052

053

054

000

## **Conditional Lagrangian Wasserstein Flow for Time Series Imputation**

Anonymous Authors<sup>1</sup>

## Abstract

Time series imputation is important for numerous real-world applications. To overcome the limitations of diffusion model-based imputation methods, e.g., slow convergence in inference, we propose a novel method for time series imputation in this work, called Conditional Lagrangian Wasserstein Flow (CLWF). Following the principle of least action in Lagrangian mechanics, we learn the velocity by minimizing the corresponding kinetic energy. Moreover, to enhance the model's performance, we estimate the gradient of a task-specific potential function using a time-dependent denoising autoencoder and integrate it into the base estimator to reduce the sampling variance. Finally, the proposed method demonstrates competitive performance compared to other state-of-the-art imputation approaches.

## 1. Introduction

Time series imputation is essential for various practical scenarios in many fields, such as transportation, environment, and medical care, etc. Deep learning-based approaches, such as RNNs, VAEs, and GANs, have been proved to be advantageous compared to traditional machine learning methods on various complex real-world multivariate time series analysis tasks (Fortuin et al., 2020). More recently, diffusion models, such as denoising diffusion probabilistic models (DDPMs) (Ho et al., 2020) and score-based generative models (SBGMs) (Song et al., 2020), have gained more and more attention in the field of time series analysis due to their powerful modelling capability (Lin et al., 2023; Meijer & Chen, 2024).

Although many diffusion model-based time series imputation approaches have been proposed and show their advantages compared to conventional deep learning models (Tashiro et al., 2021; Chen et al., 2021; 2023), they are limited to slow convergence or large computational costs. Such limitations may prevent them being applied to realworld applications. To address the aforementioned issues, in this work, we leverage the optimal transport theory (Villani et al., 2009) and Lagrangian mechanics (Arnol'd, 2013) to propose a novel method, called Conditional Lagrangian Wasserstein Flow (CLWF), for fast and accurate time series imputation.

In our method, we treat the multivariate time series imputation task as a conditional optimal transport problem, whereby the random noise is the source distribution, the missing data is the target distribution, and the observed data is the conditional information. To generate new data samples efficiently and accurately, we need to find the shortest path in the probability space according to the optimal transport theory. To this end, we first project the original source and target distributions into the Wasserstein space via sampling mini-batch OT maps. Afterwards, we construct the time-dependent intermediate samples through interpolating the source distribution and target distribution. Then according to the principle of least action in Lagrangian mechanics (Arnol'd, 2013), the optimal velocity function moving the source distribution to the target distribution is learned in a self-supervised manner by minimizing the corresponding kinetic energy. We can solve the model efficiently using flow matching in a simulation-free manner (Lipman et al., 2022; Liu et al., 2022; Albergo & Vanden-Eijnden, 2023; Tong et al., 2023).

To further improve the model's performance, we leverage the denoising affect of the time-dependent denoising autoencoder (TDAE) model which is trained on the observed time series data to estimate the gradient of task-specific potential function. By doing so, combined with the aforementioned flow model, we can formulate a new path sampler to reduce the sampling variances. Furthermore, we can interpret the gradient of the potential function as the control signal from the perspective of stochastic optimal control (SOC) in data generation (Bellman, 1966; Chen et al., 2021; Caluya & Halder, 2021; Berner et al., 2024), Consequently, the sampling procedure can be viewed as a controlled path integral (Zhang & Chen, 2022). We also explain the variance reduction effect of the new sampler using the Rao-Blackwell theorem (Casella & Robert, 1996). Moreover, we propose a resampling technique using the interpolated conditional

<sup>&</sup>lt;sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

055 samples to enhance the model's imputation performance.

Finally, CLWF is assessed on three real-word and one syn-057 thetic time series datasets for validation. The results ob-058 tained show that the proposed method achieves competitive 059 performance and admits faster convergence compared with 060 other state-of-the-art time series imputation methods.

The contributions of the paper are summarized as follows:

- We present Conditional Lagrangian Wasserstein Flow, a novel conditional generative framework based on the optimal transport theory and Lagrangian mechanics;
- We develop the efficient training and inference algorithms to solve the time series imputation problem;
- · We establish theoretical links between optimal transport, stochastic optimal control and path measures;
- · We demonstrate that the proposed method has achieved competitive performance on time series imputation tasks compared to other state-of-the-art methods.

#### 2. Preliminaries

061

062

063

064

065

066

067

068 069

070

071

073

074

075

076

077

078

079

080

081

082 083

084

085

086

087

088

089

090

091

092

093

094

095

096

In this section, we will succinctly introduce the fundamentals of stochastic differential equations, optimal transport, Shrödinger Bridge, and Lagrangian mechanics.

#### 2.1. Stochastic Differential Equations

We treat the data generation task as an initial value problem (IVP), in which  $X_0 \in \mathbb{R}^d$  is the initial data (e.g., some random noise) at the initial time t = 0, and  $X_T \in \mathbb{R}^d$  is target data at the terminal time t = T. To solve the IVP, we consider a stochastic differential equation (SDE) defined by a Borel measurable time-dependent drift function  $\mu_t : \mathbb{R}^d \times [0,T] \to \mathbb{R}^d$ , and a positive Borel measurable time-dependent diffusion function  $\sigma_t : [0,T] \to \mathbb{R}^d_{>0}$ . Accordingly, the Itô form of the SDE can be described as follows (Oksendal, 2013):

$$dX_t = \mu_t(X_t, t)dt + \sigma_t(t)dW_t, \tag{1}$$

097 where  $W_t$  is a Brownian motion/Wiener process. Note that 098 when the diffusion term is not considered, the SDE degen-099 erates to an ordinary differential equation (ODE), which 100 is typically easier to solve numerically. Nonetheless, we will use the SDE for theoretical analysis throughout the paper, as it provides a more general framework. Accordingly, The above SDE's associated forward Fokker-Planck Kol-104 mogorov (FPK) equatio (Risken & Frank, 2012) describing 105 the evolution of the marginal density  $p_t(X_t)$  reads 106

$$\frac{\partial p_t}{\partial t} + \boldsymbol{\nabla} \cdot (p_t \mu_t) = \langle D(t), \nabla^2(p_t) \rangle, \qquad (2)$$

where  $D(t) := \frac{1}{2}\sigma^{\top}(t)\sigma(t)$ ,  $\nabla^2$  represents the Hessian operator, and  $\langle \cdot, \cdot \rangle := \operatorname{trace}(\cdot^{\top}, \cdot)$  represents the Frobenius inner product

In fact, both Eq. (1) and Eq. (2) reveal the system's dynamics and act as the boundary conditions for the optimization problems introduced in later sections, each with a different focus. When the constraint is given by Eq. (1), the formalism is Lagrangian, depicting the movement of each individual particle. In contrast, when the constraint is Eq. (2), the formalism is Eulerian, representing the evolution of the population as a whole.

#### 2.2. Optimal Transport

The optimal transport (OT) problem aims to find the optimal transport plans/maps that move the source distribution to the target distribution (Villani et al., 2009; Santambrogio, 2015; Peyré et al., 2019). In the Kantorovich's formulation of the OT problem, the transport costs are minimized with respect to some probabilistic couplings/joint distributions (Villani et al., 2009; Santambrogio, 2015; Peyré et al., 2019). Let  $p_0$  and  $p_T$  be two Borel probability measures with finite second moments on the space  $\Omega \in \mathbb{R}^d$ .  $\Pi(p_0, p_T)$  denotes a set of transport plans between these two marginals. Then, the Kantorovich's OT problem is defined as follows:

$$\inf_{\pi \in \Pi(p_0, p_T)} \int_{\mathcal{X} \times \mathcal{Y}} \frac{1}{2} \|x - y\|^2 \pi(x, y) \mathrm{d}x \mathrm{d}y, \qquad (3)$$

where  $\Pi(p_0, p_T) = \{ \pi \in \mathcal{P}(\mathcal{X} \times \mathcal{Y}) : (\pi^x)_{\#} \pi = p_0, (\pi^y)_{\#} \pi = p_T \},\$ with  $\pi^x$  and  $\pi^y$  being two projections of  $\mathcal{X} \times \mathcal{Y}$  on  $\Omega$ . The minimizer of Eq. (3),  $\pi^*$ , always exists and is referred to as the OT plan.

Note that Eq. (3) can also include an entropy regularization term, the Kullback–Leibler (KL) divergence  $D_{\text{KL}}(\pi || p_0 \otimes p_T)$ . This transforms the original OT problem into the entropyregularized optimal transport (EROT) problem with Eq. (2) serving as the constraint, which frames the transport problem better in terms of convexity and stability (Cuturi, 2013). In particular, from a data generation perspective,  $p_0$  is some random initial noise and  $p_T$  is the target data distribution, and we can sample the corresponding OT plan in a mini-batch manner (Tong et al., 2023; 2024; Pooladian et al., 2023).

#### 2.3. Shrödinger Bridge

The transport problem in Sec. 2.2 can be further viewed from a distribution evolution perspective, which is particularly suitable for developing the flow-based models that model data generation process. For this reason, the Shrödinger Bridge (SB) problem is introduced herein (Léonard, 2012). Assume that  $\Omega \in C^1(\mathbb{R}^d \times [0,T]), \mathcal{P}(\Omega)$  is a probability path measure on the path space  $\Omega$ , then the goal of the SB

problem is to find the following optimal path measure:

111

112 113 114

115

116

117 118

119

120

121

122

123

124

125 126

127 128

129

130

131

132

140

141

147

148 149

150

151

152

153

154

155

156

157

$$\mathbb{P}^{*} = \underset{\mathbb{P} \in \mathcal{P}(\Omega)}{\arg\min} D_{\mathrm{KL}}(\mathbb{P} \| \mathbb{Q}),$$
subject to  $\mathbb{P}_{+} = a_{+}$  and  $\mathbb{P}_{-} = a_{-}$ 
(4)

subject to  $\mathbb{P}_0 = q_0$  and  $\mathbb{P}_T = q_T$ ,

where  $D_{\mathrm{KL}}(\mathbb{P}||\mathbb{Q}) = \begin{cases} \log\left(\frac{\mathrm{d}\mathbb{P}}{\mathrm{d}\mathbb{Q}}\right)\mathrm{d}\mathbb{P}, & \text{if } \mathbb{P} \ll \mathbb{Q}, \\ +\infty, & \text{otherwise}, \end{cases}$  and  $\mathbb{Q}$  is

a reference path measure, e.g., Brownian motion or Ornstein-Uhlenbeck process. Moreover, the distribution matching problem in Eq. (3) can be reframed as a dynamical SB problem as well (Gushchin et al., 2024; Koshizuka & Sato, 2023; Liu et al., 2024):

$$\arg\min_{\theta} \mathbb{E}_{p(X_t)} \left[ \frac{1}{2} \left\| \mu_t^{\theta}(X_t, t) \right\|^2 \right],$$
  
subject to Eq. (1) or Eq. (2), (5)

where  $\theta$  is the parameters of the variational drift function  $\mu_t$ .

#### 2.4. Lagrangian Mechanics

In this section, we formulate the data generation problem within the framework of Lagrangian mechanics (Arnol'd, 2013). Let  $p_t$  and  $\dot{p}_t := dp_t/dt$  be the density and law of the generalized coordinates  $X_t$ , respectively. Denoting the kinetic energy as  $\mathcal{K}(p_t, \dot{p}_t, t)$  and the potential energy as  $\mathcal{U}(p_t, t)$ , then the corresponding Lagrangian is given by

$$\mathcal{L}(p_t, \dot{p}_t, t) = \mathcal{K}(p_t, \dot{p}_t, t) - \mathcal{U}(p_t).$$
(6)

Further, we assume that Eq. (6) is lower semi-continuous (lsc) and strictly convex in  $\dot{p}_t$  in the Wasserstein space. Consequently,  $\mathcal{K}(x_t, \mu_t, t)$  and  $\mathcal{U}(p_t, t)$  are defined as follows, respectively:

$$\mathcal{K}(x_t, \mu_t, t) := \mathbb{E}_{p_t(x_t)} \left[ \int_0^T \frac{1}{2} \|\mu_t(x_t, t)\|^2 \mathrm{d}t \right], \quad (7)$$

$$\mathcal{U}(p_t, t) := \int_{\mathbb{R}_d} V_t(x_t) p_t(x_t) dx_t, \tag{8}$$

where  $V_t(X_t)$  is the potential function. Then the *action* in the context of Lagrangian mechanics is defined as follows:

$$\mathcal{A}(\mu_t(x_s t)) := \int_0^T \int_{\mathbb{R}_d} \mathcal{L}(x_t, \mu_t, t) dx_t dt.$$
(9)

According to *the principle of least action* (Feynman, 2005),
the shortest path is the one minimizing the action, which
is aligned with Eq. (4) in the SB theory as well. Therefore, we can leverage the Lagrangian dynamics to tackle
the OT problem for data generation. Moreover, to solve
Eq. (6), the corresponding stationary condition, i.e., the

Euler-Lagrangian equation (Arnol'd, 2013), needs to be satisfied:

$$\frac{d}{dt}\frac{\partial}{\partial \dot{p}_t}\mathcal{L}(x_t,\mu_t,t) = \frac{\partial}{\partial p_t}\mathcal{L}(p_t,\dot{p}_t,t), \qquad (10)$$

with the boundary conditions:  $\frac{dX_t}{dt} = \mu_t$ ,  $p_0 = q_0$ , and  $p_T = q_T$ .

## **3.** Conditional Lagrangian Wasserstein Flow for Time Series Imputation

In this section, building on the theory introduced in Sec. 2, we propose Conditional Lagrangian Wasserstein Flow, a novel conditional generative method for time series imputation.

#### **3.1. Time Series Imputation**

Our goal is to impute the missing time series data points based on the observations. To this end, we adopt a conditionally generative approach for time series imputation in the sample space  $\mathbb{R}^{K \times L}$ , where K represents the dimension of the multivariate time series and L represents sequence length. In our self-supervised learning approach, the total observed data  $x^{\text{obs}} \in \mathbb{R}^{K \times L}$  are partitioned into the imputation target  $x^{\text{tar}} := x^{\text{obs}} \odot M^{\text{tar}}$  and the condition  $x^{\text{cond}} := x^{\text{obs}} \odot M^{\text{cond}}$ , where  $\odot$  denotes the Hadamard product,  $M^{\text{cond}} \in \mathbb{R}^{K \times L}$  and  $M^{\text{tar}} \in \mathbb{R}^{K \times L}$  are the condition and target masks, respectively.

Consequently, the missing data points  $x^{\text{tar}}$  can be generated based on the conditions  $x^{\text{cond}}$  joint with some uninformative initial distribution  $x_0 \in \mathbb{R}^{K \times L}$  (e.g., Gaussian noise) at the initial time t = 0. Thereby, the imputation task can be described as:  $x^{\text{tar}} \sim p(x^{\text{tar}}|x_0^{\text{in}})$ , where  $x_0^{\text{in}} := \text{Concatenate}(x^{\text{cond}}, x_0) \in \mathbb{R}^{K \times L \times 2}$  is the the total input of the model.

#### 3.2. Interpolation in Wasserstein Space

To solve Eq. (7), we need to sample the intermediate variable  $X_t$  in the Wasserstein space first. To do so, the interpolation method is adopted to construct the intermediate samples (Liu et al., 2022; Albergo & Vanden-Eijnden, 2023; Tong et al., 2024). According to the OT and SB problems introduced in Sec. 2, we define the following time-differentiable interpolant:

$$I_t: \Gamma \times \Gamma \to \Gamma$$
 such that  $I_0 = X_0$  and  $I_T = X_T$ , (11)

where  $\Gamma \in \mathbb{R}^d$  is the support of the marginals  $p_0(X_0)$  and  $p_T(X_T)$ , as well as the conditional  $p(X_t|X_0, X_T, t)$ .

To implement Eq. (11), we first independently sample some random noise  $X_0 \sim \mathcal{N}(0, \sigma_0^2)$  at the initial time t = 0 and the data samples  $X_T \sim p(x^{\text{tar}})$  at the terminal time t = T, 165 respectively. Afterwards, the interpolation method is used to 166 construct the intermediate samples  $X_t \sim p(X_t|X_0, X_T, t)$ , 167 where  $t \sim \text{Uniform}(0, T)$ . More specifically, we design the 168 following sampling approach:

169

170

171

172

173

213

214

$$X_t = \frac{t}{T} (X_T + \gamma_t) + (1 - \frac{t}{T}) X_0$$
$$+ \alpha(t) \sqrt{\frac{t(T-t)}{T}} \epsilon, \quad t \in [0, T], \quad (12)$$

where  $\gamma_t \sim \mathcal{N}(0, \sigma_{\gamma}^2)$  is some random noise with variance  $\sigma_{\gamma}$  injected to the target data samples to improve the coupling's generalization property,  $\alpha(t) \geq 0$  is a timedependent scalar, and  $\epsilon \sim \mathcal{N}(0, \mathbb{I})$ .

179 Note that Eq. (12) can only allow us to generate timedependent intermediate samples in the Euclidean space but 180 not the Wasserstein space, which can lead to slow conver-181 gence as the sampling paths are not straightened. Hence, to 182 183 address this issue, we can project the samples  $X_T$  and  $X_0$ in the Wasserstein space before interpolating to strengthen 184 185 the probability flow. To this end, we leverage the method adopted in (Tong et al., 2023; 2024; Pooladian et al., 2023) 186 187 to sample the optimal mini-batch OT maps between  $X_0$ 188 and  $X_T$  first, and perform the interpolations according to Eq. (12) afterwards. Finally, we have the joint variable 189  $x_t^{\text{in}} := (x^{\text{cond}}, x_t)$  as the input for computing the velocity 190 of the Wasserstein flow. 191

#### 193 3.3. Velocity Estimation via Flow Matching

To estimate the velocity of the Wasserstein flow  $\mu_t(X_t, t)$ 195 in Eq. (1), the previous methods that require trajectory sim-196 ulation for training can result in long convergence time and 197 large computational costs (Chen et al., 2018; Onken et al., 2021). To circumvent the above issue, in this work we adopt 199 a simulation-free learning strategy based on the OT theory 200 introduce in Sec. 2.2 (Liu et al., 2022; Tong et al., 2023; Albergo & Vanden-Eijnden, 2023), which turns out to be 202 faster and more scalable to large time series datasets. 203

By drawing mini-batch interpolated samples of the source distribution and target distribution in the Wasserstein space using Eq. (12), we can now model the variational velocity function via a neural network parameterized by  $\theta$ . Then, according to Eq. (1), the target velocity can be computed as the difference between the source distribution and target distribution. Therefore, the variational velocity function  $\mu_{\theta}(x_t^{\text{in}}, t)$  can be learned trough

$$\arg\min_{\theta} \int_{0}^{T} \int_{\mathbb{R}} \left\| \frac{\mathrm{d}X_{t}}{\mathrm{d}t} - \mu_{t}^{\theta}(x_{t}, t) \right\|_{2}^{2} \mathrm{d}x_{t} \mathrm{d}t \qquad (13)$$

215  
216 
$$\approx \arg\min_{\theta} \mathbb{E}\left[\left\|\frac{x_t^{\text{tar}} - x_0}{T} - \mu_t^{\theta}(x_t^{\text{in}}, t)\right\|_2^2\right].$$
 (14)

<sup>218</sup> 219 Since Eq. (14) can be solved by drawing mini-batch samples in the Wasserstein space and performing stochastic gradient descent accordingly, the learning process operates in a simulation-free manner.

Moreover, note that Eq. (13) also obeys the principle of least action introduced in Sec. 2.4 as it minimizes the kinetic energy described in Eq. (7). Therefore, this also indicates that the geodesic that drives the particles from the source distribution to the target distribution in the OT problem described in Sec. 2 is identified, which, as a result, allows us to generate new samples with fewer simulation steps compared to standard diffusion models.

#### 3.4. Gradient of Potential Function

So far, we have demonstrated how to leverage the kinetic energy to estimate the velocity in the Lagrangian described by Eq. (6). Apart from this, we can also incorporate the prior knowledge within the task-specific potential energy into the dynamics, which enables us to further improve the data generation performance. To this end, we let  $U_t(X_t) : \mathbb{R}^d \times$  $[0, T] \to \mathbb{R}$  be the task-specific potential function depending on the generalized coordinates  $X_t$  (Yang & Karniadakis, 2020; Onken et al., 2021; Neklyudov et al., 2023b), and the dynamics (here, we assume that the particle is solely driven by the drift) of the system  $v_t(X_t, t)$  yields

$$\frac{\mathrm{d}X_t}{\mathrm{d}t} = v_t(X_t, t) = -\nabla_x U_t(X_t). \tag{15}$$

Moreover, since the data generation problem in our case can also be interpreted as a stochastic optimal control (SOC) problem (Bellman, 1966; Fleming & Rishel, 2012; Nüsken & Richter, 2021; Zhang & Chen, 2022; Holdijk et al., 2023; Berner et al., 2024), the existence of such  $U_t(X_t)$  is guaranteed by Pontryagin's Maximum Principle (PMP) (Evans, 2024). Please refer to Appendix B for further details.

To estimate  $v_t(X_t, t)$ , according to the Lagrangian in Eq. (6), we assume that the potential function takes the form  $U_t(X_t) \approx -\log \mathcal{N}(X_t | \hat{X}_t, \sigma_p^2)$ , where  $\hat{X}_t$  the estimated mean and  $\sigma_p^2$  is the pre-defined variance. As a result, the corresponding derivative is  $\nabla_x U_t(X_t) = \frac{X_t - \hat{X}_t}{\sigma_p^2}$ . In terms of practical implementation, we parameterize  $\nabla_x U(X_t)$  via a time-dependent denoising autoencoder (TDAE). More specifically, we either pre-train or jointly train the TDAE on the intermediate time series data samples  $X_t$  generated by Eq. (12). The input is perturbed with noise, while the reconstruction target remains clean, to achieve the denoising effect. Afterwards, the reconstruction  $v_t^{\phi}(X_t, t)$  parametrized by  $\phi$  depending on the predicted  $X_t$ :

$$v_t^{\phi}(X_t, t) = -\frac{s}{\sigma_p^2} \big( X_t - \text{TDAE}(X_t) \big), \qquad (16)$$

where  $\text{TDAE}(X_t)$  represents the reconstruction of the TAVE model with input  $X_t$ ,  $s := s_0 t (T - t)/T$  with  $s_0$ 

**Conditional Lagrangian Wasserstein Flow** 



Figure 1: The overall training process of Conditional Lagrangian Wasserstein Flow.

being a positive scalar, and  $\sigma_p^2$  is treated as a positive constant for simplicity.

In this manner, we can incorporate the prior knowledge learned from the accessible training data into the sampling procedure established via Eq. (14) to enhance the data generation performance.

#### 3.5. Resampling Trick

Note that during inference, the model will generate the new data whose region encompasses both  $x_t^{\text{tar}}$  and  $x_t^{\text{cond}}$ , as the new input for next function evaluation iteration. However, from the problem defined in Sec. 3.1,  $x_t^{\text{cond}}$  can be computed accurately by interpolating the condition and the initial noise via Eq. (12). Therefore, we propose the following resampling trick to update the generated intermediate samples  $\hat{x}_t$  at time t by stitching the observed data region with the generated data region:

$$\hat{x}_t = \left(\frac{t}{T}x^{\text{obs}} + (1 - \frac{t}{T})x_0\right) \odot M^{\text{cond}} + x_t^{\text{gen}} \odot M^{\text{tar}},$$
(17)

where  $x_t^{\text{gen}}$  denotes the generated intermediate data samples.

The visualization of the proposed resampling trick can be found in Appendix D.

#### 3.6. The Algorithms

We now present the proposed training and inference algorithms. In the training procedure, we minimize the flowing matching loss to learn the variational velocity function  $\mu_t^{\theta}$ using the interpolation method. To estimate the variational drift function  $v_t^{\phi}$ , we can calculate the gradient of potential function using the TDVE model. The overall training process of our method is shown in Fig. 1.

In the inference procedure, we use the ODE sampler constructed by  $\mu_t^{\theta}$  to perform the path integral. Moreover, if Algorithm 1 Training procedure

**Require:** Terminal time: T, max epoch, observed data  $X^{obs}$ , parameters:  $\theta$  and  $\phi$ . **while** epoch < max epoch **do** sample t,  $(x_0, x_T)$  **if** OT **then** sample the mini-batch OT maps; **end if** sample  $x_t$  according to Eq. (12); minimize the loss function Eq. (14); **end while if** Rao-Blackwellization **then** train a TDAE model using  $X^{obs}$  and  $X_0$ . **end if** 

we want to further reduce the sampling variances, we can use the drift function  $v_t^{\phi}$  to formulate a new sampler. In addition, we can also choose to use the resampling trick to enhance the data generation performance.

Finally, the detailed training and inference procedures are summarized in Algorithms 1 and 2, respectively.

#### 3.7. Discussion

Here, we shed some light on the proposed method's connection to stochastic optimal control, path measures, and Rao-Blackwellization.

**Stochastic optimal control.** We first following the principle of least action in Lagrangian mechanics and optimal transport theory to compute the velocity function  $\mu_t$  by minimizing the corresponding kinetic energy. To further improve the data generation performance, we leverage the marginal log density-based potential function to construct the drift function  $v_t$ . According to the stochastic optimal control theory, it suggests that  $v_t$  in fact can act as the opti-

A	Algorithm 2 Sampling procedure
ŀ	<b>Require:</b> Step number: $N$ , step size: $h_L$ ,
	sample initial noise $x_0 \sim \mathcal{N}(0, \sigma_0^2)$ , conditional informa-
	tion $x^{\text{cond}}$ .
	while $t < N$ do
	$x_t^{\text{in}} = \text{Concatenate}(x^{\text{cond}}, \hat{x}_t)$
	$\hat{x}_{t+1} = \hat{x}_t + \mu_t^{\theta}(x_t^{\text{in}}, t) \frac{T}{N}$
	if Rao-Blackwellization then
	$\hat{x}_{t+1} = \hat{x}_{t+1} + v_t^{\phi}(x_{t+1}^{\text{pred}}, t) \frac{T}{N}$
	end if
	if Resampling then
	$\hat{x}_{i+1} = \left( \frac{(t+1)T}{T} x^{\text{cond}} + (T - \frac{(t+1)T}{T}) x_0 \right)$
	$x_{t+1} = \begin{pmatrix} N & x & + (1 & N) & x_0 \end{pmatrix} \bigcirc$
	$M^{\operatorname{cond}} + \hat{x}_{t+1} \odot M^{\operatorname{tar}}$
	end if
	t = t + 1
	end while

mal control signal if we consider the data generation process
as a controlled SDE. Moreover, the optimal control signal
can be attained by solving the corresponding HJB equation
using the Hopf-Cole transformation (Evans, 2022) and the
FB-SDE theory Anderson (1982); Song et al. (2020). Please
refer to Appendix B for the detailed discussion.

Path measures. We can now establish the path integral sam-301 pler by leveraging the Random-Nikon derivative between 302 the uncontrolled and controlled path measures obtained by 303 the Girsanov theorem (Liptser & Shiryaev, 2013), and ob-304 tain the corresponding KL divergence as well. Moreover, 305 we can also derive the ELBO for the marginal density p(t)306 by solving the associated Fokker-Planck equation using the 307 Feynman-Kac formula (Karatzas & Shreve, 2014). Further 308 details can be found in Appendices B.4 and B.5, respec-309 tively.

311**Rao-Blackwellization.** If we let the sampler constructed by312 $\mu_t$  be the based sampler and the sampler constructed by  $v_t$ 313the sufficient statistic, it can be seen that the new sampler314can, according to the Rao-Blackwell theorem (Casella &315Robert, 1996), improve the data generation performance by316reducing the sampling variances. The relevant theoretical317details can be found in Appendix C.

## 4. Experiments

In the section, we present the numerical results to demonstrate the effectiveness of our approach.

#### 4.1. Datasets

319

320

322

324

325

326

327

328

329

293

We use one synthetic dataset and three public multivariate time series datasets for validation.

1) Synthetic dataset was generated by the function: x =

 $t\sin(10t + 2\pi\epsilon)$ , where  $\epsilon \sim \mathcal{N}(0, \mathbb{I})$  and  $t \in [0, 1]$  with step size of 0.01. The batch size is 200, the total number of data points is 20,000, the missing rate of the raw data is 80%. 40%, 60%, and 80% of the datapoints are masked randomly as the imputation targets, denoted as Synthetic 0.4, Synthetic 0.6 and Synthetic 0.8, respectively.

**2) PM 2.5 dataset** (Zheng et al., 2013) was collected from the air quality monitoring sites for 12 months. The missing rate of the raw data is 13%. The feature number *K* is 36 and the sequence length *L* is 36. In our experiments, only the observed datapoints are masked randomly as the imputation targets.

**3) PhysioNet dataset** (Silva et al., 2012) was collected from the intensive care unit for 48 hours. The feature number K is 35 and the sequence length L is 48. The missing rate of the raw data is 80%. 10% and 50% of the datapoints are masked randomly as the imputation targets, denoted as PhysioNet 0.1 and PhysioNet 0.5, respectively.

4) ETTh1 dataset (Zhou et al., 2021) was collected from the electric power indicators for 2 years. The feature number K is 24 and the sequence length L is 96. 25%, 37.5%, and 50% of the datapoints are masked randomly as the imputation targets, denoted as ETTh1 0.25, ETTh1 0.375 and ETTh1 0.5, respectively.

#### 4.2. Baselines

For comparison, we select the following state-of-the-art timer series imputation methods as the baselines: 1) GP-VAE (Fortuin et al., 2020), which incorporates the Gaussian Process prior into a VAE model; 2) CSDI (Tashiro et al., 2021), which is based on the conditional diffusion model; 3) CSBI (Chen et al., 2023), which adopts the Schrödinger Bridge diffusion framework; 4) DSPD-GP (Biloš et al., 2023), which combines the diffusion model with the Gaussian Process prior; 5) DLinear (Zeng et al., 2023), which utilizes the moving average kernel for decomposition; 6) LightTS (Zhang et al., 2022), which captures the temporal patterns by continuous and interval sampling; 7) Etsformer (Woo et al., 2022), which proposes to use the exponential smoothing attention and the frequency attention; 8) Times-Net (Wu et al., 2023), which extracts the complex temporal information from the transformed 2D tensors.

#### 4.3. Experimental Settings

In terms of architecture choice, both the flow model and the TDAE model are built upon Transformers (Tashiro et al., 2021). We use the ODE sampler for inference and sample the exact optimal transport maps for interpolations to achieve the optimal performance. The optimizer is Adam and the learning rate: 0.001 with linear scheduler. The maximum training epochs is 200. The mini batch size for



Figure 2: Visualization of the test imputation results on the synthetic data, green dots are the conditions, blue dots are the imputation results, and red dots are the ground truth.

training is 64. The total step number of the Euler method 345 used in CLWF is 15, while the total step numbers for other diffusion models. i.e., CSDI, CSBI, and DSPD-GP are 50, 347 as suggested in their papers. The number of the Monte Carlo 348 samples for inference is 20. The standard deviation  $\sigma_0$  for 349 the initial noise  $X_0$  is 0.1, and the standard deviation  $\sigma_{\gamma}$ 350 for the injected noise  $\gamma_t$  is 0.001. The coefficient  $\sigma_p^2$  in the 351 gradient of the potential function is 0.01. 352

#### 4.4. Overall Imputation Results

341

342 343

353

354 355

357

358

359

360

361

362

363

364

367

368

369

370

374 375

376

378

379

380

381

Tables 1 and 2 show the overall test imputation results on 356 PM 2.5, PhysioNet, and Etth1, respectively. And the results demonstrate that CLWF achieves competitive performance compared with the state-of-the-art methods in terms of RMSE and MAE.

Note that CLWF requires less simulation steps (15) and sampled paths (20) to obtain high-quality data samples, which suggests that CLWF is faster and less computational expensive, compared to other existing diffusion-based time series imputation models, such as CSDI, CSBI, and DSPD.

Table 1: Test imputation results on PM 2.5, PhysioNet 0.1, and PhysioNet 0.5 (5-trial averages). The best are in bold and the second best are underlined.

Mathad	PM	2.5	PhysioN	Net 0.1	PhysioNet 0.5	
Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
GP-VAE	43.1	26.4	0.73	0.42	0.76	0.47
CSDI	19.3	9.86	0.57	0.24	0.65	0.32
CSBI	19.0	9.80	0.55	0.23	0.63	0.31
DSPD-GP	18.3	9.70	0.54	0.22	0.68	<u>0.30</u>
CLWF	18.1	9.70	0.47	0.22	0.64	0.29

#### 4.5. Ablation Study

To further demonstrate the effectiveness of our proposed 382 method, we conduct the following ablation study experi-383 ments. 384

Table 2: Test imputation results on ETT-h1(5-trial averages). The best are in bold and the second best are underlined.

Mathad	ETT-h	1 0.25	ETT-h1	0.375	ETT-h	n1 0.5
Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
DLinear	0.541	0.402	0.577	0.404	0.506	0.347
LightTS	0.469	0.347	0.544	0.382	0.463	0.318
Etsformer	0.411	0.304	0.514	0.364	0.424	0.292
TimesNet	<u>0.262</u>	0.178	0.289	0.196	0.319	0.215
CLWF	0.197	0.128	0.263	0.171	0.323	0.205

1) Single-path-sample Results. We compare the test imputation results of CLWF and CSDI by only using one path samples. From the results shown in Table 3, it can be seen that CLWF can achieve relatively good imputation results by only using one single Monte Carlo path integral sample, which indicates that CLWF has smaller sampling variances.

Table 3: Single-sample test imputation results on PM 2.5, PhysioNet 0.1, and PhysioNet 0.5 (5-trial averages).

Mathad	PM	2.5	Physiol	Net 0.1	Physiol	Net 0.5
Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
CSDI	22.2	11.7	0.74	0.30	0.83	0.40
CLWF	18.4	10.0	0.48	0.22	0.64	0.30

2) Numbers of Diffusion Steps. We compare the test imputation results of CLWF and CSDI using varying numbers of diffusion steps From the results shown in Table 4, we can see that CLWF has better imputation performance using less simulation steps for inference compared with CSDI, which implies that CLWF has faster convergence during inference.

3) Effect of Rao-Blackwellization. We compare the test imputation results of CLWF with and without using the Rao-Blackwellzation, referred to as Base and RB, respectively. From the results shown in Table 5, it can be seen that the new sampler can further improve the time series imputation performance of the base sampler by reducing

Table 4: Test imputation results on PhysioNet 0.1 with different simulation steps (5 trials).

Mathad	5 ste	eps	10 st	eps	15 st	teps	20 st	eps
Method	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE
CSDI	0.60	0.22	0.58	0.22	0.57	0.22	0.56	0.22
CLWF	0.48	0.22	0.47	0.22	0.47	0.22	0.48	0.22

Table 5: Ablation study imputation results on Rao-Blackwellization (5-trial averages).

Mathad	PM	2.5	Physio	Net 0.1	Physio	Net 0.5
Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
Base	18.27	9.76	0.4802	0.2221	0.6476	0.2991
RB	18.08	9.71	0.4785	0.2250	0.6466	0.3003
Mathad	ETT-h	1 0.25	ETT-h	1 0.375	ETT-I	h1 0.5
Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
Base	0.1999	0.1317	0.2191	0.1422	0.2906	0.1882
RB	0.1970	0.1266	0.2185	0.1424	0.2891	0.1845

the variances/RMSEs, which is also supported by the Rao-Blackwell theorem.

Finally, please refer to Appendix E for details on the hardware and software environments used in the experiments, and to Appendix F for additional experimental results.

### 5. Related Work

387

388

389

390

395

396

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

#### 5.1. Diffusion Models

Diffusion models, such as DDPMs (Ho et al., 2020) and 416 SBGM (Song et al., 2020), are considered as the new 417 contenders to GANs on data generation tasks. But they 418 generally take relatively long time to produce high quality 419 samples. To mitigate this problem, the flowing matching 420 methods have been proposed from an OT perspective. For 421 example, ENOT uses the saddle point reformulation of the 422 423 OT problem to develop a new diffusion model (Gushchin et al., 2024) The flowing matching methods have also been 424 proposed based on the OT theory (Lipman et al., 2022; Liu, 425 2022; Liu et al., 2023; Albergo & Vanden-Eijnden, 2023; 426 Albergo et al., 2023). In particular, mini-batch couplings 427 are proposed to straighten probability flows for fast infer-428 429 ence (Pooladian et al., 2023; Tong et al., 2024; 2023).

430 The Schrödinger Bridge framework have also been applied 431 to diffusion models for improving the data generation perfor-432 mance of diffusion models. Diffusion Schrödinger Bridge 433 utilizes the Iterative Proportional Fitting (IPF) method to 434 solve the SB problem (De Bortoli et al., 2021). SB-FBSDE 435 proposes to use forward-backward (FB) SDE theory to solve 436 the SB problem through likelihood training (Chen et al., 437 2022). GSBM formulates a generalized Schrödinger Bridge 438 matching framework by including the task-specific state 439

costs for various data generation tasks (Liu et al., 2024) NLSB chooses to model the potential function rather than the velocity function to solve the Lagrangian SB problem (Koshizuka & Sato, 2023). Action Matching (Neklyudov et al., 2023a;b) leverages the principle of least action in Lagrangian mechanics to implicitly model the velocity function for trajectory inference. Another classes of diffusion models have also been proposed from an stochastic optimal control perspective by solving the HJB-PDEs (Nüsken & Richter, 2021; Zhang & Chen, 2022; Berner et al., 2024; Liu et al., 2024; Park et al., 2024).

#### 5.2. Time Series Imputation

Many diffusion-based models have been recently proposed for time series imputation (Lin et al., 2023; Meijer & Chen, 2024). For instance, CSDI (Tashiro et al., 2021) combines a conditional DDPM with a Transformer model to impute time series data. CSBI (Chen et al., 2023) adopts the FB-SDE theory to train the conditional Schrödinger Bridge model to for probabilistic time series imputation. To model the dynamics of time series from irregular sampled data, DSPD-GP (Biloš et al., 2023) uses a Gaussian process as the noise generator. TDdiff (Kollovieh et al., 2024) utilizes self guidance and learned implicit probability density to improve the time series imputation performance of the diffusion models. However, the time series imputation methods mentioned above exhibit common issues, such as slow convergence, similar to many diffusion models. Therefore, in this work, we proposed CLWF to tackle thess challenges.

## 6. Conclusion

In this work, we proposed CLWF, a novel time series imputation method based on the optimal transport theory and Lagrangian mechanics. To generate the missing time series data, following the principle of least action, CLWF learns a velocity field by minimizing the kinetic energy to move the initial random noise to the target distribution. Moreover, we can also estimate the derivative of a potential function via a TDAE model trained on the observed training data to further improve the performance of the base sampler by Rao-Blackwellization. In contrast with previous diffusion-based models, the proposed requires less simulation steps and Monet Carlo samples to produce high-quality data, which leads to fast inference. For validation, CWLF is assessed on two public datasets and achieves competitive results compared with existing methods.

### **Impact Statement**

This paper presents work whose goal is to advance the field of time series imputation and deep learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## 440 References

441

442

443

444

445

446

447

448

449 450

451

452

453

454

455

456

457

458

459

465

468

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

- Albergo, M. S. and Vanden-Eijnden, E. Building normalizing flows with stochastic interpolants. In *The Eleventh International Conference on Learning Representations*, 2023.
- Albergo, M. S., Boffi, N. M., and Vanden-Eijnden, E. Stochastic interpolants: A unifying framework for flows and diffusions. *arXiv preprint arXiv:2303.08797*, 2023.
- Anderson, B. D. Reverse-time diffusion equation models. Stochastic Processes and their Applications, 12(3):313– 326, 1982.
- Arnol'd, V. I. *Mathematical methods of classical mechanics*, volume 60. Springer Science & Business Media, 2013.
- Bellman, R. Dynamic programming. *science*, 153(3731): 34–37, 1966.
- 460 Berner, J., Richter, L., and Ullrich, K. An optimal con461 trol perspective on diffusion-based generative model462 ing. *Transactions on Machine Learning Research*, 2024.
  463 ISSN 2835-8856. URL https://openreview.
  464 net/forum?id=oYIjw37pTP.
- Bertsekas, D. *Dynamic programming and optimal control: Volume I*, volume 4. Athena scientific, 2012.
- Biloš, M., Rasul, K., Schneider, A., Nevmyvaka, Y., and
  Günnemann, S. Modeling temporal data as continuous
  functions with stochastic process diffusion. In *Interna- tional Conference on Machine Learning*, pp. 2452–2470.
  PMLR, 2023.
  - Caluya, K. F. and Halder, A. Wasserstein proximal algorithms for the schrödinger bridge problem: Density control with nonlinear drift. *IEEE Transactions on Automatic Control*, 67(3):1163–1178, 2021.
  - Casella, G. and Robert, C. P. Rao-blackwellisation of sampling schemes. *Biometrika*, 83(1):81–94, 1996.
  - Chen, R. T., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
  - Chen, T., Liu, G.-H., and Theodorou, E. Likelihood training of schrödinger bridge using forward-backward sdes theory. In *International Conference on Learning Representations*, 2022, 2022.
- Chen, Y., Georgiou, T. T., and Pavon, M. Stochastic control liaisons: Richard sinkhorn meets gaspard monge on a schrodinger bridge. *Siam Review*, 63(2):249–313, 2021.

- Chen, Y., Deng, W., Fang, S., Li, F., Yang, N. T., Zhang, Y., Rasul, K., Zhe, S., Schneider, A., and Nevmyvaka, Y. Provably convergent schrödinger bridge with applications to probabilistic time series imputation. In *International Conference on Machine Learning*, pp. 4485–4513. PMLR, 2023.
- Cuturi, M. Sinkhorn distances: Lightspeed computation of optimal transport. *Advances in neural information processing systems*, 26, 2013.
- De Bortoli, V., Thornton, J., Heng, J., and Doucet, A. Diffusion schrödinger bridge with applications to score-based generative modeling. *Advances in Neural Information Processing Systems*, 34:17695–17709, 2021.
- Domingo-Enrich, C., Han, J., Amos, B., Bruna, J., and Chen, R. T. Stochastic optimal control matching. arXiv preprint arXiv:2312.02027, 2023.
- Evans, L. C. *Partial differential equations*, volume 19. American Mathematical Society, 2022.
- Evans, L. C. An introduction to mathematical optimal control theory spring, 2024 version. *Lecture notes available at https://math.berkeley.edu/ evans/control.course.pdf*, 2024.
- Feynman, R. P. The principle of least action in quantum mechanics. In *Feynman's thesis—a new approach to quantum theory*, pp. 1–69. World Scientific, 2005.
- Fleming, W. H. and Rishel, R. W. *Deterministic and stochastic optimal control*, volume 1. Springer Science & Business Media, 2012.
- Fortuin, V., Baranchuk, D., Rätsch, G., and Mandt, S. Gpvae: Deep probabilistic time series imputation. In *International conference on artificial intelligence and statistics*, pp. 1651–1661. PMLR, 2020.
- Gozzi, F. and Russo, F. Verification theorems for stochastic optimal control problems via a time dependent fukushima– dirichlet decomposition. *Stochastic Processes and their Applications*, 116(11):1530–1562, 2006.
- Gushchin, N., Kolesov, A., Korotin, A., Vetrov, D. P., and Burnaev, E. Entropic neural optimal transport via diffusion processes. *Advances in Neural Information Processing Systems*, 36, 2024.
- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.
- Holdijk, L., Du, Y., Hooft, F., Jaini, P., Ensing, B., and Welling, M. Stochastic optimal control for collective variable free sampling of molecular transition paths. *Ad*vances in Neural Information Processing Systems, 36, 2023.

- 495 Karatzas, I. and Shreve, S. Brownian motion and stochastic Neklyudov, K., Brekelmans, R., Severo, D., and Makhzani, 496 calculus, volume 113. springer, 2014. A. Action matching: Learning stochastic dynamics from 497 samples. In International Conference on Machine Learn-498 Kollovieh, M., Ansari, A. F., Bohlke-Schneider, M., ing, pp. 25858-25889. PMLR, 2023a. Zschiegner, J., Wang, H., and Wang, Y. B. Predict, refine, 499 Neklyudov, K., Brekelmans, R., Tong, A., Atanackovic, L., synthesize: Self-guiding diffusion models for probabilis-500 tic time series forecasting. Advances in Neural Informa-Liu, Q., and Makhzani, A. A computational framework 501 tion Processing Systems, 36, 2024. for solving wasserstein lagrangian flows. arXiv preprint 502 arXiv:2310.10649, 2023b. 503 Koshizuka, T. and Sato, I. Neural lagrangian schrödinger 504 Nüsken, N. and Richter, L. Solving high-dimensional bridge: Diffusion modeling for population dynamics. In 505 The Eleventh International Conference on Learning Rephamilton-jacobi-bellman pdes using neural networks: 506 resentations, 2023. perspectives from the theory of controlled diffusions and 507 measures on path space. Partial differential equations 508 Lehmann, E. L. and Casella, G. Theory of point estimation. and applications, 2(4):48, 2021. 509 Springer Science & Business Media, 2006. 510 Oksendal, B. Stochastic differential equations: an intro-511 Léonard, C. From the schrödinger problem to the mongeduction with applications. Springer Science & Business 512 kantorovich problem. Journal of Functional Analysis, Media, 2013. 513 262(4):1879-1920, 2012. 514 Onken, D., Fung, S. W., Li, X., and Ruthotto, L. Ot-flow: 515 Lin, L., Li, Z., Li, R., Li, X., and Gao, J. Diffusion models Fast and accurate continuous normalizing flows via opti-516 for time-series applications: a survey. Frontiers of Informal transport. Proceedings of the AAAI Conference on 517 mation Technology & Electronic Engineering, pp. 1–23, Artificial Intelligence, 35(10):9223-9232, 2021. 518 2023. 519 Park, B., Choi, J., Lim, S., and Lee, J. Stochastic optimal 520 Lipman, Y., Chen, R. T., Ben-Hamu, H., Nickel, M., and control for diffusion bridges in function spaces. arXiv 521 Le, M. Flow matching for generative modeling. In The preprint arXiv:2405.20630, 2024. 522 Eleventh International Conference on Learning Repre-523 Peyré, G., Cuturi, M., et al. Computational optimal transsentations, 2022. 524 port: With applications to data science. Foundations and 525 Liptser, R. S. and Shiryaev, A. N. Statistics of random Trends® in Machine Learning, 11(5-6):355-607, 2019. 526 processes: I. General theory, volume 5. Springer Science 527 Pooladian, A.-A., Ben-Hamu, H., Domingo-Enrich, C., & Business Media, 2013. Amos, B., Lipman, Y., and Chen, R. T. Multisample 528 529 Liu, G.-H., Lipman, Y., Nickel, M., Karrer, B., Theodorou, flow matching: Straightening flows with minibatch cou-E., and Chen, R. T. Generalized schrödinger bridge 530 plings. In International Conference on Machine Learning, matching. In The Twelfth International Conference on 531 pp. 28100–28127. PMLR, 2023. Learning Representations, 2024. 532 Risken, H. and Frank, T. The Fokker-Planck Equa-533 Liu, Q. Rectified flow: A marginal preserving approach tion: Methods of Solution and Applications, volume 18. 534 to optimal transport. arXiv preprint arXiv:2209.14577, Springer Science & Business Media, 2012. 535 2022. 536 Santambrogio, F. Optimal transport for applied mathemati-537 Liu, X., Gong, C., and Liu, Q. Flow straight and fast: cians. Birkäuser, NY, 55(58-63):94, 2015. 538 Learning to generate and transfer data with rectified flow. 539 Silva, I., Moody, G., Scott, D. J., Celi, L. A., and Mark, In The Eleventh International Conference on Learning 540 R. G. Predicting in-hospital mortality of icu patients: Representations, 2022. 541 The physionet/computing in cardiology challenge 2012. 542 Liu, X., Wu, L., Ye, M., and Liu, Q. Learning diffusion In 2012 Computing in Cardiology, pp. 245-248. IEEE, 543 bridges on constrained domains. In international confer-2012. 544 ence on learning representations (ICLR), 2023. 545 Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Er-546 Meijer, C. and Chen, L. Y. The rise of diffusion models in mon, S., and Poole, B. Score-based generative modeling 547 time-series forecasting. arXiv preprint arXiv:2401.03006, through stochastic differential equations. In International 548 2024. Conference on Learning Representations, 2020. 549 10

Tashiro, Y., Song, J., Song, Y., and Ermon, S. Csdi: Con-550 551 ditional score-based diffusion models for probabilistic 552 time series imputation. Advances in Neural Information 553 Processing Systems, 34:24804–24816, 2021. 554 2013. Tong, A., Malkin, N., Fatras, K., Atanackovic, L., Zhang, 555 Y., Huguet, G., Wolf, G., and Bengio, Y. Simulation-free 556 schrödinger bridges via score and flow matching. arXiv 557 preprint arXiv:2307.03672, 2023. 558 559 Tong, A., FATRAS, K., Malkin, N., Huguet, G., Zhang, Y., 560 Rector-Brooks, J., Wolf, G., and Bengio, Y. Improving 561 and generalizing flow-based generative models with mini-562 batch optimal transport. Transactions on Machine Learn-563 ing Research, 2024. ISSN 2835-8856. URL https:// 564 openreview.net/forum?id=CD9Snc73AW. Ex-565 pert Certification. 566 567 Villani, C. et al. Optimal transport: old and new, volume 568 338. Springer, 2009. 569 570 Woo, G., Liu, C., Sahoo, D., Kumar, A., and Hoi, S. Ets-571 former: Exponential smoothing transformers for time-572 series forecasting. arXiv preprint arXiv:2202.01381, 573 2022. 574 575 Wu, H., Hu, T., Liu, Y., Zhou, H., Wang, J., and Long, M. 576 Timesnet: Temporal 2d-variation modeling for general 577 time series analysis. In The Eleventh International Con-578 ference on Learning Representations. OpenReview.net, 579 2023. 580 581 Yang, L. and Karniadakis, G. E. Potential flow generator 582 with 12 optimal transport regularity for generative models. 583 IEEE Transactions on Neural Networks and Learning 584 Systems, 33(2):528-538, 2020. 585 Yong, J. and Zhou, X. Y. Stochastic controls: Hamiltonian 586 systems and HJB equations, volume 43. Springer Science 587 & Business Media, 2012. 588 589 Zeng, A., Chen, M., Zhang, L., and Xu, Q. Are transformers 590 effective for time series forecasting? In Proceedings of 591 the AAAI conference on artificial intelligence, volume 37, 592 pp. 11121–11128, 2023. 593 594 Zhang, Q. and Chen, Y. Path integral sampler: A stochastic 595 control approach for sampling. In The Tenth Interna-596 tional Conference on Learning Representations. OpenRe-597 view.net, 2022. URL https://openreview.net/ 598 forum?id=\_uCb2ynRu7Y. 599 600 Zhang, T., Zhang, Y., Cao, W., Bian, J., Yi, X., Zheng, S., 601 and Li, J. Less is more: Fast multivariate time series 602 forecasting with light sampling-oriented mlp structures. 603 arXiv preprint arXiv:2207.01186, 2022.

- Zheng, Y., Liu, F., and Hsieh, H.-P. U-air: When urban air quality inference meets big data. In *Proceedings* of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining, pp. 1436–1444, 2013.
- Zhou, H., Zhang, S., Peng, J., Zhang, S., Li, J., Xiong, H., and Zhang, W. Informer: Beyond efficient transformer for long sequence time-series forecasting. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 11106–11115, 2021.
- Zhou, X. Y., Yong, J., and Li, X. Stochastic verification theorems within the framework of viscosity solutions. *SIAM Journal on Control and Optimization*, 35(1):243– 253, 1997.

## A. Notations

Below are some mathematical notations used throughout the paper:

- $\nabla$  denotes the Jacobian operator;
- $\nabla$  · denotes the divergence operator;
- $\nabla^2$  denotes the Hessian operator;
- $\langle \cdot, \cdot \rangle := \text{trace}(\cdot^{\top}, \cdot)$  denotes the Frobenius inner product.

## **B. Stochastic Optimal Control**

In this section, we show how the stochastic optimal control theory is related to our data generation task.

#### **B.1.** Cost Function

An SDE controlled by the deterministic control function  $u \in \mathbb{U} \subset C(\mathbb{R}^d \times [0, T]; \mathbb{R})$ , where  $\mathbb{U}$  is a set of admissible controls, reads

$$dX_t^u = (a + \sigma u)(X_t^u, t)dt + \sigma(X_t^u, t)d\overline{W}_t,$$
(18)

where  $a \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$  is the drift/advection function,  $\sigma(t) \in C^1(\mathbb{R}^d \times [0,T]);\mathbb{R}^{d \times d}$  is the diffusion function, and  $\overline{W}_t$  is the standard Brownian motion.

Therefore, the data generation task can also be interpreted as a stochastic optimal control (SOC) problem (Bellman, 1966; Fleming & Rishel, 2012; Nüsken & Richter, 2021; Zhang & Chen, 2022; Holdijk et al., 2023; Koshizuka & Sato, 2023; Domingo-Enrich et al., 2023; Berner et al., 2024) whose cost functional  $\mathcal{J}$  is defined as:

$$\mathcal{J}(u; x_{\text{init}}, t) = \mathbb{E}\left[\int_{t}^{T} h(X_{s}^{u}, u, s) \mathrm{d}s + g(X_{T}^{u}) \middle| X_{t}^{u} = x_{\text{init}}\right],\tag{19}$$

where  $h(X_s^u, u, s) := f(X_s^u, s) + \frac{1}{2} ||u(X_s^u, s)||_2^2$ , where  $f \in C^1(\mathbb{R}^d \times [t, T]; [0, \infty))$ , is the instantaneous/running cost, and  $g \in C^1(\mathbb{R}^d; \mathbb{R})$  denotes the terminal cost. The above SOC problem can be solved via dynamic programming (Bellman, 1966; Bertsekas, 2012).

#### **B.2.** Pontryagin's Maximum Principle

Consider the following optimization problem derived from Eq. 19:

$$\arg\min_{u\in\mathcal{U}}\mathcal{J}(u;x_0,t) = \arg\min_{u\in\mathcal{U}}\left\{\int_0^T h(x,u,t)\mathrm{d}t + g(X_T)\right\}$$
(20)

subject to 
$$\begin{cases} \dot{x}(t) = v(x, u, t), & 0 < t \le T, \\ x(0) = x_0, \end{cases}$$
 (21)

where  $\dot{x}(t) := \frac{dx(t)}{dt}$  and v(x, u, t) denotes the system dynamics. Accordingly, the associated Lagrangian functional is defined as

$$\mathcal{L}(x,\lambda,u,t) := \int_0^T \left\{ h(x,u,t) + \lambda(t)(\dot{x}(t) - v(x,u,t)) \right\} \mathrm{d}t + g(X_T),$$
(22)

where  $\lambda(t)$  is the Lagrangian multiplier.

5 Now, let's consider a dynamic system defined by the following Hamiltonian

$$\mathcal{H}(x,\lambda,u) := \mathcal{K} + \mathcal{U}$$
  
=  $\lambda^{\top}(t)v(x,u,t) - h(x,u,t)$  (23)

where  $\mathcal{K}$  is the kinetic energy,  $\mathcal{U}$  is the potential energy, and  $\lambda(t)$  serves as the momentum/costate here. Then, the Lagrangian functional in Eq. 22 becomes

$$\mathcal{L}(x,\lambda,u,t) := \int_0^T \left\{ -\mathcal{H}(x,\lambda,u) + \lambda(t)\dot{x}(t) \right\} \mathrm{d}t + g(X_T),$$
(24)

We now apply variations of calculus to Eq. (24) by assuming  $(x^*(t), \lambda(t), u^*(t))$  is its minimizer and adding perturbation  $\delta x, \delta u, \delta \lambda$  with  $\delta x(0) = 0$ . Now we perform the first-order Taylor expansion:

$$\delta \mathcal{L} := \mathcal{L}(x^* + \delta x, \lambda + \delta \lambda, u^* + \delta u) - \mathcal{L}(x^*, \lambda, u^*)$$
$$\approx \int_0^T \left\{ -\mathcal{H}_x \delta x - \mathcal{H}_\lambda \delta \lambda - \mathcal{H}_u \delta u + \lambda \frac{\mathrm{d}}{\mathrm{d}t} \delta x + \delta \lambda \dot{x}^* \right\} \mathrm{d}t + g_x \delta(X_T).$$
(25)

Using integration by parts, we obtain

$$\delta \mathcal{L} \approx \int_0^T \left\{ (-\mathcal{H}_x - \dot{\lambda}) \delta x - \mathcal{H}_\lambda \delta \lambda + (-\mathcal{H}_u + \dot{x}^*) \delta \lambda + \frac{\mathrm{d}}{\mathrm{d}t} (\lambda \delta x) \right\} \mathrm{d}t + g_x \delta(X_T).$$
(26)

Let  $\delta \mathcal{L} = 0$ , we see that  $(x^*(t), \lambda(t), u^*(t))$  satisfies the following *necessary optimality conditions* for the Hamiltonian system with optimal control:

$x^*(0) = x_0,$	Initial value	(27)
$\dot{x}^*(t) = \mathcal{H}_{\lambda}(x^*(t), \lambda(t), u^*(t)),$	System dynamics	(28)
$\dot{\lambda}^*(t) = -\mathcal{H}_x(x^*(t),\lambda(t),u^*(t)),$	Adjoint/costate equation	(29)
$\lambda(T) = -g_x(x^*(T)),$	Adjoint terminal value	(30)
$\mathcal{H}_u(x^*(t),\lambda(t),u^*(t)) = 0,$	Extremal	(31)

where Eq. (28) and Eq. (29) are also known as Hamilton's canonical equations.

**Theorem B.1.** *Pontryagin's Maximum Principle (PMP) (Evans, 2024).* If the  $u^*$  is the optimal solution to the optimal control problem Eq. (19), then there exists a function  $\lambda$  solution of the costate/adjoint equation for which

$$u^* = \arg\max_{u \in \mathcal{U}} \mathcal{H}(x, \lambda, u), \ 0 \le t \le T.$$
(32)

This result implies that the Hamiltonian  $\mathcal{H}$  is maximized with respect to the optimal control  $u^*$  at each time t.

#### B.3. Hamilton-Jacobi-Bellman Equation

Here, we show how to derive the expression for the optimal control. Let  $p_t \in C^{2,1}(\mathbb{R}^d \times [0,T],\mathbb{R})$  be the density, then the controlled forward Fokker-Planck Kolmogorov (FPK) equation of the controlled SDE in Eq. (18) reads

$$\frac{\partial p_t}{\partial t} + \boldsymbol{\nabla} \cdot (p_t v) = \langle D(t), \nabla^2(p_t) \rangle, \tag{33}$$

where  $v(x,t) := (a + \sigma u)(x,t)$  represents the controlled system dynamics and  $D(t) := \frac{1}{2}\sigma^{\top}(t)\sigma(t)$ .

Considering the optimization objective Eq. (20) with respect to the new constraint Eq. (33), we define formulate the following Lagrangian

$$\mathcal{L}(p, x, \psi) = \int_{0}^{T} \int_{\mathbb{R}^{d}} \left\{ h(x, u, t) p_{t}(x) - \psi(x, t) \left( \underbrace{\frac{\partial p_{t}(x)}{\partial t}}_{\text{term (1)}} + \underbrace{\nabla \cdot (vp_{t})}_{\text{term (2)}} - \underbrace{\langle D(t), \nabla^{2}(p_{t}) \rangle}_{\text{term (3)}} \right) \right\} dxdt + \int_{\mathbb{R}^{d}} g(x) p_{T}(x) dx,$$
(34)

where  $\psi(x, t)$  serves as the Lagrangian multiplier.

We now apply integration by parts to term (1) with respect to t and integration by parts to term (2) with respect to x, respectively. Then, the two-fold integration by parts are performed in term (3), and we have

$$\int_{\mathbb{R}^d} \langle D(t), \nabla^2(p_t) \rangle \psi dx = \int_{\mathbb{R}^d} \langle D(t), \nabla^2(\psi) \rangle p dx,$$
(35)

where we assume the functions have compact support such that their respective products terms vanish at both ends. Then Eq. (34) becomes

$$\mathcal{L}(p, x, \psi) = \int_0^T \int_{\mathbb{R}^d} p_t(x) \Big\{ h(x, u, t) + \frac{\partial \psi}{\partial t} + \langle \nabla \psi, v(x, t) \rangle + \langle D(t), \nabla^2(\psi) \rangle \Big\} \mathrm{d}x \mathrm{d}t \\ - \int_{\mathbb{R}^d} \psi(x, T) p_T(x) \mathrm{d}x + \int_{\mathbb{R}^d} \psi(x, 0) p_0(x) \mathrm{d}x.$$
(36)

Let the density  $p_t$  be fixed, then according to the verification theorem (Zhou et al., 1997; Gozzi & Russo, 2006), we can easily attain the optimal control that minimizes Eq. 36 with respect to u:

$$u^* = -\sigma^{\top}(t)\nabla_x\psi(x,t). \tag{37}$$

Plug Eq. (37) into Eq. (36) and let the new integral equal 0, we have

$$\mathcal{L}(p,x,\psi) = \int_0^T \int_{\mathbb{R}^d} p(x,t) \Big\{ \frac{\partial \psi}{\partial t} + \frac{1}{2} \big\| \sigma^\top(t) \nabla_x \psi(x,t) \big\|_2^2 + \langle \nabla \psi, v(x,t) \rangle + \langle D(t), \nabla^2(\psi) \rangle \Big\} \mathrm{d}x \mathrm{d}t - \int_{\mathbb{R}^d} \psi(x,T) p_T(x) \mathrm{d}x + \int_{\mathbb{R}^d} \psi(x,0) p_0(x) \mathrm{d}x.$$
(38)

And the associated minimizer  $(x^*(t), \lambda(t), u^*(t))$  satisfies the following optimality conditions:

$$\frac{\partial}{\partial p}\mathcal{L}(p,x,\psi) = 0 \tag{39}$$

$$\frac{\partial}{\partial p_T} \mathcal{L}(p, x, \psi) = 0.$$
(40)

As a result, we obtain the following partial differential equation (PDE) whose solution is the potential function  $\psi$ :

$$\frac{\partial\psi}{\partial t} + \frac{1}{2} \left\| \sigma^{\top}(t) \nabla\psi \right\|_{2}^{2} + \langle \nabla\psi, v(x,t) \rangle = -\langle D(t), \nabla^{2}(\psi) \rangle, \tag{41}$$

with the terminal condition: 
$$\psi(x,T) = g(x)$$
, (42)

where  $\frac{1}{2} \| \sigma^{\mathsf{T}}(t) \nabla \psi \|_2^2 + \langle \nabla \psi, v(x,t) \rangle$  is the Hamiltonian with  $\nabla \psi$  being the momentum. And Eq. (41) is the cerebrated Hamilton-Jacobi-Bellman (HJB) equation with the value function  $\psi = \inf_u \mathcal{J}$  being the unique viscosity solution (Zhou et al., 1997; Gozzi & Russo, 2006; Yong & Zhou, 2012; Evans, 2022).

Further, Eq. (41) can be linearized by using the Hopf-Cole transformation (Evans, 2022). To this end, we let  $\psi(x,t) = \log \tilde{p}(x)$  to have:

$$\widetilde{p} = \exp(\psi) \tag{43}$$

$$\nabla \widetilde{p} = \exp(\psi) \nabla \psi \tag{44}$$

759 We also have 

$$\langle D(t), \nabla^{2}(\tilde{p}) \rangle = \sum_{i,j=1}^{d} (D(t))_{i,j} \frac{\partial^{2}}{\partial_{x_{i}} \partial_{x_{j}}} \exp(\psi)$$

$$= \exp(\psi) \Biggl\{ \sum_{i,j=1}^{d} (D(t))_{i,j} \Biggl( \frac{\partial^{2}\psi}{\partial_{x_{i}} \partial_{x_{j}}} + \frac{\partial\psi}{\partial_{x_{i}}} \frac{\partial\psi}{\partial_{x_{j}}} \Biggr) \Biggr\}$$

$$= \exp(\psi) \Biggl\{ \langle D(t), \nabla^{2}(\psi) \rangle + \frac{1}{2} \left\| \sigma^{\top}(t) \nabla\psi \right\|_{2}^{2} \Biggr\}.$$

$$(45)$$

Combining Eq. (41), Eq. (44), and Eq. (45), we have

$$\frac{\partial \widetilde{p}}{\partial t} = \exp(\psi) \frac{\partial \psi}{\partial t} 
= -\exp(\psi) \left\{ \frac{1}{2} \left\| \sigma^{\top}(t) \nabla_x \psi \right\|_2^2 + \langle \nabla \psi, v(x,t) \rangle + \langle D(t), \nabla^2(\psi) \rangle \right\} 
= -\langle D(t), \nabla^2(\widetilde{p}) \rangle - \langle \nabla \widetilde{p}, v(x,t) \rangle,$$
(46)

which in fact is the backward Fokker-Planck Kolmogorov equation, which suggests that  $p_t$  is a reverse-time density. Then, the optimal control signal amounts to

$$u^* = -\sigma^\top(t)\nabla\log\widetilde{p}(x). \tag{47}$$

Consequently, according to Eq. 18, we have the following controlled reverse-time SDE:

$$dX_t = [f(X_t, t) - \sigma^2(t)\nabla\log\widetilde{p}_t(X, t)]dt + \sigma(t)d\overline{W}_t.$$
(48)

Recall the coupled forward-backward SDE system (Anderson, 1982; Song et al., 2020), then the above reverse-time SDE has the following forward-time counterpart:

$$dX_t = f(X_t, t)dt + \sigma(t)dW_t.$$
(49)

Further, we consider an overdamped Langevin dynamics system by letting  $f(X_t, t) := -\nabla_x \log \hat{p}_t(x)$ , where  $\hat{p}_t$  is the forward-time density. Consequently, this enables us to control the sampling process in the forward (noise-to-data) sampling process.

#### B.4. Path Sampling via Stochastic Optimal Control

Let  $\mathbb{P} \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$  be the base path measure and  $\mathbb{P}^u \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$  the associated path measure rendered by the optimal control  $u \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$ . We have the following Radon-Nikodym derivative attained by Girsanov theorem (Liptser & Shiryaev, 2013):

$$\frac{\mathbb{I}\mathbb{P}^u}{\mathrm{d}\mathbb{P}} = \exp\left\{\int_0^T u^\top(x,t)\mathrm{d}W_t - \int_0^T \frac{1}{2} \|u(x,t)\|_2^2 \mathrm{d}t\right\}$$
(50)

$$\frac{\mathrm{d}\mathbb{P}}{\mathrm{d}\mathbb{P}^u} = \exp\left\{-\int_0^T u^\top(x,t)\mathrm{d}W_t - \int_0^T \frac{1}{2}\|u(x,t)\|_2^2 \mathrm{d}t\right\},\tag{51}$$

where u satisfies the Novikov's condition:  $\mathbb{E}\left[\exp\left(\frac{1}{2}\int_0^T u^2 dt\right)\right] < \infty$ .

Noting that  $dX_t = vdt + \sigma_t dW_t$ ,  $dX_t^u = vdt + \sigma_t (dW_t + udt)$ , and  $\mathbb{E}\left[\int_0^T u^\top(x, t)dW_t\right] = 0$ , then the KL divergence between the two path measures amounts to

$$D_{\mathrm{KL}}(\mathbb{P}^{u}\|\mathbb{P}) = -\mathbb{E}_{\mathbb{P}}\left[\log\frac{\mathrm{d}\mathbb{P}^{u}}{\mathrm{d}\mathbb{P}}(X)\right] = \mathbb{E}_{\mathbb{P}}\left[\int_{0}^{T}\frac{1}{2}\|u(X,t)\|_{2}^{2}\mathrm{d}t\Big|X_{0}=x\right].$$
(52)

This result suggests that finding the optimal control enables us sample the target distribution. Furthermore, to reduce to the sampling variances, recalling the cost functional defined in Eq. 19, one can adopt the importance sampling scheme through sampling N paths from the path measure  $P^u$  and compute the average given by

$$\frac{1}{N}\sum_{i=1}^{N}\mathcal{J}(X^{i,u})\mathcal{W}(X^{i,u}),\tag{53}$$

where the importance weights  $\mathcal{W}(X^{i,u})$  given by Eq. (51) are

$$\mathcal{W}(X^{i,u}) = \exp\left\{-\int_0^T u^\top(x,t) \mathrm{d}W_t - \int_0^T \frac{1}{2} \|u(x,t)\|_2^2 \mathrm{d}t\right\}.$$
(54)

#### 825 B.5. Feynman-Kac Formula

The Feynman-Kac Formula is very powerful tool to solve parabolic PDEs.

Theorem B.2. (Feynman-Kac formula (Karatzas & Shreve, 2014)). Let  $x \in \mathbb{R}^d$  be the spatial variable and  $t \in [0,T]$  the temporal variable.

$$\frac{\partial}{\partial t}\rho(x,t) + \bar{\mu}(x,t)\frac{\partial}{\partial x}\rho(x,t) + \frac{1}{2}\bar{\sigma}^2(x,t)\frac{\partial^2}{\partial x^2}\rho(x,t) - A(x,t)\rho(x,t) + c(x,t) = 0,$$
(55)

subject to the terminal condition: 
$$\rho(x,T) = \psi(x)$$
, (56)

where  $\rho \in C^{2,1}(\mathbb{R}^d \times [0,T];\mathbb{R})$ ,  $\bar{\mu} \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$ ,  $\bar{\sigma} \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^{d\times d})$ ,  $c \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$ ,  $A \in C^1(\mathbb{R}^d \times [0,T];\mathbb{R}^d)$ , and  $\psi \in C^1(\mathbb{R}^d;\mathbb{R}^d)$  are known functions. Let  $W_t$  is a Brownian motion under path measure P and X solves the following SDE:

$$dX_t = \widetilde{\mu}(X, t) + \widetilde{\sigma}(X, t)dW_t.$$
(57)

840 Then  $\rho(x,t)$  can be represented by the Feynman-Kac formula as follow:

$$\rho(x,t) = \mathbb{E}_P \left[ \exp\left\{ -\int_t^T A(X_s,s) \mathrm{d}s \right\} \psi(X_T) + \int_t^T \exp\left\{ -\int_t^\tau A(X_s,s) \mathrm{d}s \right\} c(X_\tau,\tau) \mathrm{d}\tau \left| X_t = x \right]$$
(58)

Now we use the Feynman-Kac formula to compute the marginal distribution. To this end, we first rewrite the forward Fokker-Planck equation as follows:

$$\frac{\partial}{\partial t}p(x,t) = -\boldsymbol{\nabla} \cdot (\mu p) - \langle D(t), \nabla^2(p) \rangle$$
  
=  $-(\boldsymbol{\nabla} \cdot \mu)p - \mu \nabla p - \langle D(t), \nabla^2(p) \rangle,$  (59)

and let its coefficients match their counterparts in Eq. (55) and Eq. (56) as follows:

$$p \longrightarrow \rho$$
 (60a)

$$\mu \longrightarrow \widetilde{\mu}$$
 (60b)

$$\sigma \longrightarrow \widetilde{\sigma} \tag{60c}$$

$$\langle D(t), \nabla^2(p) \rangle \longrightarrow \frac{1}{2} \sigma^2 \frac{\partial^2}{\partial x^2} \rho$$
 (60d)

$$-\boldsymbol{\nabla} \boldsymbol{\cdot} \boldsymbol{\mu} \longrightarrow \boldsymbol{A} \tag{60e}$$

$$0 \longrightarrow c$$
 (60f)

$$g(x) \longrightarrow \psi(x).$$
 (60g)

Therefore, according to Eq. (58), we obtain the following expressions for the marginal distribution:

$$p(x,t) = \mathbb{E}_{\mathbb{P}}\left[\exp\left\{\int_{t}^{T} \nabla \cdot \mu(X_{s},s) \mathrm{d}s\right\} g(X_{T}) \Big| X_{t} = x\right]$$
(61)

$$p(x,t) = \mathbb{E}_{\mathbb{P}}\left[\exp\left\{-\int_{0}^{t} \nabla \cdot \mu(X_{s},s) \mathrm{d}s\right\}g(X_{0}) \Big| X_{0} = x\right].$$
(62)

Combining Eq. (61) with Eq. 52, we use Jensen's inequality to obtain the following ELBO:

$$\log p(x,t) \ge \mathbb{E}_{\mathbb{P}} \left[ -\int_{0}^{t} \left\{ \nabla \cdot \mu(X_{s},s) \mathrm{d}s + \log g(X_{0}) \middle| X_{0} = x \right] - \mathbb{E}_{\mathbb{P}} \left[ \log \frac{\mathrm{d}\mathbb{P}^{u}}{\mathrm{d}\mathbb{P}}(X) \right]$$

$$\approx \mathbb{E} \left[ -\int_{0}^{t} \left\{ \nabla \cdot \mu(X_{s},s) \mathrm{d}s + \log g(X_{0}) \middle| X_{0} = x \right] - \mathbb{E}_{\mathbb{P}} \left[ \log \frac{\mathrm{d}\mathbb{P}^{u}}{\mathrm{d}\mathbb{P}}(X) \right] \right]$$
(62)

$$\geq \mathbb{E}_{\mathbb{P}} \left[ -\int_{0} \left\{ \mathbf{V} \cdot \mu(X_{s}, s) + \frac{1}{2} \| u(X, s) \|_{2}^{2} \right\} \mathrm{d}s + \log g(X_{0}) | X_{0} = x \right]$$
(63)

#### **B.6. Connection to Flow Matching**

Since we have the intermediate sample  $X_t = \frac{t}{T}X_T + (1 - \frac{t}{T})X_0$ , we can directly computed the predicted terminal samples at time T using the predicted  $\hat{X}_t$  at time t without iterative function evaluations:

$$\hat{X}_T \approx \left(\hat{X}_t - (1 - \frac{t}{T})X_0\right) / (\frac{t}{T}).$$
 (64)

Assume the log-densities of  $X_0$ ,  $X_t$ , and  $X_T$  can be represented by the same function, then the terminal cost in the value function Eq. (20) is defined as  $g(X_T) := -\log p(X_T)$ . As a result, it suggests that minimizing the running cost at time talso means minimizing the terminal cost at time T.

First, following the principle of least action in Lagrangian dynamics, the uncontrolled system dynamics  $v(x_t, t)$  is learned by minimizing the associated kinetic energy.

The uncontrolled sampling process can be described as follow:

$$dx_t = v(x_t, t)dt \tag{65}$$

Since we formulate the potential energy function  $U(x_t, t) = \log p(x_t, t) = \log \mathcal{N}(x_t; \tilde{x}_t, \tilde{\sigma}_t)$ , where  $\tilde{x}_t = \text{TDAE}(x_t, t)$ attained by the reconstruction process of the TDAE model. Accordingly, the controlled sampling process can be cast as:

$$dx_t = u(x_t, t)dt = -\sigma^{\top}(t)\nabla_x \log p(x_t, t)dt = -\frac{x_t - x_t}{\tilde{\sigma}_t^2}$$
(66)

901 The corresponding sampling scheme is 902

$$x_{t+1} = x_t + v(x_t, t)\mathrm{d}t \tag{67}$$

$$x_{t+1}^u = x_{t+1} + u(x_{t+1}, t+1)d\tilde{t}.$$
(68)

## C. Rao-Blackwellization

Here we prove that the proposed sampler is, in fact, a Rao-Blackwellized trajectory sampler (Casella & Robert, 1996). We first start with the following definition:

**Definition C.1** (Sufficient statistic). A *sufficient statistic*  $\mathcal{T}$  for a parameter  $\Theta$  captures all the necessary information 911 contained in the data sample  $\mathcal{X}$  to estimate  $\Theta$ . Once  $\mathcal{T}$  is known,  $\mathcal{X}$  does not provide additional information to estimate  $\Theta$ .

912913 To determine whether a statistic is sufficient, we can apply the following theorem.

**Theorem C.2.** (*Fisher-Neyma theorem (Lehmann & Casella, 2006)*). Let probability density function of  $\mathcal{X}$  be  $p(x|\varphi)$ , 915 then the statistics  $\mathcal{T}$  are sufficient for  $\mathcal{X}$  iff  $p(x|\varphi)$  are be written in the following form:

$$p(x|\varphi) = \mathcal{F}(x)\mathcal{G}(\mathcal{T}(x);\varphi), \tag{69}$$

918 where  $\mathcal{F}(x)$  is a distribution independent of  $\theta$  and  $g(\cdot, \theta)$  captures all the dependence on  $\theta$  via sufficient statistics  $\mathcal{T}(x)$ .

Following the above theorem, in our context, we assume that the marginal distribution of  $X_t$  is the Gaussian with unknown mean and known variance:  $\mathcal{N}(X_t; m_{\varphi}(X_{t-1}), \sigma^2(t))$ . Then the joint distribution of N samples can be written and decomposed as follows:

$$p(X_{t}^{1}, X_{t}^{2}, \dots, X_{t}^{N} | \varphi) = (2\pi)^{-N/2} \sigma^{2} \exp\left(\frac{-1}{2\sigma^{2}} \sum_{i=1}^{N} (X_{t}^{i} - m_{\varphi})^{2}\right)$$

$$= (2\pi)^{-N/2} \sigma^{2} \exp\left(\frac{-1}{2\sigma^{2}} \sum_{i=1}^{N} X_{t}^{i} + \frac{m_{\varphi}}{\sigma^{2}} \sum_{i=1}^{N} X_{t}^{i} - \frac{Nm_{\varphi}^{2}}{2\sigma^{2}}\right)$$

$$= \underbrace{(2\pi)^{-N/2} \sigma^{2} \exp\left(\frac{-1}{2\sigma^{2}} \sum_{i=1}^{N} X_{t}^{i}\right)}_{\mathcal{F}(x)} \underbrace{\exp\left(\frac{m_{\varphi}}{\sigma^{2}} \sum_{i=1}^{N} X_{t}^{i} - \frac{Nm_{\varphi}^{2}}{2\sigma^{2}}\right)}_{\mathcal{G}(\mathcal{T}(x);\varphi)}.$$
(70)

The above result suggests that the trajectory sampler can be formulated as a sufficient statistic for  $\varphi$ . Consequently, for our task, we have: 1) the parameter to estimate  $\Theta := X_t$ , where  $X_t$  is the intermediate sample predicted (or its mean) at time t; 2) the data sample  $\mathcal{X} := X_0$ , where  $X_0$  is the initial noise (as well as the observation); 3) the base unbiased sampler representing the system dynamics estimated according to the principle of least action,  $S(X_0; \mu, t) := X_0 + \int_0^t \mu(X_s, s) ds =$  $X_{t-1} + \int_{t-1}^{t} \mu(X_s, s) ds = X_t$ , where  $0 < t \le T$ ; 4) the sufficient statistic representing the optimal control signal according to Pontryagin's Maximum Principle,  $\mathcal{T}(X_0; u, t) := X_0 + \int_0^{t-1} u^*(X_s, s) ds = X_{t-1}$ , where  $0 < t \le T$ . As a result, we have the new sampler  $\mathcal{S}^* := \mathbb{E}[\mathcal{S}|\mathcal{T}] = \mathbb{E}[\mathcal{T}(X_0; u, t) + \int_{t-1}^t \mu(\mathcal{T}(X_0; u, s), s) ds]$ . Then, according to the Rao-Blackwell theorem: 

**Theorem C.3.** (*Rao-Blackwell theorem* (*Casella & Robert, 1996*)). Let *S* be an unbiased estimator of some parameter  $\Theta$ , 946 and  $\mathcal{T}(\mathcal{X})$  the sufficient statistic for  $\Theta$ , then: 1)  $S^* = \mathbb{E}[S|\mathcal{T}(\mathcal{X})]$ , is an unbiased estimator for  $\Theta$ , and 2)  $\mathbb{V}[S] \ge \mathbb{V}[S^*]$ . 947 The inequality is strict unless *S* is a function of  $\mathcal{T}$ .

**Proof:** In the ODE/SDE sampling process, we have  $p(X_t|X_{t-1}, X_0) = p(X_t|X_{t-1})$ , i.e.,  $p(\Theta|\mathcal{T}, \mathcal{X}) = p(\Theta|\mathcal{T})$ . Since  $\mathcal{T}$  is a statistic of  $\mathcal{X}$  and  $\mathcal{S}$  is an estimator of  $\Theta$ , we have  $\mathbb{E}[\mathcal{S}|\mathcal{T}] = \mathbb{E}[\mathcal{S}|\mathcal{T}, \Theta]$ . We now apply the law of total expectation (*Z* and *Y* are two random variables):

$$\mathbb{E}[Z|Y] = \int zp(z|Y)dz \Longrightarrow \mathbb{E}[\mathbb{E}[Z|Y]] = \iint zp(z|y)dzp(y)dy$$
$$= \iint zp(z|y)p(y)dzdy$$
$$= \iint zp(z,y)dzdy = \int zp(z)dz = \mathbb{E}[Z],$$
(71)

to attain the following relationships:

$$\mathbb{E}[\mathcal{S}^*|\Theta] = \mathbb{E}[\mathbb{E}[\mathcal{S}|\mathcal{T}]|\Theta] = \mathbb{E}[\mathbb{E}[\mathcal{S}|\mathcal{T},\Theta]|\Theta] = \mathbb{E}[\mathcal{S}|\Theta].$$
(72)

973 Then we apply the law of total variance to attain the following relationships: 

$$\mathbb{V}[\mathcal{S}|\Theta] = \mathbb{E}[\mathbb{V}[\mathcal{S}|\mathcal{T},\Theta]|\Theta] + \mathbb{V}[\mathbb{E}[\mathcal{S}|\mathcal{T},\Theta]|\Theta]$$
$$= \mathbb{E}[\mathbb{V}[\mathcal{S}|\mathcal{T},\Theta]|\Theta] + \mathbb{V}[\mathcal{S}^*|\Theta],$$
(73)

983 where  $\mathbb{E}[\mathbb{V}[\mathcal{S}|\mathcal{T},\Theta]|\Theta] \ge 0$ , therefore  $\mathbb{V}[\mathcal{S}|\Theta] \ge \mathbb{V}[\mathcal{S}^*|\Theta]$ .

984 The results in Eq. (72) and Eq. (73) suggest that the new sampler has the same expectation as the base sampler but with 985 smaller variance (mean squared error), which is also verified by the experimental results.

#### 987 C.1. Overall Theoretical Framework

Fig. 3 visualizes the overall theoretical framework of the proposed method in this paper.

**Conditional Lagrangian Wasserstein Flow** 



Gold-6248R-3.00GHz CPU. For the software environment, the Python version is 3.9.7, the CUDA version 11.7, and the

Pytorch version is 2.0.1.

## 1045 1046

# 1047

1049

1050

1058

1066 1067

1068

Table 6: Test imputation results on PM 2.5 with different simulation steps (5 trials).

)51	Mada a	5 steps		10 steps		15 steps		20 steps	
)52	Method	RMSE	MAE	RMSE	MAE	RMSE	MAE	RMSE	MAE
)53	CSDI	$34.21 \pm 0.16$	$14.85\pm0.01$	$29.43 \pm 0.46$	$12.48\pm0.08$	$22.40 \pm 0.16$	$10.78\pm0.04$	$19.22\pm0.13$	$9.91 \pm 0.02$
54	CLWF	$18.29\pm0.002$	$9.78\pm0.004$	$18.28\pm0.003$	$9.77\pm0.005$	$18.26\pm0.006$	$9.76\pm0.004$	$18.21\pm0.002$	$9.72\pm0.004$
)55									

## F.2. Ablation study on Rao-Blackwellization.

F. Additional Experimental Results

F.1. Ablation study on simulation steps.

Table 7: Test imputation results on PM 2.5, PhysioNet 0.1, and PhysioNet 0.5 (5 trials).

1060							
1061	Mathad	PM	2.5	PhysioNet 0.1		PhysioNet 0.5	
1062	Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
1063 1064	CLWF (no RB)	$18.27 \pm 0.01$ 18.08 ± 0.02	$9.76 \pm 0.01$ 9 71 + 0 00	$0.4802 \pm 1e-4$	$0.2221 \pm 0e-4$	$0.6476 \pm 0e-4$	$0.2991 \pm 0e-4$
1065		10.00 ± 0.02	9.11 ± 0.00	0.4100 ± 10-4	$0.2200 \pm 10^{-4}$	0.0400 ± 00-4	$0.0000 \pm 0.004$

Table 8: Test imputation results on synthetic data (5-trials, values are multiplied by  $10^2$ ).

070	Mathad	Synthe	etic 0.4	Synthe	etic 0.6	Synthetic 0.8	
071	Method	RMSE	MAE	RMSE	MAE	RMSE	MAE
)72	CLWF (no RB)	$22.91 \pm 0.49$	$15.28\pm0.21$	$25.65\pm0.31$	$15.54\pm0.22$	$27.41 \pm 0.27$	$15.91\pm0.23$
073 074	CLWF (with RB)	$22.72\pm0.48$	$13.23\pm0.42$	$25.44\pm0.30$	$15.28\pm0.17$	$27.32\pm0.27$	$15.79\pm0.23$

## 1075

1078

1079

1076 F.3. Ablation study on Resampling. 1077

Table 9: Test imputation results on ETT-h1(5-trial averages). The best are in **bold** and the second best are underlined.

0	Mathad	ETT-h1 0.25		ETT-h	1 0.375	ETT-h1 0.5		
)	Wiethou	RMSE	MAE	RMSE	MAE	RMSE	MAE	
	Base	$0.1999 \pm 8e-4$	$0.1317\pm3\mathrm{e}{ ext{-}4 ext{-}$	$0.2191 \pm 2e-4$	$0.1422 \pm 4\text{e-}4$	$0.2906 \pm 1e-3$	$0.1882 \pm 2e-4$	
	RB	$0.1970 \pm 6e-4$	$0.1266 \pm 2\text{e-}4$	$0.2185 \pm 4e-4$	$0.1424\pm2\text{e-}4$	$0.2891 \pm 2e-4$	$0.1845 \pm 2e-4$	
	Resampling	$0.1976 \pm 8e-4$	$0.1257 \pm 3e-4$	$\underline{0.2165 \pm 1\text{e-}3}$	$\underline{0.1366 \pm 1\text{e-}4}$	$0.2964 \pm 1e$ -3	$0.1988 \pm 4\text{e-}4$	
	Resampling + RB	$0.1968 \pm 8e-4$	$0.1253 \pm 3e-4$	$0.2157 \pm 6e-4$	$0.1363 \pm 3e-4$	$0.2921 \pm 9e-4$	$0.1926\pm3\text{e-}4$	

1087

#### 1088 F.4. Time Efficiency 1089

1090 We report the statistics regarding the time efficiency of our method here.

Base

3.26

s/iteration

1	0	9	1
1	0	9	2

1093

1094

1095

1096 1097

1098

1099

Table 10: Inference time costs on ETT-h1 0.5.

Resampling

3.28

Resampling + RB

6.38

RB

6.36