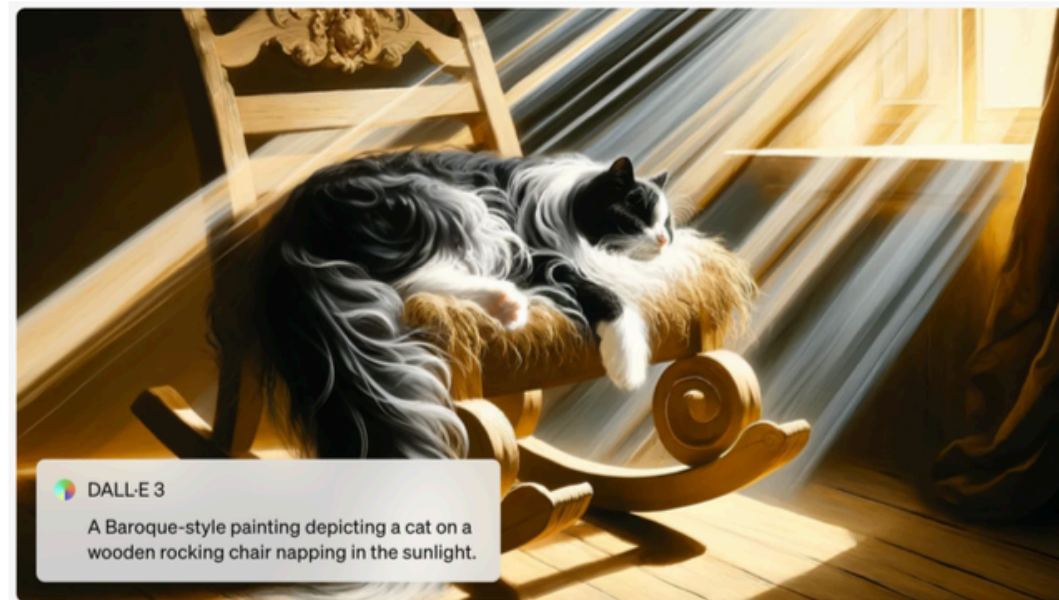


Diffusion Models

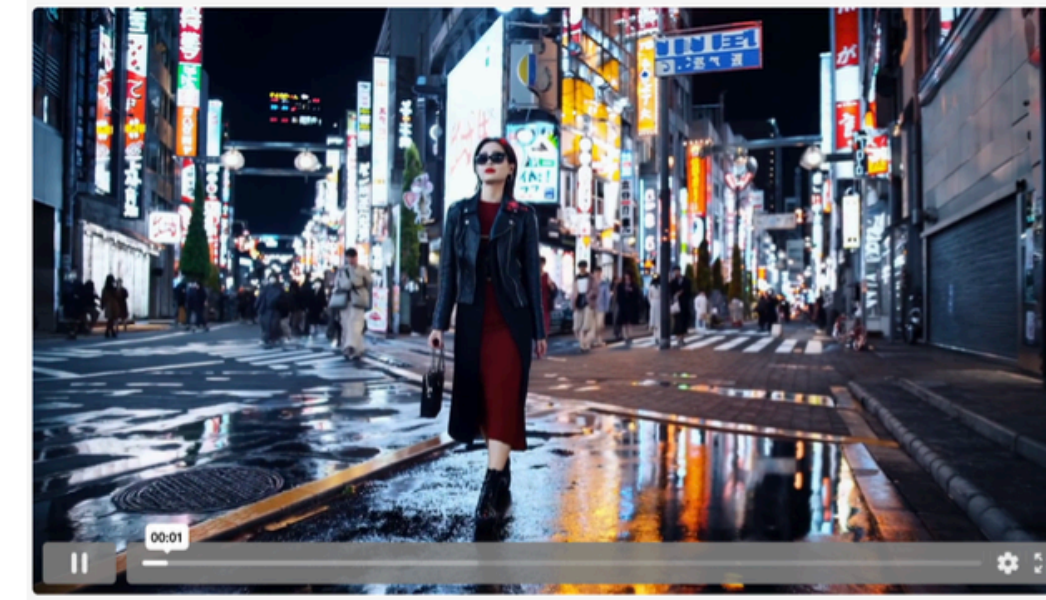
(Score-based generative models, Flow matching)

Yinbo Chen

Generative modeling

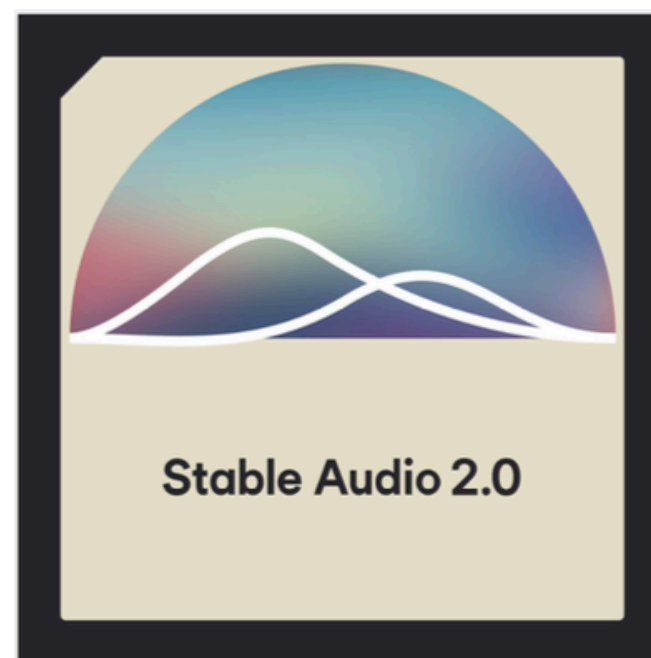


Image

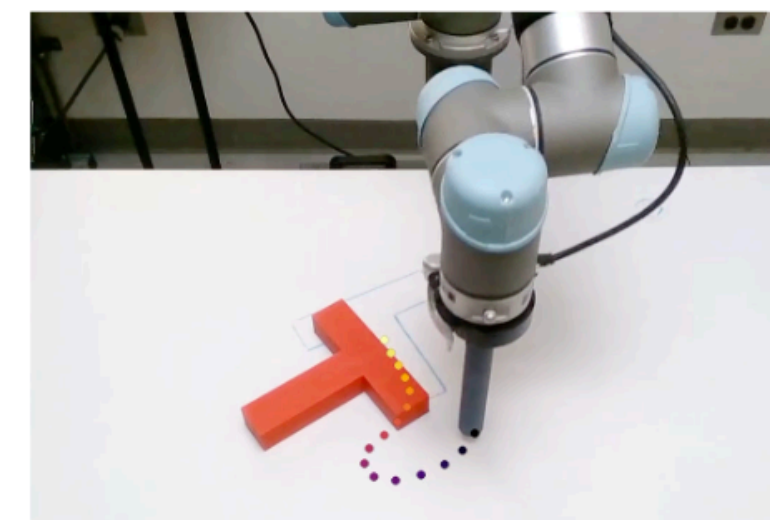


Prompt: A stylish woman walks down a Tokyo street filled with warm glowing neon and animated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a...

Video



Audio



Policy

Task definition

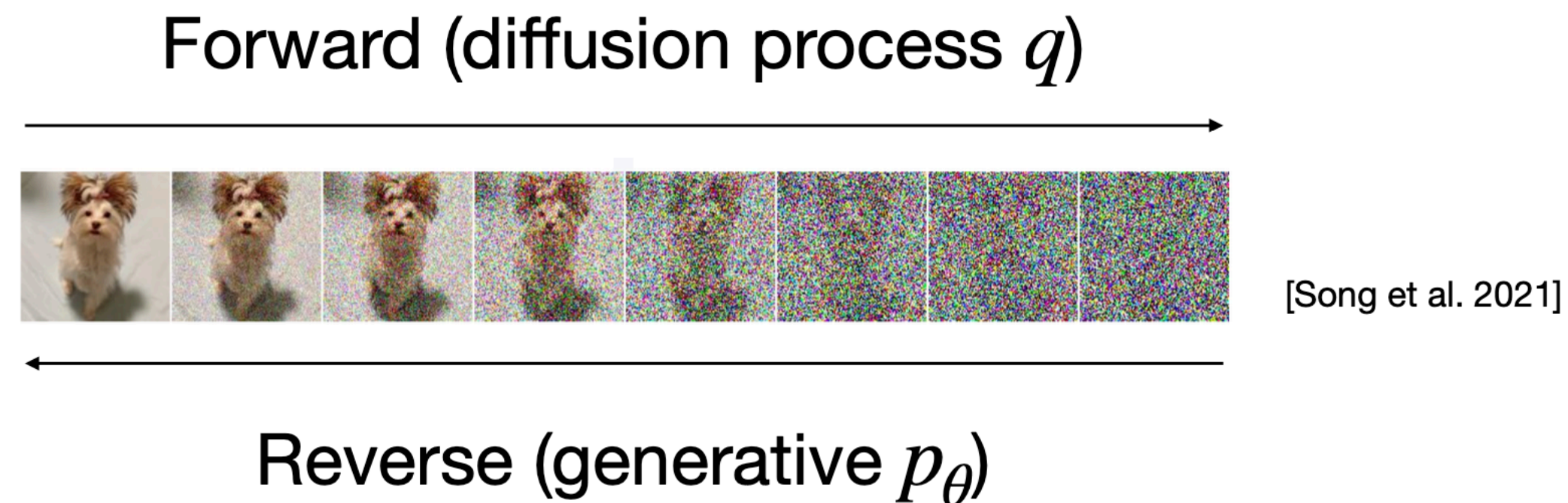
- Given a dataset of samples $x \sim q(x)$
- x is a vector in general (e.g., image: dimension = $h * w * 3$)
- Learn a model that generates samples from $p_{\theta}(x) \approx q(x)$

Overview

- **Basics:** discrete time, continuous time (SDE / ODE)
- **Applications:** acceleration, conditional generation
- Score distillation
- (Topics not covered: DreamBooth, engineering ...)

Denoising Diffusion Probabilistic Models

“Generate by denoising”



- Forward process defines joint distribution $q(x_0, \dots, x_T)$
- Learn model p_θ to reverse: sample noise x_T , denoise step by step
- (θ is the network parameters, p_θ denotes the generated distribution)

Forward process q

Gradually add noise

- Defined by a set of β_t
- $q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$
 - i.e., $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\varepsilon \quad \varepsilon \sim \mathcal{N}(0, I)$
- $\sqrt{1 - \beta_t}$ is just to preserve the variance
 - Assume $Var(x_0) = 1$, after each step Var is still 1

Forward (diffusion process q)



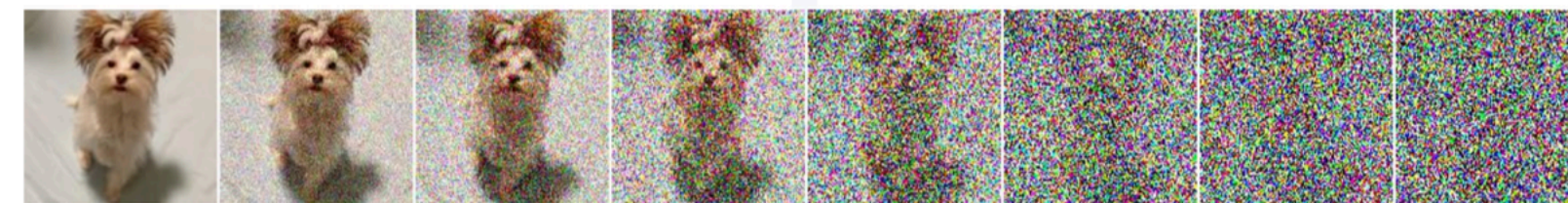
[Song et al. 2021]

Forward process q

Add noise

- Defined by a set of β_t
- $q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t}x_{t-1}, \beta_t I)$
 - i.e., $x_t = \sqrt{1 - \beta_t}x_{t-1} + \sqrt{\beta_t}\varepsilon$
- Let $\alpha_t = 1 - \beta_t$, and $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$
- $q(x_t | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)I)$

Forward (diffusion process q)

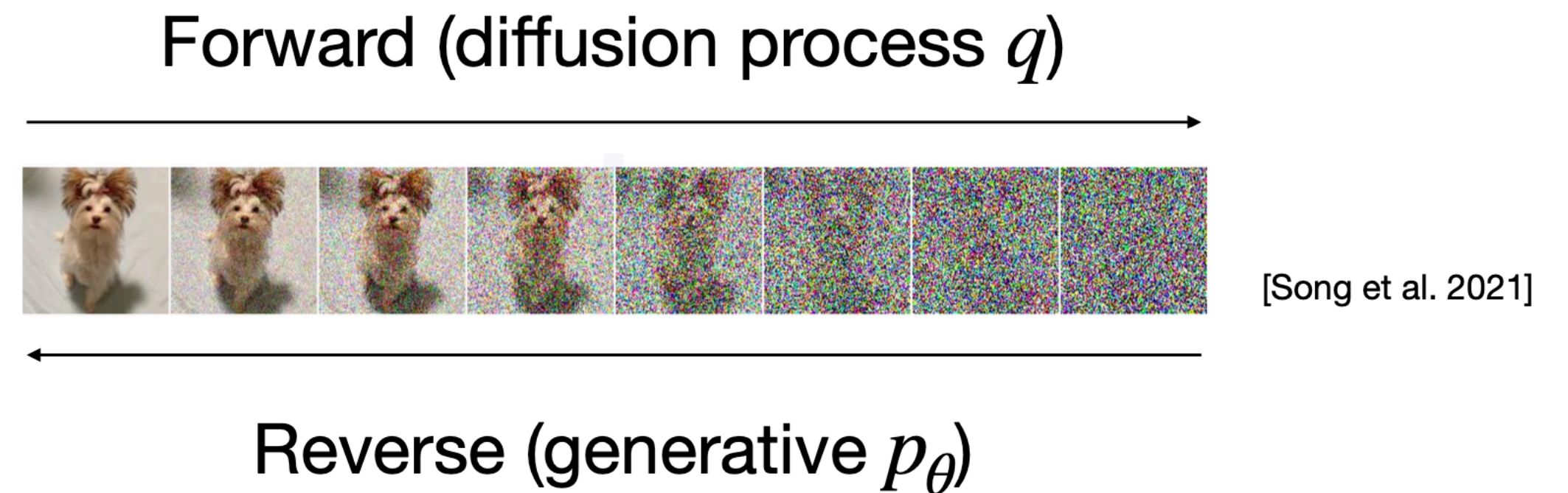


[Song et al. 2021]

Reverse process p_θ

Denoise

- How to model $q(x_{t-1} | x_t)$?
 - It is Gaussian if β_t are small enough



- $p_\theta(x_{t-1} | x_t) = \mathcal{N}(\mu_\theta(x_t, t), \sigma_t^2 I)$
 - Use mean network μ_θ to define p_θ (empirically $\sigma_t^2 = \beta_t$ is good [Ho et al. 2020])
- Model defined, how to train it?

Training loss \mathcal{L} : maximize data likelihood

Evidence lower bound (ELBO)

- **ELBO:** $\log p_{\theta}(x) \geq E_{z \sim q(z|x)} \left[\log \frac{p_{\theta}(x, z)}{q(z|x)} \right] =: \mathcal{L}$

- In diffusion models: “ x ” is x_0 , “ z ” is $x_{1 \dots T}$.

- DDPM [Ho et al. 2020] shows with $\mu_{\theta}(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_{\theta}(x_t, t) \right)$,
 $\mathcal{L} = E_{x_0, t, \varepsilon} [w_t || \varepsilon - \varepsilon_{\theta}(x_t, t) ||^2]$, where $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon$

DDPM [Ho et al. 2020]

Discard w_t in \mathcal{L} empirically works well

- DDPM [Ho et al. 2020] shows with $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}}(x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\epsilon_\theta(x_t, t))$,
 $\mathcal{L} = E_{x_0, t, \epsilon} [w_t || \epsilon - \epsilon_\theta(x_t, t) ||^2]$, where $x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon$

Algorithm 1 Training

- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on
 $\nabla_\theta ||\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)||^2$
 - 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: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
 - 5: **end for**
 - 6: **return** \mathbf{x}_0
-

Score-based generative models

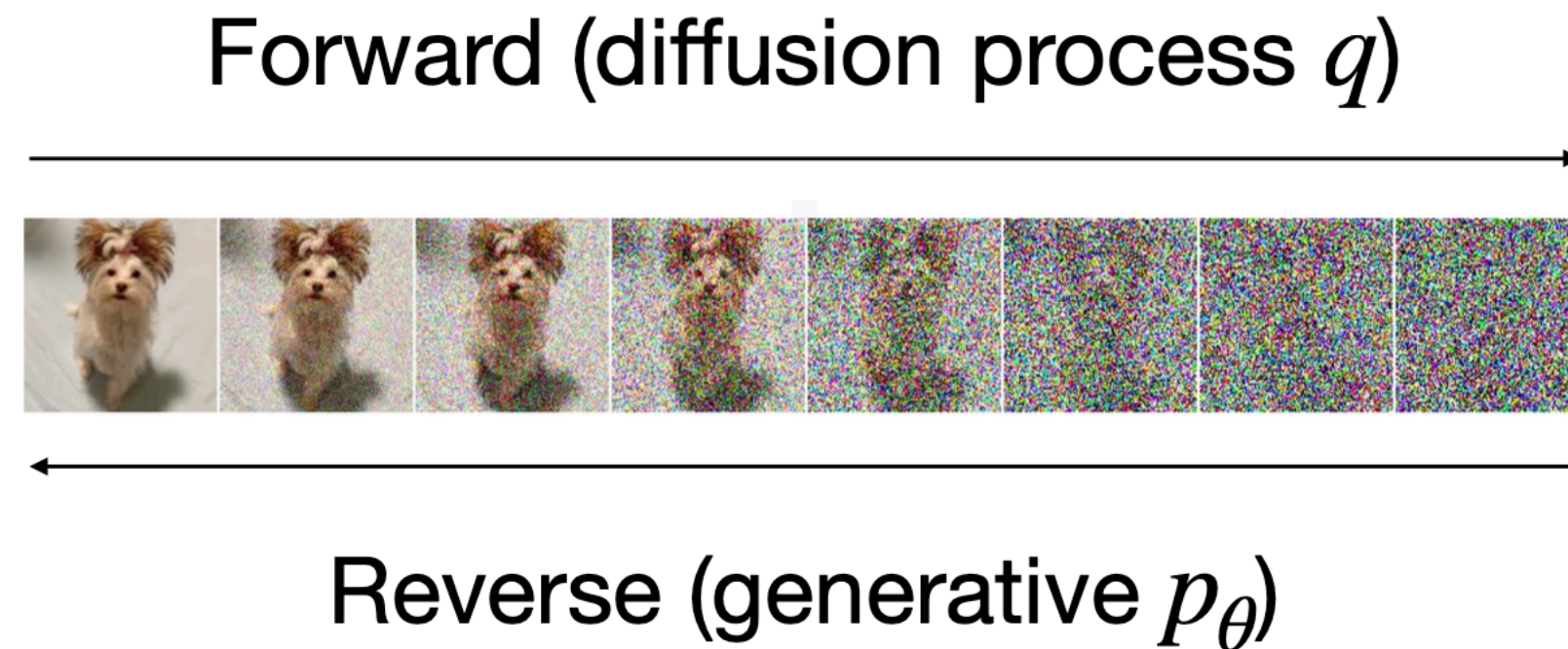
SDE in continuous time

- $x_t = \sqrt{1 - \beta_t} x_{t-1} + \sqrt{\beta_t} \varepsilon$

- Infinitely small step size \rightarrow continuous time

- SDE: $dx = -\frac{1}{2}\beta(t)x dt + \sqrt{\beta(t)} dw$

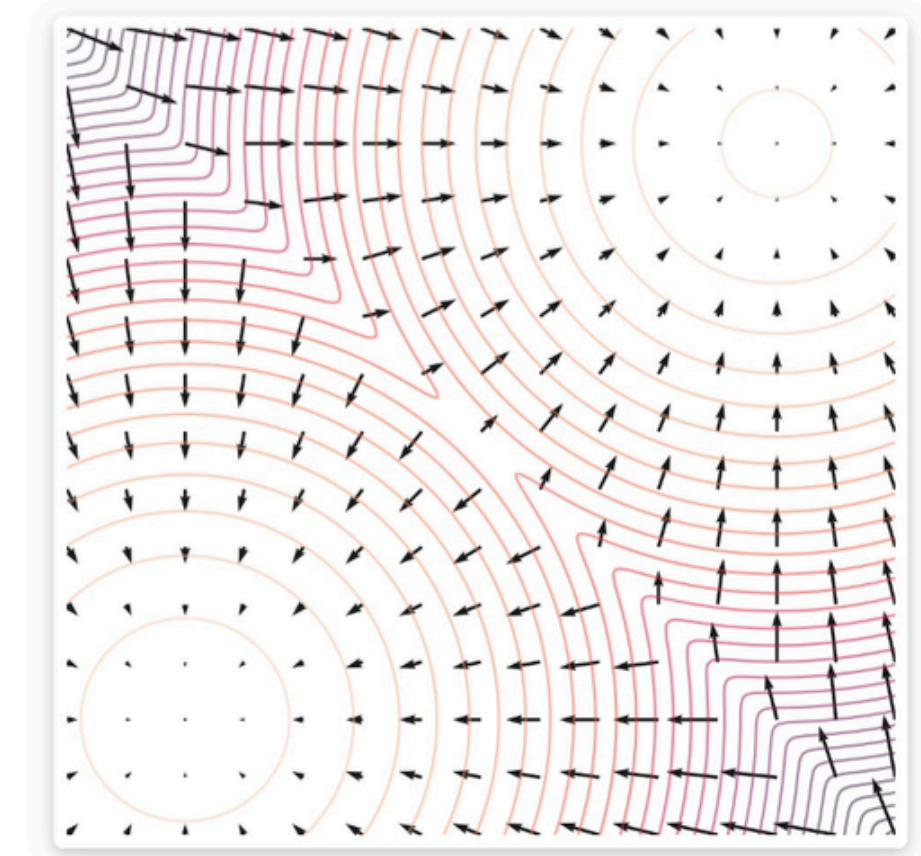
- dw : Wiener process (infinitely small noise)



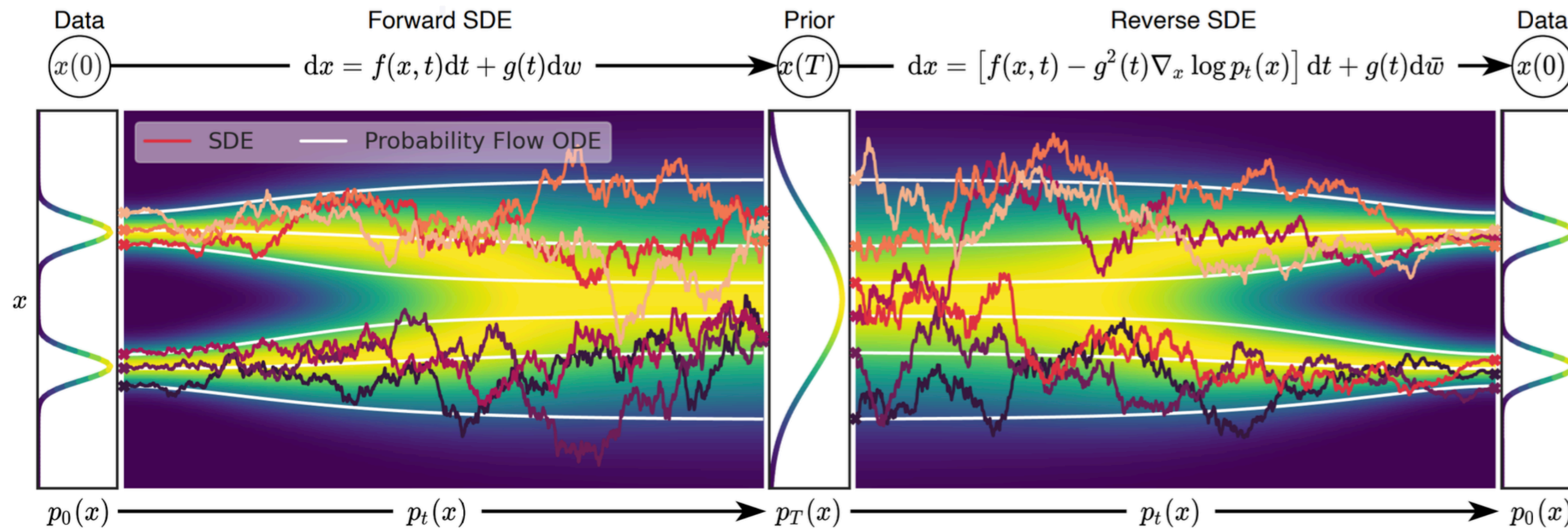
[Song et al. 2021]

Reverse the SDE

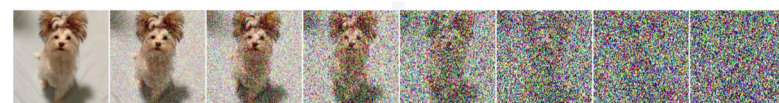
Based on score function $\nabla_x \log p(x)$



Score function (the vector field) and density function (contours) of a mixture of two Gaussians.



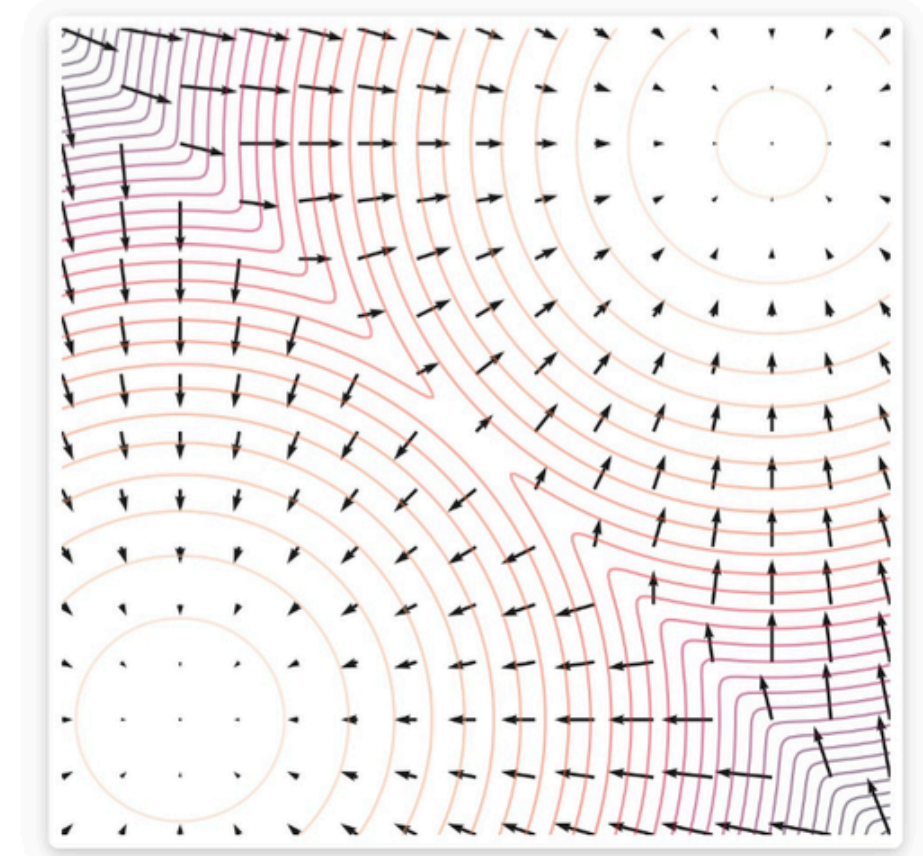
[Song et al. 2021]



How to get $\nabla_{x_t} \log q(x_t)$?

Denoising score matching

- $\min_{\theta} E_{x_0, t, \varepsilon} \left\| s_{\theta}(x_t, t) - \nabla_{x_t} \log q(x_t | x_0) \right\|^2$
 - $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \varepsilon$
 - At minimum point (not necessarily 0), $s_{\theta}(x_t, t) \approx \nabla_{x_t} \log q(x_t)$
- In diffusion, $q(x_t | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I)$
- $\nabla_{x_t} \log q(x_t | x_0) = \Sigma^{-1}(\mu - x_t) = \frac{-\varepsilon}{\sqrt{1 - \bar{\alpha}_t}}$



Score function (the vector field) and density function (contours) of a mixture of two Gaussians.

How to get $\nabla_{x_t} \log q(x_t)$?

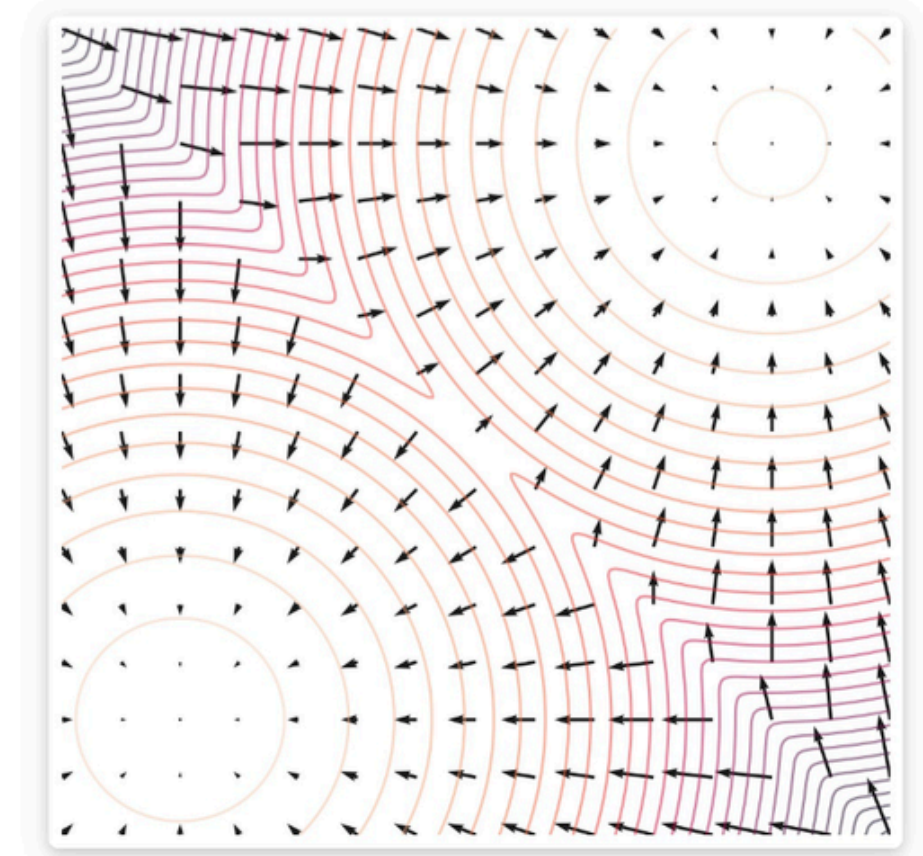
Denoising score matching

- $$\min_{\theta} E_{x_0, t, \varepsilon} \left\| s_{\theta}(x_t, t) - \frac{-\varepsilon}{\sqrt{1 - \bar{\alpha}_t}} \right\|^2$$

- Let
$$s_{\theta}(x_t, t) = \frac{-\varepsilon_{\theta}(x_t, t)}{\sqrt{1 - \bar{\alpha}_t}}$$

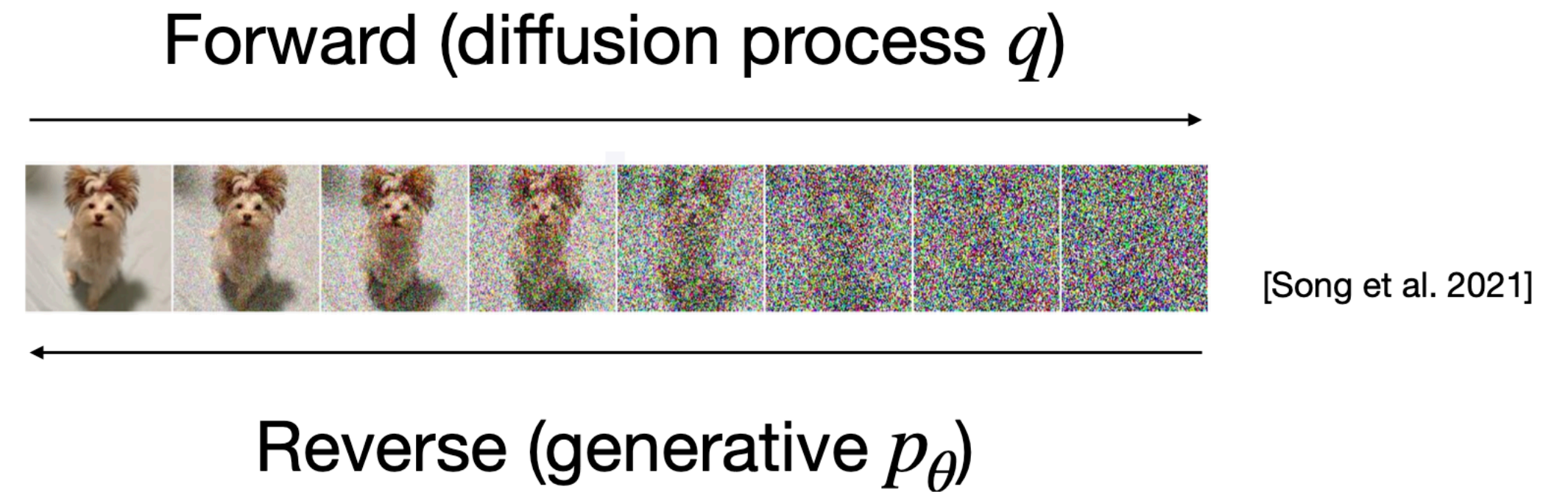
- It becomes:
$$\min_{\theta} E_{x_0, t, \varepsilon} \frac{1}{1 - \bar{\alpha}_t} \left\| \varepsilon_{\theta}(x_t, t) - \varepsilon \right\|^2$$

- Same denoising objective! (DDPM derived from ELBO)



Score function (the vector field) and density function (contours) of a mixture of two Gaussians.

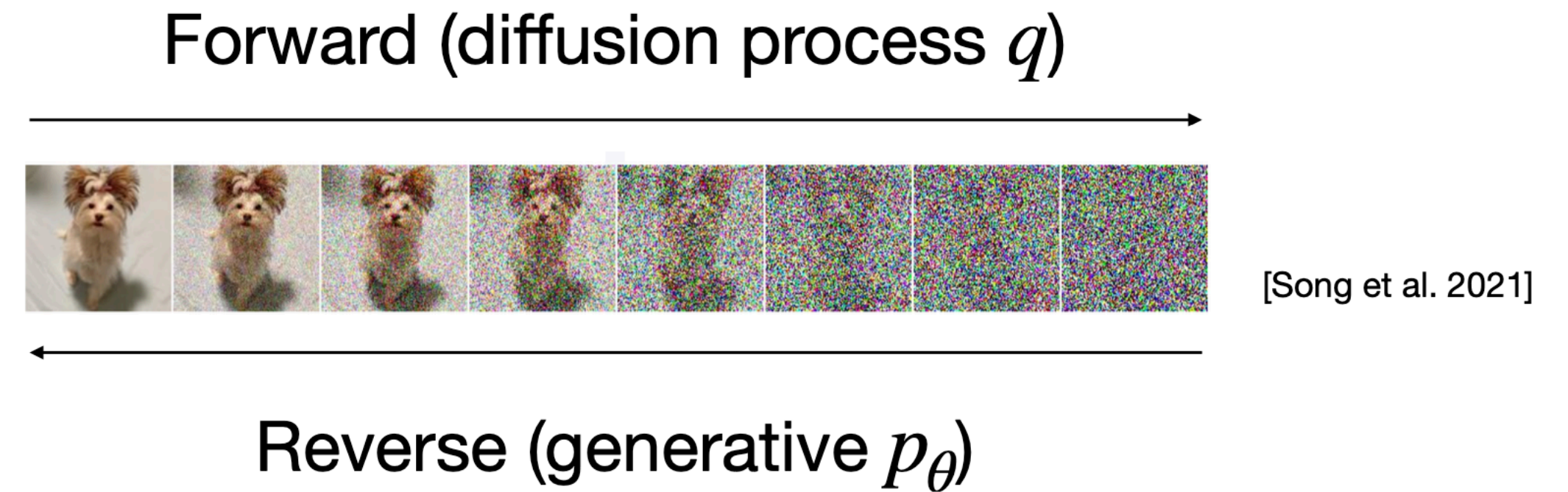
Summary



- Learn a noise predictor: $E_{x_0, t, \varepsilon} [w_t \mid \|\varepsilon_\theta(x_t, t) - \varepsilon\|^2]$
- Get score functions at all noisy dist.: $\nabla_{x_t} \log q(x_t) \approx -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \varepsilon_\theta(x_t, t)$
- Under L2 loss, the optimal $\varepsilon_\theta(x_t, t)$ is the expected ε given x_t (mean)
- In theory, w_t does not change the optimal point (but changes ELBO or not)

Generalized view

General formulation

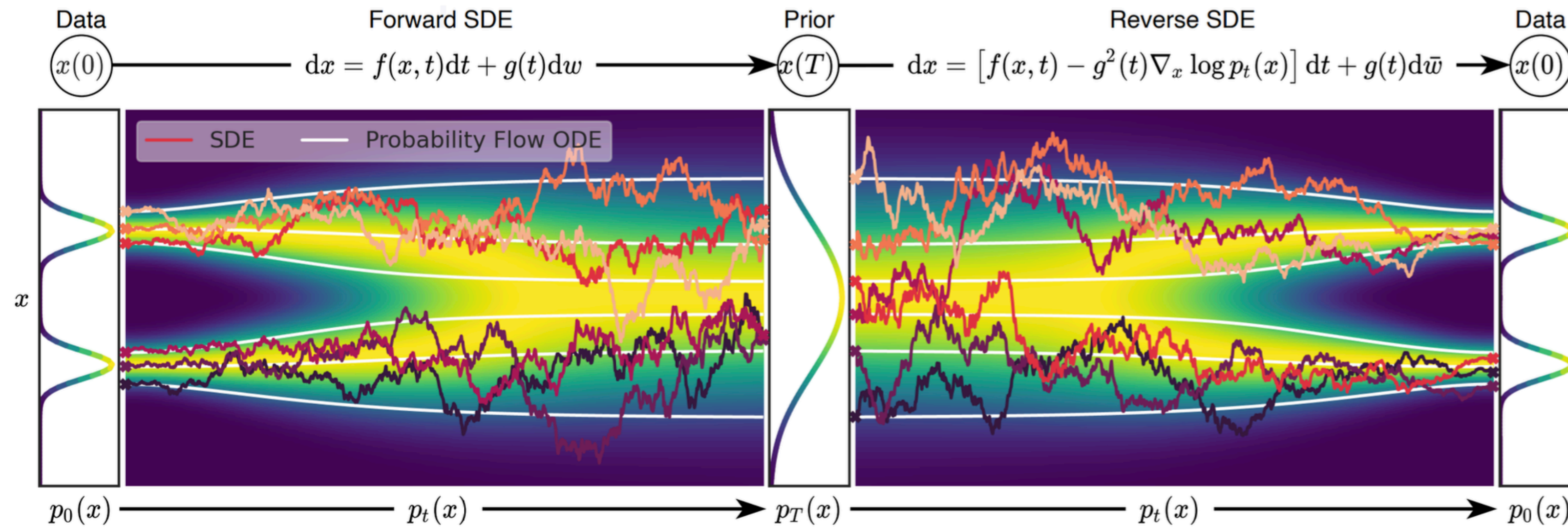


- A diffusion process is defined by $x_t = \alpha_t x_0 + \sigma_t \mathcal{E}$ (a different α than before)
 - **Signal ratio:** α_t from 1 to 0
 - **Noise ratio:** σ_t from 0 to 1
- The change rate of α_t defines the previous “ β_t ” and we don’t force variance 1

Sample prediction

- $x_t = \alpha_t x_0 + \sigma_t \varepsilon$, so $x_0 = (x_t - \sigma_t \varepsilon) / \alpha_t$
- Define $x_\theta(x_t, t)$ as $x_\theta = (x_t - \sigma_t \varepsilon_\theta) / \alpha_t$
- $||x_\theta - x_0||^2 = ||(x_t - \sigma_t \varepsilon_\theta) / \alpha_t - (x_t - \sigma_t \varepsilon) / \alpha_t||^2 = \left(\frac{\sigma_t}{\alpha_t}\right)^2 ||\varepsilon_\theta - \varepsilon||^2$
- Predict x_0 , i.e., $E_{x_0, t, \varepsilon} ||x_\theta(x_t, t) - x_0||^2$, is just a time-reweighed ε -prediction
- Furthermore, predict any $A\varepsilon + Bx_0$ from x_t learns the same score function

Probability Flow ODE (PF-ODE)



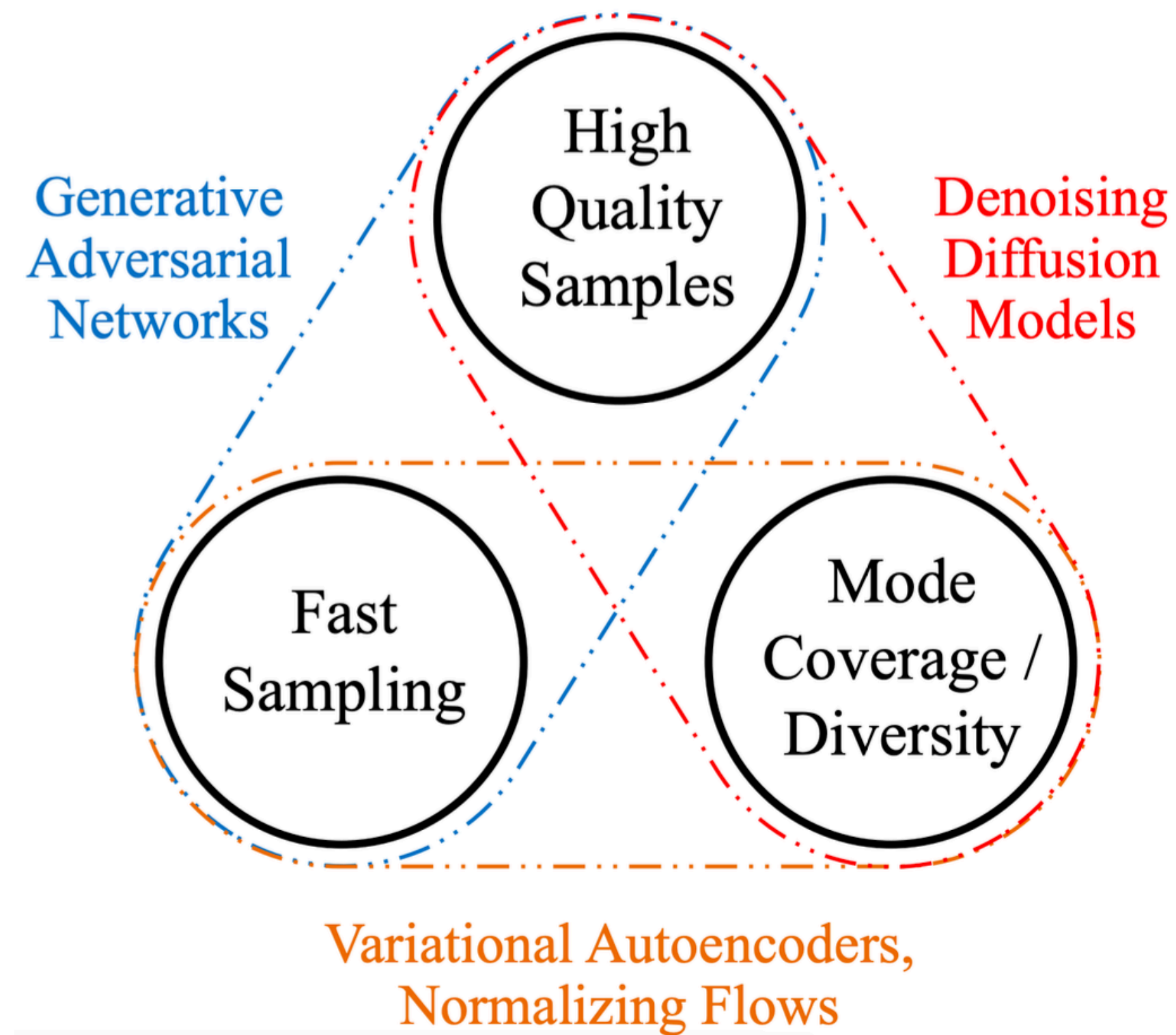
[Song et al. 2021]

- Reverse ODE: $dx = \left[f(x, t) - \frac{1}{2}g(t)^2 \nabla_x \log p_t(x) \right] dt$
- Gives the same marginal distributions $p_t(x)$ (but a different joint distribution)

Flow matching

- $x_t = \alpha_t x_0 + \sigma_t \varepsilon$
 - **Signal ratio:** $\alpha_t = 1 - t$
 - **Noise ratio:** $\sigma_t = t$
- Prediction: $\varepsilon - x_0$, i.e., $\mathcal{L} = E_{x_0, t, \varepsilon} || v_\theta(x_t, t) - (\varepsilon - x_0) ||^2$
- Is simple, is ELBO [Kingma et al., 2023] (Used in many recent SOTA models)

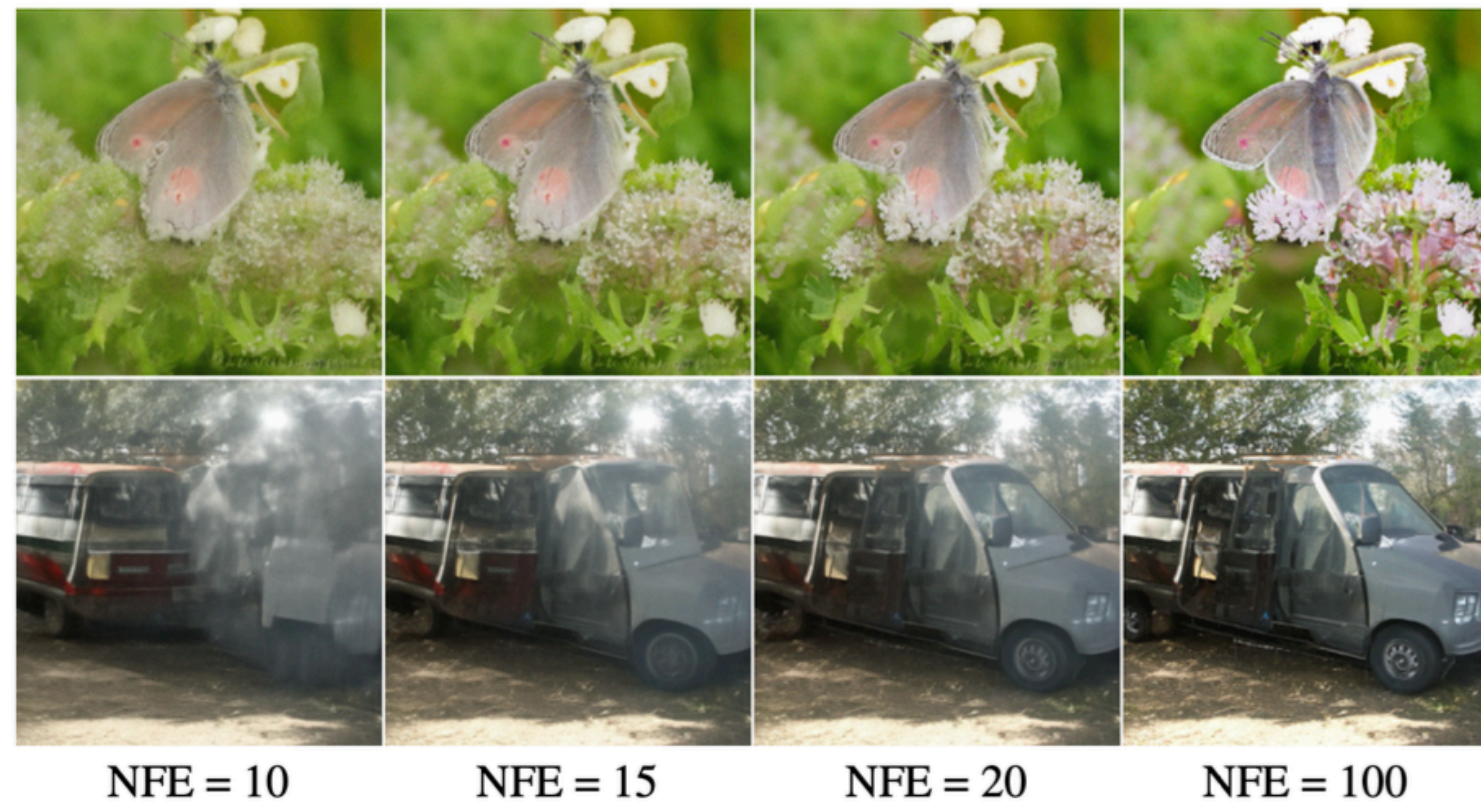
Generative learning trilemma



[Xiao et al. 2022]

Sampling acceleration

1. Better ODE/SDE solvers

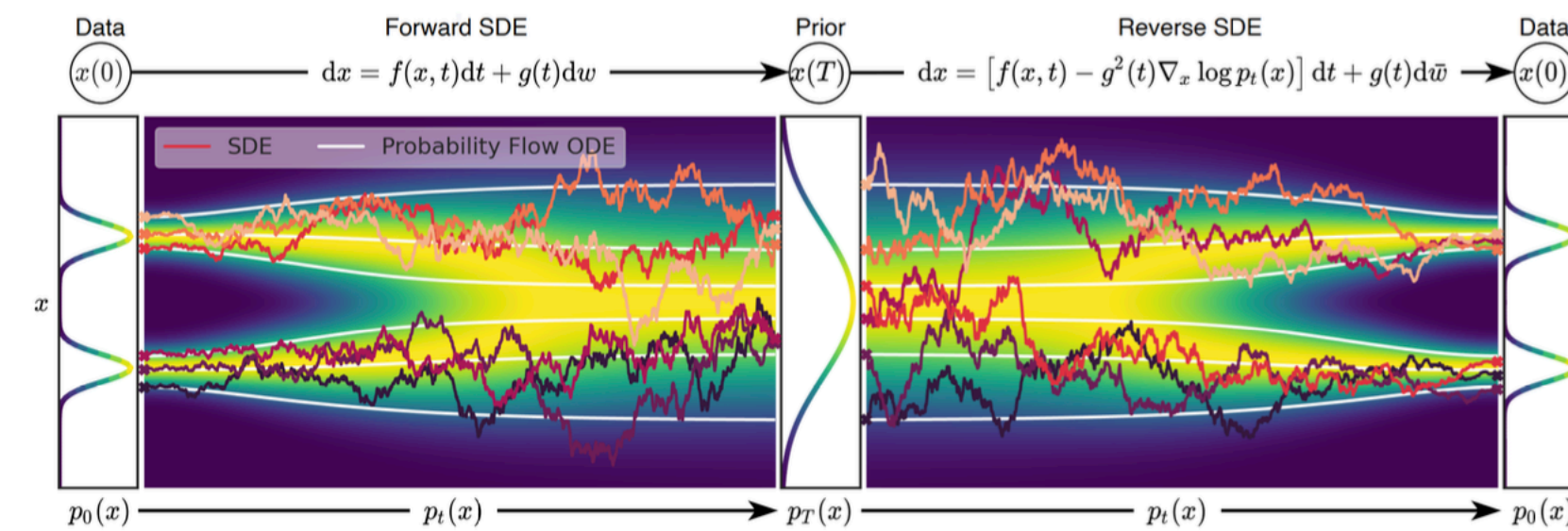


(a) DDIM [19]



(b) DPM-Solver (ours)

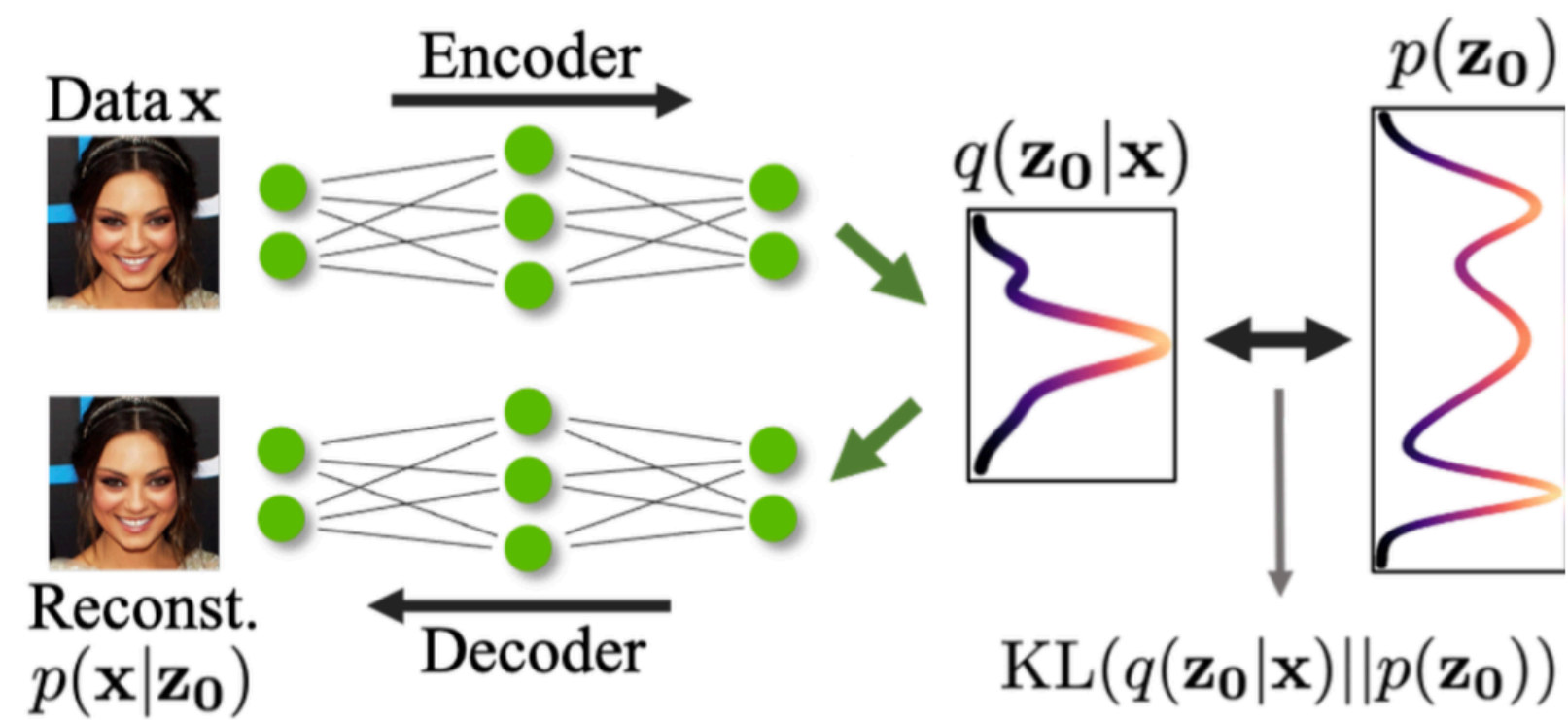
[Lu et al. 2022]



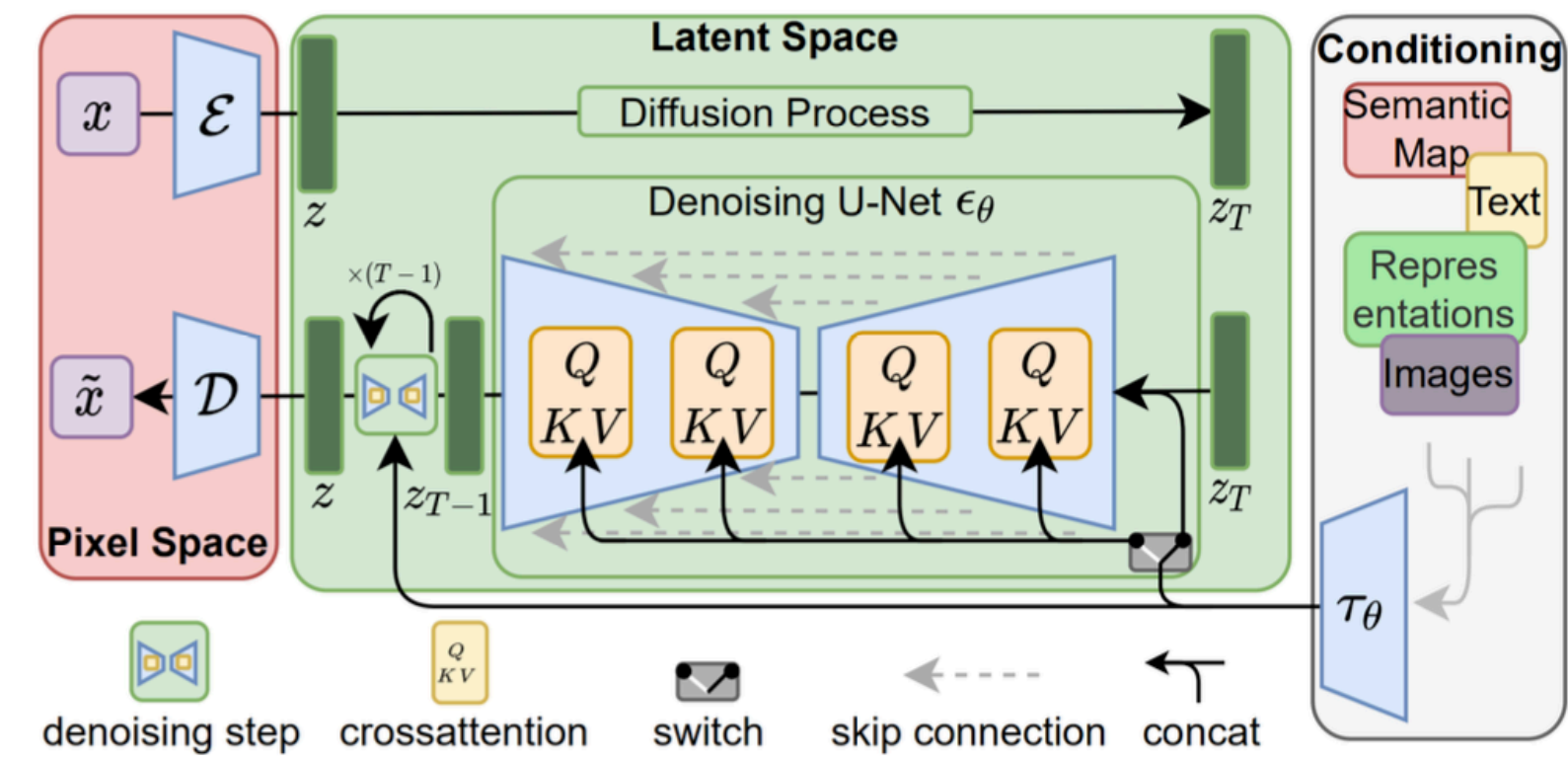
- Higher order methods:
RK4, Heun ...
- Exponential ODE Integrators:
DPM Solver ...
- ...

Sampling acceleration

2. Latent space



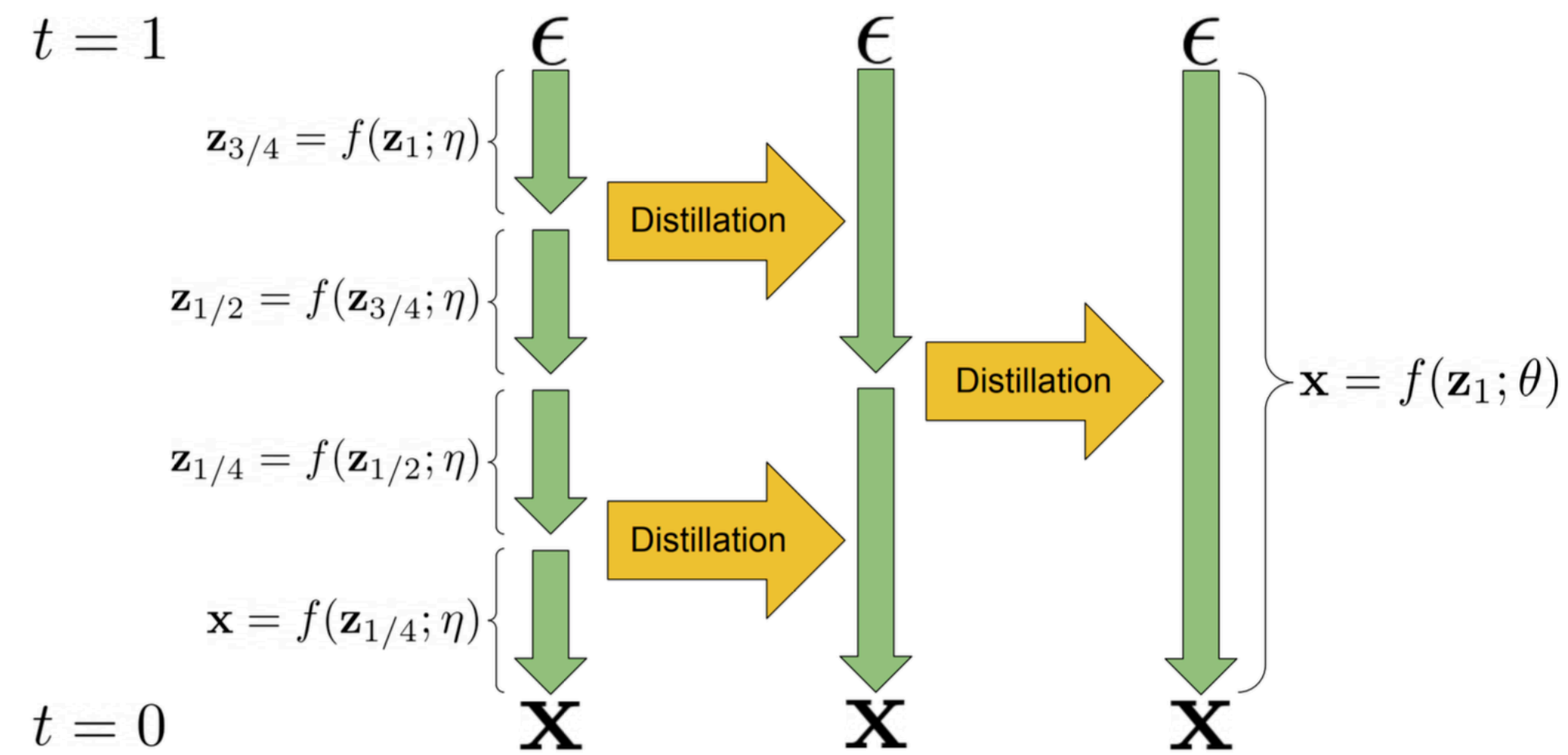
[Vahdat et al. 2021]



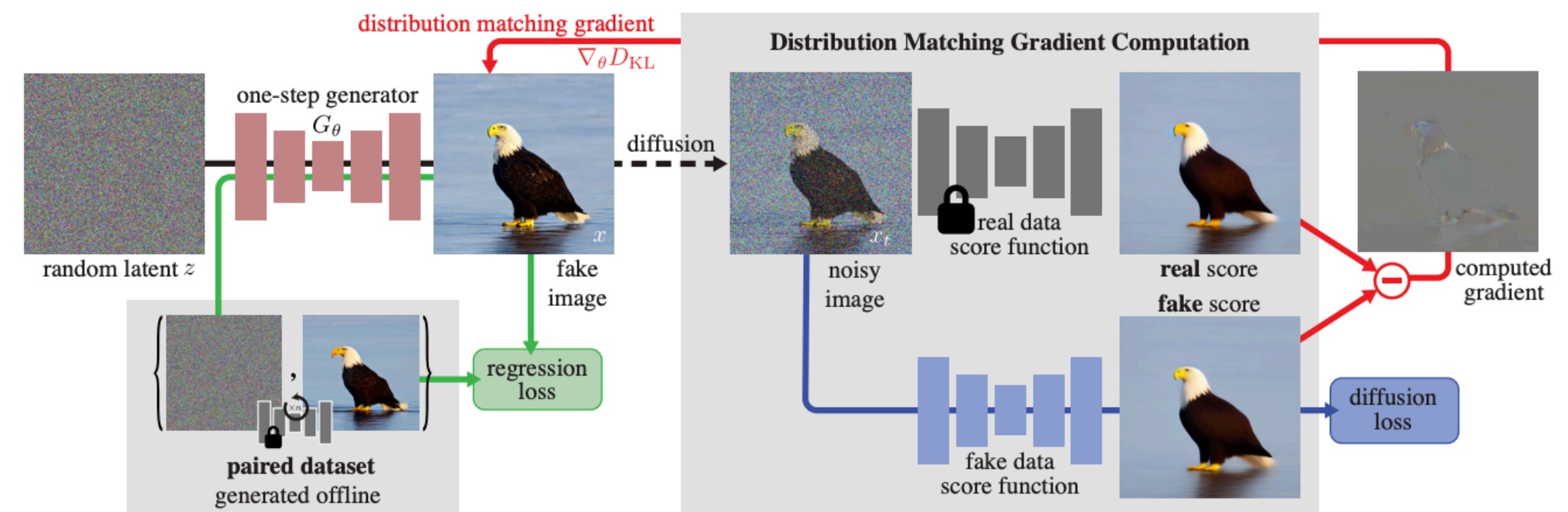
[Rombach et al. 2021]

Sampling acceleration

3. Distillation



[Salimans et al. 2022]



[Yin et al. 2024]

Conditional generation

1. Make the network conditional

- Simply, we learn $\nabla_{x_t} \log q(x_t | c)$
- Training: $E_{x_0, t, \varepsilon, c} [w_t | | \varepsilon_{\theta}(x_t, t, c) - \varepsilon | |^2]$

Conditional generation

2. Classifier-based guidance

- $q(x_t | c)q(c) = q(x_t)q(c | x_t)$. Take log and gradient to x_t :

$$\nabla_{x_t} \log q(x_t | c) = \nabla_{x_t} \log q(x_t) + \nabla_{x_t} \log q(c | x_t)$$

- $\nabla_{x_t} \log q(x_t)$: unconditional model
- $\nabla_{x_t} \log q(c | x_t)$: a separate classifier (BP gradient to the input)

Conditional generation

3. Classifier-free guidance

- $\nabla_{x_t} \log q(x_t | c) = \nabla_{x_t} \log q(x_t) + \nabla_{x_t} \log q(c | x_t)$
- $\nabla_{x_t} \log q(c | x_t) = \nabla_{x_t} \log q(x_t | c) - \nabla_{x_t} \log q(x_t)$
- $s_{\theta}^{CFG}(x_t, t, c) = s_{\theta}(x_t, t, \emptyset) + CFG \cdot (s_{\theta}(x_t, t, c) - s_{\theta}(x_t, t, \emptyset))$
- During training, condition c is randomly dropped (e.g., 10%)

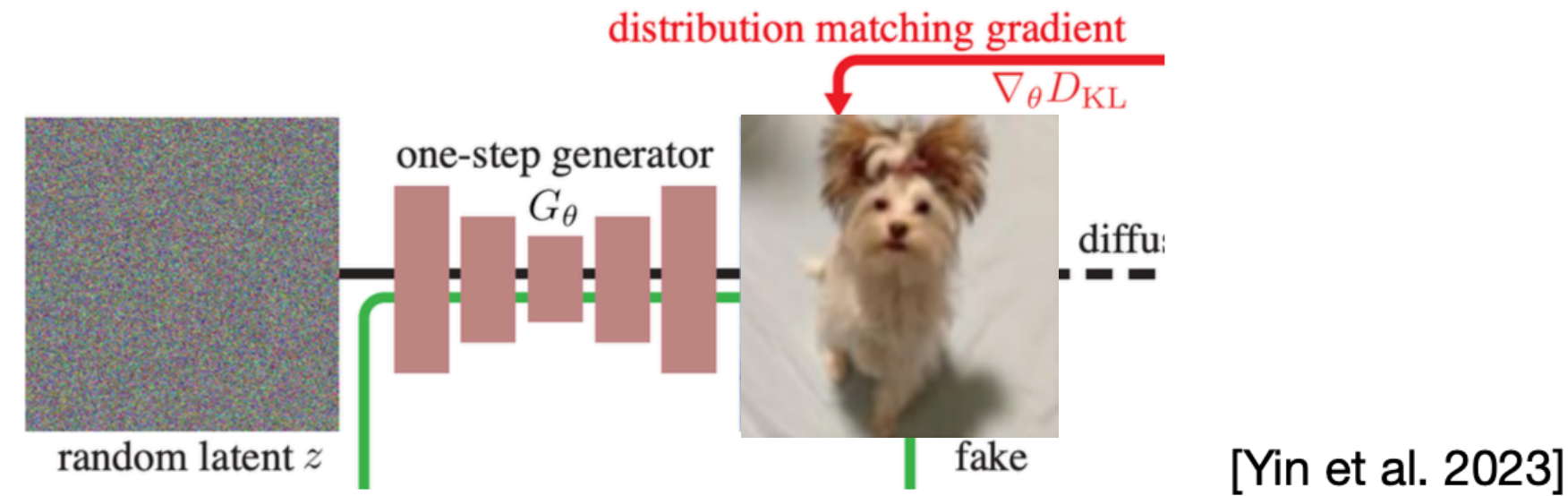
Score distillation: one-step generation

[Poole et al. 2022]

[Wang et al. 2023]

[Luo et al. 2023]

[Yin et al. 2023]



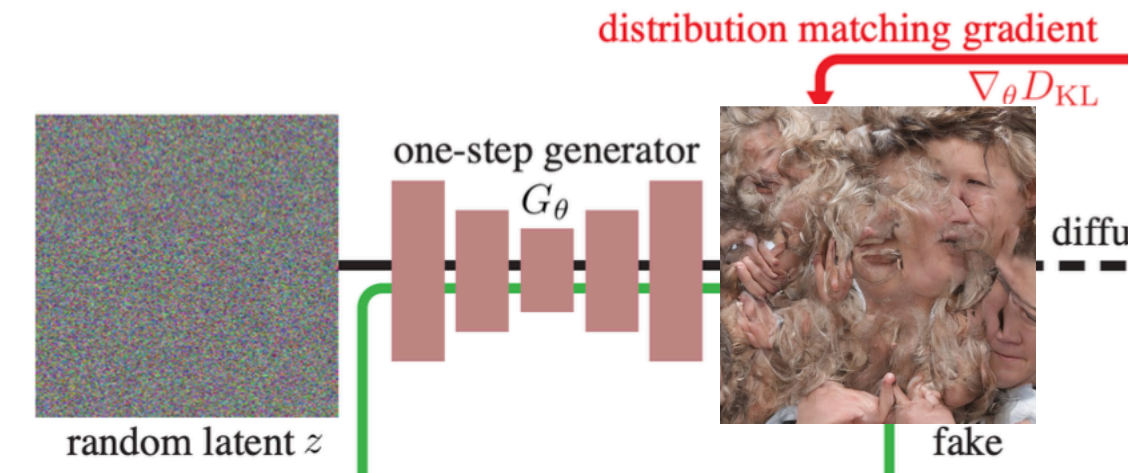
$$\begin{aligned}
 D_{KL}(p_{\text{fake}} \parallel p_{\text{real}}) &= \mathbb{E}_{x \sim p_{\text{fake}}} \left(\log \left(\frac{p_{\text{fake}}(x)}{p_{\text{real}}(x)} \right) \right) \\
 &= \mathbb{E}_{\substack{z \sim \mathcal{N}(0; \mathbf{I}) \\ x = G_{\theta}(z)}} \left[-(\log p_{\text{real}}(x) - \log p_{\text{fake}}(x)) \right].
 \end{aligned}
 \tag{1}$$

$$\nabla_{\theta} D_{KL} = \mathbb{E}_{\substack{z \sim \mathcal{N}(0; \mathbf{I}) \\ x = G_{\theta}(z)}} \left[- (s_{\text{real}}(x) - s_{\text{fake}}(x)) \nabla_{\theta} G_{\theta}(z) \right],
 \tag{2}$$

where $s_{\text{real}}(x) = \nabla_x \log p_{\text{real}}(x)$, $s_{\text{fake}}(x) = \nabla_x \log p_{\text{fake}}(x)$

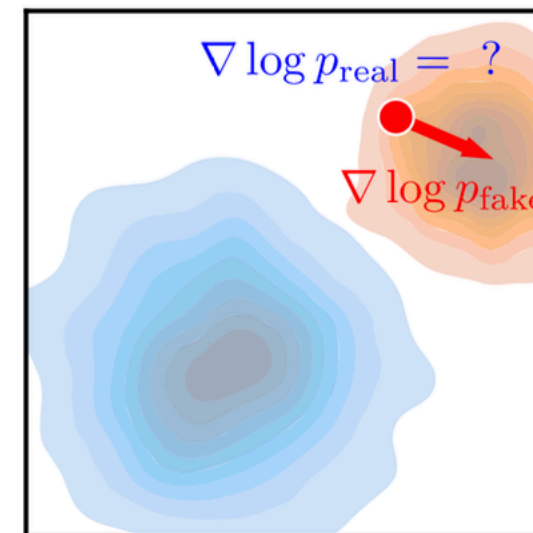
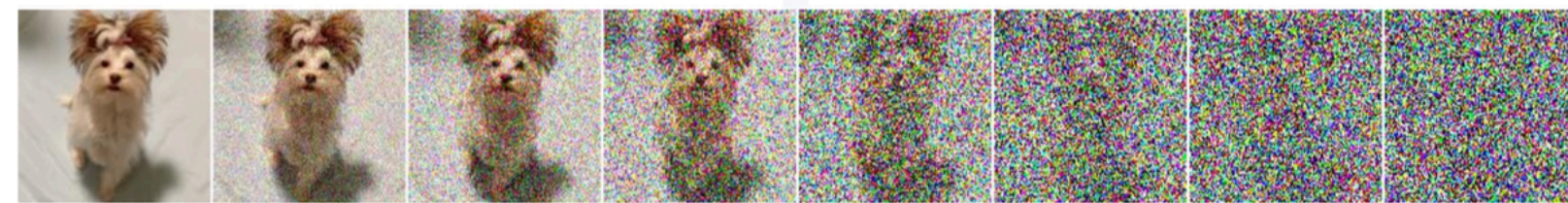
Score distillation

$$\nabla_{\theta} D_{KL} = \mathbb{E}_{\substack{z \sim \mathcal{N}(0; \mathbf{I}) \\ x = G_{\theta}(z)}} \left[- (s_{\text{real}}(x) - s_{\text{fake}}(x)) \nabla_{\theta} G_{\theta}(z) \right]$$



Is this all we need?

Issue: for a bad generated x , the score function may not be available, as it is OOD of $\nabla_{x_t} \log p(x_t)$ at any t



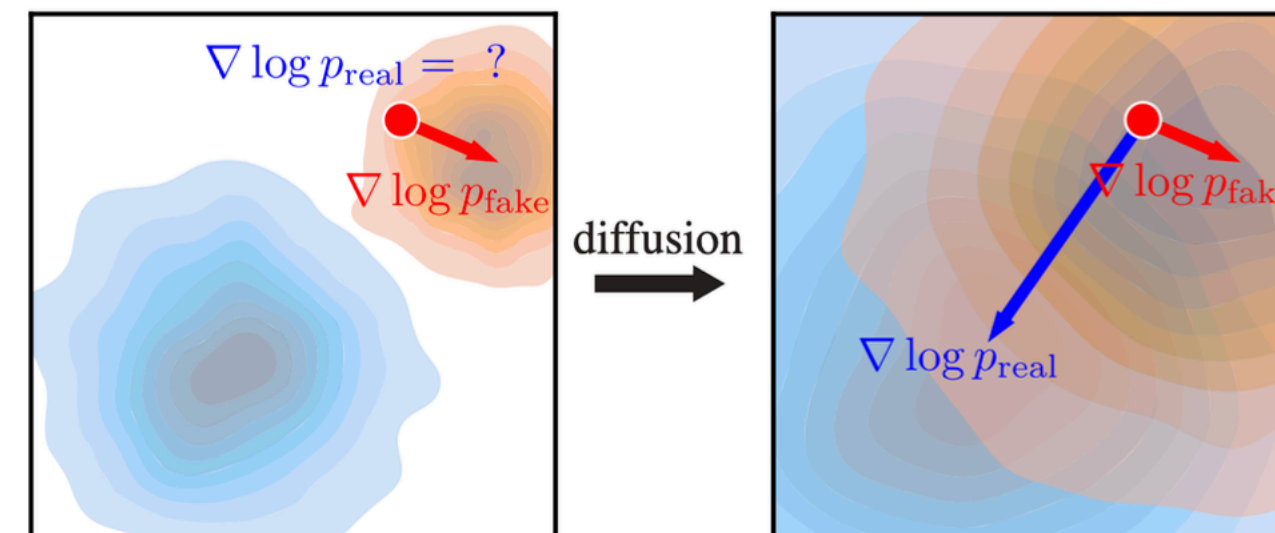
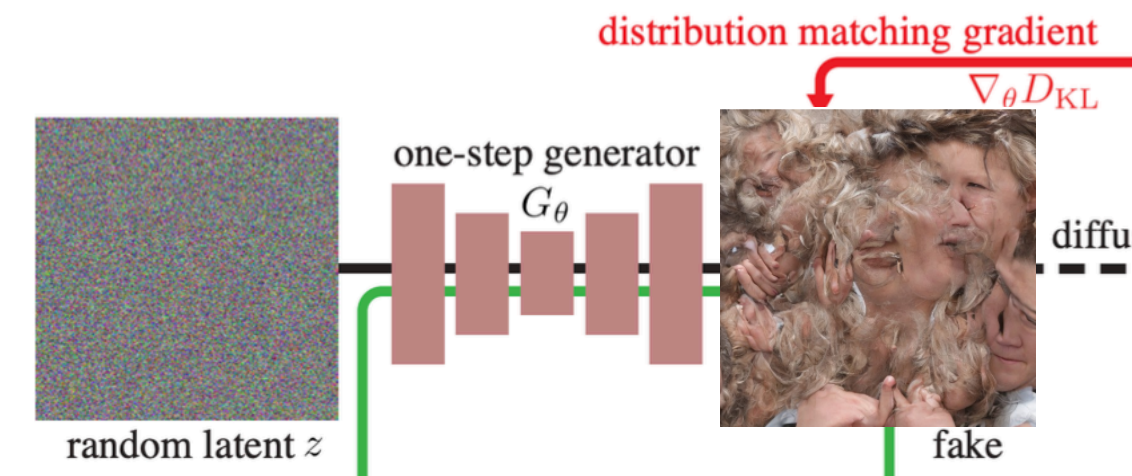
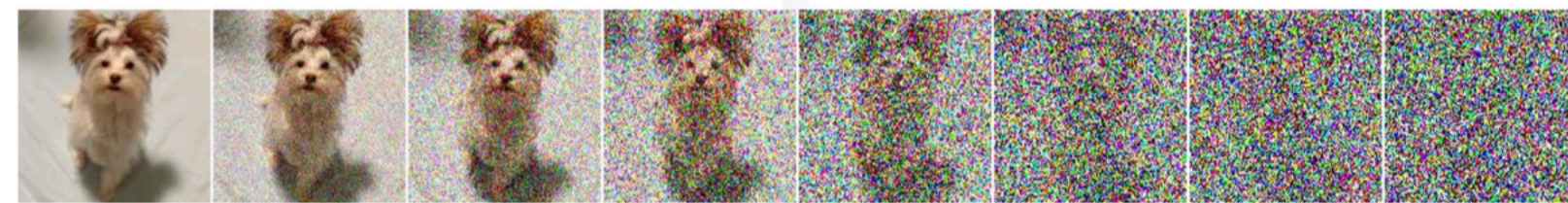
Score distillation

$$\nabla_{\theta} D_{KL} = \mathbb{E}_{\substack{z \sim \mathcal{N}(0; \mathbf{I}) \\ x = G_{\theta}(z)}} \left[- (s_{\text{real}}(x) - s_{\text{fake}}(x)) \nabla_{\theta} G_{\theta}(z) \right]$$

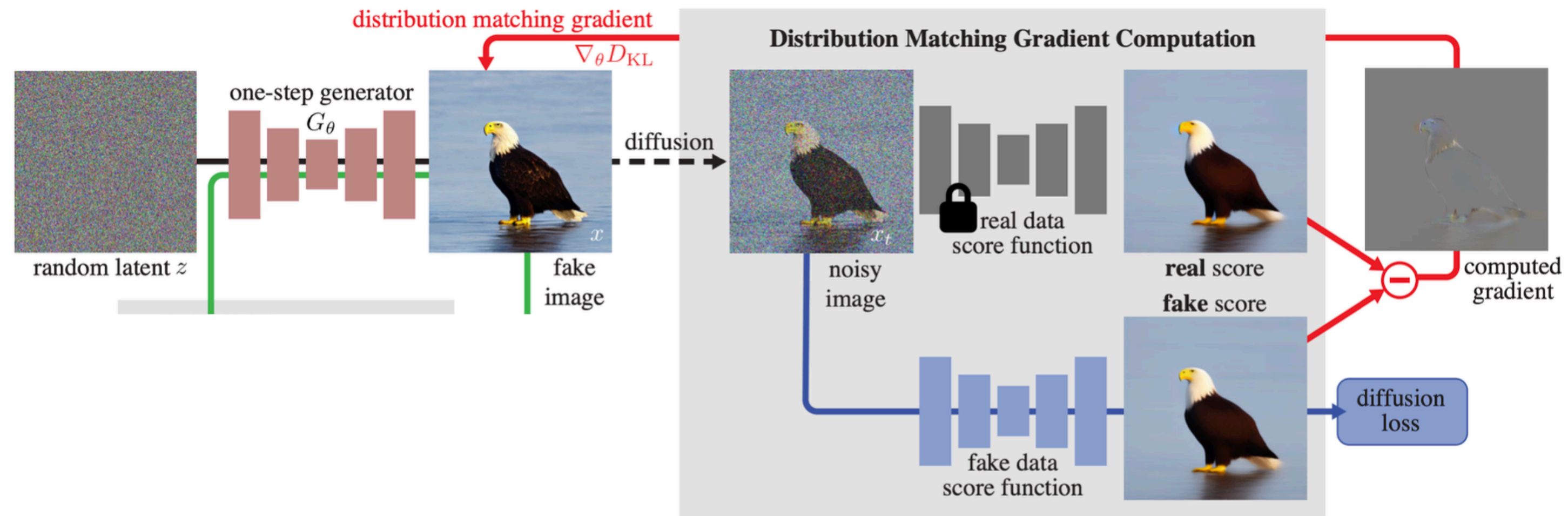


Compute on diffused distributions:

$$\nabla_{\theta} D_{KL} \simeq \mathbb{E}_{z, t, x, x_t} [w_t \alpha_t (s_{\text{fake}}(x_t, t) - s_{\text{real}}(x_t, t)) \nabla_{\theta} G_{\theta}(z)], \quad (7)$$



Score distillation



Generator: G_θ

Fake network: s_ϕ^{fake}

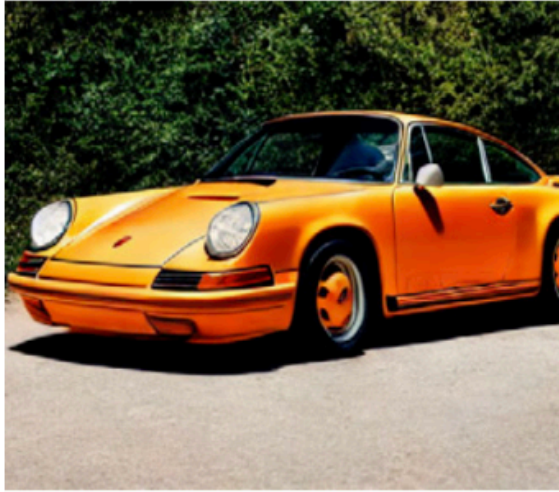
$$\nabla_\theta D_{KL} \simeq \mathbb{E}_{z,t,x,x_t} [w_t \alpha_t (s_{fake}(x_t, t) - s_{real}(x_t, t)) \nabla_\theta G_\theta(z)], \quad (7)$$

A loss to make some representation follow some distribution (with DM scores)

“macro photo of a miniature toy sloth drinking a soda, shot on a light pastel cyclorama”



“a high-resolution photo of an orange Porsche under sunshine”



“a hot air balloon in shape of a heart. Grand Canyon”



“Astronaut on a camel on mars”



DMD (ours, 1 step)
90ms

InstaFlow (1 step)
90ms

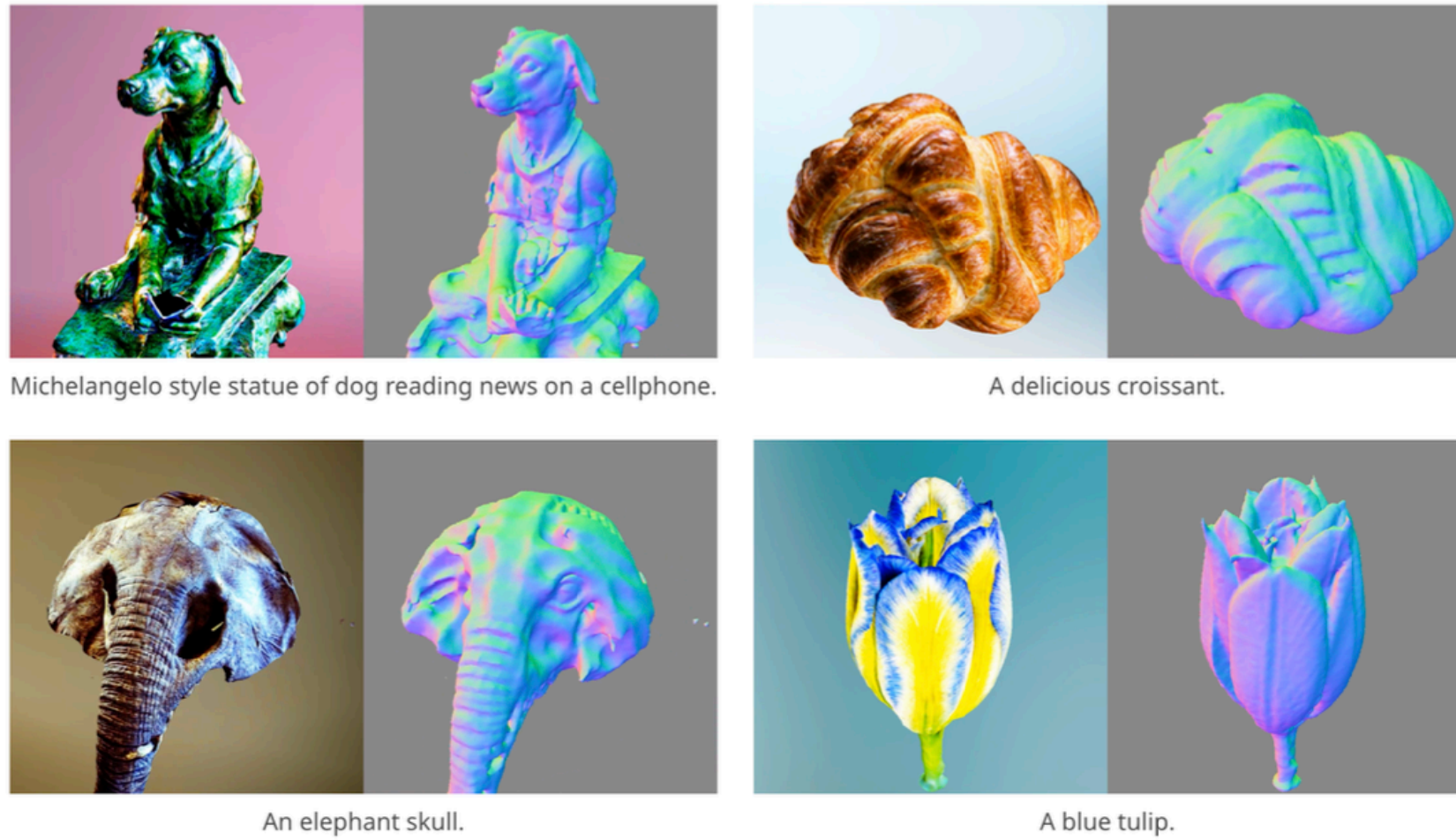
LCM (1 step)
90ms

LCM (2 steps)
120ms

DPM++ (4 steps)
260ms

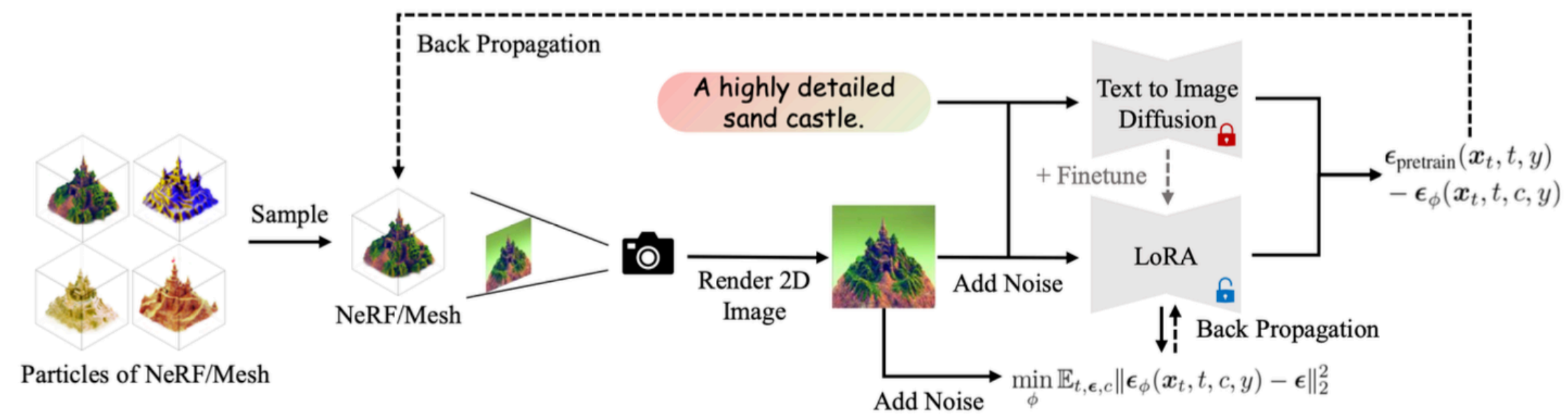
SD (50 steps)
2590ms

Score distillation: text-to-3D



Optimize a learnable 3D representation, such that the rendered views follow the text-to-image distribution (pretrained DM)

[Wang et al. 2023]



Summary

- **Basics:** discrete time, continuous time (SDE / ODE)
- **Applications:** acceleration, conditional generation
- Score distillation
- (Topics not covered: DreamBooth, engineering ...)