

UNIFIED CONTINUOUS GENERATIVE MODELS FOR DENOISING-BASED DIFFUSION

Anonymous authors

Paper under double-blind review

ABSTRACT

Recent advances in continuous generative models, encompassing multi-step processes such as diffusion and flow matching (typically requiring 8-1000 steps) and few-step methods such as consistency models (typically 1-8 steps), have yielded impressive generative performance. However, existing work often treats these approaches as distinct paradigms, leading to disparate training and sampling methodologies. We propose a unified framework for the training, sampling, and analysis of diffusion, flow matching, and consistency models. Within this framework, we derive a surrogate unified objective that, for the first time, theoretically shows that the few-step objective can be viewed as the multi-step objective plus a regularization term. Building on this framework, we introduce the Unified Continuous Generative Models Trainer and Sampler (UCGM- $\{T, S\}$), which enables efficient and stable training of both multi-step and few-step models. Empirically, our framework achieves state-of-the-art results. On ImageNet 256×256 with a 675M diffusion transformer, UCGM-T trains a multi-step model achieving 1.30 FID in 20 steps, and a few-step model achieving 1.42 FID in only 2 steps. Moreover, applying UCGM-S to REPA-E (Leng et al., 2025) improves its FID from 1.26 (at 250 steps) to 1.06 in only 40 steps, without additional cost.

1 INTRODUCTION



(a) NFE = 40, FID = 1.48.

(b) NFE = 2, FID = 1.75.

Figure 1: **Generated samples from two 675M diffusion transformers trained with our UCGM on ImageNet-1K 512×512 .** The figure showcases generated samples illustrating the flexibility of Number of Function Evaluation (NFE) and superior performance achieved by our UCGM. The left subfigure presents results with NFE = 40 (multi-step), while the right subfigure shows results with NFE = 2 (few-step). Note that the samples are sampled *without classifier-free guidance (CFG) or other guidance techniques*.

Continuous generative models, encompassing diffusion models (Ho et al., 2020; Song et al., 2020a), flow-matching models (Lipman et al., 2022; Ma et al., 2024), and consistency models (Song et al., 2023; Lu & Song, 2024), have demonstrated remarkable success in synthesizing high-fidelity data across diverse applications, including image and video generation (Peebles & Xie, 2023; Chen et al., 2024c; Ma et al., 2024; Xie et al., 2024a; Ho et al., 2022; Chen et al., 2025c).

Training and sampling of these models necessitate substantial computational resources (Karras et al., 2022; 2024b). Moreover, current research treats distinct model paradigms (diffusion models/flow matching (Karras et al., 2022) v.s. consistency models (Song et al., 2023)) independently, leading to paradigm-specific training and sampling methodologies. This fragmentation introduces two primary challenges: (a) **a deficit in unified theoretical and empirical understanding**, which constrains the

Table 1: **Existing continuous generative paradigms as special cases of our UCGM.** Prominent continuous generative models, such as Diffusion, Flow Matching, and Consistency models, can be formulated as specific parameterizations of our UCGM. The columns detail the required parameterizations for the transport coefficients $\alpha(\cdot)$, $\gamma(\cdot)$, $\hat{\alpha}(\cdot)$, $\hat{\gamma}(\cdot)$ and parameters λ , ρ , ν of UCGM. Note that $\sigma(t)$ is defined as $e^{4(2.68t-1.59)}$ in this table.

Paradigm		UCGM-based Parameterization						
Type	e.g.,	$\alpha(t) =$	$\gamma(t) =$	$\hat{\alpha}(t) =$	$\hat{\gamma}(t) =$	$\lambda \in [0, 1]$	$\rho \in [0, 1]$	$\nu \in \{1, 2\}$
Diffusion	EDM (Karras et al., 2022)	$\frac{\sigma(t)}{\sqrt{\sigma^2(t)+\frac{1}{4}}}$	$\frac{1}{\sqrt{\sigma^2(t)+\frac{1}{4}}}$	$\frac{-0.5}{\sqrt{\sigma^2(t)+\frac{1}{4}}}$	$\frac{2\sigma(t)}{\sqrt{\sigma^2(t)+\frac{1}{4}}}$	0	≥ 0	2
Flow Matching	FM (Lipman et al., 2022)	t	$1-t$	1	-1	0	≥ 0	1
Consistency	sCM (Lu & Song, 2024)	$\sin(t \cdot \frac{\pi}{2})$	$\cos(t \cdot \frac{\pi}{2})$	$\cos(t \cdot \frac{\pi}{2})$	$\sin(t \cdot \frac{\pi}{2})$	1	1	1

transfer of advancements across different paradigms; and (b) **limited cross-paradigm generalization**, as algorithms optimized for one paradigm (e.g., diffusion models) are often incompatible with others. To address these limitations, we introduce UCGM, a novel framework that establishes a unified foundation for the theoretical understanding, training and sampling of continuous generative models (diffusion, flow matching, and consistency models). Within this framework, we derive a surrogate unified objective, which not only offers a formulation equivalent to the unified objective, but also, for the first time, shows that the few-step objective can be viewed as the multi-step objective plus a self-consistency term. Within this formulation, we link the instability of few-step model training to the self-alignment term that dominates the training dynamics as $\lambda \rightarrow 1$.

The unified trainer UCGM-T is built upon a unified objective, parameterized by a consistency ratio $\lambda \in [0, 1]$. This allows a single training paradigm to flexibly produce models tailored for different inference regimes: models behave akin to multi-step diffusion or flow-matching approaches when λ is close to 0, and transition towards few-step consistency-like models as λ approaches 1. Furthermore, our unified framework supports compatibility with diverse noise schedules (e.g., linear, triangular, quadratic) without requiring algorithm-specific modifications.

Complementing UCGM-T, we propose a unified sampler UCGM-S that operates seamlessly with models trained under our objective. UCGM-S is designed to enhance and accelerate sampling from pre-trained models—including those from previous paradigms as well as ones trained via UCGM-T. The unifying power of UCGM is further demonstrated by its ability to encapsulate several major continuous generative paradigms as special instances, as summarized in Tab. 1. Moreover, as shown in Fig. 1, models trained with UCGM achieve high sample quality across a wide range of Number of Function Evaluations (NFEs).

In summary, our contributions are:

- We propose a unified framework that provides a theoretical foundation for the training and sampling of continuous generative models—including diffusion models, flow matching models, and consistency models—and derive a surrogate unified objective that, for the first time, theoretically shows that the few-step objective can be viewed as the multi-step objective plus a self-alignment term.
- We introduce a unified trainer UCGM-T, that seamlessly bridges few-step (e.g., consistency models) and multi-step (e.g., diffusion, flow matching) generative paradigms, accommodating diverse model architectures, latent autoencoders, and noise schedules. We also propose a unified sampler UCGM-S, which is compatible with our trained models and further accelerate and improve pre-trained models from existing yet distinct paradigms.
- We empirically validate the effectiveness and efficiency of UCGM. Our approach consistently matches or surpasses SOTA methods across various datasets, architectures, and resolutions, for both few-step and multi-step generation tasks (cf., the experimental results in Sec. 4).

2 PRELIMINARIES

Given a training dataset \mathcal{D} , let $p_{\text{data}}(\mathbf{x})$ represent its underlying data distribution, or $p_{\text{data}}(\mathbf{x}|\mathbf{c})$ under a condition \mathbf{c} . Continuous generative models seek to learn an estimator that gradually transforms a simple source distribution $p_{\mathbf{z}}(\mathbf{z})$ into a complex target distribution $p_{\text{data}}(\mathbf{x})$ within a continuous space. Typically, $p_{\mathbf{z}}(\mathbf{z})$ is represented by the standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. For brevity, we hereafter omit subscripts when the context is clear, and assume independence, i.e., $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x})p(\mathbf{z})$. For instance, diffusion models reverse a noising process that gradually perturbs a data sample $\mathbf{x} \sim p(\mathbf{x})$ into a noisy version $\mathbf{x}_t = \alpha(t)\mathbf{x} + \sigma(t)\mathbf{z}$, where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Over the range

$t \in [0, T]$, the perturbation intensifies with increasing t , where higher t values indicate more noise. Below, we introduce three learning paradigms for continuous generative models.

Diffusion models (Ho et al., 2020; Song et al., 2020b; Karras et al., 2022). In the widely adopted EDM method (Karras et al., 2022), the noising process is defined by setting $\alpha(t) = 1$, $\sigma(t) = t$. The training objective is given by $\mathbb{E}_{\mathbf{x}, \mathbf{z}, t} \left[\omega(t) \|\mathbf{f}_\theta(\mathbf{x}_t, t) - \mathbf{x}\|_2^2 \right]$ where $\omega(t)$ is a weighting function. The diffusion model is parameterized by $\mathbf{f}_\theta(\mathbf{x}_t, t) = c_{\text{skip}}(t)\mathbf{x}_t + c_{\text{out}}(t)\mathbf{F}_\theta(c_{\text{in}}(t)\mathbf{x}_t, c_{\text{noise}}(t))$ where \mathbf{F}_θ is a neural network, and the coefficients c_{skip} , c_{out} , c_{in} , and c_{noise} are manually designed. During sampling, EDM solves the Probability Flow Ordinary Differential Equation (PF-ODE) (Song et al., 2020b): $\frac{d\mathbf{x}_t}{dt} = [\mathbf{x}_t - \mathbf{f}_\theta(\mathbf{x}_t, t)]/t$, integrated from $t = T$ to $t = 0$.

Flow matching (Lipman et al., 2022). Flow matching models are similar to diffusion models but differ in the transport process from the source to the target distribution and in the neural network training objective. The forward transport process utilizes differentiable coefficients $\alpha(t)$ and $\gamma(t)$, such that $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$. Typically, the coefficients satisfy the boundary conditions $\alpha(1) = \gamma(0) = 1$ and $\alpha(0) = \gamma(1) = 0$. The training objective is given by $\mathbb{E}_{\mathbf{x}, \mathbf{z}, t} \left[\omega(t) \left\| \mathbf{F}_\theta(\mathbf{x}_t, t) - \left(\frac{d\alpha_t}{dt}\mathbf{z} + \frac{d\gamma_t}{dt}\mathbf{x} \right) \right\|_2^2 \right]$. Similar to diffusion models, the reverse transport process (i.e., sampling process) begins at $t = 1$ with $\mathbf{x}_1 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and solves the PF-ODE: $\frac{d\mathbf{x}_t}{dt} = \mathbf{F}_\theta(\mathbf{x}_t, t)$, integrated from $t = 1$ to $t = 0$.

Consistency models (Song et al., 2023; Lu & Song, 2024). A consistency model $\mathbf{f}_\theta(\mathbf{x}_t, t)$ is trained to map the noisy input \mathbf{x}_t directly to the corresponding clean data \mathbf{x} in one or few steps by following the sampling trajectory of the PF-ODE starting from \mathbf{x}_t . To be valid, \mathbf{f}_θ must satisfy the boundary condition $\mathbf{f}_\theta(\mathbf{x}, 0) \equiv \mathbf{x}$. Inspired by EDM (Karras et al., 2022), one approach to enforce this condition is to parameterize the consistency model as $\mathbf{f}_\theta(\mathbf{x}_t, t) = c_{\text{skip}}(t)\mathbf{x}_t + c_{\text{out}}(t)\mathbf{F}_\theta(c_{\text{in}}(t)\mathbf{x}_t, c_{\text{noise}}(t))$ with $c_{\text{skip}}(0) = 1$ and $c_{\text{out}}(0) = 0$. The training objective is defined between two adjacent time steps with a finite distance: $\mathbb{E}_{\mathbf{x}_t, t} [\omega(t)d(\mathbf{f}_\theta(\mathbf{x}_t, t), \mathbf{f}_{\theta^-}(\mathbf{x}_{t-\Delta t}, t-\Delta t))]$, where θ^- denotes $\text{stopgrad}(\theta)$, $\Delta t > 0$ is the distance between adjacent time steps, and $d(\cdot, \cdot)$ is a metric function. Discrete-time consistency models are sensitive to the choice of Δt , necessitating manually designed annealing schedules (Song & Dhariwal, 2023; Geng et al., 2024) for rapid convergence. This limitation is addressed by proposing a training objective for continuous consistency models (Lu & Song, 2024), derived by taking the limit as $\Delta t \rightarrow 0$.

3 METHODOLOGY

This section elaborates on our two primary contributions: (1) the unified framework for continuous generation models and a surrogate loss function that affords a theoretical interpretation of model behavior. (2) the concrete instantiation of the unified framework through UCGM-T (for training) and UCGM-S (for sampling).

3.1 UNIFIED FRAMEWORK FOR CONTINUOUS GENERATIVE MODELS

We first propose a unified multi-step objective for diffusion and flow-matching models, which constitute all multi-step continuous generative models. Furthermore, we extend this unified multi-step objective to encompass both few-step models and multi-step models.

Unified objective for multi-step continuous generative models. We introduce a generalized training objective below that effectively trains generative models while encompassing the formulations presented in existing studies (Karras et al., 2022; Lipman et al., 2022; Liu et al., 2022; Ho et al., 2020; Song et al., 2020a):

$$\mathcal{L}(\theta) := \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\omega(t)} \|\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{z}_t\|_2^2 \right], \quad (1)$$

where time $t \in [0, 1]$, $\omega(t)$ is the weighting function for the loss, \mathbf{F}_θ is a neural network¹ with parameters θ , $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$, and $\mathbf{z}_t = \hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}$. Here, $\alpha(t)$, $\gamma(t)$, $\hat{\alpha}(t)$, and $\hat{\gamma}(t)$ are the unified transport coefficients defined for UCGM. In this paper, we refer to equation (1) as **the multi-step objective**. Additionally, to efficiently and robustly train multi-step continuous generative models using (1), we propose the following *necessary assumption*:

¹For simplicity, unless otherwise specified, we assume that any conditioning information \mathbf{c} is incorporated into the network input. Thus, $\mathbf{F}_\theta(\mathbf{x}_t, t)$ should be understood as $\mathbf{F}_\theta(\mathbf{x}_t, t, \mathbf{c})$ when \mathbf{c} is applicable.

Assumption 1 . The coefficients function $\alpha(t), \gamma(t), \hat{\alpha}(t), \hat{\gamma}(t)$ satisfy the following constraints:

- (a) $\alpha(t) \in C^1[0, 1]$ and is non-decreasing, with $\alpha(0) = 0, \alpha(1) = 1$.
- (b) $\gamma(t) \in C^1[0, 1]$ and is non-increasing, with $\gamma(0) = 1, \gamma(1) = 0$.
- (c) $\forall t \in [0, 1], |\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)| > 0$.

Under the **Assump. 1**, diffusion and flow matching are special cases of multi-step objective (1):

- (a) **Diffusion**: following EDM (Karras et al., 2022; 2024b), by setting $\alpha(t) = 1$ and $\sigma(t) = t$, diffusion models based on EDM can be derived from (1) provided that the constraint $\gamma(t)/\alpha(t) = t$ is satisfied².
- (b) **Flow Matching**: Similarly, flow matching can be derived only when $\hat{\alpha}(t) = \frac{d\alpha(t)}{dt}$ and $\hat{\gamma}(t) = \frac{d\gamma(t)}{dt}$ (see **Sec. 2** for more technical details about EDM-based and flow-based models).

Unified objective for both multi-step and few-step models. To facilitate the interpretation of our framework, we define two prediction functions based on model F_θ as:

$$\mathbf{f}^x(F_\theta^t, \mathbf{x}_t, t) := \frac{\alpha(t) \cdot F_\theta^t - \hat{\alpha}(t) \cdot \mathbf{x}_t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} \quad \& \quad \mathbf{f}^z(F_\theta^t, \mathbf{x}_t, t) := \frac{\hat{\gamma}(t) \cdot \mathbf{x}_t - \gamma(t) \cdot F_\theta^t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)}, \quad (2)$$

where we define $F_\theta^t := F_\theta(\mathbf{x}_t, t)$. The training objective (1) thus becomes (cf., **App. F.1.1**):

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{\omega}(t)} \|\mathbf{f}^x(F_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{x}\|_2^2 \right]. \quad (3)$$

To align with the gradient of multi-step objective (1), we define a new weighting function $\hat{\omega}(t)$ in (3) as $\hat{\omega}(t) := \frac{\alpha(t) \cdot \alpha(t) \cdot \omega(t)}{(\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t))^2}$. To unify few-step models (such as consistency models) with multi-step models, we adopt a modified version of (3) by incorporating a consistency ratio $\lambda \in [0, 1]$:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{\omega}(t)} \|\mathbf{f}^x(F_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^x(F_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t)\|_2^2 \right], \quad (4)$$

where consistency models and conventional multi-steps models are special cases within the context of (4) (cf., **App. F.1.1** and **App. F.1.3**):

- (a) **Diffusion / Flow Matching**: setting $\lambda = 0$ yields diffusion and flow matching, and our unified objective (4) degrades to the objective (3), which is equivalent to the multi-step objective (1).
- (b) **Consistency Model**: setting $\lambda = 1 - \frac{\Delta t}{t}$ with $\Delta t \rightarrow 0$ recovers consistency models.

Equivalent surrogate objective for unified objective (4). Building on the unified objective (4), we derive an equivalent surrogate objective. Importantly, this surrogate not only provides an equivalent reformulation of the unified objective but also sheds light on the theoretical origin of instability in few-step models, like consistency model.

Theorem 1 (Surrogate objective for unified objective of linear case ($\alpha(t) = t, \gamma(t) = 1 - t$)). Under **Assump. 1**, let's consider a surrogate objective

$$\mathcal{G}(\theta) = \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\underbrace{\|F_\theta(\mathbf{x}_t, t) - \mathbf{z}_t\|_2^2}_{\text{Flow Matching Objective}} + \frac{\lambda}{1 - \lambda} \underbrace{\|F_\theta(\mathbf{x}_t, t) - F_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t)\|_2^2}_{\text{Self-Alignment Term}} \right], \quad (5)$$

where $\mathbf{x}_t = t \cdot \mathbf{z} + (1 - t) \cdot \mathbf{x}$, $\mathbf{z}_t = \mathbf{z} - \mathbf{x}$, $\hat{\omega}(t) = t^2 \cdot (1 - \lambda)$, $0 < \lambda < 1$. The following equation holds: $\nabla_\theta \mathcal{L}(\theta) = \nabla_\theta \mathcal{G}(\theta)$, $\forall \theta$. See **App. F.1.5** for proof and general case.

Thm. 1 establishes that optimizing the unified objective in (4) is equivalent to optimizing the surrogate objective in (5). This equivalence is useful for analysis because the surrogate, $\mathcal{G}(\theta)$, can be decomposed into two distinct components: a multi-step objective term and a self-alignment term. We can offer a physical interpretation for each component by considering the underlying function $F_\theta(\mathbf{x}_t, t)$ as a learned velocity field:

- **Flow matching objective**: This term corresponds to the learning objective of multi-step models (1). It learns the mean velocity $\mathbf{z}_t = \mathbf{z} - \mathbf{x} = \frac{\mathbf{x}_1 - \mathbf{x}_0}{1 - 0}$ of a flow trajectory.

²In EDM, with $\sigma(t) = t$, the input of neural network F_θ is $c_{in}(t)\mathbf{x}_t = c_{in}(t) \cdot (\mathbf{x} + t \cdot \mathbf{z})$. Although $c_{in}(t)$ can be manually adjusted, the coefficient before \mathbf{z} remains t times that of \mathbf{x} .

- **Self-alignment term:** This term can be considered as a regularization term, which enforces consistency of the velocity of any points within a flow trajectory, ultimately helping to straighten the learned trajectories.

Remark 1 (Analysis of instability of few-step objective (i.e. $\lambda \rightarrow 1$)). According to *Thm. 1*, as $\lambda \rightarrow 1$, the self-alignment term dominates the loss function. This term only requires the velocity to be consistent in each flow trajectory, without constraining it to match the mean velocity. Thus, while a pre-trained velocity field may initially be straightened under this objective, prolonged training with few-step objective ultimately degrades the quality of the velocity field.

3.2 INSTANTIATING THE UNIFIED FRAMEWORK FOR TRAINING (UCGM-T)

Applying the gradient identity from Lu & Song (2024)³, we derive the unified objective:

$$\mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_t, t) + 2 \cdot \frac{\Delta \mathbf{f}^{\mathbf{x}}}{B(t) - B(\lambda t)} \right\|_2^2 \right], \quad (6)$$

where the detailed derivation from (4) to (6) is provided in [App. F.1.7](#), and

$$\Delta \mathbf{f}^{\mathbf{x}} := \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t^{\theta^-}, \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\lambda t}^{\theta^-}, \mathbf{x}_{\lambda t}, \lambda t), \quad B(t) := \alpha(t) / (\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)).$$

Second-order estimator as $\lambda \rightarrow 1$. We identify that the direct estimation of the difference quotient in objective (6) is only a first-order approximation, which is susceptible to numerical precision errors. To mitigate this issue, we propose a second-order estimator:

$$\frac{\Delta \mathbf{f}^{\mathbf{x}}}{B(t) - B(\lambda t)} \approx \frac{\mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{t+\epsilon}, t + \epsilon), \mathbf{x}_{t+\epsilon}, t + \epsilon) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{t-\epsilon}, t - \epsilon), \mathbf{x}_{t-\epsilon}, t - \epsilon)}{B(t + \epsilon) - B(t - \epsilon)}.$$

See [App. F.2.3](#) for further analysis of this second-order estimator. To stabilize the training, we implement two strategies for the second-order estimation: (1) We adopt a distributive reformulation of the second-difference term to prevent direct subtraction $\Delta \mathbf{f}_t^{\mathbf{x}} = \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{t+\epsilon}, t + \epsilon), \mathbf{x}_{t+\epsilon}, t + \epsilon) \cdot \frac{1}{2\epsilon} - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{t-\epsilon}, t - \epsilon), \mathbf{x}_{t-\epsilon}, t - \epsilon) \cdot \frac{1}{2\epsilon}$. (2) we also observe that applying numerical truncation $\text{clip}(\cdot, -1, 1)$ to the second-order estimator enhances training stability (Lu & Song, 2024).

Generalized time distribution (GTD) $\text{Beta}(\theta_1, \theta_2)$. Previous studies (Yao et al., 2025; Esser et al., 2024; Song et al., 2023; Lu & Song, 2024; Karras et al., 2022; 2024b) employ non-linear functions to transform the time variable t , initially sampled from a uniform distribution $t \sim \mathcal{U}(0, 1)$. This transformation shifts the distribution of sampled times, effectively performing importance sampling and thereby accelerating the training convergence rate. For example, the `lognorm` function $f_{\text{lognorm}}(t; \mu, \sigma) = 1 / (1 + \exp(-\mu - \sigma \cdot \Phi^{-1}(t)))$ is widely used (Yao et al., 2025; Esser et al., 2024), where $\Phi^{-1}(\cdot)$ denotes the inverse Cumulative Distribution Function of the standard normal distribution.

In this work, we demonstrate that commonly used time distribution after non-linear time transformation can be well-approximated by the Beta distribution (a detailed analysis is provided in [App. F.2.1](#)). Consequently, we simplify the process by directly sampling time from a Beta distribution, i.e., $t \sim \text{Beta}(\theta_1, \theta_2)$, where θ_1 and θ_2 are parameters that control the shape of time distribution (see [App. D.1.3](#) for their settings).

Learning enhanced target score function. We additionally incorporate the enhanced target score function proposed in recent work (Tang et al., 2025) into our unified training objective in (6). This technique is not our main contribution but can be seamlessly integrated into our framework. For completeness, we provide the formulation and further analysis in [App. F.1.8](#).

An ablation study for our proposed techniques is shown in [Tab. 13](#), and the pseudocode is in [Alg. 1](#).

3.3 INSTANTIATING THE UNIFIED FRAMEWORK FOR SAMPLING (UCGM-S)

For classical iterative sampling models, such as a trained flow-matching model \mathbf{f}_{θ} , sampling from the learned distribution $p(\mathbf{x})$ involves solving the PF-ODE (Song et al., 2020b). This process typically uses numerical ODE solvers, such as the Euler or Runge-Kutta methods (Ma et al., 2024), to iteratively transform the initial Gaussian noise $\tilde{\mathbf{x}}$ into a sample from $p(\mathbf{x})$ by solving the ODE (i.e., $\frac{d\tilde{\mathbf{x}}_t}{dt} = \mathbf{f}_{\theta}(\tilde{\mathbf{x}}_t, t)$). Similarly, sampling processes in models like EDM (Karras et al., 2022; 2024b) and

³ $\nabla_{\theta} \mathbb{E}[\mathbf{F}_{\theta}^{\top} \mathbf{y}] = \frac{1}{2} \nabla_{\theta} \mathbb{E}[\|\mathbf{F}_{\theta} - \mathbf{F}_{\theta^-} + \mathbf{y}\|_2^2]$.

consistency models (Song et al., 2023) involve a comparable gradual denoising procedure. Building on these observations and our unified trainer UCGM-T, we first propose a general iterative sampling process with two stages below:

- (a) **Decomposition:** At time t , the current input $\tilde{\mathbf{x}}_t$ is decomposed into two components: $\tilde{\mathbf{x}}_t = \alpha(t) \cdot \hat{\mathbf{z}}_t + \gamma(t) \cdot \tilde{\mathbf{x}}_t$. This decomposition uses the estimation model F_θ . Specifically, the model output $F_t = F_{\theta^-}(\tilde{\mathbf{x}}_t, t)$ is computed, yielding the estimated clean component $\hat{\mathbf{x}}_t = f^x(F_t, \tilde{\mathbf{x}}_t, t)$ and the estimated noise component $\hat{\mathbf{z}}_t = f^z(F_t, \tilde{\mathbf{x}}_t, t)$.
- (b) **Reconstruction:** The next time step’s input, t' , is generated by combining the estimated components: $\tilde{\mathbf{x}}_{t'} = \alpha(t') \cdot \hat{\mathbf{z}}_t + \gamma(t') \cdot \hat{\mathbf{x}}_t$. The process then iterates to stage (a).

We then introduce two enhancement techniques below to optimize the sampling process:

- (i) **Extrapolating the estimation.** Directly utilizing the estimated $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{z}}_t$ to reconstruct the subsequent input $\tilde{\mathbf{x}}_{t'}$ can result in significant estimation errors, as the estimation model F_θ does not perfectly align with the target function F^{target} for solving the PF-ODE. Note that CFG guides a conditional model using an unconditional model, i.e., $f_\theta(\tilde{\mathbf{x}}, t) = f_\theta(\tilde{\mathbf{x}}, t) + \kappa \cdot (f_\theta^\circ(\tilde{\mathbf{x}}, t) - f_\theta(\tilde{\mathbf{x}}, t))$ where κ is the guidance ratio. This approach can be interpreted as leveraging a less accurate estimation to guide a more accurate one (Karras et al., 2024a). Extending this insight, we propose to extrapolate the next time-step estimates $\hat{\mathbf{x}}_{t'}$ and $\hat{\mathbf{z}}_{t'}$ using the previous estimates $\hat{\mathbf{x}}_t$ and $\hat{\mathbf{z}}_t$, formulated as: $\hat{\mathbf{x}}_{t'} \leftarrow \hat{\mathbf{x}}_{t'} + \kappa \cdot (\hat{\mathbf{x}}_{t'} - \hat{\mathbf{x}}_t)$ and $\hat{\mathbf{z}}_{t'} \leftarrow \hat{\mathbf{z}}_{t'} + \kappa \cdot (\hat{\mathbf{z}}_{t'} - \hat{\mathbf{z}}_t)$, where $\kappa \in [0, 1]$ is the extrapolation ratio. This extrapolation process can significantly enhance sampling quality and reduce the number of sampling steps. Notably, this technique is compatible with CFG and does not introduce additional computational overhead (see Sec. 4.2 for experimental details and App. F.1.10 for theoretical analysis).
- (ii) **Incorporating stochasticity.** During the aforementioned sampling process, the input $\tilde{\mathbf{x}}_t$ is deterministic, potentially limiting the diversity of generated samples. To mitigate this, we introduce a stochastic term ρ to $\tilde{\mathbf{x}}_t$, defined as: $\tilde{\mathbf{x}}_{t'} = \alpha(t') \cdot (\sqrt{1 - \rho} \cdot \hat{\mathbf{z}}_t + \sqrt{\rho} \cdot \mathbf{z}) + \gamma(t') \cdot \hat{\mathbf{x}}_t$, where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is a random noise vector, and ρ is the stochasticity ratio. This stochastic term acts as a random perturbation to $\tilde{\mathbf{x}}_t$, thereby enhancing the diversity of generated samples. We adopt $\rho = \lambda$ as the default configuration, with further analysis provided in App. F.1.11.

Unified sampling algorithm UCGM-S. Putting all these factors together, here we introduce a unified sampling algorithm applicable to consistency models and diffusion/flow-based models, as presented in Alg. 2. An ablation study for our proposed techniques is in Tab. 14. Extensive experiments (cf., Sec. 4) demonstrate two key features of this algorithm:

- (a) *Reduced computational resources:* It decreases the number of sampling steps required by existing models while maintaining or enhancing performance.
- (b) *High compatibility:* It is compatible with existing models, irrespective of their training objectives or noise schedules, without necessitating modifications to model architectures or tuning.

4 EXPERIMENT

This section details the experimental setup and evaluation of our proposed methodology, UCGM-{T,S}. Note that our approach relies on specific parameterizations of the transport coefficients $\alpha(\cdot)$, $\gamma(\cdot)$, $\hat{\alpha}(\cdot)$, and $\hat{\gamma}(\cdot)$, as detailed in Alg. 1 and Alg. 2. Therefore, Tab. 7 summarizes the parameterizations used in experiments, including configurations for compatibility with prior methods.

4.1 EXPERIMENTAL SETTING

Datasets. We utilize ImageNet-1K (Deng et al., 2009) at resolutions of 512×512 and 256×256 as our primary datasets, following prior studies (Karras et al., 2024b; Song et al., 2023) and adhering to ADM’s data preprocessing protocols (Dhariwal & Nichol, 2021). Additionally, CIFAR-10 (Krizhevsky et al., 2009b) at a resolution of 32×32 is employed for ablation studies.

For both 512×512 and 256×256 images, experiments are conducted using latent space generative modeling in line with previous works. Specifically: (a) For 256×256 images, we employ multiple widely-used autoencoders, including SD-VAE (Rombach et al., 2022), VA-VAE (Yao et al., 2025), and E2E-VAE (Leng et al., 2025). (b) For 512×512 images, a DC-AE (*f32c32*) (Chen et al., 2024c) with a higher compression rate is used to conserve computational resources. When utilizing SD-VAE for 512×512 images, a $2 \times$ larger patch size is applied to maintain computational parity with the 256×256 setting. Consequently, the computational burden for generating images at both 512×512

Table 2: **System-level quality comparison for multi-step generation task on class-conditional ImageNet-1K.** Notation $A \oplus B$ denotes the result obtained by combining methods A and B. \downarrow/\uparrow indicate a decrease/increase, respectively, in the metric compared to the baseline performance of the pre-trained models.

METHOD	512 × 512				256 × 256				
	NFE (↓)	FID (↓)	#Params	#Epochs	NFE (↓)	FID (↓)	#Params	#Epochs	
Diffusion & flow-matching Models									
ADM-G (Dhariwal & Nichol, 2021)	250×2	7.72	559M	388	ADM-G (Dhariwal & Nichol, 2021)	250×2	4.59	559M	396
U-ViT-H/4 (Bao et al., 2023)	50×2	4.05	501M	400	U-ViT-H/2 (Bao et al., 2023)	50×2	2.29	501M	400
DiT-XL/2 (Peebles & Xie, 2023)	250×2	3.04	675M	600	DiT-XL/2 (Peebles & Xie, 2023)	250×2	2.27	675M	1400
SiT-XL/2 (Ma et al., 2024)	250×2	2.62	675M	600	SiT-XL/2 (Ma et al., 2024)	250×2	2.06	675M	1400
MaskDiT (Zheng et al., 2023)	79×2	2.50	736M	-	MDT (Gao et al., 2023)	250×2	1.79	675M	1300
EDM2-S (Karras et al., 2024b)	63	2.56	280M	1678	REPA-XL/2 (Yu et al., 2024)	250×2	1.96	675M	200
EDM2-L (Karras et al., 2024b)	63	2.06	778M	1476	REPA-XL/2 (Yu et al., 2024)	250×2	1.42	675M	800
EDM2-XXL (Karras et al., 2024b)	63	1.91	1.5B	734	LightDiT (Yao et al., 2025)	250×2	2.11	675M	64
DiT-XL/1 \oplus Chen et al. (2024c)	250×2	2.41	675M	400	LightDiT (Yao et al., 2025)	250×2	1.35	675M	800
U-ViT-H/1 \oplus Chen et al. (2024c)	30×2	2.53	501M	400	DDT-XL/2 (Wang et al., 2025)	250×2	1.31	675M	256
REPA-XL/2 (Yu et al., 2024)	250×2	2.08	675M	200	DDT-XL/2 (Wang et al., 2025)	250×2	1.26	675M	400
DDT-XL/2 (Wang et al., 2025)	250×2	1.28	675M	-	REPA-E-XL (Leng et al., 2025)	250×2	1.26	675M	800
GANs & masked & autoregressive models									
VQGAN \oplus Esser et al. (2021)	256	18.65	227M	-	VQGAN \oplus Sun et al. (2024)	-	2.18	3.1B	300
MAGViT-v2 (Yu et al., 2023)	64×2	1.91	307M	1080	MAR-L (Li et al., 2024)	256×2	1.78	479M	800
MAR-L (Li et al., 2024)	256×2	1.73	479M	800	MAR-H (Li et al., 2024)	256×2	1.55	943M	800
VAR-d36-s (Tian et al., 2024)	10×2	2.63	2.3B	350	VAR-d30-re (Tian et al., 2024)	10×2	1.73	2.0B	350
Ours: UCGM-S sampling with models trained by prior works									
UCGM-S \oplus Karras et al. (2024b)	40 ¹²³	2.53 ^{10.03}	280M	-	UCGM-S \oplus Wang et al. (2025)	100 ⁴⁰⁰	1.27 ^{70.01}	675M	-
UCGM-S \oplus Karras et al. (2024b)	50 ¹¹³	2.04 ^{10.02}	778M	-	UCGM-S \oplus Yao et al. (2025)	100 ⁴⁰⁰	1.21 ^{10.14}	675M	-
UCGM-S \oplus Karras et al. (2024b)	40 ¹²³	1.88 ^{10.03}	1.5B	-	UCGM-S \oplus Leng et al. (2025)	80 ⁴²⁰	1.06 ^{10.20}	675M	-
UCGM-S \oplus Wang et al. (2025)	200 ¹³⁰⁰	1.25 ^{10.03}	675M	-	UCGM-S \oplus Leng et al. (2025)	20 ⁴⁸⁰	2.00 ^{10.74}	675M	-
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 0$)									
\oplus DC-AE (Chen et al., 2024c)	40	1.48	675M	800	\oplus SD-VAE (Rombach et al., 2022)	60	1.41	675M	400
\oplus DC-AE (Chen et al., 2024c)	20	1.68	675M	800	\oplus VA-VAE (Yao et al., 2025)	60	1.21	675M	400
\oplus SD-VAE (Rombach et al., 2022)	40	1.67	675M	320	\oplus E2E-VAE (Leng et al., 2025)	40	1.21	675M	800
\oplus SD-VAE (Rombach et al., 2022)	20	1.80	675M	320	\oplus E2E-VAE (Leng et al., 2025)	20	1.30	675M	800

and 256×256 resolutions remains comparable across our trained models⁴. Further details on datasets and autoencoders are provided in App. D.1.1.

Neural network architectures. We evaluate UCGM-S sampling using models trained with established methodologies. These models employ various architectures from two prevalent families commonly used in continuous generative models: (a) Diffusion Transformers, including variants such as DiT (Peebles & Xie, 2023), UViT (Bao et al., 2023), SiT (Ma et al., 2024), Lightning-DiT (Yao et al., 2025), and DDT (Wang et al., 2025). (b) UNet-based convolutional networks, including improved UNets (Karras et al., 2022; Song et al., 2020b) and EDM2-UNets (Karras et al., 2024b). For training models specifically for UCGM-T, we consistently utilize DiT as the backbone architecture. We train models of various sizes (B: 130M, L: 458M, XL: 675M parameters) and patch sizes. Notation such as XL/2 denotes the XL model with a patch size of 2. Following prior work (Yao et al., 2025; Wang et al., 2025), minor architectural modifications are applied to enhance training stability (details in App. D.1.2).

Implementation details. Our implementation is developed in PyTorch (Paszke, 2019). Training employs AdamW (Loshchilov & Hutter, 2017) for multi-step sampling models. For few-step sampling models, RAdam (Liu et al., 2019) is used to improve training stability. Consistent with standard practice in generative modeling (Yu et al., 2024; Ma et al., 2024), an exponential moving average (EMA) of model weights is maintained throughout training using a decay rate of 0.9999. All reported results utilize the EMA model. Comprehensive hyperparameters and additional implementation details are provided in App. D.1.3. Consistent with prior work (Song et al., 2020b; Ho et al., 2020; Lipman et al., 2022; Brock et al., 2018), we adopt standard evaluation protocols. The primary metric for assessing image quality is the Fréchet Inception Distance (FID) (Heusel et al., 2017), calculated on 50,000 images (FID-50K).

4.2 COMPARISON WITH SOTA METHODS FOR MULTI-STEP GENERATION

Our experiments on ImageNet-1K at 512×512 and 256×256 resolutions systematically validate the three key advantages of UCGM: (1) sampling acceleration via UCGM-S on pre-trained models, (2) ultra-efficient generation with joint UCGM-T + UCGM-S, and (3) broad compatibility.

UCGM-S: Plug-and-play sampling acceleration without additional cost. UCGM-S provides free sampling acceleration for pre-trained generative models. It reduces the required Number of

⁴Previous works often employed the same autoencoders and patch sizes for both resolutions, resulting in higher computational costs for generating 512×512 images. For example, the DiT-XL/2 model requires 524.60 GFLOPs for 512×512 generation, in contrast to 118.64 GFLOPs for 256×256 .

Table 3: System-level quality comparison for few-step generation task on class-conditional ImageNet-1K.

METHOD	512 × 512				256 × 256				
	NFE (↓)	FID (↓)	#Params	#Epochs	METHOD	NFE (↓)	FID (↓)	#Params	#Epochs
Consistency training & distillation									
sCT-M (Lu & Song, 2024)	1	5.84	498M	1837	iCT (Song & Dhariwal, 2023)	2	20.3	675M	-
	2	5.53	498M	1837	Shortcut-XL/2 (Frans et al., 2024)	1	10.6	676M	250
sCT-L (Lu & Song, 2024)	1	5.15	778M	1274		4	7.80	676M	250
	2	4.65	778M	1274		128	3.80	676M	250
sCT-XXL (Lu & Song, 2024)	1	4.29	1.5B	762	IMM-XL/2 (Zhou et al., 2025)	1×2	7.77	675M	3840
	2	3.76	1.5B	762		2×2	5.33	675M	3840
sCD-M (Lu & Song, 2024)	1	2.75	498M	1997		4×2	3.66	675M	3840
	2	2.26	498M	1997		8×2	2.77	675M	3840
sCD-L (Lu & Song, 2024)	1	2.55	778M	1434	IMM ($\omega = 1.5$)	1×2	8.05	675M	3840
	2	2.04	778M	1434		2×2	3.99	675M	3840
sCD-XXL (Lu & Song, 2024)	1	2.28	1.5B	921		4×2	2.51	675M	3840
	2	1.88	1.5B	921		8×2	1.99	675M	3840
GANs & masked & autoregressive models									
BigGAN (Brock et al., 2018)	1	8.43	160M	-	BigGAN (Brock et al., 2018)	1	6.95	112M	-
StyleGAN (Sauer et al., 2022)	1×2	2.41	168M	-	GigaGAN (Kang et al., 2023)	1	3.45	569M	-
MAGVIT-v2 (Yu et al., 2023)	64×2	1.91	307M	1080	StyleGAN (Sauer et al., 2022)	1×2	2.30	166M	-
VAR-d36-s (Tian et al., 2024)	10×2	2.63	2.3B	350	VAR-d30-re (Tian et al., 2024)	10×2	1.73	2.0B	350
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 0$)									
⊕DC-AE (Chen et al., 2024c)	32	1.55	675M	800	⊕VA-VAE (Yao et al., 2025)	16	2.11	675M	400
⊕DC-AE (Chen et al., 2024c)	16	1.81	675M	800	⊕VA-VAE (Yao et al., 2025)	8	6.09	675M	400
⊕DC-AE (Chen et al., 2024c)	8	3.07	675M	800	⊕E2E-VAE (Leng et al., 2025)	16	1.40	675M	800
⊕DC-AE (Chen et al., 2024c)	4	74.0	675M	800	⊕E2E-VAE (Leng et al., 2025)	8	2.68	675M	800
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 1$)									
⊕DC-AE (Chen et al., 2024c)	1	2.42	675M	840	⊕VA-VAE (Yao et al., 2025)	2	1.42	675M	432
⊕DC-AE (Chen et al., 2024c)	2	1.75	675M	840	⊕VA-VAE (Yao et al., 2025)	1	2.19	675M	432
⊕SD-VAE (Rombach et al., 2022)	1	2.63	675M	360	⊕SD-VAE (Rombach et al., 2022)	1	2.10	675M	424
⊕SD-VAE (Rombach et al., 2022)	2	2.11	675M	360	⊕E2E-VAE (Leng et al., 2025)	1	2.29	675M	264

Function Evaluations (NFEs) while preserving or improving generation quality, as measured by FID. Applied to 512×512 image generation, the approach demonstrates notable efficiency gains:

- (a) For the diffusion-based models, such as a pre-trained EDM2-XXL model, UCGM-S reduced NFEs from 63 to 40 (a 36.5% reduction), concurrently improving FID from 1.91 to 1.88.
- (b) When applied to the flow-based models, such as a pre-trained DDT-XL/2 model, UCGM-S achieved an FID of 1.25 with 200 NFEs, compared to the original 1.28 FID requiring 500 NFEs. This demonstrates a performance improvement achieved alongside enhanced efficiency.

This approach generalizes across different generative model frameworks and resolutions. For instance, on 256×256 resolution using the flow-based REPA-E-XL model, UCGM-S attained 1.06 FID at 80 NFEs, which surpasses the baseline performance of 1.26 FID achieved at 500 NFEs.

In summary, UCGM-S acts as a broadly applicable technique for efficient sampling, demonstrating cases where performance (FID) improves despite a reduction in sampling steps.

UCGM-T + UCGM-S: Synergistic efficiency. The combination of UCGM-T training and UCGM-S sampling yields highly competitive generative performance with minimal NFEs:

- (a) 512×512 : With a DC-AE autoencoder, our framework achieved 1.48 FID at 40 NFEs. This outperforms DiT-XL/1 ⊕ DC-AE (2.41 FID, 500 NFEs) and EDM2-XXL (1.91 FID, 63 NFEs), with comparable or reduced model size.
- (b) 256×256 : With an E2E-VAE autoencoder, we attained 1.21 FID at 40 NFEs. This result exceeds prior SOTA models like MAR-H (1.55 FID, 512 NFEs) and REPA-E-XL (1.26 FID, 500 NFEs).

Importantly, models trained with UCGM-T maintain robustness under extremely low-step sampling regimes. At 20 NFEs, the 256×256 performance degrades gracefully to 1.30 FID, a result that still exceeds the performance of several baseline models sampling with significantly higher NFEs.

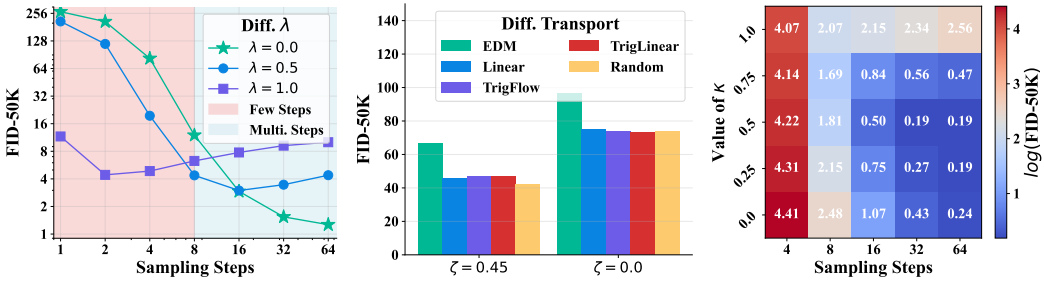
In summary, the demonstrated robustness and efficiency of UCGM-{T, S} across various scenarios underscore the high potential of our UCGM for multi-step continuous generative modeling.

4.3 COMPARISON WITH SOTA METHODS FOR FEW-STEP GENERATION

As evidenced by the results in Tab. 3, our UCGM-{T, S} framework exhibits superior performance across two key settings: $\lambda = 0$, characteristic of a multi-step regime akin to diffusion and flow-matching models, and $\lambda = 1$, indicative of a few-step regime resembling consistency models.

Few-step regime ($\lambda = 1$). Configured for few-step generation, UCGM-{T, S} achieves SOTA sample quality with minimal NFEs, surpassing existing specialized consistency models and GANs:

- (a) 512×512 : Using a DC-AE autoencoder, our model achieves an FID of 1.75 with 2 NFEs and 675M parameters. This outperforms sCD-XXL, a leading consistency distillation model, which reports 1.88 FID with 2 NFEs and 1.5B parameters.



(a) Various λ and sampling steps. (b) Different ζ and transport types. (c) Various κ and sampling steps.

Figure 2: **Ablation studies of UCGM on ImageNet-1K 256×256 .** These studies evaluate key factors of the proposed UCGM. Ablations presented in (a) and (c) utilize XL/1 models with the VA-VAE autoencoder. For the results shown in (b), B/2 models with the SD-VAE autoencoder are used to facilitate more efficient training.

(b) 256×256 : Using a VA-VAE autoencoder, our model achieves an FID of 1.42 with 2 NFEs. This is a notable improvement over IMM-XL/2, which obtains 1.99 FID with $8 \times 2 = 16$ NFEs, demonstrating higher sample quality while requiring $8 \times$ fewer sampling steps.

In summary, *these results demonstrate the capability of UCGM- $\{T, S\}$ to deliver high-quality generation with minimal sampling cost, which is advantageous for practical applications.*

Multi-step regime ($\lambda = 0$). Even when models are trained for multi-step generation, it nonetheless demonstrates competitive performance even when utilizing a moderate number of sampling steps.

(a) 512×512 : Using a DC-AE autoencoder, our model obtains an FID of 1.81 with 16 NFEs and 675M parameters. This result is competitive with or superior to existing methods such as VAR- $d30$ -s, which reports 2.63 FID with $10 \times 2 = 20$ NFEs and 2.3B parameters.

(b) 256×256 : Using an E2E-VAE autoencoder, our model achieves an FID of 1.40 with 16 NFEs. This surpasses IMM-XL/2, which obtains 1.99 FID with $8 \times 2 = 16$ NFEs, demonstrating improved quality at the same sampling cost.

In summary, *our UCGM- $\{T, S\}$ framework demonstrates versatility and high performance across both few-step ($\lambda = 1$) and multi-step ($\lambda = 0$) sampling regimes. As shown, it consistently achieves SOTA or competitive sample quality relative to existing methods, often requiring fewer sampling steps or parameters, which are important factors for efficient high-resolution image synthesis.*

4.4 ABLATION STUDY OVER THE KEY FACTORS OF UCGM

Unless otherwise specified, experiments in this section are conducted with $\kappa = 0.0$ and $\lambda = 0.0$.

Effect of λ in UCGM-T. Fig. 2a demonstrates that varying λ influences the range of effective sampling steps for trained models. For instance, with $\lambda = 1$ ⁵, optimal performance is attained at 2 sampling steps. In contrast, with $\lambda = 0.5$, optimal performance is observed at 16 steps.

Impact of ζ and transport type in UCGM. The results in Fig. 2b demonstrates that UCGM- $\{T, S\}$ is applicable with various transport types, albeit with some performance variation. Investigating these performance differences constitutes future work. The results also illustrate that the enhanced training objective (achieved with $\zeta = 0.45$ compared to $\zeta = 0.0$, per Sec. 3) consistently improves performance across all tested transport types, underscoring the efficacy of this technique.

Setting different κ in UCGM-S. Experimental results, depicted in Fig. 2c, illustrate the impact of κ on the trade-off between sampling steps and generation quality: (a) High κ values (e.g., 1.0 and 0.75) prove beneficial for extreme few-step sampling scenarios (e.g., 4 steps); (b) Moreover, mid-range κ values (0.25 to 0.5) achieve superior performance with fewer steps compared to $\kappa = 0.0$.

5 ADDITIONAL EXPERIMENTS ON LARGE-SCALE MODELS AND DATASETS

5.1 COMPARISON WITH TEXT-TO-IMAGE MODELS

We evaluate the practical efficacy of UCGM on text-to-image synthesis, with comprehensive benchmarks detailed in Tab. 4. The training efficiency is notable: fine-tuning the SANA-0.6B and SANA-1.6B backbones (batch sizes 128 and 64, respectively) for 40,000 steps required only 60 NVIDIA H800 GPU hours.

⁵For the purpose of a fair ablation study, additional stabilizing techniques were omitted for this $\lambda = 1$ case.

Table 4: **System-level comparison of UCGM against few-step text-to-image baselines.** Throughput (batch size 10) and latency (batch size 1) are evaluated on a single NVIDIA A100 GPU (BF16).

Method	NFE ↓	Throughput ↑ (samples/s)	Latency (s) ↓	#Params	GenEval ↑	DPG-Bench ↑
SDXL-DMD2 (Yin et al., 2024a)	2	2.89	0.40	0.9B	0.58	-
FLUX-Schnell (Labs, 2024)	2	0.92	1.15	12.0B	0.71	-
SANA-Sprint-0.6B (Chen et al., 2025c)	2	6.46	0.25	0.6B	0.76	81.5
SANA-Sprint-1.6B (Chen et al., 2025c)	2	5.68	0.24	1.6B	0.77	82.1
SDXL-LCM (Luo et al., 2023)	2	2.89	0.40	0.9B	0.44	-
PixArt-LCM (Chen et al., 2024b)	2	3.52	0.31	0.6B	0.42	-
PCM (Wang et al., 2024)	2	2.62	0.56	0.9B	0.55	-
SD3.5-Turbo (Esser et al., 2024)	2	1.61	0.68	8.0B	0.53	-
PixArt-DMD (Chen et al., 2024a)	1	4.26	0.25	0.6B	0.45	-
SDXL-DMD2 (Yin et al., 2024a)	1	3.36	0.32	0.9B	0.59	-
FLUX-Schnell (Labs, 2024)	1	1.58	0.68	12.0B	0.69	-
SANA-Sprint-0.6B (Chen et al., 2025c)	1	7.22	0.21	0.6B	0.72	78.6
SANA-Sprint-1.6B (Chen et al., 2025c)	1	6.71	0.21	1.6B	0.76	80.1
SDXL-LCM (Luo et al., 2023)	1	3.36	0.32	0.9B	0.28	-
PixArt-LCM (Chen et al., 2024b)	1	4.26	0.25	0.6B	0.41	-
PCM (Wang et al., 2024)	1	3.16	0.40	0.9B	0.42	-
SD3.5-Turbo (Esser et al., 2024)	1	2.48	0.45	8.0B	0.51	-
UCGM-0.6B (Ours, $\lambda = 1$)	2	6.50	0.26	0.6B	0.84	81.0
UCGM-1.6B (Ours, $\lambda = 1$)	2	5.71	0.25	1.6B	0.82	82.4
UCGM-0.6B (Ours, $\lambda = 1$)	1	7.30	0.23	0.6B	0.79	78.2
UCGM-1.6B (Ours, $\lambda = 1$)	1	6.75	0.22	1.6B	0.80	80.7

The tuning settings follow those outlined in App. D.3.

Empirical results establish a new Pareto frontier for generation quality and speed. At 2 NFE, UCGM sets a high performance standard, with our 0.6B model achieving a GenEval score of **0.84**. This significantly outperforms the 12B-parameter FLUX-Schnell (0.71) and SANA-Sprint-1.6B (0.77), demonstrating that massive parameter counts are not a prerequisite for high fidelity. This advantage persists in the challenging single-step (1 NFE) regime, where UCGM-0.6B attains a score of **0.79**—surpassing SANA-Sprint-1.6B (0.76)—while delivering the highest throughput in the benchmark at 7.30 samples/s.

Beyond quantitative metrics, UCGM proves that objective formulation outweighs system complexity. Unlike baselines like SANA-Sprint that rely on composite adversarial losses, we achieve superior fidelity using solely the singular objective in (6).

5.2 EXTENSION TO UNIFIED MULTIMODAL MODELS

We push the scalability limits of UCGM by applying it to the **20B-parameter Multi-Modal Diffusion Transformer (MM-DiT)**. The experimental outcomes, detailed in App. D.3, highlight our engineering success and underscore two distinct advantages:

- Successful training of high-capacity models:** We successfully scaled UCGM to the 20B regime using **12,976 H800 GPU hours**. Relying exclusively on public datasets (e.g., LAION-5B), our model achieves performance parity with state-of-the-art backbones on **GenEval, DPG-Bench and WISE** benchmarks, demonstrating that it handles large-scale UMMs.
- Robust distillation in few-step regimes:** Distilling 20B-parameter models presents severe stability hurdles. Our benchmarks reveal that Consistency Models (Song et al., 2023) suffer from **catastrophic collapse**, and MeanFlow (Geng et al., 2025) encounters **Out-of-Memory (OOM)** errors. In contrast, UCGM successfully distills the 20B UMM, maintaining superior stability and quality where prior arts fail.

6 CONCLUSION

We present UCGM, a unified and efficient framework for training and sampling both few-step and multi-step continuous generative models. Within this framework, we derive a surrogate unified objective that theoretically decomposes the few-step objective into the multi-step objective plus a self-alignment term. Building on this foundation, we introduce UCGM-T, which seamlessly bridges few-step (e.g., consistency models) and multi-step (e.g., diffusion, flow matching) generative paradigms, supporting diverse model architectures, latent autoencoders, and noise schedules. We further propose UCGM-S, a unified sampler compatible with our trained models, which can also accelerate and enhance pre-trained models from existing paradigms.

540 7 ETHICS STATEMENT

541 This research adheres to the *ICLR Code of Ethics* and is committed to the principles of responsible and
 542 transparent scientific inquiry. The study involves no human participants, personal or sensitive data,
 543 or any activities requiring approval from an institutional ethics review board. All datasets used are
 544 publicly accessible under appropriate licenses, with proper attribution given to their original sources.
 545 To promote openness and reproducibility, we provide our implementation code and experimental
 546 settings for verification and further development by the research community. We also declare that no
 547 conflicts of interest or external funding have influenced the design, execution, or presentation of this
 548 work.

549 8 REPRODUCIBILITY STATEMENT

550 Comprehensive details regarding the datasets, model architectures, optimization settings, and training
 551 procedures are provided in [Sec. 4.1](#) of the main paper and in [App. D](#). These materials are designed
 552 to facilitate the reliable and transparent reproduction of our results. Additionally, our source code
 553 will be made publicly available upon acceptance of the paper and **is included in the supplementary**
 554 **material**.

555 REFERENCES

- 556
 557
 558 Fan Bao, Shen Nie, Kaiwen Xue, Yue Cao, Chongxuan Li, Hang Su, and Jun Zhu. All are worth
 559 words: A vit backbone for diffusion models. In *Proceedings of the IEEE/CVF conference on*
 560 *computer vision and pattern recognition*, pp. 22669–22679, 2023.
- 561 Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale gan training for high fidelity natural
 562 image synthesis. *arXiv preprint arXiv:1809.11096*, 2018.
- 563
 564 Soravit Changpinyo, Piyush Sharma, Nan Ding, and Radu Soricut. Conceptual 12M: Pushing
 565 web-scale image-text pre-training to recognize long-tail visual concepts. In *CVPR*, 2021.
- 566
 567 Jiu hai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi Zhou,
 568 Saining Xie, Silvio Savarese, Le Xue, Caiming Xiong, and Ran Xu. Blip3-o: A Family of Fully
 569 Open Unified Multimodal Models-Architecture, Training and Dataset. *arXiv.org*, abs/2505.09568,
 570 2025a.
- 571 Jiu hai Chen, Zhiyang Xu, Xichen Pan, Yushi Hu, Can Qin, Tom Goldstein, Lifu Huang, Tianyi
 572 Zhou, Saining Xie, Silvio Savarese, et al. Blip3-o: A family of fully open unified multimodal
 573 models-architecture, training and dataset. *arXiv preprint arXiv:2505.09568*, 2025b.
- 574
 575 Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping
 576 Luo, Huchuan Lu, and Zhenguo Li. Pixart- σ : Weak-to-strong training of diffusion transformer for
 577 4k text-to-image generation. *arXiv preprint arXiv:2403.04692*, 2024a.
- 578
 579 Junsong Chen, Yue Wu, Simian Luo, Enze Xie, Sayak Paul, Ping Luo, Hang Zhao, and Zhenguo Li.
 580 Pixart- $\{\delta\}$: Fast and controllable image generation with latent consistency models. *arXiv*
 581 *preprint arXiv:2401.05252*, 2024b.
- 582
 583 Junsong Chen, Shuchen Xue, Yuyang Zhao, Jincheng Yu, Sayak Paul, Junyu Chen, Han Cai, Enze
 584 Xie, and Song Han. Sana-sprint: One-step diffusion with continuous-time consistency distillation.
 585 *arXiv preprint arXiv:2503.09641*, 2025c.
- 586
 587 Junying Chen, Zhenyang Cai, Pengcheng Chen, Shunian Chen, Ke Ji, Xidong Wang, Yunjin Yang,
 588 and Benyou Wang. Sharegpt-4o-Image: Aligning Multimodal Models with GPT-4o-Level Image
 589 Generation. *arXiv.org*, abs/2506.18095, 2025d.
- 590
 591 Junyu Chen, Han Cai, Junsong Chen, Enze Xie, Shang Yang, Haotian Tang, Muyang Li, Yao Lu, and
 592 Song Han. Deep compression autoencoder for efficient high-resolution diffusion models. *arXiv*
 593 *preprint arXiv:2410.10733*, 2024c.
- Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and
 Chong Ruan. Janus-pro: Unified multimodal understanding and generation with data and model
 scaling. *arXiv preprint arXiv:2501.17811*, 2025e.

- 594 Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao
595 Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv*
596 *preprint arXiv:2505.14683*, 2025.
- 597 Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
598 hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*,
599 pp. 248–255. Ieee, 2009.
- 600 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*
601 *in neural information processing systems*, 34:8780–8794, 2021.
- 602 Patrick Esser, Robin Rombach, and Bjorn Ommer. Taming transformers for high-resolution image
603 synthesis. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,
604 pp. 12873–12883, 2021.
- 605 Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam
606 Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for
607 high-resolution image synthesis. In *Forty-first international conference on machine learning*, 2024.
- 608 Kevin Frans, Danijar Hafner, Sergey Levine, and Pieter Abbeel. One step diffusion via shortcut
609 models. *arXiv preprint arXiv:2410.12557*, 2024.
- 610 Shanghua Gao, Pan Zhou, Ming-Ming Cheng, and Shuicheng Yan. Masked diffusion transformer is a
611 strong image synthesizer. In *Proceedings of the IEEE/CVF international conference on computer*
612 *vision*, pp. 23164–23173, 2023.
- 613 Zhengyang Geng, Ashwini Pokle, William Luo, Justin Lin, and J Zico Kolter. Consistency models
614 made easy. *arXiv preprint arXiv:2406.14548*, 2024.
- 615 Zhengyang Geng, Mingyang Deng, Xingjian Bai, J Zico Kolter, and Kaiming He. Mean flows for
616 one-step generative modeling. *arXiv preprint arXiv:2505.13447*, 2025.
- 617 Dhruba Ghosh, Hannaneh Hajishirzi, and Ludwig Schmidt. Geneval: An object-focused framework
618 for evaluating text-to-image alignment. In *Conference on Neural Information Processing Systems*
619 *(NeurIPS)*, 2023.
- 620 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
621 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the*
622 *ACM*, 63(11):139–144, 2020.
- 623 Jacky He and contributors. text-to-image-2M: A high-quality, diverse text–image training dataset.
624 <https://huggingface.co/datasets/jackyhate/text-to-image-2M>, 2024.
- 625 Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans
626 trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural*
627 *information processing systems*, 30, 2017.
- 628 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*,
629 2022.
- 630 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*
631 *neural information processing systems*, 33:6840–6851, 2020.
- 632 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J
633 Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646,
634 2022.
- 635 Peter Holderrieth and Ezra Erives. An introduction to flow matching and diffusion models. *arXiv*
636 *preprint arXiv:2506.02070*, 2025.
- 637 Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip Diffusion Models
638 with LLM for Enhanced Semantic Alignment. *arXiv.org*, abs/2403.05135, 2024.
- 639 Juan C Jiménez. Simplified formulas for the mean and variance of linear stochastic differential
640 equations. *Applied Mathematics Letters*, 49:12–19, 2015.
- 641

- 648 Minguk Kang, Jun-Yan Zhu, Richard Zhang, Jaesik Park, Eli Shechtman, Sylvain Paris, and Taesung
649 Park. Scaling up gans for text-to-image synthesis. In *Proceedings of the IEEE/CVF conference on*
650 *computer vision and pattern recognition*, pp. 10124–10134, 2023.
- 651
- 652 Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-
653 based generative models. *Advances in neural information processing systems*, 35:26565–26577,
654 2022.
- 655 Tero Karras, Miika Aittala, Tuomas Kynkäänniemi, Jaakko Lehtinen, Timo Aila, and Samuli Laine.
656 Guiding a diffusion model with a bad version of itself. *Advances in Neural Information Processing*
657 *Systems*, 37:52996–53021, 2024a.
- 658
- 659 Tero Karras, Miika Aittala, Jaakko Lehtinen, Janne Hellsten, Timo Aila, and Samuli Laine. Analyzing
660 and improving the training dynamics of diffusion models. In *Proceedings of the IEEE/CVF*
661 *Conference on Computer Vision and Pattern Recognition*, pp. 24174–24184, 2024b.
- 662 Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
663 2009a.
- 664
- 665 Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 and cifar-100 datasets. *URL: https://www.*
666 *cs.toronto.edu/kriz/cifar.html*, 6(1):1, 2009b.
- 667
- 668 Black Forest Labs. Flux. <https://github.com/black-forest-labs/flux>, 2024.
- 669
- 670 LAION. Releasing re-laion-5b: transparent iteration on laion-5b with additional safety fixes. <https://laion.ai/blog/re-laion-5b/>, 2024. Accessed: 30 aug, 2024.
- 671
- 672 Xingjian Leng, Jaskirat Singh, Yunzhong Hou, Zhenchang Xing, Saining Xie, and Liang Zheng.
673 Repa-e: Unlocking vae for end-to-end tuning with latent diffusion transformers. *arXiv preprint*
674 *arXiv:2504.10483*, 2025.
- 675
- 676 Tianhong Li, Yonglong Tian, He Li, Mingyang Deng, and Kaiming He. Autoregressive image
677 generation without vector quantization. *Advances in Neural Information Processing Systems*, 37:
678 56424–56445, 2024.
- 679 Bin Lin, Zongjian Li, Xinhua Cheng, Yuwei Niu, Yang Ye, Xianyi He, Shenghai Yuan, Wangbo Yu,
680 Shaodong Wang, Yunyang Ge, et al. Uniworld: High-resolution semantic encoders for unified
681 visual understanding and generation. *arXiv preprint arXiv:2506.03147*, 2025.
- 682
- 683 Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching
684 for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- 685
- 686 Liyuan Liu, Haoming Jiang, Pengcheng He, Weizhu Chen, Xiaodong Liu, Jianfeng Gao, and Jiawei
687 Han. On the variance of the adaptive learning rate and beyond. *arXiv preprint arXiv:1908.03265*,
688 2019.
- 689
- 690 Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and
691 transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.
- 692
- 693 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint*
694 *arXiv:1711.05101*, 2017.
- 695
- 696 Cheng Lu and Yang Song. Simplifying, stabilizing and scaling continuous-time consistency models.
697 *arXiv preprint arXiv:2410.11081*, 2024.
- 698
- 699 Simian Luo, Yiqin Tan, Longbo Huang, Jian Li, and Hang Zhao. Latent consistency models:
700 Synthesizing high-resolution images with few-step inference. *arXiv preprint arXiv:2310.04378*,
701 2023.
- 702
- 703 Nanye Ma, Mark Goldstein, Michael S Albergo, Nicholas M Boffi, Eric Vanden-Eijnden, and
704 Saining Xie. Sit: Exploring flow and diffusion-based generative models with scalable interpolant
705 transformers. In *European Conference on Computer Vision*, pp. 23–40. Springer, 2024.

- 702 Ollin Matsubara and Draw Things AI Team. Megalith-10M: A dataset of 10 million
703 public-domain photographs. [https://huggingface.co/datasets/madebyollin/
704 megalith-10m](https://huggingface.co/datasets/madebyollin/megalith-10m), 2024. CC0/Flickr-Commons images; Florence-2 captions available in the
705 *megalith-10m-florence2* variant.
- 706 Yuwei Niu, Munan Ning, Mengren Zheng, Bin Lin, Peng Jin, Jiaqi Liao, Kun-Peng Ning, Bin
707 Zhu, and Li Yuan. Wise: A World Knowledge-Informed Semantic Evaluation for Text-to-Image
708 Generation. *arXiv.org*, abs/2503.07265, 2025.
- 709 Xichen Pan, Satya Narayan Shukla, Aashu Singh, Zhuokai Zhao, Shlok Kumar Mishra, Jialiang
710 Wang, Zhiyang Xu, Jiuhai Chen, Kunpeng Li, Felix Juefei-Xu, Ji Hou, and Saining Xie. Transfer
711 between Modalities with MetaQueries. *arXiv.org*, abs/2504.06256, 2025.
- 712 A Paszke. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint
713 arXiv:1912.01703*, 2019.
- 714 Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier
715 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, et al. Scikit-learn:
716 Machine learning in python. *the Journal of machine Learning research*, 12:2825–2830, 2011.
- 717 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *Proceedings of
718 the IEEE/CVF international conference on computer vision*, pp. 4195–4205, 2023.
- 719 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
720 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-
721 ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 722 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv
723 preprint arXiv:2202.00512*, 2022.
- 724 Axel Sauer, Katja Schwarz, and Andreas Geiger. Stylegan-xl: Scaling stylegan to large diverse
725 datasets. In *ACM SIGGRAPH 2022 conference proceedings*, pp. 1–10, 2022.
- 726 Noam Shazeer. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*, 2020.
- 727 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv
728 preprint arXiv:2010.02502*, 2020a.
- 729 Yang Song and Prafulla Dhariwal. Improved techniques for training consistency models. *arXiv
730 preprint arXiv:2310.14189*, 2023.
- 731 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben
732 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint
733 arXiv:2011.13456*, 2020b.
- 734 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint
735 arXiv:2303.01469*, 2023.
- 736 Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. Roformer: Enhanced
737 transformer with rotary position embedding. *Neurocomputing*, 568:127063, 2024.
- 738 Peize Sun, Yi Jiang, Shoufa Chen, Shilong Zhang, Bingyue Peng, Ping Luo, and Zehuan Yuan.
739 Autoregressive model beats diffusion: Llama for scalable image generation. *arXiv preprint
740 arXiv:2406.06525*, 2024.
- 741 Zhicong Tang, Jianmin Bao, Dong Chen, and Baining Guo. Diffusion models without classifier-free
742 guidance. *arXiv preprint arXiv:2502.12154*, 2025.
- 743 Keyu Tian, Yi Jiang, Zehuan Yuan, Bingyue Peng, and Liwei Wang. Visual autoregressive modeling:
744 Scalable image generation via next-scale prediction. *Advances in neural information processing
745 systems*, 37:84839–84865, 2024.
- 746 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
747 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing
748 systems*, 30, 2017.

- 756 Fu-Yun Wang, Zhaoyang Huang, Alexander William Bergman, Dazhong Shen, Peng Gao, Michael
757 Lingelbach, Keqiang Sun, Weikang Bian, Guanglu Song, Yu Liu, et al. Phased consistency model.
758 *arXiv preprint arXiv:2405.18407*, 2024.
- 759
760 Shuai Wang, Zhi Tian, Weilin Huang, and Limin Wang. Ddt: Decoupled diffusion transformer. *arXiv*
761 *preprint arXiv:2504.05741*, 2025.
- 762
763 Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan Gao, Kun Yan, Sheng ming Yin, Shuai
764 Bai, Xiao Xu, Yilei Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi Wang, An Yang,
765 Bowen Yu, Chen Cheng, Dayiheng Liu, Deqing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan
766 Ni, Kai Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng Wang, Shuting Yu, Tingkun
767 Wen, Wensen Feng, Xiaoxiao Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu, Yuxuan
768 Cai, and Zenan Liu. Qwen-image technical report, 2025a. URL [https://arxiv.org/abs/
2508.02324](https://arxiv.org/abs/2508.02324).
- 769
770 Chenyuan Wu, Pengfei Zheng, Ruiran Yan, Shitao Xiao, Xin Luo, Yueze Wang, Wanli Li, Xiyan
771 Jiang, Yexin Liu, Junjie Zhou, et al. Omnigen2: Exploration to advanced multimodal generation.
772 *arXiv preprint arXiv:2506.18871*, 2025b.
- 773
774 Size Wu, Zhonghua Wu, Zerui Gong, Qingyi Tao, Sheng Jin, Qinyue Li, Wei Li, and Chen Change
775 Loy. Openuni: A simple baseline for unified multimodal understanding and generation. 2025c.
776 URL <https://arxiv.org/abs/2505.23661>.
- 777
778 Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Shuting
779 Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation. *arXiv preprint*
780 *arXiv:2409.11340*, 2024.
- 781
782 Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li,
783 Ligeng Zhu, Yao Lu, et al. Sana: Efficient high-resolution image synthesis with linear diffusion
784 transformers. *arXiv preprint arXiv:2410.10629*, 2024a.
- 785
786 Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin,
787 Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer
788 to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024b.
- 789
790 Jinheng Xie, Zhenheng Yang, and Mike Zheng Shou. Show-o2: Improved native unified multimodal
791 models. *arXiv preprint arXiv:2506.15564*, 2025.
- 792
793 Jingfeng Yao, Bin Yang, and Xinggang Wang. Reconstruction vs. generation: Taming optimization
794 dilemma in latent diffusion models. *arXiv preprint arXiv:2501.01423*, 2025.
- 795
796 Junyan Ye, Dongzhi Jiang, Zihao Wang, Leqi Zhu, Zhenghao Hu, Zilong Huang, Jun He, Zhiyuan
797 Yan, Jinghua Yu, Hongsheng Li, Conghui He, and Weijia Li. Echo-4o: Harnessing the Power of
798 GPT-4o Synthetic Images for Improved Image Generation. *arXiv*, 2025.
- 799
800 Tianwei Yin, Michaël Gharbi, Taesung Park, Richard Zhang, Eli Shechtman, Fredo Durand, and
801 William T Freeman. Improved distribution matching distillation for fast image synthesis. *arXiv*
802 *preprint arXiv:2405.14867*, 2024a.
- 803
804 Tianwei Yin, Michaël Gharbi, Richard Zhang, Eli Shechtman, Fredo Durand, William T Freeman,
805 and Taesung Park. One-step diffusion with distribution matching distillation. In *Proceedings of*
806 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 6613–6623, 2024b.
- 807
808 Lijun Yu, José Lezama, Nitesh B Gundavarapu, Luca Versari, Kihyuk Sohn, David Minnen, Yong
809 Cheng, Vighnesh Birodkar, Agrim Gupta, Xiuye Gu, et al. Language model beats diffusion-
tokenizer is key to visual generation. *arXiv preprint arXiv:2310.05737*, 2023.
- 810
811 Sihyun Yu, Sangkyung Kwak, Huiwon Jang, Jongheon Jeong, Jonathan Huang, Jinwoo Shin, and
812 Saining Xie. Representation alignment for generation: Training diffusion transformers is easier
813 than you think. *arXiv preprint arXiv:2410.06940*, 2024.
- 814
815 Biao Zhang and Rico Sennrich. Root mean square layer normalization. *Advances in Neural*
816 *Information Processing Systems*, 32, 2019.

810 Hongkai Zheng, Weili Nie, Arash Vahdat, and Anima Anandkumar. Fast training of diffusion models
811 with masked transformers. *arXiv preprint arXiv:2306.09305*, 2023.
812
813 Linqi Zhou, Stefano Ermon, and Jiaming Song. Inductive moment matching. *arXiv preprint*
814 *arXiv:2503.07565*, 2025.
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
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

864	CONTENTS	
865		
866	1 Introduction	1
867		
868	2 Preliminaries	2
869		
870	3 Methodology	3
871	3.1 Unified Framework for Continuous Generative Models	3
872	3.2 Instantiating the Unified Framework for Training (UCGM-T)	5
873	3.3 Instantiating the Unified Framework for Sampling (UCGM-S)	5
874		
875	4 Experiment	6
876	4.1 Experimental Setting	6
877	4.2 Comparison with SOTA Methods for Multi-step Generation	7
878	4.3 Comparison with SOTA Methods for Few-step Generation	8
879	4.4 Ablation Study over the Key Factors of UCGM	9
880		
881	5 Additional Experiments on Large-scale Models and Datasets	9
882	5.1 Comparison with Text-to-image Models	9
883	5.2 Extension to Unified Multimodal Models	10
884		
885	6 Conclusion	10
886		
887	7 Ethics Statement	11
888		
889	8 Reproducibility Statement	11
890		
891	A Use of LLMs	19
892		
893	B Broader Impacts	19
894		
895	C Limitations	19
896		
897	D Detailed Experiment	19
898	D.1 Detailed Experimental Setting	19
899	D.1.1 Detailed Datasets	19
900	D.1.2 Detailed Neural Architecture	19
901	D.1.3 Detailed Implementation Details	20
902	D.2 Experimental Results on Small Datasets	22
903	D.3 Experimental Results on Large-scale Unified Multimodal Models	22
904	D.4 Detailed Comparison with SOTA Methods for Multi-step Generation	25
905	D.5 Detailed Comparison with SOTA Methods for Few-step Generation	25
906	D.6 Case Studies	26
907	D.6.1 Analysis of Consistency Ratio λ	26
908	D.6.2 Analysis of Transport Types	27
909	D.7 Analysis of Pre-trained Model Tuning	28
910	D.8 Ablation study on UCGM techniques	29
911		
912	E Pseudocode	30
913	E.1 Training Algorithm for UCGM-T	30
914	E.2 Sampling Algorithm for UCGM-S	31
915		
916	F Theoretical Analysis	32
917	F.1 Main Results	32
	F.1.1 Learning Objective when $\lambda = 0$	32
	F.1.2 Closed-form Solution Analysis when $\lambda = 0$	33
	F.1.3 Learning Objective as $\lambda \rightarrow 1$	45
	F.1.4 Analysis on the Optimal Solution for $\lambda \in [0, 1]$	49
	F.1.5 Surrogate Objective for Unified Objective	50
	F.1.6 Closed-form Solution Analysis for $\lambda \in [0, 1]$	53

918		
919	F.1.7	Unified Training Objective 56
920	F.1.8	Enhanced Target Score Function 57
921	F.1.9	Unified Sampling Process 58
922	F.1.10	Extrapolating Estimation 59
923	F.1.11	Interpretation of Unified Sampling Process 60
924	F.2	Other Techniques 62
925	F.2.1	Beta Transformation 62
926	F.2.2	Kumaraswamy Transformation 62
927	F.2.3	Derivative Estimation 67
928	F.2.4	Calculation of Transport 69
929		
930		
931		
932		
933		
934		
935		
936		
937		
938		
939		
940		
941		
942		
943		
944		
945		
946		
947		
948		
949		
950		
951		
952		
953		
954		
955		
956		
957		
958		
959		
960		
961		
962		
963		
964		
965		
966		
967		
968		
969		
970		
971		

972 A USE OF LLMs

973 During the preparation of this paper, we used OpenAI’s ChatGPT to assist with language refinement,
974 including grammar correction, style polishing, and improving readability. The model was not used
975 for generating ideas, experimental design, analysis, or writing substantive technical content. All
976 scientific contributions, including theoretical derivations, method development, and experimental
977 results, are entirely the work of the authors.

979 B BROADER IMPACTS

980 This paper proposes a unified implementation and theoretical framework for recent popular continuous
981 generative models, such as diffusion models, flow matching models, and consistency models. This
982 work should provide positive impacts for the generative modeling community.

984 C LIMITATIONS

985 **Integration of training acceleration techniques.** This work does not explore the integration of
986 advanced training acceleration methods for diffusion models, such as REPA (Yu et al., 2024).

987 **Exploration of downstream applications.** The current study focuses on establishing the founda-
988 tional framework. Comprehensive exploration of its application to complex downstream generative
989 tasks, including text-to-image and text-to-video generation, is reserved for future research.

990 **Stabilization of few-step objectives.** While we theoretically decompose the few-step objective
991 into the multi-step objective and a self-alignment term, and identify the self-alignment term as the
992 source of potential instability, methods for stabilizing the few-step objective are not investigated in
993 this work. We leave this as an important direction for future research.

994 D DETAILED EXPERIMENT

995 D.1 DETAILED EXPERIMENTAL SETTING

996 D.1.1 DETAILED DATASETS

998 **Image datasets.** We conduct experiments on two datasets: CIFAR-10 (Krizhevsky et al., 2009a),
999 ImageNet-1K (Deng et al., 2009):

- 1000 (a) CIFAR-10 is a widely used benchmark dataset for image classification and generation tasks. It
1001 consists of 60,000 color images, each with a resolution of 32×32 pixels, categorized into 10
1002 distinct classes. The dataset is divided into 50,000 training images and 10,000 test images.
1003 (b) ImageNet-1K is a large-scale dataset containing over 1.2 million high-resolution images across
1004 1,000 categories.

1005 **Latent space datasets.** However, directly training diffusion transformers in the pixel space is
1006 computationally expensive and inefficient. Therefore, following previous studies (Yu et al., 2024; Ma
1007 et al., 2024), we train our diffusion transformers in latent space instead. Tab. 5 presents a comparative
1008 analysis of various Variational Autoencoder (VAE) architectures. SD-VAE is characterized by a
1009 higher spatial resolution in its latent representation (e.g., $H/8 \times W/8$) combined with a lower channel
1010 capacity (4 channels). Conversely, alternative models such as VA-VAE, E2E-VAE, and DC-AE
1011 achieve more significant spatial compression (e.g., $H/16 \times W/16$ or $H/32 \times W/32$) at the expense
1012 of an increased channel depth (typically 32 channels).

1013 A key consideration is that the computational cost of a diffusion transformer subsequently processing
1014 these latent representations is primarily dictated by their spatial dimensions, rather than their channel
1015 capacity (Chen et al., 2024c). Specifically, if the latent map is processed by a transformer by dividing
1016 it into non-overlapping patches, the cost is proportional to the number of these patches. This quantity
1017 is given by $(H/\text{Compression Ratio}/\text{Patch Size}) \times (W/\text{Compression Ratio}/\text{Patch Size})$. Here, H and
1018 W are the input image dimensions, Compression Ratio refers to the spatial compression factor of
1019 the VAE (e.g., 8, 16, 32 as detailed in Tab. 5), and Patch Size denotes the side length of the patches
1020 processed by the transformer.

1021 D.1.2 DETAILED NEURAL ARCHITECTURE

1022 Diffusion Transformers (DiTs) represent a paradigm shift in generative modeling by replacing the
1023 traditional U-Net backbone with a Transformer-based architecture. Proposed by *Scalable Diffusion*
1024 *Models with Transformers* (Peebles & Xie, 2023), DiTs exhibit superior scalability and performance
1025 in image generation tasks. In this paper, we utilize three key variants—DiT-B (130M parameters),
DiT-L (458M parameters), and DiT-XL (675M parameters).

Table 5: **Comparison of different VAE architectures in terms of latent space dimensions and channel capacity.** The table contrasts four variational autoencoder variants (SD-VAE, VA-VAE, E2E-VAE, and DC-AE) by their spatial compression ratios (latent size) and feature channel dimensions. Here, H and W denote input image height and width (e.g., 256×256 or 512×512), respectively.

	SD-VAE (both <code>ema</code> and <code>mse</code> versions) (Rombach et al., 2022)	VA-VAE (Yao et al., 2025)	E2E-VAE (Leng et al., 2025)	DC-AE (<code>f32c32</code>) (Chen et al., 2024c)
Latent Size	$(H/8) \times (W/8)$	$(H/16) \times (W/16)$	$(H/16) \times (W/16)$	$(H/32) \times (W/32)$
Channels	4	32	32	32

To improve training stability, informed by recent studies (Yao et al., 2025; Wang et al., 2025), we incorporate several architectural modifications into the DiT model: (a) SwiGLU feed-forward networks (FFN) (Shazeer, 2020); (b) RMSNorm (Zhang & Sennrich, 2019) without learnable affine parameters; (c) Rotary Positional Embeddings (RoPE) (Su et al., 2024); and (d) parameter-free RMSNorm applied to Key (K) and Query (Q) projections in self-attention layers (Vaswani et al., 2017).

D.1.3 DETAILED IMPLEMENTATION DETAILS

Experiments were conducted on a cluster equipped with 8 H800 GPUs, each with 80 GB of VRAM.

Hyperparameter configuration. Detailed hyperparameter configurations are provided in Tab. 6 to ensure reproducibility. The design of time schedules for sampling processes varies in complexity. For few-step models, typically employing 1 or 2 sampling steps, manual schedule design is straightforward. However, the time schedule \mathcal{T} utilized by our UCGM-S often comprises a large number of time points, particularly for a large number of sampling steps N . Manual design of such dense schedules is challenging and can limit the achievable performance of our UCGM- $\{T, S\}$, as prior work (Yao et al., 2025; Wang et al., 2025) has established that carefully designed schedules significantly enhance multi-step models, including flow-matching variants. To address this, we propose transforming each time point $t \in \mathcal{T}$ using a generalized Kumaraswamy transformation: $f_{\text{Kuma}}(t; a, b, c) = (1 - (1 - t^a)^b)^c$. This choice is motivated by the common practice in prior studies of applying non-linear transformations to individual time points to construct effective schedules. A specific instance of such a transformation is the `timeshift` function $f_{\text{shift}}(t; s) = \frac{st}{1+(s-1)t}$, where $s > 0$ (Yao et al., 2025). We find that the Kumaraswamy transformation, by appropriate selection of parameters a, b, c , can effectively approximate f_{shift} and other widely-used functions (cf., App. F.2.2), including the identity function $f(t) = t$ (Yu et al., 2024; Leng et al., 2025). Empirical evaluations suggest that the parameter configuration $(a, b, c) = (1.17, 0.8, 1.1)$ yields robust performance across diverse scenarios, corresponding to the "Auto" setting in Tab. 6.

Detailed implementation techniques of enhancing target score function. We enhance the target score function for conditional diffusion models by modifying the standard score $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{c})$ (Song et al., 2020b) to an enhanced version derived from the density $p_t(\mathbf{x}_t | \mathbf{c}) (p_{t, \theta}(\mathbf{x}_t | \mathbf{c}) / p_{t, \theta}(\mathbf{x}_t))^\zeta$. This corresponds to a target score of $\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{c}) + \zeta (\nabla_{\mathbf{x}_t} \log p_{t, \theta}(\mathbf{x}_t | \mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_{t, \theta}(\mathbf{x}_t))$. The objective is to guide the learning process towards distributions that yield higher quality conditional samples.

Accurate estimation of the model probabilities $p_{t, \theta}$ is crucial for the effectiveness of this enhancement. We find that using parameters from an Exponential Moving Average (EMA) of the model during training improves the stability and quality of these estimates, resulting better \mathbf{x}^* and \mathbf{z}^* in Alg. 1.

When training few-step models, direct computation of the enhanced target score gradient typically requires evaluating the model with and without conditioning (for the $p_{t, \theta}$ terms), incurring additional computational cost. To address this, we propose an efficient approximation that leverages a well-pre-trained multi-step model, denoted by parameters θ^* . Instead of computing the score gradient explicitly, the updates for the variables \mathbf{x}^* and \mathbf{z}^* (as used in Alg. 1) are calculated based on features or outputs derived from a single forward pass of the pre-trained model θ^* .

Specifically, we compute $\mathbf{F}_t = \mathbf{F}_{\theta^*}(\mathbf{x}_t, t)$, representing features extracted by the pre-trained model θ^* at time t given input \mathbf{x}_t . The enhanced updates \mathbf{x}^* and \mathbf{z}^* are then computed as follows:

- (a) For $t \in [0, s]$, the updates are: $\mathbf{x}^* \leftarrow \mathbf{x} + \zeta \cdot (\mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x})$, $\mathbf{z}^* \leftarrow \mathbf{z} + \zeta \cdot (\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z})$.
- (b) For $t \in (s, 1]$, the updates are: $\mathbf{x}^* \leftarrow \mathbf{x} + \frac{1}{2} (\mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x})$ and $\mathbf{z}^* \leftarrow \mathbf{z} + \frac{1}{2} (\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z})$.

Table 6: **Hyperparameter configurations for UCGM-{T,S} training and sampling on ImageNet-1K.** We maintain a consistent batch size of 1024 across all experiments. Training durations (epoch counts) are provided in other tables throughout the paper. The table specifies optimizer choices, learning rates, and key parameters for both UCGM-T and UCGM-S variants across different model architectures and datasets.

Task		Optimizer				UCGM-T				UCGM-S			
Resolution	VAE/AE	Model	Type	lr	(β_1, β_2)	Transport	(θ_1, θ_2)	λ	ζ	ρ	κ	\mathcal{T}	ν
Multi-step model training and sampling													
256	E2E-VAE	XL/1	AdamW	0.0002	(0.9,0.95)	Linear	(1.0,1.0)	0	0.67	0	0.5	Auto	1
	SD-VAE	XL/2	AdamW	0.0002	(0.9,0.95)	Linear	(2.4,2.4)	0	0.44	0	0.21	Auto	1
	VA-VAE	XL/1	AdamW	0.0002	(0.9,0.95)	Linear	(1.0,1.0)	0	0.47	0	0.5	Auto	1
512	DC-AE	XL/1	AdamW	0.0002	(0.9,0.95)	Linear	(1.0,1.0)	0	0.57	0	0.46	Auto	1
	SD-VAE	XL/4	AdamW	0.0002	(0.9,0.95)	Linear	(2.4,2.4)	0	0.60	0	0.4	Auto	1
Few-step model training and sampling													
256	E2E-VAE	XL/1	RAdam	0.0001	(0.9,0.999)	Linear	(0.8,1.0)	1	1.3	1	0	{1,0.5}	1
	SD-VAE	XL/2	RAdam	0.0001	(0.9,0.999)	Linear	(0.8,1.0)	1	2.0	1	0	{1,0.3}	1
	VA-VAE	XL/2	RAdam	0.0001	(0.9,0.999)	Linear	(0.8,1.0)	1	2.0	1	0	{1,0.3}	1
512	DC-AE	XL/1	RAdam	0.0001	(0.9,0.999)	Linear	(0.8,1.0)	1	1.5	1	0	{1,0.6}	1
	SD-VAE	XL/4	RAdam	0.0001	(0.9,0.999)	Linear	(0.8,1.0)	1	1.5	1	0	{1,0.5}	1

Table 7: **Comparison of different transport types employed during the sampling and training phases of our UCGM-{T,S}.** “TrigLinear” and “Random” are introduced herein specifically for ablation studies. “TrigLinear” is constructed by combining the transport coefficients of “Linear” and “TrigFlow”. “Random” represents a randomly designed transport type used to demonstrate the generality of our UCGM. Other transport types are adapted from existing methods and transformed into the transport coefficient representation used by UCGM.

	Linear	ReLinear	TrigFlow	EDM ($\sigma(t) = e^{4 \cdot (2.68t - 1.59)}$)	TrigLinear	Random
$\alpha(t)$	t	$1 - t$	$\sin(t \cdot \frac{\pi}{2})$	$\sigma(t) / \sqrt{\sigma^2(t) + 0.25}$	$\sin(t \cdot \frac{\pi}{2})$	$\sin(t \cdot \frac{\pi}{2})$
$\gamma(t)$	$1 - t$	t	$\cos(t \cdot \frac{\pi}{2})$	$1 / \sqrt{\sigma^2(t) + 0.25}$	$\cos(t \cdot \frac{\pi}{2})$	$1 - t$
$\hat{\alpha}(t)$	1	-1	$\cos(t \cdot \frac{\pi}{2})$	$-0.5 / \sqrt{\sigma^2(t) + 0.25}$	1	1
$\hat{\gamma}(t)$	-1	1	$-\sin(t \cdot \frac{\pi}{2})$	$2\sigma(t) / \sqrt{\sigma^2(t) + 0.25}$	-1	$-1 - e^{-5t}$
e.g.,	(Ma et al., 2024)	(Yao et al., 2025)	(Chen et al., 2025c)	(Karras et al., 2022)	N/A	N/A

We consistently set the time threshold $s = 0.75$. This approach allows us to incorporate the guidance from the enhanced target signal with the computational cost equivalent to a single forward evaluation of the pre-trained model θ^* per step. The enhancement ratio ζ is constrained to $[0, \infty)$ in this case.

Baselines. We compare our approach against several SOTA continuous and discrete generative models. We broadly categorize these baselines by their generation process:

- (a) Multi-step models. These methods typically synthesize data through a sequence of steps. We include various diffusion models, encompassing classical formulations like DDPM and score-based models (Song et al., 2020a; Ho et al., 2020), and advanced variants focusing on improved sampling or performance in latent spaces (Dhariwal & Nichol, 2021; Karras et al., 2022; Peebles & Xie, 2023; Zheng et al., 2023; Bao et al., 2023). We also consider flow-matching models (Lipman et al., 2022), which leverage continuous normalizing flows and demonstrate favorable training properties, along with subsequent scaling efforts (Ma et al., 2024; Yu et al., 2024; Yao et al., 2025). Additionally, we also include autoregressive models (Li et al., 2024; Tian et al., 2024; Yu et al., 2023) as the baselines, which generate data sequentially, often in discrete domains.
- (b) Few-step models. These models are designed for efficient, often single-step or few-step, generation. This category includes generative adversarial networks (Goodfellow et al., 2020), which achieve efficient one-step synthesis through adversarial training, and their large-scale variants (Brock et al., 2018; Sauer et al., 2022; Kang et al., 2023). We also evaluate consistency models (Song et al., 2023), proposed for high-quality generation adaptable to few sampling steps, and subsequent techniques aimed at improving their stability and scalability (Song & Dhariwal, 2023; Lu & Song, 2024; Zhou et al., 2025).

Crucially, we demonstrate the compatibility of UCGM-S with models pre-trained using these methods. We show how these models can be represented within the UCGM framework by defining the functions $\alpha(\cdot)$, $\gamma(\cdot)$, $\hat{\alpha}(\cdot)$, and $\hat{\gamma}(\cdot)$. Detailed parameterizations are provided in Tab. 7, with guidance for their specification presented in App. F.2.4.

D.2 EXPERIMENTAL RESULTS ON SMALL DATASETS

Since most existing few-step generation methods (Song et al., 2023; Geng et al., 2024) are limited to training models on low-resolution, small-scale datasets like CIFAR-10 (Krizhevsky et al., 2009a), we conduct our comparative experiments on CIFAR-10 to ensure fair comparison. To demonstrate the versatility of our UCGM, we employ both the "EDM" transport (see Tab. 7 for definition) and the standard 56M-parameter UNet architecture, following established practices in prior work Song et al. (2023); Geng et al. (2024).

Table 8: System-level quality comparison for few-step generation task on unconditional CIFAR-10 (32×32).

Metric	PD (Salimans & Ho, 2022)	2-RF (Liu et al., 2022)	DMD (Yin et al., 2024b)	CD (Song et al., 2023)	sCD (Lu & Song, 2024)
FID (\downarrow)	4.51	4.85	3.77	2.93	2.52
NFE (\downarrow)	2	1	1	2	2
Metric	iCT (Song & Dhariwal, 2023)	ECT (Geng et al., 2024)	sCT (Lu & Song, 2024)	IMM (Zhou et al., 2025)	UCGM
FID (\downarrow)	2.83 / 2.46	3.60 / 2.11	2.97 / 2.06	3.20 / 1.98	2.82 / 2.17
NFE (\downarrow)	1 / 2	1 / 2	1 / 2	1 / 2	1 / 2

As shown in Tab. 8, our UCGM achieves SOTA performance with just 1 NFE (Neural Function Evaluation) while maintaining competitive results for 2 NFEs. These results underscore UCGM’s robust compatibility across diverse datasets, network architectures, and transport types.

D.3 EXPERIMENTAL RESULTS ON LARGE-SCALE UNIFIED MULTIMODAL MODELS

To evaluate the scalability and efficacy of UCGM on Unified Multimodal Models (UMMs), we employ the widely adopted Multi-Modal Diffusion Transformer (MM-DiT) architecture (Esser et al., 2024; Wu et al., 2025a) as our primary backbone. Tab. 9 summarizes the performance across three benchmarks. Crucially, UCGM-S demonstrates superior efficiency, significantly outperforming the standard Euler sampler while maintaining identical NFE budgets.

For evaluation, we employ GenEval (Ghosh et al., 2023) and DPG-Bench (Hu et al., 2024) for text-to-image generation, and WISE (Niu et al., 2025) for world knowledge assessment.

Multi-step Regime ($\lambda = 0$). In the multi-step setting, our trained UCGM-20B achieves performance parity with state-of-the-art generative models. Remarkably, our model achieves these results relying exclusively on publicly available datasets, whereas many SOTA baselines depend on large-scale proprietary data. Our training corpus includes Megalith-10M (Matsubara & Team, 2024), BLIP3o-Pretrain (Chen et al., 2025a), LAION-5B (LAION, 2024), Conceptual 12M (Changpinyo et al., 2021), and text-to-image-2M (He & contributors, 2024) for pre-training, followed by fine-tuning on high-quality instruction-following datasets (BLIP3-o-60K (Chen et al., 2025a), Echo-4o-Image (Ye et al., 2025), and ShareGPT-4o-Image (Chen et al., 2025d)). We adhere to a rigorous training protocol: pre-training for 60k steps (batch size 8,192) and fine-tuning for 3k steps (batch size 1,024) on NVIDIA H800 GPUs (12,976 GPU hours in total).

The training details for UCGM-20B are as follows: we utilize the “Linear” transport type (as defined in Tab. 7), with learning rates of 1×10^{-4} for pre-training and 1×10^{-4} for fine-tuning. The model is trained using the AdamW optimizer with a cosine learning rate schedule.

Few-step Regime ($\lambda = 1$). Scaling distillation to large UMMs presents significant challenges for existing methods. As shown in Tab. 9, standard few-step techniques such as Consistency Models (CM) (Song et al., 2023) suffer from catastrophic model collapse, while MeanFlow (Geng et al., 2025) encounters prohibitive memory costs (OOM). Even when stabilized with our proposed techniques (denoted as CM* and MeanFlow*), these baselines fail to produce competitive results. In contrast, UCGM exhibits exceptional robustness, successfully distilling the 20B-parameter UMM into a few-step generator without compromising stability or requiring excessive memory overhead.

The training configuration for few-step UCGM-20B mirrors that of the multi-step variant, with the exception of a reduced learning rate of 1×10^{-5} .

Dynamic λ Strategy. To bridge the gap between generation quality and inference latency, we propose a dynamic λ training strategy. By conditioning the architecture on a scalar $\lambda \in [0, 1]$ sampled randomly during training, the model learns a continuous spectrum of generation behaviors. This design empowers users to navigate the trade-off between fidelity and speed at inference time. Empirical results confirm that this dynamic approach is highly effective: it not only matches the

Table 9: **System-level comparison of UCGM against SOTA unified multimodal models.** We report inference efficiency (NFE) and generation performance across three benchmarks. **Bold** and underline denote the best and second-best results, respectively. [†] indicates evaluation using LLM-rewritten prompts on GenEval. CM* and MeanFlow* denote baselines re-implemented with our stabilizing techniques or memory-efficient approximations (finite difference) to enable training on large-scale UMMs. All experiments were conducted on NVIDIA H800 GPUs.

Method	NFE ↓	Image Generation		
		GenEval ↑	DPG-Bench ↑	WISE ↑
Multi-step models				
Show-o (Xie et al., 2024b)	50×2	0.68	67.27	0.35
Show-o2-7B (Xie et al., 2025)	50×2	0.76	86.14	0.39
OmniGen (Xiao et al., 2024)	50×2	0.70	81.16	-
OmniGen2 (Wu et al., 2025b)	50×2	0.80 / 0.86 [†]	83.57	-
Janus-Pro (Chen et al., 2025e)	-	0.80	84.19	0.35
MetaQuery-XL (Pan et al., 2025)	30×2	0.78 / 0.80 [†]	81.10	0.55
BLIP3-o-8B (Chen et al., 2025b)	30×2 + 50×2	0.84	81.60	0.62
UniWorld-V1 (Lin et al., 2025)	28×2	0.80 / 0.84 [†]	-	0.55
OpenUni-L-512 (Wu et al., 2025c)	20×2	0.85	81.54	0.52
Bagel (Deng et al., 2025)	50×2	0.82 / 0.88 [†]	-	0.52
Qwen-Image-20B (Wu et al., 2025a)	50×2	0.87	88.32	0.62
UCGM-20B (ours, $\lambda = 0$)	20×2	0.87	86.27	0.58
Qwen-Image-20B + 20-step Euler Sampler	20×2	0.85	82.51	0.53
Qwen-Image-20B + 20-step UCGM-S	20×2	0.87	88.28	0.61
Few-step models				
OpenUni-L-512 \oplus CM (Song et al., 2023) (model collapse)	1	0.0	-	-
Qwen-Image-20B \oplus CM (Song et al., 2023) (model collapse)	1	0.0	-	-
Qwen-Image-20B \oplus CM*	8	0.51	72.17	0.24
Qwen-Image-20B \oplus CM*	4	0.51	71.27	0.22
Qwen-Image-20B \oplus CM*	2	0.44	66.39	0.19
Qwen-Image-20B \oplus CM*	1	0.01	15.41	0.04
Qwen-Image-20B \oplus MeanFlow (Geng et al., 2025) (out of memory)	-	-	-	-
Qwen-Image-20B \oplus MeanFlow*	8	0.49	83.81	0.37
Qwen-Image-20B \oplus MeanFlow*	4	0.44	83.28	0.34
Qwen-Image-20B \oplus MeanFlow*	2	0.31	80.39	0.22
Qwen-Image-20B \oplus MeanFlow*	1	0.05	62.19	0.10
OpenUni-L-512 \oplus UCGM (ours, $\lambda = 1$)	2	0.83	81.05	0.48
OpenUni-L-512 \oplus UCGM (ours, $\lambda = 1$)	1	0.76	73.27	0.42
Qwen-Image-20B\oplusUCGM (ours, $\lambda = 1$)	8	0.60	<u>84.68</u>	0.43
Qwen-Image-20B\oplusUCGM (ours, $\lambda = 1$)	4	0.62	84.23	0.41
Qwen-Image-20B\oplusUCGM (ours, $\lambda = 1$)	2	0.55	67.54	0.23
Qwen-Image-20B\oplusUCGM (ours, $\lambda = 1$)	1	0.32	56.43	0.22
Qwen-Image-20B\oplusUCGM (ours, dynamic λ)	8	0.86	86.49	0.55
Qwen-Image-20B\oplusUCGM (ours, dynamic λ)	4	<u>0.85</u>	83.92	<u>0.52</u>
Qwen-Image-20B\oplusUCGM (ours, dynamic λ)	2	0.82	81.34	0.47
Qwen-Image-20B\oplusUCGM (ours, dynamic λ)	1	0.47	57.39	0.28

flexibility of multi-step models but also secures superior performance across varying NFE budgets, highlighting the versatility of UCGM for unified multimodal generation.

We have included additional qualitative results in Fig. 3, Fig. 4, and Fig. 5 for our trained UCGM-20B.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255



Figure 3: Image generation results from UCGM-20B ($\kappa = 0.0$). Note that the model failed to generate the “START” image.

1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271



Figure 4: Image generation results from UCGM-20B ($\kappa = 0.5$). The model failed to generate the “START” image.

1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288



Figure 5: Image generation results from UCGM-20B ($\kappa = 0.9$). The “START” image was generated successfully.

1291
1292
1293
1294
1295

D.4 DETAILED COMPARISON WITH SOTA METHODS FOR MULTI-STEP GENERATION

Table 10: System-level quality comparison for multi-step generation task on class-conditional ImageNet-1K. Notation $A \oplus B$ denotes the result obtained by combining methods A and B. \downarrow/\uparrow indicate a decrease/increase, respectively, in the metric compared to the baseline performance of the pre-trained models.

METHOD	VAE/AE	Patch Size	Activation Size	NFE (\downarrow)	FID (\downarrow)	IS (\uparrow)	#Params	#Epochs
512 \times 512								
Diffusion & flow-matching models								
ADM-G (Dhariwal & Nichol, 2021)	-	-	-	250 \times 2	7.72	172.71	559M	388
U-ViT-H/4 (Bao et al., 2023)	SD-VAE (Rombach et al., 2022)	4	16 \times 16	50 \times 2	4.05	263.79	501M	400
DiT-XL/2 (Peebles & Xie, 2023)	SD-VAE (Rombach et al., 2022)	2	32 \times 32	250 \times 2	3.04	240.82	675M	600
SiT-XL/2 (Ma et al., 2024)	SD-VAE (Rombach et al., 2022)	2	32 \times 32	250 \times 2	2.62	252.21	675M	600
MaskDiT (Zheng et al., 2023)	SD-VAE (Rombach et al., 2022)	2	32 \times 32	79 \times 2	2.50	256.27	736M	-
EDM2-S (Karras et al., 2024b)	SD-VAE (Rombach et al., 2022)	-	-	63	2.56	-	280M	1678
EDM2-L (Karras et al., 2024b)	SD-VAE (Rombach et al., 2022)	-	-	63	2.06	-	778M	1476
EDM2-XXL (Karras et al., 2024b)	SD-VAE (Rombach et al., 2022)	-	-	63	1.91	-	1.5B	734
DiT-XL/1 \oplus (Chen et al., 2024c)	DC-AE (Chen et al., 2024c)	1	16 \times 16	250 \times 2	2.41	263.56	675M	400
U-ViT-H/1 \oplus (Chen et al., 2024c)	DC-AE (Chen et al., 2024c)	1	16 \times 16	30 \times 2	2.53	255.07	501M	400
REPA-XL/2 (Yu et al., 2024)	SD-VAE (Rombach et al., 2022)	2	32 \times 32	250 \times 2	2.08	274.6	675M	200
DDT-XL/2 (Wang et al., 2025)	SD-VAE (Rombach et al., 2022)	2	32 \times 32	250 \times 2	1.28	305.1	675M	-
GANs & masked & autoregressive models								
VQGAN \oplus (Esser et al., 2021)	-	-	-	256	18.65	-	227M	-
MAGViT-v2 (Yu et al., 2023)	-	-	-	64 \times 2	1.91	324.3	307M	1080
MAR-L (Li et al., 2024)	-	-	-	256 \times 2	1.73	279.9	479M	800
VAR-d36-s (Tian et al., 2024)	-	-	-	10 \times 2	2.63	303.2	2.3B	350
Ours: UCGM-S sampling with models trained by prior works								
EDM2-S (Karras et al., 2024b)	SD-VAE (Rombach et al., 2022)	-	-	40 ¹²³	2.53 ^{10.03}	-	280M	-
EDM2-L (Karras et al., 2024b)	SD-VAE (Rombach et al., 2022)	-	-	50 ¹¹³	2.04 ^{10.02}	-	778M	-
EDM2-XXL (Karras et al., 2024b)	SD-VAE (Rombach et al., 2022)	-	-	40 ¹²³	1.88 ^{10.03}	-	1.5B	-
DDT-XL/2 (Wang et al., 2025)	SD-VAE (Rombach et al., 2022)	2	32 \times 32	200 ¹³⁰⁰	1.25 ^{10.03}	-	675M	-
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 0$)								
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16 \times 16	40	1.48	-	675M	800
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16 \times 16	20	1.68	-	675M	800
Ours-XL/4	SD-VAE (Rombach et al., 2022)	4	16 \times 16	40	1.67	-	675M	320
Ours-XL/4	SD-VAE (Rombach et al., 2022)	4	16 \times 16	20	1.80	-	675M	320
256 \times 256								
Diffusion & flow-matching models								
ADM-G (Dhariwal & Nichol, 2021)	-	-	-	250 \times 2	4.59	186.70	559M	396
U-ViT-H/2 (Bao et al., 2023)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	50 \times 2	2.29	263.88	501M	400
DiT-XL/2 (Peebles & Xie, 2023)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	2.27	278.24	675M	1400
SiT-XL/2 (Ma et al., 2024)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	2.06	277.50	675M	1400
MDT (Gao et al., 2023)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	1.79	283.01	675M	1300
REPA-XL/2 (Yu et al., 2024)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	1.96	264.0	675M	200
REPA-XL/2 (Yu et al., 2024)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	1.42	305.7	675M	800
Light.DiT (Yao et al., 2025)	VA-VAE (Yao et al., 2025)	1	16 \times 16	250 \times 2	2.11	-	675M	64
Light.DiT (Yao et al., 2025)	VA-VAE (Yao et al., 2025)	1	16 \times 16	250 \times 2	1.35	-	675M	800
DDT-XL/2 (Wang et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	1.31	308.1	675M	256
DDT-XL/2 (Wang et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	250 \times 2	1.26	310.6	675M	400
REPA-E-XL (Leng et al., 2025)	E2E-VAE (Leng et al., 2025)	1	16 \times 16	250 \times 2	1.26	314.9	675M	800
GANs & masked & autoregressive models								
VQGAN \oplus (Sun et al., 2024)	-	-	-	-	2.18	-	3.1B	300
MAR-L (Li et al., 2024)	-	-	-	256 \times 2	1.78	296.0	479M	800
MAR-H (Li et al., 2024)	-	-	-	256 \times 2	1.55	303.7	943M	800
VAR-d30-re (Tian et al., 2024)	-	-	-	10 \times 2	1.73	350.2	2.0B	350
Ours: UCGM-S sampling with models trained by prior works								
DDT-XL/2 (Wang et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16 \times 16	100 ¹⁴⁰⁰	1.27 ^{70.01}	-	675M	-
Light.DiT (Yao et al., 2025)	VA-VAE (Yao et al., 2025)	1	16 \times 16	100 ¹⁴⁰⁰	1.21 ^{10.14}	-	675M	-
REPA-E-XL (Leng et al., 2025)	E2E-VAE (Leng et al., 2025)	1	16 \times 16	80 ¹⁴²⁰	1.06 ^{10.20}	-	675M	-
REPA-E-XL (Leng et al., 2025)	E2E-VAE (Leng et al., 2025)	1	16 \times 16	20 ¹⁴⁸⁰	2.00 ^{70.74}	-	675M	-
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 0$)								
Ours-XL/2	SD-VAE (Rombach et al., 2022)	2	16 \times 16	60	1.41	-	675M	400
Ours-XL/1	VA-VAE (Yao et al., 2025)	1	16 \times 16	60	1.21	-	675M	400
Ours-XL/1	E2E-VAE (Leng et al., 2025)	1	16 \times 16	40	1.21	-	675M	800
Ours-XL/1	E2E-VAE (Leng et al., 2025)	1	16 \times 16	20	1.30	-	675M	800

D.5 DETAILED COMPARISON WITH SOTA METHODS FOR FEW-STEP GENERATION

Table 11: System-level quality comparison for few-step generation task on class-conditional ImageNet-1K (512 × 512).

METHOD	VAE/AE	Patch Size	Activation Size	NFE (↓)	FID (↓)	IS	#Params	#Epochs
512 × 512								
Consistency training & distillation								
sCT-M (Lu & Song, 2024)	-	-	-	1	5.84	-	498M	1837
sCT-M (Lu & Song, 2024)	-	-	-	2	5.53	-	498M	1837
sCT-L (Lu & Song, 2024)	-	-	-	1	5.15	-	778M	1274
sCT-L (Lu & Song, 2024)	-	-	-	2	4.65	-	778M	1274
sCT-XXL (Lu & Song, 2024)	-	-	-	1	4.29	-	1.5B	762
sCT-XXL (Lu & Song, 2024)	-	-	-	2	3.76	-	1.5B	762
sCD-M (Lu & Song, 2024)	-	-	-	1	2.75	-	498M	1997
sCD-M (Lu & Song, 2024)	-	-	-	2	2.26	-	498M	1997
sCD-L (Lu & Song, 2024)	-	-	-	1	2.55	-	778M	1434
sCD-L (Lu & Song, 2024)	-	-	-	2	2.04	-	778M	1434
sCD-XXL (Lu & Song, 2024)	-	-	-	1	2.28	-	1.5B	921
sCD-XXL (Lu & Song, 2024)	-	-	-	2	1.88	-	1.5B	921
GANs & masked & autoregressive models								
BigGAN (Brock et al., 2018)	-	-	-	1	8.43	-	160M	-
StyleGAN (Sauer et al., 2022)	-	-	-	1×2	2.41	267.75	168M	-
MAGVIT-v2 (Yu et al., 2023)	-	-	-	64×2	1.91	324.3	307M	1080
VAR- <i>t</i> 36-s (Tian et al., 2024)	-	-	-	10×2	2.63	303.2	2.3B	350
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 0$)								
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16×16	32	1.55	-	675M	800
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16×16	16	1.81	-	675M	800
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16×16	8	3.07	-	675M	800
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16×16	4	74.0	-	675M	800
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 1$)								
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16×16	1	2.42	-	675M	840
Ours-XL/1	DC-AE (Chen et al., 2024c)	1	16×16	2	1.75	-	675M	840
Ours-XL/4	SD-VAE (Rombach et al., 2022)	4	16×16	1	2.63	-	675M	360
Ours-XL/4	SD-VAE (Rombach et al., 2022)	4	16×16	2	2.11	-	675M	360
256 × 256								
Consistency training & distillation								
iCT (Song & Dhariwal, 2023)	-	-	-	2	20.3	-	675M	-
Shortcut-XL/2 (Frans et al., 2024)	SD-VAE (Rombach et al., 2022)	2	16×16	1	10.6	-	676M	250
Shortcut-XL/2 (Frans et al., 2024)	SD-VAE (Rombach et al., 2022)	2	16×16	4	7.80	-	676M	250
Shortcut-XL/2 (Frans et al., 2024)	SD-VAE (Rombach et al., 2022)	2	16×16	128	3.80	-	676M	250
IMM-XL/2 (Zhou et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16×16	1×2	7.77	-	675M	3840
IMM-XL/2 (Zhou et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16×16	2×2	5.33	-	675M	3840
IMM-XL/2 (Zhou et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16×16	4×2	3.66	-	675M	3840
IMM-XL/2 (Zhou et al., 2025)	SD-VAE (Rombach et al., 2022)	2	16×16	8×2	2.77	-	675M	3840
IMM ($\omega = 1.5$)	SD-VAE (Rombach et al., 2022)	2	16×16	1×2	8.05	-	675M	3840
IMM ($\omega = 1.5$)	SD-VAE (Rombach et al., 2022)	2	16×16	2×2	3.99	-	675M	3840
IMM ($\omega = 1.5$)	SD-VAE (Rombach et al., 2022)	2	16×16	4×2	2.51	-	675M	3840
IMM ($\omega = 1.5$)	SD-VAE (Rombach et al., 2022)	2	16×16	8×2	1.99	-	675M	3840
GANs & masked & autoregressive models								
BigGAN (Brock et al., 2018)	-	-	-	1	6.95	-	112M	-
GigaGAN (Kang et al., 2023)	-	-	-	1	3.45	225.52	569M	-
StyleGAN (Sauer et al., 2022)	-	-	-	1×2	2.30	265.12	166M	-
VAR- <i>t</i> 30-re (Tian et al., 2024)	-	-	-	10×2	1.73	350.2	2.0B	350
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 0$)								
Ours-XL/1	VA-VAE (Yao et al., 2025)	1	16×16	16	2.11	-	675M	400
Ours-XL/1	VA-VAE (Yao et al., 2025)	1	16×16	8	6.09	-	675M	400
Ours-XL/1	E2E-VAE (Leng et al., 2025)	1	16×16	16	1.40	-	675M	800
Ours-XL/1	E2E-VAE (Leng et al., 2025)	1	16×16	8	2.68	-	675M	800
Ours: models trained and sampled using UCGM-{T,S} (setting $\lambda = 1$)								
Ours-XL/1	VA-VAE (Yao et al., 2025)	1	16×16	2	1.42	-	675M	432
Ours-XL/1	VA-VAE (Yao et al., 2025)	1	16×16	1	2.19	-	675M	432
Ours-XL/2	SD-VAE (Rombach et al., 2022)	2	16×16	1	2.10	-	675M	424
Ours-XL/1	E2E-VAE (Leng et al., 2025)	1	16×16	1	2.29	-	675M	264

D.6 CASE STUDIES

In this section, we provide several case studies to intuitively illustrate the technical components proposed in this paper.

D.6.1 ANALYSIS OF CONSISTENCY RATIO λ

We evaluate our approach on three synthetic benchmark datasets from `scikit-learn` (Pedregosa et al., 2011): the Two Moons (non-linear separation, see Fig. 6a), S-Curve (manifold structure, see Fig. 6b), and Swiss Roll (non-linear dimensionality reduction, see Fig. 6c). These studies yield two primary observations:

- (a) Our UCGM successfully captures the structure of the data distribution and maps initial points sampled from a Gaussian distribution to the target distribution, regardless of whether the task is few-step ($\lambda = 1$) or multi-step ($\lambda = 0$) generation.

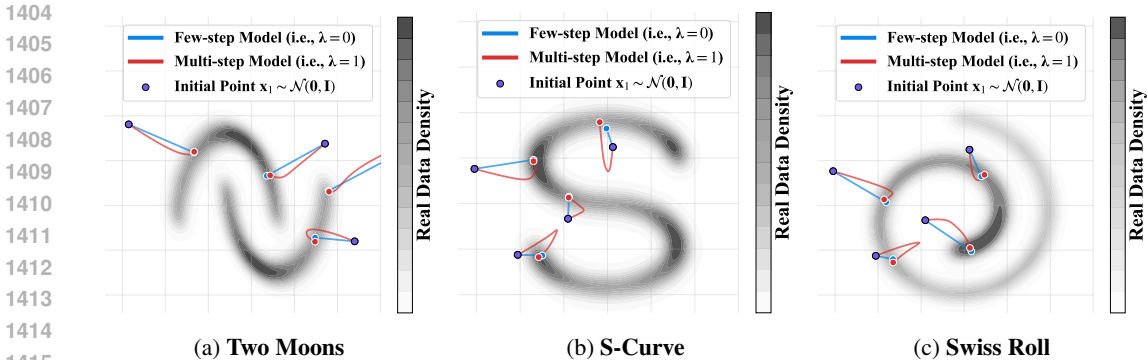


Figure 6: **Case studies of UCGM on three synthetic datasets.** These intuitive studies evaluate the ability of our UCGM to capture the latent data structure for both few-step generation ($\lambda = 1$) and multi-step generation ($\lambda = 0$) tasks.

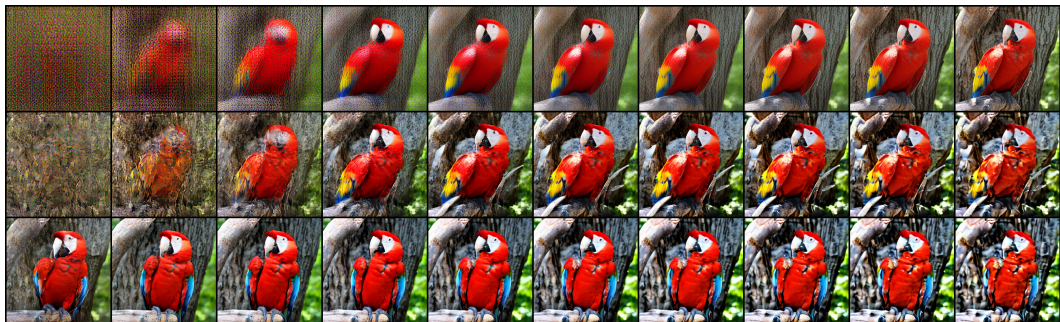


Figure 7: **Intermediate images generated during 60-step sampling from UCGM-S.** Columns display intermediate images \hat{x}_t produced at different timesteps t during a single sampling trajectory, ordered from left to right by decreasing t . Rows correspond to models trained with $\lambda \in \{0.0, 0.5, 1.0\}$, ordered from top to bottom. Note that the initial noise for generating these images is the same.

(b) Models trained for multi-step ($\lambda = 0$) and few-step ($\lambda = 1$) generation map the same initial Gaussian noise to nearly identical target data points.

To further validate these findings and explore additional properties of the consistency ratio λ , we conduct experiments on a real-world dataset (ImageNet-1K). Specifically, we trained three models with three different settings of $\lambda \in \{0.0, 0.5, 1.0\}$.

The experimental results presented in Fig. 7 demonstrate the following:

- (a) For $\lambda = 1.0$, high visual fidelity is achieved early in the sampling process. In contrast, for $\lambda = 0.0$, high visual fidelity emerges in the mid to late stages. For $\lambda = 0.5$, high-quality images appear in the mid-stage of sampling.
- (b) Despite being trained with different settings of λ values, the models produce remarkably similar generated images.

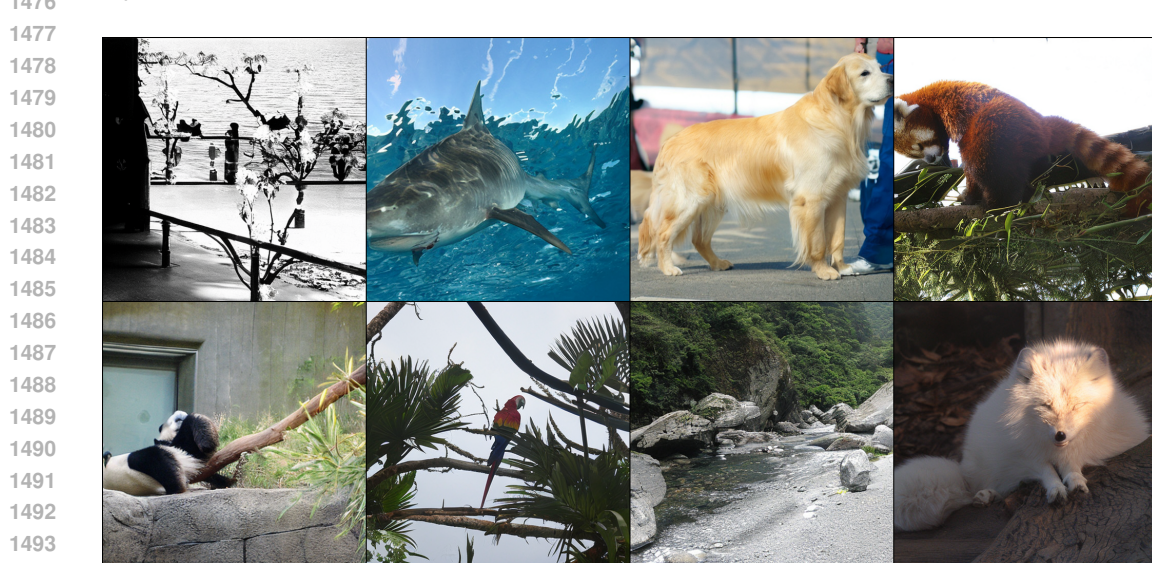
In summary, we posit that while the setting of λ affects the dynamics of the generation process, it does not substantially impact the final generated image quality. Detailed analysis of these phenomena is provided in App. F.1.1, App. F.1.3 and App. F.1.4.

D.6.2 ANALYSIS OF TRANSPORT TYPES

Generated samples, obtained using UCGM-S with two distinct pre-trained models from prior works, are presented in Fig. 9 and Fig. 8. When using the identical initial Gaussian noise for both models, the generated images exhibit notable visual similarity. This observation is unexpected, considering the models were trained independently (Karras et al., 2024b; Wang et al., 2025) using distinct algorithms, transport formulations, network architectures, and data augmentation strategies. The similarity suggests that despite these differences, the learned probability flow ODEs may be converging to similar solutions. See App. F.1.2 for a comprehensive analysis of this phenomenon.



1475 Figure 8: Visualization of generated images (512×512) from pre-trained EDM2-S (Karras et al., 2024b).



1495 Figure 9: Visualization of generated images (512×512) from pre-trained DDT-XL/2 (Wang et al., 2025).

1496 D.7 ANALYSIS OF PRE-TRAINED MODEL TUNING

1497
1498
1499
1500 Table 12: System-level quality comparison for few-step generation on class-conditional ImageNet-1K
1501 after tuning. Notation \downarrow/\uparrow indicate performance decrease/increase relative to the baseline "Generation (Gen.)"
1502 performance of the "Original (Orig.)" pre-trained models at the respective NFE. Tuning time is evaluated on a
1503 cluster with 8 NVIDIA H800 GPUs.

METHOD	#Params	Orig. Few-step Gen.		Orig. Multi-step Gen.		Tuning Efficiency		Tuned Few-step Gen.	
		NFE (\downarrow)	FID (\downarrow)	NFE (\downarrow)	FID (\downarrow)	#Epochs	Time	NFE (\downarrow)	FID (\downarrow)
REPA (Yu et al., 2024)	675M	2	177	80	1.86	0.64	\approx 13 minutes	2	1.95 ^{\downarrow175}
Lightning-DiT (Yao et al., 2025)	675M	2	217	80	1.49	0.64	\approx 10 minutes	2	2.06 ^{\downarrow215}
REPA-E (Leng et al., 2025)	675M	2	193	80	1.54	0.40	\approx 8 minutes	2	1.39 ^{\downarrow192}
DDT (Wang et al., 2025)	675M	2	191	80	1.46	0.32	\approx 11 minutes	2	1.90 ^{\downarrow189}

1509
1510 In addition to our previous studies and experiments, where we demonstrated that our UCGM-S is
1511 a plug-and-play, training-free method for accelerating the sampling process of given pre-trained
models from prior works (Yu et al., 2024; Yao et al., 2025; Leng et al., 2025; Wang et al., 2025) (cf.,

App. D.4), we have also proven that our UCGM-T is an efficient and effective unified framework for training both few-step and multi-step continuous generative models (cf., App. D.5 and App. D.4).

In this section, we evaluate the effectiveness of UCGM for tuning existing pre-trained generative models to enhance few-step generation performance. Tab. 12 presents the experimental results.

Specifically, the results demonstrate that UCGM-T facilitates the efficient conversion of continuous multi-step generative models (including diffusion and flow matching models) into high-performance few-step variants through minimal fine-tuning. For instance, the pre-trained REPA-E model (Leng et al., 2025), exhibiting 1.54 FID at 80 NFEs and 193 FID at 2 NFEs, can be efficiently tuned using UCGM-T in *only approximately 8 minutes (0.4 epoch)*. This tuning process yields a model *achieving 1.39 FID at 2 NFEs*, representing a substantial improvement in few-step generation quality with negligible tuning cost.

D.8 ABLATION STUDY ON UCGM TECHNIQUES

Tab. 13 and Tab. 14 present the ablation studies on the proposed techniques in UCGM-T and UCGM-S, conducted under the same experimental setup as Tab. 12.

For UCGM-T, we observe that removing the generalized time distribution (GTD) does not affect the performance. This is expected, since GTD is designed to generalize beyond specific time distributions, whereas our experiments are conducted under a uniform distribution. In contrast, removing both GTD and the stabilizing technique leads to a significant degradation in FID scores across all backbones, demonstrating the importance of our proposed training stabilization method.

For UCGM-S, the stochastic sampling technique, which unifies ODE and SDE samplers in a generalized formulation, does not change the quantitative performance under our setting. However, removing both the stochastic component and the extrapolation strategy results in a substantial increase in NFE, indicating that the extrapolation-based acceleration is highly effective for efficient sampling.

Table 13: Ablation study on UCGM-T techniques.

UCGM-T ($\lambda = 1$)	DDT (Wang et al., 2025)		Light.DiT (Yao et al., 2025)		REPA-E (Leng et al., 2025)	
	FID ↓	NFE ↓	FID ↓	NFE ↓	FID ↓	NFE ↓
original	1.90	2	2.06	2	1.39	2
w/o GTD	1.90	2	2.06	2	1.39	2
w/o GTD & stab.	4.75	2	13.87	2	2.45	2

Table 14: Ablation study on UCGM-S techniques.

UCGM-S	DDT (Wang et al., 2025)		Light.DiT (Yao et al., 2025)		REPA-E (Leng et al., 2025)	
	FID ↓	NFE ↓	FID ↓	NFE ↓	FID ↓	NFE ↓
original	1.27	100	1.21	100	1.06	80
w/o stoch.	1.27	100	1.21	100	1.06	80
w/o stoch. & extr.	1.26	500	1.35	500	1.26	500

E PSEUDOCODE

E.1 TRAINING ALGORITHM FOR UCGM-T

Algorithm 1 (UCGM-T). A Unified and Efficient Trainer for Few-step and Multi-step Continuous Generative Models (including Diffusion, Flow Matching, and Consistency Models)

Require: Dataset D , transport coefficients $\{\alpha(\cdot), \gamma(\cdot), \hat{\alpha}(\cdot), \hat{\gamma}(\cdot)\}$, neural network F_θ , enhancement ratio ζ , Beta distribution parameters (θ_1, θ_2) , learning rate η , $\theta^- = \theta$ only in value.

Ensure: Trained neural network F_θ for generating samples from $p(\mathbf{x})$.

```

1: repeat
2:   Sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ,  $\mathbf{x} \sim D$ ,  $t \sim \phi(t) := \text{Beta}(\theta_1, \theta_2)$ 
3:   Compute input data, such as  $\mathbf{x}_t = \alpha(t) \cdot \mathbf{z} + \gamma(t) \cdot \mathbf{x}$  and  $\mathbf{x}_{\lambda t} = \alpha(\lambda t) \cdot \mathbf{z} + \gamma(\lambda t) \cdot \mathbf{x}$ 
4:   Compute model output  $F_t = F_\theta(\mathbf{x}_t, t)$  and set  $\mathbf{z}^* = \mathbf{z}$  and  $\mathbf{x}^* = \mathbf{x}$ 
5:   if  $\zeta \in (0, 1)$  then
6:     Get enhanced  $\mathbf{x}^* = \xi(\mathbf{x}, t, \mathbf{f}^{\mathbf{x}}(F_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t), \mathbf{f}^{\mathbf{x}}(F_{\theta^-}(\mathbf{x}_t, t), \emptyset, \mathbf{x}_t, t))$  and  $\mathbf{z}^* = \xi(\mathbf{z}, t, \mathbf{f}^{\mathbf{z}}(F_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t), \mathbf{f}^{\mathbf{z}}(F_{\theta^-}(\mathbf{x}_t, t), \emptyset, \mathbf{x}_t, t))$  {Note that  $\xi(\mathbf{a}, t, \mathbf{b}, \mathbf{d}) := \mathbf{a} + (\zeta + \mathbf{1}_{t>s}(\frac{1}{2} - \zeta)) \cdot (\mathbf{b} - \mathbf{1}_{t>s} \cdot \mathbf{a} - \mathbf{d}(1 - \mathbf{1}_{t>s}))$ , where  $\mathbf{1}(\cdot)$  is the indicator function}
7:   end if
8:   if  $\lambda \in [0, 1)$  then
9:     Compute  $\mathbf{z}_t^* = \hat{\alpha}(t) \cdot \mathbf{z}^* + \hat{\gamma}(t) \cdot \mathbf{x}^*$  and  $\mathbf{z}_{\lambda t}^* = \hat{\alpha}(\lambda t) \cdot \mathbf{z}^* + \hat{\gamma}(\lambda t) \cdot \mathbf{x}^*$ 
10:    Compute  $D(t) = \alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)$ ,  $B(t) = \frac{\alpha(t)}{D(t)}$ 
11:    Let  $C(t) = \frac{\alpha(t)}{2D(t)}$ ,  $A(t) = B(t) - B(\lambda t)$  and  $\hat{\omega}(t) = C(t) \cdot A(t)$ 
12:    Let  $\Delta \mathbf{z}_t = \mathbf{z}_t^* - \mathbf{z}_{\lambda t}^*$ 
13:    Compute loss  $\mathcal{L}_t(\theta) = \|\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{z}_t^*\|_2^2 + \frac{B(\lambda t)}{\hat{\omega}(t)} \|\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) - \Delta \mathbf{z}_t\|_2^2$ 
14:    else if  $\lambda = 1$  then
15:      Compute  $\mathbf{x}_{t+\epsilon}^* = \alpha(t+\epsilon) \cdot \mathbf{z}^* + \gamma(t+\epsilon) \cdot \mathbf{x}^*$  and  $\mathbf{x}_{t-\epsilon} = \alpha(t-\epsilon) \cdot \mathbf{z}^* + \gamma(t-\epsilon) \cdot \mathbf{x}^*$ 
16:      Let  $\Delta \mathbf{f}_t^{\mathbf{x}} = \mathbf{f}^{\mathbf{x}}(F_{\theta^-}(\mathbf{x}_{t+\epsilon}, t+\epsilon), \mathbf{x}_{t+\epsilon}^*, t+\epsilon) - \mathbf{f}^{\mathbf{x}}(F_{\theta^-}(\mathbf{x}_{t-\epsilon}, t-\epsilon), \mathbf{x}_{t-\epsilon}^*, t-\epsilon)$ 
17:      Let  $\Delta B(t) = \frac{\alpha(t+\epsilon)}{\alpha(t+\epsilon)\hat{\gamma}(t+\epsilon) - \hat{\alpha}(t+\epsilon)\gamma(t+\epsilon)} - \frac{\alpha(t-\epsilon)}{\alpha(t-\epsilon)\hat{\gamma}(t-\epsilon) - \hat{\alpha}(t-\epsilon)\gamma(t-\epsilon)}$ 
18:      Compute  $F_t^{\text{target}} = F_{\theta^-}(\mathbf{x}_t, t) - 2 \cdot \text{clip}(\frac{\Delta \mathbf{f}_t^{\mathbf{x}}}{\Delta B(t)}, -1, 1)$ 
19:      Compute loss  $\mathcal{L}_t(\theta) = \|F_t - F_t^{\text{target}}\|_2^2$ 
20:    end if
21:    Update  $\theta \leftarrow \theta - \eta \nabla_\theta \int_0^1 \phi(t) \mathcal{L}_t(\theta) dt$ 
22:  until Convergence

```

E.2 SAMPLING ALGORITHM FOR UCGM-S

Algorithm 2 (UCGM-S). A Unified and Efficient Sampler for Few-step and Multi-step Continuous Generative Models (including Diffusion, Flow Matching, and Consistency Models)

Require: Initial $\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, transport coefficients $\{\alpha(\cdot), \gamma(\cdot), \hat{\alpha}(\cdot), \hat{\gamma}(\cdot)\}$, trained model F_{θ} , sampling steps N , order $\nu \in \{1, 2\}$, time schedule \mathcal{T} , extrapolation ratio κ , stochastic ratio ρ .

Ensure: Final generated sample $\tilde{\mathbf{x}} \sim p(\mathbf{x})$ and history samples $\{\tilde{\mathbf{x}}_i\}_{i=0}^N$ over generation process.

- 1: Let $N \leftarrow \lfloor (N + 1)/2 \rfloor$ if using second order sampling ($\nu = 2$) {Adjusts total steps to match first-order evaluation count}
 - 2: **for** $i = 0$ to $N - 1$ **do**
 - 3: Compute model output $F = F_{\theta^-}(\tilde{\mathbf{x}}, t_i)$, and then $\hat{\mathbf{x}}_i = f^{\mathbf{x}}(F, \tilde{\mathbf{x}}, t_i)$ and $\hat{\mathbf{z}}_i = f^{\mathbf{z}}(F, \tilde{\mathbf{x}}, t_i)$
 - 4: **if** $i \geq 1$ **then**
 - 5: Compute extrapolated estimation $\hat{\mathbf{z}} = \hat{\mathbf{z}}_i + \kappa \cdot (\hat{\mathbf{z}}_i - \hat{\mathbf{z}}_{i-1})$ and $\hat{\mathbf{x}} = \hat{\mathbf{x}}_i + \kappa \cdot (\hat{\mathbf{x}}_i - \hat{\mathbf{x}}_{i-1})$
 - 6: **end if**
 - 7: Sample $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ {An example choice of ρ for performing SDE-similar sampling is:
 $\rho = \text{clip}(\frac{|t_i - t_{i+1}| \cdot 2\alpha(t_i)}{\alpha(t_{i+1})}, 0, 1)$ }
 - 8: Compute estimated next time sample $\mathbf{x}' = \alpha(t_{i+1}) \cdot (\sqrt{1 - \rho} \cdot \hat{\mathbf{z}} + \sqrt{\rho} \cdot \mathbf{z}) + \gamma(t_{i+1}) \cdot \hat{\mathbf{x}}$
 - 9: **if** order $\nu = 2$ **and** $i < N - 1$ **then**
 - 10: Compute prediction $F' = F_{\theta}(\mathbf{x}', t_{i+1})$, $\hat{\mathbf{x}}' = f^{\mathbf{x}}(F', \mathbf{x}', t_{i+1})$ and $\hat{\mathbf{z}}' = f^{\mathbf{z}}(F', \mathbf{x}', t_{i+1})$
 - 11: Compute corrected next time sample $\mathbf{x}' = \tilde{\mathbf{x}} \cdot \frac{\gamma(t_{i+1})}{\gamma(t_i)} + \left(\alpha(t_{i+1}) - \frac{\gamma(t_{i+1})\alpha(t_i)}{\gamma(t_i)} \right) \cdot \frac{\hat{\mathbf{x}} + \hat{\mathbf{x}}'}{2}$
 - 12: **end if**
 - 13: Reset $\tilde{\mathbf{x}} \leftarrow \mathbf{x}'$
 - 14: **end for**
-

1674 F THEORETICAL ANALYSIS

1675 F.1 MAIN RESULTS

1677 F.1.1 LEARNING OBJECTIVE WHEN $\lambda = 0$

1678 Recall that $(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x})$ is a pair of latent and data variables (typically independent), and let
 1679 $t \in [0, 1]$. We have four differentiable scalar functions $\alpha, \gamma, \hat{\alpha}, \hat{\gamma}: [0, 1] \rightarrow \mathbb{R}$, the *noisy interpolant*
 1680 $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$ and $\mathbf{F}_t = \mathbf{F}_\theta(\mathbf{x}_t, t)$. We define the \mathbf{x} - and \mathbf{z} -prediction functions by

$$1681 \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\alpha(t)\mathbf{F}_t - \hat{\alpha}(t)\mathbf{x}_t}{\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)}, \quad \text{and} \quad \mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\hat{\gamma}(t)\mathbf{x}_t - \gamma(t)\mathbf{F}_t}{\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)}.$$

1684 Finally, let $\hat{w}(t) > 0$ be a weight function. We consider the \mathbf{x} - and \mathbf{z} -prediction losses

$$1685 \mathcal{L}_{\mathbf{x}}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{w}(t)} \left\| \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x} \right\|_2^2 \right],$$

$$1686 \mathcal{L}_{\mathbf{z}}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{w}(t)} \left\| \mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z} \right\|_2^2 \right].$$

1689 Recall that our unified loss function is defined by:

$$1690 \mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \frac{1}{\hat{w}(t)} \left\| \mathbf{f}^{\mathbf{x}}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t) \right\|_2^2.$$

1693 We have $\mathcal{L}(\theta) = \mathcal{L}_{\mathbf{x}}(\theta)$ when $\lambda = 0$, since $\mathbf{f}^{\mathbf{x}}(\mathbf{F}_0, \mathbf{x}_0, 0) = \mathbf{x}$. Then, we define the direct-field loss

$$1694 \mathcal{L}_{\mathbf{F}}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}), t} \left[w_{\mathbf{F}}(t) \left\| \mathbf{F}_t - (\hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}) \right\|_2^2 \right], \quad w(t) > 0.$$

1697 **Lemma 1 (Equivalence of \mathbf{x} -prediction and direct-field loss).** *For all θ ,*

$$1698 \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x} = \frac{\alpha(t)}{\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)} \left[\mathbf{F}_t - (\hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}) \right].$$

1701 Hence

$$1702 \mathcal{L}_{\mathbf{x}}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}), t} \left[\frac{\alpha(t)^2}{\hat{w}(t) (\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t))^2} \left\| \mathbf{F}_t - (\hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}) \right\|_2^2 \right],$$

1705 so $\mathcal{L}_{\mathbf{x}}$ is equivalent to $\mathcal{L}_{\mathbf{F}}$ with

$$1706 w_{\mathbf{F}}(t) = \frac{\alpha(t)^2}{\hat{w}(t) (\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t))^2}.$$

1710 *Proof.* Compute

$$1711 \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x} = \frac{\alpha(t)\mathbf{F}_t - \hat{\alpha}(t)\mathbf{x}_t}{\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)} - \mathbf{x}.$$

1714 Since $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$, the numerator becomes

$$1715 \alpha\mathbf{F}_t - \hat{\alpha}(\alpha\mathbf{z} + \gamma\mathbf{x}) - (\alpha\hat{\gamma} - \hat{\alpha}\gamma)\mathbf{x} = \alpha(t) \left[\mathbf{F}_t - (\hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}) \right].$$

1717 Dividing by $\alpha\hat{\gamma} - \hat{\alpha}\gamma$ yields the desired factorization. Substituting into $\mathcal{L}_{\mathbf{x}}$ gives the weight $w(t)$ as
 1718 above. \square

1720 **Lemma 2 (Equivalence of \mathbf{z} -Prediction and Direct-Field Loss).** *For all θ ,*

$$1721 \mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z} = \frac{\gamma(t)}{\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t)} \left[(\hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}) - \mathbf{F}_t \right].$$

1724 Hence

$$1725 \mathcal{L}_{\mathbf{z}}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}), t} \left[\frac{\gamma(t)^2}{\hat{w}(t) (\alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t))^2} \left\| \mathbf{F}_t - (\hat{\alpha}(t)\mathbf{z} + \hat{\gamma}(t)\mathbf{x}) \right\|_2^2 \right],$$

so \mathcal{L}_z is equivalent to \mathcal{L}_F with

$$w_F(t) = \frac{\gamma(t)^2}{\hat{\omega}(t) (\alpha(t) \hat{\gamma}(t) - \hat{\alpha}(t) \gamma(t))^2}.$$

Proof. Compute

$$\mathbf{f}^z(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z} = \frac{\hat{\gamma}(t) \mathbf{x}_t - \gamma(t) \mathbf{F}_t}{\alpha(t) \hat{\gamma}(t) - \hat{\alpha}(t) \gamma(t)} - \mathbf{z}.$$

Using $\mathbf{x}_t = \alpha \mathbf{z} + \gamma \mathbf{x}$, the numerator is

$$\hat{\gamma}(\alpha \mathbf{z} + \gamma \mathbf{x}) - \gamma \mathbf{F}_t - (\alpha \hat{\gamma} - \hat{\alpha} \gamma) \mathbf{z} = \gamma(t) [\hat{\alpha}(t) \mathbf{z} + \hat{\gamma}(t) \mathbf{x} - \mathbf{F}_t].$$

Dividing by $\alpha \hat{\gamma} - \hat{\alpha} \gamma$ gives the factorization. Substitution into \mathcal{L}_z yields the stated equivalence. \square

F.1.2 CLOSED-FORM SOLUTION ANALYSIS WHEN $\lambda = 0$

when $\lambda = 0$, we aim to derive the Probability Flow Ordinary Differential Equation (PF-ODE) (Song et al., 2020b) corresponding to a defined forward process from time 0 to 1.

Lemma 3 (Probability Flow ODE for the linear Gaussian forward process). *Let $p(\mathbf{x})$ be a data distribution on \mathbb{R}^d , and let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ be independent of \mathbf{x} . Let $\alpha, \gamma : [0, 1] \rightarrow \mathbb{R}$ be continuously differentiable scalar functions satisfying*

$$\alpha(0) = 0, \quad \alpha(1) = 1, \quad \gamma(0) = 1, \quad \gamma(1) = 0,$$

and assume $\gamma(t) \neq 0$ for $t \in (0, 1)$. Define the forward process

$$\mathbf{x}_t = \alpha(t) \mathbf{z} + \gamma(t) \mathbf{x}, \quad t \in [0, 1],$$

so that $\mathbf{x}_0 = \mathbf{x} \sim p(\mathbf{x})$ and $\mathbf{x}_1 = \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. Let $p_t(\mathbf{x}_t)$ denote the marginal density of \mathbf{x}_t . Then the Probability Flow ODE for this process,

$$\frac{d\mathbf{x}_t}{dt} = \mathbf{f}(\mathbf{x}_t, t) - \frac{1}{2} g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t),$$

takes the explicit form

$$\frac{d\mathbf{x}_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t - \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \quad (7)$$

Proof. We first represent the forward process \mathbf{x}_t as the solution of a SDE (Song et al., 2020b):

$$d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t) dt + g(t) d\mathbf{w}_t,$$

where \mathbf{w}_t is a standard d -dimensional Wiener process, and where $\mathbf{f}(\cdot, t)$ and $g(t)$ are to be determined so that $\mathbf{x}_t = \alpha(t) \mathbf{z} + \gamma(t) \mathbf{x}$ in law.

1. Drift term via the conditional mean. Since \mathbf{z} and \mathbf{x} are independent,

$$\mathbb{E}[\mathbf{x}_t | \mathbf{x}] = \gamma(t) \mathbf{x}.$$

Differentiating in t gives

$$\frac{d}{dt} \mathbb{E}[\mathbf{x}_t | \mathbf{x}] = \gamma'(t) \mathbf{x}. \quad (8)$$

On the other hand, we use the method of separation of variables, which is a classical method in solving PDEs, and we set the drift term as $\mathbf{f}(\mathbf{x}_t, t) = H(t) \mathbf{x}_t$ for some matrix $H(t)$, then

$$\frac{d}{dt} \mathbb{E}[\mathbf{x}_t | \mathbf{x}] = H(t) \mathbb{E}[\mathbf{x}_t | \mathbf{x}] = H(t) \gamma(t) \mathbf{x}. \quad (9)$$

Comparing (8) and (9) yields $H(t) = \gamma'(t)/\gamma(t) \mathbf{I}_d$, so

$$\mathbf{f}(\mathbf{x}_t, t) = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t.$$

2. Diffusion term via the conditional variance. The covariance of \mathbf{x}_t given \mathbf{x} is

$$\text{Var}(\mathbf{x}_t | \mathbf{x}) = \alpha(t)^2 \mathbf{I}_d.$$

For a linear SDE with drift matrix $H(t)$ and scalar diffusion $g(t)$, the covariance $\Sigma(t)$ satisfies the following Lyapunov equation (Jiménez, 2015):

$$\frac{d\Sigma(t)}{dt} = H(t)\Sigma(t) + \Sigma(t)H(t)^\top + g(t)^2 \mathbf{I}_d.$$

Substitute $\Sigma(t) = \alpha(t)^2 \mathbf{I}_d$ and $H(t) = \frac{\gamma'(t)}{\gamma(t)} \mathbf{I}_d$. Since $\frac{d}{dt}(\alpha(t)^2) = 2\alpha(t)\alpha'(t)$, we get

$$2\alpha(t)\alpha'(t)\mathbf{I}_d = 2\frac{\gamma'(t)}{\gamma(t)}\alpha(t)^2\mathbf{I}_d + g(t)^2\mathbf{I}_d.$$

Rearranging yields

$$g(t)^2 = 2\alpha(t)\alpha'(t) - 2\frac{\gamma'(t)}{\gamma(t)}\alpha(t)^2.$$

3. Probability Flow ODE. By general theory (see, e.g., de Bortoli et al.), the probability flow ODE associated with the SDE $d\mathbf{x}_t = \mathbf{f}(\mathbf{x}_t, t) dt + g(t) d\mathbf{w}_t$ is

$$\frac{d\mathbf{x}_t}{dt} = \mathbf{f}(\mathbf{x}_t, t) - \frac{1}{2}g(t)^2 \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t).$$

Substituting the expressions for \mathbf{f} and g^2 above gives

$$\frac{d\mathbf{x}_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t - \left[\alpha(t)\alpha'(t) - \frac{\gamma'(t)}{\gamma(t)}\alpha(t)^2 \right] \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t),$$

i.e.,

$$\mathbf{f}(\mathbf{x}_t, t) = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t, \quad g(t)^2 = 2\alpha(t)\alpha'(t) - 2\frac{\gamma'(t)}{\gamma(t)}\alpha(t)^2.$$

which is exactly the claimed formula (7). This result is also proved with another method in (Holdrieth & Erives, 2025) (see Proposition 1 in their section 4.2). \square

Lemma 4 (Tweedie formula (Song et al., 2020b) for the linear Gaussian model). *Under the linear Gaussian interpolation model $\mathbf{x}_t | \mathbf{x} \sim \mathcal{N}(\gamma(t)\mathbf{x}, \alpha^2(t)\mathbf{I})$, the conditional expectation of \mathbf{x} given \mathbf{x}_t is*

$$\mathbb{E}[\mathbf{x} | \mathbf{x}_t] = \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)}.$$

Proof. We write the conditional expectation by Bayes' rule:

$$\mathbb{E}[\mathbf{x} | \mathbf{x}_t] = \int \mathbf{x} p(\mathbf{x} | \mathbf{x}_t) d\mathbf{x} = \frac{1}{p_t(\mathbf{x}_t)} \int \mathbf{x} p_t(\mathbf{x}_t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x},$$

where $p_t(\mathbf{x}_t) = \int p_t(\mathbf{x}_t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$.

Since $p_t(\mathbf{x}_t | \mathbf{x}) = (2\pi\alpha^2(t))^{-d/2} \exp(-\frac{1}{2\alpha^2(t)}\|\mathbf{x}_t - \gamma(t)\mathbf{x}\|^2)$, we have

$$\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t | \mathbf{x}) = -\frac{1}{\alpha^2(t)} (\mathbf{x}_t - \gamma(t)\mathbf{x}) p_t(\mathbf{x}_t | \mathbf{x}).$$

Differentiating the marginal,

$$\nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t) = \int \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} = -\frac{1}{\alpha^2(t)} \int (\mathbf{x}_t - \gamma(t)\mathbf{x}) p_t(\mathbf{x}_t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$

Multiply by $-\alpha^2(t)$ and split:

$$-\alpha^2(t) \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t) = \mathbf{x}_t p_t(\mathbf{x}_t) - \gamma(t) \int \mathbf{x} p_t(\mathbf{x}_t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}.$$

Rearrange and divide by $\gamma(t)p_t(\mathbf{x}_t)$:

$$\frac{1}{p_t(\mathbf{x}_t)} \int \mathbf{x} p_t(\mathbf{x}_t | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} = \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} p_t(\mathbf{x}_t) / p_t(\mathbf{x}_t)}{\gamma(t)} = \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)}.$$

Hence $\mathbb{E}[\mathbf{x} | \mathbf{x}_t] = (\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)) / \gamma(t)$, as claimed. \square

Lemma 5 (Optimal predictors as conditional expectations). For each fixed t and observed \mathbf{x}_t , the pointwise minimizers $\mathbf{f}_*^{\mathbf{x}}$ and $\mathbf{f}_*^{\mathbf{z}}$ for the objective function $\mathcal{L}(\boldsymbol{\theta})$ satisfy

$$\mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t], \quad \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = \mathbb{E}[\mathbf{z} \mid \mathbf{x}_t].$$

Proof. Fix t and \mathbf{x}_t . By [Lem. 1](#) and [Lem. 2](#), we conclude that the minimizers of $\mathcal{L}(\boldsymbol{\theta})$ are equivalent to those of $\mathcal{L}_{\mathbf{x}}$ and $\mathcal{L}_{\mathbf{z}}$.

Then, up to an additive constant independent of $\mathbf{f}^{\mathbf{x}}$, the contribution of (t, \mathbf{x}_t) to $\mathcal{L}_{\mathbf{x}}$ is

$$\mathcal{J}_{\mathbf{x}}(\mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t)) = \mathbb{E}[\|\mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x}\|_2^2 \mid \mathbf{x}_t].$$

For any random vector X , the function $w \mapsto \mathbb{E}\|w - X\|^2$ is uniquely minimized at $w = \mathbb{E}[X]$. Therefore

$$\mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \arg \min_w \mathbb{E}[\|w - \mathbf{x}\|^2 \mid \mathbf{x}_t] = \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t].$$

The same argument applies to

$$\mathcal{J}_{\mathbf{z}}(\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t)) = \mathbb{E}[\|\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z}\|_2^2 \mid \mathbf{x}_t],$$

yielding

$$\mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = \mathbb{E}[\mathbf{z} \mid \mathbf{x}_t]. \quad \square$$

Theorem 2. Under the linear Gaussian interpolation model $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$, with $\mathbf{z} \sim \mathcal{N}(0, \mathbf{I})$ independent of \mathbf{x} , we have

$$\mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)}, \quad \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = -\alpha(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t).$$

Then for every t ,

$$\alpha'(t) \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) + \gamma'(t) \mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t - \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha^2(t) \right] \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t).$$

As a result, by [Lem. 3](#), we conclude:

$$\frac{d\mathbf{x}_t}{dt} = \alpha'(t) \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) + \gamma'(t) \mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t)$$

Proof. **Tweedie formula for $\mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t)$.** According to [Lem. 5](#) and [Lem. 4](#), we have

$$\mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t] = \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)}.$$

Derivation of $\mathbb{E}[\mathbf{z} \mid \mathbf{x}_t]$ for $\mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t)$. From $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$ we solve $\mathbf{z} = (\mathbf{x}_t - \gamma(t)\mathbf{x})/\alpha(t)$. Taking conditional expectation and substituting the above,

$$\begin{aligned} \mathbb{E}[\mathbf{z} \mid \mathbf{x}_t] &= \frac{1}{\alpha(t)} \left(\mathbf{x}_t - \gamma(t) \mathbb{E}[\mathbf{x} \mid \mathbf{x}_t] \right) \\ &= \frac{1}{\alpha(t)} \left(\mathbf{x}_t - \gamma(t) \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)} \right) = -\alpha(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \end{aligned}$$

Thus, according to [Lem. 5](#), we can obtain

$$\mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = -\alpha(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t).$$

Combine to obtain the claimed identity.

$$\begin{aligned} &\alpha'(t) \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) + \gamma'(t) \mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) \\ &= \alpha'(t) [-\alpha(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)] + \gamma'(t) \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)} \\ &= -\alpha(t) \alpha'(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) + \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t + \frac{\gamma'(t)}{\gamma(t)} \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) \\ &= \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t - \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha^2(t) \right] \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \end{aligned}$$

This matches the claimed formula. \square

Remark 2 (Velocity field of the flow ODE). Given \mathbf{x} and \mathbf{z} , the field $\mathbf{v}^{(\mathbf{z}, \mathbf{x})}(\mathbf{y}, t) = \alpha'(t)\mathbf{z} + \gamma'(t)\mathbf{x}$ could transport \mathbf{z} to \mathbf{x} , so the velocity field of the flow ODE can be computed as

$$\begin{aligned} \mathbf{v}^*(\mathbf{x}_t, t) &= \mathbb{E}_{(\mathbf{z}, \mathbf{x})|\mathbf{x}_t} \left[\mathbf{v}^{(\mathbf{z}, \mathbf{x})}(\mathbf{x}_t, t) | \mathbf{x}_t \right] \\ &= \mathbb{E}_{(\mathbf{z}, \mathbf{x})|\mathbf{x}_t} [\alpha'(t)\mathbf{z} + \gamma'(t)\mathbf{x} | \mathbf{x}_t] \\ &= \alpha'(t) \cdot \mathbb{E}[\mathbf{z} | \mathbf{x}_t] + \gamma'(t) \cdot \mathbb{E}[\mathbf{x} | \mathbf{x}_t] \\ &= \alpha'(t) \cdot \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) + \gamma'(t) \cdot \mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t). \end{aligned}$$

Corollary 1 (Closed-form PF-ODE for an arbitrary Gaussian mixture in \mathbb{R}^d). Let

$$p(\mathbf{x}) = \sum_{j=1}^K w_j p_j(\mathbf{x}; \mathbf{m}_j, \Sigma_j), \quad w_j > 0, \quad \sum_j w_j = 1,$$

be a Gaussian-mixture density on \mathbb{R}^d , where $p_j(\mathbf{x})$ is the density of the j -th component, and \mathbf{m}_j is the mean and Σ_j is the covariance matrix of the j -th component. In addition, let α, γ satisfy the hypotheses of [Lem. 3](#), and define the forward map

$$\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}, \quad \mathbf{x} \sim p(\mathbf{x}), \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}).$$

For each component j set

$$\boldsymbol{\mu}_j(t) = \gamma(t)\mathbf{m}_j, \quad \Sigma_j(t) = \gamma(t)^2 \Sigma_j + \alpha(t)^2 \mathbf{I}, \quad \phi_j(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_j(t), \Sigma_j(t))$$

so that

$$p_t(\mathbf{x}_t) = \sum_{j=1}^K w_j \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_j(t), \Sigma_j(t)).$$

Then the Probability-Flow ODE (7) admits the closed-form drift

$$\frac{d\mathbf{x}_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t + \left[\alpha(t)\alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \sum_{j=1}^K \frac{w_j \phi_j(\mathbf{x}_t)}{p_t(\mathbf{x}_t)} \Sigma_j(t)^{-1} (\mathbf{x}_t - \boldsymbol{\mu}_j(t)).$$

Proof. Step 1. *Affine transform of a Gaussian mixture.* Conditioned on the j -th component, $\mathbf{x} \sim \mathcal{N}(\mathbf{m}_j, \Sigma_j)$, and hence

$$\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x} \mid (j) \sim \mathcal{N}(\gamma(t)\mathbf{m}_j, \alpha(t)^2 \mathbf{I} + \gamma(t)^2 \Sigma_j).$$

Defining

$$\boldsymbol{\mu}_j(t) = \gamma(t)\mathbf{m}_j, \quad \Sigma_j(t) = \gamma(t)^2 \Sigma_j + \alpha(t)^2 \mathbf{I},$$

we conclude that the marginal of \mathbf{x}_t is

$$\begin{aligned} p_t(\mathbf{x}_t) &= \sum_{j=1}^K p_t(\mathbf{x}_t, N = j) \\ &= \sum_{j=1}^K p(N = j) p_t(\mathbf{x}_t | N = j) \\ &= \sum_{j=1}^K w_j p_t(\alpha\mathbf{z} + \gamma\mathbf{x} | N = j) \\ &= \sum_{j=1}^K w_j \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_j(t), \Sigma_j(t)). \end{aligned}$$

Step 2. *Score of the mixture.* Set

$$\phi_j(\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_j(t), \boldsymbol{\Sigma}_j(t)), \quad p_t(\mathbf{x}_t) = \sum_{j=1}^K w_j \phi_j(\mathbf{x}_t).$$

Then by the usual mixture-rule,

$$\nabla_{\mathbf{x}_t} \log p_t = \frac{1}{p_t(\mathbf{x}_t)} \sum_{j=1}^K w_j \phi_j(\mathbf{x}_t) \nabla_{\mathbf{x}_t} \log \phi_j(\mathbf{x}_t).$$

Since for each Gaussian component

$$\nabla_{\mathbf{x}_t} \log \phi_j(\mathbf{x}_t) = -\boldsymbol{\Sigma}_j(t)^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_j(t)),$$

we obtain the closed-form score

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t) = -\frac{1}{p_t(\mathbf{x}_t)} \sum_{j=1}^K w_j \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_j(t), \boldsymbol{\Sigma}_j(t)) \boldsymbol{\Sigma}_j(t)^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_j(t)).$$

Step 3. *Substitution into the PF-ODE.* By [Lem. 3](#), the Probability-Flow ODE reads

$$\frac{d\mathbf{x}_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t - \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t).$$

Substituting the expression for $\nabla \log p_t$ above (and observing that the two '-' signs cancel) yields

$$\frac{d\mathbf{x}_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t + \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \sum_{j=1}^K \frac{w_j \mathcal{N}(\mathbf{x}_t; \boldsymbol{\mu}_j(t), \boldsymbol{\Sigma}_j(t))}{p_t(\mathbf{x}_t)} \boldsymbol{\Sigma}_j(t)^{-1}(\mathbf{x}_t - \boldsymbol{\mu}_j(t)),$$

which is exactly the claimed closed-form drift. \square

Corollary 2 (Closed-form PF-ODE for a symmetric two-peak Gaussian mixture). *Let $p(x)$ be the one-dimensional, symmetric, two-peak Gaussian mixture*

$$p(x) = \frac{1}{2} \mathcal{N}(x; -m, \sigma^2) + \frac{1}{2} \mathcal{N}(x; +m, \sigma^2),$$

and let α, γ be as in [Lem. 3](#). Define

$$x_t = \alpha(t) z + \gamma(t) x, \quad \Sigma_t = \gamma(t)^2 \sigma^2 + \alpha(t)^2, \quad \mu_{\pm}(t) = \pm \gamma(t) m.$$

Then the marginal density of x_t is

$$p_t(x_t) = \frac{1}{2} \mathcal{N}(x_t; \mu_-(t), \Sigma_t) + \frac{1}{2} \mathcal{N}(x_t; \mu_+(t), \Sigma_t),$$

and the Probability-Flow ODE (7) admits the closed-form drift

$$\frac{dx_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} x_t + \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \frac{1}{\Sigma_t} \left[x_t - \gamma(t) m \tanh\left(\frac{\gamma(t) m}{\Sigma_t} x_t\right) \right].$$

Proof. Step 1. *Marginal law under the affine map.* Conditional on $x = \pm m$, one has

$$x_t = \alpha z + \gamma x \mid (x = \pm m) \sim \mathcal{N}(\pm \gamma m, \alpha^2 + \gamma^2 \sigma^2) = \mathcal{N}(\mu_{\pm}(t), \Sigma_t).$$

Since each peak has weight $\frac{1}{2}$, the marginal of x_t is $\frac{1}{2} \mathcal{N}(\mu_-, \Sigma_t) + \frac{1}{2} \mathcal{N}(\mu_+, \Sigma_t)$.

Step 2. *Score of the bimodal mixture.* Write $\phi_{\pm}(x_t) = \mathcal{N}(x_t; \mu_{\pm}(t), \Sigma_t)$, so $p_t = \frac{1}{2}(\phi_- + \phi_+)$. Then

$$\frac{d}{dx_t} \log p_t = \frac{1}{p_t} \frac{1}{2} (\phi_- \nabla \log \phi_- + \phi_+ \nabla \log \phi_+), \quad \nabla \log \phi_{\pm} = -\frac{x_t - \mu_{\pm}(t)}{\Sigma_t}.$$

Hence

$$\frac{d}{dx_t} \log p_t = -\frac{1}{2 p_t \Sigma_t} [\phi_-(x_t - \mu_-) + \phi_+(x_t - \mu_+)].$$

1998 Define

$$1999 \quad r_{\pm}(x_t) = \frac{\phi_{\pm}(x_t)}{\phi_{-}(x_t) + \phi_{+}(x_t)}, \quad \phi_{-} + \phi_{+} = 2p_t.$$

2001 Then

$$2002 \quad \frac{d}{dx_t} \log p_t = -\frac{1}{\Sigma_t} [r_{-}(x_t - \mu_{-}) + r_{+}(x_t - \mu_{+})].$$

2004 A direct computation shows

$$2005 \quad r_{+} - r_{-} = \tanh\left(\frac{\gamma m}{\Sigma_t} x_t\right), \quad r_{-}(x_t + \gamma m) + r_{+}(x_t - \gamma m) = x_t - \gamma m \tanh\left(\frac{\gamma m}{\Sigma_t} x_t\right).$$

2008 Therefore

$$2009 \quad \frac{d}{dx_t} \log p_t = -\frac{1}{\Sigma_t} \left[x_t - \gamma m \tanh\left(\frac{\gamma m}{\Sigma_t} x_t\right) \right].$$

2011 Step 3. *Substitution into the PF-ODE.* By [Lem. 3](#),

$$2012 \quad \frac{dx_t}{dt} = \frac{\gamma'}{\gamma} x_t - \left[\alpha \alpha' - \frac{\gamma'}{\gamma} \alpha^2 \right] \frac{d}{dx_t} \log p_t.$$

2015 Since $\frac{d}{dx_t} \log p_t$ carries a “-” sign, the two negatives cancel, yielding exactly

$$2016 \quad \frac{dx_t}{dt} = \frac{\gamma'}{\gamma} x_t + \left[\alpha \alpha' - \frac{\gamma'}{\gamma} \alpha^2 \right] \frac{1}{\Sigma_t} \left[x_t - \gamma m \tanh\left(\frac{\gamma m}{\Sigma_t} x_t\right) \right],$$

2019 as claimed. □

2021 **Remark 3 (OU-type schedule for the symmetric bimodal case)**. *Specialize [Cor. 2](#) to the Ornstein-Uhlenbeck-type schedule with*

$$2022 \quad \gamma(t) = e^{-st}, \quad \alpha(t) = \sqrt{1 - e^{-2st}},$$

2023 *and noise variance σ^2 in each mixture component. Then the marginal variance is*

$$2024 \quad \Sigma_t = \gamma(t)^2 \sigma^2 + \alpha(t)^2 = \sigma^2 e^{-2st} + (1 - e^{-2st}),$$

2025 *and one obtains the closed-form drift of the Probability-Flow ODE:*

$$2026 \quad \boxed{\frac{dx_t}{dt} = -s x_t + \frac{s}{\Sigma_t} \left[x_t - m e^{-st} \tanh\left(\frac{m e^{-st}}{\Sigma_t} x_t\right) \right].}$$

2033 *Proof.* We start from the general drift in [Cor. 2](#):

$$2034 \quad \frac{dx_t}{dt} = \frac{\gamma'}{\gamma} x_t + \left[\alpha \alpha' - \frac{\gamma'}{\gamma} \alpha^2 \right] \frac{1}{\Sigma_t} \left[x_t - \gamma m \tanh\left(\frac{\gamma m}{\Sigma_t} x_t\right) \right].$$

2038 We now substitute $\gamma(t) = e^{-st}$, $\alpha(t) = \sqrt{1 - e^{-2st}}$ and compute each piece in detail:

2039 Derivative of γ :

$$2040 \quad \gamma'(t) = -s e^{-st}, \quad \implies \quad \frac{\gamma'(t)}{\gamma(t)} = -s.$$

2042 Marginal variance Σ_t :

$$2043 \quad \Sigma_t = \gamma(t)^2 \sigma^2 + \alpha(t)^2 = \sigma^2 e^{-2st} + (1 - e^{-2st}).$$

2044 Square of α and its derivative:

$$2045 \quad \alpha(t)^2 = 1 - e^{-2st}, \quad \frac{d}{dt} [\alpha(t)^2] = 2s e^{-2st} \implies 2\alpha \alpha' = 2s e^{-2st} \implies \alpha \alpha' = s e^{-2st}.$$

2049 Combination term

$$2050 \quad \alpha \alpha' - \frac{\gamma'}{\gamma} \alpha^2 = s e^{-2st} - (-s)(1 - e^{-2st}) = s [e^{-2st} + 1 - e^{-2st}] = s.$$

Substitution into the general drift formula gives

$$\frac{dx_t}{dt} = -s x_t + s \frac{1}{\Sigma_t} \left[x_t - e^{-st} m \tanh\left(\frac{e^{-st} m}{\Sigma_t} x_t\right) \right].$$

Hence the final, closed-form Probability-Flow ODE is

$$\frac{dx_t}{dt} = -s x_t + \frac{s}{\Sigma_t} \left[x_t - m e^{-st} \tanh\left(\frac{m e^{-st}}{\Sigma_t} x_t\right) \right],$$

where $\Sigma_t = \sigma^2 e^{-2st} + (1 - e^{-2st})$. \square

Remark 4 (Triangular schedule for the symmetric bimodal case). *Specialize Cor. 2 to the trigonometric schedule*

$$\gamma(t) = \cos\left(\frac{\pi}{2} t\right), \quad \alpha(t) = \sin\left(\frac{\pi}{2} t\right),$$

with noise variance σ^2 in each mixture component. Then

$$\Sigma_t = \gamma(t)^2 \sigma^2 + \alpha(t)^2 = \sigma^2 \cos^2\left(\frac{\pi}{2} t\right) + \sin^2\left(\frac{\pi}{2} t\right),$$

and the closed-form drift of the Probability-Flow ODE is

$$\frac{dx_t}{dt} = -\frac{\pi}{2} \tan\left(\frac{\pi}{2} t\right) x_t + \frac{\frac{\pi}{2} \tan\left(\frac{\pi}{2} t\right)}{\Sigma_t} \left[x_t - \cos\left(\frac{\pi}{2} t\right) m \tanh\left(\frac{\cos\left(\frac{\pi}{2} t\right) m}{\Sigma_t} x_t\right) \right].$$

Proof. We begin with the general drift in Cor. 2:

$$\frac{dx_t}{dt} = \frac{\gamma'}{\gamma} x_t + \left[\alpha \alpha' - \frac{\gamma'}{\gamma} \alpha^2 \right] \frac{1}{\Sigma_t} \left[x_t - \gamma m \tanh\left(\frac{\gamma m}{\Sigma_t} x_t\right) \right].$$

For $\gamma(t) = \cos\left(\frac{\pi}{2} t\right)$, $\alpha(t) = \sin\left(\frac{\pi}{2} t\right)$,

$$\gamma'(t) = -\frac{\pi}{2} \sin\left(\frac{\pi}{2} t\right) = -\frac{\pi}{2} \alpha(t), \quad \frac{\gamma'}{\gamma} = -\frac{\pi}{2} \tan\left(\frac{\pi}{2} t\right).$$

And

$$\alpha'(t) = \frac{\pi}{2} \cos\left(\frac{\pi}{2} t\right) = \frac{\pi}{2} \gamma(t),$$

so that

$$\alpha \alpha' - \frac{\gamma'}{\gamma} \alpha^2 = \frac{\pi}{2} \alpha \gamma + \frac{\pi}{2} \frac{\alpha^3}{\gamma} = \frac{\pi}{2} \frac{\alpha}{\gamma} (\alpha^2 + \gamma^2) = \frac{\pi}{2} \tan\left(\frac{\pi}{2} t\right).$$

Substituting into the general formula immediately yields the boxed drift. \square

Remark 5 (Linear schedule for the symmetric bimodal case). *Specialize Cor. 2 to the "Linear" schedule*

$$\gamma(t) = 1 - t, \quad \alpha(t) = t, \quad t \in [0, 1].$$

Then the marginal variance is

$$\Sigma_t = \gamma(t)^2 \sigma^2 + \alpha(t)^2 = (1 - t)^2 \sigma^2 + t^2,$$

and one obtains the closed-form drift of the Probability-Flow ODE:

$$\frac{dx_t}{dt} = -\frac{x_t}{1-t} + \frac{t}{(1-t)\Sigma_t} \left[x_t - m(1-t) \tanh\left(\frac{m(1-t)}{\Sigma_t} x_t\right) \right].$$

Proof. We begin with the general drift formula from Cor. 2:

$$\frac{dx_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} x_t + \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \frac{1}{\Sigma_t} \left[x_t - \gamma(t) m \tanh\left(\frac{\gamma(t) m}{\Sigma_t} x_t\right) \right].$$

We substitute $\gamma(t) = 1 - t$ and $\alpha(t) = t$ and compute each piece:

2106 1. Derivative of γ :

$$2107 \quad \gamma'(t) = -1, \quad \implies \quad \frac{\gamma'(t)}{\gamma(t)} = -\frac{1}{1-t}.$$

2109 2. Marginal variance:

$$2110 \quad \Sigma_t = (1-t)^2 \sigma^2 + t^2.$$

2111 3. Square of α and its derivative:

$$2112 \quad \alpha(t)^2 = t^2, \quad \frac{d}{dt} [\alpha(t)^2] = 2t \implies 2\alpha\alpha' = 2t \implies \alpha(t)\alpha'(t) = t.$$

2114 4. Combination term:

$$2115 \quad \alpha\alpha' - \frac{\gamma'}{\gamma}\alpha^2 = t - \left(-\frac{1}{1-t}\right)t^2 = t + \frac{t^2}{1-t} = \frac{t}{1-t}.$$

2117 Substituting these into the general drift gives

$$2118 \quad \frac{dx_t}{dt} = -\frac{x_t}{1-t} + \frac{t}{(1-t)\Sigma_t} \left[x_t - m(1-t) \tanh\left(\frac{m(1-t)}{\Sigma_t} x_t\right) \right],$$

2119 which is the claimed closed-form Probability-Flow ODE. \square

2122 **Remark 6 (OU-type schedule for the Hermite–Gaussian $n = 1$ case).** Apply *Lem. 3* to the one-dimensional Hermite–Gaussian initial density

$$2123 \quad p_1(x) \propto x e^{-x^2/2}, \quad x > 0,$$

2124 and the OU-type schedule

$$2125 \quad \gamma(t) = e^{-st}, \quad \alpha(t) = \sqrt{1 - e^{-2st}}.$$

2126 Then the Probability–Flow ODE (7) reduces to the scalar form

$$2127 \quad \boxed{\frac{dx_t}{dt} = -\frac{s}{x_t}, \quad t \in [0, 1],}$$

2128 and integrating from $t = 1$ (with $x(1) = x_1$) to any $t \in [0, 1]$ yields the explicit solution

$$2129 \quad \boxed{x_t = \sqrt{x_1^2 + 2s(1-t)}.$$

2130 *Proof.* By *Lem. 3*, the drift of the Probability–Flow ODE is

$$2131 \quad \frac{dx_t}{dt} = \frac{\gamma'(t)}{\gamma(t)} x_t - \left[\alpha(t)\alpha'(t) - \frac{\gamma'(t)}{\gamma(t)}\alpha(t)^2 \right] \partial_{x_t} \ln p_t(x_t).$$

2132 Under $\gamma(t) = e^{-st}$ and $\alpha(t) = \sqrt{1 - e^{-2st}}$ one computes

$$2133 \quad \frac{\gamma'}{\gamma} = -s, \quad 2\alpha\alpha' = 2s e^{-2st} \implies \alpha\alpha' = s e^{-2st}, \quad -\frac{\gamma'}{\gamma}\alpha^2 = s(1 - e^{-2st}),$$

2134 hence

$$2135 \quad \alpha\alpha' - \frac{\gamma'}{\gamma}\alpha^2 = s e^{-2st} + s(1 - e^{-2st}) = s.$$

2136 Moreover, one checks that the marginal density remains $p_t(x) \propto x e^{-x^2/2}$, so $\partial_x \ln p_t(x) = \frac{1}{x} - x$. Therefore

$$2137 \quad \frac{dx_t}{dt} = -s x_t - s \left(\frac{1}{x_t} - x_t \right) = -\frac{s}{x_t}.$$

2138 Separating variables,

$$2139 \quad \frac{dx}{dt} = -\frac{s}{x} \implies \int_{x_1}^{x_t} x dx = -s \int_1^t ds \implies \frac{x_t^2 - x_1^2}{2} = -s(t-1),$$

2140 whence

$$2141 \quad x_t^2 = x_1^2 + 2s(1-t), \quad x_t = \sqrt{x_1^2 + 2s(1-t)},$$

2142 taking the positive root on $x > 0$. \square

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171

Lemma 6 (Picard–Lindelöf existence and uniqueness). *Let $v: \mathbb{R} \times [0, 1] \rightarrow \mathbb{R}$ be continuous in t and satisfy the uniform Lipschitz condition*

$$|v(x, t) - v(y, t)| \leq L|x - y|, \quad \forall x, y \in \mathbb{R}, t \in [0, 1],$$

for some constant $L < \infty$. Then for any $t_0 \in [0, 1]$ and any initial value $x(t_0) = x_0$, there exists $\delta > 0$ and a unique function

$$x \in C^1([t_0 - \delta, t_0 + \delta] \cap [0, 1])$$

solving the ODE

$$\frac{dx}{dt}(t) = v(x(t), t), \quad x(t_0) = x_0.$$

2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185

Proof. Fix $t_0 \in [0, 1]$ and $x_0 \in \mathbb{R}$. Choose $\delta > 0$ so small that $(t_0 - \delta, t_0 + \delta) \subset [0, 1]$ and $L\delta < 1$. Define the closed ball

$$B_R = \{x \in C([t_0 - \delta, t_0 + \delta], \mathbb{R}) : \|x - x_0\|_\infty \leq R\}$$

with $R > 0$ to be chosen. Consider the operator

$$(\Gamma x)(t) = x_0 + \int_{t_0}^t v(x(s), s) ds.$$

Since v is continuous on the compact set $B_R \times [t_0 - \delta, t_0 + \delta]$, it is bounded by some $M < \infty$. If we choose $R = M\delta$, then Γ maps B_R into itself:

$$\|\Gamma x - x_0\|_\infty \leq \sup_t \int_{t_0}^t |v(x(s), s)| ds \leq M\delta = R.$$

Moreover, for any $x, y \in B_R$ and any t in the interval,

$$|(\Gamma x)(t) - (\Gamma y)(t)| \leq \int_{t_0}^t |v(x(s), s) - v(y(s), s)| ds \leq L\delta \|x - y\|_\infty < \|x - y\|_\infty,$$

2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200

so Γ is a contraction. By the Banach fixed-point theorem, Γ has a unique fixed point in B_R , which is precisely the unique C^1 solution of the ODE on $[t_0 - \delta, t_0 + \delta] \cap [0, 1]$. \square

Lemma 7 (Gronwall’s inequality and no blow-up). *Let $x \in C^1([0, 1])$ satisfy*

$$|x'(t)| \leq K(1 + |x(t)|), \quad t \in [0, 1],$$

for some constant $K \geq 0$. Then

$$|x(t)| \leq (|x(1)| + 1) e^{K(1-t)} - 1, \quad \forall t \in [0, 1],$$

and in particular x does not blow up in finite time on $[0, 1]$.

2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

Proof. Define

$$y(t) = |x(t)| + 1 \geq 1.$$

Since $y(t)$ is Lipschitz, for almost every t we have

$$y'(t) = \frac{d}{dt}(|x(t)| + 1) = \text{sgn}(x(t)) x'(t),$$

and hence

$$y'(t) \geq -|x'(t)| \geq -K(1 + |x(t)|) = -K y(t).$$

Equivalently,

$$y'(t) + K y(t) \geq 0.$$

Multiply both sides by the integrating factor e^{Kt} :

$$\frac{d}{dt}(e^{Kt} y(t)) = e^{Kt}(y'(t) + K y(t)) \geq 0.$$

2214 Thus the function $t \mapsto e^{Kt}y(t)$ is non-decreasing on $[0, 1]$. For any $t \leq 1$ we then have

$$2215 e^{Kt}y(t) \leq e^{K \cdot 1}y(1) \implies y(t) \leq y(1) e^{K(1-t)} = (|x(1)| + 1) e^{K(1-t)}.$$

2216 Rewriting $y(t) = |x(t)| + 1$ gives

$$2217 |x(t)| \leq (|x(1)| + 1)e^{K(1-t)} - 1,$$

2218 as claimed. In particular $|x(t)| < \infty$ for all $t \in [0, 1]$, so no finite-time blow-up occurs. \square

2219 **Lemma 8 (Gaussian convolution preserves linear-growth bound)**. Let $p_0 \in C^1(\mathbb{R})$ be a
2220 probability density satisfying

$$2221 |\partial_x \log p_0(x)| \leq A + B|x|, \quad A, B < \infty, \quad \forall x \in \mathbb{R},$$

2222 and assume furthermore that $\|p_0\|_\infty = \sup_{x \in \mathbb{R}} p_0(x) \leq M < \infty$. For each $\sigma > 0$,
2223 define the Gaussian kernel $\phi_\sigma(u) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{u^2}{2\sigma^2}\right)$, and set $p_\sigma(x) = (p_0 * \phi_\sigma)(x) =$
2224 $\int_{\mathbb{R}} p_0(y) \phi_\sigma(x - y) dy$. Then $p_\sigma \in C^\infty(\mathbb{R})$ and there exist

$$2225 A(\sigma) = A + B M \sigma \sqrt{\frac{2}{\pi}}, \quad B(\sigma) = B,$$

2226 such that

$$2227 |\partial_x \log p_\sigma(x)| \leq A(\sigma) + B(\sigma)|x|, \quad \forall x \in \mathbb{R}.$$

2228 *Proof. Smoothness and differentiation under the integral.* Since $\phi_\sigma \in C^\infty(\mathbb{R})$ decays rapidly and
2229 $p_0 \in L^\infty(\mathbb{R})$, by dominated convergence we may differentiate under the integral to get

$$2230 p'_\sigma(x) = \int_{\mathbb{R}} p_0(y) \partial_x \phi_\sigma(x - y) dy = \int_{\mathbb{R}} p_0(y) \phi'_\sigma(x - y) dy.$$

2231 Noting $\partial_y \phi_\sigma(x - y) = -\phi'_\sigma(x - y)$, we rewrite

$$2232 p'_\sigma(x) = - \int_{\mathbb{R}} p_0(y) \partial_y \phi_\sigma(x - y) dy.$$

2233 *Integration by parts.* Integrating the above in y and using that $p_0(y)\phi_\sigma(x - y) \rightarrow 0$ as $|y| \rightarrow \infty$, we
2234 obtain

$$2235 p'_\sigma(x) = \int_{\mathbb{R}} p'_0(y) \phi_\sigma(x - y) dy = \int_{\mathbb{R}} (\partial_y \log p_0)(y) p_0(y) \phi_\sigma(x - y) dy.$$

2236 *Bounding $\partial_x \log p_\sigma$.* Hence

$$2237 |\partial_x \log p_\sigma(x)| = \frac{|p'_\sigma(x)|}{p_\sigma(x)} = \frac{|\int (\partial_y \log p_0)(y) p_0(y) \phi_\sigma(x - y) dy|}{p_\sigma(x)}$$

$$2238 \leq \frac{\int |\partial_y \log p_0(y)| p_0(y) \phi_\sigma(x - y) dy}{p_\sigma(x)} \leq \frac{\int (A + B|y|) p_0(y) \phi_\sigma(x - y) dy}{p_\sigma(x)}$$

$$2239 = A + B \frac{\int |y| p_0(y) \phi_\sigma(x - y) dy}{p_\sigma(x)}.$$

2240 *Change of variables.* Set $u = y - x$. Then

$$2241 \int |y| p_0(y) \phi_\sigma(x - y) dy = \int |x + u| p_0(x + u) \phi_\sigma(u) du \leq |x| p_\sigma(x) + \int |u| p_0(x + u) \phi_\sigma(u) du.$$

2242 Hence

$$2243 \frac{\int |y| p_0(y) \phi_\sigma(x - y) dy}{p_\sigma(x)} \leq |x| + \frac{\int |u| p_0(x + u) \phi_\sigma(u) du}{p_\sigma(x)}.$$

2244 *Using the L^∞ -bound on p_0 .* Since $p_0(x + u) \leq M$,

$$2245 \int |u| p_0(x + u) \phi_\sigma(u) du \leq M \int |u| \phi_\sigma(u) du = M \sigma \sqrt{\frac{2}{\pi}}.$$

In particular v is globally Lipschitz in x (uniformly in t) and of linear growth. By the [Lem. 6](#) together with [Lem. 7](#) to prevent finite-time blow-up, the backward ODE admits for each x_1 a unique C^1 solution on $[0, 1]$.

(2) Conservation of the CDF. Let

$$F_t(x) = \int_{-\infty}^x p_t(u) \, du \quad (\text{the CDF of } p_t).$$

Since p_t satisfies the continuity equation $\partial_t p_t + \partial_x(v p_t) = 0$, along any characteristic $t \mapsto x_t$ one computes

$$\frac{d}{dt} F_t(x_t) = \int_{-\infty}^{x_t} \partial_t p_t(u) \, du + p_t(x_t) \frac{dx_t}{dt} = -[v p_t]_{-\infty}^{x_t} + p_t(x_t) v(x_t, t) = 0,$$

using $\lim_{u \rightarrow -\infty} p_t(u) = 0$. Hence $F_t(x_t) = F_1(x_1)$ for all $t \in [0, 1]$.

(3) Quantile representation. Evaluating at $t = 0$ gives

$$F_0(x_0(x_1)) = F_1(x_1).$$

Since $F_0: \mathbb{R} \rightarrow (0, 1)$ is strictly increasing and onto, it has an inverse F_0^{-1} , and thus

$$x_0(x_1) = F_0^{-1}(F_1(x_1)).$$

(4) Monotonicity and uniqueness. If $x_1 < y_1$ then $F_1(x_1) < F_1(y_1)$, so

$$g(x_1) = F_0^{-1}(F_1(x_1)) < F_0^{-1}(F_1(y_1)) = g(y_1),$$

showing g is strictly increasing. In one dimension the strictly increasing transport between two given laws is unique, so g is the unique increasing map pushing p_1 onto p_0 . A case study presented in [Fig. 10](#) validates this theorem, considering the specific schedules discussed in [Rem. 5](#) and [Rem. 3](#). \square

Lemma 9 (Monotone transport from Gaussian to P). *Let $Z \sim N(0, 1)$ be a standard normal random variable and let X be a random variable with distribution P on \mathbb{R} , having cumulative distribution function (CDF) F_P . Define*

$$\Phi(z) = \Pr[Z \leq z], \quad F_P^{-1}(u) = \inf\{x : F_P(x) \geq u\}, \quad u \in (0, 1).$$

Then there exists a non-decreasing continuous function $g(z) = F_P^{-1}(\Phi(z))$ such that $g(Z) \stackrel{d}{=} X$ if and only if P has no atoms (i.e. F_P is continuous). Moreover, if F_P is strictly increasing then g is unique.

Proof. Existence. Since $\Phi: \mathbb{R} \rightarrow (0, 1)$ is continuous and strictly increasing, the random variable

$$U = \Phi(Z)$$

is distributed uniformly on $(0, 1)$. Hence for any $x \in \mathbb{R}$,

$$\Pr(F_P^{-1}(U) \leq x) = \Pr(U \leq F_P(x)) = F_P(x),$$

so $F_P^{-1}(U)$ has distribution P . The quantile function F_P^{-1} is non-decreasing and, by standard results on generalized inverses (see e.g. Billingsley, *Probability and Measure*), is continuous on $(0, 1)$ if and only if F_P is continuous. Therefore

$$g(z) = F_P^{-1}(\Phi(z))$$

is non-decreasing and continuous exactly when F_P is continuous, and in that case $g(Z) \stackrel{d}{=} X$.

Necessity. Suppose P has an atom at x_0 , i.e. $\Pr[X = x_0] = p > 0$. If there were a continuous non-decreasing g with $g(Z) \stackrel{d}{=} X$, then to produce a point-mass p at x_0 it would have to be constant on a set of positive Pr-mass in the continuous law of Z . But continuity of g then forces it to be constant on a strictly larger interval, yielding a mass $> p$ at x_0 , a contradiction. Thus F_P must be continuous.

Uniqueness. Let g_1, g_2 be two continuous non-decreasing functions with $g_i(Z) \stackrel{d}{=} P$. Define for $u \in (0, 1)$

$$h_i(u) = g_i(\Phi^{-1}(u)), \quad i = 1, 2.$$

Each h_i is continuous, non-decreasing, and pushes $\text{Unif}(0, 1)$ onto P . When F_P is strictly increasing, its quantile F_P^{-1} is the unique such map (classical uniqueness of quantile functions for atomless laws). Hence $h_1 \equiv h_2 \equiv F_P^{-1}$ on $(0, 1)$, and therefore $g_1 \equiv g_2$ on \mathbb{R} . \square

F.1.3 LEARNING OBJECTIVE AS $\lambda \rightarrow 1$

Lemma 10 (L^p -estimate for the difference of two absolutely continuous functions). Let $I = [a, b]$ be a compact interval and $(E, \|\cdot\|)$ a Banach space. Suppose $f, g : I \rightarrow E$ are absolutely continuous with Bochner-integrable derivatives f', g' . Fix $1 \leq p \leq \infty$. Then

$$\|f - g\|_{L^p(I; E)} \leq (b - a)^{1/p} \|f(a) - g(a)\| + \int_a^b (b - s)^{1/p} \|f'(s) - g'(s)\| ds,$$

where for $p = \infty$ one interprets $(b - s)^{1/p} = 1$. Moreover, if $1 < p < \infty$ and p' denotes the conjugate exponent $1/p + 1/p' = 1$, then by Hölder's inequality one further deduces

$$\|f - g\|_{L^p(I; E)} \leq (b - a)^{1/p} \|f(a) - g(a)\| + \left(\frac{p-1}{p}\right)^{1/p'} (b - a) \|f' - g'\|_{L^p(I; E)}.$$

Proof. Since f and g are absolutely continuous on $[a, b]$, the Fundamental Theorem of Calculus in the Bochner setting gives, for each $t \in [a, b]$,

$$f(t) - g(t) = (f(a) - g(a)) + \int_a^t (f'(s) - g'(s)) ds.$$

Set $X(s) = f'(s) - g'(s)$. Then for every $t \in [a, b]$,

$$\|f(t) - g(t)\| \leq \|f(a) - g(a)\| + \left\| \int_a^t X(s) ds \right\|.$$

We now distinguish two cases.

Case 1: $1 \leq p < \infty$. Taking the L^p -norm in the variable t over $[a, b]$ and applying Minkowski's integral inequality for Bochner integrals yields

$$\begin{aligned} \|f - g\|_{L_t^p} &\leq \|f(a) - g(a)\| \|1\|_{L^p([a, b])} + \left\| \int_a^t X(s) ds \right\|_{L_t^p} \\ &= (b - a)^{1/p} \|f(a) - g(a)\| + \left(\int_a^b \left\| \int_a^t X(s) ds \right\|^p dt \right)^{1/p} \\ &\leq (b - a)^{1/p} \|f(a) - g(a)\| + \int_a^b \|1_{[s, b]}(\cdot) X(s)\|_{L_t^p} ds. \end{aligned}$$

Here we have written $\int_a^t X(s) ds = \int_a^b 1_{[a, t]}(s) X(s) ds$ and used the fact that

$$\|1_{[s, b]}(t)\|_{L_t^p} = \left(\int_a^b 1_{[s, b]}(t) dt \right)^{1/p} = (b - s)^{1/p}.$$

Hence

$$\|f - g\|_{L^p(I; E)} \leq (b - a)^{1/p} \|f(a) - g(a)\| + \int_a^b (b - s)^{1/p} \|X(s)\| ds,$$

which is the claimed L^p -estimate.

Case 2: $p = \infty$. Taking the essential supremum in $t \in [a, b]$ in the pointwise bound $\|f(t) - g(t)\| \leq \|f(a) - g(a)\| + \int_a^t \|X(s)\| ds$ gives immediately

$$\|f - g\|_{L^\infty(I; E)} \leq \|f(a) - g(a)\| + \int_a^b \|X(s)\| ds,$$

which agrees with the above formula when $(b - s)^{1/p} = 1$.

Refinement for $1 < p < \infty$. Let p' be the conjugate exponent, $1/p + 1/p' = 1$. Applying Hölder's inequality to the integral $\int_a^b (b - s)^{1/p} \|X(s)\| ds$ gives

$$\int_a^b (b - s)^{1/p} \|X(s)\| ds \leq \left(\int_a^b (b - s)^{p'/p} ds \right)^{1/p'} \left(\int_a^b \|X(s)\|^p ds \right)^{1/p}.$$

2430 Since $p'/p = 1/(p-1)$, a direct computation yields

$$2431 \int_a^b (b-s)^{p'/p} ds = \int_0^{b-a} u^{1/(p-1)} du = \frac{p-1}{p} (b-a)^{p'}. \quad 2432$$

2433 Hence

$$2434 \left(\int_a^b (b-s)^{p'/p} ds \right)^{1/p'} = \left(\frac{p-1}{p} \right)^{1/p'} (b-a), \quad 2435$$

2436 and we arrive at

$$2437 \int_a^b (b-s)^{1/p} \|X(s)\| ds \leq \left(\frac{p-1}{p} \right)^{1/p'} (b-a) \|X\|_{L^p(I;E)}. \quad 2438$$

2439 Combining this with the previous display completes the proof of the refined estimate. \square

2440 **Lemma 11 (Uniqueness of absolutely continuous functions)**. *Let $I = [a, b]$ be a compact interval and $(E, \|\cdot\|)$ a Banach space. Suppose $f, g : I \rightarrow E$ are absolutely continuous with Bochner-integrable derivatives f', g' . If*

$$2441 f(a) = g(a) \quad \text{and} \quad f'(t) = g'(t) \quad \text{for almost every } t \in I,$$

2442 then $f(t) = g(t)$ for all $t \in I$.

2443 *Proof.* Apply **Lem. 10** (the L^p -estimate for differences) in the case $p = \infty$. Since in this case one has

$$2444 (b-s)^{1/p} = 1, \quad \|f(a) - g(a)\| = 0, \quad \|f'(s) - g'(s)\| = 0 \text{ a.e.},$$

2445 the conclusion of **Lem. 10** reads

$$2446 \|f - g\|_{L^\infty(I;E)} \leq \|f(a) - g(a)\| + \int_a^b \|f'(s) - g'(s)\| ds = 0. \quad 2447$$

2448 Hence $\|f - g\|_{L^\infty(I;E)} = 0$, which means

$$2449 \sup_{t \in I} \|f(t) - g(t)\| = 0,$$

2450 so $f(t) = g(t)$ for every $t \in I$. \square

2451 **Theorem 4 (Pathwise consistency via zero total derivative)**. *Let $p(\mathbf{x})$ be a data distribution on \mathbb{R}^d , and let $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ be independent of \mathbf{x} . Let $\alpha, \gamma : [0, 1] \rightarrow \mathbb{R}$ be C^1 scalar functions satisfying*

$$2452 \alpha(0) = 0, \alpha(1) = 1, \quad \gamma(0) = 1, \gamma(1) = 0, \quad \gamma(t) \neq 0 \forall t \in (0, 1).$$

2453 Define the forward process

$$2454 \mathbf{x}_t = \alpha(t) \mathbf{z} + \gamma(t) \mathbf{x}, \quad t \in [0, 1],$$

2455 so that $\mathbf{x}_0 = \mathbf{x} \sim p(\mathbf{x})$ and $\mathbf{x}_1 = \mathbf{z} \sim \mathcal{N}(\mathbf{0}, I)$. Let p_t be the law of \mathbf{x}_t . By **Lem. 3** the corresponding Probability Flow ODE is

$$2456 \mathbf{v}(\mathbf{x}_t, t) = \frac{d}{dt} \mathbf{x}_t = \frac{\gamma'(t)}{\gamma(t)} \mathbf{x}_t - \left[\alpha(t) \alpha'(t) - \frac{\gamma'(t)}{\gamma(t)} \alpha(t)^2 \right] \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t). \quad 2457$$

2458 Given any point \mathbf{x}_t , define

$$2459 \mathbf{g}(\mathbf{x}_t, t) = \mathbf{x}_0 = \mathbf{x}_t + \int_t^0 \mathbf{v}(\mathbf{x}_u, u) du. \quad 2460$$

2461 Let $(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{x}) \otimes \mathcal{N}(\mathbf{0}, I)$ and $t \sim \text{Unif}[0, 1]$ be all mutually independent. Write $\mathbb{E}_{(\mathbf{z}, \mathbf{x})}$ for expectation over (\mathbf{z}, \mathbf{x}) and $\mathbb{E}_{(\mathbf{z}, \mathbf{x}), t}$ for expectation over (\mathbf{z}, \mathbf{x}) and t . Suppose

$$2462 \mathbb{E}_{(\mathbf{z}, \mathbf{x})} \|\mathbf{f}(\mathbf{x}_0, 0) - \mathbf{g}(\mathbf{x}_0, 0)\| = 0, \quad \mathbb{E}_{(\mathbf{z}, \mathbf{x}), t} \left\| \frac{d}{dt} \mathbf{f}(\mathbf{x}_t, t) \right\| = 0. \quad 2463$$

Then

$$\mathbb{E}_{(\mathbf{z}, \mathbf{x}, t)} \|\mathbf{f}(\mathbf{x}_t, t) - \mathbf{g}(\mathbf{x}_t, t)\| = 0.$$

Proof. Fix a draw (\mathbf{z}, \mathbf{x}) . Along its forward trajectory $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}$, define the two curves

$$f(t) = \mathbf{f}(\mathbf{x}_t, t), \quad g(t) = \mathbf{g}(\mathbf{x}_t, t).$$

We check the hypotheses of **Lem. 11** for $f, g : [0, 1] \rightarrow \mathbb{R}^d$.

Absolute continuity. Since \mathbf{f} is C^1 in (\mathbf{x}, t) and $t \mapsto \mathbf{x}_t$ is C^1 , the composition $f(t) = \mathbf{f}(\mathbf{x}_t, t)$ is absolutely continuous, with

$$f'(t) = \frac{d}{dt} \mathbf{f}(\mathbf{x}_t, t), \quad \text{existing a.e.}$$

Also

$$g(t) = \mathbf{x}_t + \int_t^0 \mathbf{v}(\mathbf{x}_u, u) du = \mathbf{x}_0 - \int_0^t \mathbf{v}(\mathbf{x}_u, u) du$$

is the sum of a C^1 function and an absolutely continuous integral, hence itself absolutely continuous.

Coincidence of initial values. From $\mathbb{E}_{(\mathbf{z}, \mathbf{x})} \|\mathbf{f}(\mathbf{x}_0, 0) - \mathbf{g}(\mathbf{x}_0, 0)\| = 0$ we get $\mathbf{f}(\mathbf{x}_0, 0) = \mathbf{g}(\mathbf{x}_0, 0)$ almost surely, so $f(0) = g(0)$ for almost every (\mathbf{z}, \mathbf{x}) .

Coincidence of derivatives a.e. By Tonelli–Fubini,

$$0 = \mathbb{E}_{(\mathbf{z}, \mathbf{x}, t)} \left\| \frac{d}{dt} \mathbf{f}(\mathbf{x}_t, t) \right\| = \int \left(\int_0^1 \left\| \frac{d}{dt} \mathbf{f}(\mathbf{x}_t, t) \right\| dt \right) d\mathbb{P}(\mathbf{z}, \mathbf{x}).$$

Hence for almost every (\mathbf{z}, \mathbf{x}) , $\int_0^1 \|\partial_t \mathbf{f}(\mathbf{x}_t, t)\| dt = 0$, which forces $\partial_t \mathbf{f}(\mathbf{x}_t, t) = 0$ for almost all t . Thus

$$f'(t) = 0 \quad \text{for a.e. } t \in [0, 1].$$

On the other hand

$$g'(t) = \frac{d\mathbf{x}_t}{dt} - \mathbf{v}(\mathbf{x}_t, t) = \mathbf{v}(\mathbf{x}_t, t) - \mathbf{v}(\mathbf{x}_t, t) = 0, \quad \forall t \in [0, 1].$$

Conclusion by uniqueness. We have shown f, g are absolutely continuous, $f(0) = g(0)$, and $f'(t) = g'(t)$ for almost every t . By **Lem. 11**, $f(t) = g(t)$ for all $t \in [0, 1]$ (almost surely in (\mathbf{z}, \mathbf{x})). Hence $\mathbf{f}(\mathbf{x}_t, t) = \mathbf{g}(\mathbf{x}_t, t)$ a.s., and taking expectation yields $\mathbb{E}_{(\mathbf{z}, \mathbf{x}, t)} \|\mathbf{f}(\mathbf{x}_t, t) - \mathbf{g}(\mathbf{x}_t, t)\| = 0$. \square

Remark 7 (Consistency-training loss). By **Thm. 4**, to enforce $\mathbf{f}(\mathbf{x}_t, t) \approx \mathbf{g}(\mathbf{x}_t, t) = \mathbf{x}_0$ along the PF–ODE flow, we suggest two equivalent training objectives:

1. Continuous PDE-residual loss

$$\mathcal{L}_{\text{PDE}} = \mathbb{E}_{t, \mathbf{x}_t} \left\| \partial_t \mathbf{f}(\mathbf{x}_t, t) + \mathbf{v}(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t} \mathbf{f}(\mathbf{x}_t, t) \right\|^2.$$

2. Finite-difference consistency loss

$$\mathcal{L}_{\text{cons}} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{z}} \left\| \mathbf{f}(\mathbf{x}_{t+\Delta t}, t + \Delta t) - \mathbf{f}(\mathbf{x}_t, t) \right\|^2,$$

where $\mathbf{x}_t = \alpha(t)\mathbf{z} + \gamma(t)\mathbf{x}_0$ and similarly for $\mathbf{x}_{t+\Delta t}$.

Proof. We begin from the requirement that $\mathbf{f}(\mathbf{x}_t, t)$ remain constant along the flow:

$$\begin{aligned} \frac{d}{dt} \mathbf{f}(\mathbf{x}_t, t) &= (\partial_t + \mathbf{v}(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t}) \mathbf{f}(\mathbf{x}_t, t) = \partial_t \mathbf{f}(\mathbf{x}_t, t) + \underbrace{\frac{d\mathbf{x}_t}{dt}}_{=\mathbf{v}(\mathbf{x}_t, t)} \cdot \nabla_{\mathbf{x}_t} \mathbf{f}(\mathbf{x}_t, t) = 0. \end{aligned}$$

This is exactly the linear transport PDE

$$(\partial_t + \mathbf{v} \cdot \nabla) \mathbf{f}(\mathbf{x}, t) = 0.$$

2538 To train a network f to satisfy it, one may minimize the L^2 -residual over the joint law of t and \mathbf{x}_t ,
 2539 yielding

$$2540 \mathcal{L}_{\text{PDE}} = \mathbb{E}_{t, \mathbf{x}_t} \left\| \partial_t \mathbf{f}(\mathbf{x}_t, t) + \mathbf{v}(\mathbf{x}_t, t) \cdot \nabla_{\mathbf{x}_t} \mathbf{f}(\mathbf{x}_t, t) \right\|^2.$$

2542 In practice, computing the spatial gradient $\nabla_{\mathbf{x}_t} \mathbf{f}$ can be expensive. Instead, we use a small time
 2543 increment Δt and the finite-difference approximation

$$2544 \mathbf{f}(\mathbf{x}_{t+\Delta t}, t + \Delta t) - \mathbf{f}(\mathbf{x}_t, t) \approx \Delta t [\partial_t \mathbf{f} + \mathbf{v} \cdot \nabla \mathbf{f}](\mathbf{x}_t, t).$$

2545 Squaring and taking expectations over $t, \mathbf{x}_0, \mathbf{z}$ then yields the discrete consistency loss

$$2547 \mathcal{L}_{\text{cons}} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{z}} \left\| \mathbf{f}(\mathbf{x}_{t+\Delta t}, t + \Delta t) - \mathbf{f}(\mathbf{x}_t, t) \right\|^2.$$

2548 This completes the derivation of both forms of the consistency-training objective. \square

2550 Recall that $(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x})$ is a pair of latent and data variables (typically independent), and let
 2551 $t \in [0, 1]$. We have four differentiable scalar functions $\alpha, \gamma, \hat{\alpha}, \hat{\gamma}: [0, 1] \rightarrow \mathbb{R}$, the *noisy interpolant*
 2552 $\mathbf{x}_t = \alpha(t) \mathbf{z} + \gamma(t) \mathbf{x}$ and $\mathbf{F}_t = \mathbf{F}_\theta(\mathbf{x}_t, t)$. We define the \mathbf{x} - and \mathbf{z} -prediction functions by

$$2553 \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\alpha(t) \mathbf{F}_t - \hat{\alpha}(t) \mathbf{x}_t}{\alpha(t) \hat{\gamma}(t) - \hat{\alpha}(t) \gamma(t)}, \quad \text{and} \quad \mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\hat{\gamma}(t) \mathbf{x}_t - \gamma(t) \mathbf{F}_t}{\alpha(t) \hat{\gamma}(t) - \hat{\alpha}(t) \gamma(t)}.$$

2556 Since

$$\begin{aligned} 2557 \mathbf{f}^{\mathbf{x}}(\mathbf{F}_0, \mathbf{x}_0, 0) &= \frac{\alpha(0) \cdot \mathbf{F}_\theta(\mathbf{x}_0, 0) - \hat{\alpha}(0) \cdot \mathbf{x}_0}{\alpha(0) \cdot \hat{\gamma}(0) - \hat{\alpha}(0) \cdot \gamma(0)} \\ 2558 &= \frac{0 \cdot \mathbf{F}_\theta(\mathbf{x}_0, 0) - \hat{\alpha}(0) \cdot \mathbf{x}_0}{0 \cdot \hat{\gamma}(0) - \hat{\alpha}(0) \cdot 1} \\ 2559 &= \frac{\mathbf{0} - \hat{\alpha}(0) \cdot \mathbf{x}_0}{0 - \hat{\alpha}(0)} \\ 2560 &= \mathbf{x}_0 \end{aligned}$$

2565 $\mathbf{f}^{\mathbf{x}}$ satisfies the boundary condition of consistency models (Song et al., 2023) and Thm. 4. To
 2566 better understand the unified loss, let's analyze a bit further. For simplicity we use the notation
 2567 $\mathbf{f}_\theta(\mathbf{x}_t, t) := \mathbf{f}^{\mathbf{x}}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t)$, the training objective is then equal to

$$2568 \mathcal{L}(\theta) = \mathbb{E}_{t, (\mathbf{z}, \mathbf{x})} \left[\frac{1}{\hat{\omega}(t)} \left\| \mathbf{f}_\theta(\mathbf{x}_t, t) - \mathbf{f}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) \right\|_2^2 \right].$$

2570 Let $\phi_t(\mathbf{x})$ be the solution of the PF-ODE determined by the velocity field $\mathbf{v}^*(\mathbf{x}_t, t) =$
 2571 $\mathbb{E}_{(\mathbf{z}, \mathbf{x}) | \mathbf{x}_t} [\mathbf{v}^{(\mathbf{z}, \mathbf{x})}(\mathbf{x}_t, t) | \mathbf{x}_t]$ (where $\mathbf{v}^{(\mathbf{z}, \mathbf{x})}(\mathbf{y}, t) = \alpha'(t) \mathbf{z} + \gamma'(t) \mathbf{x}$) and an initial value \mathbf{x} at time
 2572 $t = 0$. Define $\mathbf{g}_\theta(\mathbf{x}, t) := \mathbf{f}_\theta(\phi_t(\mathbf{x}), t)$ that moves along the solution trajectory. When $\lambda \rightarrow 1$, the
 2573 gradient of the loss tends to

$$\begin{aligned} 2574 \lim_{\lambda \rightarrow 1} \nabla_\theta \frac{\mathcal{L}(\theta)}{2(1-\lambda)} &= \mathbb{E}_t \left[\frac{t}{\hat{\omega}(t)} \cdot \mathbb{E}_{(\mathbf{z}, \mathbf{x})} \lim_{\lambda \rightarrow 1} \left\langle \frac{\mathbf{f}_\theta(\mathbf{x}_t, t) - \mathbf{f}_\theta(\mathbf{x}_{\lambda t}, \lambda t)}{t - \lambda t}, \nabla_\theta \mathbf{f}_\theta(\mathbf{x}_t, t) \right\rangle \right] \\ 2575 &= \mathbb{E}_t \left[\frac{t}{\hat{\omega}(t)} \cdot \mathbb{E}_{(\mathbf{z}, \mathbf{x})} \left\langle \frac{d\mathbf{f}_\theta(\mathbf{x}_t, t)}{dt}, \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \right\rangle \right] \end{aligned}$$

2579 The inner expectation can be computed as:

$$\begin{aligned} 2580 &\mathbb{E}_{(\mathbf{z}, \mathbf{x}), \mathbf{x}_t} \left\langle \frac{d\mathbf{f}_\theta(\mathbf{x}_t, t)}{dt}, \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \right\rangle \\ 2581 &= \mathbb{E}_{(\mathbf{z}, \mathbf{x}), \mathbf{x}_t} \langle \partial_1 \mathbf{f}_\theta(\mathbf{x}_t, t) \cdot \mathbf{v}^{(\mathbf{z}, \mathbf{x})}(\mathbf{x}_t, t) + \partial_2 \mathbf{f}_\theta(\mathbf{x}_t, t), \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \rangle \\ 2582 &= \mathbb{E}_{(\mathbf{z}, \mathbf{x}), \mathbf{x}_t} \langle \partial_1 \mathbf{f}_\theta(\mathbf{x}_t, t) \cdot (\alpha'(t) \mathbf{z} + \gamma'(t) \mathbf{x}) + \partial_2 \mathbf{f}_\theta(\mathbf{x}_t, t), \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \rangle \\ 2583 &= \mathbb{E}_{\mathbf{x}_t} \left[\mathbb{E}_{(\mathbf{z}, \mathbf{x}) | \mathbf{x}_t} \langle \partial_1 \mathbf{f}_\theta(\mathbf{x}_t, t) \cdot (\alpha'(t) \mathbf{z} + \gamma'(t) \mathbf{x}) + \partial_2 \mathbf{f}_\theta(\mathbf{x}_t, t), \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \rangle \right] \\ 2584 &= \mathbb{E}_{\mathbf{x}_t} \langle \partial_1 \mathbf{f}_\theta(\mathbf{x}_t, t) \cdot \mathbb{E}_{(\mathbf{z}, \mathbf{x}) | \mathbf{x}_t} [\alpha'(t) \mathbf{z} + \gamma'(t) \mathbf{x} | \mathbf{x}_t] + \partial_2 \mathbf{f}_\theta(\mathbf{x}_t, t), \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \rangle \\ 2585 &= \mathbb{E}_{\mathbf{x}_t} \langle \partial_1 \mathbf{f}_\theta(\mathbf{x}_t, t) \cdot \mathbf{v}^*(\mathbf{x}_t, t) + \partial_2 \mathbf{f}_\theta(\mathbf{x}_t, t), \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \rangle \\ 2586 &= \mathbb{E}_{\mathbf{x}_t} \langle \partial_2 \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t), \nabla_\theta \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \rangle \\ 2587 &= \nabla_\theta \mathbb{E}_{\phi_t^{-1}(\mathbf{x}_t)} \frac{1}{2} \left\| \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) - \mathbf{g}_{\theta^-}(\phi_t^{-1}(\mathbf{x}_t), t) + \partial_2 \mathbf{g}_\theta(\phi_t^{-1}(\mathbf{x}_t), t) \right\|_2^2 \end{aligned}$$

Thus from the perspective of gradient, when $\lambda \rightarrow 1$ the training objective is equivalent to

$$\mathbb{E}_{\phi_t^{-1}(\mathbf{x}_t), t} \left[\frac{t}{\hat{\omega}(t)} \cdot \|\mathbf{g}_{\theta}(\phi_t^{-1}(\mathbf{x}_t), t) - \mathbf{g}_{\theta^-}(\phi_t^{-1}(\mathbf{x}_t), t) + \partial_2 \mathbf{g}_{\theta}(\phi_t^{-1}(\mathbf{x}_t), t)\|_2^2 \right]$$

which naturally leads to the solution $\mathbf{g}_{\theta}(\mathbf{x}, t) = \mathbf{x}$ (since $\mathbf{g}_{\theta}(\mathbf{x}, 0) \equiv \mathbf{x}$), or equivalently $\mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^*}(\mathbf{x}_t, t), \mathbf{x}_t, t) = \mathbf{f}_{\theta^*}(\mathbf{x}_t, t) = \phi_t^{-1}(\mathbf{x}_t)$, that is the definition of consistency function.

F.1.4 ANALYSIS ON THE OPTIMAL SOLUTION FOR $\lambda \in [0, 1]$

Below we provide some examples to illustrate the property of the optimal solution for the unified loss by considering some simple cases of data distribution.

(for simplicity define $\mathbf{f}_{\theta}(\mathbf{x}_t, t) = \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t)$)

Assume $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \Sigma)$. For $r < t$ the conditional mean

$$\mathbb{E}[\mathbf{x}_r | \mathbf{x}_t] = \gamma(r)\boldsymbol{\mu} + (\gamma(r)\gamma(t)\Sigma + \alpha(r)\alpha(t)\mathbf{I}) (\gamma(t)^2\Sigma + \alpha(t)^2\mathbf{I})^{-1} (\mathbf{x}_t - \gamma(t)\boldsymbol{\mu}),$$

denote

$$\mathbf{K}(r, t) := (\gamma(r)\gamma(t)\Sigma + \alpha(r)\alpha(t)\mathbf{I}) (\gamma(t)^2\Sigma + \alpha(t)^2\mathbf{I})^{-1},$$

using above equations we can get the optimal solution for diffusion model:

$$\mathbf{f}_{\theta^*}^{\text{DM}}(\mathbf{x}_t, t) = \mathbb{E}[\mathbf{x} | \mathbf{x}_t] = \boldsymbol{\mu} + \mathbf{K}(0, t)(\mathbf{x}_t - \gamma_t\boldsymbol{\mu}).$$

Now consider a series of t together: $t = t_T > t_{T-1} > \dots > t_1 > t_0 \approx 0$. This series could be obtained by $t_{j-1} = \lambda \cdot t_j, j = T, \dots, 0$, for instance. With an abuse of notation, denote \mathbf{x}_{t_j} as \mathbf{x}_j and $\alpha(t_j)$ as $\alpha_j, \gamma(t_j)$ as γ_j . Since $t_0 \approx 0, \mathbf{x}_0 \approx \mathbf{x}$, we could conclude the trained model $\mathbf{f}_{\theta^*}(\mathbf{x}_1, t_1) = \mathbb{E}_{\mathbf{x} | \mathbf{x}_1}[\mathbf{x} | \mathbf{x}_1]$, and consequently

$$\mathbf{f}_{\theta^*}(\mathbf{x}_{j+1}, t_{j+1}) = \mathbb{E}_{\mathbf{x}_j | \mathbf{x}_{j+1}}[\mathbf{f}_{\theta^*}(\mathbf{x}_j, t_j) | \mathbf{x}_{j+1}], j = 1, \dots, T-1.$$

Using the property of the conditional expectation, we have $\mathbb{E}_{\mathbf{x}_j}[\mathbf{f}_{\theta^*}(\mathbf{x}_j, t_j)] = \mathbb{E}_{\mathbf{x}}[\mathbf{x}], \forall j$. Using the expressions above we have

$$\mathbf{f}_{\theta^*}(\mathbf{x}_1, t_1) = \boldsymbol{\mu} + \mathbf{K}(t_0, t_1)(\mathbf{x}_1 - \gamma_1\boldsymbol{\mu})$$

and

$$\mathbf{f}_{\theta^*}(\mathbf{x}_j, t_j) = \boldsymbol{\mu} + \left[\prod_{k=1}^j \mathbf{K}(t_{k-1}, t_k) \right] \cdot (\mathbf{x}_t - \gamma_t\boldsymbol{\mu}), j = 2, \dots, T$$

Further denote $c_j = \prod_{k=1}^j \alpha_{k-1}\alpha_k + \gamma_{k-1}\gamma_k$ and assume $\Sigma = \mathbf{I}, \alpha = \sin(t), \gamma(t) = \cos(t)$. For appropriate choice of the partition scheme (e.g. even or geometric), the coefficient c_j can converge as T grows. For instance, when evenly partitioning the interval $[0, t]$, we have:

$$\lim_{T \rightarrow \infty} c(t) = \lim_{T \rightarrow \infty} \prod_{k=1}^T \alpha_{k-1}\alpha_k + \gamma_{k-1}\gamma_k = \lim_{T \rightarrow \infty} (\cos(\frac{t}{T}))^T = 1.$$

Thus the trained model can be viewed as an interpolant between the consistency model($\lambda \rightarrow 1$ or $T \rightarrow \infty$) and the diffusion model($\lambda \rightarrow 0$ or $T \rightarrow 1$):

$$\mathbf{f}_{\theta^*}(\mathbf{x}_t, t) = \boldsymbol{\mu} + c(t)(\mathbf{x}_t - \gamma(t)\boldsymbol{\mu}),$$

$$\mathbf{f}_{\theta^*}^{\text{CM}}(\mathbf{x}_t, t) = \boldsymbol{\mu} + (\mathbf{x}_t - \gamma(t)\boldsymbol{\mu}),$$

$$\mathbf{f}_{\theta^*}^{\text{DM}}(\mathbf{x}_t, t) = \boldsymbol{\mu} + \gamma(t)(\mathbf{x}_t - \gamma(t)\boldsymbol{\mu}).$$

The expression of $\mathbf{f}_{\theta^*}^{\text{CM}}$ can be obtained by first compute the velocity field $\mathbf{v}^*(\mathbf{x}_t, t) = \mathbb{E}[\alpha'(t)\mathbf{z} + \gamma'(t)\mathbf{x} | \mathbf{x}_t] = \gamma'(t)\boldsymbol{\mu}$ then solve the initial value problem of ODE to get $\mathbf{x}(0)$.

The above optimal solution can be possibly obtained by training. For example if we set the parameterization as $\mathbf{f}_{\theta}(\mathbf{x}_t, t) = (1 - \gamma_t c_t)\boldsymbol{\theta} + c_t \mathbf{x}_t$, the gradient of the loss can be computed as (let $r = \lambda \cdot t$):

$$\nabla_{\theta} \|\mathbf{f}_{\theta}(\mathbf{x}_t, t) - \mathbf{f}_{\theta^-}(\mathbf{x}_r, r)\|_2^2 = 2(1 - \gamma_t c_t) [(\alpha_t \gamma_t - \alpha_r \gamma_r)\mathbf{z} + (\gamma_r c_r - \gamma_t c_t)(\boldsymbol{\theta} - \mathbf{x})],$$

$$\nabla_{\theta} \mathbb{E}_{\mathbf{z}, \mathbf{x}} \|\mathbf{f}_{\theta}(\mathbf{x}_t, t) - \mathbf{f}_{\theta^-}(\mathbf{x}_r, r)\|_2^2 = 2(1 - \gamma_t c_t)(\gamma_r c_r - \gamma_t c_t)(\theta - \mu),$$

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta) &= \mathbb{E}_t \frac{2(1 - \gamma_t c_t)(\gamma_r c_r - \gamma_t c_t)}{\hat{\omega}(t)} (\theta - \mu) \\ &= C(\theta - \mu), \quad C = \mathbb{E}_t \frac{2(1 - \gamma_t c_t)(\gamma_r c_r - \gamma_t c_t)}{\hat{\omega}(t)}. \end{aligned}$$

Use gradient descent to update θ during training:

$$\frac{d\theta(s)}{ds} = -\nabla_{\theta} \mathcal{L}(\theta) = -C(\theta - \mu).$$

The generalization loss thus evolves as:

$$\begin{aligned} \frac{d\|\theta(s) - \mu\|^2}{ds} &= \langle \theta(s) - \mu, \frac{d\theta(s)}{ds} \rangle \\ &= \langle \theta(s) - \mu, -C(\theta(s) - \mu) \rangle \\ &= -C\|\theta(s) - \mu\|^2, \\ \implies \|\theta(s) - \mu\|^2 &= \|\theta(0) - \mu\|^2 e^{-Cs}. \end{aligned}$$

F.1.5 SURROGATE OBJECTIVE FOR UNIFIED OBJECTIVE

Proof for Thm. 1. For brevity, we omit the expectation operator \mathbb{E} in the following derivation.

Step 1. Omit the expectation operator.

$$l(\theta) = \frac{1}{\hat{\omega}(t)} \|\mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t)\|_2^2.$$

Step 2. Gradient of the loss.

$$\nabla_{\theta} l(\theta) = \frac{1}{\hat{\omega}(t)} \langle \nabla_{\theta} \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t), \Delta f \rangle, \quad (10)$$

where

$$\begin{aligned} \Delta f &= \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t) \\ &= [\mathbf{x}_t - t \cdot \mathbf{F}_{\theta}(\mathbf{x}_t, t)] - [\mathbf{x}_{\lambda t} - \lambda t \cdot \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t)] \\ &= (t - \lambda t)[(\mathbf{z} - \mathbf{x}) - \mathbf{F}_{\theta}(\mathbf{x}_t, t)] + \lambda t \cdot (\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) - \mathbf{F}_{\theta}(\mathbf{x}_t, t)). \end{aligned} \quad (11)$$

Also,

$$\nabla_{\theta} \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) = -t \cdot \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t). \quad (12)$$

Step 3. Substitute (11) and (12) into (10).

$$\begin{aligned} \nabla_{\theta} l(\theta) &= \frac{t(t - \lambda t)}{\hat{\omega}(t)} \langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{F}_{\theta}(\mathbf{x}_t, t) - (\mathbf{z} - \mathbf{x}) \rangle \\ &\quad + \frac{t\lambda t}{\hat{\omega}(t)} \langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) \rangle. \end{aligned}$$

Step 3. Use $\hat{\omega}(t) = t^2 \cdot (1 - \lambda)$.

$$\nabla_{\theta} l(\theta) = \nabla_{\theta} \|\mathbf{F}_{\theta}(\mathbf{x}_t, t) - (\mathbf{z} - \mathbf{x})\|_2^2 + \frac{\lambda}{1 - \lambda} \nabla_{\theta} \|\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t)\|_2^2.$$

This matches exactly the gradient of $\mathcal{G}(\theta)$. Hence,

$$\nabla_{\theta} \mathcal{L}(\theta) = \nabla_{\theta} \mathcal{G}(\theta).$$

□

Theorem 5 (Surrogate Loss for Unified Objective of General Case). *Define*

$$A(t) := \frac{\alpha(t)}{D(t)} - \frac{\alpha(\lambda t)}{D(\lambda t)}, \quad B(t) := \frac{\alpha(t)}{D(t)}, \quad C(t) := \frac{\alpha(t)}{2D(t)}, \quad D(t) := \alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t).$$

Let the surrogate loss be

$$\mathcal{G}(\theta) = \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\frac{C(t)}{\hat{\omega}(t)} \left(A(t) \underbrace{\|\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{z}_t\|_2^2}_{\text{Mean Velocity Alignment}} \right. \right. \\ \left. \left. + B(\lambda t) \underbrace{\|(\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t)) - (\mathbf{z}_t - \mathbf{z}_{\lambda t})\|_2^2}_{\text{Velocity Difference Consistency}} \right) \right] \quad (13)$$

Then, for all θ ,

$$\nabla_\theta \mathcal{L}(\theta) = \nabla_\theta \mathcal{G}(\theta).$$

Proof for Thm. 5. For brevity, we omit the expectation operator \mathbb{E} and weights $\hat{\omega}(t)$ in the following derivation.

$$l(\theta) = \|\mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^\mathbf{x}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t)\|_2^2 \quad (14)$$

$$\nabla_\theta l(\theta) = \langle \nabla_\theta \mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t), \Delta \mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) \rangle \quad (15)$$

In the following, we compute $\nabla_\theta \mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t)$ and $\Delta \mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t)$, respectively.

$$\begin{aligned} \nabla_\theta \mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) &= \nabla_\theta \left(\frac{\alpha(t) \cdot \mathbf{F}_\theta(\mathbf{x}_t, t) - \hat{\alpha}(t) \cdot \mathbf{x}_t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} \right) \\ &= \frac{\alpha(t)}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} \cdot \nabla_\theta \mathbf{F}_\theta(\mathbf{x}_t, t) \\ &= \frac{\alpha(t)}{D(t)} \cdot \nabla_\theta \mathbf{F}_\theta(\mathbf{x}_t, t) \end{aligned} \quad (16)$$

$$\begin{aligned} \Delta \mathbf{f}^\mathbf{x}(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) &= \frac{\alpha(t) \cdot \mathbf{F}_\theta(\mathbf{x}_t, t) - \hat{\alpha}(t) \cdot \mathbf{x}_t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} - \frac{\alpha(\lambda t) \cdot \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) - \hat{\alpha}(\lambda t) \cdot \mathbf{x}_{\lambda t}}{\alpha(\lambda t) \cdot \hat{\gamma}(\lambda t) - \hat{\alpha}(\lambda t) \cdot \gamma(\lambda t)} \\ &= \frac{\alpha(t)}{D(t)} \cdot \mathbf{F}_\theta(\mathbf{x}_t, t) - \frac{\hat{\alpha}(t)}{D(t)} \cdot \mathbf{x}_t - \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) + \frac{\hat{\alpha}(\lambda t)}{D(\lambda t)} \cdot \mathbf{x}_{\lambda t} \end{aligned} \quad (17)$$

Now, we consider to replace \mathbf{x}_t and $\mathbf{x}_{\lambda t}$ with \mathbf{z}_t and $\mathbf{z}_{\lambda t}$. Let's consider a general term (Remind $\mathbf{z}_t = \hat{\alpha}(t) \cdot \mathbf{z} + \hat{\gamma}(t) \cdot \mathbf{x}$):

$$\begin{aligned} \frac{\hat{\alpha}(s) \cdot \mathbf{x}_s}{D(s)} &= \frac{\hat{\alpha}(s) \cdot \alpha(s) \cdot \mathbf{z} + \hat{\alpha}(s) \cdot \gamma(s) \cdot \mathbf{x}}{\alpha(s) \cdot \hat{\gamma}(s) - \hat{\alpha}(s) \cdot \gamma(s)} \\ &= \frac{\alpha(s) \cdot \mathbf{z}_s - (\alpha(s) \cdot \hat{\gamma}(s) - \hat{\alpha}(s) \cdot \gamma(s)) \cdot \mathbf{x}}{\alpha(s) \cdot \hat{\gamma}(s) - \hat{\alpha}(s) \cdot \gamma(s)} \\ &= \frac{\alpha(s) \cdot \mathbf{z}_s}{\alpha(s) \cdot \hat{\gamma}(s) - \hat{\alpha}(s) \cdot \gamma(s)} - \mathbf{x} \\ &= \frac{\alpha(s) \cdot \mathbf{z}_s}{D(s)} - \mathbf{x} \end{aligned} \quad (18)$$

2754 Therefore, by substituting (18) into (17), we get:

$$\begin{aligned}
2755 & \Delta \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) \\
2756 &= \frac{\alpha(t)}{D(t)} \cdot \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \frac{\hat{\alpha}(t)}{D(t)} \cdot \mathbf{x}_t - \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t) + \frac{\hat{\alpha}(\lambda t)}{D(\lambda t)} \cdot \mathbf{x}_{\lambda t} \\
2757 &= \frac{\alpha(t)}{D(t)} \cdot \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \frac{\alpha(t)}{D(t)} \cdot \mathbf{z}_t - \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t) + \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot \mathbf{z}_{\lambda t} \\
2758 &= \left(\frac{\alpha(t)}{D(t)} - \frac{\alpha(\lambda t)}{D(\lambda t)} \right) \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t) + \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)) + \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot (\mathbf{z}_{\lambda t} - \mathbf{z}_t) \\
2759 & \tag{19}
\end{aligned}$$

2765 By substituting (19) and (16) into (15), we get:

$$\begin{aligned}
2766 & \nabla_{\theta} l(\theta) = \langle \nabla_{\theta} \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t), \Delta \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) \rangle \\
2767 &= \langle \frac{\alpha(t)}{D(t)} \cdot \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \left(\frac{\alpha(t)}{D(t)} - \frac{\alpha(\lambda t)}{D(\lambda t)} \right) \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t) \rangle \\
2768 &+ \langle \frac{\alpha(t)}{D(t)} \cdot \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)) \rangle \\
2769 &+ \langle \frac{\alpha(t)}{D(t)} \cdot \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot (\mathbf{z}_{\lambda t} - \mathbf{z}_t) \rangle \\
2770 &= \frac{\alpha(t)}{2D(t)} \cdot \left(\frac{\alpha(t)}{D(t)} - \frac{\alpha(\lambda t)}{D(\lambda t)} \right) \cdot \nabla_{\theta} \|\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t\|_2^2 \\
2771 &+ \frac{\alpha(t)}{2D(t)} \cdot \frac{\alpha(\lambda t)}{D(\lambda t)} \cdot \nabla_{\theta} \|(\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)) - (\mathbf{z}_t - \mathbf{z}_{\lambda t})\|_2^2 \\
2772 & \\
2773 & \\
2774 & \\
2775 & \\
2776 & \\
2777 & \\
2778 & \\
2779 & \\
2780 & \square
\end{aligned}$$

2781
2782
2783
2784
2785
2786
2787
2788
2789
2790
2791
2792
2793
2794
2795
2796
2797
2798
2799
2800
2801
2802
2803
2804
2805
2806
2807

F.1.6 CLOSED-FORM SOLUTION ANALYSIS FOR $\lambda \in [0, 1]$

Theorem 6 (optimal solution of surrogate objective in linear case ($\alpha(t) = t$, $\gamma(t) = 1 - t$)).
Under *Assump. 1*, let's consider a surrogate objective

$$\mathcal{G}(\theta) = \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\left\| \mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{z}_t \right\|_2^2 + \frac{\lambda}{1 - \lambda} \cdot \left\| \mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) \right\|_2^2 \right], \quad (20)$$

where $\mathbf{x}_t = t \cdot \mathbf{z} + (1 - t) \cdot \mathbf{x}$, $\mathbf{z}_t = \mathbf{z} - \mathbf{x}$, $0 < \lambda < 1$. Then the optimal solution of the surrogate objective is:

$$\mathbf{F}_{\theta^*}(\mathbf{x}_t, t) = \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\mathbf{z}_t + \frac{\lambda}{1 - \lambda} \cdot \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) \mid \mathbf{x}_t \right] \quad (21)$$

Proof. For any fixed t and \mathbf{x}_t , the objective only depends on the value $\mathbf{F}_\theta(\mathbf{x}_t, t)$ at this specific input. Therefore, we may optimize pointwise over $\mathbf{F}_\theta(\mathbf{x}_t, t)$.

Define for shorthand:

$$\mathbf{f} := \mathbf{F}_\theta(\mathbf{x}_t, t), \quad \mathbf{a} := \mathbf{z}_t, \quad \mathbf{b} := \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \quad \mathbf{f}^* := \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\mathbf{a} + \frac{\lambda}{1 - \lambda} \cdot \mathbf{b} \mid \mathbf{x}_t \right].$$

Then the surrogate objective specialized at (\mathbf{x}_t, t) is

$$\begin{aligned} \mathcal{G}_t(\mathbf{f}) &= \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{a} \right\|_2^2 + \frac{\lambda}{1 - \lambda} \left\| \mathbf{f} - \mathbf{b} \right\|_2^2 \mid \mathbf{x}_t \right] \\ &= \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{f}^* + \mathbf{f}^* - \mathbf{a} \right\|_2^2 + \frac{\lambda}{1 - \lambda} \left\| \mathbf{f} - \mathbf{f}^* + \mathbf{f}^* - \mathbf{b} \right\|_2^2 \mid \mathbf{x}_t \right] \\ &= \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{f}^* \right\|_2^2 + \left\| \mathbf{f}^* - \mathbf{a} \right\|_2^2 + 2 \langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{a} \rangle \mid \mathbf{x}_t \right] \\ &\quad + \frac{\lambda}{1 - \lambda} \cdot \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{f}^* \right\|_2^2 + \left\| \mathbf{f}^* - \mathbf{b} \right\|_2^2 + 2 \langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{b} \rangle \mid \mathbf{x}_t \right] \\ &\geq \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{f}^* \right\|_2^2 + \left\| \mathbf{f}^* - \mathbf{a} \right\|_2^2 \mid \mathbf{x}_t \right] + \frac{\lambda}{1 - \lambda} \cdot \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{f}^* \right\|_2^2 + \left\| \mathbf{f}^* - \mathbf{b} \right\|_2^2 \mid \mathbf{x}_t \right] \\ &= \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\left\| \mathbf{f} - \mathbf{f}^* \right\|_2^2 \mid \mathbf{x}_t \right] + \mathcal{G}_t(\mathbf{f}^*) \end{aligned}$$

In the following, we need to show that:

$$\mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{a} \rangle \mid \mathbf{x}_t \right] + \frac{\lambda}{1 - \lambda} \cdot \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{b} \rangle \mid \mathbf{x}_t \right] = 0 \quad (22)$$

The (22) always holds because:

$$\begin{aligned} &\mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{a} \rangle \mid \mathbf{x}_t \right] + \frac{\lambda}{1 - \lambda} \cdot \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{b} \rangle \mid \mathbf{x}_t \right] \\ &= \langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbb{E}_{\mathbf{z}, \mathbf{x}}[\mathbf{a} \mid \mathbf{x}_t] \rangle + \langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \frac{\lambda}{1 - \lambda} \cdot \mathbb{E}_{\mathbf{z}, \mathbf{x}}[\mathbf{b} \mid \mathbf{x}_t] \rangle \\ &= \langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbb{E}_{\mathbf{z}, \mathbf{x}}[\mathbf{a} + \frac{\lambda}{1 - \lambda} \cdot \mathbf{b} \mid \mathbf{x}_t] \rangle \\ &= 0 \end{aligned}$$

which completes the proof. In addition, the optimal solution of the surrogate objective also satisfies that the gradient of the surrogate objective is zero, i.e.

$$\nabla_{\theta} \mathcal{G}(\theta^*) = 0$$

□

Because it is hard to directly analyze the behavior of the optimal solution of the surrogate objective, we analyze the gradient of the surrogate objective instead.

Proposition 1 (Gradient of surrogate objective in linear case). *The gradient of the surrogate objective (20) is:*

$$\nabla_{\theta} \mathcal{G}(\theta) = 2\mathbb{E}[\langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t) + \frac{\lambda}{1-\lambda} \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)) \rangle] \quad (23)$$

Proof.

$$\nabla_{\theta} \mathcal{G}(\theta) = \nabla_{\theta} \mathbb{E} \left[\left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t \right\|_2^2 + \frac{\lambda}{1-\lambda} \cdot \left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t) \right\|_2^2 \right] \quad (24)$$

$$= 2\mathbb{E} \left[\langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t \rangle + \frac{\lambda}{1-\lambda} \cdot \langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t) \rangle \right] \quad (25)$$

$$= 2\mathbb{E} \left[\langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \frac{\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \lambda \cdot \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)}{1-\lambda} - \mathbf{z}_t \rangle \right] \quad (26)$$

$$= 2\mathbb{E} \left[\langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t) + \frac{\lambda}{1-\lambda} \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)) \rangle \right] \quad (27)$$

□

Remark 8 (Behavior of the gradient). *The gradient term shows that $\mathbf{F}_{\theta}(\mathbf{x}_t, t)$ is pushed towards the convex combination*

$$(\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t) + \frac{\lambda}{1-\lambda} \cdot (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t))$$

Therefore, λ smoothly interpolates between two regimes:

- *If $\lambda = 0$, the gradient reduces to standard flow matching gradient $\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t$, which corresponds to the gradient of the multi-step objective.*
- *If $\lambda \rightarrow 1$, the gradient becomes $\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t + t \cdot \frac{d\mathbf{F}_{\theta}(\mathbf{x}_t, t)}{dt}$, where $\mathbf{F}_{\theta}(\mathbf{x}_t, t)$ is the predicted velocity, \mathbf{z}_t is the target mean velocity, and $t \cdot \frac{d\mathbf{F}_{\theta}(\mathbf{x}_t, t)}{dt}$ acts as a temporal consistency correction, encouraging the predicted velocity field to remain coherent across time. In other words, when $\mathbf{F}_{\theta}(\mathbf{x}_t, t)$ becomes the mean velocity, there will need no correction term and the term $\frac{d\mathbf{F}_{\theta}(\mathbf{x}_t, t)}{dt}$ constantly equals to zero.*
- *If $\lambda \in (0, 1)$, the gradient simultaneously enforces accuracy of the velocity prediction and strengthens the requirement that the velocity field be consistent across all time steps. As a result, the behavior increasingly resembles a few-step consistency model as $\lambda \rightarrow 1$ because the weights of the correction term becomes increasingly larger.*

Theorem 7 (optimal solution of surrogate objective in General Case). *Under Assump. 1, let's define*

$$A(t) := \frac{\alpha(t)}{D(t)} - \frac{\alpha(\lambda t)}{D(\lambda t)}, \quad B(t) := \frac{\alpha(t)}{D(t)}, \quad C(t) := \frac{\alpha(t)}{2D(t)}, \quad D(t) := \alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t).$$

Let the surrogate loss be

$$\mathcal{G}(\theta) = \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\frac{C(t)}{\hat{\omega}(t)} \left(A(t) \left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t \right\|_2^2 + B(\lambda t) \left\| (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^{-}}(\mathbf{x}_{\lambda t}, \lambda t)) - (\mathbf{z}_t - \mathbf{z}_{\lambda t}) \right\|_2^2 \right) \right] \quad (28)$$

where $\mathbf{x}_t = \alpha(t) \cdot \mathbf{z} + \gamma(t) \cdot \mathbf{x}$, $\mathbf{z}_t = \hat{\alpha}(t) \cdot \mathbf{z} + \hat{\gamma}(t) \cdot \mathbf{x}$, $0 < \lambda < 1$. Then, the optimal solution of the surrogate objective is:

$$\mathbf{F}_{\theta^*}(\mathbf{x}_t, t) = \mathbb{E}_{\mathbf{z}, \mathbf{x}} \left[\mathbf{z}_t + \frac{B(\lambda t)}{A(t) + B(\lambda t)} (\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) - \mathbf{z}_{\lambda t}) \mid \mathbf{x}_t \right]. \quad (29)$$

Proof. The proof proceeds like the linear special case: for each fixed pair (\mathbf{x}_t, t) we optimise pointwise over the vector value $\mathbf{F}_{\theta}(\mathbf{x}_t, t)$.

For notational convenience define

$$\mathbf{f} := \mathbf{F}_{\theta}(\mathbf{x}_t, t), \quad \mathbf{a} := \mathbf{z}_t, \quad \mathbf{a}_{\lambda} := \mathbf{z}_{\lambda t}, \quad \mathbf{c} := \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t).$$

Observe that

$$\|(\mathbf{f} - \mathbf{c}) - (\mathbf{a} - \mathbf{a}_{\lambda})\|_2^2 = \|\mathbf{f} - (\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda})\|_2^2.$$

With this notation the pointwise contribution of the surrogate loss (omitting the common prefactor $1/\hat{\omega}(t)$) becomes

$$\mathcal{G}_t(\mathbf{f}) = C(t)A(t) \mathbb{E}[\|\mathbf{f} - \mathbf{a}\|_2^2 \mid \mathbf{x}_t] + C(t)B(\lambda t) \mathbb{E}[\|\mathbf{f} - (\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda})\|_2^2 \mid \mathbf{x}_t].$$

Define the weighted conditional mean

$$\mathbf{f}^* := \mathbb{E} \left[\frac{A(t)\mathbf{a} + B(\lambda t)(\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda})}{A(t) + B(\lambda t)} \mid \mathbf{x}_t \right]$$

which is well-defined under the assumptions on the coefficients. Expanding each squared term around \mathbf{f}^* gives

$$\begin{aligned} \mathcal{G}_t(\mathbf{f}) &= C(t)A(t) \mathbb{E}[\|\mathbf{f} - \mathbf{f}^*\|^2 + \|\mathbf{f}^* - \mathbf{a}\|^2 + 2\langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - \mathbf{a} \rangle \mid \mathbf{x}_t] \\ &\quad + C(t)B(\lambda t) \mathbb{E}[\|\mathbf{f} - \mathbf{f}^*\|^2 + \|\mathbf{f}^* - (\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda})\|^2 \\ &\quad \quad \quad + 2\langle \mathbf{f} - \mathbf{f}^*, \mathbf{f}^* - (\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda}) \rangle \mid \mathbf{x}_t]. \end{aligned}$$

Collecting terms yields

$$\mathcal{G}_t(\mathbf{f}) = C(t)(A(t) + B(\lambda t)) \mathbb{E}[\|\mathbf{f} - \mathbf{f}^*\|^2 \mid \mathbf{x}_t] + \mathcal{G}_t(\mathbf{f}^*) + 2\mathcal{I},$$

where the cross-term \mathcal{I} equals

$$\mathcal{I} = \mathbb{E} \left[\langle \mathbf{f} - \mathbf{f}^*, C(t)A(t)(\mathbf{f}^* - \mathbf{a}) + C(t)B(\lambda t)(\mathbf{f}^* - (\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda})) \rangle \mid \mathbf{x}_t \right].$$

By the definition of \mathbf{f}^* as the conditional expectation of the weighted target

$$C(t)A(t)\mathbf{a} + C(t)B(\lambda t)(\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda}),$$

it follows that the vector inside the inner product in \mathcal{I} has zero conditional expectation:

$$C(t)A(t) \mathbb{E}[\mathbf{a} \mid \mathbf{x}_t] + C(t)B(\lambda t) \mathbb{E}[\mathbf{a} + \mathbf{c} - \mathbf{a}_{\lambda} \mid \mathbf{x}_t] = C(t)(A(t) + B(\lambda t))\mathbf{f}^*.$$

Hence $\mathcal{I} = 0$, and we obtain the lower bound

$$\mathcal{G}_t(\mathbf{f}) \geq C(t)(A(t) + B(\lambda t)) \mathbb{E}[\|\mathbf{f} - \mathbf{f}^*\|^2 \mid \mathbf{x}_t] + \mathcal{G}_t(\mathbf{f}^*),$$

with equality at $\mathbf{f} = \mathbf{f}^*$. Therefore \mathbf{f}^* minimises the pointwise objective, and we recover (29). \square

F.1.7 UNIFIED TRAINING OBJECTIVE

Theorem 8 . Let $\lambda \in (0, 1)$ and define the scalar functions

$$\begin{aligned} A(t) &:= B(t) - B(\lambda t), & B(t) &:= \frac{\alpha(t)}{D(t)}, \\ C(t) &:= \frac{\alpha(t)}{2D(t)}, & D(t) &:= \alpha(t)\hat{\gamma}(t) - \hat{\alpha}(t)\gamma(t), & \hat{\omega}(t) &:= C(t)A(t). \end{aligned}$$

For a pair $(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x})$ and times $t, \lambda t$, we define the shorthand

$$\Delta_{\theta, \theta^-} \mathbf{f}^{\mathbf{x}} := \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t).$$

Assume θ and θ^- are two different variables and equal in value. Now we define the three functionals:

$$\begin{aligned} \mathcal{L}(\theta) &= \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\frac{1}{\hat{\omega}(t)} \left\| \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t) \right\|_2^2 \right], \\ \mathcal{G}(\theta) &= \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{z}_t \right\|_2^2 + \frac{B(\lambda t)}{\hat{\omega}(t)} \left\| (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t)) - (\mathbf{z}_t - \mathbf{z}_{\lambda t}) \right\|_2^2 \right], \\ \mathcal{N}(\theta) &= \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\frac{1}{2} \left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t) + 2 \cdot \frac{\Delta_{\theta^-, \theta^-} \mathbf{f}^{\mathbf{x}}}{A(t)} \right\|_2^2 \right]. \end{aligned}$$

Then, for all θ ,

$$\nabla_{\theta} \mathcal{L}(\theta) = \nabla_{\theta} \mathcal{G}(\theta) = \nabla_{\theta} \mathcal{N}(\theta).$$

Proof. The first equality $\nabla_{\theta} \mathcal{L}(\theta) = \nabla_{\theta} \mathcal{G}(\theta)$ is straightforward by [Thm. 5](#). Now, we prove the equality $\nabla_{\theta} \mathcal{L}(\theta) = \nabla_{\theta} \mathcal{N}(\theta)$. For brevity, we omit the expectation operator \mathbb{E} in the following derivation. Then, we compute the gradient of $\mathcal{L}(\theta)$ as follows:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta) &= \nabla_{\theta} \left[\frac{1}{\hat{\omega}(t)} \left\| \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}, \lambda t) \right\|_2^2 \right] \\ &= \frac{2}{\hat{\omega}(t)} \left\langle \nabla_{\theta} \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta}(\mathbf{x}_t, t), \mathbf{x}_t, t), \Delta_{\theta, \theta^-} \mathbf{f}^{\mathbf{x}} \right\rangle && \text{(by chain rule)} \\ &= \frac{2}{\hat{\omega}(t)} \left\langle \frac{\alpha(t)}{D(t)} \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \Delta_{\theta, \theta^-} \mathbf{f}^{\mathbf{x}} \right\rangle && \text{(by definition of } \mathbf{f}^{\mathbf{x}} \text{)} \\ &= \left\langle \nabla_{\theta} \mathbf{F}_{\theta}(\mathbf{x}_t, t), \frac{2\alpha(t)}{D(t)\hat{\omega}(t)} \Delta_{\theta, \theta^-} \mathbf{f}^{\mathbf{x}} \right\rangle. \end{aligned}$$

Since \mathbf{F}_{θ^-} does not depend on θ , we can equivalently write

$$\nabla_{\theta} \mathcal{L}(\theta) = \left\langle \nabla_{\theta} (\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_t, t)), \frac{2\alpha(t)}{D(t)\hat{\omega}(t)} \Delta_{\theta^-, \theta^-} \mathbf{f}^{\mathbf{x}} \right\rangle.$$

Next, using the identity

$$\hat{\omega}(t) = \frac{\alpha(t)}{2D(t)} A(t) \implies \frac{2\alpha(t)}{D(t)\hat{\omega}(t)} = \frac{2}{A(t)},$$

we obtain

$$\nabla_{\theta} \mathcal{L}(\theta) = \left\langle \nabla_{\theta} \left(\mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_t, t) + \frac{2}{A(t)} \Delta_{\theta^-, \theta^-} \mathbf{f}^{\mathbf{x}} \right), \frac{2}{A(t)} \Delta_{\theta^-, \theta^-} \mathbf{f}^{\mathbf{x}} \right\rangle.$$

Finally, by the standard identity $\langle \nabla f(x), f(x) \rangle = \frac{1}{2} \nabla \|f(x)\|_2^2$, we conclude

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta) &= \frac{1}{2} \nabla_{\theta} \left\| \mathbf{F}_{\theta}(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_t, t) + \frac{2}{A(t)} \Delta_{\theta^-, \theta^-} \mathbf{f}^{\mathbf{x}} \right\|_2^2 \\ &= \nabla_{\theta} \mathcal{N}(\theta). \end{aligned}$$

□

F.1.8 ENHANCED TARGET SCORE FUNCTION

Training a model directly with objective in (6) fails to produce realistic samples without Classifier-Free Guidance (CFG) (Ho & Salimans, 2022). However, while enhancing semantic information, it introduces significant computational overhead by approximately doubling the required function evaluations.

A recent approach (Tang et al., 2025) proposes modifying the target score function. Instead of the standard conditional score (Song et al., 2020b), $\nabla_{\mathbf{x}_t} \log(p_t(\mathbf{x}_t|\mathbf{c}))$, they propose an enhanced version $\nabla_{\mathbf{x}_t} \log\left(p_t(\mathbf{x}_t|\mathbf{c}) \left(\frac{p_{t,\theta}(\mathbf{x}_t|\mathbf{c})}{p_{t,\theta}(\mathbf{x}_t)}\right)^\zeta\right)$, where $\zeta \in (0, 1)$ denotes the enhancement ratio. This modification eliminates the need for CFG, enabling high-fidelity sample generation at a significantly reduced inference cost.

Inspired by this, we propose enhancing the target score function in a manner compatible with our unified training objective in (6). Specifically, we introduce a time-dependent enhancement strategy:

- (a) For $t \in [0, s]$, enhance \mathbf{x} and \mathbf{z} by applying $\mathbf{x}^* = \mathbf{x} + \zeta \cdot (\mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{x}}(\mathbf{F}_t^\emptyset, \mathbf{x}_t, t))$, $\mathbf{z}^* = \mathbf{z} + \zeta \cdot (\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{z}}(\mathbf{F}_t^\emptyset, \mathbf{x}_t, t))$. Here, $\mathbf{F}_t^\emptyset = \mathbf{F}_{\theta^-}(\mathbf{x}_t, t, \emptyset)$ and $\mathbf{F}_t = \mathbf{F}_{\theta^-}(\mathbf{x}_t, t)$.
- (b) For $t \in (s, 1]$, enhance \mathbf{x} and \mathbf{z} by applying $\mathbf{x}^* = \mathbf{x} + \frac{1}{2} (\mathbf{f}^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x})$ and $\mathbf{z}^* = \mathbf{z} + \frac{1}{2} (\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z})$. We consistently set $s = 0.75$ and see App. F.1.8 for more analysis.

An ablation study for this technique is shown in Sec. 4.4, and the pseudocode is shown in Alg. 1.

Recall that CFG proposes to modify the sampling distribution as

$$\tilde{p}_\theta(\mathbf{x}_t|\mathbf{c}) \propto p_\theta(\mathbf{x}_t|\mathbf{c})p_\theta(\mathbf{c}|\mathbf{x}_t)^\zeta,$$

Bayesian rule gives

$$p_\theta(\mathbf{c}|\mathbf{x}_t) = \frac{p_\theta(\mathbf{x}_t|\mathbf{c})p_\theta(\mathbf{c})}{p_\theta(\mathbf{x}_t)},$$

so we can further deduce

$$\begin{aligned} \tilde{p}_\theta(\mathbf{x}_t|\mathbf{c}) &\propto p_\theta(\mathbf{x}_t|\mathbf{c})p_\theta(\mathbf{c}|\mathbf{x}_t)^\zeta \\ &= p_\theta(\mathbf{x}_t|\mathbf{c})\left(\frac{p_\theta(\mathbf{x}_t|\mathbf{c})p_\theta(\mathbf{c})}{p_\theta(\mathbf{x}_t)}\right)^\zeta \\ &\propto p_\theta(\mathbf{x}_t|\mathbf{c})\left(\frac{p_\theta(\mathbf{x}_t|\mathbf{c})}{p_\theta(\mathbf{x}_t)}\right)^\zeta. \end{aligned}$$

When $t \in [0, s]$ ($s = 0.75$), inspired by above expression and a recent work (Tang et al., 2025), we choose to use below as the target score function for training

$$\nabla_{\mathbf{x}_t} \log\left(p_t(\mathbf{x}_t|\mathbf{c}) \left(\frac{p_{t,\theta}(\mathbf{x}_t|\mathbf{c})}{p_{t,\theta}(\mathbf{x}_t)}\right)^\zeta\right)$$

which equals to

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + \zeta (\nabla_{\mathbf{x}_t} \log p_{t,\theta}(\mathbf{x}_t|\mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_{t,\theta}(\mathbf{x}_t)).$$

In Thm. 2, we have shown that

$$\mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = -\alpha(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t), \text{ and } \mathbf{f}_*^{\mathbf{x}}(\mathbf{F}_t, \mathbf{x}_t, t) = \frac{\mathbf{x}_t + \alpha^2(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)},$$

so we can further deduce: For $\mathbf{f}_*^{\mathbf{z}}$ we originally want to learn:

$$\mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) = -\alpha(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t),$$

now it turns to

$$\begin{aligned} \mathbf{f}_*^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) &= -\alpha(t)\nabla_{\mathbf{x}_t} \log\left(p_t(\mathbf{x}_t|\mathbf{c}) \left(\frac{p_{t,\theta}(\mathbf{x}_t|\mathbf{c})}{p_{t,\theta}(\mathbf{x}_t)}\right)^\zeta\right) \\ &= -\alpha(t) [\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + \zeta (\nabla_{\mathbf{x}_t} \log p_{t,\theta}(\mathbf{x}_t|\mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_{t,\theta}(\mathbf{x}_t))] \\ &= -\alpha(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + \zeta (-\alpha(t)\nabla_{\mathbf{x}_t} \log p_{t,\theta}(\mathbf{x}_t|\mathbf{c}) + \alpha(t)\nabla_{\mathbf{x}_t} \log p_{t,\theta}(\mathbf{x}_t)) \\ &= -\alpha(t)\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t|\mathbf{c}) + \zeta (\mathbf{f}^{\mathbf{z}}(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{f}^{\mathbf{z}}(\mathbf{F}_t^\emptyset, \mathbf{x}_t, t)), \end{aligned}$$

thus in training we set the objective for \mathbf{f}^z as:

$$\mathbf{z}^* \leftarrow \mathbf{z} + \zeta \cdot (\mathbf{f}^z(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{f}^z(\mathbf{F}_t^\emptyset, \mathbf{x}_t, t)) .$$

Similarly, since $\mathbf{f}_*^x = \frac{\mathbf{x}_t + \alpha^2(t) \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)}{\gamma(t)}$ is also linear in the score function, we can use the same strategy to modify the training objective for \mathbf{f}^x :

$$\mathbf{x}^* \leftarrow \mathbf{x} + \zeta \cdot (\mathbf{f}^x(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{f}^x(\mathbf{F}_t^\emptyset, \mathbf{x}_t, t)) .$$

We can also derive that:

$$\mathbf{x}_t^* = \alpha(t) \cdot \mathbf{z}^* + \gamma(t) \cdot \mathbf{x}^* = \mathbf{x}_t ,$$

When $t \in (s, 1]$ ($s = 0.75$), we further slightly modify the target score function to

$$\nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t | \mathbf{c}) + \zeta (\nabla_{\mathbf{x}_t} \log p_{t, \theta}(\mathbf{x}_t | \mathbf{c}) - \nabla_{\mathbf{x}_t} \log p_t(\mathbf{x}_t)) , \zeta = 0.5$$

which corresponds to the following training objective:

$$\mathbf{x}^* \leftarrow \mathbf{x} + \frac{1}{2} (\mathbf{f}^x(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{x}) , \mathbf{z}^* \leftarrow \mathbf{z} + \frac{1}{2} (\mathbf{f}^z(\mathbf{F}_t, \mathbf{x}_t, t) - \mathbf{z}) .$$

After applying above enhanced target score matching, we can further deduce the training objective for \mathbf{f}^x as:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{\omega}(t)} \|\mathbf{f}^x(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{x}^*\|_2^2 \right] .$$

Similarly, introducing the λ , we can deduce:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{\omega}(t)} \|\mathbf{f}^x(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t, t) - \mathbf{f}^x(\mathbf{F}_\theta(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}^*, \lambda t)\|_2^2 \right] .$$

Using $\mathbf{x}_t^* = \mathbf{x}_t$, we can further deduce:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{z}, \mathbf{x}) \sim p(\mathbf{z}, \mathbf{x}), t} \left[\frac{1}{\hat{\omega}(t)} \|\mathbf{f}^x(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t^*, t) - \mathbf{f}^x(\mathbf{F}_\theta(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}^*, \lambda t)\|_2^2 \right] .$$

Enhanced target score matching for training objective (6). By following the derivation in [Thm. 8](#), we can deduce:

$$\mathcal{N}(\theta) = \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\frac{1}{2} \|\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_t, t) + 2 \cdot \frac{\Delta_{\theta^-, \theta^-} \mathbf{f}^x}{B(t) - B(\lambda t)}\|_2^2 \right] .$$

where

$$\Delta_{\theta^-, \theta^-} \mathbf{f}^x = \mathbf{f}^x(\mathbf{F}_\theta(\mathbf{x}_t, t), \mathbf{x}_t^*, t) - \mathbf{f}^x(\mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t), \mathbf{x}_{\lambda t}^*, \lambda t) .$$

Enhanced target score matching for training objective (13). By following the derivation in [Thm. 5](#), we can further deduce:

$$\mathcal{G}(\theta) = \mathbb{E}_{\mathbf{z}, \mathbf{x}, t} \left[\|\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{z}_t^*\|_2^2 + \frac{B(t)}{\hat{\omega}(t)} \|(\mathbf{F}_\theta(\mathbf{x}_t, t) - \mathbf{F}_{\theta^-}(\mathbf{x}_{\lambda t}, \lambda t)) - (\mathbf{z}_t^* - \mathbf{z}_{\lambda t}^*)\|_2^2 \right]$$

F.1.9 UNIFIED SAMPLING PROCESS

Deterministic sampling. When the stochastic ratio $\rho = 0$, let's analyze a special case where the coefficients satisfying $\hat{\alpha}(t) = \frac{d\alpha(t)}{dt}$, $\hat{\gamma}(t) = \frac{d\gamma(t)}{dt}$. Let $\Delta t = t_{i+1} - t_i$, for the core updating rule we have:

$$\begin{aligned} \mathbf{x}' &= \alpha(t_{i+1}) \cdot \hat{\mathbf{z}} + \gamma(t_{i+1}) \cdot \hat{\mathbf{x}} \\ &= (\alpha(t_i) + \alpha'(t_i) \Delta t + o(\Delta t)) \cdot \hat{\mathbf{z}} + (\gamma(t_i) + \gamma'(t_i) \Delta t + o(\Delta t)) \cdot \hat{\mathbf{x}} \\ &= (\alpha(t_i) \hat{\mathbf{z}} + \gamma(t_i) \hat{\mathbf{x}}) + (\hat{\alpha}(t_i) \hat{\mathbf{z}} + \hat{\gamma}(t_i) \hat{\mathbf{x}}) \cdot \Delta t + o(\Delta t) \\ &= (\alpha(t_i) \mathbf{f}^z(\mathbf{F}, \tilde{\mathbf{x}}, t_i) + \gamma(t_i) \mathbf{f}^x(\mathbf{F}, \tilde{\mathbf{x}}, t_i)) + (\hat{\alpha}(t_i) \mathbf{f}^z(\mathbf{F}, \tilde{\mathbf{x}}, t_i) + \hat{\gamma}(t_i) \mathbf{f}^x(\mathbf{F}, \tilde{\mathbf{x}}, t_i)) \cdot \Delta t + o(\Delta t) \\ &= (\alpha(t_i) \frac{\hat{\gamma}(t_i) \cdot \tilde{\mathbf{x}} - \gamma(t_i) \cdot \mathbf{F}(\tilde{\mathbf{x}}, t_i)}{\alpha(t_i) \cdot \hat{\gamma}(t_i) - \hat{\alpha}(t_i) \cdot \gamma(t_i)} + \gamma(t_i) \frac{\alpha(t_i) \cdot \mathbf{F}(\tilde{\mathbf{x}}, t_i) - \hat{\alpha}(t_i) \cdot \mathbf{x}_t}{\alpha(t_i) \cdot \hat{\gamma}(t_i) - \hat{\alpha}(t_i) \cdot \gamma(t_i)}) \\ &\quad + (\hat{\alpha}(t_i) \frac{\hat{\gamma}(t_i) \cdot \tilde{\mathbf{x}} - \gamma(t_i) \cdot \mathbf{F}(\tilde{\mathbf{x}}, t_i)}{\alpha(t_i) \cdot \hat{\gamma}(t_i) - \hat{\alpha}(t_i) \cdot \gamma(t_i)} + \hat{\gamma}(t_i) \frac{\alpha(t_i) \cdot \mathbf{F}(\tilde{\mathbf{x}}, t_i) - \hat{\alpha}(t_i) \cdot \mathbf{x}_t}{\alpha(t_i) \cdot \hat{\gamma}(t_i) - \hat{\alpha}(t_i) \cdot \gamma(t_i)}) \cdot \Delta t + o(\Delta t) \\ &= \tilde{\mathbf{x}} + \mathbf{F}(\tilde{\mathbf{x}}, t_i) \cdot \Delta t + o(\Delta t) \end{aligned}$$

In this case $\mathbf{F}(\cdot, \cdot)$ tries to predict the velocity field of the flow model, and we can see that the term $\tilde{\mathbf{x}} + \mathbf{F}(\tilde{\mathbf{x}}, t_i) \cdot \Delta t$ corresponds to the sampling rule of the Euler ODE solver.

Stochastic sampling. As for case when the stochastic ratio $\rho \neq 0$, follow the Euler-Maruyama numerical methods of SDE, the noise injected should be a Gaussian with zero mean and variance proportional to Δt , so when the updating rule is $\mathbf{x}' = \alpha(t_{i+1}) \cdot (\sqrt{1-\rho} \cdot \hat{\mathbf{z}} + \sqrt{\rho} \cdot \mathbf{z}) + \gamma(t_{i+1}) \cdot \hat{\mathbf{x}}$, the coefficient of \mathbf{z} should satisfy

$$\alpha(t_{i+1})\sqrt{\rho} \propto \sqrt{\Delta t}, \rho \propto \frac{\Delta t}{\alpha^2(t_{i+1})}$$

In practice, we set

$$\rho = \frac{2\Delta t \cdot \alpha(t_i)}{\alpha^2(t_{i+1})}.$$

which corresponds to $g(t) = \sqrt{2\alpha(t)}$ for the SDE $d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}$.

F.1.10 EXTRAPOLATING ESTIMATION

Theorem 9 (Local truncation error of the extrapolation estimation). Assume the sampling process uses a uniform time step size $h = t_{i+1} - t_i = t_i - t_{i-1}$. Let $\Phi(t)$ denote the virtual endpoint estimates (e.g., $\hat{\mathbf{x}}$ or $\hat{\mathbf{z}}$), assumed to be twice continuously differentiable. By setting the extrapolation coefficient $\kappa = 1$, the proposed predictor:

$$\hat{\Phi}_i^{\text{ext}} = \Phi(t_i) + \kappa \cdot (\Phi(t_i) - \Phi(t_{i-1})) \quad (30)$$

achieves a Local Truncation Error (LTE) of order $\mathcal{O}(h^2)$, effectively serving as a second-order approximation to the trajectory.

Proof. We analyze the error by expanding the exact solution $\Phi(t)$ around time t_i . Under the assumption of uniform steps, let h be the step size. The Taylor expansion for the true target value at t_{i+1} is:

$$\Phi(t_{i+1}) = \Phi(t_i) + \Phi'(t_i)h + \frac{1}{2}\Phi''(t_i)h^2 + \mathcal{O}(h^3). \quad (31)$$

Similarly, the historical value at t_{i-1} (where $t_{i-1} = t_i - h$) is expanded as:

$$\Phi(t_{i-1}) = \Phi(t_i) - \Phi'(t_i)h + \frac{1}{2}\Phi''(t_i)h^2 + \mathcal{O}(h^3). \quad (32)$$

The algorithm computes the extrapolated estimate $\hat{\Phi}_i^{\text{ext}}$ using the fixed constant κ . Substituting (32) into the extrapolation formula:

$$\begin{aligned} \hat{\Phi}_i^{\text{ext}} &= \Phi(t_i) + \kappa (\Phi(t_i) - \Phi(t_{i-1})) \\ &= \Phi(t_i) + \kappa \left(\Phi(t_i) - \left[\Phi(t_i) - \Phi'(t_i)h + \frac{1}{2}\Phi''(t_i)h^2 \right] \right) + \mathcal{O}(h^3) \\ &= \Phi(t_i) + \kappa h \Phi'(t_i) - \frac{\kappa}{2} h^2 \Phi''(t_i) + \mathcal{O}(h^3). \end{aligned} \quad (33)$$

To evaluate the local truncation error $\mathcal{E}_{loc} = \|\Phi(t_{i+1}) - \hat{\Phi}_i^{\text{ext}}\|_2$:

$$\begin{aligned} \mathcal{E}_{loc} &= \|\Phi(t_{i+1}) - \hat{\Phi}_i^{\text{ext}}\|_2 \\ &= \|(1-\kappa)h\Phi'(t_i) + \frac{1}{2}(1+\kappa)h^2\Phi''(t_i)\|_2 + \mathcal{O}(h^3) \end{aligned}$$

Thus, with uniform steps and $\kappa = 1$, the proposed extrapolation correctly captures the linear trend of the virtual endpoints, resulting in a local error of $\mathcal{O}(h^2)$ (set $\|\Phi''(\cdot)\|_2 \leq C$):

$$\mathcal{E}_{loc} = \|\Phi''(t_i)\|_2 \cdot h^2 + \mathcal{O}(h^3) \leq C \cdot h^2.$$

□

Theorem 10 (Global error of the extrapolation estimation). Under the same assumptions as in Theorem 9, suppose the sampling process uses a uniform step size h , and let $\hat{\Phi}_i$ denote the estimated virtual endpoints produced by the sampler at time t_i . When the extrapolation coefficient is set to $\kappa = 1$, the extrapolation-based update achieves a global error of order

$$\|\Phi(t_N) - \hat{\Phi}_N\|_2 = \mathcal{O}(h), \quad (34)$$

where $t_N - t_0 = Nh$ is fixed. In other words, although the local truncation error is second-order, the accumulated global error over N steps is of first order.

Proof. Define the global error at step i by

$$\mathbf{e}_i := \Phi(t_i) - \hat{\Phi}_i. \quad (35)$$

Following Algorithm 2, the extrapolated estimate used in the update is

$$\hat{\Phi}^{\text{ext}} = \hat{\Phi}_i + \kappa(\hat{\Phi}_i - \hat{\Phi}_{i-1}). \quad (36)$$

Consider the hypothetical extrapolation formed using the exact solution:

$$\Phi_{\text{true}}^{\text{ext}} = \Phi(t_i) + \kappa(\Phi(t_i) - \Phi(t_{i-1})). \quad (37)$$

The difference between the true and estimated extrapolations can be expressed directly using the global errors:

$$\begin{aligned} \Phi_{\text{true}}^{\text{ext}} - \hat{\Phi}^{\text{ext}} &= (\Phi(t_i) - \hat{\Phi}_i) + \kappa\left((\Phi(t_i) - \hat{\Phi}_i) - (\Phi(t_{i-1}) - \hat{\Phi}_{i-1})\right) \\ &= (1 + \kappa)\mathbf{e}_i - \kappa\mathbf{e}_{i-1}. \end{aligned} \quad (38)$$

Next, the local truncation error established in Theorem 9 states that, for $\kappa = 1$,

$$\Phi(t_{i+1}) - \Phi_{\text{true}}^{\text{ext}} = \mathcal{O}(h^2). \quad (39)$$

Combining the two relations yields the recursion for the global error:

$$\begin{aligned} \mathbf{e}_{i+1} &= \Phi(t_{i+1}) - \hat{\Phi}^{\text{ext}} \\ &= \underbrace{(\Phi(t_{i+1}) - \Phi_{\text{true}}^{\text{ext}})}_{\mathcal{O}(h^2)} + (\Phi_{\text{true}}^{\text{ext}} - \hat{\Phi}^{\text{ext}}) \\ &= \mathcal{O}(h^2) + (1 + \kappa)\mathbf{e}_i - \kappa\mathbf{e}_{i-1}. \end{aligned} \quad (40)$$

Setting $\kappa = 1$ gives the linear difference equation

$$\mathbf{e}_{i+1} = 2\mathbf{e}_i - \mathbf{e}_{i-1} + \mathcal{O}(h^2). \quad (41)$$

The corresponding homogeneous relation

$$\mathbf{e}_{i+1} - 2\mathbf{e}_i + \mathbf{e}_{i-1} = 0 \quad (42)$$

has characteristic polynomial $(r - 1)^2$, whose general solution is a linear function of i . Hence the homogeneous component introduces at most a linear growth factor in i but no exponential amplification.

Unrolling the recursion over N steps and noting that each step contributes an $\mathcal{O}(h^2)$ non-homogeneous term yields

$$\|\mathbf{e}_N\|_2 \leq CNh^2 + \mathcal{O}(h^3). \quad (43)$$

Because the total integration time is fixed, $Nh = t_N - t_0 = \mathcal{O}(1)$, and thus

$$\|\Phi(t_N) - \hat{\Phi}_N\|_2 = \|\mathbf{e}_N\|_2 = \mathcal{O}(h). \quad (44)$$

Therefore, although the extrapolation step achieves a second-order local truncation error, the accumulated global error across N uniform steps is of first order. \square

F.1.11 INTERPRETATION OF UNIFIED SAMPLING PROCESS

The validity of the decomposition: The decomposition of $\tilde{\mathbf{x}}_t$ is guaranteed by the design of $\mathbf{f}^{\mathbf{x}}$ and $\mathbf{f}^{\mathbf{z}}$:

$$\alpha(t) \cdot \hat{\mathbf{z}}_t + \gamma(t) \cdot \hat{\mathbf{x}}_t = \alpha(t) \cdot \mathbf{f}^{\mathbf{z}}(\mathbf{F}_{\theta^-}(\tilde{\mathbf{x}}_t, t), \tilde{\mathbf{x}}_t, t) + \gamma(t) \cdot \mathbf{f}^{\mathbf{x}}(\mathbf{F}_{\theta^-}(\tilde{\mathbf{x}}_t, t), \tilde{\mathbf{x}}_t, t) \quad (45)$$

$$= \alpha(t) \cdot \frac{\hat{\gamma}(t) \cdot \tilde{\mathbf{x}}_t - \gamma(t) \cdot \mathbf{F}_t^{\theta^-}}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} + \gamma(t) \cdot \frac{\alpha(t) \cdot \mathbf{F}_t^{\theta^-} - \hat{\alpha}(t) \cdot \tilde{\mathbf{x}}_t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} \quad (46)$$

$$= \frac{\alpha(t) \cdot \hat{\gamma}(t) \cdot \tilde{\mathbf{x}}_t - \alpha(t) \cdot \gamma(t) \cdot \mathbf{F}_t^{\theta^-} + \gamma(t) \cdot \alpha(t) \cdot \mathbf{F}_t^{\theta^-} - \gamma(t) \cdot \hat{\alpha}(t) \cdot \tilde{\mathbf{x}}_t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} \quad (47)$$

$$= \frac{\alpha(t) \cdot \hat{\gamma}(t) \cdot \tilde{\mathbf{x}}_t - \gamma(t) \cdot \hat{\alpha}(t) \cdot \tilde{\mathbf{x}}_t}{\alpha(t) \cdot \hat{\gamma}(t) - \hat{\alpha}(t) \cdot \gamma(t)} \quad (48)$$

$$= \tilde{\mathbf{x}}_t \quad (49)$$

The validity of the reconstruction: Firstly, the decomposition and reconstruction forms one DDIM step. Specifically, the decomposition gives:

$$\mathbf{x}_t = \alpha(t) \cdot \mathbf{f}^z(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t) + \gamma(t) \cdot \mathbf{f}^x(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t)$$

and the reconstruction gives:

$$\begin{aligned} \mathbf{x}_s &= \alpha(s) \cdot \mathbf{f}^z(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t) + \gamma(s) \cdot \mathbf{f}^x(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t) \\ &= \frac{\alpha(s)}{\alpha(t)} \cdot (\mathbf{x}_t - \gamma(t) \cdot \mathbf{f}^x(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t)) + \gamma(s) \cdot \mathbf{f}^x(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t) \quad (\text{by the decomposition}) \\ &= \frac{\alpha(s)}{\alpha(t)} \cdot \mathbf{x}_t + \left(\gamma(s) - \frac{\alpha(s)}{\alpha(t)} \cdot \gamma(t) \right) \cdot \mathbf{f}^x(\mathbf{F}_{\theta^-}(\mathbf{x}_t, t), \mathbf{x}_t, t) \end{aligned}$$

This actually forms one DDIM step. In the formula (11) and appendix C.1 of (Zhou et al., 2025), they show that the DDIM interpolant is self-consistent and is marginal-preserving when $\mathbf{f}^x(\mathbf{F}_{\theta^*}(\mathbf{x}_t, t), \mathbf{x}_t, t) = \mathbb{E}[\mathbf{x} | \mathbf{x}_t]$.

- For multi-step models ($\lambda = 0$), we have proved $\mathbf{f}^x(\mathbf{F}_{\theta^*}(\mathbf{x}_t, t), \mathbf{x}_t, t) = \mathbb{E}[\mathbf{x} | \mathbf{x}_t]$ in [Lem. 5](#) when $\rho = 0$ and $\kappa = 0$.
- For few-step consistency models ($\lambda \rightarrow 1$), we always set $\rho = 1$ in sampling process (see [Table 6](#)), which is the same as the sampling process of consistency model (Song et al., 2023).

Proposition 2 (Equivalence to multi-step consistency sampler when $\rho = 1$). Consider [Algorithm 2](#) and assume the per-step predictor $\mathbf{f}^x(\cdot)$ returns the same virtual-endpoint (denoised) estimate used by the multi-step consistency model. When the stochastic ratio is set to $\rho = 1$ and the sampler runs in first-order mode ($\nu = 1$), each update of [Algorithm 2](#) is identical to the corresponding update of the multi-step consistency model sampler ([Song et al., 2023](#)).

Proof. We compare a single step of [Algorithm 2](#) (with $\nu = 1$) to the standard multi-step consistency sampler update. Fix a generic step index i and use the algorithm’s notation.

Under $\nu = 1$ the algorithm does not execute the second-order correction block; hence the core update used to produce the next state \mathbf{x}' (line computing “estimated next time sample”) is

$$\mathbf{x}' = \alpha(t_{i+1}) \cdot (\sqrt{1-\rho} \hat{\mathbf{z}} + \sqrt{\rho} \mathbf{z}) + \gamma(t_{i+1}) \cdot \hat{\mathbf{x}}, \quad (50)$$

where $\hat{\mathbf{x}}$ (resp. $\hat{\mathbf{z}}$) denotes the extrapolated denoised (resp. latent) estimate computed from the current model output as in the algorithm, and $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is a freshly drawn Gaussian.

Set $\rho = 1$. Then $\sqrt{1-\rho} = 0$ and $\sqrt{\rho} = 1$, so (50) reduces to

$$\mathbf{x}' = \alpha(t_{i+1}) \mathbf{z} + \gamma(t_{i+1}) \hat{\mathbf{x}}. \quad (51)$$

Now recall the common generative parameterization used by multi-step consistency models: a noisy state at time t is written as

$$\mathbf{x}_t = \gamma(t) \mathbf{x}_0 + \alpha(t) \mathbf{z}, \quad \mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

and the consistency sampler constructs the next-step state by using the model’s estimate of the denoised endpoint \mathbf{x}_0 (denote this estimate by the same symbol $\hat{\mathbf{x}}$) together with a fresh Gaussian \mathbf{z} to form

$$\mathbf{x}_{t_{i+1}} = \gamma(t_{i+1}) \hat{\mathbf{x}} + \alpha(t_{i+1}) \mathbf{z}. \quad (52)$$

Comparing (51) and (52) we see they are algebraically identical: the next state is produced by the same deterministic combination of the denoised estimate $\hat{\mathbf{x}}$ and an independent Gaussian \mathbf{z} , scaled by the same schedule coefficients $\gamma(t_{i+1})$ and $\alpha(t_{i+1})$.

The equivalence in distribution follows because [Algorithm 2](#) draws $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ independently at each step (same as the consistency sampler), and by the proposition hypothesis \mathbf{f}^x returns the same denoised/endpoint estimate used by the consistency model. Therefore, for $\rho = 1$ and $\nu = 1$ each single-step mapping (and its randomness) produced by [Algorithm 2](#) coincides exactly with the corresponding single-step mapping of the multi-step consistency sampler. Hence the first-order variant of [Algorithm 2](#) with $\rho = 1$ is equivalent to the multi-step consistency model sampler. \square

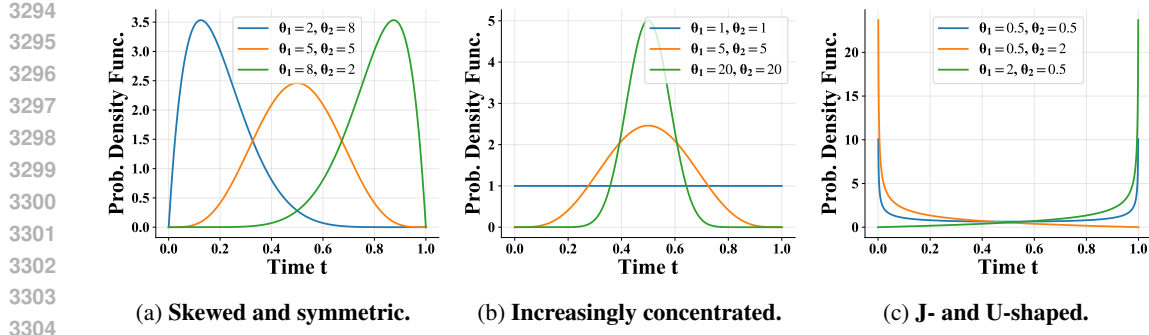


Figure 11: Probability density functions of the Beta distribution over the domain $t \in [0, 1]$ for various shape-parameter θ_1, θ_2 .

F.2 OTHER TECHNIQUES

F.2.1 BETA TRANSFORMATION

We utilize three representative cases to illustrate how the Beta transformation $f_{\text{Beta}}(t; \theta_1, \theta_2)$ generalizes time warping mechanisms for $t \in [0, 1]$.

Standard logit-normal time transformation (Yao et al., 2025; Esser et al., 2024). For $t \sim \mathcal{U}(0, 1)$, the logit-normal transformation $f_{\text{lognorm}}(t; 0, 1) = \frac{1}{1 + \exp(-\Phi^{-1}(t))}$ generates a symmetric density profile peaked at $t = 0.5$, consistent with the central maximum of the logistic-normal distribution. Analogously, the Beta transformation $f_{\text{Beta}}(t; \theta_1, \theta_2)$ (with $\theta_1, \theta_2 > 1$) produces a density peak at $t = \frac{\theta_1 - 1}{\theta_1 + \theta_2 - 2}$. When $\theta_1 = \theta_2 > 1$, this reduces to $t = 0.5$, mirroring the logit-normal case. Both transformations concentrate sampling density around critical time regions, enabling importance sampling for accelerated training. Notably, this effect can be equivalently achieved by directly sampling $t \sim \text{Beta}(\theta_1, \theta_2)$.

Uniform time distribution (Yao et al., 2025; Yu et al., 2024; Ma et al., 2024; Lipman et al., 2022). The uniform limit case emerges when $\theta_1 = \theta_2 = 1$, reducing $f_{\text{Beta}}(t; 1, 1)$ to an identity transformation. This corresponds to a flat density $p(t) = 1$, reflecting no temporal preference—a baseline configuration widely adopted in diffusion and flow-based models.

Approximately symmetrical time distribution (Song et al., 2023; Song & Dhariwal, 2023; Karras et al., 2022; 2024b). For near-symmetric configurations where $\theta_1 \approx \theta_2 > 1$, the Beta transformation induces quasi-symmetrical densities with tunable central sharpness. For instance, setting $\theta_1 = \theta_2 = 2$ yields a parabolic density peaking at $t = 0.5$, while $\theta_1 = \theta_2 \rightarrow 1^+$ asymptotically approaches uniformity. This flexibility allows practitioners to interpolate between uniform sampling and strongly peaked distributions, adapting to varying requirements for temporal resolution in training. Such approximate symmetry is particularly useful in consistency models where balanced gradient propagation across time steps is critical.

Furthermore, Fig. 11 further demonstrates the flexibility of the beta distribution.

F.2.2 KUMARASWAMY TRANSFORMATION

Lemma 12 (Piecewise monotone error). Suppose f, g are continuous and nondecreasing on $[0, 1]$, and agree at

$$0 = x_0 < x_1 < \dots < x_n = 1,$$

i.e. $f(x_j) = g(x_j)$ for $j = 0, \dots, n$. Let $\Delta_j = g(x_j) - g(x_{j-1})$. Then for every $t \in [x_{j-1}, x_j]$,

$$|f(t) - g(t)| \leq \Delta_j.$$

In particular, if each $\Delta_j \leq \frac{1}{4}$, then $\|f - g\|_{L^\infty} \leq \frac{1}{4}$.

Proof. On $[x_{j-1}, x_j]$ monotonicity gives

$$f(t) - g(t) \leq f(x_j) - g(x_{j-1}) = g(x_j) - g(x_{j-1}) = \Delta_j,$$

and similarly $g(t) - f(t) \leq \Delta_j$. \square

Theorem 11 (L^2 approximation bound of monotonic functions by generalized Kumaraswamy transformation). Let $\mathcal{G} = \{g \in C([0, 1]) : g \text{ nondecreasing, } g(0) = 0, g(1) = 1\}$, and define for $a, b, c > 0$, $f_{a,b,c}(t) = (1 - (1 - t^a)^b)^c$, $t \in [0, 1]$. Then

$$\sup_{g \in \mathcal{G}} \inf_{a,b,c > 0} \int_0^1 [f_{a,b,c}(t) - g(t)]^2 dt \leq \frac{1}{16}.$$

Proof. Let $g \in \mathcal{G}$. By continuity and the Intermediate-Value Theorem there exist

$$0 < t_1 < t_0 < t_2 < 1, \quad g(t_1) = \frac{1}{4}, \quad g(t_0) = \frac{1}{2}, \quad g(t_2) = \frac{3}{4}.$$

We will choose $(a, b, c) > 0$ so that

$$f_{a,b,c}(t_j) = g(t_j) \quad (j = 1, 0, 2),$$

and then apply the piecewise monotone [Lem. 12](#) on the partition

$$0, t_1, t_0, t_2, 1$$

to conclude $\|f_{a,b,c} - g\|_{L^\infty} \leq \frac{1}{4}$ and hence $\|f_{a,b,c} - g\|_{L^2}^2 \leq \frac{1}{16}$.

Existence via the implicit function theorem. Define

$$F : \mathbb{R}_{>0}^3 \longrightarrow \mathbb{R}^3, \quad F(a, b, c) = \begin{pmatrix} f_{a,b,c}(t_1) - \frac{1}{4} \\ f_{a,b,c}(t_0) - \frac{1}{2} \\ f_{a,b,c}(t_2) - \frac{3}{4} \end{pmatrix}.$$

Then F is C^1 , and at the ‘‘base point’’ $(a, b, c) = (1, 1, 1)$ with $(t_1, t_0, t_2) = (\frac{1}{4}, \frac{1}{2}, \frac{3}{4})$ we have $f_{1,1,1}(t) = t$ so $F(1, 1, 1) = 0$, and the Jacobian $\partial F / \partial (a, b, c)$ there is invertible. By the Implicit Function Theorem, for each fixed (t_1, t_0, t_2) near $(\frac{1}{4}, \frac{1}{2}, \frac{3}{4})$ there is a unique local solution (a, b, c) .

Global non-degeneracy of the Jacobian. In order to continue this local solution to *all* triples $0 < t_1 < t_0 < t_2 < 1$, we show $\det(\partial_{(a,b,c)} F(a, b, c))$ never vanishes.

Set

$$u(t) = 1 - (1 - t^a)^b, \quad u_j = u(t_j) \in (0, 1), \quad f_j = u_j^c.$$

Then

$$\partial_a f_j = c u_j^{c-1} \partial_a u_j, \quad \partial_b f_j = c u_j^{c-1} \partial_b u_j, \quad \partial_c f_j = u_j^c \ln u_j.$$

Hence

$$\det J = \det \begin{pmatrix} c u_1^{c-1} u_{1,a} & c u_1^{c-1} u_{1,b} & u_1^c \ln u_1 \\ c u_0^{c-1} u_{0,a} & c u_0^{c-1} u_{0,b} & u_0^c \ln u_0 \\ c u_2^{c-1} u_{2,a} & c u_2^{c-1} u_{2,b} & u_2^c \ln u_2 \end{pmatrix}.$$

Factor c from the first two columns and u_j^{c-1} from each row:

$$\det J = c^2 (u_1 u_0 u_2)^{c-1} \det \begin{pmatrix} u_{1,a} & u_{1,b} & u_1 \ln u_1 \\ u_{0,a} & u_{0,b} & u_0 \ln u_0 \\ u_{2,a} & u_{2,b} & u_2 \ln u_2 \end{pmatrix}.$$

Now

$$\begin{aligned} u_{j,b} &= -(1 - t_j^a)^b \ln(1 - t_j^a) = -(1 - u_j) \ln(1 - t_j^a), \\ u_{j,a} &= b(1 - t_j^a)^{b-1} t_j^a \ln t_j = -b(1 - u_j) \frac{t_j^a \ln t_j}{1 - t_j^a}. \end{aligned}$$

A direct—but straightforward—expansion shows

$$\det \begin{pmatrix} u_{1,a} & u_{1,b} & u_1 \ln u_1 \\ u_{0,a} & u_{0,b} & u_0 \ln u_0 \\ u_{2,a} & u_{2,b} & u_2 \ln u_2 \end{pmatrix} = c^{-2} b \frac{u_1 u_0 u_2}{(1 - u_1)(1 - u_0)(1 - u_2)} (u_0 - u_1)(u_2 - u_1)(u_2 - u_0).$$

Therefore

$$\det J(a, b, c) = b (u_1 u_0 u_2)^c \frac{(u_0 - u_1)(u_2 - u_1)(u_2 - u_0)}{(1 - u_1)(1 - u_0)(1 - u_2)} > 0,$$

since $0 < u_1 < u_0 < u_2 < 1$ and $a, b, c > 0$. Hence the Jacobian is everywhere non-zero, and the local solution by the Implicit Function Theorem extends along any path in the connected domain $\{0 < t_1 < t_0 < t_2 < 1\}$. We obtain a unique $(a, b, c) > 0$ solving

$$f_{a,b,c}(t_j) = g(t_j), \quad j = 1, 0, 2,$$

for every choice $0 < t_1 < t_0 < t_2 < 1$.

Completing the error estimate. By construction $f_{a,b,c}(0) = 0$, $f_{a,b,c}(1) = 1$, and $f_{a,b,c}(t_j) = g(t_j)$ for $j = 1, 0, 2$. On the partition

$$0, t_1, t_0, t_2, 1$$

the increments of g are each $1/4$. The piecewise monotone error [Lem. 12](#) yields $\|f_{a,b,c} - g\|_{L^\infty} \leq \frac{1}{4}$, hence

$$\int_0^1 [f_{a,b,c}(t) - g(t)]^2 dt \leq \|f - g\|_{L^\infty}^2 \leq \frac{1}{16}.$$

Since g was arbitrary in \mathcal{G} , we conclude

$$\sup_{g \in \mathcal{G}} \inf_{a,b,c > 0} \int_0^1 [f_{a,b,c}(t) - g(t)]^2 dt \leq \frac{1}{16}.$$

This completes the proof. \square

Setting and notation. Fix a positive real number $s > 0$ and consider the *shift function*

$$f_{\text{shift}}(t; s) = \frac{st}{1 + (s-1)t}, \quad t \in [0, 1].$$

For $a, b, c > 0$, define the *Kumaraswamy transform* as

$$f_{\text{Kuma}}(t; a, b, c) = \left(1 - (1 - t^a)^b\right)^c, \quad t \in [0, 1].$$

Notice that when $a = b = c = 1$ one obtains

$$f_{\text{Kuma}}(t; 1, 1, 1) = 1 - (1 - t^1)^1 = t,$$

so that the identity function appears as a special case.

We work in the Hilbert space $L^2([0, 1])$ with the inner product

$$\langle f, g \rangle = \int_0^1 f(t)g(t) dt.$$

Accordingly, we introduce the error functional

$$J(a, b, c) := \left\| f_{\text{Kuma}}(\cdot; a, b, c) - f_{\text{shift}}(\cdot; s) \right\|_2^2 \quad \text{and} \quad J_{\text{id}} := \left\| \text{id} - f_{\text{shift}}(\cdot; s) \right\|_2^2.$$

It is known that for $s \neq 1$ one has

$$\inf_{a,b,c} J(a, b, c) < J_{\text{id}}.$$

The goal is to quantify this improvement by optimally adjusting all three parameters (a, b, c) .

Quadratic approximation around the identity. Since the interesting behavior occurs near the identity $(a, b, c) = (1, 1, 1)$, we reparameterize as

$$\theta := \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix} := \begin{pmatrix} a - 1 \\ b - 1 \\ c - 1 \end{pmatrix}, \quad \text{with } \|\theta\| \ll 1.$$

Thus, we study the function

$$f_{\text{Kuma}}(t; 1 + \alpha, 1 + \beta, 1 + \gamma)$$

3456 in a small neighborhood of $(1, 1, 1)$. Writing

$$3457 \quad F(a, b, c; t) := f_{\text{Kuma}}(t; a, b, c) = \left(1 - (1 - t^a)^b\right)^c,$$

3458 a second-order Taylor expansion around $(a, b, c) = (1, 1, 1)$ gives

$$3459 \quad f_{\text{Kuma}}(t; 1 + \alpha, 1 + \beta, 1 + \gamma) = t + \sum_{i=1}^3 \theta_i g_i(t) + \frac{1}{2} \sum_{i,j=1}^3 \theta_i \theta_j h_{ij}(t) + \mathcal{O}(\|\theta\|^3), \quad (53)$$

3460 where

$$3461 \quad g_i(t) = \frac{\partial}{\partial \theta_i} f_{\text{Kuma}}(t; 1 + \theta) \Big|_{\theta=0} \quad \text{and} \quad h_{ij}(t) = \frac{\partial^2}{\partial \theta_i \partial \theta_j} f_{\text{Kuma}}(t; 1 + \theta) \Big|_{\theta=0}.$$

3462 A short calculation yields:

3463 (a) With respect to a (noting that for $b = c = 1$ one has $f_{\text{Kuma}}(t; a, 1, 1) = t^a$):

$$3464 \quad g_1(t) = \frac{\partial f_{\text{Kuma}}}{\partial a}(t; 1, 1, 1) = \frac{d}{da} t^a \Big|_{a=1} = t \ln t.$$

3465 (b) With respect to b (since for $a = 1, c = 1$ we have $f_{\text{Kuma}}(t; 1, b, 1) = 1 - (1 - t)^b$):

$$3466 \quad g_2(t) = \frac{\partial f_{\text{Kuma}}}{\partial b}(t; 1, 1, 1) = -(1 - t) \ln(1 - t).$$

3467 (c) With respect to c (noting that for $a = b = 1$ we have $f_{\text{Kuma}}(t; 1, 1, c) = t^c$):

$$3468 \quad g_3(t) = \frac{\partial f_{\text{Kuma}}}{\partial c}(t; 1, 1, 1) = t \ln t.$$

3469 Thus, we observe that

$$3470 \quad g_1(t) = g_3(t),$$

3471 which indicates an inherent redundancy in the three-parameter model. In consequence, the Gram

3472 matrix (defined below) will be of rank at most two.

3473 Next, define the difference between the identity and the shift functions:

$$3474 \quad g(t) := \text{id}(t) - f_{\text{shift}}(t; s) = t - \frac{st}{1 + (s-1)t} = (1-s) \frac{t(1-t)}{1 + (s-1)t}.$$

3475 Then, $J_{\text{id}} = \langle g, g \rangle$. Also, introduce the first-order moments and the Gram matrix:

$$3476 \quad v_i := \langle g, g_i \rangle, \quad G_{ij} := \langle g_i, g_j \rangle, \quad i, j = 1, 2, 3.$$

3477 Inserting the expansion (53) into the error functional gives

$$3478 \quad J(1 + \theta) = \|f_{\text{Kuma}}(\cdot; 1 + \theta) - f_{\text{shift}}(\cdot; s)\|_2^2 = J_{\text{id}} - 2 \sum_{i=1}^3 \theta_i v_i + \sum_{i,j=1}^3 \theta_i \theta_j G_{ij} + \mathcal{O}(\|\theta\|^3).$$

3479 Thus, the quadratic approximation (or model) of the error is

$$3480 \quad \hat{J}(\theta) := J_{\text{id}} - 2 \theta^\top v + \theta^\top G \theta.$$

3481 Since the Gram matrix G is positive semidefinite (and has a nontrivial null-space due to $g_1 = g_3$), the

3482 minimizer is determined only up to the null-space. To select the unique (minimum-norm) minimizer,

$$3483 \quad \theta^* = G^\dagger v,$$

3484 where G^\dagger denotes the Moore-Penrose pseudoinverse. The quadratic model is then minimized at

$$3485 \quad \hat{J}_{\min} = J_{\text{id}} - v^\top G^\dagger v.$$

3486 A scaling argument now shows that for any sufficiently small $\varepsilon > 0$ one has

$$3487 \quad J(1 + \varepsilon \theta^*) \leq \hat{J}(\varepsilon \theta^*) = J_{\text{id}} - \varepsilon^2 v^\top G^\dagger v < J_{\text{id}},$$

3488 so that the full nonlinear functional is improved by following the direction of θ^* .

3489 For convenience we introduce the explicit improvement factor

$$3490 \quad \rho_3(s) := \frac{v^\top G^\dagger v}{J_{\text{id}}(s)} \in (0, 1), \quad s \neq 1, \quad (54)$$

3491 so that our main bound can be written succinctly as

$$3492 \quad \min_{a,b,c>0} J(a, b, c) \leq \left(1 - \rho_3(s)\right) J_{\text{id}}(s). \quad (s > 0, s \neq 1) \quad (55)$$

3510 **Computation of the Gram matrix G .** We now compute the inner products

$$3511 \quad G_{ij} = \langle g_i, g_j \rangle, \quad i, j = 1, 2, 3.$$

3512 Since the functions g_1 and g_3 are identical, only two independent functions appear in the system. A
3513 standard fact from Beta-function calculus is that

$$3514 \quad \int_0^1 t^n \ln^2 t \, dt = \frac{2}{(n+1)^3}, \quad n > -1.$$

3515 Thus, one has

$$3516 \quad \langle g_1, g_1 \rangle = \int_0^1 t^2 \ln^2 t \, dt = \frac{2}{3^3} = \frac{2}{27},$$

$$3517 \quad \langle g_2, g_2 \rangle = \int_0^1 (1-t)^2 \ln^2(1-t) \, dt = \frac{2}{27},$$

3518 since the change of variable $u = 1 - t$ yields the same result.

$$3519 \quad \langle g_1, g_2 \rangle = - \int_0^1 t(1-t) \ln t \ln(1-t) \, dt = \frac{3\pi^2 - 37}{108}.$$

3520 It is now convenient to express the Gram matrix with an overall factor:

$$3521 \quad G = \frac{2}{27} \begin{pmatrix} 1 & r & 1 \\ r & 1 & r \\ 1 & r & 1 \end{pmatrix}, \quad r = \frac{3\pi^2 - 37}{8}.$$

3522 Since $g_1 = g_3$, it is clear that the columns (and rows) corresponding to parameters a and c are identical,
3523 so that $\text{rank}(G) = 2$. One can compute the Moore-Penrose pseudoinverse G^\dagger by eliminating one of
3524 the redundant rows/columns, inverting the resulting 2×2 block, and then re-embedding into $\mathbb{R}^{3 \times 3}$.
3525 One obtains

$$3526 \quad G^\dagger = \frac{27}{8(1-r^2)} \begin{pmatrix} 1 & -2r & 1 \\ -2r & 4 & -2r \\ 1 & -2r & 1 \end{pmatrix}.$$

3527 **Computation of the first-order moments v_i .** Recall that

$$3528 \quad g(t) = \text{id}(t) - f_{\text{shift}}(t; s) = t - \frac{st}{1 + (s-1)t}.$$

3529 This expression can be rewritten as

$$3530 \quad g(t) = (1-s)t(1-t)D_s(t), \quad \text{with } D_s(t) := \frac{1}{1 + (s-1)t}.$$

3531 Then, the first-order moments read

$$3532 \quad v_1 = v_3 = (1-s) \int_0^1 t(1-t)D_s(t) t \ln t \, dt,$$

$$3533 \quad v_2 = -(1-s) \int_0^1 t(1-t)D_s(t) (1-t) \ln(1-t) \, dt.$$

3534 These integrals can be expressed in closed form (involving logarithms and powers of $(s-1)$); in the
3535 case $s \neq 1$ at least one of the v_i is nonzero so that $\rho_3(s) > 0$.

3564 **A universal numerical improvement.** Since projecting onto the three-dimensional subspace
 3565 spanned by $\{g_1, g_2, g_3\}$ is at least as effective as projecting onto any one axis, we immediately deduce
 3566 that

$$3567 \rho_3(s) \geq \rho_1(s),$$

3568 where the one-parameter improvement factor is defined by

$$3569 \rho_1(s) := \frac{v_1(s)^2}{\langle g_1, g_1 \rangle J_{\text{id}}(s)}.$$

3572 By an elementary (albeit slightly tedious) estimate — for example, using the bounds $\frac{1}{2} \leq D_s(t) \leq 2$
 3573 valid for $|s - 1| \leq 1$ — one can show that

$$3574 \rho_1(s) \geq \frac{49}{1536}.$$

3577 Hence, one deduces that

$$3578 \rho_3(s) \geq \frac{49}{1536} \approx 0.0319, \quad \text{for } |s - 1| \leq 1.$$

3580 In particular, for $s \in [0.5, 2] \setminus \{1\}$ the optimal three-parameter Kumaraswamy transform reduces the
 3581 squared L^2 error by at least 3.19% compared with the identity mapping. Analogous bounds can be
 3582 obtained on any compact subset of $(0, \infty) \setminus \{1\}$.

3583 **Interpretation of the bound.** Inequality (55) strengthens the known qualitative result (namely, that
 3584 the three-parameter model can outperform the identity mapping) in two important respects:

- 3585 (a) Quantitative improvement: The explicit factor $\rho_3(s)$ is computable via one-dimensional integrals,
 3586 providing a concrete measure of the error reduction.
 3587 (b) Utilization of all three parameters: Even though the redundancy (i.e. $g_1 = g_3$) implies that the
 3588 Gram matrix is singular, the full three-parameter model still offers strict improvement; indeed,
 3589 one has $\rho_3(s) \geq \rho_1(s) > 0$ for $s \neq 1$. (Equality would require, hypothetically, that $v_2(s) = 0$,
 3590 which does not occur in practice.)
 3591

3592 **Summary.** For every shift parameter $s > 0$ with $s \neq 1$ there exist parameters (a, b, c) (in a
 3593 neighborhood of $(1, 1, 1)$) such that

$$3594 \left\| f_{\text{Kuma}}(\cdot; a, b, c) - f_{\text{shift}}(\cdot; s) \right\|_2^2 \leq \left(1 - \rho_3(s)\right) \left\| \text{id} - f_{\text{shift}}(\cdot; s) \right\|_2^2,$$

3596 with the improvement factor $\rho_3(s)$ defined in (54) and satisfying

$$3597 \rho_3(s) \geq 0.0319 \quad \text{on } s \in [0.5, 2] \setminus \{1\}.$$

3599 Thus, the full three-parameter Kumaraswamy transform not only beats the identity mapping but does
 3600 so by a quantifiable margin.

3601 F.2.3 DERIVATIVE ESTIMATION

3602 **Proposition 3 (Error estimates for forward and central difference quotients).** *Let $f \in C^3(I)$
 3603 where $I \subset \mathbb{R}$ is an open interval, and let $t \in I$. For $0 < \varepsilon$ small enough that $[t - \varepsilon, t + \varepsilon] \subset I$,
 3604 define the forward and central difference quotients*

$$3605 D_+ f(t) = \frac{f(t + \varepsilon) - f(t)}{\varepsilon}, \quad D_0 f(t) = \frac{f(t + \varepsilon) - f(t - \varepsilon)}{2\varepsilon}.$$

3606 Then

$$3607 D_+ f(t) = f'(t) + \frac{\varepsilon}{2} f''(t) + \frac{\varepsilon^2}{6} f^{(3)}(t + \theta_1 \varepsilon), \quad \text{for some } 0 < \theta_1 < 1,$$

$$3608 D_0 f(t) = f'(t) + \frac{\varepsilon^2}{12} \left[f^{(3)}(t + \theta_2 \varepsilon) + f^{(3)}(t - \theta_3 \varepsilon) \right], \quad \text{for some } 0 < \theta_2, \theta_3 < 1.$$

3609 In particular,

$$3610 D_+ f(t) - f'(t) = O(\varepsilon), \quad D_0 f(t) - f'(t) = O(\varepsilon^2),$$

3611 so for sufficiently small ε , the forward-difference error exceeds the central-difference error.
 3612
 3613
 3614
 3615
 3616
 3617

3618
3619
3620
3621
3622
3623
3624
3625
3626
3627
3628
3629
3630
3631
3632
3633
3634
3635
3636
3637
3638
3639
3640
3641
3642
3643
3644
3645
3646
3647
3648
3649
3650
3651
3652
3653
3654
3655
3656
3657
3658
3659
3660
3661
3662
3663
3664
3665
3666
3667
3668
3669
3670
3671

Proof. By Taylor’s theorem with Lagrange remainder, for some $0 < \theta_1 < 1$,

$$f(t + \varepsilon) = f(t) + f'(t)\varepsilon + \frac{1}{2}f''(t)\varepsilon^2 + \frac{1}{6}f^{(3)}(t + \theta_1\varepsilon)\varepsilon^3.$$

Dividing by ε gives the formula for $D_+f(t)$. Hence

$$D_+f(t) - f'(t) = \frac{1}{2}f''(t)\varepsilon + \frac{1}{6}f^{(3)}(t + \theta_1\varepsilon)\varepsilon^2 = O(\varepsilon).$$

Similarly, applying Taylor’s theorem at $t + \varepsilon$ and $t - \varepsilon$,

$$\begin{aligned} f(t + \varepsilon) &= f(t) + f'(t)\varepsilon + \frac{1}{2}f''(t)\varepsilon^2 + \frac{1}{6}f^{(3)}(t + \theta_2\varepsilon)\varepsilon^3, \\ f(t - \varepsilon) &= f(t) - f'(t)\varepsilon + \frac{1}{2}f''(t)\varepsilon^2 - \frac{1}{6}f^{(3)}(t - \theta_3\varepsilon)\varepsilon^3, \end{aligned}$$

for some $0 < \theta_2, \theta_3 < 1$. Subtracting and dividing by 2ε yields the formula for $D_0f(t)$ and

$$D_0f(t) - f'(t) = \frac{\varepsilon^2}{12}[f^{(3)}(t + \theta_2\varepsilon) + f^{(3)}(t - \theta_3\varepsilon)] = O(\varepsilon^2).$$

This completes the proof. \square

Proposition 4 . *Let $f : \mathbb{R} \rightarrow \mathbb{R}$ be differentiable, let $t \in \mathbb{R}$ and $\varepsilon > 0$. In BF16 arithmetic (1-bit sign, 8-bit exponent, 7-bit significand) with unit roundoff $\eta = 2^{-7}$, define*

$$\begin{aligned} f_{\pm} &= f(t \pm \varepsilon), \quad \Delta = f_+ - f_-, \\ E_1 &= \frac{\text{fl}(f_+) - \text{fl}(f_-)}{2\varepsilon}, \quad E_2 = \text{fl}\left(\frac{f_+}{2\varepsilon}\right) - \text{fl}\left(\frac{f_-}{2\varepsilon}\right). \end{aligned}$$

Suppose in addition that

- (1) $|\Delta| < 2^{-126}$, so that Δ (and any nearby perturbation) lies in the BF16 subnormal range;
- (2) writing $\text{fl}(f_{\pm}) = f_{\pm}(1 + \delta_{\pm})$ with $|\delta_{\pm}| \leq \eta$, one has $|f_+\delta_+ - f_-\delta_-| < 2^{-126}$, so $\tilde{f}_+ - \tilde{f}_-$ remains subnormal;
- (3) $|f_{\pm}/(2\varepsilon)| \geq 2^{-126}$, so each product $f_{\pm}/(2\varepsilon)$ lies in the normalized range;
- (4) $|f_+| + |f_-| = O(|\Delta|)$, so that any rounding in the two multiplications is not amplified by a large subtraction.

Then the “subtract-then-scale” formula E_1 may incur a relative error of order $O(1)$, whereas the “scale-then-subtract” formula E_2 retains a relative error of order $O(\eta)$.

Proof. We use two BF16 rounding models: (i) if $x \in [2^{-126}, 2^{128})$ then $\text{fl}(x) = x(1 + \delta)$, $|\delta| \leq \eta$; (ii) for any x (including subnormals), $|\text{fl}(x) - x| \leq \frac{1}{2} \text{ulp}(x)$, where $\text{ulp}_{\text{sub}} = 2^{-133}$ for subnormals.

Set $\tilde{f}_{\pm} = \text{fl}(f_{\pm}) = f_{\pm}(1 + \delta_{\pm})$, $|\delta_{\pm}| \leq \eta$.

Error in E_1 . By (1) and (2), $\tilde{f}_+ - \tilde{f}_- = \Delta + (f_+\delta_+ - f_-\delta_-)$ lies in the subnormal range. Hence

$$d = \text{fl}(\tilde{f}_+ - \tilde{f}_-) = (\tilde{f}_+ - \tilde{f}_-) + e_d, \quad |e_d| \leq \frac{1}{2} \text{ulp}_{\text{sub}} = 2^{-134}.$$

Thus

$$d = \Delta + (f_+\delta_+ - f_-\delta_-) + e_d, \quad |e_d|/|\Delta| = O(2^{-134}/|\Delta|)\mathbf{g}\eta.$$

Dividing by 2ε and rounding gives

$$E_1 = \text{fl}(d/(2\varepsilon)) = \frac{d}{2\varepsilon}(1 + \delta_q), \quad |\delta_q| \leq \eta,$$

so the relative error in E_1 can be $O(1)$.

Error in E_2 . By (3), each $f_{\pm}/(2\varepsilon)$ is normalized, so

$$g_{\pm} = \text{fl}\left(\frac{f_{\pm}}{2\varepsilon}\right) = \frac{f_{\pm}}{2\varepsilon}(1 + \delta'_{\pm}), \quad |\delta'_{\pm}| \leq \eta.$$

Subtracting and rounding (still normalized) gives

$$E_2 = \text{fl}(g_+ - g_-) = (g_+ - g_-)(1 + \delta'_d), \quad |\delta'_d| \leq \eta.$$

Since

$$g_+ - g_- = \frac{\Delta}{2\varepsilon} + \frac{f_+\delta'_+ - f_-\delta'_-}{2\varepsilon},$$

we obtain

$$E_2 = \frac{\Delta}{2\varepsilon}(1 + \delta'_d) + \frac{f_+\delta'_+ - f_-\delta'_-}{2\varepsilon}(1 + \delta'_d).$$

The second term has magnitude $\leq \eta \frac{|f_+|+|f_-|}{2\varepsilon}(1 + \eta)$, and by (4) its relative size to $\Delta/(2\varepsilon)$ is $O(\eta \frac{|f_+|+|f_-|}{|\Delta|}) = O(\eta)$.

Hence E_1 may suffer $O(1)$ relative error, while E_2 attains $O(\eta)$ relative accuracy under (1)–(4). \square

F.2.4 CALCULATION OF TRANSPORT

Transport transformation from EDM to UCGM. Take the formula (8) from EDM (Karras et al., 2022). With $\sigma_{\text{data}} = \frac{1}{2}$ and $\mathbf{n} = \sigma\mathbf{z}$, we can deduce:

$$\begin{aligned} & \mathbb{E}_{\sigma, \mathbf{x}, \mathbf{n}} \left[\lambda(\sigma) c_{\text{out}}(\sigma)^2 \left\| \mathbf{F}_\theta(c_{\text{in}}(\sigma) \cdot (\mathbf{x} + \mathbf{n}); c_{\text{noise}}(\sigma)) - \frac{1}{c_{\text{out}}(\sigma)} (\mathbf{x} - c_{\text{skip}}(\sigma) \cdot (\mathbf{x} + \mathbf{n})) \right\|_2^2 \right] \\ &= \mathbb{E}_{\sigma, \mathbf{x}, \mathbf{z}} \left[\left\| \mathbf{F}_\theta \left(\frac{1}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} (\mathbf{x} + \sigma\mathbf{z}); \frac{1}{4} \ln \sigma \right) - \frac{\sqrt{\sigma_{\text{data}}^2 + \sigma^2}}{\sigma\sigma_{\text{data}}} \left(\mathbf{x} - \frac{\sigma_{\text{data}}^2}{\sigma^2 + \sigma_{\text{data}}^2} (\mathbf{x} + \sigma\mathbf{z}) \right) \right\|_2^2 \right] \\ &= \mathbb{E}_{\sigma, \mathbf{x}, \mathbf{z}} \left[\left\| \mathbf{F}_\theta \left(\frac{1}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} (\mathbf{x} + \sigma\mathbf{z}); \frac{1}{4} \ln \sigma \right) - \frac{\sqrt{\sigma_{\text{data}}^2 + \sigma^2}}{\sigma\sigma_{\text{data}}} \left(\frac{\sigma^2}{\sigma^2 + \sigma_{\text{data}}^2} \mathbf{x} - \frac{\sigma_{\text{data}}^2}{\sigma^2 + \sigma_{\text{data}}^2} \sigma\mathbf{z} \right) \right\|_2^2 \right] \\ &= \mathbb{E}_{\sigma, \mathbf{x}, \mathbf{z}} \left[\left\| \mathbf{F}_\theta \left(\frac{1}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} (\mathbf{x} + \sigma\mathbf{z}); \frac{1}{4} \ln \sigma \right) - \left(\frac{\sigma}{\sigma_{\text{data}} \sqrt{\sigma^2 + \sigma_{\text{data}}^2}} \mathbf{x} - \frac{\sigma_{\text{data}}}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} \mathbf{z} \right) \right\|_2^2 \right] \\ &= \mathbb{E}_{\sigma, \mathbf{x}, \mathbf{z}} \left[\left\| \mathbf{F}_\theta \left(\frac{\sigma}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} \mathbf{z} + \frac{1}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} \mathbf{x}; \frac{1}{4} \ln \sigma \right) - \left(-\frac{\sigma_{\text{data}}}{\sqrt{\sigma^2 + \sigma_{\text{data}}^2}} \mathbf{z} + \frac{\sigma}{\sigma_{\text{data}} \sqrt{\sigma^2 + \sigma_{\text{data}}^2}} \mathbf{x} \right) \right\|_2^2 \right] \\ &= \mathbb{E}_{\sigma, \mathbf{x}, \mathbf{z}} \left[\left\| \mathbf{F}_\theta \left(\frac{\sigma}{\sqrt{\sigma^2 + \frac{1}{4}}} \mathbf{z} + \frac{1}{\sqrt{\sigma^2 + \frac{1}{4}}} \mathbf{x} \right) - \left(-\frac{1/2}{\sqrt{\sigma^2 + \frac{1}{4}}} \mathbf{z} + \frac{2\sigma}{\sqrt{\sigma^2 + \frac{1}{4}}} \mathbf{x} \right) \right\|_2^2 \right]. \end{aligned}$$