

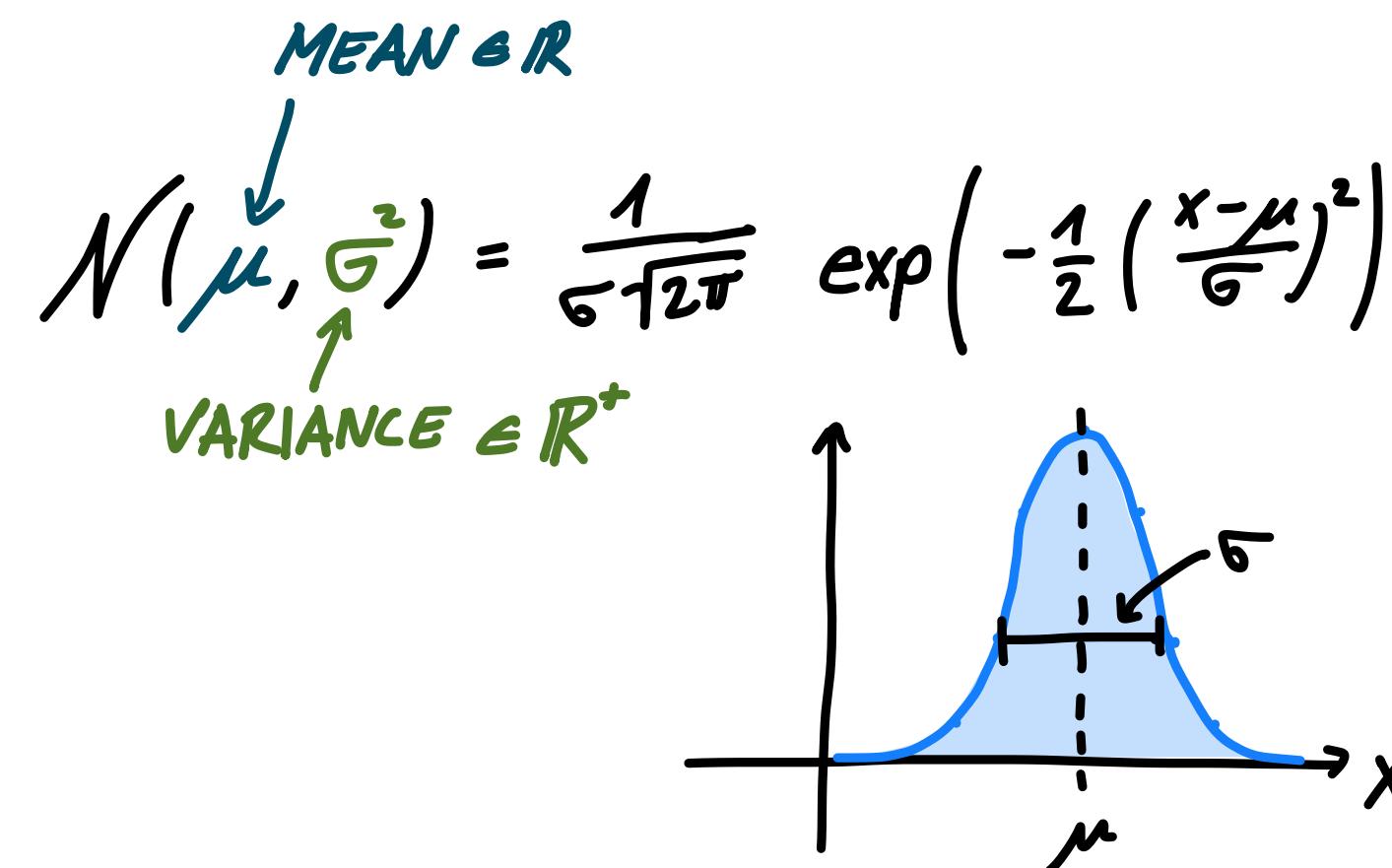


Notes on Diffusion Models

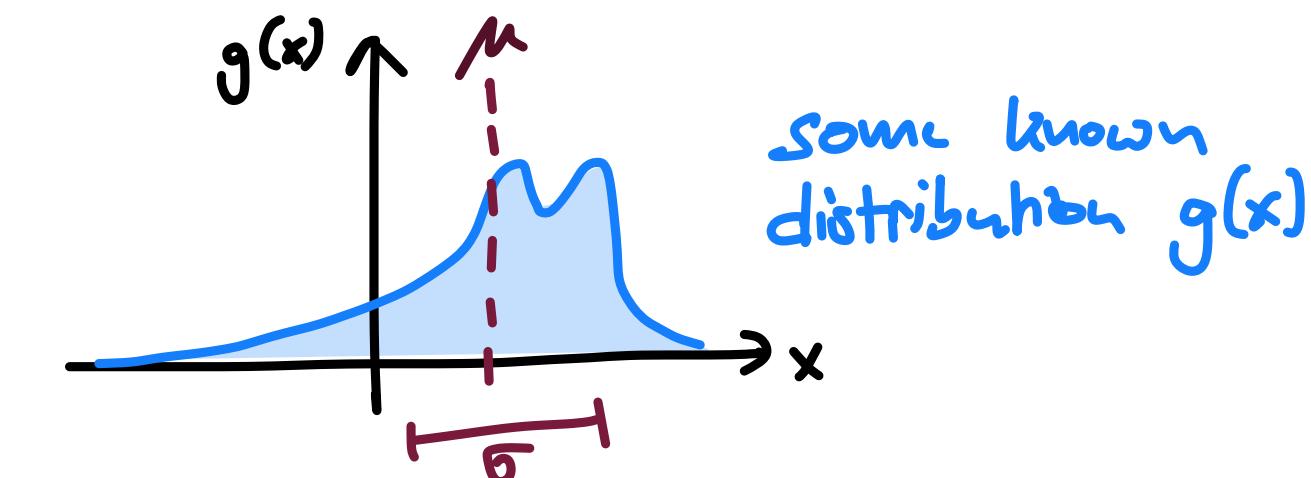
1. Sampling from a normal distribution
2. How to learn a distribution $p(x)$
3. Distribution of natural images
4. Vanilla Diffusion Model
 1. Forward Diffusion Process
 2. Denoising Process
 3. Loss Function
 4. Training and Inference
 5. KPIs
5. Improvements
 1. Song et al (DDIM)
 2. Nichol et al
 3. Dhariwal et al
 4. Nichol et al (GLIDE)
 5. Ramesh et al (DALL-E 2)
 6. Rombach et al (Stable Diffusion)
6. UNet Architecture in Depth
7. ControlNet

1 - Sampling from a normal distribution

Probability density function of a normal distribution:



In other words:



When sampling n times from g to get $\{x_i\}_{i=1..n}$ we know that the average $m = \frac{1}{n} \sum_{i=1}^n x_i$ follows a normal distribution $N\left(\mu, \frac{\sigma^2}{n}\right)$.

We are looking for a function f that randomly returns a number, such that the corresponding histogram after many executions of that function follows the pdf of the normal distribution.

Trick: Central Limit Theorem

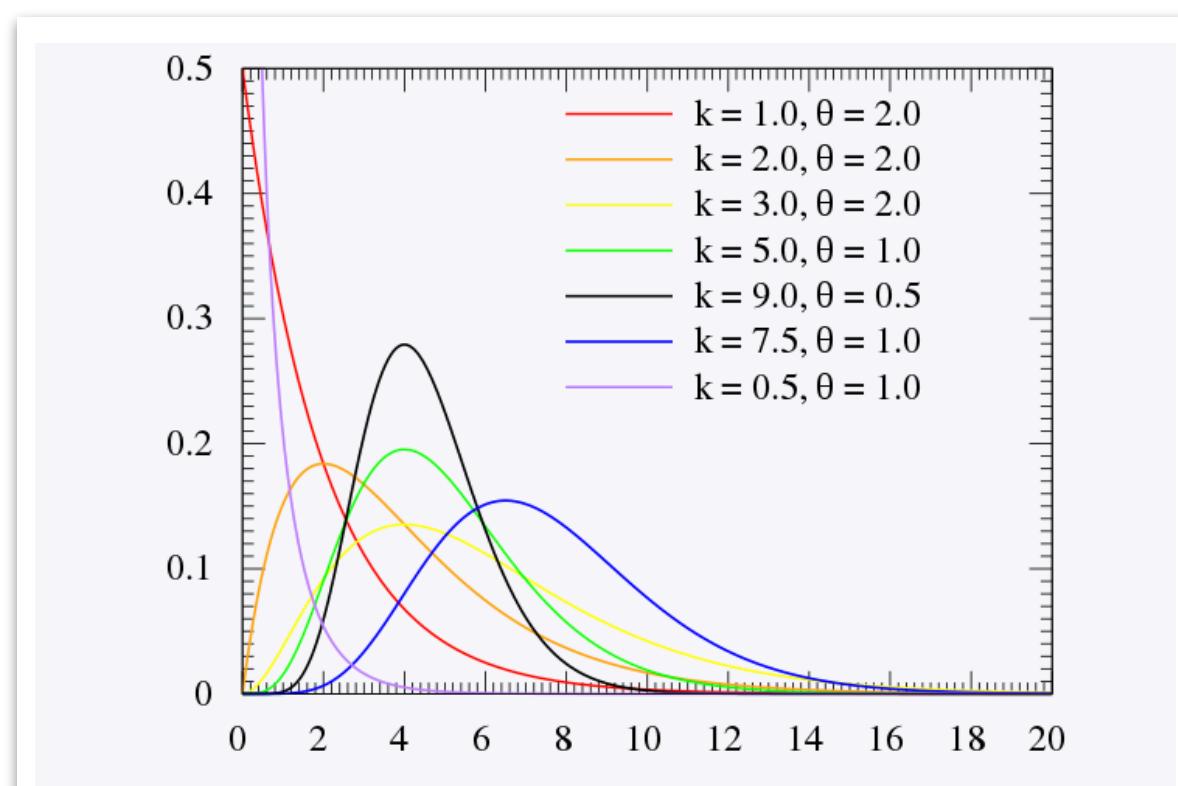
If $X_1, X_2, \dots, X_n, \dots$ are random samples drawn from a population with overall mean μ and finite variance σ^2 , and if \bar{X}_n is the sample mean of the first n samples, then the limiting form of the distribution, $Z = \lim_{n \rightarrow \infty} \left(\frac{\bar{X}_n - \mu}{\sigma_{\bar{X}}} \right)$, with $\sigma_{\bar{X}} = \sigma / \sqrt{n}$, is a standard normal distribution.^[2]

If you want to draw from $N(0,1)$, instead draw n times from g and calculate the following number:

$$r = \frac{\left(\frac{1}{n} \sum_{i=1}^n x_i \right) - \mu}{\sigma / \sqrt{n}}$$

2 - How to learn a distribution $p(x)$

Toy model: Gamma Distribution

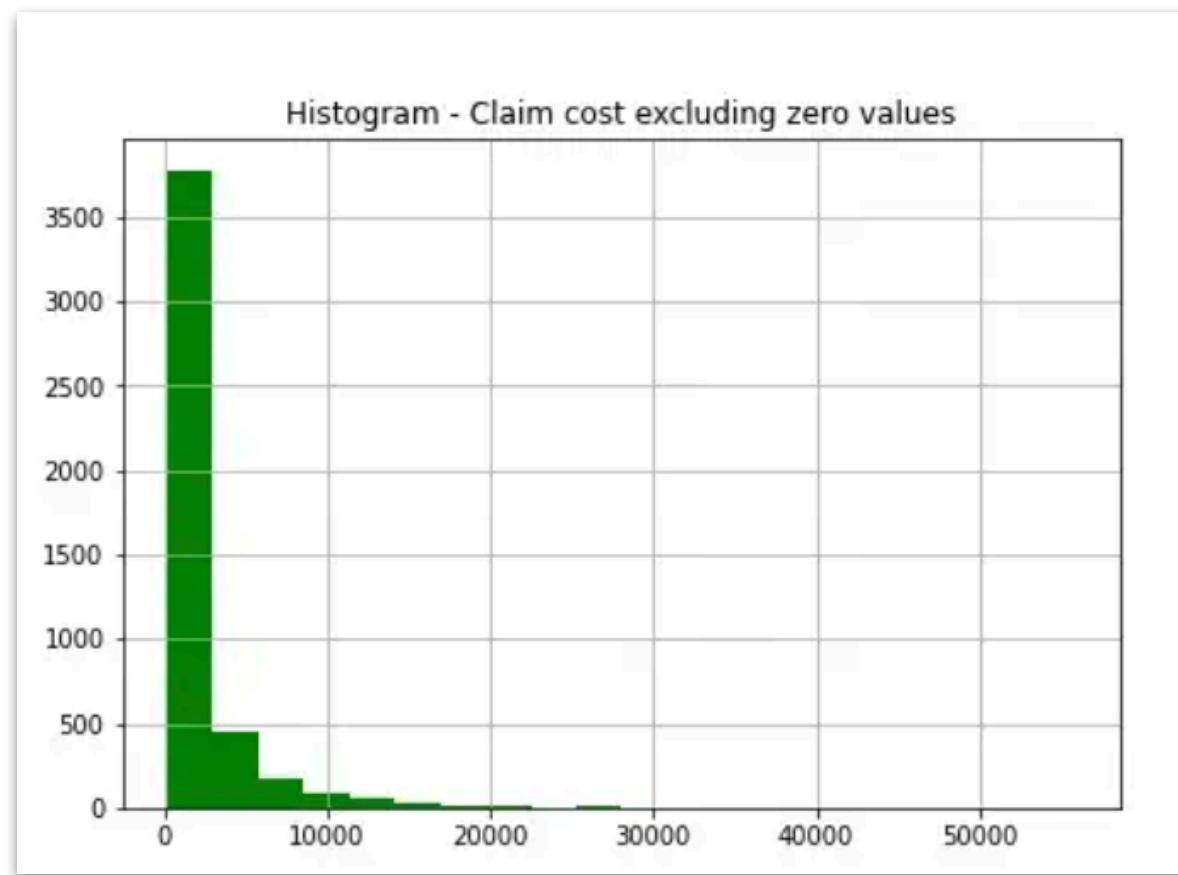


The corresponding probability density function in the shape-rate parameterization is

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1} e^{-\beta x} \beta^\alpha}{\Gamma(\alpha)} \quad \text{for } x > 0 \quad \alpha, \beta > 0,$$

where $\Gamma(\alpha)$ is the [gamma function](#). For all positive integers, $\Gamma(\alpha) = (\alpha - 1)!$.

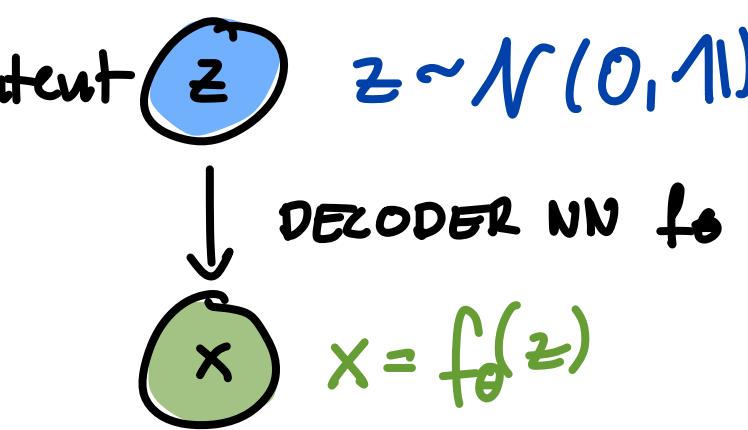
https://en.wikipedia.org/wiki/Gamma_distribution



<https://medium.com/swlh/modeling-insurance-claim-severity-b449ac426c23>

Let's learn this distribution, i.e. find a function that randomly generates values such that the corresponding histogram after many generated values follows the probability density function of the gamma distribution.

- **GRAPHICAL MODEL**



- **OPTIMIZATION VIA DATASET $\{x^{(i)}\}_{i=1..n}$**

$$\theta^* = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n p_\theta(x^{(i)})$$

$$\Leftrightarrow \theta^* = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^n \log p_\theta(x^{(i)})$$

- **MODEL DEFINITION**

$$p(z) = \mathcal{N}(0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{z^2}{2}\right)$$

$$p_\theta(x|z) = \mathcal{N}(f_\theta(z), c \mathbb{I}) \quad \text{IR}^+, \text{ hyperparameter}$$

$$= \frac{1}{\sqrt{2\pi c}} \exp\left(-\frac{(x - f_\theta(z))^2}{2c}\right)$$

Derivation of the loss function:

$$\sum_{i=1}^n \log p_\theta(x^{(i)}) = \sum_{i=1}^n \log \int p_\theta(x^{(i)}|z) p(z) dz = \dots$$

- **INTRACTABLE** \Rightarrow Integration over whole latent space
 - Trick: Knowing the mapping $x \rightarrow z$ would speed up
- \Rightarrow Introduce "encoder" to speed up training
- $$x \rightarrow \mathcal{N}(g_\phi(x), \text{diag}(h_\phi(x))) = q_\phi(z|x)$$

$$\dots = \sum_{i=1}^n \log \int p_\theta(x^{(i)}|z) p(z) \frac{q_\phi(z|x)}{q_\phi(z|x)} dz$$

$$= \sum_{i=1}^n \log \int p_\theta(x^{(i)}|z) \frac{p(z)}{q_\phi(z|x)} q_\phi(z|x) dz$$

$$= \sum_{i=1}^n \log \left\langle p_\theta(x^{(i)}|z) \frac{p(z)}{q_\phi(z|x)} \right\rangle_{z \sim q_\phi(z|x)}$$

$$\geq \sum_{i=1}^n \left\langle \log \left(p_\theta(x^{(i)}|z) \frac{p(z)}{q_\phi(z|x)} \right) \right\rangle_{z \sim q_\phi(z|x)}$$

$$= \sum_{i=1}^n \left\langle \log p_\theta(x^{(i)}|z) \right\rangle_{z \sim q_\phi(z|x)} - \left\langle \log \frac{q_\phi(z|x)}{p(z)} \right\rangle_{z \sim q_\phi(z|x)}$$

$$= \sum_{i=1}^n \left\langle \log p_\theta(x^{(i)}|z) \right\rangle_{z \sim q_\phi(z|x)} - D_{KL}(q_\phi(z|x) \parallel p(z))$$

1st term

2nd term

Jensen
inequality

2 - How to learn a distribution $p(x)$

Explicit calculation of loss terms using Gaussian distributions:

$$\begin{aligned}
 \text{1st term} &= \left\langle \log \left(\frac{1}{\sqrt{2\pi c}} \exp \left(-\frac{(x - f_\phi(z))^2}{2c} \right) \right) \right\rangle_{z \sim q_\phi(z|x)} \\
 &= \left\langle \underbrace{\log \frac{1}{\sqrt{2\pi c}}}_{:= \text{const.}} - \frac{(x - f_\phi(z))^2}{2c} \right\rangle_{z \sim q_\phi(z|x)} \\
 \text{2nd term} &= D_{KL} \left(\mathcal{N} \left(\mu_j^{(x)}, \text{diag}(\Sigma_j^{(x)})^T \right) \parallel \mathcal{N}(0, I) \right) \\
 &\quad \text{Wikipedia KL-Divergence} \\
 &= \frac{1}{2} \sum_{j=1}^d \left(\Sigma_j^{(x)} + \mu_j^{(x)2} - 1 - \log \Sigma_j^{(x)} \right)
 \end{aligned}$$

Final loss:

$$\max \stackrel{!}{=} \sum_{i=1}^n \log p_\phi(x^{(i)}) \geq \sum_{i=1}^n \left\langle \underbrace{\log \frac{1}{\sqrt{2\pi c}}}_{:= \text{const.}} - \frac{(x - f_\phi(z))^2}{2c} \right\rangle_{z \sim q_\phi(z|x)} - D_{KL}(q_\phi(z|x) \parallel p(z))$$

$$\Leftrightarrow \min \stackrel{!}{=} \sum_{i=1}^n \left(\text{REC} \left\langle [x - f_\phi(z)]^2 \right\rangle_{z \sim q_\phi(z|x)} + D_{KL}(q_\phi(z|x) \parallel p(z)) \right) = \lambda$$

reconstruction loss Latent loss

Assuming normal distributions results in the L2 norm, if we would model $p(z|x)$ differently this term would be more complicated. See this [blog post](#).

```

# Model
# #####
class VAE(torch.nn.Module):
    def __init__(self, data_dim, latent_dim, n1=512, n2=256):
        super(VAE, self).__init__()

        # encoder part
        self.fc1 = torch.nn.Linear(data_dim, n1)
        self.fc2 = torch.nn.Linear(n1, n2)
        self.fc31 = torch.nn.Linear(n2, latent_dim)
        self.fc32 = torch.nn.Linear(n2, latent_dim)

        # decoder part
        self.fc4 = torch.nn.Linear(latent_dim, n2)
        self.fc5 = torch.nn.Linear(n2, n1)
        self.fc6 = torch.nn.Linear(n1, data_dim)

    def encoder(self, x):
        h = torch.relu(self.fc1(x))
        h = torch.relu(self.fc2(h))
        return self.fc31(h), self.fc32(h) # mu, log_var

    def sampling(self, mu, log_var):
        std = torch.exp(0.5*log_var)
        eps = torch.randn_like(std)
        return eps.mul(std).add_(mu) # return z sample

    def decoder(self, z):
        h = torch.relu(self.fc4(z))
        h = torch.relu(self.fc5(h))
        # return self.fc6(h)
        # return torch.sigmoid(self.fc6(h))
        return torch.relu(self.fc6(h))

    def forward(self, x):
        mu, log_var = self.encoder(x)
        z = self.sampling(mu, log_var)
        return self.decoder(z), mu, log_var

# Loss Function
# #####
def loss_function(x_tilde, x, mu, log_var, C=1):
    n = x.size(0)
    REC = (x_tilde - x).pow(2).sum() / n
    KLD = 0.5 * torch.sum(mu.pow(2) + log_var.exp() - 1 - log_var) / n
    return C * REC + KLD, REC, KLD

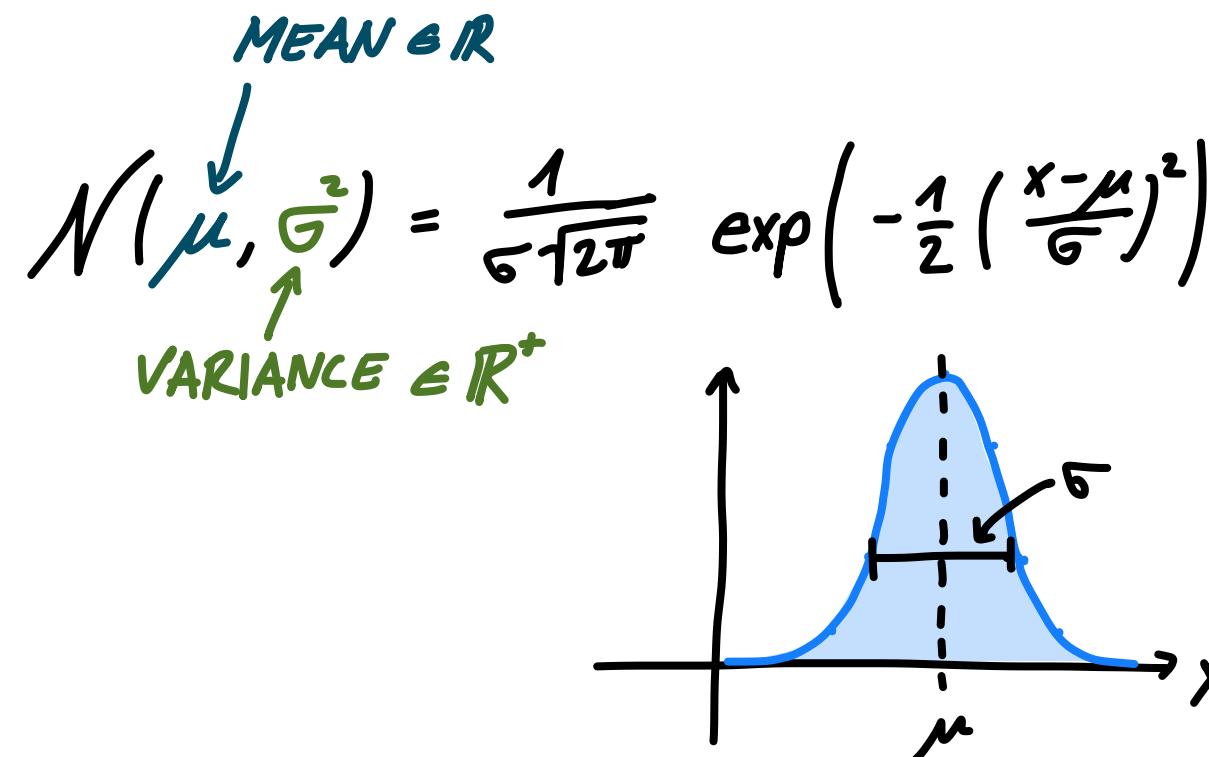
```

...PyTorch implementation, turns out we developed a "Variational Autoencoder".

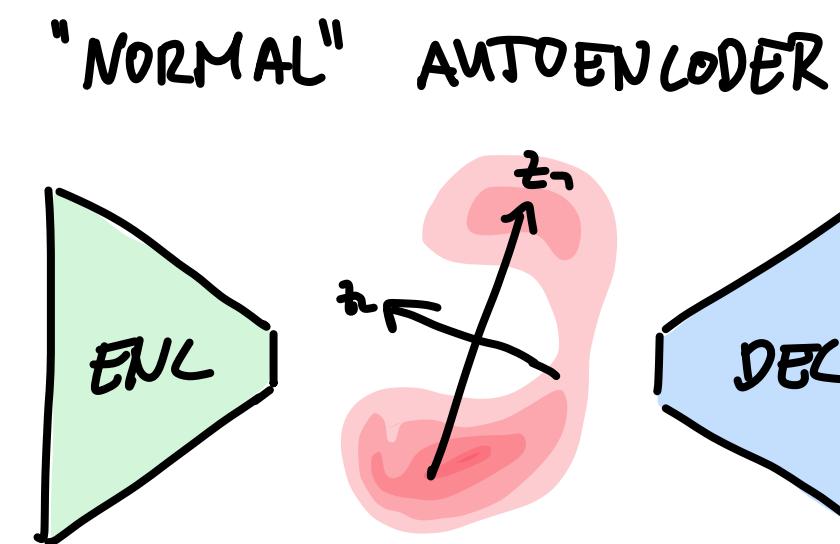
2 - How to learn a distribution $p(x)$

Some side notes:

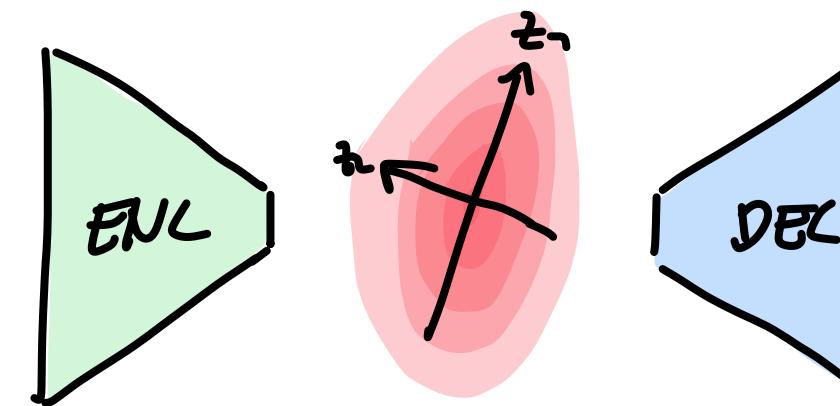
Gaussian notation: σ , σ^2 , std & variance?



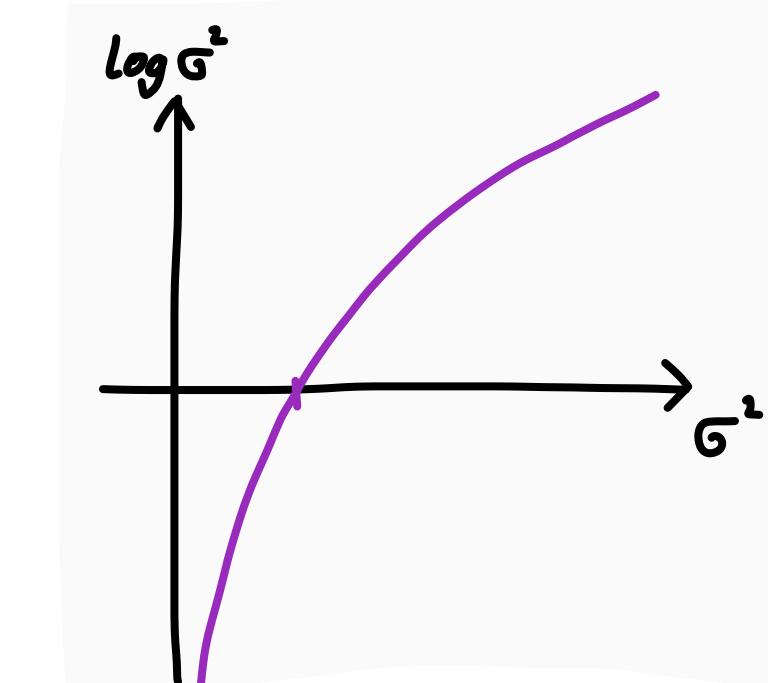
Autoencoder vs. Variational Autoencoder



VARIATIONAL AUTOENCODER
→ SAMPLING POSSIBLE

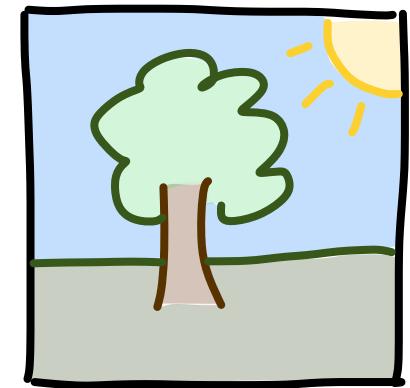


Why logvar?



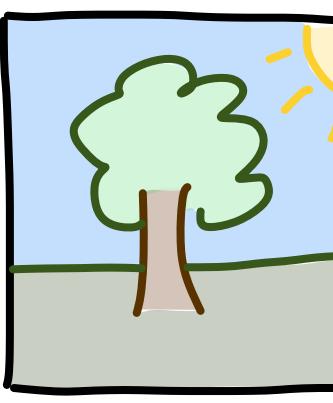
$$\log \sigma^2 \in \mathbb{R}$$
$$\sigma^2 = e^{\log \sigma^2} \in \mathbb{R}^+$$

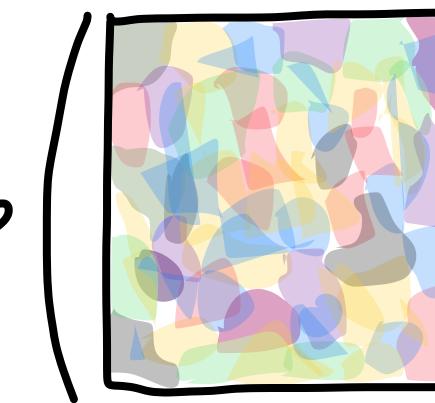
3 - Distribution of natural images


$$\hat{\mathbf{p}} \approx \begin{pmatrix} 0.53 \\ -0.81 \\ -0.62 \\ \vdots \\ 0.78 \\ 0.81 \end{pmatrix} \in \mathbb{R}^{d \times d}$$

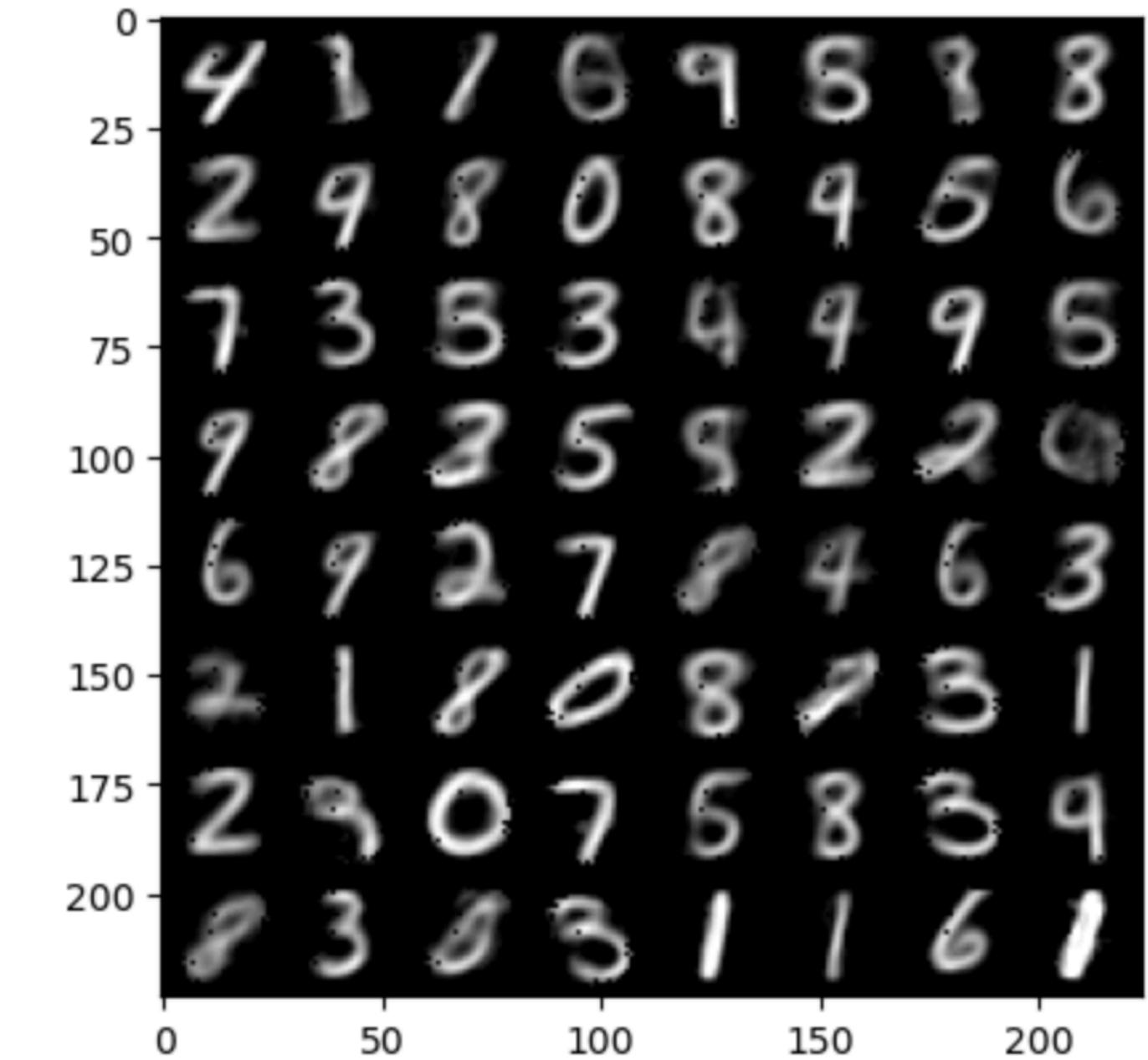
$\hat{\mathbf{p}}$: Distribution of natural images

\Rightarrow very complicated

$$P\left(\left|\hat{\mathbf{p}} - \mathbf{p}\right| \leq \epsilon\right) = \text{high}$$


$$P\left(\left|\hat{\mathbf{p}} - \mathbf{p}\right| \leq \epsilon\right) = \text{low}$$


```
n = 8
with torch.no_grad():
    z = torch.randn(n*n, LATENT_DIM).to(DEVICE)
    sample = vae.decoder(z).to(DEVICE)
    img = sample.view(n,n,28,28).cpu().numpy()
    img = np.vstack(np.dstack(img))
plt.imshow(img, cmap='gray', vmin=0, vmax=1)
plt.show()
```

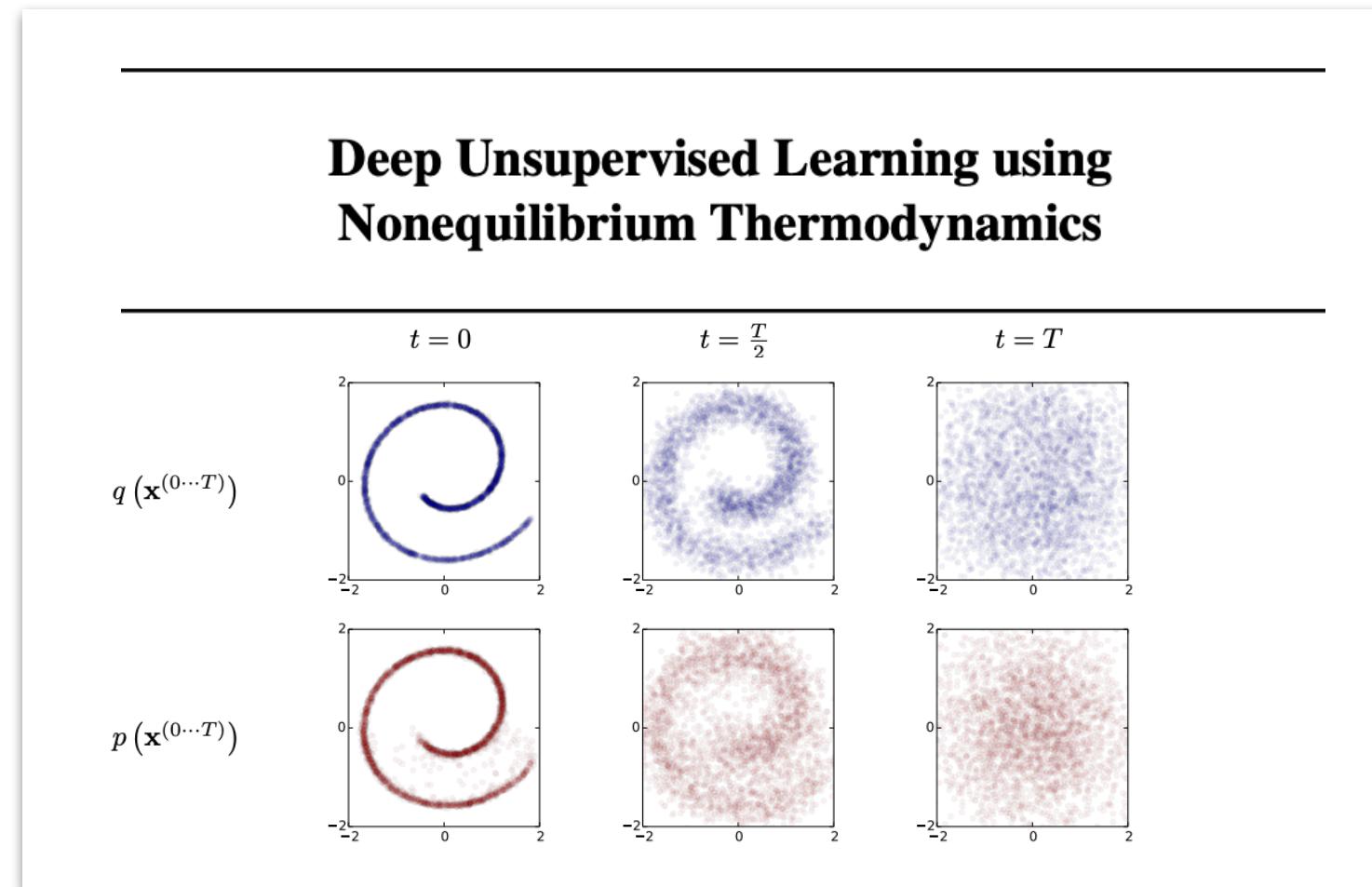


VAE from before trained on the MNIST dataset.

4 - Vanilla Diffusion Model

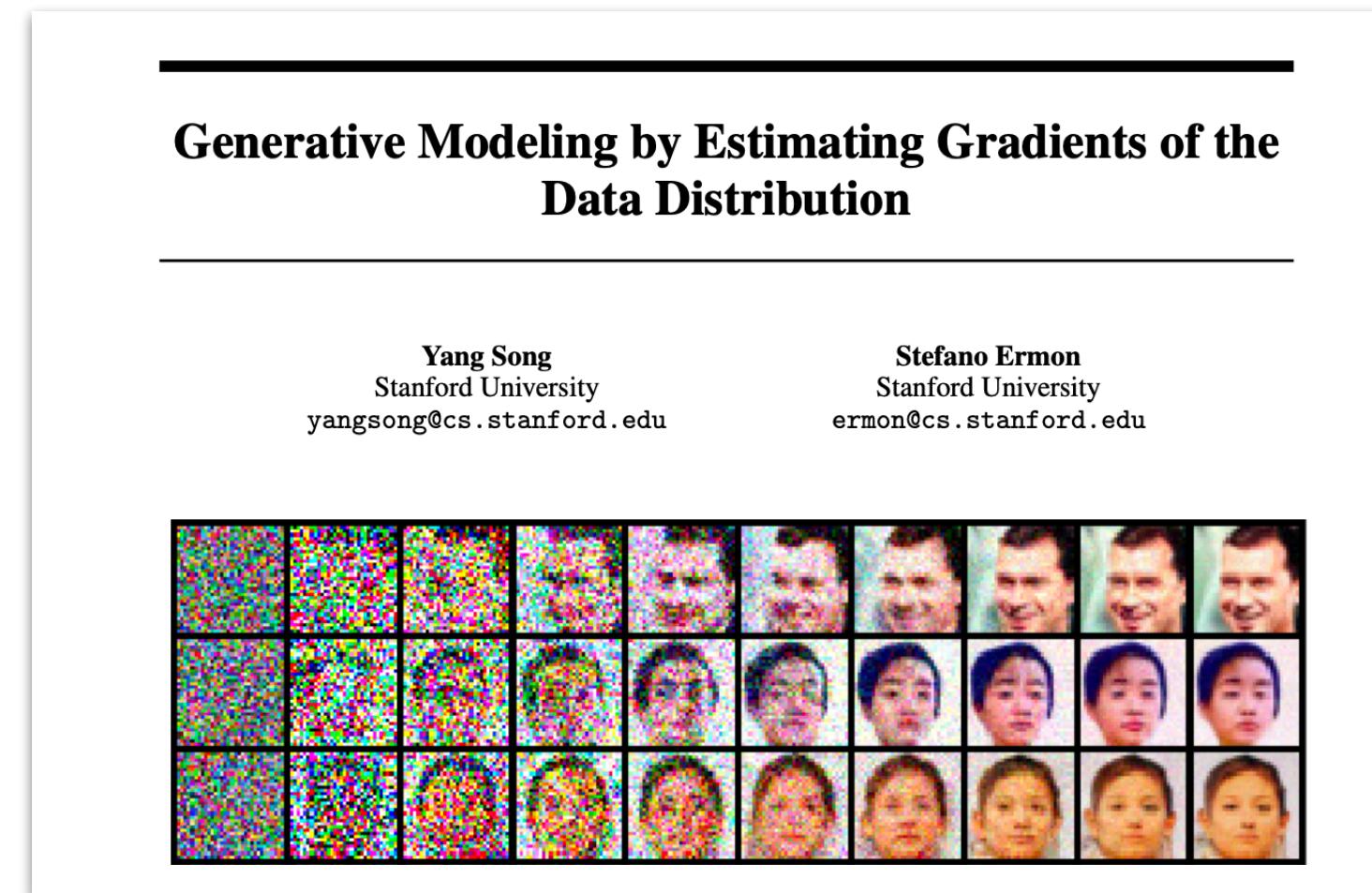
Pioneering papers:

Mar 2015: Sohl-Dickstein et al.



"The essential idea, inspired by non-equilibrium statistical physics, is to systematically and slowly destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data."

Jul 2019: Song et al.



"We introduce a new generative model where samples are produced via Langevin dynamics using gradients of the data distribution estimated with score matching."

Jul 2019: Ho et al.

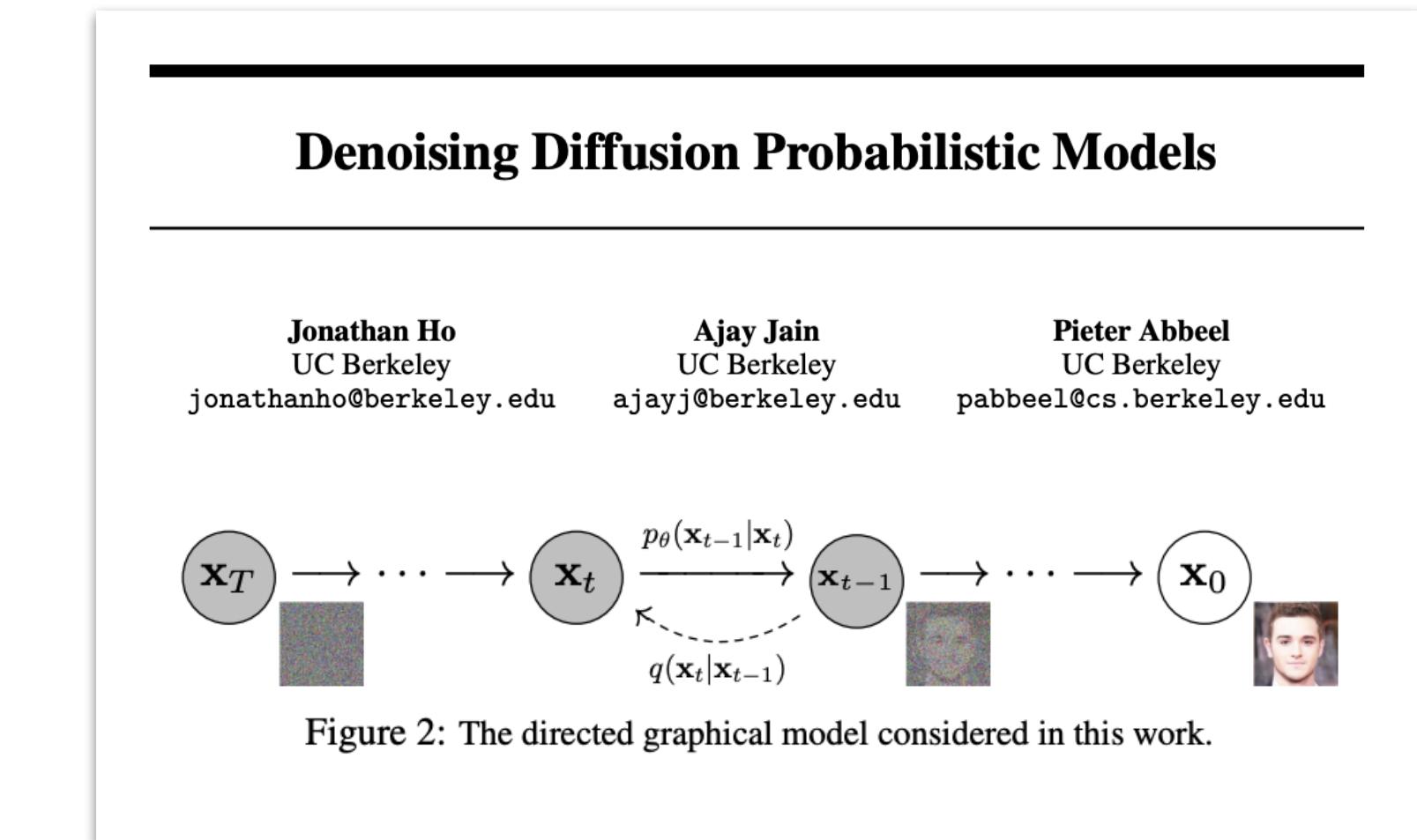


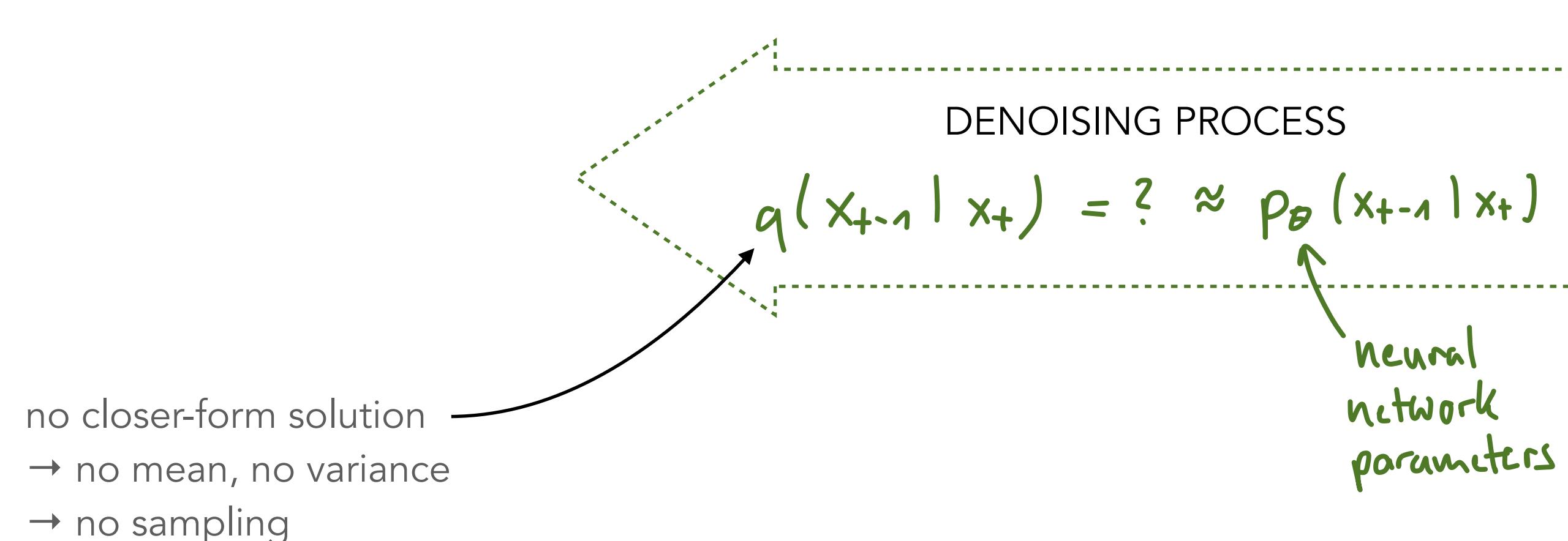
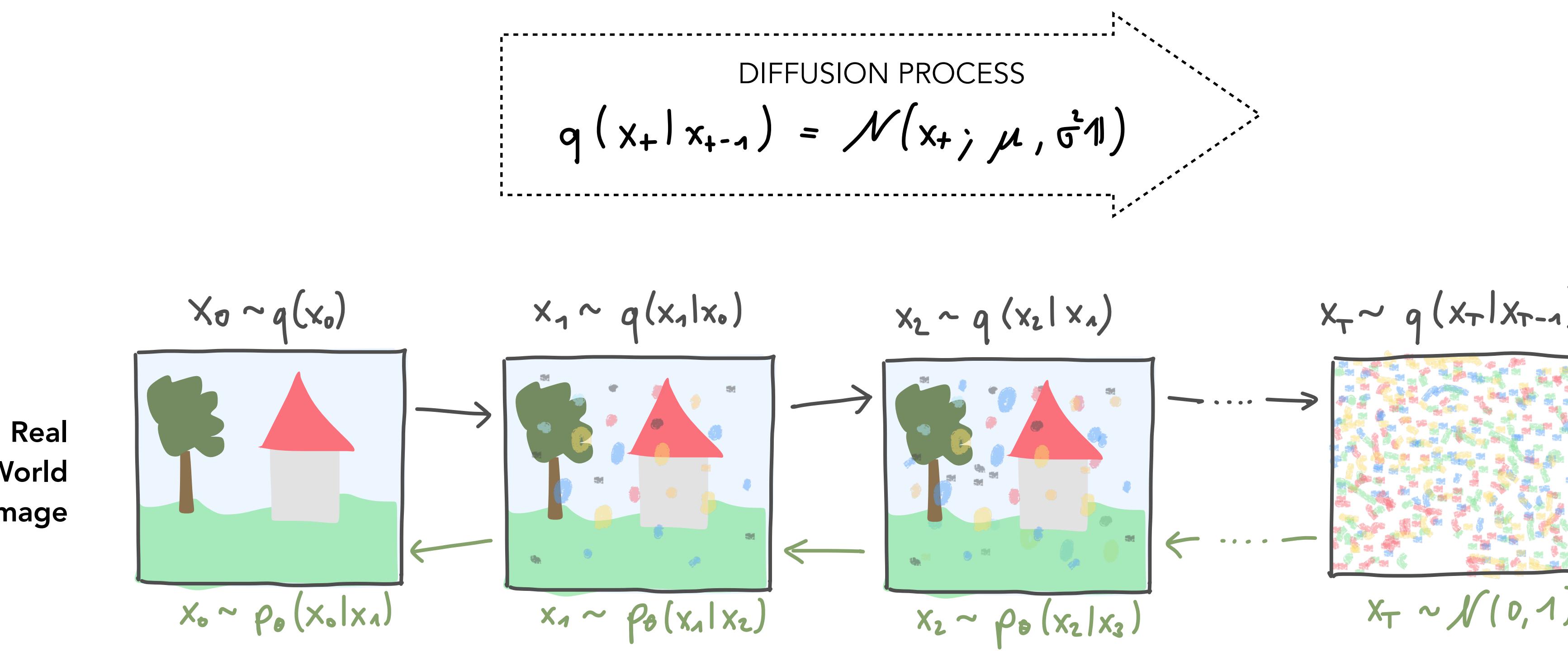
Figure 2: The directed graphical model considered in this work.

"This paper presents progress in diffusion probabilistic models [Sohl-Dickstein et al.]... When the diffusion consists of small amounts of Gaussian noise, it is sufficient to set the sampling chain transitions to conditional Gaussians too, allowing for a particularly simple neural network parameterisation."

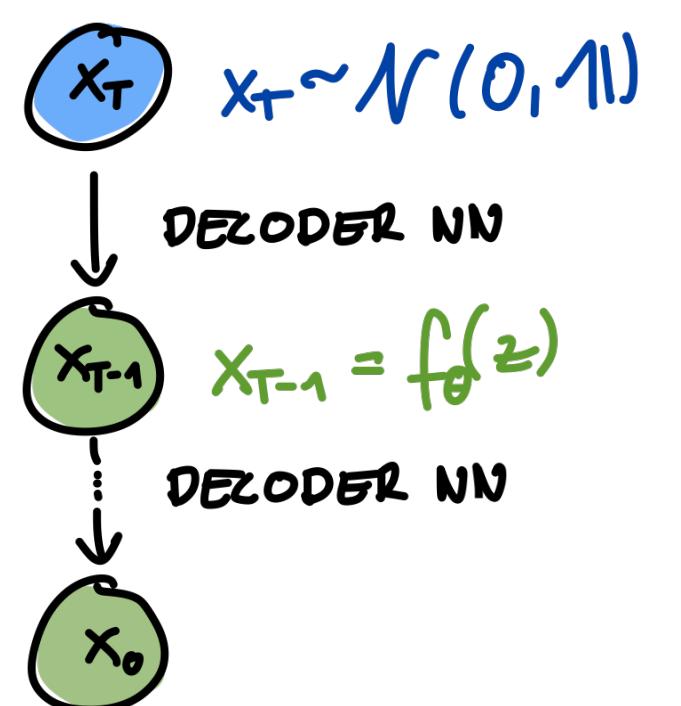


4 - Vanilla Diffusion Model

Basic idea:



Graphical Model



4 - Vanilla Diffusion Model

Bayes does not work in practice...

$$q(x_{t-1} | x_t) = \frac{q(x_t | x_{t-1}) \cdot q(x_{t-1})}{q(x_t)}$$

↑
we want to have a
closed form solution
of that distribution
so that we can sample
from

forward process
↓
 $q(x_t | x_{t-1})$

prior
↓
 $q(x_{t-1})$

↑
normalization
 $= \int q(x_t, x_{t-1}) dx_{t-1} = \int q(x_t | x_{t-1}) q(x_{t-1}) dx_{t-1}$

INTRACTABLE
↔
HIGH DIMENSIONAL
INTEGRATION

4.1 - Vanilla Diffusion Model - Forward Diffusion Process

Forward diffusion process:

Real data distribution: $x_0 \sim q(x)$

Definition of Forward Diffusion Process:

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t \mathbb{I})$$

with $\beta_t \in (0, 1)$

$$q(x_{1:T} | x_0) = \prod_{t=1}^T q(x_t | x_{t-1})$$

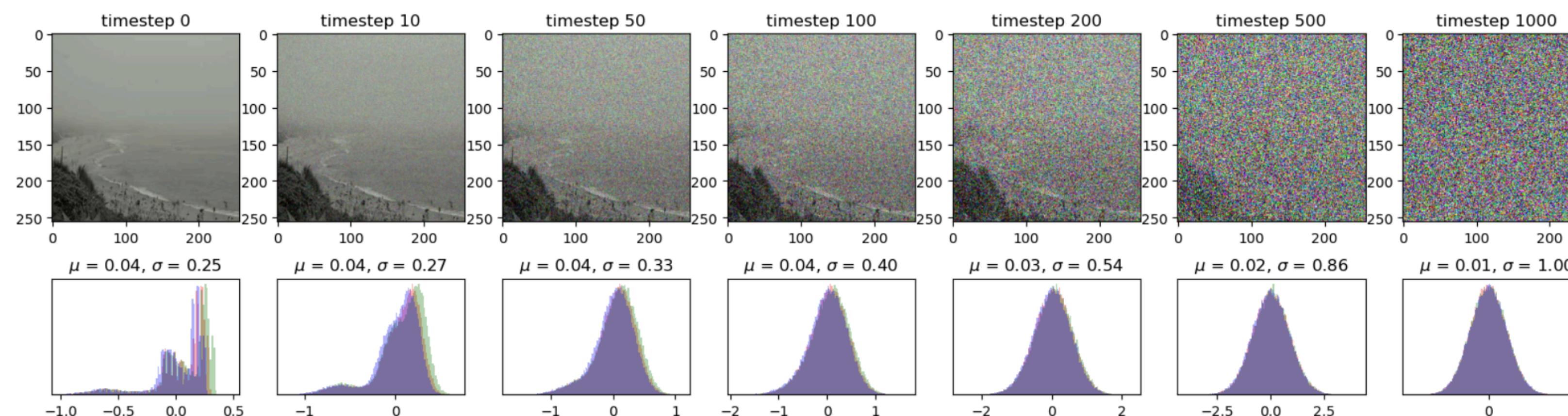
x_1, x_2, \dots, x_T

shifts μ towards 0 for large t

after a sufficiently large number of steps T we arrive at isotropic Gaussian noise.

Example:

$$T = 1000$$



Nice property (closed form for direct sampling)

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

$$\begin{aligned} x_t &= \sqrt{\bar{\alpha}_t} + \sqrt{1-\bar{\alpha}_t} \varepsilon \\ q(x_t | x_0) &= \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t} x_0, (1-\bar{\alpha}_t) \mathbb{I}) \end{aligned}$$

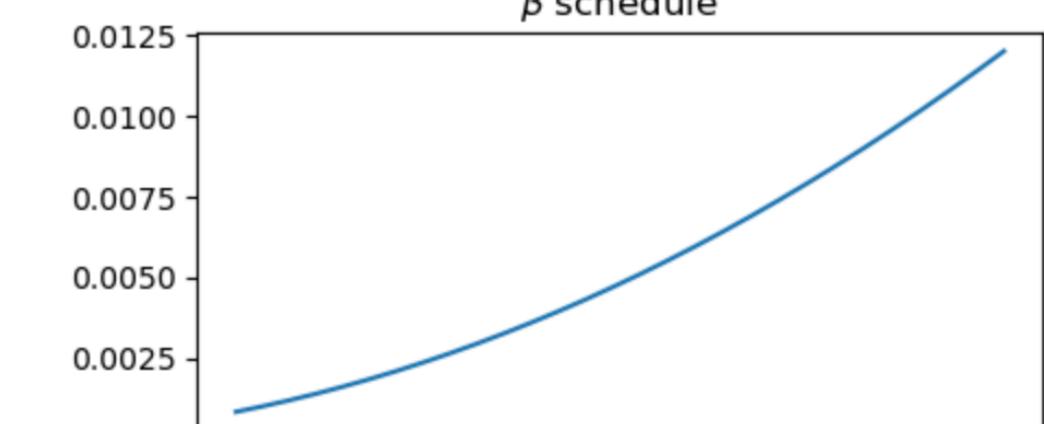
$$\alpha_t = 1 - \beta_t, \quad \varepsilon_{t-1}, \varepsilon_{t-2}, \dots \sim \mathcal{N}(0, 1)$$

$$\text{use } \mathcal{N}(0, \sigma_i^2 \mathbb{I}) \oplus \mathcal{N}(0, \sigma_i^2 \mathbb{I}) = \mathcal{N}(0, (\sigma_i^2 + \sigma_i^2) \mathbb{I})$$

$$\bar{\alpha}_t = \prod_{i=1}^t \alpha_i, \quad \varepsilon \sim \mathcal{N}(0, 1)$$

β_t increases quadratically from $\beta_1 = 0.85 \times 10^{-3}$ to $\beta_{1000} = 1.2 \times 10^{-1}$

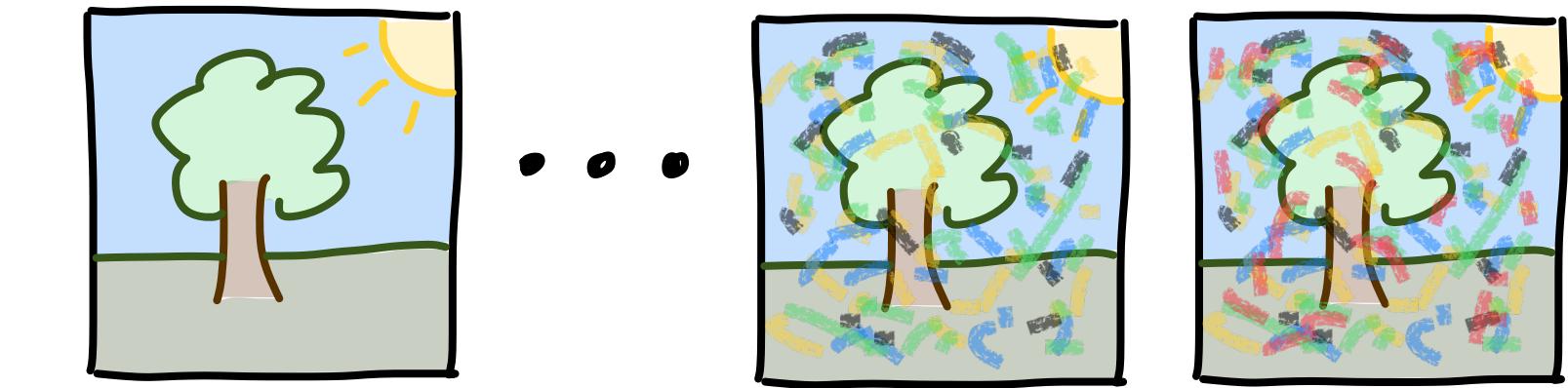
β schedule



4.1 - Vanilla Diffusion Model - Forward Diffusion Process

Important for later:

$q(x_{t-1} | x_t, x_0)$ is tractable!



$$\begin{aligned}
 q(x_{t-1} | x_t, x_0) &= \underbrace{q(x_t | x_{t-1}, x_0)}_{= q(x_t | x_{t-1})} \cdot \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} \quad \text{Bayes rule} \\
 &\propto \exp\left[-\frac{1}{2} \cdot \frac{(x_t - \sqrt{\alpha_t} x_{t-1})^2}{\beta_t}\right] \cdot \exp\left[-\frac{1}{2} \frac{(x_{t-1} - \sqrt{\bar{\alpha}_{t-1}} x_0)^2}{1 - \bar{\alpha}_{t-1}}\right] \cdot \exp\left[\frac{1}{2} \frac{(x_t - \sqrt{\alpha_t} x_0)^2}{1 - \bar{\alpha}_t}\right] \\
 &= \exp\left[-\frac{1}{2} \left(\frac{x_t^2 - 2\sqrt{\alpha_t} x_t x_{t-1} + \alpha_t x_{t-1}^2}{\beta_t} + \frac{x_{t-1}^2 - 2\sqrt{\bar{\alpha}_{t-1}} x_{t-1} x_0 + \bar{\alpha}_{t-1} x_0^2}{1 - \bar{\alpha}_{t-1}} - \frac{x_t^2 - 2\sqrt{\alpha_t} x_0 x_t + \bar{\alpha}_t x_0^2}{1 - \bar{\alpha}_t} \right)\right] \\
 &= \exp\left[-\frac{1}{2} \left(\underbrace{\left(\frac{\alpha_t}{\beta_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) x_{t-1}^2}_{1/\tilde{\beta}_t} - \underbrace{\left(\frac{2\sqrt{\alpha_t}}{\beta_t} x_t + \frac{2\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) x_{t-1}}_{2\tilde{\mu}_t/\tilde{\beta}_t} + C(x_t, x_0) \right)\right] \quad || \quad \left(\frac{x - \tilde{\mu}}{\tilde{\beta}} \right)^2 = \frac{1}{\tilde{\beta}} x^2 - \frac{2\tilde{\mu}}{\tilde{\beta}} x + \frac{\tilde{\mu}^2}{\tilde{\beta}} \\
 &= \exp\left[-\frac{1}{2} \left(\frac{x_{t-1} - \tilde{\mu}_t}{\tilde{\beta}_t} \right)^2\right] \\
 &\propto \mathcal{N}(x_{t-1}; \tilde{\mu}_t, \tilde{\beta}_t^{-1})
 \end{aligned}$$

not important

4.1 - Vanilla Diffusion Model - Forward Diffusion Process

$$\tilde{\beta}_+ = 1 / \left(\frac{\alpha_+}{\beta_+} + \frac{1}{1 - \bar{\alpha}_{t-1}} \right) = 1 / \frac{\alpha_+ (1 - \bar{\alpha}_{t-1}) + \beta_+}{\beta_+ (1 - \bar{\alpha}_{t-1})} = \frac{(1 - \bar{\alpha}_{t-1}) \beta_+}{\alpha_+ - \bar{\alpha}_+ + \beta_+} = \frac{1 - \bar{\alpha}_{t-1}}{1 - \alpha_+} \beta_+$$

\uparrow
 $1 - \alpha_+$

$$\begin{aligned}\tilde{\mu}_+ &= \left(\frac{\sqrt{\alpha_+}}{\beta_+} x_+ + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) \cdot \tilde{\beta}_+ = \left(\frac{\sqrt{\alpha_+}}{\beta_+} x_+ + \frac{\sqrt{\bar{\alpha}_{t-1}}}{1 - \bar{\alpha}_{t-1}} x_0 \right) \cdot \frac{1 - \bar{\alpha}_{t-1}}{1 - \alpha_+} \cdot \beta_+ \\ &= \frac{\sqrt{\alpha_+} (1 - \bar{\alpha}_{t-1})}{1 - \alpha_+} x_+ + \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_+}{1 - \alpha_+} x_0 \quad \underbrace{x_0}_{= \frac{1}{\sqrt{\alpha_+}} (x_+ - \sqrt{1 - \bar{\alpha}_+} \varepsilon_+)} \\ &= \dots = \frac{1}{\sqrt{\alpha_+}} \left(x_+ - \frac{1 - \bar{\alpha}_+}{\sqrt{1 - \bar{\alpha}_+}} \varepsilon_+ \right)\end{aligned}$$

$$\Rightarrow q(x_{t-1} | x_t, x_0) = \mathcal{N}(x_{t-1}; \underbrace{\frac{1}{\sqrt{\alpha_+}} \left(x_+ - \frac{1 - \bar{\alpha}_+}{\sqrt{1 - \bar{\alpha}_+}} \varepsilon_+ \right)}_{\tilde{\mu}_+}, \underbrace{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_+} \beta_+ \mathbb{I}}_{\tilde{\beta}_+})$$

4.2 - Vanilla Diffusion Model - Denoising Process

Denoising process:

Small β_t

$\Rightarrow q(x_{t-1} | x_t)$ is Gaussian but still intractable

\Rightarrow approximation $p_\theta(x_{t-1} | x_t)$ needed

Definition of denoising process:

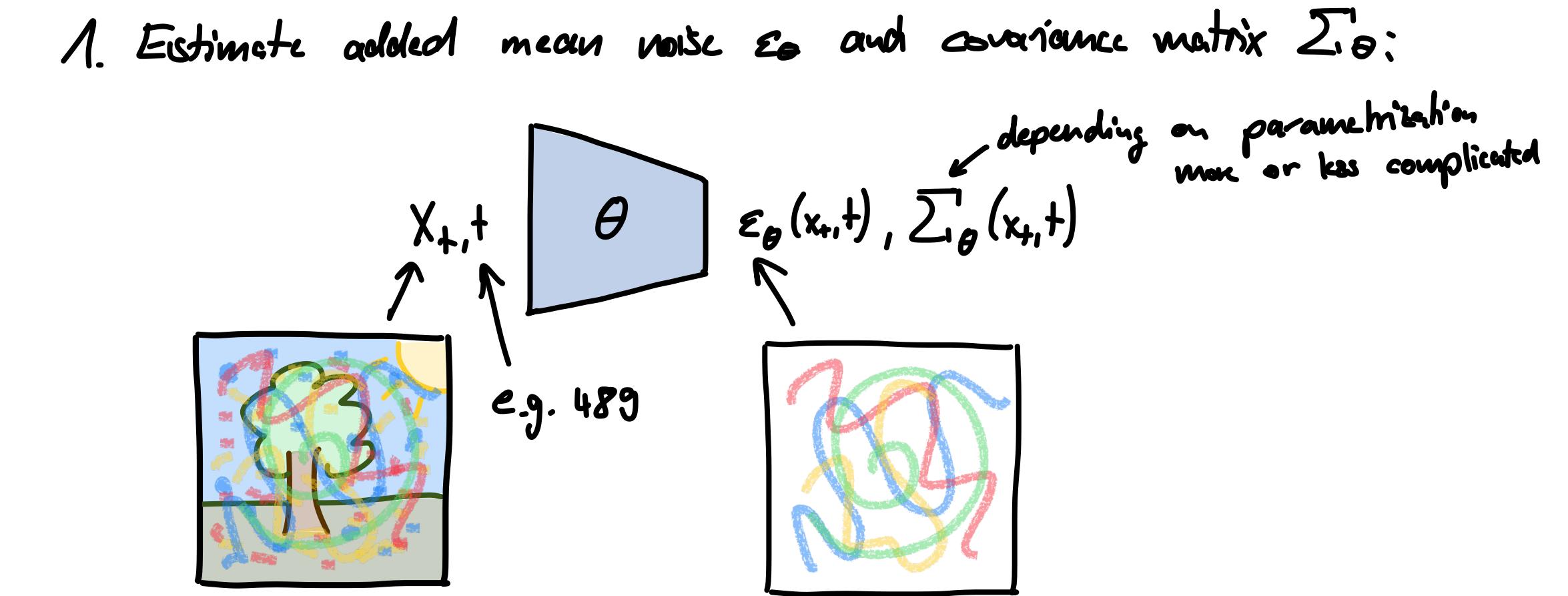
$$p_\theta(x_{t-1} | x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$$

$$\rightarrow p_\theta(x_{0:T}) = \underbrace{p_\theta(x_T)}_{= p(x_T) = \mathcal{N}(x_T; 0, 1)} \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$$

Instead of predicting μ_θ directly let's estimate the added noise $\varepsilon_\theta(x_t, t)$:

$$\mu_\theta(x_t, t) := \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \varepsilon_\theta(x_t, t) \right)$$

How one denoising step works in practice:



2. Sample a denoised sample from the estimated distribution:

$$x_{t-1} \sim \mathcal{N}\left(x_{t-1}; \underbrace{\frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \varepsilon_\theta(x_t, t) \right)}_{\mu_\theta(x_t, t)}, \Sigma_\theta(x_t, t)\right)$$

$$= \frac{1}{\sqrt{\alpha_t}} \left(\text{clear image} - \frac{1-\alpha_t}{\sqrt{1-\alpha_t}} \varepsilon_\theta(x_t, t) \right) + \dots$$

4.3 - Vanilla Diffusion Model - Loss Function

Loss function: Jansen inequality/ELBO/VLB

$$\begin{aligned} \max &\stackrel{!}{=} \left\langle \log p_{\theta}(x_0) \right\rangle_{x_0 \sim q(x)} \\ &= \left\langle \log \left(\int p_{\theta}(x_{0:T}) dx_{1:T} \right) \right\rangle_{x_0 \sim q(x)} \quad | p_{\theta}(x_0) = \int p_{\theta}(x_0|x_1) \dots p(x_{T-1}|x_T) p(x_T) dx_1 \dots dx_T \\ &= \left\langle \log \left(\int q(x_{1:T}|x_0) \cdot \frac{p_{\theta}(x_{0:T})}{q(x_{1:T}|x_0)} dx_{1:T} \right) \right\rangle_{x_0 \sim q(x)} \quad | \text{Inserting a "1"} \\ &= \left\langle \log \left\langle \frac{p_{\theta}(x_{0:T})}{q(x_{1:T}|x_0)} \right\rangle_{x_{1:T} \sim q(x_{1:T}|x_0)} \right\rangle_{x_0 \sim q(x)} \quad | \text{Jansen inequality} \\ &\geq \left\langle \log \frac{p_{\theta}(x_{0:T})}{q(x_{1:T}|x_0)} \right\rangle_{x_{0:T} \sim q(x)} \end{aligned}$$

$$\Rightarrow \boxed{\mathcal{L} = \left\langle \log \frac{q(x_{1:T}|x_0)}{p_{\theta}(x_{0:T})} \right\rangle_{x_{0:T} \sim q(x)}}$$

Also known as ELBO (Evidence Lower Bound) or VLB (Variational Lower Bound). Research "Variational Inference" for more details.

4.3 - Vanilla Diffusion Model - Loss Function

Loss function: Separation in different KL-divergence terms

$$\mathcal{L} = \left\langle \log \frac{q(x_{1:T} | x_0)}{p_\theta(x_{0:T})} \right\rangle_{x \sim q}$$

"every x_t in the equation
is sampled from q "

$$= \left\langle \log \frac{\prod_{t=1}^T q(x_t | x_{t-1})}{p_\theta(x_T) \prod_{t=1}^{T-1} p_\theta(x_{t-1} | x_t)} \right\rangle_{x \sim q}$$

$$= \left\langle -\log p_\theta(x_T) + \sum_{t=1}^{T-1} \log \frac{q(x_t | x_{t-1})}{p_\theta(x_{t-1} | x_t)} \right\rangle_{x \sim q}$$

$$= \left\langle -\log p_\theta(x_T) + \left(\sum_{t=2}^T \log \frac{q(x_t | x_{t-1})}{p_\theta(x_{t-1} | x_t)} \right) + \log \frac{q(x_0 | x_0)}{p_\theta(x_0 | x_0)} \right\rangle_{x \sim q}$$

$$= \left\langle -\log p_\theta(x_T) + \sum_{t=2}^T \log \left(\frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \cdot \frac{q(x_t | x_0)}{q(x_{t-1} | x_0)} \right) + \log \frac{q(x_0 | x_0)}{p_\theta(x_0 | x_0)} \right\rangle_{x \sim q}$$

$$= \left\langle -\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} + \sum_{t=2}^T \log \frac{q(x_t | x_0)}{q(x_{t-1} | x_0)} + \log \frac{q(x_0 | x_0)}{p_\theta(x_0 | x_0)} \right\rangle_{x \sim q}$$

$$= \left\langle -\log p_\theta(x_T) + \sum_{t=2}^T \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} + \log \frac{q(x_T | x_0)}{q(x_0 | x_0)} + \log \frac{q(x_0 | x_0)}{p_\theta(x_0 | x_0)} \right\rangle_{x \sim q}$$

$$= \left\langle \log \frac{q(x_T | x_0)}{p_\theta(x_T)} + \sum_{t=1}^T \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} - \log p_\theta(x_0 | x_0) \right\rangle_{x \sim q}$$

$$= \left\langle \log \frac{q(x_T | x_0)}{p_\theta(x_T)} \right\rangle_{x \sim q} + \sum_{t=2}^T \left\langle \log \frac{q(x_{t-1} | x_t, x_0)}{p_\theta(x_{t-1} | x_t)} \right\rangle_{x \sim q} - \left\langle \log p_\theta(x_0 | x_0) \right\rangle_{x \sim q}$$

$$= \underbrace{\left\langle D_{KL}(q(x_T | x_0) || p_\theta(x_T)) \right\rangle}_{L_T} + \underbrace{\sum_{t=2}^T \left\langle D_{KL}(q(x_{t-1} | x_t, x_0) || p_\theta(x_{t-1} | x_t)) \right\rangle}_{L_{t-1}} - \underbrace{\left\langle \log p_\theta(x_0 | x_0) \right\rangle}_{L_0}$$

$$q(x_t | x_{t-1}) = q(x_t | x_{t-1}, x_0) \stackrel{\text{Bayes}}{=} \frac{q(x_{t-1} | x_t, x_0) \cdot q(x_t | x_0)}{q(x_{t-1} | x_0)}$$

$$\log \left(\frac{q(x_2 | x_0)}{q(x_1 | x_0)} \cdot \frac{q(x_3 | x_0)}{q(x_2 | x_0)} \cdots \frac{q(x_T | x_0)}{q(x_{T-1} | x_0)} \right) = \log \frac{q(x_T | x_0)}{q(x_1 | x_0)}$$

4.3 - Vanilla Diffusion Model - Loss Function

Loss function: Explicit calculation of KL-divergence terms

- $\mathcal{L}_T = D_{KL}(q(x_T | x_0) \| p_\theta(x_T))$

$$q(x_T | x_0) = \mathcal{N}(x_T; \sqrt{\bar{\alpha}_T} x_0, (1 - \bar{\alpha}_T) \mathbb{I})$$

$$p_\theta(x_T) = p(x_T) = \mathcal{N}(x_T; 0, \mathbb{I})$$

\rightarrow independent of θ \rightarrow can be ignored in the loss

- $\mathcal{L}_0 = -\log p_\theta(x_0 | x_0)$

$$\hookrightarrow \mathcal{N}(x_0; \mu_\theta(x_0, 1), \Sigma_\theta(x_0, 1))$$

$$= \frac{1}{2} (x_0 - \mu_\theta)^T \Sigma_\theta^{-1} (x_0 - \mu_\theta) - \log(\text{Prefactor})$$

\rightarrow tractable 😊

- $\mathcal{L}_{t+1} = D_{KL}(q(x_{t+1} | x_t, x_0) \| p_\theta(x_{t+1} | x_t)) \quad \text{for } 2 \leq t \leq T$

$$q(x_{t+1} | x_t, x_0) = \mathcal{N}(x_{t+1}; \underbrace{\frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_t \right)}_{\tilde{\mu}_\theta(x_t, x_0)}, \underbrace{\frac{1 - \bar{\alpha}_{t+1}}{1 - \bar{\alpha}_t} \beta_t \mathbb{I}}_{\Sigma_\theta(x_t, t)})$$

$$p_\theta(x_{t+1} | x_t) = \mathcal{N}(x_{t+1}; \underbrace{\frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(x_t, t) \right)}_{\mu_\theta(x_t, t)}, \Sigma_\theta(x_t, t))$$

KL-Divergence between two multivariate Gaussians has a closed form solution:

$$\mathcal{L}_{t+1} = \frac{1}{2 \|\Sigma_\theta(x_t, t)\|^2} \left\| \tilde{\mu}_\theta(x_t, x_0) - \mu_\theta(x_t, t) \right\|^2 + \text{const.}$$

$$= \frac{1}{2 \|\Sigma_\theta(x_t, t)\|^2} \cdot \left\| \frac{1}{\sqrt{\bar{\alpha}_t}} \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} (\varepsilon_t - \varepsilon_\theta(x_t, t)) \right\|^2 + \text{const.}$$

$$= \frac{(1 - \alpha_t)^2}{2 \alpha_t (1 - \bar{\alpha}_t) \|\Sigma_\theta(x_t, t)\|^2} \left\| \varepsilon_t - \varepsilon_\theta(x_t, t) \right\|^2 + \text{const.}$$

\rightarrow tractable 😊

4.3 - Vanilla Diffusion Model - Loss Function

Loss Function: Simplified Loss

3.4 Simplified training objective

With the reverse process and decoder defined above, the variational bound, consisting of terms derived from Eqs. (12) and (13), is clearly differentiable with respect to θ and is ready to be employed for training. However, we found it beneficial to sample quality (and simpler to implement) to train on the following variant of the variational bound:

$$L_{\text{simple}}(\theta) := \mathbb{E}_{t, x_0, \epsilon} \left[\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2 \right] \quad (14)$$

where t is uniform between 1 and T . The $t = 1$ case corresponds to L_0 with the integral in the discrete decoder definition (13) approximated by the Gaussian probability density function times the bin width, ignoring σ_1^2 and edge effects. The $t > 1$ cases correspond to an unweighted version of Eq. (12), analogous to the loss weighting used by the NCSN denoising score matching model [55]. (L_T does not appear because the forward process variances β_t are fixed.) Algorithm 1 displays the complete training procedure with this simplified objective.

Since our simplified objective (14) discards the weighting in Eq. (12), it is a weighted variational bound that emphasizes different aspects of reconstruction compared to the standard variational bound [18, 22]. In particular, our diffusion process setup in Section 4 causes the simplified objective to down-weight loss terms corresponding to small t . These terms train the network to denoise data with very small amounts of noise, so it is beneficial to down-weight them so that the network can focus on more difficult denoising tasks at larger t terms. We will see in our experiments that this reweighting leads to better sample quality.

excerpt Ho et al.

independent of θ

$$\mathcal{L} = \langle \mathcal{L}_T + \mathcal{L}_{T-1} + \dots + \mathcal{L}_0 \rangle_{x_0 \sim q}$$

$$\Rightarrow \mathcal{L} = \langle \mathcal{L}_0 \rangle_{x_0 \sim q} + \sum_{t=2}^T \langle \mathcal{L}_{t-1} \rangle_{x_t \sim q}$$

from before:

$$\mathcal{L}_0 = \frac{1}{2} (x_0 - \mu_\theta)^T \Sigma_\theta^{-1} (x_0 - \mu_\theta) - \log(\text{Prefactor})$$

$$\mathcal{L}_{t-1} = \underbrace{\frac{(1 - \alpha_t)^2}{2\alpha_t(1 - \bar{\alpha}_t) \|\Sigma_\theta(x_t, t)\|^2}}_{\text{weighting factor;}} \|\epsilon_t - \epsilon_\theta(x_t, t)\|^2 + \text{const.}$$

will be dropped

$$\Rightarrow \mathcal{L}_{\text{simple}} = \left\langle \|\epsilon_t - \epsilon_\theta(\underbrace{\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon_t, t}_{= x_t})\|^2 \right\rangle_{x_0 \sim q, t \sim [0, \dots, T], \epsilon \sim \mathcal{N}(0, 1)}$$

4.4 - Vanilla Diffusion Model - Training and Inference

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

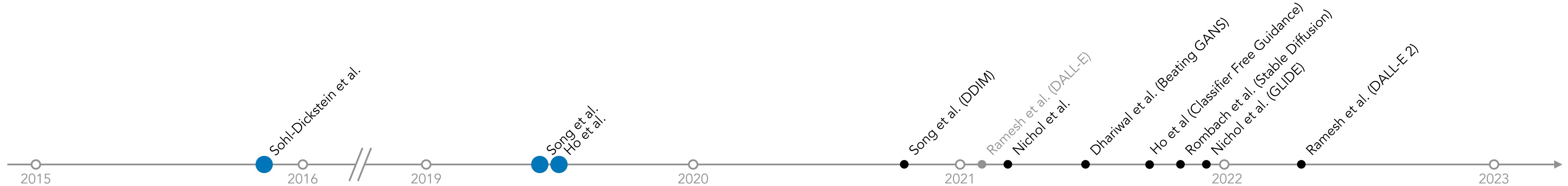
```

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  $\mathbf{x}_0$ 

```

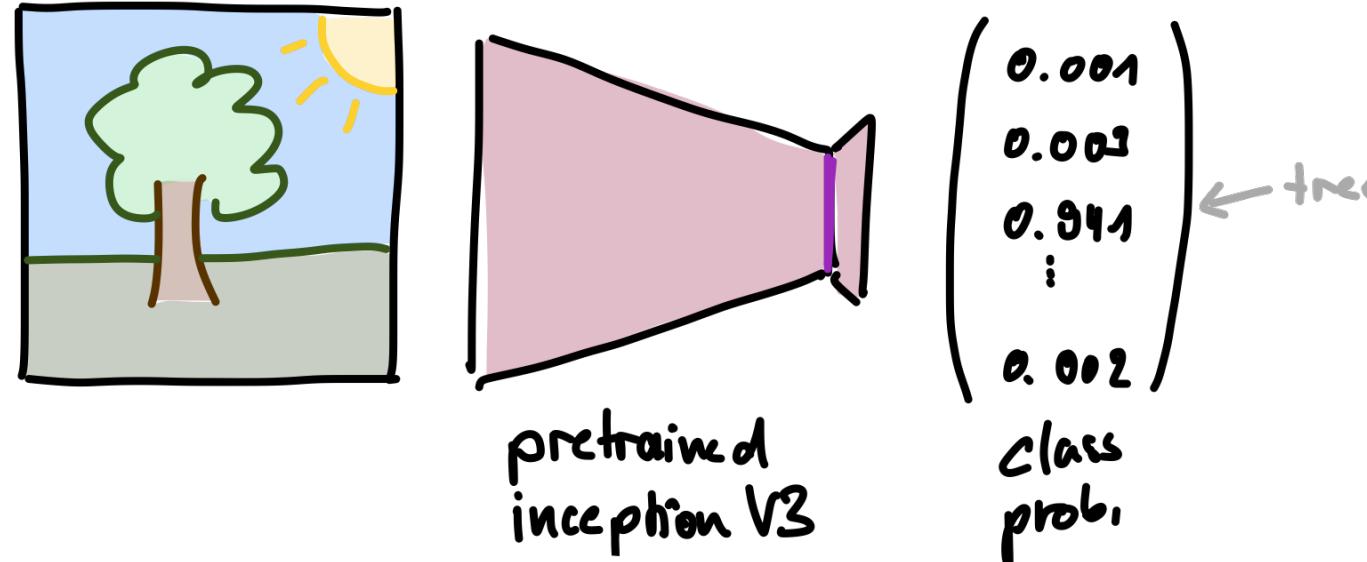
set to \tilde{p}_t or \hat{p}_t

excerpt Ho et al.

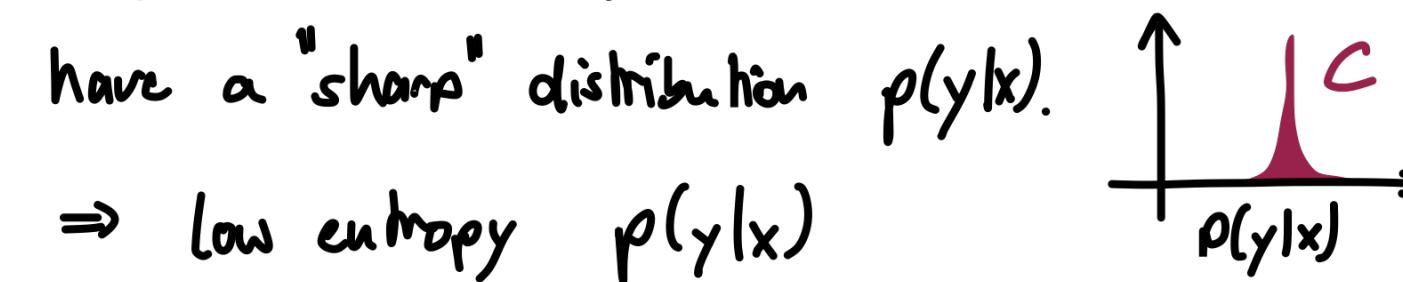


4.5 - Vanilla Diffusion Model - KPIs

Inception Score (IS)

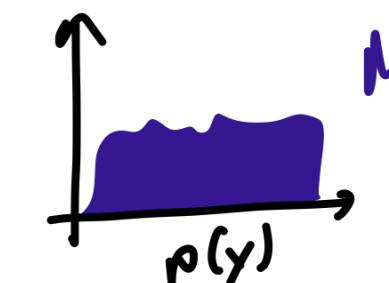


- Images with meaningful full objects should have a "sharp" distribution $p(y|x)$.



- Model should output a variety of images.

$$\Rightarrow p(y) = \int p(y|x) dx \text{ high entropy}$$



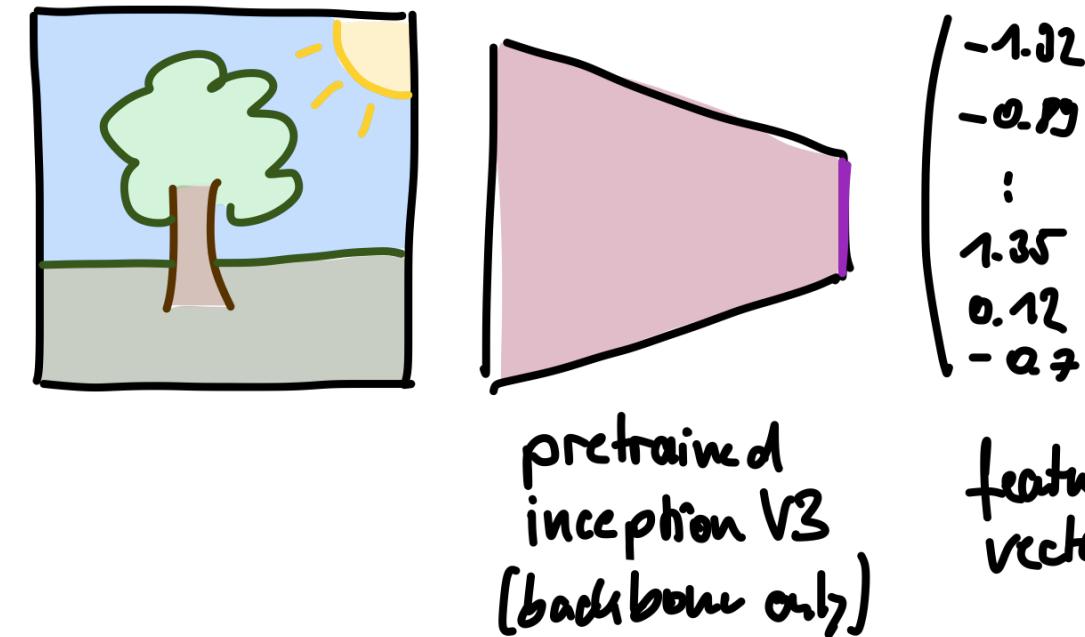
\Rightarrow Combined metric: relative entropy

$$IS = \exp \left[\left\langle D_{KL}(C || M) \right\rangle_{x \sim \text{all samples}} \right]$$

$$IS \in [1, n_{\text{classes}}]$$

bad ↑ good ↘

Frechet Inception Distance (FID)



- Downside of IS: No comparison with real images
- Feature distribution of real images: μ, Σ
- Feature distribution of syn. images: $\tilde{\mu}, \tilde{\Sigma}$
- Both distributions should be similar
- Frechet distance should be small

$$FID = \|\mu - \tilde{\mu}\|^2 + \text{Tr}(\Sigma + \tilde{\Sigma} - 2\sqrt{\Sigma \cdot \tilde{\Sigma}})$$

$$FID \in [0, \infty]$$

↑ good ↑ bad

Table 1: CIFAR10 results. NLL measured in bits/dim.

Model	IS	FID	NLL Test (Train)
Conditional			
EBM [11]	8.30	37.9	
JEM [17]	8.76	38.4	
BigGAN [3]	9.22	14.73	
StyleGAN2 + ADA (v1) [29]	10.06	2.67	
Unconditional			
Diffusion (original) [53]			≤ 5.40
Gated PixelCNN [59]	4.60	65.93	3.03 (2.90)
Sparse Transformer [7]			2.80
PixelIQN [43]	5.29	49.46	
EBM [11]	6.78	38.2	
NCSNv2 [56]			31.75
NCSN [55]	8.87 ± 0.12	25.32	
SNGAN [39]	8.22 ± 0.05	21.7	
SNGAN-DDLS [4]	9.09 ± 0.10	15.42	
StyleGAN2 + ADA (v1) [29]	9.74 ± 0.05	3.26	
Ours (L , fixed isotropic Σ)	7.67 ± 0.13	13.51	≤ 3.70 (3.69)
Ours (L_{simple})	9.46 ± 0.11	3.17	≤ 3.75 (3.72)

excerpt Ho et al.

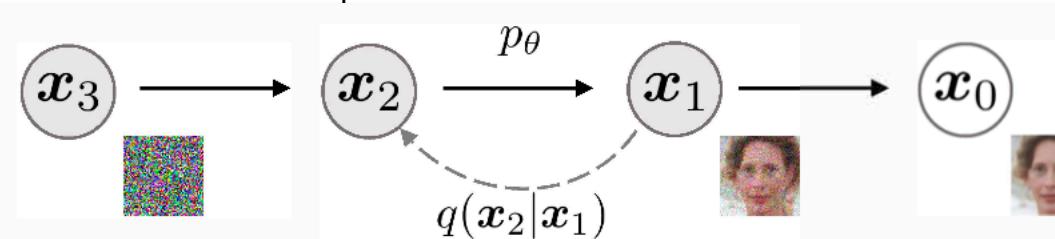
Links:
[Blogpost: IS](#)
[Blogpost: FID](#)

5.1 Improvements - Song et al (DDIM)

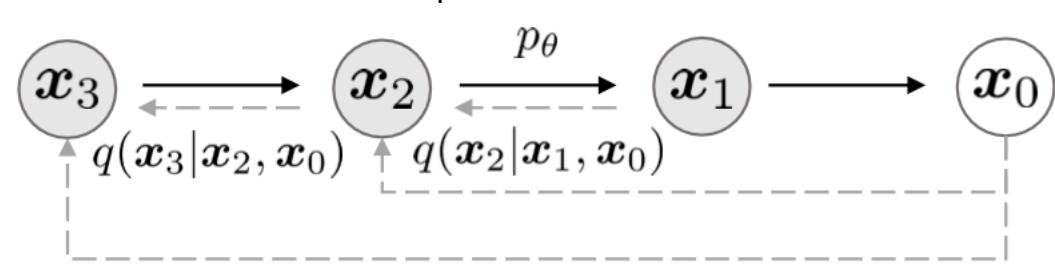
DENOISING DIFFUSION IMPLICIT MODELS

Jiaming Song, Chenlin Meng & Stefano Ermon
Stanford University
`{tsong, chenlin, ermon}@cs.stanford.edu`

Markovian forward process:



Non-Markovian forward process:



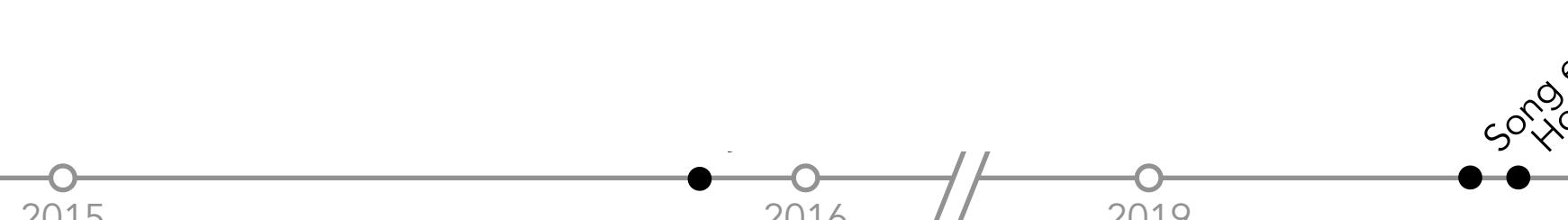
Definition of non-Markovian forward process:

$$q_\sigma(x_{1:T}|x_0) := q_\sigma(x_T|x_0) \prod_{t=2}^T q_\sigma(x_{t-1}|x_t, x_0)$$

$$q_\sigma(x_{t-1}|x_t, x_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}}x_0 + \frac{x_t - \sqrt{\alpha_t}x_0}{\sqrt{1-\alpha_{t-1}-\sigma_t^2}}, \frac{\sigma_t^2}{1-\alpha_{t-1}}\right) \quad (*)$$

\Rightarrow Bayes yields $q_\sigma(x_t|x_{t-1}, x_0)$

\Rightarrow One can prove $q_\sigma(x_t|x_0) = \mathcal{N}\left(\sqrt{\alpha_t}x_0, (1-\alpha_t)\mathbb{I}\right)$



Resulting Loss function:

$$\mathbb{E}_{\mathbf{x}_{0:T} \sim q_\sigma(\mathbf{x}_{0:T})} \left[\log q_\sigma(\mathbf{x}_T|\mathbf{x}_0) + \sum_{t=2}^T \log q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) - \sum_{t=1}^T \log p_\theta^{(t)}(\mathbf{x}_{t-1}|\mathbf{x}_t) - \log p_\theta(\mathbf{x}_T) \right]$$

- Best models for "normal" DDPM also minimise this loss
- Loss is independent of the exact time series of the forward procedure \rightarrow allows us to skip steps

Sampling Process:

$$x_t = \sqrt{\alpha_t} x_0 + \sqrt{1-\alpha_t} \varepsilon_t \rightarrow x_0 = (x_t - \sqrt{1-\alpha_t} \varepsilon_t) / \sqrt{\alpha_t} \quad \text{Inserting in (*):}$$

$$x_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left(\frac{x_t - \sqrt{1-\alpha_t} \epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } x_0\text{"}} + \underbrace{\sqrt{1-\alpha_{t-1}-\sigma_t^2} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t)}_{\text{"direction pointing to } \mathbf{x}_t\text{"}} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

Extreme cases:

$\sigma_t = 0 \Rightarrow$ Deterministic model (DDIM)

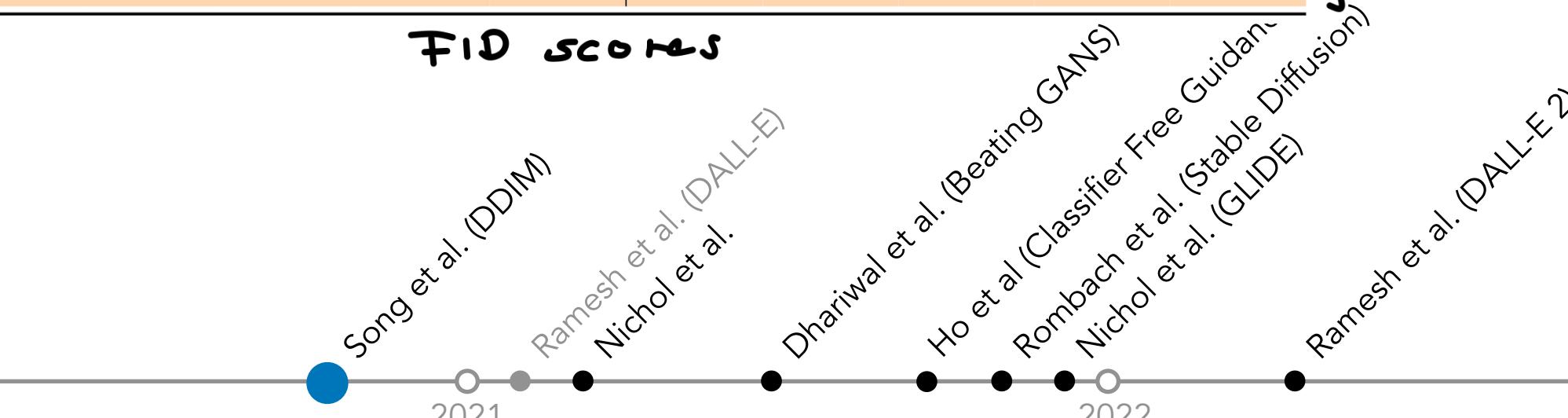
$$\hat{\alpha}_t = \sqrt{\frac{1-\alpha_{t-1}}{1-\alpha_t}} \sqrt{\frac{1-\alpha_t}{\alpha_{t-1}}}$$

\Rightarrow Markovian process of "normal" DDPMs

Experimental results:

S	CIFAR10 (32 × 32)					CelebA (64 × 64)				
	10	20	50	100	1000	10	20	50	100	1000
0.0	13.36	6.84	4.67	4.16	4.04	17.33	13.73	9.17	6.53	3.51
0.2	14.04	7.11	4.77	4.25	4.09	17.66	14.11	9.51	6.79	3.64
0.5	16.66	8.35	5.25	4.46	4.29	19.86	16.06	11.01	8.09	4.28
1.0	41.07	18.36	8.01	5.78	4.73	33.12	26.03	18.48	13.93	5.98
$\hat{\sigma}$	367.43	133.37	32.72	9.99	3.17	299.71	183.83	71.71	45.20	3.26

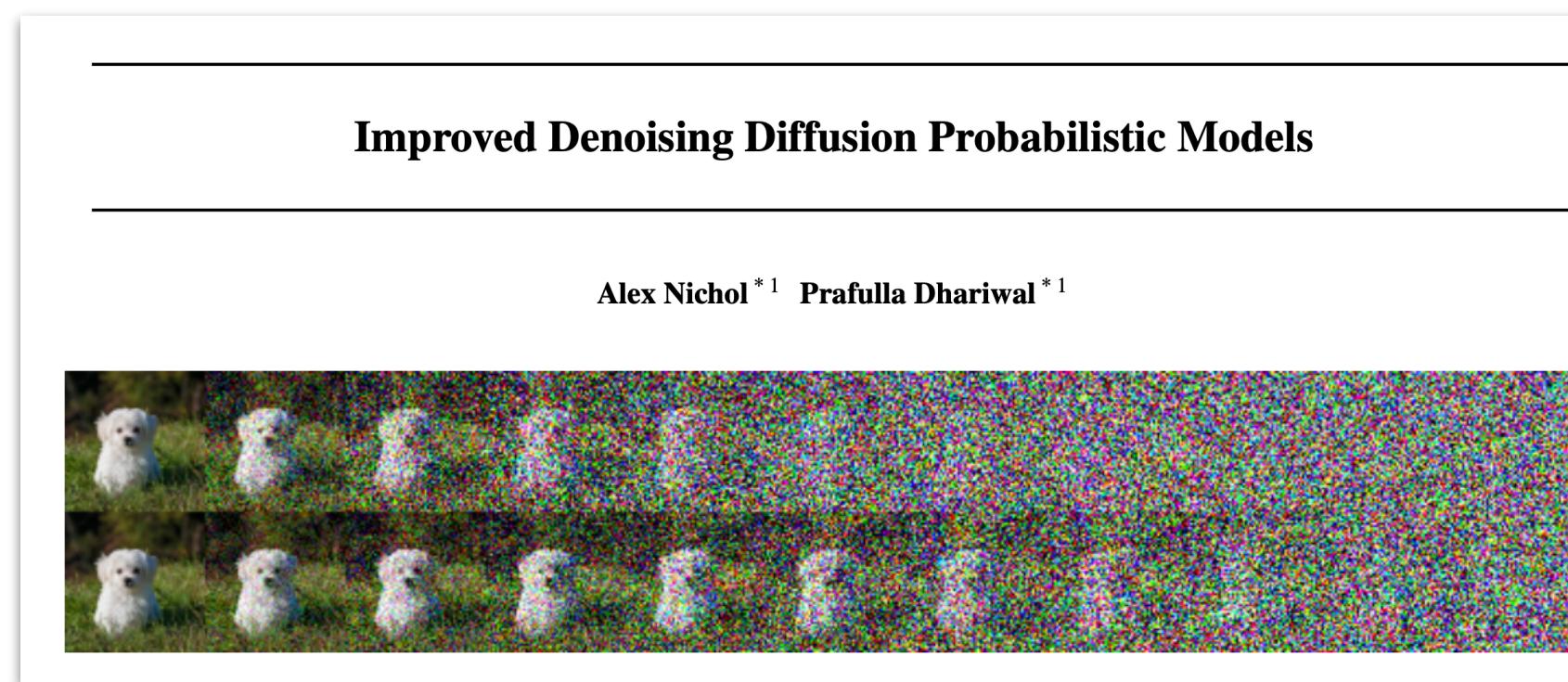
FID scores



← DDIM

} Ho et al

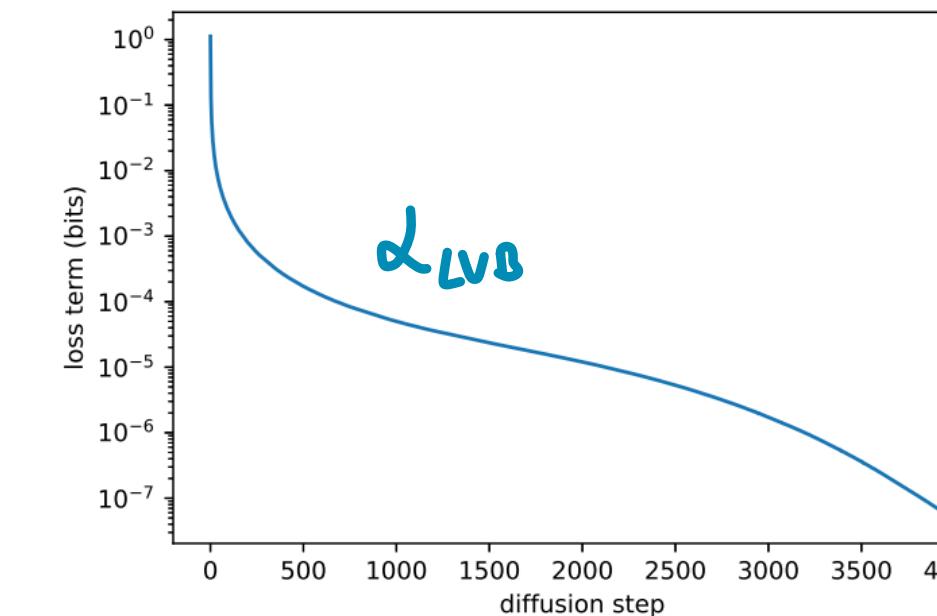
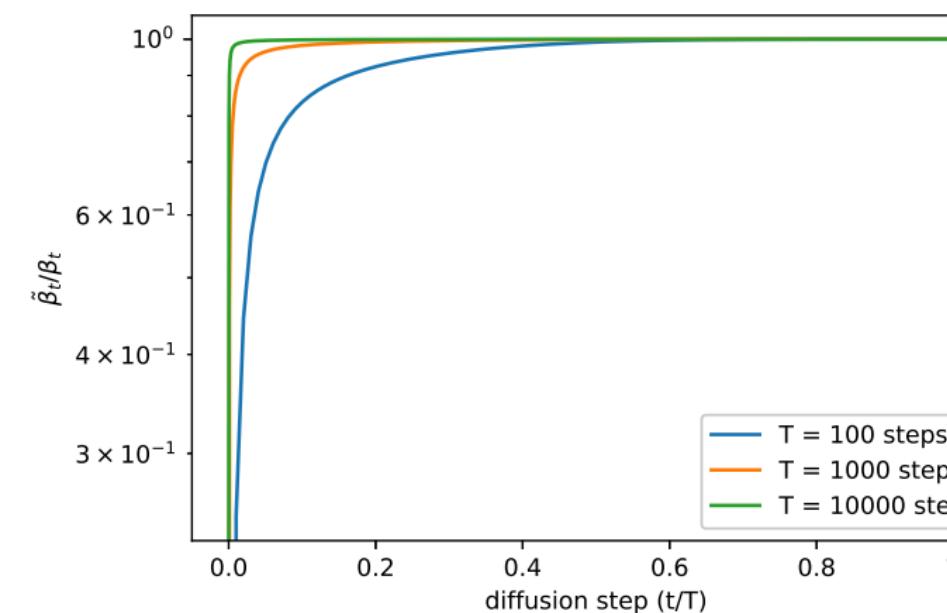
5.2 Improvements - Nichol et al



Main contributions:

1. Learning Σ_θ and hybrid loss
 2. Improved noise schedule β_t
 3. Reduced gradient noise by improved sampling
- Nice side effect: Faster sampling speed

Learning Σ_θ and hybrid loss



$$L_{\text{vlb}} := L_0 + L_1 + \dots + L_{T-1} + L_T$$

$$L_0 := -\log p_\theta(x_0|x_1)$$

$$L_{t-1} := D_{KL}(q(x_{t-1}|x_t, x_0) \parallel p_\theta(x_{t-1}|x_t))$$

$$L_T := D_{KL}(q(x_T|x_0) \parallel p(x_T))$$

$$L_{\text{simple}} = E_{t,x_0,\epsilon} [||\epsilon - \epsilon_\theta(x_t, t)||^2]$$



hyperparameter, set to β_t or $\tilde{\beta}_t$

Ho et al.: $\Sigma_\theta(x_t, t) = \overline{\sigma}_t^2 \mathbb{I}$

Recall: $q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{1-\beta_t} x_{t-1}, \beta_t \mathbb{I})$ (diffusion)

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

$\Rightarrow \tilde{\beta}_t \leq \overline{\sigma}_t^2 \leq \beta_t$

$\frac{1 - \bar{x}_{t-1} \beta_t}{1 - \bar{\alpha}_t}$

⇒ Real range is small, especially for large T

⇒ Hard to learn for a neural network (even at log-scale)

Ausatz:

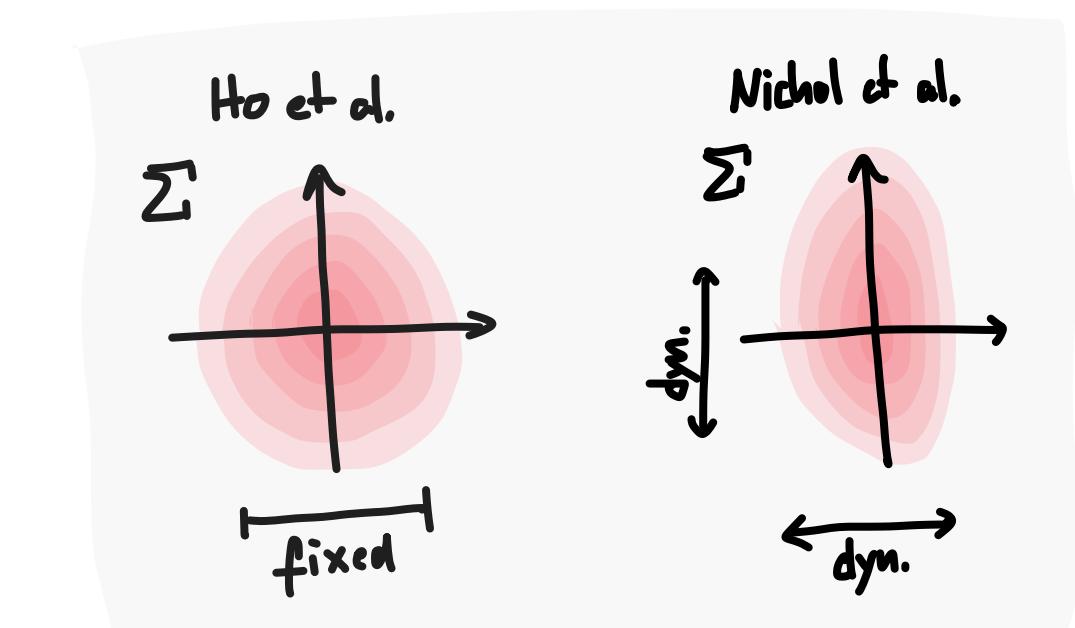
$$\Sigma_\theta(x_t, t) = \exp(v \log \beta_t + (1-v) \log \tilde{\beta}_t) \quad v \in \mathbb{R}^d$$

needs to be trained with:

$$\mathcal{L}_{\text{hybrid}} := \mathcal{L}_{\text{simple}} + \lambda \mathcal{L}_{\text{VLB}}$$

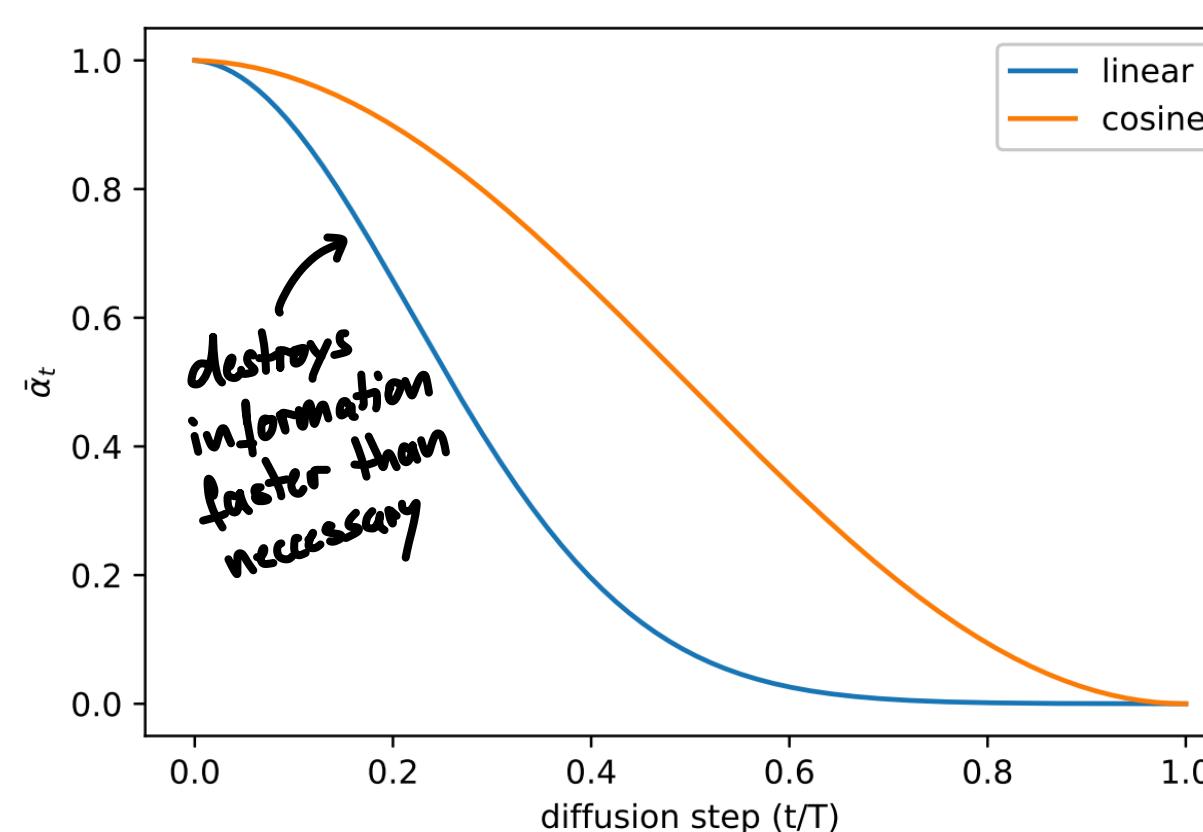
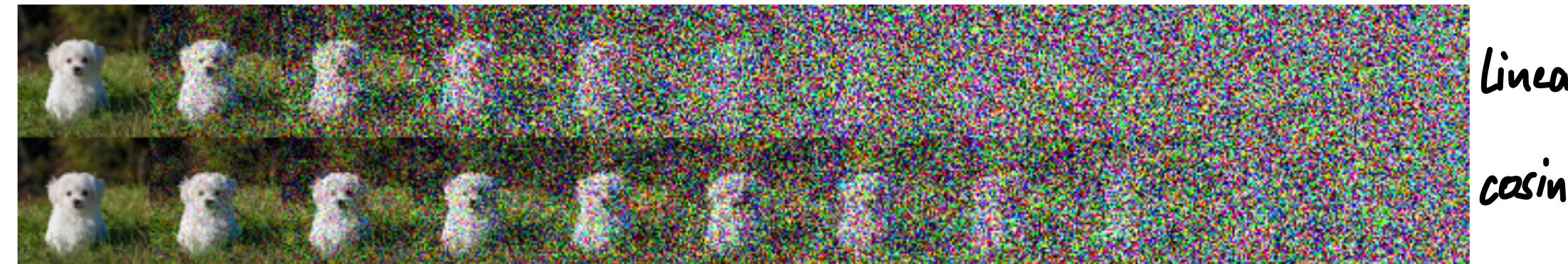
\uparrow independent of Σ_θ

$\uparrow \frac{1}{1000}$ only important for small t ; this is also where β_t and $\tilde{\beta}_t$ deviate the most!



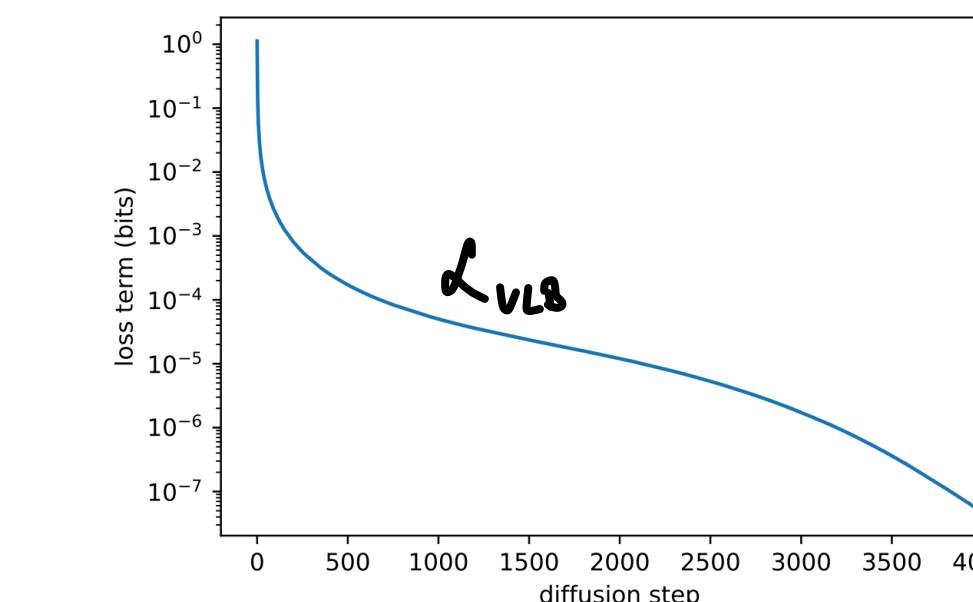
5.2 Improvements - Nichol et al

Improved noise schedule β_t



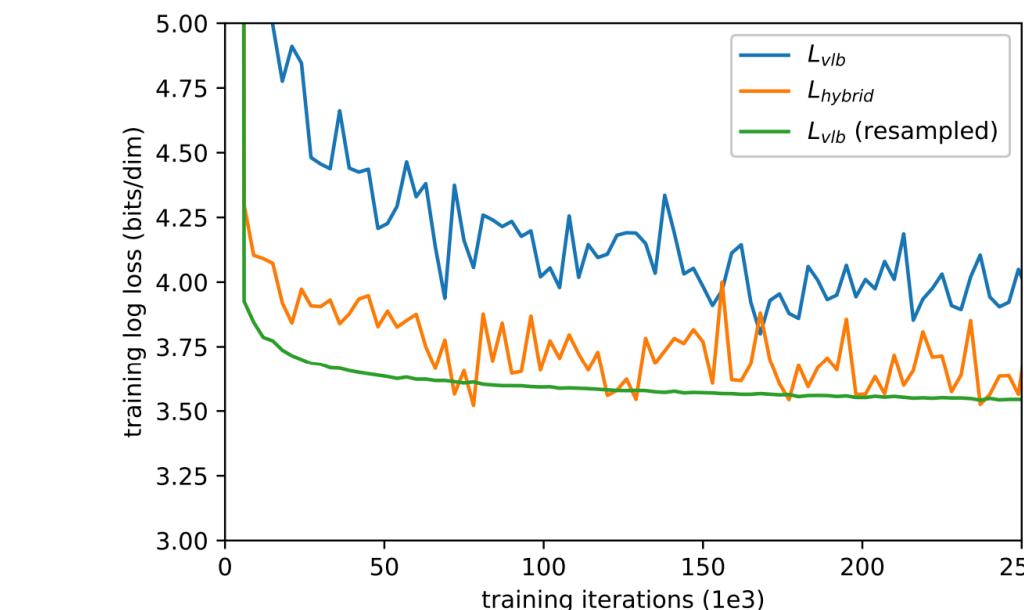
Linear
cosine

Reduced gradient noise by improved sampling



↑ spans multiple
orders of
magnitude

noisy loss landscape
⇒ when sampling uniformly



example $L_{20} \gg L_{500}$

→ rescale both so that they have comparable values but sample $t=20$ more often to compensate.

$$L_{VLB} = \left\langle \frac{\bar{\alpha}_t}{P_t} \right\rangle + \nu P_t$$

$$P_t \propto f(t^2)$$

and $\sum P_t = 1$



5.2 Improvements - Nichol et al

Nice side effect: Faster sampling speed

$S = \text{arbitrary subseries of } [1, 2, \dots, T]$

$$\beta_{S_t} := 1 - \frac{\bar{\alpha}_{S_t}}{\bar{\alpha}_{S_{t+1}}}$$

$$\tilde{\beta}_{S_t} := \frac{1 - \bar{\alpha}_{S_{t+1}}}{1 - \bar{\alpha}_{S_t}} \beta_{S_t}$$

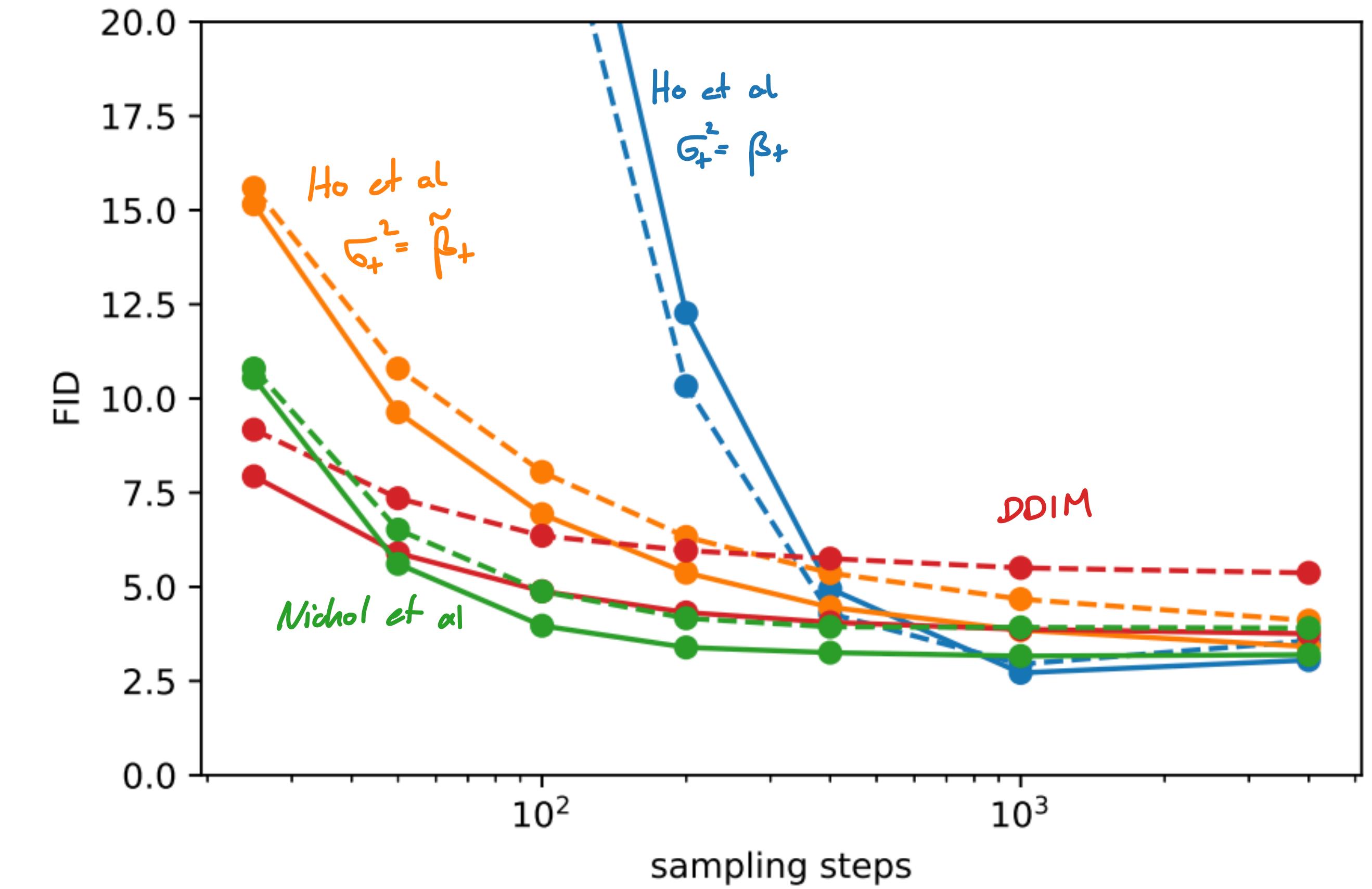
$$p(x_{S_{t+1}} | x_{S_t}) = N(x_{S_{t+1}}; \mu_\theta(x_{S_t}, S_t), \Sigma_\theta(x_{S_t}, S_t))$$

Idea:

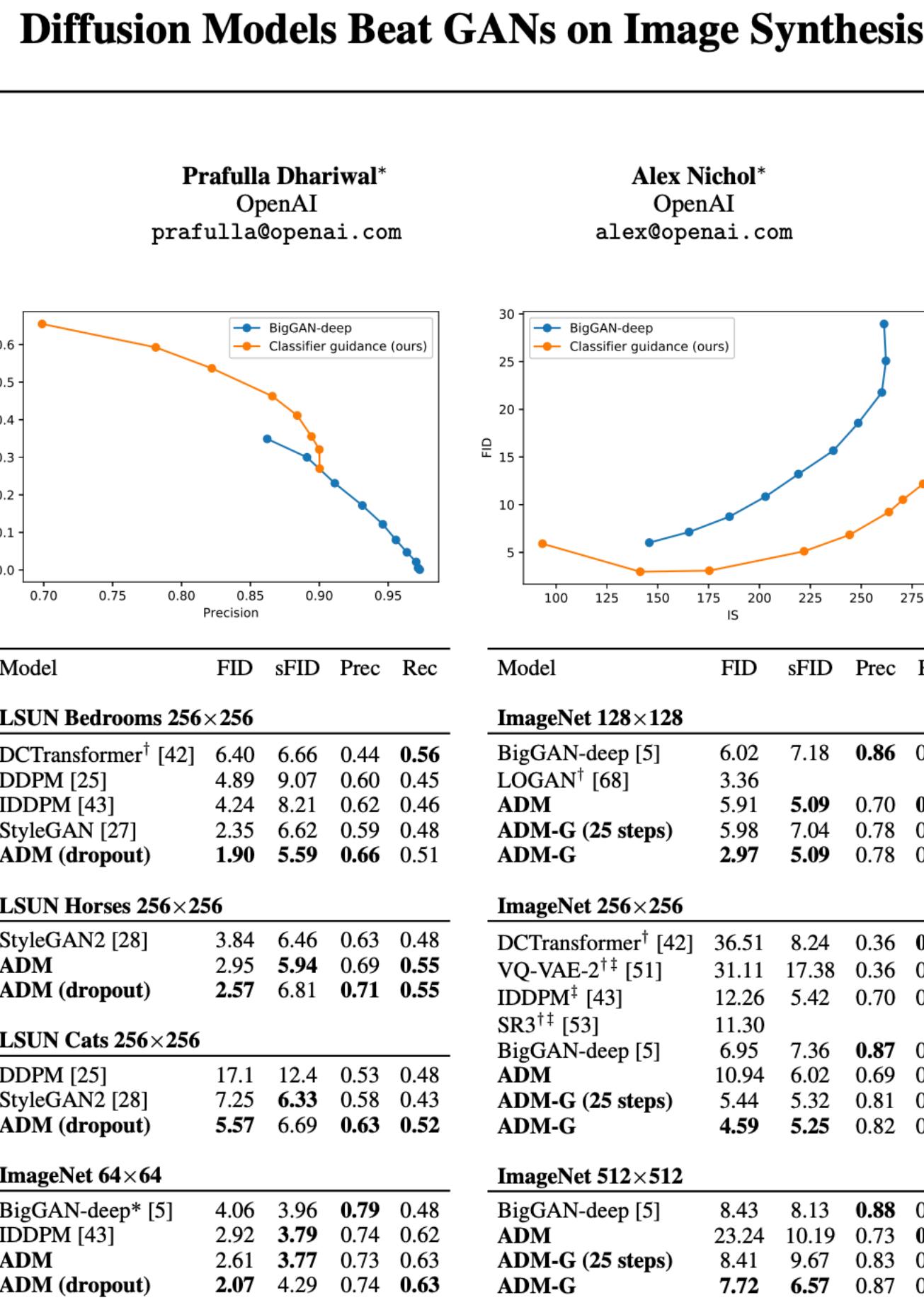
$$\beta_+ = 1 - \alpha_+$$

$$= 1 - \frac{\alpha_t \cdot \alpha_{t-1} \cdot \dots \cdot \alpha_0}{\alpha_{t-1} \cdot \alpha_{t-2} \cdot \dots \cdot \alpha_0}$$

$$= 1 - \frac{\alpha_+}{\alpha_{t-1}}$$



5.3 Improvements - Dhariwal et al



Main contributions:

1. Ablation studies on architecture design
2. Classifier guidance

Ablation studies on architecture design

- Started with UNet + global attention layer at 16x16 resolution + projection of timestamps embedded in each ResBlock.
- Ablation studies on depth vs width, number of attention heads, different resolutions of attention heads, using BigGAN Res blocks, etc...
- Conclusions:
 - More heads or few channels improve FID
 - 64 channels is best w.r.t wall clock time
 - Adaptive Gradient Normalisation Layer

Channels	Depth	Heads	Attention resolutions	BigGAN up/downsample	Rescale resblock	FID 700K	FID 1200K
160	2	1	16	✗	✗	15.33	13.21
128	4	4	32,16,8	✓		-0.21	-0.48
160	2	4	32,16,8	✓	✗	-0.54	-0.82
						-0.72	-0.66
						-1.20	-1.21
						0.16	0.25
						-3.14	-3.00

Number of heads	Channels per head	FID	Operation	FID
1		14.08	AdaGN	13.06
2		-0.50	Addition + GroupNorm	15.08
4		-0.97		
8		-1.17		
	32	-1.36		
	64	-1.03		
	128	-1.08		

In the rest of the paper, we use this final improved model architecture as our default:

- variable width with 2 residual blocks per resolution,
- multiple heads with 64 channels per head,
- attention at 32, 16 and 8 resolutions,
- BigGAN residual blocks for up and downsampling,
- and adaptive group normalisation for injecting timestep and class embeddings into residual blocks.



5.3 Improvements - Dhariwal et al

Classifier Guidance

Motivation: GANs make heavy use of class labels. So how can we embed class labels into diffusion models?

$$p_{\theta,t}(x_t|x_{t-1}, y) = \frac{p_{\theta}(x_t|x_{t-1}) p_{\phi}(y|x_t)}{Z} \quad \begin{matrix} \text{quite lengthy derivation} \\ \text{normalization constant} \\ \text{classifier} \end{matrix}$$

Recall:

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

$$\Rightarrow \log p_{\theta}(x_t|x_{t-1}) = -\frac{1}{2} (x_t - \mu)^T \Sigma^{-1} (x_t - \mu) + C$$

After a lengthy calculation:

$$\log(p_{\theta}(x_t|x_{t-1}) p_{\phi}(y|x_t)) = -\frac{1}{2} (x_t - \mu - \Sigma g)^T \Sigma^{-1} (x_t - \mu - \Sigma g) + C_2$$

$$\Rightarrow p_{\theta,t}(x_t|x_{t-1}, y) \sim \mathcal{N}(\mu + \Sigma g, \Sigma)$$

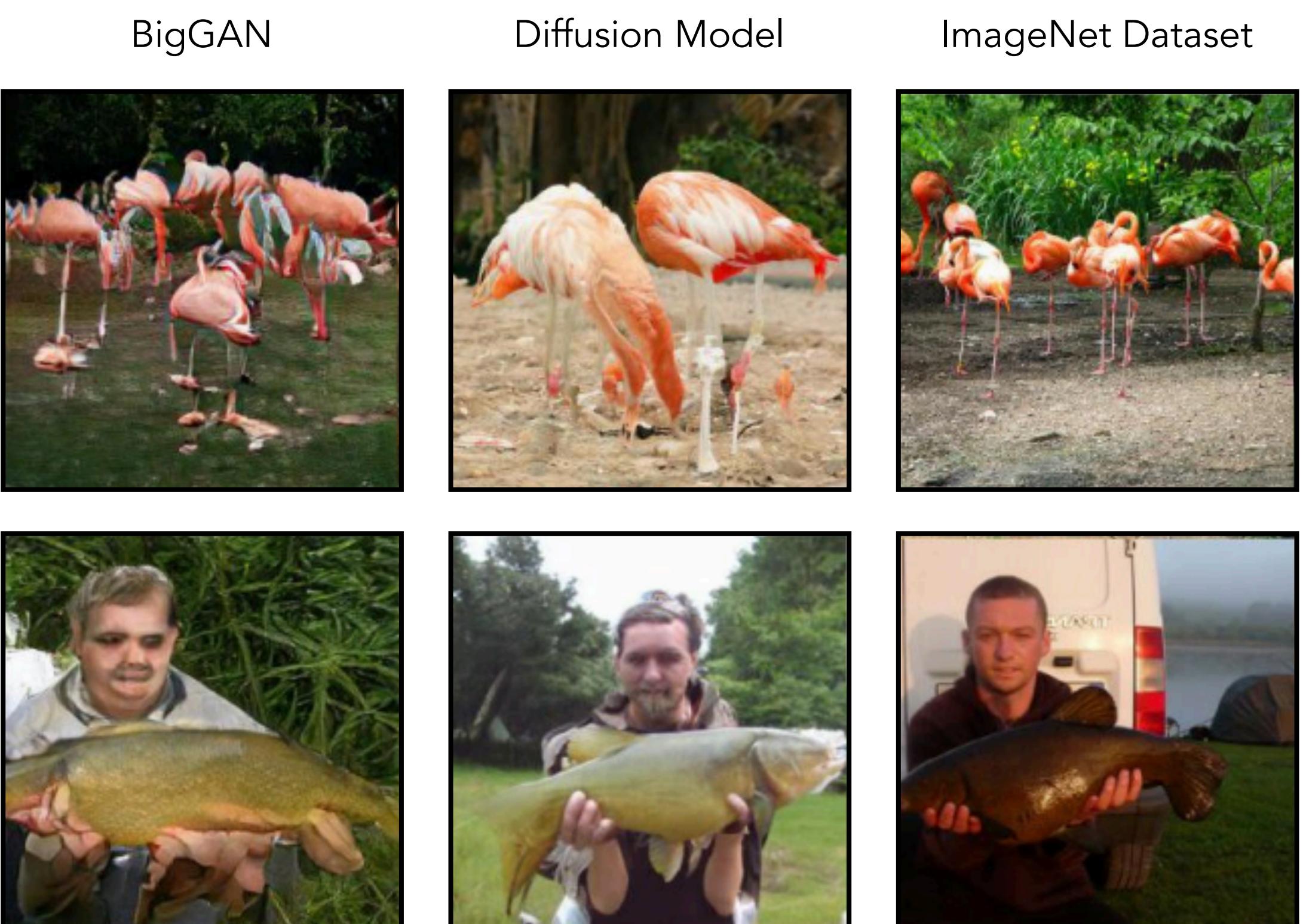
Algorithm 1 Classifier guided diffusion sampling, given a diffusion model $(\mu_{\theta}(x_t), \Sigma_{\theta}(x_t))$, classifier $p_{\phi}(y|x_t)$, and gradient scale s .

```

Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, I)$ 
for all  $t$  from  $T$  to 1 do
     $\mu, \Sigma \leftarrow \mu_{\theta}(x_t), \Sigma_{\theta}(x_t)$ 
     $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s \Sigma \nabla_{x_t} \log p_{\phi}(y|x_t), \Sigma)$ 
end for
return  $x_0$ 

```

Conditional	Guidance	Scale	FID	sFID	IS	Precision	Recall
✗	✗		26.21	6.35	39.70	0.61	0.63
✗	✓	1.0	33.03	6.99	32.92	0.56	0.65
✗	✓	10.0	12.00	10.40	95.41	0.76	0.44
✓	✗		10.94	6.02	100.98	0.69	0.63
✓	✓	1.0	4.59	5.25	186.70	0.82	0.52
✓	✓	10.0	9.11	10.93	283.92	0.88	0.32



Caution: Only works for stochastic models (e.g., DDPM). For DDIM there is an alternative solution. See the paper.

5.4 Improvements - Ho et al (Classifier Free Guidance)

Classifier-Free Diffusion Guidance

Jonathan Ho
Google Research
jonathanho@google.com

Tim Salimans
Google Research
salimans@google.com

Instead of training a separate classifier model, we choose to train an unconditional denoising diffusion model $p_\theta(\mathbf{z})$ parameterized through a score estimator $\epsilon_\theta(\mathbf{z}_\lambda)$ together with the conditional model $p_\theta(\mathbf{z}|\mathbf{c})$ parameterized through $\epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c})$. We use a single neural network to parameterize both models, where for the unconditional model we can simply input zeros for the class identifier \mathbf{c} when predicting the score, i.e. $\epsilon_\theta(\mathbf{z}_\lambda) = \epsilon_\theta(\mathbf{z}_\lambda, \mathbf{c} = 0)$. We jointly train the unconditional and conditional models simply by randomly setting \mathbf{c} to the unconditional class identifier.

Method	FID (\downarrow)	IS (\uparrow)
ADM [3]	2.07	-
CDM [6]	1.48	67.95
Ours, no guidance	1.80	53.71
Ours, with guidance		
$w = 0.1$	1.55	66.11
$w = 0.2$	2.04	78.91
$w = 0.3$	3.03	92.8
$w = 0.4$	4.30	106.2
$w = 0.5$	5.74	119.3
$w = 0.6$	7.19	131.1
$w = 0.7$	8.62	141.8
$w = 0.8$	10.08	151.6
$w = 0.9$	11.41	161
$w = 1.0$	12.6	170.1
$w = 2.0$	21.03	225.5
$w = 3.0$	24.83	250.4
$w = 4.0$	26.22	260.2

Figure 1: ImageNet 64x64 results

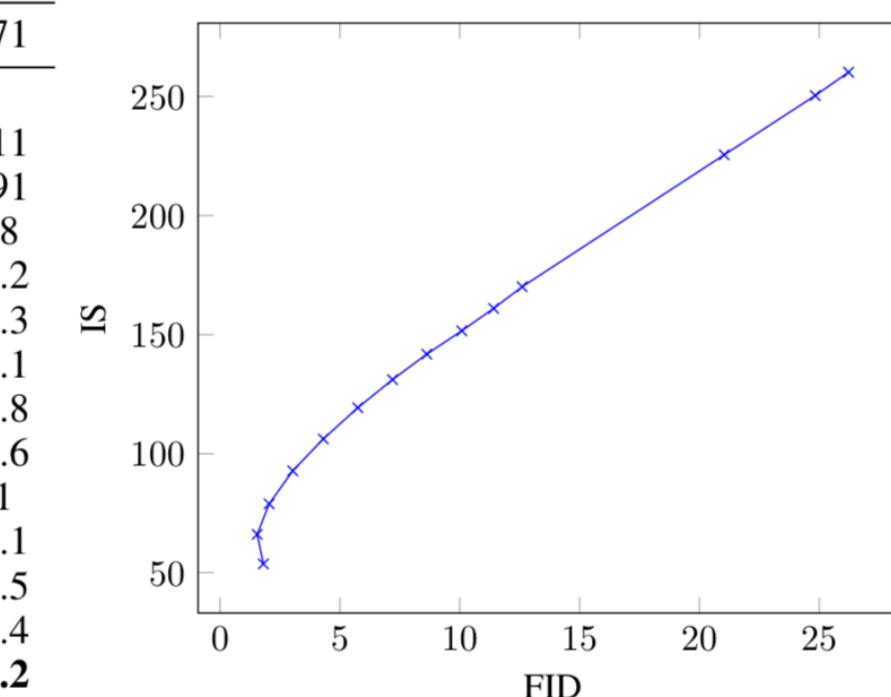
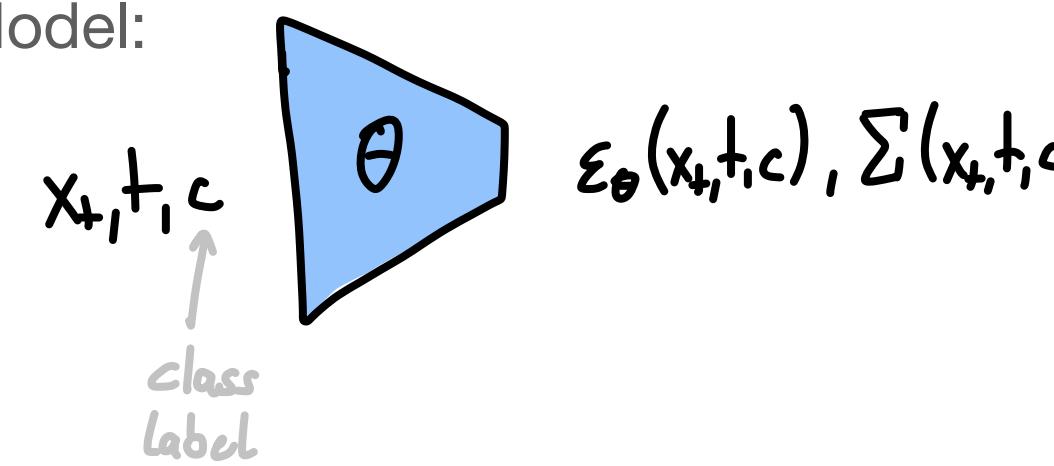
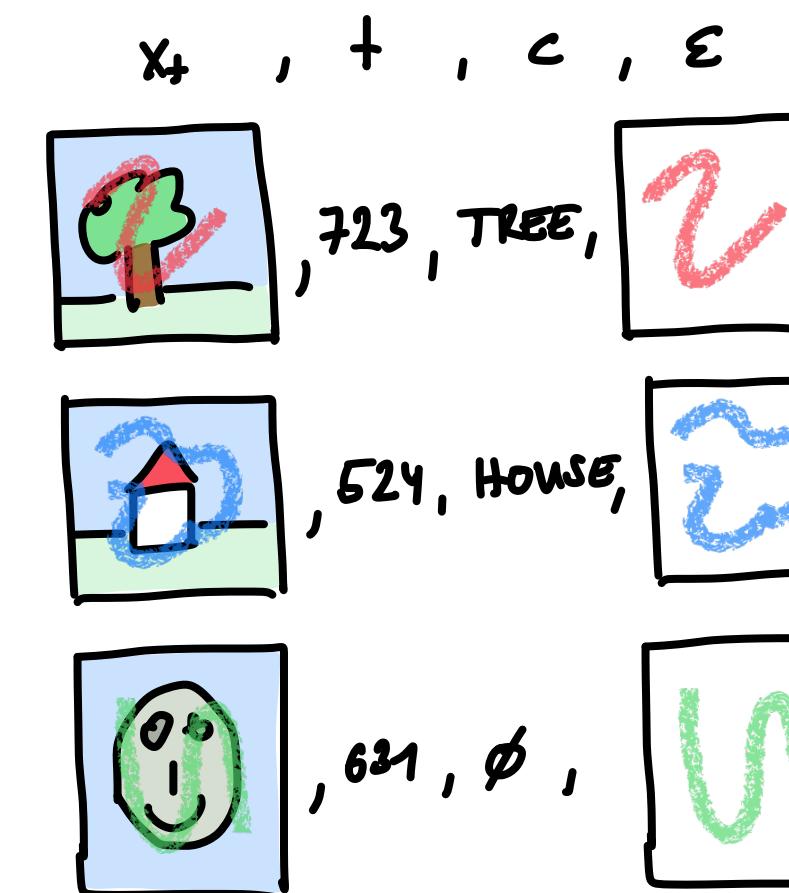


Figure 2: ImageNet 64x64 FID vs. IS

Model:

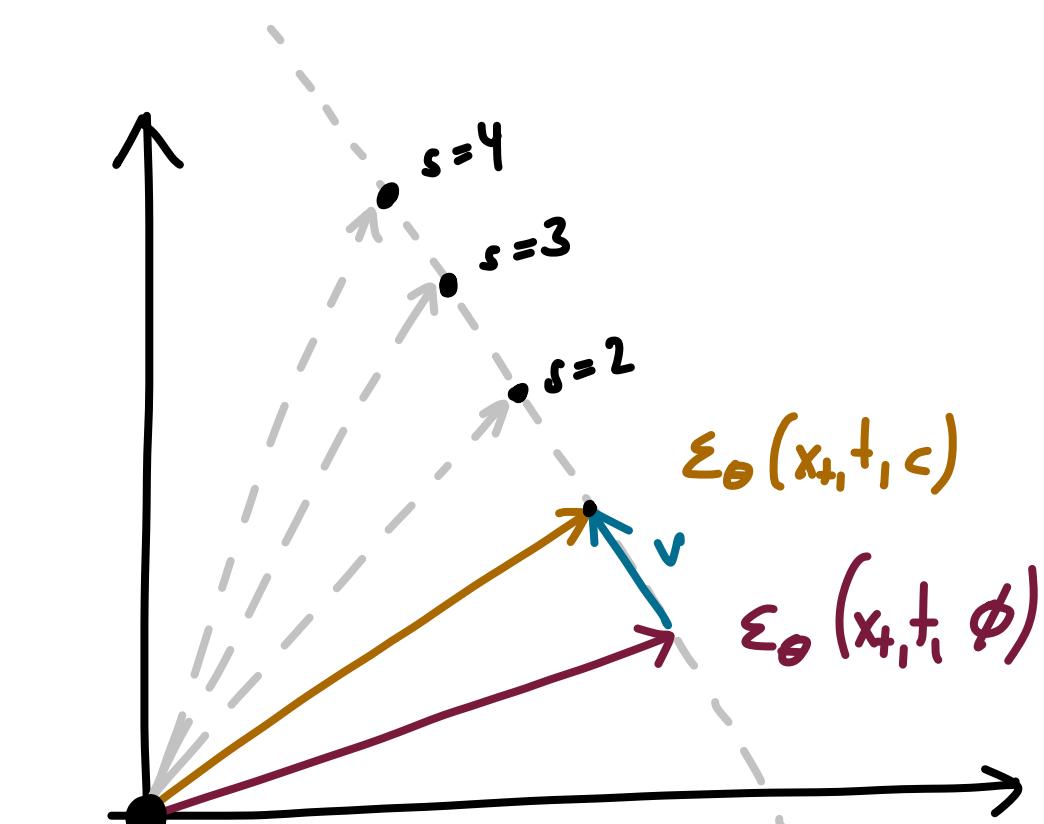


Training samples:

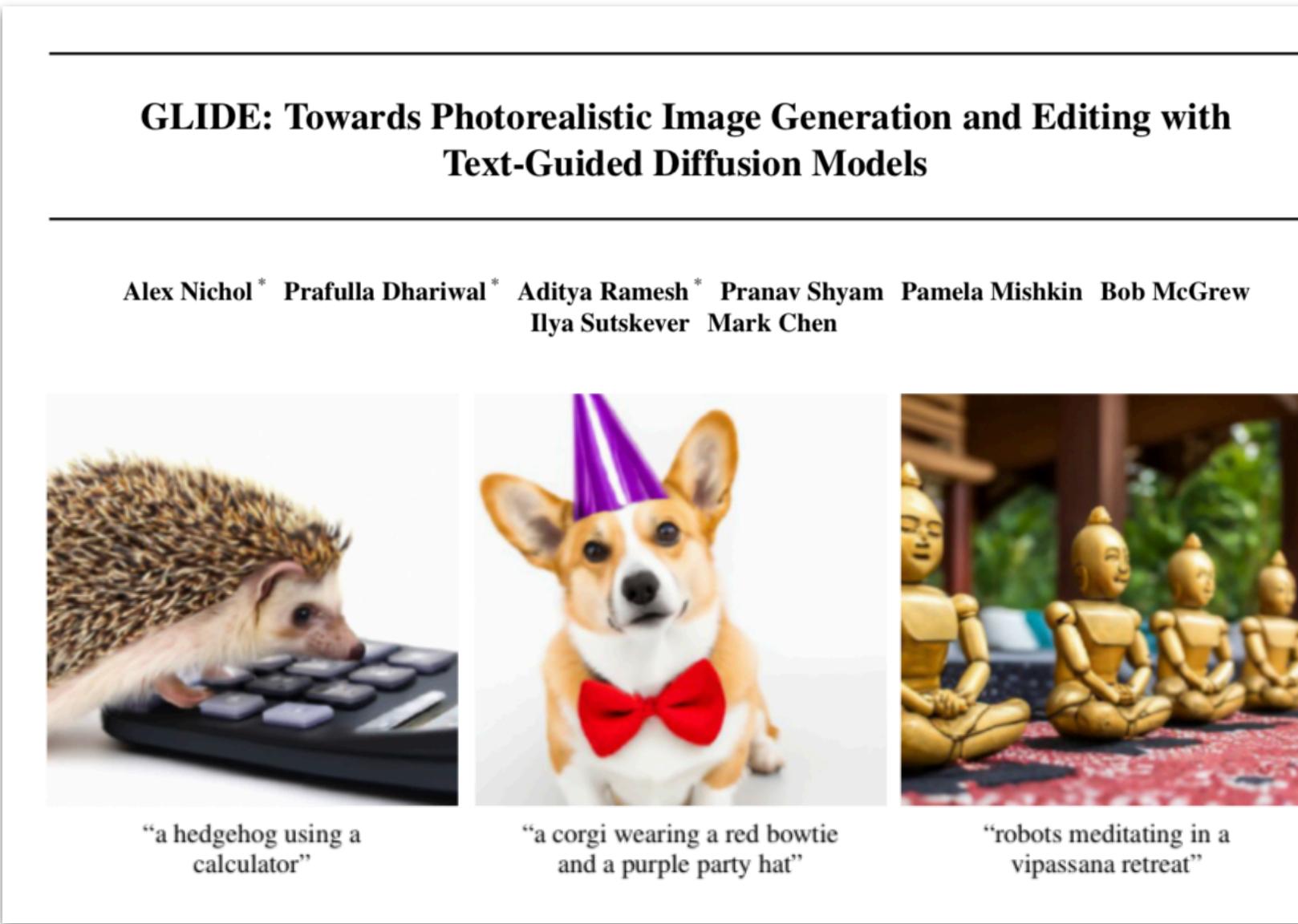


Estimated noise:

$$\begin{aligned}\tilde{\epsilon}_\theta(x_t, t, c) &= (1+w)\epsilon_\theta(x_t, t, c) - w\epsilon_\theta(x_t, t, \emptyset) \\ w &:= s-1 \\ &= \epsilon_\theta(x_t, t, \emptyset) + s[\underbrace{\epsilon_\theta(x_t, t, c) - \epsilon_\theta(x_t, t, \emptyset)}_{=: v}]\end{aligned}$$



5.4 Improvements - Nichol et al (GLIDE)



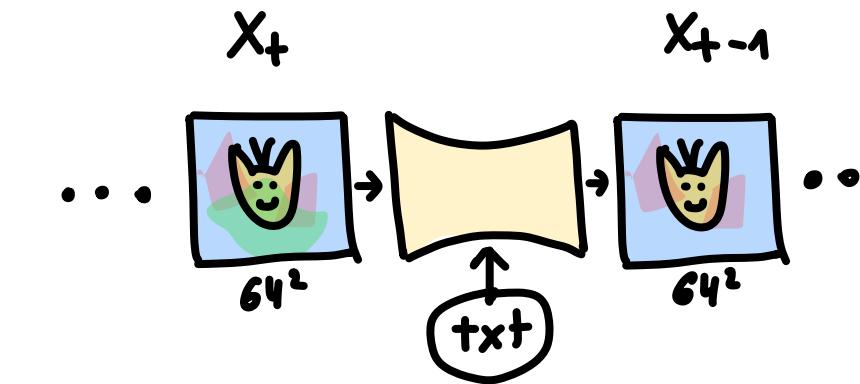
Main contributions:

1. Text conditional guidance
2. Iterative generation by (Improved) inpainting

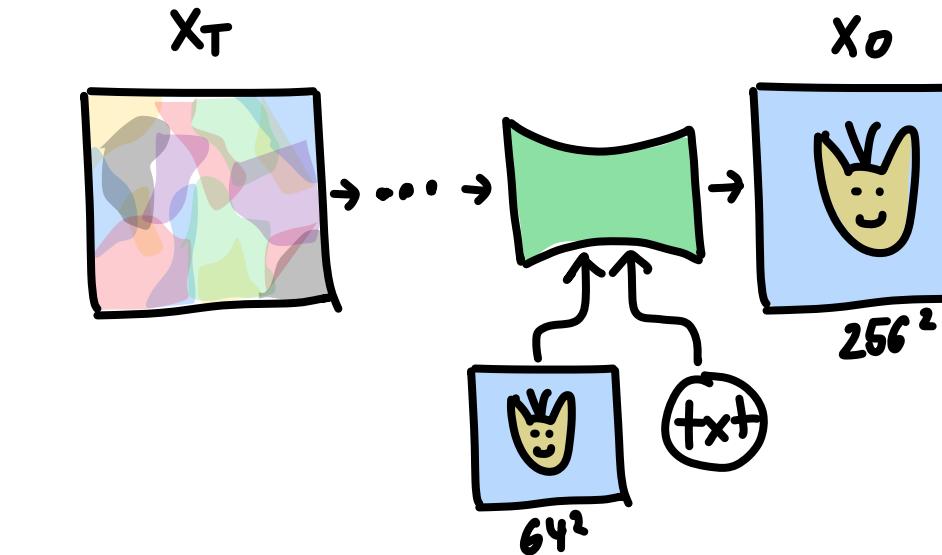


Architecture:

① text-conditional diffusion model



② text-conditional upscaling diffusion model (actually from Nichol et al.)



③ text-encoder

We adopt the ADM model architecture proposed by Dhariwal & Nichol (2021), but augment it with text conditioning information. For each noised image x_t and corresponding text caption c , our model predicts $p(x_{t-1}|x_t, c)$. To condition on the text, we first encode it into a sequence of K tokens, and feed these tokens into a Transformer model (Vaswani et al., 2017). The output of this transformer is used in two ways: first, the final token embedding is used in place of a class embedding in the ADM model; second, the last layer of token embeddings (a sequence of K feature vectors) is separately projected to the dimensionality of each attention layer throughout the ADM model, and then concatenated to the attention context at each layer.

Training:

- Pre-training on DALL-E dataset.
- Fine tuning for classifier-free guidance by randomly dropping 20% of the image descriptions.
- Authors also tried CLIP guidance but classifier-free guidance yields better results.

Improved Denoising Diffusion Probabilistic Models

C. Sample Quality on ImageNet 256 × 256

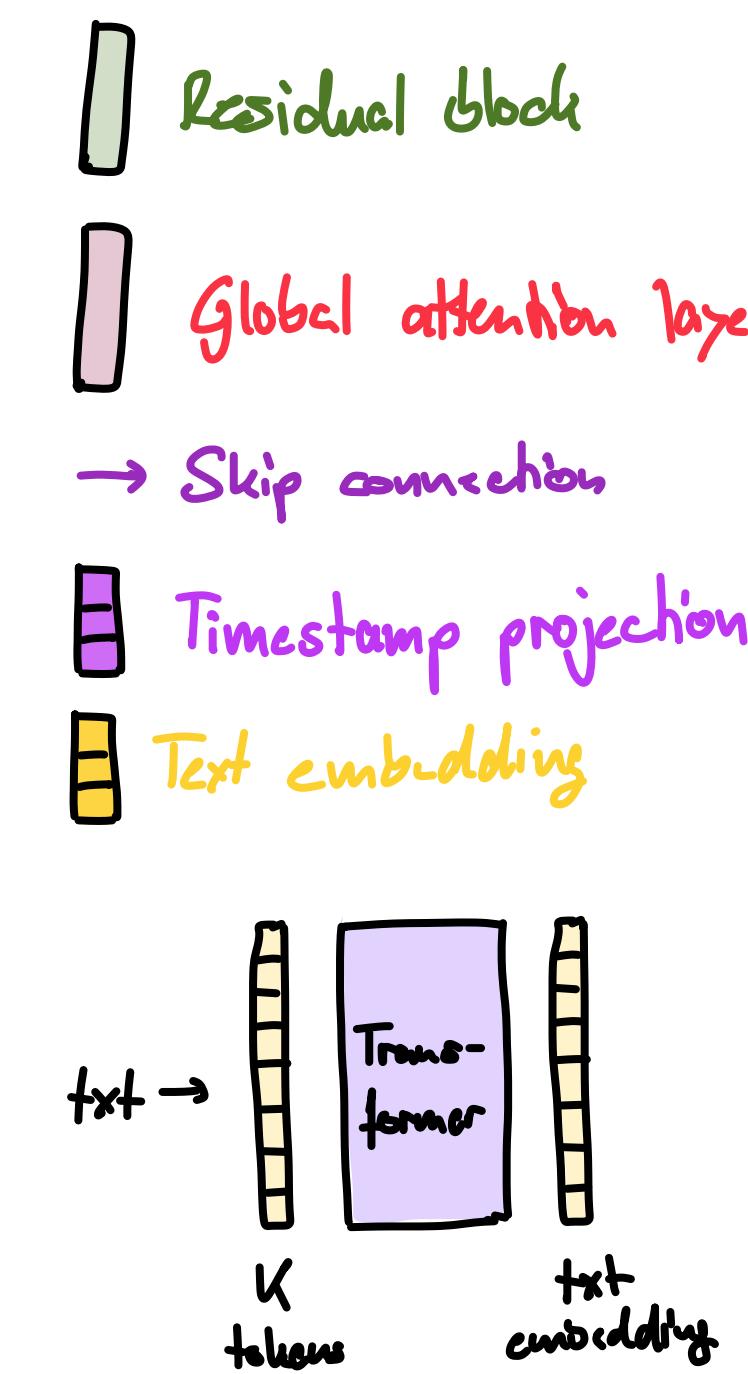
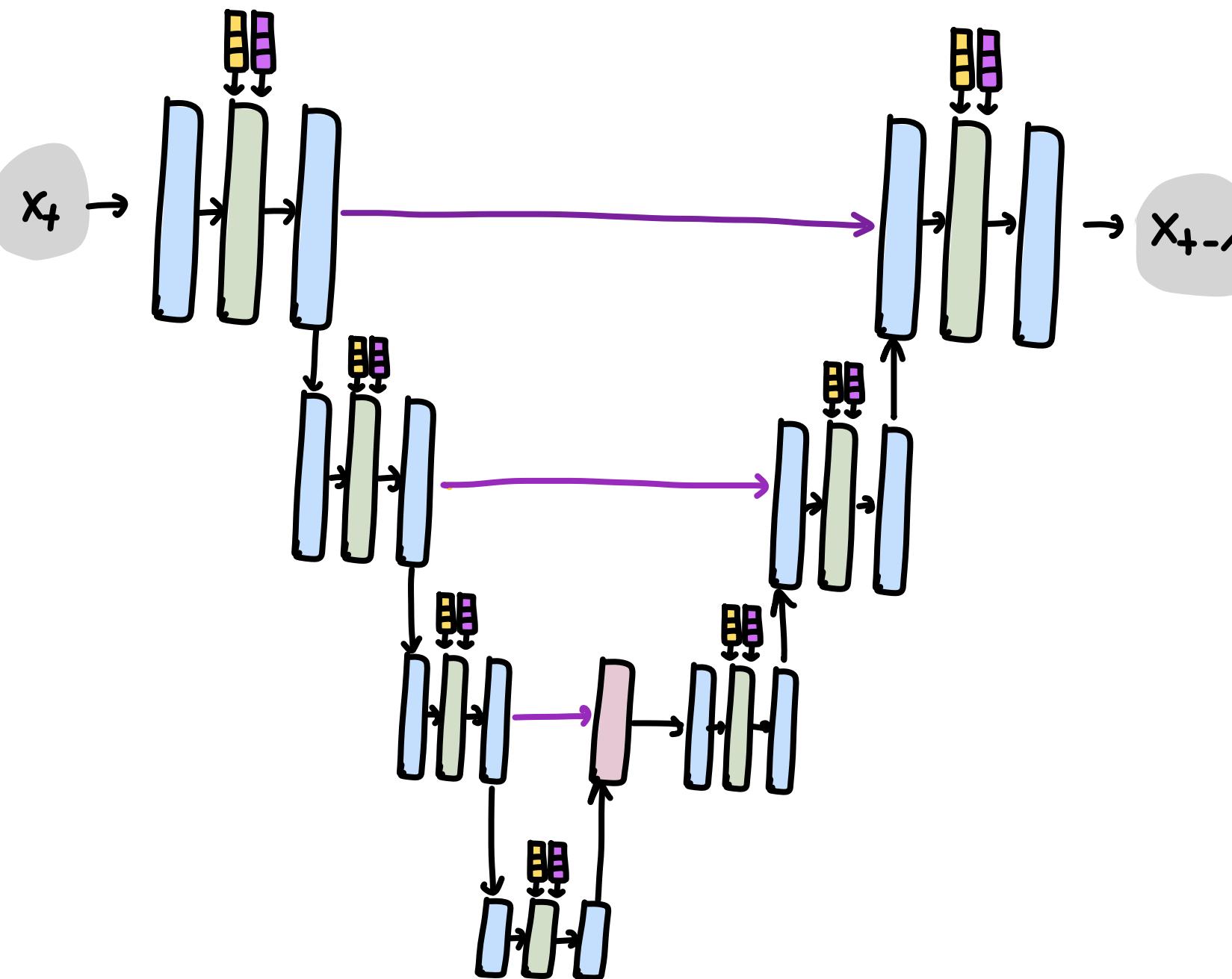
We trained two models on class conditional ImageNet 256 × 256. The first is a usual diffusion model that directly models the 256 × 256 images. The second model reduces compute by chaining a pretrained 64 × 64 model $p(x_{64}|y)$ with another upsampling diffusion model $p(x_{256}|x_{64}, y)$ to upsample images to 256 × 256. For the upsampling model, the downsampled image x_{64} is passed as extra conditioning input to the UNet. This is similar to VQ-VAE-2 (Razavi et al., 2019), which uses two stages of priors at different latent resolutions to more efficiently learn global and local features. The linear schedule worked better for 256 × 256 images, so we used that for these results. Table 5 summarizes our results. For VQ-VAE-2, we use the FIDs reported in (Ravuri & Vinyals, 2019). Diffusion models still obtain the best FIDs for a likelihood-based model, and close the gap to GANs considerably.

MODEL	FID
VQ-VAE-2 ((Razavi et al., 2019), two-stage)	38.1
Improved Diffusion (ours, single-stage)	31.5
Improved Diffusion (ours, two-stage)	12.3
BigGAN (Brock et al., 2018)	7.7
BigGAN-deep (Brock et al., 2018)	7.0

5.4 Improvements - Nichol et al (GLIDE)

Rough schema of the architecture (however, not 100% sure):

ADM Architecture with text encoding



How to find the architecture....

Nichol et al.: "We use the same model architecture as the ImageNet 64×64 model from Dhariwal & Nichol (2021), but scale the model width to 512 channels"

Dhariwal et al.: "Ho et al. [25] introduced the UNet architecture for diffusion models"

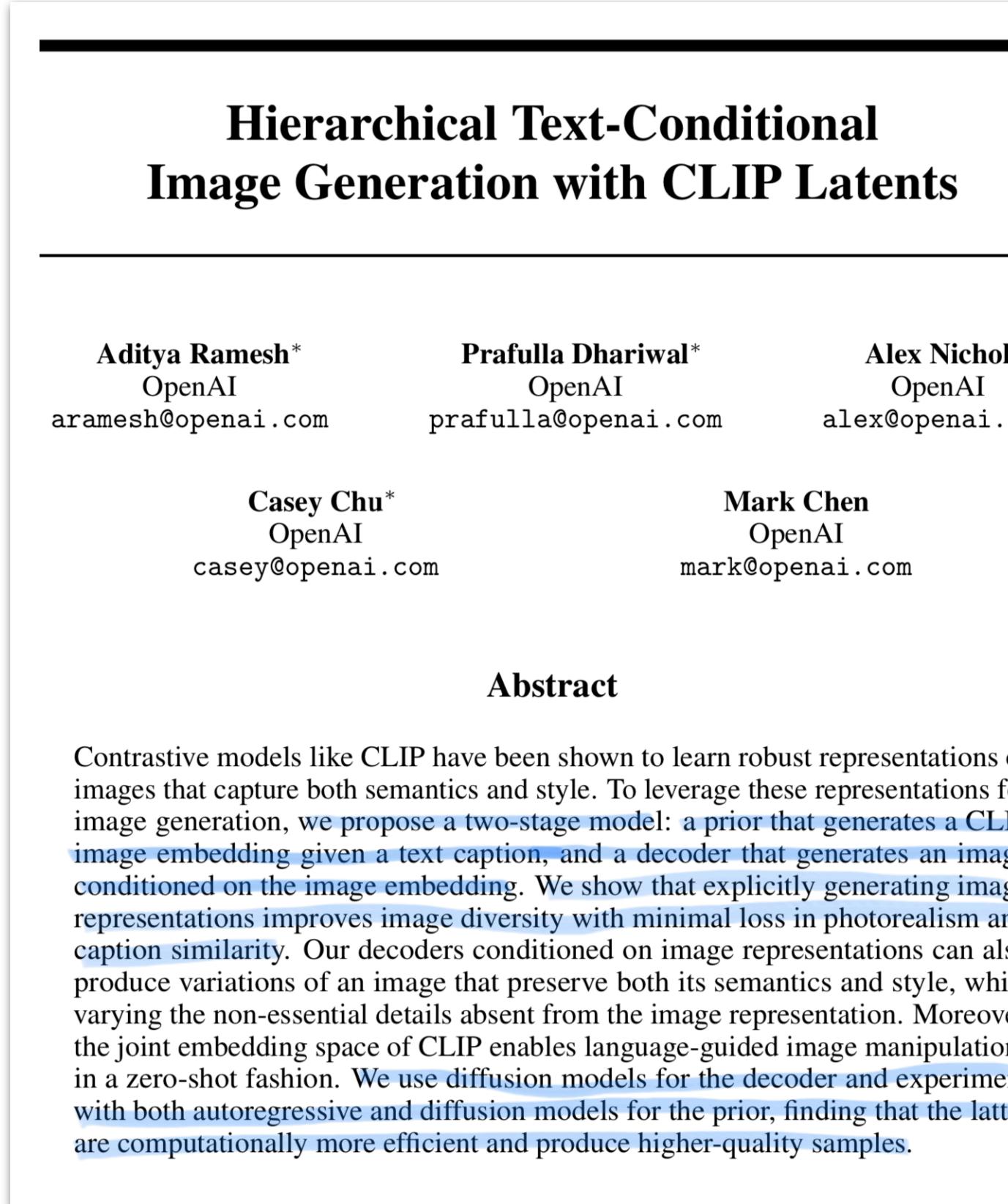
Ho et al.: "...we use a U-Net backbone similar to an unmasked PixelCNN++..."

Pixel CNN++: "PixelCNN++: Improving the PixelCNN..."



"Oil painting of a frustrated person sitting in his office, throwing around scientific articles."

5.5 Improvements - Ramesh et al (DALL-E 2)



Main contributions:

1. Leveraging powerful CLIP embeddings

Main idea:

- “CLIP embeddings are robust to image distribution shift, have impressive zero-shot capabilities, and have been fine-tuned to achieve state-of-the-art results on a wide variety of vision and language tasks.”
- “Diffusion models have emerged as a promising generative modeling framework, pushing the state-of-the-art on image and video generation tasks”

“CLIP + diffusion model = DALL-E 2”

Two staged model:

- Prior $P(z_i|y)$
- Decoder $P(x|z_i, y)$

$$\Rightarrow P(x|y) = P(x, z_i|y) = P(x|z_i, y) \cdot P(z_i|y)$$

$z_i = f(x)$
deterministic

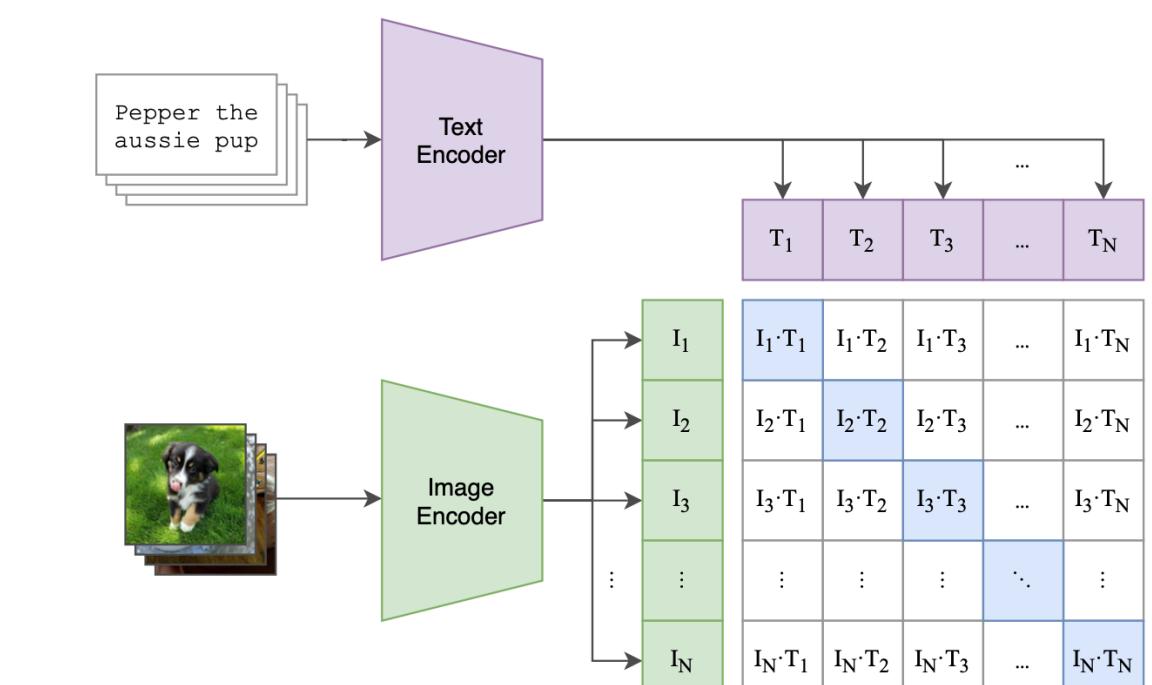
y: caption

z_i : image encoding

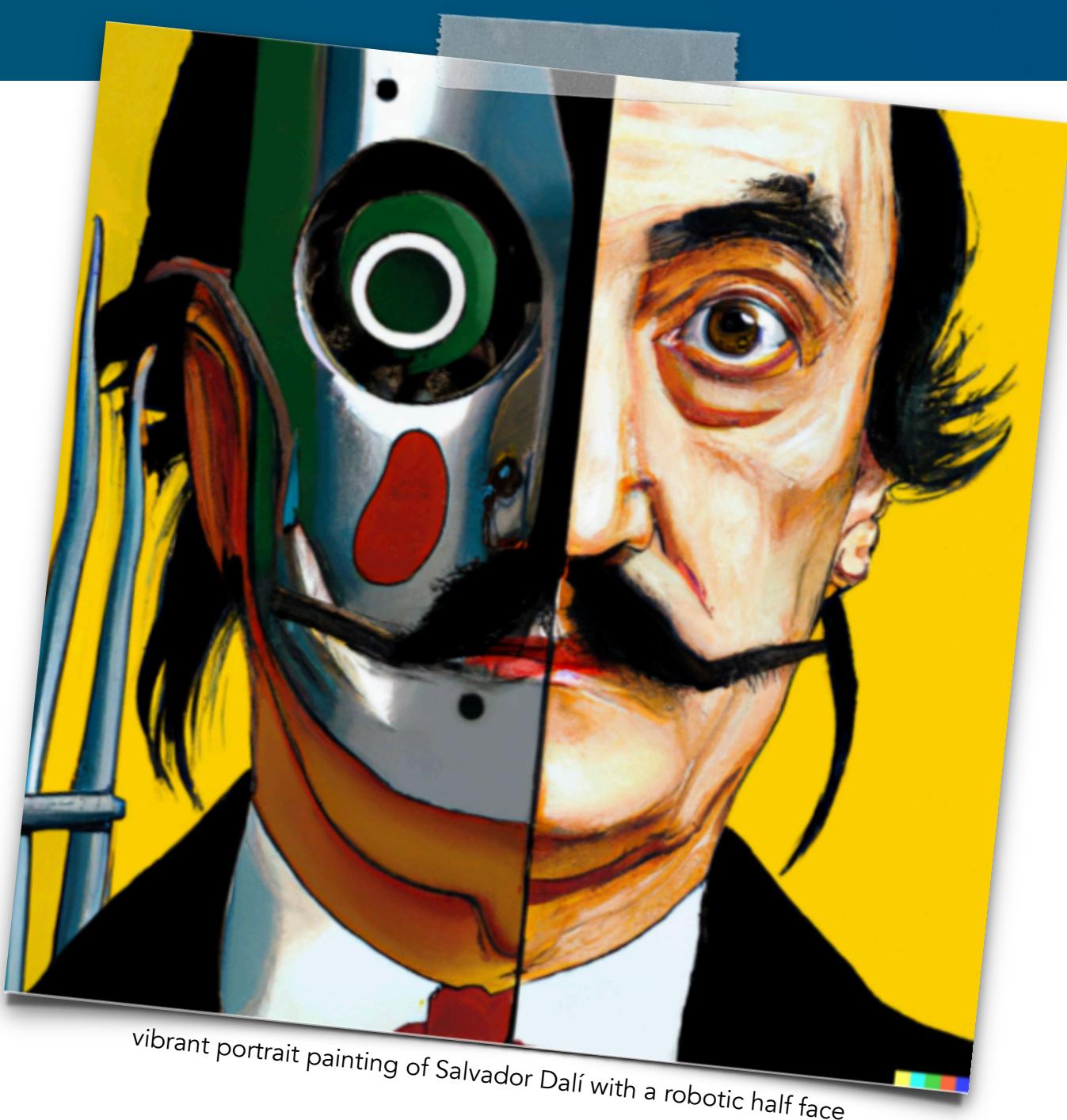
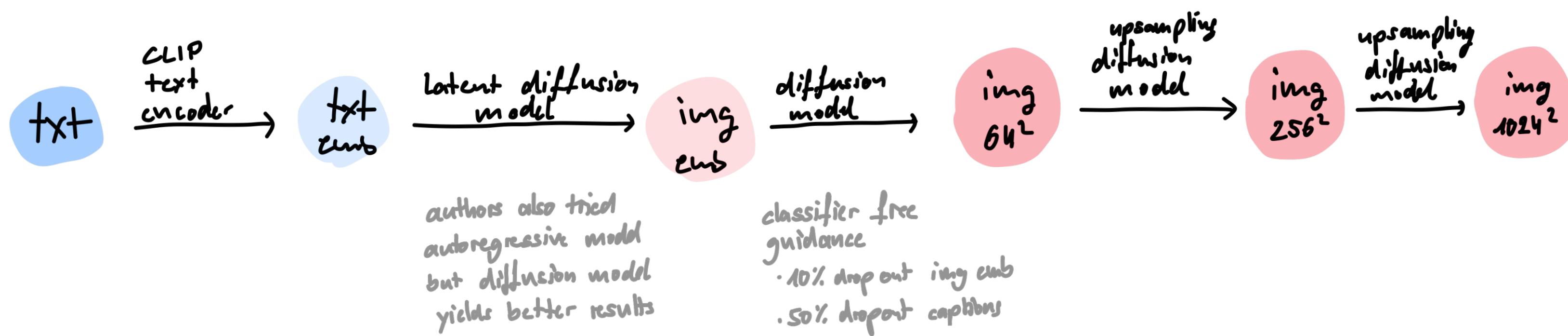
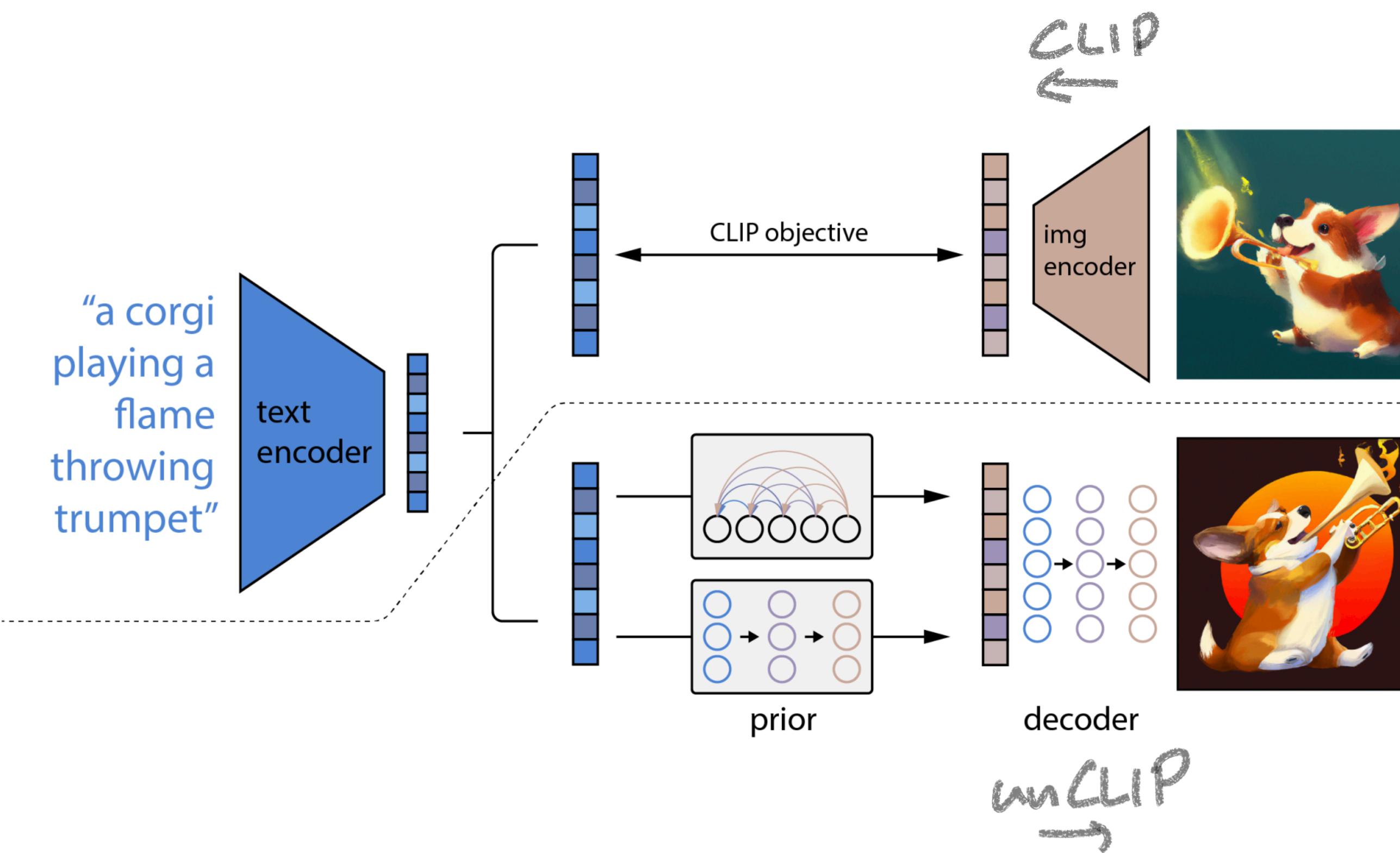
x: image

Definition
of cond. prob

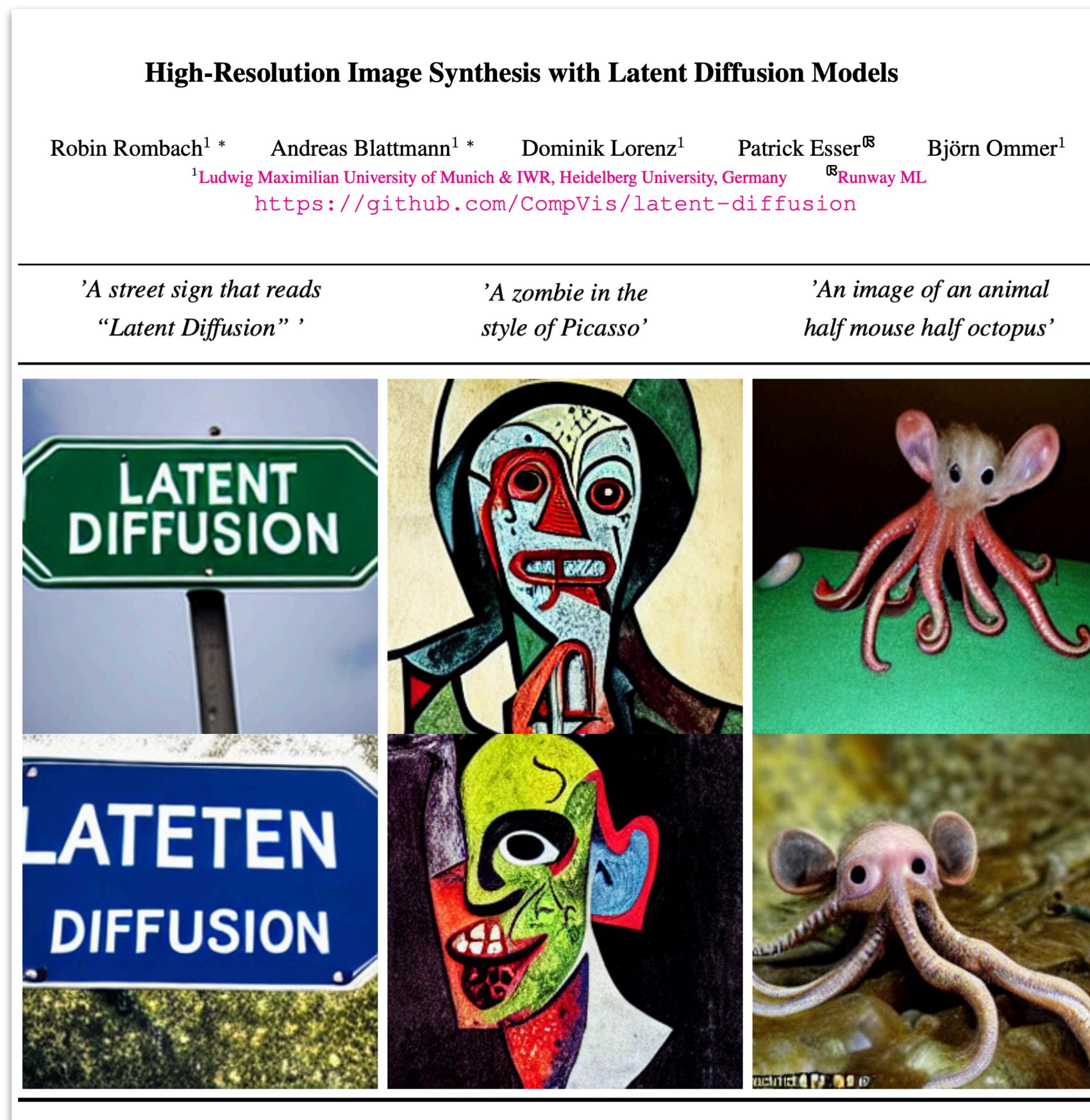
Excerpt from CLIP paper:



5.5 Improvements - Ramesh et al (DALL-E 2)

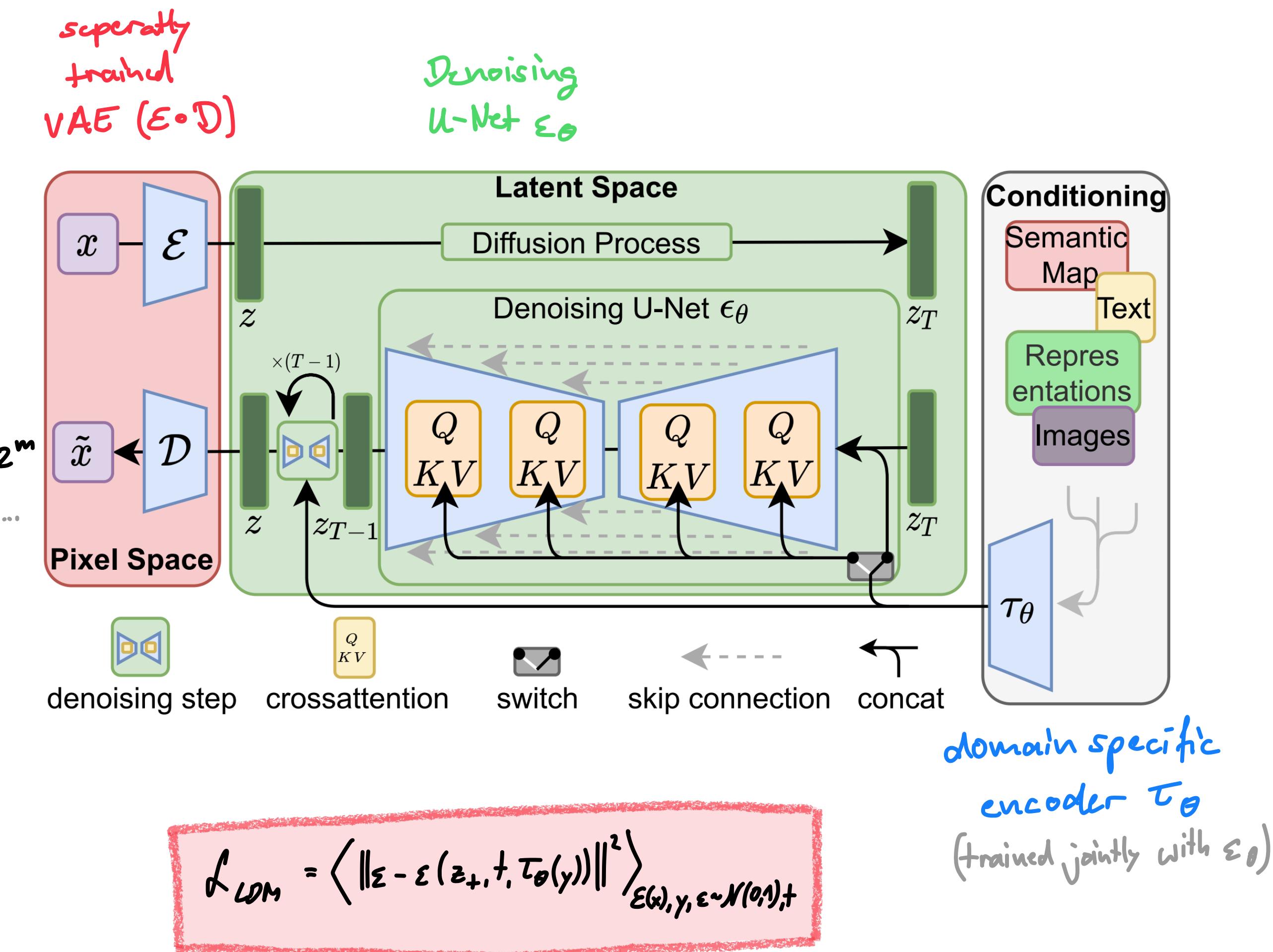


5.6 Improvements - Rombach et al (Stable Diffusion)



Main contributions:

1. Reduction of comp. requirements through latent diffusion.
2. Cross-attention layer for general purpose (multi model) conditioning.



5.6 Improvements - Rombach et al (Stable Diffusion)

About the conditioning mechanism...

Domain specific encoder:

$$\mathcal{T}_\theta: y \mapsto \mathcal{T}_\theta(y) \in \mathbb{R}^{M \times d_\tau}$$

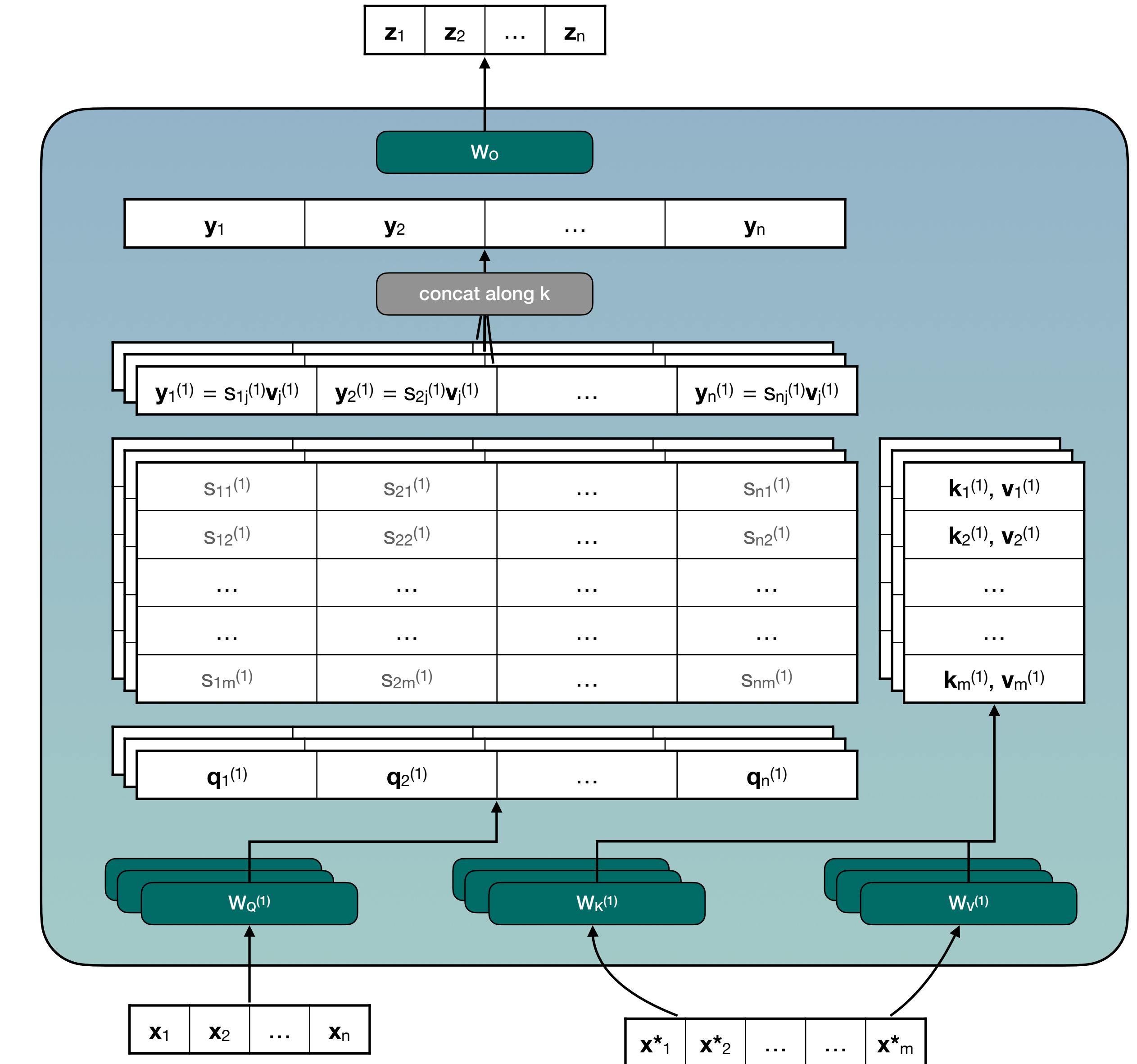
Cross-attention mechanism:

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

$$Q = W_Q^{(1)} \cdot \varphi_i(z_i) \quad K = W_K^{(1)} \cdot \mathcal{T}_\theta(y) \quad V = W_V^{(1)} \cdot \mathcal{T}_\theta(y)$$

$\uparrow d \times d_c$ $\uparrow d \times d_\tau$ $\uparrow d \times d_\tau$

$\varphi_i(z_i) \in \mathbb{R}^{N \times d_c}$ flattened intermediate representation of the UNet.



flattened rep UNet; $N \times d\phiii$

domain specific decoder; $M \times dtau$

6. UNet Architecture in Depth

...example implementation from Stable Diffusion repo...

```
● unet.time_embed
  ✓ 0.0s
    Python

layer_names = []

for i, sequential_module in enumerate(unet.input_blocks):
    print(i, sequential_module._get_name())
    for j, layer in enumerate(sequential_module):
        print(f" {i}.{j}", layer._get_name())
        layer_names.append(layer._get_name())

    print()
    for i, module in enumerate(unet.middle_block):
        layer_name = layer._get_name()
        print(f" {i}.{j}", layer_name)
        layer_names.append(layer_name)

    print()
    for i, sequential_module in enumerate(unet.output_blocks):
        print(i, sequential_module._get_name())
        for j, layer in enumerate(sequential_module):
            print(f" {i}.{j}", layer._get_name())
            layer_names.append(layer._get_name())

    print(set(layer_names))
  ✓ 0.0s
    Python
```

```
Sequential(
  (0): Linear(in_features=320, out_features=1280, bias=True)
  (1): SiLU()
  (2): Linear(in_features=1280, out_features=1280, bias=True)
)

0 TimestepEmbedSequential
  0.0 Conv2d
  1 TimestepEmbedSequential
    1.0 ResBlock
    1.1 SpatialTransformer
  2 TimestepEmbedSequential
    2.0 ResBlock
    2.1 SpatialTransformer
  3 TimestepEmbedSequential
    3.0 Downsample
  4 TimestepEmbedSequential
    4.0 ResBlock
    4.1 SpatialTransformer
  5 TimestepEmbedSequential
    5.0 ResBlock
    5.1 SpatialTransformer
  6 TimestepEmbedSequential
    6.0 Downsample
  7 TimestepEmbedSequential
    7.0 ResBlock
    7.1 SpatialTransformer
  8 TimestepEmbedSequential
    8.0 ResBlock
    8.1 SpatialTransformer
  9 TimestepEmbedSequential
    9.0 Downsample
 10 TimestepEmbedSequential
    10.0 ResBlock
 11 TimestepEmbedSequential
    11.0 ResBlock

 0.0 ResBlock
 1.0 ResBlock
 2.0 ResBlock
 3.0 ResBlock
 4.0 ResBlock
 5.0 ResBlock
 6.0 ResBlock
 7.0 ResBlock
 8.0 ResBlock
 9.0 ResBlock
10.0 ResBlock
11.0 ResBlock

 0 TimestepEmbedSequential
 1 TimestepEmbedSequential
 2 TimestepEmbedSequential
 3 TimestepEmbedSequential
 4 TimestepEmbedSequential
 5 TimestepEmbedSequential
 6 TimestepEmbedSequential
 7 TimestepEmbedSequential
 8 TimestepEmbedSequential
 9 TimestepEmbedSequential
10 TimestepEmbedSequential
11 TimestepEmbedSequential

 0.0 ResBlock
 1.0 ResBlock
 2.0 ResBlock
 3.0 ResBlock
 4.0 ResBlock
 5.0 ResBlock
 6.0 ResBlock
 7.0 ResBlock
 8.0 ResBlock
 9.0 ResBlock
10.0 ResBlock
11.0 ResBlock

 0.0 ResBlock
 1.0 ResBlock
 2.0 ResBlock
 3.0 ResBlock
 4.0 ResBlock
 5.0 ResBlock
 6.0 ResBlock
 7.0 ResBlock
 8.0 ResBlock
 9.0 ResBlock
10.0 ResBlock
11.0 ResBlock

 0.0 ResBlock
 1.0 ResBlock
 2.0 ResBlock
 3.0 ResBlock
 4.0 ResBlock
 5.0 ResBlock
 6.0 ResBlock
 7.0 ResBlock
 8.0 ResBlock
 9.0 ResBlock
10.0 ResBlock
11.0 ResBlock
```

TimestepEmbedSequential:

Simply a sequential module that passes timestamp embeddings to the layers that support it.

```
72
73     Robin Rombach, 3 months ago | 1 author (Robin Rombach)
74     class TimestepEmbedSequential(nn.Sequential, TimestepBlock):
75         """
76             A sequential module that passes timestep embeddings to the children that
77             support it as an extra input.
78         """
79
80         def forward(self, x, emb, context=None):
81             for layer in self:
82                 if isinstance(layer, TimestepBlock):
83                     x = layer(x, emb)
84                 elif isinstance(layer, SpatialTransformer):
85                     x = layer(x, context)
86                 else:
87                     x = layer(x)
88
89     return x
```

ldm/modules/diffusionmodules/openaimodels.py

The forward pass through the network:

```
def forward(self, x, timesteps=None, context=None, y=None, **kwargs):
    """
    Apply the model to an input batch.
    :param x: an [N x C x ...] Tensor of inputs.
    :param timesteps: a 1-D batch of timesteps.
    :param context: conditioning plugged in via crossattn
    :param y: an [N] Tensor of labels, if class-conditional.
    :return: an [N x C x ...] Tensor of outputs.
    """

    assert (y is not None) == (
        self.num_classes is not None
    ), "must specify y if and only if the model is class-conditional"
    hs = []
    t_emb = timestep_embedding(timesteps, self.model_channels, repeat_only=False)
    emb = self.time_embed(t_emb)

    if self.num_classes is not None:
        assert y.shape[0] == x.shape[0]
        emb = emb + self.label_emb(y)

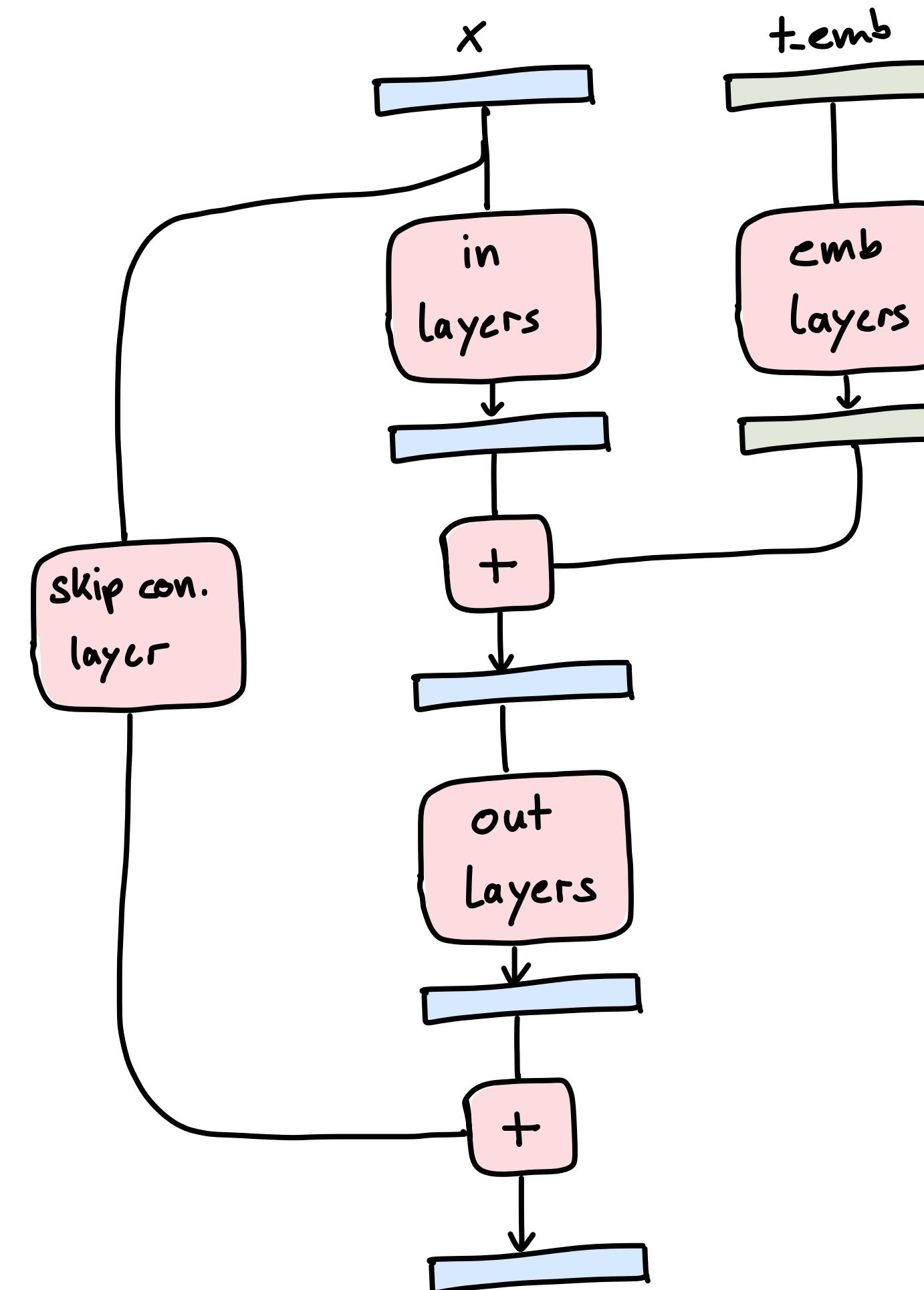
    h = x.type(self.dtype)
    for module in self.input_blocks:
        h = module(h, emb, context)
        hs.append(h)
    h = self.middle_block(h, emb, context)
    for module in self.output_blocks:
        h = th.cat([h, hs.pop()], dim=1)
        h = module(h, emb, context)
    h = h.type(x.dtype)
    if self.predict_codebook_ids:
        return self.id_predictor(h)
    else:
        return self.out(h)
```

6. UNet Architecture in Depth

ResBlock

```
unet.input_blocks[1][0] •
✓ 0.0s
Python
ResBlock(
    (in_layers): Sequential(
        (0): GroupNorm32(32, 320, eps=1e-05, affine=True)
        (1): SiLU()
        (2): Conv2d(320, 320, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    )
    (h_upd): Identity()
    (x_upd): Identity()
    (emb_layers): Sequential(
        (0): SiLU()
        (1): Linear(in_features=1280, out_features=320, bias=True)
    )
    (out_layers): Sequential(
        (0): GroupNorm32(32, 320, eps=1e-05, affine=True)
        (1): SiLU()
        (2): Dropout(p=0, inplace=False)
        (3): Conv2d(320, 320, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    )
    (skip_connection): Identity()
)
```

```
161 Robin Rombach, 3 months ago | 1 author (Robin Rombach)
162 class ResBlock(TimestepBlock):
163     """
164         A residual block that can optionally change the number of channels.
165
166         ...
167
168         def _forward(self, x, emb):
169             if self.updown:
170                 in_rest, in_conv = self.in_layers[:-1], self.in_layers[-1]
171                 h = in_rest(x)
172                 h = self.h_upd(h)
173                 x = self.x_upd(x)
174                 h = in_conv(h)
175             else:
176                 h = self.in_layers(x)
177                 emb_out = self.emb_layers(emb).type(h.dtype)
178                 while len(emb_out.shape) < len(h.shape):
179                     emb_out = emb_out[..., None]
180                 if self.use_scale_shift_norm:
181                     out_norm, out_rest = self.out_layers[0], self.out_layers[1:]
182                     scale, shift = th.chunk(emb_out, 2, dim=1)
183                     h = out_norm(h) * (1 + scale) + shift
184                     h = out_rest(h)
185                 else:
186                     h = h + emb_out
187                     h = self.out_layers(h)
188             return self.skip_connection(x) + h
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
```



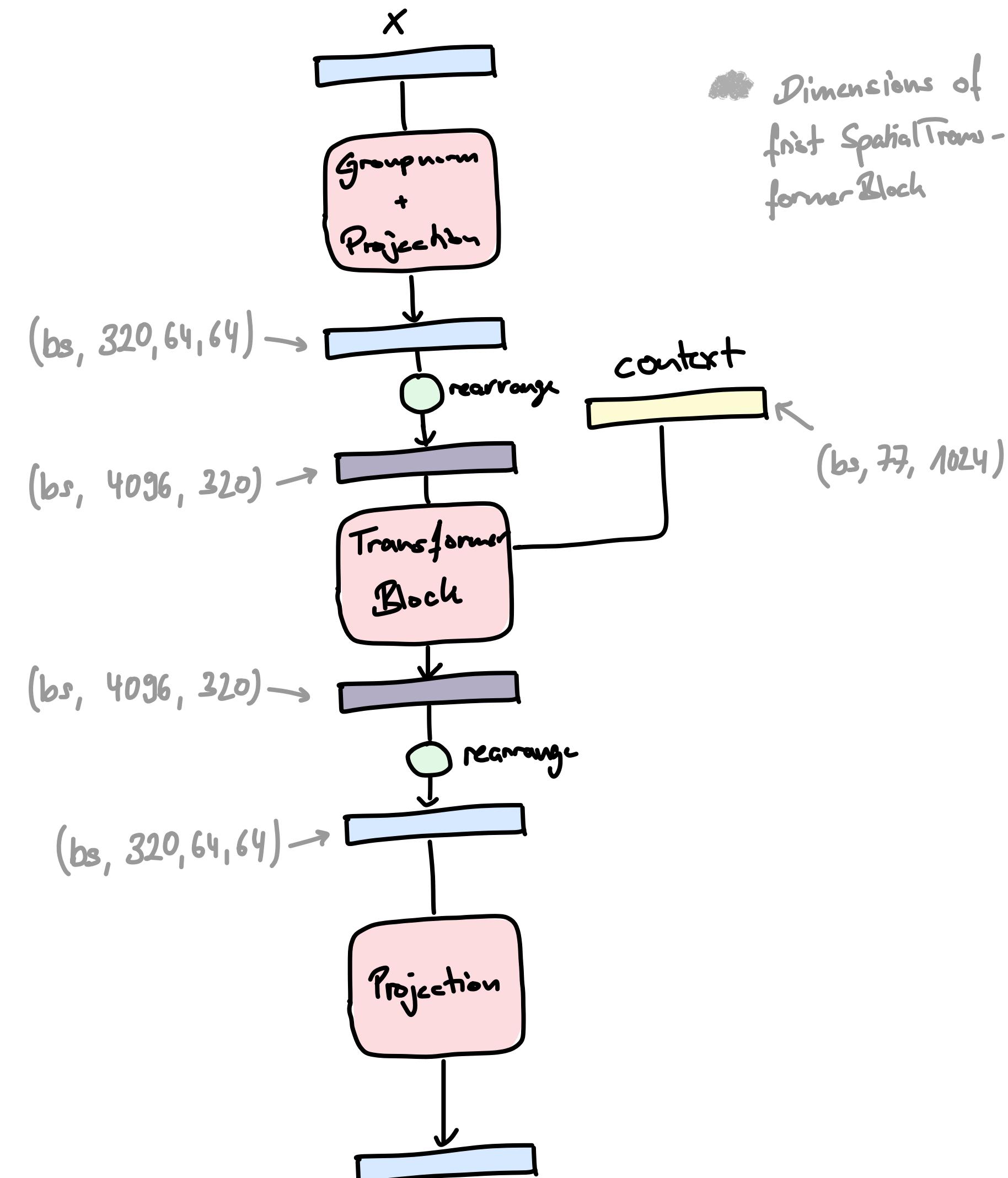
6. UNet Architecture in Depth

Spatial Transformer (Part 1)

```
unet.input_blocks[1][1]
  ✓ 0.0s
    SpatialTransformer(
      (norm): GroupNorm(32, 320, eps=1e-06, affine=True)
      (proj_in): Linear(in_features=320, out_features=320, bias=True)
      (transformer_blocks): ModuleList(
        (0): BasicTransformerBlock(
          (attn1): CrossAttention(
            (to_q): Linear(in_features=320, out_features=320, bias=False)
            (to_k): Linear(in_features=320, out_features=320, bias=False)
            (to_v): Linear(in_features=320, out_features=320, bias=False)
            (to_out): Sequential(
              (0): Linear(in_features=320, out_features=320, bias=True)
              (1): Dropout(p=0.0, inplace=False)
            )
          )
          (ff): FeedForward(
            (net): Sequential(
              (0): GEGLU(
                (proj): Linear(in_features=320, out_features=2560, bias=True)
              )
              (1): Dropout(p=0.0, inplace=False)
              (2): Linear(in_features=1280, out_features=320, bias=True)
            )
          )
          (attn2): CrossAttention(
            (to_q): Linear(in_features=320, out_features=320, bias=False)
            (to_k): Linear(in_features=1024, out_features=320, bias=False)
            (to_v): Linear(in_features=1024, out_features=320, bias=False)
            (to_out): Sequential(
              (0): Linear(in_features=320, out_features=320, bias=True)
              (1): Dropout(p=0.0, inplace=False)
            )
          )
        (norm1): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
        (norm2): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
        (norm3): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      )
      (proj_out): Linear(in_features=320, out_features=320, bias=True)
    )
```

```
278 class SpatialTransformer(nn.Module):
279     """
280         Transformer block for image-like data.
281         First, project the input (aka embedding)
282         and reshape to b, t, d.
283         Then apply standard transformer action.
284         Finally, reshape to image
285         NEW: use_linear for more efficiency instead of the 1x1 convs
286     """
287     def forward(self, x, context=None):
288         # note: if no context is given, cross-attention defaults to self-attention
289         if not isinstance(context, list):
290             context = [context]
291         b, c, h, w = x.shape
292         x_in = x
293         x = self.norm(x)
294         if not self.use_linear:
295             x = self.proj_in(x)
296             x = rearrange(x, 'b c h w -> b (h w) c').contiguous()
297             if self.use_linear:
298                 x = self.proj_in(x)
299             for i, block in enumerate(self.transformer_blocks):
300                 x = block(x, context=context[i])
301                 if self.use_linear:
302                     x = self.proj_out(x)
303                     x = rearrange(x, 'b (h w) c -> b c h w', h=h, w=w).contiguous()
304             if not self.use_linear:
305                 x = self.proj_out(x)
306         return x + x_in
```

ldm/modules/attention.py



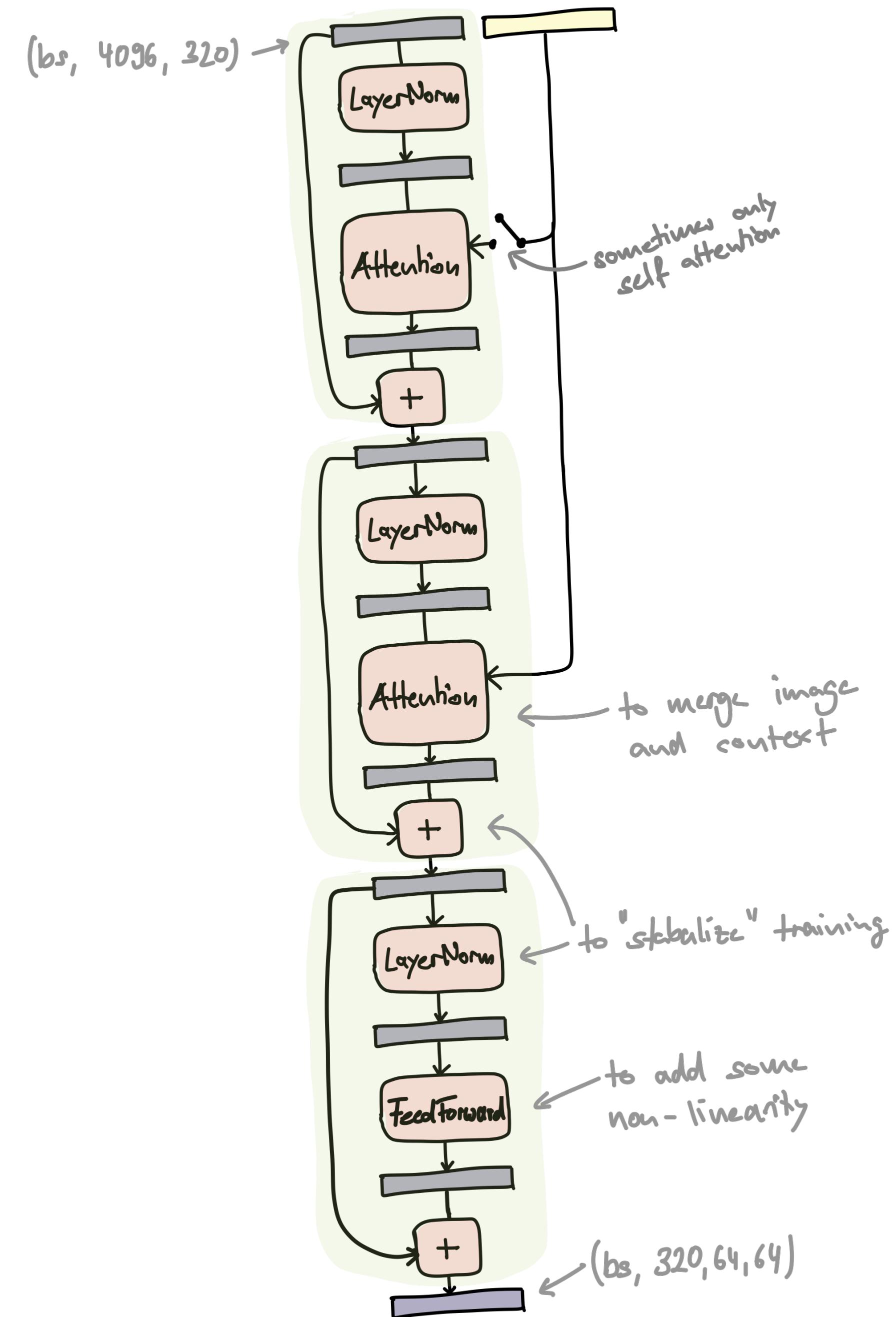
6. UNet Architecture in Depth

Spatial Transformer (Part 2)

```
unet.input_blocks[1][1]
  ✓ 0.0s
  SpatialTransformer(
    (norm): GroupNorm(32, 320, eps=1e-06, affine=True)
    (proj_in): Linear(in_features=320, out_features=320, bias=True)
    (transformer_blocks): ModuleList(
      (0): BasicTransformerBlock(
        (attn1): CrossAttention(
          (to_q): Linear(in_features=320, out_features=320, bias=False)
          (to_k): Linear(in_features=320, out_features=320, bias=False)
          (to_v): Linear(in_features=320, out_features=320, bias=False)
          (to_out): Sequential(
            (0): Linear(in_features=320, out_features=320, bias=True)
            (1): Dropout(p=0.0, inplace=False)
          )
        )
        (ff): FeedForward(
          (net): Sequential(
            (0): GELU(
              (proj): Linear(in_features=320, out_features=2560, bias=True)
            )
            (1): Dropout(p=0.0, inplace=False)
            (2): Linear(in_features=1280, out_features=320, bias=True)
          )
        )
        (attn2): CrossAttention(
          (to_q): Linear(in_features=320, out_features=320, bias=False)
          (to_k): Linear(in_features=1024, out_features=320, bias=False)
          (to_v): Linear(in_features=1024, out_features=320, bias=False)
          (to_out): Sequential(
            (0): Linear(in_features=320, out_features=320, bias=True)
            (1): Dropout(p=0.0, inplace=False)
          )
        )
      (norm1): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      (norm2): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      (norm3): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
    )
    (proj_out): Linear(in_features=320, out_features=320, bias=True)
  )
```

```
Robin Rombach, 3 months ago | 1 author (Robin Rombach)
246 class BasicTransformerBlock(nn.Module):
247
248     def _forward(self, x, context=None):
249         x = self.attn1(self.norm1(x), context=context if self.disable_self_attn else None) + x
250         x = self.attn2(self.norm2(x), context=context) + x
251         x = self.ff(self.norm3(x)) + x
252
253         return x
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
```

ldm/modules/attention.py



6. UNet Architecture in Depth

Spatial Transformer (Part 3)

```
unet.input_blocks[1][1]
  ✓ 0.0s
    SpatialTransformer(
      (norm): GroupNorm(32, 320, eps=1e-06, affine=True)
      (proj_in): Linear(in_features=320, out_features=320, bias=True)
      (transformer_blocks): ModuleList(
        (0): BasicTransformerBlock(
          (attn1): CrossAttention(
            (to_q): Linear(in_features=320, out_features=320, bias=False)
            (to_k): Linear(in_features=320, out_features=320, bias=False)
            (to_v): Linear(in_features=320, out_features=320, bias=False)
            (to_out): Sequential(
              (0): Linear(in_features=320, out_features=320, bias=True)
              (1): Dropout(p=0.0, inplace=False)
            )
          )
          (ff): FeedForward(
            (net): Sequential(
              (0): GELU(
                (proj): Linear(in_features=320, out_features=2560, bias=True)
              )
              (1): Dropout(p=0.0, inplace=False)
              (2): Linear(in_features=1280, out_features=320, bias=True)
            )
          )
        )
        (attn2): CrossAttention(
          (to_q): Linear(in_features=320, out_features=320, bias=False)
          (to_k): Linear(in_features=1024, out_features=320, bias=False)
          (to_v): Linear(in_features=1024, out_features=320, bias=False)
          (to_out): Sequential(
            (0): Linear(in_features=320, out_features=320, bias=True)
            (1): Dropout(p=0.0, inplace=False)
          )
        );
      )
      (norm1): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      (norm2): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      (norm3): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
    )
    (proj_out): Linear(in_features=320, out_features=320, bias=True)
  )
```

```
144 Dango233, 2 months ago | 2 authors (Robin Rombach and others)
145 class CrossAttention(nn.Module):
146     def __init__(self, query_dim, context_dim=None, heads=8, dim_head=64, dropout=0.):
147         super().__init__()
148         inner_dim = dim_head * heads
149         context_dim = default(context_dim, query_dim)
150
151         self.scale = dim_head ** -0.5
152         self.heads = heads
153
154         self.to_q = nn.Linear(query_dim, inner_dim, bias=False)
155         self.to_k = nn.Linear(context_dim, inner_dim, bias=False)
156         self.to_v = nn.Linear(context_dim, inner_dim, bias=False)
157
158         self.to_out = nn.Sequential(
159             nn.Linear(inner_dim, query_dim),
160             nn.Dropout(dropout)
161         )
162
163     def forward(self, x, context=None, mask=None):
164         h = self.heads
165
166         q = self.to_q(x) ❶
167         context = default(context, x) ❷
168         k = self.to_k(context)
169         v = self.to_v(context)
170
171         q, k, v = map(lambda t: rearrange(t, 'b n (h d) -> (b h) n d', h=h), (q, k, v)) ❸
172
173         # force cast to fp32 to avoid overflowing
174         if _ATTN_PRECISION == "fp32":
175             with torch.autocast(enabled=False, device_type='cuda'):
176                 q, k = q.float(), k.float()
177                 sim = einsum('b i d, b j d -> b i j', q, k) * self.scale
178             else:
179                 sim = einsum('b i d, b j d -> b i j', q, k) * self.scale
180
181         del q, k
182
183         if exists(mask):
184             mask = rearrange(mask, 'b ... -> b (...)')
185             max_neg_value = -torch.finfo(sim.dtype).max
186             mask = repeat(mask, 'b j -> (b h) () j', h=h)
187             sim.masked_fill_(~mask, max_neg_value)
188
189         # attention, what we cannot get enough of
190         sim = sim.softmax(dim=-1)
191
192         out = einsum('b i j, b j d -> b i d', sim, v)
193         out = rearrange(out, '(b h) n d -> b n (h d)', h=h)
194
195     return self.to_out(out)
```

example dimensions of first CrossAttention Block:

- $x = \text{torch.Size}([1\text{bs}, 4096, 320])$
context = $\text{torch.Size}([1\text{bs}, 77, 1024])$
 $h = 5$
 $q = \text{torch.Size}([1\text{bs}, 4096, 320])$
 $k = \text{torch.Size}([1\text{bs}, 77, 320])$
 $v = \text{torch.Size}([1\text{bs}, 77, 320])$

- $q = \text{torch.Size}([5\text{bs}, 4096, 64])$
 $k = \text{torch.Size}([5\text{bs}, 77, 64])$
 $v = \text{torch.Size}([5\text{bs}, 77, 64])$

- $\text{sim} = \text{torch.Size}([5\text{bs}, 4096, 77])$

- $\text{out} = \text{torch.Size}([5\text{bs}, 4096, 64])$

- $\text{out} = \text{torch.Size}([1\text{bs}, 4096, 320])$

- $\text{ret} = \text{torch.Size}([1\text{bs}, 4096, 320])$

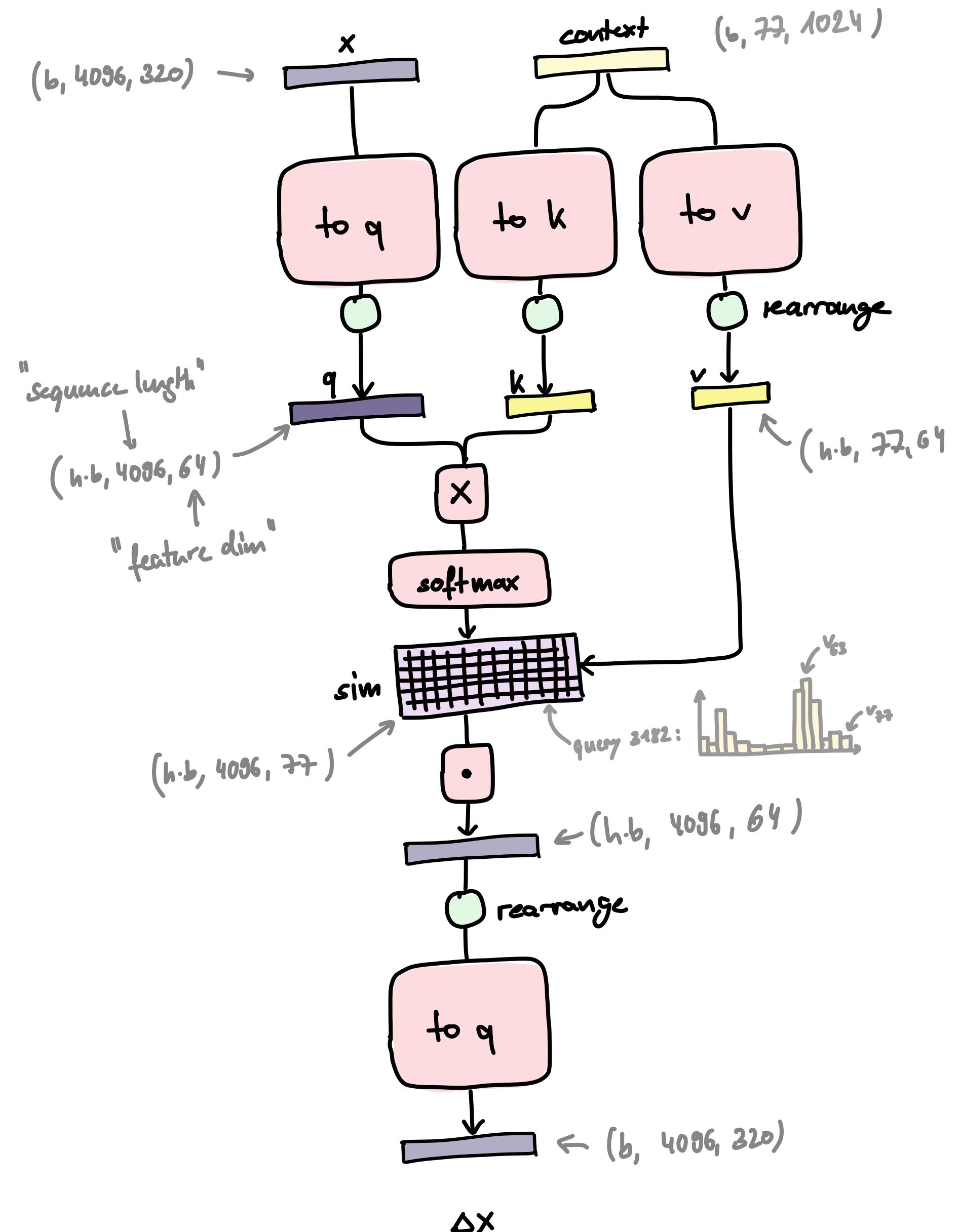
6. UNet Architecture in Depth

Spatial Transformer (Part 3)

```

unet.input_blocks[1][1]
  ✓ 0.0s
  SpatialTransformer(
    (norm): GroupNorm(32, 320, eps=1e-06, affine=True)
    (proj_in): Linear(in_features=320, out_features=320, bias=True)
    (transformer_blocks): ModuleList(
      (0): BasicTransformerBlock(
        (attn1): CrossAttention(
          (to_q): Linear(in_features=320, out_features=320, bias=False)
          (to_k): Linear(in_features=320, out_features=320, bias=False)
          (to_v): Linear(in_features=320, out_features=320, bias=False)
          (to_out): Sequential(
            (0): Linear(in_features=320, out_features=320, bias=True)
            (1): Dropout(p=0.0, inplace=False)
          )
        )
        (ff): FeedForward(
          (net): Sequential(
            (0): GEGLU(
              (proj): Linear(in_features=320, out_features=2560, bias=True)
            )
            (1): Dropout(p=0.0, inplace=False)
            (2): Linear(in_features=1280, out_features=320, bias=True)
          )
        )
      )
      (attn2): CrossAttention(
        (to_q): Linear(in_features=320, out_features=320, bias=False)
        (to_k): Linear(in_features=1024, out_features=320, bias=False)
        (to_v): Linear(in_features=1024, out_features=320, bias=False)
        (to_out): Sequential(
          (0): Linear(in_features=320, out_features=320, bias=True)
          (1): Dropout(p=0.0, inplace=False)
        )
      )
      (norm1): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      (norm2): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
      (norm3): LayerNorm((320,), eps=1e-05, elementwise_affine=True)
    )
    (proj_out): Linear(in_features=320, out_features=320, bias=True)
  )

```



7. ControlNet (Feb. 2023)

Adding Conditional Control to Text-to-Image Diffusion Models

Lvmin Zhang and Maneesh Agrawala
Stanford University



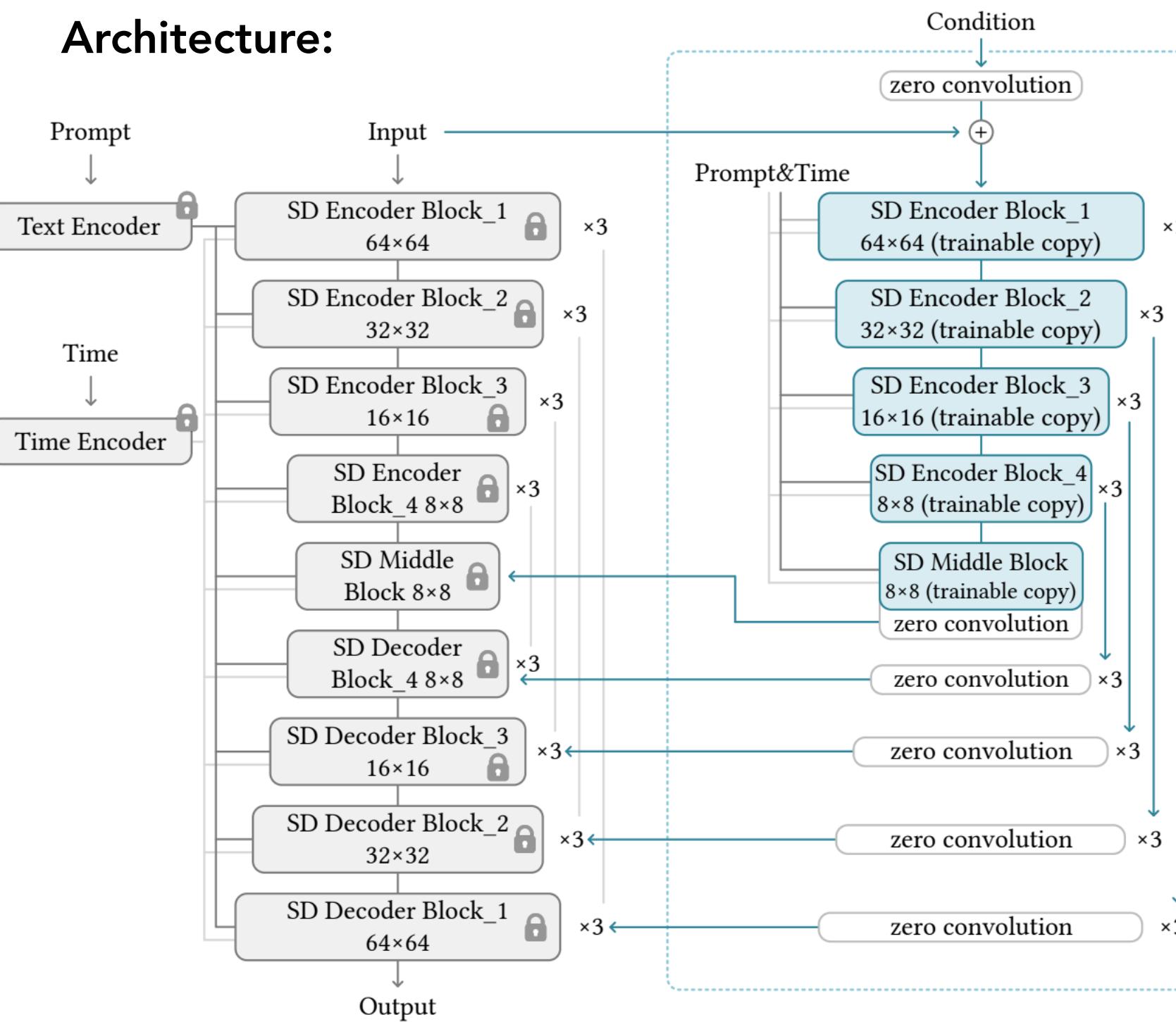
Main contributions:

1. New method to condition large text2image diffusion models through “zero convolutions”.

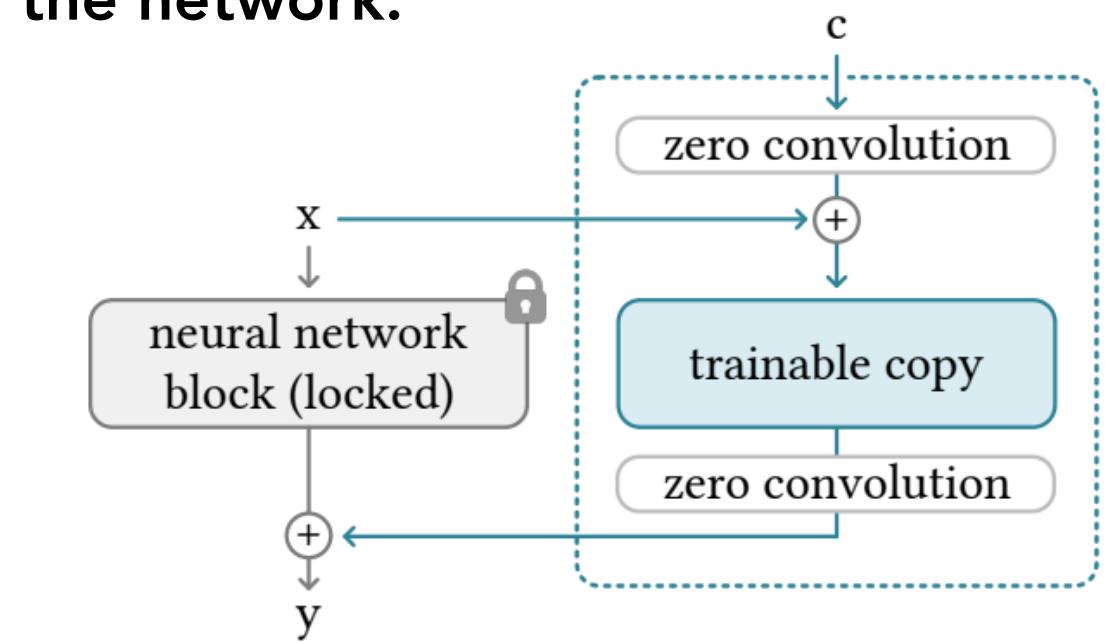
Improvement over fine-tuning:

- Network always keeps production ready state, i.e., it does not forget its general txt2image capabilities.
- Training is as fast as fine-tuning.
- No overfitting.

Architecture:



Detailed view of one “block” somewhere in the network:



Zero convolutions: Convolutional layers that are initialised with vanishing weights and biases ($W=0, B=0$).

- At the first training step: ControlNet has no effect.
- After optimisation: ControlNet effectively augments each feature vector according to the conditioning vector.