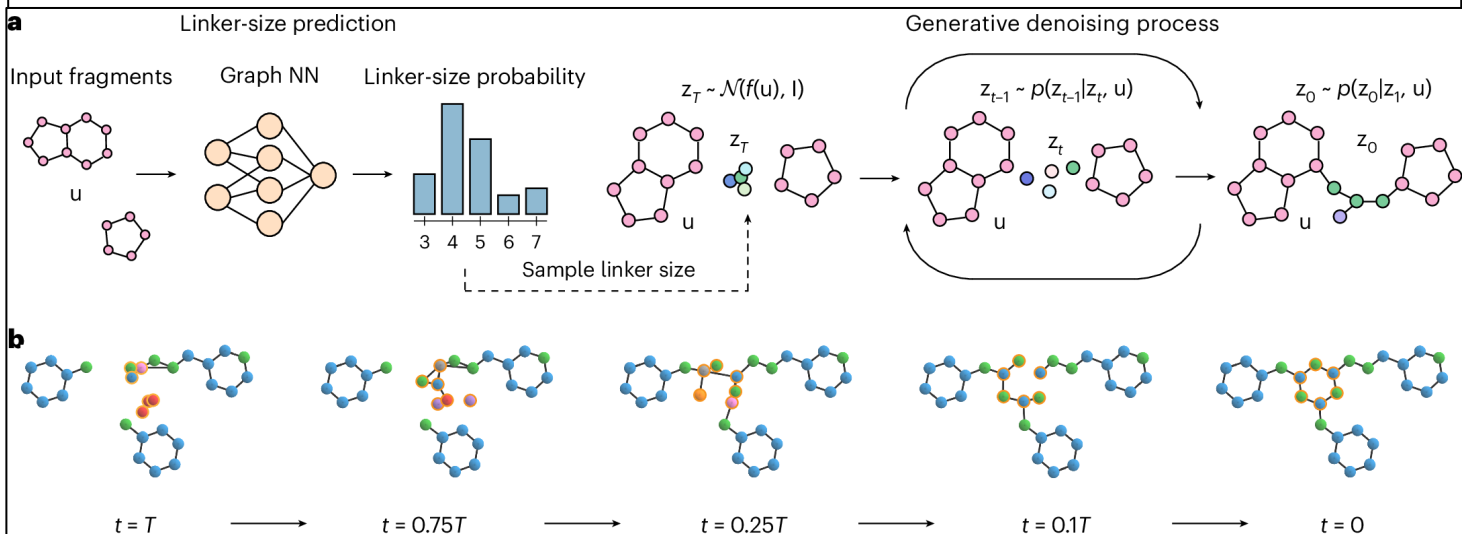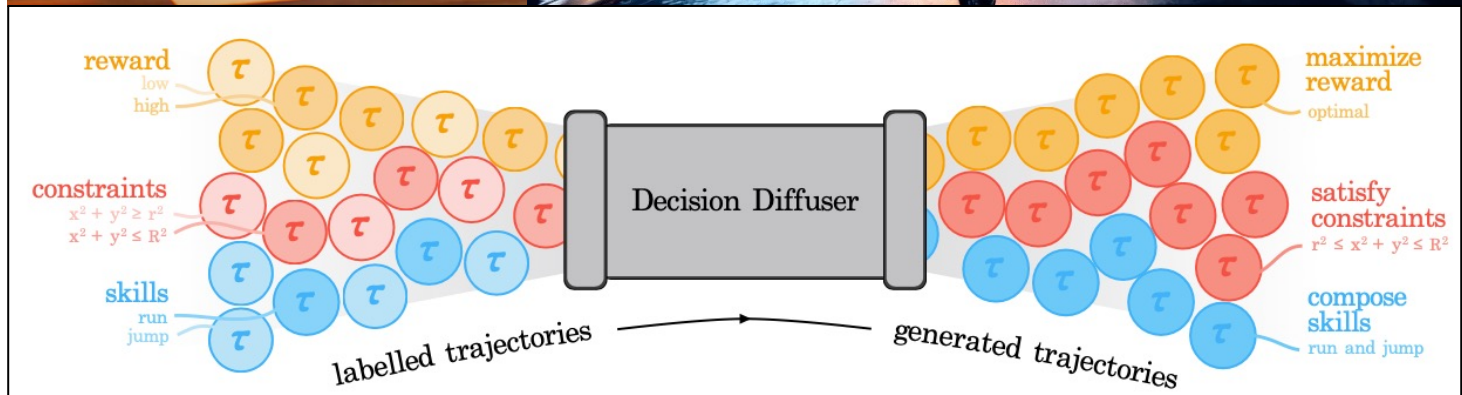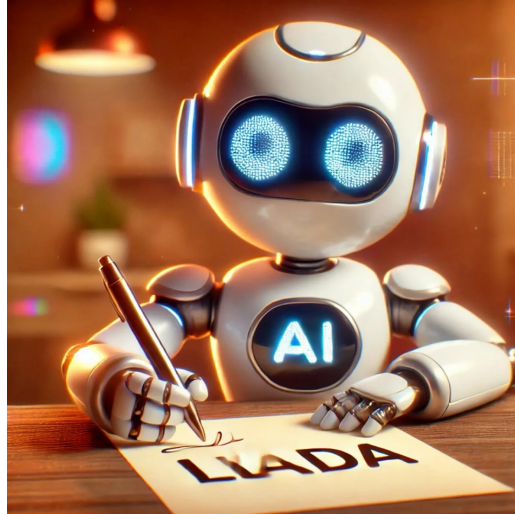# Learning Process & Sampling Complexity of Diffusion Models

Shuai Li

2025.11

# Diffusion Models

- Vision: Sora, etc.
  - SOTA result: Image, 3D, video

- Language: LLaDA

- Multi-modal Models: MMaDA

- Reinforcement Learning

- AI4Science

[1] NZYZOHZLWL, Large Language Diffusion Models, ICLR 2025 DeLTa Workshop, Oral.
[2] YTLZSTW, Multimodal Large Diffusion Language Models, NeurIPS 2025.
[3] ADGTJA, Is Conditional Generative Modeling all you need for Decision Making?, ICLR 2023.
[4] ISVSSFWBC, Equivariant 3D-conditional diffusion model for molecular linker design, Nature Machine Intelligence 2024.

# Theory Helps Training & Sampling

- Solid theoretical foundation helps efficient training & fast sampling:

- Theoretical SDE framework of diffusion family unifies training & sampling[1]

- New training paradigm with SOTA performance: Flow-matching[2]

- 10× Faster sampling algorithm: DPM-Solvers series[3], Analytic-DPM[4]

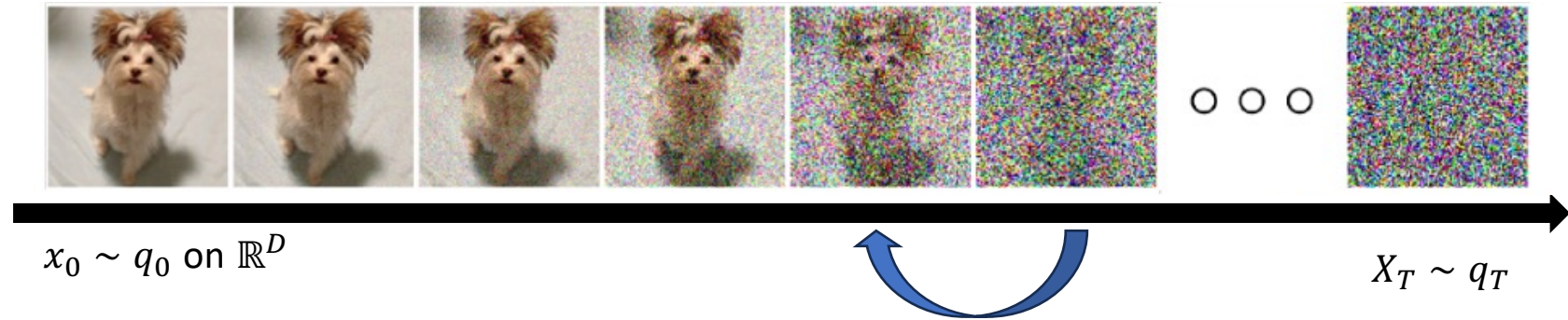[1] SDKKEP, Score-Based Generative Modeling through Stochastic Differential Equations, ICLR 2021.
[2] LG, Flow Straight and Fast: Learning to Generate and Transfer Data with Rectified Flow, ICLR 2023.
[3] LZBCLZ, Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps, NeurIPS 2022.
[4] BLZZ, Analytic-DPM: an Analytic Estimate of the Optimal Reverse Variance in Diffusion Probabilistic Models, ICLR 2022.

# Paradigm of Multi-step Diffusion Models

**Forward Process**

$x_0 \sim q_0$ on $\mathbb{R}^D$

$X_T \sim q_T$

**Core Problem 1: Training Process to Learn Denoising**

**Reverse Process**

$k+1$   $k$

$p_{t_K}$

$Y_t = X_{T-t}$

$Y_0 \sim \mathcal{N}(0, \sigma_T^2 I)$

**Core Problem 2: Sampling Complexity $K$**

# Overview

- Pretraining: Efficient Multi-manifold MoG Model
- Fine-tuning: Good Sharing Latent Guarantees Few-shot Efficiency
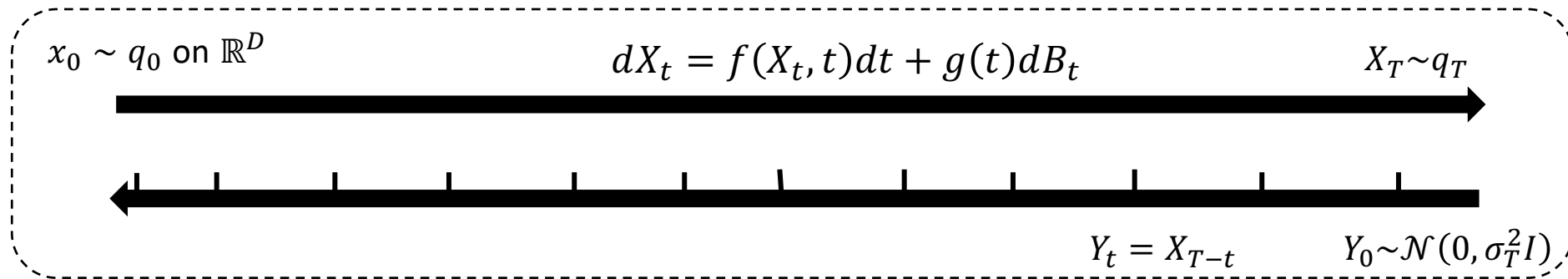- Sampling: Complexity for Multi-step Diffusion Models
- Discretization: Complexity of 1-step Models in Training Phase

# Mathematical Framework of Diffusion Models

$x_0 \sim q_0$ on $\mathbb{R}^D$        $dX_t = f(X_t, t)dt + g(t)dB_t$        $X_T \sim q_T$

$Y_t = X_{T-t}$        $Y_0 \sim \mathcal{N}(0, \sigma_T^2 I)$

> score function

> core to train DM
> unknown

- $dY_t = \left[ -f(Y_t, T-t) + \frac{1+\eta^2}{2} g^2(T-t) \nabla \log q_{T-t}(Y_t) \right] dt + \eta g(T-t) dB_t, \eta \in [0,1]$

- Score matching training objective:

> conditional distribution
> known

$$\min_{s \in \mathfrak{F}} \hat{\mathcal{L}}(s) = \frac{1}{n} \sum_{i=1}^{n} \frac{1}{T-\delta} \int_{\delta}^{T} \mathbb{E}_{X_t | X_0 = X_i} \left[ \| \nabla \log q_t(X_t | X_0) - s(X_t, t) \|_2^2 \right] dt$$

# Learning Faces Curse of Dimension

- Minimiser $s_\theta \in \text{argmin}_\Theta \hat{\mathcal{L}}(s)$ satisfies

  Estimation Error= $\frac{1}{T-\delta} \int_\delta^T \mathbb{E}_{q_t} [\|\nabla \log q_t(X_t) - s_\theta(X_t, t)\|_2^2] dt < O(n^{-1/D})$

  covering number & concentration

  $D = 3 \times 256 \times 256 \approx 2 \times 10^5$

- Good training requires training data size $n = O(10^{10^5})$ Huge!!

- Efficient training needs utilizing data structure!

# Data Structures: Existing Works

| Manifold Modeling | Latent | | # of Parameters | Estimation Error |
|---|---|---|---|---|
| Full Space [1] | General | $X$ | $O(D^{D+1})$ | $O(n^{-1/D})$ |
| Full Space [2] | Mixture of Gaussian (MoG) | $X \sim \sum_{m=1}^{M} \pi_m \mathcal{N}(\mu_m, \Sigma_m)$ | $O(MD^2)$ | $O(\frac{\sqrt{DM}}{\sqrt{n}})$ |
| Low-dim manifold [3] | General | $X = Az$, with $A \in \mathbb{R}^{D \times d}$ | $O(Dd + d^{d+1})$ | $O(n^{-\frac{2}{d}})$ |
| Multi-manifold | General | $X = \sum_{\ell=1}^{L} \pi_\ell A_\ell z_\ell$, with $A_\ell \in \mathbb{R}^{D \times d}$ | $O(LDd + \boxed{Ld^{d+1}})$ | $O(\sqrt{L} n^{-\frac{2}{d}})$ |
| Multi-manifold [4] | Gaussian | $X \sim \sum_{\ell=1}^{L} \pi_\ell \mathcal{N}(\cdot; 0, A_\ell A_\ell^\top)$ | $O(LDd)$ | $O(\frac{\sqrt{dL}}{\sqrt{n}} + \boxed{\text{Const}})$ |

[1] OAS, Diffusion Models are Minimax Optimal Distribution Estimators, ICML 2023.
[2] SCK, Learning mixtures of gaussians using the ddpm objective, NeurIPS 2023.
[3] CHZW, Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data, ICML 2023.
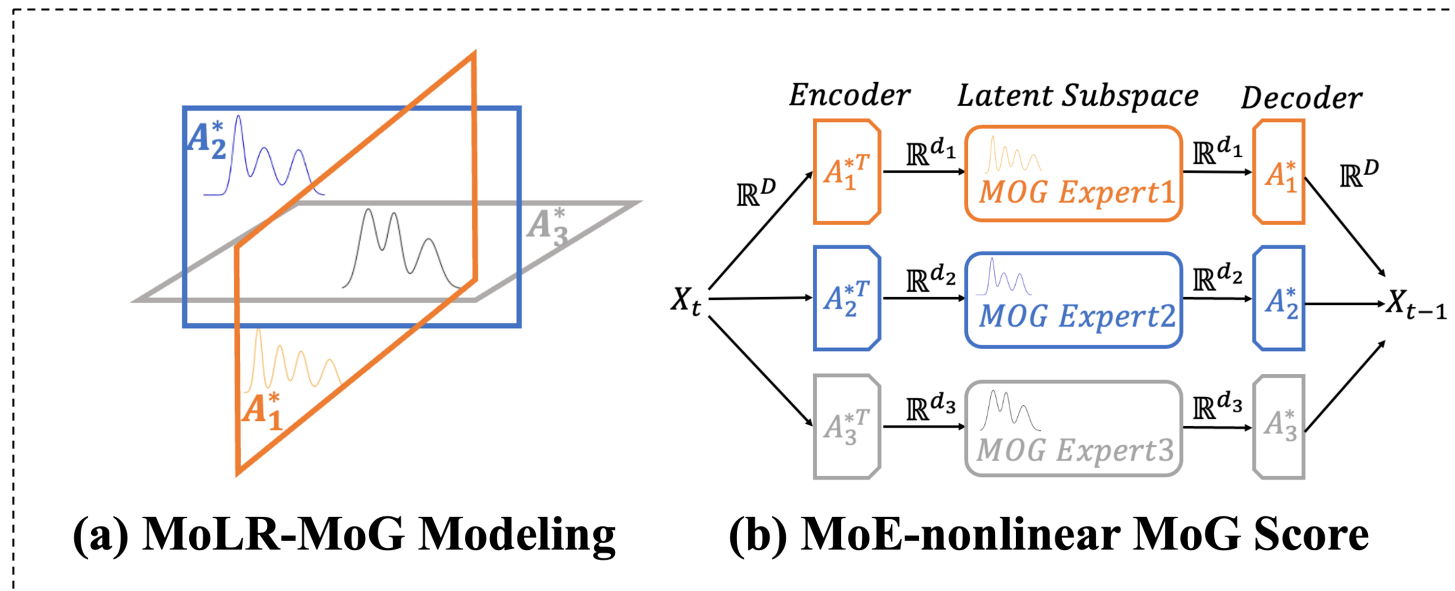[4] WZZCMQ. Diffusion models learn low-dimensional distributions via subspace clustering, NeurIPS 2024 M3L Workshop.
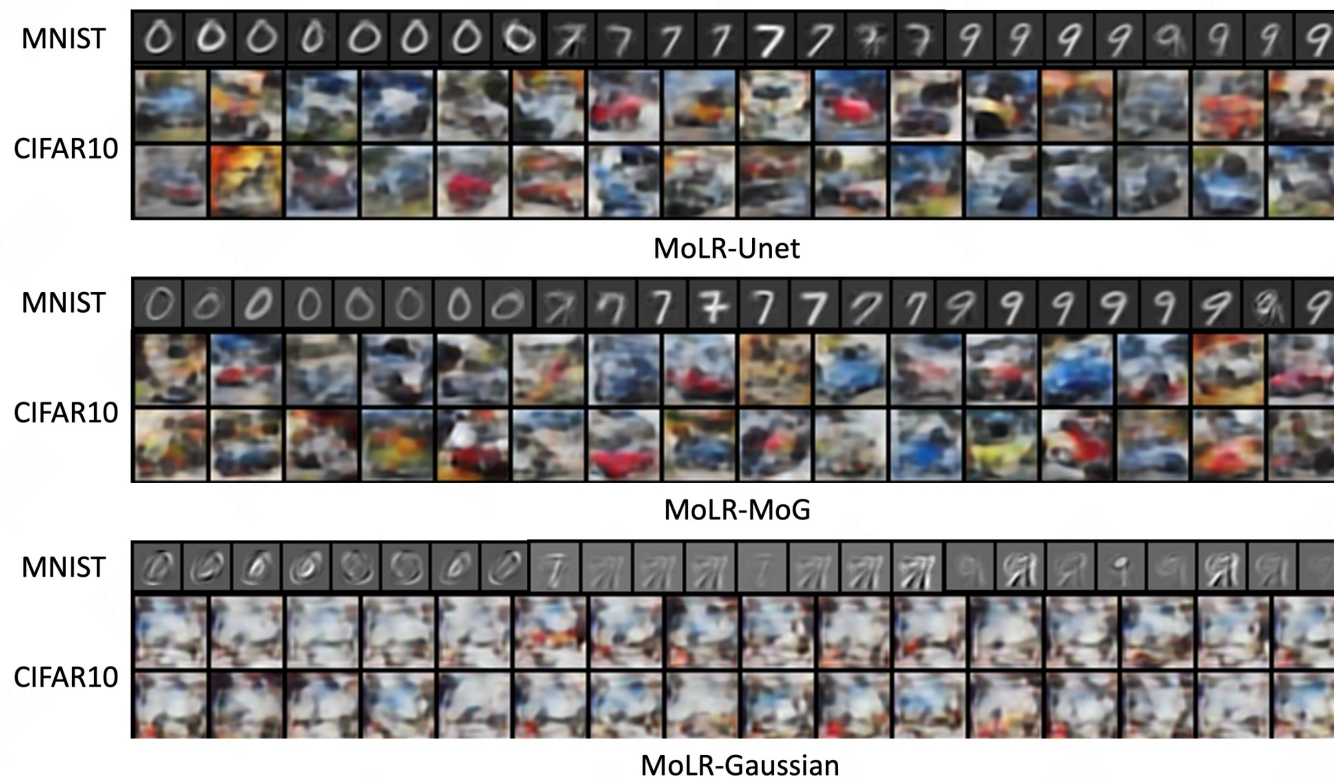
# Multi-manifold Mixture-of-Gaussian Modeling

- $X \sim \sum_{\ell=1}^{L} \pi_\ell \sum_{m=1}^{M} \pi_{\ell,m} \mathcal{N}\left(\cdot; A_\ell \mu_{\ell,m}, A_\ell \Sigma_{\ell,m} A_\ell^\top\right)$  Most general!

- Theorem. Its estimation error satisfies

$$\frac{1}{T-\delta} \int_\delta^T \mathbb{E}_{q_t}\left[\|\nabla \log q_t(X_t) - s_\theta(X_t, t)\|_2^2\right] dt < O\left(\frac{\sqrt{LM}\sqrt{dL}}{\sqrt{n}}\right)$$
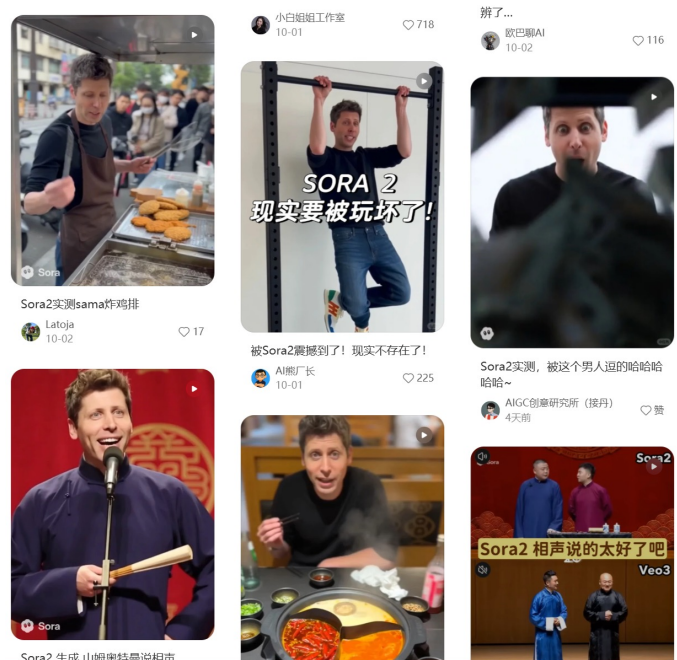


**(a) MoLR-MoG Modeling**    **(b) MoE-nonlinear MoG Score**

[1] YLJC**L**, Multi-Subspace Multi-Modal Modeling for Diffusion Models: Estimation, Convergence and Mixture of Experts (In submission)

# Much Smaller Model w/ Sufficiently Good Performance

| Latent | # of Parameters | Estimation Error | MNIST Acc/ Performance |
|--------|-----------------|------------------|------------------------|
| General | $O(LDd + Ld^{d+1})$ | $O(\sqrt{L}n^{-\frac{2}{d}})$ | 0.96 ✔ Deep NN |
| Mixture of Gaussian | $O(LDd + Ld^2)$ | $O\left(\frac{\sqrt{LM}\sqrt{dL}}{\sqrt{n}}\right)$ | 0.89 ✔ 2-layer NN |
| Gaussian | $O(LDd)$ | $O(\frac{\sqrt{dL}}{\sqrt{n}} + \text{Const})$ | 0.08 ✘ Linear NN |



MoLR-Unet

MoLR-MoG

MoLR-Gaussian

- YLJC**L**, Multi-Subspace Multi-Modal Modeling for Diffusion Models: Estimation, Convergence and Mixture of Experts (In submission)
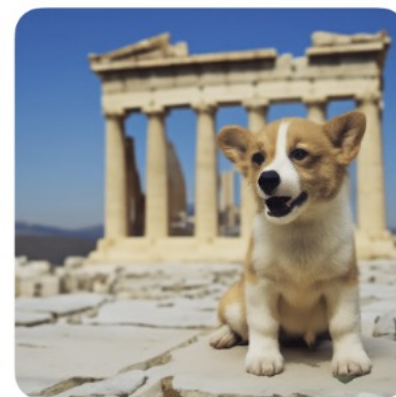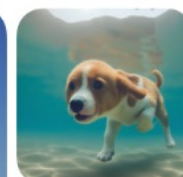
# Overview

- Pretraining: Efficient Multi-manifold MoG Model
- Fine-tuning: Good Sharing Latent Guarantees Few-shot Efficiency
- Sampling: Complexity for Multi-step Diffusion Models
- Discretization: Complexity of 1-step Models in Training Phase

Input images

**Few-shot Fine-tuning** is key to the customized creation

but no theory supports effective information sharing

# Few-shot Fine-tuning

- Pretrain w/ large source data (2.3 Billion): $\{X_{s,i}\}_{i=1}^{n_s} \sim q_0^s$ on $\mathbb{R}^D$

- $\min_{s \in \text{Source } \mathfrak{P}} \hat{\mathcal{L}}_s(s) = \frac{1}{n_s} \sum_{i=1}^{n_s} \frac{1}{T-\delta} \int_\delta^T \mathbb{E}_{X_t|X_0=X_{s,i}} \left[ \|\nabla \log q_t^s(X_t|X_0) - s(X_t, t)\|_2^2 \right] dt$

  e.g. 882M

- Estimation error $O\left(n_s^{-\frac{2}{d}}\right)$ — Tolerable!

- Fine-tune with limited target data (~10 images): $\{X_{\text{ta},i}\}_{i=1}^{n_{\text{ta}}} \sim q_0^{\text{ta}}$

- $\min_{s \in \text{Target } \mathfrak{P}} \hat{\mathcal{L}}_{\text{ta}}(s) = \frac{1}{n_{\text{ta}}} \sum_{i=1}^{n_{\text{ta}}} \frac{1}{T-\delta} \int_\delta^T \mathbb{E}_{X_t|X_0=X_{\text{ta},i}} \left[ \|\nabla \log q_t^{\text{ta}}(X_t|X_0) - s(X_t, t)\|_2^2 \right] dt$

  e.g. 1.5M 0.17%

- Estimation error $O\left(n_{\text{ta}}^{-\frac{2}{d}}\right)$ — Meaningless!
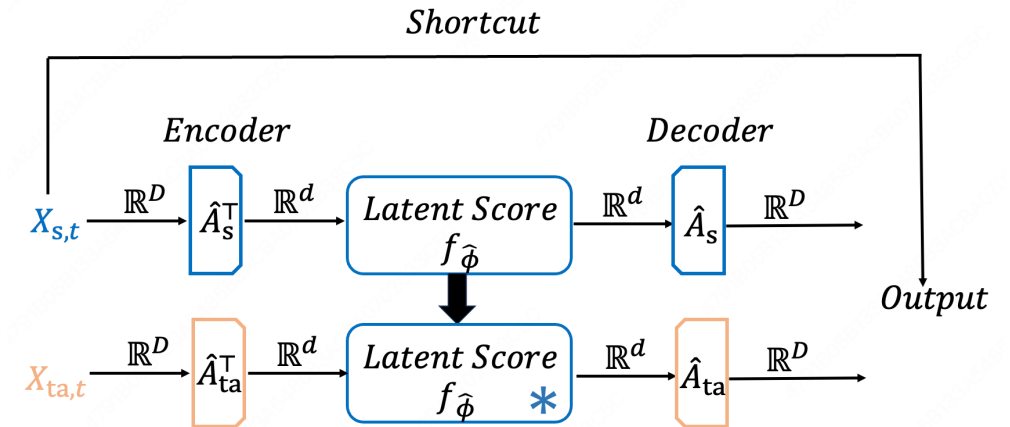
# Information-sharing Model Design

- Empirical works share most parameters and fine-tune key parameters

- Assumption. The source and target data admit linear structure and share latent space $X_s = A_s z$ and $X_{ta} = A_{ta} z, z \in \mathbb{R}^d$



- Then the score function is

$$\nabla \log q_t^{ta}(X) = A_{ta} \nabla \log q_t^{\text{Latent}}(A_{ta}^\top X) - \frac{1}{\sigma_t^2}(I_D - A_{ta}A_{ta}^\top)X$$

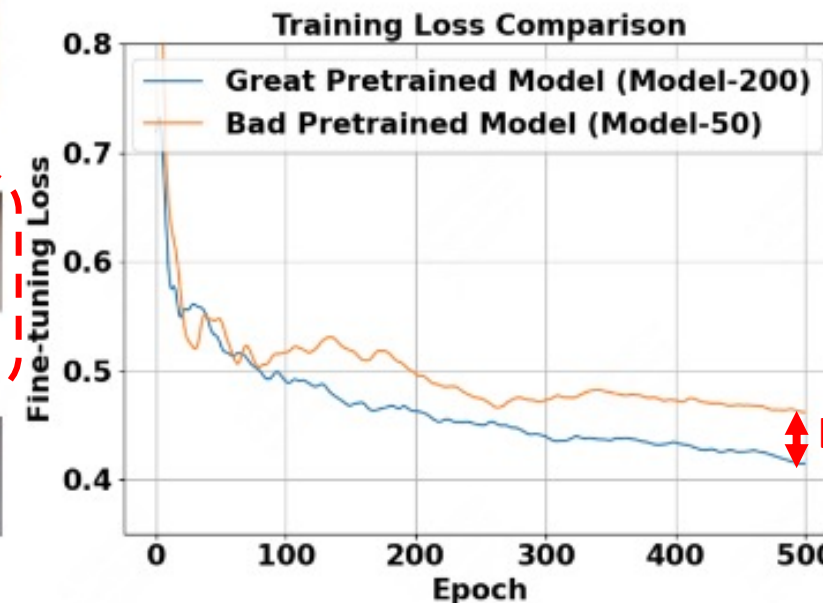Shared Latent Score

# Bad Latent Leads to Large Estimation Error



Strange Generation!

(i) Target Dataset

(ii) Model-50 (Underfitting Bad Pretrained Model)

(iii) Model-200 (Good Pretrained Model)

Training Loss Comparison

Large Gap!

- **Theorem**. W/ bad latent

$$\frac{1}{T-\delta}\int_{\delta}^{T} \mathbb{E}_{q_t^{\mathsf{ta}}}\left[\|\nabla \log q_t^{\mathsf{ta}}(X_t) - s_\theta(X_t, t)\|_2^2\right] dt \geq \text{Const}$$

[1] YLJC**L**, Evaluating the Role of Great Pre-trained Diffusion Models in Few-shot Phase: Warm-up and Acceleration (In submission).

# Bad Latent Suffers Bad Local Minima



Fine-tuning Results based on Great Pre-trained Models (SD3 Medium)

Fine-tuning Results based on *Overfitting* Bad Pre-trained Models (SD3 Medium with 1k overfitting steps)

A *cat* on top of a wooden floor

A *cat* in a chef outfit

A *cat* with a city in the background

A *cat* wearing a yellow shirt

A *cat* in a police outfit

Prompt cat but results in dog figure

Bad latent fails to fit target feature!

- **Theorem**. W/ bad latent, $\exists\, s_\theta^{\text{few-shot}} \neq s_{\theta^*}^{\text{few-shot}}$ s.t. $\dfrac{\partial s_\theta^{\text{few-shot}}}{\partial \theta} \approx 0$

[1] YLJC**L**, Evaluating the Role of Great Pre-trained Diffusion Models in Few-shot Phase: Warm-up and Acceleration (In submission)

# Good Latent Secures Efficiency

- **Theorem**. The estimation error of few-shot diffusion model is

$$\frac{1}{T-\delta} \int_{\delta}^{T} \mathbb{E}_{q_t^{\text{ta}}} \left[ \left\| \nabla \log q_t^{\text{ta}}(X_t) - s_{\hat{A}_{\text{ta}}, \hat{\phi}}(X_t, t) \right\|_2^2 \right] dt \leq O\left( n_{\text{ta}}^{-\frac{1}{2}} + n_s^{-\frac{2}{d}} \right)$$

Guarantee good latent

- $O\left( n_{\text{ta}}^{-\frac{1}{2}} \right)$ explains why $5 - 8$ images are enough for few-shot fine-tuning

Table 1: The requirement of $n_{ta}$ in popular datasets. We use latent dimension in Pope et al. (2021).

| Dataset | CIFAR-10 | CIFAR-100 | CelebA | MS-COCO | ImageNet |
|---|---|---|---|---|---|
| Dataset Size | $6 \times 10^4$ | $6 \times 10^4$ | $2 \times 10^5$ | $3.3 \times 10^5$ | $1.2 \times 10^6$ |
| Latent Dimension | 25 | 22 | 24 | 37 | 43 |
| The Requirement of $n_{ta}$ | 6 | 8 | 8 | 5 | 5 |

- YHCJW**L**, Few-shot diffusion models escape the curse of dimensionality. NeurIPS 2024.

# Good Latent Leads to Good Landscape

- Theorem. With a good shared latent, the landscape of the few-shot optimization is $\kappa$-strongly convex w/ convergence rate
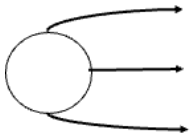
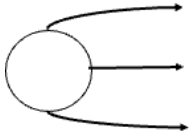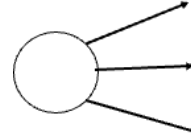$$\left\| \hat{A}_{\text{ta}}^{(i)} \hat{A}_{\text{ta}}^{(i)\top} - A_{\text{ta}} A_{\text{ta}}^\top \right\|_F \leq \left( \frac{\kappa - 1}{\kappa + 1} \right)^i \|A_{\text{ta}}\|_F \left\| \hat{A}_{\text{ta}}^{(0)} - A_{\text{ta}} \right\|_F$$

# Overview

- Pretraining: Efficient Multi-manifold MoG Model
- Fine-tuning: Good Sharing Latent Guarantees Few-shot Efficiency
- Sampling: Complexity for Multi-step Diffusion Models
- Discretization: Complexity of 1-step Models in Training Phase

# Common Forward Processes

$$dX_t = f(X_t, t)dt + g(t)dB_t \qquad T$$



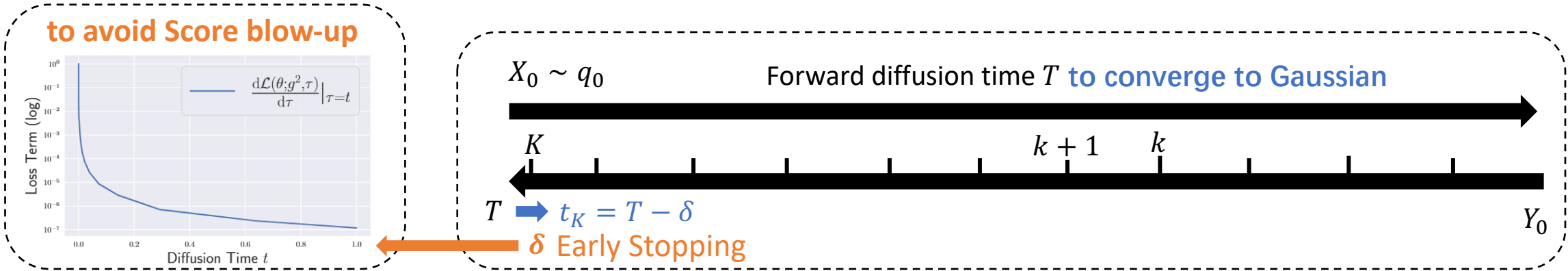|  |  | Trajectory | Forward Distribution |  |
|---|---|---|---|---|
| Variance Preserving (VP) [1] | $f(X_t, t) = -\frac{1}{2}X_t$ <br> $g(t) = 1$ | | $\mathcal{N}(0, I_D)$ | stability.ai <br> Midjourney |
| Variance Exploding (VE–SMLD) [2] | $f(X_t, t) = 0$ <br> $g(t) = \sqrt{2}$ | | $\mathcal{N}(0, TI_D)$ | Stanford University |
| Variance Exploding (VE–EDM) [3] | $f(X_t, t) = 0$ <br> $g(t) = \sqrt{2t}$ | | $\mathcal{N}(0, T^2 I_D)$ | |
| Rectified Flow (RF) [4] | $X_t = (1-t)X_0 + tZ$ <br> $t \in [0,1]$ | | $\mathcal{N}(0, I_D)$ | |

[1] HJA, Denoising diffusion probabilistic models, NeurIPS 2020.
[2] SE, Generative modeling by estimating gradients of the data distribution, NeurIPS 2019.
[3] KAAL, Elucidating the Design Space of Diffusion-Based Generative Models, NeurIPS 2022.
[4] LG, Flow Straight and Fast: Learning to Generate and Transfer Data with Rectified Flow, ICLR 2023.

# Sampling Complexity: Objective



to avoid Score blow-up

$X_0 \sim q_0$     Forward diffusion time $T$ **to converge to Gaussian**

$K$        $k+1$   $k$

$T \Rightarrow t_K = T - \delta$

$\delta$ Early Stopping

$Y_0$

- Objective:

  With accurate score $\|\nabla \log q_t(X) - s_\theta(X, t)\|_2^2 \leq \epsilon_{\text{score}}^2$

  Minimize sample complexity $K$ s.t.

  $$\text{KL}\left(p_{t_K}, q_\delta\right) \leq \epsilon_{\text{KL}}^2 \text{ and } W_2^2(q_0, q_\delta) \leq \epsilon_{W_2}^2$$

# Sample Complexity: General Guarantee for Reverse SDE

- Theorem. Sample complexity can be divided by

$$\text{KL}\big(p_{t_K}, q_\delta\big) \leq \underset{\substack{\text{Convergence of} \\ \text{Forward Process}}}{\text{KL}\big(\mathcal{N}(0, \sigma_T^2), q_T\big)} + \sum_{k=0}^{K-1} \mathbb{E}_{q_{t_k}(x)} \underset{\text{Discretization}}{\text{KL}\left(p_{t_{k+1}|t_k}(\cdot \,|x), q_{t_{k+1}|t_k}(\cdot \,|x)\right)}$$

$$\leq D^2 m_T / \sigma_T^2 + D^2 (T/\delta)^{\frac{1}{a}} / K \leq \tilde{O}\big(\epsilon_{\text{KL}}^2\big)$$

- Then the sample complexity requires $K = O(D^2 (T/\delta)^{\frac{1}{a}} / \epsilon_{\text{KL}}^2)$ where $\delta$ satisfies

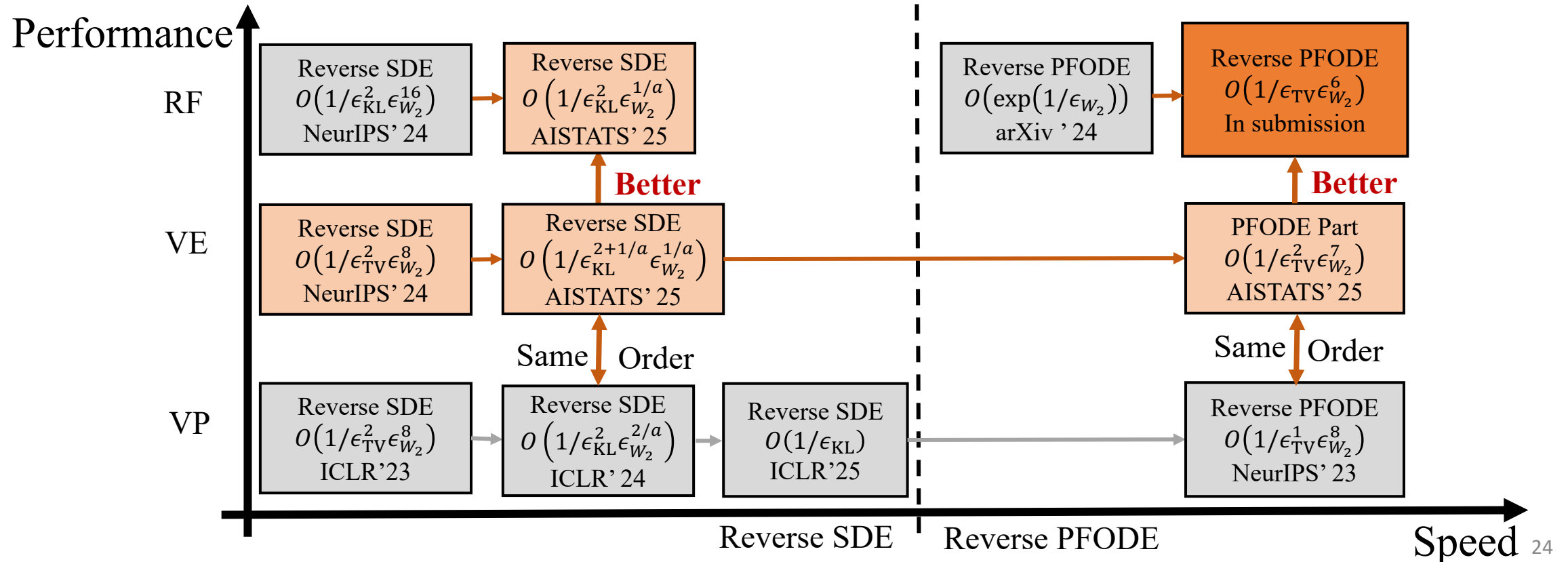$$W_2^2(q_0, q_\delta) \leq \sigma_\delta^2 \leq \epsilon_{W_2}^2$$

# Sample Complexities

| | $m_T$ | $\sigma_T^2$ | $T$: $\mathrm{KL}(\mathcal{N}(0,\sigma_T^2),q_T)$ $\leq \frac{m_T}{\sigma_T^2} \leq \epsilon_{\mathrm{KL}}^2$ | $\sigma_\delta^2$ | $\delta$: $\mathrm{W}_2^2(q_0,q_\delta) \leq$ $\sigma_\delta^2 \leq \epsilon_{W_2}^2$ | $K$: $O(D^2(T/\delta)^{\frac{1}{a}}/\epsilon_{\mathrm{KL}}^2)$ |
|---|---|---|---|---|---|---|
| **VP** | $e^{-T}$ | $1-e^{-2T}$ | $\log(1/\epsilon_{\mathrm{KL}})$✓ | $\delta$ | $\epsilon_{W_2}^2$ ✗ | $O\left(D^2/\epsilon_{\mathrm{KL}}^2\epsilon_{W_2}^{2/a}\right)$ |
| **VE (SMLD)** | $1$ | $T$ | $1/\epsilon_{\mathrm{KL}}^2$ ✗ | $\delta$ | $\epsilon_{W_2}^2$ ✗ | $O\left(D^2/\epsilon_{\mathrm{KL}}^{2+2/a}\epsilon_{W_2}^{2/a}\right)$ |
| **VE (EDM)** | $1$ | $T^2$ | $1/\epsilon_{\mathrm{KL}}$ ✗ | $\delta^2$ | $\epsilon_{W_2}$ ✓ | $O\left(D^2/\epsilon_{\mathrm{KL}}^{2+1/a}\epsilon_{W_2}^{1/a}\right)$ |
| **RF** | $1$ | $1$ | $1$✓ | $\delta^2$ | $\epsilon_{W_2}$ ✓ | $O\left(D^2/\epsilon_{\mathrm{KL}}^2\epsilon_{W_2}^{1/a}\right)$ |

- VP better in $T$ and VE (EDM) better in $\delta$

- RF better in both $T$ and $\delta$ and thus has a better complexity

# Results Extend to PRODE

[1] YWJL, Leveraging drift to improve sample complexity of variance exploding diffusion models. NeurIPS 2024.
[2] YJL, The Polynomial Iteration Complexity for Variance Exploding Diffusion Models: Elucidating SDE and ODE Samplers. AISTATS 2025.
[3] YZJCL, Elucidating Rectified Flow with Deterministic Sampler: Polynomial Discretization Complexity for Multi and One-step Models. Arxiv.

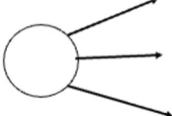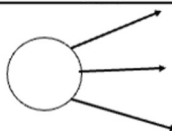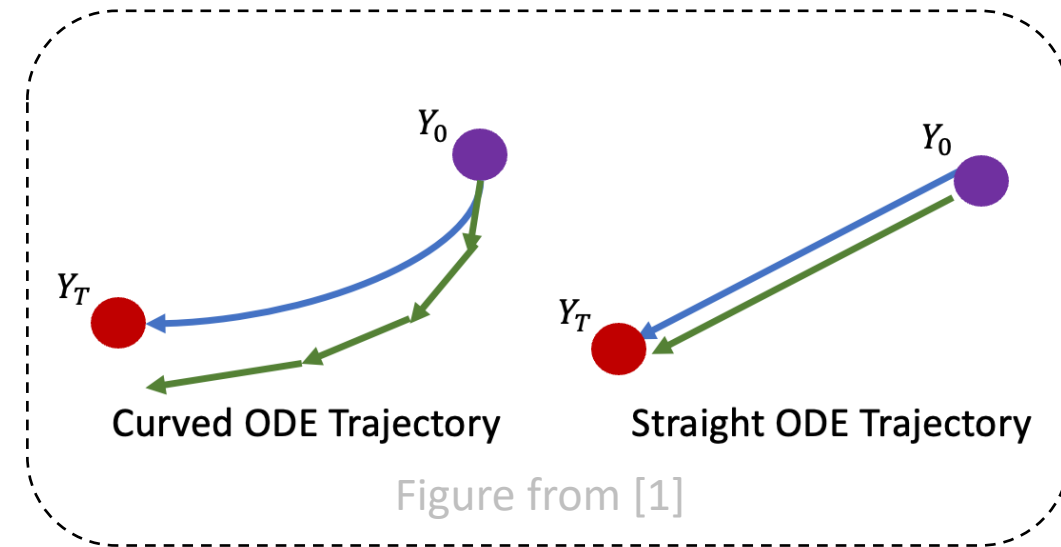- Reverse SDE generate diverse samples while PFODE generate fast

# Overview

- Pretraining: Efficient Multi-manifold MoG Model
- Fine-tuning: Good Sharing Latent Guarantees Few-shot Efficiency
- Sampling: Complexity for Multi-step Diffusion Models
- Discretization: Complexity of 1-step Models in Training Phase

# Linear Trajectory & PFODE Achieve 1-step Generation

- PFODE generate deterministically compared to reverse SDE

- VE-EDM and RF have linear trajectory

| | | |
|---|---|---|
| Variance Exploding (VE−EDM) [3] | $f(X_t, t) = 0$ $g(t) = \sqrt{2t}$ | |
| Rectified Flow (RF) [4] | $X_t = (1-t)X_0 + tZ$ $t \in [0,1]$ | |



Curved ODE Trajectory     Straight ODE Trajectory

Figure from [1]

[1] LZMPL, InstaFlow: One Step is Enough for High-Quality Diffusion-Based Text-to-Image Generation, ICLR 2024.

# 1-Step Mapping Function from Multi-step

- For PFODE reverse process of multi-step diffusion models

$$\mathrm{d}Y_t = v(Y_t, t)\mathrm{d}t, Y_0 \sim q_T$$

  the corresponding 1-step mapping function (by integral) is

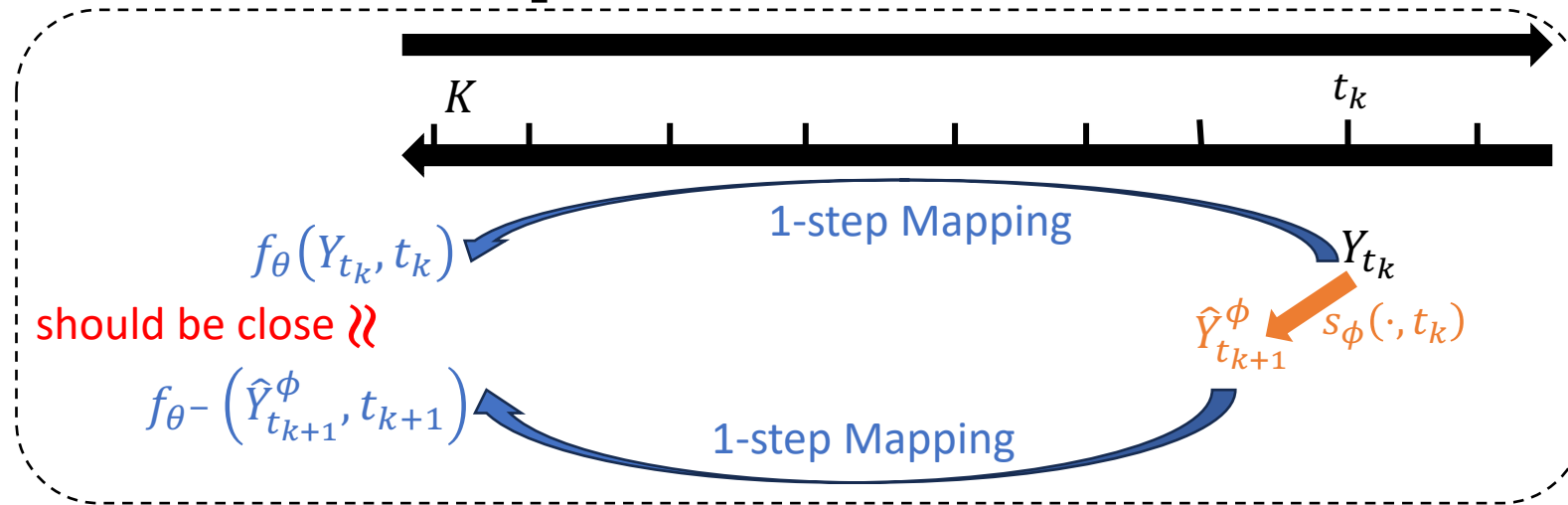$$f(Y_t, t) = Y_{T-\delta} = X_\delta \approx X_0, \forall t \in [0, T - \delta]$$

to avoid Score blow-up

- Use NN $f_\theta(Y_t, t)$ to approximate 1-step mapping function $f$

# What is a Good Optimization Objective?

- Consistency distillation to learn good 1-step mapping [1]

$$\mathcal{L}_{\mathrm{CD}}^K(\boldsymbol{\theta}, \boldsymbol{\theta}^-; \boldsymbol{\phi}) := \mathbb{E}_{X_0}\left[\mathbb{E}_{Y_{t_k}|X_0}\left\|\boldsymbol{f}_{\boldsymbol{\theta}}(Y_{t_k}, t_k) - \boldsymbol{f}_{\boldsymbol{\theta}^-}\left(\hat{Y}_{t_{k+1}}^{\phi}, t_{k+1}\right)\right\|_2^2\right]$$



- Minimize $K$ s.t. $W_2^2\left(f_{\theta}(\mathcal{N}(0, \sigma_T^2 I_d), 0; K), q_0\right) \leq \epsilon_{W_2}^2$

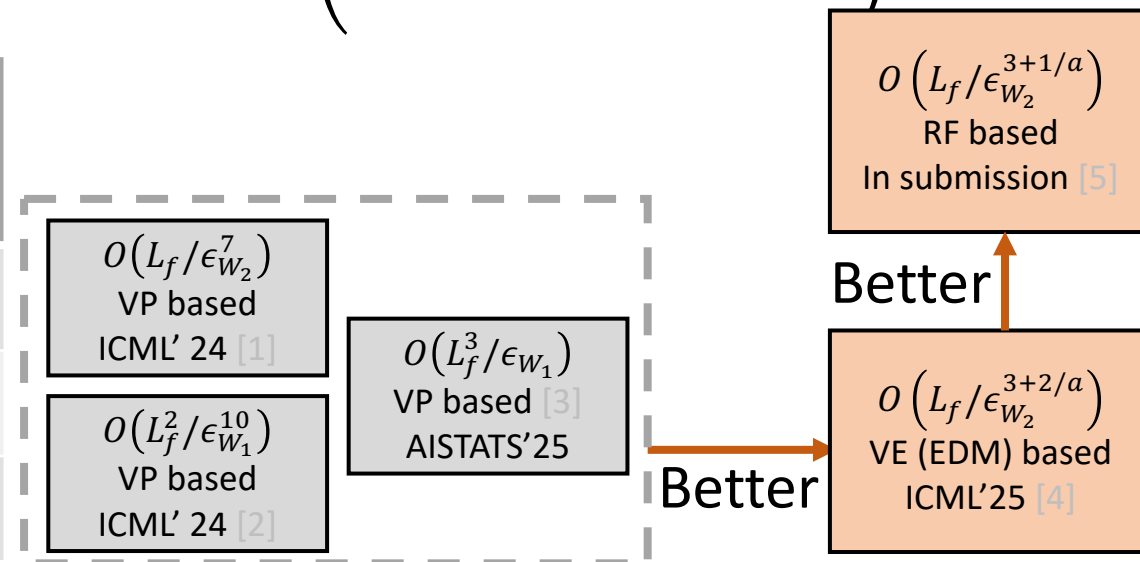[1] SDCS, Consistency Models, ICML 2023.

# Similar Balance

[1] LCF, Sampling is as easy as keeping the consistency: convergence guarantee for consistency models , ICML 2024
[2] DCWY, Theory of consistency diffusion models: Distribution estimation meets fast sampling, ICML 2024
[3] LHW, Towards a mathematical theory for consistency training in diffusion models, AISTATS 2025
[4] YJV**L**, Improved Discretization Complexity Analysis of Consistency Models: Variance Exploding Forward Process and Decay Discretization Scheme, ICML 2025
[5] YZJC**L**, Elucidating Rectified Flow with Deterministic Sampler: Polynomial Discretization Complexity for Multi and One-step Models, Arxiv.

- **Theorem**. For 1-step generation models,

$$W_2^2\left(f_\theta\left(\mathcal{N}(0, \sigma_T^2 I_d), T - \delta\right), q_0\right) \leq \frac{m_T}{\sigma_T^2} + \frac{L_f^2 (T/\delta)^{\frac{2}{a}}}{K^2 \delta^4} + \sigma_\delta^2 \leq \epsilon_{W_2}^2$$

- Then it requires discretization complexity $K = O\left(L_f (T/\delta)^{\frac{1}{a}}/(\delta^2 \epsilon_{W_2})\right)$

| | $T$: $\dfrac{m_T}{\sigma_T^2} \leq \epsilon_{W_2}^2$ | $\delta$: $\sigma_\delta^2 \leq \epsilon_{W_2}^2$ | $K$: $O(L_f(T/\delta)^{\frac{1}{a}}/(\delta^2 \epsilon_{W_2}))$ |
|---|---|---|---|
| VP | $\log(1/\epsilon_{W_2})$ ✔ | $\epsilon_{W_2}^2$ ✗ | $O\left(L_f/\epsilon_{W_2}^{5+2/a}\right)$ |
| VE (EDM) | $1/\epsilon_{W_2}$ ✗ | $\epsilon_{W_2}$ ✔ | $O\left(L_f/\epsilon_{W_2}^{3+2/a}\right)$ |
| RF | $1$ ✔ | $\epsilon_{W_2}$ ✔ | $O\left(L_f/\epsilon_{W_2}^{3+1/a}\right)$ |

# Conclusions

- Pretraining: Efficient Multi-manifold MoG Model
  - Empirical: Much less parameters with good enough performance
  - Theoretical: Estimation error escape the curse of dimensionality
- Fine-tuning: Good Sharing Latent Guarantees Few-shot Efficiency
  - Model the sharing scheme between pretraining and few-shot fine-tuning
  - Show effect of latent quality on estimation and optimization
- Sampling: Complexity for Multi-step Diffusion Models
  - Unified framework for sampling complexities of VP, VE, RF models
- Discretization: Complexity of 1-step Models in Training Phase
  - Support good performances of RF models

# Future Work

- Pretraining Phase
  - SOTA Results with Multi-manifold MoG Modeling and Fewer Parameters
  - Global Optimization Guarantee and Generalization Mechanism
- Few-shot Fine-tuning Phase
  - Multi-task Meta-learning and Few-shot Fine-tuning Framework and Analysis
- Sampling Process of Multi-step Diffusion Models
  - Conditional Generation: Analysis of influence additional guidance
- Learning Process of 1-Step Generative Models
  - With the simplified MoG latent of Multi Subspace MoG modeling, better training and SOTA Results

# Thanks!



## Shuai Li

- Associate Professor
- Shanghai Jiao Tong University
- Research: RL/ML theory
- https://shuaili8.github.io/

Students:
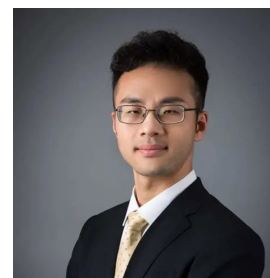


Ruofeng Yang   Zhijie Wang   Yongcan Li   Zhaoyu Zhu

Collaborators:



Bo Jiang   Baoxiang Wang   Cheng Chen   Ruinan Jin

# Questions?