

# Elucidating the Design Choice of Probability Paths in Flow Matching for Forecasting

**Soon Hoe Lim**

*Department of Mathematics, KTH Royal Institute of Technology  
Nordita, KTH Royal Institute of Technology and Stockholm University*

*shlim@kth.se*

**Yijin Wang**

*International Computer Science Institute*

*yijin.wang@berkeley.edu*

**Annan Yu**

*Center for Applied Mathematics, Cornell University*

*ay262@cornell.edu*

**Emma Hart**

*Department of Mathematics, Emory University*

*emmahart@lbl.gov*

**Michael W. Mahoney**

*Department of Statistics, University of California at Berkeley  
International Computer Science Institute  
Lawrence Berkeley National Laboratory*

*mmahoney@stat.berkeley.edu*

**Xiaoye S. Li**

*Lawrence Berkeley National Laboratory*

*xsli@lbl.gov*

**N. Benjamin Erichson**

*International Computer Science Institute  
Lawrence Berkeley National Laboratory*

*erichson@icsi.berkeley.edu*

**Reviewed on OpenReview:** <https://openreview.net/forum?id=JApMDLwbLR>

## Abstract

Flow matching has recently emerged as a powerful paradigm for generative modeling and has been extended to probabilistic time series forecasting. However, the impact of the specific choice of probability path model on forecasting performance, particularly for high-dimensional spatio-temporal dynamics, remains under-explored. In this work, we demonstrate that forecasting spatio-temporal data with flow matching is highly sensitive to the selection of the probability path model. Motivated by this insight, we propose a novel probability path model designed to improve forecasting performance. Our empirical results across various dynamical system benchmarks show that our model achieves faster convergence during training and improved predictive performance compared to existing probability path models. Importantly, our approach is efficient during inference, requiring only a few sampling steps. This makes our proposed model practical for real-world applications and opens new avenues for probabilistic forecasting.

## 1 Introduction

Generative modeling has achieved remarkable success in recent years, especially for generating high-dimensional objects by learning mappings from simple, easily-sampled reference distributions,  $\pi_0$ , to complex target distributions,  $\pi_1$ . In particular, diffusion models have pushed the capabilities of generating realistic samples across various data modalities, including images (Ho et al., 2020; Song et al., 2020b; Karras et al., 2022),

videos (Ho et al., 2022; Blattmann et al., 2023; Gupta et al., 2023), and spatio-temporal scientific data like climate and weather patterns (Pathak et al., 2024; Kohl et al., 2024). Despite their impressive performance, diffusion models often come with high computational costs during training and inference. Additionally, they typically assume a Gaussian reference distribution, which may not be optimal for all data types and can limit modeling flexibility.

One promising alternative is flow matching, where the mappings are learned via a stochastic process that transforms  $\pi_0$  into  $\pi_1$  through an ordinary differential equation (ODE), approximating its marginal vector flow (Lipman et al., 2022; Albergo et al., 2023; Liu et al., 2022; Tong et al., 2023b; Pooladian et al., 2023; Lipman et al., 2024). While score-based models (Song & Ermon, 2019; Song et al., 2020a;b; Ho et al., 2020) are specific instances of flow matching, with Gaussian transition densities, the general framework allows for a broader class of interpolating paths. This flexibility can lead to deterministic sampling schemes that are faster and require fewer steps (Zhang & Chen, 2022). Recent work has demonstrated the remarkable capabilities of flow matching models for generating high-dimensional images (Esser et al., 2024) and discrete data (Gat et al., 2024).

Building on this, flow matching in latent space has recently been applied to forecasting spatio-temporal data (Davtyan et al., 2023) (predicting future frames in videos). This approach leverages latent representations to capture the complex dynamics inherent in temporal data. However, spatio-temporal forecasting, especially for videos and dynamical systems data, presents unique challenges. A video prediction model capable of generalizing to new, unseen scenarios must implicitly “understand” the scene: detecting and classifying objects, learning how they move and interact, estimating their 3D shapes and positions, and modeling the physical laws governing the environment (Battaglia et al., 2016). Similarly, accurate weather forecasting requires capturing intricate physical processes and interactions across multiple scales (Dueben & Bauer, 2018; Schultz et al., 2021).

We observe that, in the context of spatio-temporal forecasting, the performance of flow matching is highly sensitive to the choice of the probability path model, an important topic which has not been widely explored within a unified framework. Different probability paths can significantly impact the accuracy and convergence of forecasting models, particularly when dealing with complex dynamical systems characterized by partial differential equations (PDEs) and chaotic behaviors. Motivated by this, we propose a novel probability path model tailored for probabilistic forecasting of dynamical systems, with a particular focus on PDE datasets relevant to scientific applications. Our model leverages the continuous dynamics intrinsic to spatio-temporal data by interpolating between consecutive sequential samples. This approach ensures better alignment with the constructed flow, leading to improved predictive performance, more stable training, and greater inference efficiency. Existing probability path models often fail to fully capture the continuous nature of spatio-temporal data, resulting in a misalignment with flow-based methods and suboptimal outcomes. Our proposed model addresses these limitations directly.

Building on previous approaches, we provide a theoretical framework and efficient algorithms tailored to probabilistic forecasting of high-dimensional spatio-temporal dynamics using flow matching in latent space. Within this framework, we demonstrate that our probability path model outperforms existing models across several forecasting tasks involving PDEs, achieving faster convergence during training and requiring fewer sampling steps during inference. These advances enhance the practicality of flow matching approaches for real-world applications, particularly in scenarios where computational resources and time are critical constraints. Our main contributions are the following.

- **Theoretical Framework and Efficient Algorithms:** We present a theoretical framework and efficient algorithms for applying flow matching in latent space to the probabilistic forecasting of dynamical systems (see Algorithms 1-2), extending the approach of (Lipman et al., 2022) and (Davtyan et al., 2023). Our approach is specifically tailored for time series data, enabling effective modeling of complex temporal dependencies inherent in dynamical systems.
- **Novel Probability Path Model:** Using dynamical optimal transport and the Schrödinger’s bridge perspective (see Theorem 3.2), we motivate and propose a new probability path model (see Section 3.3) specifically designed for modeling dynamical systems data. We provide theoretical insights to show that the variance of the vector field (VF) generating our proposed path can be lower than

that of the optimal transport (OT) VF proposed by (Lipman et al., 2022) for sufficiently correlated spatio-temporal samples (see Theorem 3.3). We further provide intuitions to understand why our model leads to smoother training loss curve and faster convergence compared to other models.

- **Empirical Validation:** We provide extensive empirical results to demonstrate that our proposed probability path model can outperform other models on several forecasting tasks involving PDEs (see Section 5). Our results demonstrate that the proposed probability path model outperforms existing flow matching models, achieving faster convergence during training and improved predictive performance across several evaluation metrics, while requiring fewer sampling steps.

## 2 Flow Matching for Probabilistic Forecasting

In this section, we first introduce the objective of probabilistic forecasting, and then we discuss how flow matching can be used for learning conditional distributions in latent space.

**Probabilistic forecasting framework.** Suppose that we are given a training set of  $n$  trajectories, with each trajectory of length  $m$ ,  $S_n = \{(\mathbf{x}^{1:m})^{(i)}\}_{i=1,\dots,n}$ , where  $(\mathbf{x}^{1:m})^{(i)} = ((x^1)^{(i)}, \dots, (x^m)^{(i)})$ , (with the  $(x^l)^{(i)} \in \mathbb{R}^d$ ), coming from an underlying continuous-time dynamical system. For simplicity, we denote the trajectories as  $\mathbf{x}^{1:m} = (x^1, \dots, x^m)$  unless there is a need to specify the corresponding  $n$ . The trajectories are observed at arbitrary time points  $\mathbf{t}_{1:m} = (t_1, \dots, t_m)$  such that  $x^i := x(t_i) \in \mathbb{R}^d$  and  $(x(t))_{t \in [t_1, t_m]}$  are the observed states of the ground truth system. In practice, we may have access to only few trajectories, i.e.,  $n$  is small or even  $n = 1$ , and the trajectories themselves may be observed at different time stamps.

The goal of probabilistic forecasting is to predict the distribution of the upcoming  $l$  elements given the first  $k$  elements, where  $m = l + k$ , i.e.:  $q(x^{k+1}, \dots, x^{k+l} | x^1, \dots, x^k) = \prod_{i=1}^l q(x^{k+i} | x^1, \dots, x^{k+i-1})$ . We propose to model each one-step predictive conditional distribution via a probability density path. Instead of using score-based diffusion models to specify the path, we choose flow matching, a simpler method to train generative models. With flow matching, we directly work with probability paths, and we can simply avoid reasoning about (forward) diffusion processes altogether. Moreover, we shall work in a latent space.

**Flow matching in latent space.** Let  $z^\tau = \mathcal{E}(x^\tau)$  for  $\tau = 1, \dots, m$ , where  $\mathcal{E}$  denotes a pre-trained encoder that maps from the data space to a lower dimensional latent space. Working in the latent space, our goal is to approximate the ground truth distribution  $q(z^\tau | x^1, \dots, x^{\tau-1})$  by the parametric distribution  $p(z^\tau | z^{\tau-1})$ , which can then be decoded as  $x^\tau = \mathcal{D}(z^\tau)$ . The latent dynamics can be modeled by an ODE:

$$\dot{Z}_t = u_t(Z_t), \quad (1)$$

where  $u_t$  is the (velocity) VF describing the instantaneous rate of change of the state at time  $t$ . Learning the dynamics of the system is equivalent to approximating the VF  $u_t$ . This can be done by regressing a neural network using the mean squared error (MSE) loss.

Following the idea of flow matching, we infer the dynamics of the system generating  $\mathbf{z}$  from the collection of latent observables by learning a time-dependent VF  $v_t : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ ,  $t \in [0, 1]$ , such that the ODE

$$\dot{\phi}_t(Z) = v_t(\phi_t(Z)), \quad \phi_0(Z) = Z, \quad (2)$$

defines a time-dependent diffeomorphic map (called a flow),  $\phi_t(Z) : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ , that pushes a reference distribution  $p_0(Z)$  towards the distribution  $p_1(Z) \approx q(Z)$  along some probability density path  $p_t(Z)$  and the corresponding VF  $u_t(Z)$ . In other words,  $p_t = [\phi_t]_* p_0$ , where  $[\cdot]_*$  denotes the push-forward operation. Here,  $q$  is the ground truth distribution,  $p$  denotes a probability density path, i.e.,  $p : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}_{>0}$ , and  $\int p_t(Z) dZ = 1$ . We also write  $Z_t = \phi_t(Z)$ ; and thus the ODE can be written as  $\dot{Z}_t = v_t(Z_t)$ ,  $Z_0 = Z$ . Typically the reference distribution  $p_0$  is chosen to be the standard Gaussian (Lipman et al., 2022; Liu et al., 2022).

In other words, the main goal of flow matching is to learn a deterministic coupling between  $p$  and  $q$  by learning a VF  $v_t$  such that the solution to the ODE (2) satisfies  $Z_0 \sim p$  and  $Z_1 \sim q$ . When  $\mathbf{Z} = (Z_t)_{t \in [0,1]}$  solves Eq. (2) for a given function  $v_t$ , we say that  $\mathbf{Z}$  is a flow with the VF  $v_t$ . If we have such a VF, then  $(Z_0, Z_1)$  is a coupling of  $(p, q)$ . If we can sample from  $p$ , then we can generate approximate samples from

the coupling by sampling  $Z_0 \sim p$  and numerically integrating Eq. (2). This can be viewed as a continuous normalizing flow (Chen et al., 2018).

If one were given a pre-defined probability path  $p_t(Z)$  and the corresponding VF  $u_t(Z)$  that generates the path, then one could parametrize  $v_t(Z)$  with a neural network  $v_t^\theta(Z)$ , with  $\theta$  the learnable parameter, and solve the least square regression by minimizing the flow matching loss:

$$\min_{\theta} \mathcal{L}_{fm}(\theta) := \mathbb{E}_{t, p_t(Z)} \omega(t) \|v_t^\theta(Z) - u_t(Z)\|^2, \quad (3)$$

where  $t \in \mathcal{U}[0, 1]$ ,  $Z \sim p_t(Z)$  and  $\omega(t) > 0$  is a weighting function. The weighting function weights the importance of the  $L^2$  loss at different times  $t$  (noise level), balancing the importance of high frequency and low frequency components. Following the standard framework (Lipman et al., 2022), we take  $\omega(t) = 1$ . Note that we choose to include  $\omega(t)$  in Eq. (3) for generality as different choice of weighting function corresponds to different parametrization of the network output (see (Anonymous, 2025) for more details).

However, we do not have prior knowledge for choosing  $p_t$  and  $u_t$ , and there are many choices that can satisfy  $p_1 \approx q$ . Moreover, we do not have access to a closed form  $u_t$  that generates the desired  $p_t$ . We shall follow the approach of (Lipman et al., 2022) and construct a target probability path by mixing simpler conditional probability paths. This probability path is the marginal probability path:

$$p_t(Z) = \int p_t(Z|\tilde{Z})q(\tilde{Z})d\tilde{Z}, \quad (4)$$

obtained by marginalizing the conditional probability density paths  $p_t(Z|\tilde{Z})$  over observed latent trajectories  $\tilde{Z}$ , with  $p_0(Z|\tilde{Z}) = p(Z)$  and  $p_1(Z|\tilde{Z}) = \mathcal{N}(Z|\tilde{Z}, \epsilon^2 I)$  for a small  $\epsilon > 0$ . Doing so gives us a marginal probability  $p_1$  which is a mixture distribution that closely approximates  $q$ . Then, assuming that  $p_t(Z) > 0$  for all  $Z$  and  $t$ , we can also define a marginal VF as:

$$u_t(Z) = \int u_t(Z|\tilde{Z}) \frac{p_t(Z|\tilde{Z})q(\tilde{Z})}{p_t(Z)} d\tilde{Z}, \quad (5)$$

where  $u_t(Z|\tilde{Z})$  is a conditional VF (conditioned on the latent trajectory  $\tilde{Z}$ ). It turns out that this way of mixing the conditional VFs leads to the correct VF for generating the marginal probability path (4). We can then break down the intractable marginal VF into simpler conditional VFs which depends on a single sample.

To deal with the intractable integrals in Eq. (4)-(5) which complicates computation of an unbiased estimator of  $\mathcal{L}_{fm}$ , we shall minimize the conditional loss proposed by (Lipman et al., 2022):

$$\min_{\theta} \mathcal{L}_{cfm}(\theta) := \mathbb{E}_{t, p_t(Z|\tilde{Z}), q(\tilde{Z})} \omega(t) \|v_t^\theta(Z) - u_t(Z|\tilde{Z})\|^2, \quad (6)$$

where  $t \in \mathcal{U}[0, 1]$ ,  $\tilde{Z} \sim q(\tilde{Z})$ ,  $Z \sim p_t(Z|\tilde{Z})$ , and  $u_t(Z|\tilde{Z})$  is the VF defined *per sample*  $\tilde{Z}$  that generates the conditional probability path  $p_t(Z|\tilde{Z})$ . Importantly, one can show that the solution of (6) is guaranteed to converge to the same result in (3); see Theorem D.1 in App. D. Therefore, the conditional flow matching loss can match the pre-defined target probability path, constructing the flow that pushes  $p_0$  towards  $p_1$ . Since both the probability path and VF are defined per sample, we can sample unbiased estimates of the conditional loss efficiently, particularly so with suitable choices of conditional probability paths and VFs.

### 3 Probability Path Models for Probabilistic Forecasting

In this section, we describe the family of probability paths that we consider for flow matching, and we propose an improved model tailored for probabilistic forecasting of spatio-temporal data.

#### 3.1 Common Probability Path Models

The family of Gaussian conditional probability paths gives us tractable choices to work with since the relevant quantities in Eq. (6) and thus the conditional flow can be defined explicitly. Therefore, we will work with

Table 1: Choices of probability density paths that we study in this paper. For VE-diffusion,  $\sigma_t$  is increasing in  $t$ ,  $\sigma_0 = 0$ . For VP-diffusion,  $\beta =$  noise scale.

Model	$a_t$	$b_t$	$c_t^2$
VE-diff. (Song et al., 2020b; Lipman et al., 2022)	1	0	$\sigma_{1-t}^2$
VP-diff. (Song et al., 2020b; Lipman et al., 2022)	$e^{-\frac{1}{2}T(1-t)}$	0	$1 - e^{-T(1-t)}, T(t) = \int_0^t \beta(s) ds$
OT-VF (Liu et al., 2022)	$t$	0	$(1 - (1 - \epsilon_{min})t)^2, \epsilon_{min} \geq 0$
SI (Chen et al., 2024)	$1 - t$	$t$ or $t^2$	$\epsilon^2 t(1-t)^2, \epsilon > 0$
<b>Ours</b>	$1 - t$	$t$	$\sigma_{min}^2 + \sigma^2 t(1-t), \sigma_{min}, \sigma \geq 0$

Gaussian probability paths. Moreover, we are going to solve (6) over the dataset of all transition pairs  $\mathcal{D}_{pair} = \{(z^{\tau-1}, z^\tau)\}_{\tau=2, \dots, m}$ , and use a pair of points for  $\tilde{Z}$ , setting  $\tilde{Z} = (Z_0, Z_1) \in \mathcal{D}_{pair}$ . In particular, we consider the following class of models for the probability paths:

$$p_t(Z|\tilde{Z} := (Z_0, Z_1)) = \mathcal{N}(Z|a_t Z_0 + b_t Z_1, c_t^2 I), \quad (7)$$

where  $a_t, b_t$  and  $c_t$  are differentiable time-dependent functions on  $[0, 1]$ , and  $I$  denotes the identity.

Table 1 provides five different choices of probability paths, including our proposed choice, that we study here. The optimal transport (OT) VF model was initially proposed by (Lipman et al., 2022), and setting  $\epsilon_{min} = 0$  gives us the rectified flow model of (Liu et al., 2022), which proposed connecting data and noise on a straight line. The stochastic interpolant (SI) model in Table 1 is the one considered by (Chen et al., 2024). The VE- and VP-diffusion conditional VFs (derived with Theorem D.2) coincide with the VFs governing the Probability Flow ODE for the VE and VP diffusion paths proposed in (Song et al., 2020b). It has been shown that combining diffusion conditional VFs with the flow matching objective leads to a training alternative that is more stable and robust when compared to existing score matching approaches (Lipman et al., 2022).

As remarked in (Lipman et al., 2022), there are many choices of VFs that generate a given probability path. We shall use the simplest VF that generates flow whose map is affine linear. Let  $p_t(Z|\tilde{Z})$  be the Gaussian probability path defined in Eq. (7) and consider the flow map  $\psi_t$  defined as  $\psi_t(Z) := a_t Z_0 + b_t Z_1 + c_t Z$  with  $c_t > 0$ . Then the unique VF that defines  $\psi_t$  is (see Theorem D.2 and the proof in App. D):

$$u_t(Z|\tilde{Z}) = \frac{c_t'}{c_t}(Z - (a_t Z_0 + b_t Z_1)) + a_t' Z_0 + b_t' Z_1, \quad (8)$$

where prime denotes derivative with respect to  $t$ , and  $u_t(Z|\tilde{Z})$  generates the Gaussian path  $p_t(Z|\tilde{Z})$ .

In view of this, minimizing the conditional loss becomes:

$$\min_{\theta} \mathcal{L}_{cfm}(\theta) := \mathbb{E}_{t, z^\tau, z^{\tau-1}, Z} \omega(t) \left\| v_t^\theta(Z) - \frac{c_t'}{c_t}(Z - (a_t z^\tau + b_t z^{\tau-1})) - a_t' z^\tau - b_t' z^{\tau-1} \right\|^2, \quad (9)$$

where  $t \sim \mathcal{U}[0, 1]$ ,  $Z \sim p_t(Z|z^\tau, z^{\tau-1})$  and  $\mathbf{z} \sim q(\mathbf{z})$ . We refer to this as the Flow Matching loss parametrization and shall work with this parametrization. There are other parametrizations: most popular ones are the Score Matching loss, Score Flow loss and DDPM loss. See App. C for a comparison of different loss parametrizations and App. B for connections to SDE based generative models.

In practice, it may be beneficial to learn the vector field  $v_t^\theta$  using additional context information, in which case (9) becomes:

$$\min_{\theta} \mathcal{L}_{cfm}(\theta) := \mathbb{E}_{t, z^\tau, z^{\tau-1}, Z, C} \omega(t) \left\| v_t^\theta(Z|C) - \frac{c_t'}{c_t}(Z - (a_t z^\tau + b_t z^{\tau-1})) - a_t' z^\tau - b_t' z^{\tau-1} \right\|^2, \quad (10)$$

where  $C$  represents the context information. The choice of  $C$  is task-dependent. For spatio-temporal tasks, we adopt the sparse conditioning scheme of Davtyan et al. (2023), choosing  $C$  to be past frames as follows. Given a sampled frame  $z^\tau$ , we sample another index  $c$  uniformly from  $\{1, \dots, \tau - 2\}$  and use  $z^c$ , which we call context frame, together with the previous frame  $z^{\tau-1}$ , as the two conditioning frames. Thus, in this case  $C = (z^{\tau-1}, z^c, \tau - c)$  (see Algorithm 1). While conditioning on as many past frames as possible

is desirable to improve prediction accuracy, sparse conditioning is sufficient to achieve favorable trade-off between computational efficiency and accuracy (Davtyan et al., 2023). Datasets with more challenging dynamics may require the use of more context frames.

### 3.2 Optimal Probability Paths

We consider the problem of selecting the optimal probability paths within the class of the Gaussian probability paths. We shall make use of the Schrödinger bridge (Léonard, 2013; Chen et al., 2021b) perspective, and seek to find the *optimal* stochastic processes that evolve a given measure into another, subject to marginal constraints and based on a *prior belief*. This optimal process describes a novel probability path model that we propose for probabilistic forecasting of spatio-temporal data.

To be more precise, let  $\nu, \nu'$  be two given probability measures and let  $\mathbb{Q}$  be the path measure of an arbitrary stochastic process. The Schrödinger bridge (SB) is the solution to the following constrained minimization problem over all path measures  $\mathbb{P}$  (which are absolutely continuous with respect to  $\mathbb{Q}$ ) of stochastic processes on the finite time interval  $[0, 1]$ :

$$\min_{\mathbb{P}_0=\nu, \mathbb{P}_1=\nu'} D_{KL}(\mathbb{P}||\mathbb{Q}), \quad (11)$$

where  $D_{KL}$  denotes the Kullback–Leibler divergence and  $\mathbb{P}_t$  denotes the time marginal of  $\mathbb{P}$  at time  $t$ . Typically  $\nu$  and  $\nu'$  are the (empirical) marginal distributions of an unknown continuous-time dynamics observed at the initial and terminal times, and  $\mathbb{Q}$  is the path measure of a prior (or reference) process that represents our belief of the dynamics before observing the data. The solution  $\mathbb{P}^*$  to the problem can then be interpreted as the optimal dynamics that conforms to the prior belief  $\mathbb{Q}$  while respecting the data marginals  $\mathbb{P}_0^* = \nu$ ,  $\mathbb{P}_1^* = \nu'$ . In other words, the SB is the (path measure of the) finite-time diffusion which admits as initial and terminal distributions the two distributions of interest and is the closest in KL divergence to (the path measure of) a reference diffusion. Recent work (Shi et al., 2024; Pooladian & Niles-Weed, 2024) focuses on developing improved algorithms to solve SB problems for general classes of reference diffusions. SBs have also been used to formulate generative models by interpolating distributions on a finite time interval (Wang et al., 2021; De Bortoli et al., 2021; Chen et al., 2021a; Peluchetti, 2023; Tong et al., 2023a; Liu et al., 2023; Gottwald & Reich, 2024).

We will consider Gaussian prior processes for  $\mathbb{Q}$ . The Gaussian probability paths that we consider can be formulated via the differential equation:  $dZ_t = (\dot{a}_t Z_0 + \dot{b}_t Z_1 + \dot{c}_t \xi)dt$ , where  $\xi \sim \mathcal{N}(0, I)$ . Note that  $\mathbb{E}[\dot{Z}_t|Z_0, Z_1] = \dot{a}_t Z_0 + \dot{b}_t Z_1 =: \alpha_t$ . We shall take  $\mathbb{Q}$  to be the path measure of the linear SDE with drift  $\alpha_t$ :

$$dY_t = \alpha_t dt + \omega dW_t, \quad t \in [0, 1], \quad (12)$$

where  $\omega > 0$  is a regularity parameter and  $W_t$  is the standard Wiener process. While many choices for  $\mathbb{Q}$  exists, the rationale behind our choice is that the Gaussian process  $Y_t$  admits a minimal representation  $Y_t = \mathbb{E}[Z_t|Z_0, Z_1] + \omega W_t$  that incorporates Brownian motion as a reference noise to characterize the stochasticity surrounding the conditional estimation of  $Z_t$ . Importantly, such choice allows interpretation of the SBs as generalized dynamical optimal transport (DOT) (Chewi et al., 2024) between two (not necessarily Gaussian) measures. In particular, adapting Theorem 1 in (Bunne et al., 2023) to our setting, we have:

**Proposition 3.1.** *Consider the SB problem with  $Y_t$  as the reference process:*

$$\min_{\mathbb{P}_0=\nu, \mathbb{P}_1=\nu'} D_{KL}(\mathbb{P}||\mathbb{Y}), \quad (13)$$

where  $\mathbb{Y}$  is the path measure induced by  $(Y_t)_{t \in [0,1]}$ . Then, (13) is equivalent to  $\inf_{(\rho_t, v_t)} \mathbb{E} \left[ \int_0^1 C(\rho_t, v_t) dt \right]$ , with

$$C(\rho_t, v_t) := \frac{\|v_t\|^2}{2\omega^2} + \frac{\omega^2}{8} \|\nabla \log \rho_t\|^2 - \frac{1}{2} \langle \alpha_t, \nabla \log \rho_t \rangle, \quad (14)$$

where the infimum is taken over all pairs  $(\rho_t, v_t)$  such that  $\rho_0 = \nu$ ,  $\rho_1 = \nu'$ ,  $\rho_t$  absolutely continuous, and satisfies  $\partial_t \rho_t = -\nabla \cdot ((v_t + \alpha_t)\rho_t)$ .

The DOT problem (14) can be seen as a generalized version of DOT with quadratic cost (minimizing a "kinetic energy"); see Eq. (4.41a)-(4.41c) in (Lipman et al., 2024). The GDOT (14) not only minimizes the

quadratic cost, but also minimizes the Fisher information and maximizes the expected alignment of the score function with the drift  $\alpha_t$ .

In our case, given a pair  $(Z_0, Z_1)$  of data points, we are interested in constructing SBs where the marginal constraints are Gaussians centered around  $Z_0$  and  $Z_1$ , which we denote as  $\xi_0 := \mathcal{N}(Z_0, \sigma_{min}^2 I)$  and  $\xi_1 := \mathcal{N}(Z_1, \sigma_{min}^2 I)$  respectively, for some given  $\sigma_{min} \geq 0$ . These are the Gaussian SBs:

$$\min_{\mathbb{P}_0=\xi_0, \mathbb{P}_1=\xi_1} D_{KL}(\mathbb{P} \parallel \mathbb{Y}). \quad (15)$$

It turns out that the solution of these Gaussian SBs admits a closed-form expression (Bunne et al., 2023).

**Theorem 3.2.** *The solution  $\mathbb{P}^*$  to the Gaussian SB (15) is the path measure of a Markov Gaussian process with the marginal variable  $X_t \sim \mathcal{N}(\mu_t, \sigma_t^2 I)$ , where*

$$\mu_t = (1 - t + a_t - a_0 - t(a_1 - a_0))Z_0 + (t + b_t - b_0 - t(b_1 - b_0))Z_1, \quad (16)$$

$$\sigma_t^2 = \sigma_{min}^2 + \sigma^2 t(1 - t), \quad (17)$$

with  $\sigma^2 = \sqrt{4\sigma_{min}^4 + \omega^4} - 2\sigma_{min}^2 > 0$ .

See App. D.3 for the proof. In particular, imposing the boundary constraints  $a_0 = b_1 = 1, a_1 = b_0 = 0$  on the functions  $a_t$  and  $b_t$ , we have  $\mu_t = a_t Z_0 + b_t Z_1$  in Theorem 3.2. Note that  $\sigma_t^2$  is independent of  $a_t$  and  $b_t$ .

### 3.3 A Novel Probability Path Model

Our proposed probability path model is a stochastic extension of the straight-line trajectory connecting consecutive latent samples  $(Z_0, Z_1)$ , i.e.,  $\mu_t = (1 - t)Z_0 + tZ_1$ . This choice of  $a_t$  and  $b_t$  can be motivated by the following optimality principle. Given a pair of data points  $(Z_0, Z_1) \in \mathbb{R}^d \times \mathbb{R}^d$ , the linear interpolation path  $\mu_t$  arises naturally as the solution to the following variational problem:

$$\min_{\mu: [0,1] \rightarrow \mathbb{R}^d} \int_0^1 \|\dot{\mu}_t\|^2 dt \quad \text{subject to } \mu_0 = Z_0, \mu_1 = Z_1. \quad (18)$$

This is a classical energy minimization problem, where the goal is to find the smoothest curve connecting the endpoints  $Z_0$  and  $Z_1$ . It can be solved using the Euler-Lagrange equations under the given boundary conditions. This formulation corresponds to minimizing the kinetic energy of the curve among all smooth interpolations between  $Z_0$  and  $Z_1$ . The solution, which is constant-velocity motion along the straight line, has a natural physical interpretation: it is the least effort way to move from  $Z_0$  to  $Z_1$  in Euclidean space.

More generally, when the boundary conditions are probability distributions (e.g., Gaussians centered at  $Z_0$  and  $Z_1$ ), this linear interpolation corresponds to the displacement interpolation in the 2-Wasserstein space. In the case where the marginals are Gaussian and the cost function is quadratic, the interpolating path of distributions is also Gaussian with linearly interpolated means, minimizing an upper bound on the expected kinetic energy among all stochastic paths with prescribed marginals at  $t = 0$  and  $t = 1$ ; see Section 4.7 in (Lipman et al., 2024) for more details. Alternative interpolants are possible, but the linear form admits analytical tractability and natural interpretation in terms of velocity fields induced by diffusion processes. Importantly, this leads to trajectories that are generally easier to sample with ODE solvers, leading to ODE simulations with smaller errors (Liu et al., 2022).

For the stochastic extension, we build on the deterministic backbone and relax the boundary constraints to Gaussians centered around  $Z_0$  and  $Z_1$ . We shall look for the "most likely" probability paths connecting the boundary distributions as follows. First, we consider the noise perturbed paths  $Z_t = \mu_t + \omega W_t$ , where  $W_t$  is the standard Wiener process, as our reference process. Then, we take the solution to the Schrödinger bridge (15) with respect to the reference process as our proposed probability path model.

More precisely, in view of Theorem 3.2, we propose to choose  $a_t = 1 - t, b_t = t, c_t^2 = \sigma_{min}^2 + \sigma^2 t(1 - t)$ , in which case we have the probability path described by:

$$p_t(Z|\tilde{Z}) = \mathcal{N}(Z|tZ_1 + (1 - t)Z_0, (\sigma_{min}^2 + \sigma^2 t(1 - t))I), \quad (19)$$

which transports a Gaussian distribution centered around  $Z_0$  with variance  $\sigma_{min}^2$  at  $t = 0$  to a Gaussian distribution centered around  $Z_1$  at  $t = 1$  with variance  $\sigma_{min}^2$ . Here  $\sigma_{min}, \sigma \geq 0$  are tunable parameters. In the case when  $\sigma_{min} = 0$ , it describes a Brownian bridge that interpolates between  $Z_0$  and  $Z_1$  (Gasbarra et al., 2007). To ensure numerical stability when sampling  $t \sim \mathcal{U}[0, 1]$ , it is beneficial to use a small  $\sigma_{min} > 0$ . Note that  $\sigma^2$  is a scale factor determining the magnitude of fluctuations around the path interpolating between  $Z_0$  and  $Z_1$ . The variance  $c_t^2$  is minimum with the value of  $\sigma_{min}^2$  at the endpoints  $t = 0$  and  $t = 1$ , and the maximum variance is  $\sigma_{min}^2 + \sigma^2/4$  which occurs in the middle of the path at  $t = 1/2$ .

The variance schedule is designed to balance exploration and stability. Low variance at the start ensures stable initialization, preventing the trajectory from deviating too far from the initial distribution. High variance in the middle allows the model to explore diverse paths in the latent space, avoiding mode collapse and enhancing diversity in the generated trajectories. Low variance at the end sharpens the trajectory, ensuring accurate reconstruction of the desired output. This strategy is inspired by findings in diffusion models that utilize a forward noising process and a backward denoising process (Ho et al., 2020; Song et al., 2020b), where such variance patterns have been shown to effectively manage the trade-off between exploration and refinement.

The corresponding VF that defines the flow is then given by (applying Theorem D.2):

$$u_t(Z|\tilde{Z}) = Z_1 - Z_0 + \frac{\sigma^2}{2} \frac{1 - 2t}{\sigma_{min}^2 + \sigma^2 t(1 - t)} \epsilon, \quad (20)$$

where  $\epsilon := Z - (tZ_1 + (1 - t)Z_0)$ . Under reasonable assumptions, the variance of this VF is lower than the variance of the OT-VF of (Lipman et al., 2022) with  $\epsilon_{min} := 0$  (rectified flow). Here the variance is taken with respect to the randomness in the samples  $z^\tau$  and the Gaussian samples drawn during gradient descent updates.

More precisely, denote the VF that corresponds to our model as  $u_t$  and the VF that corresponds to the rectified flow model as  $\tilde{u}_t$ . For a given  $\tau$ , they generate the probability path  $Z_t = tz^\tau + (1 - t)z^{\tau-1} + c_t\xi$  and  $\tilde{Z}_t = tz^{\tau-1} + (1 - t)\eta$  respectively, where  $\xi, \eta \sim \mathcal{N}(0, I)$ ,  $c_t = \sqrt{\sigma_{min}^2 + \sigma^2 t(1 - t)}$  and  $t \in [0, 1]$ . Applying Eq. (8), we have:

$$u_t(Z_t|z^{\tau-1}, z^\tau) = z^\tau - z^{\tau-1} + c_t'\xi, \quad (21)$$

$$\tilde{u}_t(\tilde{Z}_t|z^{\tau-1}) = z^{\tau-1} - \eta. \quad (22)$$

**Theorem 3.3.** *Suppose that  $(z^\tau)_{\tau=1, \dots, m}$ , with the  $z^\tau \in \mathbb{R}^d$ , is a discrete-time stochastic process with nonzero correlation in time and let  $t \in [0, 1]$  be given. If  $\text{Cov}(z^{\tau-1}, z^\tau) \geq \frac{1}{2} \left( \left( \frac{\sigma^4}{4\sigma_{min}^2} - 1 \right) I + \text{Var}(z^\tau) \right)$ , then  $\text{Var}(\tilde{u}_t(\tilde{Z}_t|z^{\tau-1})) \geq \text{Var}(u_t(Z_t|z^{\tau-1}, z^\tau))$ .*

See App. D.4 for a proof and related discussions. Theorem 3.3 implies that if the consecutive latent variables  $z^\tau, z^{\tau-1}$  are sufficiently correlated and  $\sigma$  is chosen small enough relative to  $\sigma_{min}$ , then the variance of the VF that corresponds to our probability path model is lower than that corresponds to the rectified flow model.

## 4 An Efficient Probabilistic Forecasting Algorithm

In this section, we present efficient algorithms for training and inferencing the flow matching model.

Recently (Davtyan et al., 2023) proposed an efficient algorithm for latent flow matching for the task of video prediction, using the probability path generated by the OT-VF of (Lipman et al., 2022). To enable efficient training, we shall follow (Davtyan et al., 2023) and leverage the iterative nature of sampling from the learned flow and use a single random conditioning element from the past at each iteration. However, our method differs from (Davtyan et al., 2023) as we use different probability paths and target VFs.

**Training.** Recall that  $\tilde{Z} = (Z_0, Z_1)$  denote training samples from  $\mathcal{D}_{pair} = \{(z^{\tau-1}, z^\tau)\}_{\tau=2, \dots, m}$ . In other words, we set  $Z_1$  to be the target element and  $Z_0$  to be the reference element chosen to be the previous element before the target element. In this way, our target probability path model maps a distribution centered



around a previous state to the distribution of the current states, which is more natural from the point of view of probabilistic forecasting whose goal is to obtain an ensemble of forecasts. Note that this differs from the approach of [Davtyan et al. \(2023\)](#), where  $\tilde{Z} = Z_1$  (i.e., they do not use a reference element to define their probability path, whereas we use a pair of elements  $(Z_0, Z_1)$ ). Algorithm 1 summarizes the training procedure of our method.

Both the autoencoder and the VF neural network can also be jointly trained in an end-to-end manner, but our empirical results show that separating the training can lead to improved performance. Moreover, doing so allows us to better assess the impact of using different probability paths.

---

**Algorithm 1** Flow matching for spatio-temporal data

---

**Input:** Dataset of sequences  $D$ , number of iterations  $M$

**for**  $i$  in range(1,  $M$ ) **do**

    Sample a sequence  $\mathbf{x}$  from the dataset  $D$

    Encode it with a pre-trained encoder to obtain  $\mathbf{z}$

    Choose a random target element  $z^\tau, \tau \in \{3, \dots, |\mathbf{z}|\}$ , from  $\mathbf{z}$

    Sample a step  $t \sim U[0, 1]$

    Sample a noisy observation  $Z \sim p_t(Z | z^\tau, z^{\tau-1})$ , where  $p_t$  is given by Eq. (7)

    Compute  $u_t(Z | z^\tau, z^{\tau-1})$  according to (20)

    Sample a condition frame  $z^c, c \in \{1, \dots, \tau - 2\}$

    Update the parameters  $\theta$  via gradient descent

$$\nabla_{\theta} \|v_t^{\theta}(Z | z^{\tau-1}, z^c, \tau - c) - u_t(Z | z^{\tau}, z^{\tau-1})\|^2$$

**end for**

**Return:** A learned VF,  $v_t^{\theta^*}$

---

**Inference.** We use an ODE sampler during inference to generate forecasts. The ODE sampler is described as follows. Let  $(Y_i^\tau)_{i=1, \dots, N-1}$  denote the generation process, where  $N - 1$  is the number of integration steps and the superscript  $\tau$  denotes the time index for which the generation/forecast is intended for. Given the previous elements  $(x^1, \dots, x^{T-1})$  of a time series sample, in order to generate the next element (i.e., the  $T$ -th element), we start with sampling the initial condition  $Y_0^T$  from  $\mathcal{N}(z^{T-1}, \sigma_{sam}^2 I)$  for some small  $\sigma_{sam} \geq 0$ , where  $z^{T-1} = \mathcal{E}(x^{T-1})$ . This is in contrast to the procedure of ([Davtyan et al., 2023](#)), which simply uses a mean-zero Gaussian sample instead. We then use an ODE solver to integrate the learned VF along the time interval  $[0, 1]$  to obtain  $Y_{N-1}^T$ . We use  $\mathcal{D}(Y_{N-1}^T)$  as an estimate of  $x^T$ , and forecasting is done autoregressively.

Algorithm 2 summarizes this procedure when the sampling is done using the forward Euler scheme. Note that we can also use computationally more expensive numerical schemes such as the Runge-Kutta (RK) schemes.

---

**Algorithm 2** One-step ahead forecasting with forward Euler

---

**Input:** A sequence  $(x^1, \dots, x^{T-1})$  containing the previous elements, number of integration steps  $N - 1$ , grid  $s_0 = 0 < s_1 < \dots < s_{N-1} = 1$ , a learnt VF  $v_s^{\theta^*}$  for  $s \in [0, 1]$

Set  $\Delta s_n = s_{n+1} - s_n$  for  $n = 0, \dots, N - 2$

Sample  $Y_0^T \sim \mathcal{N}(\mathcal{E}(x^{T-1}), \sigma_{sam}^2 I)$ ,  $\sigma_{sam} \geq 0$

**for**  $n$  in range(0,  $N - 1$ ) **do**

    Sample  $c \sim \mathcal{U}(2, \dots, T - 1)$

$y^{T-c} = \mathcal{E}(x^{T-c})$

$Y_{n+1}^T = Y_n^T + \Delta s_n v_{s_n}^{\theta^*}(Y_n^T | Y_0^T, y^{T-c}, T - c)$

**end for**

**Return:** An estimate of  $x^T$ ,  $\hat{x}^T = \mathcal{D}(Y_{N-1}^T)$

---

## 5 Empirical Results

In this section, we present our main empirical results to elucidate the design choice of probability paths within the flow matching framework (comparison with other generative modeling frameworks is not our focus here). We focus on PDE dynamics forecasting tasks here (additional details and results can be found in App. E-F). We test the performance of our probability path model, i.e., Eq. (19) with  $a_t = 1 - t$ ,  $b_t = t$  and  $c_t = \sqrt{\sigma_{min}^2 + \sigma^2 t(1-t)}$  on these tasks. We pick  $\sigma_{min} = 0.001$ , and treat  $\sigma$  and  $\sigma_{sam}$  as tunable parameters. We compare our proposed model with four other models of probability paths:

- **RIVER** (Davtyan et al., 2023): RIVER uses the OT-VF model in Table 1, i.e.,  $a_t = 0$ ,  $b_t = t$ ,  $c_t = 1 - (1 - \epsilon_{min})t$ , choosing  $\epsilon_{min} = 10^{-7}$ .
- **VE-diffusion** in Table 1: We use  $\sigma_t = \sigma_{min} \sqrt{\left(\frac{\sigma_{max}}{\sigma_{min}}\right)^{2t} - 1}$  with  $\sigma_{min} = 0.01$ ,  $\sigma_{max} = 0.1$ , and sample  $t$  uniformly from  $[0, 1 - \epsilon]$  with  $\epsilon = 10^{-5}$ , following (Song et al., 2020b).
- **VP-diffusion** in Table 1: We use  $\beta(s) = \beta_{min} + s(\beta_{max} - \beta_{min})$  where  $\beta_{min} = 0.1$ ,  $\beta_{max} = 20$  and  $t$  is sampled from  $\mathcal{U}[0, 1 - \epsilon]$  with  $\epsilon = 10^{-5}$ , following (Song et al., 2020b). Thus,  $T(s) = s\beta_{min} + \frac{1}{2}s^2(\beta_{max} - \beta_{min})$ .
- **The stochastic interpolant (SI) path** in Table 1: We consider the path proposed by (Chen et al., 2024) and use the suggested choice of  $a_t = 1 - t$ ,  $b_t = t^2$  and  $c_t = \epsilon(1 - t)\sqrt{t}$  (see Eq. (2) in (Chen et al., 2024) and note that  $Var((1 - t)W_t) = (1 - t)^2 t$  for the standard Wiener process  $W_t$ ). We also consider the choice with  $b_t = t$  instead. We choose  $\epsilon = 0.01$  for both choices. This is a path that is similar to ours, but with the variance  $c_t^2$  chosen such that the maximum occurs at  $t = 1/\sqrt{3}$  instead of at the middle of the path at  $t = 1/2$ . We shall see that different forms of variance can lead to vastly different performance in the considered tasks.

**Evaluation metrics.** We evaluate the models using the following metrics. First, we use the standard mean squared error (MSE) and the relative Frobenius norm error (RFNE) to measure the difference between predicted and true snapshots. Second, we compute metrics such as the peak signal-to-noise ratio (PSNR), and the structural similarity index measure (SSIM) to further quantify the quality and similarity of the generated snapshots (Wang et al., 2004). Third, we use the Pearson correlation coefficient to assess the correlation between predicted and true snapshots.

**Training details.** We use an autoencoder (AE) to embed the training data into a low-dimensional latent space, which enables the model to capture the most relevant features of the data while reducing dimensionality (Azencot et al., 2020); see App. F for further discussion of the motivation. We then train a flow matching model in this latent space. Training generative models in latent space has also been considered by (Vahdat et al., 2021) for score matching models and by (Dao et al., 2023) for flow matching models. To train the AE, we minimize a loss function that consists of reconstruction error, in terms of MSE, between the input data and its reconstructed version from the latent space. The choice of the AE architecture is tailored to the complexity of the dataset (see App. F for details).

### 5.1 Probabilistic Forecasting of Dynamical Systems

We evaluate the performance of our proposed probability path model on challenging dynamical systems to demonstrate its effectiveness in forecasting complex continuous dynamics. Specifically, we consider the following tasks (for details see App. F.2):

- **Fluid Flow Past a Cylinder (FPC):** This task involves forecasting the vorticity of a fluid flowing past a cylinder. The model conditions on the first 5 frames and predicts the subsequent 20 frames at a resolution of  $64 \times 64$  with 1 channel representing vorticity.
- **Shallow-Water Equation (SWE):** This dataset models the dynamics of shallow-water equations (Takamoto et al., 2022), capturing essential aspects of geophysical fluid flows. We use the first 5 frames for conditioning and predict the next 15 frames at a resolution of  $128 \times 128$  with 1 channel representing horizontal flow velocity.

- **Diffusion-Reaction Equation (DRE):** This dataset models the dynamics of a 2D diffusion-reaction equation (Takamoto et al., 2022). We use the first 5 frames for conditioning and predict 15 future frames at a  $128 \times 128$  with 2 channels representing velocity in the  $x$  and  $y$  directions.
- **Incompressible Navier-Stokes Equation (NSE):** As a more challenging benchmark, we consider forecasting the dynamics of a 2D incompressible Navier-Stokes equation (Takamoto et al., 2022). We use the first 5 frames for conditioning and predict the next 20 frames at a resolution of  $512 \times 512$  with 2 channels representing velocity in the  $x$  and  $y$  directions.

Table 2 summarizes the performance of our model compared to other probability path models across all tasks. It can be seen that our probability path model achieves the lowest test MSE and RFNE across all tasks, indicating more accurate forecasts. Moreover, the higher PSNR and SSIM scores indicate that our model better preserves spatial structures in the predictions. Despite the similarity of our proposed model with the stochastic interpolant of (Chen et al., 2024), in that both models use consecutive samples to define the path, our model outperforms the stochastic interpolant model, suggesting that choosing the maximum variance to occur at the middle of the path is a better choice. Importantly, our model is highly efficient during inference time since it requires very few sampling steps; this is significantly lower compared to other models.

Figure 1 shows the Pearson correlation coefficients of the predicted snapshots over time for all models. Our model’s predictions shows a slower decay of correlation coefficients compared to other models, indicating better temporal consistency and long-term predictive capability. Correlation coefficients about 95% indicate performance on par with physics-based numerical simulators.

Table 2: Results for forecasting dynamical systems using different probability path models for flow matching. Results are averaged over 5 generations obtained with 9 sampling steps ( $N = 10$ ) using RK4. For our model, we use  $\sigma_{min} = 0.001$  and  $\sigma_{sam} = 0$  for all the considered tasks.

Task	Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
Flow past Cylinder	RIVER	1.86e-03	4.48e-02	44.30	0.99
	VE-diffusion	2.29e-01	4.74e-01	27.36	0.55
	VP-diffusion	3.58e-03	6.09e-02	42.37	0.98
	SI ( $b_t = t^2$ )	3.40e-03	6.10e-02	41.81	0.98
	SI ( $b_t = t$ )	1.19e-02	9.96e-02	39.64	0.99
	<b>Ours (<math>\sigma = 0.01</math>)</b>	<b>3.79e-04</b>	<b>2.30e-02</b>	<b>48.88</b>	<b>1.00</b>
Shallow-Water	RIVER	9.18e-04	1.49e-01	34.92	0.92
	VE-diffusion	1.32e-02	5.66e-01	28.10	0.55
	VP-diffusion	1.39e-03	1.81e-01	34.07	0.87
	SI ( $b_t = t^2$ )	1.05e-03	1.53e-01	35.59	0.89
	SI ( $b_t = t$ )	6.74e-04	1.29e-01	36.08	0.93
	<b>Ours (<math>\sigma = 0.1</math>)</b>	<b>6.60e-04</b>	<b>1.28e-01</b>	<b>36.10</b>	<b>0.93</b>
Diffusion-Reaction	RIVER	2.87e-03	2.28e-01	38.12	0.82
	VE-diffusion	1.04e-01	1.71	32.98	0.36
	VP-diffusion	2.02e-02	7.03e-01	34.98	0.52
	SI ( $b_t = t^2$ )	6.17e-02	8.62e-01	<b>45.68</b>	0.76
	SI ( $b_t = t$ )	3.72e-04	1.18e-01	34.24	0.89
	<b>Ours (<math>\sigma = 0</math>)</b>	<b>3.56e-04</b>	<b>1.16e-01</b>	34.34	<b>0.90</b>
Navier-Stokes	RIVER	2.84e-02	8.67e-01	<b>30.75</b>	0.63
	VE-diffusion	1.58e-01	2.31	26.90	0.33
	VP-diffusion	2.09e-01	2.48	27.96	0.30
	SI ( $b_t = t^2$ )	1.27e-03	2.66e-01	30.73	0.90
	SI ( $b_t = t$ )	1.13e-03	2.54e-01	30.66	0.93
	<b>Ours (<math>\sigma = 0.1</math>)</b>	<b>1.13e-03</b>	<b>2.53e-01</b>	30.66	<b>0.93</b>

Figure 2 compares the training loss curves of our model with others for the FPC and the SWE task. Our method leads to faster convergence during training and smoother loss curves. This suggests that our model requires fewer iterations to generate high-quality samples when compared to other models. We find that there is no significant difference in training time between the models when trained for the same number of epochs.

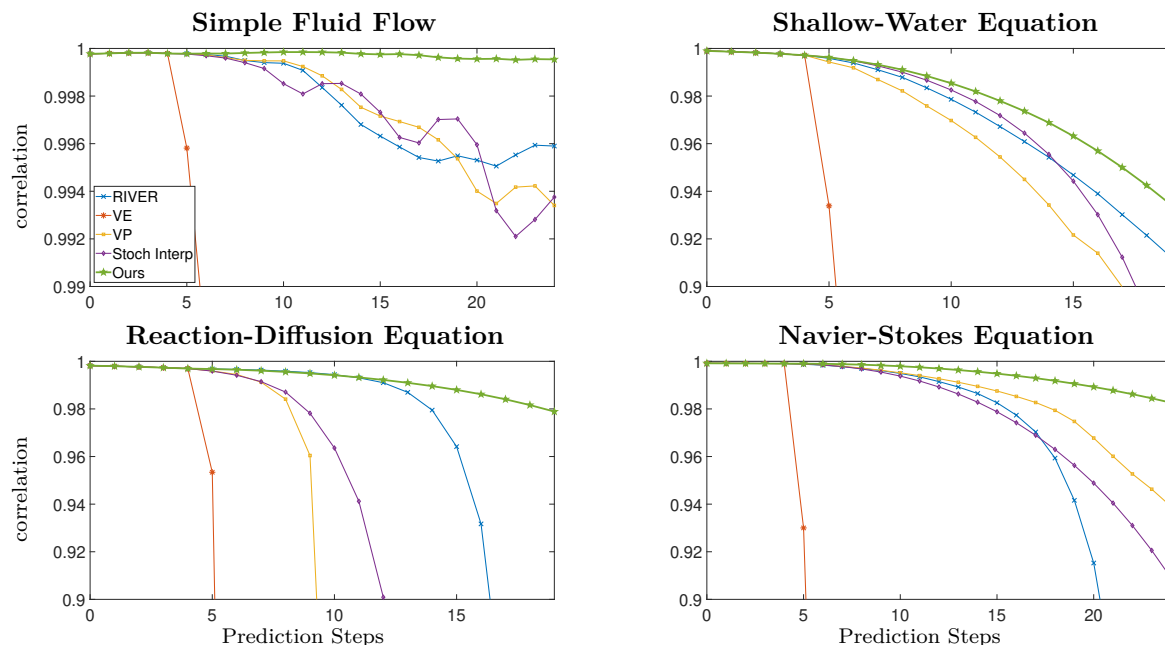


Figure 1: Pearson correlation coefficient to assess the correlation between predicted and true snapshots at various prediction steps for different probability path models. Our probability path model shows the best performance on all three tasks. Note that the first 5 snapshots are the conditioning snapshots.

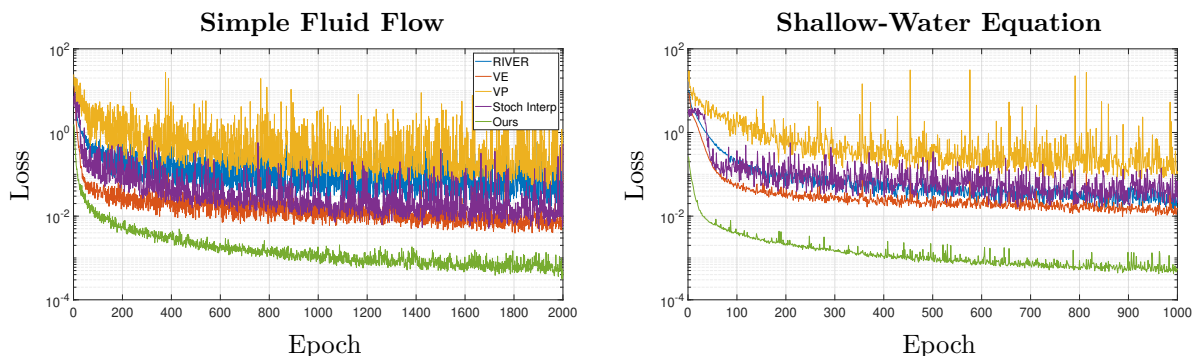


Figure 2: Training loss for different models of probability path for the fluid flow past a cylinder task. Our model leads to fastest convergence and smoothest loss curve among all models.

## 5.2 Discussion

Our empirical results showed that our proposed model consistently outperformed other models across several forecasting tasks involving different PDEs. Our model shows improved training efficiency, with faster convergence reducing the computational resources and time required for model training. Moreover, our model is efficient during inference time since it only requires a few sampling steps, making it practical for real-world applications where computational efficiency is crucial. Additionally, the model maintained better temporal consistency, as indicated by a slower decay of Pearson correlation coefficients over longer prediction horizons. These findings validate the effectiveness of our approach in modeling complex dynamical systems. Our focus in this paper is on near- to mid-term forecasting (15–20 steps), which is a standard set-up in spatio-temporal video and PDE benchmarks. While our framework is compatible with longer rollouts in principle, scaling to significantly longer horizons remains a long-standing challenge, as for most forecasting tasks, due to compounding errors and distributional drift, especially in high-dimensional dynamics. This is an important direction for future work.

Our ablation study (see App. E.2) further validated the advantages of our proposed probability path model. We found that larger  $\sigma$  values not only contributed to smoother training loss curves but also enhanced the overall stability and efficiency of the model. Moreover, we saw that  $\sigma_{\text{sam}} = 0$  can be fixed without compromising accuracy. We also provide studies to understand the effect of using context frames and the choice of  $\sigma_{\text{min}}$  on test performance in App. E.2. In particular, in Table 5 we see noticeable improvement in test performance using  $\sigma_{\text{min}} > 0$ , showing the advantages of going beyond the commonly used Brownian bridge model (see (Tong et al., 2023b;a)). Importantly, the fact that our model achieved improved performance even with the simplest sampler (Euler scheme) and a minimal number of sampling steps (as few as one for the FPC task; see Figure 4) validates its practical applicability, especially in scenarios where computational resources and time are limited.

Lastly, we provide an expanded discussion to position the proposed probability path relative to the stochastic interpolant (SI) path. While the SI path offers simplicity and analytical tractability, it lacks principled guidance for choosing  $b_t$ , which is an important design choice. Our proposed probability path is grounded in the dynamical SB framework, which, in our formulation, models the most likely stochastic evolution between two Gaussian distributions centered around time-adjacent samples (with a minimal variance  $\sigma_{\text{min}}^2$ ) under prior dynamics. From this perspective, the path is not simply a heuristic interpolation and the choice of  $b_t = t$  is justified, whereas the inclusion of a non-zero  $\sigma_{\text{min}}^2$  is important in the context of probabilistic forecasting. Note that when  $\sigma_{\text{min}} = 0$ , our proposed model simplifies to the SI model with  $b_t = t$ . A key implication of this viewpoint is that the variance structure of the path plays a crucial role: it reflects the uncertainty in matching the two endpoints under the prior dynamics. Our proposed path introduces a variance profile that is symmetric and peaks at the midpoint  $t = 1/2$ , better reflecting the intrinsic uncertainty in interpolating between endpoints. In contrast, the preferred SI path suggested in (Chen et al., 2024) (with  $b_t = t^2$ ) has an asymmetric variance peak at  $t = 1/\sqrt{3}$ , which may misalign with the data geometry. Empirically, we see that our probability path model mostly outperforms the SI models across the considered tasks (see Table 2). As expected, when  $\sigma_{\text{min}}$  is very close to zero, the test performance of our model is only marginally better than that of the SI model with  $b_t = t$  for most tasks. As discussed earlier (see also App. E.2), using bigger values of  $\sigma_{\text{min}} > 0$  can not only improve test performance but also help improve training stability.

## 6 Conclusion

In this work, we investigated the use of flow matching in latent space for probabilistic forecasting of spatio-temporal dynamics, providing a theoretical framework and efficient algorithms. We demystified the critical role of the probability path design in this setting and proposed an improved probability path model. Our model is theoretically motivated via the SB and dynamical optimal transport perspective. It leverages the inherent continuity and correlation in the spatio-temporal data, leading to more stable training and faster convergence. Our empirical evaluations on several PDE forecasting tasks demonstrated that our model performs favorably when compared to existing models. While we focus on the flow matching approach, we leave comprehensive comparisons with more computationally demanding approaches, such as score matching (Song et al., 2020b) and bridge matching (De Bortoli et al., 2021), to future work.

Our findings position flow matching as a promising alternative to diffusion-based generative models in PDE forecasting. While diffusion models have shown strong performance, they are typically not simulation-free and require many iterative sampling steps. They rely on simulating reverse SDEs or ODEs, often using 25-100+ steps for generation. Therefore, they can be computationally costly and difficult to tune for high-dimensional spatio-temporal systems. In contrast, flow matching enables training without simulating stochastic processes, by directly learning a continuous vector field aligned with a designed probability path and few-step sampling, offering improved inference efficiency. Our results show that careful design of the probability path, especially those tailored to continuous-time dynamics, can significantly improve training stability and forecast performance. These insights suggest that flow-based methods not only offer theoretical elegance via connections to optimal transport and Schrödinger bridges, but also hold practical advantages for scalable and controllable forecasting in complex dynamical systems.

## Acknowledgments

The computations were enabled by resources provided by the National Academic Infrastructure for Supercomputing in Sweden (NAISS), partially funded by the Swedish Research Council through grant agreement no. 2022-06725 (NAISS 2024/5-269).

SHL would like to acknowledge support from the Wallenberg Initiative on Networks and Quantum Information (WINQ) and the Swedish Research Council (VR/2021-03648). AY was supported by the SciAI Center, funded by the Office of Naval Research under Grant Number N00014-23-1-2729. EH was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Department of Energy Computational Science Graduate Fellowship under Award Number DE-SC0024386. NBE would like to acknowledge LBL’s LDRD initiative for providing partial support. XSL was supported in part by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research’s Applied Mathematics SciML program under Contract No. DE-AC02-05CH11231 at Lawrence Berkeley National Laboratory.

## 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.
- Juan Miguel Lopez Alcaraz and Nils Strodthoff. Diffusion-based time series imputation and forecasting with structured state space models, 2023.
- Anonymous. Diffusion models and Gaussian flow matching: Two sides of the same coin. In *ICLR Blogposts 2025*, 2025.
- Omri Azencot, N Benjamin Erichson, Vanessa Lin, and Michael W Mahoney. Forecasting sequential data using consistent Koopman autoencoders. In *International Conference on Machine Learning*, pp. 475–485. PMLR, 2020.
- Peter Battaglia, Razvan Pascanu, Matthew Lai, Danilo Jimenez Rezende, et al. Interaction networks for learning about objects, relations and physics. *Advances in Neural Information Processing Systems*, 29, 2016.
- Yoshua Bengio, Nicolas Boulanger-Lewandowski, and Razvan Pascanu. Advances in optimizing recurrent networks. In *2013 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 8624–8628. IEEE, 2013.
- Marin Biloš, Kashif Rasul, Anderson Schneider, Yuriy Nevmyvaka, and Stephan Günnemann. Modeling temporal data as continuous functions with stochastic process diffusion, 2023.
- Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendeleevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023.
- George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time Series Analysis: Forecasting and Control*. John Wiley & Sons, 2015.
- Charlotte Bunne, Ya-Ping Hsieh, Marco Cuturi, and Andreas Krause. The Schrödinger bridge between Gaussian measures has a closed form. In *International Conference on Artificial Intelligence and Statistics*, pp. 5802–5833. PMLR, 2023.
- Salva Rühling Cachay, Bo Zhao, Hailey Joren, and Rose Yu. Dyffusion: A dynamics-informed diffusion model for spatiotemporal forecasting, 2023.
- Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. *Advances in Neural Information Processing Systems*, 31, 2018.

- Tianrong Chen, Guan-Horng Liu, and Evangelos A Theodorou. Likelihood training of Schrödinger bridge using forward-backward sdes theory. *arXiv preprint arXiv:2110.11291*, 2021a.
- Yifan Chen, Mark Goldstein, Mengjian Hua, Michael S Albergo, Nicholas M Boffi, and Eric Vandenberg. Probabilistic forecasting with stochastic interpolants and Föllmer processes. *arXiv preprint arXiv:2403.13724*, 2024.
- Yongxin Chen, Tryphon T Georgiou, and Michele Pavon. Stochastic control liaisons: Richard Sinkhorn meets Gaspard Monge on a Schrodinger bridge. *SIAM Review*, 63(2):249–313, 2021b.
- Sinho Chewi, Jonathan Niles-Weed, and Philippe Rigollet. Statistical optimal transport. *arXiv preprint arXiv:2407.18163*, 2024.
- Tim Colonius and Kunihiko Taira. A fast immersed boundary method using a nullspace approach and multi-domain far-field boundary conditions. *Computer Methods in Applied Mechanics and Engineering*, 197(25-28):2131–2146, 2008.
- Quan Dao, Hao Phung, Binh Nguyen, and Anh Tran. Flow matching in latent space. *arXiv preprint arXiv:2307.08698*, 2023.
- Aram Davtyan, Sepehr Sameni, and Paolo Favaro. Efficient video prediction via sparsely conditioned flow matching. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 23263–23274, 2023.
- Valentin De Bortoli, James Thornton, Jeremy Heng, and Arnaud Doucet. Diffusion Schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709, 2021.
- Tim Dockhorn, Arash Vahdat, and Karsten Kreis. Score-based generative modeling with critically-damped Langevin diffusion. *arXiv preprint arXiv:2112.07068*, 2021.
- Peter D Dueben and Peter Bauer. Challenges and design choices for global weather and climate models based on machine learning. *Geoscientific Model Development*, 11(10):3999–4009, 2018.
- N Benjamin Erichson, Lionel Mathelin, J Nathan Kutz, and Steven L Brunton. Randomized dynamic mode decomposition. *SIAM Journal on Applied Dynamical Systems*, 18(4):1867–1891, 2019.
- N Benjamin Erichson, Lionel Mathelin, Zhewei Yao, Steven L Brunton, Michael W Mahoney, and J Nathan Kutz. Shallow neural networks for fluid flow reconstruction with limited sensors. *Proceedings of the Royal Society A*, 476(2238):20200097, 2020.
- Patrick Esser, Sumith Kulal, Andreas Blattmann, Rahim Entezari, Jonas Müller, Harry Saini, Yam Levi, Dominik Lorenz, Axel Sauer, Frederic Boesel, et al. Scaling rectified flow transformers for high-resolution image synthesis. In *Forty-first International Conference on Machine Learning*, 2024.
- Lawrence C Evans. *Partial Differential Equations*, volume 19. American Mathematical Society, 2022.
- Dario Gasbarra, Tommi Sottinen, and Esko Valkeila. Gaussian bridges. In *Stochastic Analysis and Applications: The Abel Symposium 2005*, pp. 361–382. Springer, 2007.
- Itai Gat, Tal Remez, Neta Shaul, Felix Kreuk, Ricky TQ Chen, Gabriel Synnaeve, Yossi Adi, and Yaron Lipman. Discrete flow matching. *arXiv preprint arXiv:2407.15595*, 2024.
- Georg A Gottwald and Sebastian Reich. Localized Schrödinger bridge sampler. *arXiv preprint arXiv:2409.07968*, 2024.
- Agrim Gupta, Lijun Yu, Kihyuk Sohn, Xiuye Gu, Meera Hahn, Li Fei-Fei, Irfan Essa, Lu Jiang, and José Lezama. Photorealistic video generation with diffusion models. *arXiv preprint arXiv:2312.06662*, 2023.
- Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.

- Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022.
- Yang Hu, Xiao Wang, Lirong Wu, Huatian Zhang, Stan Z Li, Sheng Wang, and Tianlong Chen. FM-TS: Flow matching for time series generation. *arXiv preprint arXiv:2411.07506*, 2024.
- Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models. *Advances in Neural Information Processing Systems*, 35:26565–26577, 2022.
- Diederik P Kingma and Ruiqi Gao. Understanding the diffusion objective as a weighted integral of ELBOs. *arXiv preprint arXiv:2303.00848*, 2023.
- Georg Kohl, Liwei Chen, and Nils Thuerey. Benchmarking autoregressive conditional diffusion models for turbulent flow simulation. In *ICML 2024 AI for Science Workshop*, 2024.
- Aditi S Krishnapriyan, Alejandro F Queiruga, N Benjamin Erichson, and Michael W Mahoney. Learning continuous models for continuous physics. *Communications Physics*, 6(1):319, 2023.
- Christian Léonard. A survey of the Schrödinger problem and some of its connections with optimal transport. *arXiv preprint arXiv:1308.0215*, 2013.
- Yaron Lipman, Ricky TQ Chen, Heli Ben-Hamu, Maximilian Nickel, and Matt Le. Flow matching for generative modeling. *arXiv preprint arXiv:2210.02747*, 2022.
- Yaron Lipman, Marton Havasi, Peter Holderrieth, Neta Shaul, Matt Le, Brian Karrer, Ricky TQ Chen, David Lopez-Paz, Heli Ben-Hamu, and Itai Gat. Flow matching guide and code. *arXiv preprint arXiv:2412.06264*, 2024.
- Guan-Horng Liu, Yaron Lipman, Maximilian Nickel, Brian Karrer, Evangelos A Theodorou, and Ricky TQ Chen. Generalized Schrödinger bridge matching. *arXiv preprint arXiv:2310.02233*, 2023.
- 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.
- Kai Lv, Liang Yuan, and Xiaoyu Ni. Learning autoencoder diffusion models of pedestrian group relationships for multimodal trajectory prediction. *IEEE Transactions on Instrumentation and Measurement*, 73:1–12, 2024. doi: 10.1109/TIM.2024.3375973.
- Dongwei Lyu, Rie Nakata, Pu Ren, Michael W Mahoney, Arben Pitarka, Nori Nakata, and N Benjamin Erichson. WaveCastNet: An AI-enabled wavefield forecasting framework for earthquake early warning. *arXiv preprint arXiv:2405.20516*, 2024.
- Caspar Meijer and Lydia Y Chen. The rise of diffusion models in time-series forecasting. *arXiv preprint arXiv:2401.03006*, 2024.
- Jaideep Pathak, Yair Cohen, Piyush Garg, Peter Harrington, Noah Brenowitz, Dale Durran, Morteza Mardani, Arash Vahdat, Shaoming Xu, Karthik Kashinath, et al. Kilometer-scale convection allowing model emulation using generative diffusion modeling. *arXiv preprint arXiv:2408.10958*, 2024.
- Stefano Peluchetti. Diffusion bridge mixture transports, Schrödinger bridge problems and generative modeling. *Journal of Machine Learning Research*, 24(374):1–51, 2023.
- Aram-Alexandre Pooladian and Jonathan Niles-Weed. Plug-in estimation of Schrödinger bridges. *arXiv preprint arXiv:2408.11686*, 2024.
- Aram-Alexandre Pooladian, Heli Ben-Hamu, Carles Domingo-Enrich, Brandon Amos, Yaron Lipman, and Ricky TQ Chen. Multisample flow matching: Straightening flows with minibatch couplings. *arXiv preprint arXiv:2304.14772*, 2023.



- Hao Qu, Yongshun Gong, Meng Chen, Junbo Zhang, Yu Zheng, and Yilong Yin. Forecasting fine-grained urban flows via spatio-temporal contrastive self-supervision. *IEEE Transactions on Knowledge and Data Engineering*, 35(8):8008–8023, 2023. doi: 10.1109/TKDE.2022.3200734.
- Alejandro Queiruga, N Benjamin Erichson, Liam Hodgkinson, and Michael W Mahoney. Stateful ODE-nets using basis function expansions. *Advances in Neural Information Processing Systems*, 34:21770–21781, 2021.
- Kashif Rasul, Calvin Seward, Ingmar Schuster, and Roland Vollgraf. Autoregressive denoising diffusion models for multivariate probabilistic time series forecasting. In *International Conference on Machine Learning*, pp. 8857–8868. PMLR, 2021.
- Pu Ren, Rie Nakata, Maxime Lacour, Ilan Naiman, Nori Nakata, Jialin Song, Zhengfa Bi, Osman Asif Malik, Dmitriy Morozov, Omri Azencot, et al. Learning physics for unveiling hidden earthquake ground motions via conditional generative modeling. *arXiv preprint arXiv:2407.15089*, 2024.
- Simo Särkkä and Arno Solin. *Applied Stochastic Differential Equations*, volume 10. Cambridge University Press, 2019.
- Martin G Schultz, Clara Betancourt, Bing Gong, Felix Kleinert, Michael Langguth, Lukas Hubert Leufen, Amirpasha Mozaffari, and Scarlet Stadler. Can deep learning beat numerical weather prediction? *Philosophical Transactions of the Royal Society A*, 379(2194):20200097, 2021.
- Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. *Advances in Neural Information Processing Systems*, 28, 2015.
- Yuyang Shi, Valentin De Bortoli, Andrew Campbell, and Arnaud Doucet. Diffusion Schrödinger bridge matching. *Advances in Neural Information Processing Systems*, 36, 2024.
- Gaurav Shrivastava and Abhinav Shrivastava. Video prediction by modeling videos as continuous multi-dimensional processes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7236–7245, 2024.
- Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations*, 2023.
- Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020a.
- Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *Advances in Neural Information Processing Systems*, 32, 2019.
- Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020b.
- Yang Song, Conor Durkan, Iain Murray, and Stefano Ermon. Maximum likelihood training of score-based diffusion models. *Advances in Neural Information Processing Systems*, 34:1415–1428, 2021.
- Makoto Takamoto, Timothy Praditia, Raphael Leiteritz, Daniel MacKinlay, Francesco Alesiani, Dirk Pflüger, and Mathias Niepert. PDEBench: An extensive benchmark for scientific machine learning. *Advances in Neural Information Processing Systems*, 35:1596–1611, 2022.
- Ella Tamir, Najwa Laabid, Markus Heinonen, Vikas Garg, and Arno Solin. Conditional flow matching for time series modelling. In *ICML 2024 Workshop on Structured Probabilistic Inference & Generative Modeling*, 2024.

- Yusuke Tashiro, Jiaming Song, Yang Song, and Stefano Ermon. CSDI: Conditional score-based diffusion models for probabilistic time series imputation. 2021.
- Jakub M Tomczak. Latent variable models. In *Deep Generative Modeling*, pp. 57–127. Springer, 2021.
- Alexander Tong, Nikolay Malkin, Kilian Fatras, Lazar Atanackovic, Yanlei Zhang, Guillaume Huguët, Guy Wolf, and Yoshua Bengio. Simulation-free Schrödinger bridges via score and flow matching. *arXiv preprint arXiv:2307.03672*, 2023a.
- Alexander Tong, Nikolay Malkin, Guillaume Huguët, Yanlei Zhang, Jarrid Rector-Brooks, Kilian Fatras, Guy Wolf, and Yoshua Bengio. Improving and generalizing flow-based generative models with minibatch optimal transport. *arXiv preprint arXiv:2302.00482*, 2023b.
- Arash Vahdat, Karsten Kreis, and Jan Kautz. Score-based generative modeling in latent space. *Advances in Neural Information Processing Systems*, 34:11287–11302, 2021.
- Gefei Wang, Yuling Jiao, Qian Xu, Yang Wang, and Can Yang. Deep generative learning via Schrödinger bridge. In *International Conference on Machine Learning*, pp. 10794–10804. PMLR, 2021.
- Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004.
- Haomin Wen, Youfang Lin, Yutong Xia, Huaiyu Wan, Qingsong Wen, Roger Zimmermann, and Yuxuan Liang. DiffSTG: Probabilistic spatio-temporal graph forecasting with denoising diffusion models, 2024.
- Tijin Yan, Hongwei Zhang, Tong Zhou, Yufeng Zhan, and Yuanqing Xia. ScoreGrad: Multivariate probabilistic time series forecasting with continuous energy-based generative models, 2021.
- Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. *arXiv preprint arXiv:2204.13902*, 2022.

## Appendix

This appendix is organized as follows. In App. A, we provide a detailed discussion of related work and connect them to our studies. In App. B, we provide some remarks on the connection of flow matching models to other generative models. In App. C, we provide commonly used loss parametrizations and compare them with our flow matching loss. In App. D, we provide theorems and proofs to justify the discussions in Section 2-3 (see Theorem D.1-D.2 and Theorem 3.2-3.3), as well as additional discussions. In App. E, we provide additional empirical results. In App. F, we provide the missing experimental details.

## A Related Work

**Modeling and forecasting dynamical systems.** Traditionally dynamical systems arising in scientific applications have been studied through mathematical models based on physical laws, such as ODEs and PDEs (Evans, 2022). These methods often require significant domain knowledge and strong assumptions, limiting their flexibility in data-driven scenarios where governing equations are unknown. Statistical and machine learning (ML) methods have emerged as powerful alternatives to study these systems. Classical approaches include autoregressive integrated moving average (ARMA) and state-space models such as Kalman filters (Box et al., 2015), which have long been used for time series forecasting but struggle with highly nonlinear dynamics. ML methods such as recurrent neural networks (RNNs) (Bengio et al., 2013) and variants for spatio-temporal data (Shi et al., 2015; Lyu et al., 2024) are capable of learning complex temporal dependencies but they are challenging to train. Neural ODEs (Chen et al., 2018; Queiruga et al., 2021) offer a structured approach to continuous-time modeling by integrating neural networks with ODE solvers. However, these models still face challenges with generalizing to new systems, and learning continuous physical dynamics (Krishnapriyan et al., 2023).

**Generative modeling for time series forecasting.** Generative models, particularly score-based diffusion models and flow-based approaches, have gained significant attention for learning complex data distributions in high-dimensional spaces such as images (Song et al., 2020b; Karras et al., 2022; Esser et al., 2024), videos (Davtyan et al., 2023; Shrivastava & Shrivastava, 2024), and dynamical systems (Pathak et al., 2024; Kohl et al., 2024; Ren et al., 2024). In the context of time series forecasting, diffusion models like TimeGrad (Rasul et al., 2021; Meijer & Chen, 2024) incorporate conditioning on previous time steps into both the forecasting process and the loss function. Building upon TimeGrad, CSDI (Tashiro et al., 2021) enhances performance in imputation tasks by replacing the RNN encoder with a transformer. ScoreGrad (Yan et al., 2021) adapts this framework to a score-based diffusion model for multivariate probabilistic forecasting. Methods such as DSPD and CSPD (Biloš et al., 2023), based on DDPM and SDE respectively, model dynamics as continuous functions and introduce time-correlated noise functions. Another relevant work, SSSD<sup>S4</sup> (Alcaraz & Strodtzoff, 2023), uses state-space models (S4) to encode time series and performs diffusion over the temporal domain instead of across multivariate components. While these models demonstrate strong performance, they often require many sampling steps during generation. Related to our approach, (Chen et al., 2024) proposes an SDE framework utilizing stochastic interpolants (Albergo et al., 2023) for probabilistic forecasting, but their sampler also requires many steps during generation. Another recent work is (Tamir et al., 2024), which introduces a flow matching framework for time series modeling within the data space; however, it concentrates on small ODE datasets and does not address forecasting tasks. Meanwhile, (Hu et al., 2024) proposes a model based on rectified flow with ODE-based straight line transport for efficient time series generation.

Applications of diffusion models to spatio-temporal data have primarily focused on video generation (Singer et al., 2023). For spatio-temporal forecasting on sensor grids, graph neural networks (GNNs) are widely employed in tasks such as traffic prediction (Qu et al., 2023; Lv et al., 2024; Wen et al., 2024) and air quality forecasting (Wen et al., 2024). Another recent work is DYffusion (Cachay et al., 2023), which utilizes a non-Markovian diffusion process to propagate images temporally, similar to DDIM. This method iteratively refines initial predictions at larger time steps by conditioning them on predictions at smaller time steps.

## B Connection to SDE Based Generative Model Frameworks

In this section, we provide some remarks on the connection between flow matching models and SDE based generative models (Song et al., 2020b).

Consider the following continuous-time Gaussian latent variable model (Tomczak, 2021):  $Z_t = \mathcal{E}(X_t)$ ,

$$Z_t = A_t Z_0 + B_t Z_1 + L_t \epsilon, \quad t \in [0, 1], \quad (23)$$

$$X_t = \mathcal{D}(Z_t), \quad (24)$$

where  $t$  is the continuous variable,  $X_0 \in \mathbb{R}^d$  represent data samples,  $Z \in \mathbb{R}^p$  is the latent variable,  $\epsilon \sim \mathcal{N}(0, I)$  is independent of the random variables  $X_0, Z_0, Z_1$ . Here  $A_t, B_t$  and  $L_t \geq 0$  are pre-specified coefficients which are possibly matrix-valued and time-dependent,  $\mathcal{D}$  and  $\mathcal{E}$  denote the decoder and encoder map respectively, and  $\mathcal{D} \circ \mathcal{E} = I$ . Note that  $Z_0$  and  $Z_1$  are initial and terminal point of the path  $(Z_t)_{t \in [0,1]}$  in the latent space.

The above latent variable model can be identified (up to equivalence in law for each  $t$ ) with the linear SDE of the form:

$$d\hat{Z}_t = F_t \hat{Z}_t dt + H_t Z_1 dt + G_t dW_t, \quad \hat{Z}_0 = Z_0, \quad t \in [0, 1], \quad (25)$$

where  $(W_t)_{t \in [0,1]}$  is the standard Wiener process. By matching the moments, we obtain

$$F_t = \dot{A}_t A_t^{-1}, \quad (26)$$

$$H_t = \dot{B}_t - \dot{A}_t A_t^{-1} B_t, \quad (27)$$

$$G_t G_t^T = \dot{L}_t L_t^T + L_t \dot{L}_t^T - \dot{A}_t A_t^{-1} L_t L_t^T - L_t L_t^T A_t^{-T} (\dot{A}_t)^T, \quad (28)$$

where the overdot denotes derivative with respect to  $t$  and  $A^T$  denotes the transpose of  $A$ .

Under the above formulation, various existing generative models such as DDPM (Ho et al., 2020), VP-SDE and VE-SDE of (Song et al., 2020b;a), the critically damped SDE of (Dockhorn et al., 2021), the flow matching models in (Lipman et al., 2022; Tong et al., 2023b; Liu et al., 2022) and the stochastic interpolants of (Albergo et al., 2023) can be recovered, and new models can be derived.

The following proposition establishes the connection between flow matching using our proposed probability path model, the Gaussian latent variable model (23) and the linear SDE model (25).

**Proposition B.1.** *For every  $t \in [0, 1]$ , the  $Z_t$  defined in Eq. (19) can be identified, up to equivalence in law, with the  $Z_t$  generated by the latent variable model (23) with  $A_t = (1-t)I$ ,  $B_t = tI$ ,  $L_t = \sqrt{\sigma_{min}^2 + \sigma^2 t(1-t)}I$ . For  $t \in [0, 1)$ , it can also be identified with the solution  $\hat{Z}_t$  of the linear SDE (25) with  $F_t = -I/(1-t)$ ,  $H_t = (1 + \frac{t}{1-t})I$  and  $G_t = \sqrt{\sigma^2 + \frac{2\sigma_{min}^2}{1-t}}I$ . Moreover,  $\lim_{t \rightarrow 1} \hat{Z}_t =^d Z_1 + \sigma_{min}\epsilon$ , where  $\epsilon \sim \mathcal{N}(0, I)$  and  $=^d$  denotes equivalence in distribution.*

*Proof.* The identification follows from matching the moments of  $Z_t$  and  $\hat{Z}_t$ , i.e., applying Eq. (26)-(28).

To prove that  $\lim_{t \rightarrow 1} \hat{Z}_t =^d Z_1 + \sigma_{min}\epsilon$ , we use the explicit solution of the SDE:

$$\hat{Z}_t = \Phi(t, 0)Z_0 + \int_0^t \Phi(t, s)H_s Z_1 ds + \int_0^t \Phi(t, s)G_s dW_s,$$

where  $\Phi(t, s)$  is the fundamental solution of the homogeneous equation  $d\Phi(t, s) = F_t \Phi(t, s)dt$  with  $\Phi(s, s) = I$ . For our  $F_t = -I/(1-t)$ , we have  $\Phi(t, s) = \exp(-\int_s^t \frac{1}{1-u} du)I = (\frac{1-t}{1-s})I$ . Substituting this and the formula for  $H_t$  into the solution, we obtain  $\hat{Z}_t = (1-t)Z_0 + \int_0^t (\frac{1-t}{1-s})(1 + \frac{s}{1-s})Z_1 ds + \int_0^t (\frac{1-t}{1-s})\sqrt{\sigma^2 + \frac{2\sigma_{min}^2}{1-s}}I dW_s$ .

Now, let us examine each term as  $t \rightarrow 1$ . First,  $(1-t)Z_0 \rightarrow 0$  as  $t \rightarrow 1$  and  $\int_0^t (\frac{1-t}{1-s})(1 + \frac{s}{1-s})Z_1 ds = tZ_1 \rightarrow Z_1$  as  $t \rightarrow 1$ . It remains to deal with the stochastic integral term  $M_t := \int_0^t (\frac{1-t}{1-s})\sqrt{\sigma^2 + \frac{2\sigma_{min}^2}{1-s}}I dW_s$ . Note that

$M_t$  is an Itô integral that has zero mean, i.e.  $\mathbb{E}M_t = 0$ , and using Itô's formula (Särkkä & Solin, 2019),

$$\mathbb{E}M_t^2 = \int_0^t ((1-t)/(1-s))^2 \left( \sigma^2 + \frac{2\sigma_{min}^2}{1-s} \right) ds \quad (29)$$

$$= (1-t)^2 \left[ \sigma^2 \int_0^t \frac{1}{(1-s)^2} ds + 2\sigma_{min}^2 \int_0^t \frac{1}{(1-s)^3} ds \right] \quad (30)$$

$$= (1-t)^2 \left[ \sigma^2 \left( \frac{t}{1-t} \right) + \sigma_{min}^2 \left( \frac{1}{(1-t)^2} - 1 \right) \right] \quad (31)$$

$$= (1-t)t\sigma^2 + \sigma_{min}^2 t(2-t), \quad (32)$$

which tends to  $\sigma_{min}^2$  as  $t \rightarrow 1$ . Combining the above results,  $\lim_{t \rightarrow 1} \hat{Z}_t \sim \mathcal{N}(Z_1, \sigma_{min}^2 I)$ . □

## C On Different Loss Parametrizations

In this section, we list popular choices of loss parametrization considered in the literature and connect them to our flow matching loss. We refer to (Kingma & Gao, 2023) for a more comprehensive discussion. Recall that the Gaussian path that we consider is:  $Z_t = a_t Z_0 + b_t Z_1 + c_t \xi$ , where  $\xi \sim \mathcal{N}(0, I)$ . In general, these loss parametrizations take the form of:

$$\mathcal{L}(\theta) := \mathbb{E}_{t, p_t(Z|\tilde{Z}), q(\tilde{Z})} \omega(t) \|m_t^\theta(Z) - m_t(Z, \tilde{Z})\|^2, \quad (33)$$

where  $\omega(t) > 0$  is a weighting function,  $m_t(Z, \tilde{Z})$  is the object (conditioned on  $\tilde{Z}$ ) to be learnt and  $m_t^\theta$  is a neural network model used to learn the object of interest. Depending on which object one would like to learn/match, we have different loss parametrizations.

**Flow matching.** The flow matching loss that we focus in this paper is:

$$\mathcal{L}_{cfm}(\theta) := \mathbb{E}_{t, p_t(Z|\tilde{Z}), q(\tilde{Z})} \omega(t) \|v_t^\theta(Z) - u_t(Z|\tilde{Z})\|^2, \quad (34)$$

where one aims to learn the flow generating vector field:

$$u_t(Z|\tilde{Z}) = \frac{c_t'}{c_t} (Z - (a_t Z_0 + b_t Z_1)) + a_t' Z_0 + b_t' Z_1, \quad (35)$$

**Score matching.** The score matching loss is:

$$\mathcal{L}_{sm}(\theta) := \mathbb{E}_{t, p_t(Z|\tilde{Z}), q(\tilde{Z})} \lambda(t) \|s_t^\theta(Z) - \nabla \log p_t(Z|\tilde{Z})\|^2, \quad (36)$$

where  $\lambda(t) > 0$  is a weighting function and one aims to learn the score function:

$$\nabla \log p_t(Z|\tilde{Z}) = \frac{a_t Z_0 + b_t Z_1 - Z}{c_t^2}. \quad (37)$$

If  $\lambda(t) = c_t^2$ , then this reduces to the original score matching loss (Song & Ermon, 2019), whereas if  $\lambda(t) = \beta(1-t)$ , this becomes the score flow loss (Song et al., 2021).

**Noise matching.** The noise matching loss is:

$$\mathcal{L}_{nm}(\theta) := \mathbb{E}_{t, p_t(Z|\tilde{Z}), q(\tilde{Z})} \|\epsilon_t^\theta(Z) - \epsilon_t(Z|\tilde{Z})\|^2, \quad (38)$$

where one aims to learn the noise:

$$\epsilon_t(Z|\tilde{Z}) = \frac{Z - (a_t Z_0 + b_t Z_1)}{c_t}. \quad (39)$$

## D Theoretical Results and Proofs

In this section, we provide theorems and proofs to justify the discussions in Section 2 and Section 3.

### D.1 Connecting Flow Matching with Conditional Flow Matching

The following theorem justifies the claim that minimizing  $\mathcal{L}_{fm}$  is equivalent to minimizing  $\mathcal{L}_{cfm}$ .

**Theorem D.1.** *If the conditional vector field  $u_t(Z|\tilde{Z})$  generates the conditional probability path  $p_t(Z|\tilde{Z})$ , then the marginal vector field  $u_t$  in Eq. (5) generates the marginal probability path  $p_t$  in Eq. (4). Moreover, if  $p_t(Z) > 0$  for all  $t, Z$ , then  $\mathcal{L}_{fm}$  and  $\mathcal{L}_{cfm}$  are equal up to a constant independent of  $\theta$ .*

*Proof.* The proof is a straightforward extension of the proofs of Theorem 1-2 in (Lipman et al., 2022) from conditioning on data samples to conditioning on latent samples and allowing an arbitrary weighting function  $\omega(t)$ .

Suppose that the conditional vector field  $u_t(Z|\tilde{Z})$  generates the conditional probability path  $p_t(Z|\tilde{Z})$ , we would like to show that the marginal vector field  $u_t$  in Eq. (5) generates the marginal probability path  $p_t$  in Eq. (4). To show this, it suffices to verify that  $p_t$  and  $u_t$  satisfy the continuity equation:

$$\frac{d}{dt}p_t(Z) + \text{div}(p_t(Z)u_t(Z)) = 0, \quad (40)$$

where the divergence operator,  $\text{div}$ , is defined with respect to the latent variable  $Z = (Z^1, \dots, Z^d)$ , i.e.,  $\text{div} = \sum_{i=1}^d \frac{\partial}{\partial Z^i}$ .

We begin by taking derivative of  $p_t(Z)$  with respect to time:

$$\frac{d}{dt}p_t(Z) = \frac{d}{dt} \int p_t(Z|\tilde{Z})q(\tilde{Z})d\tilde{Z} \quad (41)$$

$$= \int \frac{d}{dt}p_t(Z|\tilde{Z})q(\tilde{Z})d\tilde{Z} \quad (42)$$

$$= - \int \text{div}(u_t(Z|\tilde{Z})p_t(Z|\tilde{Z}))q(\tilde{Z})d\tilde{Z} \quad (43)$$

$$= -\text{div} \left( \int u_t(Z|\tilde{Z})p_t(Z|\tilde{Z})q(\tilde{Z})d\tilde{Z} \right) \quad (44)$$

$$= -\text{div}(u_t(Z)p_t(Z)). \quad (45)$$

In the third line, we use the fact that  $u_t(\cdot|\tilde{Z})$  generates  $p_t(\cdot|\tilde{Z})$ . In the last line, we use Eq. (5). In the second and forth line above, the exchange of integration and differentiation can be justified by assuming that the integrands satisfy the regularity conditions of the Leibniz rule.

Next, we would like to show that if  $p_t(Z) > 0$  for all  $t, Z$ , then  $\mathcal{L}_{fm}$  and  $\mathcal{L}_{cfm}$  are equal up to a constant independent of  $\theta$ . We follow (Lipman et al., 2022) and assume that  $q(Z)$  and  $p_t(Z|\tilde{Z})$  are decreasing to zero sufficiently fast as  $\|Z\| \rightarrow 0$ , that  $u_t, v_t, \nabla_{\theta}v_t$  are bounded, so that all integrals exist and exchange of integration order is justified via Fubini's theorem.

Using the bilinearity of the 2-norm, we decompose the squared losses as:

$$\|v_t^{\theta}(Z) - u_t(Z)\|^2 = \|v_t^{\theta}(Z)\|^2 - 2\langle v_t^{\theta}(Z), u_t(Z) \rangle + \|u_t(Z)\|^2, \quad (46)$$

$$\|v_t^{\theta}(Z) - u_t(Z|\tilde{Z})\|^2 = \|v_t^{\theta}(Z)\|^2 - 2\langle v_t^{\theta}(Z), u_t(Z|\tilde{Z}) \rangle + \|u_t(Z|\tilde{Z})\|^2. \quad (47)$$

Now,

$$\mathbb{E}_{p_t(Z)} \|v_t^\theta(Z)\|^2 = \int \|v_t^\theta(Z)\|^2 p_t(Z) dZ \quad (48)$$

$$= \int \int \|v_t^\theta(Z)\|^2 p_t(Z|\tilde{Z}) q(\tilde{Z}) d\tilde{Z} dZ \quad (49)$$

$$= \mathbb{E}_{q(\tilde{Z}), p_t(Z|\tilde{Z})} \|v_t^\theta(Z)\|^2, \quad (50)$$

where we use Eq. (4) in the second equality above and exchange the order of integration in the third equality.

Next, we compute:

$$\mathbb{E}_{p_t(Z)} \langle v_t^\theta(Z), u_t(Z) \rangle = \int \left\langle v_t^\theta(Z), \frac{\int u_t(Z|\tilde{Z}) p_t(Z|\tilde{Z}) q(\tilde{Z}) d\tilde{Z}}{p_t(Z)} \right\rangle p_t(Z) dZ \quad (51)$$

$$= \int \left\langle v_t^\theta(Z), \int u_t(Z|\tilde{Z}) p_t(Z|\tilde{Z}) q(\tilde{Z}) d\tilde{Z} \right\rangle dZ \quad (52)$$

$$= \int \int \langle v_t^\theta(Z), u_t(Z|\tilde{Z}) \rangle p_t(Z|\tilde{Z}) q(\tilde{Z}) d\tilde{Z} dZ \quad (53)$$

$$= \mathbb{E}_{q(\tilde{Z}), p_t(Z|\tilde{Z})} \langle v_t^\theta(Z), u_t(Z|\tilde{Z}) \rangle, \quad (54)$$

where we first plug in Eq. (5) and then exchange the order the integration in order to arrive at the last equality.

Finally, noting that  $u_t$  are  $\omega(t)$  independent of  $\theta$  (and are thus irrelevant for computing the loss gradients), we have proved the desired result.  $\square$

## D.2 Identifying the Vector Field that Generates the Gaussian Paths

Similar to Theorem 3 in (Lipman et al., 2022), we have the following result, which identifies the unique vector field that generates the Gaussian probability path.

**Theorem D.2.** *Let  $p_t(Z|\tilde{Z})$  be the Gaussian probability path defined in Eq. (7) and consider the flow map  $\psi_t$  defined as  $\psi_t(Z) = a_t Z_0 + b_t Z_1 + c_t Z$  with  $c_t > 0$ . Then the unique vector field that defines  $\psi_t$  is:*

$$u_t(Z|\tilde{Z}) = \frac{c'_t}{c_t} (Z - (a_t Z_0 + b_t Z_1)) + a'_t Z_0 + b'_t Z_1, \quad (55)$$

where prime denotes derivative with respect to  $t$ , and  $u_t(Z|\tilde{Z})$  generates the Gaussian path  $p_t(Z|\tilde{Z})$ .

*Proof.* Let  $p_t(Z|\tilde{Z})$  be the Gaussian probability path defined in Eq. (7) and consider the flow map  $\psi_t$  defined as  $\psi_t(Z) = a_t Z_0 + b_t Z_1 + c_t Z$ . We would like to show that the unique vector field that defines  $\psi_t$  is:

$$u_t(Z|\tilde{Z}) = \frac{c'_t}{c_t} (Z - (a_t Z_0 + b_t Z_1)) + a'_t Z_0 + b'_t Z_1, \quad (56)$$

and  $u_t(Z|\tilde{Z})$  generates the Gaussian path  $p_t(Z|\tilde{Z})$ .

We denote  $w_t = u_t(Z|\tilde{Z})$  for notational simplicity. Then,

$$\frac{d}{dt} \psi_t(Z) = w_t(\psi_t(Z)). \quad (57)$$

Since  $\psi_t$  is invertible (as  $c_t > 0$ ), we let  $Z = \psi^{-1}(Y)$  and obtain

$$\psi'_t(\psi^{-1}(Y)) = w_t(Y), \quad (58)$$

where the prime denotes derivative with respect to  $t$  and we have used the apostrophe notation for the derivative to indicate that  $\psi'_t$  is evaluated at  $\psi^{-1}(Y)$ .

Inverting  $\psi_t(Z)$  gives:

$$\psi_t^{-1}(Y) = \frac{Y - \mu_t(\tilde{Z})}{c_t}, \quad (59)$$

where  $\mu_t(\tilde{Z}) := a_t Z_0 + b_t Z_1$ .

Differentiating  $\psi_t$  with respect to  $t$  gives  $\psi'_t(Z) = c'_t Z + \mu'_t(\tilde{Z})$ .

Plugging the last two equations into Eq. (58), we obtain:

$$w_t(Y) = \frac{c'_t}{c_t}(Y - \mu_t(\tilde{Z})) + \mu'_t(\tilde{Z}) \quad (60)$$

which is the result that we wanted to show. □

### D.3 Solution to the Gaussian Schrodinger Bridge Problem

We prove Theorem 3.2 in this subsection. First, we recall Theorem 3 from (Bunne et al., 2023). Theorem 3.2 then follows from an application of the theorem.

Let  $\xi_0 \sim \mathcal{N}(\mu_0, \Sigma_0)$  and  $\xi_1 \sim \mathcal{N}(\mu_1, \Sigma_1)$  be two given Gaussian random variables. From now on, by abusing the notation we will also denote the measures of these Gaussians by  $\xi_0$  and  $\xi_1$  respectively.

Let  $\sigma \geq 0$  and set:

$$D_\sigma := (4\Sigma_0^{1/2}\Sigma_1\Sigma_0^{1/2} + \sigma^4 I)^{1/2} \quad (61)$$

$$C_\sigma := \frac{1}{2}(\Sigma_0^{1/2}D_\sigma\Sigma_0^{-1/2} - \sigma^2 I). \quad (62)$$

Consider the following Schrodinger bridges with Gaussian marginal constraints:

$$\min_{\mathbb{P}_0=\xi_0, \mathbb{P}_1=\xi_1} D_{KL}(\mathbb{P}||\mathbb{Q}), \quad (63)$$

where  $\mathbb{Q}$  is the path measure of the linear SDE:

$$dY_t = (c_t Y_t + \alpha_t)dt + g_t dW_t := f_t dt + g_t dW_t. \quad (64)$$

Here,  $c_t : \mathbb{R}^+ \rightarrow \mathbb{R}$ ,  $\alpha_t : \mathbb{R}^+ \rightarrow \mathbb{R}^d$  and  $g_t : \mathbb{R}^+ \rightarrow \mathbb{R}^+$  are smooth functions, and  $W_t$  is the standard Wiener process.

The following theorem is a simplified version of Theorem 3 in (Bunne et al., 2023).

**Theorem D.3** (Theorem 3 in (Bunne et al., 2023)). *The solution  $\mathbb{P}^*$  to the Gaussian Schrodinger bridges is (the path measure of) a Markov Gaussian process whose marginal variable  $X_t \sim \mathcal{N}(\mu_t, \Sigma_t)$ , where*

$$\mu_t = (\tau_t - r_t \tau_1)\mu_0 + r_t \mu_1 + \zeta(t) - r_t \zeta(1), \quad (65)$$

$$\Sigma_t = (\tau_t - r_t \tau_1)^2 \Sigma_0 + r_t^2 \Sigma_1 + r_t(\tau_t - r_t \tau_1)(C_{\sigma^*} + C_{\sigma^*}^T) + \kappa(t, t)(1 - \rho_t)I. \quad (66)$$

In the above,

$$\begin{aligned} \tau_t &= \exp\left(\int_0^t c_s ds\right), \quad \kappa(t, t') = \tau_t \tau_{t'} \int_0^t \tau_s^{-2} g_s^2 ds, \quad r_t = \kappa(t, 1)/\kappa(1, 1), \\ \sigma^* &= \sqrt{\tau_1^{-1} \kappa(1, 1)}, \quad \zeta(t) = \tau_t \int_0^t \tau_s^{-1} \alpha_s ds, \quad \rho_t = \frac{\int_0^t \tau_s^{-2} g_s^2 ds}{\int_0^1 \tau_s^{-2} g_s^2 ds}. \end{aligned}$$

Now we prove Theorem 3.2.



*Proof of Theorem 3.2.* We shall apply Theorem D.3 to the Gaussian SB (15) by setting  $\mu_0 := Z_0$ ,  $\mu_1 = Z_1$ ,  $\Sigma_0 = \Sigma_1 := \sigma_{min}^2 I$ ,  $\alpha_t := \dot{a}_t Z_0 + \dot{b}_t Z_1$ ,  $c_t := 0$ , and  $g_t := \omega$ . Then, we have  $\tau_t = 1$ ,  $\kappa(t, t') = \omega^2 t$ ,  $r_t = t$ ,  $\sigma^* = \omega$ ,  $\zeta(t) = \int_0^t \alpha_s ds = \left(\int_0^t \dot{a}_s ds\right) Z_0 + \left(\int_0^t \dot{b}_s ds\right) Z_1 = (a_t - a_0)Z_0 + (b_t - b_0)Z_1$ , and  $\rho_t = t$ . Also,  $D_{\sigma^*} = (4\sigma_{min}^4 + \omega^4)^{1/2} I$  and  $C_{\sigma^*} = C_{\sigma^*}^T = \frac{1}{2}((4\sigma_{min}^4 + \omega^4)^{1/2} - \omega^2)I$ .

Thus,

$$\mu_t = (1-t)Z_0 + tZ_1 + (a_t - a_0)Z_0 + (b_t - b_0)Z_1 - t[(a_1 - a_0)Z_0 + (b_1 - b_0)Z_1] \quad (67)$$

and

$$\Sigma_t = [((1-t)^2 + t^2)\sigma_{min}^2 + t(1-t)((4\sigma_{min}^4 + \omega^4)^{1/2} - \omega^2) + \omega^2 t(1-t)]I \quad (68)$$

$$= \sigma_{min}^2 + \left(\sqrt{4\sigma_{min}^4 + \omega^4} - 2\sigma_{min}^2\right) \cdot t(1-t). \quad (69)$$

Theorem 3.2 then follows from Theorem D.3 with the above formula for  $\mu_t$  and  $\Sigma_t$ .  $\square$

#### D.4 Comparing the Variance of Different Vector Field Models

We begin with providing intuitions for understanding our proposed probability path model. We expect that our model improves upon the other considered models, as it takes advantage of the inherent continuity and correlation in the spatio-temporal data. Intuitively, for time series samples whose underlying dynamics are continuous and obey a physical law, a Gaussian sample is typically further from the time series samples, so the distance between a Gaussian sample and a time series sample should generally be larger than the distance between consecutive time series samples (which can be highly correlated). Therefore, using a probability path that connects consecutive time series samples could lead to faster convergence and more stable training, when compared to using a path that simply connects a time series sample to a Gaussian sample, since the resulting flow model in the former case can better capture the true system dynamics with less effort.

Moreover, if the consecutive samples are sufficiently correlated, then the variance of the VF corresponding to our proposed probability path model can be lower than the variance of the VF corresponding to the other choices of probability paths. Precise statements capturing this are the contents of Theorem 3.3, which focuses on comparison between our proposed model and the optimal transport VF model of (Lipman et al., 2022).

Before proving Theorem 3.3, we start with the following lemma.

**Lemma D.4.** *Let  $A, B, C, D$  be random vectors where  $C$  and  $D$  are independent, both  $A$  and  $B$  are independent of  $C$  and  $D$  (but  $A$  and  $B$  could be dependent). If  $Cov(A, B) \geq (Var(C) - Var(D) + Var(B))/2$ , then*

$$Var(A + D) \geq Var(A - B + C) = Var(B - A + C), \quad (70)$$

where  $A \geq B$  means that  $A - B$  is positive semidefinite.

*Proof.* We compute:

$$Var(A + D) = Var(A - B + C + B + D - C) \quad (71)$$

$$= Var(A - B + C) + Var(B + D) + Var(C) + 2Cov(A - B + C, B + D) - 2Cov(A - B + C, C) - 2Cov(B + D, C) \quad (72)$$

$$= Var(A - B + C) + Var(B + D) - Var(C) + 2Cov(A - B + C, B) + 2Cov(A - B + C, D) \quad (73)$$

$$= Var(A - B + C) + Var(B + D) - Var(C) + 2Cov(A, B) - 2Var(B) \quad (74)$$

$$= Var(A - B + C) - Var(C) + Var(D) + 2Cov(A, B) - Var(B), \quad (75)$$

where we have simply rearranged the terms in the first equality, used the formula  $Var(A + B + C) = Var(A) + Var(B) + Var(C) + 2Cov(A, B) + 2Cov(A, C) + 2Cov(B, C)$ , bilinearity of covariance, the facts that  $Cov(A, A) = Var(A)$  and  $Var(cA) = c^2 Var(A)$  for a scalar  $c$ , as well as the assumption that both  $A$  and  $B$  are independent of  $C, D$  in the last four equalities.

Therefore, if  $-Var(C) + Var(D) + 2Cov(A, B) - Var(B) \geq 0$ , then we have  $Var(A + D) \geq Var(A - B + C)$ .  $\square$

Now we prove Theorem 3.3.

*Proof.* Note that  $c'_t = \frac{\sigma^2(1-2t)}{2\sqrt{\sigma_{min}^2 + \sigma^2 t(1-t)}}$ ,  $Var(c'_t \xi) = (c'_t)^2 I$  and  $Var(-\eta) = I$ . Therefore, using these and applying Lemma D.4 with  $A := z^{\tau-1}$ ,  $B := z^\tau$ ,  $C := c'_t \xi$  and  $D := -\eta$ , allow us to establish the claim that  $Var(\tilde{u}_t(Z_t|z^{\tau-1})) \geq Var(u_t(Z_t|z^{\tau-1}, z^\tau))$  if

$$Cov(z^{\tau-1}, z^\tau) \geq \frac{1}{2} \left( \frac{\sigma^4(1-2t)^2 I}{4(\sigma_{min}^2 + \sigma^2 t(1-t))} + Var(z^\tau) - I \right). \quad (76)$$

Since the function  $f(t) := \frac{\sigma^4(1-2t)^2}{4(\sigma_{min}^2 + \sigma^2 t(1-t))}$  is maximized at the endpoints  $t = 0, 1$  with the maximum value of  $\sigma^4/4\sigma_{min}^2$ , the desired result stated in the theorem follows.  $\square$

Lastly, we provide some discussions following Theorem 3.3 (using the notations there).

**Discussions.** Let us denote  $v_t^\theta(Z) := v_t^\theta(z^\tau, z^{\tau-1}, \xi)$  to show the explicit dependence of the vector field neural net on the random samples  $t, z^\tau, z^{\tau-1}$  and  $\xi \sim \mathcal{N}(0, I)$  drawn during each update of gradient descent during training. During each gradient descent update, our model involves computation of

$$\nabla_\theta \mathcal{L}_{cfm}(\theta; t, \xi, z^\tau, z^{\tau-1}) = 2\nabla_\theta v_t^\theta(z^\tau, z^{\tau-1}, \xi)^T \cdot (v_t^\theta(z^\tau, z^{\tau-1}, \xi) - u_t(Z_t|z^\tau, z^{\tau-1})) \quad (77)$$

$$=: 2\nabla_\theta v_t^\theta(z^\tau, z^{\tau-1}, \xi)^T \cdot \Delta_t^\theta(z^{\tau-1}, z^\tau, \xi), \quad (78)$$

with  $t \sim \mathcal{U}[0, 1]$ ,  $\xi \sim \mathcal{N}(0, I)$  and the latent samples  $z^\tau, z^{\tau-1}$  drawn randomly.

Similarly, for the rectified flow model let us denote  $\tilde{v}_t^\theta(Z) := v_t^\theta(z^{\tau-1}, \eta)$  to show the explicit dependence of the vector field neural net on the random samples  $t, z^{\tau-1}$  and  $\eta \sim \mathcal{N}(0, I)$  drawn during each update of gradient descent during training. Each update of gradient descent using the rectified flow model involves computation of

$$\nabla_\theta \tilde{\mathcal{L}}_{cfm}(\theta; t, \eta, z^{\tau-1}) = 2\nabla_\theta \tilde{v}_t^\theta(z^{\tau-1}, \eta)^T \cdot (\tilde{v}_t^\theta(z^{\tau-1}, \eta) - \tilde{u}_t(\tilde{Z}_t|z^{\tau-1})) \quad (79)$$

$$=: 2\nabla_\theta \tilde{v}_t^\theta(z^{\tau-1}, \eta)^T \cdot \tilde{\Delta}_t^\theta(z^{\tau-1}, \eta), \quad (80)$$

with  $t \sim \mathcal{U}[0, 1]$ ,  $\eta \sim \mathcal{N}(0, I)$  and the latent sample  $z^{\tau-1}$  drawn randomly.

Lower gradient variance results in smoother training loss curve and potentially faster convergence, so it is useful to compare the variances of the loss gradient for the two models. However, the variances are highly dependent on  $\nabla_\theta v_t^\theta$ ,  $\nabla_\theta \tilde{v}_t^\theta$  and their covariance with the other random vectors appearing in Eq. (77) and Eq. (79), making such comparison challenging without strong assumptions. Heuristically, the difference in the variances of the loss gradient during each update for the two models is primarily determined by the difference between  $Var(\Delta_t^\theta(z^{\tau-1}, z^\tau, \xi))$  and  $Var(\tilde{\Delta}_t^\theta(z^{\tau-1}, \eta))$  if  $\nabla_\theta v_t^\theta$  and  $\nabla_\theta \tilde{v}_t^\theta$  are relatively stable. In this case, we have  $Var(\Delta_t^\theta(z^{\tau-1}, z^\tau, \xi)) \leq Var(\tilde{\Delta}_t^\theta(z^{\tau-1}, \eta))$  if we suppose the assumptions in Theorem 3.3,  $Var(\tilde{v}_t^\theta) \geq Var(v_t^\theta)$  and  $Cov(v_t^\theta, u_t) \geq Cov(\tilde{v}_t^\theta, \tilde{u}_t)$ .

The implications of Theorem 3.3 together with the heuristics above could partially explain why our probability path model leads to smoother loss curve and faster convergence (see Figure 2) compared to the RIVER method of (Davtyan et al., 2023). On the other hand, the dependence of the lower bound in the theorem on  $\sigma$  and  $\sigma_{min}$  suggests that using values of  $\sigma$  that is relatively large enough might not keep the variance of the vector field low, which could partially explain the phenomenon displayed in Figure 2, where using  $\sigma = 0.1$  and  $\sigma_{min} = 0.001$  leads to large loss fluctuations.

## E Additional Empirical Results

In this section, we provide additional experimental results.

## E.1 Visualization of Flow Patterns and Dynamics

Figure 3 provides visual results of the predicted snapshots by our model for each task. The visualizations highlight our model’s ability to capture complex flow patterns and dynamics.

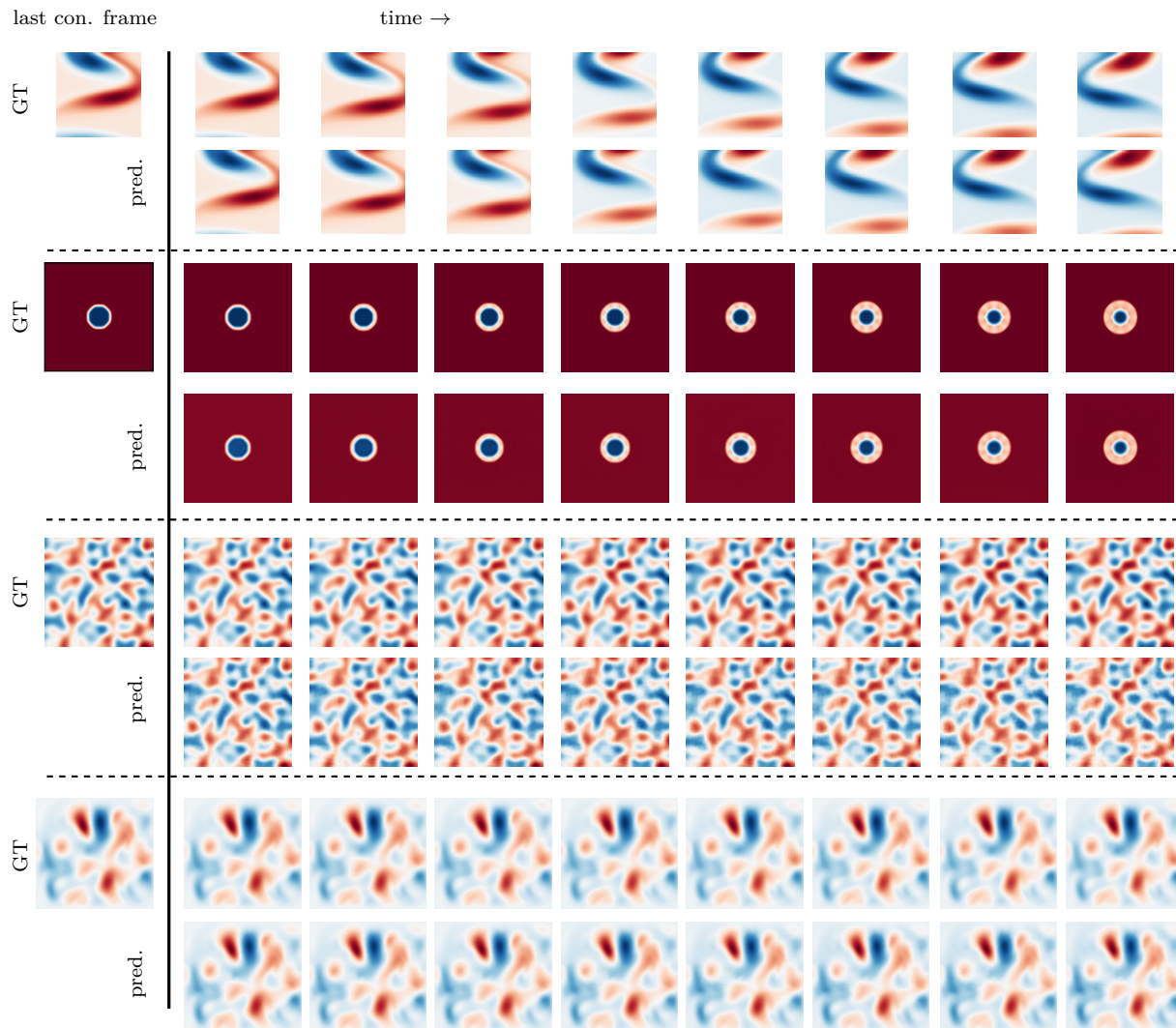


Figure 3: Visualization of predicted frames using our model of probability path for the considered tasks. From top to bottom: fluid flow past a cylinder, shallow-water equation, diffusion-reaction equation, and incompressible Navier-Stokes equation. In each case, GT indicates the ground truth frames and pred. indicates the predicted frames.

## E.2 Ablation Studies

To further assess our model, we conducted a detailed ablation study focusing on the impact of various hyperparameters for the fluid flow (FPC) task. Specifically, we study the impact of the values of  $\sigma$ , the choice of sampler, and the number of sampling steps during inference. For small  $\sigma_{sam}$ , we find that this parameter has negligible impact on test performance, so we fixed  $\sigma_{sam} = 0$  for all experiments in this section.

**Impact of  $\sigma$  on training stability.** Figure 5 illustrates the effect of different  $\sigma$  values on the training loss curve for our method on the fluid flow past a cylinder task. We observed that larger values of  $\sigma$  (e.g.,  $\sigma = 0.1$ ) resulted in smoother loss curves and more stable convergence during training.

**Effect of  $\sigma$ , sampler choice, and sampling steps on test performance.** Table 6 investigates how different values of  $\sigma$ , the choice of sampler (Euler or RK4), and the number of sampling steps affect test performance. It can be seen, that even with as few as one sampling step ( $N = 2$ ) using the Euler scheme, our model perform reasonably well. However, increasing the number of sampling steps or employing the more computationally intensive RK4 sampler can help to lead to better results. From Figure 4, we see that our model leads to the smallest test MSE for both samplers at all sampling steps. In particular, lowest MSEs can be achieved using as few as one sampling step for both sampling schemes. Moreover, while using the RK4 sampler can lead to a lower test MSE for all models, the performance gap is much smaller for our model, showing robustness of our model to sampling choice.

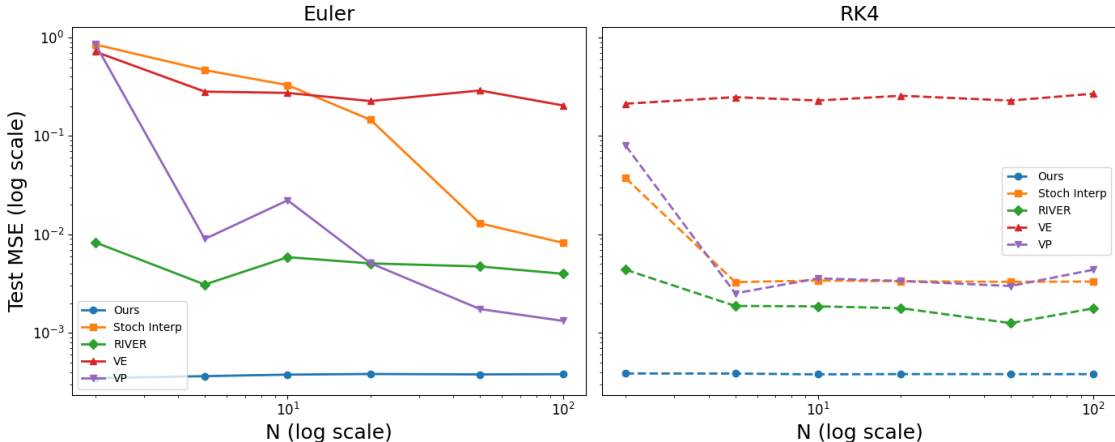


Figure 4: Test MSE vs. number of sampling steps  $N - 1$  (we consider  $N = 2, 5, 10, 20, 50, 100$ ) for the five models with the Euler sampler (left) and the RK4 sampler (right).

**Effect of context frames on test performance.** Table 3 shows the test performance in terms of the considered evaluation metrics in the case when the neural networks are trained without using the context frames, in which case  $C$  is null in (10) and we are minimizing the loss in (9). We see that without using the proposed context frames, the test performance of all models degrades across the evaluation metrics, showing the benefits of using the additional context information. The degradation is particularly significant for RIVER, VE-diffusion and VP-diffusion.

Table 3: Ablation study to assess the impact of context frames on test performance of the considered models, given that the same pre-trained autoencoder is used. Results are averaged over 5 generations.

Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	1.21	1.13	20.84	0.25
VE-diffusion	1.01	1.07	23.21	0.23
VP-diffusion	1.04	1.06	21.18	0.24
Stochastic interpolant	5.37e-03	7.04e-02	41.36	0.97
Ours ( $\sigma = 0.01, \sigma_{sam} = 0$ )	<b>8.42e-04</b>	<b>3.18e-02</b>	<b>46.75</b>	<b>0.99</b>

We further assess the impact of the random conditioning frames  $z^c$ , considering the case when the neural networks are trained with the reference frames but without using random conditioning frames, in which case  $C = z^{\tau-1}$  in (10). Table 4 shows the results. We see that in this case RIVER and VE-diffusion achieve lower test performance across all evaluation metrics, which shows that these models indeed benefit from the use of conditional frames. Interestingly, our model, the stochastic interpolant model and the VP-diffusion model remain relatively robust, with comparable or even improved test performance. Despite this, our model outperforms the other models across all evaluation metrics regardless of whether the conditioning frames are involved. For consistency and fair comparison with RIVER, we retain the conditioning scheme by default.

Table 4: Ablation study to assess the impact of the random conditioning frames on test performance of the considered models, given that the same pre-trained autoencoder is used. Results are averaged over 5 generations.

Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	4.68e-03	6.19e-02	43.26	0.97
VE-diffusion	3.96e-01	6.04e-01	26.40	0.53
VP-diffusion	1.25e-03	3.67e-02	45.73	0.99
Stochastic interpolant	2.33e-03	4.90e-02	43.99	0.98
Ours ( $\sigma = 0.01, \sigma_{sam} = 0$ )	<b>3.21e-04</b>	<b>2.16e-02</b>	<b>49.23</b>	<b>1.00</b>

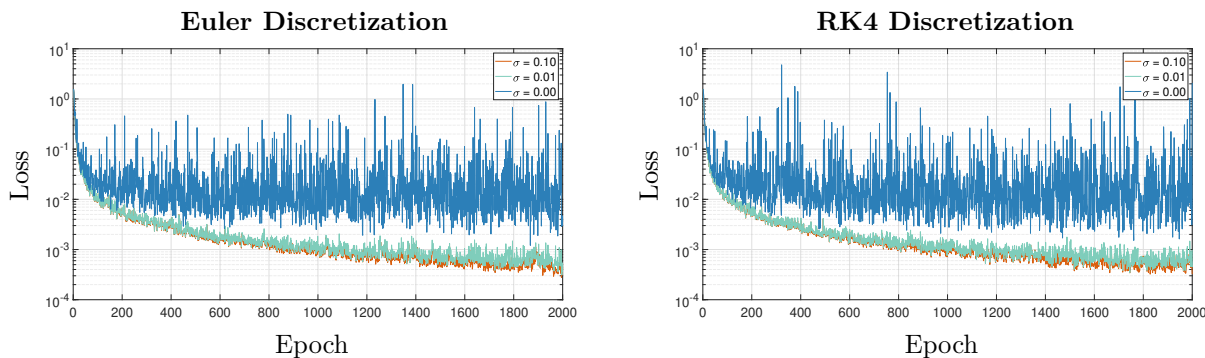


Figure 5: Training loss for different values of  $\sigma$  using our probability path model. The left subplot shows results for the Euler sampler, and the right subplot for the RK4 sampler. We see that the loss curve is sensitive to the choice of  $\sigma$ , with larger values of  $\sigma$  giving smoother loss curves.

**Sensitivity to  $\sigma_{min}$ .** Table 5 shows that using  $\sigma_{min} > 0$  (while fixing the other hyperparameters) leads to noticeable improvement in test performance of our model across the evaluation metrics. This is in line with the observation that using  $\sigma_{min} > 0$  improves training stability, since this alleviates the singularity of the target VF (20) at  $t = 0, 1$ . In fact, we see that using  $\sigma_{min} > 0.001$  (the default value that we use) leads to better test results for the FPC task.

Table 5: Sensitivity analysis of our model to  $\sigma_{min}$ . Results are averaged over 5 generations.

$\sigma_{min}$	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
0.0	1.98e-03	4.41e-02	44.74	0.99
0.001	3.79e-04	2.30e-02	48.88	1.00
0.01	3.78e-04	2.26e-02	49.11	1.00
0.1	3.17e-04	2.12e-02	49.52	1.00
1.0	3.46e-04	2.17e-02	49.44	1.00

## F Experimental Details

In this section, we provide the experimental details for the tasks considered in Section 5.

### F.1 On the Choice of Gaussian Reference Processes

Our choice of a Gaussian reference process is not directly derived from the characteristics of a particular spatio-temporal dataset, but rather follows a modeling convention that is widely adopted, particularly when physical or statistical smoothness assumptions are in play. Gaussian processes are a standard modeling choice in these domains for several reasons:

Table 6: Ablation study for the fluid flow past a cylinder task. Results are averaged over 5 generations.

$\sigma$	sampler	$N$	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
0.0	Euler	5	7.02e-04	2.96e-02	47.15	0.99
0.01	Euler	5	3.63e-04	4.14e-03	48.89	1.00
0.1	Euler	5	2.97e-03	4.11e-02	43.99	0.98
0.0	Euler	10	7.31e-04	3.01e-02	47.02	0.99
0.01	Euler	10	3.76e-04	2.31e-02	48.76	1.00
0.1	Euler	10	2.85e-03	5.04e-02	44.12	0.98
0.0	RK4	10	3.89e-04	2.26e-02	49.19	1.00
0.01	RK4	10	3.79e-04	2.30e-02	48.88	1.00
0.1	RK4	10	6.56e-03	7.83e-02	40.68	0.97
0.0	Euler	20	7.40e-04	3.03e-02	46.98	0.99
0.01	Euler	20	3.82e-04	2.33e-02	48.71	1.00
0.1	Euler	20	2.78e-03	4.98e-02	44.19	0.98
0.0	RK4	20	3.88e-04	2.26e-02	49.19	1.00
0.01	RK4	20	6.53e-04	2.80e-02	47.63	0.99
0.1	RK4	20	6.51e-04	7.80e-02	40.70	0.97

- They arise as solutions to linear stochastic partial differential equations (e.g., the heat equation with additive white noise), which are common in spatio-temporal physical systems.
- They offer desirable analytical properties such as closed-form marginals, smooth sample paths, and tractable likelihoods.
- In many practical settings (e.g., climate modeling, diffusion, spatial statistics), Gaussian priors are used because they act as regularizing priors on functions or fields that evolve over space and time.

## F.2 Details on the Datasets

**Fluid flow past a cylinder (FPC).** We use the fluid flow past a stationary cylinder at a Reynolds number of 100 as a simple test problem. This fluid flow is a canonical problem in fluid dynamics characterized by a periodically shedding wake structure (Erichson et al., 2020; 2019). The flow dynamics are governed by the two-dimensional incompressible Navier–Stokes equations:

$$\frac{\partial \mathbf{u}}{\partial t} + (\mathbf{u} \cdot \nabla) \mathbf{u} = -\frac{1}{\rho} \nabla p + \nu \nabla^2 \mathbf{u},$$

$$\nabla \cdot \mathbf{u} = 0,$$

where  $\mathbf{u} = (u, v)$  is the velocity field,  $p$  is the pressure,  $\rho$  is the fluid density, and  $\nu$  is the kinematic viscosity. The vorticity field  $\omega$  is obtained from the velocity field via:

$$\omega = \nabla \times \mathbf{u},$$

providing insights into the rotational characteristics of the flow.

For simulating the data, the Immersed Boundary Projection Method (IBPM) has been used (Colonius & Taira, 2008). The flow tensor has dimensions  $199 \times 449 \times 151$ , representing 151 temporal snapshots on a  $449 \times 199$  spatial grid. We crop and spatially subsample the data which results in a  $64 \times 64$  spatial field.

**Shallow-water equation (SWE).** The shallow-water equations, derived from the compressible Navier–Stokes equations, can be used for modeling free-surface flow problems. We consider the 2D equation used in (Takamoto et al., 2022), which is the following system of hyperbolic PDEs:

$$\partial_t h + \nabla h \mathbf{u} = 0, \quad \partial_t h \mathbf{u} + \nabla \left( \mathbf{u}^2 h + \frac{1}{2} g_r h^2 \right) = -g_r h \nabla b, \quad (81)$$

where  $\mathbf{u} = u, v$  being the velocities in the horizontal and vertical direction respectively,  $h$  describes the water depth, and  $b$  describes a spatially varying bathymetry.  $h\mathbf{u}$  can be interpreted as the directional momentum components and  $g_r$  describes the gravitational acceleration. The mass and momentum conservation properties can hold across shocks in the solution and thus challenging datasets can be generated. This equation finds application in modeling tsunamis and flooding events.

We use the dataset generated and provided by PDEBench (Takamoto et al., 2022). The data file (2D\_rdb\_NA\_NA.h5) can be downloaded from [https://github.com/pdebench/PDEBench/tree/main/pdebench/data\\_download](https://github.com/pdebench/PDEBench/tree/main/pdebench/data_download). The data sample is a series of 101 frames at a  $128 \times 128$  pixel resolution and come with 1 channel. The simulation considered in (Takamoto et al., 2022) is a 2D radial dam break scenario. On a square domain  $\Omega = [-2.5, 2.5]^2$ , the water height is initialized as a circular bump in the center of the domain:

$$h(t = 0, x, y) = \begin{cases} 2, & \text{for } r < \sqrt{x^2 + y^2}, \\ 1, & \text{for } r \geq \sqrt{x^2 + y^2}, \end{cases}$$

with the radius  $r$  randomly sampled from  $\mathcal{U}(0.3, 0.7)$ . The dataset is simulated with a finite volume solver using the *PyClaw* package. We apply standardization and then normalization to the range of  $[-1, 1]$  to preprocess the simulated data.

**Incompressible Navier-Stokes equation (NSE).** The Navier-Stokes equation is the incompressible version of the compressible Navier-Stokes equation, and it can be used to model hydromechanical systems, turbulent dynamics and weather. We use the inhomogeneous version of the equation (which includes a vector field forcing term  $\mathbf{u}$ ) considered by (Takamoto et al., 2022):

$$\nabla \cdot \mathbf{v} = 0, \quad \rho(\partial_t \mathbf{v} + \mathbf{v} \cdot \nabla \mathbf{v}) = -\nabla p + \eta \Delta \mathbf{v} + \mathbf{u}, \quad (82)$$

where  $\rho$  is the mass density,  $\mathbf{v}$  is the fluid velocity,  $p$  is the gas pressure and  $\eta$  is shear viscosity. The initial conditions  $\mathbf{v}_0$  and inhomogeneous forcing parameters  $\mathbf{u}$  are each drawn from isotropic Gaussian random fields with truncated power-law decay  $\tau$  of the power spectral density and scale  $\sigma$ , where  $\tau_{v_0} = -3$ ,  $\sigma_{v_0} = 0.15$ ,  $\tau_u = -1$ ,  $\sigma_u = 0.4$ . The domain is taken to be the unit square  $\Omega = [0, 1]^2$  and the viscosity  $\eta = 0.01$ . The equation is numerically simulated using *Phiflow*. Boundary conditions are taken to be Dirichlet to clamp the field velocity to zero at the perimeter.

We use the dataset generated and provided by PDEBench (Takamoto et al., 2022). The data file (ns\_incom\_inhom\_2d\_512-0.h5) can be downloaded from [https://github.com/pdebench/PDEBench/tree/main/pdebench/data\\_download](https://github.com/pdebench/PDEBench/tree/main/pdebench/data_download). The data sample is a series of 1000 frames at a  $512 \times 512$  pixel resolution and come with 2 channels. We do not apply any data preprocessing procedure here.

**Diffusion-reaction equation (DRE).** We use the 2D extension of diffusion-reaction equation of (Takamoto et al., 2022) which describes two non-linearly coupled variables, namely the activator  $u = u(t, x, y)$  and the inhibitor  $v = v(t, x, y)$ . The equation is given by:

$$\partial_t u = D_u \partial_{xx} u + D_u \partial_{yy} u + R_u, \quad (83)$$

$$\partial_t v = D_v \partial_{xx} v + D_v \partial_{yy} v + R_v, \quad (84)$$

where  $D_u$  and  $D_v$  are the diffusion coefficient for the activator and inhibitor respectively,  $R_u = R_u(u, v)$  and  $R_v = R_v(u, v)$  are the activator and inhibitor reaction function respectively. The domain of the simulation includes  $x \in (-1, 1)$ ,  $y \in (-1, 1)$ ,  $t \in (0, 5]$ . This equation can be used for modeling biological pattern formation.

The reaction functions for the activator and inhibitor are defined by the Fitzhugh-Nagumo equation as:  $R_u(u, v) = u - u^3 - k - v$ ,  $R_v(u, v) = u - v$ , where  $k = 5 \times 10^{-3}$ , and the diffusion coefficients for the activator and inhibitor are  $D_u = 1 \times 10^{-3}$  and  $D_v = 5 \times 10^{-3}$  respectively. The initial condition is generated as standard Gaussian noise  $u(0, x, y) \sim \mathcal{N}(0, 1.0)$  for  $x \in (-1, 1)$  and  $y \in (-1, 1)$ . We take a no-flow Neumann boundary condition:  $D_u \partial_x u = 0$ ,  $D_v \partial_x v = 0$ ,  $D_u \partial_y u = 0$ , and  $D_v \partial_y v = 0$  for  $x, y \in (-1, 1)^2$ .

We use a downsampled version of the dataset provided by PDEBench (Takamoto et al., 2022). The data file (2D\_diff-react\_NA\_NA.h5) can be downloaded from <https://github.com/pdebench/PDEBench/tree/>

[main/pdebench/data\\_download](#). The data sample is a series of 101 frames at a  $128 \times 128$  pixel resolution and come with 2 channels. The sample frames are generated using the finite volume method for spatial discretization, and the time integration is performed using the built-in fourth order Runge-Kutta method in the *scipy* package. We do not apply any data preprocessing procedure here.

### F.3 Details on Pre-Training the Autoencoder

We provide details on pre-training the autoencoder here. The choice of first pre-training an autoencoder is motivated by the computational challenges of working directly with the high-dimensional spatial resolution of PDE datasets. Training directly in the ambient space requires substantial GPU memory and computational resources, making it impractical for large-scale or high-resolution datasets. By leveraging a latent-space representation, we achieve significant dimensionality reduction while preserving the essential structure of the data, enabling efficient training and inference with standard hardware configurations. For these datasets, latent-space modeling provides a critical balance between computational efficiency, scalability, and performance.

We use the same architecture for the encoder and decoder for all the tasks, with the architecture parameters chosen based on the complexity of the task.

**The encoder.** The encoder first applies a 2D convolution (`conv_in`) to the input frame, which reduces the number of channels from `in_channels` to `mid_channels`, and processes the spatial dimensions. Then, a series of ResidualBlock layers, which progressively process and downsample the feature map, making it smaller in spatial dimensions but more enriched in terms of features, are applied. After the residual blocks, the feature map undergoes an attention process via a multi-head attention layer. This layer helps the encoder focus on important parts of the input, learning relationships between spatial positions in the image. For the post-attention step, the feature map is further processed by residual blocks and normalized, preparing it for the final convolution. The output of the encoder is obtained by applying a final 2D convolution (`out_conv`), which maps the processed feature map to the desired number of output channels (`out_channels`).

**The decoder.** The decoder takes the encoded feature map and transforms it back into an output with similar spatial dimensions as the input. Similar to the encoder, the decoder starts with a convolution that adjusts the number of channels from `in_channels` to `mid_channels`. Then, an attention mechanism (similar to the encoder) is applied to focus on important aspects of the encoded features. Next, a series of ResidualBlock layers, combined with UpBlock layers, are used to progressively increase the spatial dimensions of the feature map (upsampling), undoing the compression applied by the encoder. After the upsampling, the output is normalized and passed through a final convolution (`out_conv`), mapping the internal feature representation to the desired number of output channels (`out_channels`).

Table 7 summarizes the architecture parameters used for the considered tasks.

Table 7: Parameters chosen for the encoder (decoder) architecture.

Task	Fluid flow	Shallow-water eq.	Navier-Stokes eq.	Diffusion-reaction eq.
<code>in_channels</code>	1 (1)	1 (1)	2 (2)	2 (2)
<code>out_channels</code>	1 (1)	1 (1)	2 (2)	2 (2)
<code>mid_channels</code>	64 (128)	128 (256)	128 (256)	128 (256)

**Training details.** We train the autoencoder using AdamW with batch size of 32, no weight decay and  $\beta = (0.9, 0.999)$ . We use the cosine learning rate scheduler with warmup. For the fluid flow past a cylinder task, we train for 2000 epochs and use learning rate of 0.001. For the Navier-Stokes task, we train for 500 epochs and use learning rate of 0.0001. For the other two tasks we train for 5000 epochs and use learning rate of 0.0005. Our implementation is in PyTorch, and all experiments are run on an NVIDIA A100-SXM4 GPU with 40 GB VRAM belonging to an internal SLURM cluster.

**Impact of the autoencoder (AE) size on test performance.** We consider using pre-trained AEs with increasing number of middle channels (`mid_channels`) in the encoder and decoder, resulting in AEs with



374,605, 1,485,589, 5,916,709 and 23,615,557 trainable parameters respectively. Table 8 shows the effect of the size of AE on the test performance for all models. We see that there is an optimal size of AE that leads to the best test performance for RIVER, the stochastic interpolant model and our model. For VE-diffusion, the best test performance is achieved with the smallest AE, suggesting that the model might overfit when AE increases in size or suffer from optimization issues in bigger latent spaces. For VP-diffusion, the test results are relatively insensitive to the size of the AE used.

Table 8: Ablation study to assess the impact of the size of the pre-trained autoencoder (AE) on test performance using the fluid flow (FPC) task. The number of channels in the decoder is in parenthesis. Results are averaged over 5 generations.

Model	mid_channels	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	16 (32)	2.54e-03	5.43e-02	42.32	0.98
	32 (64)	2.56e-03	4.89e-02	<b>44.55</b>	<b>0.99</b>
	64 (128)	<b>1.86e-03</b>	<b>4.48e-02</b>	44.30	<b>0.99</b>
	128 (256)	2.69e-03	5.52e-02	42.48	0.98
VE-diffusion	16 (32)	<b>1.91e-01</b>	<b>3.99e-01</b>	<b>29.45</b>	<b>0.68</b>
	32 (64)	4.28e-01	6.19e-01	27.20	0.49
	64 (128)	2.29e-01	4.74e-01	27.36	0.55
	128 (256)	5.12e-01	7.22e-01	26.65	0.42
VP-diffusion	16 (32)	<b>2.02e-03</b>	5.11e-02	42.32	<b>0.98</b>
	32 (64)	2.84e-03	5.11e-02	<b>44.14</b>	<b>0.98</b>
	64 (128)	3.58e-03	6.09e-02	42.37	<b>0.98</b>
	128 (256)	2.03e-03	<b>4.97e-02</b>	42.97	<b>0.98</b>
Stoch. interp.	16 (32)	1.44e-02	1.22e-01	36.35	0.92
	32 (64)	1.09e-02	9.01e-02	41.15	0.95
	64 (128)	<b>3.40e-03</b>	<b>6.10e-02</b>	<b>41.81</b>	<b>0.98</b>
	128 (256)	2.41e-02	1.36e-01	38.02	0.90
Ours ( $\sigma = 0.01, \sigma_{sam} = 0$ )	16 (32)	2.49e-03	5.70e-02	41.30	0.98
	32 (64)	<b>2.87e-04</b>	<b>1.93e-02</b>	<b>50.62</b>	<b>1.00</b>
	64 (128)	3.79e-04	2.30e-02	48.88	<b>1.00</b>
	128 (256)	3.53e-04	2.24e-02	48.96	<b>1.00</b>

#### F.4 Details on Training the Flow Matching Models

**Architecture.** The vector field regressor is a transformer-based model designed to process latent vector fields and predict refined outputs with spatial and temporal dependencies. It uses key parameters like `depth` and `mid_depth`, which control the number of transformer encoder layers in the input, middle, and output stages. The `state_size` and `state_res` parameters define the number of channels and spatial resolution of the input data, while `inner_dim` sets the embedding dimension for processing. The model uses learned positional encodings and a time projection to incorporate spatial and temporal context into the input, which can include `input_latents`, `reference_latents`, and `conditioning_latents`. The input is projected into the inner dimension and passed through a series of transformer layers, with intermediate outputs from the input blocks concatenated with the output layers to refine predictions. Finally, the model projects the processed data back to the original spatial resolution and channel size using BatchNorm, producing the final vector field output.

Table 9 summarizes the architecture parameters used for the considered tasks.

**Training details.** For all the considered tasks, we train the regressor using AdamW with batch size of 32, learning rate of 0.00005, no weight decay and  $\beta = (0.9, 0.999)$ . We use the cosine learning rate scheduler with warmup. For the fluid flow past cylinder, we train for 2000 epochs, for the shallow-water equation and diffusion-reaction task we train for 1000 epochs, and for the Navier-Stokes task we train for 100 epochs. Our implementation is in PyTorch, and all experiments are run on a single NVIDIA A100-SXM4 GPU with 40 GB VRAM belonging to an internal SLURM cluster. For the Navier-Stokes task, due to GPU memory

Table 9: Parameters chosen for the vector field neural network.

Parameter	Fluid flow	Shallow-water eq.	Navier-Stokes eq.	Diffusion-reaction eq.
<code>state_size</code>	4	4	8	4
<code>state_res</code>	[8,8]	[16, 16]	[64, 64]	[16, 16]
<code>inner_dim</code>	512	512	512	512
<code>depth</code>	4	4	4	4
<code>mid_depth</code>	5	5	5	5

constraints in the experiments, we use gradient accumulation to simulate the desired batch size while training with a smaller per-step batch.

### F.5 Details on the Evaluation Metrics

In addition to the standard mean squared error and relative Frobenius norm error (RFNE), we use the Pearson correlation coefficient to measure the linear relationship between the forecasted frames and the target frames. The range of this coefficient is  $[-1, 1]$ , with zero implying no correlation. Correlations of  $-1$  or  $+1$  imply an exact linear relationship. Positive correlations imply that as  $x$  increases, so does  $y$ . Negative correlations imply that as  $x$  increases,  $y$  decreases. In addition, we use peak signal-to-noise ratio (PSNR) to evaluate the quality of signal representation against corrupting noise, and structural similarity index measure (SSIM) (Wang et al., 2004) to assess perceptual results. The presented results are computed by averaging over batch size and number of sample generations.

### F.6 Standard Deviations for the Presented Results

Table 10-14 provide the standard deviation of the results presented in the main paper.

Table 10: Standard deviation results for the fluid flow past a cylinder task using different choices of probability paths for flow matching. Results are averaged over 5 generations.

Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	1.04e-03	2.70e-02	1.46	4.92e-03
VE-diffusion	5.22e-02	3.34e-01	6.02e-01	3.14e-02
VP-diffusion	7.09e-03	3.96e-02	4.26e-01	3.30e-03
SI ( $b_t = t^2$ )	6.35e-05	3.53e-02	5.39e-02	4.03e-04
SI ( $b_t = t$ )	1.48e-04	8.70e-02	2.97e-02	5.30e-04
Ours ( $\sigma = 0.01$ , $\sigma_{sam} = 0$ , RK4)	4.26e-06	5.29e-03	3.61e-02	3.63e-05

Table 11: Standard deviation for the results of ablation study for the fluid flow past a cylinder task. Results are averaged over 5 generations.

$\sigma$	sampler	$N$	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
0.0	Euler	5	1.00e-05	1.25e-02	4.37e-02	9.76e-05
0.01	Euler	5	6.71e-06	4.14e-03	5.17e-02	4.81e-05
0.1	Euler	5	4.88e-05	4.11e-02	5.39e-02	3.91e-04
0.0	Euler	10	1.24e-05	1.29e-02	4.28e-02	8.27e-05
0.01	Euler	10	8.10e-06	4.37e-03	6.42e-02	6.21e-05
0.1	Euler	10	4.32e-05	4.02e-02	5.13e-02	3.06e-04
0.0	RK4	10	3.45e-06	7.56e-03	2.57e-02	2.46e-05
0.01	RK4	10	4.26e-06	5.29e-03	3.61e-02	3.63e-05
0.1	RK4	10	8.62e-06	5.88e-02	2.69e-02	2.61e-04
0.0	Euler	20	8.14e-06	1.31e-02	3.11e-02	5.02e-05
0.01	Euler	20	4.25e-06	4.51e-03	3.66e-02	2.77e-05
0.1	Euler	20	2.12e-05	3.96e-02	2.21e-02	1.03e-04
0.0	RK4	20	1.54e-06	7.49e-03	1.84e-02	1.26e-05
0.01	RK4	20	3.99e-06	1.29e-02	1.74e-02	2.95e-05
0.1	RK4	20	2.84e-05	5.85e-02	8.10e-03	8.33e-05

Table 12: Standard deviation results for the shallow-water equation task using different choices of probability paths for flow matching. Results are averaged over 5 generations.

Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	2.28e-05	8.67e-02	8.39e-02	1.23e-03
VE-diffusion	7.29e-04	3.29e-01	2.17e-01	1.68e-02
VP-diffusion	2.55e-04	1.10e-01	3.62e-01	1.56e-02
SI ( $b_t = t^2$ )	2.97e-06	1.03e-01	1.86e-02	2.03e-04
SI ( $b_t = t$ )	1.37e-06	7.18e-02	4.11e-03	2.12e-04
Ours ( $\sigma = 0.1, \sigma_{sam} = 0, \text{RK4}$ )	1.55e-06	7.06e-02	5.98e-03	2.11e-04

Table 13: Standard deviation results for the diffusion-reaction equation task using different choices of probability paths for flow matching. Results are averaged over 5 generations.

Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	5.56e-04	2.54e-01	7.88e-01	1.37e-02
VE-diffusion	1.51e-02	1.11	5.89e-01	1.22e-02
VP-diffusion	1.55e-03	6.15e-01	1.06	7.08e-03
SI ( $b_t = t^2$ )	6.44e-04	1.15	8.60e-02	6.98e-04
SI ( $b_t = t$ )	8.15e-07	5.01e-02	5.67e-03	8.25e-05
Ours ( $\sigma = 0, \sigma_{sam} = 0, \text{RK4}$ )	8.47e-07	4.77e-02	6.22e-03	9.53e-05

Table 14: Standard deviation results for the Navier-Stokes equation task using different choices of probability paths for flow matching. Results are averaged over 5 generations.

Model	Test MSE ( $\downarrow$ )	Test RFNE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )
RIVER	7.52e-04	8.31e-01	9.34e-02	3.80e-03
VE-diffusion	6.53e-04	1.21	1.19e-01	1.61e-03
VP-diffusion	3.88e-03	1.62	1.64e-01	5.74e-03
SI ( $b_t = t^2$ )	5.11e-06	3.27e-02	2.37e-02	3.34e-04
SI ( $b_t = t$ )	4.87e-07	1.09e-02	1.98e-03	6.22e-05
Ours ( $\sigma = 0.1, \sigma_{sam} = 0, \text{RK4}$ )	4.81e-07	1.06e-02	2.30e-03	6.01e-05