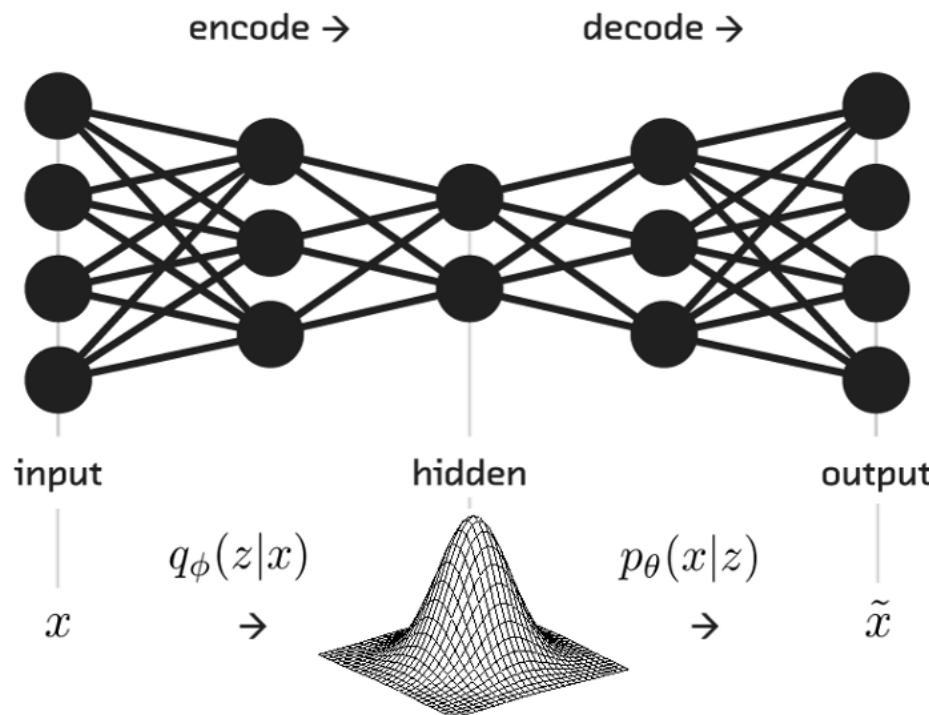


Diffusion Model Review

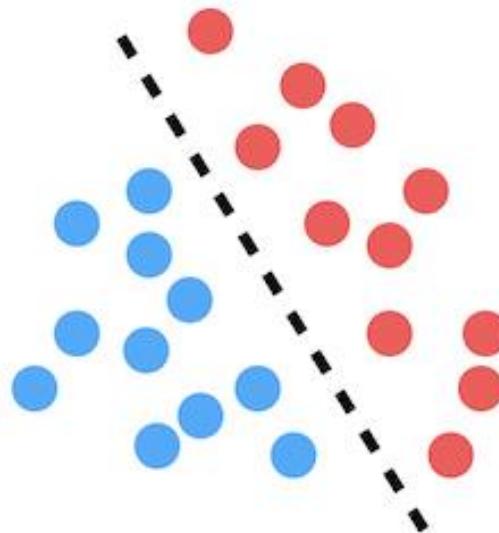
CAO Hanqun

Generative Model

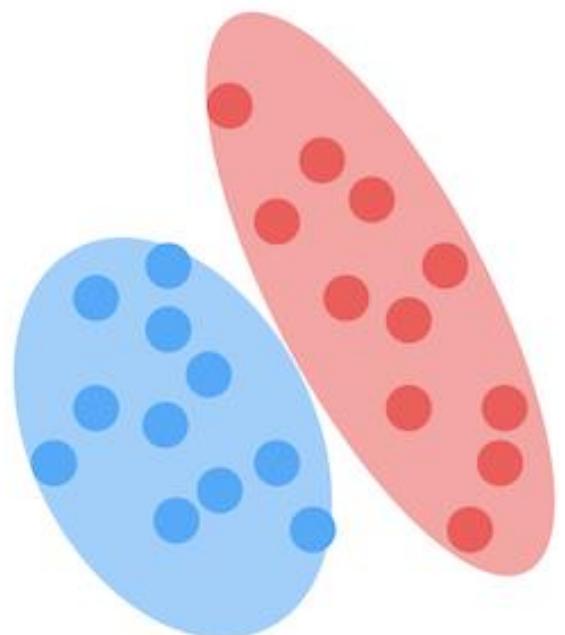
- Given data x , generating samples following distribution $p(x)$



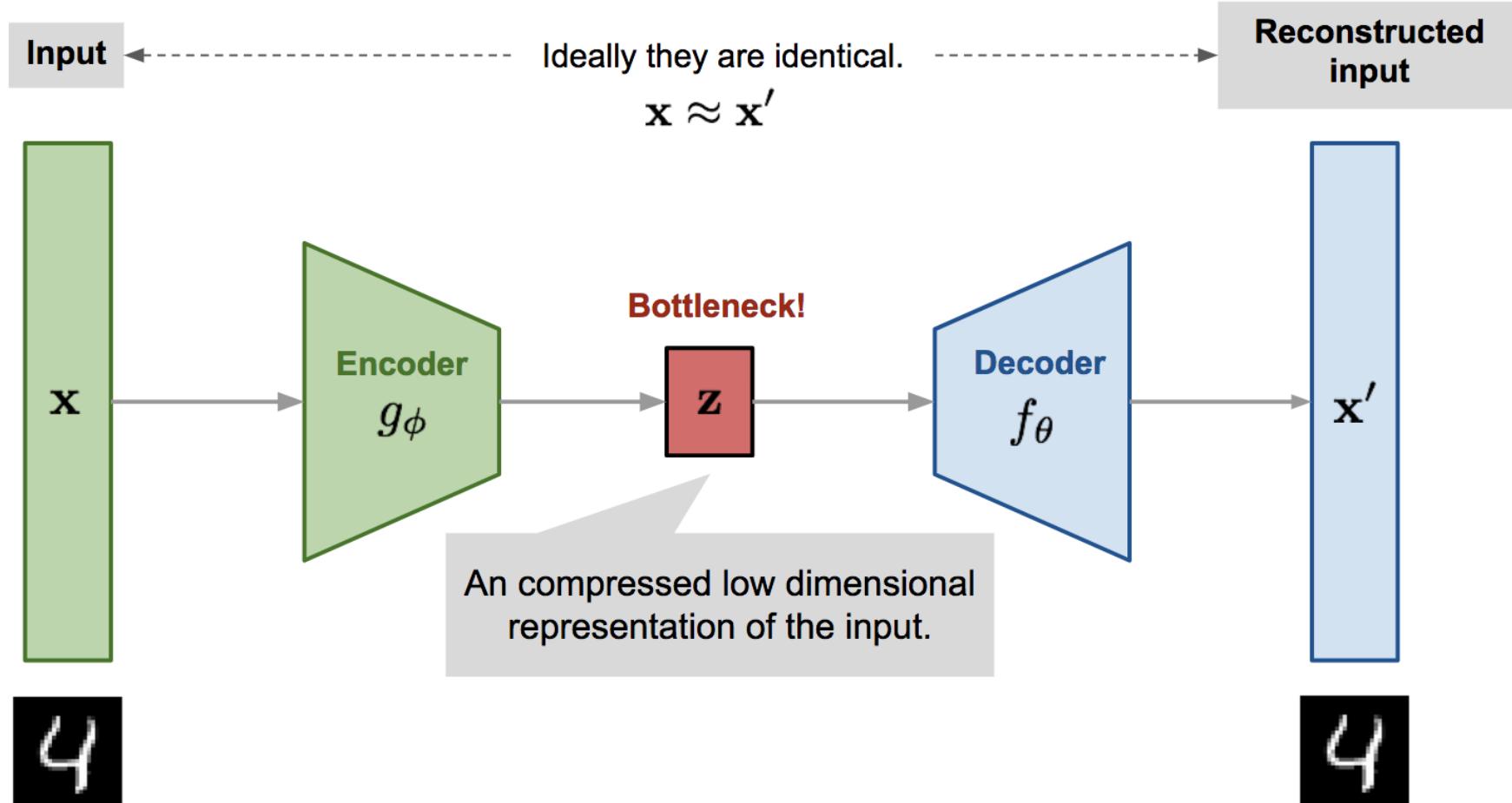
Discriminative



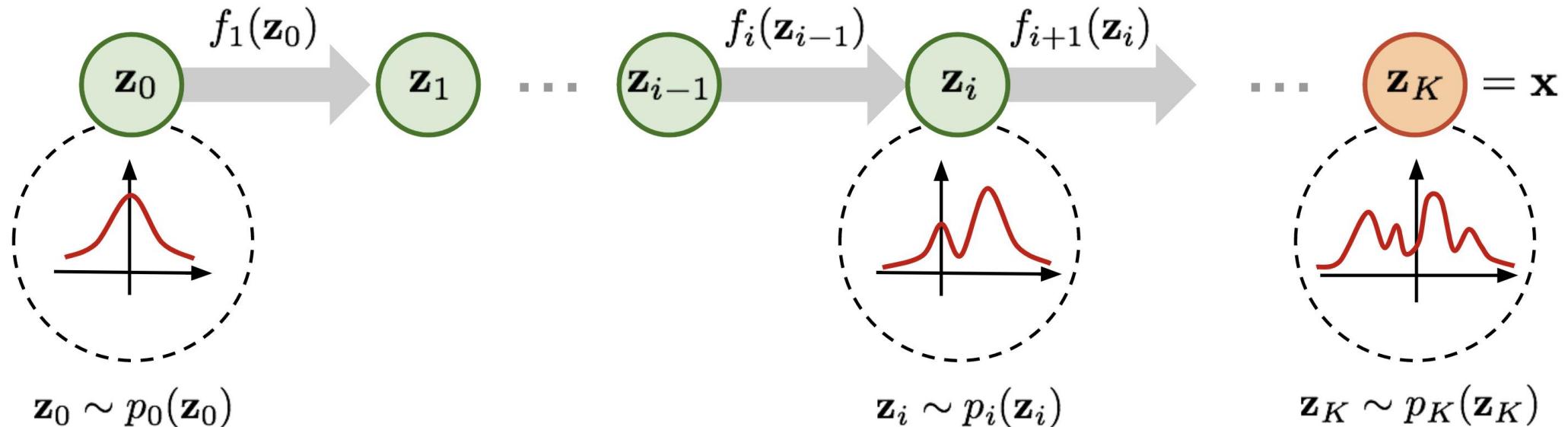
Generative



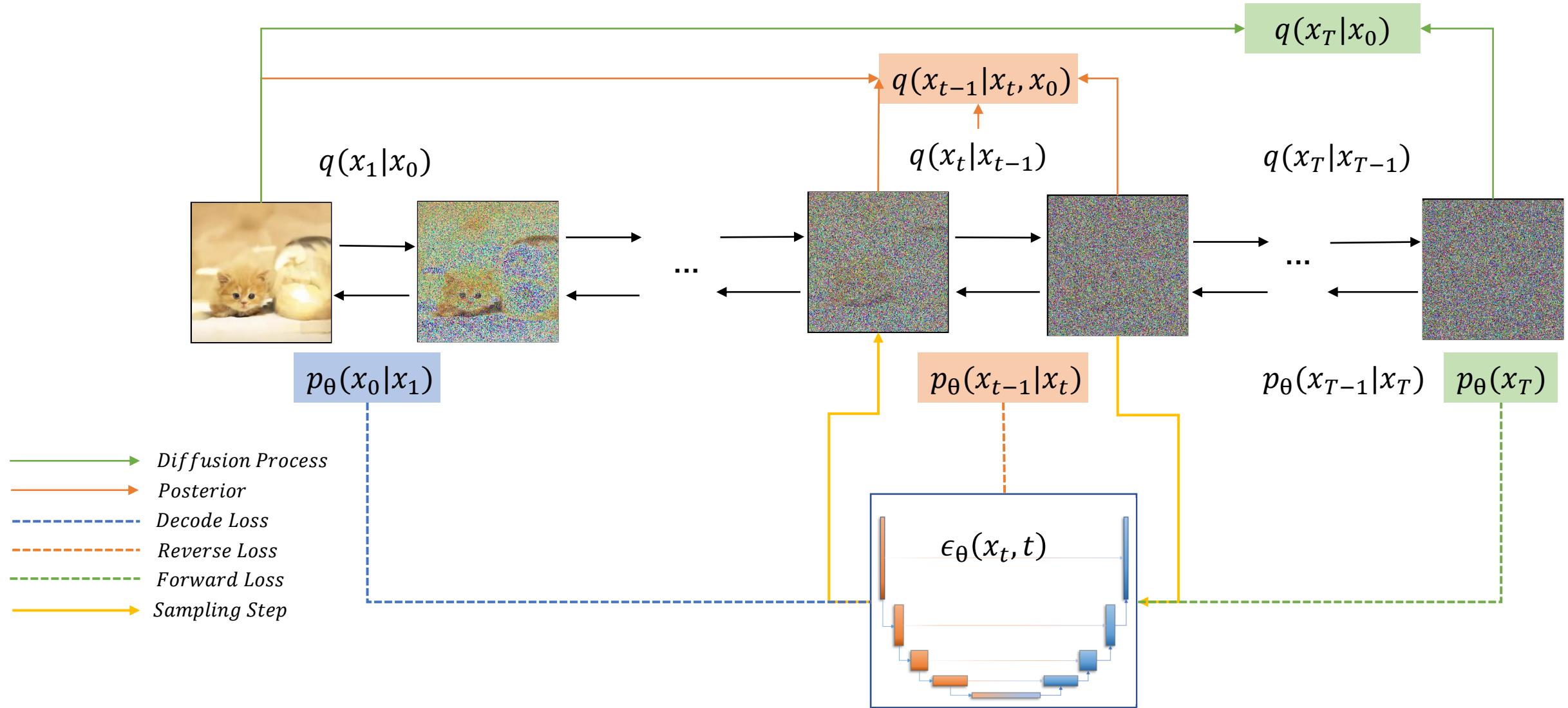
From VAE to NF, to Diffusion Model



From VAE to NF, to Diffusion Model



From VAE to NF, to Diffusion Model



Diffusion Model Math Formulation

- Forward Process $F(x, \sigma) = F_T(x_{T-1}, \sigma_T) \cdots \circ F_t(x_{t-1}, \sigma_t) \cdots \circ F_1(x_0, \sigma_1)$

- Reverse Process $R(x, \sigma) = R_1(x_1, \sigma_1) \cdots \circ R_t(x_t, \sigma_t) \cdots \circ R_T(x_T, \sigma_T)$

- Transition Kernels $x_t = F_t(x_{t-1}, \sigma_t), \quad x_{t-1} = R_t(x_t, \sigma_t)$

$$\tilde{x}_0 = [R_1(x_1, \sigma_1) \cdots \circ R_t(x_t, \sigma_t) \cdots \circ R_T(x_T, \sigma_T)] \circ [F_T(x_{T-1}, \sigma_T) \cdots \circ F_t(x_{t-1}, \sigma_t) \cdots \circ F_1(x_0, \sigma_1)]$$

- Noising & Denoising Process

$$\mathbb{E}_{F(x_0, \sigma)} [-\log R(x_T, \tilde{\sigma})]$$

- Training Objective based on MLE

Landmark Works: From discrete to Continuous

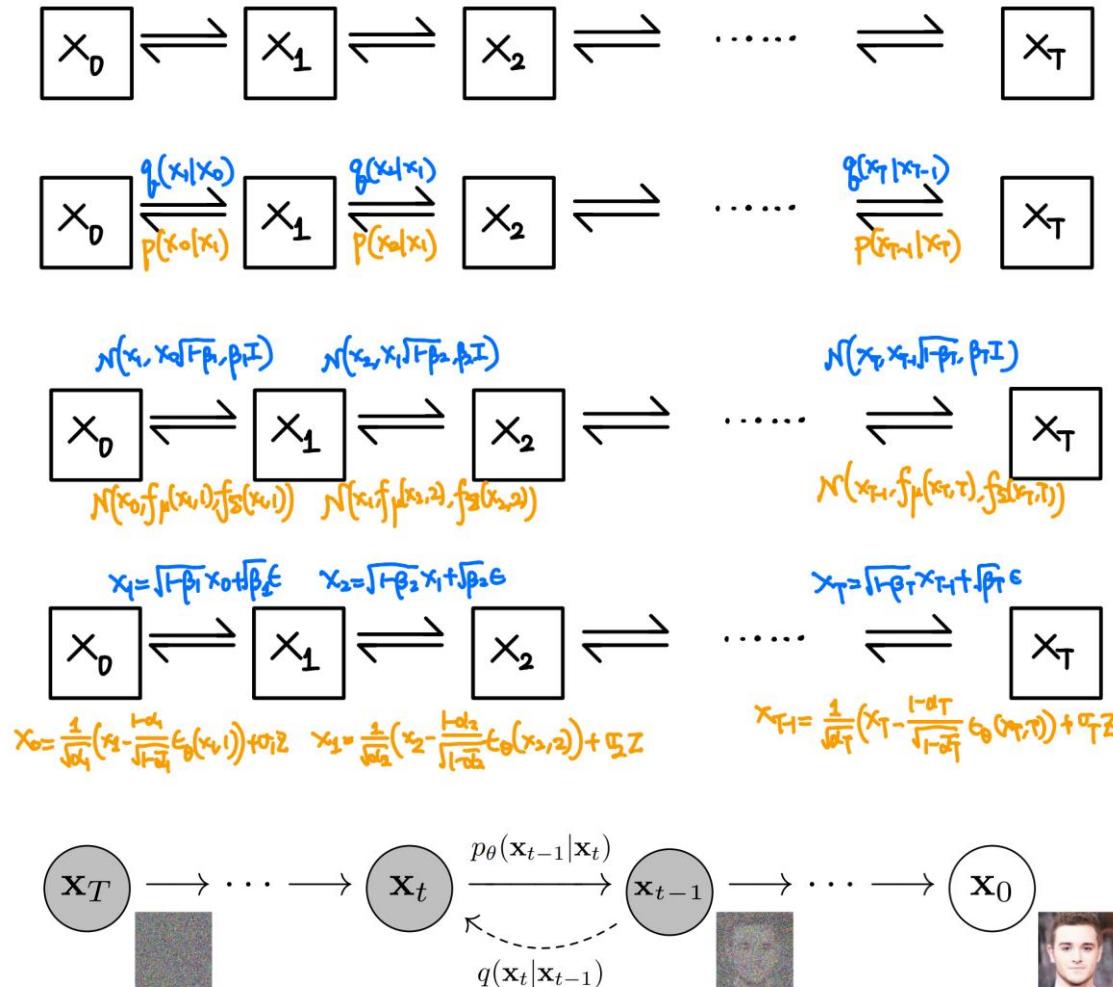


Figure 2: The directed graphical model considered in this work.

Algorithm 1 Training

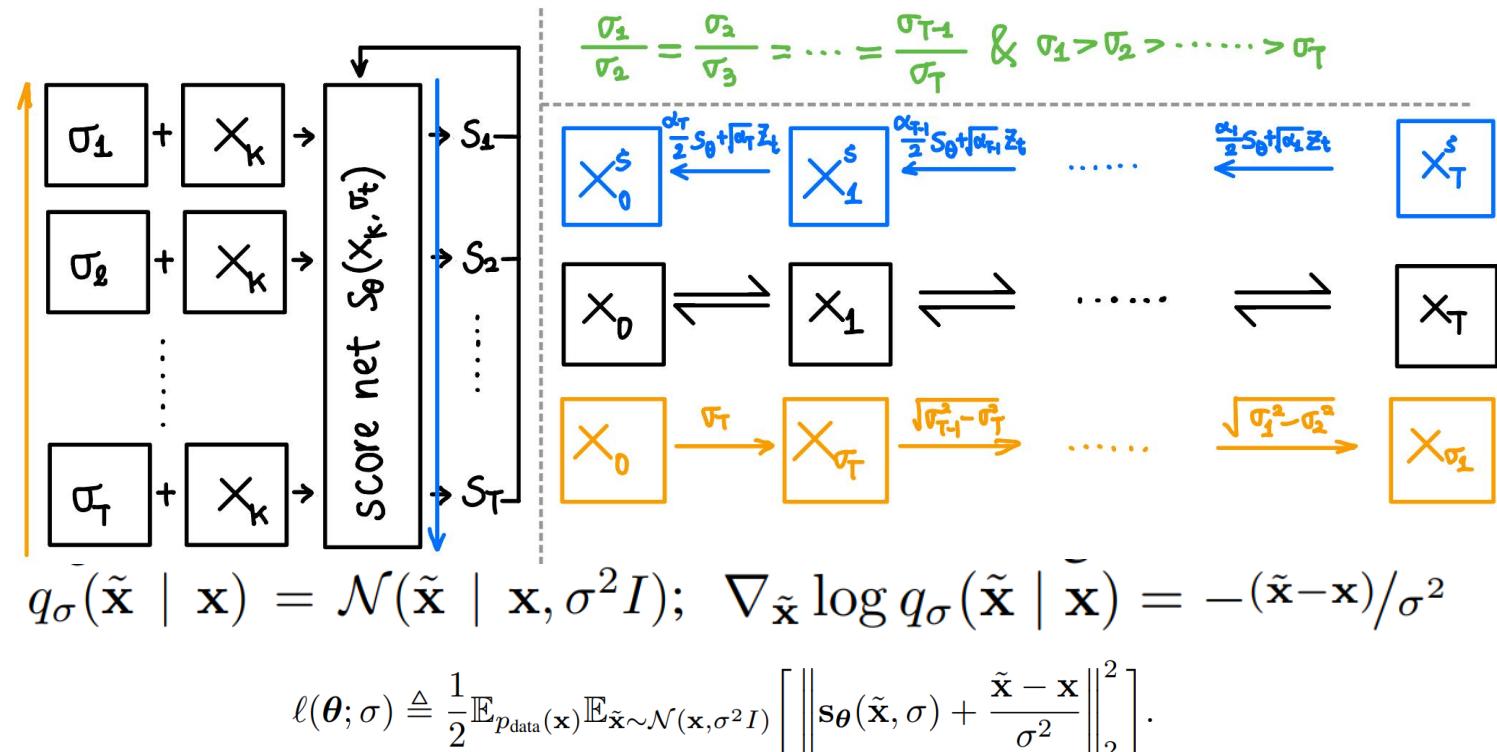
- 1: **repeat**
 - 2: $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
 - 3: $t \sim \text{Uniform}(\{1, \dots, T\})$
 - 4: $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
 - 5: Take gradient descent step on

$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$
 - 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
-

Landmark Works: From discrete to Continuous



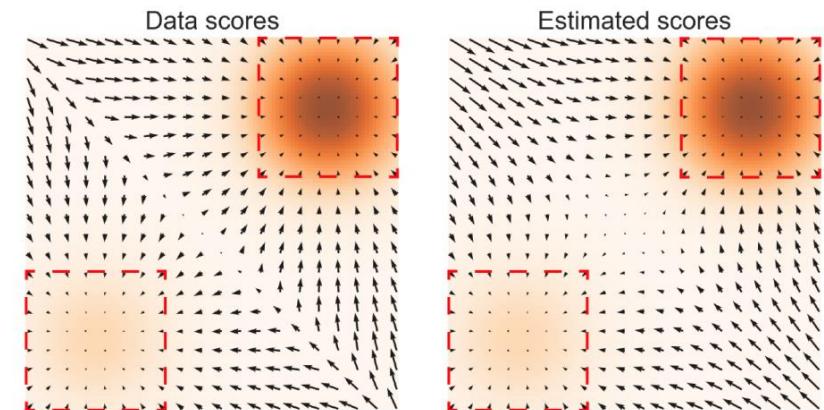
Method	Loss
\mathcal{L}_{ESM}	$\frac{1}{2} \mathbb{E}[\mathbf{s}_\theta(Y_s, s) - \nabla \log q(Y_s) _\Lambda^2]$
\mathcal{L}_{ISM}	$\mathbb{E}[\frac{1}{2} \mathbf{s}_\theta(Y_s, s) _\Lambda^2 + \nabla \cdot (\Lambda^\top \mathbf{s}_\theta)]$
\mathcal{L}_{SSM}	$\mathbb{E}[\frac{1}{2} \mathbf{s}_\theta(Y_s, s) _\Lambda^2 + v^\top \nabla(\Lambda^\top \mathbf{s}_\theta)v]$
\mathcal{L}_{DSM}	$\frac{1}{2} \mathbb{E}[\mathbf{s}_\theta(Y_s, s) - \nabla \log q(Y_s Y_0) _\Lambda^2]$

Algorithm 1 Annealed Langevin dynamics [1]

```

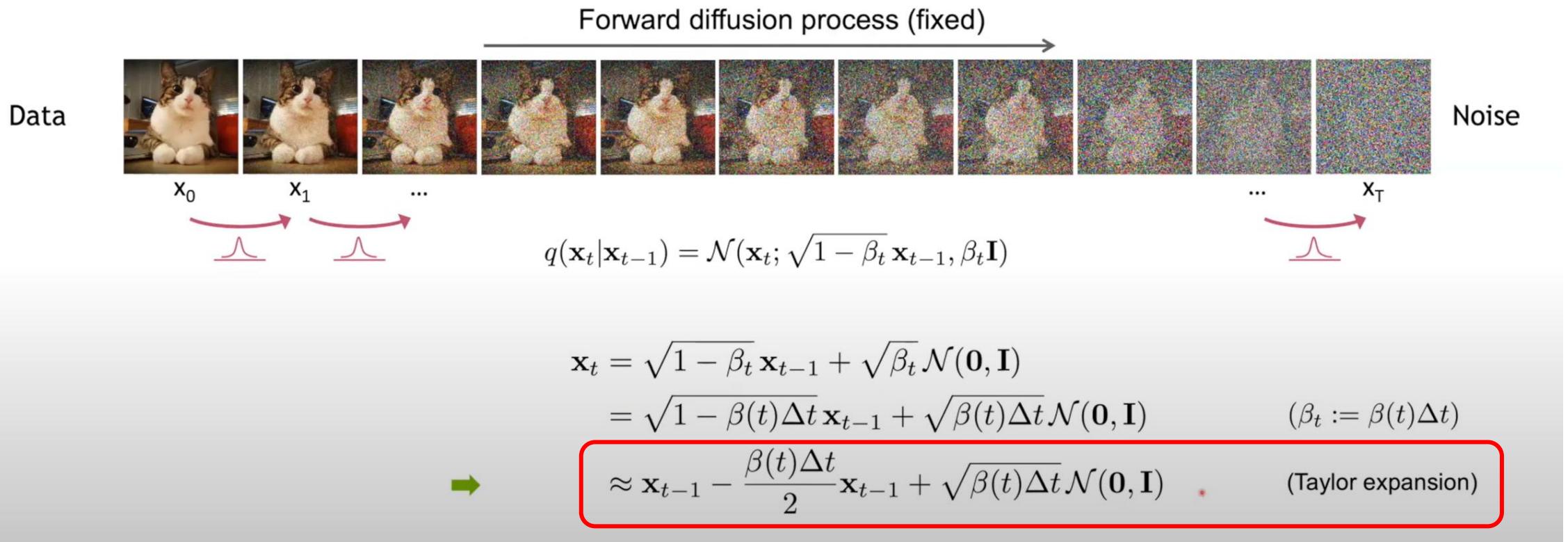
Require:  $\{\sigma_i\}_{i=1}^L, \epsilon, T.$ 
1: Initialize  $\mathbf{x}_0$ 
2: for  $i \leftarrow 1$  to  $L$  do
3:    $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$   $\triangleright \alpha_i$  is the step size.
4:   for  $t \leftarrow 1$  to  $T$  do
5:     Draw  $\mathbf{z}_t \sim \mathcal{N}(0, I)$ 
6:      $\mathbf{x}_t \leftarrow \mathbf{x}_{t-1} + \alpha_i \mathbf{s}_\theta(\mathbf{x}_{t-1}, \sigma_i) + \sqrt{2\alpha_i} \mathbf{z}_t$ 
7:    $\mathbf{x}_0 \leftarrow \mathbf{x}_T$ 
8: if denoise  $\mathbf{x}_T$  then
9:   return  $\mathbf{x}_T + \sigma_T^2 \mathbf{s}_\theta(\mathbf{x}_T, \sigma_T)$ 
10: else
11:   return  $\mathbf{x}_T$ 

```



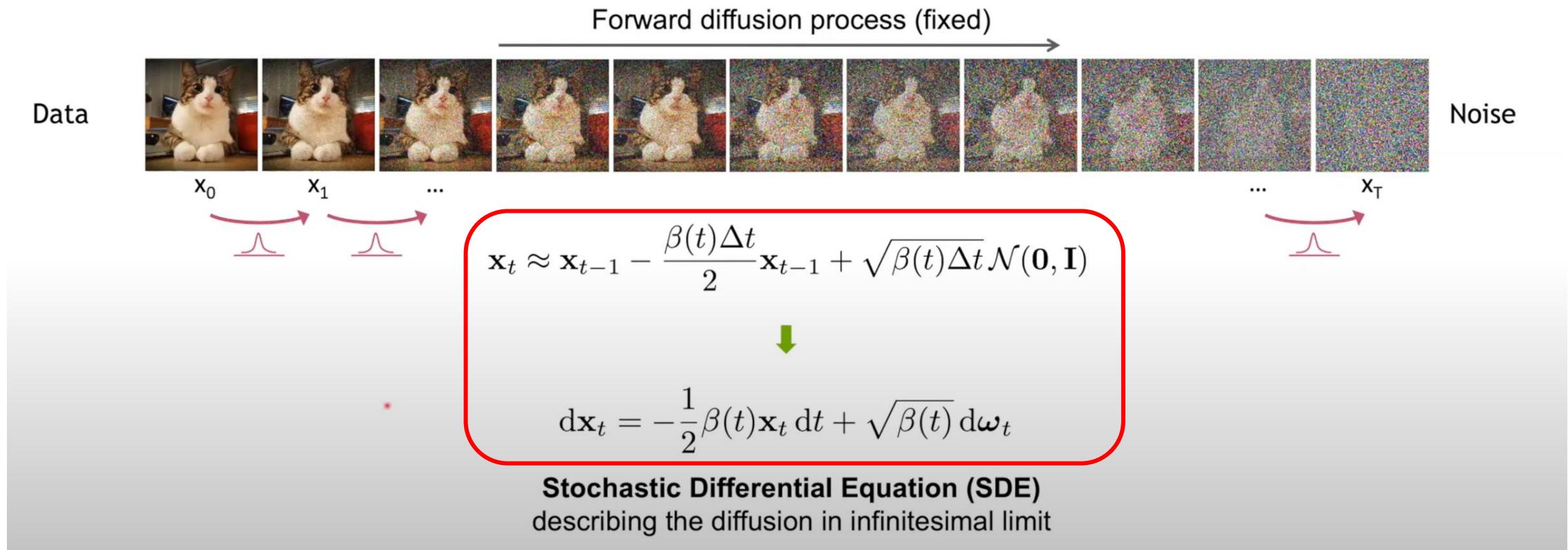
Landmark Works: From discrete to Continuous

Consider the limit of many small steps:

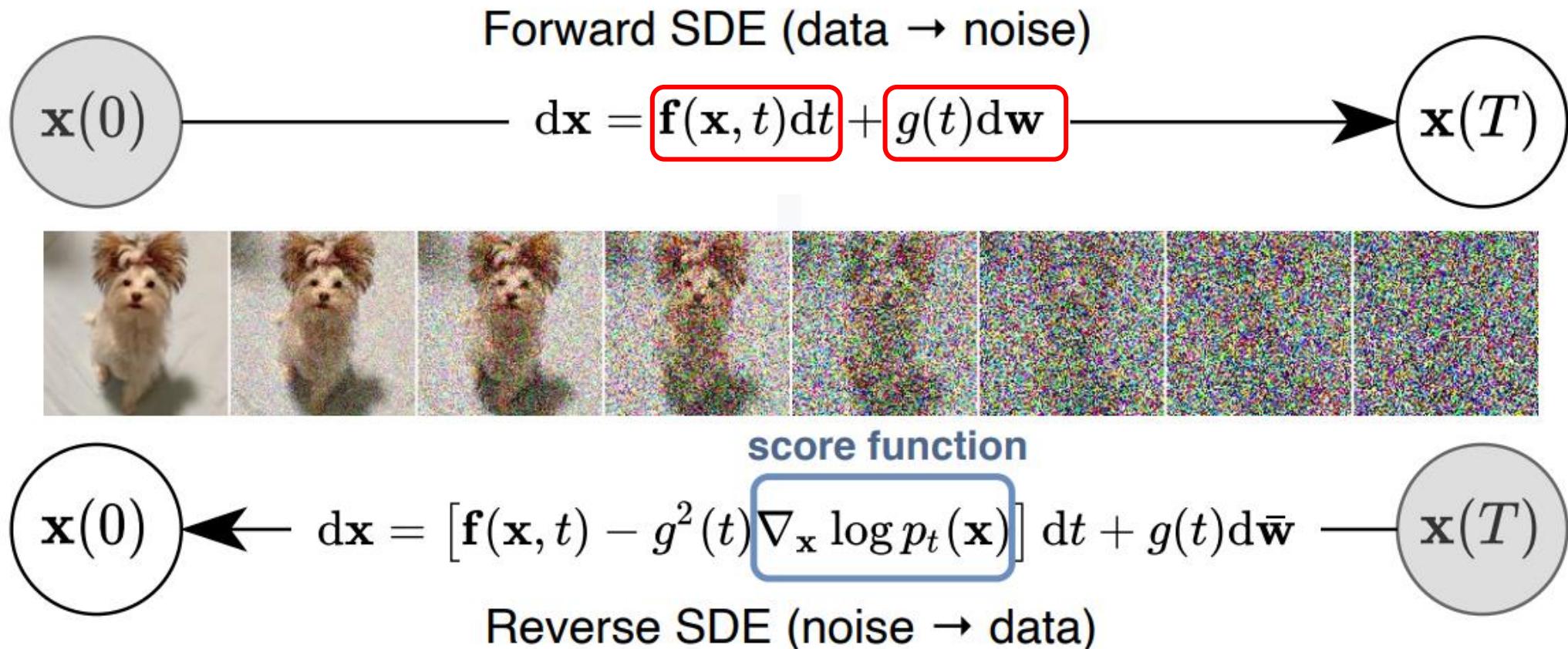


Landmark Works: From discrete to Continuous

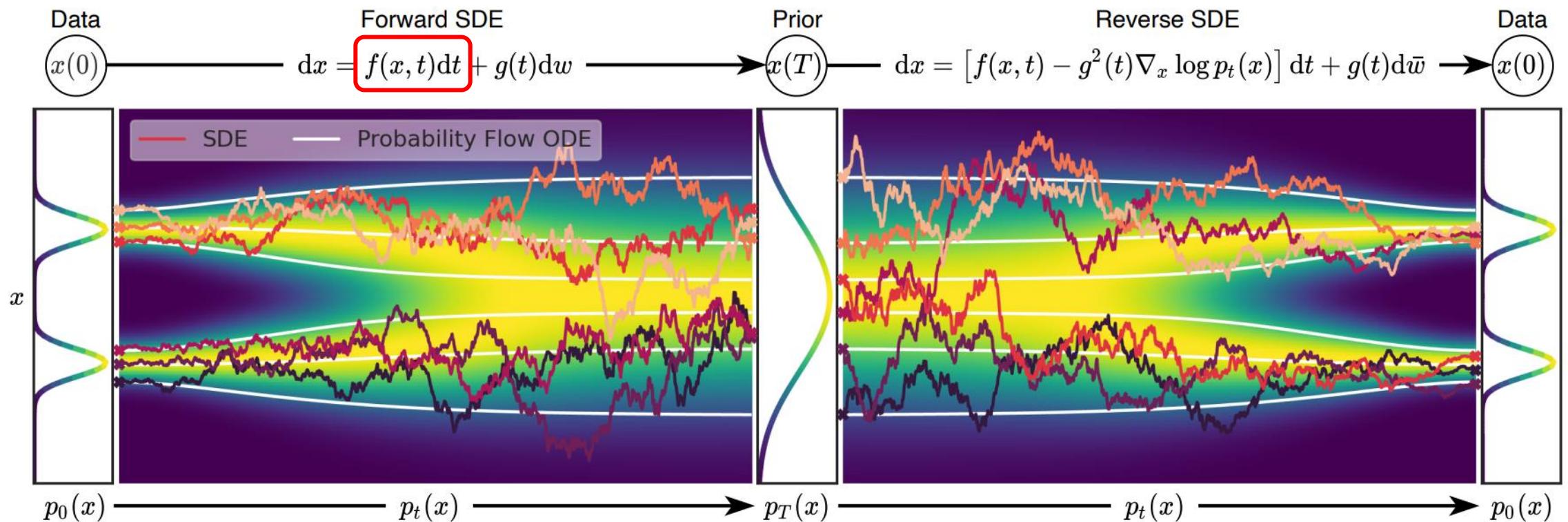
Consider the limit of **many small steps**:



Landmark Works: From discrete to Continuous



Landmark Works: From discrete to Continuous



- SDE: Higher Performance
- ODE: Higher Speed

Landmark Works: From discrete to Continuous

$$p_{0t}(\mathbf{x}(t) \mid \mathbf{x}(0)) = \begin{cases} \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0), [\sigma^2(t) - \sigma^2(0)]\mathbf{I}), & (\text{VE SDE}) \\ \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0)e^{-\frac{1}{2}\int_0^t \beta(s)ds}, \mathbf{I} - \mathbf{I}e^{-\int_0^t \beta(s)ds}) & (\text{VP SDE}) \\ \mathcal{N}(\mathbf{x}(t); \mathbf{x}(0)e^{-\frac{1}{2}\int_0^t \beta(s)ds}, [1 - e^{-\int_0^t \beta(s)ds}]^2\mathbf{I}) & (\text{sub-VP SDE}) \end{cases}$$

Algorithm 2 PC sampling (VE SDE)

```

1:  $\mathbf{x}_N \sim \mathcal{N}(\mathbf{0}, \sigma_{\max}^2 \mathbf{I})$ 
2: for  $i = N - 1$  to 0 do
3:    $\mathbf{x}'_i \leftarrow \mathbf{x}_{i+1} + (\sigma_{i+1}^2 - \sigma_i^2)\mathbf{s}_{\theta*}(\mathbf{x}_{i+1}, \sigma_{i+1})$ 
4:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:    $\mathbf{x}_i \leftarrow \mathbf{x}'_i + \sqrt{\sigma_{i+1}^2 - \sigma_i^2}\mathbf{z}$ 
6:   for  $j = 1$  to  $M$  do
7:      $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
8:      $\mathbf{x}_i \leftarrow \mathbf{x}_i + \epsilon_i \mathbf{s}_{\theta*}(\mathbf{x}_i, \sigma_i) + \sqrt{2\epsilon_i}\mathbf{z}$ 
9:   return  $\mathbf{x}_0$ 
```

Algorithm 3 PC sampling (VP SDE)

1: $\mathbf{x}_N \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 2: for $i = N - 1$ to 0 do 3: $\mathbf{x}'_i \leftarrow (2 - \sqrt{1 - \beta_{i+1}})\mathbf{x}_{i+1} + \beta_{i+1}\mathbf{s}_{\theta*}(\mathbf{x}_{i+1}, i + 1)$ 4: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 5: $\mathbf{x}_i \leftarrow \mathbf{x}'_i + \sqrt{\beta_{i+1}}\mathbf{z}$	Predictor
6: for $j = 1$ to M do 7: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 8: $\mathbf{x}_i \leftarrow \mathbf{x}_i + \epsilon_i \mathbf{s}_{\theta*}(\mathbf{x}_i, i) + \sqrt{2\epsilon_i}\mathbf{z}$	Corrector
9: return \mathbf{x}_0	

-- Sampling along the reverse trajectory is actually finding the numerical solution of differential equations.

Landmark Works: Label_Conditional Diffusion

Diffusion Models Beat GANs on Image Synthesis

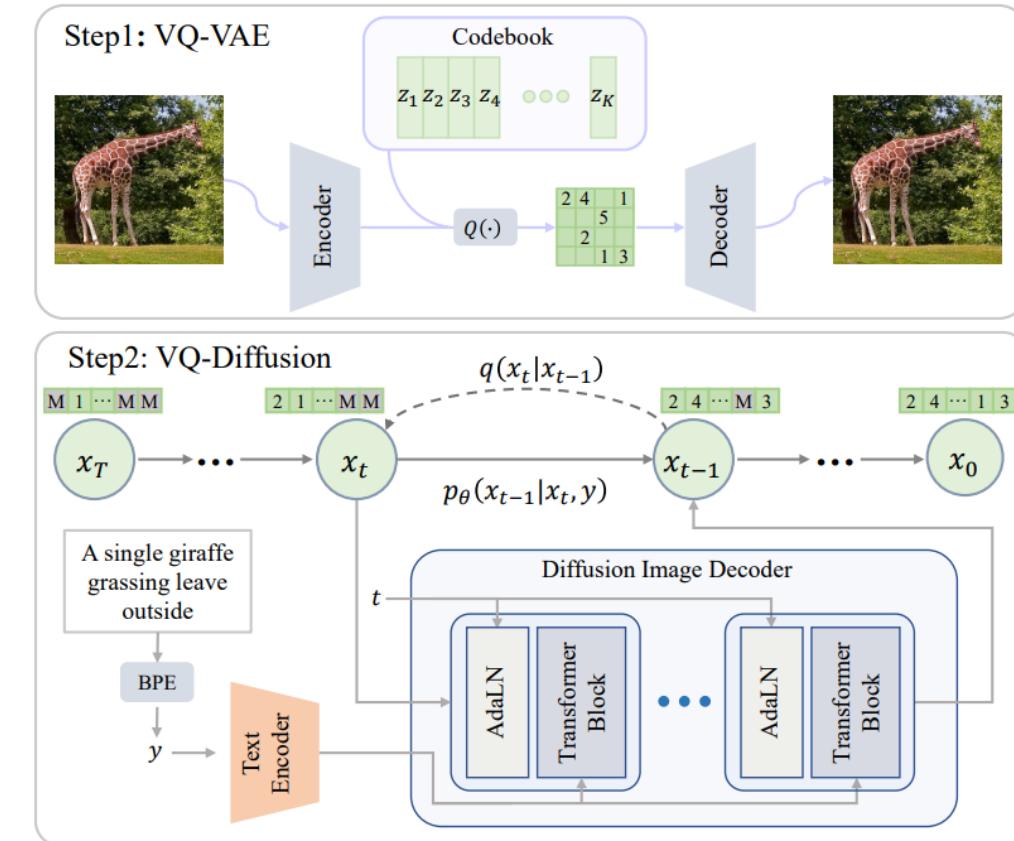
Algorithm 1 Classifier guided diffusion sampling, given a diffusion model $(\mu_\theta(x_t), \Sigma_\theta(x_t))$, classifier $p_\phi(y|x_t)$, and gradient scale s .

```

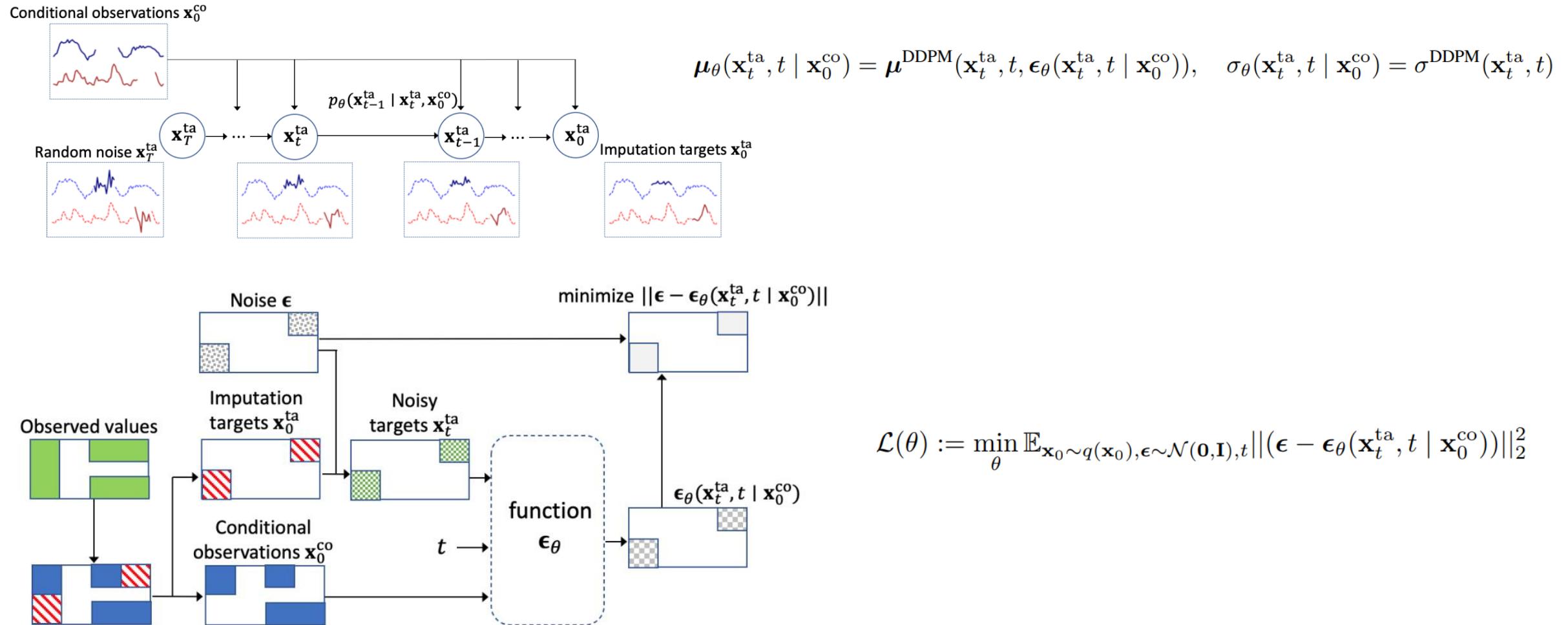
Input: class label  $y$ , gradient scale  $s$ 
 $x_T \leftarrow$  sample from  $\mathcal{N}(0, \mathbf{I})$ 
for all  $t$  from  $T$  to 1 do
     $\mu, \Sigma \leftarrow \mu_\theta(x_t), \Sigma_\theta(x_t)$ 
     $x_{t-1} \leftarrow$  sample from  $\mathcal{N}(\mu + s\Sigma \nabla_{x_t} \log p_\phi(y|x_t),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

SCORE-BASED GENERATIVE MODELING THROUGH STOCHASTIC DIFFERENTIAL EQUATIONS

$$d\mathbf{x} = \{\mathbf{f}(\mathbf{x}, t) - g(t)^2[\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) + \nabla_{\mathbf{x}} \log p_t(\mathbf{y} | \mathbf{x})]\}dt + g(t)d\bar{\mathbf{w}}.$$

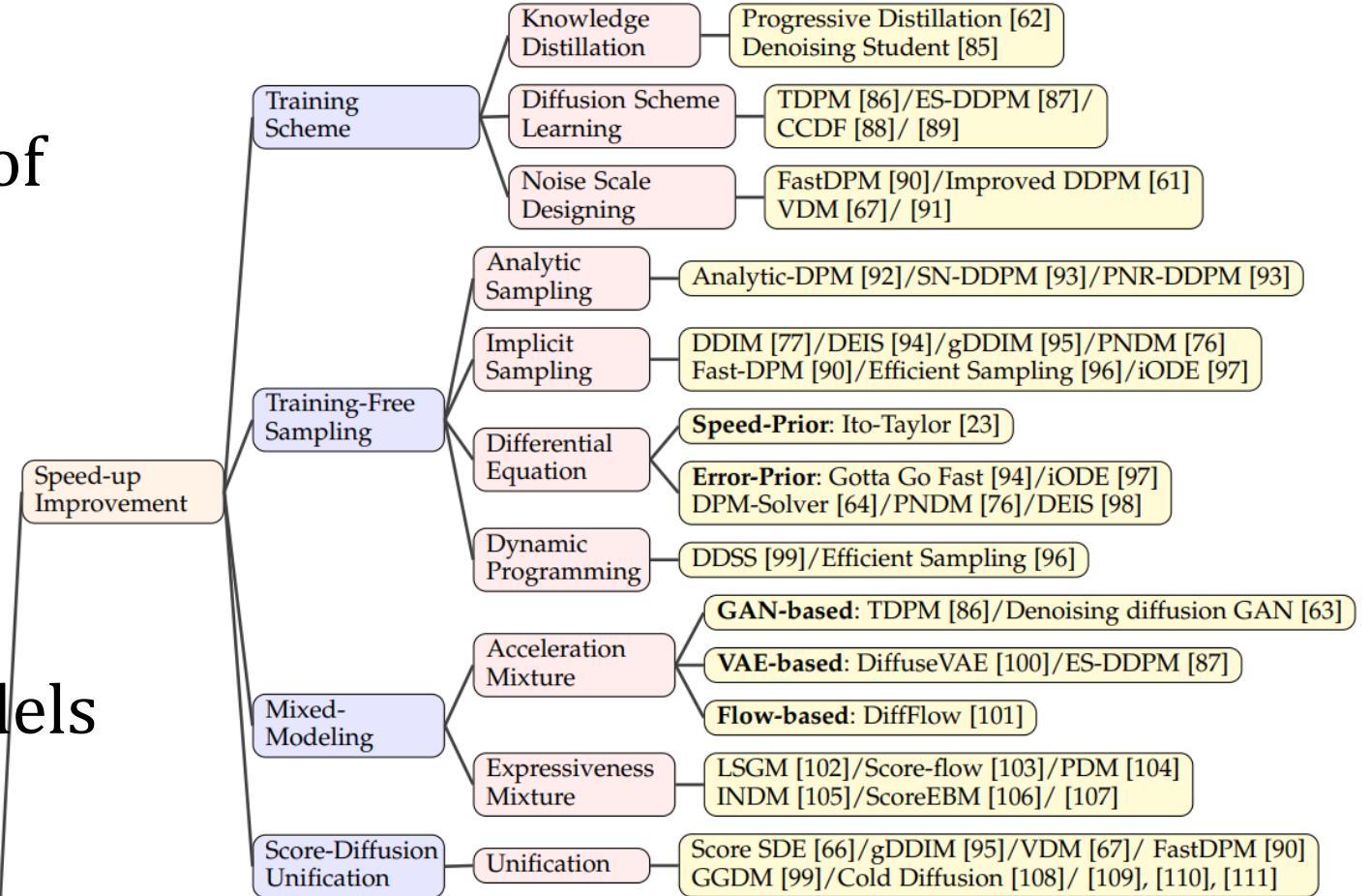


Landmark Works: Data_Conditional Diffusion

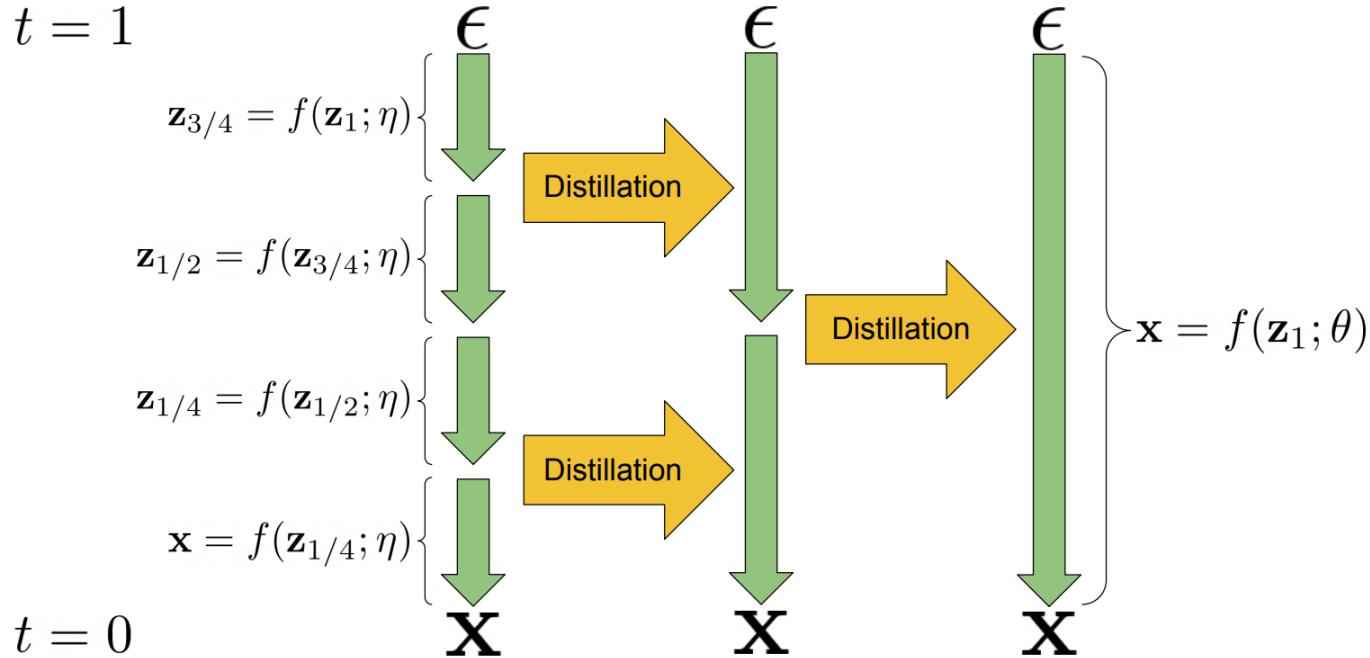


Diffusion Needs Improvements

- Gaussian-based diffusion sampler takes thousands of steps to sample.
1. Change diffusion pattern
 2. Advanced sampling
 3. Seek help from other models
 4. Unified Framework



Training Scheme: Distillation



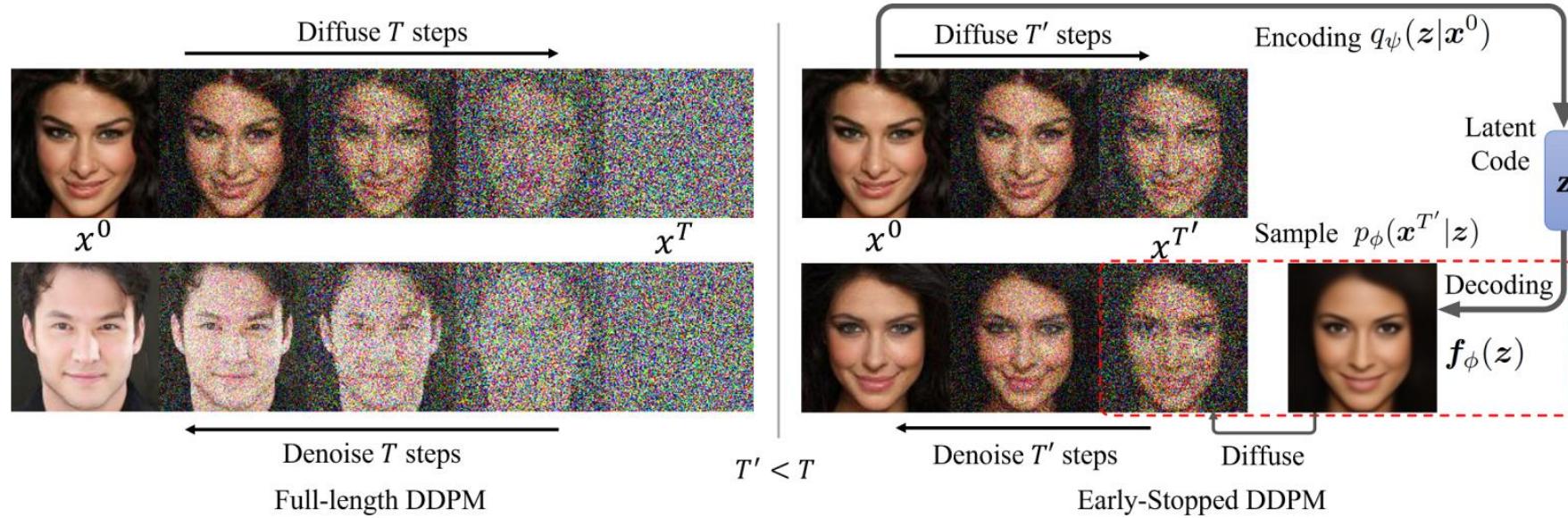
Algorithm 2 Progressive distillation

```

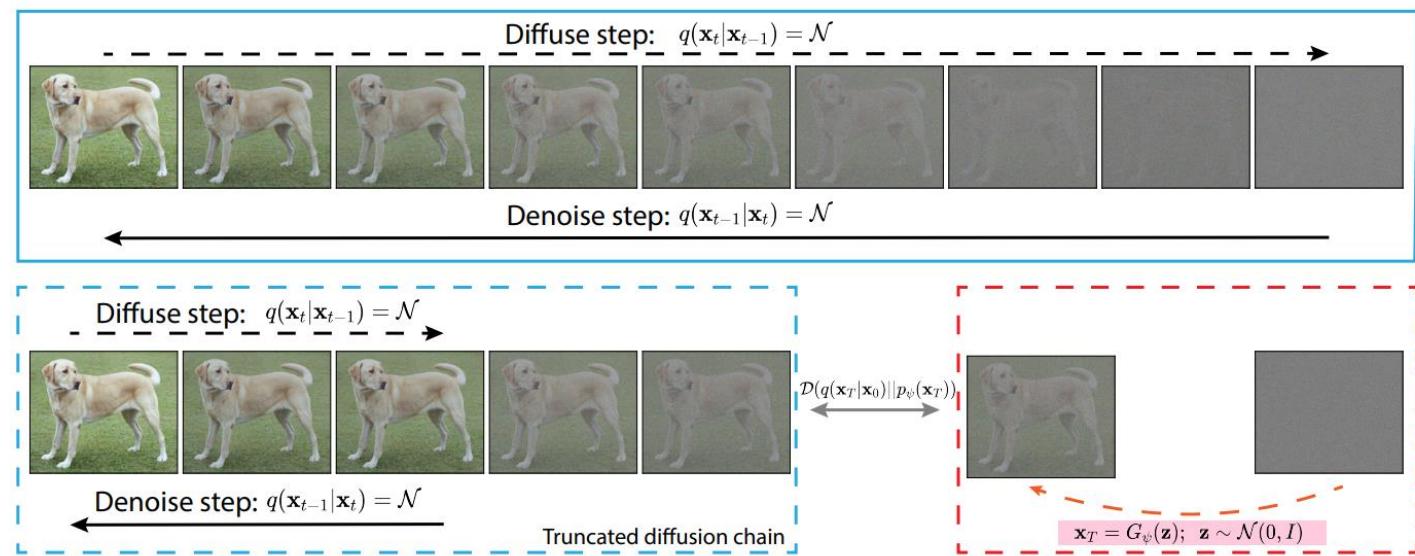
Require: Trained teacher model  $\hat{\mathbf{x}}_\eta(\mathbf{z}_t)$ 
Require: Data set  $\mathcal{D}$ 
Require: Loss weight function  $w()$ 
Require: Student sampling steps  $N$ 
for  $K$  iterations do
     $\theta \leftarrow \eta$                                  $\triangleright$  Init student from teacher
    while not converged do
         $\mathbf{x} \sim \mathcal{D}$ 
         $t = i/N, i \sim \text{Cat}[1, 2, \dots, N]$ 
         $\epsilon \sim N(0, I)$ 
         $\mathbf{z}_t = \alpha_t \mathbf{x} + \sigma_t \epsilon$ 
        # 2 steps of DDIM with teacher
         $t' = t - 0.5/N, t'' = t - 1/N$ 
         $\mathbf{z}_{t'} = \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_t) + \frac{\sigma_{t'}}{\sigma_t} (\mathbf{z}_t - \alpha_t \hat{\mathbf{x}}_\eta(\mathbf{z}_t))$ 
         $\mathbf{z}_{t''} = \alpha_{t''} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}) + \frac{\sigma_{t''}}{\sigma_{t'}} (\mathbf{z}_{t'} - \alpha_{t'} \hat{\mathbf{x}}_\eta(\mathbf{z}_{t'}))$ 
         $\tilde{\mathbf{x}} = \frac{\mathbf{z}_{t''} - (\sigma_{t''}/\sigma_t) \mathbf{z}_t}{\alpha_{t''} - (\sigma_{t''}/\sigma_t) \alpha_t}$        $\triangleright$  Teacher  $\hat{\mathbf{x}}$  target
         $\lambda_t = \log[\alpha_t^2/\sigma_t^2]$ 
         $L_\theta = w(\lambda_t) \|\tilde{\mathbf{x}} - \hat{\mathbf{x}}_\theta(\mathbf{z}_t)\|_2^2$ 
         $\theta \leftarrow \theta - \gamma \nabla_\theta L_\theta$ 
    end while
     $\eta \leftarrow \theta$                                  $\triangleright$  Student becomes next teacher
     $N \leftarrow N/2$   $\triangleright$  Halve number of sampling steps
end for

```

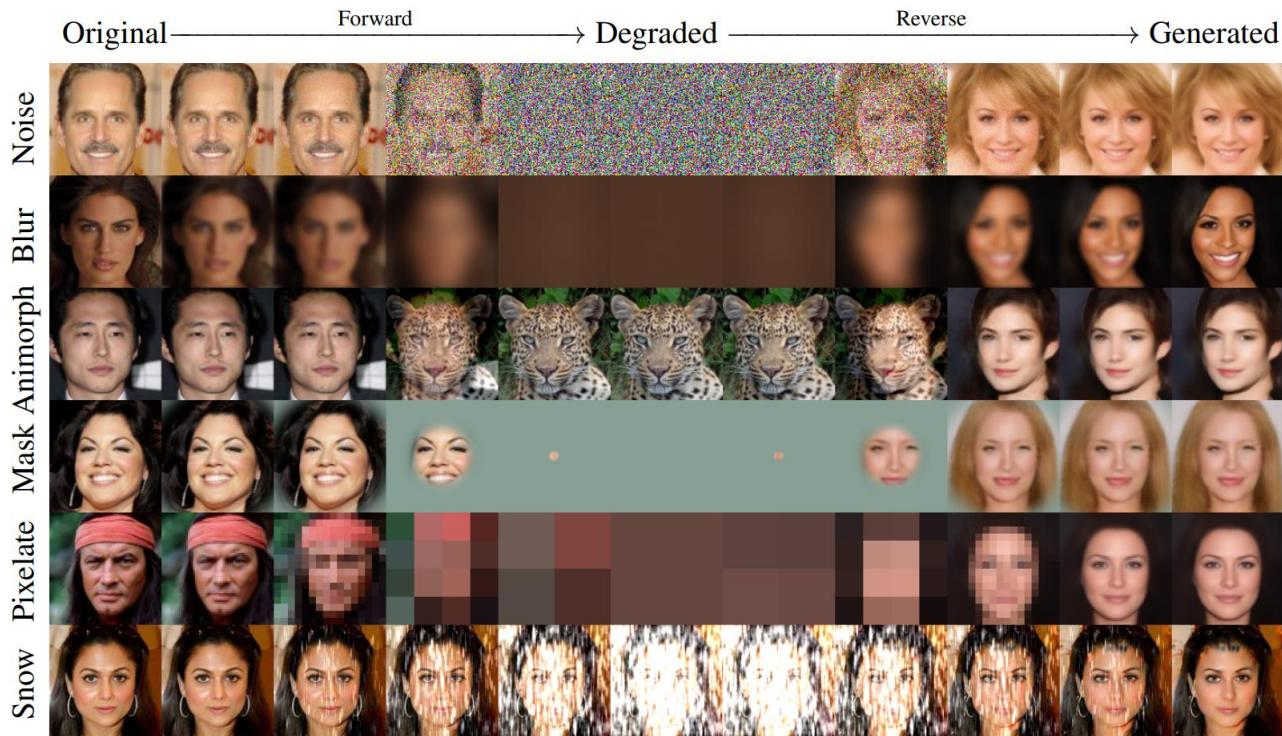
Training Scheme: Stop!



Early Stop!



Training Scheme: Exploring the distribution



Algorithm 1 Naive Sampling

Input: A degraded sample x_t
for $s = t, t - 1, \dots, 1$ **do**
 $\hat{x}_0 \leftarrow R(x_s, s)$
 $x_{s-1} = D(\hat{x}_0, s - 1)$
end for
Return: x_0

Algorithm 2 Improved Sampling for Cold Diffusion

Input: A degraded sample x_t
for $s = t, t - 1, \dots, 1$ **do**
 $\hat{x}_0 \leftarrow R(x_s, s)$
 $x_{s-1} = x_s - D(\hat{x}_0, s) + D(\hat{x}_0, s - 1)$
end for

Training Scheme: Design the Noise

- FastDPM: link noise with time t with bijective map

$$\mathcal{R}(t) = (\Delta\beta)^{\frac{t}{2}} \Gamma\left(\hat{\beta} + 1\right)^{\frac{1}{2}} \Gamma\left(\hat{\beta} - t + 1\right)^{-\frac{1}{2}}.$$

- VDM: link noise with time t with bijective map

$$\sigma_t^2 = \text{sigmoid}(\gamma_{\boldsymbol{\eta}}(t)) \quad \alpha_t^2 = \text{sigmoid}(-\gamma_{\boldsymbol{\eta}}(t))$$

guide the training with SNR:

$$\text{SNR}(t) = \exp(-\gamma_{\boldsymbol{\eta}}(t)) \quad \mathcal{L}_T(\mathbf{x}) = \frac{T}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), i \sim U\{1, T\}} [(\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{x} - \hat{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{z}_t; t)\|_2^2],$$

$$\mathcal{L}_{\infty}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_{\text{SNR}_{\min}}^{\text{SNR}_{\max}} \|\mathbf{x} - \tilde{\mathbf{x}}_{\boldsymbol{\theta}}(\mathbf{z}_v, v)\|_2^2 dv,$$

- Improved DDPM: learn noise scale by new loss function

$$L_{\text{hybrid}} = L_{\text{simple}} + \lambda L_{\text{vlb}}$$

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:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
```

Advanced Sampling: Analytic Solver

Theorem 1. (Score representation of the optimal solution to Eq. (4), proof in Appendix A.2)

The optimal solution $\mu_n^*(\mathbf{x}_n)$ and σ_n^{*2} to Eq. (4) are

$$\begin{aligned}\mu_n^*(\mathbf{x}_n) &= \tilde{\mu}_n \left(\mathbf{x}_n, \frac{1}{\sqrt{\bar{\alpha}_n}} (\mathbf{x}_n + \bar{\beta}_n \nabla_{\mathbf{x}_n} \log q_n(\mathbf{x}_n)) \right), \\ \sigma_n^{*2} &= \lambda_n^2 + \left(\sqrt{\frac{\bar{\beta}_n}{\alpha_n}} - \sqrt{\bar{\beta}_{n-1} - \lambda_n^2} \right)^2 \left(1 - \bar{\beta}_n \mathbb{E}_{q_n(\mathbf{x}_n)} \frac{\|\nabla_{\mathbf{x}_n} \log q_n(\mathbf{x}_n)\|^2}{d} \right),\end{aligned}$$

- Analytical DPM: Start from vlb optimization to explore analytical solutions
- Extended Analytical DPM: Suppose that reverse mean is independent of reverse noise, finding optimal mean and noise

Theorem 3.2. (Optimal solution to optimization solely w.r.t. mean) For any covariance Σ_n , the optimal mean to problem (12) is always Eq. (4), i.e.,

$$\mu_n^*(\mathbf{x}_n) = \tilde{\mu}_n \left(\mathbf{x}_n, \frac{1}{\sqrt{\bar{\alpha}_n}} (\mathbf{x}_n - \sqrt{\bar{\beta}_n} \mathbb{E}_{q(\mathbf{x}_0|\mathbf{x}_n)} [\epsilon_n]) \right),$$

which is irrelevant to Σ_n .

Theorem 3.3. (Optimal solution to optimization w.r.t. covariance solely) Suppose $\Sigma_n(\mathbf{x}_n) = \text{diag}(\sigma_n(\mathbf{x}_n)^2)$. For any mean $\mu_n(\mathbf{x}_n)$ that is parameterized by a noise prediction network $\hat{\epsilon}_n(\mathbf{x}_n)$ as in Eq. (6), the optimal covariance $\tilde{\sigma}_n^*(\mathbf{x}_n)^2$ to problem (13) is

$$\begin{aligned}\tilde{\sigma}_n^*(\mathbf{x}_n)^2 &= \sigma_n^*(\mathbf{x}_n)^2 + \gamma_n^2 \underbrace{\frac{\bar{\beta}_n}{\bar{\alpha}_n} (\hat{\epsilon}_n(\mathbf{x}_n) - \mathbb{E}_{q(\mathbf{x}_0|\mathbf{x}_n)} [\epsilon_n])^2}_{\text{error}} \\ &= \lambda_n^2 \mathbf{1} + \gamma_n^2 \frac{\bar{\beta}_n}{\bar{\alpha}_n} \mathbb{E}_{q(\mathbf{x}_0|\mathbf{x}_n)} [(\epsilon_n - \hat{\epsilon}_n(\mathbf{x}_n))^2], \quad (14)\end{aligned}$$

where $\sigma_n^*(\mathbf{x}_n)^2$ is the optimal covariance to the joint optimization problem in Theorem 3.1, $\epsilon_n = \frac{\mathbf{x}_n - \sqrt{\bar{\alpha}_n} \mathbf{x}_0}{\sqrt{\bar{\beta}_n}}$ is the noise used to generate \mathbf{x}_n from \mathbf{x}_0 , $(\cdot)^2$ is the element-wise square, $\mathbf{1}$ is the vector of ones and $\gamma_n = \sqrt{\bar{\alpha}_{n-1}} - \sqrt{\bar{\beta}_{n-1} - \lambda_n^2} \sqrt{\frac{\bar{\alpha}_n}{\bar{\beta}_n}}$.

Advanced Sampling: Implicit/non-Gaussian

$$\left\{ \begin{array}{l} q(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\beta}_t \mathbf{I}), \\ \text{where } \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) := \frac{\sqrt{\alpha_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \quad \text{and} \quad \tilde{\beta}_t := \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \\ q_\sigma(\mathbf{x}_{t-1} | \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\left(\sqrt{\alpha_{t-1}} \mathbf{x}_0 + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\alpha_t} \mathbf{x}_0}{\sqrt{1 - \alpha_t}}, \sigma_t^2 \mathbf{I}\right). \end{array} \right.$$

- Implicit Sampling

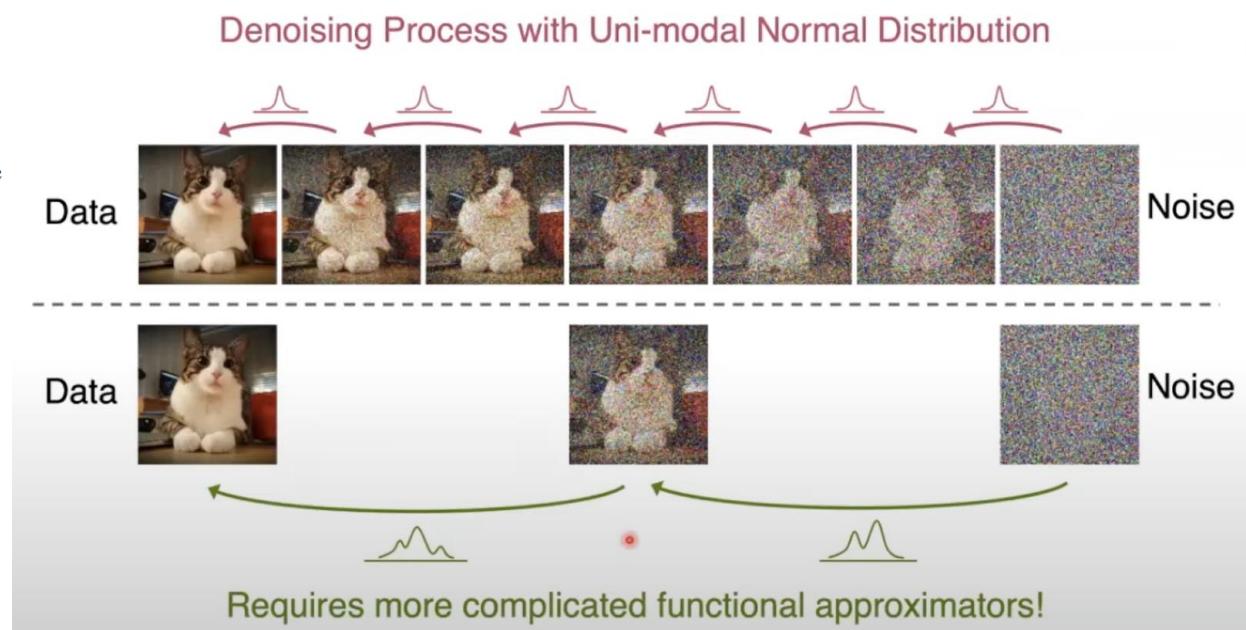
$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } \mathbf{x}_0\text{"}} + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t)}_{\text{"direction pointing to } \mathbf{x}_t\text{"}} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

- Implicit Sampling as continuous neural ODE

$$\frac{\mathbf{x}_{t-\Delta t}}{\sqrt{\alpha_{t-\Delta t}}} = \frac{\mathbf{x}_t}{\sqrt{\alpha_t}} + \left(\sqrt{\frac{1 - \alpha_{t-\Delta t}}{\alpha_{t-\Delta t}}} - \sqrt{\frac{1 - \alpha_t}{\alpha_t}} \right) \epsilon_\theta^{(t)}(\mathbf{x}_t)$$

$$d\bar{\mathbf{x}}(t) = \epsilon_\theta^{(t)} \left(\frac{\bar{\mathbf{x}}(t)}{\sqrt{\sigma^2 + 1}} \right) d\sigma(t),$$

Comparison between DDPM and DDIM



Advanced Sampling: Differential Equation Solver

- Difference between ODE-based and SDE-based models
- Corresponding solver for ODE and SDE
- Difference between higher order and lower order solvers

Improved Techniques:

1. Mid-point solver
2. Multi-methods solver
3. Threshold
4. Dynamic step size adjustment

Algorithm 2 Our stochastic sampler with $\sigma(t) = t$ and $s(t) = 1$.

```

1: procedure STOCHASTICSAMPLER( $D_\theta(\mathbf{x}; \sigma)$ ,  $t_{i \in \{0, \dots, N\}}$ ,  $\gamma_{i \in \{0, \dots, N-1\}}$ ,  $S_{\text{noise}}$ )
2:   sample  $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, t_0^2 \mathbf{I})$ 
3:   for  $i \in \{0, \dots, N-1\}$  do
4:     sample  $\epsilon_i \sim \mathcal{N}(\mathbf{0}, S_{\text{noise}}^2 \mathbf{I})$ 
5:      $\hat{t}_i \leftarrow t_i + \gamma_i t_i$ 
6:      $\hat{\mathbf{x}}_i \leftarrow \mathbf{x}_i + \sqrt{\hat{t}_i^2 - t_i^2} \epsilon_i$ 
7:      $\mathbf{d}_i \leftarrow (\hat{\mathbf{x}}_i - D_\theta(\hat{\mathbf{x}}_i; \hat{t}_i)) / \hat{t}_i$ 
8:      $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - \hat{t}_i) \mathbf{d}_i$ 
9:     if  $t_{i+1} \neq 0$  then
10:       $\mathbf{d}'_i \leftarrow (\mathbf{x}_{i+1} - D_\theta(\mathbf{x}_{i+1}; t_{i+1})) / t_{i+1}$ 
11:       $\mathbf{x}_{i+1} \leftarrow \hat{\mathbf{x}}_i + (t_{i+1} - \hat{t}_i) (\frac{1}{2} \mathbf{d}_i + \frac{1}{2} \mathbf{d}'_i)$ 
12:   return  $\mathbf{x}_N$ 

```

Algorithm 1 Dynamic step size extrapolation for solving Reverse Diffusion Processes

Require: $s_\theta, \epsilon_{rel}, \epsilon_{abs}, h_{init} = 0.01, r = 0.9, \theta = 0.9$ ▷ for images: $\epsilon_{abs} = \frac{y_{max} - y_{min}}{256}$

```

 $t \leftarrow 1$ 
 $h \leftarrow h_{init}$ 
Initialize  $\mathbf{x}$ 
 $\mathbf{x}'_{prev} \leftarrow \mathbf{x}$ 
while  $t > 0$  do
  Draw  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
   $\mathbf{x}' \leftarrow \mathbf{x} - h f(\mathbf{x}, t) + h g(t)^2 s_\theta(\mathbf{x}, t) + \sqrt{h} g(t) \mathbf{z}$  ▷ Euler-Maruyama
   $\tilde{\mathbf{x}} \leftarrow \mathbf{x} - h f(\mathbf{x}', t - h) + h g(t - h)^2 s_\theta(\mathbf{x}', t - h) + \sqrt{h} g(t - h) \mathbf{z}$ 
   $\mathbf{x}'' \leftarrow \frac{1}{2}(\mathbf{x}' + \tilde{\mathbf{x}})$  ▷ Improved Euler (SDE version)
   $\delta \leftarrow \max(\epsilon_{abs}, \epsilon_{rel} \max(|\mathbf{x}'|, |\mathbf{x}'_{prev}|))$  ▷ Element-wise operations
   $E_2 \leftarrow \frac{1}{\sqrt{n}} \|(\mathbf{x}' - \mathbf{x}'') / \delta\|_2$ 
  if  $E_2 \leq 1$  then
     $\mathbf{x} \leftarrow \mathbf{x}''$ 
     $t \leftarrow t - h$ 
     $\mathbf{x}'_{prev} \leftarrow \mathbf{x}'$ 
     $h \leftarrow \min(t, \theta h E_2^{-r})$  ▷ Accept
  return  $\mathbf{x}$  ▷ Extrapolation
▷ Dynamic step size update

```

Advanced Sampling: Dynamic Sampling

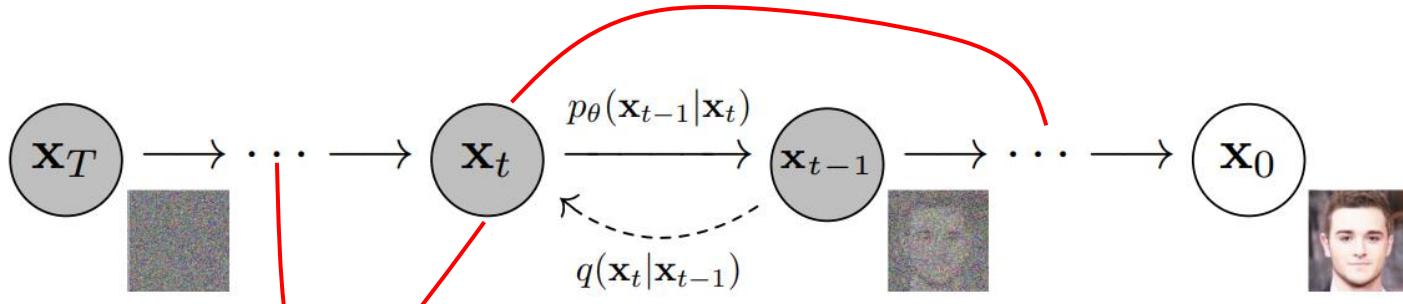


Figure 2: The directed graphical model considered in this work.

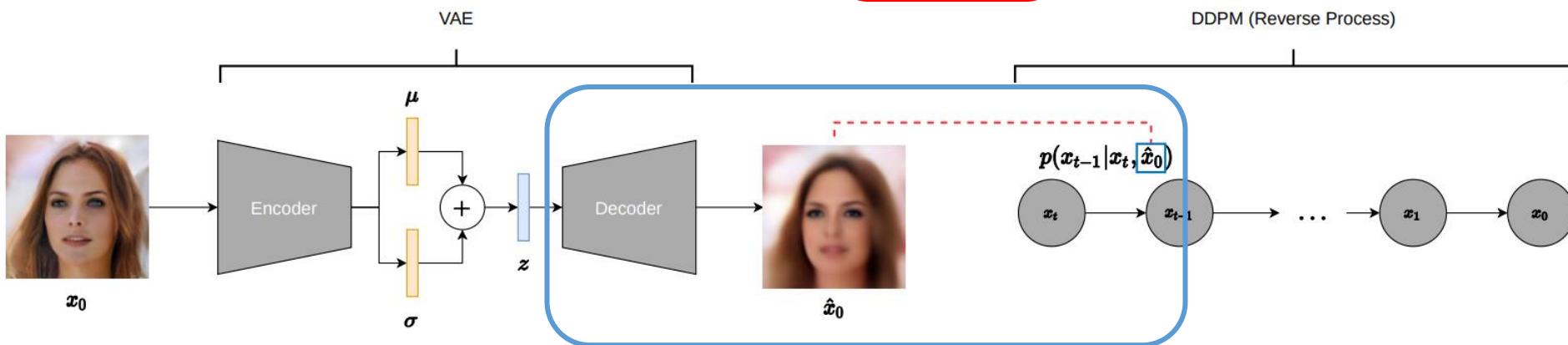
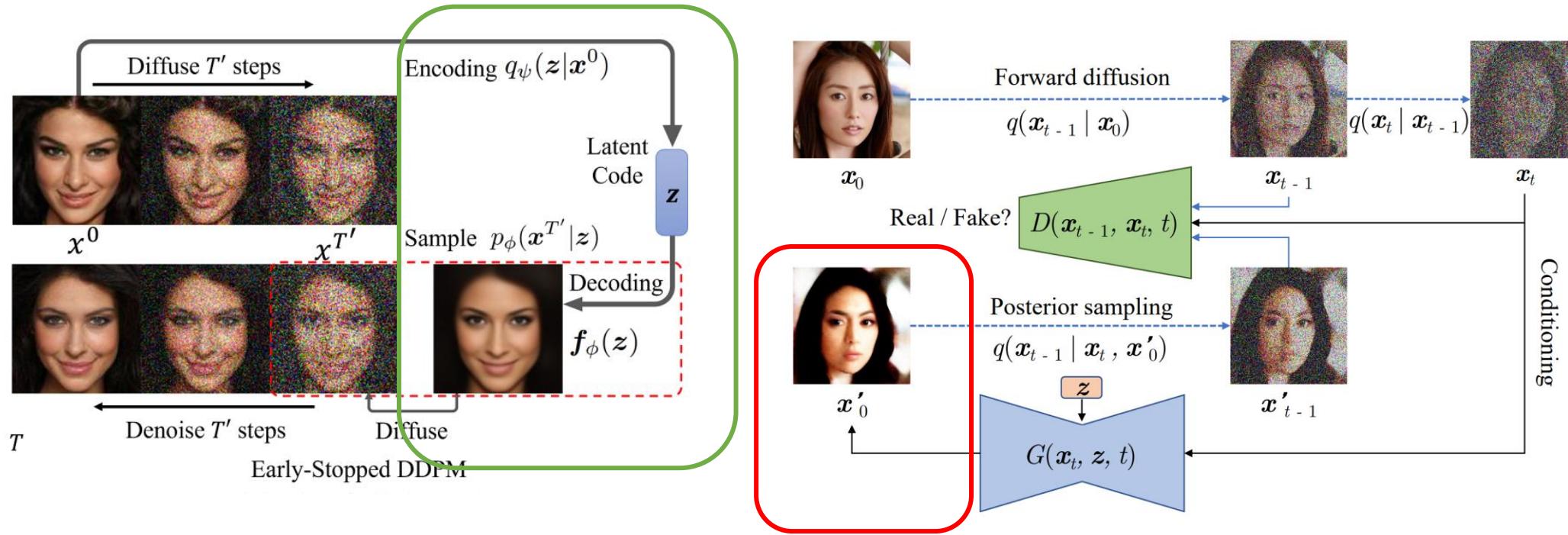
We can construct an ELBO path that entails $K \in \mathbb{N}$ refinement steps. I.e., for any K , and any given path of inference timesteps $0 = t'_0 < t'_1 < \dots < t'_{K-1} < t'_K = 1$, one can derive a corresponding ELBO

$$-L_{\text{ELBO}} = \mathbb{E}_q D_{\text{KL}}(q(\mathbf{x}_1|\mathbf{x}_0) \| p_\theta(\mathbf{x}_1)) + \sum_{i=1}^K L(t'_i, t'_{i-1}) \quad (16)$$

where

$$L(t, s) = \begin{cases} -\mathbb{E}_q \log p_\theta(\mathbf{x}_t|\mathbf{x}_0) & s = 0 \\ \mathbb{E}_q D_{\text{KL}}(q(\mathbf{x}_s|\mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_s|\mathbf{x}_t)) & s > 0 \end{cases} \quad (17)$$

Mixed-Modeling: Acceleration



Mixed-Modeling: Expressiveness

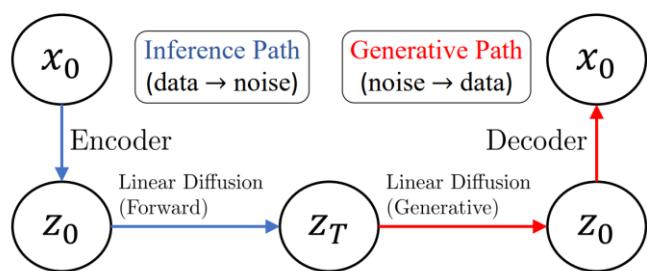


Figure 3: LSGM.

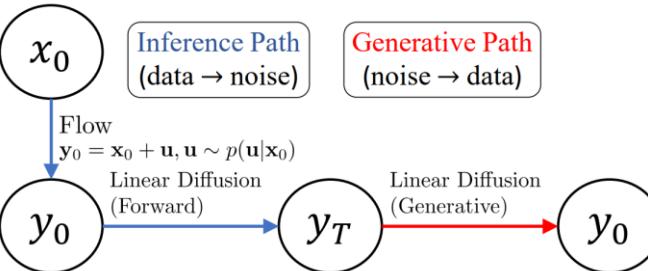


Figure 4: ScoreFlow.

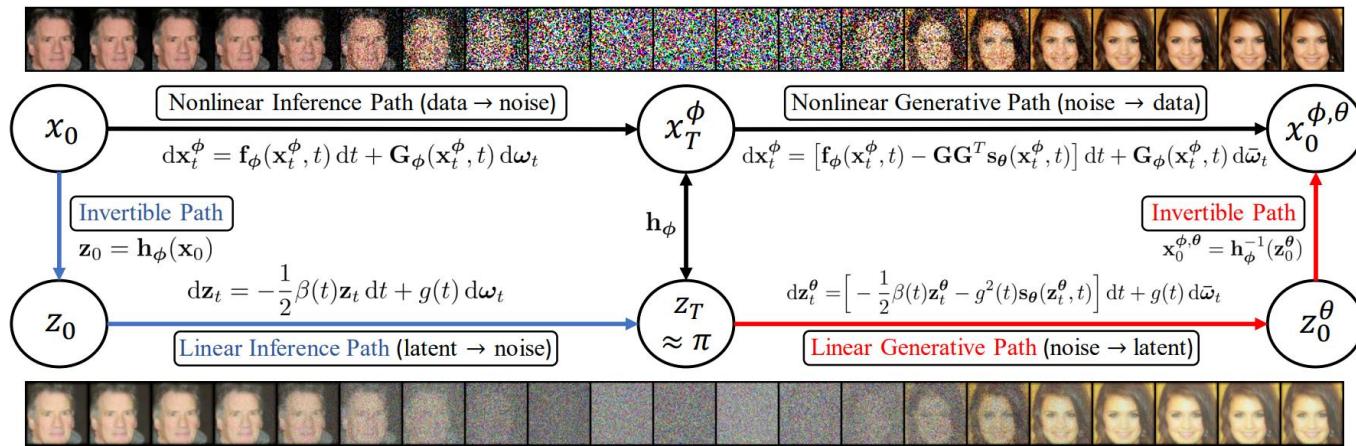
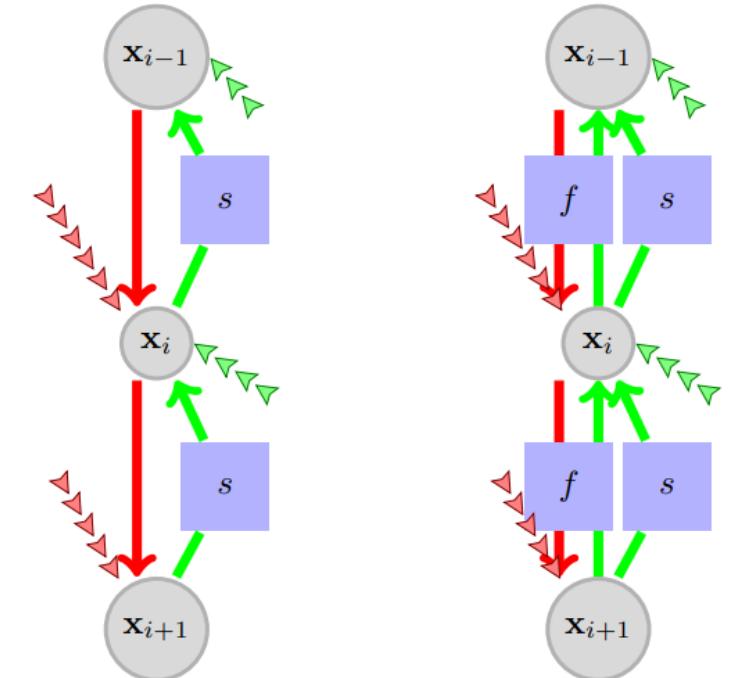


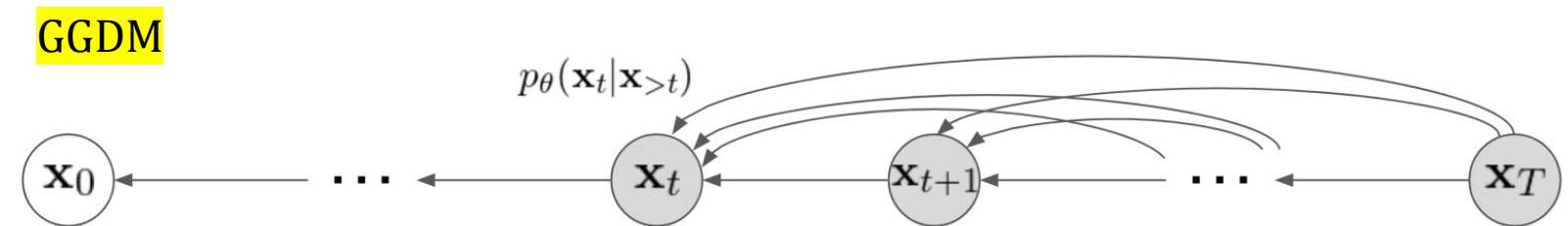
Figure 2: PDM attains a ladder structure between a data space and a latent space by the invertibility of normalizing flow. A score network estimates data score following blue arrows, and a model generates images following red arrows.



Diffusion Models

DiffFlow

Unification: Unified Framework

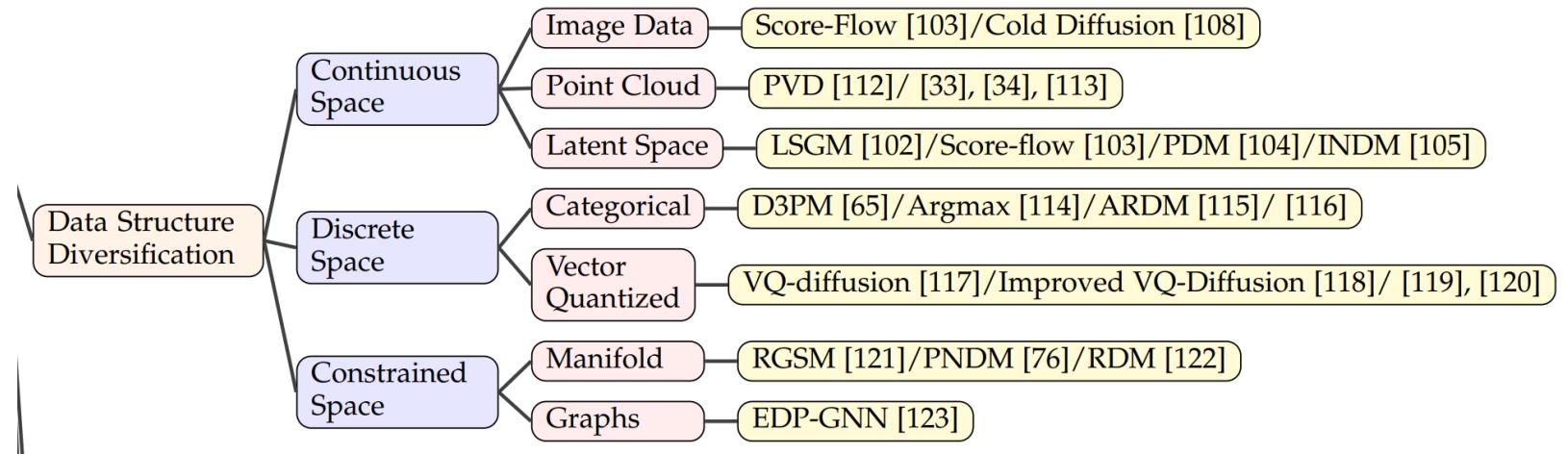


- ScoreSDE: Unified SDE Framework
- gDDIM: Generalized Implicit diffusion model
- VDM: Express the whole pipeline w.r.t. signal-to-noise ratio
- FastDPM: Express the whole pipeline w.r.t time-based function
- GGDM: Consider a fully-dependent diffusion framework
- Cold Diffusion: Summarize diffusion as Degrader and Reconstructor

Diffusion Needs Improvements

- Traditional RGB-image is not applicable for all fields.

1. Continuous Space
2. Discrete Space
3. Constrained Space



Diversification: Continuous

3D Point Cloud

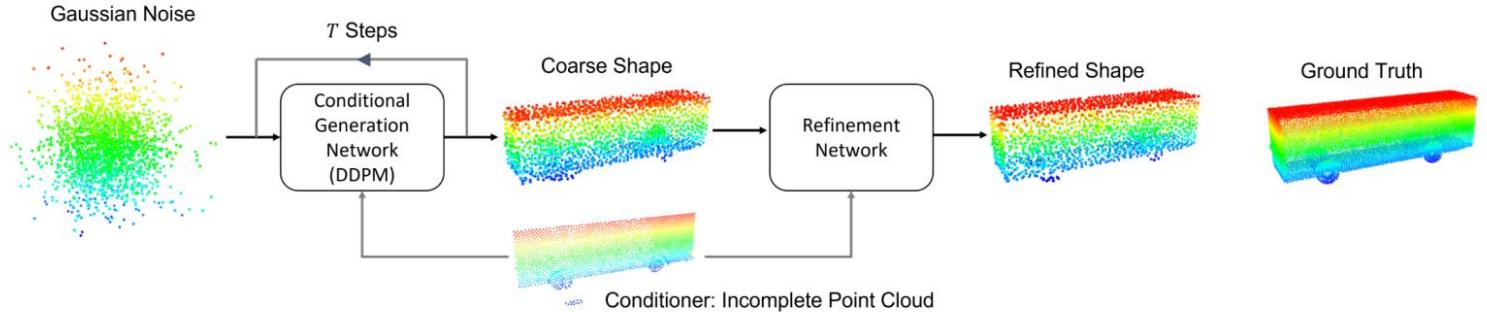
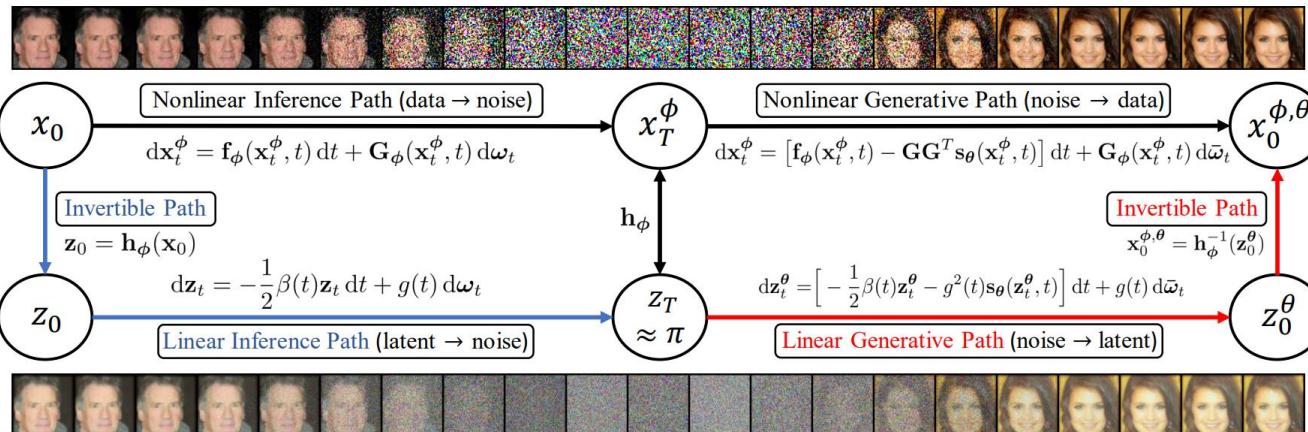


Figure 1: Our **Conditional Point Diffusion-Refinement (PDR)** paradigm first moves a Gaussian noise step by step towards a coarse completion of the partial observation through a diffusion model (DDPM). Then it refines the coarse point cloud by one step to obtain a high quality point cloud.

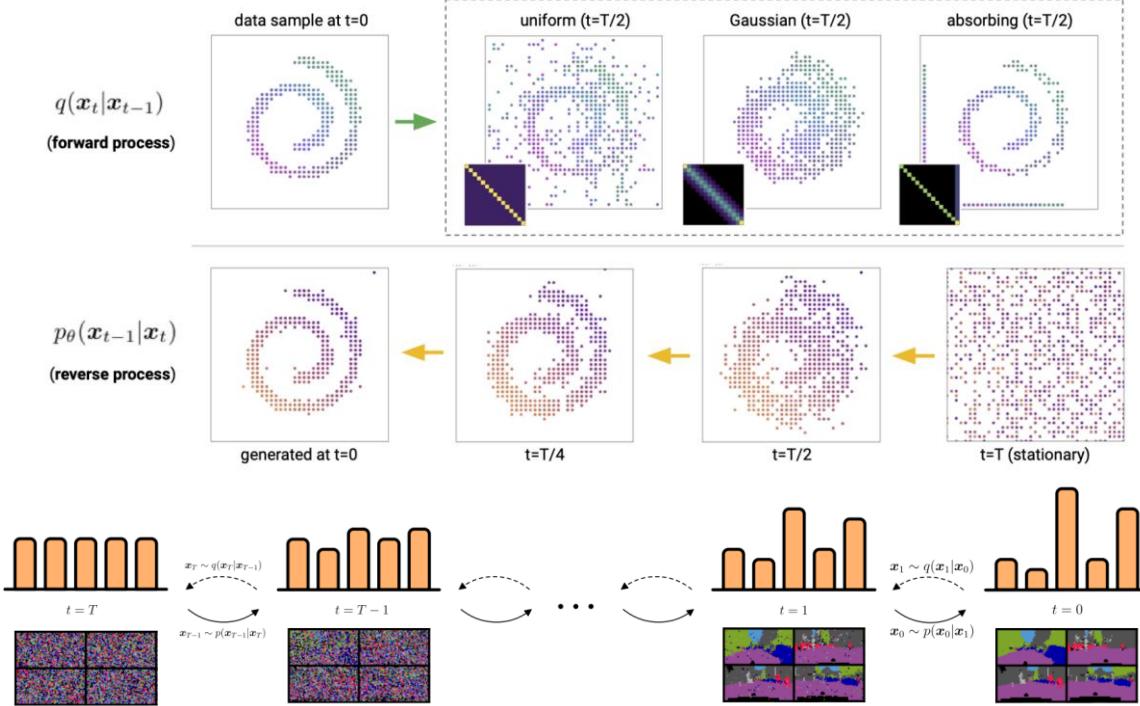


Latent Space

Figure 2: PDM attains a ladder structure between a data space and a latent space by the invertibility of normalizing flow. A score network estimates data score following blue arrows, and a model generates images following red arrows.

Diversification: Discrete

Categorical/Multinomial



Vector-Quantized

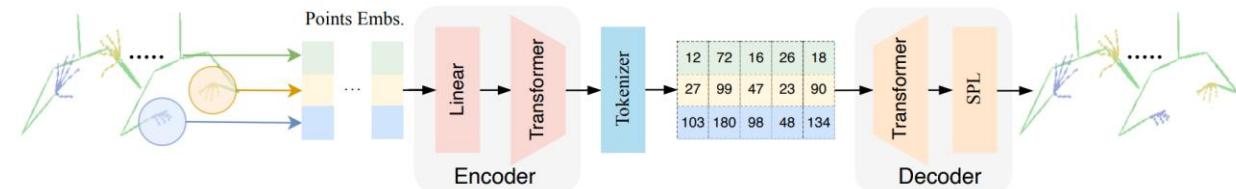
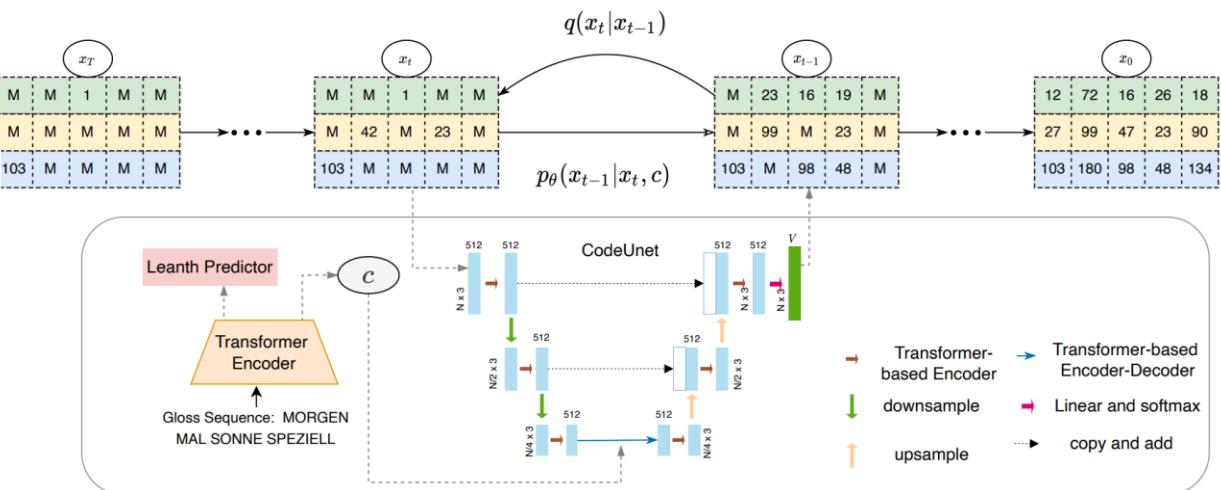
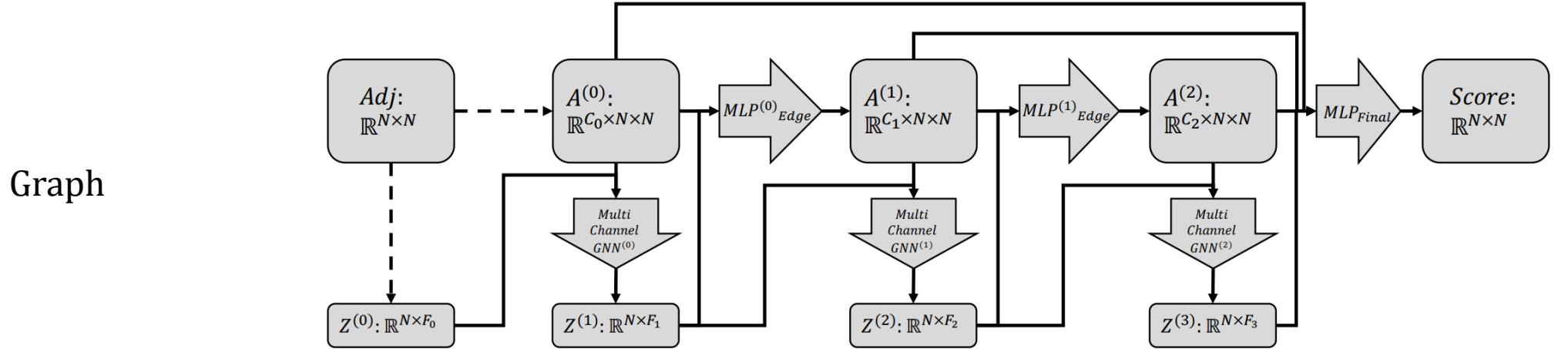


Figure 1: The architecture of the first stage model Pose-VQVAE for learning the discrete latent codes.



Diversification: Constrained



Graph

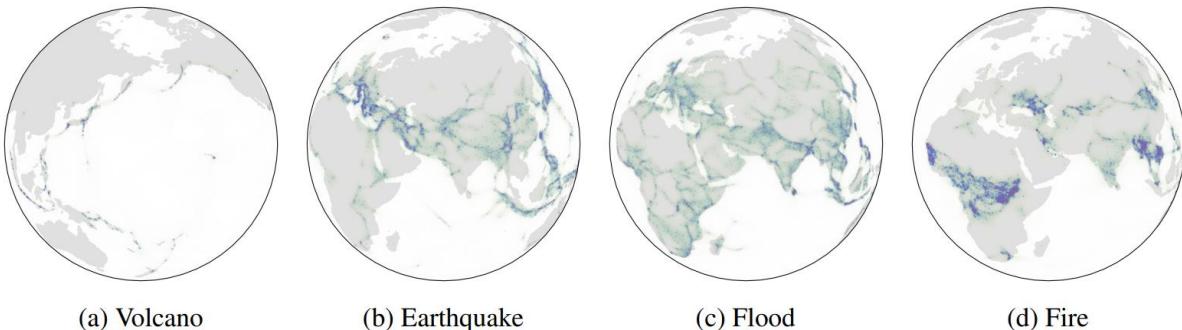
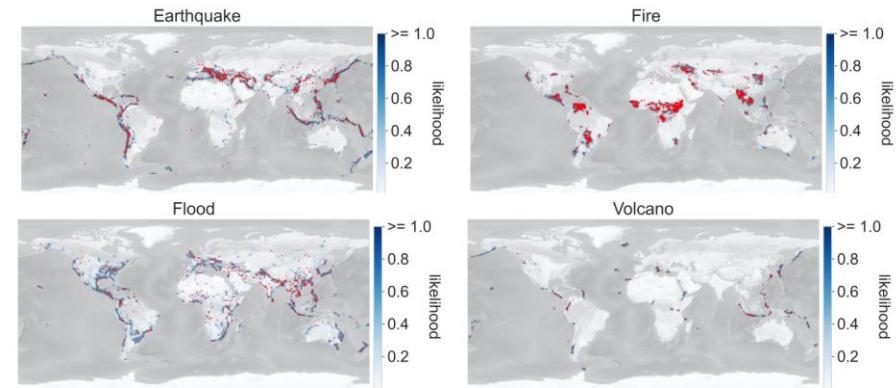


Figure 1: Trained score-based generative models on earth sciences data. The learned density is colored green-blue. Blue and red dots represent training and testing datapoints, respectively.



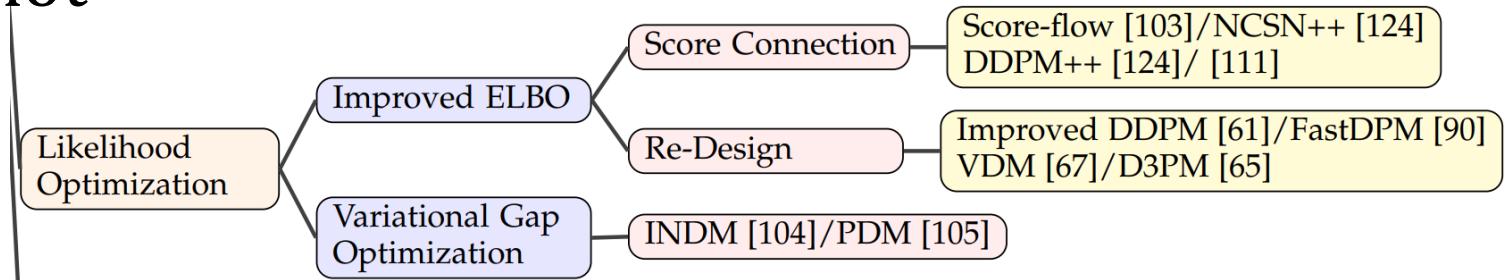
Manifold

Figure 1: Density of models trained on earth datasets. Red dots are samples from the test set.

Diffusion Needs Improvements

- Likelihood Optimization achieved by ELBO may not work sometimes.

1. Advanced ELBO Design
2. Variational Gap Optimization



Likelihood Optimization: Advanced ELBO

Theorem 3. Let $p_{0t}(\mathbf{x}' \mid \mathbf{x})$ denote the transition distribution from $p_0(\mathbf{x})$ to $p_t(\mathbf{x})$ for the SDE in Eq. (1). With the same notations and conditions in Theorem 1, we have

$$-\log p_{\theta}^{\text{SDE}}(\mathbf{x}) \leq \mathcal{L}_{\theta}^{\text{SM}}(\mathbf{x}) = \mathcal{L}_{\theta}^{\text{DSM}}(\mathbf{x}), \quad (11)$$

Score Connection: Represent ELBO based on score-matching loss

where $\mathcal{L}_{\theta}^{\text{SM}}(\mathbf{x})$ is defined as

$$-\mathbb{E}_{p_{0T}(\mathbf{x}' \mid \mathbf{x})}[\log \pi(\mathbf{x}')] + \frac{1}{2} \int_0^T \mathbb{E}_{p_{0t}(\mathbf{x}' \mid \mathbf{x})} \left[2g(t)^2 \nabla_{\mathbf{x}'} \cdot \mathbf{s}_{\theta}(\mathbf{x}', t) + g(t)^2 \|\mathbf{s}_{\theta}(\mathbf{x}', t)\|_2^2 - 2\nabla_{\mathbf{x}'} \cdot \mathbf{f}(\mathbf{x}', t) \right] dt,$$

and $\mathcal{L}_{\theta}^{\text{DSM}}(\mathbf{x})$ is given by

$$\begin{aligned} & -\mathbb{E}_{p_{0T}(\mathbf{x}' \mid \mathbf{x})}[\log \pi(\mathbf{x}')] + \frac{1}{2} \int_0^T \mathbb{E}_{p_{0t}(\mathbf{x}' \mid \mathbf{x})} \left[g(t)^2 \|\mathbf{s}_{\theta}(\mathbf{x}', t) - \nabla_{\mathbf{x}'} \log p_{0t}(\mathbf{x}' \mid \mathbf{x})\|_2^2 \right] dt \\ & - \frac{1}{2} \int_0^T \mathbb{E}_{p_{0t}(\mathbf{x}' \mid \mathbf{x})} \left[g(t)^2 \|\nabla_{\mathbf{x}'} \log p_{0t}(\mathbf{x}' \mid \mathbf{x})\|_2^2 + 2\nabla_{\mathbf{x}'} \cdot \mathbf{f}(\mathbf{x}', t) \right] dt. \end{aligned}$$

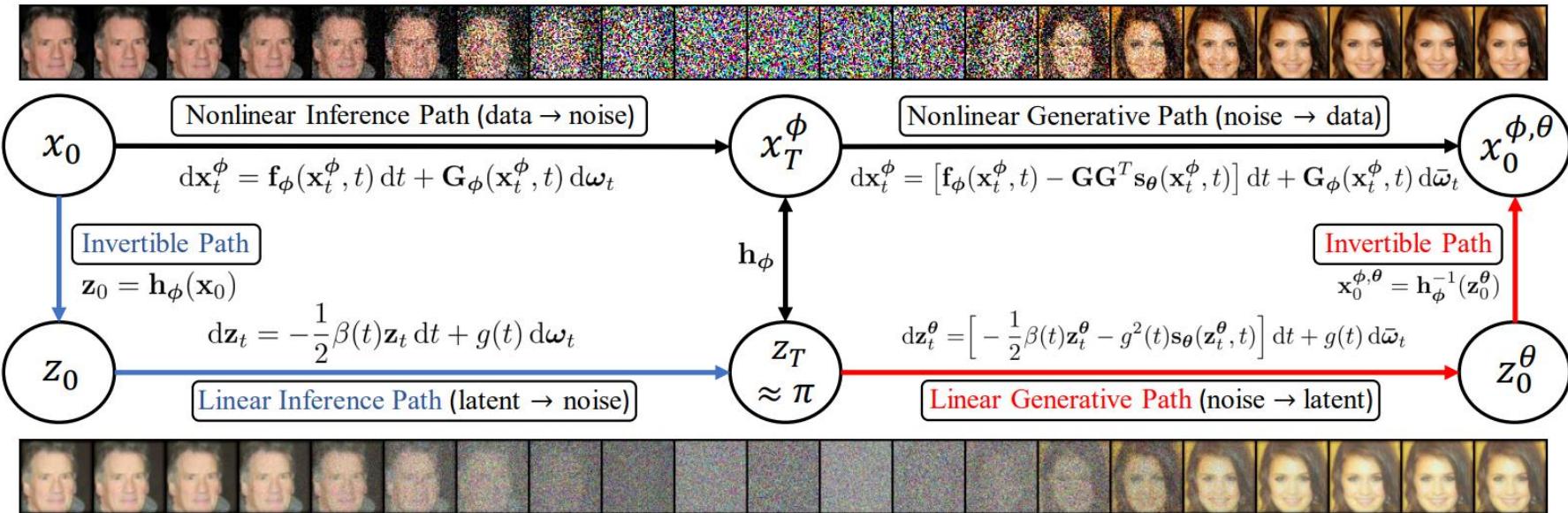
$$L_{\text{hybrid}} = L_{\text{simple}} + \lambda L_{\text{vlb}}$$

$$\mathcal{L}_T(\mathbf{x}) = \frac{T}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I}), i \sim U\{1, T\}} \left[(\text{SNR}(s) - \text{SNR}(t)) \|\mathbf{x} - \hat{\mathbf{x}}_{\theta}(\mathbf{z}_t; t)\|_2^2 \right],$$

$$\mathcal{L}_{\infty}(\mathbf{x}) = \frac{1}{2} \mathbb{E}_{\epsilon \sim \mathcal{N}(0, \mathbf{I})} \int_{\text{SNR}_{\min}}^{\text{SNR}_{\max}} \|\mathbf{x} - \tilde{\mathbf{x}}_{\theta}(\mathbf{z}_v, v)\|_2^2 dv,$$

Re-Design: design vlb loss from a different perspective to obtain better convergence

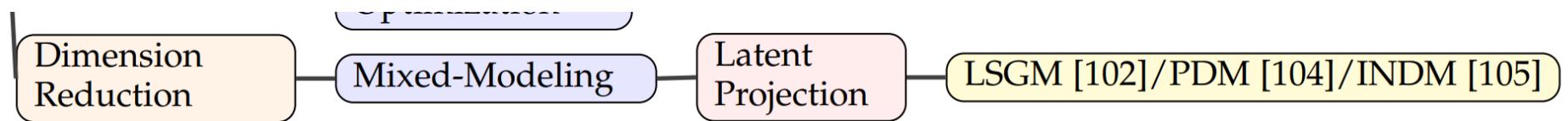
Likelihood Optimization: Variational Gap



$$\begin{aligned} \text{Gap}_d(\mu_\phi^d, \nu_{\phi,\theta}^d) &= \text{NELBO} - \text{NLL} \\ &= D_{KL}(\mu_\phi^d \| \nu_{\phi,\theta}^d) - D_{KL}(p_r \| p_{\phi,\theta}), \end{aligned}$$

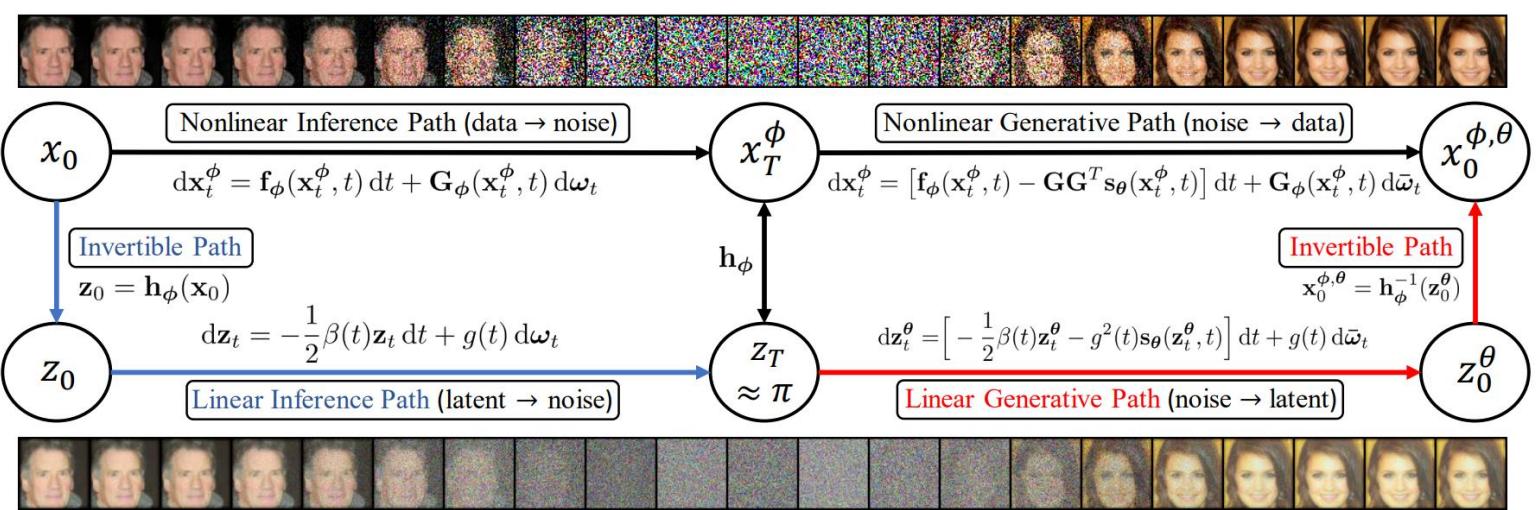
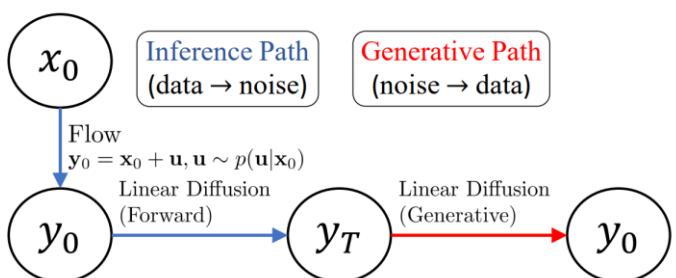
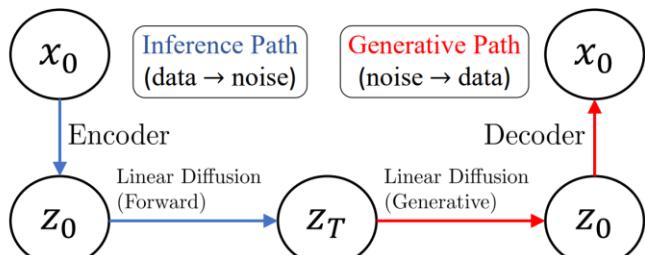
$$\phi^* = \arg \min_{\phi} \text{Gap}_d(\mu_\phi^d, \nu_{\phi,\theta}^d) = \arg \min_{\phi} D_{KL}(\mu_\phi^d \| \nu_{\phi,\theta}^d) = \arg \min_{\phi} D_{KL}(p_r \| p_{\phi,\theta}).$$

Diffusion Needs Improvements

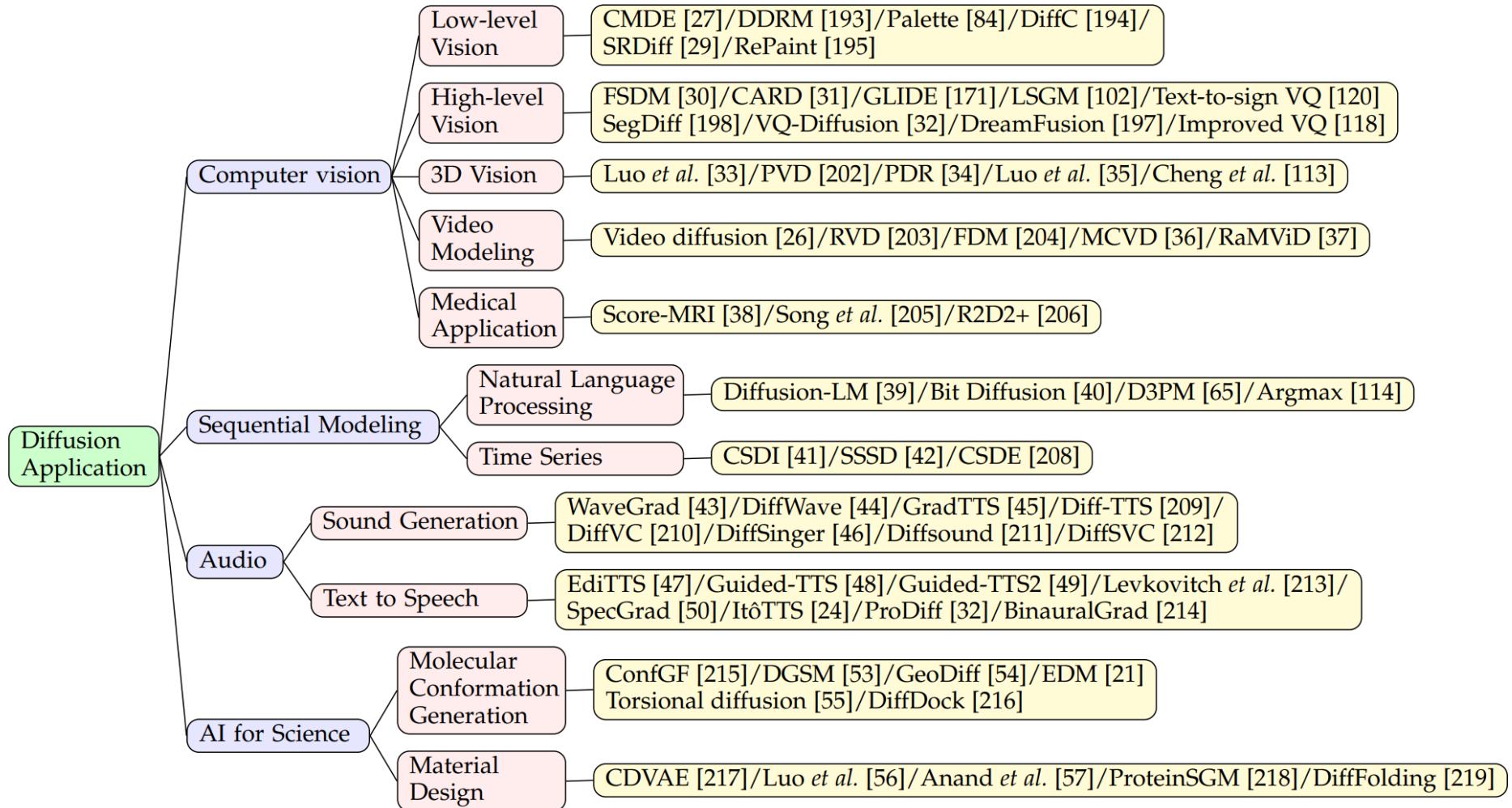


- Gradually transferred latent space cannot reduce the dimension
1. NF/VAE will help you

Dimension Reduction: Latent Projection



Diffusion Has Applications



Low-Level Vision:

CMDE [27]	2021	RGB-Image	SDE	Inpainting, Super-Resolution, Edge to image translation
DDRM [193]	2022	RGB-Image	Diffusion	Super-Resolution, Deblurring, Inpainting, Colorization
Palette [84]	2022	RGB-Image	Diffusion	Colorization, Inpainting, Uncropping, JPEG Restoration
DiffC [194]	2022	RGB-Image	SDE	Compression
SRDiff [29]	2021	RGB-Image	Diffusion	Super-Resolution -
RePaint [195]	2022	RGB-Image	Diffusion	Inpainting, Super-resolution, Edge to Image Translation

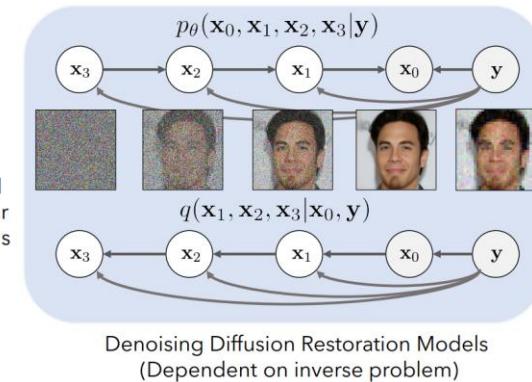
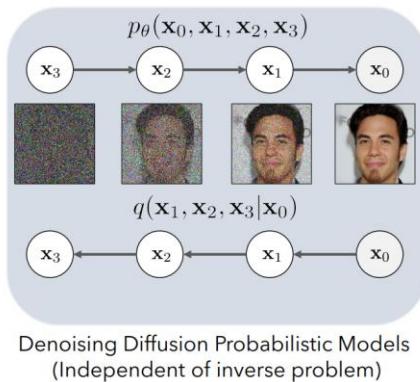
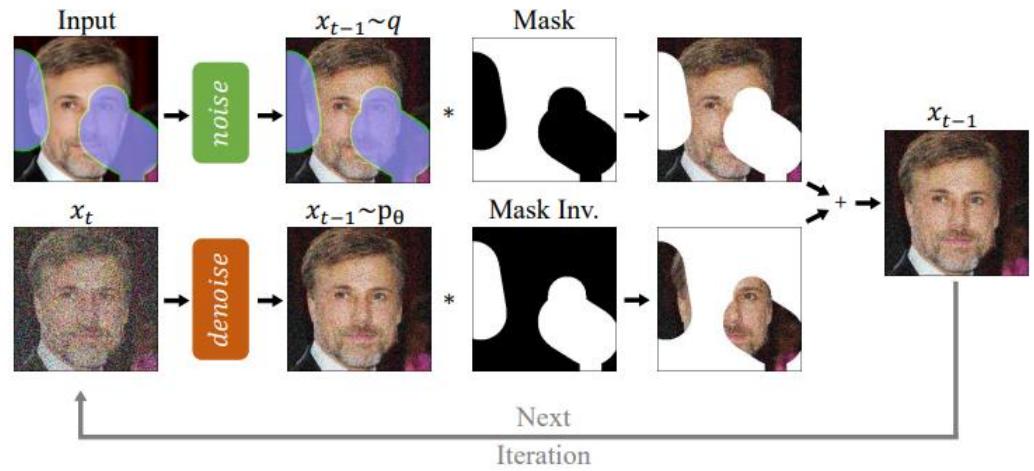
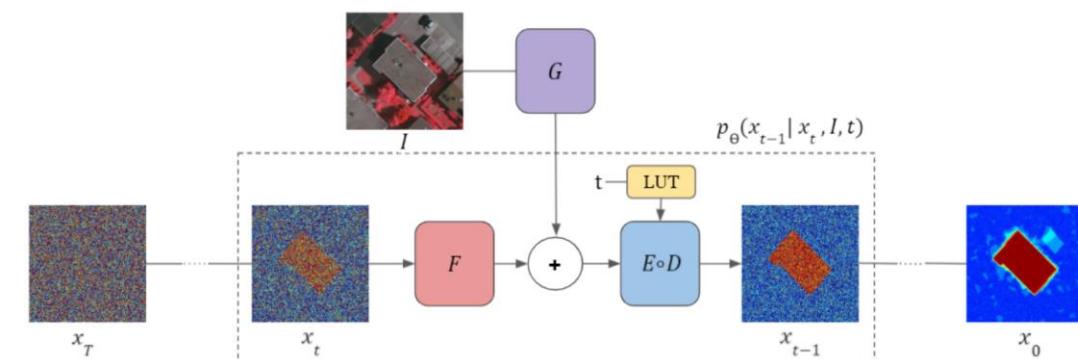


Figure 2. Illustration of our DDRM method for a specific inverse problem (super-resolution + denoising). We can use unsupervised DDPM models as a good solution to the DDRM objective.



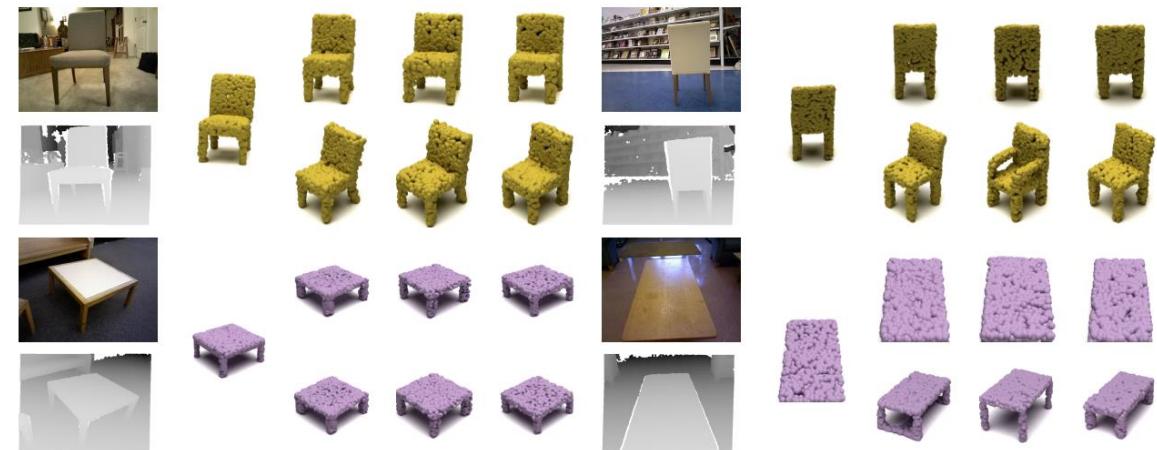
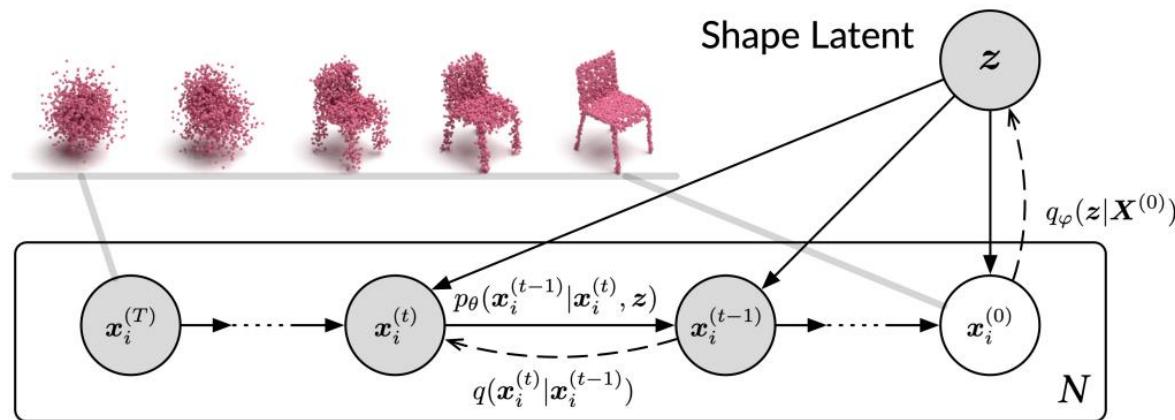
High-Level Vision:

FSDM [30]	2022	RGB-Image	Diffusion	Few-shot Generation
CARD [31]	2022	RGB-Image	Diffusion	Conditional Generation
GLIDE [171]	2022	RGB-Image	Diffusion	Conditional Generation
LSGM [102]	2022	RGB-Image	SDE	UnConditional & Conditional Generation
SegDiff [198]	2022	RGB-Image	Diffusion	Segmentation
VQ-Diffusion [117]	2022	VQ Data	Diffusion	Text-to-Image Synthesis
DreamFusion [197]	2023	VQ Data	Diffusion	Text-to-Image Synthesis
Text-to-Sign VQ [120]	2022	VQ Data	Diffusion	Conditional Pose Generation
Improved VQ-Diff [118]	2022	VQ Data	Diffusion	Text-to-Image Synthesis



3D Vision:

Luo's Model [33]	2021	Point Cloud	Diffusion	Point Cloud Generation
PVD [202]	2022	Point Cloud	Diffusion	Point Cloud Generation, Point-Voxel representation
PDR [34]	2022	Point Cloud	Diffusion	Point Cloud Completion
Cheng's Model [113]	2022	Point Cloud	Diffusion	Point Cloud Generation
Luo's Model [35]	2022	Point Cloud	Score	Point Cloud Denoising



Video Modeling:

VDM [26]	2022	Video	Diffusion	Text-Conditioned Video Generation
RVD [203]	2022	Video	Diffusion	Video Forecasting, Video compression
FDM [204]	2022	Video	Diffusion	Video Forecasting, Long-range Video modeling
MCVD [36]	2022	Video	Diffusion	Video Prediction, Video Generation, Video Interpolation
RaMViD [37]	2022	Video	SDE	Conditional Generation



Figure 1: The first two and last two frames of a video are given and our model does fill in the missing frames very accurate and with much detail.

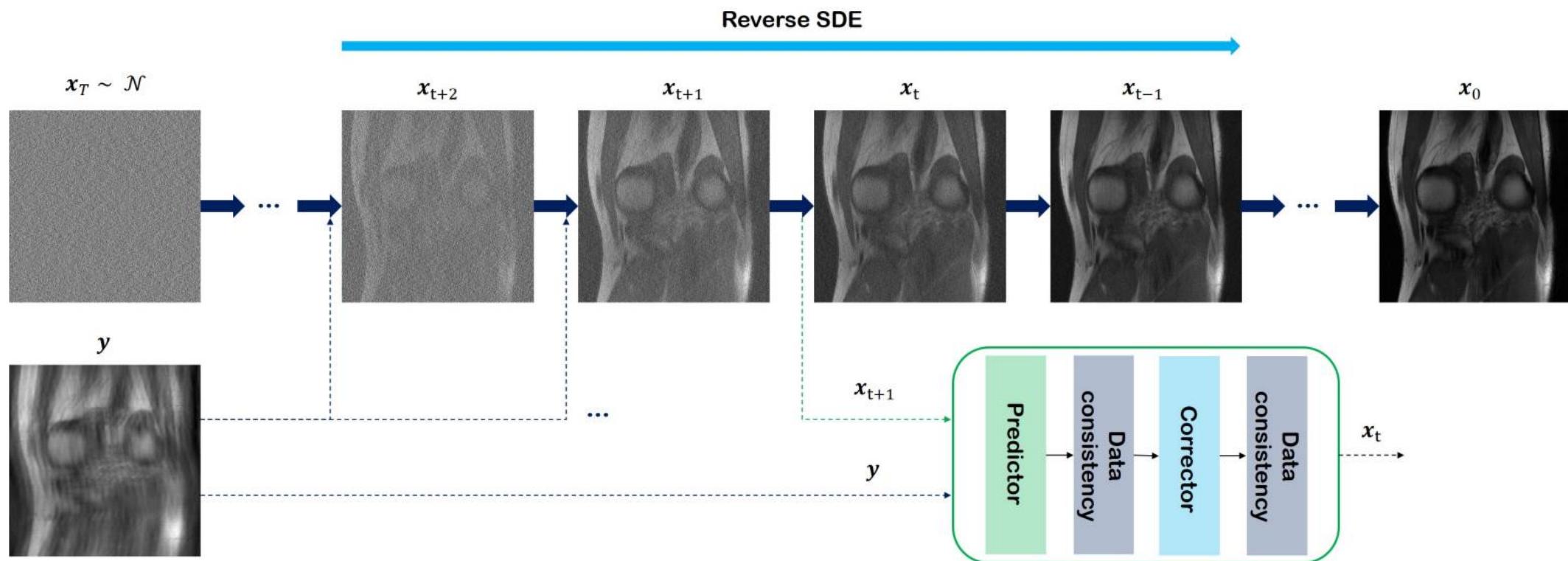
Medical Image Processing:

Score-MRI [38]
Song's Model [205]
R2D2+ [206]

2022 | MRI
2022 | MRI, CT
2022 | MRI

SDE
SDE
SDE

MRI Reconstruction
MRI Reconstruction, CT Reconstruction
MRI Denoising



Sequence Modeling: NLP

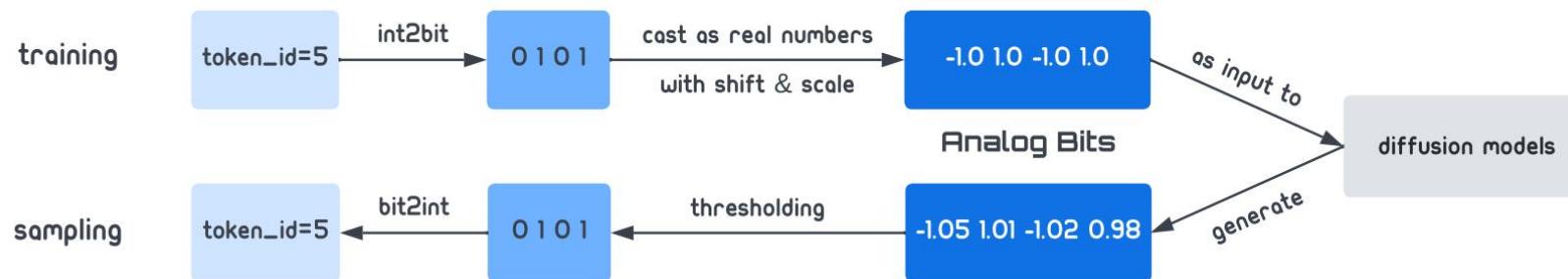
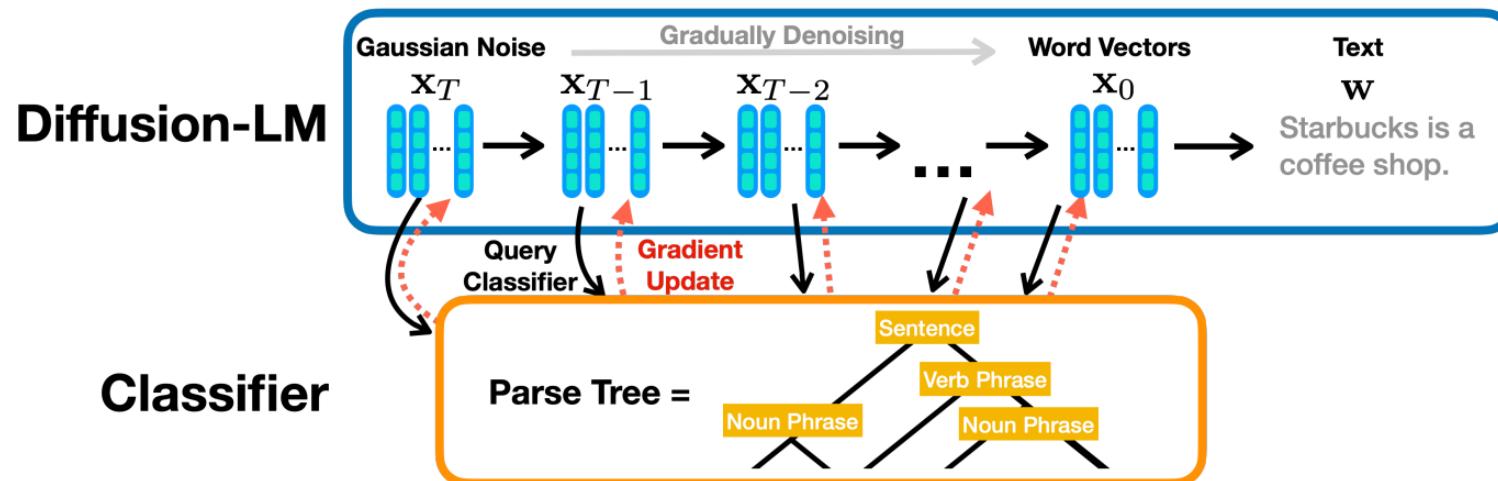
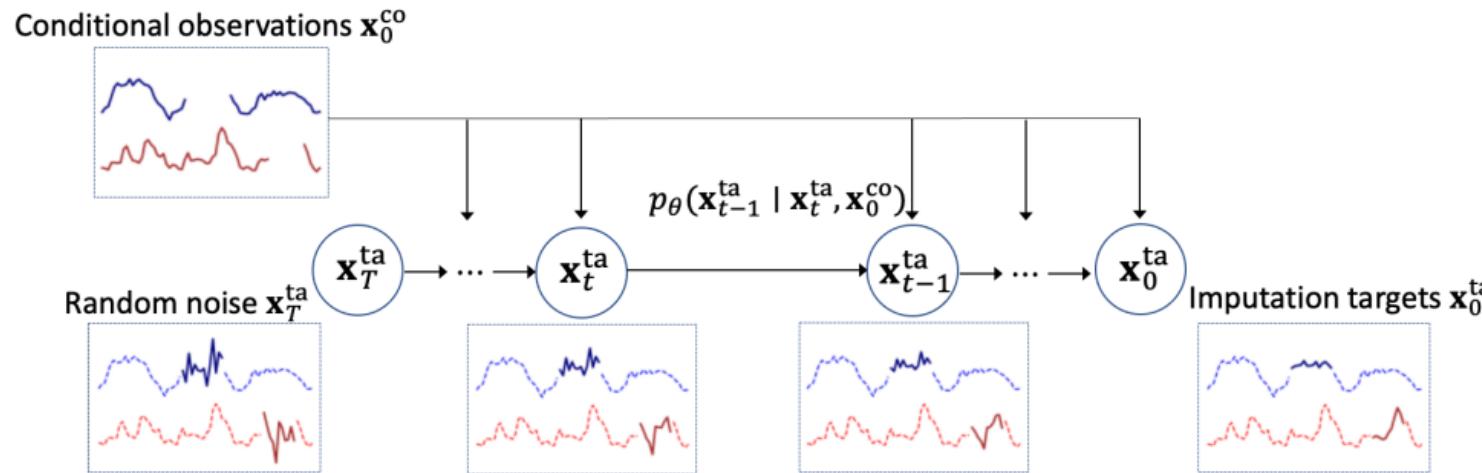


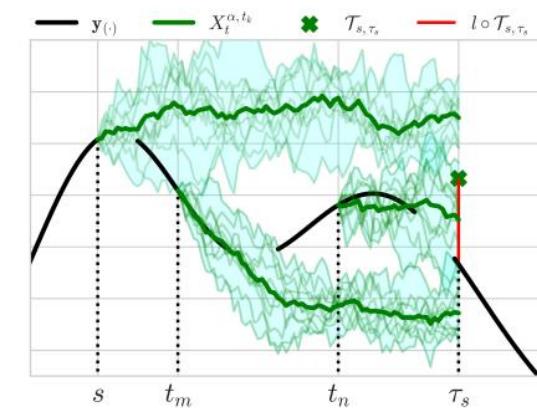
Figure 1: Bit Diffusion: modeling discrete data using continuous diffusion models with analog bits.



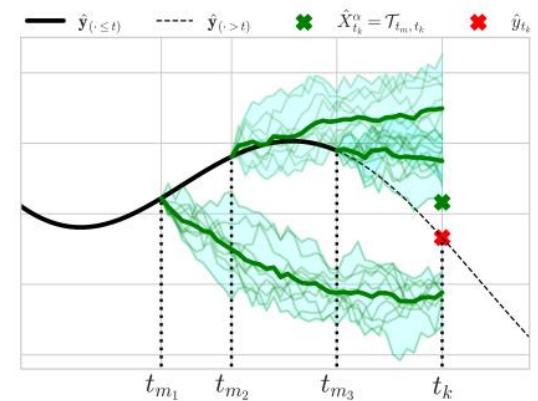
Sequence Modeling: Time Series



Conditioning matters!!!



(a) Network Training with the MFcond loss



(b) Network Inference

Audio: Sound Generation

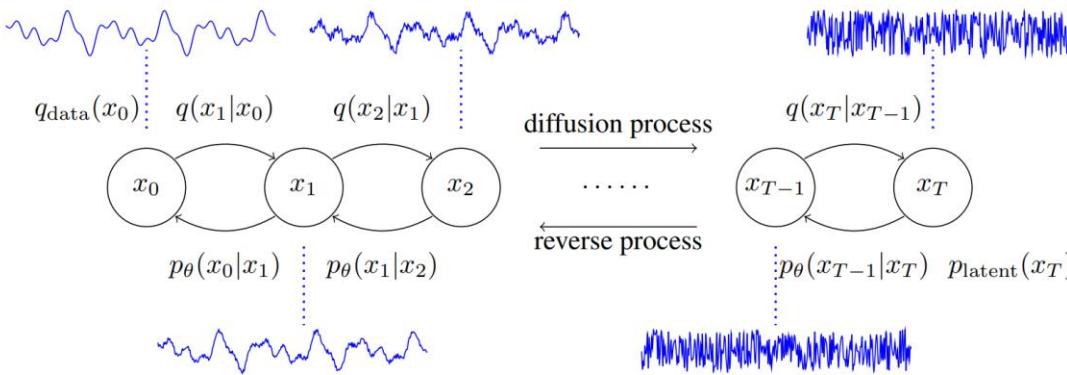


Figure 1: The diffusion and reverse process in diffusion probabilistic models. The reverse process gradually converts the white noise signal into speech waveform through a Markov chain $p_\theta(x_{t-1}|x_t)$.

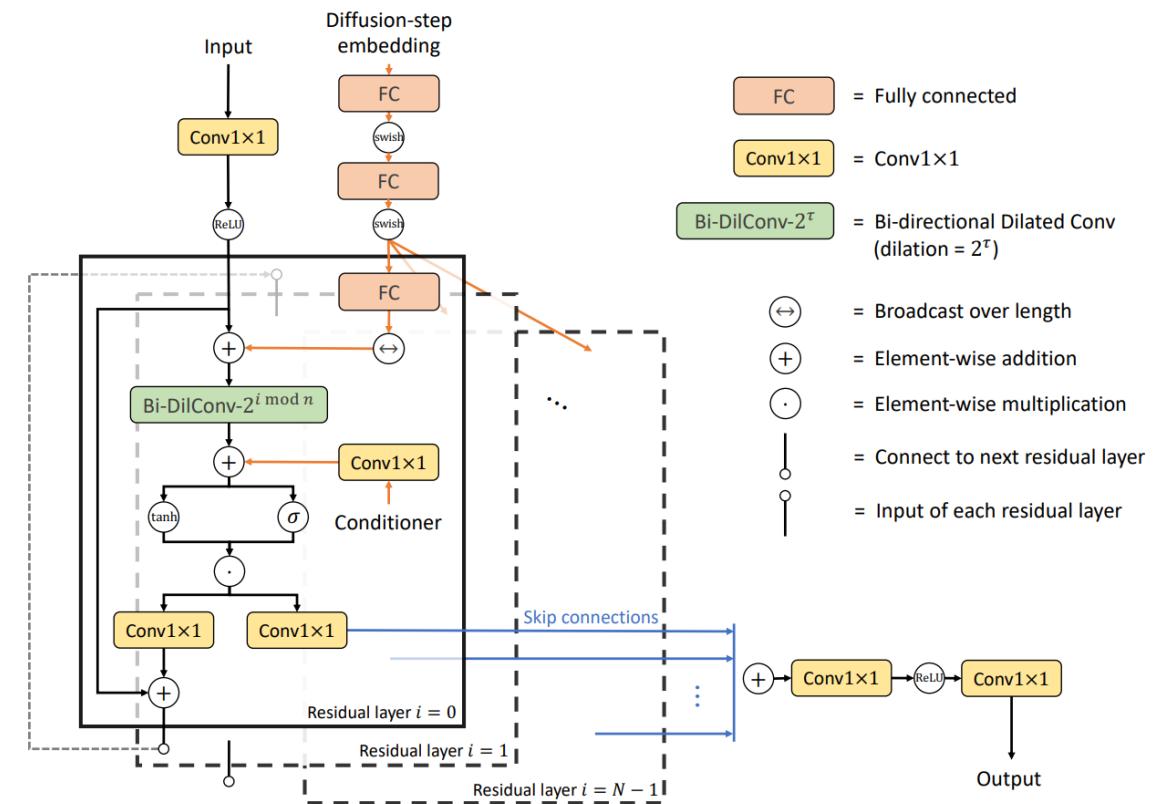


Figure 2: The network architecture of DiffWave in modeling $\epsilon_\theta : \mathbb{R}^L \times \mathbb{N} \rightarrow \mathbb{R}^L$.

Audio: Text to Speech

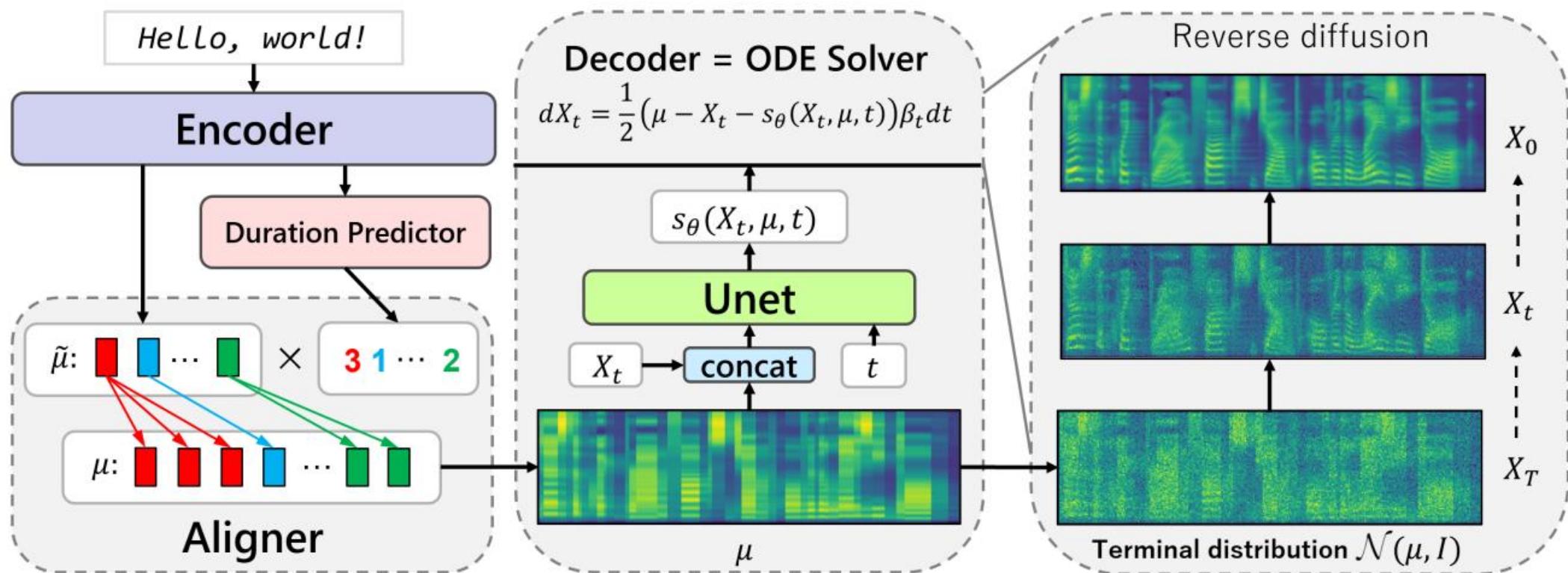
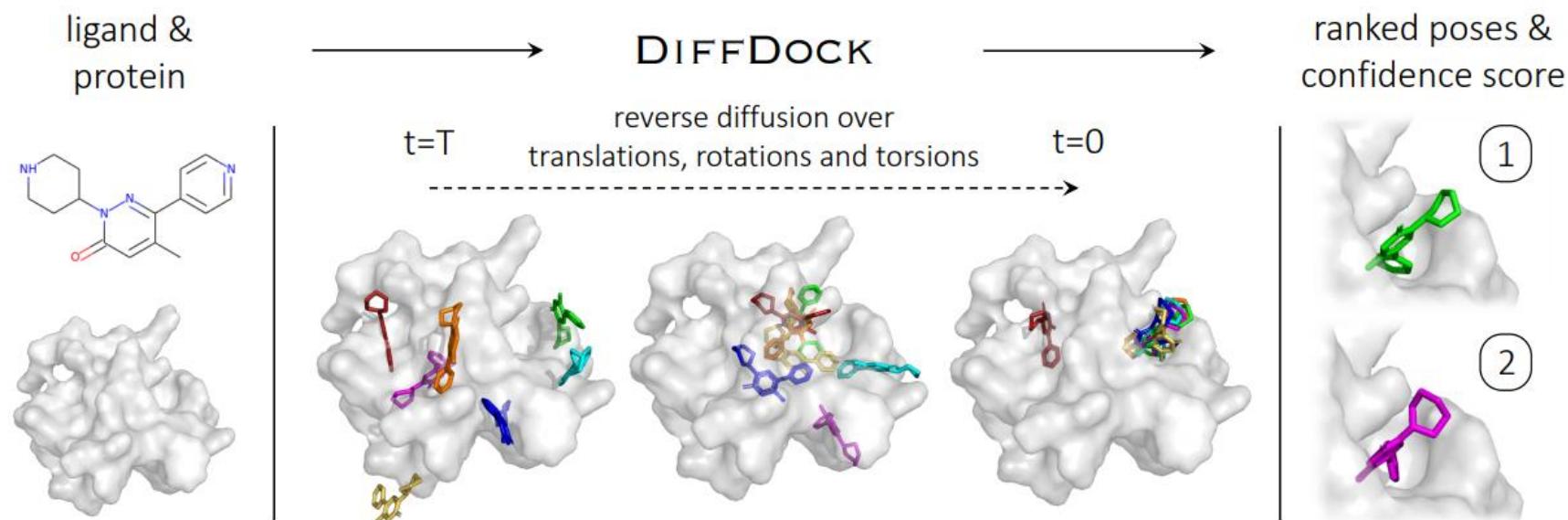
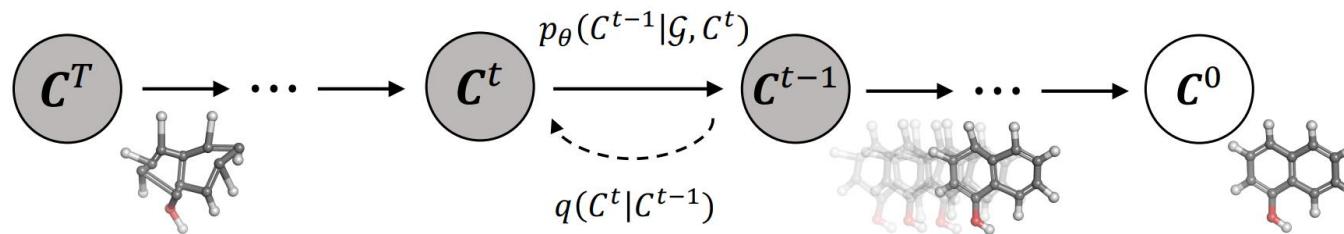
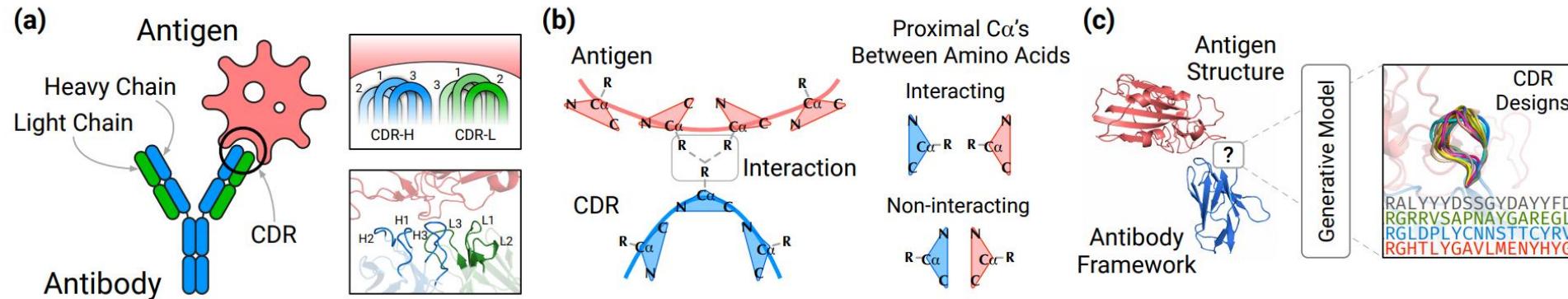


Figure 2. Grad-TTS inference scheme.

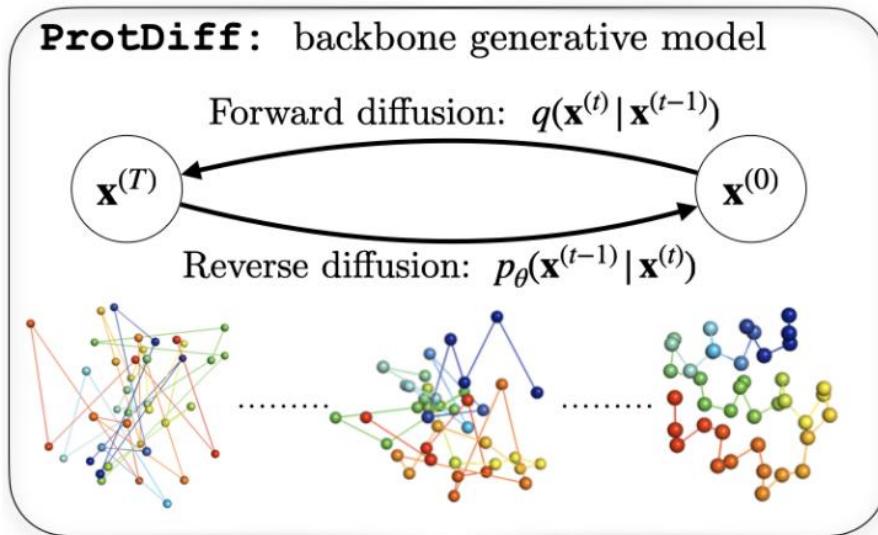
AI for Science: Molecular Generation



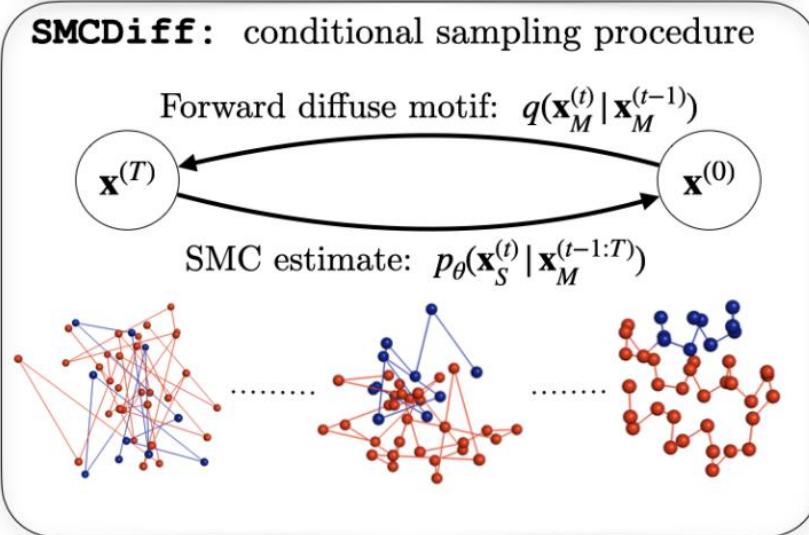
AI for Science: Protein Modeling



ProtDiff: backbone generative model



SMCDiff: conditional sampling procedure



Further Directions:

- Attention on diffusion model class:
 - Prior distribution, transition kernel, sampling algorithm, and diffusion schemes
- Training objective & evaluation metric:
 - Evaluation mismatch, Improved objective for MLE
- Application and inductive bias:
 - Inductive bias, more practice

Materials: Thanks for listening

Paper: <https://arxiv.org/abs/2209.02646>

GitHub: <https://github.com/chq1155/A-Survey-on-Generative-Diffusion-Model>

Be a contributor, Involve in diffusion research,
conduct diffusion applications!



Thanks for listening
and discussion!