

# Flow Map Distillation Without Data

Shangyuan Tong<sup>\*</sup>  
MIT

sytong@csail.mit.edu

Nanye Ma<sup>\*</sup>  
NYU

nm3607@nyu.edu

Saining Xie<sup>†</sup>  
NYU

saining.xie@nyu.edu

Tommi Jaakkola<sup>†</sup>  
MIT

tommi@csail.mit.edu

## Abstract

*State-of-the-art flow models achieve remarkable quality but require slow, iterative sampling. To accelerate this, flow maps can be distilled from pre-trained teachers, a procedure that conventionally requires sampling from an external dataset. We argue that this data-dependency introduces a fundamental risk of **Teacher-Data Mismatch**, as a static dataset may provide an incomplete or even misaligned representation of the teacher’s full generative capabilities. This leads us to question whether this reliance on data is truly necessary for successful flow map distillation. In this work, we explore a data-free alternative that samples only from the prior distribution—a distribution the teacher is guaranteed to follow by construction, thereby circumventing the mismatch risk entirely. To demonstrate the practical viability of this philosophy, we introduce a principled framework that learns to predict the teacher’s sampling path while actively correcting for its own compounding errors to ensure high fidelity. Our approach surpasses all data-based counterparts and establishes a new state-of-the-art by a significant margin. Specifically, distilling from SiT-XL/2+REPA, our method reaches an impressive FID of **1.45** on ImageNet 256×256, and **1.49** on ImageNet 512×512, both with only 1 sampling step. We hope our work establishes a more robust paradigm for accelerating generative models and motivates the broader adoption of flow map distillation without data.*

Project page: [data-free-flow-distill.github.io](https://data-free-flow-distill.github.io)

## 1. Introduction

Diffusion models [22, 28, 73, 78] and flow models [1, 18, 40, 43, 56, 91, 92] have revolutionized high-fidelity synthesis [23, 60, 62, 86], yet their reliance on numerically integrating an Ordinary Differential Equation (ODE) creates a significant computational bottleneck. To resolve this latency, flow maps [5], which learn the solution operator of the ODE directly, offer a principled path to acceleration, bypassing iterative solving by taking large “jumps”

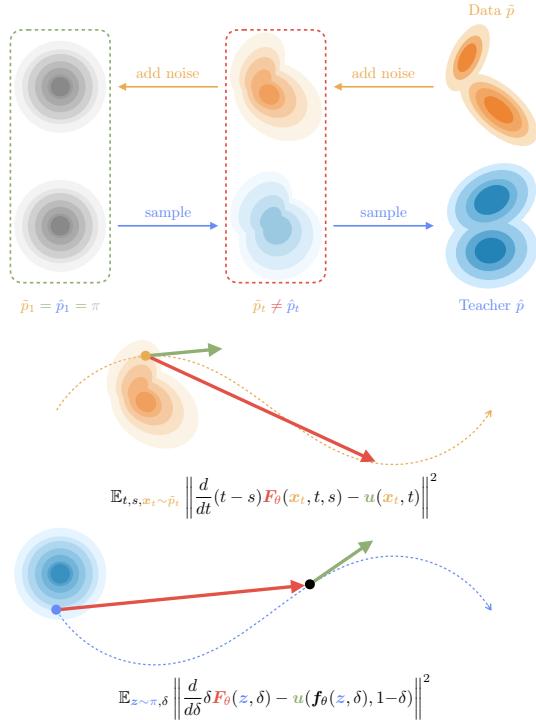


Figure 1. **Teacher-Data Mismatch and the data-free alternative.** (Top) Conventional data-based distillation relies on intermediate distributions ( $\tilde{p}_t$ ) derived from a static dataset, which could be misaligned with the teacher’s true generative distributions ( $\hat{p}_t$ ). (Bottom) The data-free paradigm, in contrast, samples only from the prior ( $\pi$ ), the single distribution with guaranteed alignment, thereby circumventing the mismatch risk by construction.

along the generative trajectory. While flow maps can be trained from scratch [4, 12, 14, 79], a more flexible alternative is to distill them from powerful, pre-trained teacher models [3, 32, 64, 65, 79]. This modular strategy allows for the compression of state-of-the-art models, which are often the product of advanced training [38, 100, 102] and post-training [39, 41, 84, 90, 101] techniques.

We observe that the dominant and most successful flow map distillation approaches are *data-based*, relying on samples from an external dataset to train the student. We argue

<sup>\*</sup>Equal contribution, <sup>†</sup>Equal advising

that this tacitly accepted dependency introduces a fundamental risk of **Teacher-Data Mismatch**. As illustrated in Fig. 1, a static dataset may provide an incomplete or misaligned representation of the teacher’s true generative capabilities. This discrepancy arises frequently in practice: when a teacher generalizes beyond its original training set [26, 27, 54, 58, 70, 75, 96, 99]; when post-hoc fine-tuning [39, 41, 84, 90, 101] shifts the teacher’s distribution away from the original data; or when the teacher’s proprietary training data is simply unavailable [6, 35, 62, 71, 87]. In these scenarios, forcing a student to match the teacher on a misaligned dataset fundamentally constrains its potential.

Fortunately, this mismatch is not inevitable. We observe that while the teacher’s generative paths may diverge from the dataset, they are, by definition, anchored to the prior distribution. As shown in Fig. 1, the prior serves as the single point of guaranteed alignment: it is the shared origin for the teacher’s generation and the termination point for any noising process. This insight leads us to question whether the common reliance on data is truly necessary. We posit that we can instead build a robust, *data-free* alternative by sampling only from the prior, thereby circumventing the mismatch risk entirely by construction.

To operationalize this philosophy, we introduce a principled framework designed to track the teacher’s dynamics purely from the prior. Our method takes a sample from the prior and a scalar integration interval, predicting where the flow should jump. We show that optimality is achieved when the model’s *generating velocity*, the rate at which it traverses its own path, aligns with the teacher’s instantaneous velocity. Nevertheless, like any autonomous numerical solver, this prediction process is susceptible to compounding errors. To mitigate this, we propose a correction mechanism that further aligns the model’s *noising velocity*, the marginal velocity of the noising flow implied by the student’s predicted distribution, back to the teacher. We name our proposal **FreeFlow**, emphasizing its defining characteristic as a completely data-free distillation framework for flow maps.

We validate our approach through extensive experiments on ImageNet [63]. Distilling from a SIT-XL/2+REPA [100] teacher, FreeFlow establishes a new state-of-the-art, reaching an impressive FID of **1.45** on  $256 \times 256$  and **1.49** on  $512 \times 512$  with 1 function evaluation (1-NFE), significantly surpassing all data-based baselines. Furthermore, by leveraging its nature as a fast, consistent proxy, FreeFlow enables efficient inference-time scaling [51, 72], allowing for the search of optimal noise samples in a single step. Ultimately, our findings confirm that an external dataset is not an essential requirement for high-fidelity flow map distillation, and the risk of Teacher-Data Mismatch can be avoided entirely without compromising performance. We believe this work provides a more robust foundation for generative model acceleration and motivates a shift toward the data-free paradigm.

## 2. Preliminaries

**Diffusion and flow.** Diffusion models [22, 28, 73, 78] and flow models [1, 18, 40, 43, 56, 91, 92] are trained to reverse a reference noising process that transports the data distribution  $p \equiv p_0$  to a easy-to-sample prior distribution  $\pi \equiv p_1$  like  $\mathcal{N}(\mathbf{0}, \mathbf{I})$ . We denote the interpolating distributions in between as  $p_t$  and their samples  $\mathbf{x}_t$ , indexed by time  $t \in [0, 1]$ . For the linear interpolation scheme [40, 43, 50] that we utilize throughout the paper, given a pair of terminal samples  $\mathbf{x} \sim p$  and  $\mathbf{z} \sim \pi$ , we construct the noising process from its interpolants,  $\mathbf{x}_t = \mathbf{I}_t(\mathbf{x}, \mathbf{z}) := (1 - t)\mathbf{x} + t\mathbf{z}$ , which in turn defines a conditional velocity, pointing in the direction from prior to data,  $\mathbf{u}(\mathbf{x}_t, t | \mathbf{x}, \mathbf{z}) := -\partial_t \mathbf{I}_t(\mathbf{x}, \mathbf{z}) = \mathbf{x} - \mathbf{z}$ . Taking the expectation over  $p$  and  $\pi$ , we arrive at the marginal instantaneous velocity,  $\mathbf{u} : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$ , a vector field that dictates how the samples evolve, which governs the noising process with the following ODE:

$$d\mathbf{x}(t) = -\mathbf{u}(\mathbf{x}(t), t)dt, \quad (1)$$

where  $\mathbf{x}(t) \in \mathbb{R}^d$  denotes the state of the system. In practice, the typically unknown  $\mathbf{u}$  can be well approximated by a model  $\mathbf{g}_\psi$  with parameters  $\psi$ , trained with denoising score matching [77, 83] or conditional flow matching [40, 56]:

$$\mathbb{E}_{\mathbf{x}, \mathbf{z}, t} \|\mathbf{g}_\psi(\mathbf{I}_t(\mathbf{x}, \mathbf{z}), t) - \mathbf{u}(\mathbf{I}_t(\mathbf{x}, \mathbf{z}), t | \mathbf{x}, \mathbf{z})\|^2. \quad (2)$$

For sampling, we need to solve Eq. (1) by integrating the flow backward in time:

$$\phi_{\mathbf{u}}(\mathbf{x}_t, t, s) = \mathbf{x}_t + \int_t^s -\mathbf{u}(\mathbf{x}(\tau), \tau)d\tau, \quad (3)$$

where  $\phi_{\mathbf{u}} : D_{\text{flow}} \rightarrow \mathbb{R}^d$ ,  $D_{\text{flow}} = \{(\mathbf{y}, \zeta, \xi) \mid \mathbf{y} \in \mathbb{R}^d, \zeta \in [0, 1], \xi \in [0, \zeta]\}$ , denotes the generating flow equipped with underlying velocity field  $\mathbf{u}$ , and  $t - s$  is the integration time interval. The standard sampling procedure of flow models corresponds to calculating  $\phi_{\mathbf{u}}(\mathbf{z}, 1, 0)$ ,  $\mathbf{z} \sim \pi$ . In practice, we resort to some numerical solver, like Euler [74, 78] and Heun methods [28]. Since the underlying trajectories typically exhibit complicated structure and curvature [88, 93], such numerical integration procedures often require dozens or even hundreds of NFEs for a single generation.

**Flow maps.** Instead of approximating the instantaneous velocity  $\mathbf{u}$ , a flow map model [3, 5, 12, 14, 32, 64, 65, 79],  $f_\theta$  parameterized by  $\theta$ , is trained to directly approximate  $\phi_{\mathbf{u}}$ . Existing works dissect and utilize key properties of  $\phi_{\mathbf{u}}$  to construct their training objective, typically via the local dynamics described by  $\mathbf{u}$ . For example, in Mean-Flow [14], the network  $\mathbf{F}_\theta : D_{\text{flow}} \rightarrow \mathbb{R}^d$  represents the average velocity  $\mathbf{f}$  travels over its path:  $\mathbf{f}_\theta(\mathbf{x}_t, t, s) = \mathbf{x}_t + (t - s)\mathbf{F}_\theta(\mathbf{x}_t, t, s)$ . At optimality, we know from Eq. (3)

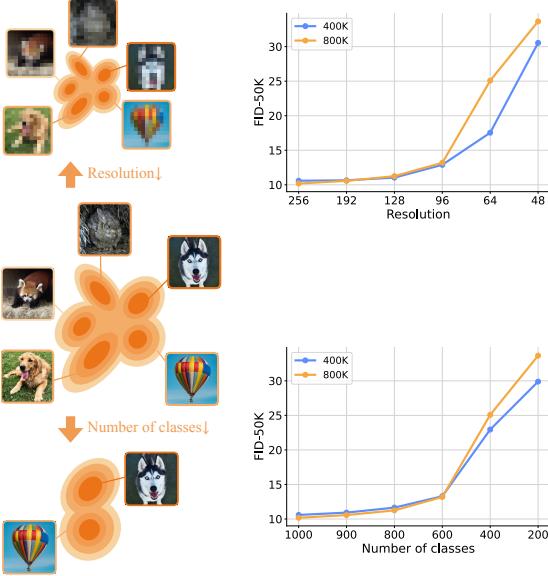


Figure 2. **Impact of Teacher-Data Mismatch.** With a fixed teacher model, increasing augmentation induces a more severe mismatch between teacher and data, degrading student performance.

that  $(t-s)\mathbf{F}_{\theta^*}(\mathbf{x}_t, t, s) = \int_t^s -\mathbf{u}(\mathbf{x}(\tau), \tau) d\tau$ . Differentiating both sides w.r.t.  $t$  leads to

$$\mathbf{F}_{\theta^*}(\mathbf{x}_t, t, s) + (t-s)\frac{d}{dt}\mathbf{F}_{\theta^*}(\mathbf{x}_t, t, s) = \mathbf{u}(\mathbf{x}_t, t), \quad (4)$$

where  $\frac{d}{dt}\mathbf{F}$  is the total derivative that can be further expanded into  $\nabla_{\mathbf{x}_t}\mathbf{F}\frac{d\mathbf{x}_t}{dt} + \partial_t\mathbf{F}$ . This identity then motivates the following practical training objective:

$$\mathbb{E}_{t,s,\mathbf{x}_t} \|\mathbf{F}_{\theta}(\mathbf{x}_t, t, s) - \text{sg}(\mathbf{u}_{\text{MF}})\|^2, \quad (5)$$

where  $\text{sg}(\cdot)$  denotes the stop-gradient operation, and  $\mathbf{u}_{\text{MF}} = \mathbf{u}(\mathbf{x}_t, t) - (t-s)\frac{d}{dt}\mathbf{F}_{\theta}(\mathbf{x}_t, t, s)$ . The decision to drop the remaining gradients is mostly empirical and common in prior literature [12, 14, 45, 76], as it results in faster, less resource-demanding, and often more stable training. Here,  $\mathbf{u}$  is either co-trained from scratch [14] together with Eq. (2), or a pre-trained flow model [57, 64, 79] in a procedure known as flow map distillation. In this paper, we focus on the second case, since this modular approach allows for easier incorporation of advanced training [38, 100, 102] and post-training [39, 41, 84, 90, 101] techniques.

### 3. With or Without Data?

The goal of flow map distillation is to create a student  $f_{\theta}$  that faithfully reproduces the full generative process of the given  $\phi_u$ , just with fewer NFEs. Intuitively, existing methods learn the teacher sampling dynamics at a series of intermediate states  $\mathbf{x}_t$  (note the expectation over  $\mathbf{x}_t$  in Eq. (5)). Our discussion begins from this very point, with a closer inspection of a foundational, yet largely unexamined, element: the underlying distribution from which these states  $\mathbf{x}_t$  are drawn.

### 3.1. The Risk of Teacher-Data Mismatch

The distribution of  $\mathbf{x}_t$  is conventionally formed by sampling from a *data-noising* distribution, which we denote as  $\tilde{p}_t$ . This is the set of all interpolants  $\mathbf{I}_t(\tilde{\mathbf{x}}, \mathbf{z})$  generated by taking a data point  $\tilde{\mathbf{x}} \sim \tilde{p}$ <sup>1</sup> and a prior sample  $\mathbf{z} \sim \pi$ . Although tacitly accepted, this practice implicitly assumes that  $\tilde{p}_t$  is a suitable representation of the states the teacher model follows over its sampling trajectory.

We note that the teacher defines a distinct set of intermediate states via Eq. (3). The set of all  $\mathbf{x}_t$  along these solution paths constitutes the true *teacher-generating* distribution, which we denote as  $\hat{p}_t$ . For the student to perfectly reproduce the teacher, it should be trained to match the teacher’s dynamics over  $\hat{p}_t$ . The central problem that we identify in this paper, which we term the **Teacher-Data Mismatch**, is that these two distributions are not equivalent:  $\tilde{p}_t \neq \hat{p}_t$ .

By training on  $\tilde{p}_t$ , the student is compelled to learn the teacher’s dynamics only on trajectories that are anchored to the static dataset  $\tilde{p}$ . Any generative behavior of the teacher that starts from  $\pi$  and evolves through states not well-represented by  $\tilde{p}_t$  will be systematically ignored during training. Consequently, even a perfectly converged student is not guaranteed to reproduce the teacher’s outputs, as it has fundamentally been trained to distill the wrong process.

To examine and validate the impact of the discussed mismatch, we design a controlled experiment on ImageNet, where we introduce deliberately misaligned  $\tilde{p}_t$  distributions by applying data augmentations during the training of conventional flow map distillation. As shown in Fig. 2, the quality of the learned flow map is highly sensitive to the fidelity and representativeness of the distillation dataset: stronger augmentation leads to a larger discrepancy between  $\hat{p}_t$  and  $\tilde{p}_t$ , and, in turn, results in a more significant degradation in student performance.

This mismatch is not merely a theoretical curiosity; it manifests in several common and critical scenarios. First, when a powerful teacher model has generalized beyond its training set [26, 27, 54, 58, 70, 75, 96, 99], or even when it simply employs the widely adopted Classifier-Free Guidance (CFG) [21], its generative distribution  $\hat{p}_0$  will contain novel, extrapolated samples not present in  $p$ , causing its trajectories  $\hat{p}_t$  to necessarily diverge from the data-noising paths  $p_t$ . Second, if the teacher has been altered by post-hoc fine-tuning [39, 41, 84, 90, 101], its generating flow is deliberately modified, again forcing  $\hat{p}_t$  to diverge from the original data-noising distribution. Third, a teacher model may be released publicly while its massive, proprietary training data is not [6, 35, 62, 71, 87]. In this case,  $p$  is simply unavailable, and any proxy dataset used will almost certainly create a severe mismatch.

<sup>1</sup>We use  $\tilde{p}$  to denote the dataset available at distillation time. As we will discuss, it may differ from  $p$ .



Figure 3. Selected samples from FreeFlow-XL/2 model at  $512 \times 512$  resolution with 1-NFE. More uncurated results are in App. D.

### 3.2. Towards a Data-Free Alternative

A straightforward remedy to the Teacher-Data Mismatch is to directly sample from  $\hat{p}_t$  during training. This would involve sampling  $z \sim \pi$ , integrating the teacher model from  $t = 1$  to a random time  $t$  to get  $x_t = \phi_u(z, 1, t)$ , and then using this  $x_t$  in the distillation loss. Just like the case of knowledge distillation [20, 43, 47, 104], where the model is trained to learn the fully integrated outcomes at  $t = 0$ , obtaining reference trajectories on-the-fly is prohibitively costly, whereas pre-computing them offline scales poorly. In short, for high-dimensional or complex conditional tasks, generating enough samples to adequately represent the underlying distribution simply becomes intractable.

This apparent impasse leads us to re-examine the properties of these two distributions. While  $\hat{p}_t$  and  $\tilde{p}_t$  diverge for  $t \in [0, 1)$ , they are, by construction, identical at  $t = 1$ . The data-noising process terminates at the prior distribution (*i.e.*,  $\tilde{p}_1 \equiv \pi$ ), and the teacher’s generative process begins at the same prior (*i.e.*,  $\hat{p}_1 \equiv \pi$ ). This observation provides a crucial foothold. The prior  $\pi$  is the one distribution we can sample from that is guaranteed to be on-distribution for the teacher’s generative process, completely circumventing the risk of Teacher-Data Mismatch.

This insight motivates our central question: Is the commonly followed data-dependency truly necessary for flow map distillation? We argue that it is not. In the following, we explore a data-free alternative, building a new flow map distillation objective governed only by the prior distribution.

## 4. Flow Map Distillation Without Data

Our exploration now moves from motivation to mechanism. A flow map model can be generally understood as directly modeling a segment of a full generative trajectory, and the core principle for training such a model is to enforce consistency with  $u$  at some point along this segment, ensuring the learned dynamics are locally correct. The segment’s two key points provide natural candidates: a sampled start-point,  $x_t$ , and a predicted end-point,  $x_s = f_\theta(x_t, t, s)$ .

This perspective provides a clear lens through which to view the distillation process. In the conventional, data-based setting, the start-point  $x_t$  is drawn from a series of data-

noising distributions  $\tilde{p}_t$ . It is thus natural to constrain the model by perturbing this start-point, which corresponds to differentiating the optimal condition,  $(t - s)F_{\theta^*}(x_t, t, s) = \int_1^s -u(x(\tau), \tau) d\tau$ , with respect to  $t$ . This operation leads to the MeanFlow identity [14] in Eq. (4), which effectively enforces consistency at the start of the segment.

In our data-free investigation, however, we only sample our start-point from the prior  $\pi$ , which fixes  $x_t = z$  at  $t = 1$ . Consequently, perturbing the start-time  $t$ , and in turn the start-point, is no longer a meaningful operation. Thus, we consider the symmetrical alternative: if we cannot enforce consistency by perturbing the sampled start-point, we can instead do so by perturbing the predicted end-point. This provides a different path to ensuring the student’s local dynamics are correct, and it corresponds to differentiating the optimal condition with respect to the end time  $s$ .

To formalize this, we first simplify our notation to reflect this prior-anchored ( $t = 1$ ) view. That is: (1) We define the integration duration as  $\delta = t - s = 1 - s$ , where  $\delta \in [0, 1]$ ; (2) The flow map  $f_\theta(z, \delta) : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$  approximates the true flow  $\phi_u(z, 1, 1 - \delta)$ ; (3) The average velocity  $F_\theta(z, \delta) : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$  is linked by the parameterization  $f_\theta(z, \delta) = z + \delta F_\theta(z, \delta)$ . The optimal condition, anchored at  $t = 1$ , thus reduced to:

$$\delta F_{\theta^*}(z, \delta) = \int_1^{1-\delta} -u(x(\tau), \tau) d\tau. \quad (6)$$

Following our exposition, we differentiate both sides of Eq. (6) w.r.t.  $\delta$  (equivalent to differentiating w.r.t.  $-s$ ):

$$F_{\theta^*}(z, \delta) + \delta \partial_\delta F_{\theta^*}(z, \delta) = u(f_{\theta^*}(z, \delta), 1 - \delta). \quad (7)$$

Eq. (7) differs from Eq. (4) in subtle ways: (1) The time derivative of  $F_\theta$  is just a partial derivative, as  $z$  does not depend on  $\delta$ ; (2)  $u$  is evaluated at a state predicted by  $f_\theta$ .

The identity defined in Eq. (7) provides a sufficient condition for optimality, which motivates the following loss:

$$\mathbb{E}_{z, \delta} \| F_\theta(z, \delta) - \text{sg}(u_{\text{target}}) \|^2, \quad (8)$$

where  $u_{\text{target}} = u(f_\theta(z, \delta), 1 - \delta) - \delta \partial_\delta F_\theta(z, \delta)$ . Remarkably, we verify that Eq. (8) is formulated by only sampling

from the prior  $\pi$ , without any reliance on an external dataset  $\tilde{p}$  and thus free from the risks of Teacher-Data Mismatch. Hence, it achieves the goal of our exploration.

#### 4.1. Predict With Generating Flows

We now analyze our proposed objective in Eq. (8). Note that  $\partial_\delta \mathbf{f}_\theta(\mathbf{z}, \delta)$ , the model prediction’s rate of change with respect to the integration time, is the velocity with which the model travels along its generating flow. Thus, the optimality of the student is equivalent to the alignment between the model’s *generating velocity* and the underlying velocity. Indeed, it is easy to see that the loss value of Eq. (8) is the same as  $\mathbb{E}_{\mathbf{z}, \delta} \|\partial_\delta \mathbf{f}_\theta(\mathbf{z}, \delta) - \mathbf{u}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)\|^2$ , which evaluates to 0 if and only if  $\partial_\delta \mathbf{f} = \mathbf{u}$ . Intuitively, the student is analogous to an autonomous ODE solver, which uses its current estimated state to query the derivative function and compute the next state. The student learns to “ride” the teacher’s vector field, starting from  $\pi$  and extending outward, step by step, based entirely on its own evolving predictions.

In practice,  $\partial_\delta \mathbf{F}_\theta$  can be calculated easily and efficiently via Jacobian-vector product (JVP) with forward-mode automatic differentiation, barring some advanced computation kernels, which currently require customized solutions [45]. Still, it is desirable to work with a more flexible loss function with no such limitations, which is why we derive a discrete-time alternative (detail in App. A.1) that numerically approximates  $\partial_\delta \mathbf{F}_\theta$  with finite differences.

Consequently, we abstract away the computation detail, and use the general notation  $\mathbf{v}_G(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)$  to denote the student’s generating velocity  $\partial_\delta \mathbf{f}_\theta(\mathbf{z}, \delta)$ . The understanding that Eq. (8) aligns  $\mathbf{v}_G$  and  $\mathbf{u}$  can also be observed from its optimization gradients:

$$\nabla_\theta \mathbb{E}_{\mathbf{z}, \delta} \left[ \mathbf{F}_\theta(\mathbf{z}, \delta)^\top \text{sg}\left( \Delta_{\mathbf{v}_G, \mathbf{u}}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta) \right) \right], \quad (9)$$

where  $\Delta_{\mathbf{v}_G, \mathbf{u}}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)$  is the difference between  $\mathbf{v}_G(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)$  and  $\mathbf{u}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)$ . Explicitly writing out the gradients makes it easier to adopt techniques like gradient weighting/normalization explored in Sec. 5.1. Note that, if  $\mathbf{F}_\theta$  is further parameterized, we should replace the first term in Eq. (9) with the actual network output to ensure effective gradient control by only modifying  $\Delta_{\mathbf{v}_G, \mathbf{u}}$ .

Lastly, we note that advanced sampling techniques can be easily incorporated in  $\mathbf{u}$ , e.g., for classifier-free guidance (CFG) [21], we could simply replace  $\mathbf{u}(\mathbf{x}_t, t \mid c)^2$ <sup>2</sup> with  $\mathbf{u}_\gamma(\mathbf{x}_t, t \mid c) = \gamma \cdot \mathbf{u}(\mathbf{x}_t, t \mid c) + (1 - \gamma) \cdot \mathbf{u}(\mathbf{x}_t, t \mid c = \emptyset)$ , where  $c$  is condition with  $\emptyset$  referring to a null input, and  $\gamma$  is the guidance strength. Furthermore, the model can be trained on a range of  $\gamma$  values, which enables the ability to effortlessly change the guidance strength at inference time.

**The Challenge of Error Accumulation.** In practice, the student model  $\mathbf{f}_\theta$  is only a learned approximation, not a

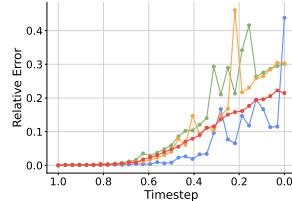


Figure 4. Approximation errors accumulate as the prediction proceeds from noise to data.

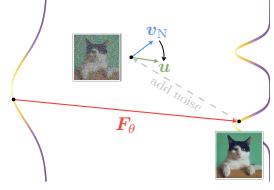


Figure 5. Correction objective in Eq. (11) aligns the student’s noising velocity  $v_N$  with  $u$ .

mathematically perfect one. At any  $\delta$ , its prediction  $\mathbf{f}_\theta(\mathbf{z}, \delta)$  may contain a small approximation error, placing it slightly off the true teacher trajectory  $\phi_u(\mathbf{z}, 1, 1 - \delta)$ . Because the objective is self-referential, this small deviation influences the target used for subsequent steps. The student queries the teacher at its current and potentially slightly erroneous state, and the resulting velocity target, while correct for that state<sup>3</sup>, may not guide the student back toward the true path. Such errors can compound as the integration proceeds from  $\delta = 0$  to  $\delta = 1$ . In Fig. 4, we measure the relative differences between the student’s predicted trajectory and the teacher’s true sampling path, which empirically quantifies and confirms such a phenomenon as the student progressively diverges from the teacher when  $\delta$  increases.

#### 4.2. Correct With Noising Flows

The problem identified above is that the student is trained to predict the next state based on its current one, but it has no means to correct its own deviations and pull the trajectory back towards the teacher’s true path. Drawing inspiration from Song et al. [78], we seek to correct the marginal distributions of the student solutions, analogous to a predictor-corrector method for solving ODEs [2]. Additionally, the correction objective cannot reintroduce the data-dependency we have worked to remove, meaning that we do not consider objectives like GANs [15] that rely on an external dataset.

Variational Score Distillation [85] was originally proposed as a training procedure to distill distributions from pre-trained diffusion models by minimizing the Integral KL divergence [48]. We slightly adapt it to our setting, where we use  $q \equiv q_0$  to denote the marginal distribution of clean samples generated by the model,  $\int \mathbf{f}_\theta(\mathbf{z}, 1) d\pi$ . Specifically, it has been shown that  $q = p$  if and only if their IKL divergence is 0, which is defined as

$$D_{\text{IKL}}(q \parallel p) := \int_0^1 \mathbb{E}_{\mathbf{x}_r \sim q_r} \left[ \log \frac{q_r(\mathbf{x}_r)}{p_r(\mathbf{x}_r)} \right] dr, \quad (10)$$

where  $q_r$  and  $p_r$  are the marginal interpolating distributions, following the same noising process constructed by  $\mathcal{I}$ .

The optimization gradient of Eq. (10) w.r.t.  $\theta$  is  $\mathbb{E}_{r, \mathbf{x}_r} \left[ (\nabla_\theta \mathbf{x}_r)^\top (\nabla_{\mathbf{x}_r} \log q_r(\mathbf{x}_r) - \nabla_{\mathbf{x}_r} \log p_r(\mathbf{x}_r)) \right]$ . As

<sup>2</sup>We omit  $c$  everywhere else in the paper for simplicity.

<sup>3</sup>This assumes the teacher is perfect, which is not true in practice.

the score functions are interchangeable with the marginal velocities [50], we can optimize with the following gradient instead for correcting the student’s prediction:

$$\nabla_{\theta} \mathbb{E}_{z, n, r} \left[ F_{\theta}(z, 1)^{\top} \text{sg} \left( \Delta_{v_N, u}(\mathbf{I}_r(f_{\theta}(z, 1), n), r) \right) \right], \quad (11)$$

where  $n$  is sampled from the prior  $\pi$  like  $z$ . We verify that Eq. (11) is also formulated by only sampling from  $\pi$ , free from the risks of Teacher-Data Mismatch. Here,  $v_N$  denotes the marginal velocity of the noising flow constructed from the generated distribution  $q$  with the interpolating function  $\mathbf{I}$ , and  $\Delta_{v_N, u}$  is the difference between  $v_N$  and  $u$ . We illustrate the high-level understanding of this mechanism in Fig. 5.

Since  $v_N$  is unknown, we approximate it with another online network  $g_{\psi}$ <sup>4</sup>, full-parameter [48, 95, 97, 98, 107, 108] or LoRA [24, 53, 85], with loss in Eq. (2). More specifically, for a pair of samples  $f_{\theta}(z, 1)$  and  $n$ , the conditional noising velocity is  $-\partial_r \mathbf{I}_r(f_{\theta}(z, 1), n)$ , and we arrive at  $v_N$  by taking the expectation over  $z \sim \pi$  and  $n \sim \pi$ :

$$\mathbb{E}_{z, n, r} \| g_{\psi}(\mathbf{I}_r(f_{\theta}(z, 1), n), r) + \partial_r \mathbf{I}_r(f_{\theta}(z, 1), n) \|^2. \quad (12)$$

We highlight the similarity of the gradient forms between Eq. (9) and Eq. (11). Consequently, we identify that the optimality of the student is also equivalent to the alignment between the model’s *noising velocity* and the underlying velocity. That is, Eq. (10) evaluates to 0 if and only if  $\Delta_{v_N, u} = 0$ . Such a velocity alignment perspective offers a series of new understandings, which provide the essential reasoning behind the practical design choices discussed later in Sec. 5.1. We further note that a comprehensive correction procedure should correct the full predicted trajectory of the student, rather than only the end sample considered in Eq. (10). However, we do not find such a design helpful in the experiment settings we considered in Sec. 5.

## 5. Experiments

We empirically validate our proposed method on ImageNet [63] at  $256 \times 256$  and  $512 \times 512$  resolutions using FID-50K [19], with implementation details provided in App. B.

### 5.1. Design Decisions

We analyze each design choice through targeted qualitative and quantitative studies, presenting key findings in the main text and deferring full analyses to the App. A. Unless specified otherwise, we adopt the DiT-B/2 architecture [55], use the pre-trained SiT-B/2 [50] as  $u$  and student initialization, train the model for 400K iterations (roughly equivalent to 80 epochs with a batch size of 256) with uniformly sampling

<sup>4</sup>Further parameterizations on  $g_{\psi}$  are permissible.

$r$ sampling; Eq. (11)	FID $\downarrow$	$r$ range; Eq. (11)	FID $\downarrow$
LogitNormal(-0.4, 1.6)	6.24	[0, 0.6]	91.82
LogitNormal(0.0, 1.6)	5.95	[0, 0.7]	24.62
LogitNormal(0.4, 1.6)	5.78	[0, 0.8]	9.00
LogitNormal(0.8, 1.6)	5.63	[0, 0.9]	6.64
LogitNormal(1.2, 1.6)	5.78	[0, 1.0]	6.02

(a) <b><math>r</math> sampling.</b> More emphasis on higher noise levels leads to better results when aligning $v_N$ and $u$ .	(b) <b><math>r</math> range.</b> With a uniform distribution over $r$ , dropping higher noise levels leads to worse results.
interval; Eq. (11)	FID $\downarrow$
[0, 0.5]	7.09
[0, 0.6]	5.72
[0, 0.7]	5.63
[0, 0.8]	6.44
[0, 0.9]	7.41
[0, 1.0]	8.65

objective	$k$	FID $\downarrow$
Eq. (9)	0.0	11.91
Eq. (9)	0.5	11.71
Eqs. (9) and (11)	1.0	12.40
Eqs. (9) and (11)	0.0	43.53
Eqs. (9) and (11)	0.5	10.58
Eqs. (9) and (11)	1.0	5.58

(c) **Guidance interval.** Compared to teacher sampling, a more aggressive guidance interval is better.

(d) **Gradient weighting.** With both objectives, stronger decay on  $\Delta_{v_G, u}$  leads to better training.

Table 1. Empirical investigations of various design decisions.

$\gamma \in [1, 2]$ , and evaluate with the best  $\gamma = 2$ . We begin with designs specific to training with one of Eqs. (9) and (11), followed by studies on how to properly combine the two.

**Sampling of  $r$  in Eq. (11).** Traditionally, within the diffusion framework and from the divergence minimization perspective, the sampling of  $r$  in Eq. (11) often follows a uniform distribution over the discretized steps designed for generation [59, 85, 98]. Here, we provide a new perspective on the Eq. (11), which drives our proposal on the sampling distribution of  $r$ . Recall that the optimality of our correction objective is  $\Delta_{v_N, u} = 0$ , the alignment between the noising velocity of the model’s generated distribution and the underlying velocity. The velocity fields induce a pair of continuity equations  $\partial_t p_t(x) = -\nabla_x \cdot (p_t(x)u(x, t))$  and  $\partial_t q_t(x) = -\nabla_x \cdot (q_t(x)v_N(x, t))$ , which dictate the evolution of  $p$  and  $q$ . We note that  $p_1 = q_1 = \pi$  by construction, and the gap between  $p_0$  and  $q_0$  can be understood as the time-integrated accumulation of the differences in their corresponding probability fluxes. This understanding suggests we place a greater emphasis on higher noise levels, and we empirically validate this intuition in Tabs. 1a and 1b.

**Handling guidance in Eq. (11).** The usual treatment for including guidance in Eq. (11) is the same as in Eq. (9): replacing it with the guided velocity  $u_{\gamma}$ . However, we highlight that there is a subtle but major difference between  $v_G$  and  $v_N$ . In a traditional flow model training with Eq. (2), we essentially train the model to learn a dataset’s noising velocity, and use it as the generating velocity during sampling, *i.e.*, the two velocities are identical as they describe the same process. However, this equivalence no longer holds in the presence of techniques like CFG [29, 103], and the distinction is especially prominent at high noise levels. We resort to dropping the guidance application at high noise levels, in a similar fashion to the guidance interval [34] used

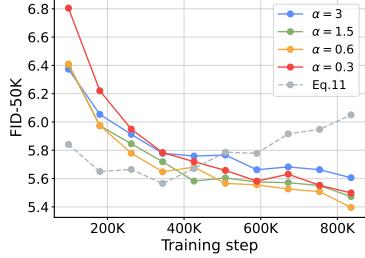
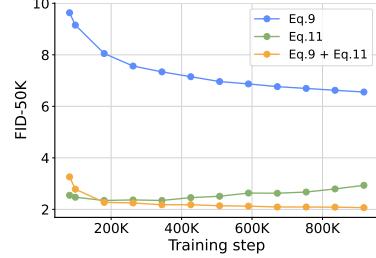
Figure 6. Performance is robust across  $\alpha$ .

Figure 7. Synergy between Eqs. (9) and (11).

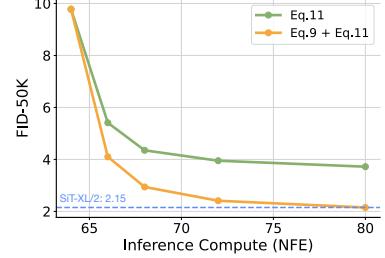


Figure 8. Inference-time scaling.

for flow sampling. Furthermore, as demonstrated in Tab. 1c, we stress that their empirical behaviors are different, and one typically needs to limit the interval significantly more aggressively (the pre-trained SiT-B/2 does not benefit from guidance interval with  $\gamma = 2$ ) in our correction objective.

**Adaptive gradient balancing between Eqs. (9) and (11).** We now discuss how to fuse the training signals from Eqs. (9) and (11). First, in a similar manner to prior works [12, 14], we decide to split the training batch between the prediction and correction objectives (75% and 25%, respectively), since correction is slightly more expensive compute-wise. Then, we adopt an adaptive gradient balancing strategy [10], where the correction gradients are scaled by some dynamic weight  $\lambda$  before concatenating with the prediction gradients. With the form similarity between Eq. (9) and Eq. (11), and the observation that both optimizations lack aleatoric uncertainty [31]<sup>5</sup>, we design  $\lambda = \alpha \frac{\mathbb{E}\|\Delta_{v_G, u}\|}{\mathbb{E}\|\Delta_{v_N, u}\| + \epsilon}$ , where the expectation is taken over the mini-batch and  $\epsilon = 10^{-6}$  is used for numerical stability. We show that the model performance is robust across a wide range of  $\alpha$  in Fig. 6.

**Gradient norm manipulations in Eq. (9).** We note that we can change the magnitude of  $\Delta_{v_G, u}$  freely without changing the optimal solution, which can also be understood as changing the loss metrics [13]. Concretely, we can scale  $\Delta_{v_G, u}$  with per-sample positive weights. Prior works [13, 14, 76] mostly explored scaling with  $1/(\|\Delta_{v_G, u}\|^2 + \epsilon)^k$ , where  $\epsilon = 10^{-4}$  is used for numerical stability and we vary the power term  $k$ . Since the actual weighting applied depends on the original norm of  $\Delta_{v_G, u}$ , which changes with the dimension of data  $d$ , we first divide it by  $\sqrt{d}$  so that it is dimension invariant before calculating the weight. We note that such a weighting design corresponds to applying a power-law decay on  $\Delta_{v_G, u}$ , and a larger  $k$  indicates a stronger decay, effectively dampening the gradient contributions. In Tab. 1d, we find that the model prefers weightings with a stronger decay when training with both Eqs. (9) and (11). We hypothesize that this is because their signals may not always agree with each other in practice, and by applying a dampener on  $\Delta_{v_G, u}$ , we mitigate the conflict between the two objectives and promote a more harmonious joint optimization.

<sup>5</sup> $\Delta_{v_G, u}$  and  $\Delta_{v_N, u}$  represent the quality of alignments, unlike Eq. (2).

## 5.2. Main Results

**Comparisons with prior work.** In Tab. 2, we benchmark our approach against existing proposals for learning fast flows on class-conditional ImageNet [63] generation at both  $256 \times 256$  and  $512 \times 512$  resolutions. We highlight three key findings from our main results. (1) **Our method achieves state-of-the-art performance by a significant margin.** Distilling from SiT-XL/2+REPA [100], our method reaches an impressive FID of **1.45** on  $256 \times 256$  and **1.49** on  $512 \times 512$ , greatly outperforming prior proposals. (2) **Competitive performance is realized very early in training.** Our model surpasses the final performance of many strong baselines after only 100K iterations ( $\approx 20$  epochs), demonstrating exceptional training efficiency. (3) **Our student faithfully reproduces the teacher’s full capabilities with 1-NFE.** Our distilled model consistently stays within 10% of the teacher’s original performance with only a single step, even when the teacher employs advanced training techniques like REPA [100], and sampling techniques like guidance intervals [34]. This is a critical advantage over training fast flows from scratch: it allows us to seamlessly inherit the benefits of complex teacher training recipes, which can be complicated to adopt [36], simply by distilling the final product. Crucially, we emphasize that our results are achieved using *only* the pre-trained teacher model, without requiring access to a single sample from an external dataset, real or synthetic, thus validating the effectiveness of the data-free paradigm.

**Synergy between prediction and correction.** While theory suggests that either learning the flow trajectories (Eq. (9)) or matching the marginal distributions (Eq. (11)) could suffice for generation, we find that neither is robust in isolation. The prediction objective, when used alone, falls victim to the error accumulation identified in Sec. 4.1, plateauing at suboptimal fidelity (blue line in Fig. 7). Conversely, training only with the correction objective (Eq. (11)) leads to gradual mode collapse and performance degradation (green line in Fig. 7; also gray baseline in Fig. 6). Fig. 7 illustrates the powerful synergy realized by our framework on a SiT-XL/2 teacher. By combining both signals at their optimal settings, we achieve performance strictly superior to either independent component. The prediction signals construct

Table 2. **Class-conditional generation on ImageNet  $256 \times 256$  and  $512 \times 512$ .** \* indicates the use of AutoGuidance [29]. Methods marked with  $\dagger$  are initialized from pretrained models. Crucially, unlike other listed distillation baselines, ours is constructed to be entirely data-free.

Class-Conditional ImageNet $256 \times 256$					Class-Conditional ImageNet $512 \times 512$									
Method	Epochs	#Params	NFE $\downarrow$	FID $\downarrow$	Method	Epochs	#Params	NFE $\downarrow$	FID $\downarrow$					
<i>Teacher Diffusion / Flow Models</i>														
SiT-XL/2 [50]	1400	675M	$250 \times 2$	2.06	SiT-XL/2 [50]	600	675M	$250 \times 2$	2.62					
SiT-XL/2+REPA [100]	800	675M	434	1.37	SiT-XL/2+REPA [100]	400	675M	460	1.37					
<i>Fast Flow from scratch</i>														
Shortcut-XL/2 [12]	250	675M	1 128	10.60 3.80	EDM2-S* [29]	1678	280M	$63 \times 2$	1.34					
IMM-XL/2 [106]	3840	675M	$1 \times 2$ $8 \times 2$	7.77 1.99	EDM2-XXL [30]	734	1.5B	82	1.40					
STEI [42]	1420 $\dagger$	675M	1 8	7.12 1.96	EDM2-XXL* [29]			$63 \times 2$	1.25					
MeanFlow-XL/2 [14]	240 1000	676M	1 2	3.43 2.20	<i>Fast Flow from scratch</i>									
DMF-XL/2 [36]	880 $\dagger$	675M	1 4	2.16 1.51	sCT-XXL [45]	761 $\dagger$	1.5B	1 2	4.29 3.76					
<i>Fast Flow by distillation</i>														
<i>Teacher: SiT-XL/2 (FID = 2.06)</i>														
SDEI [42]	20	675M	8	2.46	DMF-XL/2 [36]	540 $\dagger$	675M	1 4	2.12 1.68					
FACM [57]	–	675M	2	2.07	<i>Teacher: EDM2-S* (FID = 1.34)</i>									
FreeFlow-XL/2	20 300	678M	1 1	2.24 <b>1.69</b>	AyF-S [64]	80	280M	1 4	3.32 1.70					
<i>Teacher: SiT-XL/2+REPA (FID = 1.37)</i>					<i>Teacher: EDM2-XXL (FID = 1.40)</i>									
FACM [57]	–	675M	2	1.52	sCD-XXL [45]	320		1 2	2.28 1.88					
$\pi$ -Flow [7]	448	675M	1 2	2.85 1.97	sCD-XXL+VSD [45]	32	1.5B	1 2	2.16 1.89					
FreeFlow-XL/2	20 300	678M	1 1	1.84 <b>1.45</b>	<i>Teacher: SiT-XL/2 (FID = 2.62)</i>									
<i>Teacher: SiT-XL/2+REPA (FID = 1.37)</i>														
FreeFlow-XL/2	20 200	678M	1	3.01 <b>2.25</b>	FreeFlow-XL/2	20 200	678M	1	2.11 <b>1.49</b>					

the generative path, while the correction signals act as a stabilizer to rectify compounding errors, ensuring consistent improvement throughout training.

**Inference-time scaling.** The recently proposed inference-time scaling framework [51, 72] offers a promising avenue to trade additional compute for generation quality. However, existing search strategies typically require the full integration of  $\phi_u$  for every candidate, making the search process prohibitively expensive. We propose a more efficient alternative: by distilling the teacher into a flow map, we create a fast proxy that retains the teacher’s mapping from noise to data. This allows us to conduct the expensive search using the cheap, one-step student, transferring only the optimal noise to the teacher for final generation. We investigate a Best-of- $N$  search with an oracle verifier [51], employing our student (trained for only 20 epochs) to guide the fixed SiT-XL/2 teacher. As shown in Fig. 8, this approach drastically improves the teacher’s sampling quality. Crucially, the results highlight the benefit of our prediction objective (Eq. (9)): while the correction-only model (Eq. (11)) yields improvements, it is notably less efficient at identifying transferable noise candidates due to a lack of guaranteed trajectory alignment. Our combined objective, by enforcing strict

consistency with  $\phi_u$ , enables a much more effective search. With a total budget of only 80 NFEs, our method outperforms the teacher’s standard classifier-free guidance sampling at 128 NFEs. This result demonstrates a powerful practical trade-off: by shifting a fraction of the inference burden to a short distillation phase, we enable the compute-efficient deployment of large-scale diffusion models.

## 6. Conclusion

In this work, we challenge the conventional reliance on external datasets for flow map distillation. We identify a fundamental vulnerability in this practice, the Teacher-Data Mismatch, and argue that a static dataset is an inherently unreliable proxy for a teacher’s full generative capabilities. This data-dependency is not only risky but also unnecessary. Our principled, data-free alternative, formulated as a predictor-corrector framework, resolves this mismatch by construction as it samples only from the prior. Our strong empirical results, which establish a new state-of-the-art, confirm the practical viability and strength of this data-free paradigm. We believe this work provides a more robust foundation for accelerating generative models and hope it motivates a broader exploration of flow map distillation without data.

## Acknowledgments

We are grateful to Kaiming He for valuable discussions and feedback on the manuscript. This work was partly supported by the Google TPU Research Cloud (TRC) program and the Google Cloud Research Credits program (GCP19980904). ST and TJ acknowledge support from the Machine Learning for Pharmaceutical Discovery and Synthesis (MLPDS) consortium, the DTRA Discovery of Medical Countermeasures Against New and Emerging (DOMANE) threats program, the NSF Expeditions grant (award 1918839) Understanding the World Through Code. SX acknowledges support from the MSIT IITP grant (RS-2024-00457882) and the NSF award IIS-2443404.

## 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, 2
- [2] Eugene L Allgower and Kurt Georg. *Numerical continuation methods: an introduction*. Springer Science & Business Media, 2012. 5
- [3] David Berthelot, Arnaud Autef, Jierui Lin, Dian Ang Yap, Shuangfei Zhai, Siyuan Hu, Daniel Zheng, Walter Talbott, and Eric Gu. Tract: Denoising diffusion models with transitive closure time-distillation. *arXiv preprint arXiv:2303.04248*, 2023. 1, 2, 18
- [4] Nicholas M Boffi, Michael S Albergo, and Eric Vanden-Eijnden. How to build a consistency model: Learning flow maps via self-distillation. *arXiv preprint arXiv:2505.18825*, 2025. 1
- [5] Nicholas Matthew Boffi, Michael Samuel Albergo, and Eric Vanden-Eijnden. Flow map matching with stochastic interpolants: A mathematical framework for consistency models. *Transactions on Machine Learning Research*, 2025. 1, 2, 14, 18
- [6] Siyu Cao, Hangting Chen, Peng Chen, Yiji Cheng, Yutao Cui, Xinchi Deng, Ying Dong, Kipper Gong, Tianpeng Gu, Xiusen Gu, et al. Hunyuaniimage 3.0 technical report. *arXiv preprint arXiv:2509.23951*, 2025. 2, 3
- [7] Hansheng Chen, Kai Zhang, Hao Tan, Leonidas Guibas, Gordon Wetzstein, and Sai Bi. pi-flow: Policy-based few-step generation via imitation distillation. *arXiv preprint arXiv:2510.14974*, 2025. 8, 18
- [8] Earl A Coddington, Norman Levinson, and T Teichmann. Theory of ordinary differential equations, 1956. 13
- [9] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in neural information processing systems*, 34:8780–8794, 2021. 16, 17
- [10] Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 12873–12883, 2021. 7
- [11] 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. 14
- [12] Kevin Frans, Danijar Hafner, Sergey Levine, and Pieter Abbeel. One step diffusion via shortcut models. *arXiv preprint arXiv:2410.12557*, 2024. 1, 2, 3, 7, 8
- [13] Zhengyang Geng, Ashwini Pokle, William Luo, Justin Lin, and J Zico Kolter. Consistency models made easy. *arXiv preprint arXiv:2406.14548*, 2024. 7, 18
- [14] 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. 1, 2, 3, 4, 7, 8, 18
- [15] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 5
- [16] Jiatao Gu, Shuangfei Zhai, Yizhe Zhang, Lingjie Liu, and Joshua M Susskind. Boot: Data-free distillation of denoising diffusion models with bootstrapping. In *ICML 2023 Workshop on Structured Probabilistic Inference & Generative Modeling*, 2023. 18
- [17] Jonathan Heek, Emiel Hoogeboom, and Tim Salimans. Multistep consistency models. *arXiv preprint arXiv:2403.06807*, 2024. 18
- [18] Eric Heitz, Laurent Belcour, and Thomas Chambon. Iterative  $\alpha$ -(de) blending: A minimalist deterministic diffusion model. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–8, 2023. 1, 2
- [19] 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, 15, 16
- [20] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015. 4
- [21] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. 3, 5
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020. 1, 2, 13
- [23] Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3d. In *International conference on machine learning*, pages 8867–8887. PMLR, 2022. 1
- [24] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 6
- [25] Aapo Hyvärinen and Peter Dayan. Estimation of non-normalized statistical models by score matching. *Journal of Machine Learning Research*, 6(4), 2005. 15
- [26] Zahra Kadkhodaie, Florentin Guth, Eero P Simoncelli, and Stéphane Mallat. Generalization in diffusion models arises from geometry-adaptive harmonic representations. *arXiv preprint arXiv:2310.02557*, 2023. 2, 3

- [27] Mason Kamb and Surya Ganguli. An analytic theory of creativity in convolutional diffusion models. *arXiv preprint arXiv:2412.20292*, 2024. 2, 3
- [28] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in neural information processing systems*, 35:26565–26577, 2022. 1, 2, 14, 16, 17
- [29] Tero Karras, Miika Aittala, Tuomas Kynkänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine. Guiding a diffusion model with a bad version of itself. *Advances in Neural Information Processing Systems*, 37:52996–53021, 2024. 6, 8
- [30] Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24174–24184, 2024. 8, 16
- [31] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? *Advances in neural information processing systems*, 30, 2017. 7
- [32] Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong He, Yuki Mitsuji, and Stefano Ermon. Consistency trajectory models: Learning probability flow ODE trajectory of diffusion. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 2, 18
- [33] Ba Jimmy Kingma, Diederik P. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 16, 17
- [34] Tuomas Kynkänniemi, 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. 6, 7
- [35] 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. 2, 3
- [36] Kyungmin Lee, Sihyun Yu, and Jinwoo Shin. Decoupled meanflow: Turning flow models into flow maps for accelerated sampling. *arXiv preprint arXiv:2510.24474*, 2025. 7, 8, 18
- [37] Sangyun Lee, Yilun Xu, Tomas Geffner, Giulia Fanti, Karsten Kreis, Arash Vahdat, and Weili Nie. Truncated consistency models. *arXiv preprint arXiv:2410.14895*, 2024. 18
- [38] Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng. Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint arXiv:2504.10483*, 2025. 1, 3
- [39] Junzhe Li, Yutao Cui, Tao Huang, Yinping Ma, Chun Fan, Miles Yang, and Zhao Zhong. Mixgrpo: Unlocking flow-based grpo efficiency with mixed ode-sde. *arXiv preprint arXiv:2507.21802*, 2025. 1, 2, 3
- [40] 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, 15
- [41] Jie Liu, Gongye Liu, Jiajun Liang, Yangguang Li, Jiaheng Liu, Xintao Wang, Pengfei Wan, Di Zhang, and Wanli Ouyang. Flow-grpo: Training flow matching models via online rl. *arXiv preprint arXiv:2505.05470*, 2025. 1, 2, 3
- [42] Wenze Liu and Xiangyu Yue. Learning to integrate diffusion odes by averaging the derivatives. *arXiv preprint arXiv:2505.14502*, 2025. 8, 18
- [43] 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, 4, 15
- [44] Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, et al. Instaflow: One step is enough for high-quality diffusion-based text-to-image generation. In *The Twelfth International Conference on Learning Representations*, 2023. 18
- [45] Cheng Lu and Yang Song. Simplifying, stabilizing and scaling continuous-time consistency models. *arXiv preprint arXiv:2410.11081*, 2024. 3, 5, 8, 18
- [46] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787, 2022. 14
- [47] Eric Luhman and Troy Luhman. Knowledge distillation in iterative generative models for improved sampling speed. *arXiv preprint arXiv:2101.02388*, 2021. 4, 18
- [48] Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff-instruct: A universal approach for transferring knowledge from pre-trained diffusion models. *Advances in Neural Information Processing Systems*, 36:76525–76546, 2023. 5, 6, 15, 18
- [49] Weijian Luo, Zemin Huang, Zhengyang Geng, J Zico Kolter, and Guo-jun Qi. One-step diffusion distillation through score implicit matching. *Advances in Neural Information Processing Systems*, 37:115377–115408, 2024. 18
- [50] Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant transformers. *arXiv preprint arXiv:2401.08740*, 2024. 2, 6, 8, 15, 16, 17
- [51] Nanye Ma, Shangyuan Tong, Haolin Jia, Hexiang Hu, Yu-Chuan Su, Mingda Zhang, Xuan Yang, Yandong Li, Tommi Jaakkola, Xuhui Jia, et al. Inference-time scaling for diffusion models beyond scaling denoising steps. *arXiv preprint arXiv:2501.09732*, 2025. 2, 8
- [52] Chenlin Meng, Robin Rombach, Ruiqi Gao, Diederik Kingma, Stefano Ermon, Jonathan Ho, and Tim Salimans. On distillation of guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14297–14306, 2023. 18
- [53] Thuan Hoang Nguyen and Anh Tran. Swiftbrush: One-step text-to-image diffusion model with variational score distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7807–7816, 2024. 6

- [54] Matthew Niedoba, Berend Zwartsenberg, Kevin Murphy, and Frank Wood. Towards a mechanistic explanation of diffusion model generalization. *arXiv preprint arXiv:2411.19339*, 2024. 2, 3
- [55] William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 4195–4205, 2023. 6, 16
- [56] Stefano Peluchetti. Non-denoising forward-time diffusions. *arXiv preprint arXiv:2312.14589*, 2023. 1, 2
- [57] Yansong Peng, Kai Zhu, Yu Liu, Pingyu Wu, Hebei Li, Xiaoyan Sun, and Feng Wu. Flow-anchored consistency models. *arXiv preprint arXiv:2507.03738*, 2025. 3, 8, 18
- [58] Jakiw Pidstrigach. Score-based generative models detect manifolds. *Advances in Neural Information Processing Systems*, 35:35852–35865, 2022. 2, 3
- [59] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022. 6
- [60] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr, 2021. 1
- [61] Yuxi Ren, Xin Xia, Yanzuo Lu, Jiacheng Zhang, Jie Wu, Pan Xie, Xing Wang, and Xuefeng Xiao. Hyper-sd: Trajectory segmented consistency model for efficient image synthesis. *Advances in Neural Information Processing Systems*, 37: 117340–117362, 2024. 18
- [62] 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. 1, 2, 3, 16
- [63] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115:211–252, 2015. 2, 6, 7
- [64] Amirrejiba Sabour, Sanja Fidler, and Karsten Kreis. Align your flow: Scaling continuous-time flow map distillation. *arXiv preprint arXiv:2506.14603*, 2025. 1, 2, 3, 8, 18
- [65] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022. 1, 2, 18
- [66] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Advances in neural information processing systems*, 29, 2016. 16
- [67] Tim Salimans, Thomas Mensink, Jonathan Heek, and Emiel Hoogeboom. Multistep distillation of diffusion models via moment matching. *Advances in Neural Information Processing Systems*, 37:36046–36070, 2024. 18
- [68] Axel Sauer, Frederic Boesel, Tim Dockhorn, Andreas Blattmann, Patrick Esser, and Robin Rombach. Fast high-resolution image synthesis with latent adversarial diffusion distillation. In *SIGGRAPH Asia 2024 Conference Papers*, pages 1–11, 2024.
- [69] Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion distillation. In *European Conference on Computer Vision*, pages 87–103. Springer, 2024. 18
- [70] Christopher Scarvelis, Haitz Sáez de Ocáriz Borde, and Justin Solomon. Closed-form diffusion models. *arXiv preprint arXiv:2310.12395*, 2023. 2, 3
- [71] Team Seedream, Yumpeng Chen, Yu Gao, Lixue Gong, Meng Guo, Qiushan Guo, Zhiyao Guo, Xiaoxia Hou, Weilin Huang, Yixuan Huang, et al. Seedream 4.0: Toward next-generation multimodal image generation. *arXiv preprint arXiv:2509.20427*, 2025. 2, 3
- [72] Raghav Singhal, Zachary Horvitz, Ryan Teehan, Mengye Ren, Zhou Yu, Kathleen McKeown, and Rajesh Ranganath. A general framework for inference-time scaling and steering of diffusion models. *arXiv preprint arXiv:2501.06848*, 2025. 2, 8
- [73] 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, 2
- [74] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2, 14
- [75] 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, 3
- [76] Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. In *The Twelfth International Conference on Learning Representations*, 2024. 3, 7, 18
- [77] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in neural information processing systems*, 32, 2019. 2, 15
- [78] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. 1, 2, 5, 14
- [79] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *Proceedings of the 40th International Conference on Machine Learning*, pages 32211–32252, 2023. 1, 2, 3, 18
- [80] Endre Süli and David F Mayers. *An introduction to numerical analysis*. Cambridge university press, 2003. 14
- [81] Joshua Tian Jin Tee, Kang Zhang, Hee Suk Yoon, Dhananjaya Nagaraja Gowda, Chanwoo Kim, and Chang D Yoo. Physics informed distillation for diffusion models. *arXiv preprint arXiv:2411.08378*, 2024. 18
- [82] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. 16
- [83] Pascal Vincent. A connection between score matching and denoising autoencoders. *Neural computation*, 23(7):1661–1674, 2011. 2

- [84] Bram Wallace, Meihua Dang, Rafael Rafailov, Linqi Zhou, Aaron Lou, Senthil Purushwalkam, Stefano Ermon, Caiming Xiong, Shafiq Joty, and Nikhil Naik. Diffusion model alignment using direct preference optimization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8228–8238, 2024. 1, 2, 3
- [85] Zhengyi Wang, Cheng Lu, Yikai Wang, Fan Bao, Chongxuan Li, Hang Su, and Jun Zhu. Prolificdreamer: High-fidelity and diverse text-to-3d generation with variational score distillation. *Advances in Neural Information Processing Systems*, 36:8406–8441, 2023. 5, 6
- [86] Joseph L Watson, David Juergens, Nathaniel R Bennett, Brian L Tripp, Jason Yim, Helen E Eisenach, Woody Ahern, Andrew J Borst, Robert J Ragotte, Lukas F Milles, et al. De novo design of protein structure and function with rfdiffusion. *Nature*, 620(7976):1089–1100, 2023. 1
- [87] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng-ming Yin, Shuai Bai, Xiao Xu, Yilei Chen, et al. Qwen-image technical report. *arXiv preprint arXiv:2508.02324*, 2025. 2, 3
- [88] Zhisheng Xiao, Karsten Kreis, and Arash Vahdat. Tackling the generative learning trilemma with denoising diffusion GANs. In *International Conference on Learning Representations*, 2022. 2
- [89] Sirui Xie, Zhisheng Xiao, Diederik Kingma, Tingbo Hou, Ying Nian Wu, Kevin P Murphy, Tim Salimans, Ben Poole, and Ruiqi Gao. Em distillation for one-step diffusion models. *Advances in Neural Information Processing Systems*, 37:45073–45104, 2024. 18
- [90] Jiazheng Xu, Xiao Liu, Yuchen Wu, Yuxuan Tong, Qinkai Li, Ming Ding, Jie Tang, and Yuxiao Dong. Imageward: Learning and evaluating human preferences for text-to-image generation. *Advances in Neural Information Processing Systems*, 36:15903–15935, 2023. 1, 2, 3
- [91] Yilun Xu, Ziming Liu, Max Tegmark, and Tommi Jaakkola. Poisson flow generative models. *Advances in Neural Information Processing Systems*, 35:16782–16795, 2022. 1, 2
- [92] Yilun Xu, Ziming Liu, Yonglong Tian, Shangyuan Tong, Max Tegmark, and Tommi Jaakkola. Pfgm++: Unlocking the potential of physics-inspired generative models. In *International Conference on Machine Learning*, pages 38566–38591. PMLR, 2023. 1, 2
- [93] Yilun Xu, Shangyuan Tong, and Tommi S. Jaakkola. Stable target field for reduced variance score estimation in diffusion models. In *The Eleventh International Conference on Learning Representations*, 2023. 2
- [94] Yanwu Xu, Yang Zhao, Zhisheng Xiao, and Tingbo Hou. 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*, pages 8196–8206, 2024. 18
- [95] Yilun Xu, Weili Nie, and Arash Vahdat. One-step diffusion models with  $f$ -divergence distribution matching. *arXiv preprint arXiv:2502.15681*, 2025. 6, 18
- [96] Mingyang Yi, Jiacheng Sun, and Zhenguo Li. On the generalization of diffusion model. *arXiv preprint arXiv:2305.14712*, 2023. 2, 3
- [97] Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and Bill Freeman. Improved distribution matching distillation for fast image synthesis. *Advances in neural information processing systems*, 37:47455–47487, 2024. 6, 15, 18
- [98] Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman, and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6613–6623, 2024. 6, 15, 16, 18
- [99] TaeHo Yoon, Joo Young Choi, Sehyun Kwon, and Ernest K Ryu. Diffusion probabilistic models generalize when they fail to memorize. In *ICML 2023 workshop on structured probabilistic inference & generative modeling*, 2023. 2, 3
- [100] Siyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and Saining Xie. Representation alignment for generation: Training diffusion transformers is easier than you think. *arXiv preprint arXiv:2410.06940*, 2024. 1, 2, 3, 7, 8, 17, 18
- [101] Yinan Zhang, Eric Tzeng, Yilun Du, and Dmitry Kislyuk. Large-scale reinforcement learning for diffusion models. In *European Conference on Computer Vision*, pages 1–17. Springer, 2024. 1, 2, 3
- [102] Boyang Zheng, Nanye Ma, Shengbang Tong, and Saining Xie. Diffusion transformers with representation autoencoders. *arXiv preprint arXiv:2510.11690*, 2025. 1, 3
- [103] Candi Zheng and Yuan Lan. Characteristic guidance: Non-linear correction for diffusion model at large guidance scale. In *International Conference on Machine Learning*, pages 61386–61412. PMLR, 2024. 6
- [104] Hongkai Zheng, Weili Nie, Arash Vahdat, Kamyr Azizzadehsheli, and Anima Anandkumar. Fast sampling of diffusion models via operator learning. In *International conference on machine learning*, pages 42390–42402. PMLR, 2023. 4, 18
- [105] Kaiwen Zheng, Yuji Wang, Qianli Ma, Huayu Chen, Jintao Zhang, Yogesh Balaji, Jianfei Chen, Ming-Yu Liu, Jun Zhu, and Qinsheng Zhang. Large scale diffusion distillation via score-regularized continuous-time consistency. *arXiv preprint arXiv:2510.08431*, 2025. 18
- [106] Linqi Zhou, Stefano Ermon, and Jiaming Song. Inductive moment matching. *arXiv preprint arXiv:2503.07565*, 2025. 8
- [107] Mingyuan Zhou, Huangjie Zheng, Yi Gu, Zhendong Wang, and Hai Huang. Adversarial score identity distillation: Rapidly surpassing the teacher in one step. *arXiv preprint arXiv:2410.14919*, 2024. 6, 18
- [108] Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation. In *Forty-first International Conference on Machine Learning*, 2024. 6, 18

## A. Extended Technical Discussion

We continue our technical discussion in Secs. 4 and 5, dividing this section into three parts.

### A.1. More on the Prediction Objective

**Optimality of training with Eq. (8).** Recall that we have already established that the loss value of Eq. (8) equals  $\mathbb{E}_{z,\delta} \|\partial_\delta f_\theta(z, \delta) - u(f_\theta(z, \delta), 1 - \delta)\|^2$ . Thus, at optimality, we have  $\partial_\delta f_\theta(z, \delta) = u(f_\theta(z, \delta), 1 - \delta)$  for all  $z \sim \pi$  and  $\delta \in [0, 1]$ . Furthermore, we note that the parametric relationship  $f_\theta(z, \delta) = z + \delta F_\theta(z, \delta)$ , which satisfies  $f_\theta(z, 0) = z$  for all  $z \sim \pi$  by design. We show that, assuming standard regularity conditions on  $u$ , we ensure  $\phi_u$  is uniquely defined via an initial value problem by the Picard-Lindelöf theorem [8], thus providing a well-defined criterion for  $f_\theta$ . This is a direct result of the fundamental theory of ODE.

**Proposition A.1.** *Let  $\phi_u(x_t, t, s)$  be defined as in Eq. (3), and we assume that there exists some  $L > 0$  for all  $y, y' \in \mathbb{R}^d$  and  $r \in [0, 1]$ ,  $\|u(y, r) - u(y', r)\| \leq L \|y - y'\|$ . We further assume that  $f_\theta(z, \delta)$  is continuously differentiable for  $\delta \in [0, 1]$ . Then, we have  $f_\theta(z, \delta) = \phi_u(z, 1, 1 - \delta)$  for all  $z \sim \pi, \delta \in [0, 1]$ , if and only if, (i)  $f_\theta(z, 0) = z$  for all  $z \sim \pi$ , and (ii)  $\partial_\delta f_\theta(z, \delta) = u(f_\theta(z, \delta), 1 - \delta)$  for all  $z \sim \pi, \delta \in [0, 1]$ .*

*Proof.* By assuming  $u(y, r)$  continuous in  $r$  and Lipschitz continuous in  $y$ , we could apply the Picard-Lindelöf theorem [8] to guarantee the existence and uniqueness of its solution  $\phi_u$ , which is

$$\phi_u(x_t, t, s) = x_t + \int_t^s -u(x(\tau), \tau) d\tau.$$

$\implies$ :

(i). Since the equality  $f_\theta(z, \delta) = \phi_u(z, 1, 1 - \delta)$  holds for all  $\delta \in [0, 1]$ , it must hold for  $\delta = 0$ . We know that  $\phi_u(z, 1, 1) = z$ . Thus, for any  $z \sim \pi$ ,

$$f_\theta(z, 0) = \phi_u(z, 1, 1) = z.$$

(ii). Given that  $f_\theta(z, \delta) = \phi_u(z, 1, 1 - \delta)$  for all  $\delta \in [0, 1]$ , and both  $f_\theta(z, \delta)$  and  $\phi_u(z, 1, 1 - \delta)$  are continuously differentiable with respect to  $\delta$  on  $[0, 1]$ , their derivatives with respect to  $\delta$  must be equal. That is, for all  $z \sim \pi, \delta \in [0, 1]$ ,

$$\partial_\delta f_\theta(z, \delta) = \frac{d}{d\delta} \phi_u(z, 1, 1 - \delta) = u(\phi_u(z, 1, 1 - \delta), 1 - \delta) = u(f_\theta(z, \delta), 1 - \delta).$$

$\iff$ :

Since  $f_\theta(z, 0) = z$  for all  $z \sim \pi$ , and  $\partial_\delta f_\theta(z, \delta) = u(f_\theta(z, \delta), 1 - \delta)$  for all  $z \sim \pi, \delta \in [0, 1]$ , we know that  $f_\theta(z, \delta)$  is a solution to the initial value problem on the interval  $\delta \in [0, 1]$ :

$$\frac{d}{d\delta} f_\theta(z, \delta) = u(f_\theta(z, \delta), 1 - \delta), \quad f_\theta(z, 0) = z.$$

By definition,  $\phi_u(z, 1, 1 - \delta)$  is also a solution to the same IVP on  $[0, 1]$ . Thus, we arrive at the equality between  $f_\theta(z, \delta)$  and  $\phi_u(z, 1, 1 - \delta)$ .  $\square$

**Discrete-time objective.** In Sec. 4.1, we skipped over the development of a more flexible discrete-time training objective for Eq. (8), which does not require efficient JVP implementations. First, recall that our default continuous-time objective is

$$\mathbb{E}_{z,\delta} \left\| F_\theta(z, \delta) + \text{sg}\left(\delta \partial_\delta F_\theta(z, \delta) - u(f_\theta(z, \delta), 1 - \delta)\right)\right\|^2. \quad (13)$$

Specifically, we can discretize the time horizon and approximate the partial derivative  $\partial_\delta F_\theta$  using finite differences. Following the common practice in the sampling procedure of flow models [22], we divide the time horizon into  $N$  intervals with  $N + 1$  boundary points:  $1 = t_1 > t_2 > \dots > t_{N+1} = 0$ . We note that the generating velocity at point  $f_\theta(z, 1 - t_i)$ , or equivalently  $f_\theta(z, t_1 - t_i)$  where  $1 \leq i \leq N$ , can be approximated by  $\frac{1}{t_i - t_{i+1}}(f_\theta(z, t_1 - t_{i+1}) - f_\theta(z, t_1 - t_i))$ . Short-handing  $t_a - t_b$  to  $\delta_{a,b}$ , we have the following discrete-time objective:

$$\mathbb{E}_{z,i} \left\| F_\theta(z, \delta_{1,i+1}) + \text{sg}\left(\delta_{1,i} \cdot \frac{F_\theta(z, \delta_{1,i+1}) - F_\theta(z, \delta_{1,i})}{\delta_{i,i+1}} - u(f_\theta(z, \delta_{1,i}), t_i)\right)\right\|^2, \quad (14)$$

which approximates Eq. (8) as  $\delta_{i,i+1} \rightarrow 0$ . Just like Eq. (13) equals to  $\mathbb{E}_{\mathbf{z},\delta} \|\partial_\delta \mathbf{f}_\theta(\mathbf{z}, \delta) - \mathbf{u}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)\|^2$ , we note that the loss value of Eq. (14) is the same as  $\mathbb{E}_{\mathbf{z},i} \|(\mathbf{f}_\theta(\mathbf{z}, \delta_{1,i+1}) - \mathbf{f}_\theta(\mathbf{z}, \delta_{1,i})) / \delta_{i,i+1} - \mathbf{u}(\mathbf{f}_\theta(\mathbf{z}, \delta_{1,i}), t_i)\|^2$ , where  $(\mathbf{f}_\theta(\mathbf{z}, \delta_{1,i+1}) - \mathbf{f}_\theta(\mathbf{z}, \delta_{1,i})) / \delta_{i,i+1}$  is the model’s generating velocity, approximating  $\partial_\delta \mathbf{f}_\theta(\mathbf{z}, \delta)$ . If we use  $v_G$  to further simplify the notations, we arrive at the following optimization gradients in Eq. (9):

$$\nabla_\theta \mathbb{E}_{\mathbf{z},\delta} \left[ \mathbf{F}_\theta(\mathbf{z}, \delta)^\top \text{sg} \underbrace{\left( v_G(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta) - \mathbf{u}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta) \right)}_{\Delta_{v_G, \mathbf{u}}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)} \right], \quad (15)$$

**Error analysis.** With a discrete-time objective, we essentially train the student to trace a numerically integrated trajectory by an ODE solver. More concretely, the loss function inherently mimics a specific solver, whose error rate depends on the technique deployed. As written, Eq. (14) can be understood as using an Euler method [74, 78] for solving the flow, specified at time steps  $t_1, t_2, \dots, t_N$ . Similar to existing fast-solver literature, we could readily adapt the objective using higher-order finite differences to mimic more precise, higher-order solvers [28, 46]. The classical theory from numerical analysis [80] indicates that a local truncation error at  $t_i$  bounded by  $\mathcal{O}((\delta_{i,i+1})^{p+1})$  leads to the global error rate is  $\mathcal{O}((\delta_{\max})^p)$ , where  $\delta_{\max} := \max_i \delta_{i,i+1} = \max_i (t_i - t_{i+1})$  is the maximum step size. This component is the discretization error, which is entirely independent of the network’s approximation error (the inevitable training inaccuracy discussed in Sec. 4.1). Together, these two errors represent the total deviation of the student’s trajectory from the true teacher flow defined by  $\mathbf{u}$ . The total error can also be shown [5] to directly control the Wasserstein distance between the generative distribution of the student and the teacher.

**Number of discretization steps.** The above discussion suggests that one should take as small a step as possible,  $(t_i - t_{i+1}) \rightarrow 0$ , and as precise a solver as possible, for the discretization error to be minimal. This means that we should use the continuous-time setup if allowed, and for the discrete-time setting, we should choose to use a large number of discretization steps  $N$ . However, this is not the trend we observe in practice. Surprisingly, larger  $N$  after a certain point actually degrades the performance. We attribute this behavior to the accumulated approximation error (discussed in Sec. 4.1) made by the model. Since  $\mathbf{u}$  is evaluated at the model’s prediction, the goodness of the trajectory implicitly depends on the quality of  $\mathbf{f}_\theta$  itself, and such a problem exacerbates as  $N$  becomes larger. We further note that incorporating higher-order solvers in  $\mathbf{u}$  like the Heun method [28] is helpful, as it significantly reduces the discretization error. Fortunately, in Fig. 9, we observe that while the model is sensitive to the number of discretization steps and solver choices when training with Eq. (9) alone, the corrective signal in Eq. (11) effectively renders such a decision unimportant, making our final proposal robust against  $N$  and the precision of solvers.

**Sampling of  $\delta$ .** Another important aspect of the training algorithm is the time sampling of  $\delta$ . On one hand, the correctness of the model propagates from small jumps  $\delta = 0$  to large jumps  $\delta = 1$ . On the other hand, the high-frequency features of the images only emerge at a lower noise level (large  $\delta$ ). To balance these two notions, we investigate a mixture of a logit-normal distribution [11, 28] and a fixed value of  $\delta = 0$  (effectively a dropout on  $\delta$ ). Results are presented in Tab. 3. Note that for the discrete-time setting, we apply an additional step of the floor operation with respect to our predefined set of discrete time steps.

**Confident region warmup.** We observe that naively optimizing with Eq. (9) brings about instability early on in training, mainly with the JVP approach. We hypothesize that it is because  $\mathbf{u}$  is evaluated at a state predicted by  $\mathbf{f}_\theta$ , which can be out-of-distribution for  $\mathbf{u}$  at initialization. We propose to add noise to the predicted state before feeding it to  $\mathbf{u}$ . Specifically, we replace  $\mathbf{u}(\mathbf{f}_\theta(\mathbf{z}, \delta), 1 - \delta)$ , with  $\mathbf{u}(\mathbf{I}_{t_c|1-\delta}(\mathbf{f}_\theta(\mathbf{z}, \delta), \mathbf{n}), t_c)$ , where  $\mathbf{I}_{t_c|t}$  is the generalized interpolating function, a transition kernel that takes samples from the noise level of  $t$  to  $t_c$ .

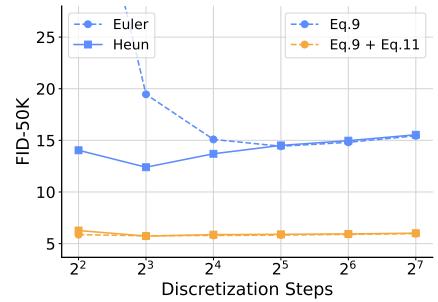


Figure 9. **Effects of discretization steps and solver type.** Our correction objective in Eq. (11) makes the model robust against both.

% of $\delta = 0$	$\delta$ sampling	FID $\downarrow$
10%	LogitNormal(0, 1)	12.40
0%	LogitNormal(0, 1)	47.09
30%	LogitNormal(0, 1)	13.19
50%	LogitNormal(0, 1)	14.30
10%	LogitNormal(-0.8, 1)	13.78
10%	LogitNormal(-0.4, 1)	12.98
10%	LogitNormal(0, 1.2)	12.60
0%	LogitNormal(-0.8, 1.6)	13.43
0%	LogitNormal(-0.4, 1.6)	12.98
0%	LogitNormal(0, 1.6)	12.59

Table 3. **Ablation of  $\delta$  sampling.** A mixture of a logit-normal and a fixed  $\delta=0$  works well.

For the linear interpolation scheme [40, 43, 50], we have  $\mathbf{I}_{t_c|t}(\mathbf{x}_t, \mathbf{n}) = \frac{1-t_c}{1-t}\mathbf{x}_t + \sqrt{t_c^2 - \frac{(1-t_c)^2}{(1-t)^2}t^2}\mathbf{n}$ , assuming  $t_c > t$ ; if  $t_c \leq t$ , it simply does nothing and returns  $\mathbf{x}_t$ . At the start of training, we linearly decrease  $t_c$  from 1 to 0 over a short warmup of 10K steps. Intuitively, this procedure ensures that  $\mathbf{u}$  always operates in a region where it is confident. After the inclusion of this warmup period, we do not find any instability of training  $f_\theta$  with a reasonable sampling distribution of  $\delta$ .

## A.2. More on the Correction Objective

**From minimizing IKL to aligning the noising velocities.** We know from prior works [48, 98] that minimizing Eq. (10) w.r.t.  $\theta$  gives us the following gradient:

$$\mathbb{E}_{r, \mathbf{x}_r} \left[ (\nabla_\theta \mathbf{x}_r)^\top (\nabla_{\mathbf{x}_r} \log q_r(\mathbf{x}_r) - \nabla_{\mathbf{x}_r} \log p_r(\mathbf{x}_r)) \right], \quad (16)$$

where  $\nabla_{\mathbf{x}_r} \log q_r(\mathbf{x}_r)$  and  $\nabla_{\mathbf{x}_r} \log p_r(\mathbf{x}_r)$  are the score functions [25, 77] of  $q_r$  and  $p_r$  respectively. It has also been shown that there exists a direct bijective translation between the score functions and the marginal velocities [50]. Specifically, for the linear interpolation scheme [40, 43, 50] and our definition of the conditional velocity, we have:

$$\begin{aligned} \mathbf{u}(\mathbf{x}_r, r) &= \frac{r}{1-r} \nabla_{\mathbf{x}_r} \log p_r(\mathbf{x}_r) + \frac{1}{1-r} \mathbf{x}_r \\ \nabla_{\mathbf{x}_r} \log p_r(\mathbf{x}_r) &= \frac{1-r}{r} \mathbf{u}(\mathbf{x}_r, r) - \frac{1}{r} \mathbf{x}_r. \end{aligned}$$

Additionally, recall that the student predicts the data sample as  $\mathbf{f}_\theta(\mathbf{z}, 1)$ . With average velocity parameterization, the intermediate state is thus  $\mathbf{x}_r = \mathbf{I}_r(\mathbf{f}_\theta(\mathbf{z}, 1), \mathbf{n}) = (1-r)(\mathbf{z} + \mathbf{F}_\theta(\mathbf{z}, 1)) + r\mathbf{n}$ . Thus, we could rewrite Eq. (16) as

$$\mathbb{E}_{r, \mathbf{x}_r} \left[ \frac{(1-r)^2}{r} (\nabla_\theta \mathbf{F}_\theta(\mathbf{z}, 1))^\top (\mathbf{v}_N(\mathbf{I}_r(\mathbf{f}_\theta(\mathbf{z}, 1), \mathbf{n}), r) - \mathbf{u}(\mathbf{I}_r(\mathbf{f}_\theta(\mathbf{z}, 1), \mathbf{n}), r)) \right],$$

where  $\mathbf{v}_N$  is the marginal noising velocity induced by the student's generated distribution  $q$ , similar to  $\mathbf{u}$  with  $p$ . Dropping the weighting  $\frac{(1-r)^2}{r}$ , as it does not change the optimal solution and provides us with easier-to-control gradients, we arrive at Eq. (11). Note that a different interpolation scheme does not change our investigation.

$$\nabla_\theta \mathbb{E}_{\mathbf{z}, \mathbf{n}, r} \left[ \mathbf{F}_\theta(\mathbf{z}, 1)^\top \text{sg} \left( \underbrace{\mathbf{v}_N(\mathbf{I}_r(\mathbf{f}_\theta(\mathbf{z}, 1), \mathbf{n}), r)}_{\Delta_{\mathbf{v}_N, \mathbf{u}}(\mathbf{I}_r(\mathbf{f}_\theta(\mathbf{z}, 1), \mathbf{n}), r)} - \mathbf{u}(\mathbf{I}_r(\mathbf{f}_\theta(\mathbf{z}, 1), \mathbf{n}), r) \right) \right]. \quad (17)$$

**Learning rate.** Our correction objective defined in Eq. (11) assumes access to  $\mathbf{v}_N$ , which is the velocity of the noising flow starting from the generated distribution, and we approximate it by training an auxiliary model  $\mathbf{g}_\psi$  concurrently with Eq. (2). The quality of the optimization signal  $\Delta_{\mathbf{v}_N, \mathbf{u}}$  depends on the quality of this approximation. Empirically in Tab. 4, we confirm this intuition and observe that the algorithm benefits from having a larger learning rate on  $\mathbf{g}_\psi$  compared to  $f_\theta$ . We note that it is also possible to adopt a two time-scale update rule [19, 97], where we train multiple iterations on  $\mathbf{g}_\psi$  before updating  $f_\theta$ , but we do not explore this option given the overhead. In a related vein, we find that while training with only the prediction objective in Eq. (9) permits a wide range of learning rates, incorporating the correction objective in Eq. (11) prefers a smaller one on  $f_\theta$  (adopted  $3 \times 10^{-5}$  compared to  $1 \times 10^{-4}$  in SiT).

learning rate for $\mathbf{g}_\psi$	FID ↓
$3 \times 10^{-5}$	8.28
$6 \times 10^{-5}$	5.77
$8 \times 10^{-5}$	5.72
$1 \times 10^{-4}$	5.63

Table 4. **Ablation of  $\mathbf{g}_\psi$ 's lr.** A higher lr compared to  $f_\theta$ 's one ( $3 \times 10^{-5}$ ) is better.

**Additional note on sampling of  $r$  in Eq. (11).** Recall that our understanding of Eq. (11) with a PDE perspective through the continuity equation in Sec. 5.1 leads to our design of sampling  $r$  more in the higher noise levels. We would like to make a brief note that, similar to Eq. (9), if  $\mathbf{F}_\theta$  is further parameterized, we should replace the first term in Eq. (11) with the actual network output. Additionally, one needs to ensure  $\Delta_{\mathbf{v}_N, \mathbf{u}}$  across different  $r$  values are roughly on the same scale first, so that their actual contributions can be precisely controlled via the sampling of  $r$ . Concretely, for our case of the linear interpolation and velocity parameterization [40, 43, 50],  $\Delta_{\mathbf{v}_N, \mathbf{u}}$  is really just the difference between the direct outputs of two neural networks, which does not require further manipulations, assuming the network outputs are of consistent scale. In comparison, another

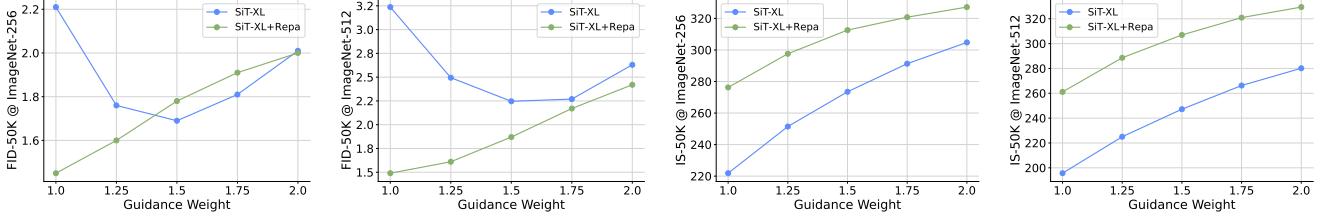


Figure 10. Performances across different guidance strength parameter  $\gamma$  in terms of FID (↓) and Inception Score (↑).

commonly used setup is the EDM parameterization [28], whose velocity at  $(\mathbf{x}_r = \mathbf{f}_\theta(\mathbf{z}, 1) + r\mathbf{n}, r)$ ,  $\mathbf{n} \sim \pi$  is of the form  $\frac{c_{\text{skip}}(r)-1}{r}\mathbf{x}_r + \frac{c_{\text{out}}(r)}{r}\mathbf{G}$ , where  $\mathbf{G}$  is the actual network (the auxiliary or teacher model). Notice that the outputs of  $\mathbf{G}$  are by design of constant norm across  $r$ , so the magnitude of  $\Delta_{\mathbf{v}_{\mathbf{N}}, \mathbf{u}}$  is proportional to  $\frac{c_{\text{out}}(r)}{r}$ . In Yin et al. [98], an additional weighting of  $r$  is applied, which makes the overall gradient contributions from Eq. (11) proportional to  $c_{\text{out}}(r)$ . Substituting in the actual terms used, we have  $\frac{0.5r}{\sqrt{0.5^2+r^2}}$ , which heavily downplays the effect of small  $r$  (lower noise levels).

### A.3. More on How to Combine Both

**Additional schedule on  $\alpha$ .** We set  $\alpha = 0$  for the first 10K steps (recall that the prediction objective has a warmup of 10K steps), and follow it with a linear warmup of another 10K steps before  $\alpha$  settles at a reference value  $\alpha_{\text{ref}}$ . While the model performance is shown to be robust across a wide range of  $\alpha$  in Fig. 6, we further notice a general trend that larger  $\alpha$  learns faster in the beginning of training, and smaller  $\alpha$  converges better as the training continues. Hence, for our REPA distillation tasks, we try out a simple inverse square root decay [30, 33], and arrive at the following schedule on  $\alpha$ :

$$\alpha = \alpha_{\text{ref}} \frac{\text{clip}\left(\frac{n-T_{\text{delay}}}{T_{\text{warmup}}}, 0, 1\right)}{\sqrt{\max(n/T_{\text{decay}}, 1)}},$$

where  $n$  is the current training iteration,  $T_{\text{delay}} = T_{\text{warmup}} = 10K$ , and  $T_{\text{decay}}$  is the hyperparameter that controls the decay rate with  $\infty$  indicating no decay applied ( $\alpha$  stays at the constant value  $\alpha_{\text{ref}}$  after warmup). We also would like to clarify that, considering the adopted 75-25 split, an  $\alpha$  value of 0.3 really means that the gradient contribution of the correction objective is around 10% of that of the prediction objective.

**Additional empirical results on prediction and correction synergy.** In addition to presenting the model progress in Fig. 7, we list the final performances of training with only prediction Eq. (9), only correction Eq. (11), and both in Tab. 5. Specifically, all models are distilled from SiT-XL/2 [50], and we compare them in terms of FID [19] and Inception Score [66] at 1.5M iterations (300 epochs). For each method, we select the optimal  $\gamma$  based on the FID performances. Recall from Fig. 7, at this point in training, although Eq. (9) makes progress, its absolute performance still lags far behind the other two configurations because of the significant error accumulation. In contrast, Eq. (11) has already suffered from mode collapse, and its diversity continues to deteriorate.

## B. Implementation Details

**Training.** We follow the standard practice and train our models in the latent space of the VAE used in Rombach et al. [62]. The model architecture we use is based on a standard DiT [55] with  $2\times 2$  patches. Recall that the standard input to a flow model for ImageNet is  $\mathbf{f}_\theta(\mathbf{x}_t, t, c)$  ( $c$  is the class label), and we need to include two additional scalar inputs: jump duration  $\delta$  and guidance strength  $\gamma$ . That is, during both training and inference, the model follows  $\mathbf{f}_\theta(\mathbf{z}, 1, c, \delta, \gamma)$ . Thus, we add a few layers to handle these additional conditions. For both, we follow the standard design for including scalar input. Specifically, we use a 256-dimensional frequency embedding [9, 82] followed by a two-layer SiLU-activated MLP with the same dimensionality as the model’s hidden size. We then add all four embeddings from  $t$ ,  $c$ ,  $\delta$ , and  $\gamma$  together (the newly added ones,  $\delta$  and  $\gamma$ , are initialized at 0) and feed the sum to each block. The same goes for  $\mathbf{g}_\psi$ , where it is modified to take input  $\mathbf{g}_\psi(\mathbf{x}_t, t, c, \gamma)$  as it needs to track different noising flows from different  $\gamma$  of  $\mathbf{f}_\theta$ . No further architecture changes are necessary for stable and

objective	FID ↓	IS ↑
Eq. (9)	5.78	257.02
Eq. (11)	3.19	258.15
Eqs. (9) and (11)	1.69	273.49

Table 5. Synergy between Eqs. (9) and (11). Together, they achieve performance that neither could attain in isolation.

Table 6. Detailed experimental configurations of our main results.

<b>Task</b>	SiT-XL/2 [50]		SiT-XL/2+REPA [100]					
teacher model resolution	256×256	512×512	256×256	512×512				
<b>General</b>								
iterations (epochs)	1.5M (300)	1M (200)	1.5M (300)	1M (200)				
batch size		256						
optimizer		Adam [33]						
optimizer betas		(0.9, 0.99)						
optimizer eps		$1 \times 10^{-8}$						
weight decay		0.0						
dropout		0.0						
$f_\theta$ learning rate		constant $3 \times 10^{-5}$						
$g_\psi$ learning rate		constant $1 \times 10^{-4}$						
EMA decay		0.9999						
<b>Network</b>								
params (M)		678						
FLOPs (G)	119	525	119	525				
depth		28						
hidden dim		1152						
heads		16						
patch size		2×2						
change from teacher	additional input for $\delta$ and $\gamma$ in $f_\theta$ ; additional input for $\gamma$ in $g_\psi$							
<b>Training</b>								
<i>Specific to Eq. (9)</i>								
confident region warmup duration		10K						
$\delta$ type		discrete; uniform; N=8						
$\delta$ sampling	LogitNormal(0, 1)		LogitNormal(-0.4, 1.2)					
% of $\delta = 0$		10%						
$u$ type		Heun solver [28]						
$k$		1						
<i>Specific to Eq. (11)</i>								
$r$ sampling		LogitNormal(0.8, 1.6)						
guidance interval	[0, 0.4]	[0, 0.5]	[0, 0.3]	[0, 0.3]				
<i>Relevant to both Eqs. (9) and (11)</i>								
split between Eqs. (9) and (11)		75% : 25%						
$\gamma$ range		[1, 2]						
$\alpha_{\text{ref}}$	0.3		0.6					
$\alpha_{\text{ref}}$ schedule ( $T_{\text{delay}}, T_{\text{warmup}}, T_{\text{decay}}$ )	(10K, 10K, $\infty$ )		(10K, 10K, 25K)					

effective training. For each of our reported entries listed in Tab. 2, we present their implementation details in Tab. 6. All model trainings are done with an internal JAX codebase on TPU.

**Evaluation.** We observe small performance variations between TPU-based FID evaluation and GPU-based FID evaluation (ADM’s TensorFlow evaluation suite [9]<sup>6</sup>). To ensure a fair comparison with the baseline methods, we convert all of our models into PyTorch, sample all of our models on GPU, and obtain FID scores using the ADM evaluation suite for reporting the final results of our XL-size models in Tabs. 2 and 5 and Fig. 10. Additionally, since our model is trained on a range of guidance strengths  $\gamma \sim \mathcal{U}(1, 2)$ , we can efficiently sweep for an optimal value during inference. We report the best FID in Tab. 2, and provide the complete performance curves in Fig. 10.

<sup>6</sup><https://github.com/openai/guided-diffusion/tree/main/evaluations>

## C. Related Work

Among the existing distillation approaches, BOOT [16] stands as the most closely related precursor, sharing the distinct operational characteristic of being data-free. However, our works diverge fundamentally in their conceptual positioning and the identified imperative for removing data. BOOT frames the data-free property primarily as a practical/logistical advantage, emphasizing the benefits of bypassing the storage and privacy burdens associated with massive, proprietary training sets. In contrast, we argue that the exclusion of data is not merely a convenience but a theoretical necessity for ensuring distributional fidelity. We elevate the data-free paradigm from a strategy of efficiency to one of correctness, presenting it as the rigorous solution to the identified Teacher-Data Mismatch. Our proposed method also differs significantly from BOOT and its subsequent improvements, which are not necessarily data-free in nature. In particular, Gu et al. [16] focuses on the specific signal-ODE parameterization, which requires a separate loss just to enforce the boundary condition  $f_\theta(z, 0) = z$ . In comparison, our prediction objective stems from the properties of average velocity [14], which satisfy the boundary condition by design. Tee et al. [81] and the Lagrangian objective in Boffi et al. [5] consider more general ODE formulations. Their optimization involves costly computations of the gradients over the partial derivatives  $\partial_\delta f_\theta$ , whereas ours does not (the partial derivatives in Eq. (8) are placed inside the stop-gradient operation). This improved training efficiency also originates from the average velocity perspective and our deduced identity in Eq. (7). Furthermore, we introduce an auxiliary correction objective for the accumulated prediction errors, pushing the model performance beyond the current state-of-the-art. In doing so, we believe our work finally completes the picture, validating the data-free paradigm as a robust and promising foundation for the future of generative model acceleration.

Our contribution sits within a much broader body of literature dedicated to accelerating diffusion and flow models. These techniques generally fall into two categories based on their distillation targets. The first category operates at the **trajectory level**, attempting to compress the complex ODE integration into fewer steps by directly mimicking the sampling path or its solution operator [47, 104]. The foundational work of Progressive Distillation [52, 65] further established the viability of this direction through an iterative strategy that progressively halves the required sampling steps. This paradigm was significantly expanded by Consistency Models [3, 45, 76, 79], which enforce a property of self-consistency along the trajectory, allowing the model to map arbitrary intermediate states directly to the data origin. More recent approaches [5, 7, 13, 14, 17, 32, 36, 37, 42–44, 57, 61, 64] have further refined this objective by formulating direct matching conditions between the student’s transport map and the teacher’s vector field. The second category operates at the **distribution level** [48, 49, 67–69, 89, 94, 95, 97, 98, 107, 108], where the student is trained to match the teacher’s marginal distribution directly, often utilizing adversarial or score-based objectives without strictly adhering to the teacher’s specific trajectory. We make contributions in both directions by proposing an efficient algorithm for distilling trajectories without data and elucidating additional design spaces for better distribution matching objectives. Together, our proposed predictor-corrector framework can be seen as combining the strengths of the two categories [45, 105], achieving superior quality while maintaining desired diversity, all without reliance on external data.

## D. Additional Visual Results

In Figs. 11 to 18, we present additional uncurated samples generated by FreeFlow-XL/2 at  $512 \times 512$  resolution with only 1-NFE. Again, we emphasize that, during training, we only make use of the teacher model (SiT-XL/2+REPA [100]), without querying any samples from ImageNet.



Figure 11. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “robin” (15)



Figure 12. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “Siberian husky” (250)



Figure 13. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “barn” (425)



Figure 14. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “desktop computer” (527)

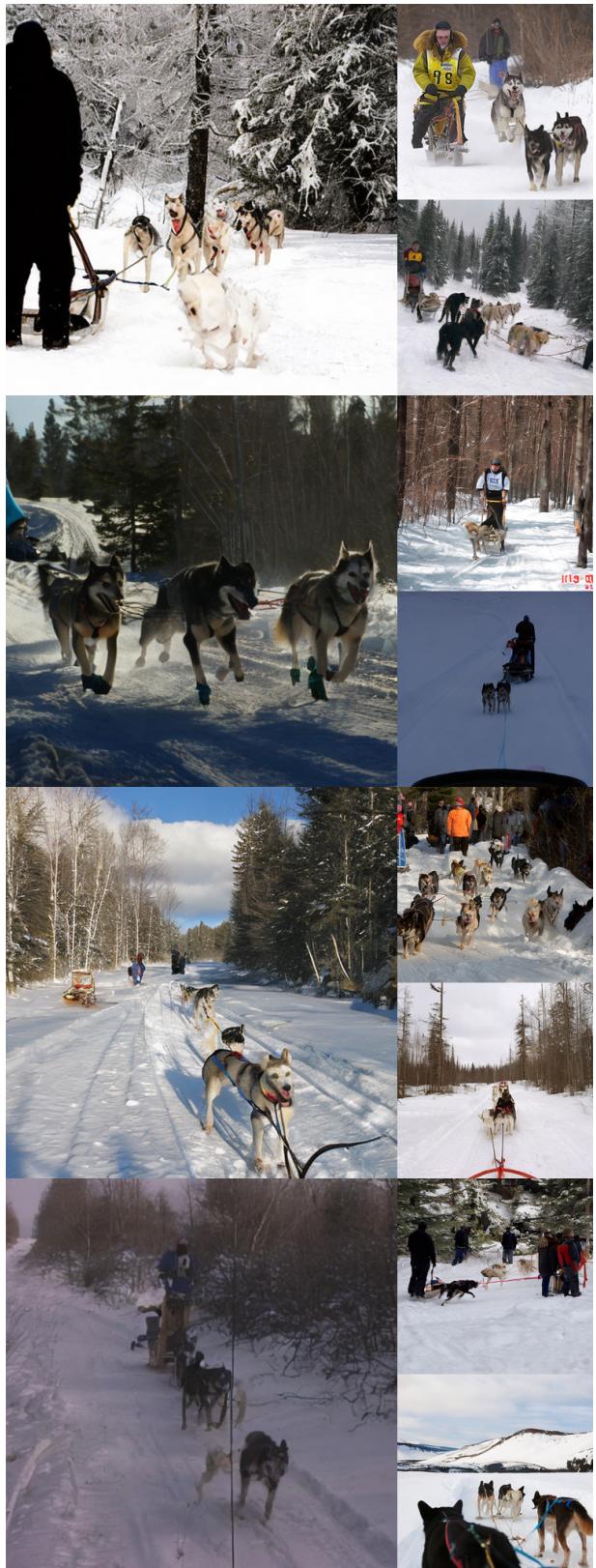


Figure 15. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “dogsled” (537)



Figure 16. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “fire truck” (555)

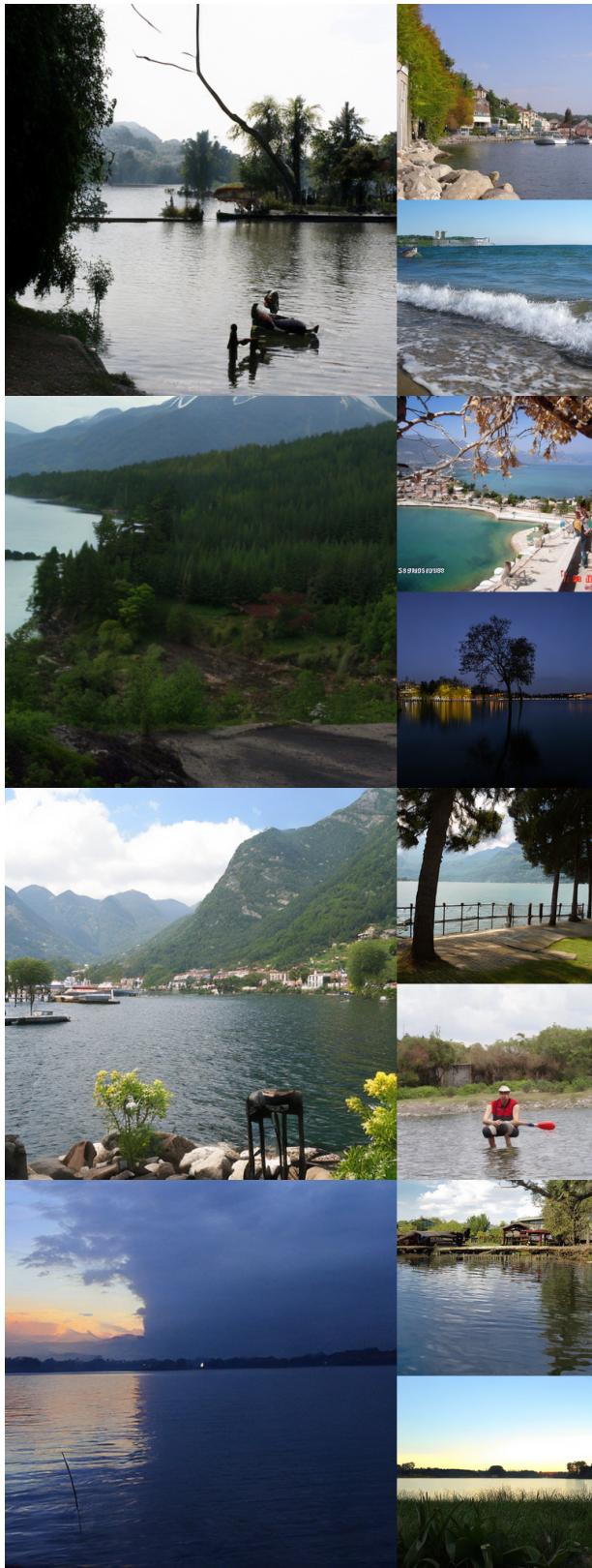


Figure 17. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “lakeside” (975)



Figure 18. Uncurated  $512 \times 512$  samples by FreeFlow, 1-NFE.  
Class label = “volcano” (980)