

---

# Modeling Continuous Stochastic Processes with Dynamic Normalizing Flows

---

Ruizhi Deng<sup>1,2\*</sup> Bo Chang<sup>1</sup> Marcus A. Brubaker<sup>1,3,4</sup> Greg Mori<sup>1,2</sup> Andreas M. Lehrmann<sup>1</sup>  
<sup>1</sup>Borealis AI <sup>2</sup>Simon Fraser University <sup>3</sup>York University <sup>4</sup>Vector Institute

## Abstract

Normalizing flows transform a simple base distribution into a complex target distribution and have proved to be powerful models for data generation and density estimation. In this work, we propose a novel type of normalizing flow driven by a differential deformation of the Wiener process. As a result, we obtain a rich time series model whose observable process inherits many of the appealing properties of its base process, such as efficient computation of likelihoods and marginals. Furthermore, our continuous treatment provides a natural framework for irregular time series with an independent arrival process, including straightforward interpolation. We illustrate the desirable properties of the proposed model on popular stochastic processes and demonstrate its superior flexibility to variational RNN and latent ODE baselines in a series of experiments on synthetic and real-world data.

## 1 Introduction

Expressive models for sequential data form the statistical basis for downstream tasks in a wide range of domains, including computer vision, robotics, and finance. Recent advances in deep generative architectures, especially the concept of reversibility, have led to tremendous progress in this area and created a new perspective on many of the long-standing limitations that are typical in traditional approaches based on structured decompositions (e.g., state-space models).

We argue that the power of a time series model depends on its properties in the following areas: (1 – Resolution) Common time series models are discrete with respect to time. As a result, they make the implicit assumption of a uniformly spaced temporal grid, which precludes their application from asynchronous tasks with a separate arrival process. (2 – Structural assumptions) The expressiveness of a temporal model is determined by the dependencies and shapes of its variables. In particular, the topological structure should be rich enough to capture the dynamics of the underlying process but sparse enough to allow for robust learning and efficient inference. (3 – Generation) A good time series model must be able to generate unbiased samples from the true underlying process in an efficient way. (4 – Inference) Given a trained model, it should support standard inference tasks, such as interpolation, forecasting, and likelihood calculation.

Recently, deep generative modeling has enabled vastly increased flexibility while keeping generation and inference tractable, owing to novel techniques like amortized variational inference [29, 12], reversible generative models [43, 30], and networks based on differential equations [9, 36].

In this work, we approach the modeling of continuous and irregular time series with a *reversible generative model for stochastic processes*. Our approach builds upon ideas from normalizing flows; however, instead of a static base distribution, we transform a dynamic base process into an observable one. In particular, we introduce the *continuous-time flow process (CTFP)*, a novel type of generative model

---

\*Work developed during an internship at Borealis AI. Correspond to [wsdmdeng@gmail.com](mailto:wsdmdeng@gmail.com).

that decodes the base continuous Wiener process into a complex observable process using a dynamic instance of normalizing flows. The resulting observable process is thus continuous in time. In addition to the appealing properties of static normalizing flows (e.g., efficient sampling and exact likelihood), this also enables a series of inference tasks that are typically unattainable in time series models with complex dynamics, such as interpolation and extrapolation at arbitrary timestamps. Furthermore, to overcome the simple covariance structure of the Wiener process, we augment the reversible mapping with latent variables and optimize this latent CTFP variant using variational optimization. Our approach is illustrated in Figure 1.

**Contributions.** In summary, we propose the *continuous-time flow process (CTFP)*, a novel generative model for continuous stochastic processes. It has the following appealing properties: (1) it induces flexible and consistent joint distributions on arbitrary and irregular time grids, with easy-to-compute density and an efficient sampling procedure; (2) the stochastic process generated by CTFP is guaranteed to have continuous sample paths, making it a natural fit for data with continuously-changing dynamics; (3) CTFP can perform interpolation and extrapolation conditioned on given observations. We validate our model and its latent variant on various stochastic processes and real-world datasets and show superior performance to state-of-the-art methods, including the variational recurrent neural network (VRNN) [12] and latent ordinary differential equation (latent ODE) [44].

## 2 Related Work

The following sections discuss the relevant literature on statistical models for sequential data and put it in context with the proposed approach.

**Early Work.** Among the most popular traditional time series models are latent variable models following the state-space equations [16], including the well-known variants with discrete and linear state-space [2, 27]. In the non-linear case, exact inference is typically intractable and we need to resort to approximate techniques [26, 24, 7, 8, 47]. Our CTFP can be viewed as a form of a continuous-time extended Kalman filter where the nonlinear observation process is noiseless and invertible and the temporal dynamics are a Wiener process. The final result, however, is more expressive than a Wiener process but retains some of its appealing properties like closed-form likelihood, interpolation, and extrapolation. Tree-based variants of non-linear Markov models have been proposed in [35]. An augmentation with switching states increases the expressiveness of state-space models but introduces additional challenges for learning [17] and inference [1]. Marginalization over an expansion of the state-space equations in terms of non-linear basis functions extends classical Gaussian processes [42] to Gaussian process dynamical models [25].

**Variational Sequence Models.** Following its success on image data, many works extended the variational autoencoder (VAE) [29] to sequential data [5, 12, 18, 37]. While RNN-based variational sequence models [12, 5] can model distributions over irregular timestamps, those timestamps have to be discrete and thus the models lack the notion of continuity. As a result, they are not suitable for modeling sequential data that have continuous underlying dynamics. Furthermore, it is not straightforward to perform interpolation at arbitrary timestamps using those models.

Latent ODEs [44] use an ODE-RNN as encoder and propagate a latent variable along a time interval using a neural ODE. This formulation ensures that the latent trajectory is continuous in time. However, decoding of the latent variables to observations is done at each time step independently. As a result, there is no guarantee that sample paths are continuous, which causes problems similar to the ones observed in variational sequence models. Neural stochastic differential equations (neural SDEs) [36] replace the deterministic latent trajectory of a latent ODE with a latent stochastic process but do also not generate continuous sample paths.

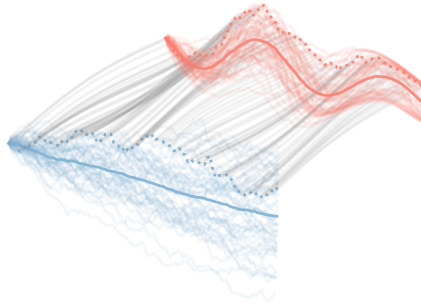


Figure 1: **Overview.** Wiener processes are continuous stochastic processes with appealing properties but limited flexibility. We propose to learn a complex observed process (red) through a differential deformation (grey) of the base Wiener process (blue), thereby preserving the advantages of the base process.

Recently, Qin et al. [41] proposed the recurrent neural process model. However, members of the neural process family [28, 19, 20, 46] only model the conditional distribution of data given observations and are not generic generative models.

**Normalizing Flows in Time Series Modeling.** Multiple recent works [33, 38, 45] apply reversible generative models to sequential data and show promise at capturing complex distributions. Mehrasa et al. [38] and Shchur et al. [45] use normalizing flows to model the distribution of inter-arrival time between events in temporal point processes. Kumar et al. [33] generate video frames using conditional normalizing flows. However, these models only use normalizing flows to model probability distributions in real space. In contrast, our model extends the domain of normalizing flows from distributions in real space to continuous-time stochastic processes.

### 3 Background

Our model is built upon the study of stochastic processes and recent advances in normalizing flow research. The following sections introduce the necessary background in these areas.

#### 3.1 Stochastic Processes

A stochastic process can be defined as a collection of random variables that are indexed by time. An example of a continuous stochastic process is the *Wiener process*. The  $d$ -dimensional Wiener process  $\mathbf{W}_\tau$  can be characterized by the following properties: (1)  $\mathbf{W}_0 = 0$ ; (2)  $\mathbf{W}_t - \mathbf{W}_s \sim \mathcal{N}(0, (t - s)\mathbf{I}_d)$  for  $s \leq t$ , and  $\mathbf{W}_t - \mathbf{W}_s$  is independent of past values of  $\mathbf{W}_{s'}$  for all  $s' \leq s$ . The joint density of  $(\mathbf{W}_{\tau_1}, \dots, \mathbf{W}_{\tau_n})$  can be written as the product of the conditional densities:  $p(\mathbf{w}_{\tau_1}, \dots, \mathbf{w}_{\tau_n}) = \prod_{i=1}^n p_{\mathbf{W}_{\tau_i} | \mathbf{W}_{\tau_{i-1}}}(\mathbf{w}_{\tau_i} | \mathbf{w}_{\tau_{i-1}})$  for  $0 \leq \tau_1 < \dots < \tau_n \leq T$ .

The conditional distribution of  $p_{\mathbf{W}_t | \mathbf{W}_s}$ , for  $s < t$ , is multivariate Gaussian; its conditional density is

$$p_{\mathbf{W}_t | \mathbf{W}_s}(\mathbf{w}_t | \mathbf{w}_s) = \mathcal{N}(\mathbf{w}_t; \mathbf{w}_s, (t - s)\mathbf{I}_d), \quad (1)$$

where  $\mathbf{I}_d$  is a  $d$ -dimensional identity matrix. This equation also provides a way to sample from  $(\mathbf{W}_{\tau_1}, \dots, \mathbf{W}_{\tau_n})$ . Furthermore, given  $\mathbf{W}_{t_1} = \mathbf{w}_{t_1}$  and  $\mathbf{W}_{t_2} = \mathbf{w}_{t_2}$ , the conditional distribution of  $\mathbf{W}_t$  for  $t_1 \leq t \leq t_2$  is also Gaussian:

$$p_{\mathbf{W}_t | \mathbf{w}_{t_1}, \mathbf{w}_{t_2}}(\mathbf{w}_t | \mathbf{w}_{t_1}, \mathbf{w}_{t_2}) = \mathcal{N}\left(\mathbf{w}_t; \mathbf{w}_{t_1} + \frac{t - t_1}{t_2 - t_1}(\mathbf{w}_{t_2} - \mathbf{w}_{t_1}), \frac{(t_2 - t)(t - t_1)}{t_2 - t_1} \mathbf{I}_d\right). \quad (2)$$

This is known as the Brownian bridge. An important property of the Wiener process is that the sample paths are continuous in time with probability one. This property allows our models to generate continuous sample paths and perform interpolation and extrapolation tasks.

#### 3.2 Normalizing Flows

*Normalizing flows* [43, 13, 31, 14, 39, 30, 3, 10, 32, 40] are reversible generative models that allow both density estimation and sampling. If our interest is to estimate the density function  $p_{\mathbf{X}}$  of a random vector  $\mathbf{X} \in \mathbb{R}^d$ , then normalizing flows assume  $\mathbf{X} = f(\mathbf{Z})$ , where  $f : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a bijective function, and  $\mathbf{Z} \in \mathbb{R}^d$  is a random vector with a simple density function  $p_{\mathbf{Z}}$ . The probability density function can be evaluated using the change of variables formula:

$$\log p_{\mathbf{X}}(\mathbf{x}) = \log p_{\mathbf{Z}}(g(\mathbf{x})) + \log \left| \det \left( \frac{\partial g}{\partial \mathbf{x}} \right) \right|, \quad (3)$$

where we denote the inverse of  $f$  by  $g$  and  $\partial g / \partial \mathbf{x}$  is the Jacobian matrix of  $g$ . Sampling from  $p_{\mathbf{X}}$  can be done by first drawing a sample from the simple distribution  $\mathbf{z} \sim p_{\mathbf{Z}}$ , and then apply the bijection  $\mathbf{x} = f(\mathbf{z})$ .

Chen et al. [9], Grathwohl et al. [21] proposed the *continuous normalizing flow*, which uses the *neural ordinary differential equation* (neural ODE) to model a flexible bijective mapping. Given  $\mathbf{z} = \mathbf{h}(t_0)$  sampled from the base distribution  $p_{\mathbf{Z}}$ , it is mapped to  $\mathbf{h}(t_1)$  based on the mapping defined by the ODE:  $d\mathbf{h}(t)/dt = f(\mathbf{h}(t), t)$ . The change in log-density is computed by the *instantaneous change of variables formula* [9]:

$$\log p_{\mathbf{X}}(\mathbf{h}(t_1)) = \log p_{\mathbf{Z}}(\mathbf{h}(t_0)) - \int_{t_0}^{t_1} \text{tr} \left( \frac{\partial f}{\partial \mathbf{h}(t)} \right) dt. \quad (4)$$

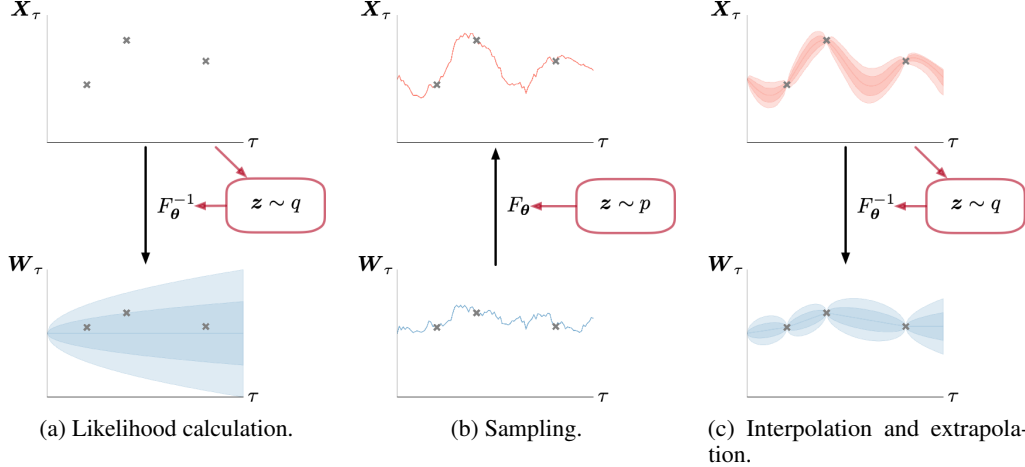


Figure 2: **(Latent) Continuous-Time Flow Processes (CTFPs)**. (a) Likelihood calculation. Given an irregular time series  $\{\mathbf{x}_{\tau_i}\}$ , the inverse flow  $F_\theta^{-1}$  maps the observed process to a set of Wiener points  $\{\mathbf{w}_{\tau_i}\}$  for which we can compute the likelihood according to Equation 7. (b) Sampling. Given a set of timestamps  $\{\tau_i\}$ , we sample a Wiener process and use the forward flow  $F_\theta$  to obtain a sample of the observed process. (c) Interpolation and extrapolation. In order to compute the density at an unobserved point  $\mathbf{x}_\tau$ , we compute the left-sided (extrapolation; Equation 1) or two-sided (interpolation; Equation 2) conditional density of its Wiener point  $\mathbf{w}_\tau$  and adjust for the flow (Equation 11). **Notes:** The effect of the latent variables  $\mathbf{Z}$  in our latent CTFP model is indicated by red boxes. The shaded areas represent 70% and 95% confidence intervals.

One potential disadvantage of the neural ODE model is that it preserves the topology of the input space, and there are classes of functions that cannot be represented by neural ODEs. Dupont et al. [15] proposed the augmented neural ODE (ANODE) model to address this limitation. Note that the original formulation of ANODE is not a generative model and it does not support the computation of likelihoods  $p_{\mathcal{X}}(\mathbf{x})$  or sampling from the target distribution  $\mathbf{x} \sim p_{\mathcal{X}}$ . In this work, we formulate a modified version of ANODE that can be used as a (conditional) generative model.

## 4 Model

We define our proposed *continuous-time flow process (CTFP)* in Section 4.1. In Section 4.2, a generative variant of ANODE is presented as a component to implement CTFP. Since the proposed stochastic process is continuous in time, it enables interpolation and extrapolation at arbitrary time points, as described in Section 4.3. Finally, richer covariance structures are enabled by the latent CTFP model presented in Section 4.4.

### 4.1 Continuous-Time Flow Process

Let  $\{(\mathbf{x}_{\tau_i}, \tau_i)\}_{i=1}^n$ , denote a sequence of irregularly spaced time series data. We assume the time series to be an (incomplete) realization of a continuous stochastic process  $\{\mathbf{X}_\tau\}_{\tau \in [0, T]}$ . In other words, this stochastic process induces a joint distribution of  $(\mathbf{X}_{\tau_1}, \dots, \mathbf{X}_{\tau_n})$ . Our goal is to model  $\{\mathbf{X}_\tau\}_{\tau \in [0, T]}$  such that the log-likelihood of the observations

$$\mathcal{L} = \log p_{\mathbf{X}_{\tau_1}, \dots, \mathbf{X}_{\tau_n}}(\mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n}) \quad (5)$$

is maximized. We define the continuous-time flow process (CTFP)  $\{F_\theta(\mathbf{W}_\tau; \tau)\}_{\tau \in [0, T]}$  such that

$$\mathbf{X}_\tau = F_\theta(\mathbf{W}_\tau; \tau), \quad \forall \tau \in [0, T], \quad (6)$$

where  $F_\theta(\cdot; \tau) : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is an invertible mapping parametrized by the learnable parameters  $\theta$  for every  $\tau \in [0, T]$ , and  $\mathbf{W}_\tau$  is a  $d$ -dimensional Wiener process.

The log-likelihood in Equation 5 can be rewritten using the change of variables formula. Let  $\mathbf{w}_{\tau_i} = F_\theta^{-1}(\mathbf{x}_{\tau_i}; \tau_i)$ , then

$$\mathcal{L} = \sum_{i=1}^n \left[ \log p_{\mathbf{W}_{\tau_i} | \mathbf{W}_{\tau_{i-1}}}(\mathbf{w}_{\tau_i} | \mathbf{w}_{\tau_{i-1}}) - \log \left| \det \frac{\partial F_\theta(\mathbf{w}_{\tau_i}; \tau_i)}{\partial \mathbf{w}_{\tau_i}} \right| \right], \quad (7)$$

where  $\tau_0 = 0$ ,  $\mathbf{W}_0 = 0$ , and  $p_{\mathbf{W}_{\tau_i}|\mathbf{W}_{\tau_{i-1}}}$  is defined in Section 3.1. Figure 2a shows an example of the likelihood calculation. Sampling from a CTFP is straightforward: given the timestamps  $\tau_i$ , we first sample a realization of the Wiener process  $\{\mathbf{w}_{\tau_i}\}_{i=1}^n$ , then map them to  $\mathbf{x}_{\tau_i} = F_{\theta}(\mathbf{w}_{\tau_i}; \tau_i)$ . Figure 2b illustrates this procedure.

The normalizing flow  $F_{\theta}(\cdot; \tau)$  transforms a simple base distribution induced by  $\mathbf{W}_{\tau}$  on an arbitrary time grid into a more complex shape in the observation space. It is worth noting that given a continuous realization of  $\mathbf{W}_{\tau}$ , as long as  $F_{\theta}(\cdot; \tau)$  is implemented as a continuous mapping, the resulting trajectory  $\mathbf{x}_{\tau}$  is also continuous.

## 4.2 Generative ANODE

In principle, any normalizing flow model indexed by time  $\tau$  could be used as  $F_{\theta}(\cdot; \tau)$  in Equation 6. We proceed with the continuous normalizing flow and ANODE, because it has a free-form Jacobian and efficient trace estimator [15, 21]. In particular, we consider the following instantiation of ANODE as a generative model: For any  $\tau \in [0, T]$  and  $\mathbf{w}_{\tau} \in \mathbb{R}^d$ , we map  $\mathbf{w}_{\tau}$  to  $\mathbf{x}_{\tau}$  by solving the following initial value problem:

$$\frac{d}{dt} \begin{pmatrix} \mathbf{h}_{\tau}(t) \\ a_{\tau}(t) \end{pmatrix} = \begin{pmatrix} f_{\theta}(\mathbf{h}_{\tau}(t), a_{\tau}(t), t) \\ g_{\theta}(a_{\tau}(t), t) \end{pmatrix}, \quad \begin{pmatrix} \mathbf{h}_{\tau}(t_0) \\ a_{\tau}(t_0) \end{pmatrix} = \begin{pmatrix} \mathbf{w}_{\tau} \\ \tau \end{pmatrix}, \quad (8)$$

where  $\mathbf{h}_{\tau}(t) \in \mathbb{R}^d$ ,  $t \in [t_0, t_1]$ ,  $f_{\theta} : \mathbb{R}^d \times \mathbb{R} \times [t_0, t_1] \rightarrow \mathbb{R}^d$ , and  $g_{\theta} : \mathbb{R} \times [t_0, t_1] \rightarrow \mathbb{R}$ . Then  $F_{\theta}$  in Equation 6 is defined as the solution of  $\mathbf{h}_{\tau}(t)$  at  $t = t_1$ :

$$F_{\theta}(\mathbf{w}_{\tau}; \tau) := \mathbf{h}_{\tau}(t_1) = \mathbf{h}_{\tau}(t_0) + \int_{t_0}^{t_1} f_{\theta}(\mathbf{h}_{\tau}(t), a_{\tau}(t), t) dt. \quad (9)$$

Note that the index  $t$  represents the independent variable in the initial value problem and should not be confused with  $\tau$ , the timestamp of the observation.

Using Equation 4, the log-likelihood  $\mathcal{L}$  can be calculated as follows:

$$\mathcal{L} = \sum_{i=1}^n \left[ \log p_{\mathbf{W}_{\tau_i}|\mathbf{W}_{\tau_{i-1}}}(\mathbf{h}_{\tau_i}(t_0) | \mathbf{h}_{\tau_{i-1}}(t_0)) - \int_{t_0}^{t_1} \text{tr} \left( \frac{\partial f_{\theta}(\mathbf{h}_{\tau_i}(t), a_{\tau_i}(t), t)}{\partial \mathbf{h}_{\tau_i}(t)} \right) dt \right], \quad (10)$$

where  $\mathbf{h}_{\tau_i}(t_0)$  is obtained by solving the ODE in Equation 8 backwards from  $t = t_1$  to  $t = t_0$ , and the trace of the Jacobian can be estimated by Hutchinson's trace estimator [23, 21].

## 4.3 Interpolation and Extrapolation with CTFP

Time-indexed normalizing flows and the Brownian bridge allow us to define conditional distributions on arbitrary timestamps. They also permit the CTFP model to perform interpolation and extrapolation given partial observations, which are of great importance in time series modeling.

Interpolation means that we can model the conditional distribution  $p_{\mathbf{X}_{\tau}|\mathbf{X}_{\tau_i}, \mathbf{X}_{\tau_{i+1}}}(\mathbf{x}_{\tau} | \mathbf{x}_{\tau_i}, \mathbf{x}_{\tau_{i+1}})$  for all  $\tau \in [\tau_i, \tau_{i+1}]$  and  $i = 1, \dots, n-1$ . This can be done by mapping the values  $\mathbf{x}_{\tau}$ ,  $\mathbf{x}_{\tau_i}$  and  $\mathbf{x}_{\tau_{i+1}}$  to  $\mathbf{w}_{\tau}$ ,  $\mathbf{w}_{\tau_i}$  and  $\mathbf{w}_{\tau_{i+1}}$ , respectively. After that, Equation 2 can be applied to obtain the conditional density of  $p_{\mathbf{W}_{\tau}|\mathbf{W}_{\tau_i}, \mathbf{W}_{\tau_{i+1}}}(\mathbf{w}_{\tau} | \mathbf{w}_{\tau_i}, \mathbf{w}_{\tau_{i+1}})$ . Finally, we have

$$\log p_{\mathbf{X}_{\tau}|\mathbf{X}_{\tau_i}, \mathbf{X}_{\tau_{i+1}}}(\mathbf{x}_{\tau} | \mathbf{x}_{\tau_i}, \mathbf{x}_{\tau_{i+1}}) = \log p_{\mathbf{W}_{\tau}|\mathbf{W}_{\tau_i}, \mathbf{W}_{\tau_{i+1}}}(\mathbf{w}_{\tau} | \mathbf{w}_{\tau_i}, \mathbf{w}_{\tau_{i+1}}) - \log \left| \det \frac{\partial \mathbf{x}_{\tau}}{\partial \mathbf{w}_{\tau}} \right|. \quad (11)$$

Extrapolation can be done in a similar fashion using Equation 1. This allows the model to predict continuous trajectories into the future, given past observations. Figure 2c shows a visualization of interpolation and extrapolation using CTFP.

## 4.4 Latent Continuous-Time Flow Process

The CTFP model inherits the Markov property from the Wiener process, which is a strong assumption and limits its ability to model stochastic processes with complex temporal dependencies. In order to enhance the expressive power of the CTFP model, we augment it with a latent variable  $\mathbf{Z} \in \mathbb{R}^m$ ,

whose prior distribution is an isotropic Gaussian  $p_{\mathbf{z}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; 0, \mathbf{I}_m)$ . As a result, the data distribution can be approximated by a diverse collection of CTFP models conditioned on sampled latent variables  $\mathbf{z}$ .

The generative model in Equation 6 is augmented to  $\mathbf{X}_\tau = F_\theta(\mathbf{W}_\tau; \mathbf{Z}, \tau), \forall \tau \in [0, T]$ , which induces the conditional distribution  $\mathbf{X}_{\tau_1}, \dots, \mathbf{X}_{\tau_n} | \mathbf{Z}$ . Similar to the initial value problem in Equation 8, we define  $F_\theta(\mathbf{w}_\tau; \mathbf{z}, \tau) = \mathbf{h}_\tau(t_1)$ , where

$$\frac{d}{dt} \begin{pmatrix} \mathbf{h}_\tau(t) \\ \mathbf{a}_\tau(t) \end{pmatrix} = \begin{pmatrix} f_\theta(\mathbf{h}_\tau(t), \mathbf{a}_\tau(t), t) \\ g_\theta(\mathbf{a}_\tau(t), t) \end{pmatrix}, \quad \begin{pmatrix} \mathbf{h}_\tau(t_0) \\ \mathbf{a}_\tau(t_0) \end{pmatrix} = \begin{pmatrix} \mathbf{w}_\tau \\ (\mathbf{z}, \tau)^\top \end{pmatrix}. \quad (12)$$

Depending on the sample of the latent variable  $\mathbf{z}$ , the CTFP model has different gradient fields and thus different output distributions.

For simplicity of notation, the subscripts of density functions are omitted from now on. For the augmented generative model, the log-likelihood becomes  $\mathcal{L} = \log \int_{\mathbb{R}^m} p(\mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n} | \mathbf{z}) p(\mathbf{z}) d\mathbf{z}$ , which is intractable to evaluate. Following the variational autoencoder approach [29], we introduce an approximate posterior distribution of  $\mathbf{Z} | \mathbf{X}_{\tau_1}, \dots, \mathbf{X}_{\tau_n}$ , denoted by  $q(\mathbf{z} | \mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n})$ . The implementation of the approximate posterior distribution is an ODE-RNN encoder [44]. With the approximate posterior distribution, we can derive an importance-weighted autoencoder (IWAE) [6] lower bound of the log-likelihood on the right-hand side of the inequality:

$$\begin{aligned} \mathcal{L} &= \log \mathbb{E}_{\mathbf{z} \sim q} \left[ \frac{p(\mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n} | \mathbf{z}) p(\mathbf{z})}{q(\mathbf{z} | \mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n})} \right] \\ &\geq \mathbb{E}_{\mathbf{z}_1, \dots, \mathbf{z}_K \sim q} \left[ \log \left( \frac{1}{K} \sum_{k=1}^K \frac{p(\mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n} | \mathbf{z}_k) p(\mathbf{z}_k)}{q(\mathbf{z}_k | \mathbf{x}_{\tau_1}, \dots, \mathbf{x}_{\tau_n})} \right) \right] =: \mathcal{L}_{\text{IWAE}}, \end{aligned} \quad (13)$$

where  $K$  is the number of samples from the approximate posterior distribution.

## 5 Experiments

In this section, we apply our models on synthetic data generated from common continuous-time stochastic processes and complex real-world datasets. The proposed CTFP and latent CTFP models are compared against two baseline models: latent ODEs [44] and variational RNNs (VRNNs) [12]. The latent ODE model with the ODE-RNN encoder is designed specifically to model time series data with irregular observation times. VRNN is a popular variational filtering model that demonstrated superior performance on structured sequential data.

For VRNNs, we append the time gap between two observations as an additional input to the neural network. Both latent CTFP and latent ODE models use ODE-RNN [44] as the inference network; GRU [11] is used as the RNN cell in latent CTFP, latent ODE, and VRNN models. All three latent variable models have the same latent dimension and GRU hidden state dimension. Please see the supplementary materials for details about our experimental setup and model implementations.

### 5.1 Synthetic Datasets

We simulate three irregularly-sampled time series datasets; all of them are univariate. **Geometric Brownian motion (GBM)** is a continuous-time stochastic process widely used in mathematical finance. It satisfies the following stochastic differential equation:  $dX_\tau = \mu X_\tau d\tau + \sigma X_\tau dW_\tau$ , where  $\mu$  and  $\sigma$  are the drift term and variance term, respectively. The timestamps of the observations are in the range between 0 and  $T = 30$  and are sampled from a homogeneous Poisson point process with an intensity of  $\lambda_{\text{train}} = 2$ . To further evaluate the model's capacity to capture the dynamics of GBM, we test the model with observation time-steps sampled from Poisson point processes with intensities of  $\lambda_{\text{test}} = 2$  and  $\lambda_{\text{test}} = 20$ . **Ornstein-Uhlenbeck process (OU Process)** is another type of widely used continuous-time stochastic process. The OU process satisfies the following stochastic differential equation:  $dX_\tau = \theta(\mu - X_\tau) d\tau + \sigma dW_\tau$ . We use the same set of observation intensities as in our GBM experiments to sample observation timestamps in the training and test sets. **Mixture of OUs.** To demonstrate the latent CTFP's capability to model sequences sampled from different continuous-time stochastic processes, we train our models on a dataset generated by mixing the sequences sampled from two different OU processes with different values of  $\theta, \mu, \sigma$ , and

Table 1: **Quantitative Evaluation (Synthetic Data)**. We show test negative log-likelihood on three synthetic stochastic processes across different models. Below each process, we indicate the intensity of the Poisson point process from which the timestamps for the test sequences were sampled. “Ground Truth” refers to the closed-form negative log-likelihood of the true underlying data generation process. [GBM: geometric Brownian motion; OU: Ornstein–Uhlenbeck process; M-OU: mixture of OUs.]

Model	GBM		OU		M-OU
	$\lambda_{\text{test}} = 2$	$\lambda_{\text{test}} = 20$	$\lambda_{\text{test}} = 2$	$\lambda_{\text{test}} = 20$	$\lambda_{\text{test}} = (2, 20)$
Latent ODE [44]	3.826	5.935	3.066	3.027	2.690
VRNN [12]	3.762	3.492	<b>2.729</b>	<b>1.939</b>	1.415
CTFP ( <b>ours</b> )	<b>3.107</b>	<b>1.929</b>	2.902	1.941	1.408
Latent CTFP ( <b>ours</b> )	<b>3.107</b>	1.930	2.902	<b>1.939</b>	<b>1.392</b>
Ground Truth	3.106	1.928	2.722	1.888	1.379

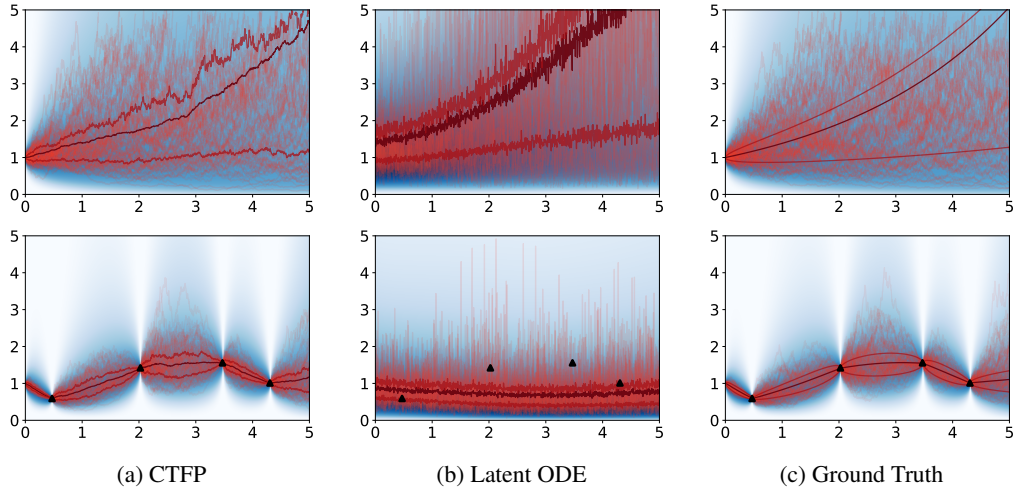


Figure 3: **Comparison between CTFP and latent ODE on the GBM data**. We consider the generation and interpolation tasks for (a) CTFP, (b) latent ODE, and (c) ground truth. In each subfigure, the upper panel shows samples generated from the model and the lower panel shows results for interpolation. The observed points for interpolation are marked by black triangles. In addition to the sample trajectories (red) and the marginal density (blue), we also show the sample-based estimates (closed-form for ground truth) of the inter-quartile range (dark red) and mean (brown) of the marginal density.

different observation intensities. We defer the details of the parameters of the synthetic dataset to the supplementary materials.

**Results.** The results are presented in Table 1. We report the exact negative log-likelihood (NLL) per observation for CTFP. For latent ODE, latent CTFP, and VRNN, we report the (upper bound of) NLL estimated by the IWAE bound [6] in Equation 13, using  $K = 25$  samples of latent variables. We also show the NLL of the test set computed with the ground truth density function.

The results on the test set sampled from the GBM indicate that the CTFP model can recover the true data generation process as the NLL estimated by CTFP is close to the ground truth. In contrast, latent ODE and VRNN models fail to recover the true data distribution. On the M-OU dataset, the latent CTFP models show better performance than the other models. Moreover, latent CTFP outperforms CTFP by 0.016 nats, indicating its ability to leverage the latent variables.

Although trained on samples with an observation intensity of  $\lambda_{\text{train}} = 2$ , CTFP can better adapt to samples with a bigger observation intensity (and thus denser time grid) of  $\lambda_{\text{test}} = 20$ . We hypothesize that the superior performance of CTFP models when  $\lambda_{\text{test}} = 20$  is due to their capability to model continuous stochastic processes, whereas the baseline models do not have the notion of continuity. We further corroborate this hypothesis in an ablation study where the base Wiener process is replaced with i.i.d. Gaussian random variables, such that the base process is no longer continuous in time (see the supplementary materials).

Table 2: **Quantitative Evaluation (Real-World Data)**. We show test negative log-likelihood on Mujoco-Hopper, Beijing Air-Quality Dataset (BAQD) and PTB Diagnostic Database (PTBDB) across different models. For CTFP, the reported values are exact; for the other four models, we report IWAE bounds using  $K = 125$  samples. Latent CTFP-ET stands for a latent CTFP model trained and evaluated with the exact Jacobian trace. Lower values correspond to better performance. Standard deviations are based on 5 independent runs.

Model	Mujoco-Hopper [44]	BAQD [4]	PTBDB [49]
Latent ODE [44]	$24.775 \pm 0.010$	$2.789 \pm 0.011$	$-0.818 \pm 0.009$
VRNN [12]	$9.113 \pm 0.018$	$0.604 \pm 0.007$	<b><math>-1.999 \pm 0.008</math></b>
CTFP ( <b>ours</b> )	$-16.249 \pm 0.034$	$-2.361 \pm 0.020$	$-1.324 \pm 0.028$
Latent CTFP ( <b>ours</b> )	<b><math>-31.397 \pm 0.063</math></b>	<b><math>-6.894 \pm 0.046</math></b>	<b><math>-1.999 \pm 0.010</math></b>
Latent CTFP-ET ( <b>ours</b> )	<b><math>-21.934 \pm 0.029</math></b>	<b><math>-2.894 \pm 0.046</math></b>	<b><math>-1.999 \pm 0.010</math></b> <sup>1</sup>

Figure 3 provides a qualitative comparison between CTFP and latent ODE trained on the GBM data, both on the generation task (upper panels) and the interpolation task (lower panels). The results in the upper panels show that CTFP can generate continuous sample paths and accurately estimate the marginal mean and quantiles. In contrast, the sample paths generated by latent ODE are more volatile and discontinuous due to its lack of continuity. For the interpolation task, the results of CTFP are consistent with the ground truth in terms of both point estimation and uncertainty estimation. For latent ODE on the interpolation task, Figure 3b shows that the latent variables from the variational posterior shift the density to the region where the observations lie. However, although latent ODE is capable of performing interpolation, there is no guarantee that the (reconstructed) sample paths pass through the observed points (triangular marks in Figure 3b), as discussed in Section 2. In addition to these difficulties with the interpolation task, the qualitative comparison between samples further highlights the importance of our models’ continuity when generating samples of continuous dynamics.

## 5.2 Real-World Datasets

We also evaluate our models on real-world datasets with continuous and complex dynamics. The following three datasets are considered: **Mujoco-Hopper** [44] consists of 10,000 sequences that are simulated by a “Hopper” model from the DeepMind Control Suite in a MuJoCo environment [48]. **PTB Diagnostic Database** (PTBDB) [4] consists of excerpts of ambulatory electrocardiography (ECG) recordings. Each sequence is one-dimensional and the sampling frequency of the recordings is 125 Hz. **Beijing Air-Quality Dataset** (BAQD) [49] is a dataset consisting of multi-year recordings of weather and air quality data across different locations in Beijing. The variables in consideration are temperature, pressure, and wind speed, and the values are recorded once per hour. We segment the data into sequences, each covering the recordings of a whole week. Please refer to the supplementary materials for additional details about data preprocessing.

Similar to our synthetic data experiment settings, we compare the CTFP and latent CTFP models against latent ODE and VRNN. It is worth noting that the latent ODE model in the original work [44] uses a fixed output variance and is evaluated using mean squared error (MSE); we adapt the model to our tasks with a predicted output variance (see supplementary materials). We further study the effect of using RealNVP [14] as the invertible mapping  $F_{\theta}(\cdot; \tau)$ . This experiment can be regarded as an ablation study and results are presented in the supplementary materials as well.

**Results.** The results are shown in Table 2. We report the exact negative log-likelihood (NLL) per observation for CTFP and the (upper bound of) NLL estimated by the IWAE bound, using  $K = 125$  samples of latent variables, for latent ODE, latent CTFP, and VRNN. For each setting, the mean and standard deviation of five evaluation runs are reported. The latent CTFP models trained and evaluated with Hutchinson’s trace estimator [23, 21] on multi-dimensional data could lead to a biased IWAE estimation and unstable training due to the variance of Hutchinson’s trace estimator and the log-sum-exp operation. Therefore, we also report results trained and estimated with the exact Jacobian trace, denoted by latent CTFP-ET. The evaluation results show that the latent CTFP model outperforms VRNN and latent ODE models on real-world datasets, indicating that CTFP is better at modeling irregular time series data with continuous dynamics. Table 2 also suggests that the latent CTFP model consistently outperforms the CTFP model, demonstrating that with the latent variables,

<sup>1</sup>Hutchinson’s trace estimator using Rademacher-distributed noise yields the same results as computation of the exact trace for one-dimensional data.



the latent CTFP model is more expressive and able to capture the data distribution better. We defer additional experimental results with different observation processes but similar conclusions to the supplementary materials.

## 6 Conclusion

In summary, we propose the continuous-time flow process (CTFP), a reversible generative model for stochastic processes, and its latent variant. It maps a simple continuous-time stochastic process, i.e., the Wiener process, into a more complicated process in the observable space. As a result, many desirable mathematical properties of the Wiener process are retained, including the efficient sampling of continuous paths, likelihood evaluation on arbitrary timestamps, and inter-/extrapolation given observed data. Our experimental results demonstrate the superior performance of the proposed models on various datasets.

## Broader Impact Statement

Time series models could be applied to a wide range of applications, including natural language processing, recommendation systems, traffic prediction, medical data analysis, forecasting, and others. Our research improves over the existing models on a particular type of data: irregular time series data.

There are opportunities for applications using the proposed models for beneficial purposes, such as weather forecasting, pedestrian behavior prediction for self-driving cars, and missing healthcare data interpolation or prediction. We encourage practitioners to understand the impacts of using CTFP in particular real-world scenarios.

One potential risk is that the capability of interpolation and extrapolation can be used in malicious ways. An adversary might be able to use the proposed model to infer private information given partial observations, which leads to privacy concerns. We would encourage further research to address this risk using tools like differential privacy.

## Funding Transparency Statement

This work was conducted at Borealis AI and partly supported by Mitacs through the Mitacs Accelerate program.

## References

- [1] D. Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.
- [2] L.E. Baum and T. Petrie. Statistical inference for probabilistic functions of finite state markov chains. In *The Annals of Mathematical Statistics*, 1966.
- [3] Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Joern-Henrik Jacobsen. Invertible residual networks. In *International Conference on Machine Learning*, pages 573–582, 2019.
- [4] R Boussejot, D Kreiseler, and A Schnabel. Nutzung der ekg-signaldatenbank cardiodat der ptb über das internet. *Biomedizinische Technik/Biomedical Engineering*, 40(s1):317–318, 1995.
- [5] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*, 2015.
- [6] Yuri Burda, Roger Grosse, and Ruslan Salakhutdinov. Importance weighted autoencoders. In *International Conference on Learning Representations*, 2016.
- [7] O. Cappé, E. Moulines, and T. Rydén. *Hidden Markov Models and Dynamical Systems*. Springer, 2005.
- [8] O. Cappé, S.J. Godsill, and E. Moulines. An overview of existing methods and recent advances in sequential monte carlo. In *Proceedings of the IEEE*, 2007.
- [9] Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in neural information processing systems*, pages 6571–6583, 2018.
- [10] Tian Qi Chen, Jens Behrmann, David K Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for invertible generative modeling. In *Advances in Neural Information Processing Systems*, pages 9913–9923, 2019.

- [11] Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, 2014.
- [12] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua Bengio. A recurrent latent variable model for sequential data. In *Advances in neural information processing systems*, pages 2980–2988, 2015.
- [13] Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- [14] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. In *International Conference on Learning Representations*, 2017.
- [15] Emilien Dupont, Arnaud Doucet, and Yee Whye Teh. Augmented neural ODEs. In *Advances in Neural Information Processing Systems*, pages 3134–3144, 2019.
- [16] James Durbin and Siem Jan Koopman. *Time Series Analysis by State Space Methods*. Oxford University Press, 2012.
- [17] E. Fox, E. Sudderth, M. Jordan, and A. Willsky. Non-parametric bayesian learning of switching linear dynamical systems. In *NeurIPS*, 2008.
- [18] Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural models with stochastic layers. In *Advances in neural information processing systems*, pages 2199–2207, 2016.
- [19] Marta Garnelo, Dan Rosenbaum, Christopher Maddison, Tiago Ramalho, David Saxton, Murray Shanahan, Yee Whye Teh, Danilo Rezende, and SM Ali Eslami. Conditional neural processes. In *International Conference on Machine Learning*, pages 1704–1713, 2018.
- [20] Marta Garnelo, Jonathan Schwarz, Dan Rosenbaum, Fabio Viola, Danilo J Rezende, SM Eslami, and Yee Whye Teh. Neural processes. *arXiv preprint arXiv:1807.01622*, 2018.
- [21] Will Grathwohl, Ricky T. Q. Chen, Jesse Bettencourt, and David Duvenaud. Scalable reversible generative models with free-form continuous dynamics. In *International Conference on Learning Representations*, 2019.
- [22] Junxian He, Daniel Spokoyny, Graham Neubig, and Taylor Berg-Kirkpatrick. Lagging inference networks and posterior collapse in variational autoencoders. *arXiv preprint arXiv:1901.05534*, 2019.
- [23] Michael F Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 19(2):433–450, 1990.
- [24] K. Ito and K. Xiong. Gaussian filters for nonlinear filtering problems. In *IEEE Trans. on Automatic Control*, 2000.
- [25] J.M.Wang, D.J.Fleet, and A.Hertzmann. Gaussian process dynamical models for human motion. In *PAMI*, 2008.
- [26] S.J. Julier and J.K. Uhlmann. A new extension of the kalman filter to nonlinear systems. In *Aerospace/Defense Sensing, Simulation and Controls*, 1997.
- [27] R.E. Kalman. A new approach to linear filtering and prediction problems. In *Journal of Basic Engineering*, 1960.
- [28] Hyunjik Kim, Andriy Mnih, Jonathan Schwarz, Marta Garnelo, Ali Eslami, Dan Rosenbaum, Oriol Vinyals, and Yee Whye Teh. Attentive neural processes. In *International Conference on Learning Representations*, 2019.
- [29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [30] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1x1 convolutions. In *Advances in Neural Information Processing Systems*, pages 10215–10224, 2018.
- [31] Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling. Improved variational inference with inverse autoregressive flow. In *Advances in neural information processing systems*, pages 4743–4751, 2016.
- [32] Ivan Kobyzev, Simon Prince, and Marcus A Brubaker. Normalizing flows: Introduction and ideas. *arXiv preprint arXiv:1908.09257*, 2019.
- [33] Manoj Kumar, Mohammad Babaeizadeh, Dumitru Erhan, Chelsea Finn, Sergey Levine, Laurent Dinh, and Durk Kingma. Videoflow: A conditional flow-based model for stochastic video generation. *arXiv preprint arXiv:1903.01434*, 2019.
- [34] Jean-François Le Gall. *Brownian motion, martingales, and stochastic calculus*, volume 274. Springer, 2016.

- [35] A.M. Lehrmann, P. Gehler, and S. Nowozin. Efficient nonlinear markov models for human motion. In *CVPR*, 2014.
- [36] Xuechen Li, Ting-Kam Leonard Wong, Ricky TQ Chen, and David Duvenaud. Scalable gradients for stochastic differential equations. *arXiv preprint arXiv:2001.01328*, 2020.
- [37] Rui Luo, Weinan Zhang, Xiaojun Xu, and Jun Wang. A neural stochastic volatility model. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [38] Nazanin Mehrasa, Ruizhi Deng, Mohamed Osama Ahmed, Bo Chang, Jiawei He, Thibaut Durand, Marcus Brubaker, and Greg Mori. Point process flows. *arXiv preprint arXiv:1910.08281*, 2019.
- [39] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017.
- [40] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *arXiv preprint arXiv:1912.02762*, 2019.
- [41] Shenghao Qin, Jiacheng Zhu, Jimmy Qin, Wenshuo Wang, and Ding Zhao. Recurrent attentive neural process for sequential data. *arXiv preprint arXiv:1910.09323*, 2019.
- [42] C.E. Rasmussen. *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [43] Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International Conference on Machine Learning*, pages 1530–1538, 2015.
- [44] Yulia Rubanova, Tian Qi Chen, and David K Duvenaud. Latent ordinary differential equations for irregularly-sampled time series. In *Advances in Neural Information Processing Systems*, pages 5321–5331, 2019.
- [45] Oleksandr Shchur, Marin Biloš, and Stephan Günnemann. Intensity-free learning of temporal point processes. *International Conference on Learning Representations (ICLR)*, 2020.
- [46] Gautam Singh, Jaesik Yoon, Youngsung Son, and Sungjin Ahn. Sequential neural processes. In *Advances in Neural Information Processing Systems*, pages 10254–10264, 2019.
- [47] Simo Särkkä. On unscented kalman filtering for state estimation of continuous-time nonlinear systems. In *IEEE Trans. on Automatic Control*, 2007.
- [48] Yuval Tassa, Yotam Doron, Alistair Muldal, Tom Erez, Yazhe Li, Diego de Las Casas, David Budden, Abbas Abdolmaleki, Josh Merel, Andrew Lefrancq, Timothy Lillicrap, and Martin Riedmiller. DeepMind control suite. Technical report, DeepMind, January 2018. URL <https://arxiv.org/abs/1801.00690>.
- [49] Shuyi Zhang, Bin Guo, Anlan Dong, Jing He, Ziping Xu, and Song Xi Chen. Cautionary tales on air-quality improvement in beijing. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 473(2205):20170457, 2017.

## A Finite-Dimensional Distribution of CTFP

Equation 7 in Section 4 is the log density of the distribution obtained by applying the normalizing flow models to the finite-dimensional distribution of Wiener process on a given time grid. A natural question that would arise is why the distribution described by Equation 7 necessarily matches the finite-dimensional distribution of  $\mathbf{X}_\tau = F_\theta(\mathbf{W}_\tau, \tau)$ . In other words, it is left to close the gap between the distributions of samples obtained by two different ways to justify Equation 7: (1) first getting a sample path of  $\mathbf{X}_\tau$  by applying the transformation defined by  $F_\theta$  to a sample of  $\mathbf{W}_\tau$  and then obtaining the finite-dimensional observation of  $\mathbf{X}_\tau$  on the time grid; (2) first obtaining the finite-dimensional sample of  $\mathbf{W}_\tau$  and applying the normalizing flows to this finite-dimensional distribution. To justify the finite-dimensional distribution of CTFP, we choose to work with the canonical Wiener space  $(\Omega, \Sigma)$  equipped with the unique Wiener measure  $\mu_{\mathbf{W}}$  where  $\Omega = C([0, +\infty), \mathbb{R}^d)$  is the set of continuous functions from  $[0, +\infty)$  to  $\mathbb{R}^d$ ,  $\Sigma$  is the Borel  $\sigma$ -algebra generated by all the cylinder sets of  $C([0, +\infty), \mathbb{R}^d)$ , and  $\mathbf{W}_\tau(\omega) = \omega(\tau)$  for  $\omega \in \Omega$ . We refer the reader to Chapter 2 of [34] for more details. Given a time grid  $0 < \tau_1 < \tau_2 < \dots < \tau_n$ , the distribution of observations of Wiener process on this discrete time grid is called the finite-dimensional distribution of  $\mathbf{W}_\tau$ . It is a push-forward measure on  $(\mathbb{R}^{d \times n}, \mathcal{B}(\mathbb{R}^{d \times n}))$  induced by the projection mapping  $\pi_{\tau_1, \tau_2, \dots, \tau_n} : (\Omega, \Sigma) \rightarrow ((\mathbb{R}^{d \times n}, \mathcal{B}(\mathbb{R}^{d \times n})))$  on this grid where  $\mathcal{B}(\cdot)$  denotes the Borel  $\sigma$ -algebra. Therefore, for each Borel (measurable) set  $B$  of  $\mathbb{R}^{d \times n}$ , the finite-dimensional distribution of  $B$  is  $\mu_{\mathbf{W}} \circ \pi^{-1}(B) = \mu_{\mathbf{W}}(\{\omega | (\mathbf{W}_{\tau_1}(\omega) \dots \mathbf{W}_{\tau_n}(\omega)) \in B\})$ . We drop the subscript of  $\pi$  for the simplicity of notation. We base the justification on the following two propositions.

**Proposition 1.** *Let  $F_\theta(\cdot, \cdot)$  be defined as Equations 8 and 9 in Section 4.2. The mapping from  $(\Omega, \Sigma, \mu_{\mathbf{W}})$  to  $(\Omega, \Sigma)$  defined by  $\omega(\tau) \rightarrow F_\theta(\omega(\tau), \tau)$  is measurable and therefore induces a pushforward measure  $\mu_{\mathbf{W}} \circ F_\theta^{-1}$ .*

*Proof.* As  $F_\theta$  is continuous in both  $\omega$  and  $\tau$ , it is easy to show  $F_\theta(\omega(\tau), \tau)$  is also continuous in  $\tau$  for each  $\omega$  continuous in  $\tau$ . As  $F_\theta(\cdot, \tau)$  is invertible for each  $\tau$ ,  $F_\theta(\cdot, \tau)$  is an homeomorphism between  $\mathbb{R}^d$  and  $\mathbb{R}^d$ . Therefore, the pre-image of each Borel set of  $\mathbb{R}^d$  under  $F_\theta(\cdot, \tau)$  for each  $\tau$  is also Borel. As a result, the pre-image of each cylinder set of  $C([0, +\infty), \mathbb{R}^d)$  under the mapping defined by  $F_\theta(\cdot, \cdot)$  is also a cylinder set, which is enough to show the mapping is measurable.  $\square$

This proposition shows  $\mathbf{X}_\tau$  is a stochastic process also defined in the space of continuous functions as Wiener process. It provides a solid basis to for defining finite-dimensional distribution of  $\mathbf{X}_\tau$  on  $\mathbb{R}^{d \times n}$  in a similar ways as Wiener process using projection. The two sampling methods mentioned above can be characterized by two different mappings from  $(\Omega, \Sigma, \mu_{\mathbf{W}})$  to  $(\mathbb{R}^{d \times n}, \mathcal{B}(\mathbb{R}^{d \times n}))$ : (1) applying transformation defined by  $F_\theta$  to a function in  $C([0, +\infty), \mathbb{R}^d)$  and then applying the projection  $\pi$  to the transformed function given a time grid; (2) applying the projection to a continuous function on a time grid and applying the transformation defined by  $F_\theta(\cdot, \tau)$  for each  $\tau$  individually. We can check the pushforward measures induced by the two mappings agree on every Borel set of  $\mathbb{R}^{d \times n}$  as their pre-images are the same in  $(\Omega, \Sigma, \mu_{\mathbf{W}})$ . Therefore we have the following proposition:

**Proposition 2.** *Given a finite subset  $\{\tau_1, \tau_2, \dots, \tau_n\} \subset (0, +\infty)$ , the finite-dimensional distribution of  $\mathbf{X}_\tau$  is the same as the distribution of  $(F_\theta(\mathbf{W}_{\tau_1}, \tau_1), \dots, F_\theta(\mathbf{W}_{\tau_n}, \tau_n))$ , where  $(\mathbf{W}_{\tau_1}, \dots, \mathbf{W}_{\tau_n})$  is a  $n \times d$ -dimensional random variable with finite-dimensional distribution of  $\mathbf{W}_\tau$ .*

*Proof.* It suffices to check that given the fixed time grid, for each Borel set  $B \subset \mathbb{R}^{d \times n}$ , the preimage of  $B$  is the same under the two mappings. They are both  $\{\omega | (F_\theta(\mathbf{W}_{\tau_1}(\omega), \tau_1), F_\theta(\mathbf{W}_{\tau_2}(\omega), \tau_2), \dots, F_\theta(\mathbf{W}_{\tau_n}(\omega), \tau_n)) \in B\}$ .  $\square$

## B Experiment Setup and Model Architecture Details

We describe the details on synthetic dataset generation, real-world dataset pre-processing, model architecture as well as training and evaluation settings in this section.

### B.1 Synthetic Dataset Details

For the geometric Brownian motion (GBM), we sample 10000 trajectories from a GBM with the parameters of  $\mu = 0.2$  and a variance of  $\sigma = 0.5$  in the interval of  $[0, 30]$ . The timestamps of the

observations are sampled from a homogeneous Poisson point process with an intensity of  $\lambda_{\text{train}} = 2$ . We evaluate the model on the observations timestamps sampled from two homogeneous Poisson processes separately with intensity values of  $\lambda_{\text{test}} = 2$  and  $\lambda_{\text{test}} = 20$ .

For the Ornstein–Uhlenbeck (OU) process, the parameters of the process we sample trajectories from are  $\theta = 2$ ,  $\mu = 1$ , and  $\sigma = 10$ . We also sample 10000 trajectories and use the same set of observation intensity values,  $\lambda_{\text{train}}$  and  $\lambda_{\text{test}}$ , to sample observation timestamps from homogeneous Poisson processes for training and test.

For the mixture of OU processes (MOU), we sample 5000 sequences from each of two different OU processes and mix them to obtain 10000 sequences. One OU process has the parameters of  $\theta = 2$ ,  $\mu = 1$ , and  $\sigma = 10$  and the observation timestamps are sampled from a homogeneous Poisson process with  $\lambda_{\text{train}} = 2$ . The other OU process has the parameters of  $\theta = 1.0$ ,  $\mu = 2.0$ , and  $\sigma = 5.0$  with observation timestamps sampled with  $\lambda_{\text{train}} = 20$ .

For the 10000 trajectories of each dataset, we use 7000 trajectories for training and 1000 trajectories for validation. We test the model on 2000 trajectories for each value of  $\lambda_{\text{test}}$ . To test the model with  $\lambda_{\text{test}} = 20$  on GBM and OU process, we also use 2000 sequences.

## B.2 Real-World Dataset Details

As mentioned in Section 5.2 of the paper, we compare our models against the baselines on three datasets: Mujoco-Hopper, Beijing Air-Quality dataset (BAQD), and PTB Diagnostic Database(PTBDB). The three datasets can be downloaded using the following links:

- <http://www.cs.toronto.edu/~rtqichen/datasets/HopperPhysics/training.pt>
- <https://www.kaggle.com/shayanfazeli/heartbeat/download>
- <https://archive.ics.uci.edu/ml/datasets/Beijing+Multi-Site+Air-Quality+Data>

We pad all sequences into the same length for each dataset. The sequence length of the Mujoco-Hopper dataset is 200 and the sequence length of BAQD is 168. The maximum sequence length in the PTBDB dataset is 650. We rescale the indices of sequences to real numbers in the interval of  $[0, 120]$  and take the rescaled values as observation timestamps for all datasets. To make the sequences asynchronous or irregularly-sampled, we sample observation timestamps  $\{\tau_i\}_{i=1}^n$  from a homogeneous Poisson process with an intensity of 2 that is independent of the data. For each sampled timestamp, the value of the closest observation is taken as its corresponding value. The timestamps of all sampled sequences are shifted by a value of 0.2 since  $\mathbf{W}_0 = 0$  deterministically for the Wiener process and there’s no variance for the CTFP model’s prediction at  $\tau = 0$ .

## B.3 Model Architecture Details

To ensure a fair comparison, we use the same values for hyper-parameters including the latent variable and hidden state dimensions across all models. Likewise, we keep the underlying architectures as similar as possible and use the same experimental protocol across all models.

For CTFP and Latent CTFP, we use a one-block augmented neural ODE module that maps the base process to the observation process. For the augmented neural ODE model, we use an MLP model consisting of 4 hidden layers of size 32–64–64–32 for the model in Equation 8 and Equation 12. In practice, the implementation of  $g$  in the two equations is optional and its representation power can be fully incorporated into  $f$ . This architecture is used for both synthetic and real-world datasets. For the latent CTFP and latent ODE models appearing in Section 5, we use the ODE-RNN model as the recognition network. For synthetic datasets, the ODE-RNN model consists of a one-layer GRU cell with a hidden dimension of 20 (the rec-dims parameter in its original implementation) and a one-block neural ODE module that has a single hidden layer of size 100, and it outputs a 10-dimensional latent variable. The same architecture is used by both latent ODE and latent CTFP models. For real-world datasets, the ODE-RNN architecture uses a hidden state of dimension 20 in the GRU cell and an MLP with a 128-dimensional hidden layer in the neural ODE module. The ODE-RNN model produces a 64-dimensional latent variable. For the generation network of the latent ODE (V2) model, we use an ODE function with one hidden layer of size 100 for synthetic datasets

and 128 for real-world datasets. The decoder network has 4 hidden layers of size 32–64–64–32; it maps a latent trajectory to outputs of Gaussian distributions at different time steps.

The VRNN model is implemented using a GRU network. The hidden state of the VRNN models is 20-dimensional for synthetic and real-world datasets. The dimension of the latent variable is 64 for real-world datasets and 10 for synthetic datasets. We use an MLP of 4 hidden layers of size 32–64–64–32 for the decoder network, an MLP with one hidden layer that has the same dimension as the hidden state for the prior proposal network, and an MLP with two hidden layers for the posterior proposal network. For synthetic data sampled from Geometric Brownian Motion, we apply an exponential function to the samples of all models. Therefore the distribution predicted by latent ODE and VRNN at each timestamp is a log-normal distribution.

## B.4 Training and Evaluation Settings

For synthetic data, we train all models using the IWAE bound with 3 samples and a flat learning rate of  $5 \times 10^{-4}$  for all models. We also consider models trained with or without the aggressive training scheme proposed by He et al. [22] for latent ODE and latent CTFP. We choose the best-performing model among the ones trained with or without the aggressive scheme based IWAE bound, estimated with 25 samples on the validation set for evaluation. The batch size is 100 for CTFP models and 25 for all the other models. For experiments on real-world datasets, we did a hyper-parameter search on learning rates over two values of  $5 \times 10^{-4}$  and  $10^{-4}$ , and whether using the aggressive training schemes for latent CTFP and latent ODE models. We report the evaluation results of the best-performing model based on IWAE bound estimated with 125 samples.

## C Ablation Study Results

### C.1 Additional Experiment Results on Real-world Datasets

We provide additional experiment results on real-world datasets using different intensity value  $\lambda$ s of 1 and 5 to sample observation processes in Table 1 below.

Table 3: Ablation Study on Time Interval for Real-World Data

Model	Negative Log-Likelihood					
	Mujoco-Hopper		BAQD		PTBDB	
	$\lambda_{\text{test}} = 1$	$\lambda_{\text{test}} = 5$	$\lambda_{\text{test}} = 1$	$\lambda_{\text{test}} = 5$	$\lambda_{\text{test}} = 1$	$\lambda_{\text{test}} = 5$
Latent ODE	25.082 ± 0.011	24.599 ± 0.004	2.948 ± 0.006	2.686 ± 0.006	-0.633 ± 0.006	-0.892 ± 0.009
VRNN	10.553 ± 0.010	8.543 ± 0.008	0.044 ± 0.007	-1.016 ± 0.001	<b>-1.552 ± 0.011</b>	<b>-2.545 ± 0.005</b>
CTFP	-10.152 ± 0.084	-23.241 ± 0.057	-1.255 ± 0.022	-3.784 ± 0.035	-1.028 ± 0.028	-1.824 ± 0.014
Latent CTFP	<b>-30.469 ± 0.079</b>	<b>-33.412 ± 0.035</b>	<b>-7.276 ± 0.061</b>	<b>-6.226 ± 0.016</b>	<b>-1.552 ± 0.010</b>	<b>-2.533 ± 0.008</b>

### C.2 I.I.D. Gaussian as Base Process

In this experiment, we replace the base Wiener process with I.I.D Gaussian random variables and keep the other components of the models unchanged. This model and its latent variant are named CTFP-IID-Gaussian and latent CTFP-IID-Gaussian. As a result, the trajectories sampled from CTFP-IID-Gaussian are not continuous and we use this experiment to study the continuous property of models and its impact on modeling irregular time series data with continuous dynamics. The results are presented in Table 4 and Table 5.

The results show that CTFP consistently outperforms CTFP-IID-Gaussian, and latent CTFP outperforms latent CTFP-IID-Gaussian. The results corroborate our hypothesis that the superior performance of CTFP models can be partially attributed to the continuous property of the model. Moreover, latent CTFP-IID-Gaussian shows similar but slightly better performance than latent ODE models. The results comply with our hypothesis as the models are very similar and both models have no notion of continuity in the decoder. We believe the performance gain of latent CTFP-IID-Gaussian comes from the use of (dynamic) normalizing flow which is more flexible than Gaussian distributions used by latent ODE.

Table 4: Comparison between CTFP, CTFP-IID-Gaussian, latent CTFP, and latent CTFP-IID-Gaussian on synthetic datasets. We report NLL per observation.

Model	GBM		OU		M-OU
	$\lambda_{\text{test}} = 2$	$\lambda_{\text{test}} = 20$	$\lambda_{\text{test}} = 2$	$\lambda_{\text{test}} = 20$	$\lambda_{\text{test}} = (2, 20)$
Latent ODE [44]	3.826	5.935	3.066	3.027	2.690
CTFP-IID-Gaussian	4.952	4.094	3.025	3.024	2.716
Latent CTFP-IID-Gaussian	3.945	5.072	3.017	3.000	2.689
CTFP ( <b>ours</b> )	<b>3.107</b>	<b>1.929</b>	2.902	1.941	1.408
Latent CTFP ( <b>ours</b> )	<b>3.107</b>	1.930	2.902	<b>1.939</b>	<b>1.392</b>
Ground Truth	3.106	1.928	2.722	1.888	1.379

Table 5: Comparison Between CTFP, CTFP-IID-Gaussian, latent CTFP, and latent CTFP-IID-Gaussian on real-world datasets. We report NLL per observation.

Model	Mujoco-Hopper [44]	BAQD [4]	PTBDB [49]
Latent ODE [44]	$24.775 \pm 0.010$	$2.789 \pm 0.011$	$-0.818 \pm 0.009$
CTFP-IID-Gaussian	$22.023 \pm 0.010$	$3.398 \pm 0.006$	$-0.375 \pm 0.003$
Latent CTFP-IID-Gaussian	$17.397 \pm 0.007$	$1.471 \pm 0.005$	$-1.436 \pm 0.005$
CTFP ( <b>ours</b> )	$-16.249 \pm 0.034$	$-2.361 \pm 0.020$	$-1.324 \pm 0.028$
Latent CTFP ( <b>ours</b> )	<b><math>-31.397 \pm 0.063</math></b>	<b><math>-6.894 \pm 0.046</math></b>	<b><math>-1.999 \pm 0.010</math></b>

### C.3 CTFP-RealNVP

In this experiment, we replace the continuous normalizing flow in CTFP model with another popular choice of normalizing flow model, RealNVP [14]. This variant of CTFP is named CTFP-RealNVP and its latent version is called latent CTFP-RealNVP. Note that the trajectories sampled from CTFP-RealNVP model are still continuous. We evaluate CTFP-RealNVP and latent CTFP-RealNVP models on datasets with high dimensional data, Mujoco-Hopper, and BAQD. The results are shown in Table 6.

Table 6: Comparison between CTFP, CTFP-RealNVP, and their latent variants on Mujoco-Hopper and BAQD datasets. We report NLL per observation.

Model	Mujoco	BAQD
CTFP-RealNVP	$-23.061 \pm 0.000$	$-5.099 \pm 0.002$
Latent CTFP-RealNVP	$-23.602 \pm 0.001$	$-5.109 \pm 0.005$
CTFP	$-16.249 \pm 0.034$	$-2.361 \pm 0.020$
Latent CTFP	<b><math>-31.397 \pm 0.063</math></b>	<b><math>-6.894 \pm 0.046</math></b>

The table indicates that CTFP-RealNVP outperforms CTFP. However, when incorporating the latent variable, the latent CTFP-RealNVP performs significantly worse than latent CTFP. The worse performance might be because RealNVP cannot make full use of the information in the latent variable due to its structural constraints as we discussed in Section 4.2.

## D Additional Details for Latent ODE Models on Mujoco-Hopper Data

The original latent ODE paper focuses on point estimation and uses the mean squared error as the performance metric [44]. When applied to our problem setting and evaluated using the log-likelihood, the model performs unsatisfactorily. In Table 7, the first row shows the negative log-likelihood on the Mujoco-Hopper dataset. The inferior NLL of the original latent ODE is potentially caused by the use of a fixed output variance of  $10^{-6}$ , which magnifies even a small reconstruction error.

To mitigate this issue, we propose two modified versions of the latent ODE model. For the first version (V1), given a pretrained (original) latent ODE model, we do a logarithmic scale search for the output variance and find the value that gives the best performance on the validation set. The second version (V2) uses an MLP to predict the output mean and variance. Both modified versions have much better performance than the original model, as shown in Table 7, rows 2–3. It also shows that the

second version of the latent ODE model (V2) outperforms the first one (V1) on the Mujoco-Hopper dataset. Therefore, we use the second version (V2) for all the experiments in the main text.

Table 7: Comparison of different version of latent ODE models on Mujoco-Hopper Datasets.

Model	NLL
Latent ODE (original)	$4 \times 10^7 \pm 9 \times 10^5$
Latent ODE (V1)	$45.874 \pm 0.001$
Latent ODE (V2)	$24.775 \pm 0.010$
VRNN	$9.113 \pm 0.018$
CTFP	$-16.249 \pm 0.034$
Latent CTFP	<b><math>-31.397 \pm 0.063</math></b>

## E Qualitative Sample for VRNN Model

We sample trajectories from the VRNN model [12] trained on Geometric Brownian Motion (GBM) by running the model on a dense time grid and show the trajectories in Figure 4. We compare the trajectories sampled from the model with trajectories sampled from GBM. As we can see, the sampled trajectories from VRNN are not continuous in time.

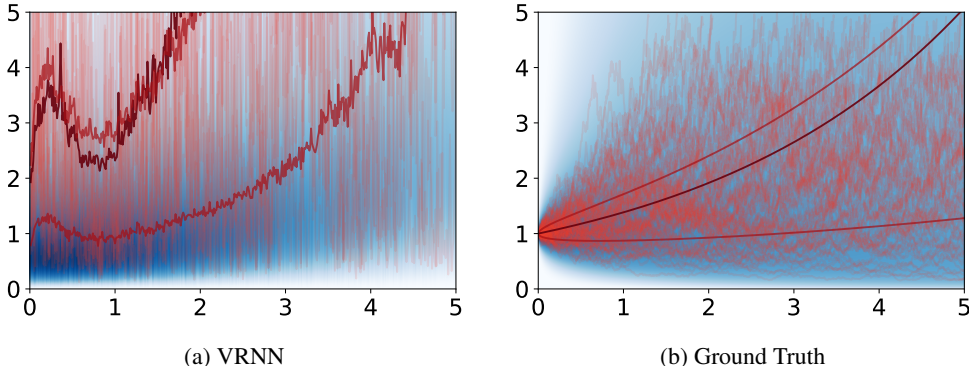


Figure 4: Sample trajectories and marginal density estimation by VRNN (a). We compare the results with sample trajectories and marginal density with ground truth (b). In addition to the sample trajectories (red) and the marginal density (blue), we also show the sample-based estimates (closed-form for ground truth) of the inter-quartile range (dark red) and mean (brown) of the marginal density.

We also use VRNN to estimate the marginal density of  $\mathbf{X}_\tau$  for each  $\tau \in (0, 5]$  and show the results in Figure 4. It is not straightforward to use VRNN model for marginal density estimation. For each timestamp  $\tau \in (0, 5]$ , we get the marginal density of  $\mathbf{X}_\tau$  by running VRNN on a time grid with two timestamps, 0 and  $\tau$ : at the first step, the input to VRNN model is  $\mathbf{x}_0 = 1$  and we can get prior distributions of the latent variable  $\mathbf{Z}_\tau$ . Note that a sampled trajectory from GBM is always 1 when  $\tau = 0$ . Conditioned on the sampled latent codes  $\mathbf{z}_0$  and  $\mathbf{z}_\tau$ , VRNN proposes  $p(\mathbf{x}_\tau | \mathbf{x}_0, \mathbf{z}_\tau, \mathbf{z}_0)$  at the second step. We average the conditional density over 125 samples of  $\mathbf{Z}_\tau$  and  $\mathbf{Z}_0$  to estimate the marginal density.

The marginal density estimated using a time grid with two timestamps is not consistent with the trajectories sampled on a different dense time grid. The results indicate that the choice of time grid has a great impact on the distribution modeled by VRNN and the distributions modeled by VRNN on different time grids can be inconsistent. In contrast, our proposed CTFP models do not have such problems.