

---

# What does gradient variance tell us about the optimality of vector fields?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We investigate the underlying cause of common failure modes behind the state-  
2 of-the-art fixed-point generative models, Rectified Flows and Schrödinger Bridge  
3 Matching. First, we introduce a gradient-variance diagnostic that directly measures  
4 how learned vector fields induce transport errors after integration. Contrary to  
5 common belief, we show that low-variance (i.e., well-aligned) interpolants can  
6 still yield high transport error. We also demonstrate that standard neural network  
7 architectures (e.g., MLPs or CNNs) are fundamentally limited, failing to exactly  
8 represent even Gaussian-to-Gaussian transport. While prior work on Rectified  
9 Flows notes that strict convergence to optimal couplings is not guaranteed, it  
10 overlooks a key limitation: with noiseless interpolants, further iterations do not  
11 improve transport. If the pairings are already straight, the process stagnates; if  
12 they are deterministic but not straight, the vector field simply memorizes these  
13 pairings without achieving true optimal transport. We prove the existence of such  
14 memorizing vector fields.

## 15 1 Introduction

16 The current state of the art in generative modeling involves learning dynamics between a known  
17 source distribution—such as a standardized Gaussian—and a target distribution from which many  
18 samples are available. When the learned dynamics are based on an ordinary differential equation  
19 (ODE), the model effectively learns a vector field connecting the two distributions (Lipman et al.,  
20 2022; Tong et al., 2023b,a). When the learned dynamics are based on a stochastic differential equation  
21 (SDE), either the score function (Ho et al., 2020; Song et al., 2020), or both the score and the vector  
22 field (Albergo et al., 2023), are learned, depending on the choice of SDE. Ultimately, an SDE can be  
23 rewritten as an ODE—referred to in the literature as the *probability flow ODE* (Song et al., 2020).

24 These methods perform very well in practice across a wide range of data modalities: images (Ho  
25 et al., 2022a; Balaji et al., 2022; Rombach et al., 2022), video (Ho et al., 2022b; Wang et al., 2024;  
26 Zhou et al., 2022), audio (Huang et al., 2023; Kong et al., 2020; Liu et al., 2023; Ruan et al., 2023),  
27 and molecular data (Hooeboom et al., 2022; Xu et al., 2022). However, their main limitations are the  
28 high computational cost of the generation of samples (as integration over an ODE is required) (Song  
29 et al., 2020; Tong et al., 2023b; Albergo et al., 2023), and their inability to learn optimal transport  
30 maps, which are often essential for unpaired data translation tasks (Shi et al., 2024).

31 To address computational inefficiency, new models have been proposed. Consistency models acceler-  
32 ate sampling for score-based approaches (Song et al., 2023; Salimans & Ho, 2022; Kim et al., 2023),  
33 while Rectified Flows (Liu, 2022; Roy et al., 2024; Lee et al., 2024) aim to learn straight vector fields  
34 that can be integrated in a single step. Rectified Flows iteratively update couplings and retrain vector  
35 fields to straighten transport paths, but repeated rectifications can accumulate errors, and it remains  
36 theoretically unclear whether one or two rectifications suffice under general conditions. Input-Convex

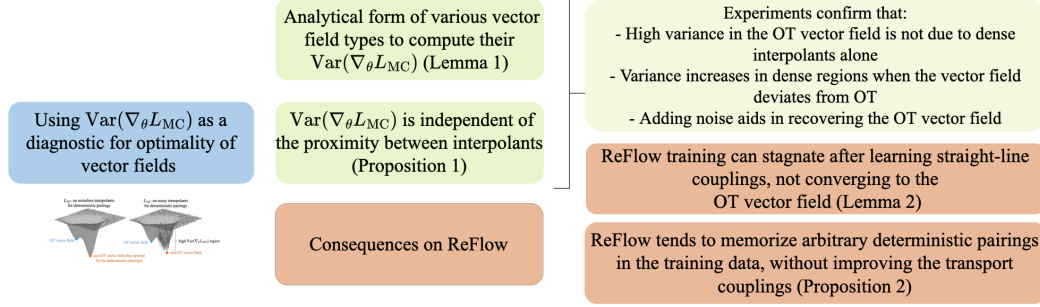


Figure 1: Flow chart of our main results.

Neural Networks (ICNN)-based parameterizations can, in theory, guarantee optimal couplings for noiseless interpolants, but are difficult to optimize in practice (Makkuva et al., 2020; Huang et al., 2020).

Schrödinger Bridge Matching (SBM) (Shi et al., 2024; Peluchetti, 2023; De Bortoli et al., 2024) extends these ideas by learning both forward and backward vector fields using noisy interpolants, approximating entropic optimal transport. This bidirectional approach can improve stability and mitigate error accumulation, but requires training two networks per iteration.

In this paper, we address fundamental limitations of iterative generative models by analyzing how gradient variance reveals suboptimality in learned vector fields. Our theoretical and empirical study uncovers why standard neural architectures struggle to represent even simple transports and how repeated rectifications can lead to memorization rather than improvement.

#### Contributions:

- We introduce a gradient-variance diagnostic for evaluating the optimality of learned vector fields.
- We prove that zero-loss is unattainable for Gaussian transport in Lemma 1 and bound variance for various types of pairings (optimal, straight, and random) in Proposition 1.
- We prove ReFlow with noiseless interpolants stagnate when arriving at straight couplings in Lemma 2, and *memorize* when trained on deterministic couplings in Proposition 2.
- We empirically validate our findings on synthetic and CIFAR-10 data.

## 2 Background

**Notation.** Let  $\mathbb{R}^d$  denote  $d$ -dimensional Euclidean space. We use  $\mathcal{P}(\mathbb{R}^d)$  for the set of probability distributions on  $\mathbb{R}^d$ . The source and target distributions are  $\pi_0$  and  $\pi_1$ , with  $X_0 \sim \pi_0$  and  $X_1 \sim \pi_1$ . A transport map is  $T : \mathbb{R}^d \rightarrow \mathbb{R}^d$ , and  $T_{\#}\pi_0$  denotes the pushforward of  $\pi_0$  by  $T$ . Interpolants are written  $X_t = I(X_0, X_1, t)$ , with  $t \in [0, 1]$ . We write  $\mathbb{E}[\cdot]$  for expectation,  $\text{Var}[\cdot]$  for variance, and  $\nabla_{\theta}$  for gradients with respect to parameters  $\theta$ . Bold symbols (e.g.,  $\mathbf{x}$ ) denote vectors. All other notation is defined in context. For  $k \in \mathbb{R}$ ,  $R_{k^{\circ}}$  denotes a  $d \times d$  rotation matrix corresponding to a counterclockwise rotation by  $k$  degrees in the plane of interest. All other notation is defined in context.

**Optimal Transport and Entropy-Regularized OT** Let  $\pi_0, \pi_1$  be probability measures on  $\mathbb{R}^d$ . The Monge problem (Monge, 1781) seeks a transport map  $T$  minimizing:

$$\inf_T \int \|T(x) - x\|^2 d\pi_0(x) \quad \text{s.t.} \quad T_{\#}\pi_0 = \pi_1, \quad (1)$$

where  $T_{\#}\pi_0$  denotes the pushforward of  $\pi_0$  under  $T$ . The Kantorovich (Kantorovitch, 1958) relaxation introduces couplings  $\pi \in \Pi(\pi_0, \pi_1)$  and solves:

$$\mathcal{W}_2^2(\pi_0, \pi_1) = \inf_{\pi \in \Pi(\pi_0, \pi_1)} \mathbb{E}_{(X_0, X_1) \sim \pi} [\|X_1 - X_0\|^2], \quad (2)$$

where  $\Pi(\pi_0, \pi_1)$  denotes the set of joint distributions with marginals  $\pi_0$  and  $\pi_1$ . Entropy-regularized Optimal Transport (eOT) adds a Kullback–Leibler divergence penalty  $\mathcal{W}_2^\epsilon(\pi_0, \pi_1) = \inf_{\pi \in \Pi(\pi_0, \pi_1)} \mathbb{E}_\pi[\|X_1 - X_0\|^2] + \epsilon D_{\text{KL}}(\pi \| \pi_0 \otimes \pi_1)$ .

When  $\pi_0$  is absolutely continuous, the Monge and Kantorovich problems admit the same deterministic optimal plan. A dynamic formulation describes OT as evolving a path  $\{X_t\}_{t \in [0,1]}$  connecting  $X_0 \sim \rho_0$  and  $X_1 \sim \rho_1$ . For convex costs, the optimal path is given by the straight-line interpolant  $X_t = (1-t)X_0 + tX_1$  (McCann, 1997).

**Conditional Flow Matching** Given the interpolants  $X_t = I(X_0, X_1, t) = \alpha_t X_0 + \beta_t X_1 + \gamma_t \epsilon$  where  $\epsilon \sim \mathcal{N}(0, I)$ , the CFM objective Lipman et al. (2022) is:

$$\mathcal{L}_{CFM}(v) = \mathbb{E}_{t \sim \mathcal{U}(0,1), (X_0, X_1) \sim \pi, \epsilon} [\|(X_1 - X_0) - v(X_t, t)\|^2] \quad (3)$$

The learned vector field  $v$  generates flows via the ODE:

$$\frac{d}{dt} X_t = v(X_t, t) \quad (4)$$

**Rectified Flows (or ReFlow)** Liu et al. (2022) propose an iterative procedure to straighten transport paths. At each iteration  $k$ , a vector field  $v^{(k)}$  is trained using the CFM loss  $\mathcal{L}_{CFM}$  on the current coupling  $(X_0^{(k)}, X_1^{(k)})$ . The updated coupling is then generated via  $X_1^{(k+1)} = \text{ODE-Solve}[v^{(k)}](X_0^{(k)}, t=1)$ , and the process repeats until convergence. A coupling  $(X_0^{(k)}, X_1^{(k)})$  is considered straight if

$$\mathbb{E} [X_0^{(k)} - X_1^{(k)} \mid X_t^{(k)} = x] = X_0^{(k)} - X_1^{(k)}. \quad (5)$$

where  $X_t^{(k)} = (1-t)X_0^{(k)} + tX_1^{(k)}$ . A coupling is considered straight if, for deterministic couplings  $(X_0, X_1)$ , the mapping  $(X_t, t) \mapsto (X_0, X_1)$  is injective. Although some works claim one or two rectifications suffice to obtain straight couplings (Lee et al., 2024) (through one iteration) (Roy et al., 2024) (through two iterations, though the paper got retracted later on), counterexamples (see Counter-Example 1) shows this is not guaranteed under basic assumptions, after one rectification.

**Schrödinger Bridge Matching** The Schrödinger Bridge seeks the entropic interpolation between  $\pi_0$  and  $\pi_1$ :

$$\inf_{v^+, v^- : \mathbb{P}_0^{v^+, v^-} = \pi_0, \mathbb{P}_1^{v^+, v^-} = \pi_1} D_{KL}(\mathbb{P}^{v^+, v^-} \| \mathbb{P}^{\text{ref}}) \quad (6)$$

where  $\mathbb{P}^{\text{ref}}$  is a reference process. SBM Shi et al. (2024) learns bidirectional vector fields via:

$$\mathcal{L}_{SBM} = \mathbb{E}_t [\|v^+(X_t, t) - (X_1 - X_0)\|^2 + \|v^-(X_t, t) - (X_0 - X_1)\|^2]. \quad (7)$$

We will refer to one iteration of SBM in only one direction as CFM+SI (CFM and stochastic interpolants), and CFM (when no noise is added to the interpolants).

### 3 Gradient Variance Analysis

**Why study gradient variance?** Analyzing  $\text{Var}[\nabla_\theta L_{\text{MC}}(v, T, I)]$  across different *choices of pairings*  $T$  (e.g. random vs. optimal vs. straight), *interpolants*  $I$  (noiseless vs. stochastic), and *vector-field classes*  $v_\theta$  provides insight into which solutions are favored by the loss landscape. Although a loss may have multiple minima, optimization often prefers those with lower gradient variance, even if they are suboptimal in terms of transport cost. Figure 2 illustrates this effect: for deterministic couplings and noiseless interpolants, the lowest variance and loss minimizer can be non-OT, while introducing stochasticity can shift preference toward more optimal solutions.

**Surprising implications for interpreting gradient variance.** When the vector field is OT, straight pairings can exhibit higher gradient variance than random pairings (see Figure 4). Conversely, with a non-OT vector field, the loss can display substantial variance for pairings that are themselves OT (see Figure 5). Furthermore, for straight non-OT pairings, the most stable solution is the corresponding optimal (non-OT) vector field, Figure 6.

**Clarification.** We have vector fields that are optimal in the OT sense, and others that are optimal for a certain type of pairing (yielding the smallest error values); we will refer to the latter as optimal non-OT. The proof for the following statements can be found in Appendix B.

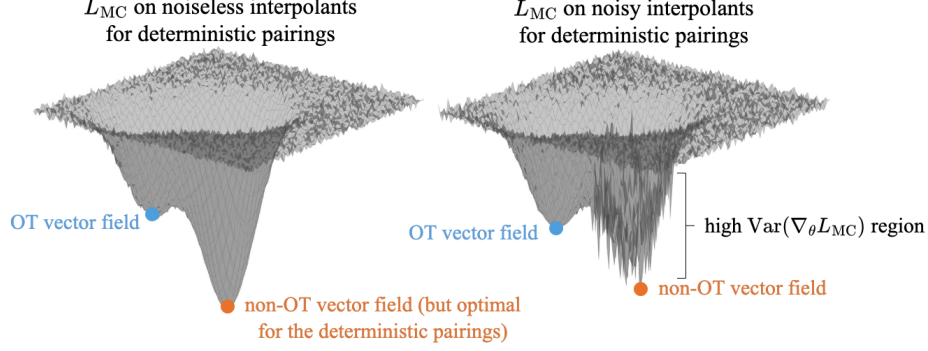


Figure 2: Schematic representing a hypothetical loss ( $L_{MC}$ ) landscape. The gradient variance of the loss acts as an indicator of solution quality. The schematic illustrates how the variance of the loss gradient reveals information about the optimality of the vector field under different interpolant types.

**Definition** We begin by rewriting the CFM objective from Equation 3:

$$L(v, T, I) = \mathbb{E}_{X_0 \sim \pi_0} \left[ \|T(X_0) - X_0 - v(X_t, t)\|^2 \right], \quad X_t = I(X_0, T(X_0), t), \quad (8)$$

where  $T$  is a (possibly random or deterministic) map satisfying  $T_{\#} \pi_0 = \pi_1$ . A Monte Carlo approximation of this loss, using samples  $\{X_0^{(s)}\}_{s=1}^S$ , is given by:

$$L_{MC}(v, T, I) = \frac{1}{S} \sum_{s=1}^S \left\| T(X_0^{(s)}) - X_0^{(s)} - v(X_t^{(s)}, t) \right\|^2, \quad X_t^{(s)} = I(X_0^{(s)}, T(X_0^{(s)}), t). \quad (9)$$

where  $I(X_0, X_1, t)$  will be the interpolant of our choice.

The gradient variance will have the following formulation:

$$\text{Var}[\nabla_{\theta} L_{MC}(v, T, I)] = \text{Var} \left[ \nabla_{\theta} \frac{1}{S} \sum_{s=1}^S \left\| T(X_0^{(s)}) - X_0^{(s)} - v(X_t^{(s)}, t) \right\|^2 \right], \quad (10)$$

where  $X_t^{(s)} = I(X_0^{(s)}, T(X_0^{(s)}), t)$

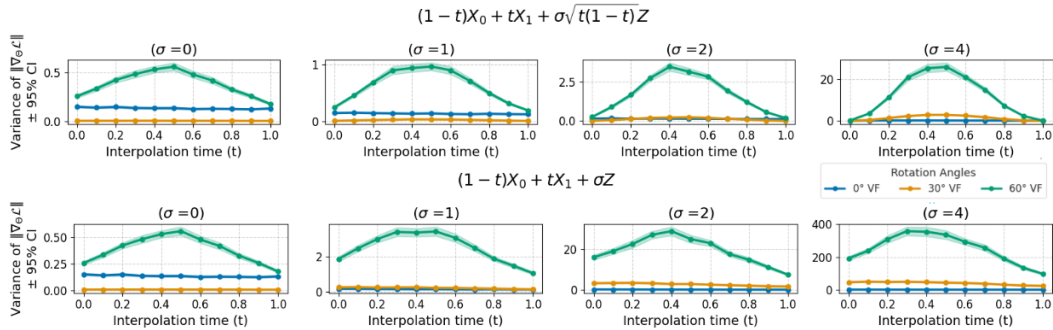


Figure 3: Gradient variance (with 95% confidence intervals) for Gaussian transport paired  $(X_0, R_{30^\circ} X_0 + \mu)$  under different noise levels ( $\sigma$  of two interpolants mentioned in the title of the figures). We compare vector fields inducing  $0^\circ$  (OT, blue),  $30^\circ$  (non-OT, optimal for the pairing, orange), and  $60^\circ$  (non-OT, non-optimal, green) rotations. Shaded regions show confidence intervals calculated over 100 bootstrap samples. As noise increases ( $\sigma = 0 \rightarrow 4$ ), the optimal transport (OT) field exhibits significantly reduced variance ( $p < 0.01$ , paired t-test) while maintaining lower transport error. Confidence intervals narrow with higher  $\sigma$ , indicating more stable gradient estimates despite increased stochasticity.



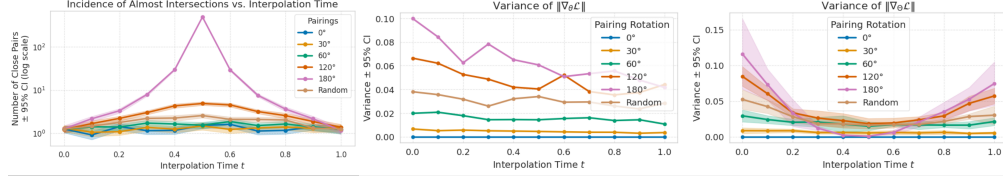


Figure 4: Gradient variance over time for several pairing types:  $(X_0, \Sigma^{1/2} \mathbf{R}_{k^\circ} X_0 + \mu)$  with  $k \in \{0, 30, 60, 120, 180\}$ , and  $(X_0, X_1)$  with  $X_1 \sim \mathcal{N}(\mu, \Sigma)$ , over the OT vector field. Shaded regions indicate 95% confidence intervals computed via 100 bootstrap samples. Notably, variance does not increase in regions where interpolant trajectories come closer together: for example, with  $180^\circ$  rotation, all interpolants intersect at  $t = 1/2$ , yet variance is lowest there. Random pairings exhibit lower variance under the OT field than structured pairings with  $120^\circ$  (straight) or  $180^\circ$  (not-straight) rotations, highlighting that variance is not simply a function of interpolant density.

### 3.1 Gaussian setting

We begin our analysis of gradient variance in the loss landscape by considering the analytically tractable case of Gaussian-to-Gaussian optimal transport (see experimental details in Appendix E, and proofs in Appendix B). This setting allows us to derive the exact OT vector field and study its properties in detail. To ensure the precision of our empirical and theoretical analysis, we use a tailored parameterization of this vector field, since standard architectures such as MLPs, CNNs, and Transformers fail to represent it exactly. While neural networks are universal approximators and can, in theory, represent matrix inverse series to arbitrary precision, in our setting, it is more informative to analyze the structure and gradient variance of the true optimal parameterization.

In Lemma 1, we derive this closed-form expression for the optimal vector field between two Gaussian distributions with different means and covariances, and we show that a multilayer perceptron (MLP) cannot approximate it without error.

**Lemma 1.** *Let  $X_0 \sim \mathcal{N}(0, \mathbf{I}_d)$  and  $X_1 \sim \mathcal{N}(\mu, \mathbf{M}_d)$ , where  $\mathbf{M}_d$  is a positive definite and symmetric matrix. The OT vector field is given by*

$$\hat{v}_{OT}(X_t, t, \hat{\Theta}, \hat{\theta}) = \hat{\theta} + \hat{\Theta}[\mathbf{I}_d + t\hat{\Theta}]^{-1}(X_t - t\hat{\theta}),$$

and the rotating vector field:

$$\hat{v}_{rOT}(X_t, t, \hat{\Theta}_{rOT}, \hat{\theta}) = \hat{\theta} + \hat{\Theta}_{rOT}[\mathbf{I}_d + t\hat{\Theta}_{rOT}]^{-1}(X_t - t\hat{\theta}),$$

where  $X_t = (1-t)X_0 + tX_1$ ,  $\hat{\Theta} = \mathbf{M}_d^{1/2} - \mathbf{I}_d$ ,  $\hat{\Theta}_{rOT} = \mathbf{M}_d^{1/2} \mathbf{R}_{k^\circ} - \mathbf{I}_d$ ,  $\hat{\theta} = \mu$ , and  $\mathbf{R}_{k^\circ}$  is a rotation matrix with  $\mathbf{R}_{k^\circ} \neq -\mathbf{R}_{k^\circ}$ . Furthermore, the function  $\hat{v}_{OT}$  cannot be exactly represented by an MLP, CNN, Transformer architecture when given concatenated inputs  $[X_t, t]$ .

An intuitive reason why this vector field cannot be represented with zero error by typical neural network parameterizations is the difficulty in capturing the matrix inverse term  $[\mathbf{I} + t\hat{\Theta}]^{-1}$ . This inverse can be expressed as an infinite series:  $\mathbf{I} - t\hat{\Theta} + (t\hat{\Theta})^2 - \dots$ , which standard networks struggle to represent precisely. For a complete proof, see Appendix B. Motivated by this, we experimented with augmenting the network input using higher powers of  $t$  on CIFAR-10, essentially giving the model a polynomial handle on time. While sinusoidal time embeddings yielded the lowest gradient variance, using only the linear term in  $t$  led to the best FID (see Appendix F.2 for full results).

**Under noiseless interpolant:** We now characterize the variance of the gradients over the noiseless interpolant  $x_t = (1-t)x_0 + tx_1$ , as this is the commonly used type of interpolant in ReFlow architectures.

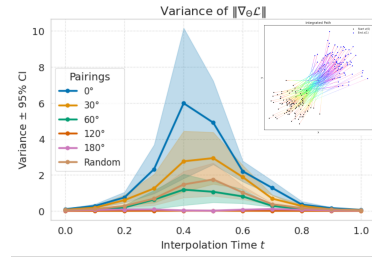


Figure 5: Gradient variance over time for a vector field rotating by  $120^\circ$ , under various pairing types including rotated and random couplings. In the top right, the two Gaussians are shown along with sample trajectories, color-coded by integration time. Optimal pairings exhibit the highest variance, while random pairings do not yield the maximum variance.

**Proposition 1** (Variance of transport between gaussian). Let  $\pi_0 \sim \mathcal{N}(0, \mathbf{I}_d)$ , and  $\pi_1 \sim \mathcal{N}(\mu, \mathbf{M}_d)$ , where  $\mathbf{M}_d$  is a positive definite and symmetric matrix. Let the following push forward maps:  $T_{OT}(X_0) = \mathbf{M}_d^{1/2}X_0 + \mu$ ,  $T_{rOT}(X_0) = \mathbf{M}_d^{1/2}\mathbf{R}X_0 + \mu$ , and  $T_{rand}(X_0) = X_1$ , where  $\mathbf{R}$  is a rotation matrix. Let the linear transformation be  $v(X_t, t) = \theta + \Theta[\mathbf{I}_d + t\Theta]^{-1}(X_t - t\theta)$ , where  $X_t = (1 - t)X_0 + tX_1$ . Let:

$$L_{MC}(\Theta, \theta, T, v) = \frac{1}{N} \sum_{s=1}^N \|T(X_0^{(s)}) - X_0^{(s)} - v(X_t^{(s)}, t, \Theta, \theta)\|^2. \quad (11)$$

Then, for optimal,  $(\hat{\Theta}, \hat{\theta})$ , and rotating  $(\hat{\Theta}_{rOT}, \hat{\theta})$  we have  $\text{Var}[\nabla_{\theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, T_{OT})] = \text{Var}[\nabla_{\Theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, T_{OT})] = \text{Var}[\nabla_{\theta}(\hat{\Theta}, \hat{\theta}_{rOT}, v_{rOT}, T_{rOT})] = \text{Var}[\nabla_{\Theta} L_{MC}(\hat{\Theta}, \hat{\theta}_{rOT}, v_{rOT}, T_{rOT})] = 0$ , and  $\text{Var}[\nabla_{\theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, v_{OT})] = \frac{4}{N} \|\mathbf{M}_d(\mathbf{R} - \mathbf{I})\|^2$ , and  $\text{Var}[\nabla_{\Theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, v_{OT})] = \frac{8}{N} \text{tr}((\mathbf{M}_d^{1/2}(\mathbf{R} - \mathbf{I}))^2)$ . Moreover, for the maps  $T_{rOT}$  and  $T_{rand}$  and  $v_{OT}$  vector field, the variance of the gradients of the optimal vector field does not increase in regions where the interpolant trajectories come closer together.

Some takeaways from this proposition are that low gradient variance does not necessarily indicate transport optimality, especially for straight or rotated pairings, as can be observed from our empirical findings. For example, when  $v$  and  $T$  are aligned (both rotate by the same degree), the loss and its gradients will be zero. Additionally, while straight pairings are known to have zero variance, under the optimal vector, they will exhibit variance in the gradients (see in Figures 3, 5, and 6).

**Remark 1.** Once Rectified Flows obtain straight pairings, the training objective favors the corresponding straight vector field—i.e., the one that satisfies the pairing exactly—regardless of whether it corresponds to optimal transport. As a result, the model may converge to a suboptimal solution that perfectly fits the pairings but does not minimize transport cost.

**Variance under stochastic interpolants:** In the noiseless case, once we reach the correct straight pairing, there exists a vector field that achieves zero loss on the optimization problem. However, with a noisy interpolant, the situation changes significantly. As shown in Figures 3, adding noise to the interpolant causes the optimal vector field to exhibit lower gradient variance than the vector field optimal for the (suboptimal) fixed straight pairing.

To conclude this section, we highlight the following key takeaways:

- **(A)** Variance does not arise from intersections or regions where the interpolating lines come close together, and **(B)** random pairings exhibit lower variance than structured couplings, such as a  $120^\circ$  rotation (see Figure 4).
- **(A)** The variance of the gradients can be high in regions where interpolants come close if the learned vector field is suboptimal, and **(B)** even a vector field that rotates by  $180^\circ$  (causing all interpolants to meet at  $t = 1/2$ ) can still be integrated over, since our integration methods are discretized and effectively jump over the intersection point (see Figure 5).
- Adding noise to the interpolants aids in avoiding the learned vector field collapsing to a straight/deterministic pairing structure (see Figure 3, and Figure 6).

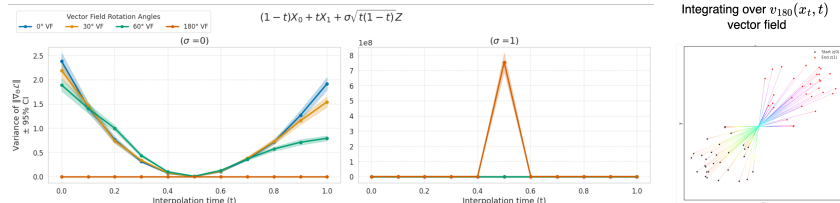


Figure 6: Integration over a vector field performing a  $180^\circ$  rotation. Although all interpolation paths intersect at  $t = 1/2$ , and  $\mathbb{E}[X_1 - X_0 | X_t = x_{1/2}] \neq X_1 - X_0$ , we still recover  $180^\circ$ -rotated couplings. This is because the numerical integration procedure discretizes time and does not evaluate the vector field precisely at the point of intersection, effectively bypassing it.

Table 1: Comparison of FID, NLL, and sample variance for CFM+SI (first iteration of SBM) and CFM (Forward and Backward). The backward vector field exhibits approximately 10 times lower variance in integrated samples compared to the forward vector field, despite both endpoints being standardized. This highlights an inherent asymmetry in sampling stability between forward and backward flows.

	Direction	FID↓	Scaled -NLL↓	Variance $\pm$ Std. Error
<b>CMF+SI</b>	Backward	–	1.423564	$0.01647 \pm 0.00024$
	Forward	4.250	–	$0.36179 \pm 1.23 \times 10^{-5}$
<b>CFM</b>	Backward	–	1.426562	$0.01414 \pm 0.00018$
	Forward	4.199	–	$0.36177 \pm 1.23 \times 10^{-5}$

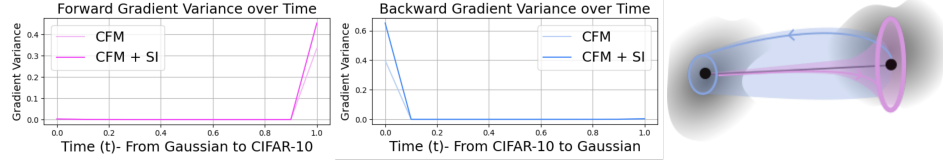


Figure 7: Comparison of gradient variance for the forward and backward passes in CFM and CFM+SI (first iteration of SBM). The schematic on the right offers intuition for how gradient variance may influence sample variance during integration. Notably, the backward-pass gradient variance peaks near the CIFAR-10 endpoint, yet the resulting sample variance is lower compared to the forward pass. This is possible, for instance, if the forward field increases linearly while the backward field decreases linearly. These effects are quantitatively reflected in Table 4.

### 3.2 General Case: First iteration of SBM

To test whether our gradient-variance diagnostic extends beyond Gaussians, we trained four U-Net models to learn the transport map from  $\pi_0 = \mathcal{N}(0, I_d)$  to  $\pi_1 \sim \text{CIFAR-10}$ , using both noiseless and noisy interpolants (see for experimental details Appendix F). In all settings (Figure 7), the backward field’s gradient variance peaks at  $t = 0$ , where CIFAR-10 appears as input. Conversely, the forward field’s variance peaks at  $t = 1$ , when CIFAR-10 appears as the output.

To assess generation stability, we repeatedly integrated 64 fixed datapoints across 100 training checkpoints, sampling at regular intervals near the FID/NLL optima (Table 1). We report the resulting sample variances in the table’s fourth column:

$$\text{Var}[x_1^{\text{gen}}] \quad \text{and} \quad \text{Var}[x_0^{\text{gen}}],$$

which show consistently higher variability in the forward direction, suggesting reduced robustness during forward-time sampling.

**Observation 1.** *The observed difference in sample variance across vector fields optimized near the optima may potentially impact subsequent iterations of SBM (Shi et al., 2024). However, precisely characterizing how this affects theoretical bounds is non-trivial and is left as a promising avenue for future research.*

## 4 Rectified Flow Dynamics: Minimizer Memorizes

### 4.1 Theoretical Results

So far, we have seen that noiseless interpolants let ReFlow converge training to straight, but not necessarily optimal, couplings, with zero gradient variance. We now formalize why ReFlow can both (A) stagnate once it hits a straight coupling (Lemma 2), and (B) memorize arbitrary deterministic pairings when trained on a finite dataset (Proposition 2).

**Lemma 2** (Idempotence of Rectified Flows). *Let  $\pi_0, \pi_1$  be distributions on  $\mathbb{R}^D$  with densities, and  $(Z_0, Z_1)$  their straight-line coupling via linear interpolant  $Z_t = (1 - t)Z_0 + tZ_1$ . Subsequent Rectified Flow iterations with this noise-free interpolant yield identical couplings:*

$$\text{ReFlow}^{(k)}(Z_0, Z_1) = (Z_0, Z_1) \quad \forall k \geq 1$$

This stagnation stems from the bijective relationship between interpolants  $(Z_t, t)$  and initial pairs  $(Z_0, Z_1)$ , given by our assumption of the straightness of the deterministic couplings.

Some recent works have suggested that 1-Reflow might be sufficient to recover straight pairings Lee et al. (2024), and a now-retracted claim Roy et al. (2024) proposed that 2-Reflow is enough. However, no formal proof has been established to date. Under mild assumptions (see Appendix A), we provide a simple counter example demonstrating that 1-Reflow is not sufficient to guarantee straight paths.

**Counter Example 1** (1-Reflow May Fail Under Mild Assumptions). *Even if the learned transport map  $T(x_0)$  is injective (e.g., under standard Lipschitz and growth conditions; see Appendix A), the straight-line interpolant  $I(x_0, T(x_0), t)$  need not be. For instance, a rotation-based map like  $T(x_0) = R_{180^\circ} x_0 + 5$  (which can be realized by a continuous vector field) leads to overlapping interpolants, breaking injectivity of  $(x_t, t) \mapsto x_0$ . As a result, a new ReFlow step cannot reconstruct  $(x_0, x_1)$  and fails to straighten the coupling. See Appendix D for details and a visual example.*

Relaxing the straightness assumption to consider general deterministic couplings and finite training datasets leads to a nuanced but equally critical limitation:

**Proposition 2** (The Minimizer Memorizes). *Let  $\pi_0$  and  $\pi_1$  be two probability densities on  $\mathbb{R}^D$ , and let  $T : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be a deterministic transport map pushing  $\pi_0$  to  $\pi_1$  (so if  $Z_0 \sim \pi_0$ , then  $Z_1 = T(Z_0) \sim \pi_1$ ). Suppose we draw a finite dataset of  $N$  i.i.d. samples  $\{Z_0^{(i)}, Z_1^{(i)} = T(Z_0^{(i)})\}_{i=1}^N$ , and for each  $i$  we also sample  $m$  time-points  $\{t^{(i,j)}\}_{j=1}^m \subset [0, 1]$  (e.g. uniformly). Let  $Z_t^{(i,j)} = (1 - t^{(i,j)}) Z_0^{(i)} + t^{(i,j)} Z_1^{(i)}$ . Define the empirical loss over this doubly indexed dataset by*

$$L_{MC}^{\det}(v_\Theta) = \frac{1}{Nm} \sum_{i=1}^N \sum_{j=1}^m \left\| (Z_1^{(i)} - Z_0^{(i)}) - v_\Theta(Z_t^{(i,j)}, t^{(i,j)}) \right\|^2.$$

Then there exists a (deterministic) vector field  $v$  attaining zero loss:  $L_{MC}^{\det}(v) = 0$ .

## 4.2 Impact of memorization

In this subsection, we discuss the implications of Lemma 2 and Proposition 2.

**Remark 2** (Condition for Recovering the Same Couplings). *Defining a look-up vector field over the training data—or approximating it via the Universal Approximation Theorem (UAT)—does not guarantee that integration will recover the original pairings. To ensure this, the integration process must traverse the specific time steps  $\{t^{(i,j)}\}_{j=1}^m$  associated with each training pair  $X_{(i)}$ . These are the only points along the interpolation  $x_t = (1 - t)x_0 + tx_1$  where we are guaranteed that  $v(x_t, t) = x_1 - x_0$ . While this condition might appear restrictive, the proposition remains valid for any natural number  $m$ .*

This insight arose from empirical observations. For example, when training a neural network on pairs  $(X_0, T(X_0))$  with  $T(X_0) = R_{180^\circ} X_0 + \mu$ , we found that—despite all interpolants intersecting at  $t = 1/2$ —the model successfully learned the rotation and preserved the pairings during integration. See Figure 8 for a visualization.

**Remark 3** (Noise Breaks Proof Assumptions). *A key advantage of using noisy interpolants is that they break the bijection between  $(x_t, t)$  and  $(x_0, T(x_0))$ , thereby violating the key steps of the proofs of Lemma 2 and Proposition 2. This disruption prevents the model from becoming stuck in suboptimal or deterministic straight pairings. As shown in Figure 3, this can lead to learning vector fields that are closer to the optimal solution.*

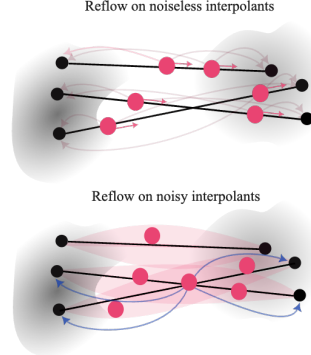


Figure 8: Schematic showing the two types of interpolants. When training with deterministic pairings, sampling at an intersection has zero probability mass. As a result, the mapping  $(Z_t, t) \mapsto (Z_0, T(Z_0), t)$  becomes injective, making it straightforward to define a vector field at these points by  $v(Z_t, t) = Z_0 - T(Z_0)$ . In contrast, for noisy interpolants, this injectivity no longer holds.

Table 2: Comparison of CFM and CFM with stochastic interpolants (CFM+SI) across low and high dimensions. CFM+SI consistently generates samples that better match the target distribution (as indicated by lower MMD and Sinkhorn distances), without resorting to memorization, as evidenced by its strong performance on both generated and true data metrics. This highlights CFM+SI’s superior generalization and sample quality compared to CFM. Experiments done over 10 seeds, the mean is presented, for full table check Appendix E.

Dimension		3				50			
		Gen	Mem	True	Data	Gen	Mem	True	Data
CFM	LogProb	4.0150	4.0890	4.1330	4.0155	54.8299	53.6502	52.5094	53.6244
	MMD	0.0034	$1.758 \times 10^{-6}$	0.0014	0.0032	0.0021	$9.089 \times 10^{-6}$	0.0020	0.0019
	Sinkhorn	0.0730	$1.411 \times 10^{-5}$	0.0637	0.0790	15.1900	0.0045	14.3221	15.7400
CFM+SI	LogProb	4.1270	4.0960	4.1330	4.0155	54.7220	53.8890	52.5094	53.6244
	MMD	0.0018	$3.105 \times 10^{-5}$	0.0014	0.0032	0.0020	$6.09 \times 10^{-5}$	0.0020	0.0019
	Sinkhorn	0.0680	$3.557 \times 10^{-4}$	0.0637	0.0790	15.1689	0.0304	14.3221	15.7400

**Remark 4** (Why Doesn’t CFM Memorize?). Although CFM does not explicitly use noisy interpolants, it avoids memorization for a different reason. At iteration 0,  $(x_0, x_1)$  are drawn independently from continuous distributions (e.g.,  $\pi_0 \sim \mathcal{N}(\mu_0, \Sigma)$  and  $\pi_1 \sim \mathcal{N}(\mu_1, \Sigma)$ ), so for dimensions  $D > 2$  interpolants don’t cross (see Proposition 3, in the Appendix C). However, the noise from  $x_0$  acts like the noise from the stochastic interpolants, breaking the bijection between  $(x_t, t)$  and  $(x_0, x_1)$  and invalidating the assumptions needed for the memorization argument to apply. In this sense, randomness in the pairings plays a similar role to the explicit noise described in Remark 3.

While Rectified Flows can help straighten trajectories, they risk collapsing onto memorized pairings at the first iteration, when deterministic couplings are formed. This memorization undermines generalization and highlights a limitation of relying solely on early deterministic transport.

### 4.3 Memorisation for CFM and CFM + SI

To probe CFM’s tendency to memorize deterministic pairings, we train two variants—one using noiseless interpolants and one with added noise—on an CFM model that maps  $\pi_0 = \mathcal{N}(0, I_d)$  to a target Gaussian mixture  $\pi_1$  (more details for the experiments in Appendix E). We evaluate in both low ( $d = 3$ ) and high ( $d = 50$ ) dimensions using three metrics (log-likelihood under the true mixture, MMD, and Sinkhorn distance) and four comparisons: (i) *Gen*: generated held-out vs. true samples, (ii) *Mem*: generated vs. training pair integrations, (iii) *True*: true vs. true reference, and (iv) *Data*: training pair vs. true samples. This setup isolates whether noiseless CFM simply memorizes its finite dataset, and whether stochastic interpolants improve generalization without sacrificing sample quality.

As shown in Table 4, CFM tends to memorize the deterministic pairings it is trained on, reflected in very low memorization distances but weaker generalization. In contrast, CFM+SI demonstrates consistently lower distances to the true distribution across both dimensions, suggesting superior generalization and less reliance on memorized pairings.

## 5 Conclusion

We analyzed flow-based models through the lens of gradient variance and showed that deterministic interpolants—especially in Rectified Flows—can lead to memorization and stagnation. Noise, whether explicit (SI) or implicit (random pairings), breaks this effect, improving both generalization and sample quality. Our theoretical results and experiments on Gaussians and CIFAR-10 demonstrate that variance-sensitive optimization prefers suboptimal flows unless stochasticity is introduced. These findings challenge the reliability of straight interpolants and underscore the importance of noise in guiding models toward better transport.

**Limitations.** Our experiments are limited to CIFAR-10 as the only real-world dataset; extending to more diverse or higher-resolution datasets may reveal new effects. We prioritized foundational settings to clarify misconceptions about interpolants and gradient variance. Our analysis of SBM remains preliminary, and further theoretical work is needed to understand its temporal asymmetries in its gradient variance as shown in Section 3.2.

## References

- Michael S Albergo, Nicholas M Boffi, and Eric Vanden-Eijnden. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Vladimir I Arnold. *Ordinary differential equations*. Springer Science & Business Media, 1992.
- Yogesh Balaji, Seungjun Nah, Xun Huang, Arash Vahdat, Jiaming Song, Karsten Kreis, Miika Aittala, Timo Aila, Samuli Laine, Jaakko Lehtinen, et al. ediff-i: Text-to-image diffusion models with an ensemble of expert denoisers. *arXiv preprint arXiv:2211.01324*, 2022.
- Valentin De Bortoli, Iryna Korshunova, Andriy Mnih, and Arnaud Doucet. Schrodinger bridge flow for unpaired data translation. *Advances in Neural Information Processing Systems*, 37: 103384–103441, 2024.
- Jean Feydy, Thibault Séjourné, François-Xavier Vialard, Shun-ichi Amari, Alain Trounev, and Gabriel Peyré. Interpolating between optimal transport and mmd using sinkhorn divergences. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 2681–2690, 2019.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.
- Jonathan Ho, Chitwan Saharia, William Chan, David J Fleet, Mohammad Norouzi, and Tim Salimans. Cascaded diffusion models for high fidelity image generation. *Journal of Machine Learning Research*, 23(47):1–33, 2022a.
- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:1–13, 2022b.
- Emiel Hoogeboom, Victor Garcia Satorras, Clément Vignac, and Max Welling. Equivariant diffusion for molecule generation in 3d. In *International conference on machine learning*, pp. 8867–8887. PMLR, 2022.
- Chin-Wei Huang, Ricky TQ Chen, Christos Tsirigotis, and Aaron Courville. Convex potential flows: Universal probability distributions with optimal transport and convex optimization. *arXiv preprint arXiv:2012.05942*, 2020.
- Haohe Huang, Hao Liang, Karan Parmar, Andrew Zhu, Harsh Misra, Stephan Schneider, and Jesse Engel. Audioldm: Text-to-audio generation with latent diffusion models. *arXiv preprint arXiv:2301.12503*, 2023.
- Leonid Kantorovitch. On the translocation of masses. *Management science*, 5(1):1–4, 1958.
- Dongjun Kim, Chieh-Hsin Lai, Wei-Hsiang Liao, Naoki Murata, Yuhta Takida, Toshimitsu Uesaka, Yutong He, Yuki Mitsufuji, and Stefano Ermon. Consistency trajectory models: Learning probability flow ode trajectory of diffusion. *arXiv preprint arXiv:2310.02279*, 2023.
- Zhifeng Kong, Wei Ping, Jiaji Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2020.
- Sangyun Lee, Zinan Lin, and Giulia Fanti. Improving the training of rectified flows. *arXiv preprint arXiv:2405.20320*, 2024.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Qiang Liu. Rectified flow: A marginal preserving approach to optimal transport. *arXiv preprint arXiv:2209.14577*, 2022.
- Xingchao Liu, Chengyue Gong, and Qiang Liu. Flow straight and fast: Learning to generate and transfer data with rectified flow. *arXiv preprint arXiv:2209.03003*, 2022.

342 Zalan Liu, Yi Ren, Joel Shor, Shuo Huang, Bowen Li, Andros Belles, Vinay Suresh, and Rif A  
343 Saurous. Audioldm: a language modeling approach to audio generation. *IEEE/ACM Transactions*  
344 *on Audio, Speech, and Language Processing*, 2023.

345 Ashok Makkuva, Amirhossein Taghvaei, Sewoong Oh, and Jason Lee. Optimal transport mapping via  
346 input convex neural networks. In *International Conference on Machine Learning*, pp. 6672–6681.  
347 PMLR, 2020.

348 Robert J McCann. A convexity principle for interacting gases. *Advances in mathematics*, 128(1):  
349 153–179, 1997.

350 Gaspard Monge. Mémoire sur la théorie des déblais et des remblais. *Mem. Math. Phys. Acad. Royale*  
351 *Sci.*, pp. 666–704, 1781.

352 Stefano Peluchetti. Diffusion bridge mixture transports, schrödinger bridge problems and generative  
353 modeling. *Journal of Machine Learning Research*, 24(374):1–51, 2023.

354 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
355 resolution image synthesis with latent diffusion models. *Proceedings of the IEEE/CVF Conference*  
356 *on Computer Vision and Pattern Recognition*, pp. 10684–10695, 2022.

357 Saptarshi Roy, Vansh Bansal, Purnamrita Sarkar, and Alessandro Rinaldo. 2-rectifications are  
358 enough for straight flows: A theoretical insight into wasserstein convergence. *arXiv e-prints*, pp.  
359 arXiv–2410, 2024.

360 Yixiong Ruan, Janne Koh, Raphael Hofmann, Beat Gfeller, Marco Tagliasacchi, and Jesse Engel.  
361 Audiogen: Textually guided audio generation. *arXiv preprint arXiv:2209.15352*, 2023.

362 Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. *arXiv*  
363 *preprint arXiv:2202.00512*, 2022.

364 Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion schrödinger  
365 bridge matching. *Advances in Neural Information Processing Systems*, 36, 2024.

366 Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben  
367 Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint*  
368 *arXiv:2011.13456*, 2020.

369 Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. *arXiv preprint*  
370 *arXiv:2303.01469*, 2023.

371 Alexander Tong, Nikolay Malkin, Kilian Fatras, Lazar Atanackovic, Yanlei Zhang, Guillaume Huguët,  
372 Guy Wolf, and Yoshua Bengio. Simulation-free schrödinger bridges via score and flow matching.  
373 *arXiv preprint arXiv:2307.03672*, 2023a.

374 Alexander Tong, Nikolay Malkin, Guillaume Huguët, Yanlei Zhang, Jarrod Rector-Brooks, Kilian  
375 Fatras, Guy Wolf, and Yoshua Bengio. Conditional flow matching: Simulation-free dynamic  
376 optimal transport. *arXiv preprint arXiv:2302.00482*, 2(3), 2023b.

377 Yikai Wang, Xinzhou Wang, Zilong Chen, Zhengyi Wang, Fuchun Sun, and Jun Zhu. Vidu4d: Single  
378 generated video to high-fidelity 4d reconstruction with dynamic gaussian surfels. In *Advances in*  
379 *Neural Information Processing Systems*, 2024.

380 Minkai Xu, Lantao Yu, Yang Song, Chence Shi, Stefano Ermon, and Jian Tang. Geodiff: A geometric  
381 diffusion model for molecular conformation generation. *arXiv preprint arXiv:2203.02923*, 2022.

382 Lijun Zhou, Jianfeng Tan, Sheng Gu, Yue Zhang, Chunyuan Wang, Xiaohan Wu, Yifan Zhao, Wenhao  
383 Li, Yichao Wu, Minh Hoai, et al. Magvit: Masked generative video transformer. *arXiv preprint*  
384 *arXiv:2212.05199*, 2022.



## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: We reference our claims in the introduction of the paper.

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: We have separated a paragraph for limitations in the conclusion.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory assumptions and proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)



Justification: Results in our paper build on top of each other, and when they do not, we always present them in the body of the proposition or reference them.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: We don't propose any new method, and the experimental details can be found in the appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: We will provide our code in the supplementary materials.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental setting/details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We provide detail in the Appendix section.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment statistical significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We have computed confidence intervals for most our empirical results. For FID results, we mention them more for reference, as we don't try to benchmark against any method.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).

- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes]

Justification: Presented in the Appendix.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: Anonymity is preserved, we don't include human subjects, or sensitive data. All items use are publicly available.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: The paper is quite theoretical, and explains why two state-of-the-art models don't perform as expected. We don't propose solutions for it.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.

- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: We don't release any data models.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: The main package we use is torchcfm, and we cite the paper on multiple occasions, and we mention it in the Appendix also.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.

- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: We introduce no asset; we just study empirically and theoretically already established methods.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: no human subjects

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional review board (IRB) approvals or equivalent for research with human subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: No crowdsourcing was used, and no human participants.

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.

- 694                   • We recognize that the procedures for this may vary significantly between institutions  
695                   and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the  
696                   guidelines for their institution.  
697                   • For initial submissions, do not include any information that would break anonymity (if  
698                   applicable), such as the institution conducting the review.

699   **16. Declaration of LLM usage**

700   Question: Does the paper describe the usage of LLMs if it is an important, original, or  
701   non-standard component of the core methods in this research? Note that if the LLM is used  
702   only for writing, editing, or formatting purposes and does not impact the core methodology,  
703   scientific rigorousness, or originality of the research, declaration is not required.

704   Answer: [NA]

705   Justification: Used LLMs only for editing.

706   Guidelines:

- 707                   • The answer NA means that the core method development in this research does not  
708                   involve LLMs as any important, original, or non-standard components.  
709                   • Please refer to our LLM policy (<https://neurips.cc/Conferences/2025/LLM>)  
710                   for what should or should not be described.

## A Assumptions

**Assumption 1.** Suppose that  $v(x, t)$  satisfies the following conditions:

- **Lipschitz continuity in  $x$ :** There exists  $L > 0$  such that for all  $x, y \in \mathbb{R}^D$  and  $t \in [0, 1]$ ,

$$\|v(x, t) - v(y, t)\| \leq L\|x - y\|. \quad (12)$$

- **Linear growth:** There exists  $C > 0$  such that for all  $x \in \mathbb{R}^D$  and  $t \in [0, 1]$ ,

$$\|v(x, t)\| \leq C(1 + \|x\|). \quad (13)$$

With Assumption 1 in place we can conclude via Arnold (1992), that the ODE  $dx_t = v(x_t, t)dt$  admits unique solutions for all initial conditions  $x_0$ , and so that  $T(x_0) = x_0 + \int_0^1 v(x_t, t)dt$  is an injective map.

## B Proofs - Gaussian Setting:

**Extra-Lemma 1.** Let  $q(x) = \frac{1}{x}$ , with  $x \in \mathbb{R}^d$ . There is no finite parameterization of a Multilayer Perceptron (MLP) that can represent this function with zero error.

*Proof.* Assume, for the sake of contradiction, that there exists a finite MLP that represents  $q(x)$  exactly. To avoid issues at  $x = 0$  where  $q(x)$  becomes undefined, we restrict the domain of  $x$  to  $\mathbb{R} \setminus [a, b]$ , where  $0 \notin [a, b]$  and  $a, b \in \mathbb{R}$ . This ensures that  $q(x)$  is Lipschitz continuous on this domain.

Next, consider the scaling property of  $q(x)$ : for any scalar  $c \neq 0$ ,  $q(cx) = \frac{1}{c}q(x)$ . This property implies that  $q(x)$  scales inversely with  $c$ .

However, standard MLPs do not inherently possess this scaling property. For an MLP to exhibit this behavior, its activation functions would need to satisfy  $\sigma(cx) = \frac{1}{c}\sigma(x)$ , which is not true for commonly used activation functions such as ReLU or tanh.

Since the desired function  $q(x)$  exhibits a scaling property that standard MLPs cannot replicate, we conclude that no finite MLP can represent  $q(x)$  exactly.  $\square$

**Lemma 1.** Let  $X_0 \sim \mathcal{N}(0, I_d)$  and  $X_1 \sim \mathcal{N}(\mu, M_d)$ , where  $M_d$  is a positive definite and symmetric matrix. The OT vector field is given by

$$\hat{v}_{OT}(X_t, t, \hat{\Theta}, \hat{\theta}) = \hat{\theta} + \hat{\Theta}[I_d + t\hat{\Theta}]^{-1}(X_t - t\hat{\theta}),$$

and the rotating vector field:

$$\hat{v}_{rOT}(X_t, t, \hat{\Theta}_{rOT}, \hat{\theta}) = \hat{\theta} + \hat{\Theta}_{rOT}[I_d + t\hat{\Theta}_{rOT}]^{-1}(X_t - t\hat{\theta}),$$

where  $X_t = (1 - t)X_0 + tX_1$ ,  $\hat{\Theta} = M_d^{1/2} - I_d$ ,  $\hat{\Theta}_{rOT} = M_d^{1/2}R_{k^\circ} - I_d$ ,  $\hat{\theta} = \mu$ , and  $R_{k^\circ}$  is a rotation matrix with  $R_{k^\circ} \neq -R_{k^\circ}$ . Furthermore, the function  $\hat{v}_{OT}$  cannot be exactly represented by an MLP, CNN, Transformer architecture when given concatenated inputs  $[X_t, t]$ .

*Proof.* The optimal transport map from  $\mathcal{N}(0, I_d)$  to  $\mathcal{N}(\mu, M_d)$  is  $T_{OT}(x) = \mu + M_d^{1/2}x$ . Thus, if we denote the optimal coupling by  $X_0^*$  and  $X_1^*$ , we have  $X_1^* = T(X_0^*) = \mu + M_d^{1/2}X_0^*$ . The displacement interpolation is defined as  $X_t = (1 - t)X_0^* + tX_1^*$ . Substitute the expression for  $X_1^*$  into the interpolation:

$$X_t = (1 - t)X_0^* + t(\mu + M_d^{1/2}X_0^*) = \left[(1 - t)I_d + tM_d^{1/2}\right]X_0^* + t\mu. \quad (14)$$

Solving for  $X_0^*$  gives

$$X_0^* = \left[(1 - t)I_d + tM_d^{1/2}\right]^{-1}(X_t - t\mu).$$

Now, substitute  $X_0^*$  back into the transport map to get

$$X_1^* = \mu + M_d^{1/2} X_0^* = \mu + M_d^{1/2} \left[ (1-t)I_d + t M_d^{1/2} \right]^{-1} (X_t - t\mu). \quad (15)$$

The optimal displacement (vector field) is defined as  $v(X_t, t) = X_1^* - X_0^*$ . Therefore, we have

$$v_{OT}(X_t, t) = \mu + M_d^{1/2} \left[ (1-t)I_d + t M_d^{1/2} \right]^{-1} (X_t - t\mu) \quad (16)$$

$$- \left[ (1-t)I_d + t M_d^{1/2} \right]^{-1} (X_t - t\mu) \quad (17)$$

$$= \mu + \left( M_d^{1/2} - I_d \right) \left[ (1-t)I_d + t M_d^{1/2} \right]^{-1} (X_t - t\mu). \quad (18)$$

This is the desired expression for the vector field.

For  $v_{rOT}$ , the computation is analogous, with the key difference that the transport map is now given by  $T_{rOT}(x) = \mu + M_d^{1/2} R X_0$ . This defines a valid pushforward map, since for any rotation matrix  $R_{k^\circ}$ , if  $Z \sim \mathcal{N}(0, I_d)$ , then  $RZ \sim \mathcal{N}(0, I_d)$  as well, due to the rotational invariance of the standard Gaussian.

Proceeding as before, we solve for  $X_0^*$ , substitute it back into the transport map, and derive the corresponding vector field. This yields the expression for  $v_{rOT}$ . The derivation follows the same steps as in the case of  $v_{OT}$ , with the only difference being that  $M_d^{1/2}$  is replaced by  $M_d^{1/2} R$ .

Also we can't have  $R = -R$  ( $180^\circ$  rotation), because then the inverse  $\left[ (1-t)I_d - t M_d^{1/2} \right]^{-1} = \left[ I_d - t(M_d^{1/2} + I) \right]^{-1}$  is not computable when  $t(M_d^{1/2} + I) = I$ , (for example for  $M_d = I$ , and  $t = 1/2$ ). For a more complete approach see Liu (2022).

We proceed to prove that this parameterization cannot be represented with zero error by an MLP. Representing the function  $\frac{\theta}{1+t\theta}$  with zero error is equivalent to representing  $\frac{1}{\theta}$  with zero error. By applying Extra-Lemma 1, we conclude that such a representation is not possible.  $\square$

**Proposition 1.** Let  $\pi_0 \sim \mathcal{N}(0, I_d)$ , and  $\pi_1 \sim \mathcal{N}(\mu, M_d)$ , where  $M_d$  is a positive definite and symmetric matrix. Let the following push forward maps:  $T_{OT}(X_0) = M_d^{1/2} X_0 + \mu$ ,  $T_{rOT}(X_0) = M_d^{1/2} R X_0 + \mu$ , and  $T_{rand}(X_0) = X_1$ , where  $R$  is a rotation matrix. Let the linear transformation be  $v(X_t, t) = \theta + \Theta [I_d + t\Theta]^{-1} (X_t - t\theta)$ , where  $X_t = (1-t)X_0 + tX_1$ . Let:

$$L_{MC}(\Theta, \theta, T, v) = \frac{1}{N} \sum_{s=1}^N \|T(X_0^{(s)}) - X_0^{(s)} - v(X_t^{(s)}, t, \Theta, \theta)\|^2. \quad (19)$$

Then, for optimal,  $(\hat{\Theta}, \hat{\theta})$  we have  $\text{Var}[\nabla_{\theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, T_{OT})] = \text{Var}[\nabla_{\Theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, T_{OT})] = \text{Var}[\nabla_{\theta}(\hat{\Theta}, \hat{\theta}, v_{rOT}, T_{rOT})] = \text{Var}[\nabla_{\Theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{rOT}, T_{rOT})] = 0$ , and  $\text{Var}[\nabla_{\theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, v_{OT})] = \frac{4}{N} \|M_d 1\|^2$ , and  $\text{Var}[\nabla_{\Theta} L_{MC}(\hat{\Theta}, \hat{\theta}, v_{OT}, v_{OT})] = \frac{8}{N} \text{tr}((M_d^{1/2}(R - I))^2)$ . Moreover, for the maps  $T_{rOT}$  and  $T_{rand}$  and  $v_{OT}$  vector field, the variance of the gradients of the optimal vector field does not increase in regions where the interpolant trajectories come closer together.

*Proof.* We break this proof into the computation of variance over time and the value at time zero to simplify computations.

**1. Way 1: for  $t = 0$  but closed form is nice:** We compute gradient variance in two ways one relies on the fact that when  $T$  is deterministic there is a direct bijection between  $X_t$  and  $X_0$ , this of course



wont give us the values in terms of  $t$ , but the formulas are somehow nicer and easier to work with.  
 We are basically looking at the variance for  $(X_0, 0)$

$$\text{Var} \left[ \nabla_{\boldsymbol{\theta}} \frac{1}{N} \sum \left\| \mathbf{M}_d^{1/2} X_0^{(n)} - X_0^{(n)} + \mu - \boldsymbol{\Theta} X_0^{(n)} - \boldsymbol{\theta} \right\|^2 \right] \quad (20)$$

$$= \text{Var} \left[ -\frac{2}{N} \sum \mathbf{1}^\top \left( \mu + (\mathbf{M}_d^{1/2} - \boldsymbol{\Theta} - \mathbf{I}) X_0^{(n)} - \boldsymbol{\theta} \right) \right] \quad (21)$$

$$= \frac{4}{N} \text{Var} \left[ \mathbf{1}^\top (\mathbf{M}_d^{1/2} - \boldsymbol{\Theta} - \mathbf{I}) X_0^{(n)} \right] \quad (22)$$

$$= \frac{4}{N} \mathbf{1}^\top (\mathbf{M}_d^{1/2} - \boldsymbol{\Theta} - \mathbf{I}) (\mathbf{M}_d^{1/2} - \boldsymbol{\Theta} - \mathbf{I})^\top \mathbf{1}. \quad (23)$$

$$\text{Var} \left[ \nabla_{\boldsymbol{\Theta}_i} \frac{1}{N} \sum \left\| \mathbf{M}_d^{1/2} X_0^{(n)} - X_0^{(n)} - \mu - \boldsymbol{\Theta} X_0^{(n)} - \boldsymbol{\theta} \right\|^2 \right] \quad (24)$$

$$= \text{Var} \left[ -\frac{2}{N} \sum X_0^\top \left( \mu + (\mathbf{M}_d^{1/2} - \mathbf{I} - \boldsymbol{\Theta}) X_0^{(n)} - \boldsymbol{\theta} \right) \right]. \quad (25)$$

**Gradient variance for  $T_{OT}$  and  $v_{OT}$**  For the rotated case:

Note that our vector field is optimal when  $\hat{\boldsymbol{\Theta}} = \mathbf{M}_d^{1/2} - \mathbf{I}$  and  $\hat{\boldsymbol{\theta}} = \mu$ . Then we have:

$$\text{Var} \left[ \nabla_{\boldsymbol{\theta}} L_{MC}(v(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}}, X_0)) \right] = 0, \quad (26)$$

and

$$\text{Var} \left[ \nabla_{\boldsymbol{\Theta}} L_{MC}^{OT}(v(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}}, X_0)) \right] = \frac{4}{N} \text{Var} [X_0^\top (\mu - \hat{\boldsymbol{\theta}})] = \frac{4 \|\mu - \hat{\boldsymbol{\theta}}\|^2}{N} = 0. \quad (27)$$

**Gradient variance for  $T_{rOT}$  and  $v_{OT}$**  For the rotated case:

$$\text{Var} \left[ \nabla_{\boldsymbol{\theta}} \frac{1}{N} \sum \left\| \mathbf{M}_d^{1/2} \mathbf{R} X_0^{(n)} - X_0^{(n)} - \mu - \boldsymbol{\Theta} X_0^{(n)} - \boldsymbol{\theta} \right\|^2 \right] \quad (28)$$

$$= \text{Var} \left[ -\frac{2}{N} \sum \mathbf{1}^\top \left( \mu + (\mathbf{M}_d^{1/2} \mathbf{R} - \boldsymbol{\Theta} - \mathbf{I}) X_0^{(n)} - \boldsymbol{\theta} \right) \right] \quad (29)$$

$$= \frac{4}{N} \text{Var} \left[ \mathbf{1}^\top (\mathbf{M}_d^{1/2} \mathbf{R} - \boldsymbol{\Theta} - \mathbf{I}) X_0^{(n)} \right]. \quad (30)$$

At the optimum, where  $\hat{\boldsymbol{\Theta}} = \mathbf{M}_d^{1/2} - \mathbf{I}$ , we obtain:

$$\frac{4}{N} \text{Var} \left[ \mathbf{1}^\top (\mathbf{M}_d^{1/2} \mathbf{R} - \mathbf{M}_d^{1/2}) X_0 \right] = \frac{4}{N} \mathbf{1}^\top (\mathbf{M}_d^{1/2} (\mathbf{R} - \mathbf{I})) (\mathbf{M}_d^{1/2} (\mathbf{R} - \mathbf{I}))^\top \mathbf{1} \quad (31)$$

$$= \frac{4}{N} \mathbf{1}^\top (\mathbf{M}_d^{1/2} (\mathbf{R} - \mathbf{I})) (\mathbf{R} - \mathbf{I})^\top \mathbf{M}_d^{1/2 \top} \mathbf{1} \quad (32)$$

$$= \frac{4}{N} \mathbf{1}^\top \mathbf{M}_d^{1/2} (2\mathbf{I} - \mathbf{R}^\top - \mathbf{R}) \mathbf{M}_d^{1/2 \top} \mathbf{1}. \quad (33)$$

For the variance of the gradient w.r.t.  $\boldsymbol{\Theta}$  at the optimum, we get:

$$\text{Var} \left[ \nabla_{\boldsymbol{\Theta}} L_{MC}^{rOT}(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}}) \right] = \frac{4}{N} \text{Var} \left[ X_0^\top \mathbf{M}_d^{1/2} (\mathbf{R} - \mathbf{I}) X_0 \right] \quad (34)$$

$$= \frac{8}{N} \text{tr} \left( (\mathbf{M}_d^{1/2} (\mathbf{R} - \mathbf{I}))^2 \right). \quad (35)$$

The last equality follows from the known variance formula of a quadratic form for Gaussian random vectors, specifically when  $X_0 \sim \mathcal{N}(0, \mathbf{I})$ .

783 **Gradient variance for  $T_{rOT}$  and  $v_{rOT}$**  For the rotated case:

$$\text{Var} \left[ \nabla_{\boldsymbol{\theta}} \frac{1}{N} \sum \left\| \mathbf{M}_d^{1/2} \mathbf{R} X_0^{(n)} - X_0^{(n)} - \mu - \boldsymbol{\Theta} X_0^{(n)} - \boldsymbol{\theta} \right\|^2 \right] \quad (36)$$

$$= \text{Var} \left[ -\frac{2}{N} \sum \mathbf{1}^\top \left( \mu + (\mathbf{M}_d^{1/2} \mathbf{R} - \boldsymbol{\Theta} - \mathbf{I}) X_0^{(n)} - \boldsymbol{\theta} \right) \right] \quad (37)$$

$$(38)$$

784 At the optimum, where  $\hat{\boldsymbol{\Theta}} = \mathbf{M}_d^{1/2} \mathbf{R} - \mathbf{I}$ , we obtain:

$$\frac{4}{N} \text{Var} \left[ \mathbf{1}^\top (\mathbf{M}_d^{1/2} \mathbf{R} - \mathbf{M}_d^{1/2} \mathbf{R}) X_0 \right] = 0 \quad (39)$$

785 For the variance of the gradient w.r.t.  $\boldsymbol{\Theta}$  at the optimum, we get:

$$\text{Var} \left[ \nabla_{\boldsymbol{\Theta}} L_{MC}^{rOT}(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}}) \right] = \frac{4}{N} \text{Var} \left[ X_0^\top \mathbf{M}_d^{1/2} (\mathbf{R} - \mathbf{R}) X_0 \right] \quad (40)$$

$$= 0 \quad (41)$$

786 **Way 2 - (for any  $t$ , closed form is not nice): Gradient variance over time** Before analyzing  
787 different kinds of pairings that we want to analyze, first, we want a true formulation of the gradients:

$$\text{Var}[\nabla_{\boldsymbol{\theta}} L_{MC}^{type}(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}})] \quad (42)$$

$$= \text{Var} \left[ \nabla_{\boldsymbol{\theta}} \frac{1}{N} \sum [|(T_{type}(X_0^{(n)}) - X_0^{(n)} - v_{\boldsymbol{\Theta}, \boldsymbol{\theta}}(X_t^{(n)}, t))|^2] \right] \quad (43)$$

$$= \text{Var} \left[ -\frac{2}{N} \sum (T_{type}(X_0^{(n)}) - X_0^{(n)} - \boldsymbol{\theta} - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} (X_t^{(n)} - t\boldsymbol{\theta}))^\top (-\mathbf{1} - t\boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} \mathbf{1}) \right] \quad (44)$$

$$= \text{Var} \left[ -\frac{2}{N} \sum (T_{type}(X_0^{(n)}) - X_0^{(n)} - \boldsymbol{\theta} - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} (X_t^{(n)} - t\boldsymbol{\theta}))^\top ([I_d + t\boldsymbol{\Theta}]^{-1} \mathbf{1}) \right] \quad (45)$$

788 For now, let  $\mathbf{C}_{t, \boldsymbol{\Theta}} = [I_d + t\boldsymbol{\Theta}]^{-1}$ . We arrive at a more simplified:

$$\text{Var}[\nabla_{\boldsymbol{\theta}} L_{MC}^{type}(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}})] = \frac{4}{N} \text{Var}[(T_{type}(X_0) - X_0 - \boldsymbol{\theta} - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} (X_t - t\boldsymbol{\theta}))^\top \mathbf{C}_{t, \boldsymbol{\Theta}} \mathbf{1}] \quad (46)$$

$$= \frac{4}{N} \text{Var}[(T_{type}(X_0) - X_0 - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} X_t)^\top \mathbf{C}_{t, \boldsymbol{\Theta}} \mathbf{1}] \quad (47)$$

$$= (\mathbf{C}_{t, \boldsymbol{\Theta}} \mathbf{1})^\top \frac{4}{N} \text{Var}[(T_{type}(X_0) - X_0 - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} X_t)^\top \mathbf{C}_{t, \boldsymbol{\Theta}} \mathbf{1}] \quad (48)$$

789 Moving forward to gradients with respect to  $\boldsymbol{\Theta}$  we get:

$$\text{Var}[\nabla_{\boldsymbol{\Theta}} L_{MC}^{type}(\hat{\boldsymbol{\Theta}}, \hat{\boldsymbol{\theta}})] \quad (49)$$

$$= \text{Var} \left[ \nabla_{\boldsymbol{\Theta}} \frac{1}{N} \sum [|(T_{type}(X_0^{(n)}) - X_0^{(n)} - v_{\boldsymbol{\Theta}, \boldsymbol{\theta}}(X_t^{(n)}, t^{(n)}))|^2] \right] \quad (50)$$

$$= \frac{4}{N} \text{Var} \left[ (T_{type}(X_0) - X_0 - \boldsymbol{\theta} - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} (X_t - t\boldsymbol{\theta}))^\top ((t\mathbf{C}_{t, \boldsymbol{\Theta}} - \mathbf{I}) \mathbf{C}_{t, \boldsymbol{\Theta}} (X_t - t\boldsymbol{\theta})) \right] \quad (51)$$

$$(52)$$

790 **CASE 1 (OT):** This case is nice because we just need to prove that:

$$\text{Var}[(T_{OT}(X_0) - X_0 - \boldsymbol{\theta} - \boldsymbol{\Theta} [I_d + t\boldsymbol{\Theta}]^{-1} (X_t - t\boldsymbol{\theta}))^\top] = \mathbf{0}, \quad (53)$$

791 which actually happens once we replace  $\Theta = (M^{1/2} - I)$ , and  $\theta = \mu$ . We will show this for  
 792 completion:

$$M^{1/2}X_0 - X_0 - \Theta[I_d + t\Theta]^{-1}X_t \quad (54)$$

$$= (M^{1/2} - I)(I_d + t(M^{1/2} - I))^{-1}(X_0 + t(M^{1/2} - I)X_0 - tM^{1/2}X_0 - (1 - t)X_0) \quad (55)$$

$$= (M^{1/2} - I)(I_d + t(M^{1/2} - I))^{-1}\mathbf{0} = \mathbf{0} \quad (56)$$

793 **Remark** If our network is not quite at the optima so  $M^{1/2}X_0 - X_0 - \Theta[I_d + t\Theta]^{-1}X_t = \epsilon$  we can  
 794 notice that the variance of  $\theta$  and  $\Theta$  decrease in time because  $\frac{1}{1+tx}$  and  $\frac{1+tx+t}{(1+tx)^2}$  are decreasing in  $t$ .

795 **No connection between gradient variance and crossings** To show the lack of connection between  
 796 crossings and variance of gradients, it is enough to offer counter-examples. In the rotation case,  
 797 we will look at the variance when the rotation matrix does a  $180^\circ$  degree rotation, which means all  
 798 interpolation lines should meet at  $t = 1/2$ . As we will show, there is no peak in variance at time  
 799  $t = 1/2$ . In the Gaussian case, we will observe a very similar behaviour.

## 800 Case 2 (rOT):

801 1. For  $\theta$  the variance of the gradients will be:

$$\text{Var}[\nabla_{\theta} L_{\text{MC}}^{\text{type}}(\hat{\Theta}, \hat{\theta})] = \frac{4}{N} \mathbf{1}^T C_{t,\hat{\Theta}}^T \text{Var}[(M^{1/2}R - I - \hat{\Theta}C_{t,\hat{\Theta}}(t(M^{1/2}R - I) - I))X_0]C_{t,\hat{\Theta}} \mathbf{1} \quad (57)$$

$$= \frac{4}{N} \mathbf{1}^T C_{t,\hat{\Theta}}^T A^T A C_{t,\hat{\Theta}} \mathbf{1}, \quad (58)$$

802 where  $A = M^{1/2}R - I - (M^{1/2} - I)[I + t(M^{1/2} - I)]^{-1}(t(M^{1/2}R - I) + I)$ , with  $A^T A \approx$   
 803  $M + I - R - R^T + 2(1 - t)(R - I)^T(M^{1/2} - I)(R - I)$ , for small  $M^{1/2} - I$  (approximation  
 804 obtained via Taylor expansion).

805 **Counter-example** For  $R = -I$  and  $M^{1/2} = 2I$  (which is  $180^\circ$  rotation so all interpolants meet at  
 806 one point) we have  $A = -3I - [I - tI]^{-1}(-3It + I)$  so  $A^T A = (\frac{16}{1-t})^2 I$ . The variance in this case  
 807 will be increasing for  $t \in [0, 1]$ . This aligns with our empirical results, however, it does not align with  
 808 what is believed in literature, that the variance would peak at  $t = 1/2$ , because crossings increase  
 809 variance.

810 2. For  $\Theta$ , the variance of the gradients will be:

$$\text{Var}[\nabla_{\Theta} L_{\text{MC}}^{rOT}(\hat{\Theta}, \hat{\theta})] \quad (59)$$

$$= \frac{4}{N} \text{Var}[(T_{rOT}(X_0) - X_0 - \theta - \Theta[I_d + t\Theta]^{-1}(X_t - t\theta))^T ((tC_{t,\Theta} - I)C_{t,\Theta}(X_t - t\theta))] \quad (60)$$

$$= \frac{4}{N} \text{Var}[(M^{1/2}R - I - (M^{1/2} - I)[I_d + t(M^{1/2} - I)]^{-1}(tM^{1/2}R + (1 - t)I)X_0]^T \quad (61)$$

$$(tC_{t,\Theta} - I)C_{t,\Theta}(tM^{1/2} + (1 - t)I)X_0] \quad (62)$$

$$= \frac{4}{N} \text{Var}[(AX_0)^T (tC_{t,\Theta} - I)C_{t,\Theta}(tM^{1/2}R + (1 - t)I)X_0] = \frac{4}{N} \text{Var}[X_0^T A^T B X_0] = 2\text{Tr}((AB)^2) \quad (63)$$

811 **Counter-example** Let  $M^{1/2} = 2I$ , and  $R = -I$ . Then  $A = \frac{4}{1+t}I$ , and  $B = \frac{(3t-1)}{(1+t)^2}$  so  
 812  $\text{Var}[\nabla_{\Theta} L_{\text{MC}}^{rOT}(\hat{\Theta}, \hat{\theta})] = \frac{4}{N} \text{Var}[\frac{16}{(1+t)^2} \frac{(3t-1)}{(1+t)^2} X_0^T X_0] = \frac{8d}{N} \frac{(16(3t-1))^2}{(1+t)^8}$  which is decreasing for  
 813  $t \in [1/9, 1/3]$  and increasing  $t \in [1/3, 1]$ .

## 814 Case 3 (Random):

815 1. For  $\theta$  the variance of the gradients will be:

$$\text{Var}[\nabla_{\theta} L_{\text{MC}}^{\text{type}}(\hat{\Theta}, \hat{\theta})] = \frac{4}{N} \mathbf{1}^T C_{t, \hat{\Theta}}^T \text{Var}[(X_1 - X_0 - \hat{\Theta} C_{t, \hat{\Theta}}(tX_1 - (1-t)X_0))] C_{t, \hat{\Theta}} \mathbf{1} \quad (64)$$

$$= \frac{4}{N} \mathbf{1}^T C_{t, \hat{\Theta}}^T [\text{Var}[(I - \hat{\Theta} C_{t, \hat{\Theta}} t)X_1] + \text{Var}[(-I - \hat{\Theta} C_{t, \hat{\Theta}}(1-t)I)X_0]] C_{t, \hat{\Theta}} \mathbf{1}, \quad (65)$$

816 because of independence of  $X_0$ , and  $X_1$ . Continuing:

$$\text{Var}[(I - \hat{\Theta} C_{t, \hat{\Theta}} t)X_1] = (I - t(M^{1/2} - I)(I + t(M^{1/2} - I))^{-1})^T M \quad (66)$$

$$(I - t(M^{1/2} - I)(I + t(M^{1/2} - I))^{-1}) \quad (67)$$

$$= (I + t(M^{1/2} - I))^{-1} M (I + t(M^{1/2} - I))^{-1}, \quad (68)$$

817 and,

$$\text{Var}[(I - \hat{\Theta} C_{t, \hat{\Theta}}(1-t))X_0] = (I - (1-t)(M^{1/2} - I)(I + t(M^{1/2} - I))^{-1})^T \quad (69)$$

$$(I - (1-t)(M^{1/2} - I)(I + t(M^{1/2} - I))^{-1}). \quad (70)$$

818 **Counter-example:** For  $M^{1/2} = 2I$  we have  $\text{Var}[(I - \hat{\Theta} C_{t, \hat{\Theta}} t)X_1] = 4 \frac{1}{(1+t)^2} I$  which is decreasing

819 for  $t \in [0, 1]$ . We have  $\text{Var}[(I - \hat{\Theta} C_{t, \hat{\Theta}}(1-t))X_0] = 4 \frac{t^2}{(1+t)^2} I$ . That is  $4 \frac{t^2+1}{(1+t)^2}$  which is decreasing

820 for  $t \in [0, 1]$ . We have the overall variance:

$$\text{Var}[\nabla_{\theta} L_{\text{MC}}^{\text{type}}(\hat{\Theta}, \hat{\theta})] = 4 \mathbf{1}^T \frac{1}{1+t} \frac{t^2+1}{(1+t)^2} \frac{1}{1+t} I \mathbf{1} = 4d \frac{t^2+1}{(1+t)^4}, \quad (71)$$

821 which is decreasing in  $t \in [0, 1]$ .

## 822 C Proofs - Reflow

823 **Lemma 2.** Let  $\pi_1$  and  $\pi_2$  be two distributions on  $\mathbb{R}^N$  admitting densities, and let  $(Z_0, Z_1)$  be their  
824 straight-line coupling. If we apply Rectified Flows again using the noise-free interpolant, we recover  
825 the same coupling.

826 *Proof.* Since  $(Z_0, Z_1)$  is the straight coupling, for each  $t \in [0, 1]$  the interpolation

$$Z_t = (1-t)Z_0 + tZ_1$$

827 is  $\sigma(Z_0, Z_1)$ -measurable. Define

$$v(z, t) := \mathbb{E}[Z_1 - Z_0 \mid Z_t = z],$$

828 so that by the Doob–Dynkin lemma there is a (deterministic) function  $v$  satisfying

$$v(Z_t, t) = \mathbb{E}[Z_1 - Z_0 \mid Z_t] = Z_1 - Z_0 \quad \text{almost surely.}$$

829 Hence if we re-solve the same linear ODE

$$\frac{dX_t}{dt} = v(X_t, t) \quad \text{with} \quad X_0 = Z_0,$$

830 the unique solution is exactly

$$X_t = Z_0 + t(Z_1 - Z_0) = Z_t.$$

831 Because the noise-free Rectified Flow minimizes the mean-squared drift error,

$$\mathbb{E} \|v(Z_t, t) - (Z_1 - Z_0)\|^2 = 0,$$

832 no further change occurs. Thus the pairing remains  $(Z_0, Z_1)$ .  $\square$

833 Moreover

**Proposition 3** (Interpolating lines don't meet in high dimensions). *Let  $x_0, x_1 \sim \pi_0(x)$  and  $y_0, y_1 \sim \pi_1(y)$ , where  $\pi_0$  and  $\pi_1$  are probability distributions on  $\mathbb{R}^d$  admitting a density. Define the linear interpolants  $l_i(t_i) = (1 - t_i)x_i + t_i y_i$  for  $i \in \{0, 1\}$ .*

A. *For  $d \geq 2$ , the lines  $l_0$  and  $l_1$  cross at  $t = t_i = t_j \in (0, 1)$  with probability 0.*

B. *For  $d > 2$  the two lines intersect for  $t_i, t_j \in (0, 1)$  with probability 0.*

*Proof.* Throughout this proof, we use the following theorem: Let  $X$  be a random variable with a continuous probability distribution in  $\mathbb{R}^n$ , and let  $A$  be a lower-dimensional subset of  $\mathbb{R}^n$ . Since  $\mathbb{P}$  admits a density, it follows that it is absolutely continuous wrt to the Lebesgue measure  $\lambda$  then as  $\lambda(A) = 0$  by absolute continuity, we have that  $\mathbb{P}(X \in A) = 0$ .

**A.** Suppose the lines  $l_0(t)$  and  $l_1(t)$  intersect at some  $t = \hat{t} \in (0, 1)$ . This implies

$$l_0(\hat{t}) = l_1(\hat{t}),$$

which expands to

$$(1 - \hat{t})x_0 + \hat{t}y_0 = (1 - \hat{t})x_1 + \hat{t}y_1.$$

Rearranging terms, we find

$$y_1 = \frac{1 - \hat{t}}{\hat{t}}(x_0 - x_1) + y_0.$$

To compute the probability of such an intersection, note that  $y_1$  must lie exactly on the affine subspace defined by the above equation, which is a one-dimensional line segment in  $\mathbb{R}^d$ .

The joint probability of the points  $(x_0, x_1, y_0, y_1)$  can be written as

$$\mathbb{P}(l_0(t) \text{ intersects } l_1(t) \text{ for } t \in (0, 1)) = \mathbb{P}(x_0, x_1, y_0) \cdot \mathbb{P}(y_1 = \frac{1 - \hat{t}}{\hat{t}}(x_0 - x_1) + y_0 \mid x_0, x_1, y_0).$$

Since  $\pi_1(y)$  is continuous and differentiable, the probability density of  $y_1$  lying on any lower-dimensional subspace (e.g., a line segment) in  $\mathbb{R}^d$ , with  $d > 2$ , is zero. Therefore,

$$\mathbb{P}(y_1 = \frac{1 - \hat{t}}{\hat{t}}(x_0 - x_1) + y_0 \mid x_0, x_1, y_0) = 0,$$

which implies

$$\mathbb{P}(l_0(t) \text{ intersects } l_1(t) \text{ at } t = \hat{t} \text{ for } t \in (0, 1)) = 0.$$

Hence, the lines  $l_0(t)$  and  $l_1(t)$  intersect with probability zero for  $t \in (0, 1)$ . **B.** Suppose the lines  $l_0(t_0)$  and  $l_1(t_1)$  intersect at some  $t_0 = \hat{t}_0, t_1 = \hat{t}_1$  for  $t_0, t_1 \in (0, 1)$ . This implies

$$l_0(\hat{t}_0) = l_1(\hat{t}_1),$$

which expands to

$$(1 - \hat{t}_0)x_0 + \hat{t}_0 y_0 = (1 - \hat{t}_1)x_1 + \hat{t}_1 y_1.$$

Rearranging terms, we find

$$y_1 = \frac{(1 - \hat{t}_0)x_0 - (1 - \hat{t}_1)x_1 + \hat{t}_0 y_0}{\hat{t}_1}.$$

which for  $\hat{t}_0, \hat{t}_1 \in (0, 1)$  we have that  $y_1$  would belong to a 2D surface. Just like before we have:

$$\mathbb{P}(l_0(t_0) \text{ intersects } l_1(t_1) \text{ for } \hat{t}_i \in (0, 1)) = \mathbb{P}(x_0, x_1, y_0) \cdot \mathbb{P}\left(y_1 = \frac{(1 - \hat{t}_0)x_0 - (1 - \hat{t}_1)x_1 + \hat{t}_0 y_0}{\hat{t}_1} \mid x_0, x_1, y_0\right) = 0.$$

□

**Proposition 2.** Let  $\pi_0$  and  $\pi_1$  be two probability densities on  $\mathbb{R}^D$ , and let  $T : \mathbb{R}^D \rightarrow \mathbb{R}^D$  be a deterministic transport map pushing  $\pi_0$  to  $\pi_1$  (so if  $Z_0 \sim \pi_0$ , then  $Z_1 = T(Z_0) \sim \pi_1$ ). Suppose we draw a finite dataset of  $N$  i.i.d. samples  $\{Z_0^{(i)}, Z_1^{(i)} = T(Z_0^{(i)})\}_{i=1}^N$ , and for each  $i$  we also sample  $m$  time-points  $\{t^{(i,j)}\}_{j=1}^m \subset [0, 1]$  (e.g. uniformly). Let  $Z_t^{(i,j)} = (1 - t^{(i,j)}) Z_0^{(i)} + t^{(i,j)} Z_1^{(i)}$ . Define the empirical loss over this “doubly indexed” dataset by

$$L_{\text{MC}}^{\text{det}}(v_\Theta) = \frac{1}{Nm} \sum_{i=1}^N \sum_{j=1}^m \left\| (Z_1^{(i)} - Z_0^{(i)}) - v_\Theta(Z_t^{(i,j)}, t^{(i,j)}) \right\|^2.$$

Then there exists a (deterministic) vector field  $v$  attaining zero loss:  $L_{\text{MC}}^{\text{det}}(v) = 0$ .

*Proof.* Our finite dataset is

$$\mathcal{D} = \left\{ (Z_0^{(i)}, Z_1^{(i)}, t^{(i,j)}, Z_t^{(i,j)}) : i = 1, \dots, N; j = 1, \dots, m \right\}.$$

We may view this as drawing from the discrete joint “empirical” measure

$$\hat{P} = \frac{1}{Nm} \sum_{i=1}^N \sum_{j=1}^m \delta_{(Z_0^{(i)}, Z_1^{(i)}, t^{(i,j)}, Z_t^{(i,j)})}.$$

For each datapoint  $(Z_t^{(i,j)}, t^{(i,j)})$  we know

$$Z_1^{(i)} - Z_0^{(i)} = \mathbb{E}[Z_1 - Z_0 \mid Z_0 = Z_0^{(i)}, Z_1 = Z_1^{(i)}, t = t^{(i,j)}, Z_t = Z_t^{(i,j)}].$$

Let  $\mathcal{Z} = \bigcup_{i \in N, t \in [0,1]} (Z_t^{(i)}, t)$ . Let  $\mathcal{A} = \bigcap_{i \in N, t \in [0,1]} (Z_t^{(i)}, t)$ .

The probability of sampling  $t^{(i,j)} = a \in \mathcal{A}$  is 0, because  $t$  lives in a two-dimensional distribution and the points of intersection are a countable union (maximum  $N!$  intersection) of one-dimensional spaces.

$$\mathbb{E}[Z_1 - Z_0 \mid Z_0 = Z_0^{(i)}, Z_1 = Z_1^{(i)}, t = t^{(i,j)}, Z_t = Z_t^{(i,j)}] = \mathbb{E}[Z_1 - Z_0 \mid t = t^{(i,j)}, Z_t = Z_t^{(i,j)}].$$

Hence we can define the exact conditional-expectation vector field

$$v(z, t) := \mathbb{E}_{\hat{P}}[Z_1 - Z_0 \mid Z_t = z, t],$$

which on each of our training points  $(z, t) = (Z_t^{(i,j)}, t^{(i,j)})$  satisfies

$$v(Z_t^{(i,j)}, t^{(i,j)}) = Z_1^{(i)} - Z_0^{(i)}.$$

Substituting this  $v$  into the Monte Carlo loss gives, term by term,

$$\left\| Z_1^{(i)} - Z_0^{(i)} - v(Z_t^{(i,j)}, t^{(i,j)}) \right\|^2 = \left\| Z_1^{(i)} - Z_0^{(i)} - (Z_1^{(i)} - Z_0^{(i)}) \right\|^2 = 0.$$

Averaging over all  $i, j$  yields  $L_{\text{MC}}^{\text{det}}(v) = 0$ .

Thus the loss can be driven exactly to zero on the finite sample by choosing  $v$  to interpolate the known displacements  $Z_1^{(i)} - Z_0^{(i)}$  at the sampled intermediate points  $(Z_t^{(i,j)}, t^{(i,j)})$ .  $\square$

## D Counter Example for straightness after 1 iteration

**Proposition 4** (Limitations of ReFlow Iterations on Noiseless Interpolants). Let  $\pi_0, \pi_1 \subset \mathbb{R}^D$ , with  $D > 2$ . Let  $(Z_0, Z_1) = 1\text{-ReFlow}(X_0, X_1)$  denote the coupling obtained after one ReFlow iteration. Under the assumption that the interpolant  $p(z, t) = (1 - t)z + tT(z)$  is injective in  $z$  for each  $t$ , there exists a vector field  $\hat{v}_1(x_t, t)$  such that performing a second ReFlow iteration using  $(Z_0, Z_1)$  yields the same coupling:

$$2\text{-ReFlow}(X_0, X_1) = (Z_1, Z_0).$$

Moreover, this second flow  $\hat{v}_1(x_t, t)$  generates straight-line paths and achieves zero loss. Therefore, further ReFlow iterations do not alter the couplings.

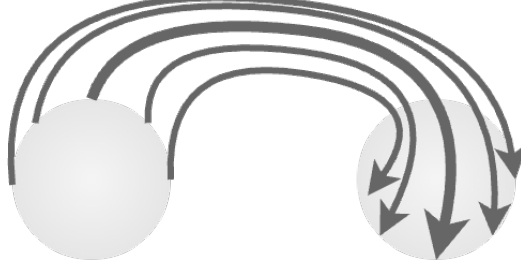


Figure 9: 180° rotation being realized by a continuous vector field.

885 *Annotated Proof.* Let  $v_0(x_t, t)$  be the learned vector field after CFM, with Assumption 1. From  
 886 Arnold (1992), the transport map  $T(z_0) = z_0 + \int_0^1 v(z_t, t)dt$  with  $z_0 \in \pi_0$  is injective.

887 **Step 1: Injectivity of  $T$**

- 888 • *Assumption:*  $T$  is injective (guaranteed by properties of  $v$ , e.g., Lipschitz, linear growth).
- 889 • *Potential Failure:* If  $T$  is not injective, the argument fails immediately.

890 **Step 2: Injectivity of the Interpolant  $p(z, t)$**

891 We claim that  $p(z, t) = (1 - t)z + tT(z)$  is also injective in  $z$  for each  $t \in [0, 1]$ .

- 892 • *Proof (by contradiction):* Suppose  $p(z_0^{(1)}, \hat{t}) = p(z_0^{(2)}, \hat{t})$  for some  $\hat{t}$  and  $z_0^{(1)} \neq z_0^{(2)}$ .
- 893 • *Rearranging:*

$$T(z_0^{(2)}) - T(z_0^{(1)}) = \frac{1 - \hat{t}}{\hat{t}}(z_0^{(1)} - z_0^{(2)})$$

- 894 • This statement doesn't really contradict with  $T$  injectivity.
- 895 • *Failure Point:* As shown in Counterexample 1, this is **not always true**. For example, if  $T$  is  
 896 a rotation plus translation,  $p(z, t)$  can fail to be injective even if  $T$  is injective.

897 One might argue that such a transport map is unlikely to be learned in practice; however, this is  
 898 not the point. Our argument is purely theoretical: injectivity of  $T$  does not imply injectivity of  
 899  $I(x_0, T(x_0), t)$ , and thus 1-ReFlow is insufficient even under standard regularity assumptions.

900 **Step 3: Construction of Inverse and New Vector Field**

901 Assuming injectivity, we can define an inverse  $f^{-1}(z_t, t) = z_0$ , and then set  $v_1(z_t, t) =$   
 902  $T(f^{-1}(z_t, t)) - f^{-1}(z_t, t)$ .

- 903 • *Dependency:* This construction **only works if  $f^{-1}$  exists**, i.e., if  $p(z, t)$  is injective.
- 904 • *Failure Point:* If  $p(z, t)$  is not injective,  $f^{-1}$  is not well-defined, and this construction fails.

905 **Step 4: Straight-Line Paths and Zero Loss**

906 Because  $z_1 - z_0$  is uniquely determined by  $z_t$ , we have  $\int_0^1 \mathbb{E}[\text{Var}(Z_1 - Z_0 | Z_t)] = 0$ , implying  
 907 straight paths and zero loss.

- 908 • *Dependency:* Again, relies on the injectivity of  $p(z, t)$ .

909

□

910 **Counter Example 1.** Let  $\pi_0 \sim \mathcal{N}(0, I_d)$  and  $\pi_1 \sim \mathcal{N}(5, I_d)$ , and let  $T(x_0) = R_{180^\circ}x_0 + 5$ , where  
 911  $R_{180^\circ}$  is a  $180^\circ$  rotation.  $T$  is injective, but the interpolant  $I(z_0, T(z_0), t) = (1-t)z_0 + tT(z_0)$  is  
 912 **not** injective (distinct  $z_0$  can map to the same  $x_t$  for some  $t$ ). Thus,  $f(x_t, t) = x_0$  is not well-defined,  
 913 and the construction of  $v_{\text{new}}(x_t, t)$  fails.

Step	Assumption Needed	Where It Can Fail
1	$T$ injective	Pathological $T$
2	$p(z, t)$ injective for $t$	Nonlinear $T$ (e.g., rotations)
3	$f^{-1}$ exists	$p(z, t)$ not injective
4	Unique $z_0$ for each $z_t$	$p(z, t)$ not injective

Table 3: Summary of dependencies and failure points in the proof.

914 **Final Note.** This proposition is only valid under the additional assumption that the interpolant  
 915  $p(z, t)$  is injective in  $z$  for each  $t$ . The counterexample demonstrates that this is not always the case,  
 916 even when  $T$  is injective. Therefore, care must be taken before applying this argument in general  
 917 settings.  $\square$

## 918 E Synthetic Experiments Details

919 **Experiment Setting** We evaluate generative models on synthetic datasets in dimensions 3 and 50.  
 920 Each dataset is constructed by sampling from a Gaussian Mixture Model (GMM) with randomly  
 921 initialized means and covariances, following our implementation in `generate_datasets.py`. The  
 922 source distribution is standard normal, and the target is the GMM. We compare Conditional Flow  
 923 Matching (CFM) and CFM with Stochastic Interpolation (CFM+SI), both implemented as neural  
 924 ODEs with time-conditioned MLP (Three-layer MLP, width 64, SELU activations) vector fields.  
 925 Models are trained and evaluated on both in-sample (training) and out-of-sample (test) data. All  
 926 metrics are computed as described below. We have used resources from Feydy et al. (2019) to  
 927 compute the distances.

### 928 Metric Descriptions

- 929 • **Log Probability (LogProb):** Measures the average log-likelihood of generated samples  
 930 under the target GMM distribution. Lower values indicate a better fit the modes.
- 931 • **Maximum Mean Discrepancy (MMD):** A kernel-based statistical distance between two  
 932 distributions, here computed using a Gaussian kernel. Lower values indicate better sample  
 933 quality.
- 934 • **Sinkhorn Distance:** An entropy-regularized approximation of the Wasserstein (optimal  
 935 transport) distance between empirical distributions, computed using the Sinkhorn algorithm.  
 936 Lower values indicate closer distributions.

937 **Note on Log-Likelihood Values** To improve readability and avoid confusion, we report and plot  
 938 the **positive** values of the log-likelihood (LogProb) throughout this paper, rather than the conventional  
 939 negative log-likelihood (NLL). This allows for a more intuitive comparison, where lower values  
 940 indicate better model performance.

941 It is important to note that log-likelihood (LogProb) primarily rewards models that generate samples  
 942 close to the high-density regions (modes) of the target distribution, rather than accurately capturing  
 943 the overall shape or support of the distribution. As a result, models that concentrate samples around  
 944 the modes can achieve high log-likelihood scores even if they do not match the full distribution  
 945 well. This should be kept in mind when interpreting LogProb values alongside other distance-based  
 946 metrics.

947 **How to read the table?** This table summarizes several metrics for model evaluation. For readers  
 948 unfamiliar with these results, here is how to interpret them:



Table 4: Comparison of CFM and CFM+SI across dimensions 3 and 50.

Dimension		3				50			
		Gen	Mem	True	Data	Gen	Mem	True	Data
CFM	LogProb	4.0150 $\pm 0.0032$	4.0156 $\pm 0.0032$	4.1330 $\pm 0.031$	4.0155 $\pm 0.035$	54.8299 $\pm 0.015$	53.6502 $\pm 0.26$	52.5094 $\pm 0.0833$	53.6244 $\pm 0.014$
	MMD	0.0034 $\pm 0.0005$	$1.758 \times 10^{-6}$ $\pm 0.0001$	0.0014 $\pm 0.0001$	0.0032 $\pm 0.0002$	0.0021 $\pm 0.0001$	$9.089 \times 10^{-6}$ $\pm 0.0001$	0.0020 $\pm 2e-09$	0.0019 $\pm 0.0001$
	Sinkhorn	0.0730 $\pm 0.0054$	$1.411 \times 10^{-5}$ $\pm 0.002$	0.0637 $\pm 0.002$	0.0790 $\pm 0.004$	15.1900 $\pm 0.162$	0.0045 $\pm 0.110$	14.3221 $\pm 0.110$	15.7400 $\pm 0.130$
CFM+SI	LogProb	4.1270 $\pm 0.0010$	4.0960 $\pm 0.016$	4.1330 $\pm 0.0008$	4.0155 $\pm 0.0013$	54.7220 $\pm 0.14$	53.8890 $\pm 0.16$	52.5094 $\pm 0.11$	53.6244 $\pm 0.26$
	MMD	0.0018 $\pm 0.00048$	$3.105 \times 10^{-5}$ $\pm 0.0003$	0.0014 $\pm 0.0003$	0.0032 $\pm 0.0002$	0.0020 $\pm 0.0001$	$6.09 \times 10^{-5}$ $\pm 0.0001$	0.0020 $\pm 3e-09$	0.0019 $\pm 0.0001$
	Sinkhorn	0.0680 $\pm 0.0015$	$3.557 \times 10^{-4}$ $\pm 0.007$	0.0637 $\pm 0.007$	0.0790 $\pm 0.004$	15.1689 $\pm 0.170$	0.0304 $\pm 0.220$	14.3221 $\pm 0.220$	15.7400 $\pm 0.230$

949 Consider the case for dimension 3. For the **Gen** (generated) column, we want the log-likelihood  
950 value to be closer to the **True** value rather than the **Data** value. If the generated value is closer to the  
951 data, it indicates memorization, rather than true generalization. For example, in CFM, the generated  
952 value is much closer to the data than the true value, suggesting memorization. This discrepancy  
953 arises because negative log-likelihood (NLL) tends to favor models that sample near the training data  
954 modes, rather than those that capture the full distribution.

955 Similarly, in the **Mem** (memorization) column, a value close to the data again indicates overfitting to  
956 the training points.

957 For the **MMD** and **Sinkhorn** metrics, these measure distances between newly generated trajectories  
958 (starting from training points) and the pairings from previous iterations. In both MMD and Sinkhorn,  
959 we observe that CFM+SI memorizes roughly ten times less than standard CFM, indicating better  
960 generalization.

## 961 E.1 Asymmetry in Gaussian mixtures

962 **Experiment Settings** We systematically investigated the Conditional Flow Matching (CFM) and  
 963 CFM+SI approaches for learning mappings between mixtures of Gaussian distributions (GMM-to-  
 964 GMM). In each experiment, both the source and target distributions were Gaussian mixtures, with the  
 965 number of components for each varied across the grid:  $\{1, 2, 4, 8, 16, 32, 64\}$ . The neural network  
 966 architecture for the vector field was a multilayer perceptron (MLP) with configurable width ( $w$ , e.g.,  
 967 64 by default) and input dimension ( $d$ , e.g., 10). Each cell in the results grid corresponds to a specific  
 968 pair of source and target GMM component counts, allowing us to analyze the effect of distribution  
 969 complexity on learning dynamics and gradient variance. Training was performed for up to 50,000  
 970 epochs using a batch size of 128 and a learning rate of  $10^{-3}$ , with integration performed via Neural  
 ODEs.

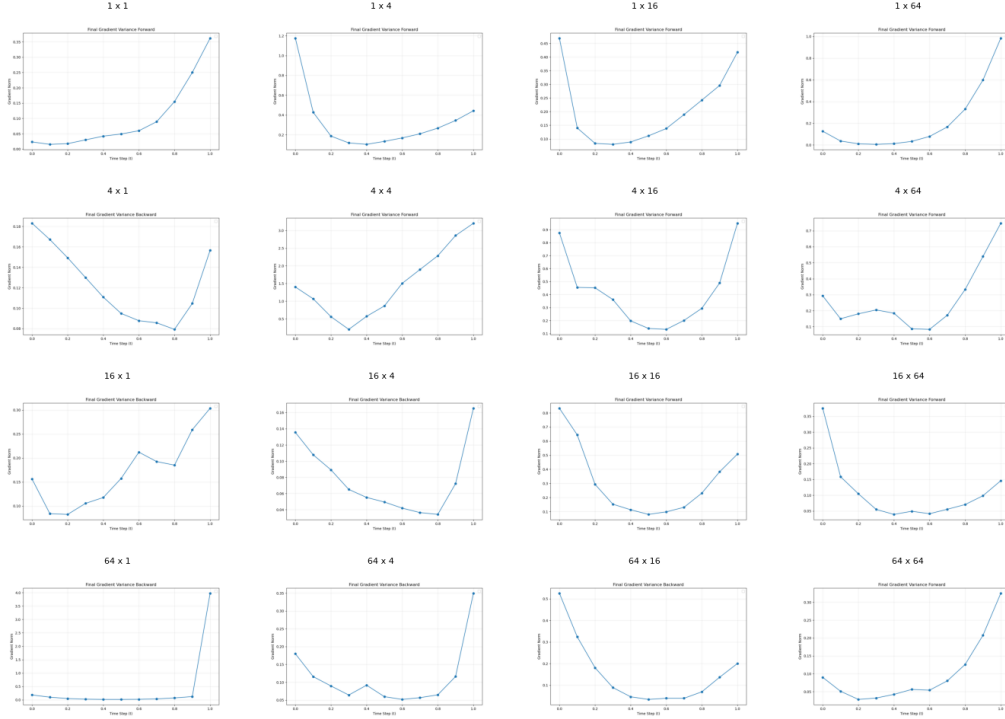


Figure 10: CFM in 10 dimensions. Each title indicates the number of Gaussian components in the source distribution multiplied by the number in the target distribution (e.g.,  $4 \times 16$  means 4 modes at source, 16 at target).

971

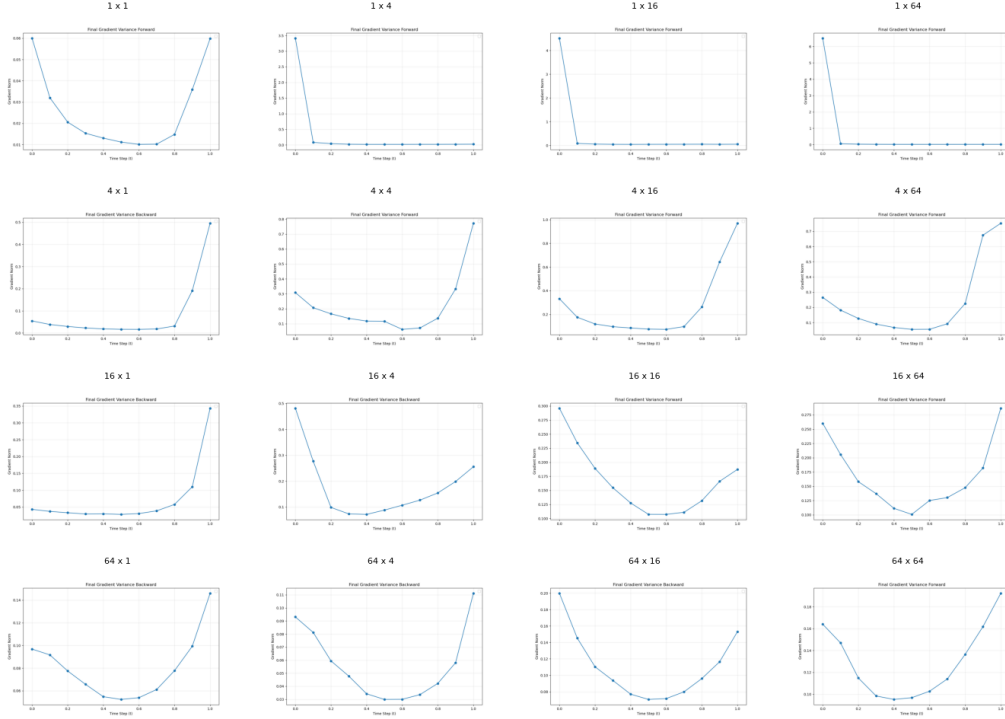


Figure 11: CFM combined with stochastic interpolants for mixture-to-mixture transport in 10 dimensions. Each panel is labeled as source modes  $\times$  target modes.

We can notice that the variance of the gradients looks more stable in the case of the Figure 11 (stochastic interpolants), than in the case of Figure 10 (noisless).

## F CIFAR-10 Experiment Details

**Model architecture.** All experiments used a U-Net-based neural network (UNetModelWrapper) with the following configuration: input shape  $(3, 32, 32)$ , base channels 128, 2 residual blocks per level, channel multipliers  $[1, 2, 2, 2]$ , attention at  $16 \times 16$  resolution (4 heads, 64 head channels), and dropout rate 0.1. The model is wrapped in a Neural ODE solver (Euler method).

**Training.** Models were trained on the CIFAR-10 training set, using random horizontal flips and normalization to  $[-1, 1]$ . Optimization used Adam with learning rate  $2 \times 10^{-4}$ , batch size 128, gradient clipping at 1.0, and a linear warmup over the first 5,000 steps. Each run used 400,001 steps (unless otherwise noted), with exponential moving average (EMA) of model weights (0.9999 decay). Checkpoints were saved every 20,000 steps. All experiments used 4 data loader workers and CUDA if available.

**Flow objectives.** We used Conditional Flow Matching (CFM, `-model cfm`), Schrödinger Bridge Matching (SBM, `-model sbm`), and other variants.

**Bidirectional setup.** Both forward (Gaussian  $\rightarrow$  CIFAR-10) and backward (CIFAR-10  $\rightarrow$  Gaussian) models were trained independently with identical hyperparameters. For the forward model, the source is standard Gaussian noise and the target is real images; for the backward model, the roles are swapped.

**Hardware and runtime.** All CIFAR-10 experiments were conducted on a compute cluster equipped with NVIDIA A10 GPUs (24 GB VRAM, CUDA 12.2). Each training run was allocated a single A10 GPU and typically ran for 24 hours to reach 240,000 optimization steps. These resources enabled efficient training of both forward and backward models at the scale reported in the main text.

## 995 **F.1 FID and Interpolant Details**

996 For Figure 7, we trained four vector fields (two forward and two backward), each using a different  
 997 interpolant:

- 998 • **CFM:** The interpolant is deterministic,

$$x_t = (1 - t)x_0 + tx_1.$$

- 999 • **CFM+SI:** The interpolant includes stochastic noise,

$$x_t = (1 - t)x_0 + tx_1 + \sigma\sqrt{t(1 - t)}Z,$$

1000 where  $Z \sim \mathcal{N}(0, I)$  and  $\sigma = 0.01$ .

1001 The models were evaluated as follows:

- 1002 • CFM forward FID: **4.199**
- 1003 • CFM backward scaled NLL: **1.426**
- 1004 • CFM+SI forward FID: **4.250**
- 1005 • CFM+SI backward scaled NLL: **1.423**

## 1006 **F.2 Does adding powers of $t$ help models learn better?**

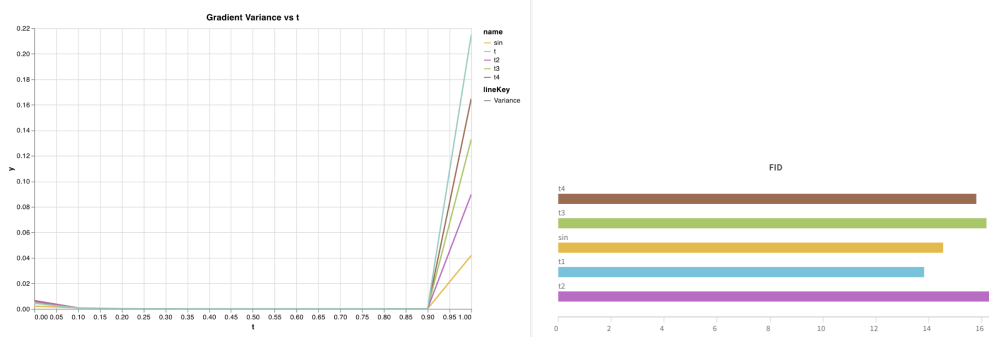


Figure 12: Variance not necessarily correlated with performance. I was wondering if it is worth saying that sinusoidal time embedding is better at representing function s like  $q(x) = \frac{1}{1+x}$ , than polynomials. (Taylor vs Fourier)