

000 001 002 003 004 005 SOFLOW: SOLUTION FLOW MODELS FOR ONE-STEP 006 GENERATIVE MODELING 007 008 009

010 **Anonymous authors**
 011 Paper under double-blind review
 012
 013
 014
 015
 016
 017
 018
 019
 020
 021
 022
 023
 024

ABSTRACT

011 The multi-step denoising process in diffusion and Flow Matching models causes
 012 major efficiency issues, which motivates research on few-step generation. We
 013 present Solution Flow Models (SoFlow), a framework for one-step generation
 014 from scratch. By analyzing the relationship between the velocity function and
 015 the solution function of the velocity Ordinary Differential Equation (ODE), we
 016 propose a flow matching loss and a solution consistency loss to train our mod-
 017 els. The flow matching loss allows our models to provide estimated velocity fields
 018 for Classifier-Free Guidance (CFG) during training, which improves generation
 019 performance. Notably, our consistency loss does not require the calculation of
 020 the Jacobian-Vector Product (JVP), a common requirement in recent works that
 021 is not well-optimized in deep learning frameworks like PyTorch. Experimental
 022 results indicate that, when trained from scratch using the same diffusion trans-
 023 former (DiT) architecture and with an equal number of training epochs, our mod-
 024 els achieve better FID-50K scores compared to MeanFlow models on the Im-
 ageNet 256x256 dataset.



033 Figure 1: Generated one-step samples on Imagenet 256×256 dataset by our Solution Flow Models.
 034

035 1 INTRODUCTION

036 Diffusion models (Sohl-Dickstein et al., 2015; Song & Ermon, 2019; Ho et al., 2020; Song et al.,
 037 2020) and Flow Matching models (Lipman et al., 2022; Liu et al., 2022; Albergo & Vanden-Eijnden,
 038 2022) have emerged as foundational frameworks in generative modeling. Diffusion models operate
 039 by systematically adding noise to data and then learning a denoising process to generate high-quality
 040 samples. Flow Matching offers a more direct alternative, modeling the velocity fields that transport
 041 a simple prior distribution to a complex data distribution. Despite their power and success across
 042 various generative tasks (Esser et al., 2024; Ma et al., 2024; Polyak et al., 2024), both approaches
 043 rely on an iterative, multi-step sampling process, which hinders their generation efficiency.

044 Addressing this latency is becoming a key area of research. Consistency Models (Song et al., 2023;
 045 Song & Dhariwal, 2023; Geng et al., 2024; Lu & Song, 2024) and related techniques (Kim et al.,
 046 2023; Wang et al., 2024a; Frans et al., 2024; Heek et al., 2024) have gained prominence by enabling
 047 rapid, few-step generation. These methods learn a direct mapping from any point on a generative
 048 trajectory to a consistent, “clean” output, bypassing the need for iterative refinement. However, this
 049 paradigm introduces significant challenges. Consistency models trained from scratch often fail to
 050 leverage Classifier-Free Guidance (CFG) for enhancing sample quality, and they are further hindered
 051 by instability caused by changing optimization targets. Recent works (Peng et al., 2025; Geng
 052 et al., 2025) address this instability by incorporating a flow matching loss. However, this approach
 053 introduces a new computational bottleneck: it relies on Jacobian-Vector Product (JVP) calculations,
 which are much less optimized than forward propagation in frameworks such as PyTorch.

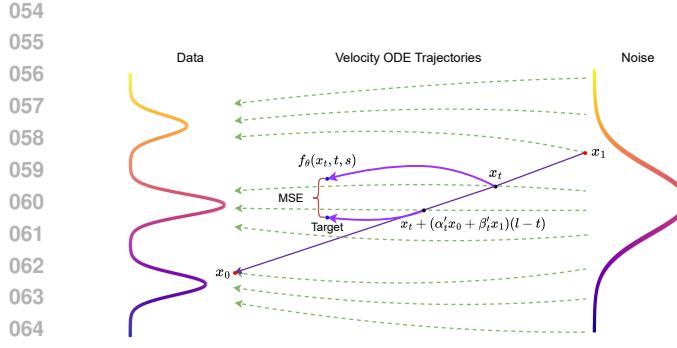


Figure 2: Illustration of solution consistency loss: The plot shows a straight-line flow matching trajectory defined by a data-noise pair (x_0, x_1) , where x_t is the intermediate point given by $x_t = \alpha_t x_0 + \beta_t x_1$ (α_t, β_t are C^1 functions with $\alpha_0 = \beta_1 = 1, \alpha_1 = \beta_0 = 0$). Given three time points $s < l < t$, the mean square error is computed between our model $f_\theta(x_t, t, s)$ and the stop-gradient target $f_{\theta^-}(x_t + (\alpha'_t x_0 + \beta'_t x_1)(l-t), l, s)$.

In this work, we introduce a new approach for one-step generation that avoids these limitations. Instead of relying on iterative ODE solvers, we propose to directly learn the solution function of the velocity ODE defined by Flow Matching. We denote this function as $f(x_t, t, s)$, which explicitly maps a state x_t at time t to its evolved state x_s at any other time s . To learn it with a parameterized model $f_\theta(x_t, t, s)$, we first analyze the properties that can make a neural network into a valid solution function. Based on this analysis, we formulate a training objective comprising a flow matching loss and a solution consistency loss. The resulting model, $f_\theta(x_t, t, s)$, not only accommodates the flow matching objective and CFG naturally but also eliminates the need for expensive JVP calculations during training. Our experimental results demonstrate the effectiveness of this approach, showing that our models achieve superior FID-50K scores compared to MeanFlow models on the ImageNet 256×256 dataset, using the same diffusion transformer (DiT) architecture and training duration.

2 RELATED WORKS

Diffusion, Flow Matching and Stochastic Interpolants Diffusion models (Sohl-Dickstein et al., 2015; Song et al., 2020; Kingma et al., 2021; Karras et al., 2022) and Flow Matching (Lipman et al., 2022; Liu et al., 2022) are widely adopted generative modeling frameworks. These approaches either progressively corrupt data with noise and train a neural network to denoise it, or learn a velocity field that governs a transformation from the data distribution to a simple prior. They have been scaled successfully for image generation (Rombach et al., 2022; Saharia et al., 2022; Podell et al., 2023) and video generation (Ho et al., 2022; Brooks et al., 2024) tasks. Stochastic interpolants (Albergo & Vanden-Eijnden, 2022; Albergo et al., 2023) build upon these concepts by explicitly defining stochastic trajectories between the data and prior distributions and aligning their associated velocity fields to enable effective distributional transport.

Few-step Generative Models Reducing the number of sampling steps has become an active research direction, driven by the need for faster generation and better theoretical insights. Prior efforts fall into two main streams: (i) distillation-based approaches that compress pre-trained multi-step models into few-step generators (Salimans & Ho, 2022; Sauer et al., 2024; Geng et al., 2023; Luo et al., 2023b; Yin et al., 2024; Zhou et al., 2024); and (ii) from-scratch training, which learns fast samplers without teachers, most prominently through the consistency-training family (CMs, iCT, ECT, sCT) (Song et al., 2023; Song & Dhariwal, 2023; Geng et al., 2024; Lu & Song, 2024).

Consistency Models (CMs) collapse noised inputs directly to clean data, enabling one- or few-step generation (Song et al., 2023). While first applied mainly in distillation (Luo et al., 2023a; Geng et al., 2024), later works established that they can also be trained from scratch via consistency training (Song & Dhariwal, 2023). Building on this, iCT (Song & Dhariwal, 2023) simplifies objectives and improves stability with robust losses and teacher-free training, while ECT (Geng et al., 2024) introduces a continuous-time ODE formulation that unifies CMs and diffusion. The recent sCT framework (Lu & Song, 2024) further stabilizes continuous-time consistency training, scaling CMs up to billion-parameter regimes.

Beyond fixed start/end settings, newer approaches extend CMs to arbitrary timestep transitions, aligning better with continuous-time dynamics. Shortcut models (Frans et al., 2024) condition on noise level and step size to flexibly support both one- and few-step sampling with shared weights. MeanFlow (Geng et al., 2025) introduces interval-averaged velocities and analyzes the relationship

108 between averaged and instantaneous velocities. IMM (Zhou et al., 2025) matches moments across
 109 transitions in a single-stage objective, reducing variance and avoiding multi-stage distillation. Flow-
 110 Anchored CMs (Peng et al., 2025) regularize shortcut learning with a flow matching anchor, im-
 111 proving stability and generalization. On the distillation side, Align Your Flow (Sabour et al., 2025)
 112 unifies CM and FM into continuous-time flow maps effective across arbitrary step counts, further
 113 enhanced with autoguidance and adversarial fine-tuning. Finally, Transition Models (Wang et al.,
 114 2025) reformulate generation around exact finite-interval dynamics, unifying few- and many-step
 115 regimes and achieving monotonic quality gains as step budgets increase.

117 3 PRELIMINARY: FLOW MATCHING

119 We first introduce the setting of Flow Matching. Flow Matching models learn a velocity field that
 120 transforms a known prior distribution like standard Gaussian distribution to the data distribution.
 121 More precisely, we denote the data distribution as $p(x_0)$, and the prior distribution as $p(x_1)$, which
 122 is standard Gaussian distribution $N(0, I)$ in our setting.

123 A general noising process is defined as $x_t = \alpha_t x_0 + \beta_t x_1$, where $x_t \in \mathbb{R}^n, t \in [0, 1]$, α_t and β_t are
 124 continuously differentiable functions satisfying the boundary conditions $\alpha_0 = 1, \alpha_1 = 0, \beta_0 = 0$,
 125 and $\beta_1 = 1$. Then, the marginal velocity field associated with the noising process is defined as

$$126 \quad v(x_t, t) = \mathbb{E}_{p(x_0, x_1 | x_t)} [\alpha'_t x_0 + \beta'_t x_1], \quad (1)$$

128 which is a conditional expectation of the conditional velocity field. α'_t and β'_t denote the derivatives
 129 of α_t and β_t . According to the Flow Matching framework, given the marginal velocity field $v(x_t, t)$,
 130 new samples can be generated by first sampling $x_1 \sim p(x_1)$, and then solving the initial value
 131 problem (IVP) of the following velocity ordinary differential equation (ODE) from $t = 1$ to $t = 0$:

$$132 \quad \frac{dX(t)}{dt} = v(X(t), t), X(1) = x_1. \quad (2)$$

135 To construct a generative model based on this principle, a straightforward approach is to approxi-
 136 mate the marginal velocity field using a neural network v_θ , parameterized by θ , which is trained by
 137 minimizing the following mean square objective:

$$138 \quad \mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, x_t} \|v_\theta(x_t, t) - v(x_t, t)\|^2. \quad (3)$$

140 However, directly optimizing a neural network with this loss is impractical since the real velocity
 141 field $v(x, t)$ is a conditional expectation given x_t . To overcome this challenge, a conditional variant
 142 of the flow matching loss is introduced (Lipman et al., 2022; Liu et al., 2022):

$$143 \quad \mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, x_0, x_1, x_t} \|v_\theta(x_t, t) - (\alpha'_t x_0 + \beta'_t x_1)\|^2. \quad (4)$$

144 Minimizing \mathcal{L}_{CFM} is equivalent to minimizing the original objective \mathcal{L}_{FM} . This loss function is
 145 tractable since we only need to sample data-noise pairs to construct targets for the network.

147 4 SOLUTION FLOW MODELS

149 4.1 FORMULATIONS

151 The main purpose of this work is to investigate how can we train a neural network that can directly
 152 solve the velocity ODE, thereby eliminating the reliance on numerical ODE solvers in Flow Match-
 153 ing models. We study this problem under the assumption that the velocity field $v(x_t, t) \in \mathbb{R}^n$, with
 154 $x_t \in \mathbb{R}^n$ and $t \in [0, 1]$, is continuously differentiable and globally Lipschitz continuous with respect
 155 to x_t , i.e. $\|v(x_t, t) - v(x'_t, t)\|_2 \leq L\|x_t - x'_t\|_2, \forall x_t, x'_t \in \mathbb{R}^n, t \in [0, 1]$.

156 According to the existence and uniqueness theorem of ODEs, these two assumptions guarantee that
 157 given initial condition $X(t) = x_t$, where $x_t \in \mathbb{R}^n, t \in [0, 1]$ can be arbitrarily chosen, a unique
 158 solution $X(s), s \in [0, t]$ to the velocity ODE $\frac{dX(s)}{ds} = v(X(s), s)$ exists. Since the initial condition
 159 can be varied, we denote this unique solution by notation $f(x_t, t, s)$. Then we immediately have two
 160 equations by ODE's definition, for any $x_t \in \mathbb{R}^n, 0 \leq s \leq t \leq 1$:

$$161 \quad f(x_t, t, t) = x_t, \quad (5)$$

$$\partial_3 f(x_t, t, s) = v(f(x_t, t, s), s), \quad (6)$$

where $\partial_3 f(x_t, t, s) \in \mathbb{R}^n$ is the partial derivative function of $f(x_t, t, s)$ with respect to the third variable. Equivalently, the two equations can also be written as the following integral equation:

$$f(x_t, t, s) = x_t + \int_t^s v(f(x_t, t, u), u) du. \quad (7)$$

Since $f(x_t, t, s)$ maps an initial value x_t at time t to the unique solution of the velocity ODE at time s , we denote it as the solution function in this paper. Owing to the continuous differentiability of the velocity field $v(x_t, t)$, the solution function $f(x_t, t, s)$ is also continuously differentiable. To realize one-step generation, we need to train a model $f_\theta(x_t, t, s)$ to approximate the ground truth solution function $f(x_t, t, s)$, which is determined uniquely by the velocity field $v(x_t, t)$.

Under the setting discussed above, the following two conditions are sufficient to ensure that $f_\theta(x_t, t, s) = f(x_t, t, s)$ for all $x_t \in \mathbb{R}^n, 0 \leq s \leq t \leq 1$:

$$f_\theta(x_t, t, t) = x_t, \quad (8)$$

$$\partial_1 f_\theta(x_t, t, s) v(x_t, t) + \partial_2 f_\theta(x_t, t, s) = 0. \quad (9)$$

where $\partial_1 f_\theta(x_t, t, s) \in \mathbb{R}^{n \times n}$ is the Jacobian matrix function of $f_\theta(x_t, t, s)$ with respect to the first variable, $v(x_t, t)$ is multiplied by it as a matrix-vector product, and $\partial_2 f_\theta(x_t, t, s) \in \mathbb{R}^n$ is the partial derivative vector function of $f_\theta(x_t, t, s)$ with respect to the second variable.

This can be proven by expanding the following derivative for $0 \leq s \leq l \leq t \leq 1$:

$$\frac{d}{dl}(f_\theta(f(x_t, t, l), l, s)) = \partial_1 f_\theta(f(x_t, t, l), l, s) \underbrace{\partial_3 f(x_t, t, l)}_{=v(f(x_t, t, l), l)} + \partial_2 f_\theta(f(x_t, t, l), l, s) = 0. \quad (10)$$

The expression is equal to 0 due to the first condition in Eq. 8. Thus, we know that the model $f_\theta(x_t, t, s)$ is indeed the true solution function $f(x_t, t, s)$ by:

$$f_\theta(x_t, t, s) = f_\theta(f(x_t, t, t), t, s) = f_\theta(f(x_t, t, s), s, s) = f(x_t, t, s), \quad (11)$$

where we use the property that $f(x_t, t, t) = f_\theta(x_t, t, t) = x_t$ and $f_\theta(f(x_t, t, l), l, s)$ is invariant with respect to l , allowing us to set $l = t$ and $l = s$ to finish the proof.

4.2 LEARNING OBJECTIVES

Now we consider how to build efficient objectives for neural networks to learn the ground truth solution operator, according to the two conditions mentioned in the previous section. To ensure the first boundary condition of Eq. 8 is satisfied, we adopt the following parameterization:

$$f_\theta(x_t, t, s) = a(t, s)x_t + b(t, s)F_\theta(x_t, t, s), \quad (12)$$

where $a(t, s)$ and $b(t, s)$ are continuously differentiable scalar functions satisfying $a(t, t) = 1$ and $b(t, t) = 0$ for all $t \in [0, 1]$, and $F_\theta(x, t, s)$ denotes a raw neural network. Following the notations in the last section, we use $\partial_1 a(t, s), \partial_2 a(t, s), \partial_1 b(t, s), \partial_2 b(t, s)$ to denote the partial derivatives of $a(t, s)$ and $b(t, s)$ with respect to the first and second variables. In our experiments, we choose two specific parameterizations, including the Euler parameterization and the triangular parameterization:

$$f_\theta(x_t, t, s) = x_t + (s - t)F_\theta(x_t, t, s), \quad (13)$$

$$f_\theta(x_t, t, s) = \cos \frac{\pi}{2}(s - t)x_t + \sin \frac{\pi}{2}(s - t)F_\theta(x_t, t, s). \quad (14)$$

Obviously, these two parameterizations satisfy the boundary condition $f_\theta(x_t, t, t) = x_t$. Then, we construct two loss functions for training using Eq. 9.

Flow Matching Loss Firstly, we consider a special situation of Eq. 9 when $t = s$. Recall that the boundary condition Eq. 8 gives $f_\theta(x_t, t, t) = x_t, \forall x_t \in \mathbb{R}^n, t \in [0, 1]$, we have:

$$\partial_1 f_\theta(x_t, t, t) = I_n, 0 = \left. \frac{df_\theta(x_t, l, l)}{dl} \right|_{l=t} = \partial_2 f_\theta(x_t, t, t) + \partial_3 f_\theta(x_t, t, t) \quad (15)$$

where $I_n \in \mathbb{R}^{n \times n}$ is the identity matrix. Then Eq. 8 can be simplified to $v(x_t, t) = \partial_3 f_\theta(x_t, t, t)$. Under our parameterization mentioned above, we have

$$v(x_t, t) = \partial_3 f_\theta(x_t, t, t) = \partial_2 a(t, t)x_t + \partial_2 b(t, t)F_\theta(x_t, t, t) + \underbrace{b(t, t)}_0 \partial_3 F_\theta(x_t, t, t). \quad (16)$$

where the complex term containing $\partial_3 F_\theta(x_t, t, t)$ is canceled since $b(t, t) = 0$ by our choice of parameterization. Now we obtain a flow matching loss for our neural networks:

$$L_{\text{FM}}(\theta) = \mathbb{E}_{t, x_0, x_1, x_t} \left[\frac{w_{\text{FM}}(t, \text{MSE})}{D} \|\partial_2 a(t, t)x_t + \partial_2 b(t, t)F_\theta(x_t, t, t) - (\alpha'_t x_0 + \beta'_t x_1)\|_2^2 \right], \quad (17)$$

where D is the data dimension, $\alpha'_t x_0 + \beta'_t x_1$ is used to replace the intractable marginal velocity field $v(x_t, t) = \mathbb{E}_{p(x_0, x_1 | x_t)}[\alpha'_t x_0 + \beta'_t x_1]$ during the training process following the standard Flow Matching framework. In addition, this velocity term can also be provided by a teacher model in distillation situations, but in this paper, we focus on the situation of training from scratch.

Previous works (Geng et al., 2024; 2025) have demonstrated that choosing an adaptive weighting function is beneficial for few-step generative models. Following their approach, we choose

$$w_{\text{FM}}(t, \text{MSE}) = \frac{1}{|\partial_2 b(t, t)|(\text{MSE} + \epsilon)^p} \quad (18)$$

as our weighting function, where $|\partial_2 b(t, t)|$ is used to balance the raw network's gradients across time, and MSE represents the original mean squared error. Here ϵ is a smoothing factor to prevent excessively small values, and p is a factor that determines how robust the loss is. For $p = 0$, the objective degenerates to the mean squared error. For $p > 0$, this factor will penalize the data points with large errors in a data batch to make the objective more robust.

The original Flow Matching framework samples t uniformly, while more recent works (Esser et al., 2024; Geng et al., 2025) suggest sampling t from a logit-normal distribution. We also sample t from $\sigma(\mathcal{N}(\mu_{\text{FM}}, \sigma_{\text{FM}}^2))$, where $\sigma(\cdot)$ and \mathcal{N} represent the sigmoid function and normal distribution.

Solution Consistency Loss We now consider how to build a training target for the $s < t$ situation using Eq. 8. According to the Taylor's expansion, we have an approximation equation:

$$\frac{f_\theta(x_t, t, s) - f_\theta(x_t + v(x_t, t)(l - t), l, s)}{t - l} = (\partial_1 f_\theta(x_t, t, s)v(x_t, t) + \partial_2 f_\theta(x_t, t, s)) + o(1), \quad (19)$$

where $l \in (s, t)$ is close to t . Thus, we can adopt the following objective:

$$L_{\text{SCM}}(\theta) = \mathbb{E}_{\substack{t, l, s, \\ x_0, x_1, x_t}} \left[\frac{w_{\text{SCM}}(t, l, s, \text{MSE})}{D} \|f_\theta(x_t, t, s) - f_{\theta^-}(x_t + (\alpha'_t x_0 + \beta'_t x_1)(l - t), l, s)\|_2^2 \right], \quad (20)$$

where D is the data dimension, θ^- means taking the stop-gradient operation to the parameters, and $(\alpha'_t x_0 + \beta'_t x_1)$ is again used to replace the intractable $v(x_t, t) = \mathbb{E}_{p(x_0, x_1 | x_t)}[\alpha'_t x_0 + \beta'_t x_1]$ following the common practice. The adaptive weighting is chosen as follows:

$$w_{\text{SCM}}(t, l, s, \text{MSE}) = \frac{1}{(t - l)|b(t, s)|} \times \frac{1}{(\frac{\text{MSE}}{(t - l)^2} + \epsilon)^p}, \quad (21)$$

where the first term ensures the gradient magnitude of the raw network $F_\theta(x_t, t, s)$ is stable, and the second term is again used to provide a more robust loss by scaling down the coefficient for data points with large mean square errors when $p > 0$. Besides, we divide the mean square error in the adaptive term by $(t - l)^2$ since its magnitude is proportional to $(t - l)^2$.

As for the sampling method of t, l, s during training, we first sample t and s from two logit-normal distributions, $\sigma(\mathcal{N}(\mu_t, \sigma_t^2))$ and $\sigma(\mathcal{N}(\mu_s, \sigma_s^2))$, respectively. Here, s is clamped to ensure $s < t - 10^{-4}$. We then determine l using the following method:

$$l = t + (s - t) \times r(k, K), \quad (22)$$

where k, K represent the current and total training steps, respectively. The function $r(k, K)$ represents a monotonically decreasing schedule that gradually moves l towards t throughout training.

To avoid numerical issues, l is clamped to ensure $l < t - 10^{-4}$. This schedule decreases from an initial value r_{init} to an end value r_{end} . In our ablation studies, we test exponential, cosine, linear, and constant schedules. For more implementation details, please refer to Appendix A.

The total training loss is a combination of the flow matching and solution consistency losses: $L(\theta) = \lambda L_{\text{FM}}(\theta) + (1 - \lambda)L_{\text{SCM}}(\theta)$. The parameter λ controls the balance between them by determining the fraction of a data batch dedicated to computing $L_{\text{FM}}(\theta)$, while the remaining fraction is used for $L_{\text{SCM}}(\theta)$. We perform ablation studies to determine the optimal value of λ .

4.3 CLASSIFIER-FREE GUIDANCE

Classifier-Free Guidance (CFG) (Ho & Salimans, 2022) is a standard technique in diffusion models for enhancing conditional generation. The models are trained with randomly dropped conditions to mix conditional and unconditional data. During inference, CFG is applied by linearly combining predictions from the label-conditional and unconditional models to enhance generation quality.

To apply CFG for our models, we first introduce the ground-truth guided marginal velocity field:

$$v_g(x_t, t, c) = wv(x_t, t | c) + (1 - w)v(x_t, t), \quad (23)$$

where $v(x_t, t) = \mathbb{E}_{p(x_0, x_1 | x_t)}[\alpha'_t x_0 + \beta'_t x_1]$ denotes the unconditional marginal velocity field, c represents conditions (e.g., class labels), $v(x_t, t | c) = \mathbb{E}_{p(x_0, x_1 | x_t, c)}[\alpha'_t x_0 + \beta'_t x_1]$ denotes the conditional marginal velocity field, and w is the CFG strength. To ensure our conditional model serves as the solution function for the velocity ODE defined by this guided field (which depends purely on the data distribution and is model-independent), it should satisfy the following equations:

$$f_\theta(x_t, t, t, c) = x_t, \quad (24)$$

$$\partial_1 f_\theta(x_t, t, s, c)v_g(x_t, t, c) + \partial_2 f_\theta(x_t, t, s, c) = 0. \quad (25)$$

Analogous to the unconditional setting, we define a guided flow matching loss $L_{\text{FM}}^g(\theta)$ as:

$$\mathbb{E}_{t, x_t, c} \left[\frac{w_{\text{FM}}(t, \text{MSE})}{D} \|\partial_2 a(t, t)x_t + \partial_2 b(t, t)F_\theta(x_t, t, t, c) - v_g(x_t, t, c)\|_2^2 \right], \quad (26)$$

and a guided solution consistency loss $L_{\text{SCM}}^g(\theta)$ as:

$$\mathbb{E}_{t, l, s, x_t, c} \left[\frac{w_{\text{SCM}}(t, l, s, \text{MSE})}{D} \|f_\theta(x_t, t, s, c) - f_{\theta^-}(x_t + v_g(x_t, t, c)(l - t), l, s, c)\|_2^2 \right], \quad (27)$$

where D is the data dimension, θ^- denotes the stop-gradient operator applied to the targets, and the adaptive weighting functions remain consistent with the unconditional case. The total guided training loss $L^g(\theta)$ is defined as their linear combination $\lambda L_{\text{FM}}^g(\theta) + (1 - \lambda)L_{\text{SCM}}^g(\theta)$.

However, since the unconditional velocity field is generally intractable during conditional training, $v_g(x_t, t, c)$ is not directly accessible. To address this, we train the network to concurrently predict the unconditional velocity field. Specifically, we randomly replace the condition c with an empty label ϕ (with a probability of 0.1) and update the model using the unconditional flow matching loss (Eq. 17) and the unconditional solution consistency loss (Eq. 20). Consequently, we can approximate $v(x_t, t)$ using the model prediction $v_{\text{uncond}} = \partial_2 a(t, t)x_t + \partial_2 b(t, t)F_{\theta^-}(x_t, t, t, \phi)$.

For data points with a non-empty condition c , we compute the guided losses by substituting the velocity term $v_g(x_t, t, c)$ in Eq. 26 and Eq. 27 with the estimator $w(\alpha'_t x_0 + \beta'_t x_1) + (1 - w)v_{\text{uncond}}$, where w is the CFG strength. Here, $\alpha'_t x_0 + \beta'_t x_1$ is used to replace the intractable conditional marginal velocity field $v(x_t, t | c)$, similar to the unconditional formulation.

Notably, the term $w(\alpha'_t x_0 + \beta'_t x_1) + (1 - w)v_{\text{uncond}}$ typically exhibits higher variance than the original term $\alpha'_t x_0 + \beta'_t x_1$, primarily due to the scaling effect of w . This increased variance can hinder training convergence and degrade final performance. We observe that the model learns the guided velocity field via Eq. 27 when conditioned on c . Therefore, we can obtain a model-predicted guided velocity v_{guided} via $\partial_2 a(t, t)x_t + \partial_2 b(t, t)F_{\theta^-}(x_t, t, t, c)$, and then employ

$$v_{\text{mix}} = m(w(\alpha'_t x_0 + \beta'_t x_1) + (1 - w)v_{\text{uncond}}) + (1 - m)v_{\text{guided}} \quad (28)$$

to approximate the target $v_g(x_t, t, c)$, where $0 < m \leq 1$ acts as a velocity mixing ratio. Since the stochastic term $\alpha'_t x_0 + \beta'_t x_1$ is the dominant source of variance, using a small m effectively mitigates this issue, as the model-predicted v_{guided} possesses significantly lower variance.

Regarding the inference process of our model, it is straightforward since $f_\theta(x_t, t, s) = a(t, s)x_t + b(t, s)F_\theta(x_t, t, s)$ is a neural solution function to the velocity ODE, we only need to sample $x_1 \sim \mathcal{N}(0, I)$ and apply this operator with $t = 1$ and $s = 0$ to obtain the clean data in a single step. Furthermore, by selecting different t and s , our model can also achieve multi-step sampling.

5 EXPERIMENTS

5.1 SETTINGS

We conduct our major experiments on the ImageNet 256×256 dataset (Deng et al., 2009). SoFlow models operate within the latent space of a pre-trained VAE¹ (Rombach et al., 2022), which is a common practice in recent works (Peebles & Xie, 2023b; Geng et al., 2025). The tokenizer converts 256×256 images into a 32×32×4 latent representation. We assess generation quality using the Fréchet Inception Distance (FID) (Heusel et al., 2017), computed over a set of 50,000 generated samples. To evaluate computational efficiency, we report the number of function evaluations (NFE), with a particular focus on the single-step (1-NFE) scenario. All models presented were trained from scratch. In addition to the standard time variable t , our model incorporates an additional time variable s . We provide this to the networks by feeding the positional embeddings (Vaswani et al., 2017) of their difference, $s - t$. For more implementation details, please refer to Appendix A.

5.2 ABLATION STUDY

Table 1: **Ablation studies on ImageNet 256×256 class-conditional generation.** All models have 131M parameters and are trained with a batch size of 256 for 400K iterations from scratch.

$l \rightarrow t$ schedule	FID-50K (1-NFE)	flow matching data ratio	FID-50K (1-NFE)	loss weighting coefficient p	FID-50K (1-NFE)
Exponential	59.97	0%	66.55	0.0	69.11
Cosine	60.10	25%	61.53	0.5	60.59
Linear	60.64	50%	60.24	1.0	59.97
Constant	60.13	75%	59.97	1.5	64.25
	(a)		(b)		(c)
noising schedule & parameterization	FID-50K (1-NFE)	CFG guidance strength w	FID-50K (1-NFE)	velocity mix ratio m	FID-50K (1-NFE)
Linear, Euler	12.89	1.5	35.70	0.25	12.89
Linear, Triangular	23.57	2.0	22.37	0.5	15.67
Triangular, Euler	14.26	2.5	15.68	0.75	17.18
Triangular, Triangular	19.87	3.0	12.89	1.0	19.18
	(d)		(e)		(f)

We conduct various ablation studies to determine the optimal hyperparameters for SoFlow models. Following the methodology of Geng et al. (2025), we utilize the DiT-B/4 architecture in our experiments, which features a “base” sized diffusion transformer with a patch size of 4. We train our models for 400K iterations. For performance reference, the original DiT-B/4 (Peebles & Xie, 2023b) achieved an FID-50K of 68.4 with 250-NFE sampling. Meanflow (Geng et al., 2025) reports an FID-50K of 61.06 with 1-NFE sampling, and they claim to have reproduced SiT-B/4 (Ma et al., 2024) with an FID-50K of 58.9 using 250-NFE sampling. The FID-50K scores mentioned here are without classifier-free guidance (CFG). Our ablation studies are conducted in two groups: a first group of three experiments without classifier-free guidance (CFG) and a second group of three with CFG. Table 1 presents our results, which are analyzed as follows:

$l \rightarrow t$ Schedule We compare exponential, cosine, linear, and constant schedules in our experiments (Table 1a). The similar FID-50K results show that the speed of the $l \rightarrow t$ transition during the training process has a relatively small influence on performance. We choose the exponential schedule following previous works (Song & Dhariwal, 2023; Geng et al., 2024).

¹SD-VAE: <https://huggingface.co/stabilityai/sd-vae-ft-mse>

378 **Flow Matching Data Ratio** Although the solution consistency loss alone is sufficient to enable
 379 one-step generation (corresponding to a 0% ratio in Table 1b), experimental results show that incor-
 380 porating a large ratio of flow matching loss is effective for improving performance.
 381

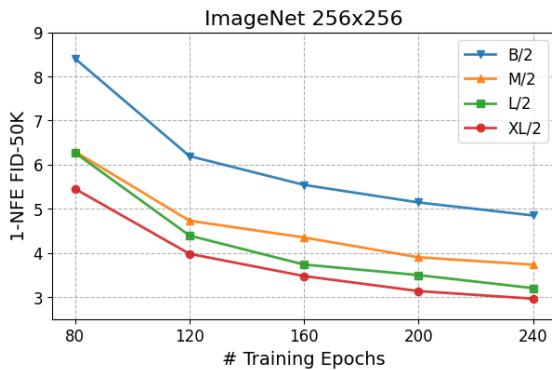
382 **Loss Weighting Coefficient p** We observe that the choice of loss metric has a significant influence
 383 on the performance of one-step generation, which aligns with previous works (Song & Dhariwal,
 384 2023; Geng et al., 2024; 2025). As shown in Table 1c, the original mean square loss with $p = 0$
 385 performs much worse than $p = 0.5$ or $p = 1$ situations.
 386

387 **Noising Schedule and Parameterization** Our method is compatible with different flow matching
 388 noising schedules and solution function parameterizations. As presented in Table 1d, experiments
 389 show that the linear noising schedule $x_t = (1-t)x_0 + tx_1$ along with the Euler parameterization
 390 $f_\theta(x_t, t, s) = x_t + (s-t)F_\theta(x_t, t, s)$ performs best, compared to the triangular noising schedule
 391 $x_t = \cos(\frac{\pi}{2}t)x_0 + \sin(\frac{\pi}{2}t)x_1$ and its corresponding parameterization $f_\theta(x_t, t, s) = \cos \frac{\pi}{2}(s-t)x_t +$
 392 $\sin \frac{\pi}{2}(s-t)F_\theta(x_t, t, s)$.
 393

393 **CFG Guidance Strength w** Unlike diffusion models, our model enables CFG during the training
 394 stage rather than the inference stage, which allows for image generation with 1 NFE. Experiments
 395 in Table 1e demonstrate that CFG can also greatly improve the generation quality for our model.
 396

397 **Velocity Mix Ratio m** Experiments (Table 1f) show that a relatively small m is beneficial for
 398 performance. This can be explained by the fact that the CFG guidance strength w amplifies the
 399 randomness in the velocity term, and a smaller m can suppress this randomness by partially replacing
 400 the velocity term with the model’s prediction of the guided velocity field.
 401

402 5.3 COMPARING WITH OTHER WORKS



404 Figure 3: **FID-50K performance of our**
 405 **models.** The plot shows the FID-50K
 406 performance of our models with vary-
 407 ing numbers of parameters, all trained
 408 from scratch on the ImageNet 256×256
 409 dataset. We apply CFG during train-
 410 ing and report the FID-50K scores of gen-
 411 erated images using a 1-NFE sam-
 412 pling pro-
 413 cess. The results consistently demon-
 414 strate that the performance of our models
 415 improves as the model size increases.
 416

417 **ImageNet 256×256 Results** We begin by analyzing the 1-NFE FID-50K results of our model
 418 across various model sizes and training durations, as shown in Figure 3. Our model’s 1-NFE per-
 419 formance gradually improves with an increase in model parameters, which is consistent with previous
 420 observations on diffusion transformers (Peebles & Xie, 2023b; Ma et al., 2024).
 421

422 Next, we compare our model’s one-step generation performance against previous models, with the
 423 results summarized in Table 2. To make relatively fair comparisons, we train our model with the
 424 same batch size and number of iterations as the Meanflow models (Geng et al., 2025). Experimental
 425 results show that our model consistently outperforms Meanflow models across all evaluated model
 426 sizes when trained from scratch. Specifically, for smaller models with DiT-B/2 and DiT-M/2 ar-
 427 chitectures, our model demonstrates significant improvements, achieving FID-50K scores of 4.85
 428 and 3.73, respectively. Our model also achieves superior FID-50K values for larger architectures,
 429 namely DiT-L/2 and DiT-XL/2, reaching 3.20 and 2.96, respectively. Notably, our models employ
 430 Classifier-Free Guidance (CFG) during training, which enables generation with exactly 1-NFE dur-
 431 ing inference, similar to Meanflow models.
 432

method	epochs	params	NFE	FID↓
Generative Adversarial Networks				
BigGAN (Brock et al., 2018)	-	112M	1	6.95
StyleGAN-XL (Sauer et al., 2022)	-	166M	1	2.30
GigaGAN (Kang et al., 2023)	-	569M	1×2	3.45
Masked and Autoregressive Models				
Mask-GIT (Chang et al., 2022)	555	227M	8	6.18
MagViT-v2 (Yu et al., 2023)	1080	307M	64	1.78
LlamaGen-XL (Sun et al., 2024)	300	775M	576	2.62
VAR (Tian et al., 2024)	350	2.0B	10	1.80
MAR (Li et al., 2024)	800	943M	64	1.55
RandAR-XL (Pang et al., 2025)	300	775M	256	2.22
Multi-step Diffusion Models				
LDM-4-G (Rombach et al., 2022)	170	395M	250×2	3.60
MDTv2 (Gao et al., 2023)	700	676M	250×2	1.63
DiT-XL/2 (Peebles & Xie, 2023a)	1400	675M	250×2	2.27
SiT-XL/2 (Ma et al., 2024)	1400	675M	250×2	2.06
FlowDCN-XL/2 (Wang et al., 2024b)	400	675M	250×2	2.00
SiT-REPA-XL/2 (Yu et al., 2024)	800	675M	250×2	1.42
Few-step Consistency Models				
iCT-XL/2 [†] (Song & Dhariwal, 2023)	-	675M	1 / 2	34.24 / 20.30
Shortcut-XL/2 (Frans et al., 2024)	250	675M	1 / 4	10.60 / 7.8
IMM-XL/2 (Zhou et al., 2025)	3840	675M	1×2 / 2×2	7.77 / 3.99
MeanFlow-B/2 (Geng et al., 2025)	240	131M	1	6.17
MeanFlow-M/2 (Geng et al., 2025)	240	308M	1	5.01
MeanFlow-L/2 (Geng et al., 2025)	240	459M	1	3.84
MeanFlow-XL/2 (Geng et al., 2025)	240	675M	1 / 2	3.43 / 2.93
SoFlow-B/2	240	131M	1 / 2	4.85 / 4.24
SoFlow-M/2	240	308M	1 / 2	3.73 / 3.42
SoFlow-L/2	240	459M	1 / 2	3.20 / 2.90
SoFlow-XL/2	240	675M	1 / 2	2.96 / 2.66

Table 2: **FID-50K results for class-conditional generation on ImageNet-256×256.** ×2 denotes an NFE of 2 per sampling step incurred by CFG. [†] Results as reported in (Zhou et al., 2025).

method	NFE	FID
iCT (Song & Dhariwal, 2023)	1	2.83
ECT (Geng et al., 2024)	1	3.60
sCT (Lu & Song, 2024)	1	2.97
IMM (Zhou et al., 2025)	1	3.20
MeanFlow (Geng et al., 2025)	1	2.92
SoFlow	1	2.90

Table 3: **FID-50K results for unconditional generation on CIFAR-10.**

CIFAR-10 Results In Table 3, we report unconditional generation results on the CIFAR-10 dataset, where performance is measured by the FID-50K metric with 1-NFE sampling. For our model, we adopt the U-Net architecture (Ronneberger et al., 2015) developed from (Song et al., 2020), aligning with prior works. Our method is applied directly to the pixel space, with a resolution of 32×32. For more implementation details, please refer to Appendix A. Our method achieves competitive performance compared to prior approaches on this dataset.

6 CONCLUSION

We have presented SoFlow, a simple yet effective framework for one-step generative modeling. Our approach directly learns the solution function of the velocity ODE, enabling single-step sampling without iterative solvers. By leveraging a bi-time formulation and a hybrid training objective combining a flow matching and a solution consistency loss, SoFlow naturally support CFG during training and avoid JVP that are not well-optimized on deep learning frameworks like PyTorch. Our method demonstrates competitive performance on class-conditioned Imagenet 256×256 generation task, outperforming Meanflow models when trained from scratch under the same settings.

486 REFERENCES
487

- 488 Michael S Albergo and Eric Vanden-Eijnden. Building normalizing flows with stochastic inter-
489 polants. *arXiv preprint arXiv:2209.15571*, 2022.
- 490 Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying
491 framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- 492 Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural
493 image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- 494 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe
495 Taylor, Troy Luhman, Eric Luhman, et al. Video generation models as world simulators. 2024.
496 *URL https://openai.com/research/video-generation-models-as-world-simulators*, 3:1, 2024.
- 497 Huiwen Chang, Han Zhang, Lu Jiang, Ce Liu, and William T Freeman. Maskgit: Masked generative
498 image transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
499 Recognition*, 2022.
- 500 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
501 hierarchical image database. 2009.
- 502 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
503 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers
504 for high-resolution image synthesis. In *Forty-first international conference on machine learning*,
505 2024.
- 506 Kevin Frans, Danijar Hafner, Sergey Levine, and Pieter Abbeel. One step diffusion via shortcut
507 models. *arXiv preprint arXiv:2410.12557*, 2024.
- 508 Shanghua Gao, Pan Zhou, Ming-Ming Cheng, and Shuicheng Yan. Mdtv2: Masked diffusion trans-
509 former is a strong image synthesizer. *arXiv preprint arXiv:2303.14389*, 2023.
- 510 Zhengyang Geng, Ashwini Pokle, and J Zico Kolter. One-step diffusion distillation via deep equi-
511 librium models. *Advances in Neural Information Processing Systems*, 36:41914–41931, 2023.
- 512 Zhengyang Geng, Ashwini Pokle, William Luo, Justin Lin, and J Zico Kolter. Consistency models
513 made easy. *arXiv preprint arXiv:2406.14548*, 2024.
- 514 Zhengyang Geng, Mingyang Deng, Xingjian Bai, J Zico Kolter, and Kaiming He. Mean flows for
515 one-step generative modeling. *arXiv preprint arXiv:2505.13447*, 2025.
- 516 Jonathan Heek, Emiel Hoogeboom, and Tim Salimans. Multistep consistency models. *arXiv
517 preprint arXiv:2403.06807*, 2024.
- 518 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter.
519 Gans trained by a two time-scale update rule converge to a local nash equilibrium. 2017.
- 520 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint
521 arXiv:2207.12598*, 2022.
- 522 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in
523 neural information processing systems*, 33:6840–6851, 2020.
- 524 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
525 Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–
526 8646, 2022.
- 527 Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung
528 Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference
529 on computer vision and pattern recognition*, pp. 10124–10134, 2023.
- 530 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
531 based generative models. *Advances in neural information processing systems*, 35:26565–26577,
532 2022.

- 540 Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka,
 541 Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency trajectory models: Learning proba-
 542 bility flow ode trajectory of diffusion. *arXiv preprint arXiv:2310.02279*, 2023.
- 543 Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Ad-*
 544 *vances in neural information processing systems*, 34:21696–21707, 2021.
- 545 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
 546 *arXiv:1412.6980*, 2014.
- 547 Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image
 548 generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:
 549 56424–56445, 2024.
- 550 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
 551 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 552 Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei
 553 Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*,
 554 2019.
- 555 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
 556 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- 557 Cheng Lu and Yang Song. Simplifying, stabilizing and scaling continuous-time consistency models.
 558 *arXiv preprint arXiv:2410.11081*, 2024.
- 559 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models: Synthe-
 560 sizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*, 2023a.
- 561 Weijian Luo, Tianyang Hu, Shifeng Zhang, Jiacheng Sun, Zhenguo Li, and Zhihua Zhang. Diff-
 562 instruct: A universal approach for transferring knowledge from pre-trained diffusion models.
 563 *Advances in Neural Information Processing Systems*, 36:76525–76546, 2023b.
- 564 Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and Sain-
 565 ing Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant
 566 transformers. In *European Conference on Computer Vision*, pp. 23–40. Springer, 2024.
- 567 Ziqi Pang, Tianyuan Zhang, Fujun Luan, Yunze Man, Hao Tan, Kai Zhang, William T Freeman, and
 568 Yu-Xiong Wang. Randar: Decoder-only autoregressive visual generation in random orders. In
 569 *Proceedings of the Computer Vision and Pattern Recognition Conference*, pp. 45–55, 2025.
- 570 William Peebles and Saining Xie. Scalable diffusion models with transformers. 2023a.
- 571 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of*
 572 *the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023b.
- 573 Yansong Peng, Kai Zhu, Yu Liu, Pingyu Wu, Hebei Li, Xiaoyan Sun, and Feng Wu. Flow-anchored
 574 consistency models. *arXiv preprint arXiv:2507.03738*, 2025.
- 575 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
 576 Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image
 577 synthesis. *arXiv preprint arXiv:2307.01952*, 2023.
- 578 A Polyak, A Zohar, A Brown, A Tjandra, A Sinha, A Lee, A Vyas, B Shi, CY Ma, CY Chuang, et al.
 579 Movie gen: A cast of media foundation models, 2025. URL <https://arxiv.org/abs/2410.13720>,
 580 pp. 51, 2024.
- 581 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
 582 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
 583 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 584 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed-
 585 ical image segmentation. In *International Conference on Medical image computing and computer-*
 586 *assisted intervention*, pp. 234–241. Springer, 2015.

- 594 Amirmojtaba Sabour, Sanja Fidler, and Karsten Kreis. Align your flow: Scaling continuous-time
 595 flow map distillation. *arXiv preprint arXiv:2506.14603*, 2025.
- 596
- 597 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar
 598 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic
 599 text-to-image diffusion models with deep language understanding. *Advances in neural informa-*
 600 *tion processing systems*, 35:36479–36494, 2022.
- 601 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv*
 602 *preprint arXiv:2202.00512*, 2022.
- 603
- 604 Axel Sauer, Katja Schwarz, and Andreas Geiger. Stylegan-xl: Scaling stylegan to large diverse
 605 datasets. In *ACM SIGGRAPH 2022 conference proceedings*, pp. 1–10, 2022.
- 606 Axel Sauer, Dominik Lorenz, Andreas Blattmann, and Robin Rombach. Adversarial diffusion dis-
 607 tillation. In *European Conference on Computer Vision*, pp. 87–103. Springer, 2024.
- 608
- 609 Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised
 610 learning using nonequilibrium thermodynamics. In *International conference on machine learn-*
 611 *ing*, pp. 2256–2265. pmlr, 2015.
- 612 Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. *arXiv*
 613 *preprint arXiv:2310.14189*, 2023.
- 614
- 615 Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution.
 616 *Advances in neural information processing systems*, 32, 2019.
- 617 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
 618 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*
 619 *arXiv:2011.13456*, 2020.
- 620
- 621 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. 2023.
- 622
- 623 Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan.
 624 Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint*
 625 *arXiv:2406.06525*, 2024.
- 626
- 627 Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling:
 628 Scalable image generation via next-scale prediction. *Advances in neural information processing*
 629 *systems*, 37:84839–84865, 2024.
- 630
- 631 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
 632 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural informa-*
 633 *tion processing systems*, 30, 2017.
- 634
- 635 Fu-Yun Wang, Zhengyang Geng, and Hongsheng Li. Stable consistency tuning: Understanding and
 636 improving consistency models. *arXiv preprint arXiv:2410.18958*, 2024a.
- 637
- 638 Shuai Wang, Zexian Li, Tianhui Song, Xubin Li, Tiezheng Ge, Bo Zheng, and Limin Wang. Explor-
 639 ing dcn-like architecture for fast image generation with arbitrary resolution. *Advances in Neural*
 640 *Information Processing Systems*, 37:87959–87977, 2024b.
- 641
- 642 Zidong Wang, Yiyuan Zhang, Xiaoyu Yue, Xiangyu Yue, Yangguang Li, Wanli Ouyang, and
 643 Lei Bai. Transition models: Rethinking the generative learning objective. *arXiv preprint*
 644 *arXiv:2509.04394*, 2025.
- 645
- 646 Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman,
 647 and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of*
 648 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6613–6623, 2024.
- 649
- 650 Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong
 651 Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion-
 652 tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023.

- 648 Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and
649 Saining Xie. Representation alignment for generation: Training diffusion transformers is easier
650 than you think. *arXiv preprint arXiv:2410.06940*, 2024.
- 651
- 652 Linqi Zhou, Stefano Ermon, and Jiaming Song. Inductive moment matching. *arXiv preprint*
653 *arXiv:2503.07565*, 2025.
- 654 Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity
655 distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation.
656 In *Forty-first International Conference on Machine Learning*, 2024.
- 657
- 658
- 659
- 660
- 661
- 662
- 663
- 664
- 665
- 666
- 667
- 668
- 669
- 670
- 671
- 672
- 673
- 674
- 675
- 676
- 677
- 678
- 679
- 680
- 681
- 682
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690
- 691
- 692
- 693
- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701

702 **A IMPLEMENTATION DETAILS**

704 We first provide a detailed description of the different scheduling strategies used to determine the
 705 intermediate time variable l during training. Recall that l is computed from t and s as follow:

$$706 \quad l = t + (s - t) \times r(k, K),$$

707 where k, K represent the current and total training steps, and the function $r(k, K)$ is a monotonically
 708 decreasing function that controls the progression of l from an initial value near t towards t itself as
 709 training advances. Below we specify the four schedules compared in our ablation studies:

710 **Exponential schedule:**

$$712 \quad r(k, K) = r_{\text{init}} \times \left(\frac{r_{\text{end}}}{r_{\text{init}}} \right)^{\frac{k}{K}},$$

714 **Cosine schedule:**

$$715 \quad r(k, K) = r_{\text{end}} + (r_{\text{init}} - r_{\text{end}}) \times \frac{1}{2} \left(1 + \cos \left(\pi \cdot \frac{k}{K} \right) \right),$$

717 **Linear schedule:**

$$718 \quad r(k, K) = r_{\text{init}} + (r_{\text{end}} - r_{\text{init}}) \times \frac{k}{K},$$

720 **Constant schedule:**

$$721 \quad r(k, K) = r_{\text{end}}.$$

722 In all cases, l is clamped to satisfy $l < t - 10^{-4}$ to ensure numerical stability, s is also clamped to
 723 satisfy $s < t - 10^{-4}$ at the same time. The initial value r_{init} and end value r_{end} are hyperparameters
 724 controlling the starting and final relative position between l and t .

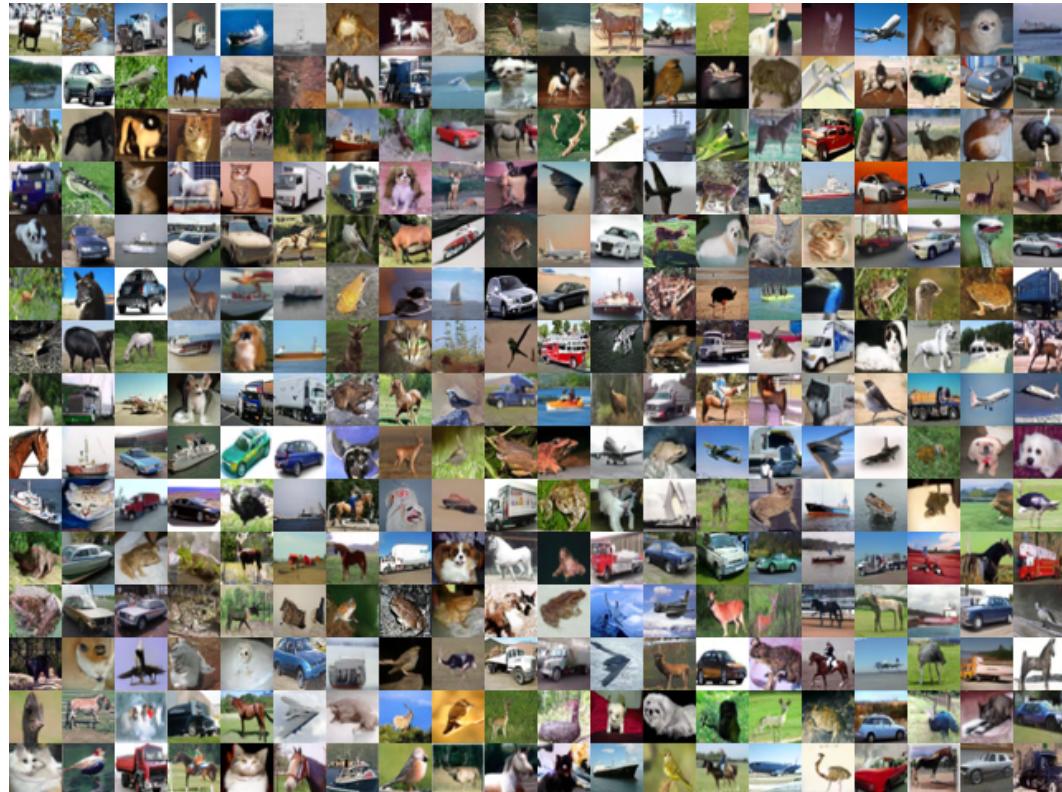
725 We now detail the training configurations for our model on two benchmark datasets. All experiments
 726 are run on NVIDIA H100 GPUs.

728 **ImageNet 256×256** We employ the Adam optimizer (Kingma & Ba, 2014) with a constant learning
 729 rate of 1×10^{-4} and betas set to $(0.9, 0.99)$, without learning rate decay or weight decay.
 730 Following standard practice, we evaluate model performance using Exponential Moving Average
 731 (EMA) with a decay rate of 0.9999. For time sampling, the logit-normal distribution parameters are
 732 set as: $\mu_{\text{FM}} = -0.2, \sigma_{\text{FM}} = 1.0; \mu_t = 0.2, \sigma_t = 0.8; \mu_s = -1.0, \sigma_s = 0.8$. The flow matching
 733 data ratio is 75% and the adaptive loss weight coefficient $p = 1.0$. The schedule parameters are
 734 $r_{\text{init}} = \frac{1}{10}$ and $r_{\text{end}} = \frac{1}{500}$. We use a linear noising schedule with Euler parameterization. The CFG
 735 strength w is set to 2.5, 2.25, 2.0, and 2.0 for the B/2, M/2, L/2, and XL/2 models, respectively,
 736 while the velocity mix ratio m is set to 0.25 for all models. Following common practice, we decay
 737 the CFG strength to 0 gradually for high-noise areas, where $t > 0.75$. Finally, architectural details
 738 are provided in Table 4.

739 Table 4: Architectural configurations on ImageNet 256×256.

741 architectures	SoFlow-B/4	SoFlow-B/2	SoFlow-M/2	SoFlow-L/2	SoFlow-XL/2
743 parameters (M)	131	131	308	459	676
744 FLOPs (G)	5.6	23.1	54.0	119.0	119.0
745 depth	12	12	16	24	28
746 hidden dimension	768	768	1024	1024	1152
747 attention heads	12	12	16	16	16
748 patch size	4×4	2×2	2×2	2×2	2 × 2
749 training epochs	80	240	240	240	240

751 **CIFAR-10** Training uses a batch size of 1024 for 800K iterations, consistent with MeanFlow. We
 752 adopt RAdam (Liu et al., 2019) with a learning rate of 1×10^{-4} , following (Song & Dhariwal, 2023;
 753 Geng et al., 2024). Time sampling parameters are: $\mu_{\text{FM}} = -0.9, \sigma_{\text{FM}} = 1.6; \mu_t = -0.9, \sigma_t = 1.6;$
 754 $\mu_s = -4.0, \sigma_s = 1.6$. The flow matching data ratio is 0% with $p = 0.5$. Schedule parameters are
 755 set to smaller values than ImageNet: $r_{\text{init}} = \frac{1}{50}$ and $r_{\text{end}} = \frac{1}{3200}$. Linear noising schedule and Euler
 756 parameterization are used.

756
757
758
759
760
761
B MORE VISUAL SAMPLES775
776
Figure 4: We show curated samples from class-conditioned generation by the DiT-XL/2 model with
777
1-NFE on the ImageNet 256×256 dataset.808
809
Figure 5: We show uncurated samples from unconditional generation by the UNet model with 1-
NFE on the CIFAR-10 dataset.

810 **C PSEUDO-CODE FOR GUIDED TRAINING PROCESS**
811

812 **Algorithm 1** Train Solution Flow Models with CFG
813

814 **Input:** model $f_\theta(x_t, t, s, c)$, flow matching data ratio λ , CFG strength w , velocity mix ratio m
815 Sample a data batch: $x_0, c \sim p_{data}(x_0, c)$
816 Split data for two losses by λ : $(x_0^{\text{FM}}, x_0^{\text{SCM}}) = x_0, (c^{\text{FM}}, c^{\text{SCM}}) = c$
817 Sample $t_{\text{FM}}, t_{\text{SCM}}, l_{\text{SCM}}, s_{\text{SCM}}$ according to subsection 4.2
818 Get $x_t^{\text{FM}}, v_t^{\text{FM}}, x_t^{\text{SCM}}, v_t^{\text{SCM}}$ by standard Flow Matching framework
819 Compute $v_{\text{mix}}^{\text{FM}}$ and $v_{\text{mix}}^{\text{SCM}}$ by Eq. 28 with w and m
820 Randomly replace $v_{\text{mix}}^{\text{FM}}$ and $v_{\text{mix}}^{\text{SCM}}$ by v_t^{FM} and v_t^{SCM} with CFG drop rate 0.1
821 Compute $L_{\text{FM}}(\theta)$ according to Eq. 26 by replacing $v_g(x_t, t, c)$ term with $v_{\text{mix}}^{\text{FM}}$
822 Compute $L_{\text{SCM}}(\theta)$ according to Eq. 27 by replacing $v_g(x_t, t, c)$ term with $v_{\text{mix}}^{\text{SCM}}$
823 Update f_θ via gradient descent according to $L(\theta) = \lambda L_{\text{FM}}(\theta) + (1 - \lambda) L_{\text{SCM}}(\theta)$
824

825

826
827 **D LARGE LANGUAGE MODEL USAGE**
828

829 We only adopt large language models to polish the writing.
830

831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863