

---

# Understanding Generalization in Diffusion Distillation via Probability Flow Distance

---

Anonymous Authors<sup>1</sup>

## Abstract

Diffusion distillation provides an effective approach for learning lightweight and few-steps diffusion models with efficient generation. However, evaluating their generalization remains challenging: theoretical metrics are often impractical for high-dimensional data, while no practical metrics rigorously measure generalization. In this work, we bridge this gap by introducing probability flow distance (PFD), a theoretically grounded and computationally efficient metric to measure generalization. Specifically, PFD quantifies the distance between distributions by comparing their noise-to-data mappings induced by the probability flow ODE. Using PFD under the diffusion distillation setting, we empirically uncover several key generalization behaviors, including: (1) quantitative scaling behavior from memorization to generalization, (2) epoch-wise double descent training dynamics, and (3) bias-variance decomposition. Beyond these insights, our work lays a foundation for generalization studies in diffusion distillation and bridges them with diffusion training.

## 1. Introduction

While diffusion models (Song et al., 2021c; Lipman et al., 2023; Liu et al., 2023; Ho et al., 2020) have revolutionized generative AI, their real-world deployment remains limited by substantial computational cost. Knowledge distillation (Hinton et al., 2015), which transfers capabilities from a high-capacity teacher to a significantly smaller student, has therefore become a standard approach for practical diffusion-model deployment. Distillation has shown impressive gains in efficiency through both lightweight student architectures (Chen et al., 2025; Kim et al., 2024; Hu et al., 2026) and shorter sampling trajectories (Yin et al., 2024b;a; Salimans

& Ho, 2022; Song et al., 2023; Meng et al., 2023; Sauer et al., 2024; Xie et al., 2024). Beyond efficiency, recent study shows that distillation can also mitigate memorization (Borkar et al., 2026), a known issue in diffusion model training that raises copyright and privacy concerns (Gu et al., 2025; Carlini et al., 2023; Somepalli et al., 2023a;b). In some cases, diffusion distillation even leads to surprising improvements in generation quality (Ma et al., 2025).

However, despite the remarkable performance of distilled models, a comprehensive understanding of their generalization ability and working mechanisms remains elusive. Specifically, existing metrics for evaluating the generalizability of diffusion distillation suffer from significant shortcomings: common empirical metrics like Fréchet inception distance (FID) (Heusel et al., 2017), Inception Score (IS) (Salimans et al., 2016) focus on generation quality, but they cannot distinguish between memorization and generalization, where both can yield high-quality outputs. Neural Network Divergence (NND) (Arora et al., 2017; Gulrajani et al., 2019) proposed to measure the generalizability for generative adversarial networks (GANs) (Goodfellow et al., 2020). However, it requires a large amount of data for evaluation and is not suitable for diffusion distillation. Although recent works measure generalization by evaluating the likelihood of generated samples that are copied from the training data (Zhang et al., 2024; Yoon et al., 2023), this can be misleading, as pure noise may be misclassified as generalized output.

On the other hand, other approaches aim to measure generalization by comparing the distance between the student distribution and the teacher distribution. While distributional distance such as Kullback-Leibler divergence (KL) (Chen et al., 2023e; Nie et al., 2024; Li et al., 2023), total variation (TV) (Chen et al., 2023c; Li et al., 2024b;a; Yang et al., 2024), and 2-Wasserstein distance ( $W_2$ ) (Gao & Zhu, 2025; Bortoli, 2022; Chen et al., 2023a; Gao et al., 2025b) are theoretically appealing, they are often computationally expensive and thus impractical for diffusion distillation. (We defer a detailed discussion in Section 2.4.) As summarized in Table 1, existing metrics are not both accurate and efficient for evaluating diffusion distillation in practice, highlighting the need for a generalization metric

---

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

that is both theoretically grounded and practically tractable.

Moreover, developing a principled evaluation framework within the distillation setting is imperative. Such a framework is crucial not only for deepening our insight into the mechanisms underlying diffusion distillation but also for providing systematic guidance in designing more effective architectures, training strategies, and practical benchmarks.

**Our Contribution.** In this work, we propose a novel metric, termed the probability flow distance (PFD), that can faithfully evaluate the generalization ability of distilled diffusion models. Specifically, our PFD metric quantifies distributional differences by leveraging the backward probability flow ODE (PF-ODE) (Song et al., 2021c), which is widely used in the sampling process of diffusion models. Unlike practical metrics such as FID, which cannot distinguish between memorization and generalization, our PFD provides a theoretically grounded measure of distance between distributions, offering a more reliable assessment of generalization. Compared to theoretical metrics like the Wasserstein distance, PFD is computationally more efficient by leveraging the benign properties of PF-ODE. Moreover, by leveraging PFD, our analysis reveals several intriguing generalization phenomena that offer new insights into the learning behavior of diffusion distillation, which could potentially be informative for training diffusion models in general.

- **Scaling behavior from memorization to generalization.**

Our metric quantitatively characterizes the scaling behavior of diffusion distillation in the transition from memorization to generalization. Specifically, we demonstrate that generalization follows a universal scaling behavior governed by  $N/\sqrt{|\theta|}$ , where  $N$  is the training dataset size and  $|\theta|$  is the number of model parameters. In contrast, prior studies (Zhang et al., 2024; Yoon et al., 2023) have only considered the effects of model capacity or dataset size in isolation, without capturing their joint influence on generalization.

- **Understanding generalization in learning dynamics.**

Our PFD metric reveals and reconciles several intriguing generalization phenomena within the learning dynamics. Specifically, under regimes of sufficient training data, we identify an “**epochwise double descent**” behavior in diffusion distillation, where the generalization error initially decreases, subsequently increases, and finally decreases again toward convergence. While similar phenomena have been observed in overparameterized supervised models, we provide the first empirical validation of this behavior in the context of diffusion distillation. Conversely, when training data is limited, we observe early learning, where models initially generalize but later transition to memorization until convergence.

- **The bias-variance trade-off of generalization errors.**

Finally, we show that our PFD metric naturally introduces

a bias–variance decomposition of the generalization error, extending classical statistical learning theory to diffusion models. Empirically, we observe a trade-off consistent with supervised learning: increasing model capacity reduces bias but increases variance, yielding a characteristic U-shaped generalization error curve. This finding could potentially help us identify overfitting and memorization in the distillation and training of diffusion models.

## 2. Introduction of Probability Flow Distance

In this section, we propose a new metric called probability flow distance (PFD), which is designed to quantify the distance between two arbitrary probability distributions. The design of PFD is motivated by the PF-ODE, which we first review in Section 2.1. We then formally define PFD in Section 2.2 and present its empirical estimation with theoretical guarantees in Section 2.3.

### 2.1. The Mapping Induced by PF-ODE

In general, PF-ODE (Song et al., 2021c) is a class of ordinary differential equations (ODE) that aims to reverse a forward process, where Gaussian noise is progressively added to samples drawn from an underlying distribution, denoted as  $p_{\text{data}}$ . The forward process and the PF-ODE can be described as follows:

- *Forward process.* Given a sample  $\mathbf{x}_0 \stackrel{i.i.d.}{\sim} p_{\text{data}}(\mathbf{x})$ , the forward process progressively corrupts it by adding Gaussian noise. This process can be characterized by the stochastic differential equation (SDE)  $d\mathbf{x}_t = f(t)\mathbf{x}_t dt + g(t)d\mathbf{w}_t$ , where  $t \in [0, T]$  is the time index,  $\{\mathbf{w}_t\}_{t \in [0, T]}$  is a standard Wiener process, and  $f(t), g(t) : \mathbb{R}_+ \rightarrow \mathbb{R}$  are the drift and diffusion functions that control the noise schedule. In this work, we adopt the noise schedule proposed by elucidated diffusion models (EDM) (Karras et al., 2022a), where  $f(t) = 0$  and  $g(t) = \sqrt{2t}$ . Substituting this into the SDE and integrating both sides, we obtain  $\mathbf{x}_t = \mathbf{x}_0 + \int_0^t \sqrt{2\tau} d\mathbf{w}_\tau$ . For ease of exposition, we use  $p_t(\mathbf{x}_t)$  to denote the distribution of the noisy image  $\mathbf{x}_t$  for each  $t \in [0, T]$ . In particular, it is worth noting that  $p_0(\mathbf{x}) = p_{\text{data}}(\mathbf{x})$  and  $p_T(\mathbf{x}) \rightarrow \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$  as  $T \rightarrow +\infty$ .

- *Probability flow ODE.* According to (Song et al., 2021c), the PF-ODE can transform a noise sample  $\mathbf{x}_T$  back into a clean data sample  $\mathbf{x}_0$ . Specifically, under the EDM noise scheduler, the PF-ODE admits the following form:

$$d\mathbf{x}_t = -t \nabla \log p_t(\mathbf{x}_t) dt, \quad (1)$$

where  $\nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t)$  (or simply  $\nabla \log p_t(\mathbf{x}_t)$ ) denotes the *score function* of the distribution  $p_t(\mathbf{x}_t)$  at time  $t \in [0, T]$ . According to (Song et al., 2021c), the backward PF-ODE (1) and the forward SDE have the same distribution at each

timestep  $t$ . In practice, since the score function  $\log p_t(\mathbf{x}_t)$  is unknown, in diffusion models we approximate it using a neural network  $\mathbf{s}_\theta(\mathbf{x}_t, t)$  and employ a numerical solver to generate samples from Equation (1). Additional details are provided in Section A.4.

The backward PF-ODE induces a *unique* mapping  $\Phi_{p_{\text{data}}}$  from  $\mathbf{x}_T$  to  $\mathbf{x}_0$ . Integrating both sides of (1) from  $T$  to 0, the mapping  $\Phi_{p_{\text{data}}}$  can be expressed as:

$$\Phi_{p_{\text{data}}}(\mathbf{x}_T) := \mathbf{x}_T - \int_T^0 t \nabla \log p_t(\mathbf{x}_t) dt. \quad (2)$$

Previous work (Song et al., 2021c) shows that, when  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$  and  $T \rightarrow +\infty$ , the random variable  $\Phi_{p_{\text{data}}}(\mathbf{x}_T)$  is distributed according to  $p_{\text{data}}(\mathbf{x})$ . Consequently, if the data distribution  $p_{\text{data}}$  is known, the score function  $\nabla \log p_t(\mathbf{x}_t)$  is explicitly available, and the backward PF-ODE induces a deterministic and unique mapping from the Gaussian distribution to  $p_{\text{data}}$  (we formally prove uniqueness in the identity property of our Theorem 1).

## 2.2. Definition of Probability Flow Distance

Based on the above setup, we define a metric to measure the distance between any two distributions as follows.

**Definition 1** (Probability flow distance (PFD)). *For any two given distributions  $p$  and  $q$  of the same dimension, we define their distribution distance as*

$$\text{PFD}(p, q) := \left( \mathbb{E}_{\mathbf{x}_T} \left[ \left\| \Psi \circ \Phi_p(\mathbf{x}_T) - \Psi \circ \Phi_q(\mathbf{x}_T) \right\|_2^2 \right] \right)^{1/2}. \quad (3)$$

Here,  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$ ,  $\Phi_p$  and  $\Phi_q$  denote the mappings between the noise and image spaces for distributions  $p$  and  $q$ , respectively, as defined in (2), and  $\Psi(\cdot)$  represents an image descriptor.

Intuitively, when the image descriptor  $\Psi(\cdot)$  is an *identity* mapping, PFD measures the distance between two distributions  $p$  and  $q$  by comparing their respective noise-to-image mappings  $\Phi_p(\cdot)$  and  $\Phi_q(\cdot)$  starting from the same Gaussian noise input  $\mathbf{x}_T$ . Small values of PFD indicate that the two distributions generate similar data from identical noise and, therefore, are close to each other.

However, in practice, measuring distributional distance alone does *not* fully capture the perceptual quality of natural images. To incorporate perceptual quality, we employ an image descriptor  $\Psi(\cdot)$ , which is typically implemented using a pre-trained neural network such as DINOv2 (Oquab et al., 2024) or Self-Supervised Copy Detection Descriptor (SSCD) (Pizzi et al., 2022). Measuring distances in such feature spaces is a common practice in existing metrics for generative models (Heusel et al., 2017; Salimans et al., 2016; Stein et al., 2024), as it tends to align better with

human perception (Stein et al., 2024). As shown by our ablation studies in Section E.2, the PFD metric with an image descriptor effectively captures differences in perceptual quality while models are generalizing, where a smaller quantity indicates better perceptual quality. Therefore, for all experiments in Sections 3 and 4, we use an image descriptor (e.g., SSCD) for  $\Psi(\cdot)$ .

Nonetheless, for simplicity and analytical tractability, we assume the image descriptor  $\Psi(\cdot)$  to be an identity mapping in the following theoretical analysis. Specifically, under Definition 1, we show that PFD satisfies the axioms of a metric (Definition 2.15 in (Rudin, 2021)).

**Theorem 1.** *For any two distributions  $p$  and  $q$ , the PFD satisfies the following properties:*

- (Positivity)  $\text{PFD}(p, q) > 0$  for any  $p \neq q$ .
- (Identity Property)  $\text{PFD}(p, q) = 0$  if and only if  $p = q$ .
- (Symmetry)  $\text{PFD}(p, q) = \text{PFD}(q, p)$ .
- (Triangle Inequality)  $\text{PFD}(p, q) \leq \text{PFD}(p, p') + \text{PFD}(p', q)$  for all  $p'$ .

We defer the proof to Section B. Note that Theorem 1 establishes the theoretical validity of PFD as a metric for measuring the distance between any two probability distributions.

## 2.3. Empirical Estimation of PFD

In practice, the expectation in (3) is intractable due to the complexity of the underlying distributions. Thus, we approximate the PFD using finite samples:

$$\hat{\text{PFD}}(p, q) = \left( \frac{1}{M} \sum_{i=1}^M \left\| \Phi_p(\mathbf{x}_T^{(i)}) - \Phi_q(\mathbf{x}_T^{(i)}) \right\|_2^2 \right)^{1/2}. \quad (4)$$

Here,  $\hat{\text{PFD}}(p, q)$  is the empirical version of  $\text{PFD}(p, q)$  computed over  $M$  independent samples  $\{\mathbf{x}_T^{(i)}\}_{i=1}^M \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$  with  $T \rightarrow \infty$ .

Specifically, our finite-sample approximation relies on two key assumptions: (i) the score functions are smooth at all timesteps, and (ii) the score functions of two distributions remain uniformly close within a bounded region of the input space, which can be described as follows.

**Assumption 1.** *Let  $p$  and  $q$  be two distributions with the same dimension, where we assume:*

- (i) (Lipschitz score functions) *There exists a constant  $L > 0$  such that for all  $\mathbf{x}_1, \mathbf{x}_2$  and  $t \in [0, T]$ , it holds that*

$$\left\| \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_1) - \nabla_{\mathbf{x}} \log p_t(\mathbf{x}_2) \right\|_2 \leq L \left\| \mathbf{x}_1 - \mathbf{x}_2 \right\|_2, \quad (5)$$

and similarly for  $q_t$ .

(ii) (Uniform Closeness) For all  $t \in [0, T]$ , there exists a constant  $\epsilon > 0$  such that

$$\|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) - \nabla_{\mathbf{x}} \log q_t(\mathbf{x})\|_2 \leq \epsilon. \quad (6)$$

The Lipschitz continuity of the score function is a common assumption widely adopted in the theoretical analysis of score functions in diffusion models (Block et al., 2020; Lee et al., 2022; Chen et al., 2023d;a; Zhu et al., 2023; Chen et al., 2023b). More recently, this property has been rigorously established under the assumption that the data distribution is a mixture of Gaussians (Liang et al., 2024). The uniform closeness assumption holds when  $p$  and  $q$  both follow Assumption 1 (i) and have support on compact domains, which is often the case for image distributions. Under Assumption 1, the concentration of the empirical estimate  $\hat{\text{PFD}}(p, q)$  to  $\text{PFD}(p, q)$  can be characterized as follows.

**Theorem 2.** Let  $p$  and  $q$  be two distributions satisfying Assumption 1, and let  $\hat{\text{PFD}}(p, q)$  denote the empirical estimate of  $\text{PFD}(p, q)$  using  $M$  independent samples, as defined in (4). Then, for any  $\gamma > 0$ , this empirical estimate satisfies the following bound:

$$|\hat{\text{PFD}}(p, q) - \text{PFD}(p, q)| \leq \gamma \text{ whenever } M \geq \frac{\kappa^4(L, \epsilon)}{2\gamma^4} \log \frac{2}{\eta}, \quad (7)$$

with probability at least  $1 - \eta$ . Here,  $\kappa(L, \epsilon) := \exp\left(\frac{LT_\xi^2}{2}\right) \xi + \frac{\epsilon}{L} \left(\exp\left(\frac{LT_\xi^2}{2}\right) - 1\right)$  is a constant, with a numerical constant  $\xi > 0$  and a finite timestep  $T_\xi$  depending only on  $\xi$ .

We defer the proof to Section B. Given the score functions of both distributions are smooth and uniformly close, our result in Theorem 2 guarantees that  $\text{PFD}(p, q)$  can be approximated to arbitrary precision by its empirical estimate  $\hat{\text{PFD}}(p, q)$  with high probability, given a finite number of samples.

#### 2.4. Advantages of PFD over Existing Metrics.

As summarized in Table 1, we conclude this section by highlighting the advantages of the proposed PFD over commonly used theoretical metrics for measuring distributional distances, including both density-based and sample-based methods. A comparison with practical evaluation metrics is presented at the end of Section 3.

**Sampling efficiency.** We compare the sampling efficiency of PFD, FID, and  $W_2$  on Gaussian distributions, as shown in Figure 6a, using the experimental settings described in Section C.1. With the same number of estimated samples ( $M = 4096$ ), PFD attains a relative error of approximately  $4 \times 10^{-3}$ , whereas FID and  $W_2$  achieve only about  $2 \times 10^{-2}$ . Consequently, to reach the same level of relative error, PFD requires substantially fewer samples.

Table 1. Metrics comparison.

	Efficient	Accurate
<i>Theoretical distances</i>		
Density-based (KL, TV)	✗	✓
Sample-based ( $W_2$ , MMD)	✗	✓
<i>Practical metrics</i>		
FID, IS, NND, etc.	✓	✗
<b>PFD (Ours)</b>	✓	✓

**Computational efficiency.** Even with the same number of estimated samples  $M$ , PFD requires less computation time than other methods:

- **Comparison with sample-based distances.** Sample-based distances such as the Wasserstein distance and Maximum Mean Discrepancy (MMD) incur  $O(M^2)$  computational complexity. In contrast, PFD requires only  $O(M)$  computation, making it substantially more efficient.
- **Comparison with density-based distances.** Density-based distances, such as KL divergence, total variation distance, and Jensen–Shannon divergence, require approximating probability densities through computationally intensive methods like the Skilling–Hutchinson trace estimator (Song et al., 2021c; Skilling, 1989; Hutchinson, 1989), which become prohibitively expensive in high-dimensional settings for evaluating diffusion models. In contrast, PFD directly estimates the distributional distance using the score function inherently learned by diffusion models. Furthermore, density-based metrics are theoretically ill-suited for image data, as probability densities remain undefined outside the image manifold (Loaiza-Ganem et al., 2024).

### 3. Quantifying Generalization Errors via PFD

In this section, we leverage the PFD metric in Section 2 to rigorously define and evaluate the generalization error of diffusion models under distillation settings. Specifically, this metric enables us to distinguish between memorization and generalization behaviors, as well as analyze the transition from memorization to generation (MtoG).

This MtoG transition has been explored in recent studies for training diffusion models (Yoon et al., 2023; Zhang et al., 2024; Kadkhodaie et al., 2024; Bonnaire et al., 2025), which highlight two learning regimes of diffusion models depending on dataset size and model capacity: (i) **Memorization regime:** Large models trained on small datasets memorize the empirical distribution  $p_{\text{emp}}(\mathbf{x})$  of the training data, yielding poor generalization and no novel samples. (ii) **Generalization regime:** For fixed model capacity, as the number of training samples increases, the model transitions into generalization, approximating the true data distribution

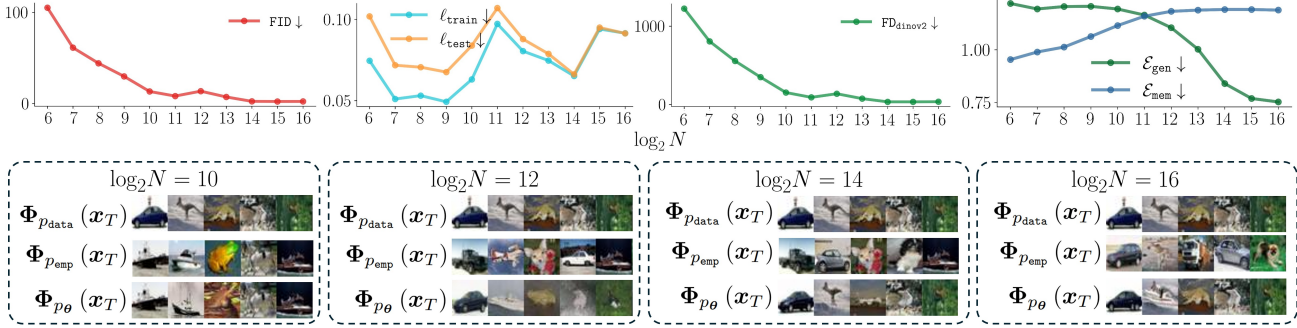


Figure 1. Comparison of practical metrics on the MtoG transition. The top figure plots multiple evaluation metrics as functions of  $\log_2 N$ . The bottom figure visualizes the generation when  $N = 2^{10}, 2^{12}, 2^{14}, 2^{16}$ , sampled from the  $p_{\text{data}}$  (top row), the  $p_{\text{emp}}$  (middle row), and  $p_{\theta}$  (bottom row). The same column shared the same initial noise  $\mathbf{x}_T$ .

$p_{\text{data}}(\mathbf{x})$  and generating new samples.

However, existing approaches (Yoon et al., 2023; Zhang et al., 2024; Alaa et al., 2022) quantify generalization solely by measuring the distance from a generated sample to its nearest neighbor in the training dataset. While effective for identifying memorization, these sample-level measures are inadequate for capturing true generalization because they do not quantify distributional distance. Simply showing that an output is distinct from training data is insufficient; for instance, such metrics could erroneously classify a sample of pure noise as a valid generated sample. To address these limitations, we leverage the PFD metric to assess generalization at the distributional level. Specifically, we quantify how closely the distribution learned via diffusion models,  $p_{\theta}$ , approximates the underlying distribution  $p_{\text{data}}(\mathbf{x})$  and how closely it aligns with the empirical distribution  $p_{\text{emp}}(\mathbf{x})$ . Based on this approach, we formally define generalization and memorization errors as follows.

**Definition 2** (Generalization and Memorization Errors). Consider a diffusion model  $s_{\theta}$  trained on a finite dataset  $\mathcal{D} = \{\mathbf{y}^{(i)}\}_{i=1}^N$ , where each sample  $\mathbf{y}^{(i)}$  is drawn i.i.d. from the underlying distribution  $p_{\text{data}}(\mathbf{x})$ . Denote the learned distribution induced by a diffusion model  $s_{\theta}$  as  $p_{\theta}(\mathbf{x})$ . Using the PFD metric, we can formally define the generalization and memorization errors as follows:

$$\mathcal{E}_{\text{gen}}(\theta) := \text{PFD}(p_{\theta}, p_{\text{data}}), \quad \mathcal{E}_{\text{mem}}(\theta) := \text{PFD}(p_{\theta}, p_{\text{emp}}), \quad (8)$$

where the empirical distribution  $p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{y}^{(i)})$ , with  $\delta(\cdot)$  denoting the Dirac delta function.

Here, given access to  $p_{\text{emp}}(\mathbf{x})$ , the memorization error  $\mathcal{E}_{\text{mem}}(\theta)$  can be exactly computed (see Section D). We further show that  $\mathcal{E}_{\text{mem}}(\theta)$  coincides with metrics introduced in (Yoon et al., 2023; Zhang et al., 2024).

**Evaluation protocol for generalization.** Under the diffusion distillation setting (see Section 4), we study generalization behaviors of a student model distilled from a teacher model. Specifically, a pretrained diffusion model  $s_{\theta_t}(\mathbf{x}_t)$

with parameters  $\theta_t$  serves as the teacher and induces a distribution  $p_{\theta_t}$ , which we treat as the underlying data distribution, i.e.,  $p_{\text{data}} = p_{\theta_t}$ . We employ the most straightforward distillation strategy by directly training a student model  $s_{\theta}$  on samples drawn from  $p_{\theta_t}$ . Generalization behaviors of the student model is then evaluated by comparing the student-induced distribution  $p_{\theta}$  with  $p_{\theta_t}$  using the generalization error metrics defined in Definition 2.

In our experiments for the rest of the paper, both the teacher and student models adopt the U-Net architecture (Ronneberger et al., 2015). The teacher model  $s_{\theta_t}$  is trained on the CIFAR-10 dataset (Krizhevsky et al., 2009) with a fixed model architecture (UNet-10 introduced in Section C.2). The student model  $s_{\theta}$  is trained on samples generated by the teacher, with the number of distillation samples varying from  $N = 2^6$  to  $N = 2^{16}$ , using the same training hyperparameters but different model sizes. For evaluating the generalization error in (8), we compute the PFD between the teacher and student models using  $M = 10^4$  samples drawn from shared initial noise, as defined in (4). Similarly, for the memorization error, we compute the PFD between the student model and the empirical distribution of the training data. Additional details for the evaluation protocol and ablation studies are provided in Section C.3 and Section E.

**Comparison with practical metrics for evaluating generalization.** Before we use the proposed metrics  $\mathcal{E}_{\text{gen}}$  and  $\mathcal{E}_{\text{mem}}$  for revealing the generalization properties of diffusion models in Section 4, we conclude this section by demonstrating their advantages over commonly used practical metrics, such as FID,  $\text{FD}_{\text{DINOv2}}$  (Stein et al., 2024), training and testing loss  $\ell_{\text{train}}, \ell_{\text{test}}$  (see Equation (10)). We defer a more comprehensive comparison with other metrics such as IS, NND, KID (Bińkowski et al., 2018), CMMD (Jayasumana et al., 2024), Precision, and Recall (Kynkäänniemi et al., 2019) to Section C.4.

As shown in Figure 1, we compare several metrics for capturing the MtoG transition under distillation. Among them, only the proposed metric  $\mathcal{E}_{\text{gen}}$  consistently tracks this tran-

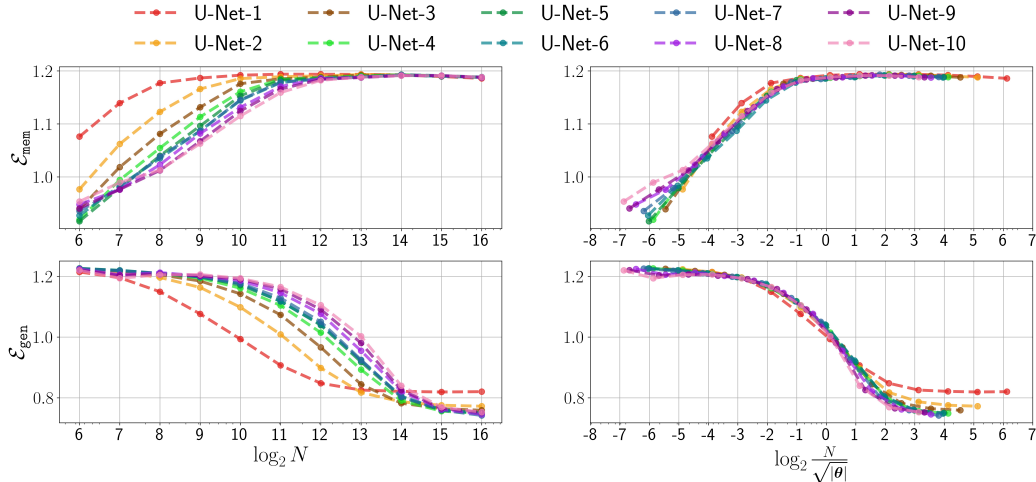


Figure 2. **Scaling behavior in the MtoG transition.**  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  plotted against Left:  $\log_2(N)$  for a range of U-Net architectures (U-Net-1 to U-Net-10). Right: the same metrics plotted against  $\log_2(N/\sqrt{|\theta|})$ , where  $|\theta|$  is the number of model parameters.

sition as the number of distillation samples increases. This conclusion is supported by the qualitative results in the bottom row of Figure 1: when  $N \geq 2^{10}$ , increasing  $N$  makes the distilled model’s generations (bottom row) visually closer to the teacher distribution (top row), indicating improved generalization. Among all metrics, only PFD captures this decreasing trend.

From a theoretical perspective, FID,  $\text{FD}_{\text{DINOv2}}$  relies on a Gaussian assumption for feature distributions, while  $\ell_{\text{train}}$  and  $\ell_{\text{test}}$  provide only upper bounds on the negative log-likelihood of the learned distribution  $p_{\theta}$  (Song et al., 2021b); as a result, none of these metrics can accurately capture generalization. In contrast, PFD is both theoretically well-founded and empirically validated as a reliable metric for measuring generalization.

## 4. Findings of Key Generalization Behaviors

Based on the evaluation protocol in Section 3, this section reveals several key generalization behaviors in diffusion distillation: (i) MtoG scaling behaviors with model capacity and distillation dataset size (Section 4.1), (ii) double descent in learning dynamics (Section 4.2), and (iii) the bias-variance trade-off of generalization error (Section 4.3).

### 4.1. Scaling Behaviors of Generalization

We investigate the scaling behavior of the distilled models with respect to both model capacity  $|\theta|$  and distillation data size  $N$ , using the metrics  $\mathcal{E}_{\text{gen}}$  and  $\mathcal{E}_{\text{mem}}$ . We evaluate ten U-Net architectures over CIFAR-10 dataset, with model sizes ranging from 0.9M to 55.7M parameters (U-Net-1 to U-Net-10). For each model, we compute  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  across varying distillation dataset sizes, following the distillation

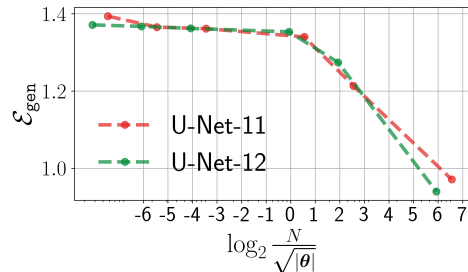


Figure 3. **Scaling behavior over ImageNet dataset.** U-Net-11, U-Net-12 contains 124.2M and 295.9M parameters. Distilling data size  $N$  ranging from  $2^6$  to  $2^{20}$ .

setting outlined in Section 3. In addition, we conduct experiments on the ImageNet dataset (Deng et al., 2009). Results for CIFAR-10 and ImageNet are reported in Figure 2 and Figure 3, respectively. Additional experimental details are provided in Section C.5. From these results, we observe the following:

**MtoG transitions governed by the ratio  $N/\sqrt{|\theta|}$ .** As shown in Figure 2 (left), for a fixed model capacity  $|\theta|$ , our metrics reveal a clear transition from memorization to generalization as the number of distillation samples  $N$  increases, consistent with prior experiments training diffusion models from scratch. (Zhang et al., 2024; Yoon et al., 2023). Moreover, in contrast to prior work that focuses solely on the effect of distillation sample size  $N$ , our results in Figure 2 (right) for CIFAR-10 and Figure 3 for ImageNet reveal a consistent quantitative scaling behavior when using our metric, governed by the ratio  $N/\sqrt{|\theta|}$  between data size and model capacity. Remarkably, both  $\mathcal{E}_{\text{gen}}$  and  $\mathcal{E}_{\text{mem}}$  metrics exhibit near-identical MtoG transition curves across models of varying sizes when plotted against this ratio. As such, analogous to the empirical scaling laws observed in large

language models (Kaplan et al., 2020), this predictable trend provides practical guidance for the development of diffusion models, particularly when scaling up model size, data, or compute to achieve optimal performance gains.

## 4.2. Generalization across Learning Dynamics

Building on the findings in Section 4.1, we further examine the generalization behavior across different training regimes. Under the distillation setting for CIFAR-10 dataset in Section 3, we analyze the learning dynamics of a UNet model with fixed model capacity (UNet-10 introduced in Section C.2) distilled with  $N = 2^6, 2^{12},$  and  $2^{16}$ , corresponding to the memorization, transition, and generalization regimes in Section 4.1, respectively. The model is distilled using stochastic gradient descent (SGD) for 500 epochs, during which we track  $\mathcal{E}_{\text{mem}}, \mathcal{E}_{\text{gen}}, \ell_{\text{train}},$  and  $\ell_{\text{test}}$  at each epoch. The results in Figure 4 reveal several notable generalization behaviors that align with phenomena previously observed in the training of overparameterized deep models (Zhang et al., 2017; Nakkiran et al., 2021):

**Epochwise double descent of the generalization error in the generalization regime.** In contrast, as shown in Figure 4 (right), distilling in the generalization regime ( $N = 2^{16}$ ) reveals a clear instance of the *double descent* phenomenon (Nakkiran et al., 2021) in the generalization error. Specifically, the error initially decreases, then increases during intermediate distilling epochs, and finally decreases again as distilling approaches convergence. Notably, this non-monotonic behavior is not captured by the standard training and test losses  $\ell_{\text{train}}$  and  $\ell_{\text{test}}$ , both of which decrease monotonically throughout distilling. This implies that extended distilling can improve generalization performance in the generalization regime.

**Epochwise early learning behavior in memorization and transition regimes.** As shown in Figure 4 (left & middle), in both the memorization ( $N = 2^6$ ) and transition ( $N = 2^{12}$ ) regimes, the generalization error initially decreases during distilling but reaches its minimum at an early epoch, after which it begins to increase again. This *early learning* (or early generalization) phenomenon becomes more salient as the distillation sample size increases from the memorization to the transition regime. As shown in the visualization at the bottom of Section 3, the model at Epoch 85 clearly exhibits generalization, whereas the model at Epoch 500 fails to generalize. This is also corroborated by the divergence of training loss  $\ell_{\text{train}}$  and test loss  $\ell_{\text{test}}$  at the top of the figure. It is worth mentioning that, although early learning behavior has been theoretically and visually demonstrated in previous works (Li et al., 2024d; 2023; Bonnaire et al., 2025), our work develops the first framework for precisely capturing this phenomenon.

## 4.3. Bias-variance Trade-off of the Generalization Error

In statistical learning theory, bias-variance trade-off is a classical yet fundamental concept in supervised learning which helps us understand and analyze the sources of prediction error in the model (Kohavi et al., 1996; Hastie et al., 2009; Yang et al., 2020; Belkin et al., 2019). Specifically, bias-variance decomposition expresses the expected generalization error as the sum of two components: (i) the *bias term*, which quantifies the discrepancy between the expected model prediction and the true function: high bias indicates systematic error or underfitting; and (ii) the *variance term*, which measures the prediction variability of the model across different training sets: high variance reflects sensitivity to data fluctuations or overfitting.

However, in unsupervised learning settings such as diffusion models, the notion of generalization error was not well-defined prior to our work, in contrast to the well-established definitions in supervised learning. As a result, bias-variance decomposition in this context remains largely unexplored. In this work, we address this gap through the generalization error measure  $\mathcal{E}_{\text{gen}}$  (see Equation (8)), which admits a bias-variance decomposition analogous to that in the supervised setting, as we detail below.

**Definition 3** (Bias-Variance Decomposition of  $\mathcal{E}_{\text{gen}}$ ). *Based on the same setup as Definition 2, we can decompose  $\mathcal{E}_{\text{gen}}$  in Equation (8) as*

$$\mathbb{E}_{\mathcal{D}} [\mathcal{E}_{\text{gen}}^2 (p_{\theta(\mathcal{D})})] = \mathcal{E}_{\text{bias}}^2 + \mathcal{E}_{\text{var}} \quad (9)$$

where  $p_{\theta(\mathcal{D})}$  denotes the distribution induced by a diffusion model  $\theta(\mathcal{D})$  trained on a given training dataset  $\mathcal{D}$  sampled from  $p_{\text{data}}$ . Specifically, the bias and variance terms are defined as:

$$\begin{aligned} \mathcal{E}_{\text{bias}} &:= \mathbb{E}_{\mathbf{x}_T} (\|\Psi \circ \Phi_{p_{\text{data}}}(\mathbf{x}_T) - \overline{\Psi \circ \Phi_{p_{\theta}}}(\mathbf{x}_T)\|_2^2)^{1/2}, \\ \mathcal{E}_{\text{var}} &:= \mathbb{E}_{\mathcal{D}} \mathbb{E}_{\mathbf{x}_T} [\|\Psi \circ \Phi_{p_{\theta(\mathcal{D})}}(\mathbf{x}_T) - \overline{\Psi \circ \Phi_{p_{\theta}}}(\mathbf{x}_T)\|_2^2], \end{aligned}$$

with  $\overline{\Psi \circ \Phi_{p_{\theta}}}(\cdot) := \mathbb{E}_{\mathcal{D}} [\Psi \circ \Phi_{p_{\theta(\mathcal{D})}}(\cdot)]$ .

Intuitively, our definitions of the bias term  $\mathcal{E}_{\text{bias}}$  and the variance term  $\mathcal{E}_{\text{var}}$  are both well-justified: (i)  $\mathcal{E}_{\text{bias}}$  quantifies the systematic error between the learned distribution  $p_{\theta}$  and the ground-truth distribution  $p_{\text{data}}$ ; and (ii)  $\mathcal{E}_{\text{var}}$  captures the variability of model predictions across different training sets by measuring the distance between  $p_{\theta}$  and the mean  $\overline{p_{\theta}}$  which can be empirically estimated by averaging over multiple datasets  $\mathcal{D}$  sampled from  $p_{\text{data}}$ . Experimental results, following the protocol in Section 3, are shown in Figure 5, with detailed settings in Section C.7.

In Figure 5 (a), when diffusion models are distilled in the generalization regime, the resulting generalization decomposition aligns with classical bias-variance theory from supervised learning: as model complexity increases, the

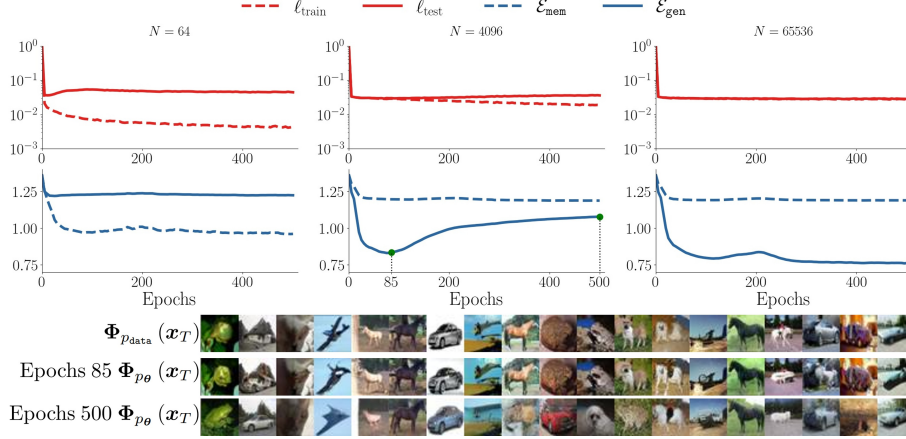


Figure 4. **Training dynamics of diffusion models in different regimes.** The top figure plots  $\mathcal{E}_{\text{mem}}$ ,  $\mathcal{E}_{\text{gen}}$ ,  $\ell_{\text{train}}$ ,  $\ell_{\text{test}}$  over training epochs for different different dataset sizes:  $N = 2^6$  (left),  $2^{12}$  (middle),  $2^{16}$  (right). The bottom figure visualizes the generation when  $N = 2^{12}$ . The top row shows samples from the underlying distribution  $\Phi_{p_{\text{data}}}(x_T)$ , while the middle and bottom rows display outputs from the trained diffusion model  $\Phi_{p_{\theta}}(x_T)$  at epoch 85 and 500, respectively.

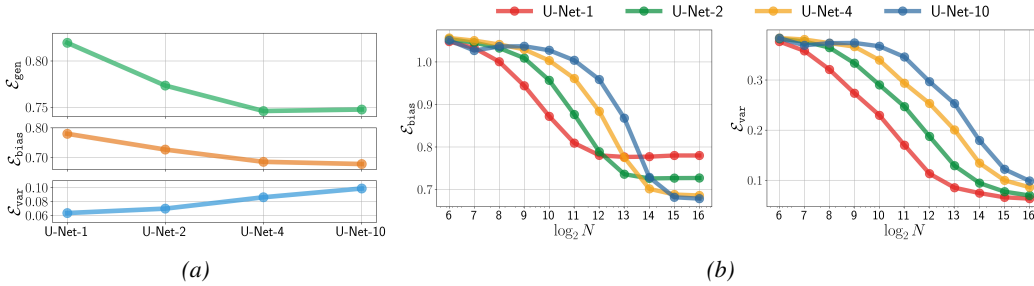


Figure 5. **Bias-Variance Trade-off.** (a) plots the generalization error  $\mathcal{E}_{\text{gen}}$ , bias  $\mathcal{E}_{\text{bias}}$ , and variance  $\mathcal{E}_{\text{var}}$  across different network architectures with a fixed distillation sample size of  $N = 2^{16}$ . (b) shows  $\mathcal{E}_{\text{bias}}$  and  $\mathcal{E}_{\text{var}}$  as functions of the number of distillation samples  $N$  for various network architectures.

bias term  $\mathcal{E}_{\text{bias}}$  decreases while the variance term  $\mathcal{E}_{\text{var}}$  increases, resulting in a U-shaped generalization error curve. Additionally, Figure 5 (b) further illustrates the effect of the distillation sample size  $N$  and number of parameters  $|\theta|$ : increasing  $N$  reduces both  $\mathcal{E}_{\text{bias}}$  and  $\mathcal{E}_{\text{var}}$ , thereby lowering the generalization error  $\mathcal{E}_{\text{gen}}$ , as expected; In contrast, increasing  $|\theta|$  consistently increases  $\mathcal{E}_{\text{var}}$ , and its effect on  $\mathcal{E}_{\text{bias}}$  depends on the size of  $N$ : it decreases  $\mathcal{E}_{\text{bias}}$  when  $N \geq 2^{15}$  but increases it when  $N \leq 2^{11}$ .

## 5. Discussion & Conclusion

While PFD is designed to evaluate generalization in distillation settings, we believe its impact extends beyond distillation and is highly relevant to diffusion model training:

### Distillation is closely aligned with training from scratch.

The teacher-student framework has been widely adopted to understand learning phenomena observed in real-world scenarios, both empirically (Betzalel et al., 2024; Lee et al., 2021; Saglietti et al., 2022) and theoretically (Advani et al.,

2020; Goldt et al., 2019; d’Ascoli et al., 2020). In this work, the phenomena we observe, such as the MtoG transition (left column of Figure 2) and early learning dynamics (middle column of Figure 4), are consistent with observations reported in prior studies of models trained from scratch. This alignment suggests that distillation provides a tractable and controllable setting for understanding training in real-world.

### Training from scratch is increasingly close to distillation.

Modern diffusion models are often trained on large-scale datasets constructed from heterogeneous sources. Due to the scarcity of sufficiently high-quality real-world data, large-scale datasets such as LAION-5B (Schuhmann et al., 2022) inevitably contain a substantial fraction of synthetic data generated by foundation models. As a result, training from scratch increasingly resembles a teacher-student setting, where models implicitly learn from other pretrained models rather than purely from natural data.

Taken together, these facts motivate the use of PFD beyond explicit distillation settings, as a general tool for analyzing generalization in real-world large-scale diffusion models.

## 6. Impact Statement

This work advances the theoretical and empirical understanding of generalization in diffusion models by introducing Probability Flow Distance (PFD). Improved understanding of generalization can help mitigate risks associated with memorization, including privacy leakage and copyright concerns. We do not anticipate direct negative societal impacts arising from this work; however, as with all advances in generative modeling, downstream applications should be developed and deployed responsibly.

## References

Advani, M. S., Saxe, A. M., and Sompolinsky, H. High-dimensional dynamics of generalization error in neural networks. *Neural Networks*, 132:428–446, 2020.

Aithal, S. K., Maini, P., Lipton, Z., and Kolter, J. Z. Understanding hallucinations in diffusion models through mode interpolation. In *Advances in Neural Information Processing Systems*, volume 37, pp. 134614–134644, 2024.

Alaa, A., Van Breugel, B., Saveliev, E. S., and Van Der Schaar, M. How faithful is your synthetic data? sample-level metrics for evaluating and auditing generative models. In *International Conference on Machine Learning*, pp. 290–306. PMLR, 2022.

Ambrogioni, L. In search of dispersed memories: Generative diffusion models are associative memory networks. *Entropy*, 26(5):381, 2024.

Arora, S., Ge, R., Liang, Y., Ma, T., and Zhang, Y. Generalization and equilibrium in generative adversarial nets (gans). In *International Conference on Machine Learning*, pp. 224–232. PMLR, 2017.

Bao, F., Nie, S., Xue, K., Cao, Y., Li, C., Su, H., and Zhu, J. All are worth words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023.

Belkin, M., Hsu, D., Ma, S., and Mandal, S. Reconciling modern machine-learning practice and the classical bias-variance trade-off. *Proceedings of the National Academy of Sciences*, 116(32):15849–15854, 2019.

Berthelot, D., Autef, A., Lin, J., Yap, D. A., Zhai, S., Hu, S., Zheng, D., Talbott, W., and Gu, E. Tract: Denoising diffusion models with transitive closure time-distillation. *arXiv preprint arXiv:2303.04248*, 2023.

Betzalel, E., Penso, C., and Fetaya, E. Evaluation metrics for generative models: An empirical study. *Machine Learning and Knowledge Extraction*, 6(3):1531–1544, 2024.

Bhatia, R. *Matrix analysis*, volume 169. Springer Science & Business Media, 2013.

Bińkowski, M., Sutherland, D. J., Arbel, M., and Gretton, A. Demystifying mmd gans. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=r1lUOzWCW>.

Block, A., Mroueh, Y., and Rakhlin, A. Generative modeling with denoising auto-encoders and langevin sampling. *arXiv preprint arXiv:2002.00107*, 2020.

Bonnaire, T., Urfin, R., Biroli, G., and Mezard, M. Why diffusion models don’t memorize: The role of implicit dynamical regularization in training. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025. URL <https://openreview.net/forum?id=BSZqpqqgM0>.

Borkar, J., Chadha, K., Mireshghallah, N., Zhang, Y., Velićhe, I.-E., Mitra, A., Smith, D. A., Xu, Z., and Garcia-Olano, D. Memorization dynamics in knowledge distillation for language models. *arXiv preprint arXiv:2601.15394*, 2026.

Bortoli, V. D. Convergence of denoising diffusion models under the manifold hypothesis. *Transactions on Machine Learning Research*, 2022. ISSN 2835-8856. URL <https://openreview.net/forum?id=MhK5aXo3gB>. Expert Certification.

Carlini, N., Hayes, J., Nasr, M., Jagielski, M., Sehwag, V., Tramèr, F., Balle, B., Ippolito, D., and Wallace, E. Extracting training data from diffusion models. In *32nd USENIX security symposium (USENIX Security 23)*, pp. 5253–5270, 2023.

Chen, H., Lee, H., and Lu, J. Improved analysis of score-based generative modeling: User-friendly bounds under minimal smoothness assumptions. In *International Conference on Machine Learning*, pp. 4735–4763. PMLR, 2023a.

Chen, J., Hu, D., Huang, X., Coskun, H., Sahni, A., Gupta, A., Goyal, A., Lahiri, D., Singh, R., Idelbayev, Y., et al. Snapgen: Taming high-resolution text-to-image models for mobile devices with efficient architectures and training. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 7997–8008, 2025.

Chen, M., Huang, K., Zhao, T., and Wang, M. Score approximation, estimation and distribution recovery of diffusion models on low-dimensional data. In *Proceedings of the 40th International Conference on Machine Learning*, 2023b.

- 495 Chen, S., Chewi, S., Lee, H., Li, Y., Lu, J., and Salim, A.  
 496 The probability flow ODE is provably fast. In *Thirty-*  
 497 *seventh Conference on Neural Information Processing*  
 498 *Systems*, 2023c. URL [https://openreview.net/](https://openreview.net/forum?id=KD6MFeWSAd)  
 499 [forum?id=KD6MFeWSAd](https://openreview.net/forum?id=KD6MFeWSAd).
- 500 Chen, S., Chewi, S., Li, J., Li, Y., Salim, A., and Zhang,  
 501 A. Sampling is as easy as learning the score: theory for  
 502 diffusion models with minimal data assumptions. In *The*  
 503 *Eleventh International Conference on Learning Represen-*  
 504 *tations*, 2023d. URL [https://openreview.net/](https://openreview.net/forum?id=zyLVMgsZ0U_)  
 505 [forum?id=zyLVMgsZ0U\\_](https://openreview.net/forum?id=zyLVMgsZ0U_).
- 506 Chen, S., Daras, G., and Dimakis, A. Restoration-  
 507 degradation beyond linear diffusions: A non-asymptotic  
 508 analysis for ddim-type samplers. In *International Con-*  
 509 *ference on Machine Learning*, pp. 4462–4484. PMLR,  
 510 2023e.
- 511 Chen, Z. On the interpolation effect of score smoothing.  
 512 *arXiv preprint arXiv:2502.19499*, 2025.
- 513 Chong, M. J. and Forsyth, D. Effectively unbiased fid and  
 514 inception score and where to find them. In *Proceedings*  
 515 *of the IEEE/CVF Conference on Computer Vision and*  
 516 *Pattern Recognition*, pp. 6070–6079, 2020.
- 517 Coddington, E. A. and Levinson, N. *Theory of ordinary*  
 518 *differential equations*. McGraw-Hill New York, 1955.
- 519 d’Ascoli, S., Sagun, L., and Biroli, G. Triple descent and the  
 520 two kinds of overfitting: Where & why do they appear?  
 521 In *Advances in Neural Information Processing Systems*,  
 522 volume 33, pp. 3058–3069, 2020.
- 523 Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei,  
 524 L. Imagenet: A large-scale hierarchical image database.  
 525 In *2009 IEEE conference on computer vision and pattern*  
 526 *recognition*, pp. 248–255. Ieee, 2009.
- 527 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn,  
 528 D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer,  
 529 M., Heigold, G., Gelly, S., Uszkoreit, J., and Houlsby,  
 530 N. An image is worth 16x16 words: Transformers for  
 531 image recognition at scale. In *International Conference*  
 532 *on Learning Representations*, 2021. URL [https://](https://openreview.net/forum?id=YicbFdNTTy)  
 533 [openreview.net/forum?id=YicbFdNTTy](https://openreview.net/forum?id=YicbFdNTTy).
- 534 Gao, R., Hoogeboom, E., Heek, J., Bortoli, V. D., Murphy,  
 535 K. P., and Salimans, T. Diffusion models and gaussian  
 536 flow matching: Two sides of the same coin. In *The Fourth*  
 537 *Blogpost Track at ICLR 2025*, 2025a. URL [https://](https://openreview.net/forum?id=C8Yyg9wy0s)  
 538 [openreview.net/forum?id=C8Yyg9wy0s](https://openreview.net/forum?id=C8Yyg9wy0s).
- 539 Gao, X. and Zhu, L. Convergence analysis for general  
 540 probability flow ODEs of diffusion models in wasserstein  
 541 distances. In *The 28th International Conference on Ar-*  
 542 *tificial Intelligence and Statistics*, 2025. URL [https://](https://openreview.net/forum?id=EkO8rb3liX)  
 543 [openreview.net/forum?id=EkO8rb3liX](https://openreview.net/forum?id=EkO8rb3liX).
- 544 Gao, X., Nguyen, H. M., and Zhu, L. Wasserstein con-  
 545 vergence guarantees for a general class of score-based  
 546 generative models. *Journal of Machine Learning Re-*  
 547 *search*, 26(43):1–54, 2025b. URL [http://jmlr.](http://jmlr.org/papers/v26/24-0902.html)  
 548 [org/papers/v26/24-0902.html](http://jmlr.org/papers/v26/24-0902.html).
- 549 Goldt, S., Advani, M., Saxe, A. M., Krzakala, F., and Zde-  
 borová, L. Dynamics of stochastic gradient descent for  
 two-layer neural networks in the teacher-student setup.  
 In *Advances in Neural Information Processing Systems*,  
 volume 32, 2019.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B.,  
 Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y.  
 Generative adversarial networks. *Communications of the*  
*ACM*, 63(11):139–144, 2020.
- Gu, X., Du, C., Pang, T., Li, C., Lin, M., and Wang, Y.  
 On memorization in diffusion models. *arXiv preprint*  
*arXiv:2310.02664*, 2023.
- Gu, X., Du, C., Pang, T., Li, C., Lin, M., and Wang,  
 Y. On memorization in diffusion models. *Transac-*  
*tions on Machine Learning Research*, 2025. ISSN 2835-  
 8856. URL [https://openreview.net/forum?](https://openreview.net/forum?id=D3DBqvSDBj)  
[id=D3DBqvSDBj](https://openreview.net/forum?id=D3DBqvSDBj).
- Gulrajani, I., Raffel, C., and Metz, L. Towards GAN  
 benchmarks which require generalization. In *In-*  
*ternational Conference on Learning Representations*,  
 2019. URL [https://openreview.net/forum?](https://openreview.net/forum?id=HkxKH2AcFm)  
[id=HkxKH2AcFm](https://openreview.net/forum?id=HkxKH2AcFm).
- Hastie, T., Tibshirani, R., Friedman, J., et al. The elements  
 of statistical learning, 2009.
- Hessel, J., Holtzman, A., Forbes, M., Bras, R. L., and Choi,  
 Y. CLIPScore: a reference-free evaluation metric for  
 image captioning. In *Empirical Methods in Natural Lan-*  
*guage Processing*, 2021.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and  
 Hochreiter, S. Gans trained by a two time-scale update  
 rule converge to a local nash equilibrium. In *Advances*  
*in Neural Information Processing Systems*, volume 30,  
 2017.
- Hinton, G., Vinyals, O., and Dean, J. Distilling  
 the knowledge in a neural network. *arXiv preprint*  
*arXiv:1503.02531*, 2015.
- Ho, J. and Salimans, T. Classifier-free diffusion guidance.  
*arXiv preprint arXiv:2207.12598*, 2022.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion prob-  
 abilistic models. In *Advances in Neural Information*  
*Processing Systems*, volume 33, pp. 6840–6851, 2020.

- 550 Hoeffding, W. Probability inequalities for sums of bounded  
551 random variables. *The Collected Works of Wassily Ho-*  
552 *effding*, pp. 409–426, 1994.
- 553 Hu, D., Gupta, A., Gabidolla, M., Sahni, A., Coskun, H., Li,  
554 Y., Idelbayev, Y., Mahmood, A., Lebedev, A., Lahiri, D.,  
555 et al. Snappgen++: Unleashing diffusion transformers for  
556 efficient high-fidelity image generation on edge devices.  
557 *arXiv preprint arXiv:2601.08303*, 2026.
- 559 Hutchinson, M. F. A stochastic estimator of the trace of the  
560 influence matrix for laplacian smoothing splines. *Com-*  
561 *munications in Statistics-Simulation and Computation*,  
562 18(3):1059–1076, 1989.
- 563 Jayasumana, S., Ramalingam, S., Veit, A., Glasner, D.,  
564 Chakrabarti, A., and Kumar, S. Rethinking fid: Towards a  
565 better evaluation metric for image generation. In *Proce-*  
566 *edings of the IEEE/CVF Conference on Computer Vision*  
567 *and Pattern Recognition*, pp. 9307–9315, 2024.
- 569 Jolicœur-Martineau, A., Piché-Taillefer, R., Mitliagkas,  
570 I., and des Combes, R. T. Adversarial score match-  
571 ing and improved sampling for image generation. In  
572 *International Conference on Learning Representations*,  
573 2021. URL [https://openreview.net/forum?](https://openreview.net/forum?id=eLfqMl3z3lq)  
574 [id=eLfqMl3z3lq](https://openreview.net/forum?id=eLfqMl3z3lq).
- 575 Kadkhodaie, Z., Guth, F., Simoncelli, E. P., and Mallat, S.  
576 Generalization in diffusion models arises from geometry-  
577 adaptive harmonic representations. In *The Twelfth In-*  
578 *ternational Conference on Learning Representations*,  
579 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=ANvmVS2Yr0)  
580 [id=ANvmVS2Yr0](https://openreview.net/forum?id=ANvmVS2Yr0).
- 582 Kamb, M. and Ganguli, S. An analytic theory of creativ-  
583 ity in convolutional diffusion models. *arXiv preprint*  
584 *arXiv:2412.20292*, 2024.
- 585 Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B.,  
586 Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and  
587 Amodi, D. Scaling laws for neural language models.  
588 *arXiv preprint arXiv:2001.08361*, 2020.
- 590 Karras, T., Aittala, M., Aila, T., and Laine, S. Elucidating  
591 the design space of diffusion-based generative models.  
592 In *Advances in Neural Information Processing Systems*,  
593 volume 35, pp. 26565–26577, 2022a.
- 594 Karras, T., Aittala, M., Aila, T., and Laine, S. Eluci-  
595 dating the design space of diffusion-based generative  
596 models. [https://github.com/NVlabs/edm/](https://github.com/NVlabs/edm/tree/main)  
597 [tree/main](https://github.com/NVlabs/edm/tree/main)>, 2022b.
- 599 Karras, T., Aittala, M., Lehtinen, J., Hellsten, J., Aila, T.,  
600 and Laine, S. Analyzing and improving the training  
601 dynamics of diffusion models. In *Proceedings of the*  
602 *IEEE/CVF Conference on Computer Vision and Pattern*  
603 *Recognition*, pp. 24174–24184, 2024.
- 604 Kim, B.-K., Song, H.-K., Castells, T., and Choi, S. Bk-  
sdm: A lightweight, fast, and cheap version of stable  
diffusion. In *European Conference on Computer Vision*,  
pp. 381–399. Springer, 2024.
- Kohavi, R., Wolpert, D. H., et al. Bias plus variance de-  
composition for zero-one loss functions. In *International*  
*Conference on Machine Learning*, volume 96, pp. 275–  
283. Citeseer, 1996.
- Krizhevsky, A., Hinton, G., et al. Learning multiple layers  
of features from tiny images. 2009.
- Kynkäänniemi, T., Karras, T., Laine, S., Lehtinen, J., and  
Aila, T. Improved precision and recall metric for assess-  
ing generative models. In *Advances in Neural Information*  
*Processing Systems*, volume 32, 2019.
- Lee, H., Lu, J., and Tan, Y. Convergence for score-based  
generative modeling with polynomial complexity. In  
*Advances in Neural Information Processing Systems*, vol-  
ume 35, pp. 22870–22882, 2022.
- Lee, S., Goldt, S., and Saxe, A. Continual learning in  
the teacher-student setup: Impact of task similarity. In  
*International Conference on Machine Learning*, pp. 6109–  
6119. PMLR, 2021.
- Li, G., Huang, Y., Efimov, T., Wei, Y., Chi, Y., and Chen, Y.  
Accelerating convergence of score-based diffusion mod-  
els, provably. In *Proceedings of the 41st International*  
*Conference on Machine Learning*, 2024a.
- Li, G., Wei, Y., Chen, Y., and Chi, Y. Towards non-  
asymptotic convergence for diffusion-based generative  
models. In *The Twelfth International Conference on*  
*Learning Representations*, 2024b. URL [https://](https://openreview.net/forum?id=4VGEeER6W9)  
[openreview.net/forum?id=4VGEeER6W9](https://openreview.net/forum?id=4VGEeER6W9).
- Li, G., Wei, Y., Chi, Y., and Chen, Y. A sharp convergence  
theory for the probability flow odes of diffusion models.  
*arXiv preprint arXiv:2408.02320*, 2024c.
- Li, P., Li, Z., Zhang, H., and Bian, J. On the generalization  
properties of diffusion models. In *Advances in Neural*  
*Information Processing Systems*, volume 36, pp. 2097–  
2127, 2023.
- Li, X., Dai, Y., and Qu, Q. Understanding generaliz-  
ability of diffusion models requires rethinking the hid-  
den gaussian structure. In *The Thirty-eighth Annual*  
*Conference on Neural Information Processing Systems*,  
2024d. URL [https://openreview.net/forum?](https://openreview.net/forum?id=Sk2duBGvrK)  
[id=Sk2duBGvrK](https://openreview.net/forum?id=Sk2duBGvrK).
- Liang, J., Huang, Z., and Chen, Y. Low-dimensional adapta-  
tion of diffusion models: Convergence in total variation.  
*arXiv preprint arXiv:2501.12982*, 2025.

- 605 Liang, Y., Shi, Z., Song, Z., and Zhou, Y. Unraveling the  
606 smoothness properties of diffusion models: A gaussian  
607 mixture perspective. *arXiv preprint arXiv:2405.16418*,  
608 2024.
- 609 Lipman, Y., Chen, R. T. Q., Ben-Hamu, H., Nickel, M., and  
610 Le, M. Flow matching for generative modeling. In *The*  
611 *Eleventh International Conference on Learning Represen-*  
612 *tations*, 2023. URL [https://openreview.net/](https://openreview.net/forum?id=PqvMRDCJT9t)  
613 [forum?id=PqvMRDCJT9t](https://openreview.net/forum?id=PqvMRDCJT9t).
- 614 Liu, X., Gong, C., and qiang liu. Flow straight and fast:  
615 Learning to generate and transfer data with rectified flow.  
616 In *The Eleventh International Conference on Learning*  
617 *Representations*, 2023. URL [https://openreview.](https://openreview.net/forum?id=XVjTTLnw5z)  
618 [net/forum?id=XVjTTLnw5z](https://openreview.net/forum?id=XVjTTLnw5z).
- 619 Loiza-Ganem, G., Ross, B. L., Hosseinzadeh, R., Caterini,  
620 A. L., and Cresswell, J. C. Deep generative models  
621 through the lens of the manifold hypothesis: A survey  
622 and new connections. *Transactions on Machine Learn-*  
623 *ing Research*, 2024. ISSN 2835-8856. URL [https://](https://openreview.net/forum?id=a90WpmSi0I)  
624 [openreview.net/forum?id=a90WpmSi0I](https://openreview.net/forum?id=a90WpmSi0I). Sur-  
625 vey Certification, Expert Certification.
- 626 Lukoianov, A., Yuan, C., Solomon, J., and Sitzmann, V.  
627 Locality in image diffusion models emerges from data  
628 statistics. In *The Thirty-ninth Annual Conference on Neu-*  
629 *ral Information Processing Systems*, 2025. URL [https:](https://openreview.net/forum?id=skunuOdav0)  
630 [//openreview.net/forum?id=skunuOdav0](https://openreview.net/forum?id=skunuOdav0).
- 631 Luo, C. Understanding diffusion models: A unified perspec-  
632 tive. *arXiv preprint arXiv:2208.11970*, 2022.
- 633 Ma, X., Yu, R., Liu, S., Fang, G., and Wang, X. Diffusion  
634 model is effectively its own teacher. In *Proceedings of*  
635 *the Computer Vision and Pattern Recognition Conference*,  
636 pp. 12901–12911, 2025.
- 637 Meng, C., Rombach, R., Gao, R., Kingma, D., Ermon, S.,  
638 Ho, J., and Salimans, T. On distillation of guided diffu-  
639 sion models. In *Proceedings of the IEEE/CVF conference*  
640 *on computer vision and pattern recognition*, pp. 14297–  
641 14306, 2023.
- 642 Nakkiran, P., Kaplun, G., Bansal, Y., Yang, T., Barak, B.,  
643 and Sutskever, I. Deep double descent: Where bigger  
644 models and more data hurt. *Journal of Statistical Mechan-*  
645 *ics: Theory and Experiment*, 2021(12):124003, 2021.
- 646 Nichol, A. Q. and Dhariwal, P. Improved denoising diffusion  
647 probabilistic models. In *International Conference on*  
648 *Machine Learning*, pp. 8162–8171. PMLR, 2021.
- 649 Nie, S., Guo, H. A., Lu, C., Zhou, Y., Zheng, C., and Li,  
650 C. The blessing of randomness: SDE beats ODE in  
651 general diffusion-based image editing. In *The Twelfth*  
652 *International Conference on Learning Representations*,  
653 2024. URL [https://openreview.net/forum?](https://openreview.net/forum?id=DesYwmUG00)  
654 [id=DesYwmUG00](https://openreview.net/forum?id=DesYwmUG00).
- 655 Niedoba, M., Zwartsenberg, B., Murphy, K., and Wood, F.  
656 Towards a mechanistic explanation of diffusion model  
657 generalization. *arXiv preprint arXiv:2411.19339*, 2024.
- 658 Oquab, M., Darcet, T., Moutakanni, T., Vo, H. V.,  
659 Szafraniec, M., Khalidov, V., Fernandez, P., HAZIZA,  
D., Massa, F., El-Nouby, A., Assran, M., Ballas, N.,  
Galuba, W., Howes, R., Huang, P.-Y., Li, S.-W., Misra,  
I., Rabbat, M., Sharma, V., Synnaeve, G., Xu, H., Je-  
gou, H., Mairal, J., Labatut, P., Joulin, A., and Bo-  
janowski, P. DINOv2: Learning robust visual features  
without supervision. *Transactions on Machine Learn-*  
*ing Research*, 2024. ISSN 2835-8856. URL [https://](https://openreview.net/forum?id=a68SUt6zFt)  
[openreview.net/forum?id=a68SUt6zFt](https://openreview.net/forum?id=a68SUt6zFt). Fea-  
tured Certification.
- Pham, B., Raya, G., Negri, M., Zaki, M. J., Ambrogioni,  
L., and Krotov, D. Memorization to generalization:  
Emergence of diffusion models from associative memory  
networks. In *New Frontiers in Associative Memories*,  
2025. URL [https://openreview.net/forum?](https://openreview.net/forum?id=IWZnhP3YgK)  
[id=IWZnhP3YgK](https://openreview.net/forum?id=IWZnhP3YgK).
- Pizzi, E., Roy, S. D., Ravindra, S. N., Goyal, P., and Douze,  
M. A self-supervised descriptor for image copy detection.  
In *Proceedings of the IEEE/CVF Conference on Com-*  
*puter Vision and Pattern Recognition*, pp. 14532–14542,  
2022.
- Poole, B., Jain, A., Barron, J. T., and Mildenhall, B. Dream-  
fusion: Text-to-3d using 2d diffusion. In *The Eleventh*  
*International Conference on Learning Representations*,  
2023. URL [https://openreview.net/forum?](https://openreview.net/forum?id=FjNys5c7VyY)  
[id=FjNys5c7VyY](https://openreview.net/forum?id=FjNys5c7VyY).
- Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G.,  
Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J.,  
et al. Learning transferable visual models from natural  
language supervision. In *International Conference on*  
*Machine Learning*, pp. 8748–8763. PmLR, 2021.
- Ronneberger, O., Fischer, P., and Brox, T. U-net: Con-  
volutional networks for biomedical image segmentation.  
In *Medical Image Computing and Computer-Assisted*  
*Intervention–MICCAI 2015: 18th International Confer-*  
*ence, Munich, Germany, October 5-9, 2015, proceedings,*  
*part III 18*, pp. 234–241. Springer, 2015.
- Rudin, W. Principles of mathematical analysis. 2021.
- Saglietti, L., Mannelli, S., and Saxe, A. An analytical  
theory of curriculum learning in teacher-student networks.  
In *Advances in Neural Information Processing Systems*,  
volume 35, pp. 21113–21127, 2022.

- 660 Sajjadi, M. S., Bachem, O., Lucic, M., Bousquet, O., and  
661 Gelly, S. Assessing generative models via precision and  
662 recall. In *Advances in Neural Information Processing*  
663 *Systems*, volume 31, 2018.
- 664 Salimans, T. and Ho, J. Progressive distillation for fast  
665 sampling of diffusion models. In *International Confer-*  
666 *ence on Learning Representations*, 2022. URL <https://openreview.net/forum?id=TIIXIpzhoI>.
- 669 Salimans, T., Goodfellow, I., Zaremba, W., Cheung, V., Rad-  
670 ford, A., and Chen, X. Improved techniques for training  
671 gans. In *Advances in Neural Information Processing*  
672 *Systems*, volume 29, 2016.
- 674 Sauer, A., Lorenz, D., Blattmann, A., and Rombach, R. Ad-  
675 versarial diffusion distillation. In *European Conference*  
676 *on Computer Vision*, pp. 87–103. Springer, 2024.
- 678 Schuhmann, C., Beaumont, R., Vencu, R., Gordon, C.,  
679 Wightman, R., Cherti, M., Coombes, T., Katta, A., Mullis,  
680 C., Wortsman, M., et al. Laion-5b: An open large-scale  
681 dataset for training next generation image-text models.  
682 *Advances in neural information processing systems*, 35:  
683 25278–25294, 2022.
- 684 Skilling, J. The eigenvalues of mega-dimensional matrices.  
685 *Maximum Entropy and Bayesian Methods: Cambridge,*  
686 *England, 1988*, pp. 455–466, 1989.
- 688 Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and  
689 Goldstein, T. Diffusion art or digital forgery? investigat-  
690 ing data replication in diffusion models. In *Proceedings*  
691 *of the IEEE/CVF Conference on Computer Vision and*  
692 *Pattern Recognition*, pp. 6048–6058, 2023a.
- 694 Somepalli, G., Singla, V., Goldblum, M., Geiping, J., and  
695 Goldstein, T. Understanding and mitigating copying  
696 in diffusion models. *arXiv preprint arXiv:2305.20086*,  
697 2023b.
- 699 Song, J., Meng, C., and Ermon, S. Denoising diffu-  
700 sion implicit models. In *International Conference on*  
701 *Learning Representations*, 2021a. URL <https://openreview.net/forum?id=StlgIarCHLP>.
- 704 Song, Y., Durkan, C., Murray, I., and Ermon, S. Maxi-  
705 mum likelihood training of score-based diffusion models.  
706 In *Advances in Neural Information Processing Systems*,  
707 volume 34, pp. 1415–1428, 2021b.
- 709 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A.,  
710 Ermon, S., and Poole, B. Score-based generative mod-  
711 eling through stochastic differential equations. In *Inter-*  
712 *national Conference on Learning Representations*,  
713 2021c. URL <https://openreview.net/forum?id=PXTIG12RRHS>.
- 714 Song, Y., Dhariwal, P., Chen, M., and Sutskever, I. Consis-  
tency models. In *International Conference on Machine*  
*Learning*, 2023.
- Stein, G., Cresswell, J., Hosseinzadeh, R., Sui, Y., Ross, B.,  
Villicroze, V., Liu, Z., Caterini, A. L., Taylor, E., and  
Loaiza-Ganem, G. Exposing flaws of generative model  
evaluation metrics and their unfair treatment of diffusion  
models. In *Advances in Neural Information Processing*  
*Systems*, volume 36, 2024.
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna,  
Z. Rethinking the inception architecture for computer  
vision. In *Proceedings of the IEEE Conference on Com-*  
*puter Vision and Pattern Recognition*, pp. 2818–2826,  
2016.
- Tinaz, B., Fabian, Z., and Soltanolkotabi, M. Emergence  
and evolution of interpretable concepts in diffusion mod-  
els through the lens of sparse autoencoders. In *Sec-*  
*ond Workshop on Visual Concepts*, 2025. URL <https://openreview.net/forum?id=plCImt7QaW>.
- Vastola, J. Generalization through variance: how noise  
shapes inductive biases in diffusion models. In *The Thir-*  
*teenth International Conference on Learning Represent-*  
*tations*, 2025. URL <https://openreview.net/forum?id=7lUdo8Vuqa>.
- Wang, B. An analytical theory of power law spectral bias in  
the learning dynamics of diffusion models. *arXiv preprint*  
*arXiv:2503.03206*, 2025.
- Wang, B. and Vastola, J. J. The unreasonable effectiveness  
of gaussian score approximation for diffusion models and  
its applications. *arXiv preprint arXiv:2412.09726*, 2024.
- Wang, F.-Y., Huang, Z., Bergman, A., Shen, D., Gao, P.,  
Lingelbach, M., Sun, K., Bian, W., Song, G., Liu, Y.,  
et al. Phased consistency models. *Advances in neural*  
*information processing systems*, 37:83951–84009, 2024a.
- Wang, P., Zhang, H., Zhang, Z., Chen, S., Ma, Y., and Qu, Q.  
Diffusion models learn low-dimensional distributions via  
subspace clustering. *arXiv preprint arXiv:2409.02426*,  
2024b.
- Xie, S., Xiao, Z., Kingma, D. P., Hou, T., Wu, Y. N., Murphy,  
K. P., Salimans, T., Poole, B., and Gao, R. EM distilla-  
tion for one-step diffusion models. In *The Thirty-eighth*  
*Annual Conference on Neural Information Processing*  
*Systems*, 2024. URL <https://openreview.net/forum?id=rafVvthuxD>.
- Xu, Y., Zhao, Y., Xiao, Z., and Hou, T. Ufogen: You forward  
once large scale text-to-image generation via diffusion  
gans. In *Proceedings of the IEEE/CVF Conference on*  
*Computer Vision and Pattern Recognition (CVPR)*, pp.  
8196–8206, 2024.

- 715 Yang, R., Wang, Z., Jiang, B., and Li, S. Leveraging drift to  
716 improve sample complexity of variance exploding diffusion  
717 models. In *Advances in Neural Information Processing*  
718 *Systems*, volume 37, pp. 107662–107702, 2024.
- 719 Yang, R., Li, Y., Jiang, B., Chen, C., and Li, S. Multi-  
720 subspace multi-modal modeling for diffusion models:  
721 Estimation, convergence and mixture of experts. In  
722 *The Fourteenth International Conference on Learning*  
723 *Representations*, 2026. URL [https://openreview.](https://openreview.net/forum?id=MPWIM6rxxU)  
724 [net/forum?id=MPWIM6rxxU](https://openreview.net/forum?id=MPWIM6rxxU).
- 725 Yang, Z., Yu, Y., You, C., Steinhardt, J., and Ma, Y. Rethink-  
726 ing bias-variance trade-off for generalization of neural  
727 networks. In *International Conference on Machine Learn-*  
728 *ing*, pp. 10767–10777. PMLR, 2020.
- 729 Ye, Z., Zhu, Q., Tao, M., and Chen, M. Provable separations  
730 between memorization and generalization in diffusion  
731 models. In *The Fourteenth International Conference*  
732 *on Learning Representations*, 2026. URL [https://](https://openreview.net/forum?id=42gfTZzyvV)  
733 [openreview.net/forum?id=42gfTZzyvV](https://openreview.net/forum?id=42gfTZzyvV).
- 734 Yin, T., Gharbi, M., Park, T., Zhang, R., Shechtman, E., Du-  
735 rand, F., and Freeman, B. Improved distribution matching  
736 distillation for fast image synthesis. *Advances in neural*  
737 *information processing systems*, 37:47455–47487, 2024a.
- 738 Yin, T., Gharbi, M., Zhang, R., Shechtman, E., Durand, F.,  
739 Freeman, W. T., and Park, T. One-step diffusion with  
740 distribution matching distillation. In *Proceedings of the*  
741 *IEEE/CVF conference on computer vision and pattern*  
742 *recognition*, pp. 6613–6623, 2024b.
- 743 Yoon, T., Choi, J. Y., Kwon, S., and Ryu, E. K. Diffusion  
744 probabilistic models generalize when they fail to memo-  
745 rize. In *ICML 2023 Workshop on Structured Probabilistic*  
746 *Inference & Generative Modeling*, 2023.
- 747 Zhang, C., Bengio, S., Hardt, M., Recht, B., and Vinyals,  
748 O. Understanding deep learning requires rethinking gen-  
749 eralization. In *International Conference on Learning*  
750 *Representations*, 2017. URL [https://openreview.](https://openreview.net/forum?id=Sy8gdB9xx)  
751 [net/forum?id=Sy8gdB9xx](https://openreview.net/forum?id=Sy8gdB9xx).
- 752 Zhang, H., Zhou, J., Lu, Y., Guo, M., Wang, P., Shen, L., and  
753 Qu, Q. The emergence of reproducibility and consistency  
754 in diffusion models. In *Proceedings of the 41st Interna-*  
755 *tional Conference on Machine Learning*, volume 235 of  
756 *Proceedings of Machine Learning Research*, pp. 60558–  
757 60590. PMLR, 2024. URL [https://proceedings.](https://proceedings.mlr.press/v235/zhang24cn.html)  
758 [mlr.press/v235/zhang24cn.html](https://proceedings.mlr.press/v235/zhang24cn.html).
- 759 Zhang, Y., Chen, T., Wang, Z., Wu, Y. N., Zhou, M., and  
760 Leong, O. Restoration score distillation: From corrupted  
761 diffusion pretraining to one-step high-quality generation.  
762 *arXiv preprint arXiv:2505.13377*, 2025.
- 763 Zhang, Z., Li, X., Li, X., Shi, L., Wu, M., Tao, M., and  
764 Qu, Q. Generalization of diffusion models arises with  
765 a balanced representation space. In *The Fourteenth*  
766 *International Conference on Learning Representations*,  
767 2026. URL [https://openreview.net/forum?](https://openreview.net/forum?id=57TheGgNAN)  
768 [id=57TheGgNAN](https://openreview.net/forum?id=57TheGgNAN).
- 769 Zhu, Z., Locatello, F., and Cevher, V. Sample complexity  
770 bounds for score-matching: Causal discovery and gen-  
771 erative modeling. In *Advances in Neural Information*  
772 *Processing Systems*, volume 36, pp. 3325–3337, 2023.

The appendix is organized as follows. We first discuss related work in Section A. Next, we provide detailed proofs for Section 2 in Section B. Experimental settings and additional discussions for Section 3 and Section 4 are presented in Section C. We then offer further discussion related to  $\mathcal{E}_{\text{mem}}$  in Section D. Finally, ablation studies for PFD are included in Section E.

## A. Related Works

In this section, we briefly review related work on generalization metrics for diffusion models, discuss diffusion model generalizability, and cover the fundamentals of training diffusion models.

### A.1. Generalization Metrics for Diffusion Models

Generalization metrics quantify the distance between the learned distribution and the underlying data distribution in diffusion models. To measure this distributional gap, theoretical works commonly employ metrics such as Kullback-Leibler (KL) divergence (Chen et al., 2023e; Nie et al., 2024; Li et al., 2023), total variation (TV) (Chen et al., 2023c; Li et al., 2024b;a; Yang et al., 2024; Li et al., 2024c; Liang et al., 2025), and Wasserstein distance (Gao & Zhu, 2025; Bortoli, 2022; Chen et al., 2023a; Gao et al., 2025b). However, these metrics are practically inefficient for diffusion models. Practical metrics focus on various perspective, including negative log-likelihood (NLL) (Song et al., 2021c), image generation quality: Fréchet inception distance (FID) (Heusel et al., 2017), inception score (IS) (Salimans et al., 2016),  $FD_{\text{dinov2}}$  (Stein et al., 2024), maximum mean discrepancy (MMD) (Bińkowski et al., 2018), CLIP maximum mean discrepancy (CMMD) (Jayasumana et al., 2024); alignment:  $CLIP_{\text{score}}$  (Hessel et al., 2021), and  $\text{precision}$ ,  $\text{recall}$  (Sajjadi et al., 2018; Kynkäänniemi et al., 2019). However, these practical metrics are not explicitly designed to evaluate the generalizability of diffusion models. Thus, there is a need for a generalization metric that are both theoretical grounded and practically efficient for diffusion models. To address this gap, we propose PFD, a novel generalization metric that is theoretically proven to be a valid distributional distance and can be efficiently approximated by its empirical version using a polynomial number of samples. In practice, PFD requires fewer samples for estimation and is the only existing metric that explicitly quantifies generalization in diffusion models.

### A.2. Diffusion Model Generalizability

Recent works have shown that diffusion models transition from memorization to generalization as the number of training samples increases (Yoon et al., 2023; Zhang et al., 2024). With sufficient data, models trained with different architectures, loss functions, and even disjoint datasets can reproduce each other’s outputs, indicating a strong convergence toward the underlying data distribution (Zhang et al., 2024; Kadkhodaie et al., 2024). To explain this strong generalization, (Kadkhodaie et al., 2024) attribute it to the emergence of a geometric-adaptive harmonic basis, while others argue that generalization arises from interpolation across the data manifold (Aithal et al., 2024; Chen, 2025). In contrast, (Vastola, 2025) explain generalization through inductive biases in the noise, whereas (Zhang et al., 2026) attribute it to the emergence of a balanced representation space. Theoretical insights by (Li et al., 2023) provide generalization bounds using KL-divergence under simplified models, (Ye et al., 2026) establish a dual theoretical separation explaining memorization and generalization. More recent efforts focus on characterizing the learned noise-to-image mapping for generalized diffusion models, either through Gaussian parameterizations (Li et al., 2024d; Wang & Vastola, 2024), mixture of low rank Gaussian parameterizations (Wang et al., 2024b), multi-subspace multi-modal modeling (Yang et al., 2026) or patch-wise optimal score functions (Niedoba et al., 2024; Kamb & Ganguli, 2024), and analyses based on statistical properties of image datasets (Lukoianov et al., 2025). In parallel, generalization has also been studied through the lens of associative memory (Pham et al., 2025; Ambrogioni, 2024) and sparse autoencoders (Tinaz et al., 2025). However, despite these theoretical analyses and qualitative insights, prior work lacks a quantitative framework for measuring generalizability. In this paper, we propose PFD, a metric that enables such quantitative evaluation. Using this measure, we uncover further insights into the generalization behavior of diffusion models, as discussed in Section 4.

### A.3. Diffusion Model Distillation

Recent work on diffusion distillation has demonstrated substantial efficiency gains through both lightweight student architectures (Chen et al., 2025; Kim et al., 2024; Hu et al., 2026) and reduced sampling trajectories (Yin et al., 2024b;a; Salimans & Ho, 2022; Song et al., 2023; Meng et al., 2023; Sauer et al., 2024; Xie et al., 2024). Representative approaches include progressive distillation (Salimans & Ho, 2022; Berthelot et al., 2023), adversarial distillation (Jolicoeur-Martineau

et al., 2021; Xu et al., 2024; Sauer et al., 2024), and consistency-based distillation methods (Song et al., 2023; Wang et al., 2024a). Beyond architectural and sampling efficiency, Meng et al. (Meng et al., 2023) distill the classifier-free guidance score (Ho & Salimans, 2022) to further accelerate guidance at inference time. Moreover, diffusion distillation has been successfully applied beyond efficiency improvements, including text-to-3D generation (Poole et al., 2023), learning from corrupted data distributions (Zhang et al., 2025), and even achieving unexpected gains in generation quality (Ma et al., 2025). In this paper, we adopt the most straightforward distillation paradigm, where the student is trained on data generated by a teacher model. Nevertheless, PFD is orthogonal to the choice of distillation strategy and can be naturally extended to more advanced techniques, such as progressive, adversarial, or consistency-based distillation frameworks, offering a flexible foundation for future extensions.

#### A.4. Training Diffusion Models

To enable sampling via the PF-ODE (1), we train a neural network  $s_\theta(\mathbf{x}_t, t)$  to approximate the score function  $\nabla \log p_t(\mathbf{x}_t)$  using denoising score matching loss (Song et al., 2021c):

$$\min_{\theta} \ell(\theta) = \frac{1}{N} \sum_{i=1}^N \int_0^T \lambda_t \mathbb{E}_{\epsilon \sim \mathcal{N}(0, T^2 \mathbf{I}_n)} \left[ \left\| s_\theta(\mathbf{x}^{(i)} + t\epsilon, t) + \epsilon/t \right\|_2^2 \right] dt, \quad (10)$$

$\lambda_t$  denotes a scalar weight for the loss at  $t$ . Given the learned score function, the corresponding noise-to-image mapping is:

$$\Phi_{p_\theta}(\mathbf{x}_T) = \mathbf{x}_T - \int_T^0 t s_\theta(\mathbf{x}_t, t) dt. \quad (11)$$

Although alternative training objectives exist, such as predicting noise  $\mathbf{x}_T$  (Ho et al., 2020), clean image  $\mathbf{x}_0$  (Karras et al., 2022a), rectified flow  $\mathbf{x}_T - \mathbf{x}_0$  (Liu et al., 2023) or other linear combinations of  $\mathbf{x}_0$  and  $\mathbf{x}_T$  (Salimans & Ho, 2022), prior works (Luo, 2022; Gao et al., 2025a) have shown that it is still possible to recover an approximate score function  $s_\theta(\mathbf{x}_t, t)$  from these methods.

## B. Proof in Section 2

*Proof of Theorem 1.* It is trivial to show  $\text{PFD}(p, q) > 0$  for any  $p \neq q$  and  $\text{PFD}(p, q) = \text{PFD}(q, p)$ , and thus we omit the proof.

- Proof of  $p = q \Leftrightarrow \text{PFD}(p, q) = 0$  :

– ( $\Rightarrow$ ) If  $p = q$ ,  $\nabla \log p_t(\mathbf{x}_t) = \nabla \log q_t(\mathbf{x}_t)$ , thus:

$$d\mathbf{x}_t = -t (\nabla \log p_t(\mathbf{x}_t) - \nabla \log q_t(\mathbf{x}_t)) dt = 0 \quad (12)$$

Thus,  $\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)$  is the solution of the ODE function Equation (12) with initial  $\mathbf{x}_T = \mathbf{0}$ . Thus  $\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T) = \mathbf{0}$  for all  $\mathbf{x}_T$ . Thus  $\text{PFD}(p, q) = 0$

– ( $\Leftarrow$ ) If  $\text{PFD}(p, q) = 0$  and  $\Phi_p, \Phi_q$  are continuous function w.r.t  $\mathbf{x}_T$ , then we have  $\Phi_p(\mathbf{x}_T) = \Phi_q(\mathbf{x}_T)$  for all  $\mathbf{x}_T$ . If  $\mathbf{x}_0 = \Phi(\mathbf{x}_T)$ , from the transformation of probability identities, we have:

$$p(\mathbf{x}_0) = \frac{\partial}{\partial [\mathbf{x}_0]_1} \cdots \frac{\partial}{\partial [\mathbf{x}_0]_n} \int_{\{\epsilon | \Phi(\epsilon) \leq \mathbf{x}_0\}} p_{\mathcal{N}}(\epsilon) d^n \epsilon, \quad (13)$$

where  $[\mathbf{x}_0]_i$  denotes the  $i$ -th element of  $\mathbf{x}_0$ ,  $\mathbf{f}(\epsilon) \leq \mathbf{x}_0$  denotes the element wise less or equal.  $p_{\mathcal{N}}(\cdot)$  is the probability

density function (PDF) of Gaussian distribution  $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I}_n)$ . Thus, for all  $\mathbf{x}_0$  we have:

$$\begin{aligned}
 p(\mathbf{x}_0) - q(\mathbf{x}_0) &= \frac{\partial}{\partial[\mathbf{x}_0]_1} \cdots \frac{\partial}{\partial[\mathbf{x}_0]_n} \int_{\{\epsilon | \Phi_p(\epsilon) \leq \mathbf{x}_0\}} p_{\mathcal{N}}(\epsilon) d^n \epsilon \\
 &\quad - \frac{\partial}{\partial[\mathbf{x}_0]_1} \cdots \frac{\partial}{\partial[\mathbf{x}_0]_n} \int_{\{\epsilon | \Phi_q(\epsilon) \leq \mathbf{x}_0\}} p_{\mathcal{N}}(\epsilon) d^n \epsilon, \\
 &= \frac{\partial}{\partial[\mathbf{x}_0]_1} \cdots \frac{\partial}{\partial[\mathbf{x}_0]_n} \int_{\{\epsilon | \Phi_p(\epsilon) \leq \mathbf{x}_0\}} p_{\mathcal{N}}(\epsilon) d^n \epsilon \\
 &\quad - \frac{\partial}{\partial[\mathbf{x}_0]_1} \cdots \frac{\partial}{\partial[\mathbf{x}_0]_n} \int_{\{\epsilon | \Phi_p(\epsilon) \leq \mathbf{x}_0\}} p_{\mathcal{N}}(\epsilon) d^n \epsilon, \\
 &= 0,
 \end{aligned} \tag{14}$$

so  $p = q$ .

- Proof of  $\text{PFD}(p, q) \leq \text{PFD}(p, p') + \text{PFD}(p', q)$ :

$$\begin{aligned}
 &\text{PFD}(p, q) \\
 &= \left( \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(0, T^2 \mathbf{I})} \left[ \|\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)\|_2^2 \right] \right)^{1/2} \\
 &\leq \left( \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(0, T^2 \mathbf{I})} \left[ \left( \|\Phi_p(\mathbf{x}_T) - \Phi_{p'}(\mathbf{x}_T)\|_2 + \|\Phi_{p'}(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)\|_2 \right)^2 \right] \right)^{1/2} \\
 &\leq \left( \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(0, T^2 \mathbf{I})} \left[ \|\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)\|_2^2 \right] \right)^{1/2} \\
 &\quad + \left( \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(0, T^2 \mathbf{I})} \left[ \|\Phi_p(\mathbf{x}_T) - \Phi_{p'}(\mathbf{x}_T)\|_2^2 \right] \right)^{1/2} \\
 &= \text{PFD}(p, p') + \text{PFD}(p', q)
 \end{aligned} \tag{15}$$

□

**Lemma 1.** Under Assumption 1, for all  $\mathbf{x}_T \in \mathcal{N}(\mathbf{0}, T^2 \mathbf{I}_n)$ , as  $T \rightarrow \infty$ , we have:

$$\|\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)\|_2 \leq \exp\left(\frac{LT_\xi^2}{2}\right) \xi + \frac{\epsilon}{L} \left( \exp\left(\frac{LT_\xi^2}{2}\right) - 1 \right), \tag{16}$$

where  $\xi$  is a numerical constant and a finite timestep  $T_\xi$  depending only on  $\xi$ .

*Proof of Lemma 1.* Let  $\phi_t, t \in [0, T]$  denotes the ODE trajectory:

$$\begin{aligned}
 \phi_t &= \mathbf{x}_t^p - \mathbf{x}_t^q, \\
 \mathbf{x}_t^p &= \mathbf{x}_T - \int_T^t \tau \nabla_{\mathbf{x}} \log p_\tau(\mathbf{x}_\tau^p) d\tau, \\
 \mathbf{x}_t^q &= \mathbf{x}_T - \int_T^t \tau \nabla_{\mathbf{x}} \log q_\tau(\mathbf{x}_\tau^q) d\tau,
 \end{aligned} \tag{17}$$

From the definition,  $\phi_0 = \Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)$ . Because  $\lim_{T \rightarrow \infty} \phi_t = \mathbf{x}_T - \mathbf{x}_T = \mathbf{0}$ , from the  $\epsilon - \delta$  definition of the limit, given  $\mathbf{x}_T$ , and a constant  $\xi$ , there exists a finite  $T_\xi$  related to  $\xi$  such that:

$$\|\phi_t\|_2 \leq \xi \quad \text{for all } t \geq T_\xi. \tag{18}$$

As  $t \leq T_\xi$ , we have:

$$\begin{aligned}
 \frac{d\phi_t}{dt} &= -t (\nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t^p) - \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t^q)), \\
 \|\phi_{T_\xi}\|_2 &\leq \xi.
 \end{aligned} \tag{19}$$

Apply Assumption 1 to Equation (19), we could obtain the following integral inequality w.r.t  $\|\phi_t\|_2$ :

$$\begin{aligned} \frac{d\|\phi_t\|_2}{dt} &\leq \left\| \frac{d\phi_t}{dt} \right\|_2 \\ &\leq t \|\nabla_{\mathbf{x}} \log p_t(\mathbf{x}_t^p) - \nabla_{\mathbf{x}} \log q_t(\mathbf{x}_t^q)\|_2 \\ &\leq t(\epsilon + L\|\phi_t\|_2), \\ \|\phi_{T_\xi}\|_2 &\leq \xi, \quad 0 \leq t \leq T_\xi, \end{aligned} \quad (20)$$

where the first inequality comes from the fact that  $\frac{d\|\phi_t\|_2}{dt} \leq \left\| \frac{d\phi_t}{dt} \right\|_2$ . From Grönwall's inequality (Coddington & Levinson, 1955), we could solve  $\|\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)\|_2 = \|\phi_0\|_2 \leq \exp\left(\frac{LT_\xi^2}{2}\right)\xi + \frac{\epsilon}{L} \left( \exp\left(\frac{LT_\xi^2}{2}\right) - 1 \right)$ .  $\square$

*Proof of Theorem 2.* Let  $\mathbf{X} := \|\Phi_p(\mathbf{x}_T) - \Phi_q(\mathbf{x}_T)\|_2^2$ . From Lemma 1,

$$0 \leq \mathbf{X} \leq \kappa^2(L, \epsilon),$$

with  $\kappa(L, \epsilon) := \exp\left(\frac{LT_\xi^2}{2}\right)\xi + \frac{\epsilon}{L} \left( \exp\left(\frac{LT_\xi^2}{2}\right) - 1 \right)$ . From Hoeffding's inequality (Hoeffding, 1994), we have:

$$\mathbb{P} \left( \left| \mathbb{E}[\mathbf{X}] - \frac{1}{M} \sum_{i=1}^M \mathbf{X}_i \right| \geq \gamma \right) \leq 2 \exp \left( -\frac{2M\gamma^2}{\kappa^4(L, \epsilon)} \right), \quad (21)$$

with  $M$  samples to achieve  $\gamma$  accuracy. Thus, we could guarantee  $\mathbb{P} \left( \left| \mathbb{E}[\mathbf{X}] - \frac{1}{M} \sum_{i=1}^M \mathbf{X}_i \right| \leq \gamma \right)$  with probability  $\eta$ , when:

$$M \geq \frac{\kappa^4(L, \epsilon)}{2\gamma^2} \log \frac{2}{\eta}. \quad (22)$$

Because

$$\left| \text{PFD}(p, q) - \hat{\text{PFD}}(p, q) \right| = \left| \sqrt{\mathbb{E}[\mathbf{X}]} - \sqrt{\frac{1}{M} \sum_{i=1}^M \mathbf{X}_i} \right| \quad (23)$$

$$\leq \sqrt{\left| \mathbb{E}[\mathbf{X}] - \frac{1}{M} \sum_{i=1}^M \mathbf{X}_i \right|}. \quad (24)$$

We could guarantee that  $\mathbb{P} \left( \left| \text{PFD}(p, q) - \hat{\text{PFD}}(p, q) \right| \leq \gamma \right)$  with probability  $\eta$ , when:

$$M \geq \frac{\kappa^4(L, \epsilon)}{2\gamma^4} \log \frac{2}{\eta}. \quad (25)$$

$\square$

**Example 1.** The Wasserstein-2 distance  $W_2(\cdot, \cdot)$  is the lower bound of the probability flow distance, i.e.,

$$W_2(p, q) \leq \text{PFD}(p, q), \quad (26)$$

Specifically, let  $p$  and  $q$  be multivariate Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ ,  $\mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ , respectively, where  $\boldsymbol{\mu}_1, \boldsymbol{\mu}_2 \in \mathbb{R}^n$  and  $\boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2 \in \mathbb{R}^{n \times n}$ . The PFD is given by

$$\text{PFD}(p, q) = \left( \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2 + \left\| \boldsymbol{\Sigma}_1^{1/2} - \boldsymbol{\Sigma}_2^{1/2} \right\|_F \right)^{1/2}, \quad (27)$$

under this case, the equality  $\text{PFD}(p, q) = W_2(p, q)$  holds when  $\boldsymbol{\Sigma}_1 \boldsymbol{\Sigma}_2 = \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1$ .

Proof of Example 1. Proof of  $W_2(p, q) \leq \text{PFD}(p, q)$ . From the definition of Wasserstein-2 distance:

$$W_2(p, q) = \inf_{\gamma \in \Gamma(p, q)} \left( \mathbb{E}_{(\mathbf{x}_p, \mathbf{x}_q) \sim \gamma} \|\mathbf{x}_p - \mathbf{x}_q\|_2^2 \right)^{1/2}, \quad (28)$$

where  $\Gamma(p, q)$  is the set of all couplings of  $p$  and  $q$ . As proofed by (Song et al., 2021c), the noise-to-image mapping  $\Phi_p$  and  $\Phi_q$  pushes the Gaussian distribution  $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I}_n)$  to the  $p$  and  $q$  distribution respectively. Thus we could find the coupling  $\gamma_{\text{PFD}} := (\Phi_p, \Phi_q)_{\#} \mathcal{N}(\mathbf{0}, T^2 \mathbf{I}_n)$ , i.e., the pushforward of  $\mathcal{N}(\mathbf{0}, T^2 \mathbf{I}_n)$  by  $(\Phi_p, \Phi_q)$ , such that

$$\text{PFD}(p, q) = \left( \mathbb{E}_{(\mathbf{x}_p, \mathbf{x}_q) \sim \gamma_{\text{PFD}}} \|\mathbf{x}_p - \mathbf{x}_q\|_2^2 \right)^{1/2} \geq W_2(p, q) \quad (29)$$

When distribution  $p(\mathbf{x})$  is Gaussian distribution  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$  with  $\boldsymbol{\mu} \in \mathbb{R}^n$  and  $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$ . Thus  $p_t(\mathbf{x})$  is  $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma} + \sigma_t^2 \mathbf{I}_n)$ , thus the score function could be calculated as,

$$\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) = (\boldsymbol{\Sigma} + t^2 \mathbf{I}_n)^{-1} (\boldsymbol{\mu} - \mathbf{x}). \quad (30)$$

By plugging in Equation (30) to Equation (2), we could obtain the ODE equation w.r.t  $\mathbf{x}$ :

$$d\mathbf{x} = -t (\boldsymbol{\Sigma} + t^2 \mathbf{I}_n)^{-1} (\boldsymbol{\mu} - \mathbf{x}) dt, \quad (31)$$

The above ODE equation has a close form solution:

$$\mathbf{x}_t = \boldsymbol{\mu} + \mathbf{U} \text{diag} \left( \left[ \sqrt{\frac{\lambda_1 + t^2}{\lambda_1 + T^2}}, \dots, \sqrt{\frac{\lambda_n + t^2}{\lambda_n + T^2}} \right] \right) \mathbf{U}^\top (\mathbf{x}_T - \boldsymbol{\mu}) \quad (32)$$

where  $\mathbf{U}, \lambda_k, k \in [n]$  are singular value decomposition of  $\boldsymbol{\Sigma}$ ,  $\boldsymbol{\Sigma} = \mathbf{U} \text{diag}([\lambda_1, \dots, \lambda_n]) \mathbf{U}^\top$ .  $\text{diag}(\cdot)$  convert a vector in  $\mathbb{R}^n$  into diagonal matrix  $\mathbb{R}^{n \times n}$ , and  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I}_n)$ . Let  $\mathbf{x}_T = T\boldsymbol{\epsilon}$  with  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ . As  $t = 0$  and  $T \rightarrow \infty$ , we have:

$$\mathbf{x}_t = \left( \mathbf{I}_n - \mathbf{U} \text{diag} \left( \left[ \sqrt{\frac{\lambda_1 + t^2}{\lambda_1 + T^2}}, \dots, \sqrt{\frac{\lambda_n + t^2}{\lambda_n + T^2}} \right] \right) \mathbf{U}^\top \right) \boldsymbol{\mu}, \quad (33)$$

$$+ \mathbf{U} \text{diag} \left( \left[ T \sqrt{\frac{\lambda_1 + t^2}{\lambda_1 + T^2}}, \dots, T \sqrt{\frac{\lambda_n + t^2}{\lambda_n + T^2}} \right] \right) \mathbf{U}^\top \mathbf{x}_T, \quad (34)$$

$$= \boldsymbol{\mu} + \mathbf{U} \text{diag} \left( \left[ \sqrt{\lambda_1}, \dots, \sqrt{\lambda_n} \right] \right) \mathbf{U}^\top \mathbf{x}_T, \quad (35)$$

$$= \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \mathbf{x}_T = \Phi(\mathbf{x}_T). \quad (36)$$

Thus, plugging in Definition 1, we have:

$$\text{PFD}(p, q) = \left( \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})} \left[ \|\Phi_1(\mathbf{x}_T) - \Phi_2(\mathbf{x}_T)\|_2^2 \right] \right)^{1/2} \quad (37)$$

$$= \left( \mathbb{E}_{\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})} \left[ \left\| \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_1^{1/2} \mathbf{x}_T - \boldsymbol{\mu}_2 - \boldsymbol{\Sigma}_2^{1/2} \mathbf{x}_T \right\|_2^2 \right] \right)^{1/2} \quad (38)$$

$$= \left( \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2 + \left\| \boldsymbol{\Sigma}_1^{1/2} - \boldsymbol{\Sigma}_2^{1/2} \right\|_F^2 \right)^{1/2} \quad (39)$$

$$= \left( \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2 + \text{Tr} \left( \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2 - 2\boldsymbol{\Sigma}_1^{1/2} \boldsymbol{\Sigma}_2^{1/2} \right) \right)^{1/2} \quad (40)$$

From Wasserstein-2 distance for Gaussian distribution  $p, q$  has close form solution:

$$W_2(p, q) = \left( \|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2\|_2^2 + \text{Tr} \left( \boldsymbol{\Sigma}_1 + \boldsymbol{\Sigma}_2 - 2 \left( \boldsymbol{\Sigma}_1^{1/2} \boldsymbol{\Sigma}_2 \boldsymbol{\Sigma}_1^{1/2} \right)^{1/2} \right) \right)^{1/2}. \quad (41)$$

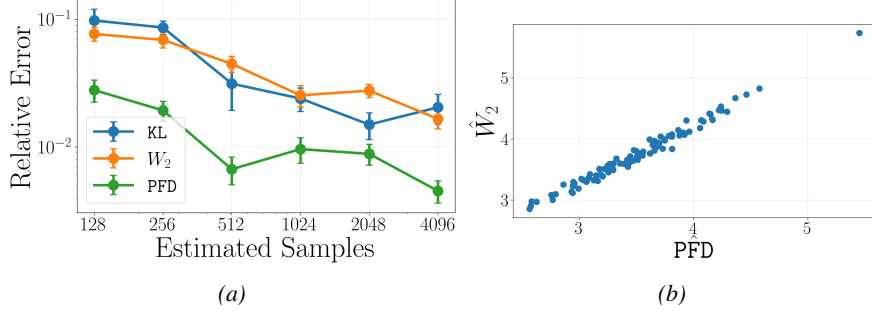


Figure 6. **Comparison of different metrics on synthetic datasets.** The figure illustrates (a) the sample efficiency of KL,  $W_2$ , and PFD under a Gaussian distribution, and (b) the correlation between  $W_2$  and PFD under a mixture of Gaussians.

From Lemma 2, we have  $W_2(p, q) \leq \text{PFD}(p, q)$ . And specifically,  $W_2(p, q) = \text{PFD}(p, q)$  when  $\Sigma_1 \Sigma_2 = \Sigma_2 \Sigma_1$ .  $\square$

**Lemma 2.** Given two positive semi-definite matrix  $\Sigma_1, \Sigma_2 \in \mathbb{R}^{n \times n}$ ,

$$0 \leq \text{Tr} \left( \Sigma_1^{1/2} \Sigma_2^{1/2} \right) \leq \text{Tr} \left( \left( \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} \right)^{1/2} \right). \quad (42)$$

*Proof of Lemma 2.* Because  $\Sigma_1, \Sigma_2$  are positive semi-definite matrix,  $\text{Tr} \left( \Sigma_1^{1/2} \Sigma_2^{1/2} \right) \geq 0$  and

$$\text{Tr} \left( \left( \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} \right)^{1/2} \right) = \text{Tr} \left( \sqrt{\left( \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} \right) \left( \Sigma_1^{1/2} \Sigma_2 \Sigma_1^{1/2} \right)^\top} \right) = \left\| \Sigma_1^{1/2} \Sigma_2^{1/2} \right\|_*, \quad (43)$$

where  $\|\cdot\|_*$  is the nuclear norm (trace norm). From trace norm inequality ((Bhatia, 2013) Chapter IV, Section 2), for a random matrix  $M$ ,  $\text{Tr}(M) \leq \|M\|_*$ . Thus, we have:

$$\text{Tr} \left( \Sigma_1^{1/2} \Sigma_2^{1/2} \right) \leq \left\| \Sigma_1^{1/2} \Sigma_2^{1/2} \right\|_*. \quad (44)$$

$\square$

## C. Experiments

In this section, we provide experimental details and additional discussion of the main results.

### C.1. Comparison of different metrics on synthetic datasets.

To support the claim that PFD is both sample-efficient and a meaningful distributional metric, we conduct numerical experiments on synthetic datasets, comparing PFD with the KL divergence (KL) and the 2-Wasserstein distance ( $W_2$ ). Specifically, in Figure 6a, we evaluate the sample efficiency of these methods under multivariate Gaussian distributions. We set the dimension of the Gaussian distributions to 5, randomly generate their means and variances, and repeat the experiment 10 times. By varying the number of samples used for estimation from 128 to 4096, denoted by  $M$ , we report the relative error  $|\mathcal{E} - \hat{\mathcal{E}}|/|\mathcal{E}|$ , where  $\mathcal{E}$  denotes the ground-truth value and  $\hat{\mathcal{E}}$  its empirical estimate. In Figure 6b, we examine the correlation between estimates of  $W_2$  and PFD on mixtures of Gaussian distributions. We consider mixtures of 5 Gaussian components, each with 5 dimensions. Both metrics are estimated using  $M = 4096$  samples, and the experiment is repeated 100 times to generate the plotted results.

As shown in Figure 6a, PFD exhibits significantly better sampling efficiency than both KL and  $W_2$ . With  $M = 4096$  samples, PFD achieves a relative error of approximately  $4 \times 10^{-2}$ , whereas KL and  $W_2$  incur errors on the order of  $2 \times 10^{-1}$ . Moreover, the computational complexity of PFD is  $O(M)$ , compared to  $O(M^2)$  for  $W_2$ . Consequently, under the same sampling budget, PFD is substantially more efficient in both estimation accuracy and computational cost.

As illustrated in Figure 6b, when measuring the distance between two mixture-of-Gaussian distributions, the estimated PFD exhibits a strong linear correlation with  $\hat{W}_2$  (with correlation coefficient 0.992). This result indicates that PFD

Table 2. U-Net architectures details.

Name	Dimensions for encoder blocks	Number of residual blocks	Number of parameters $ \theta $
U-Net-1	[32, 32, 32]	4	0.9M
U-Net-2	[64, 64, 64]	4	3.5M
U-Net-3	[96, 96, 96]	4	7.9M
U-Net-4	[128, 128, 128]	4	14.0M
U-Net-5	[80, 160, 160]	4	17.1M
U-Net-6	[160, 160, 160]	3	17.8M
U-Net-7	[160, 160, 160]	4	21.8M
U-Net-8	[192, 192, 192]	4	31.3M
U-Net-9	[224, 224, 224]	4	42.7M
U-Net-10	[256, 256, 256]	4	55.7M

captures meaningful distributional distance, which align with 2-Wasserstein distance, even for more complex, multimodal distributions.

### C.2. Network Architecture Details

In this subsection, we provide details of the U-Net architectures, as summarized in Table 2. The U-Net follows an encoder-decoder design, where the encoder comprises multiple encoder blocks. The column **Dimensions for encoder blocks** indicates the feature dimensions of each encoder block, while **Number of residual blocks** specifies how many residual blocks are used within each encoder block. The decoder is symmetric to the encoder. For further architectural details, please refer to (Karras et al., 2022b). By varying the encoder block dimensions and the number of residual blocks, we scale the U-Net model from 0.9M to 55.7M parameters.

### C.3. Evaluation Protocol

In this subsection, we provide details of the evaluation protocol introduced in Section 3, as well as the comparison between the synthetic dataset from the teacher model and the real dataset.

**Experiment settings for evaluation protocol.** The teacher model  $\theta_t$  and the student model  $\theta$  share a similar U-Net architecture (Ronneberger et al., 2015) with different numbers of parameters, as introduced in Section C.2. The teacher model, with UNet-10 architecture, is trained on the CIFAR-10 dataset (Krizhevsky et al., 2009) using the EDM noise scheduler (Karras et al., 2022a), with a batch size of 128 for 1,000 epochs. The student model <sup>1</sup> is trained using the variance-preserving (VP) noise scheduler (Ho et al., 2020), under the same training hyperparameters. We use one A40 GPU with 48 GB video random access memory (VRAM) for all experiments. We generated three subsets of initial noise  $\{\mathbf{x}_{\text{train},T}^{(i)}\}_{i=1}^N, \{\mathbf{x}_{\text{gen},T}^{(i)}\}_{i=1}^M, \{\mathbf{x}_{\text{test},T}^{(i)}\}_{i=1}^M \stackrel{\text{iid}}{\sim} \mathcal{N}(0, T^2 \mathbf{I}_n)$ . The training and test datasets are produced using the teacher model:

$$\mathcal{D} := \{\mathbf{x}_{\text{train}}^{(i)}\}_{i=1}^N = \{\Phi_{p_{\theta_t}}(\mathbf{x}_{\text{train},T}^{(i)})\}_{i=1}^N, \quad \mathcal{D}_{\text{test}} := \{\mathbf{x}_{\text{test}}^{(i)}\}_{i=1}^M = \{\Phi_{p_{\theta_t}}(\mathbf{x}_{\text{test},T}^{(i)})\}_{i=1}^M.$$

To evaluate the student model, we generate an evaluation dataset from itself:

$$\mathcal{D}_{\text{gen}} := \{\mathbf{x}_{\text{gen}}^{(i)}\}_{i=1}^M = \{\Phi_{p_{\theta}}(\mathbf{x}_{\text{gen},T}^{(i)})\}_{i=1}^M.$$

All samples are generated using the second-order Heun solver (Karras et al., 2022a) with 18 sampling steps. We vary the number of training samples  $N$  from  $2^6$  to  $2^{16}$  in powers of two.  $M$  is set to 50,000 for the experiments in Section C.4, and 10,000 for the rest.

**Experiment settings for validating the synthetic dataset with real real-world dataset.** We evaluate FID and  $\mathcal{E}_{\text{mem}}$  for diffusion models with UNet-4 architecture, trained separately on the synthetic dataset  $\mathcal{D}$  and CIFAR-10 training dataset. We keep the number of training dataset  $N$  the same for these two settings, ranging from  $2^6$  to  $2^{15}$ , with a power of 2. Then we

<sup>1</sup>The architecture of the student model varies across experiments and will be described in detail for each specific case.

## Probability Flow Distance

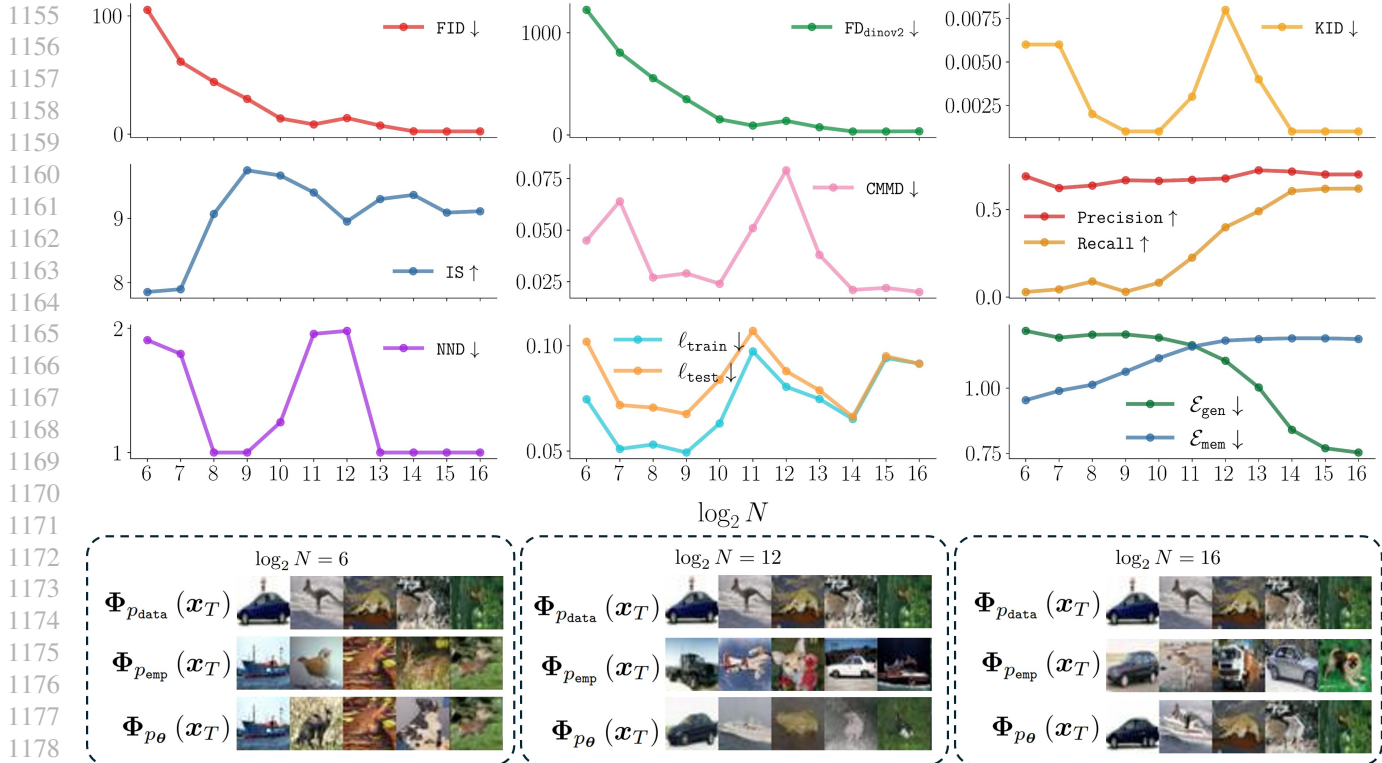


Figure 7. Comparison of practical metrics on the MtoG transition. The top figure plots multiple evaluation metrics as functions of  $\log_2 N$ . The bottom figure visualizes the generation under three numbers of training samples ( $2^6, 2^{12}, 2^{16}$ ). For each setting, the figure shows generations from the underlying distribution (top row), empirical data distribution (middle row), and the learned distribution from the diffusion model (bottom row). Each column corresponds to the same initial noise.

evaluate the FID between  $\mathcal{D}_{\text{gen}}$  and  $\mathcal{D}_{\text{test}}$  (CIFAR-10 test dataset) for the synthetic (real-world) setting, with  $M = 10000$ . To evaluate  $\mathcal{E}_{\text{mem}}$ , we use the initial noise  $\{\mathbf{x}_{\text{gen}}^{(i)}\}_{i=1}^M$ .

### C.4. Comparison with Practical Metrics for Generalization Evaluation

In this subsection, we expand upon the experiment presented in Section 3, which compares our proposed metric with practical metrics for evaluating generalization. We compare  $\mathcal{E}_{\text{gen}}$  and  $\mathcal{E}_{\text{mem}}$  with well-used generative model metrics, including FID,  $FD_{\text{DINOv2}}$ , KID, CMMD, Precision, Recall, NND, IS. We also including the training and testing loss  $\ell_{\text{train}}, \ell_{\text{test}}$  (Equation (10)) as comparison. We evaluating their ability in capturing the MtoG transition, under the evaluation protocol proposed in Section 3.

Metric	Dataset(s)
FID, $FD_{\text{DINOv2}}$ , KID, CMMD, Precision, Recall, NND	$\mathcal{D}_{\text{gen}}$ vs. $\mathcal{D}_{\text{test}}$
$FID_{\text{train}}, FD_{\text{DINOv2, train}}$	$\mathcal{D}$ vs. $\mathcal{D}_{\text{test}}$
IS	$\mathcal{D}_{\text{gen}}$
$\ell_{\text{train}}$	$\mathcal{D}$
$\ell_{\text{test}}$	$\mathcal{D}_{\text{test}}$
$\mathcal{E}_{\text{mem}}, \mathcal{E}_{\text{gen}}$	$\{\mathbf{x}_{\text{gen}, T}^{(i)}\}_{i=1}^M$

Table 3. Datasets used to evaluate each metric.

We use UNet-10 for the student model in this experiment. We summarized datasets used by these metrics in Table 3. Results are shown in Figure 7, summarized into one sentence, only  $\mathcal{E}_{\text{gen}}$  and  $\mathcal{E}_{\text{mem}}$  could quantitatively capture this transition. We

include detailed discussions below:

**Results discussions.** Figure 7 (bottom) is consistent with prior empirical observations (Zhang et al., 2024; Yoon et al., 2023): In the memorization regimes ( $N = 2^6$ ),  $p_\theta$  tends to memorize the empirical distribution  $p_{\text{emp}}$ , resulting in similar generation between  $\Phi_{p_{\text{emp}}}(\mathbf{x}_T)$  and  $\Phi_{p_\theta}(\mathbf{x}_T)$ ; in the transition regime ( $N = 2^{12}$ ), the model lacks sufficient capacity to memorize and the sample complexity is inadequate for generalization, leading to poor-quality generations  $\Phi_{p_\theta}(\mathbf{x}_T)$ ; in the generalization regimes ( $N = 2^{16}$ ),  $p_\theta$  captures the underlying distribution  $p_{\text{data}}$ , and the generations  $\Phi_{p_{\text{data}}}(\mathbf{x}_T)$  and  $\Phi_{p_\theta}(\mathbf{x}_T)$  are closely aligned.

As shown in Figure 7 (top), when  $N$  increases,  $\mathcal{E}_{\text{mem}}$  consistently increases and  $\mathcal{E}_{\text{gen}}$  consistently decreases. This aligns with our intuition: as sample complexity grows, models tend to generalize and memorize less. In contrast, all other metrics fail to capture this transition effectively. The reasons can be summarized as follows:

- **FID,  $\text{FD}_{\text{DINOv2}}$ , KID, IS, and CMMD are sensitive to generation quality.** Image quality metrics, including FID,  $\text{FD}_{\text{DINOv2}}$ , KID, IS, and CMMD, show degradation in performance at  $N = 2^{12}$ . This drop is primarily due to degraded visual quality in the generated samples, as visualize in Figure 7 (bottom-middle). However, at this sample complexity, the generated data still captures low-level features such as colors and structures from the underlying distribution. This is evident from the visual similarity between  $\Phi_{p_{\text{data}}}(\mathbf{x}_T)$  and  $\Phi_{p_\theta}(\mathbf{x}_T)$ , suggesting the model have some generalizability. In comparison, only  $\mathcal{E}_{\text{gen}}$  decreases consistently around  $N = 2^{12}$ , indicating it captures generalizability better than others despite visual degradation.
- **FID,  $\text{FD}_{\text{DINOv2}}$  and Recall are sensitive to diversity.** The monotonic trends for FID,  $\text{FD}_{\text{DINOv2}}$  and Recall are due to their sensitivity to the diversity of  $\mathcal{D}_{\text{gen}}$ , rather than their ability to measure generalizability. At small  $N$ , the model memorizes the training samples, resulting in  $\mathcal{D}_{\text{gen}}$  closely resembling  $\mathcal{D}$  and exhibiting significantly lower diversity than  $\mathcal{D}_{\text{test}}$ , since  $N \ll M$ . Under these conditions, FID,  $\text{FD}_{\text{DINOv2}}$  are large because they are biased towards the diversity of the evaluation samples (as proved in (Chong & Forsyth, 2020)). Meanwhile, Recall is low because the the support of  $\mathcal{D}_{\text{test}}$  is limited, reducing the probability that samples drawn from  $\mathcal{D}_{\text{gen}}$  lie within the support of  $\mathcal{D}_{\text{test}}$ . In contrast,  $\mathcal{E}_{\text{gen}}$  measures generalizability by directly quantifying the distance between the generation from the learned distribution and the underlying distribution and is less affected by the diversity of the generated samples.
- **NND and  $\ell$  fail to capture the generalizability.** The NND, originally designed for assessing the generalization of GANs, is sensitive to image quality and increases during the transition regime. Additionally, it produces identical values across a wide range of sample sizes (e.g.,  $N = 2^8, 2^9, 2^{13}, 2^{14}, 2^{15}, 2^{16}$ ), making it unreliable for evaluating generalization in diffusion models. Similarly, neither the training loss  $\ell_{\text{train}}$  nor the test loss  $\ell_{\text{test}}$  exhibits a consistent decreasing trend as  $N$  increases, indicating that these losses do not directly reflect either memorization or generalization. While the loss gap  $\ell_{\text{test}} - \ell_{\text{train}}$  does tend to decrease with larger  $N$ , it cannot serve as a robust generalization metric either. This is because even a randomly initialized model  $\theta$  can exhibit a small loss gap.

In conclusion,  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  are the only metrics that could capture the MtoG transition for diffusion models. They evaluate the generalization (memorization) by directly measuring the distance between the learned distribution by the diffusion model and the underlying (empirical) distribution. Unlike other metrics, they are less affected by the quality or diversity of the evaluating samples.

### C.5. Scaling Behaviors of the MtoG Transition

In this subsection, we provide detailed experimental settings for Section 4.1, along with additional experiments to further investigate the MtoG transition across more architectures (e.g., Transformer-based models (Bao et al., 2023)). We also investigate the scaling behavior of the MtoG transition under the DINOv2 descriptor.

**Experiment settings.** The detailed architectures of the student models, from U-Net-1 to U-Net-10, are provided in Section C.2, with model sizes ranging from 0.9M to 55.7M parameters. We scale up the architectures by increasing the dimensionality of the encoder blocks and the number of residual blocks. For the ImageNet experiments, we adopt the U-Net architectures proposed in (Karras et al., 2024), referred to as U-Net-11 and U-Net-12. These models contain 124.2M and 295.9M parameters, respectively.

**MtoG transition between U-Net and transformer architecture.** To further investigate the impact of network architecture, we compare the U-Net architecture with the transformer-based UViT (Bao et al., 2023). Specifically, we use the U-Net-9

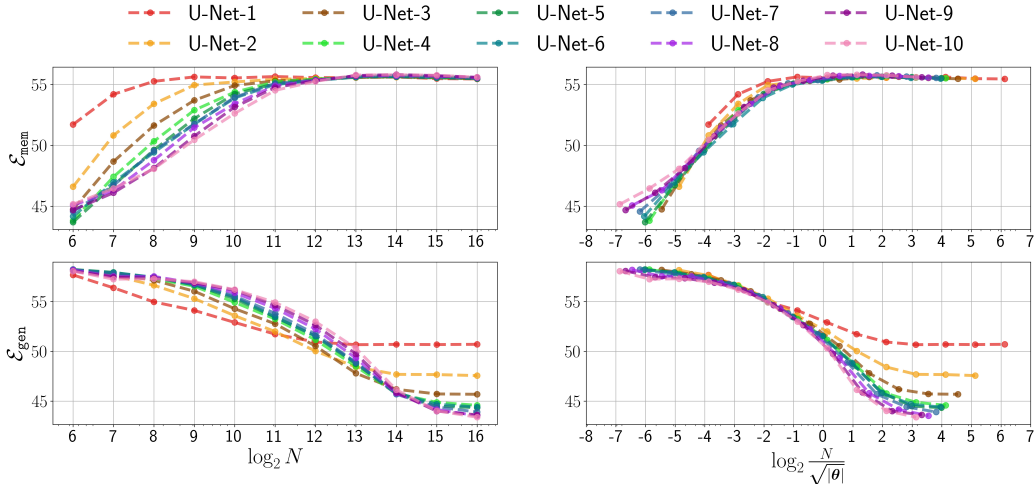


Figure 8. Scaling behavior in the MtoG transition under DINOv2 descriptor.  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  plotted against  $\log_2(N)$  for a range of U-Net architectures (U-Net-1 to U-Net-10). Right: the same metrics plotted against  $\log_2(N/\sqrt{|\theta|})$ , where  $|\theta|$  is the number of model parameters.

from Table 2, containing 42.7M parameters, and design the UViT model with comparable parameters of 44.2M. Both models are trained for 1000 epochs. Using the same experimental setup described in Section 4.1, we plot the MtoG transition curves for both U-Net and UViT, as shown in Figure 9.

As illustrated in Figure 9, with a similar number of parameters and the same training data sizes, UViT exhibits a higher  $\mathcal{E}_{\text{mem}}$  in the memorization regime ( $2^6 \leq N \leq 2^{10}$ ) and a higher  $\mathcal{E}_{\text{gen}}$  in the generalization regime ( $2^{11} \leq N \leq 2^{15}$ ), suggesting a lower model capacity compared to U-Net under these conditions. However, when provided with sufficient training data ( $N = 2^{16}$ ), UViT achieves a lower  $\mathcal{E}_{\text{gen}}$ , demonstrating better generalization performance. This observation is consistent with prior findings on transformer architectures in classification tasks: transformer-based models, lacking the inductive biases inherent to CNNs, tend to generalize poorly when trained on limited data (Dosovitskiy et al., 2021).

**Scaling behavior of the MtoG transition under the DINOv2 descriptor.** The scaling behavior under the DINOv2 descriptor is shown in Figure 8. Both  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  exhibit trends consistent with those observed under the SSCD descriptor (see Figure 2). The only difference is that, under the DINOv2 descriptor, models with varying parameter sizes show greater differentiation in the generalization regime compared to those under the SSCD descriptor. Further discussion on this can be found in the ablation study on image descriptors in Section E.2.

### C.6. Early Learning and Double Descent in Learning Dynamics

In this subsection, we build on the discussion from Section 4.2. In Figure 4, we evaluate  $\ell_{\text{train}}$  and  $\ell_{\text{test}}$  across the three training regimes. Notably, the gap  $\ell_{\text{test}} - \ell_{\text{train}}$  emerges as a practical heuristic for identifying the training regime: In the memorization regime, the gap increases steadily with training; In the transition regime, the gap remains near zero during early training (when generalization improves) and increases for further training (when generalization degrades); in the generalization regime, the gap remains close to zero throughout training. While  $\ell_{\text{test}} - \ell_{\text{train}}$  is not a strict measure of generalization, it proves to be a useful empirical indicator of training regimes for diffusion models. Practically, by setting aside a test dataset to estimate this gap, we can more effectively identify the training regime for diffusion models.

**Training dynamics of diffusion models under the DINOv2 descriptor.** The training dynamics under the DINOv2 descriptor are shown in Figure 10. Both  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  exhibit trends consistent with those observed under the SSCD descriptor for  $N = 64$  and  $N = 4096$  (see Figure 4). For  $N = 65536$ ,  $\mathcal{E}_{\text{gen}}$  still displays a double descent pattern under the DINOv2 descriptor; however, instead of a rise between the two drops, the curve remains relatively flat.

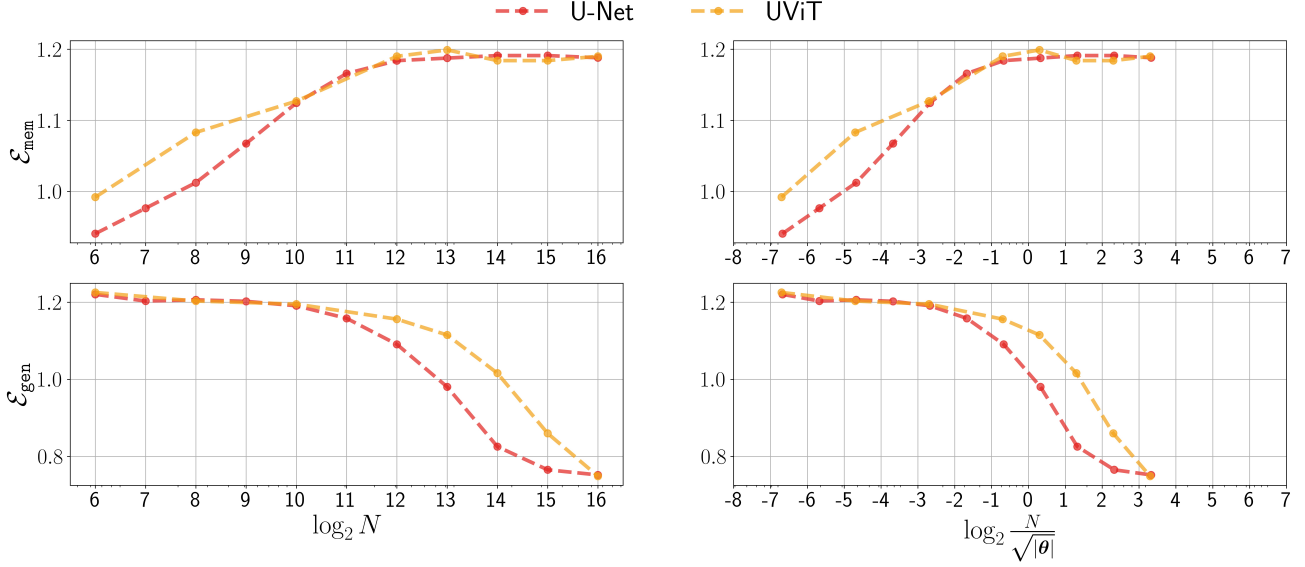


Figure 9. Comparison of scaling behavior between UNet and Transformer architectures in the MtoG transition.  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  plotted against  $\log_2(N)$  for U-Net architecture (U-Net-9) and UViT architecture. Right: the same metrics plotted against  $\log_2(N/\sqrt{|\theta|})$ , where  $|\theta|$  is the number of model parameters.

### C.7. Bias-Variance Decomposition of Generalization Error

To approximate  $\overline{\Psi \circ \Phi_{p_\theta}(\cdot)}$ , we independently sample two training datasets,  $\mathcal{D}_1$  and  $\mathcal{D}_2$ , for each specified number of training samples  $N$ . We then train two student models,  $\theta(\mathcal{D}_1)$  and  $\theta(\mathcal{D}_2)$ , using these datasets. The quantity  $\overline{\Psi \circ \Phi_{p_\theta}(\cdot)}$  is approximated as follows:

$$\overline{\Psi \circ \Phi_{p_\theta}(\cdot)} \approx \frac{1}{2}(\Psi \circ \Phi_{p_{\theta(\mathcal{D}_1)}}(\cdot) + \Psi \circ \Phi_{p_{\theta(\mathcal{D}_2)}}(\cdot)). \quad (45)$$

**Bias-Variance Decomposition of Generalization Error under the DINOv2 Descriptor.** The bias-variance decomposition under the DINOv2 descriptor is shown in Figure 11. Overall, the results are consistent with those observed under the SSCD descriptor, with two differences: (1) for  $N = 65536$ ,  $\mathcal{E}_{\text{gen}}$  does not exhibit a U-shaped curve under the DINOv2 descriptor; and (2)  $\mathcal{E}_{\text{bias}}$  for U-Net-1 and U-Net-2 does not decrease monotonically, instead, it first decreases and then increases.

### D. Further Discussions of $\mathcal{E}_{\text{mem}}$

In this section, we present the mathematical formulation for estimating  $\mathcal{E}_{\text{mem}}$  and compare it with the existing memorization metric.

**Empirically estimate  $\mathcal{E}_{\text{mem}}$ .** As described in Definition 1 and Definition 2, estimating  $\mathcal{E}_{\text{mem}}$  requires access to the mapping  $\Phi_{p_{\text{emp}}}(\cdot)$ . According to Equation (2), this mapping is determined by the score function of the empirical distribution, denoted as  $\nabla \log \hat{p}_t(\mathbf{x}_t)$ . Based on prior works (Karras et al., 2022a; Zhang et al., 2024; Gu et al., 2023), the score function of the empirical distribution has a closed-form expression:

$$\nabla \log \hat{p}_t(\mathbf{x}_t) = \frac{1}{T^2} \left( \frac{\mathbb{E}_{\mathbf{x} \sim p_{\text{emp}}}[\mathcal{N}(\mathbf{x}_t; \mathbf{x}, T^2 \mathbf{I}_n) \cdot \mathbf{x}]}{\mathbb{E}_{\mathbf{x} \sim p_{\text{emp}}}[\mathcal{N}(\mathbf{x}_t; \mathbf{x}, T^2 \mathbf{I}_n)]} - \mathbf{x}_t \right), \quad (46)$$

where  $p_{\text{emp}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \delta(\mathbf{x} - \mathbf{y}^{(i)})$  corresponds to the empirical distribution over the training dataset  $\mathbf{y}^{(i)}_{i=1}^N$ . This formulation allows us to numerically compute  $\nabla \log \hat{p}_t(\mathbf{x}_t)$  for any given  $t$ . Subsequently, we can use a numerical solver to estimate the integral in Equation (2), thereby enabling the estimation of  $\mathcal{E}_{\text{mem}}$ .

**Comparison between existing memorization metric and  $\mathcal{E}_{\text{mem}}$ .** Previous works (Yoon et al., 2023; Zhang et al., 2024)

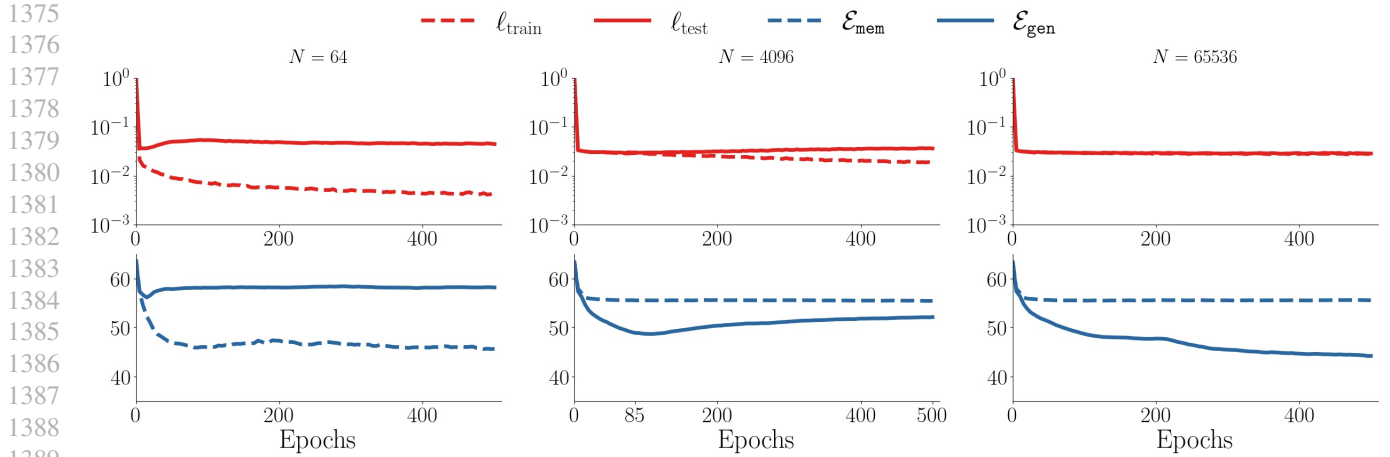


Figure 10. Training dynamics of diffusion models under DINOv2 descriptor in different regimes. The figure plots  $\mathcal{E}_{\text{mem}}, \mathcal{E}_{\text{gen}}, \ell_{\text{train}}, \ell_{\text{test}}$  over training epochs for different dataset sizes:  $N = 2^6$  (left),  $2^{12}$  (middle),  $2^{16}$  (right).

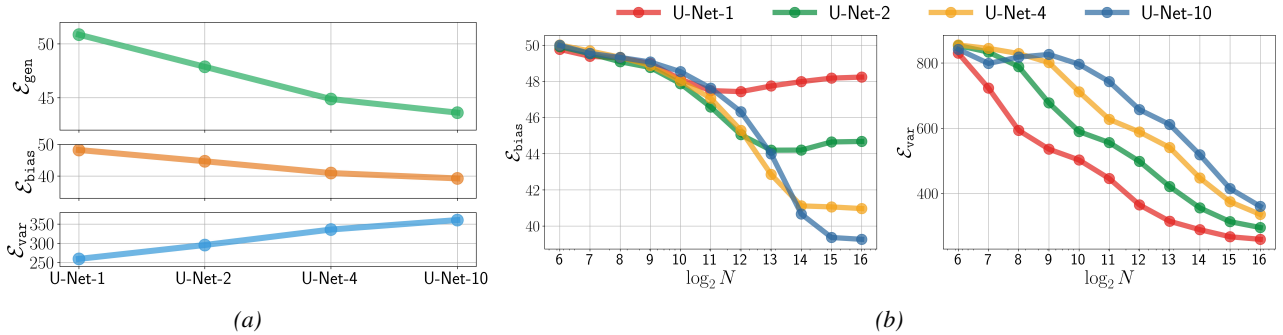


Figure 11. Bias-Variance Trade-off under DINOv2 descriptor. (a) plots the generalization error  $\mathcal{E}_{\text{gen}}$ , bias  $\mathcal{E}_{\text{bias}}$ , and variance  $\mathcal{E}_{\text{var}}$  across different network architectures with a fixed training sample size of  $N = 2^{16}$ . (b) shows  $\mathcal{E}_{\text{bias}}$  and  $\mathcal{E}_{\text{var}}$  as functions of the number of training samples  $N$  for various network architectures.

define memorization metric as:

$$\text{M Distance}(p_\theta) := \mathbb{E}_{\mathbf{x}_T} \left[ \min_{\mathbf{x} \sim p_{\text{emp}}} \|\Psi(\mathbf{x}) - \Psi \circ \Phi_{p_\theta}(\mathbf{x}_T)\|_2 \right], \quad (47)$$

A generated sample  $\Phi_{p_\theta}(\mathbf{x}_T)$  is a memorized sample if it is close enough to one of the sample  $\mathbf{x}$  from  $p_{\text{emp}}$ . It is easy to show that  $\mathcal{E}_{\text{mem}}$  is a more strict metric than M Distance, i.e. " $\mathcal{E}_{\text{mem}}(p^\theta) = 0$ " is a sufficient but not necessary condition for " $\text{M Distance}(p^\theta) = 0$ ". We propose  $\mathcal{E}_{\text{mem}}$  in order to unify the definition of memorization and generalization.

## E. Ablation Study

In this section, we present ablation studies on the evaluation protocol, examining the effects of different noise schedulers and sampling methods (Section E.1), image descriptors (Section E.2), sample sizes for evaluation (Section E.3), and teacher models (Section E.4).

### E.1. Sampling Methods

In this subsection, we present ablation studies on various noise schedulers and sampling strategies. Specifically, we evaluate the performance of the following methods: Variance Preserving (VP) (Song et al., 2021c), Variance Exploding (VE) (Song et al., 2021c), iDDPM (Nichol & Dhariwal, 2021) + DDIM (Song et al., 2021a), and EDM (Karras et al., 2022a). The specific form of  $f(t), g(t)$  used in each approach are detailed in Table 1 of (Karras et al., 2022a). Additionally, each method

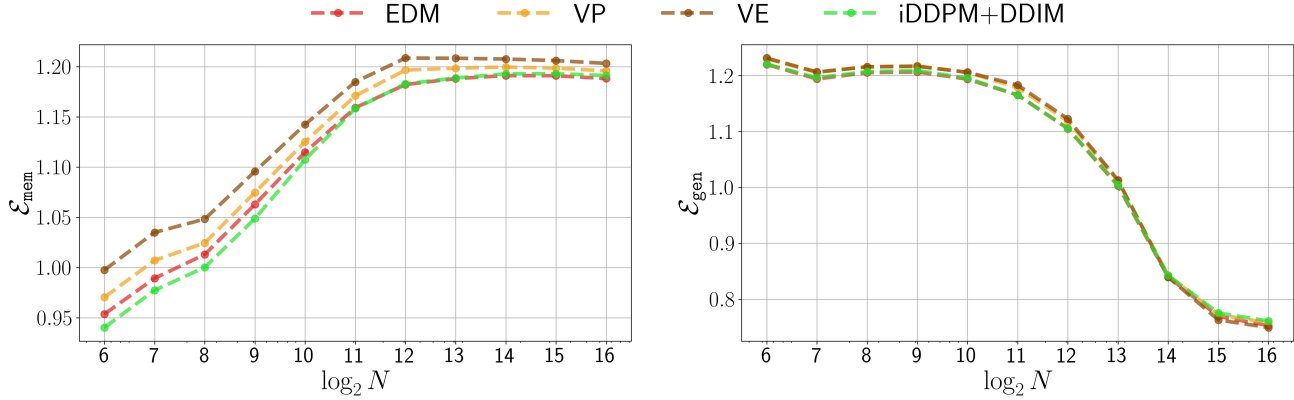


Figure 12. **Comparison of different sampling methods.**  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  plotted against  $\log_2(N)$  for different sampling methods, including: EDM, VP, VE, iDDPM+DDIM.

also differs in its choice of ODE solver and timestep discretization strategy. For sampling, we use 256 steps for VP, 1000 for VE, 100 for iDDPM + DDIM, and 18 for EDM. All experiments are conducted under the evaluation protocol described in Section 3, where we estimate the  $\mathcal{E}_{\text{gen}}$  under different training samples  $N$ . The student models use the UNet-10 architecture. During the ablation study, both the teacher and student models use the same sampling method<sup>2</sup> as specified above.

As shown in Figure 12, different samplers yield highly consistent results, demonstrating that PFD can be extended to various noise schedules, i.e., different choices of  $f(t)$  and  $g(t)$ .

## E.2. Image Descriptors

In this subsection, we present ablation studies on the image descriptor  $\Psi$  used in Equation (3). The descriptors evaluated include DINOv2 (Oquab et al., 2024), InceptionV3 (Szegedy et al., 2016), CLIP (Radford et al., 2021), SSCD (Pizzi et al., 2022), and the identity function. All experiments follow the evaluation protocol described in Section 3, where we estimate both  $\mathcal{E}_{\text{mem}}$  and  $\mathcal{E}_{\text{gen}}$  across varying training sample sizes  $N$  and different student model architectures: U-Net-1, U-Net-2, U-Net-4, and U-Net-10.

As shown in Figure 13, different feature embeddings reveal a consistent trend in the memorization-to-generalization (MtoG) transition across various U-Net architectures. With limited training samples, smaller models exhibit lower generalization scores. Conversely, with sufficient training data, larger models tend to have lower generalization scores. When comparing with  $\mathcal{E}_{\text{gen}}$  measured in pixel space (i.e., using the identity function as the descriptor), we observe that  $\mathcal{E}_{\text{gen}}$  values are nearly identical across diffusion architectures when  $N \geq 2^{15}$ . In this regime, all models have learned low-level image features such as color and structure; however, only the larger models capture high-level perceptual details. Because pixel-space measurements fail to reflect these high-level features, they yield similar  $\mathcal{E}_{\text{gen}}$  values regardless of model size. Therefore, it is better to evaluate  $\mathcal{E}_{\text{gen}}$  in a feature space, which better captures perceptual differences between models.

Different feature descriptors mainly differ in the generalization regime. Specifically,  $\mathcal{E}_{\text{gen}}$  varies the most across architectures when using the DINOv2 descriptor, and the least when using the SSCD descriptor. This is because each descriptor capture different aspects of the image. SSCD focuses on detecting duplicate content and is more sensitive to low-frequency features, while DINOv2 emphasizes perceptual quality and captures high-frequency features. Diffusion models with limited capacity tend to learn low-frequency information first, as it is more easier to learn (Wang, 2025). As a result, under the SSCD descriptor, different architectures show more similar  $\mathcal{E}_{\text{gen}}$  values, since they are all primarily capturing the same low-frequency information in the early training stages.

## E.3. Evaluation Sample Number

In this subsection, we present ablation studies on the number of samples  $M$  used by  $\hat{\text{PFD}}$  to approximate PFD, as defined in Equation (4). All experiments follow the evaluation protocol described in Section 3, where we estimate  $\mathcal{E}_{\text{gen}}$  across varying

<sup>2</sup>Note that the noise scheduler used for sampling could differ from that used during training.

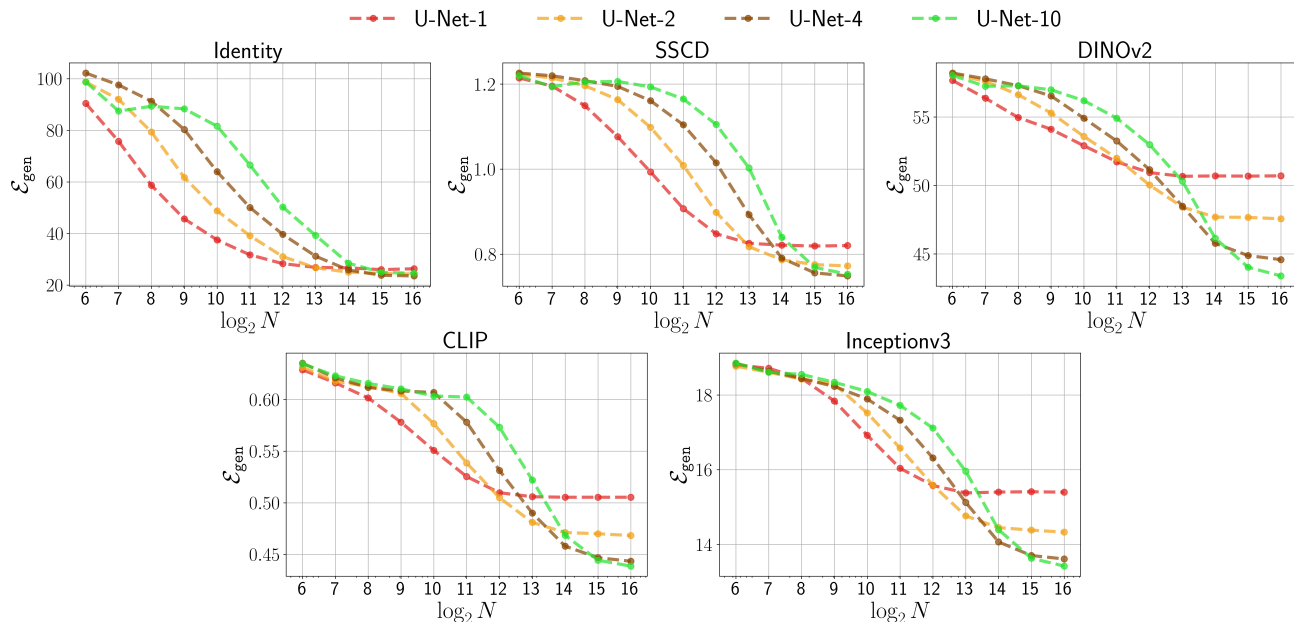


Figure 13. Comparison between different image descriptors.  $\mathcal{E}_{\text{gen}}$  plotted against  $\log_2(N)$  for a range of U-Net architectures (U-Net-1, U-Net-2, U-Net-4, U-Net-10) using different image descriptors, including identity function, SSCD, DINOv2, CLIP, Inceptionv3.

training sample sizes  $N$  and different student model architectures: U-Net-1, U-Net-2, U-Net-4, and U-Net-10. We vary  $M \in \{10, 32, 100, 316, 1000, 3163, 10000\}$ , and for each setting, generate 5 independent sets of  $\{\mathbf{x}_{\text{gen},T}^{(i)}\}_{i=1}^M$  initial noise estimate  $\mathcal{E}_{\text{gen}}$ , computing both the mean and variance.

As shown in Figure 14, the variance of  $\mathcal{E}_{\text{gen}}$  approaches zero as  $M$  increases to 10,000, indicating that when  $M \geq 10000$ , the empirical estimate of  $\mathcal{E}_{\text{gen}}$  converges to its value over the underlying distribution. This result holds consistently across different model architectures.

#### E.4. Teacher Model Architecture

We end this section by examining how different teacher models affect the evaluation protocol. Specifically, we consider three types of diffusion models: EDM, Rectified Flow (Rect) (Liu et al., 2023), and UViT. Using the CIFAR-10 dataset, we train three teacher models, one for each of these diffusion types. For each teacher model, we then evaluate all three diffusion models as student models. We report their corresponding  $\mathcal{E}_{\text{gen}}$  values. Both teacher and student models use the same sampling method, the second-order Heun solver with 18 steps.

As shown in Figure 15, the  $\mathcal{E}_{\text{gen}}$  is approximately 0.7 when both the student and teacher models are selected from EDM or UViT. However,  $\mathcal{E}_{\text{gen}}$  increases to around 0.8 when either the student or teacher model is Rect. According to its original paper, Rect has the poorest generative performance among the three, as measured by FID. This suggests that the teacher model should possess strong generative performance to serve as an underlying distribution that is close to the real-world data distribution. Therefore, in this paper, we adopt EDM as the teacher model, as it achieves the lowest FID among the three models.

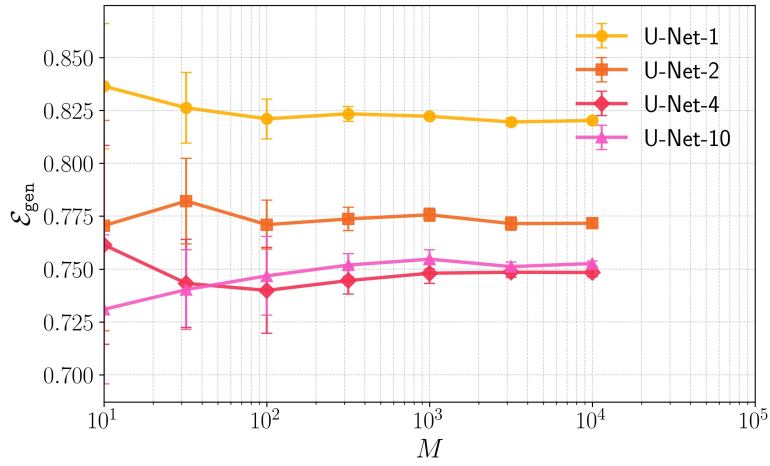


Figure 14. **Comparison across evaluation sample sizes.** The mean and variance of  $\mathcal{E}_{\text{gen}}$  are plotted against the number of evaluation samples  $M$  for various U-Net architectures (U-Net-1, U-Net-2, U-Net-4, U-Net-10), with a fixed number of training samples  $N = 2^{16}$ .

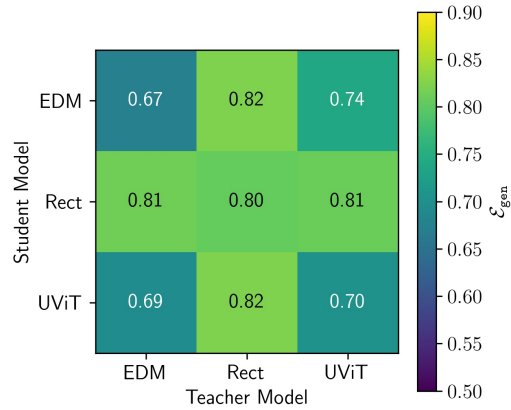


Figure 15. **Comparison of different teacher models.** The figure shows the  $\mathcal{E}_{\text{gen}}$  values for various student models (EDM, Rect, UViT) trained using different teacher models (EDM, Rect, UViT), with a fixed training data size of  $N = 2^{16}$ .