

From Navigation to Refinement: Revealing the Two-Stage Nature of Flow-based Diffusion Models through Oracle Velocity

Haoming Liu^{1,2} Jinnuo Liu¹ Yanhao Li¹ Liuyang Bai¹ Yunkai Ji¹ Yuanhe Guo^{1,2}
Shenji Wan^{1,2} Hongyi Wen^{1,2}

¹Center for Data Science, New York University Shanghai ²New York University

Abstract

Flow-based diffusion models have emerged as a leading paradigm for training generative models across images and videos. However, their memorization-generalization behavior remains poorly understood. In this work, we revisit the flow matching (FM) objective and study its marginal velocity field, which admits a closed-form expression, allowing exact computation of the oracle FM target. Analyzing this oracle velocity field reveals that flow-based diffusion models inherently formulate a two-stage training target: an early stage guided by a mixture of data modes, and a later stage dominated by the nearest data sample. The two-stage objective leads to distinct learning behaviors: the early navigation stage generalizes across data modes to form global layouts, whereas the later refinement stage increasingly memorizes fine-grained details. Leveraging these insights, we explain the effectiveness of practical techniques such as timestep-shifted schedules, classifier-free guidance intervals, and latent space design choices. Our study deepens the understanding of diffusion model training dynamics and offers principles for guiding future architectural and algorithmic improvements.

1. Introduction

Diffusion models [7, 14, 33, 34] have emerged as a powerful class of generative methods, capable of synthesizing high-fidelity samples across diverse domains, such as images [5, 8, 19, 20, 30] and videos [3, 27, 29, 39]. These models learn complex data distributions by progressively transforming a prior distribution (e.g., Gaussian) into the data distribution through an interpolation process parameterized by denoising or transport directions. Various formulations have been proposed to characterize the diffusion process, including probabilistic models [14, 34], score-based ODEs/SDEs [36–38], and flow matching [1, 22, 23]. These

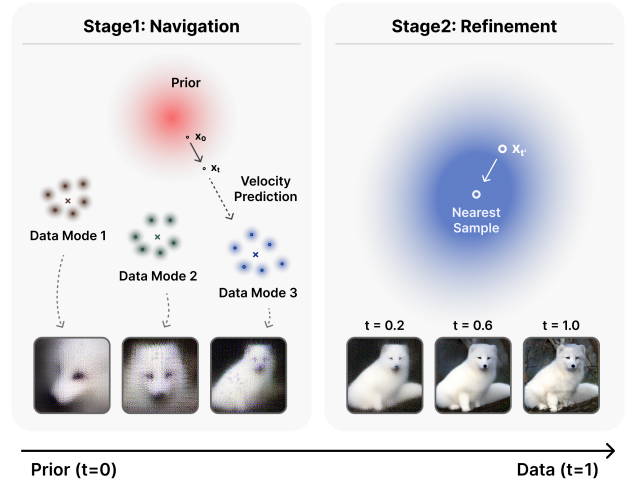


Figure 1. Illustration of the two stages in flow-based diffusion models. In the navigation stage (near prior), the target is guided by a mixture of multiple data samples, forming global layouts. In the refinement stage (near data), the target is dominated by the nearest data sample, refining fine-grained visual details.

perspectives give rise to diverse training objectives (e.g., noise, score, or velocity prediction) and correspond to different parameterizations of the probability flow ODE (PF-ODE) [21, 38]. Recent advances have largely converged on the flow matching formulation, where the model is trained to predict the velocity field under a linear schedule (also known as the canonical linear flow or rectified flow [23]). Thanks to its simplicity and stable training dynamics, this canonical formulation has become the de facto standard for training state-of-the-art diffusion models.

Meanwhile, a growing line of research has sought to understand the training and inference behaviors of diffusion models, particularly on the balance between **memorization** (i.e., the tendency to reproduce training samples) and **generalization** (i.e., the ability to synthesize novel samples).

Earlier works analyze this phenomenon from diverse viewpoints, such as by deriving quantitative metrics to characterize model behavior [11, 43] and by analyzing the underlying factors that lead to memorization/generalization [15, 32, 42]. In particular, most works primarily: (1) focus on the model behaviors when sampling from scratch, and (2) verify their findings on datasets of small scales/resolutions (e.g., FFHQ [16], CIFAR-10 [17]). While such insights can be informative in low-data regimes, such memorization behaviors become increasingly unlikely as the data scales to ImageNet-level [6] or beyond. Song *et al.* [35] recently observe the memorization/generalization divergence on ImageNet-scale diffusion models when resuming sampling from different temporal ranges. In general, resuming from an earlier timestep (near the prior) leads to novel samples, whereas later ones tend to reproduce training images. Hence, a key question arises: *what underlying principles govern the balance between memorization and generalization in diffusion models trained on large-scale data?*

To address this, we trace the origin of model behavior back to its training objective. We first revisit flow matching (FM) [22, 23] and its gradient-equivalent proxy, conditional flow matching (CFM), where the former is generally regarded as intractable due to the unavailable ground-truth velocity field. In this work, we refine this conventional view by showing that the marginal velocity field of rectified flow admits a closed-form expression under a Gaussian prior and a finite training set (Sec. 2). This allows us to compute the **oracle velocity** (defined by the FM objective) at any location in the sample space. Through the lens of the oracle velocity field, the effective training target of flow-based diffusion models exhibits two distinctive stages (Sec. 3). In the early stage that is closer to the prior (termed as the **navigation** stage), the oracle velocity contains combined information of multiple data points, and guides the model toward a mixture of relevant data modes. In the later stage (termed as the **refinement** stage), the velocity field is dominated by a single data point. In addition, we identify *data dimensionality* and *sample size* as the two key factors for determining the point of stage transition over time.

We next analyze the model’s behavior in light of this two-stage structure (Sec. 4). Overall, the navigation stage primarily establishes the global image layout, whereas the refinement stage concentrates on polishing fine-grained visual details. We hypothesize that the model’s **generalization** capability arises mainly from the **navigation** stage, while its **memorization** behavior roots from the **refinement** stage. Moreover, we observe that the learning difficulty differs across stages: the refinement stage poses a greater challenge to learn, while the navigation stage is comparatively easier.

We leverage our stage-level insights to understand why several empirically effective techniques succeed (Sec. 5). Specifically, we investigate how the two-stage nature sheds

light on: (1) the effect of the timestep shift factor for designing non-uniform inference-time schedules, (2) the role of classifier-free guidance (CFG) intervals in balancing navigation and refinement, and (3) the influence of latent space structure on model learning. Finally, we discuss how these findings relate to broader practices in diffusion model training and highlight promising yet underexplored directions for future improvement.

2. Oracle Velocity from Empirical Mixture

2.1. Preliminaries on Flow Matching

We consider two endpoint distributions: a simple prior distribution p_{prior} serving as the source, and a target data distribution p_{data} as the destination. The goal is to learn a continuous transport map that evolves the source into the target. This process induces a family of intermediate marginal densities $\{p_t\}_{t \in [0,1]}$ defined over a normalized time interval $t \in [0, 1]$ ¹. Despite the formulation differences between diffusion models and flow matching, their intrinsic dynamics are both governed by the *continuity equation*:

$$\frac{\partial p_t(x)}{\partial t} + \nabla_x \cdot (u_t(x)p_t(x)) = 0, \quad (1)$$

which couples the evolving probability path p_t with the underlying velocity field $u_t(x)$ that transports the probability mass along the path. The flow matching (FM) objective learns a neural network $v_t(x_t; \theta)$ to regress $u_t(x_t)$:

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, p_t(x_t)} \|v_t(x_t; \theta) - u_t(x_t)\|^2. \quad (2)$$

However, the FM objective is generally intractable when a closed-form u_t is unavailable, so a common practice is to construct a conditional probability path $p_t(x_t | x_1)$ based on a particular data sample x_1 . This is also known as the conditional flow matching (CFM) objective:

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, q(x_1), p_t(x_t|x_1)} \|v_t(x_t; \theta) - u_t(x_t | x_1)\|^2, \quad (3)$$

which has been shown to share identical gradients with the FM objective. Meanwhile, the probability path p_t can be constructed in various ways. One predominant approach is to randomly sample $x_0 \sim p_{\text{prior}}$, $x_1 \sim p_{\text{data}}$, and interpolate linearly via time-dependent scaling factors α_t and σ_t :

$$x_t = \alpha_t x_1 + \sigma_t x_0. \quad (4)$$

Hence, the conditional training target is given by:

$$u_t(x_t | x_1) = \dot{\alpha}_t x_1 + \dot{\sigma}_t x_0. \quad (5)$$

For rectified flow, we set $(\alpha_t, \sigma_t) = (t, 1 - t)$ and adopts $u_t(x_t | x_1) = x_1 - x_0$ as the training target.

¹The time conventions used in diffusion and flow matching works are highly inconsistent; here we follow Lipman *et al.* [22], let $t = 0$ corresponds to the prior distribution and $t = 1$ to the data distribution.

2.2. Closed-Form Oracle under Gaussian Prior

While the original Flow Matching (FM) objective (Eq. 2) is intractable under unknown probability paths, it admits a closed-form oracle velocity field under a few mild conditions in the context of flow-based diffusion models. Specifically, we assume (1) a Gaussian prior distribution, (2) a finite dataset $\{x_1^{(i)}\}_{i=1}^N$ approximating p_{data} , and (3) a linear interpolation path as in rectified flow. Under these assumptions, we can explicitly compute the conditional expectation $\mathbb{E}[u_t(x_t | x_1) | x_t]$, which yields the following closed-form expression for the oracle velocity field.

Theorem 2.1 (Closed-Form Oracle under Gaussian Prior). *Let the data distribution be represented as an empirical mixture over finite dataset $\{x_1^{(i)}\}_{i=1}^N$ and consider the linear interpolation $x_t = \alpha_t x_1^{(i)} + \sigma_t x_0$ with $x_0 \sim \mathcal{N}(0, I)$ and uniformly sampled $x_1^{(i)}$. Then, the oracle velocity field $u_t^*(x_t, t) := \mathbb{E}[u_t(x_t | x_1) | x_t]$ admits the closed form:*

$$u_t^*(x_t, t) = A_t \sum_{i=1}^N \gamma_i(x_t, t) x_1^{(i)} + B_t x_t, \quad (6)$$

where the coefficients $A_t = \dot{\alpha}_t - \frac{\alpha_t \dot{\sigma}_t}{\sigma_t}$, $B_t = \frac{\dot{\sigma}_t}{\sigma_t}$, and the normalized posterior weights $\gamma_i(x_t, t)$ are given by:

$$\gamma_i(x_t, t) = \frac{\exp\left(-\frac{\|x_t - \alpha_t x_1^{(i)}\|^2}{2\sigma_t^2}\right)}{\sum_{j=1}^N \exp\left(-\frac{\|x_t - \alpha_t x_1^{(j)}\|^2}{2\sigma_t^2}\right)}. \quad (7)$$

The proof can be given by applying Bayes' rule to the path marginal (a Gaussian mixture) and taking the conditional expectation given x_t ; the full derivation is provided in the Appendix. For class-conditional generation, the oracle velocity can be computed within each class-specific subset of samples $\{x_1^{(i)}\}_{i \in \mathcal{I}_y}$, denoted as $u_t^*(x_t, t | y)$.

3. The Two-Stage Training Target

Given the closed-form expression in Eq. 6, Flow Matching training can be viewed as a supervised learning problem, where the network is tasked with predicting the target label $u_t^*(x_t, t)$ for each input pair (x_t, t) . Interestingly, this (conceptually infinite) training dataset comprises samples with two distinct characteristics, separated by the timestep t . Roughly speaking, for x_t samples with $t \in [0.0, 0.1]$, the training target $u_t^*(x_t, t)$ is influenced by multiple data points from the target distribution, whereas for larger t , the target becomes dominated by a single nearest data point.

Specifically, we consider a class-conditional image generation setting on the ImageNet [6] dataset, where we train flow-based diffusion models within the latent space of VAE [41] and SD-VAE [31]. Under rectified flow [23], the noisy latents x_t are constructed from randomly sampled prior/data pairs (x_0, x_1) following Eq. 4. Then, we can

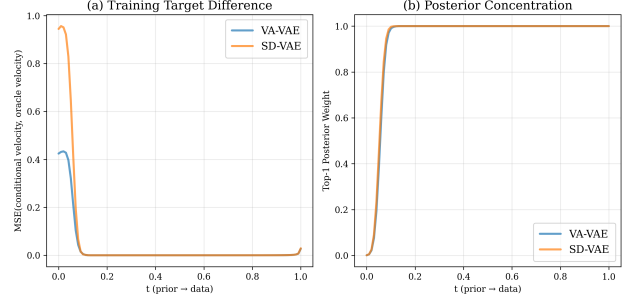


Figure 2. (a) MSE between u_t^* and the CFM target $(x_1 - x_0)$ across timesteps; (b) Average top-1 posterior weight $\gamma_i(x_t, t)$ showing rapid concentration after $t = 0.1$; both plots reveal a clear two-stage behavior emerging in the oracle training target.

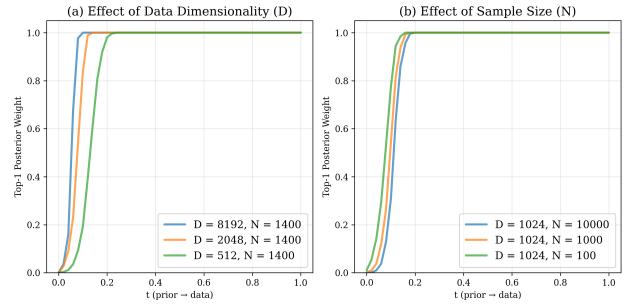


Figure 3. Plots of top-1 posterior weight under varying conditions. (a) Higher data dimensionality accelerates saturation. (b) A larger sample size delays the transition. Zoom in for the details.

compute the class-conditioned oracle velocity $u_t^*(x_t, t | y)$ and contrast against the CFM target $(x_1 - x_0)$. As shown in Fig. 2(a), the MSE between the noisy and oracle targets reveals that the divergence is concentrated in the early interval $t \in [0.0, 0.1]$, while the two targets closely align for later timesteps ($t > 0.1$). A minor discrepancy is also observed as $t \rightarrow 1$, which arises from the rapidly shrinking $2\sigma_t^2$ term in the denominator of γ_i . Meanwhile, the two-stage target can be further validated through the top-1 posterior weight among $\gamma_i(x_t, t)$ (Fig. 2(b)), where the normalized top-1 weight rapidly saturates to 1 beyond $t = 0.1$, suggesting that the oracle velocity field has collapsed to a single dominant data point. Together, these results demonstrate that the oracle training target inherently exhibits a two-stage nature: *in the early stage ($t \in [0.0, 0.1]$), the oracle target corresponds to the mixture of multiple data samples, whereas in the later stage ($t \in (0.1, 1.0]$), it rapidly switches to a near-deterministic target akin to the CFM objective (Eq. 3).*

Here, the rapid saturation of the top-1 posterior mass roots from the posterior weighting defined in Eq. 7. More specifically, we first recall the conditional case (i.e., considering a single data sample), where the marginal distribution at x_t follows a Gaussian whose variance is scaled by some temporal factors. In contrast, under the oracle setting, the marginal distribution at x_t is a mixture of N Gaussians,

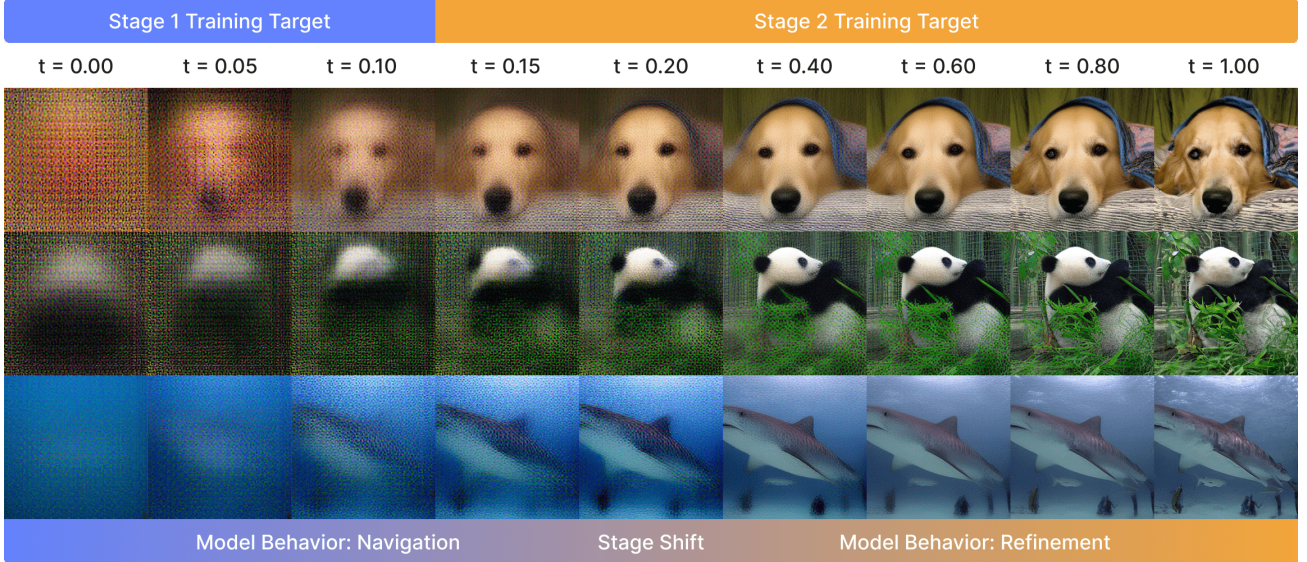


Figure 4. Intermediate predictions of a LightningDiT-XL/1 [41] model across timesteps. Overall, early stages primarily **navigate** global layout, while later stages **refine** fine-grained details. Notably, the empirically observed stage transition ($t \approx 0.2$) in model behavior lags slightly behind the training target shift ($t \approx 0.1$). Zoom in for the best view. Better view with color.

where N denotes the number of samples. Considering the expression of $\gamma_i(x_t, t)$ (Eq. 7) in a D -dimensional space, as the squared distances scale with D and enter the exponent divided by $2\sigma_t^2$, even modest differences in distance translate into exponentially large differences in weight as D grows (for fixed σ_t). As a result, once σ_t becomes small compared to the typical inter-sample distance, the posterior $\gamma_i(x_t, t)$ becomes sharply peaked on the nearest sample. To sum up, *the temporal positioning and the sharpness of the stage split are governed by the interplay between data dimensionality D , the sample size N , the time-dependent scaling factor σ_t , and the geometric spread of the dataset.*

We verify our claim in Fig. 3, where the top-1 posterior weight is plotted under varying dimensionalities and sample sizes. Overall, we observe that increasing the data dimension results in faster top-1 posterior weight concentration, whereas increasing the number of samples mitigates this effect. In more concrete settings (e.g., ImageNet data at 256^2 resolution), we typically have latent dimensions $D = 4096$ or 8192 and the class-level sample size $N \approx 1400$, thereby yielding the observed top-1 posterior saturation around $t = 0.1$. This implies that *the effective training target naturally differs across datasets of varying dimensionality and sample size, even under the same rectified flow objective.* We also note that our demonstration adopts synthetic unit-Gaussian samples to approximate mean- and variance-normalized data; while real data distributions may deviate slightly, their behavior remains similar as in Fig. 2. The influence of latent space structure on this phenomenon is further analyzed in Sec. 5.3.

4. Model Behaviors under Two-Stage Target

4.1. Observing Empirical Stage Transitions

We next examine whether a trained flow-based model exhibits analogous stage-specific behaviors characterized by the two-stage target. Specifically, we adopt a LightningDiT-XL/1 [41] model and visualize its intermediate predictions, which are calculated by taking a single Euler step from an intermediate timestep $t' \in [0, 1]$ to $t = 1$ based on the velocity prediction at $x_{t'}$. While diffusion models are designed to progressively denoise noisy inputs, our study aims to identify the transition in model behavior by analyzing the temporal evolution of intermediate predictions. As demonstrated in Fig. 8, several observations arise. First, the model’s prediction at the prior $t = 0$ collapses to a coarse class-mean (e.g., predominantly blue for shark, black-and-white for panda). As the trajectory progresses, the early stage primarily navigates the global image layout, which stabilizes around $t \approx 0.2$. The later stage focuses exclusively on refining local visual details with minimal semantic or structural deviation, which complies with the consistent training target in this stage. In addition, we also find that *the stage transition in model behavior occurs slightly later than the stage split implied by the oracle target.* One potential explanation for this is that the model may require additional temporal margin to rectify its accumulated prediction errors after shifting to the consistent, CFM-like training target.

4.2. Relations to Memorization/Generalization

Moving forward, we further interpret the model’s **generalization** (i.e., the ability to generate novel samples) and

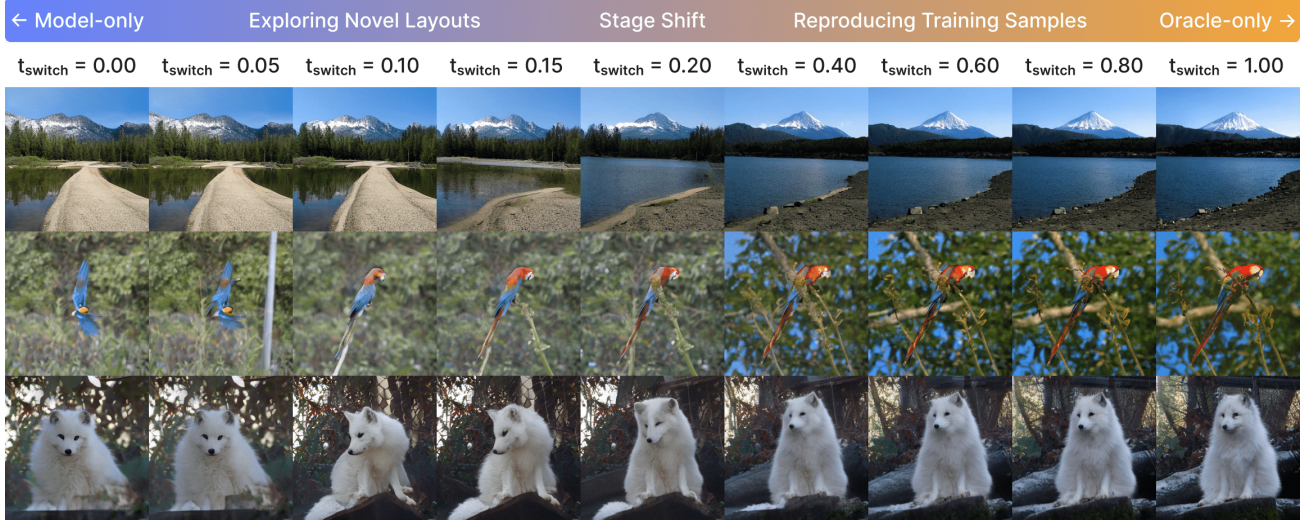


Figure 5. Mixed sampling results with switch point t_{switch} . Oracle u_t^* is used before t_{switch} and LightningDiT-XL/1 [41] afterward. Overall, early switching yields diverse novel outputs (generalization), while late switching reproduces training samples (memorization).

memorization (i.e., the tendency to reproduce training samples) behaviors through the lens of stage-level insights. We first consider a mixed sampling scheme. Specifically, we start from a random Gaussian prior and take Euler steps based on the oracle velocity u_t^* in the early stage. Beyond a certain threshold t_{switch} , we switch to take Euler steps based on the model’s velocity predictions. The oracle velocity in the first stage ensures that the intermediate states remain precisely on the interpolated distribution, isolating the effect of imperfect model predictions. By varying t_{switch} , we observe distinct behaviors of generalization or memorization. The qualitative results are presented in Fig. 5. When the oracle velocity is applied throughout the entire trajectory, the process deterministically retrieves a random training sample; when $t_{\text{switch}} \in (0.2, 1.0]$, the model can largely replicate the training trajectory, producing images that closely resemble the training images. We attribute such **memorization** behaviors to the trivial and consistent training target in this regime. In contrast, when $t_{\text{switch}} \in [0.0, 0.2]$, the strong prior corruption prevents the model from inferring the original training trajectory, thus deviating from the training instance and exhibiting **generalization** capability. These results also demonstrate that memorization-from-scratch (i.e., sampling from the prior and recovering exact training samples) is highly unlikely when the dataset size substantially exceeds the effective memorization capacity of the model; in contrast, resuming from training trajectories in the refinement stage is very likely to trigger memorization behaviors.

Beyond training trajectories, we investigate whether the model’s refinement capability generalizes to unseen data. We construct x_t by combining random Gaussian priors and unseen image latents from the ImageNet validation split, resuming the sampling process at t_{resume} . As illustrated in

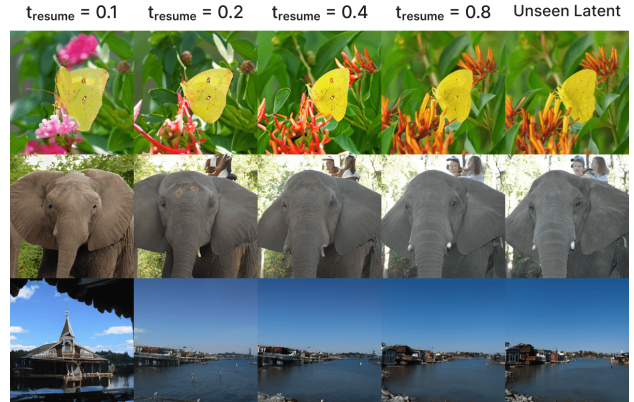


Figure 6. Qualitative results for refinement generalization. When resuming from $t_{\text{resume}} \geq 0.2$ on validation image latents, the model largely preserves global structure while improvising fine details.

Fig. 6, for $t_{\text{resume}} \geq 0.2$ the global layout is maintained, while still introducing slight variations in fine details. Based on the results on seen/unseen data, we can conclude that: **generalization** stems from the early **navigation** stage, where recovering training trajectories is difficult; whereas **memorization** arises in the late **refinement** stage, where the targets are trivial and consistent.

4.3. Learning Differs in the Two Stages

Our goal is to assess how the model’s learning difficulty varies across timesteps and whether it aligns with the two-stage structure. As shown in Fig. 7(a), we plot the training MSE across timesteps by comparing its predicted velocity (without classifier-free guidance [13]) to both the noisy conditional target $(x_1 - x_0)$ and the class-conditioned oracle target $u_t^*(x_t, t | y)$. Consistent with the discrepancy

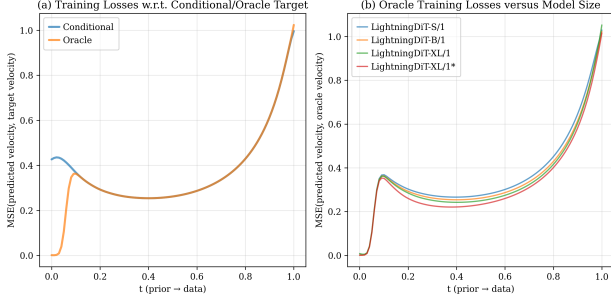


Figure 7. Training loss trends across timesteps. (a) Training losses (MSE) with respect to conditional and oracle targets. (b) Oracle training losses under varying model sizes. Most models are trained for 100 epochs on the ImageNet [6] dataset. * indicates the model is trained for 800 epochs ($8\times$ training compute).

between these targets, *the loss divergence is concentrated in the navigation stage*, while both curves align closely in the refinement stage where their effective forms coincide. Notably, the model achieves near-perfect oracle loss when $t \rightarrow 0$, consistent with the trivial class-mean prediction under extreme noise in Fig. 4. The oracle loss then rises steadily through the navigation stage, reflecting the increasing challenge of steering toward specific data modes. At the transition into the refinement stage, we observe a noticeable bump in loss, followed by a dip in the mid-refinement region where the intermediate predictions are typically smoothed images. The loss climbs again as $t \rightarrow 1$, possibly due to the diversity of potential fine-detail refinements near the data manifold. Importantly, while sharing similar overall trends, *the oracle loss plot is inherently dependent on the latent space structures*; we provide further comparison between VA-VAE [41] and SD-VAE [31] in Fig. 9.

Fig. 7(b) investigates how model capacity and training compute influence the oracle loss across timesteps. Interestingly, *all models, regardless of parameter size or training duration, exhibit nearly identical performance in the early navigation stage*, where the observed model behavior is to predict coarse image layouts from heavily Gaussian-corrupted inputs. In contrast, *the divergence emerges in the vast refinement stage: larger models or those trained with substantially more compute achieve noticeably lower oracle losses*. Such superiority can also be captured by the generation Fréchet inception distance (gFID) [12] results (Tab.1). These results suggest that additional capacity and optimization primarily benefit the refinement of high-frequency details in generated images and are measurable by existing metrics. In addition, we highlight that the improvements to the navigation stage remain limited, possibly due to the difficulty of inferring global structure from near-prior noise, which needs subtle measurements beyond the oracle losses.

To further validate if the navigation performance is insensitive to model capacity, we conduct a mixed-generation

Sampling Interval	Stage 1 Model	Stage 2 Model	gFID@50K ↓
[0.0, 0.1] + [0.1, 1.0]	XL	XL	2.94
	Base	XL	3.71
	XL	Base	11.26
	Base	Base	12.45
[0.0, 0.2] + [0.2, 1.0]	XL	XL	2.60
	Base	XL	4.47
	XL	Base	9.24
	Base	Base	12.01

Table 1. gFID performance under different model combinations. We adopt the Base and XL variants of LightningDiT [41]. “Stage 1 Model” is used on the first sampling interval, while “Stage 2 Model” is used on the second. We evenly assign 25 uniform sampling steps to each sub-interval for fair comparison (no CFG).

Timestep Shift	Percentage of $t \in [0.0, 0.2]$	gFID@50K ↓
$s = 4.0$	6%	18.89
$s = 2.0$	12%	14.82
$s = 1.0$ (uniform)	22%	12.99
$s = 0.7$	28%	12.46
$s = 0.5$	34%	12.23
$s = 0.3$	46%	12.66
$s = 0.1$	72%	19.91

Table 2. Generation performance with different timestep shift factors (NFE = 50, no CFG). We adopt a LightningDiT-B/1 [41] model for evaluation. Overall, a moderate increase in early (navigation) steps gives the best trade-off in sampling quality.

experiment (Tab. 1). In this setup, we explicitly swap different models for Stage 1 and Stage 2 during sampling while keeping the total NFE quota (i.e., number of function evaluations in sampling) fixed. Consistent with the oracle-loss trends, replacing the Stage 1 model with a smaller-capacity variant yields little degradation in gFID, indicating that *coarse layout prediction under near-prior noise is largely unaffected by model size*. In contrast, substituting the Stage 2 model with a weaker one leads to a substantial drop in generation quality, confirming that the refinement on high-frequency details is the primary beneficiary of additional capacity and training compute.

5. Elucidating Stage-aware Practices

Building on the two-stage perspective developed in earlier sections, we now revisit several widely used empirical practices through stage-related insights. Our goal is to clarify why these techniques work and how their effects can be better understood and optimized by explicitly considering the distinct roles of navigation and refinement.

5.1. Optimizing Timestep Schedule

A common practice in diffusion sampling is to use a uniform timestep schedule over $t \in [0, 1]$. However, a natural question arises: *what is the optimal allocation of computation between the navigation and refinement stages under*

CFG Interval	gFID@50K ↓	CFG Interval	gFID@50K ↓
None	12.99	[0.0, 1.0]	10.79
[0.0, 0.1]	6.33	[0.0, 0.2]	6.62
[0.1, 0.2]	5.21	[0.0, 0.4]	8.03
[0.2, 0.3]	7.70	[0.0, 0.6]	9.19
[0.3, 0.4]	9.39	[0.0, 0.8]	10.14
[0.4, 0.5]	10.39	[0.1, 0.3]	3.54
[0.5, 0.6]	10.59	[0.1, 0.4]	4.16
[0.6, 0.7]	11.06	[0.1, 0.5]	2.82
[0.7, 0.8]	11.38	[0.1, 0.6]	2.80
[0.8, 0.9]	11.64	[0.1, 0.7]	2.86
[0.9, 1.0]	12.20	[0.1, 0.8]	2.97

Table 3. Generation performance under different CFG intervals (CFG factor $\omega = 2.5$). We adopt a LightningDiT-B/1 [41] model for evaluation. Optimal intervals concentrate in the early-mid **refinement** stage, while excessive guidance in the earliest **navigation** steps can degrade the sampling quality.

a fixed NFE budget? We narrow down our discussion to *timestep shifting* [8], an existing technique for constructing non-uniform timestep schedules. Specifically, timestep shifting transforms a uniform timestep t_n with a smooth monotonic mapping: $t_m = \frac{s t_n}{1 + (s-1)t_n}$, where s is a shift factor. With $s < 1$, we allocate more steps to the early (navigation) timesteps, while $s > 1$ biases the schedule toward later (refinement) timesteps. While this technique was originally proposed to address noise-level imbalance in high-resolution image generation, this approach naturally provides a flexible knob to control stage-wise sampling step allocation. The results in Tab. 2 indicate that biasing the timestep distribution slightly toward earlier steps leads to noticeably better samples, suggesting that *modestly emphasizing the navigation phase is beneficial for inference*.

5.2. Optimizing CFG Intervals

Classifier-free guidance (CFG) [13] modulates the conditional signal during sampling by amplifying the difference between conditional and unconditional predictions. Prior work has shown that applying CFG only on a selected sub-interval, rather than across all timesteps, leads to improved generation quality [18]. We revisit this idea in the context of flow-based diffusion models with results in Tab. 3. Our ablations reveal three consistent trends. First, *the most effective single short interval (of width 0.1) lies in the transition stage between navigation and refinement*. Second, when the CFG interval is expanded, it becomes beneficial to exclude the very initial segment [0.0, 0.1], likely because amplified guidance at extremely noisy states interferes with the formation of stable global layouts. Third, *the overall optimal CFG ranges tend to span the early and mid refinement stage*. This aligns with the prediction dynamics in Fig. 8, where velocity norms peak, suggesting that these timesteps correspond to high-confidence refinement in which CFG most effectively enhances sample fidelity.

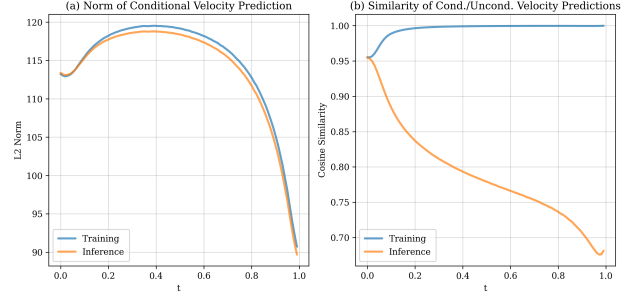


Figure 8. Analysis of model prediction trends. (a) Norm of velocity predictions peaks around $t=0.4$. (b) Conditional and unconditional velocity predictions diverge rapidly during inference.

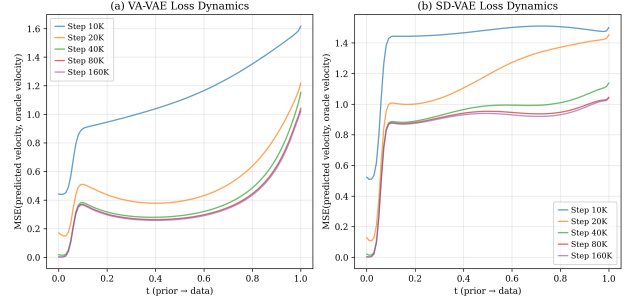


Figure 9. Convergence of oracle loss under different latent spaces: VA-VAE [41] losses decrease smoothly in the refinement stage; whereas SD-VAE [31] exhibits wavy patterns. The gFID dynamics are attached in the Appendix, where VA-VAE converges faster.

5.3. Influence of Latent Space Structure

We next examine how the choice of latent space affects the model’s dynamics under a two-stage oracle target. Fig. 9 compares the oracle loss dynamics under VA-VAE [41] and SD-VAE [31] latent spaces during training. In the navigation stage, both spaces exhibit a gradual oracle loss increase, suggesting that *the growing difficulty of coarse-layout prediction is agnostic to the latent space*. In the refinement stage, however, the loss trends diverge substantially. While the absolute magnitudes are not directly comparable, certain patterns arise from the shapes of the oracle loss curves. For VA-VAE, whose latent space is aligned with semantic structure through DINO-based VF loss [28, 41], the oracle loss converges smoothly toward a parabola-like shape as training progresses. This indicates that a well-organized latent manifold, where semantic modes are arranged coherently. In comparison, SD-VAE is trained purely for low-level reconstruction, so it converges to a flatter curve with noticeable waviness, suggesting that its latent manifold is less structured and the model struggles to perform consistent denoising across stages. Overall, *latent spaces with high-level concept alignment possess the properties that are preferred by diffusion modeling, as such alignment tends to induce clearer mode organization under Gaussian marginals and thus facilitates both navigation and refinement*.

6. Related Work

6.1. Diffusion-based Generative Models

Diffusion models learn data distributions by through scheduled interpolations to a simple prior (typically Gaussian). This process is governed by the probability flow ODE (PF-ODE) [38] and can be interpreted through several complementary viewpoints [21]. From a probabilistic perspective, DDPM-style methods define a discrete-time forward noising process and learn to invert it via noise/data prediction targets [7, 14]. The score matching perspective instead characterizes the dynamics in continuous time using ODE/SDE formulations [36, 38]. Another line of work sources from continuous normalizing flows (CNFs) [22] and stochastic interpolants [1, 2] and unifies to flow matching (FM), where the model directly learns the velocity field governing the transport between prior and data. Recent work favors the rectified flow formulation [23] for its simplicity and stable optimization; whereas emerging research explores new generative paradigms for few-step diffusion, such as ShortCut models [9], MeanFlow [10], and AlphaFlow [44].

6.2. Memorization versus Generalization

Understanding the balance between memorization and generalization is a central theme in generative modeling. One line of work focuses on measuring model behaviors through quantitative metrics, such as the effective model memorization (EMM) factor [11] and the probability flow distance (PFD) [43]. Another line of work investigates the factors influencing the memorization-generalization trade-off. For example, Zhang *et al.* [42] observe distinctive generalization and memorization regimes during training; Bonnaire *et al.* [4] show that diffusion models possess separate generalization and memorization timescales, with larger datasets widening the effective generalization window during training; Shi *et al.* [32] identify that declining entropy in recursively generated training data triggers a memorization-dominated collapse in diffusion models, which can be alleviated through an entropy-based selection strategy; Niedoba *et al.* [26] reveal that network denoisers generalize through localized denoising operations; Kadkhodaie *et al.* [15] find that two independently trained denoisers converge to nearly identical score functions. However, most existing studies are based on small datasets, low resolutions, and sampling-from-scratch analysis, making it difficult to extrapolate their conclusions to modern diffusion models at a larger scale. Recent work by Song *et al.* [35] observes this memorization-generalization trade-off at ImageNet scale, yet without probing the underlying causes. Our work advances this discussion by analyzing flow-based diffusion models through the oracle velocity in closed form. This enables us to dissect the two-stage behaviors between navigation and refinement.

7. Discussion

7.1. Improving Navigation Capability

Our findings indicate that the navigation performance under the FM/CFM objective remains largely unchanged even when model capacity is scaled. This suggests an inherent limitation in the near-prior regime and points to two possible explanations: 1) the prediction task in this early stage may be intrinsically trivial, causing models of different capacities to converge to similar solutions; 2) the supervision in this regime is too weak to benefit from larger model capacity, suggesting that extra guidance may be required. In addition, standard metrics like gFID provide little insight into intermediate prediction quality, which highlights the need for new measures that assess navigated layout fidelity and better capture model behaviors in the early stage.

7.2. The Hidden Mechanism behind Scaling

As the stage transition depends on data dimensionality and dataset size, scaling along these axes directly reshapes the underlying navigation-refinement dynamics. For example, higher dimensionality amplifies the curse of dimensionality, while increasing the dataset size exposes the model to a broader set of global layouts and fine-grained details. In practice, state-of-the-art diffusion systems scale data, model capacity, and training compute simultaneously; the findings from the oracle velocity offer intuition for why scaling systematically stabilizes navigation, enriches denoising, and ultimately improves generative quality.

7.3. Oracle Velocity for Training

A natural question is whether the closed-form oracle can directly serve as a training target. Our experiments show that oracle-supervised training is feasible and yields convergence behavior similar to CFM. However, this approach remains impractical at scale: computing the oracle during training introduces substantial data-scaling bottlenecks and heavy I/O overhead that grow quickly with dataset size. Nonetheless, the closed-form oracle remains a valuable analytical tool for analyzing model behaviors.

8. Conclusion

In this work, we introduce a principled framework for understanding diffusion models. By deriving the closed-form oracle velocity, we show that the effective training target is inherently two-stage: a multi-sample guided navigation regime near the prior and a single-sample dominated refinement regime near the data. This perspective explains the divergence between memorization and generalization, clarifies what diffusion models learn and how they sample, and sheds light on widely adopted empirical practices. Our findings provide an oracle-driven viewpoint for advancing the development of next-generation diffusion models.

References

- [1] Michael Samuel Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic interpolants. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 8
- [2] Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023. 8
- [3] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023. 1
- [4] Tony Bonnaire, Raphaël Urfin, Giulio Biroli, and Marc Mezard. 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. 8
- [5] Google DeepMind. Nano banana (gemini 2.5 flash image). <https://ai.google.dev/gemini-api/docs/image-generation>, 2025. 1
- [6] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 2, 3, 6, 11
- [7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 1, 8
- [8] Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024. 1, 7
- [9] Kevin Frans, Danijar Hafner, Sergey Levine, and Pieter Abbeel. One step diffusion via shortcut models. In *The Thirteenth International Conference on Learning Representations*, 2025. 8
- [10] Zhengyang Geng, Mingyang Deng, Xingjian Bai, J Zico Kolter, and Kaiming He. Mean flows for one-step generative modeling. *arXiv preprint arXiv:2505.13447*, 2025. 8
- [11] Xiangming Gu, Chao Du, Tianyu Pang, Chongxuan Li, Min Lin, and Ye Wang. On memorization in diffusion models, 2024. 2, 8
- [12] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
- [13] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 5, 7, 12
- [14] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1, 8
- [15] Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization in diffusion models arises from geometry-adaptive harmonic representations. In *The Twelfth International Conference on Learning Representations*, 2024. 2, 8
- [16] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019. 2
- [17] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). URL <http://www.cs.toronto.edu/kriz/cifar.html>, 2010. 2
- [18] Tuomas Kynkäänniemi, Miika Aittala, Tero Karras, Samuli Laine, Timo Aila, and Jaakko Lehtinen. Applying guidance in a limited interval improves sample and distribution quality in diffusion models. *Advances in Neural Information Processing Systems*, 37:122458–122483, 2024. 7
- [19] Black Forest Labs. Flux.1. <https://github.com/black-forest-labs/flux>, 2023. 1, 15
- [20] Black Forest Labs, Stephen Batifol, Andreas Blattmann, Frederic Boesel, Saksham Consul, Cyril Diagne, Tim Dockhorn, Jack English, Zion English, Patrick Esser, et al. Flux.1 kontext: Flow matching for in-context image generation and editing in latent space. *arXiv preprint arXiv:2506.15742*, 2025. 1
- [21] Chieh-Hsin Lai, Yang Song, Dongjun Kim, Yuki Mitsufuji, and Stefano Ermon. The principles of diffusion models. *arXiv preprint arXiv:2510.21890*, 2025. 1, 8
- [22] Yaron Lipman, Ricky T. Q. Chen, Heli Ben-Hamu, Maximilian Nickel, and Matthew Le. Flow matching for generative modeling. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 2, 8
- [23] Xingchao Liu, Chengyue Gong, and qiang liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. In *The Eleventh International Conference on Learning Representations*, 2023. 1, 2, 3, 8
- [24] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. 12
- [25] E. A. Nadaraya. On estimating regression. *Theory of Probability & Its Applications*, 9(1):141–142, 1964. 11
- [26] Matthew Niedoba, Berend Zwartsenberg, Kevin Patrick Murphy, and Frank Wood. Towards a mechanistic explanation of diffusion model generalization. In *Forty-second International Conference on Machine Learning*, 2025. 8
- [27] OpenAI. Sora: A text-to-video generation model. <https://openai.com/index/sora>, 2024. 1
- [28] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 7, 12
- [29] Xiangyu Peng, Zangwei Zheng, Chenhui Shen, Tom Young, Xinying Guo, Binluo Wang, Hang Xu, Hongxin Liu, Mingyan Jiang, Wenjun Li, et al. Open-sora 2.0: Training a commercial-level video generation model in \$200 k. *arXiv preprint arXiv:2503.09642*, 2025. 1
- [30] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and

- Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 1
- [31] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 3, 6, 7, 12
- [32] Lianghe Shi, Meng Wu, Huijie Zhang, Zekai Zhang, Molei Tao, and Qing Qu. A closer look at model collapse: From a generalization-to-memorization perspective. In *The Impact of Memorization on Trustworthy Foundation Models: ICML 2025 Workshop*, 2025. 2, 8
- [33] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International conference on machine learning*, pages 2256–2265. PMLR, 2015. 1
- [34] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *International Conference on Learning Representations*, 2021. 1
- [35] Kiwhan Song, Jaeyeon Kim, Sitan Chen, Yilun Du, Sham Kakade, and Vincent Sitzmann. Selective underfitting in diffusion models. *arXiv preprint arXiv:2510.01378*, 2025. 2, 8
- [36] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019. 1, 8
- [37] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *Advances in neural information processing systems*, 33:12438–12448, 2020.
- [38] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 1, 8
- [39] Team Wan, Ang Wang, Baole Ai, Bin Wen, Chaojie Mao, Chen-Wei Xie, Di Chen, Feiwu Yu, Haiming Zhao, Jianxiao Yang, et al. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*, 2025. 1
- [40] Geoffrey S. Watson. Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)*, 26 (4):359–372, 1964. 11
- [41] Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruction vs. generation: Taming optimization dilemma in latent diffusion models. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 15703–15712, 2025. 3, 4, 5, 6, 7, 12, 13, 14
- [42] Huijie Zhang, Jinfan Zhou, Yifu Lu, Minzhe Guo, Peng Wang, Liyue Shen, and Qing Qu. The emergence of reproducibility and generalizability in diffusion models. *arXiv preprint arXiv:2310.05264*, 2023. 2, 8
- [43] Huijie Zhang, Zijian Huang, Siyi Chen, Jinfan Zhou, Zekai Zhang, Peng Wang, and Qing Qu. Understanding generalization in diffusion models via probability flow distance. In *High-dimensional Learning Dynamics 2025*, 2025. 2, 8
- [44] Huijie Zhang, Aliaksandr Siarohin, Willi Menapace, Michael Vasilkovsky, Sergey Tulyakov, Qing Qu, and Ivan Skorokhodov. Alphaflow: Understanding and improving meanflow models. *arXiv preprint arXiv:2510.20771*, 2025. 8

A. Proof of Theorem 2.1

The Flow Matching (FM) objective (Eq. 8) is given by:

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, p_t(x_t)} \|v_t(x_t; \theta) - u_t(x_t)\|^2. \quad (8)$$

The marginal velocity field $u_t(x_t)$ in Eq. 8 is generally intractable under unknown probability paths. However, we want to show that it admits a closed-form solution under the following conditions: (i) a Gaussian prior distribution, (ii) a finite dataset $\{x_1^{(i)}\}_{i=1}^N$ approximating p_{data} , and (iii) a linear interpolation path as in rectified flow. Formally, our goal is to obtain a closed-form expression for the following conditional expectation:

$$u_t^*(x_t, t) := \mathbb{E}_{x_1 \sim p_{\text{data}}} [u_t(x_t | x_1) | x_t]. \quad (9)$$

Under linear flows (i.e., probability paths constructed via linear interpolation), we have:

$$x_t = \alpha_t x_1 + \sigma_t x_0 \implies x_0 = \frac{x_t - \alpha_t x_1}{\sigma_t}. \quad (10)$$

Recall that the conditional velocity in CFM is given by:

$$u_t(x_t | x_1) = \dot{\alpha}_t x_1 + \dot{\sigma}_t x_0. \quad (11)$$

Substituting Eq. 10 into Eq. 11 writes $u_t(x_t | x_1)$ as a function of x_t and x_1 :

$$u_t(x_t | x_1) = \left(\dot{\alpha}_t - \frac{\alpha_t \dot{\sigma}_t}{\sigma_t} \right) x_1 + \frac{\dot{\sigma}_t}{\sigma_t} x_t. \quad (12)$$

Taking the conditional expectation given x_t yields:

$$u_t^*(x_t, t) := \mathbb{E}[u_t(x_t | x_1) | x_t] \quad (13)$$

$$= \left(\dot{\alpha}_t - \frac{\alpha_t \dot{\sigma}_t}{\sigma_t} \right) \mathbb{E}[x_1 | x_t] + \frac{\dot{\sigma}_t}{\sigma_t} x_t, \quad (14)$$

Given a finite dataset $\{x_1^{(i)}\}_{i=1}^N$, we essentially approximate the true data distribution $p_{\text{data}}(x_1)$ via an empirical mixture:

$$p_{\text{data}}(x_1) \approx \frac{1}{N} \sum_{i=1}^N \delta(x_1 - x_1^{(i)}), \quad (15)$$

where $\delta(\cdot)$ is the Dirac delta function with $\delta(0) = \infty$ and zero elsewhere. Accordingly, the empirical probability path marginal $\tilde{p}_t(x_t)$ is given by a Gaussian mixture:

$$\tilde{p}_t(x_t) = \frac{1}{N} \sum_{i=1}^N \mathcal{N}(x_t; \alpha_t x_1^{(i)}, \sigma_t^2 I). \quad (16)$$

By Bayes' rule, the posterior $p(x_1^{(i)} | x_t)$ is proportional to $p(x_t | x_1^{(i)})$ up to a common normalizing factor that ensures the probabilities sum to one. Since all mixture components share the same Gaussian covariance $\sigma_t^2 I$ and uniform

prior weight $1/N$, their normalization constants cancel out when computing the posterior weights. Let $\gamma_i(x_t, t)$ denote the resulting normalized weighting function that reflects the relative contribution of each data sample $x_1^{(i)}$ to the current point x_t , we have:

$$\gamma_i(x_t, t) = \frac{\exp\left(-\frac{\|x_t - \alpha_t x_1^{(i)}\|^2}{2\sigma_t^2}\right)}{\sum_{j=1}^N \exp\left(-\frac{\|x_t - \alpha_t x_1^{(j)}\|^2}{2\sigma_t^2}\right)}. \quad (17)$$

Hence, the posterior mean is given by:

$$\mathbb{E}[x_1 | x_t] = \sum_{i=1}^N \gamma_i(x_t, t) x_1^{(i)}. \quad (18)$$

This is also known as the Nadaraya-Watson estimator [25, 40]. Combining Eq. 14 and 18, we reach the closed form:

$$u_t^*(x_t, t) = A_t \sum_{i=1}^N \gamma_i(x_t, t) x_1^{(i)} + B_t x_t, \quad (19)$$

where $A_t = \dot{\alpha}_t - \frac{\alpha_t \dot{\sigma}_t}{\sigma_t}$, $B_t = \frac{\dot{\sigma}_t}{\sigma_t}$. We refer to this closed-form expression of the marginal velocity field under the linear probability path construction as the *oracle velocity field*. Moreover, the oracle velocity field can also be evaluated conditionally; for instance, under a class-conditional generation setting, it can be computed within each class-specific subset, denoted as $u_t^*(x_t, t | y)$. \square

B. Oracle-Supervised Training

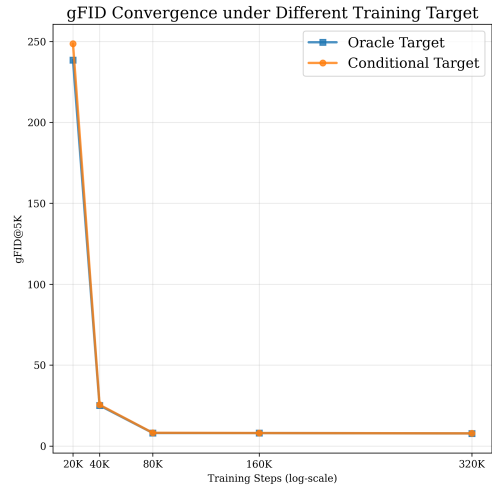


Figure 10. Convergence of gFID@5K when training rectified flow models with oracle/conditional target on a 100-class ImageNet [6] subset. Oracle supervision offers a slight advantage in early training (due to more accurate guidance on near-prior timesteps), while the later training dynamics largely overlap (as the oracle target collapses to the conditional target on most timesteps).

C. Implementation Details

architecture			
LightningDiT [41] variants	Small	Base	XL
depth	12	12	28
hidden dim	384	768	1152
heads	6	12	16
image size	256		
patch size	1 (VA-VAE), 2 (SD-VAE)		
latent size	16 × 16		
training			
epochs	{ 100, 800 }		
optimizer	AdamW [24] ($\beta_1, \beta_2 = 0.9, 0.995$)		
batch size	512		
learning rate	1e-4		
learning rate schedule	constant		
weight decay	0		
max gradient norm	1.0		
ema decay	0.9999		
time sampler	Uniform[0, 1]		
class token drop (for CFG)	0.1		
sampling			
ODE solver	Euler		
ODE steps	50		
time steps	uniform / stage-wise uniform / shifted		
CFG [13] scale	{ 1.0, 2.5 }		

Table 4. Implementation details.

D. Illustration of Timestep Shift

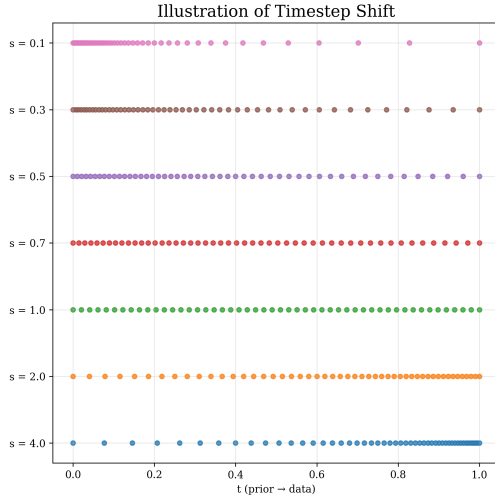


Figure 11. Illustration of timestep shift mapping $t_m = \frac{s t_n}{1 + (s-1)t_n}$, where s is a shift factor and t_n is the uniform sampling schedule. Intuitively, with $s < 1$, we allocate more steps to the early (navigation) timesteps, while $s > 1$ biases the schedule toward later (refinement) timesteps. The best gFID is achieved with $s = 0.5$.

E. Impact of Latent Space

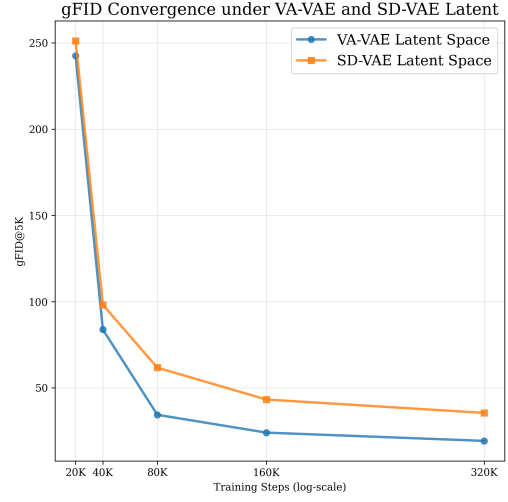


Figure 12. Convergence of gFID@5K when training rectified flow models under different latent spaces. We use LightningDiT-B/1 for VA-VAE [41] and LightningDiT-B/2 for SD-VAE [31] to align the training resolution to 16^2 . The training in the VA-VAE latent space converges faster, indicating a better latent space structure.

F. Quantifying Memorization Behaviors

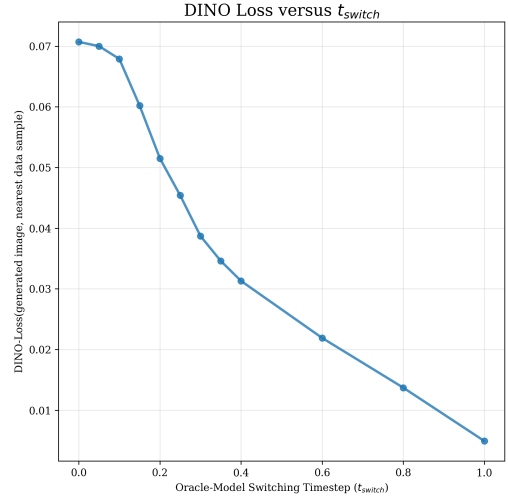


Figure 13. Quantitative results for oracle-model mixed generation. We report the DINO [28] loss (MSE of DINO self-attention maps, reflecting structural dissimilarity) between the generated images and the nearest training sample, measured across different switching timesteps (t_{switch}). Overall, we observe: (1) a sharp decline emerges after $t_{\text{switch}} \approx 0.1$ (i.e., the shift of training target); (2) when the loss falls below roughly 0.05, the generated layouts become closely aligned with those of the training samples (Fig. 14).

G. Additional Qualitative Results



Figure 14. Mixed sampling results with switch point t_{switch} . Oracle u_t^* is used before t_{switch} and LightningDiT-XL/1 [41] afterward. Overall, early switching yields diverse novel outputs (generalization), while late switching reproduces training samples (memorization). Despite minor variations across sampling trajectories, the empirical stage transition (i.e., reverting to training-like layouts) emerges around $t = 0.2$, slightly lagging behind the shift in the training target. Zoom in for the best view. Better view with color.

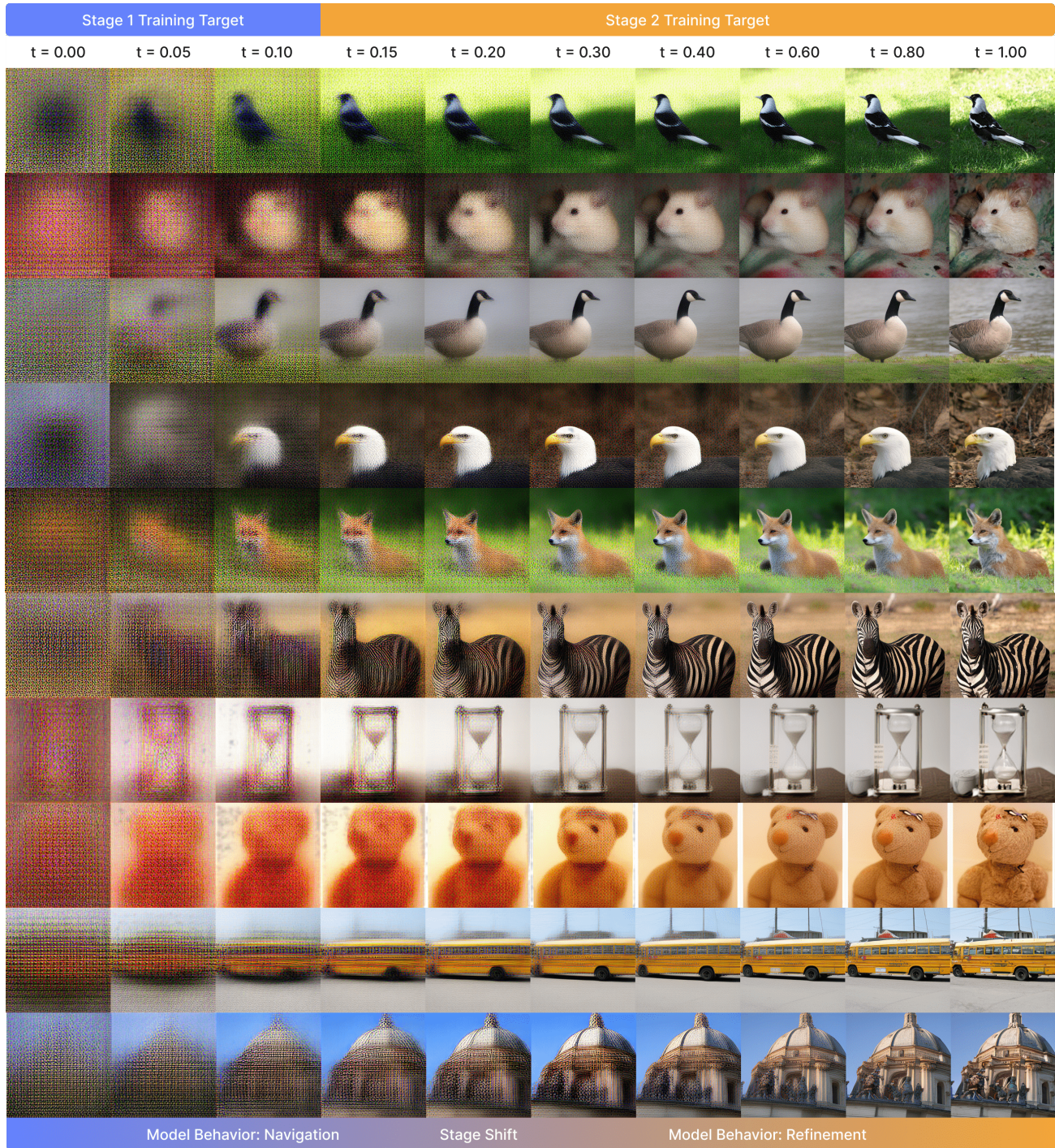


Figure 15. Intermediate predictions of a LightningDiT-XL/1 [41] model across timesteps. Overall, early stages primarily navigate global layout via smoothed predictions, while later stages refine fine-grained details. Zoom in for the best view. Better view with color.

H. Two-Stage Behaviors in Flux.1



Figure 16. Qualitative illustration of two-stage behavior in `Flux.1[dev]` [19]. Specifically, we first generate a reference latent z_{gt} via text-to-image sampling. Then, we re-noise it by interpolating with Gaussian noise at a chosen t_{resume} and resume sampling. Owing to Flux’s higher-dimensional latent space, the stage transition appears earlier than models trained on 256^2 ImageNet data, and the model can reliably recover nearly identical images even after $\sim 90\%$ noise corruption. We also note that Flux employs a non-uniform, resolution-aware timestep schedule and a different time convention; all t_{resume} values shown here are converted to our convention for consistency.