

ICCV 2025 Tutorial

Learning Deep Low-Dimensional Models from High-Dimensional Data: Theory to Practice

Sam Buchanan, Yi Ma, **Qing Qu**, Liyue Shen, Peihao Wang, Zhihui Zhu
October 23, 2025

EECS, University of Michigan

This Tutorial: The Outline

Session I: Introduction of Basic Low-D Models

Session II: Understanding Low-D Structures in Representations

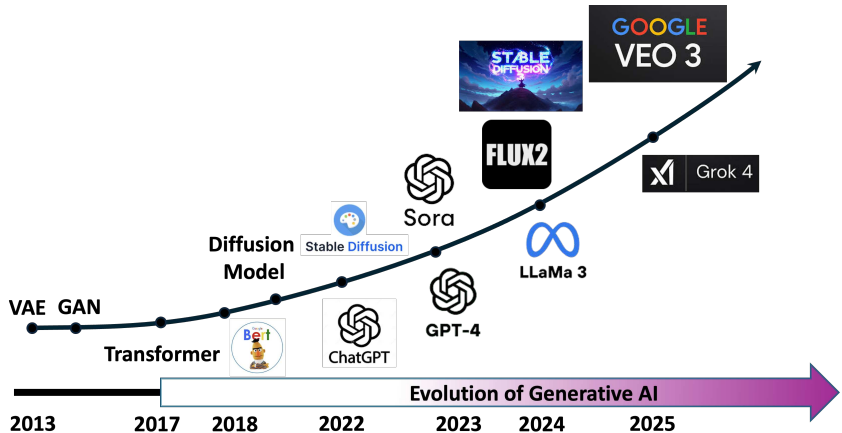
Session III: Understanding Low-D Structures in Generative Models

- **Lecture III-1: Low-Dimensional Models for Understanding Generalization in Diffusion Models**
- Lecture III-2: Exploring Low-Dimensional Structures for Controlling Diffusion Models

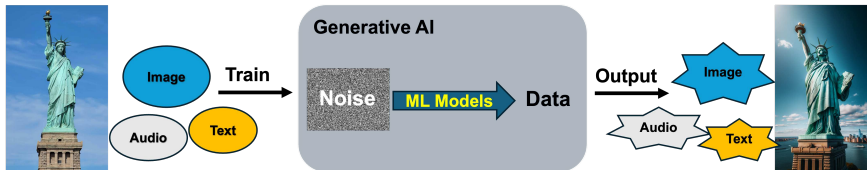
Session IV: Designing Deep Networks for Pursuing Low-D Structures

- Lecture IV-1: ReduNet: A White-box Deep Network from the Principle of Maximizing Rate Reduction
- Lecture IV-2: White-Box Transformers via Sparse Rate Reduction

The Emergence and Revolution of Generative AI

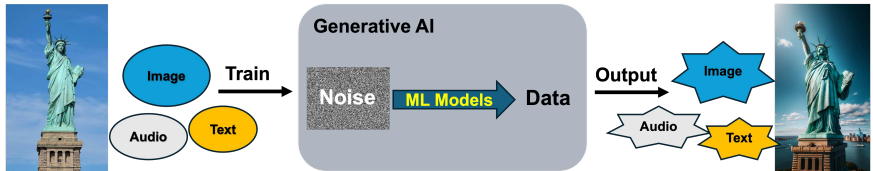


The Family of Generative Models¹

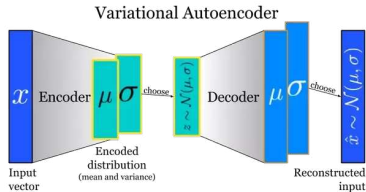


¹Images credited to Prof. Mengdi Wang

The Family of Generative Models¹



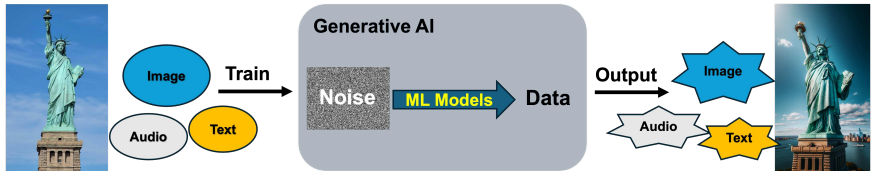
Generative models in the past:



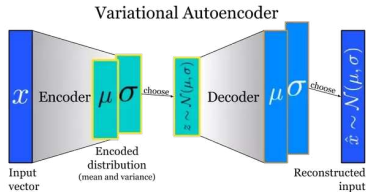
(a) VAE (Kingma & Welling, 2013):
poor generation quality.

¹Images credited to Prof. Mengdi Wang

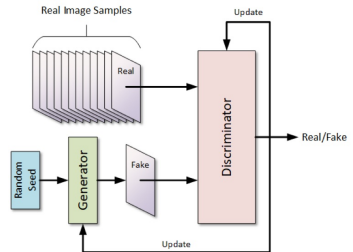
The Family of Generative Models¹



Generative models in the past:



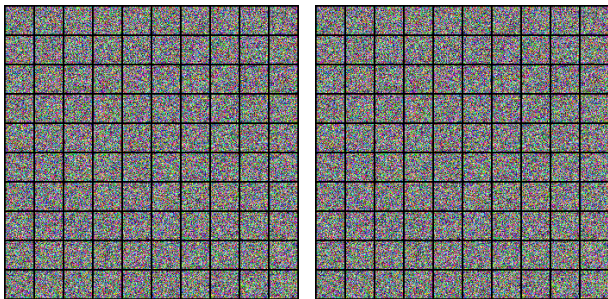
(a) VAE (Kingma & Welling, 2013):
poor generation quality.



(b) GAN (Goodfellow et al. 2014):
unstable to train on large datasets.

¹Images credited to Prof. Mengdi Wang

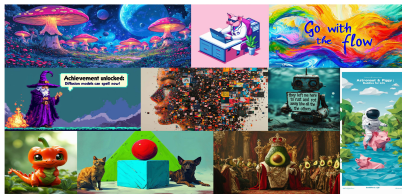
A Revolution by Diffusion Models²



(Sohl-Dickstein et al. 2015, Song and Ermon 2019, Ho et al. 2020)

²<https://yang-song.net/blog/2021/score/>

Many Commercial Applications



Text-to-Image Generation
(Stable Diffusion 3.5, Stability AI)

Autoregression: ✓ High quality ✓ Arbitrary-length ✓ KV caching ✗ Not Parallelizable

Generation steps

There are three categories of the average

There are three categories of the average rate

There are three categories of the average rate of...

Diffusion:

✗ Lower quality ✗ Fixed-length ✗ No KV caching ✓ Parallelizable

the reusability

will continue to

the

Repeal the reusability cuts and the law will continue to reduce the

Repeal the reusability cuts and prove the law will continue to reduce the deficit.

Block Diffusion:

✓ High quality ✓ Arbitrary-length ✓ KV caching ✓ Parallelizable

On September 17, we be

On September 17, 2016, we will be giving the release of

On September 17, 2016, we will be giving the beta-release of the to our server testing ...

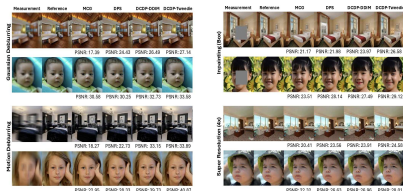
Diffusion LLM

(Block Diffusion, Arriola et al.'25)

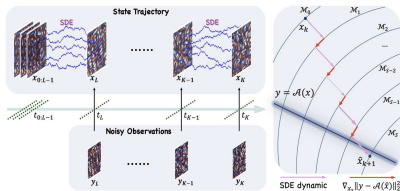


Video Generation
(VEO, Google)

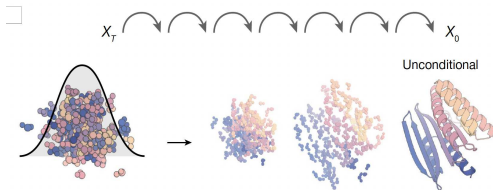
Revolutionizing Scientific Discovery



Inverse Problems
(DCDP, Li et al.'24)

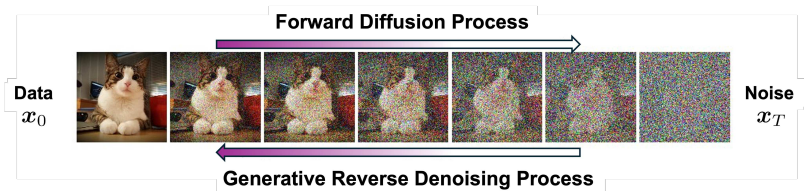


Data Assimilation
(FlowDAS, Chen et al.'25)



Protein Design
(RFDiffusion, Watson et al.'23)

What are Diffusion Models?



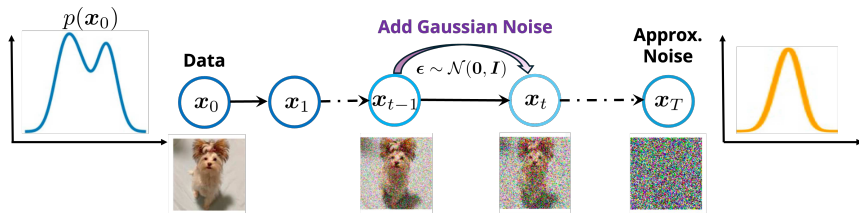
- **Forward process:** progressively adding noise to an image x_0 ;³

$$x_t = \alpha_t x_0 + \beta_t \epsilon, \quad \epsilon \sim \mathcal{N}(\mathbf{0}, I).$$

- **Backward process:** starting from a random noise ϵ , progressively denoising to generate an image x_0 .

³Here, α_t and β_t are some pre-defined noise scales.

Forward Process: Progressively Adding Noise



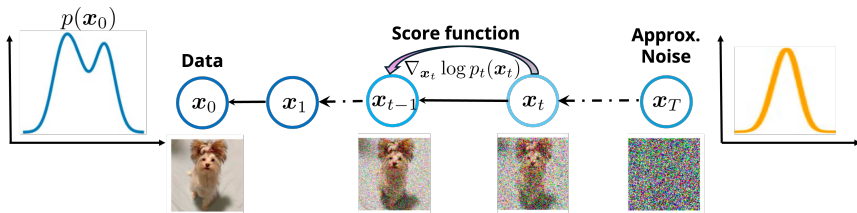
Forward stochastic differential equation (SDE):

$$dx = f(x, t)dt + g(t) \cdot \underset{\text{Brownian}}{dw}.$$

- $f(\cdot, t) : \mathbb{R}^d \rightarrow \mathbb{R}^d$ and $g(\cdot) : \mathbb{R} \rightarrow \mathbb{R}$ are pre-defined *diffusion* and *drift* functions, respectively.⁴

⁴Here, $f(x, t) = \frac{d\log\alpha_t}{dt}x$ and $g(t) = \frac{d\beta_t^2}{dt} - 2\beta_t^2 \frac{d\log\alpha_t}{dt}$.

Generative Backward Process: Progressive Denoising



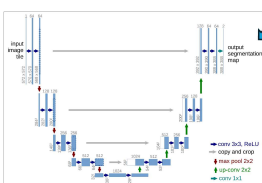
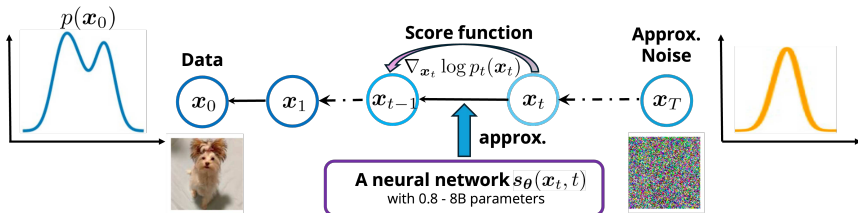
Backward probability flow **ODE** (Song et al., 2020):

$$dx = \left[f(x, t) - \frac{1}{2}g(t)^2 \cdot \underbrace{\left[\nabla_x \log p_t(x) \right]}_{\text{score function}} \right] dt.$$

Deterministic, much faster with slightly inferior sample quality.⁵

⁵For example, EDM (Karras et al., 2022), DPM-solver (Lu et al., 2022).

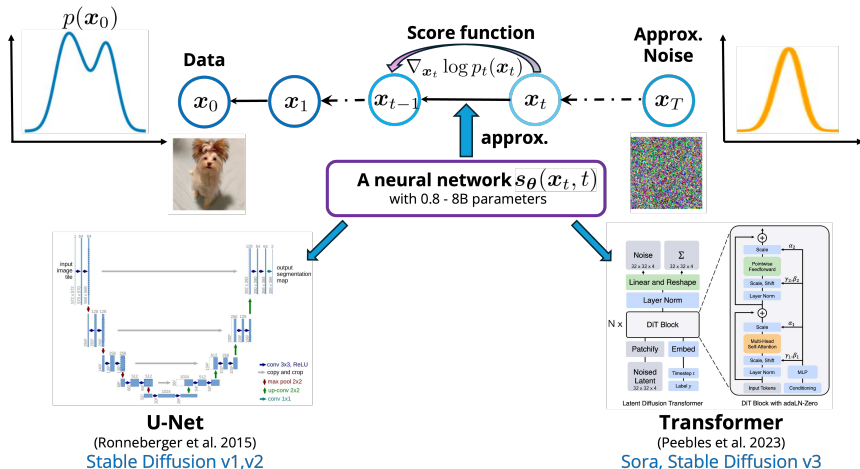
How to Estimate the Score Function?



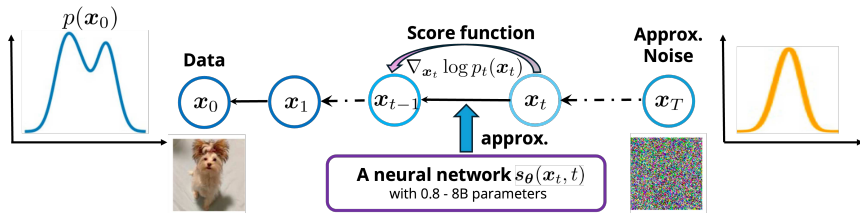
U-Net

(Ronneberger et al. 2015)
Stable Diffusion v1,v2

How to Estimate the Score Function?



How do We Learn the Neural Network?

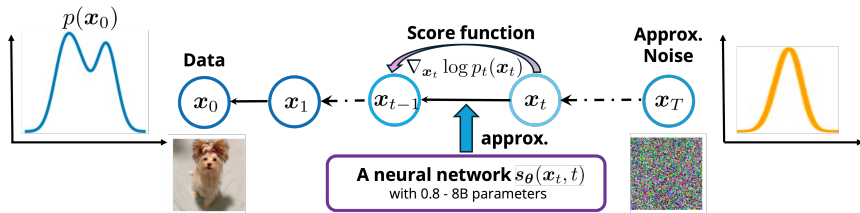


Training loss: we can learn the denoiser $s_\theta(x_t, t)$ simply by solving⁶

$$\min_{\theta} \mathcal{L}(\theta) := \mathbb{E}_{t \sim \mathcal{U}[0,1], x_0 \sim p(x_0)} \left[\beta_t^2 \mathbb{E}_{x_t \sim p(x_t|x_0)} \|\nabla_{x_t} \log p(x_t) - s_\theta(x_t, t)\|_2^2 \right]$$

⁶This can be achieved by sampling $x_0 \sim p(x_0)$, $t \sim \mathcal{U}[0, 1]$, and $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$, to run stochastic gradient descent on $\mathcal{L}(\theta)$ to optimize the network parameters θ .

How do We Learn the Neural Network?

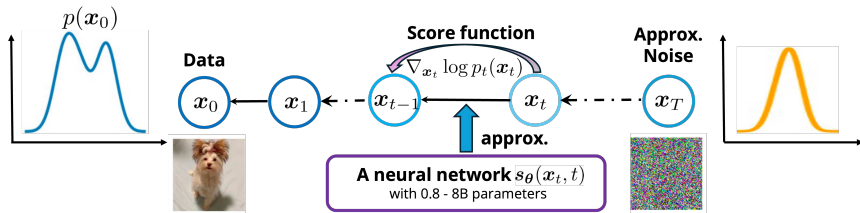


Training loss: we can learn the denoiser $s_\theta(x_t, t)$ simply by solving⁶

$$\begin{aligned} \min_{\theta} \mathcal{L}(\theta) &:= \mathbb{E}_{t \sim \mathcal{U}[0,1], x_0 \sim p(x_0)} \left[\beta_t^2 \|\nabla_{x_t} \log p(x_t) - s_\theta(x_t, t)\|_2^2 \right] \\ &\quad x_t \sim p(x_t | x_0) \\ &= \mathbb{E}_{t \sim \mathcal{U}[0,1], x_0 \sim p(x_0)} \left[\|\epsilon + \beta_t s_\theta(x_t, t)\|_2^2 \right] + \text{const.} \\ &\quad \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \end{aligned}$$

⁶This can be achieved by sampling $x_0 \sim p(x_0)$, $t \sim \mathcal{U}[0, 1]$, and $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, to run stochastic gradient descent on $\mathcal{L}(\theta)$ to optimize the network parameters θ .

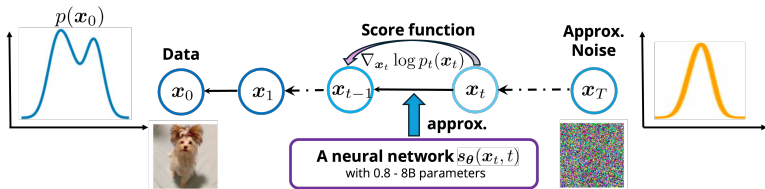
How do We Learn the Neural Network?



Training loss: we can learn the denoiser $s_\theta(x_t, t)$ simply by solving⁷

$$\begin{aligned} \min_{\theta} \mathcal{L}(\theta) &:= \mathbb{E}_{t \sim \mathcal{U}[0,1], x_0 \sim p(x_0)} \left[\beta_t^2 \|\nabla_{x_t} \log p(x_t) - s_\theta(x_t, t)\|_2^2 \right] \\ &\quad x_t \sim p(x_t | x_0) \\ &\approx \sum_{i=1}^N \mathbb{E}_{t \sim \mathcal{U}[0,1], \epsilon \sim \mathcal{N}(\mathbf{0}, I)} \left[\|\epsilon + \beta_t s_\theta(x_t^{(i)}, t)\|_2^2 \right] + \text{const.} \end{aligned}$$

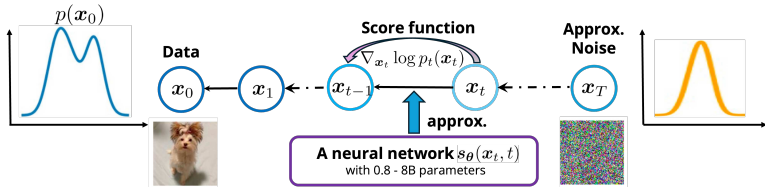
⁷This can be achieved by sampling $x_0 \sim p(x_0)$, $t \sim \mathcal{U}[0, 1]$, and $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$, to run stochastic gradient descent on $\mathcal{L}(\theta)$ to optimize the network parameters θ .



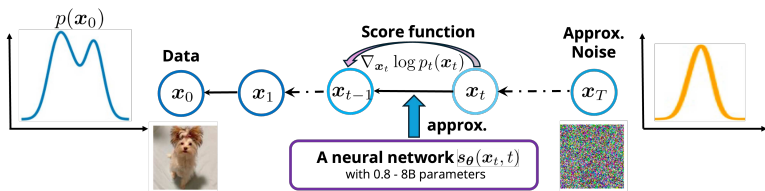
Session III: Exploring Low-D Structures in Generative Models

- **Lecture III-1: Low-Dimensional Models for Understanding Generalization in Diffusion Models**
- Lecture III-2: Exploring Low-Dimensional Structures for Controlling Diffusion Models

Generalization of Diffusion Models



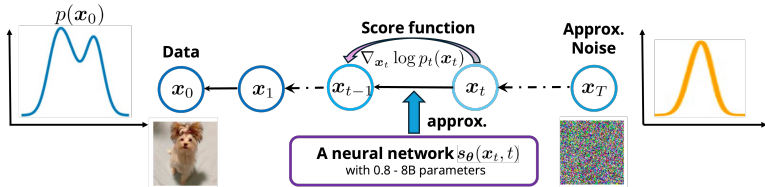
Generalization of Diffusion Models



Generalization: in practice, diffusion models are trained using **empirical loss** of $\mathcal{L}(\theta)$ with **finite** training samples $\mathcal{S} = \{x^{(i)}\}$:

- **Question I:** As diffusion models are trained to fit the training samples, why and when can they **generate new sensible** samples without curse of dimensionality?

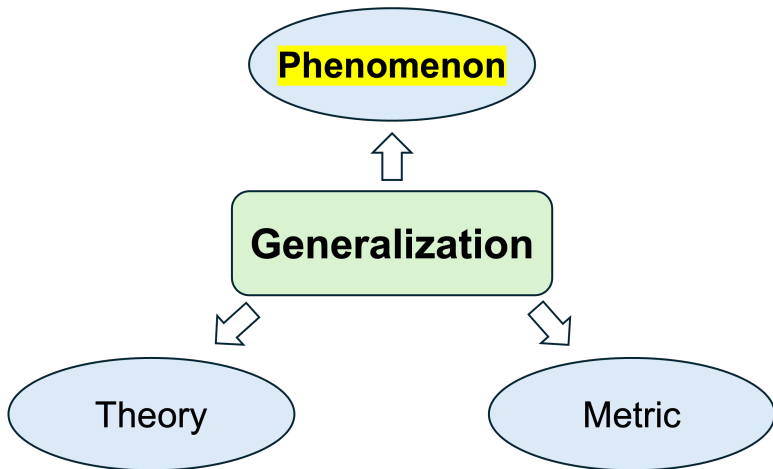
Generalization of Diffusion Models



Generalization: in practice, diffusion models are trained using **empirical loss** of $\mathcal{L}(\theta)$ with **finite** training samples $\mathcal{S} = \{x^{(i)}\}$:

- **Question I:** As diffusion models are trained to fit the training samples, why and when can they **generate new sensible** samples without curse of dimensionality?
- **Question II:** Without **ground-truth** data distribution, how can we measure generalization behavior in terms of training samples and model capacity?

1. **Phenomenon of Generalization** via Model Reproducibility
2. **Theory of Generalization** via Low-dimensional Models
3. **Quantifying Generalization** via Probability Flow Distance
4. **Conclusion & Discussion**



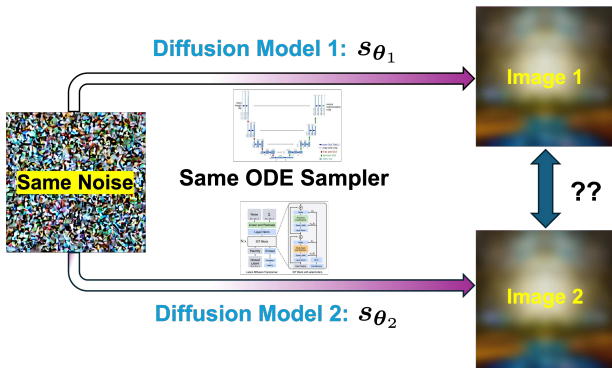
⁸H. Zhang*, J. Zhou*, Y. Lu, M. Guo, P. Wang, L. Shen, and Q. Qu. The Emergence of Reproducibility and Consistency in Diffusion Models. ICML, 2024. (NeurIPS'23 Workshop on Diffusion Models, **Best Paper Award**)

Phenomenon of Generalization via Model Reproducibility

Reproducibility in Diffusion Models

Question for Audience

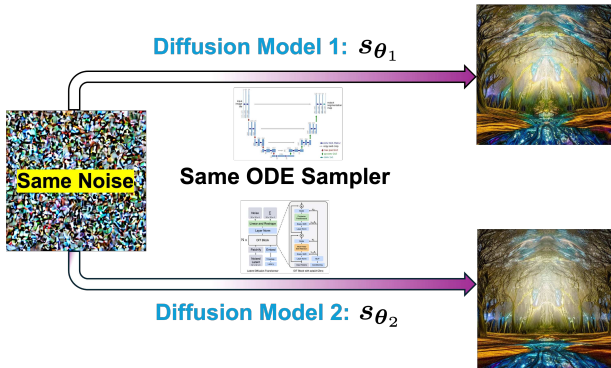
Starting from the **same noise input**, how are the generated data samples from various diffusion models related to each other?



Reproducibility in Diffusion Models

Question for Audience

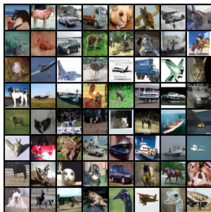
Starting from the **same noise input**, how are the generated data samples from various diffusion models related to each other?



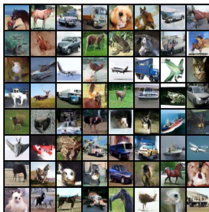
Reproducibility in Diffusion Model

Question for Audience

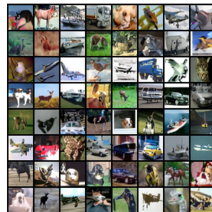
Starting from the **same noise input**, how are the generated data samples from various diffusion models related to each other?



(a) DDPMv4



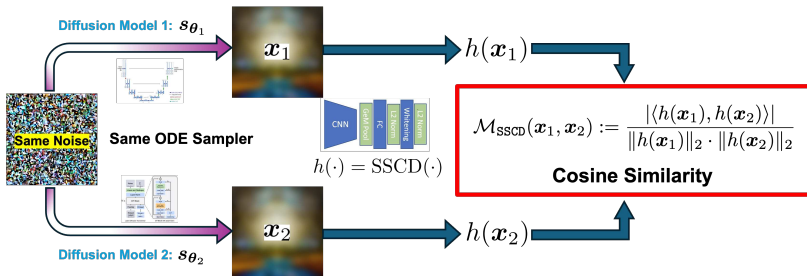
(b) CT



(c) U-ViT

Training on the same dataset, sampling by an ODE deterministic sampler.

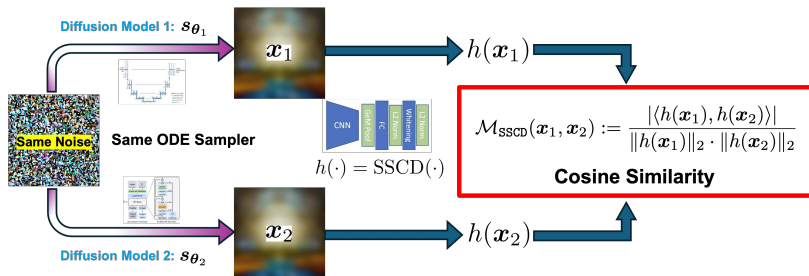
How to Measure Reproducibility Quantitatively?



Self-supervised copy detection (SSCD) similarity $\mathcal{M}_{\text{SSCD}}(\cdot, \cdot)$.

- Here, $h(\cdot) = \text{SSCD}(\cdot)$ represents a neural image descriptor for copy detection. (Pizzi et al.'22, Somepalli et al.'23)

How to Measure Reproducibility Quantitatively?

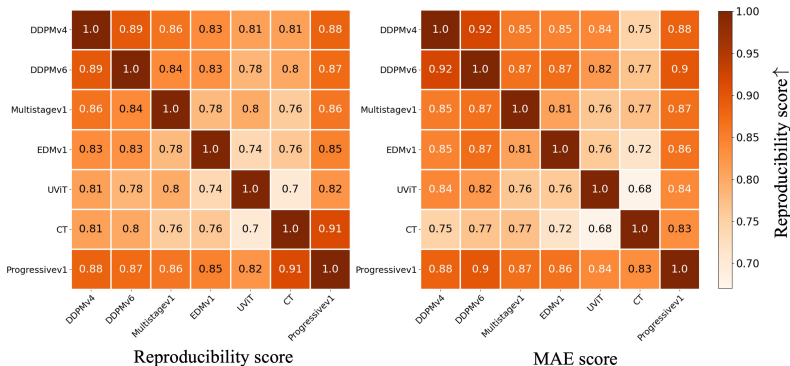


Reproducibility (RP) Score:

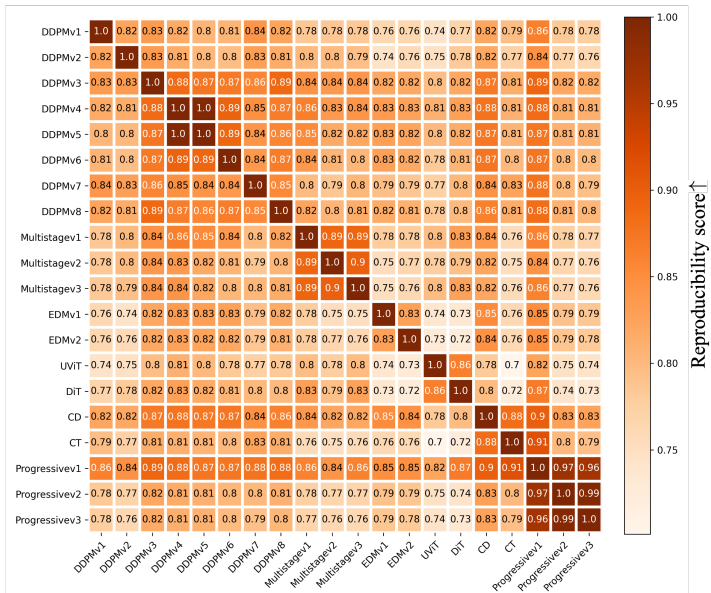
$$\text{RP Score} := \mathbb{P}(\mathcal{M}_{\text{SSCD}}(x_1, x_2) > 0.6).$$

- A *probability measure* of the similarity between models.
- We sample 10^4 random noise pairs to estimate the probability.

Quantitative Analysis of Diffusion Models (Cifar10)



- **Network architectures.** Transformer (U-ViT) vs U-Nets.
- **Training loss.** Consistency loss (CT), EDMv1, and others.
- **Sampling procedures.** DPM (DDPMv4), EDMv1, vs CT.
- **Perturbation kernels.** VP (DDPMv4), sub-VP(DDPMv6), EDMv1.



Reproducibility is Rare in Other Generative Models

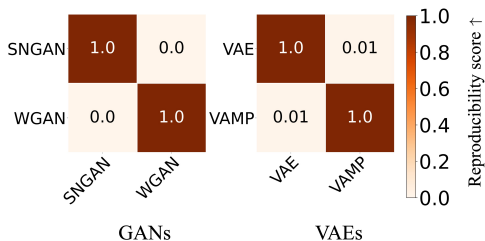


Figure 4: Reproducibility for GANs and VAEs.

- Before this work, only for VAE with a factorized prior distribution over the latent variables (Khemakhem et al. 2020).

Reproducibility is Rare in Other Generative Models

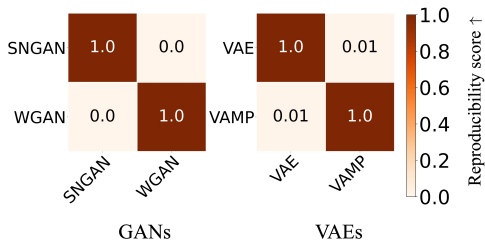
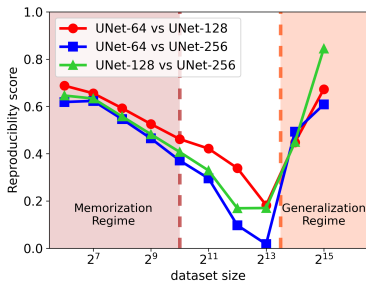


Figure 4: Reproducibility for GANs and VAEs.

- Before this work, only for VAE with a factorized prior distribution over the latent variables (Khemakhem et al. 2020).
- **Prevalent phenomenon in diffusion model!**

Reproducibility Manifests in Two Different Regimes

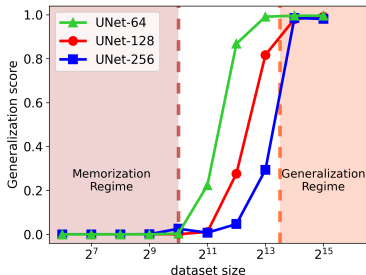


Reproducibility (RP) Score:

$$\text{RP Score} := \mathbb{P}(\mathcal{M}_{\text{SSCD}}(x_1, x_2) > 0.6).$$

Higher implies better reproducibility between two diffusion models.

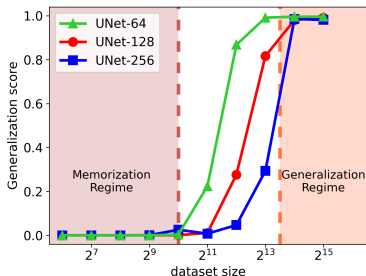
Reproducibility Manifests in Two Different Regimes



Generalization (GL) score is defined to

measure the difference between a **newly generated sample** x
and the **whole training dataset** $\mathcal{S} = \{x^{(i)}\}_{i=1}^N$.

Reproducibility Manifests in Two Different Regimes

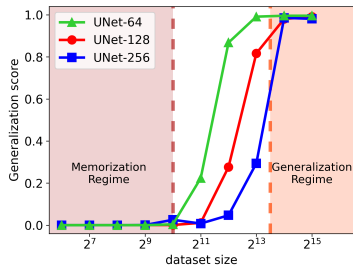
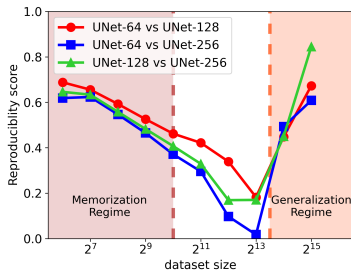


Generalization (GL) score (or perhaps memorization score?)

$$\text{GL Score} := 1 - \mathbb{P} \left(\max_{i \in [N]} [\mathcal{M}_{\text{SSCD}}(\mathbf{x}, \mathbf{x}^{(i)})] > 0.6 \right),$$

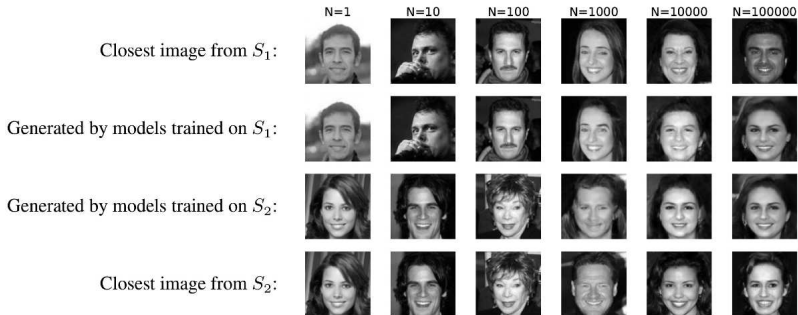
is also a *probability* measure. Higher implies better generalizability.

From “Memorization” to “Generalization”



Reproducibility manifests in **two distinct regimes**, with a **strong** correlation with model’s **generalizability**.

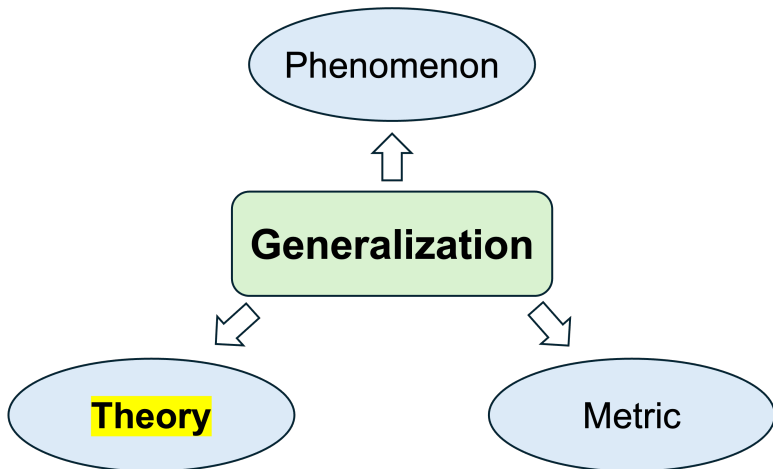
Complementary Results from Concurrent Work⁹



Non-overlapping subsets from the same distribution:

- the same model is trained on disjoint parts S_1, S_2 of the dataset
- the image is generated from the same initial noise

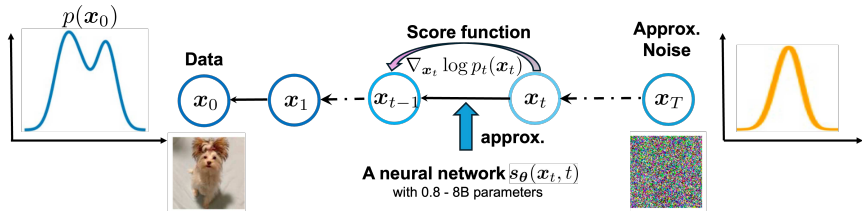
⁹Z Kadkhodaie, et al.'24 "Generalization in diffusion models arises from geometry-adaptive harmonic representation." (ICLR'24, **Outstanding Paper Award**)



¹⁰P. Wang*, H. Zhang*, Z. Zhang, S. Chen, Y. Ma, and Q. Qu. Diffusion Model Learns Low-Dimensional Distributions via Subspace Clustering. Arxiv Preprint arXiv:2409.02426, 2024.

Theory of Generalization via Low-dimensional Models

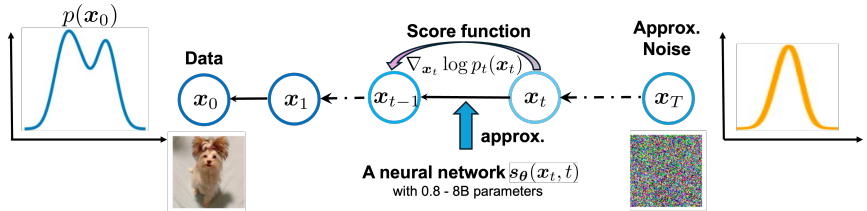
Why Does Reproducibility Manifest in Distinct Regimes?



Backward ODE sampler:

$$dx = \left[f(x, t) - \frac{1}{2} g(t)^2 \cdot \underbrace{\left[\nabla_x \log p_t(x) \right]}_{\text{score function}} \right] dt.$$

Why Does Reproducibility Manifest in Distinct Regimes?

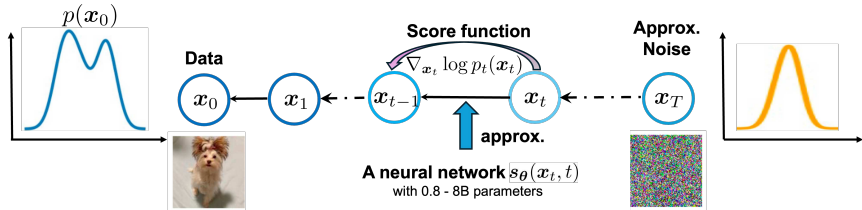


Backward ODE sampler:

$$dx = \left[f(x, t) - \frac{1}{2}g(t)^2 \cdot \underbrace{\nabla_x \log p_t(x)}_{\text{score function}} \right] dt.$$

- **Score approximation.** Reproducibility implies that s_θ can well approximate the score function $\nabla_x \log p_t(x)$.

Why Does Reproducibility Manifest in Distinct Regimes?



Backward ODE sampler:

$$dx = \left[f(x, t) - \frac{1}{2}g(t)^2 \cdot \underbrace{\nabla_x \log p_t(x)}_{\text{score function}} \right] dt.$$

- **Score approximation.** Reproducibility implies that s_θ can well approximate the score function $\nabla_x \log p_t(x)$.
- **Learning different scores in distinct regimes.** But for which $p(x_0)$ are we learning the score function?

Definition

Given a training dataset $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ of N -samples, the **empirical distribution** $p_{\text{emp}}(\mathbf{x})$ of \mathcal{S} can be characterized by the **multi-delta distribution**:

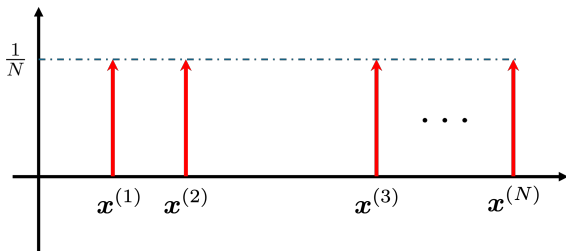
$$p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}^{(i)}).$$

Learning Empirical Distribution in Memorization Regime

Definition

Given a training dataset $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N$ of N -samples, the **empirical distribution** $p_{\text{emp}}(\mathbf{x})$ of \mathcal{S} can be characterized by the **multi-delta distribution**:

$$p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}^{(i)}).$$



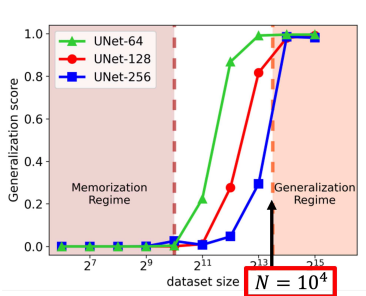
The curse of dimensionality: for image dataset (e.g., CelebA, Cifar),

$$p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}^{(i)}) \approx p_{\text{data}}(\mathbf{x}),$$

to be ε -close,¹¹ we need **extremely large** $N \geq (L/\varepsilon)^d$

¹¹We can draw this conclusion by a simple covering argument, the image dimension $d = 32 \times 32 = 1024$ for Cifar. See also recent work by Li et al., 2024.

Interpolation/Extrapolation of True Data Distribution



The curse of dimensionality: for image dataset (e.g., CelebA, Cifar),

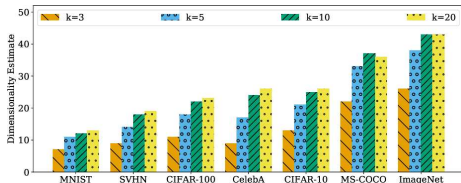
$$p_{\text{emp}}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x^{(i)}) \approx p_{\text{data}}(x),$$

to be ε -close,¹¹ we need **extremely large** $N \geq (L/\varepsilon)^d$!

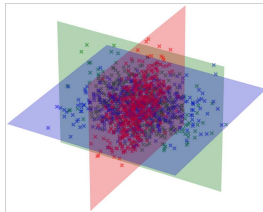
¹¹We can draw this conclusion by a simple covering argument, the image dimension $d = 32 \times 32 = 1024$ for Cifar. See also recent work by Li et al., 2024.

The Intrinsic Low-Dimensionality of Data¹²

The low-dim of model reflects the intrinsic dimension of our data:



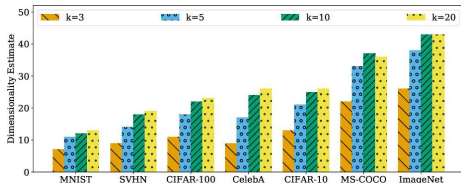
Intrinsic dimension of image datasets.



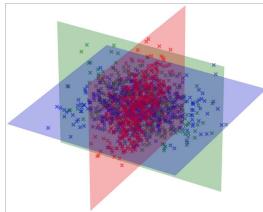
¹²Image credit: P. Pope et al., ICLR'2021.

The Intrinsic Low-Dimensionality of Data¹²

The low-dim of model reflects the intrinsic dimension of our data:



Intrinsic dimension of image datasets.



The blessing of dimensionality: the intrinsic data dimension r is **much lower** than the ambient dimension d , i.e., $r \ll d$.

¹²Image credit: P. Pope et al., ICLR'2021.

- Denoising autoencoder (DAE) formulation:

$$\min_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}) := \sum_{i=1}^N \int_0^1 \lambda_t \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, I_n)} \left[\left\| \mathbf{x}_{\boldsymbol{\theta}}(\mathbf{x}_t, t) - \mathbf{x}^{(i)} \right\|^2 \right] dt$$

Intrinsic Low-Dimensionality of the Model

- Denoising autoencoder (DAE) formulation:

$$\min_{\theta} \ell(\theta) := \sum_{i=1}^N \int_0^1 \lambda_t \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I_n)} \left[\left\| \mathbf{x}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}^{(i)} \right\|^2 \right] dt$$

- Tweedie's formula, $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \beta_t \epsilon$:

$$\underbrace{\mathbf{x}_{\theta}(\mathbf{x}_t, t)}_{\text{neural networks, e.g., U-Net}} \approx \underbrace{\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]}_{\text{posterior mean}} = (\mathbf{x}_t + \beta_t^2 \underbrace{\nabla_{\mathbf{x}} \log p_t(\mathbf{x})}_{\text{score function}}) / \alpha_t.$$

Intrinsic Low-Dimensionality of the Model

- Denoising autoencoder (DAE) formulation:

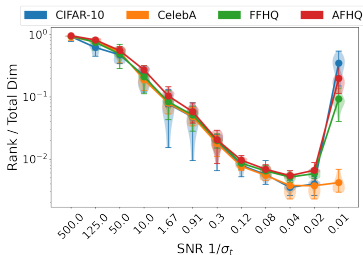
$$\min_{\theta} \ell(\theta) := \sum_{i=1}^N \int_0^1 \lambda_t \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I_n)} \left[\left\| \mathbf{x}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}^{(i)} \right\|^2 \right] dt$$

- Tweedie's formula, $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \beta_t \epsilon$:

$$\underbrace{\mathbf{x}_{\theta}(\mathbf{x}_t, t)}_{\text{neural networks, e.g., U-Net}} \approx \underbrace{\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]}_{\text{posterior mean}} = (\mathbf{x}_t + \beta_t^2 \underbrace{[\nabla_{\mathbf{x}} \log p_t(\mathbf{x})]}_{\text{score function}}) / \alpha_t.$$

- Evaluate the **rank ratio** of the Jacobian $\mathbf{J}_{\theta, t}(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \mathbf{x}_{\theta}(\mathbf{x}_t)$.

Intrinsic Low-Dimensionality of the Model



- Denoising autoencoder (DAE) formulation:

$$\min_{\theta} \ell(\theta) := \sum_{i=1}^N \int_0^1 \lambda_t \mathbb{E}_{\epsilon \sim \mathcal{N}(\mathbf{0}, I_n)} \left[\left\| \mathbf{x}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}^{(i)} \right\|^2 \right] dt$$

- Tweedie's formula, $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \beta_t \epsilon$:

$$\underbrace{\mathbf{x}_{\theta}(\mathbf{x}_t, t)}_{\text{neural networks, e.g., U-Net}} \approx \underbrace{\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t]}_{\text{posterior mean}} = (\mathbf{x}_t + \beta_t^2 \underbrace{[\nabla_{\mathbf{x}} \log p_t(\mathbf{x})]}_{\text{score function}}) / \alpha_t.$$

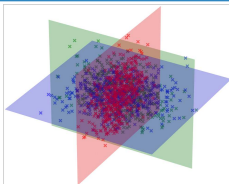
- Evaluate the **rank ratio** of the Jacobian $\mathbf{J}_{\theta, t}(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \mathbf{x}_{\theta}(\mathbf{x}_t)$.

Data Assumption

Mixture of low-rank Gaussian (MoLRG):

$$p_{\text{data}}(\mathbf{x}) = \frac{1}{K} \sum_{i \in [K]} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma_i) \text{ with } \Sigma_i = \mathbf{U}_i \mathbf{U}_i^\top,$$

where K is the number of clusters, and $\mathbf{U}_i \in \mathbb{R}^{d \times r}$ is the low-rank basis for the i th cluster with $r \ll d$, with $\mathbf{U}_i \perp \mathbf{U}_j (i \neq j)$.



¹³Chen et al. Score Approximation, Estimation and Distribution Recovery of Diffusion Models on Low-Dimensional Data. *ICML*, 2023.

Study Generalization under Low-Dimensional Models

Data Assumption

Mixture of low-rank Gaussian (MoLRG):

$$p_{\text{data}}(\mathbf{x}) = \frac{1}{K} \sum_{i \in [K]} \mathcal{N}(\mathbf{x}; \mathbf{0}, \Sigma_i) \text{ with } \Sigma_i = \mathbf{U}_i \mathbf{U}_i^\top,$$

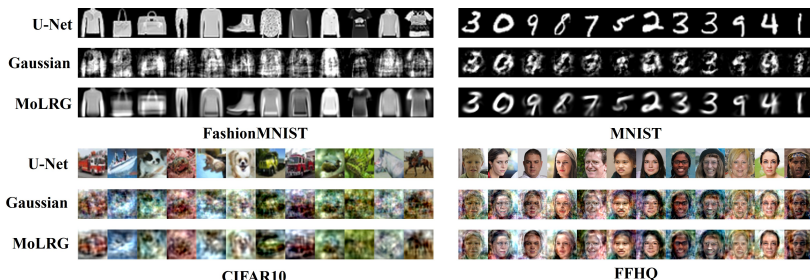
where K is the number of clusters, and $\mathbf{U}_i \in \mathbb{R}^{d \times r}$ is the low-rank basis for the i th cluster with $r \ll d$, with $\mathbf{U}_i \perp \mathbf{U}_j (i \neq j)$.

Lemma 1. Suppose that $p_{\text{data}}(\mathbf{x}_0) \sim \text{MoLRG}$. For all $t > 0$,

$$\mathbb{E}[\mathbf{x}_0 | \mathbf{x}_t] = \frac{\alpha_t}{\alpha_t^2 + \beta_t^2} \sum_{k=1}^K w_k \mathbf{U}_k^* \mathbf{U}_k^{*\top} \mathbf{x}_t,$$

where $w_k = \frac{\pi_k \exp(\phi_t \|\mathbf{U}_k^{*\top} \mathbf{x}_t\|^2)}{\sum_{k=1}^K \pi_k \exp(\phi_t \|\mathbf{U}_k^{*\top} \mathbf{x}_t\|^2)}$ and $\phi_t := \alpha_t^2 / (2\beta_t^2(\alpha_t^2 + \beta_t^2))$.

MoLRG Serves as a Good Approximation of Image Distributions



DAE-generated images with different parameterizations:

- **1st row:** U-Net parameterization;
- **2nd row:** optimal denoiser of a single Gaussian;
- **3rd row:** optimal denoiser of MoLRG (see Lemma 1).

A Simple Case Study: Single Low-rank Gaussian $K = 1$

Theorem (Equivalence to PCA)

Suppose that

- The distribution $p(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0; \mathbf{0}, \mathbf{U}_g \mathbf{U}_g^\top)$ with $\mathbf{U}_g \in \mathcal{O}^{d \times r}$;
- For each $t \in [0, 1]$, we parameterize the denoiser $\mathbf{x}_U(\mathbf{x}_t, t)$:

$$\mathbf{x}_U(\mathbf{x}_t, t) = \frac{\alpha_t}{\alpha_t^2 + \beta_t^2} \cdot \mathbf{U} \mathbf{U}^\top \mathbf{x}_t$$

Let $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(i)} & \dots & \mathbf{x}^{(N)} \end{bmatrix}$ be the training data matrix, then

A Simple Case Study: Single Low-rank Gaussian $K = 1$

Theorem (Equivalence to PCA)

Suppose that

- The distribution $p(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0; \mathbf{0}, U_g U_g^\top)$ with $U_g \in \mathcal{O}^{d \times r}$;
- For each $t \in [0, 1]$, we parameterize the denoiser $\mathbf{x}_U(\mathbf{x}_t, t)$:

$$\mathbf{x}_U(\mathbf{x}_t, t) = \frac{\alpha_t}{\alpha_t^2 + \beta_t^2} \cdot U U^\top \mathbf{x}_t$$

Let $\mathbf{X} = \begin{bmatrix} \mathbf{x}^{(i)} & \dots & \mathbf{x}^{(N)} \end{bmatrix}$ be the training data matrix, then

- The training loss can be reduced to the **PCA problem**:

$$\max_U \|\mathbf{U}^\top \mathbf{X}\|_F^2, \quad \text{s.t.} \quad \mathbf{U}^\top \mathbf{U} = \mathbf{I}_r.$$

A Simple Case Study: Single Low-rank Gaussian $K = 1$

Theorem (Equivalence to PCA)

Suppose that

- The distribution $p(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0; \mathbf{0}, U_g U_g^\top)$ with $U_g \in \mathcal{O}^{d \times r}$;
- For each $t \in [0, 1]$, we parameterize the denoiser $\mathbf{x}_U(\mathbf{x}_t, t)$:

$$\mathbf{x}_U(\mathbf{x}_t, t) = \frac{\alpha_t}{\alpha_t^2 + \beta_t^2} \cdot U U^\top \mathbf{x}_t$$

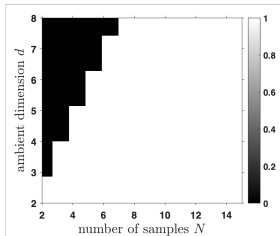
Let $\mathbf{X} = [\mathbf{x}^{(i)} \quad \dots \quad \mathbf{x}^{(N)}]$ be the training data matrix, then

- The training loss can be reduced to the **PCA problem**:

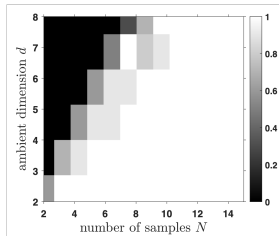
$$\max_U \|\mathbf{U}^\top \mathbf{X}\|_F^2, \quad \text{s.t.} \quad \mathbf{U}^\top \mathbf{U} = \mathbf{I}_r.$$

- Thus, it holds for the global solution \mathbf{U}_\star w.h.p. that
 - (i) If $N \geq r$, we have $\|\mathbf{U}_\star \mathbf{U}_\star^\top - U_g U_g^\top\|_F < \delta$;
 - (ii) If $N < r$, we have $\|\mathbf{U}_\star \mathbf{U}_\star^\top - U_g U_g^\top\|_F \geq \sqrt{r - N} - \delta$.

Phase Transitions on MoLRG with Parameterized Networks

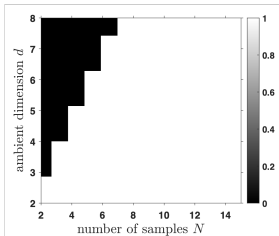


(a) PCA

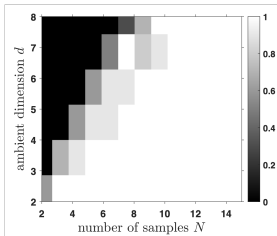


(b) DM ($K=1$)

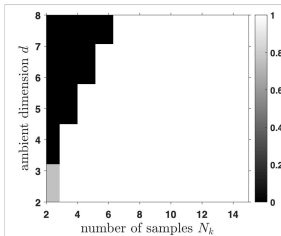
Phase Transitions on MoLRG with Parameterized Networks



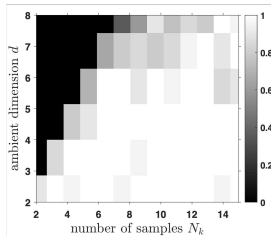
(a) PCA



(b) DM ($K=1$)



(c) Subspace Clustering



(d) DM ($K=2$)

Study of Multiple Low-Dim Subspaces $K > 1$

Theorem (Equivalence to Subspace Clustering)

Suppose that

- $p(x_0)$ is MoLRG with $K > 1$;
- If we parameterize the DAE network

$$\mathbf{x}_\theta(\mathbf{x}_t, t) = \frac{\alpha_t}{\alpha_t^2 + \beta_t^2} \sum_{k=1}^K w_k(\theta; \mathbf{x}_t) \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t.$$

Theorem (Equivalence to Subspace Clustering)

Suppose that

- $p(x_0)$ is MoLRG with $K > 1$;
- If we parameterize the DAE network

$$\mathbf{x}_\theta(\mathbf{x}_t, t) = \frac{\alpha_t}{\alpha_t^2 + \beta_t^2} \sum_{k=1}^K w_k(\theta; \mathbf{x}_t) \mathbf{U}_k \mathbf{U}_k^\top \mathbf{x}_t.$$

Then the training problem is equivalent to **subspace clustering**

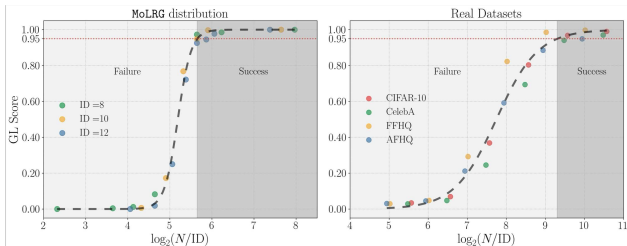
$$\max_{\theta} \frac{1}{N} \sum_{k=1}^K \sum_{i \in \mathcal{C}_k(\theta)} \|\mathbf{U}_k^\top \mathbf{x}^{(i)}\|^2 \quad \text{s.t.} \quad [\mathbf{U}_1, \dots, \mathbf{U}_K] \in \mathcal{O}^{n \times dK},$$

with $\mathcal{C}_k(\theta) := \{i \in [N] : \|\mathbf{U}_k^\top \mathbf{x}^{(i)}\| \geq \|\mathbf{U}_l^\top \mathbf{x}^{(i)}\|, \forall l \neq k\}$ for $k \in [K]$. The sample complexity is **linear** between N and d .

Phase Transition of Generalization: $N \geq c \cdot \text{ID}$

Training U-Net model with fixed capacity on synthetic & real dataset:

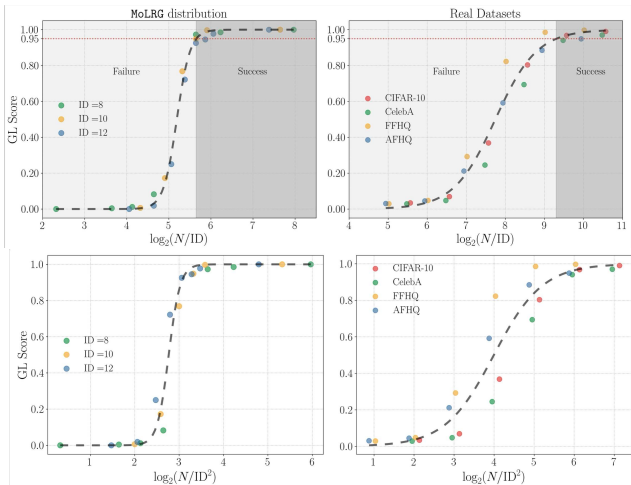
Fitting N vs. ID



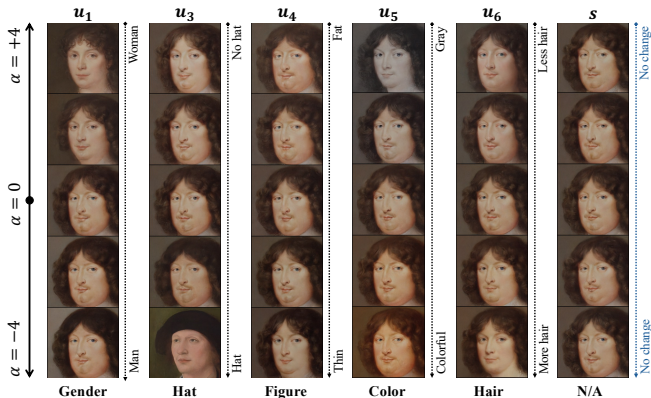
Phase Transition of Generalization: $N \geq c \cdot \text{ID}$

Training U-Net model with fixed capacity on synthetic & real dataset:

Fitting N vs. ID

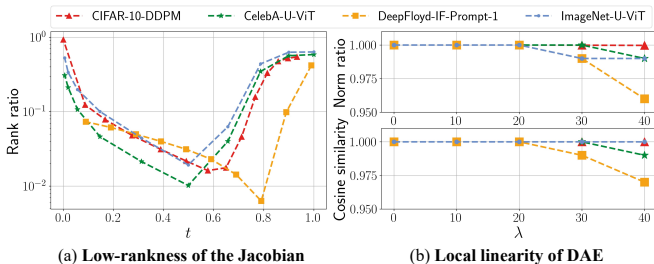


Low-Dimensional Bases are Concept Steering Vectors



$$\begin{aligned}
 x_{\theta,t}(x_t + \lambda u_i) &\approx x_{\theta,t}(x_t) + \lambda J_{\theta,t}(x_t) u_i \\
 &= x_{\theta,t}(x_t) + \lambda \sum_{j=1}^r \sigma_j u_j u_j^\top u_i \\
 &= \hat{x}_{0,t} + \lambda \sigma_i u_i.
 \end{aligned}$$

Inductive Bias Towards “Simple” Solutions¹⁴



The trained network via Adam tends to have simple structures:

- **Low-rankness** of the Jacobian:

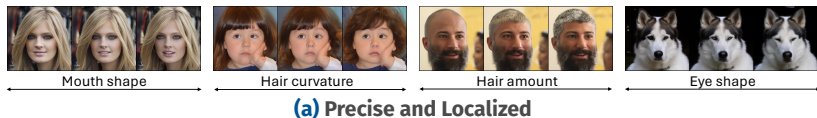
$$\mathbf{J}_{\theta,t}(\mathbf{x}_t) = \mathbf{U}\Sigma\mathbf{U}^\top = \sum_{i=1}^r \sigma_i \mathbf{u}_i \mathbf{u}_i^\top.$$

- **Local linearity** of the DAE:

$$\mathbf{x}_{\theta,t}(\mathbf{x}_t + \lambda\Delta\mathbf{x}) \approx \mathbf{x}_{\theta,t}(\mathbf{x}_t) + \lambda\mathbf{J}_{\theta,t}(\mathbf{x}_t) \cdot \Delta\mathbf{x}$$

¹⁴X. Li, Y. Dai, Q. Qu. Understanding Generalizability of Diffusion Models Requires Rethinking the Hidden Gaussian Structure. *NeurIPS*, 2024.

Low-rank Controllable Image Editing (LOCO Edit)¹⁵

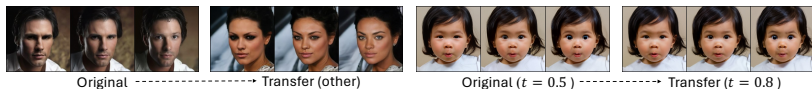


¹⁵S. Chen*, H. Zhang*, M. Guo, Y. Lu, P. Wang, and Q. Qu. Exploring Low-Dimensional Subspaces in Diffusion Models for Controllable Image Editing. NeurIPS, 2024.

Low-rank Controllable Image Editing (LOCO Edit)¹⁵



(a) Precise and Localized



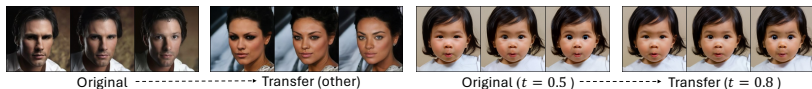
(b) Homogeneity & Transferability

¹⁵S. Chen*, H. Zhang*, M. Guo, Y. Lu, P. Wang, and Q. Qu. Exploring Low-Dimensional Subspaces in Diffusion Models for Controllable Image Editing. NeurIPS, 2024.

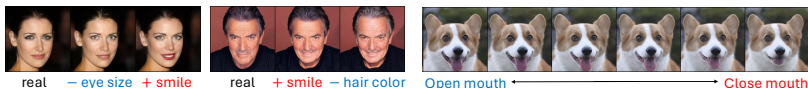
Low-rank Controllable Image Editing (LOCO Edit)¹⁵



(a) Precise and Localized



(b) Homogeneity & Transferability

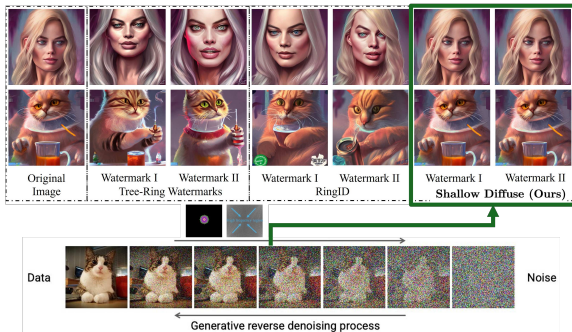


(c) Composability & Disentanglement

(d) Linearity

¹⁵S. Chen*, H. Zhang*, M. Guo, Y. Lu, P. Wang, and Q. Qu. Exploring Low-Dimensional Subspaces in Diffusion Models for Controllable Image Editing. NeurIPS, 2024.

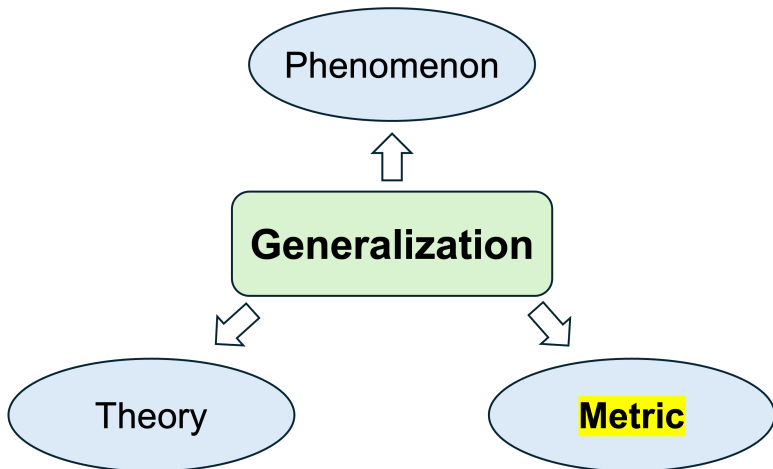
Robust and Invisible Watermarking via Shallow Diffuse¹⁶



Key idea: Inject the watermark Δx in the **Null Space** of $J_{\theta,t}(x_t)$:

$$x_{\theta,t}(x_t^{\mathcal{W}}) = x_{\theta,t}(x_t) + \boxed{\lambda \underset{\approx 0}{J_{\theta,t}(x_t)} \cdot \Delta x} \approx x_{\theta,t}(x_t)$$

¹⁶W. Li, H. Zhang, Q. Qu. Shallow Diffuse: Robust and Invisible Watermarking through Low-Dimensional Subspaces in Diffusion Models. NeurIPS'25, (**spotlight, top 3%**)



¹⁷H. Zhang, Z. Huang, S. Chen, J. Zhou, Z. Zhang, P. Wang, Q. Qu. Understanding Generalization in Diffusion Models via Probability Flow Distance. *Arxiv Preprint arXiv:2505.20123*, 2025.

Quantifying Generalization via Probability Flow Distance

How to Quantify Generalization for GenAI Models?

In practice, measuring generalization errors is hard **without the groundtruth** distribution $p_{\text{data}}(\boldsymbol{x})$. Previously, we defined GL Score:

$$\text{GL Score} := 1 - \mathbb{P} \left(\max_{i \in [N]} \left[\mathcal{M}_{\text{SSCD}}(\boldsymbol{x}, \boldsymbol{x}^{(i)}) \right] > 0.6 \right).$$

- It measures the **dissimilarity** between generated \boldsymbol{x} and the whole training dataset $\mathcal{S} = \{\boldsymbol{x}^{(i)}\}_{i=1}^N$.

How to Quantify Generalization for GenAI Models?

In practice, measuring generalization errors is hard **without the groundtruth** distribution $p_{\text{data}}(x)$. Previously, we defined GL Score:

$$\text{GL Score} := 1 - \mathbb{P} \left(\max_{i \in [N]} [\mathcal{M}_{\text{SSCD}}(x, x^{(i)})] > 0.6 \right).$$

- It measures the **dissimilarity** between generated x and the whole training dataset $\mathcal{S} = \{x^{(i)}\}_{i=1}^N$.
- However, it **alone cannot** correctly quantify generalization: A **randomly drawn** sample x can be misspecified!

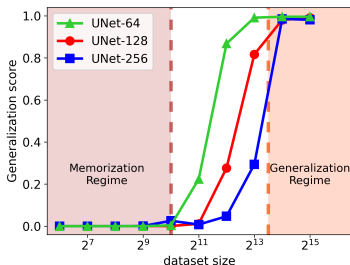
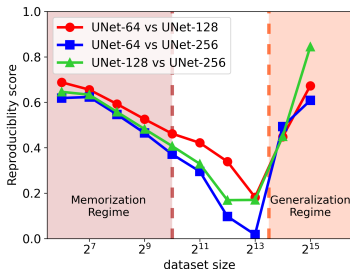
How to Quantify Generalization for GenAI Models?

In practice, measuring generalization errors is hard **without the groundtruth** distribution $p_{\text{data}}(x)$. Previously, we defined GL Score:

$$\text{GL Score} := 1 - \mathbb{P} \left(\max_{i \in [N]} [\mathcal{M}_{\text{SSCD}}(x, x^{(i)})] > 0.6 \right).$$

- It measures the **dissimilarity** between generated x and the whole training dataset $\mathcal{S} = \{x^{(i)}\}_{i=1}^N$.
- However, it **alone cannot** correctly quantify generalization: A **randomly drawn** sample x can be misspecified!
- Instead, we can predict generalization by **combining** GL Score and RP Score.

Reproducibility vs. Generalization



Reproducibility (RP) Score:

$$\text{RP Score} := \mathbb{P}(\mathcal{M}_{\text{SSCD}}(x_1, x_2) > 0.6).$$

Higher implies better reproducibility between two diffusion models.

Quantify Generalization Error $s_\theta \approx \nabla \log p_{\text{data}}$?

Without p_{data} , we study generalization under **teacher-student** setup:

Quantify Generalization Error $s_\theta \approx \nabla \log p_{\text{data}}$?

Without p_{data} , we study generalization under **teacher-student** setup:

- **Teacher model** s_{θ_*} : a large pre-trained model to generate training dataset $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N$; serve as a proxy of p_{data} .

Quantify Generalization Error $s_\theta \approx \nabla \log p_{\text{data}}$?

Without p_{data} , we study generalization under **teacher-student** setup:

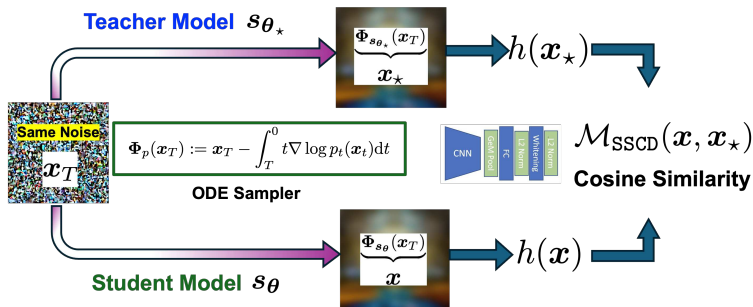
- **Teacher model** s_{θ_*} : a large pre-trained model to generate training dataset $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N$; serve as a proxy of p_{data} .
- **Student model** s_θ : a smaller model trained on \mathcal{S} .

Quantify Generalization Error $s_\theta \approx \nabla \log p_{\text{data}}$?

Without p_{data} , we study generalization under **teacher-student** setup:

- **Teacher model** s_{θ_*} : a large pre-trained model to generate training dataset $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N$; serve as a proxy of p_{data} .
- **Student model** s_θ : a smaller model trained on \mathcal{S} .
- **Evaluation**: We collect noise $\mathcal{T} = \{\mathbf{x}_T^{(i)}\}_{i=1}^M$ distinct from training set \mathcal{S} to compare reproducibility between s_θ and s_{θ_*} !

Quantify Generalization Error $s_\theta \approx \nabla \log p_{\text{data}}$?



Without p_{data} , we study generalization under **teacher-student** setup:

- **Teacher model** s_{θ_*} : a large pre-trained model to generate training dataset $\mathcal{S} = \{x^{(i)}\}_{i=1}^N$; serve as a proxy of p_{data} .
- **Student model** s_θ : a smaller model trained on \mathcal{S} .
- **Evaluation**: We collect noise $\mathcal{T} = \{x_T^{(i)}\}_{i=1}^M$ distinct from training set \mathcal{S} to compare reproducibility between s_θ and s_{θ_*} !

Quantifying Generalization via Probability Flow Distance (PFD)

Definition (Probability Flow Distance)

Let p_θ denote the distribution induced by s_θ , and similarly p_{θ_\star} for s_{θ_\star} . We measure the distance between p_θ and p_{θ_\star} by

$$\text{PFD}(p_\theta, p_{\theta_\star}) := \left(\mathbb{E}_{\mathbf{x}_T} \left[\left\| h \circ \Phi_{p_\theta}(\mathbf{x}_T) - h \circ \Phi_{p_{\theta_\star}}(\mathbf{x}_T) \right\|_2^2 \right] \right)^{\frac{1}{2}},$$

where $\mathbf{x}_T \sim \mathcal{N}(0, T^2 \mathbf{I})$, and $\Phi_p(\mathbf{x}_T) = \mathbf{x}_T - \int_T^0 t \nabla \log p_t(\mathbf{x}_t) dt$.

Quantifying Generalization via Probability Flow Distance (PFD)

Definition (Probability Flow Distance)

Let p_θ denote the distribution induced by s_θ , and similarly p_{θ_\star} for s_{θ_\star} . We measure the distance between p_θ and p_{θ_\star} by

$$\text{PFD}(p_\theta, p_{\theta_\star}) := \left(\mathbb{E}_{\mathbf{x}_T} \left[\left\| h \circ \Phi_{p_\theta}(\mathbf{x}_T) - h \circ \Phi_{p_{\theta_\star}}(\mathbf{x}_T) \right\|_2^2 \right] \right)^{\frac{1}{2}},$$

where $\mathbf{x}_T \sim \mathcal{N}(0, T^2 \mathbf{I})$, and $\Phi_p(\mathbf{x}_T) = \mathbf{x}_T - \int_T^0 t \nabla \log p_t(\mathbf{x}_t) dt$.

- In practice, the expectation can be well approximated by empirical mean over finite samples $\{\mathbf{x}_T^{(i)}\}_{i=1}^M \stackrel{i.i.d.}{\sim} \mathcal{N}(0, T^2 \mathbf{I}_n)$.

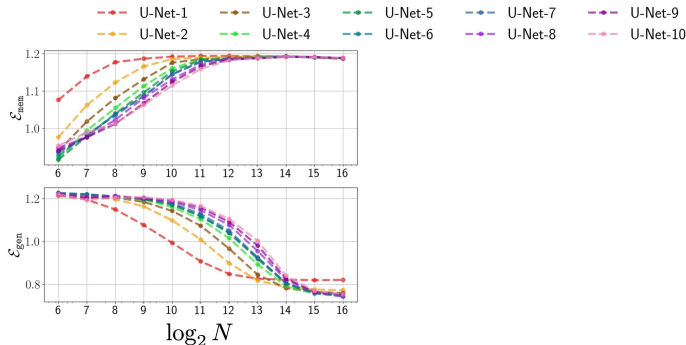
Finding I: Empirical Scaling Behavior with $\log_2(N/\sqrt{|\theta|})$

Generalization and Memorization Errors

Suppose s_θ is trained on $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N \stackrel{i.i.d.}{\sim} p_{\theta_*}$, and let $p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}^{(i)})$. We define

$$\mathcal{E}_{\text{gen}}(\theta) := \text{PFD}(p_\theta, p_{\theta_*}), \quad \mathcal{E}_{\text{mem}}(\theta) := \text{PFD}(p_\theta, p_{\text{emp}}).$$

Finding I: Empirical Scaling Behavior with $\log_2(N/\sqrt{|\theta|})$

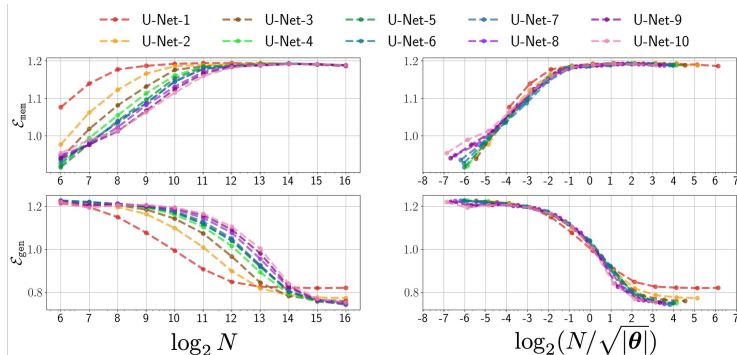


Generalization and Memorization Errors

Suppose s_θ is trained on $\mathcal{S} = \{\mathbf{x}^{(i)}\}_{i=1}^N \stackrel{i.i.d.}{\sim} p_{\theta_*}$, and let $p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{x}^{(i)})$. We define

$$\mathcal{E}_{\text{gen}}(\theta) := \text{PFD}(p_\theta, p_{\theta_*}), \quad \mathcal{E}_{\text{mem}}(\theta) := \text{PFD}(p_\theta, p_{\text{emp}}).$$

Finding I: Empirical Scaling Behavior with $\log_2(N/\sqrt{|\theta|})$

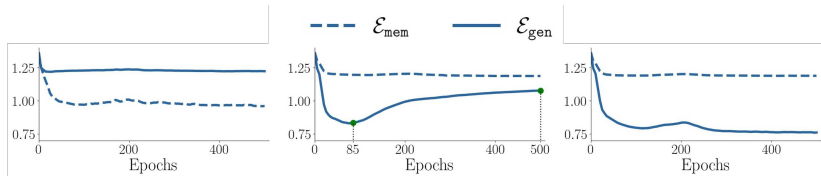


Generalization and Memorization Errors

Suppose s_θ is trained on $\mathcal{S} = \{x^{(i)}\}_{i=1}^N \stackrel{i.i.d.}{\sim} p_{\theta_*}$, and let $p_{\text{emp}}(x) = \frac{1}{N} \sum_{i=1}^N \delta(x - x^{(i)})$. We define

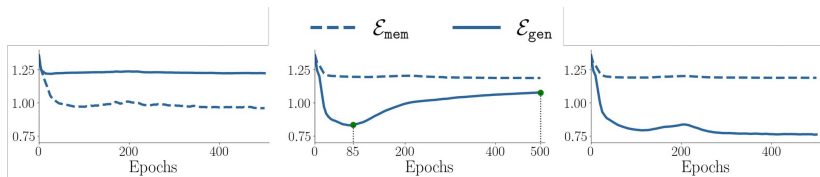
$$\mathcal{E}_{\text{gen}}(\theta) := \text{PFD}(p_\theta, p_{\theta_*}), \quad \mathcal{E}_{\text{mem}}(\theta) := \text{PFD}(p_\theta, p_{\text{emp}}).$$

Finding II: Epoch-wise Double Descent vs. Early Learning



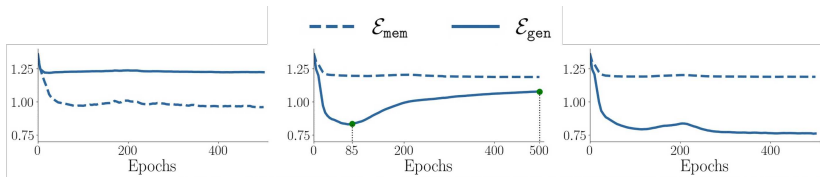
- **Generalization regime (right):** This reveals a clear **epoch-wise double descent** of the generalization error.

Finding II: Epoch-wise Double Descent vs. Early Learning



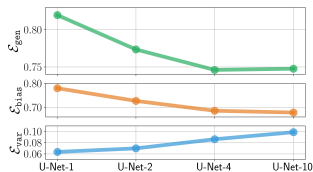
- **Generalization regime (right):** This reveals a clear **epoch-wise double descent** of the generalization error.
- **Transition regime (middle):** The **early generalization** becomes salient as the training sample size increases.

Finding II: Epoch-wise Double Descent vs. Early Learning

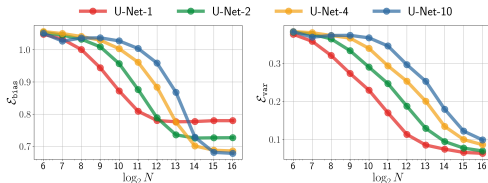


- **Generalization regime (right):** This reveals a clear **epoch-wise double descent** of the generalization error.
- **Transition regime (middle):** The **early generalization** becomes salient as the training sample size increases.

Finding III: Bias-Variance Decomposition



(a)

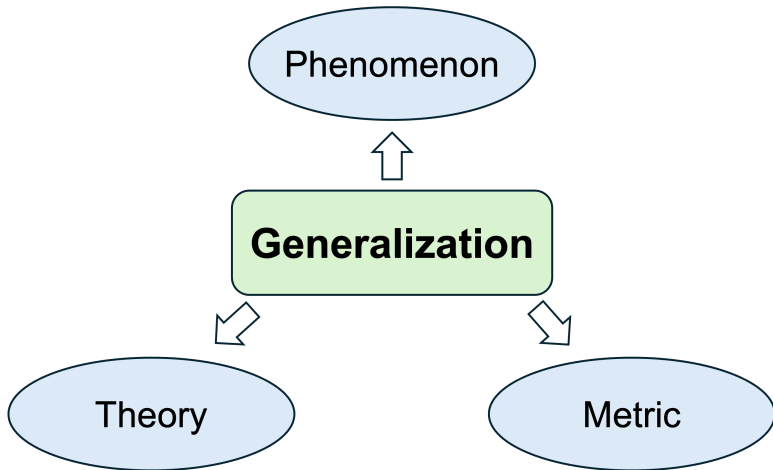


(b)

There exists a decomposition of \mathcal{E}_{gen} , such that

$$\mathbb{E}_{\mathcal{D}} [\mathcal{E}_{\text{gen}}^2(p_{\theta(\mathcal{D})})] = \mathcal{E}_{\text{bias}}^2 + \mathcal{E}_{\text{var}},$$

where $p_{\theta(\mathcal{D})}$ denotes the distribution induced by a diffusion model $\theta(\mathcal{D})$ trained on a given training dataset $\mathcal{D} \sim p_{\text{data}}$.



Conclusion & Discussion

Take-Home Message

- **Phenomenon:** DMs exhibit **unique reproducibility** transition from **memorization** to **generalization** regimes.

Take-Home Message

- **Phenomenon:** DMs exhibit **unique reproducibility** transition from **memorization** to **generalization** regimes.
- **Theory:** We can learn **low-dimensional data distribution** via DMs without the curse of dimensionality.

Take-Home Message

- **Phenomenon:** DMs exhibit **unique reproducibility** transition from **memorization** to **generalization** regimes.
- **Theory:** We can learn **low-dimensional data distribution** via DMs without the curse of dimensionality.
- **Metric:** We can quantify generalization via reproducibility with intriguing **generalization properties**.

Major References

1. H. Zhang*, J. Zhou*, Y. Lu, M. Guo, P. Wang, L. Shen, and Q. Qu. [The Emergence of Reproducibility and Consistency in Diffusion Models](#). *ICML*, 2024. (NeurIPS'23 Workshop on Diffusion Models, **Best Paper Award**)
2. P. Wang*, H. Zhang*, Z. Zhang, S. Chen, Y. Ma, and Q. Qu. [Diffusion Model Learns Low-Dimensional Distributions via Subspace Clustering](#). *Arxiv Preprint arXiv:2409.02426*, 2024.
3. S. Chen*, H. Zhang*, M. Guo, Y. Lu, P. Wang, and Q. Qu. [Exploring Low-Dimensional Subspaces in Diffusion Models for Controllable Image Editing](#). *NeurIPS*, 2024.
4. X. Li, Y. Dai, Q. Qu. [Understanding Generalizability of Diffusion Models Requires Rethinking the Hidden Gaussian Structure](#). *NeurIPS*, 2024.
5. H. Zhang*, Y. Lu*, I. Alkhouri, S. Ravishankar, D. Song, Q. Qu. [Improving Efficiency of Diffusion Models via Multi-Stage Framework and Tailored Multi-Decoder Architectures](#). *CVPR*, 2024.

6. W. Li, H. Zhang, Q. Qu. [Shallow Diffuse: Robust and Invisible Watermarking through Low-Dimensional Subspaces in Diffusion Models](#). *NeurIPS*, 2025 (**spotlight, top 3.2 %**).
7. H. Zhang, Z. Huang, S. Chen, J. Zhou, Z. Zhang, P. Wang, Q. Qu. [Understanding Generalization in Diffusion Models via Probability Flow Distance](#). *Arxiv Preprint arXiv:2505.20123*, 2025.
8. X. Li*, Z. Zhang*, X. Li, S. Chen, Z. Zhu, P. Wang, and Q. Qu. [Understanding Representation Dynamics of Diffusion Models via Low-Dimensional Modeling](#). *NeurIPS*, 2025.
9. X. Li, R. Wang, Q. Qu. [Towards Understanding the Mechanisms of Classifier-Free Guidance](#). *NeurIPS*, 2025 (**spotlight, top 3.2 %**).
10. Lianghe Shi, Meng Wu, Huijie Zhang, Zekai Zhang, Molei Tao, Qing Qu. [A Closer Look at Model Collapse: From a Generalization-to-Memorization Perspective](#). *Neural Information Processing Systems (NeurIPS'25)*, 2025. (**spotlight, top 3.2%**)
11. S. Chen, Y. Zhang, S. Liu, Q. Qu. [The Dual Power of Interpretable Token Embeddings: Jailbreaking Attacks and Defenses for Diffusion Model Unlearning](#). *Arxiv Preprint arXiv:2504.21307*, 2025.

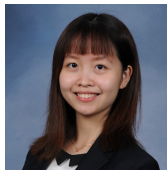
Acknowledgement



Huijie Zhang



Peng Wang



Siyi Chen



Zekai Zhang



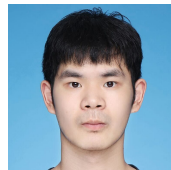
Xiang Li



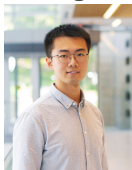
Minzhe Guo



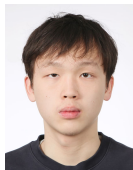
Yifu Lu



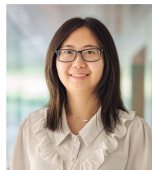
Jinfan Zhou



Lianghe Shi



Zijian Huang



Liyue Shen



Yi Ma

Acknowledgement



Thank You!



SCAN ME