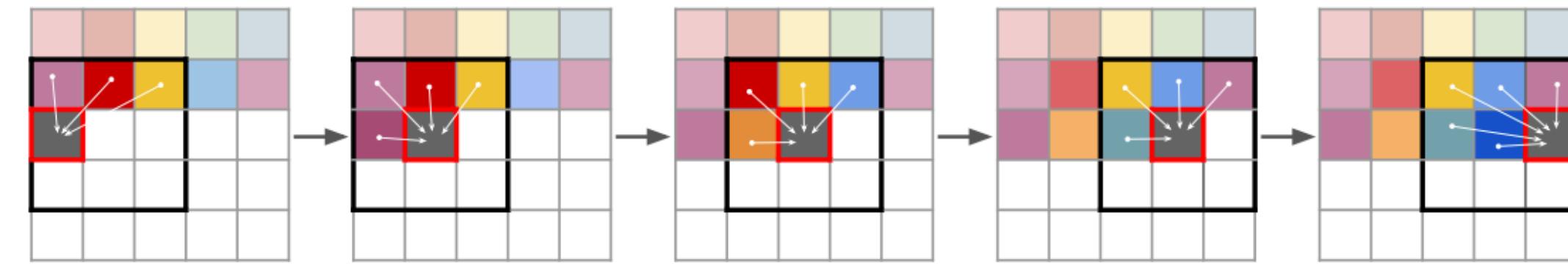


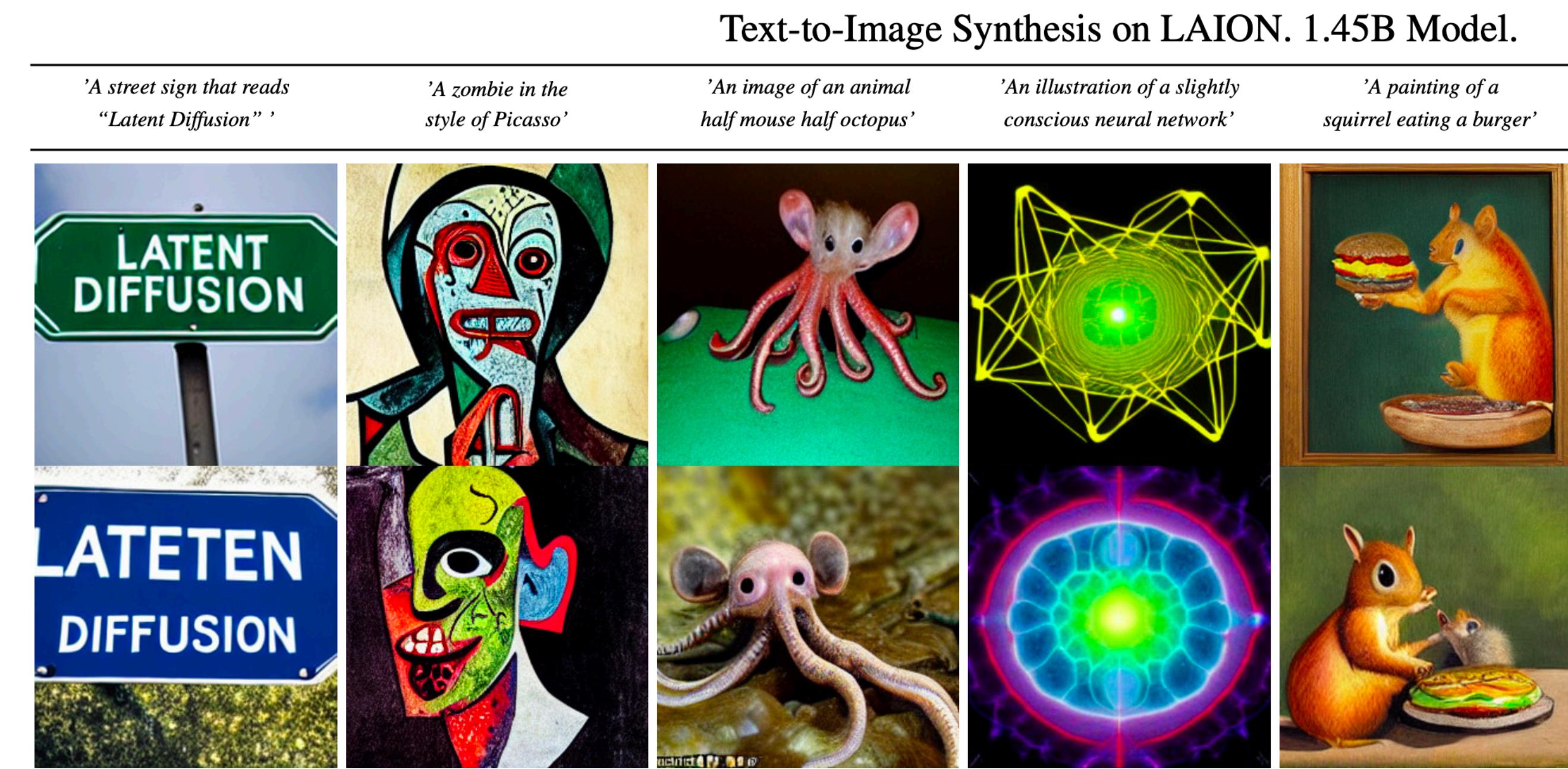
Figure 3. Sliding attention window.

Stable Diffusion's experiments

Text-to-image

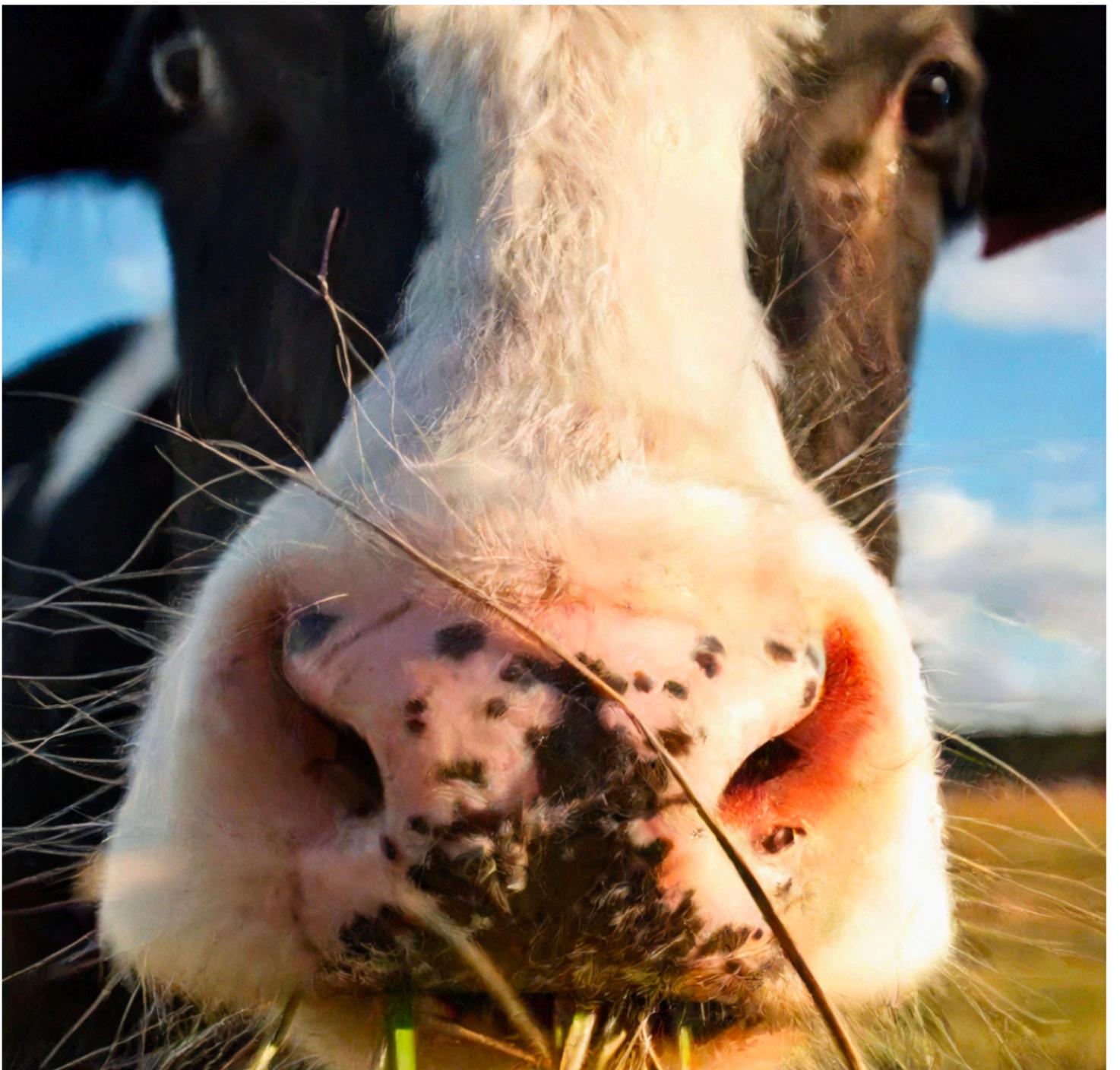
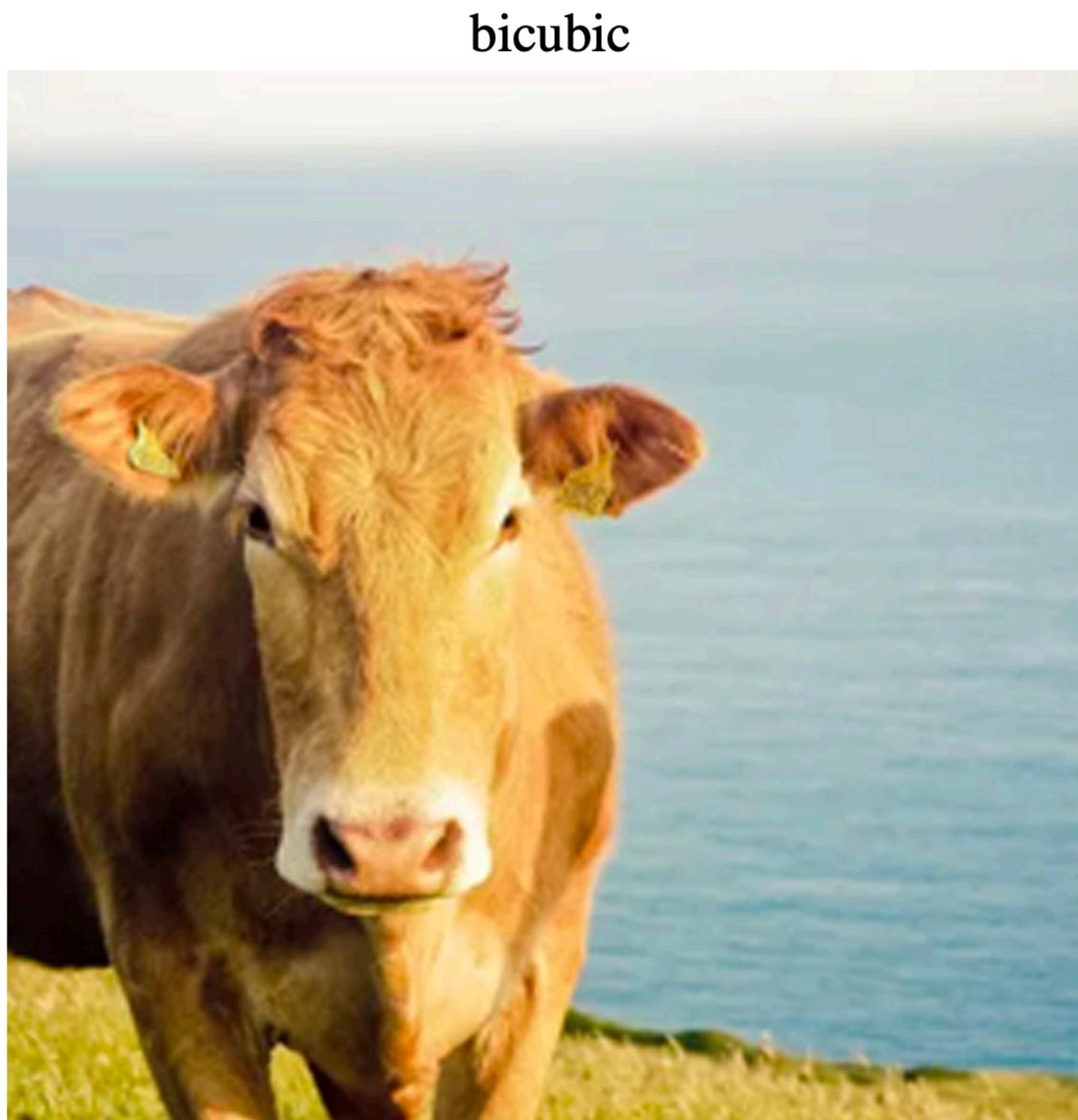
Text-Conditional Image Synthesis			
Method	FID ↓	IS↑	Nparams
CogView [†] [17]	27.10	18.20	4B
LAFITE [†] [109]	26.94	<u>26.02</u>	75M
GLIDE* [59]	<u>12.24</u>	-	6B
Make-A-Scene* [26]	11.84	-	4B
<i>LDM-KL-8</i>	23.31	20.03 ± 0.33	1.45B
<i>LDM-KL-8-G*</i>	12.63	30.29 ± 0.42	1.45B
<i>250 DDIM steps</i>			

Table 2. Evaluation of text-conditional image synthesis on the 256×256 -sized MS-COCO [51] dataset: with 250 DDIM [84] steps our model is on par with the most recent diffusion [59] and autoregressive [26] methods despite using significantly less parameters. [†]/*:Numbers from [109]/[26]



Stable Diffusion's experiments

Super-resolution (1024x1024)



Conclusion

- Saw diffusion in latent space
 - With simultaneously training - derived non-trivial training objective for this approach
 - End-to-end
 - Iterative optimization
 - Saw what models in fixed latent space can achieve using crossattention with prompts