

---

# Learning interpretable continuous-time models of latent stochastic dynamical systems

---

Lea Duncker<sup>1</sup> Gergő Böhner<sup>1</sup> Julien Boussard<sup>2</sup> Maneesh Sahani<sup>1</sup>

## Abstract

We develop an approach to learn an interpretable semi-parametric model of a latent continuous-time stochastic dynamical system, assuming noisy high-dimensional outputs sampled at uneven times. The dynamics are described by a nonlinear stochastic differential equation (SDE) driven by a Wiener process, with a drift evolution function drawn from a Gaussian process (GP) conditioned on a set of learnt fixed points and corresponding local Jacobian matrices. This form yields a flexible nonparametric model of the dynamics, with a representation corresponding directly to the interpretable portraits routinely employed in the study of nonlinear dynamical systems. The learning algorithm combines inference of continuous latent paths underlying observed data with a sparse variational description of the dynamical process. We demonstrate our approach on simulated data from different nonlinear dynamical systems.

## 1. Introduction

Many dynamical systems with intrinsic noise may be modelled in continuous time using the framework of stochastic differential equations (SDE). However identifying a good SDE model from intermittent observations of the process is challenging, particularly if the dynamical process is nonlinear and the observations are indirect and noisy. A common response is to assume a latent process that operates in discretised time, often called a state-space model. This approach has been applied in contexts ranging from modelling human motion (Wang et al., 2006) to solving control problems (Eleftheriadis et al., 2017). However, it assumes that observations, and the critical phenomena of the dynamics, can be

---

<sup>1</sup>Gatsby Computational Neuroscience Unit, University College London, London, United Kingdom <sup>2</sup>Stanford University, Palo Alto, California, USA. Correspondence to: Lea Duncker <duncker@gatsby.ucl.ac.uk>.

accurately modelled using a discrete time grid.

A further challenge when the goal is to gain insight into a physical or biological system whose parametric description is unknown, is to obtain an *interpretable* model of the dynamics from observed data, whether modelled in discrete or continuous time. State-space models that rely on nonparametric or flexibly parametrised descriptions of dynamics, for example using Gaussian process (GP) priors or recurrent neural networks (RNN), may be effective at prediction but inevitably leave interpretation to a second analytic stage, posing its own challenges.

In this paper, we consider continuous-time latent SDE models of the form

$$d\mathbf{x} = \mathbf{f}(\mathbf{x})dt + \sqrt{\Sigma} d\mathbf{w}$$
$$\mathbb{E}_{y|\mathbf{x}}[\mathbf{y}(t_i)] = g(\mathbf{C}\mathbf{x}(t_i) + \mathbf{d}), \quad i = 1, \dots, T, \quad (1)$$

where the temporal evolution of a latent variable  $\mathbf{x} \in \mathbb{R}^K$  is described by a nonlinear SDE with dynamical evolution function  $\mathbf{f} : \mathbb{R}^K \mapsto \mathbb{R}^K$  and incremental noise covariance  $\Sigma$  shaping the Wiener noise process  $\mathbf{w}(t)$ . Note that the nonlinear SDE induces a non-Gaussian prior on  $\mathbf{x}(t)$  with no easy access to finite marginal distributions. The latent state is observed indirectly through noisy measurements  $\mathbf{y}_i \in \mathbb{R}^N$  at unevenly spaced time points  $t_i$ . The measurements are distributed with a known parametric form and generalized linear dependence; that is the expected value is  $g(\mathbf{C}\mathbf{x} + \mathbf{d})$  with inverse-link function  $g$  and parameters  $\mathbf{C} \in \mathbb{R}^{N \times K}$  and  $\mathbf{d} \in \mathbb{R}^N$ . We seek to infer latent paths  $\mathbf{x}(t)$  along with the dynamical parameters and an interpretable representation of the dynamical mapping  $\mathbf{f}$ .

What do we mean by interpretable? The properties of dynamical systems are frequently analyzed by characterizing dynamical fixed points and local behaviour near these points (Sussillo & Barak, 2013). When  $\mathbf{f}$  is a learnt, general function, fixed points must be found numerically (Golub & Sussillo, 2018). This makes it difficult to propagate uncertainty about  $\mathbf{f}$  to the number and location of fixed points, and to the local dynamics around them. Our approach is to develop a non-parametric Gaussian-process model for  $\mathbf{f}$  conditioned on the learnt locations of fixed points and associated local Jacobians. Thus, we implicitly integrate out the details of

$\mathbf{f}$ , while optimising directly over the components of the interpretable dynamical portrait.

The paper is organised as follows: In section 2 we review background material on the related Gaussian Process State-Space Model (GP-SSM) and previous work on GP approximations to SDEs (Archambeau et al., 2007; 2008). We also briefly review the inducing point approach for GP models. In section 3 we make use of GP priors to represent the unknown nonlinear dynamics  $\mathbf{f}$ , incorporating interpretable structure by conditioning the GP on fixed points and local Jacobian matrices of the system. We derive a Variational Bayes algorithm for approximate inference and parameter learning in section 4. Finally, we demonstrate the performance of our algorithm on a number of nonlinear dynamical system examples in section 5.

## 2. Background

### 2.1. Gaussian Process State-Space-Model

A discrete-time analogue of the model in (1) is the GP-SSM, where the latent state evolution over a fixed step size is modelled as

$$\mathbf{x}_{\ell+1} = \mathbf{f}(\mathbf{x}_\ell) + \boldsymbol{\epsilon}_\ell \quad (2)$$

where  $\boldsymbol{\epsilon}_\ell \sim \mathcal{N}(\boldsymbol{\epsilon}_\ell | 0, D)$ . There have been a range of approaches for performing approximate inference in this model, based on Assumed Density Filtering (Deisenroth et al., 2009; Ramakrishnan et al., 2011), Expectation Propagation (Deisenroth & Mohamed, 2012), variational inference (Frigola et al., 2014), or recurrent recognition networks (Eleftheriadis et al., 2017). The model in (1) requires a different treatment for latent path inference, as it maintains the continuous-time structure of the system of interest.

### 2.2. Gaussian Process Approximation to SDEs

The problem of performing approximate inference in continuous-time SDE models has been considered previously, with the two main approaches being Expectation Propagation (Cseke et al., 2016) and variational inference (Archambeau et al., 2007; 2008). We review the latter approach in this section, as our Variational Bayes algorithm in section 4 extends this work.

Archambeau et al. (2007; 2008) consider the model in (1) under linear Gaussian observations. The authors derive an approximate inference algorithm based on a variational Gaussian approximation to the posterior process on  $\mathbf{x}(t)$  under the constraint that the approximate process has Markov structure, as is the case for the true posterior process. The most general way to construct such an approximation is via a linear time-varying SDE of the form

$$d\mathbf{x} = (-\mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t)) dt + \sqrt{\boldsymbol{\Sigma}} d\mathbf{w} \quad (3)$$

The instantaneous marginal distributions of this approximation at any time  $t$  are Gaussian, with means  $\mathbf{m}_x(t)$  and covariances  $\mathbf{S}_x(t)$  that evolve in time according to the ordinary differential equations (ODEs):

$$\begin{aligned} \frac{d\mathbf{m}_x}{dt} &= -\mathbf{A}(t)\mathbf{m}_x + \mathbf{b}(t) \\ \frac{d\mathbf{S}_x}{dt} &= -\mathbf{A}(t)\mathbf{S}_x - \mathbf{S}_x\mathbf{A}(t)^\top + \boldsymbol{\Sigma} \end{aligned} \quad (4)$$

Archambeau et al. (2007; 2008) derive a lower bound to the marginal log-likelihood – often called the variational free energy or evidence lower bound – whose maximisation with respect to  $q_x$  is equivalent to minimising the Kullback-Leibler (KL) divergence between the approximate and true posterior process. The free energy has the form

$$\mathcal{F} = \sum_i \langle \log p(\mathbf{y}_i | \mathbf{x}_i) \rangle_{q_x} - \text{KL}[q_x(\mathbf{x}) \| p(\mathbf{x})] \quad (5)$$

The first term is the expected log-likelihood under the approximation and only depends on the marginal distributions  $q_x(\mathbf{x}(t))$ . The second term is the KL-divergence between the continuous-time approximate posterior process and the prior process. Archambeau et al. (2007) show that this term can be written as

$$\text{KL}[q_x(\mathbf{x}) \| p(\mathbf{x})] = \int_{\mathcal{T}} dt \langle (\mathbf{f} - \mathbf{f}_q)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{f} - \mathbf{f}_q) \rangle_q \quad (6)$$

where both  $\mathbf{f}$  and  $\mathbf{f}_q$  are evaluated at  $\mathbf{x}(t)$ , and  $\mathbf{f}_q(\mathbf{x}(t)) = -\mathbf{A}(t)\mathbf{x}(t) + \mathbf{b}(t)$ . Note that the noise covariance  $\boldsymbol{\Sigma}$  is deliberately chosen to be equal for the SDEs in  $q_x$  and  $p$ , as this term would diverge otherwise.

To maximise  $\mathcal{F}$  with respect to  $\mathbf{m}_x(t)$  and  $\mathbf{S}_x(t)$ , subject to the constraint that the approximate posterior process has Markov structure according to equation (3), one can find the stationary points of the Lagrangian

$$\mathcal{L} = \mathcal{F} - \mathcal{C}_1 - \mathcal{C}_2 \quad (7)$$

with

$$\begin{aligned} \mathcal{C}_1 &= \int_{\mathcal{T}} dt \text{Tr} \left[ \Psi \left( \frac{d\mathbf{S}_x}{dt} + \mathbf{A}\mathbf{S}_x + \mathbf{S}_x\mathbf{A}^\top - \boldsymbol{\Sigma} \right) \right] \\ \mathcal{C}_2 &= \int_{\mathcal{T}} dt \boldsymbol{\lambda}^\top \left( \frac{d\mathbf{m}_x}{dt} + \mathbf{A}\mathbf{m}_x - \mathbf{b} \right) \end{aligned} \quad (8)$$

where  $\Psi$  and  $\boldsymbol{\lambda}$  are Lagrange multipliers. Archambeau et al. (2007; 2008) derive a smoothing algorithm that involves iterating fixed point updates of this Lagrangian. These are either closed form, or require solving ODEs forward and backward in time, thus achieving linear time complexity. In section 4, we will modify this original algorithm in order to improve its numerical stability, and show how to incorporate it in an efficient Variational Bayes algorithm.

### 2.3. Sparse Gaussian Processes using inducing points

In later sections of the paper, we will make use of the sparse variational inducing point approach of Titsias (2009). The key idea of inducing point approaches is to condition a GP  $\zeta(\mathbf{x}) \sim \mathcal{GP}(0, \kappa(\mathbf{x}, \mathbf{x}'))$  on so-called inducing variables  $\mathbf{u} \in \mathbb{R}^M$ . These can be thought of as pseudo-observations of the function at  $M$  locations  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_M] \in \mathbb{R}^{K \times M}$ . An augmented prior for the GP and inducing variables can be written as

$$\mathbf{u} \sim \mathcal{N}(\mathbf{u}|0, \mathbf{K}_{zz}), \zeta|\mathbf{u} \sim \mathcal{GP}(\mu_{\zeta|\mathbf{u}}(\mathbf{x}), \nu_{\zeta|\mathbf{u}}(\mathbf{x}, \mathbf{x}')) \quad (9)$$

The conditioned GP mean and covariance function are

$$\begin{aligned} \mu_{\zeta|\mathbf{u}}(\mathbf{x}) &= \mathbf{K}_{\cdot z}(\mathbf{x})\mathbf{K}_{zz}^{-1}\mathbf{u} \\ \nu_{\zeta|\mathbf{u}}(\mathbf{x}, \mathbf{x}') &= \kappa(\mathbf{x}, \mathbf{x}') - \mathbf{K}_{\cdot z}(\mathbf{x})\mathbf{K}_{zz}^{-1}\mathbf{K}_{z \cdot}(\mathbf{x}) \end{aligned} \quad (10)$$

Where  $[\mathbf{K}_{zz}]_{ij} = \kappa(\mathbf{z}_i, \mathbf{z}_j)$ , and  $[\mathbf{K}_{\cdot z}(\mathbf{x})]_i = \kappa(\mathbf{x}, \mathbf{z}_i)$ . The computational complexity of building the mean and covariance in (10) is linear in the number of  $\mathbf{x}$  input points and cubic only in the number of inducing points  $M$ . If we were to integrate over the inducing variables in this augmented prior, we would recover the original model. However, the inducing variables can also be kept in the model as auxiliary variables, which may be incorporated into approaches for variational inference (Titsias, 2009).

### 3. Interpretable priors on nonlinear dynamics

Similarly to the GP-SSM work, we wish to model  $\mathbf{f}$  using the framework of GPs. GPs can represent a flexible class of nonlinear dynamics. However, it may be difficult to interpret the inferred function with respect to studying the underlying dynamical system that generated the observed data. As stated above, standard analysis approaches for nonlinear dynamical systems rely on identifying local fixed points  $\mathbf{s}_i$ , where  $\mathbf{f}(\mathbf{s}_i) = \mathbf{0}$ , and the locally-linearised dynamics around them, given by the Jacobians  $\nabla_{\mathbf{x}}\mathbf{f}(\mathbf{x})|_{\mathbf{x}=\mathbf{s}_i}$  (Sussillo & Barak, 2013). This strategy motivates our approach to interpretability.

#### 3.1. A Gaussian Process prior for dynamics

In order to arrive at a modelling framework that makes fixed points and Jacobian matrices readily available for analysis, we introduce a GP prior conditioned directly on these parameters (see also Bohner & Sahani, 2018). The fixed point locations and Jacobians around them can be viewed as further hyperparameters specifying the prior mean and covariance function of the GP, which we will denote by  $\boldsymbol{\theta} = \{\mathbf{f}_s^{(i)}, \mathbf{J}_s^{(i)}\}_{i=1}^L$ , where  $L$  denotes the total number of fixed point locations. With  $\mathbf{f}_s^{(i)} = \mathbf{f}(\mathbf{s}_i) = \mathbf{0}$  and  $[\mathbf{J}_s^{(i)}]_{k,m} = \frac{\partial f_k(\mathbf{x})}{\partial x_m}|_{\mathbf{x}=\mathbf{s}_i}$ . We can hence write a GP prior conditioned on the fixed points and Jacobians for each di-

mension in  $\mathbf{f}$ , using the fact that a GP and its derivative process are still jointly distributed as a GP.

The Variational Bayes approach in section 4 will make use of a sparse variational approximation for  $\mathbf{f}$  using inducing variables, as in Titsias (2009). To make later notation more compact, we therefore directly introduce the augmented model including inducing variables drawn from the conditioned GP prior here. We denote the joint covariance matrix between inducing variables, fixed points and Jacobians as

$$\mathbf{K}_{zz}^{\boldsymbol{\theta}} = \begin{bmatrix} \mathbf{K}_{zz} & \mathbf{K}_{zs} & \mathbf{K}_{zs}^{\nabla_2} \\ \mathbf{K}_{sz} & \mathbf{K}_{ss} & \mathbf{K}_{ss}^{\nabla_2} \\ \mathbf{K}_{sz}^{\nabla_1} & \mathbf{K}_{ss}^{\nabla_1} & \mathbf{K}_{ss}^{\nabla_1 \nabla_2} \end{bmatrix} = \begin{bmatrix} \mathbf{K}_{zz} & \tilde{\mathbf{K}}_{zs} \\ \tilde{\mathbf{K}}_{sz} & \tilde{\mathbf{K}}_{ss} \end{bmatrix} \quad (11)$$

where the superscript  $\nabla_i$  denotes the derivative of the covariance function with respect to its  $i$ th input argument such that  $[\mathbf{K}_{ss}^{\nabla_2}]_{ij} = \frac{\partial}{\partial \mathbf{s}} \kappa(\mathbf{z}_i, \mathbf{s})|_{\mathbf{s}=\mathbf{s}_j}$ . The conditional prior on the inducing variables given  $\boldsymbol{\theta}$  can then be written as

$$\mathbf{u}_k|\boldsymbol{\theta} = \mathcal{N}\left(\mathbf{u} \left| \tilde{\mathbf{K}}_{zs}\tilde{\mathbf{K}}_{ss}^{-1}\mathbf{v}_k^{\boldsymbol{\theta}}, \mathbf{K}_{zz} - \tilde{\mathbf{K}}_{zs}\tilde{\mathbf{K}}_{ss}^{-1}\tilde{\mathbf{K}}_{sz} \right.\right) \quad (12)$$

where  $\mathbf{v}_k^{\boldsymbol{\theta}} = [f_{s,k}^{(1)}, \dots, f_{s,k}^{(L)}, \mathbf{J}_{k,\cdot}^{(1)}, \dots, \mathbf{J}_{k,\cdot}^{(L)}]^{\top}$  collects the fixed-point and derivative observations relating to  $f_k$ . Finally, for the conditional prior on  $f_k$ , given the inducing variables and  $\boldsymbol{\theta}$ , we have

$$f_k|\mathbf{u}_k, \boldsymbol{\theta} \sim \mathcal{GP}\left(\mu_{f|u}^{\boldsymbol{\theta}}(\mathbf{x}), \nu_{f|u}^{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}')\right) \quad (13)$$

with

$$\mu_{f|u}^{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{a}_z^{\boldsymbol{\theta}}(\mathbf{x}) \begin{bmatrix} \mathbf{u}_k \\ \mathbf{v}_k^{\boldsymbol{\theta}} \end{bmatrix} \quad (14)$$

$$\nu_{f|u}^{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}') = \kappa(\mathbf{x}, \mathbf{x}') - \mathbf{a}_z^{\boldsymbol{\theta}}(\mathbf{x})\mathbf{K}_{zz}^{\boldsymbol{\theta}}\mathbf{a}_z^{\boldsymbol{\theta}}(\mathbf{x}')^{\top}$$

$$\mathbf{a}_z^{\boldsymbol{\theta}}(\mathbf{x}) = [\mathbf{K}_{\cdot z}(\mathbf{x}) \quad \mathbf{K}_{\cdot s}(\mathbf{x}) \quad \mathbf{K}_{\cdot s}^{\nabla_2}(\mathbf{x})] \mathbf{K}_{zz}^{\boldsymbol{\theta}}^{-1} \quad (15)$$

#### 3.2. Automatic selection of the number of fixed points

When the generative SDE dynamics are unknown, so are the number of fixed points in the system. We therefore take the general approach of introducing more fixed points than expected, and ‘pruning’ by hyperparameter optimisation. In particular, we include noise variance parameters for each fixed point, representing uncertainty about the zero-value of the function at the fixed point location. We hence have

$$\mathbf{f}_i^s = \mathbf{f}(\mathbf{s}_i) + \alpha_i \boldsymbol{\epsilon} = \mathbf{0} + \alpha_i \boldsymbol{\epsilon} \quad (16)$$

with  $\boldsymbol{\epsilon} \sim \mathcal{N}(0, I)$ . The variance parameters  $\alpha_i$  will enter our model simply via an added diagonal matrix to the  $\mathbf{K}_{ss}$  block in (11). When the  $\alpha_i$  are optimised, the uncertainty for superfluous fixed points will grow, while that of the fixed points the system is actually using will shrink. When the uncertainty for a fixed point is large, conditioning on it in the GP prior for  $\mathbf{f}$  will have negligible effect on prediction.

## 4. Variational inference and learning

We can derive an efficient Variational Bayes (VB) algorithm (Attias, 2000) for variational inference and learning in the model in (1) by maximising a variational free energy. We assume that our full variational distribution factorises as

$$q(\mathbf{x}, \mathbf{f}, \mathbf{u}) = q_x(\mathbf{x})q_{f,u}(\mathbf{f}, \mathbf{u}) \quad (17)$$

Following Titsias (2009), we choose  $q_{f,u}(\mathbf{f}, \mathbf{u}) = \prod_{k=1}^K p(f_k | \mathbf{u}_k, \boldsymbol{\theta}) q_u(\mathbf{u}_k)$ . The variational approximation of the posterior over the inducing variables are chosen to be of the form  $q_u(\mathbf{u}_k) = \mathcal{N}(\mathbf{u}_k | \mathbf{m}_u^k, \mathbf{S}_u^k)$ . The marginal variational distribution  $q_f(\mathbf{f}) = \prod_k \int d\mathbf{u}_k p(f_k | \mathbf{u}_k, \boldsymbol{\theta}) q_u(\mathbf{u}_k)$  is also a GP. The resulting expression for the variational free energy is of the form:

$$\mathcal{F}^* = \langle \mathcal{F} \rangle_{q_f} - \sum_{k=1}^K \text{KL}[q_u(\mathbf{u}_k) || p(\mathbf{u}_k | \boldsymbol{\theta})] \quad (18)$$

The VB algorithm iterates over an inference step, where the distribution  $q_x$  over the latent path is updated, a learning step, where  $q_{f,u}$  and the parameters in the output mapping are updated, and a hyperparameter learning step, where the kernel hyperparameters and fixed point locations are updated.

### 4.1. Inference

Our inference approach extends the work of Archambeau et al. (2007; 2008) to a wider class of observation models and to a nonparametric Bayesian treatment of the dynamics  $\mathbf{f}$  under the conditioned sparse GP prior introduced in section 3.

After using integration by parts on the Lagrangian in (7) (exchanging  $\mathcal{F}$  for  $\mathcal{F}^*$ ), we take variational derivatives with respect to  $\mathbf{m}_x(t)$  and  $\mathbf{S}_x(t)$ . Since our model has a rotational non-identifiability with respect to the latents  $\mathbf{x}$ , we fix  $\boldsymbol{\Sigma} = \mathbf{I}$  without loss of generality. We arrive at the following set of fixed point equations:

$$\frac{d\boldsymbol{\Psi}}{dt} = \mathbf{A}(t)^\top \boldsymbol{\Psi}(t) + \boldsymbol{\Psi}(t) \mathbf{A}(t) - \frac{\partial \mathcal{F}^*}{\partial \mathbf{S}_x} \odot \mathbb{P} \quad (19)$$

$$\frac{d\boldsymbol{\lambda}}{dt} = \mathbf{A}(t)^\top \boldsymbol{\lambda}(t) - \frac{\partial \mathcal{F}^*}{\partial \mathbf{m}_x} \quad (20)$$

$$\mathbf{A}(t) = \left\langle \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right\rangle_{q_x q_f} + 2\boldsymbol{\Psi}(t) \quad (21)$$

$$\mathbf{b}(t) = \langle \mathbf{f}(\mathbf{x}) \rangle_{q_x q_f} + \mathbf{A}(t) \mathbf{m}_x(t) - \boldsymbol{\lambda}(t) \quad (22)$$

with  $\mathbb{P}_{ij} = \frac{1}{2}$  for  $i \neq j$  and 1 otherwise and  $\odot$  denotes the Hadamard product. In contrast to previous work, we explicitly take the symmetric variations of  $\mathbf{S}_x(t)$  into account, which leads to slightly modified equations in (19) compared to the work in Archambeau et al. (2007; 2008), and thus to

improved numerical stability in practice. As a result, we can work with the fixed point updates (21) and (22) directly, without introducing a learning rate parameter that blends the updates with the previous value of the variational parameters  $\mathbf{A}$  and  $\mathbf{b}$ , as was done by Archambeau et al. (2007; 2008).

The inference algorithm involves solving the set of coupled ODEs in (4) and (19)-(22) using the conditions  $\mathbf{m}_x(0) = \mathbf{m}_{x,0}$ ,  $\mathbf{S}_x(0) = \mathbf{S}_{x,0}$  and  $\boldsymbol{\lambda}(T) = 0$ ,  $\boldsymbol{\Psi}(T) = 0$ . In principle, it is possible to use any ODE solver to do this. In this work, we choose to solve (4) using the forward Euler method with fixed step size  $\Delta t$  to obtain  $\mathbf{m}_x$  and  $\mathbf{S}_x$  evaluated on an evenly spaced grid. Similarly, we then solve (19) and (20) backwards in time to obtain evaluations of  $\boldsymbol{\lambda}$  and  $\boldsymbol{\Psi}$ . The solutions from the ODEs can then be used with equations (21) and (22) to obtain evaluations of  $\mathbf{A}$  and  $\mathbf{b}$  on the same time-grid used for solving the ODEs.

Evaluating the expectations of the terms involving  $\mathbf{f}$  with respect to  $q_x$  and  $q_f$  only involves computing Gaussian expectations of covariance functions and their derivatives. These can be computed analytically for choices such as an exponentiated quadratic covariance function. We update the initial state values  $\mathbf{m}_{x,0}$  and  $\mathbf{S}_{x,0}$  using the same procedure as that described by Archambeau et al. (2008). Given the function evaluations on the inference time-grid, we use linear interpolation to obtain function evaluations of  $\mathbf{m}_x$  and  $\mathbf{S}_x$  at arbitrary time points. Further details on the inference algorithm are given in the supplementary material.

### 4.2. Learning

#### 4.2.1. DYNAMICS

The only terms in (18) that depend on parameters in  $\mathbf{f}$  are the expected KL-divergence between the prior and approximate posterior processes and the KL-divergence relating to the inducing variables for  $\mathbf{f}$ , which are jointly quadratic in the inducing variables and Jacobians. Thus, given  $\mathbf{m}_x(t)$ ,  $\mathbf{S}_x(t)$ ,  $\mathbf{A}(t)$  and  $\mathbf{b}(t)$ , we can find closed form updates for the Jacobians and variational parameters relating to  $\mathbf{f}$ . For  $\mathbf{S}_u^k$  the update is of the form

$$\mathbf{S}_u^k = \left( \boldsymbol{\Omega}_u^{-1} + \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{[u,u]} \right)^{-1} \quad (23)$$

with  $\boldsymbol{\Omega}_u = \mathbf{K}_{zz} - \tilde{\mathbf{K}}_{zs} \tilde{\mathbf{K}}_{ss}^{-1} \tilde{\mathbf{K}}_{sz}$  and where the operation  $[X]_{[u,u]}$  selects the first  $M \times M$  block of  $\mathbf{X}$ . The inducing variable means and Jacobians around the fixed-point locations can be updated jointly as

$$\begin{bmatrix} \mathbf{m}_u^1 & \cdots & \mathbf{m}_u^K \\ \mathbf{J}_1 & \cdots & \mathbf{J}_K \end{bmatrix} = \mathbf{B}_1^{-1} (\mathbf{B}_2 - \mathbf{B}_3) \quad (24)$$



with

$$\begin{aligned}
 \mathbf{B}_1 &= \left( \tilde{\Omega} + \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{[uj,uj]} \right) \\
 \mathbf{B}_2 &= \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{[:,uj]}^\top \langle \mathbf{f}_q \rangle_{q_x}^\top \\
 \mathbf{B}_3 &= \int_{\mathcal{T}} dt [\langle \nabla_x \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{[:,uj]}^\top \mathbf{S}_x \mathbf{A}^\top \\
 \tilde{\Omega} &= \begin{bmatrix} \Omega_u^{-1} & -\Omega_u^{-1} \mathbf{G} \\ -\mathbf{G}^\top \Omega_u^{-1} & \mathbf{G}^\top \Omega_u^{-1} \mathbf{G} \end{bmatrix}, \mathbf{G} = [\tilde{\mathbf{K}}_{zs} \tilde{\mathbf{K}}_{ss}^{-1}]_{[j,j]}
 \end{aligned}$$

where  $[X]_{[uj,uj]}$  selects the first  $M \times M$  and last  $LK \times LK$  block of  $X$ ,  $[X]_{[:,uj]}$  selects the first  $M$  and last  $LK$  columns of  $X$ , and  $[X]_{[j,j]}$  selects the last  $LK \times LK$  block of  $X$ . The one-dimensional integrals can be computed efficiently using, for instance, Gauss-Legendre quadrature. Detailed derivations are given in the supplementary material, where we also provide closed form updates for the sparse variational GP approach for modelling  $\mathbf{f}$  without further conditioning on fixed points and Jacobians.

#### 4.2.2. OUTPUT MAPPING AND HYPERPARAMETERS

The only term that depends on the parameters  $\mathbf{C}$  and  $\mathbf{d}$  in (18) is the expected log-likelihood. Depending on the choice of likelihood, the optimal update may be available in closed form. Otherwise, parameter updates may be found by direct optimisation of the variational free energy. Similarly, the covariance function hyperparameters and fixed point locations are learnt by maximising the variational free energy. The inducing point locations can also be included here, though we chose to hold them fixed on a grid for all examples shown in this paper.

#### 4.3. Computational Complexity

The main costs of the algorithm come from evaluating GP predictions at a set of input points, and solving the ODEs (19)-(20). Computing the GP predictions using the basic sparse inducing-point approach scales cubically in the number of inducing points, the number of fixed points and the number of entries in the Jacobians, but scales linearly in the number of input points. Solving the ODEs using simple forward Euler integration achieves linear time complexity. In principle, adaptive ODE solvers could achieve a lower cost. Similarly, recent advances in scalable sparse Gaussian Process methods could improve on the cubic dependence on the inducing points. Hence, the cubic cost relating to the candidate number of fixed points could be viewed as the intrinsic cost of our description of the dynamics, while the costs relating to integrating ODEs and computing sparse Gaussian Process predictions could be improved.

## 5. Experiments

In this section, we apply our algorithm to data generated from different nonlinear dynamical systems. In all experiments, we choose an exponentiated quadratic covariance function in the prior over the dynamics  $\mathbf{f}$  and initialise the inducing point means and Jacobian matrices at zero. Each fixed point observation’s uncertainty is initialised with a standard deviation of 0.1. We generate  $\mathbf{C}$  and  $\mathbf{d}$  by drawing their entries from Gaussian distributions unless otherwise stated, and initialise our algorithm at these parameter values. For inference, we solve the ODEs (19)-(22) using the forward Euler method with  $\Delta t = 1\text{ms}$ . Unless stated otherwise, the link function is the identity  $g(z) = z$ .

### 5.1. Double-well dynamics

We first demonstrate our method on the classic one-dimensional double-well example, where the latent SDE evolves with drift  $f(x) = 4x(1 - x^2)$ . We simulate data on 20 trials with multivariate Gaussian outputs of dimensionality  $N = 15$  with unknown variances 0.25, and observe the output process at 20 randomly sampled time-points per trial. We chose 8 evenly spaced inducing points in  $(-3, 3)$  for  $f$ . While the true dynamics have three fixed points, we condition the prior on  $f$  on four fixed points and use the method outlined in section 3.2 to automatically select the correct number. The results are summarised in Figure 1, demonstrating that our algorithm can successfully perform inference and interpretable learning of the SDE path and dynamics, respectively, and does not move away from the good initial location for the model parameters  $\mathbf{C}$  and  $\mathbf{d}$ .

### 5.2. Van der Pol’s oscillator

Our next example examines a two-dimensional system where the dynamics contain a limit cycle around an unstable fixed point. The dynamics are given by

$$f_1(\mathbf{x}) = \rho\tau \left( x_1 - \frac{1}{3}x_1^3 - x_2 \right), \quad f_2(\mathbf{x}) = \frac{\tau}{\rho}x_1 \quad (25)$$

with a time constant  $\tau$ . We generate data from (1) using these dynamics with  $\rho = 2, \tau = 15, N = 20$  output dimensions and Gaussian measurement noise with unknown variances 2.25 on 20 repeated trials. We use  $5 \times 5$  inducing points evenly spaced in  $(-2, 2)$ . The results are summarised in Figure 2, demonstrating that our description of the dynamics successfully captures the limit cycle of the generative dynamics.

### 5.3. Neural population dynamics

This example demonstrates our algorithm under multivariate point-process observations. We model the intensity functions of the  $n$ th output process as  $\eta_n(t) =$

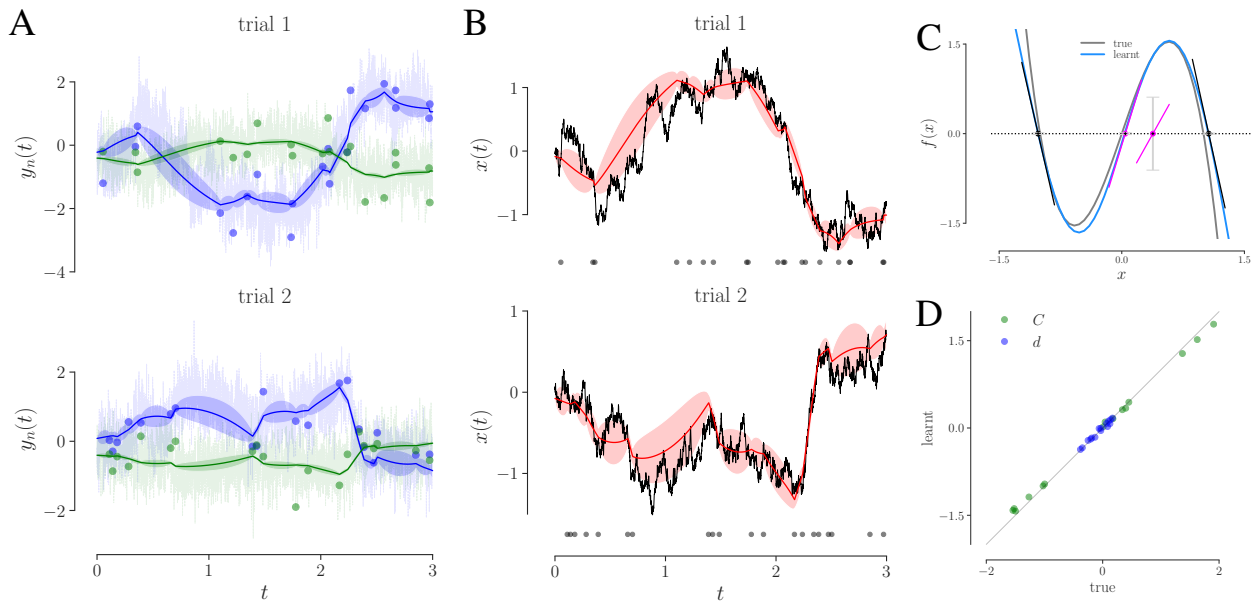


Figure 1. Double-well dynamics. A: Two example dimensions of the output process on two different trials. The dots represent the observed data-points of the noisy output processes plotted in faint lines. The solid blue/green traces are the inferred posterior means with  $\pm 1$  posterior standard deviation tubes around them. B: True and inferred latent SDE trajectory for the same example trials as in A. The red traces represent the posterior means with  $\pm 1$  posterior standard deviation tubes around them, black traces show the true latent SDE path. The black dots indicate the times when observations of  $\mathbf{y}$  were made. C: True and learnt dynamics together with the learnt fixed-point locations and tangent lines. Stable fixed points are shown in black, unstable ones in magenta. The uncertainty about the fixed point observation is illustrated using grey error bars representing  $\pm 1$  standard deviation. Only the additional fourth fixed point is associated with high uncertainty. D: True vs. learnt model parameters  $\mathbf{C}$  and  $\mathbf{d}$ .

$\exp(\sum_{k=1}^K C_{nk}x_k(t) + d_n)$ . Conditioned on the intensity function, the  $\phi(n)$  observed event-times  $\mathbf{t}^{(n)}$  are generated by a Poisson process with log-likelihood

$$\log p(\mathbf{t}^{(n)}|\eta_n) = -\int_{\mathcal{T}} \eta_n(t)dt + \sum_{i=1}^{\phi(n)} \log \eta_n(t_i^{(n)}) \quad (26)$$

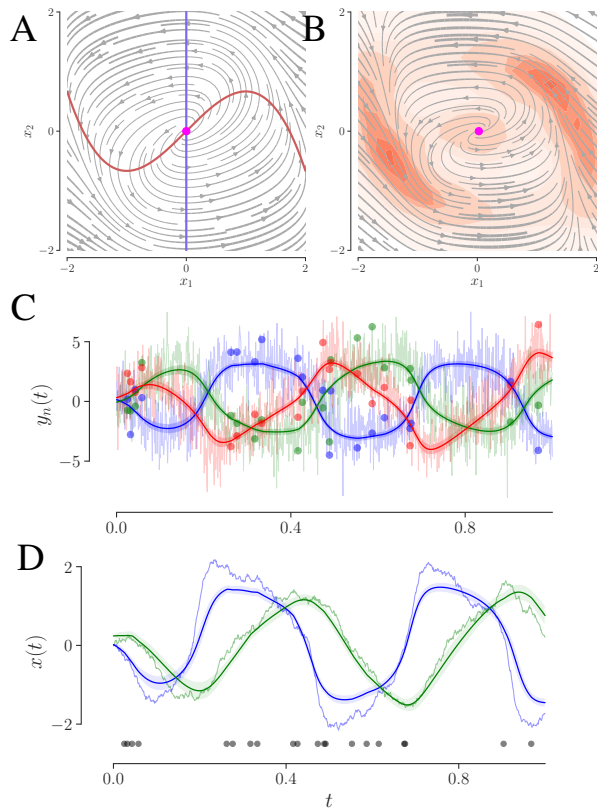
In contrast to the Gaussian observation case, the first term in the log-likelihood above is continuous in  $\eta_n(t)$  and the absence of events is also informative towards the underlying intensity of the process.

An interesting application for this setting lies in neural data analysis, where data may be available as a set of spike times of a population of simultaneously recorded neurons jointly embedded in a circuit involved in performing a computation. In fact, studying neural population activity as a dynamical system has gained increasing traction in the field of neuroscience in recent years (Macke et al., 2011; Shenoy et al., 2013; Pandarinath et al., 2018), and data analysis methods that can obtain such descriptions are thus of great interest.

We simulate a two-dimensional latent SDE using the dynamics  $f_k(\mathbf{x}) = -x_k + \sigma_k(w_{k1}x_1 - w_{k2}x_2 - z_k)$  for  $k = 1, 2$ , where  $\sigma_k(x) = (1 + \exp(-b_kx))^{-1}$ . Depending on the choice of parameters  $b_k$ ,  $w_{kj}$  and  $z_k$  the dynamical sys-

tem will exhibit different properties. We explore the two regimes where the system either has two stable and one unstable fixed points (Figure 3C left) or exhibits a single stable spiral (Figure 3C right).

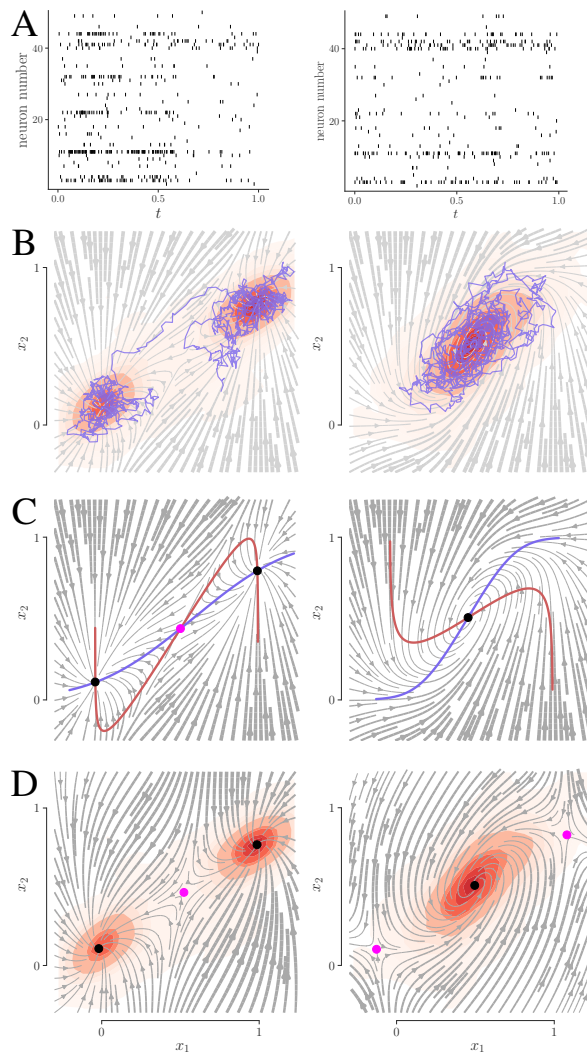
We simulate data from 50 neurons on 25 trials for each of the two parameter regimes for  $b_k$ ,  $w_{kj}$  and  $z_k$ . Figure 3A shows example neural spike trains under the two regimes. Figure 3B illustrates sample paths through the latent space under the different dynamical regimes, together with the density of latent locations visited across all trials. In both settings, we initialise our algorithm with three fixed points and inducing points placed on an evenly spaced  $4 \times 4$  grid in  $(-0.25, 1.25)$ , and hold the parameters relating to the output mapping constant. Figure 3D shows the estimated flow fields in both settings, together with the location of the fixed points and their stability as indicated by the eigenvalues of the Jacobian matrices. In both settings, our method successfully recovers the main qualitative distinguishing features of the dynamics. In the regime where the dynamics are conditioned on three fixed points but the generative system only contains one, the two additional fixed points will either be associated with higher uncertainty or move to regions where no or little data was observed, as indicated by the superimposed density plots.



**Figure 2.** Van der Pol’s oscillator. A: Streamline plot of the true dynamics together with nullclines and the unstable fixed point. B: Density plot of the locations visited by the latents across all trials used for learning in red, and streamline plot of the learnt dynamics with the location of the learnt fixed point. The eigenvalues of the learnt Jacobian matrix indicate that the fixed point is unstable. C: Three example dimensions of the output process. The dots represent the observed data-points of the noisy output process. The solid traces show the the posterior means with  $\pm 1$  standard deviation tubes around them. D: The true latent SDE path together with the posterior mean  $\pm 1$  posterior standard deviation of each latent dimension. Black dots represent the locations where the 20 measurements of the output process were made.

#### 5.4. Multistable chemical reaction dynamics

This example is based on the dynamical system in [Ganapathisubramanian \(1991\)](#), which describes nonlinear dynamics of two species of iodine in the iodate-AS(III) system under imperfect mixing by coupled first-order ODEs. We use these ODEs to describe  $f$  and generate data according to (1) with high-dimensional Gaussian observations representing spectroscopic measurements, which can approximately be described as a linear mapping from concentrations based on the  $I^-$  and  $IO_3^-$  absorption spectra provided in [Kireev & Shnyrev \(2015\)](#). We simulate data on 20 trials with different initial conditions, collecting 50 unevenly spaced samples from 13 spectroscopy measurements on each trial. Figure 4



**Figure 3.** Neural population dynamics. Left: simulations with parameter settings  $b_1 = 1.9, b_2 = 0.5, z_1 = 3, z_2 = 3.9, w_{11} = 10, w_{12} = 5, w_{21} = 9, w_{22} = 3$ . Right: simulations with parameter settings  $b_1 = 0.4, b_2 = 0.6, z_1 = 1.7, z_2 = 7, w_{11} = 20, w_{12} = 16, w_{21} = 21, w_{22} = 6$ . A: Raster plot of the observed spike times for a population of 50 neurons for an example trial. B: Example paths through the two-dimensional latent space on the same trial as A, together with a density plot of latent locations visited across all trials that were used for learning the dynamics, shown in red. C: Streamline plots of the true dynamics together with their fixed points and nullclines for each latent dimension. Stable fixed points are black, unstable ones are magenta. D: Same density plots as in B together with streamline plots of the learnt dynamics and learnt fixed points. The fixed point stability is shown as indicated by the eigenvalues of the learnt Jacobian matrices.

shows observed data and the latent SDE path on an example trial, as well as the true and estimated dynamical portraits.

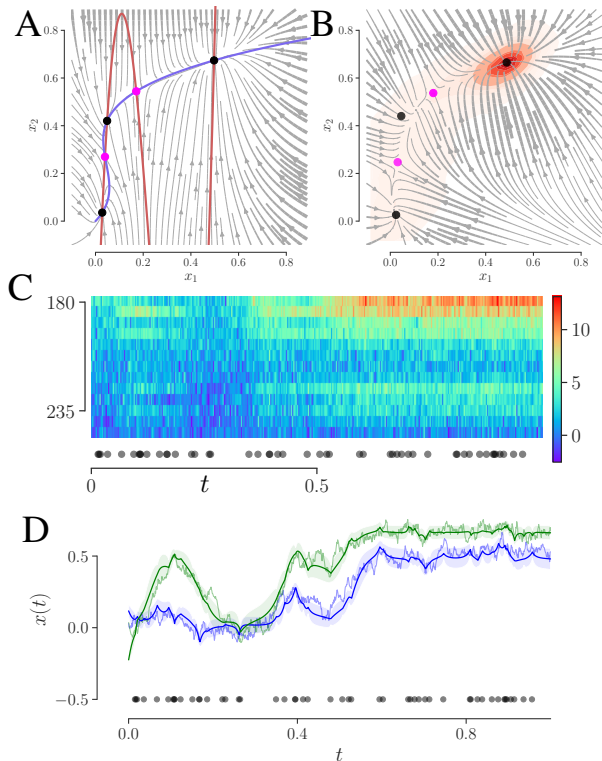


Figure 4. Multistable chemical reaction dynamics. A: Streamline plot of the concentration dynamics for two species of iodine, together with nullclines and fixed points. Stable fixed points are black, unstable ones are magenta. B: Learnt dynamics and fixed points with stability determined by eigenvalues of learnt Jacobians. The red contour plot illustrates the density of latent path locations across all trials used for training. C: Example spectroscopy measurements (output process) across light wavelengths (nm). D: Example true latent path together with the inferred posterior mean and  $\pm 1$  standard deviation tubes for each latent dimension on the same trial as C. The black dots indicate the time points at which measurements were taken.

## 6. Discussion

We have introduced a flexible and general variational Bayesian framework for the interpretable modelling of a continuous-time latent dynamical process from intermittent observations. Using a suitable GP prior, we integrate over a nonparametric description of the system dynamics, conditioned on its fixed points and associated local Jacobian matrices, thus both avoiding the need to assume a specific parametric dynamical form and directly obtaining a meaningful portrait of the dynamical structure. The approach applies to a variety of multivariate observation models, with many updates available in closed form.

The effectiveness of the approach is demonstrated using

data simulated from a number of realistic but known nonlinear dynamical systems describing physical, biological and chemical phenomena. In each case, it was possible to recover a meaningful description of fixed points and nearby dynamics even when data were sparse; and an inferred dynamical model that approximated the true systems well over large regions of the state space.

A similar prior over dynamics could be adopted within a discrete-time model such as the GP-SSM, albeit with a less natural interpretation of the local Jacobians. However, real-world systems evolve in continuous time, and in some contexts available observations do not arrive at discrete sample times. Retaining a continuous-time model means that the variational posterior over latents can be described by a system of coupled ODEs. While the solution of these may incur a discretisation error, this is a numerical issue related to the choice of ODE solver, rather than the assumption of a discretised model. Indeed, the ODE solution can exploit an adaptive step size in a way that would be impractical within a discrete-time model.

Our work also differs from other GP-based approaches to time series modelling, where each dimension of the process  $x_k(t)$  is modelled via an independent GP (Damianou et al., 2011; Duncker & Sahani, 2018). In this case, the prior on  $\mathbf{x}(t)$  evaluated at any finite set of points can be described by a multivariate Gaussian distribution, which greatly simplifies the inference. However, this cannot capture correlations across the dimensions of the latent process and thus comes at a loss of the descriptive power.

The variational inference approach for SDEs from Archambeau et al. (2007; 2008) relied on a Gaussian observation model and known dynamics (Archambeau et al., 2007), or a known parameterisation of the dynamics (Archambeau et al., 2008), both of which are restrictive. Here, we have extended the inference approach to handle a wider class of observation models, as well as a nonparametric GP description of the dynamics. Batz et al. (2018) also use a GP to model the drift function of an SDE. However, they consider the setting where dense or sparse observations of the SDE path are directly available, while we treat the entire SDE output as latent. Furthermore, the interpretable nonparametric representation of the SDE dynamics in terms of their fixed points and local Jacobian matrices is novel.

While we have demonstrated our algorithm in the setting of unevenly sampled multivariate Gaussian and multivariate point process observations, the inference approach extends readily to other stochastic processes typically considered challenging to model, such as marked point processes. We therefore expect this approach to have diverse applications, ranging from neuroscience to chemistry and finance.



## Acknowledgements

This work was funded by the Simons Foundation (SCGB 323228, 543039; MS) and the Gatsby Charitable Foundation.

## References

- Archambeau, C., Cornford, D., Opper, M., and Shawe-Taylor, J. Gaussian process approximations of stochastic differential equations. *Journal of machine learning research*, 1:1–16, 2007.
- Archambeau, C., Opper, M., Shen, Y., Cornford, D., and Shawe-taylor, J. S. Variational inference for diffusion processes. In *Advances in Neural Information Processing Systems*, pp. 17–24, 2008.
- Attias, H. A variational bayesian framework for graphical models. In *Advances in neural information processing systems*, pp. 209–215, 2000.
- Batz, P., Ruttor, A., and Opper, M. Approximate bayes learning of stochastic differential equations. *Physical Review E*, 98(2):022109, 2018.
- Bohner, G. and Sahani, M. Empirical fixed point bifurcation analysis. *arXiv preprint arXiv:1807.01486*, 2018.
- Cseke, B., Schnoerr, D., Opper, M., and Sanguinetti, G. Expectation propagation for continuous time stochastic processes. *Journal of Physics A: Mathematical and Theoretical*, 49(49):494002, 2016.
- Damianou, A., Titsias, M. K., and Lawrence, N. D. Variational gaussian process dynamical systems. In *Advances in Neural Information Processing Systems*, pp. 2510–2518, 2011.
- Deisenroth, M. and Mohamed, S. Expectation propagation in gaussian process dynamical systems. In *Advances in Neural Information Processing Systems*, pp. 2609–2617, 2012.
- Deisenroth, M. P., Huber, M. F., and Hanebeck, U. D. Analytic moment-based gaussian process filtering. In *Proceedings of the 26th annual international conference on machine learning*, pp. 225–232. ACM, 2009.
- Duncker, L. and Sahani, M. Temporal alignment and latent gaussian process factor inference in population spike trains. In Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., and Garnett, R. (eds.), *Advances in Neural Information Processing Systems 31*, pp. 10466–10476. Curran Associates, Inc., 2018.
- Eleftheriadis, S., Nicholson, T., Deisenroth, M., and Hensman, J. Identification of gaussian process state space models. In *Advances in Neural Information Processing Systems*, pp. 5309–5319, 2017.
- Frigola, R., Chen, Y., and Rasmussen, C. E. Variational gaussian process state-space models. In *Advances in Neural Information Processing Systems*, pp. 3680–3688, 2014.
- Ganapathisubramanian, N. Tristability in the iodate–as (iii) chemical system arising from a model of stirring and mixing effects. *The Journal of chemical physics*, 95(4):3005–3008, 1991.
- Golub, M. and Sussillo, D. Fixedpointfinder: A tensorflow toolbox for identifying and characterizing fixed points in recurrent neural networks. *Journal of Open Source Software*, 3(31):1003, 2018.
- Kireev, S. and Shnyrev, S. Study of molecular iodine, iodate ions, iodide ions, and triiodide ions solutions absorption in the uv and visible light spectral bands. *Laser Physics*, 25(7):075602, 2015.
- Macke, J. H., Buesing, L., Cunningham, J. P., Byron, M. Y., Shenoy, K. V., and Sahani, M. Empirical models of spiking in neural populations. In *Advances in neural information processing systems*, pp. 1350–1358, 2011.
- Pandarathna, C., O’Shea, D. J., Collins, J., Jozefowicz, R., Stavisky, S. D., Kao, J. C., Trautmann, E. M., Kaufman, M. T., Ryu, S. I., Hochberg, L. R., et al. Inferring single-trial neural population dynamics using sequential autoencoders. *Nature methods*, pp. 1, 2018.
- Quiñonero-Candela, J. and Rasmussen, C. E. A unifying view of sparse approximate gaussian process regression. *Journal of Machine Learning Research*, 6(Dec):1939–1959, 2005.
- Ramakrishnan, N., Ertin, E., and Moses, R. L. Assumed density filtering for learning gaussian process models. In *Statistical Signal Processing Workshop (SSP), 2011 IEEE*, pp. 257–260. IEEE, 2011.
- Shenoy, K. V., Sahani, M., and Churchland, M. M. Cortical control of arm movements: a dynamical systems perspective. *Annual review of neuroscience*, 36:337–359, 2013.
- Snelson, E. and Ghahramani, Z. Sparse gaussian processes using pseudo-inputs. In *Advances in neural information processing systems*, pp. 1257–1264, 2006.
- Sussillo, D. and Barak, O. Opening the black box: low-dimensional dynamics in high-dimensional recurrent neural networks. *Neural computation*, 25(3):626–649, 2013.

Titsias, M. Variational learning of inducing variables in sparse gaussian processes. In *Artificial Intelligence and Statistics*, pp. 567–574, 2009.

Turner, R., Deisenroth, M., and Rasmussen, C. State-space inference and learning with gaussian processes. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pp. 868–875, 2010.

Wang, J., Hertzmann, A., and Fleet, D. J. Gaussian process dynamical models. In *Advances in neural information processing systems*, pp. 1441–1448, 2006.

## A. Variational Lower Bound

We derive a variational lower bound to the marginal log-likelihood of our model using Jensen’s inequality

$$\begin{aligned} \log p(\mathbf{y}|\boldsymbol{\theta}) &= \log \int d\mathbf{x}d\mathbf{f}d\mathbf{u} p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\mathbf{f})p(\mathbf{f}|\mathbf{u},\boldsymbol{\theta})p(\mathbf{u}|\boldsymbol{\theta}) \\ &\geq \int d\mathbf{x}d\mathbf{f}d\mathbf{u} q(\mathbf{x},\mathbf{f},\mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\mathbf{f})p(\mathbf{f}|\mathbf{u},\boldsymbol{\theta})p(\mathbf{u}|\boldsymbol{\theta})}{q(\mathbf{x},\mathbf{f},\mathbf{u})} \\ &\stackrel{\text{def}}{=} \mathcal{F}^* \end{aligned}$$

where  $p(\mathbf{f}|\mathbf{u},\boldsymbol{\theta})p(\mathbf{u}|\boldsymbol{\theta}) = \prod_k p(f_k|\mathbf{u}_k,\boldsymbol{\theta})p(\mathbf{u}_k|\boldsymbol{\theta})$ . Choosing a factorised variational distribution of the form

$$q(\mathbf{x},\mathbf{f},\mathbf{u}) = q_x(\mathbf{x}) \prod_k p(f_k|\mathbf{u}_k,\boldsymbol{\theta})q_u(\mathbf{u}_k)$$

we can rewrite the bound as

$$\begin{aligned} \mathcal{F}^* &= \int d\mathbf{x}d\mathbf{f}d\mathbf{u} q(\mathbf{x},\mathbf{f},\mathbf{u}) \log \frac{p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\mathbf{f}) \prod_k p(\mathbf{u}_k|\boldsymbol{\theta})}{q_x(\mathbf{x}) \prod_k q_u(\mathbf{u}_k)} \\ &= \langle \log p(\mathbf{y}|\mathbf{x}) \rangle_{q_x} - \langle \text{KL}[q_x(\mathbf{x})\|p(\mathbf{x}|\mathbf{f})] \rangle_{q_f} \\ &\quad - \sum_k \text{KL}[q_u(\mathbf{u}_k)\|p(\mathbf{u}_k)] \end{aligned}$$

where

$$q_f(\mathbf{f}) = \prod_k \int d\mathbf{u}_k p(f_k|\mathbf{u}_k,\boldsymbol{\theta})q_u(\mathbf{u}_k)$$

and  $q_x(\mathbf{x})$  is described by (3) and (4). We can derive the Kullback-Leibler divergence between the distributions over SDE paths  $q_x(\mathbf{x})$  and  $p(\mathbf{x}|\mathbf{f})$  by discretising time in steps of  $\Delta t$ . The discretised paths have Markovian structure with

$$\begin{aligned} p(\mathbf{x}_{t+1}|\mathbf{x}_t,\mathbf{f}) &= \mathcal{N}(\mathbf{x}_{t+1}|\mathbf{x}_t + \mathbf{f}(\mathbf{x}_t)\Delta t, \boldsymbol{\Sigma}\Delta t) \\ q_x(\mathbf{x}_{t+1}|\mathbf{x}_t) &= \mathcal{N}(\mathbf{x}_{t+1}|\mathbf{x}_t + \mathbf{f}_q(\mathbf{x}_t)\Delta t, \boldsymbol{\Sigma}\Delta t) \end{aligned}$$

We can hence write

$$\begin{aligned} \text{KL}[q_x(\mathbf{x})\|p(\mathbf{x})] &= \sum_{t=1}^{T-1} \int d\mathbf{x}_t q(\mathbf{x}_t) \int d\mathbf{x}_{t+1} q(\mathbf{x}_{t+1}|\mathbf{x}_t) \log \frac{q(\mathbf{x}_{t+1}|\mathbf{x}_t)}{p(\mathbf{x}_{t+1}|\mathbf{x}_t)} \\ &= \frac{1}{2} \sum_{t=1}^{T-1} \Delta t \langle (\mathbf{f} - \mathbf{f}_q)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{f} - \mathbf{f}_q) \rangle_{q_x} \end{aligned}$$

Taking the limit as  $\Delta t \rightarrow 0$ , we obtain

$$\text{KL}[q_x(\mathbf{x})\|p(\mathbf{x})] = \frac{1}{2} \int_{\mathcal{T}} dt \langle (\mathbf{f} - \mathbf{f}_q)^\top \boldsymbol{\Sigma}^{-1} (\mathbf{f} - \mathbf{f}_q) \rangle_{q_x}$$

## B. Inference Details

### B.1. Lagrangian

The full Lagrangian, after applying integration by parts to the constraints in (8), has the form

$$\mathcal{L} = \mathcal{F}^* - \mathcal{C}_1 - \mathcal{C}_2$$

$$\begin{aligned} \mathcal{C}_1 &= \int_{\mathcal{T}} dt \left( \text{Tr} \left[ \Psi (\mathbf{A} \mathbf{S}_x + \mathbf{S}_x \mathbf{A}^\top - I) - \frac{d\Psi}{dt} \mathbf{S}_x \right] \right) \\ &\quad + \text{Tr} [\Psi(T) \mathbf{S}_x(T)] - \text{Tr} [\Psi(0) \mathbf{S}_x(0)] \\ \mathcal{C}_2 &= \int_{\mathcal{T}} dt \left( \lambda^\top (\mathbf{A} \mathbf{m}_x - \mathbf{b}) - \frac{d\lambda}{dt} \mathbf{m}_x \right) \\ &\quad + \lambda(T)^\top \mathbf{m}_x(T) - \lambda(0)^\top \mathbf{m}_x(0) \end{aligned}$$

For the variational free energy term  $\mathcal{F}^*$ , we have from before

$$\mathcal{F} = \sum_i \langle \log p(\mathbf{y}_i | \mathbf{x}_i) \rangle_{q_x} - \text{KL}[q_x(\mathbf{x}) \| p(\mathbf{x})]$$

and

$$\mathcal{F}^* = \langle \mathcal{F} \rangle_{q_f} - \sum_{k=1}^K \text{KL}[q_u(\mathbf{u}_k) \| p(\mathbf{u}_k | \boldsymbol{\theta})]$$

The Kullback-Leibler divergences can be evaluated as

$$\begin{aligned} \text{KL}[q(\mathbf{u}_k) \| p(\mathbf{u}_k | \boldsymbol{\theta})] &= \frac{1}{2} \left( \text{Tr} [\boldsymbol{\Omega}_u^{-1} \mathbf{S}_u^k] - M + \log \frac{|\boldsymbol{\Omega}_u|}{|\mathbf{S}_u^k|} \right) \\ &\quad + (\boldsymbol{\mu}_u^k - \mathbf{m}_u^k)^\top \boldsymbol{\Omega}_u^{-1} (\boldsymbol{\mu}_u^k - \mathbf{m}_u^k) \end{aligned}$$

with

$$\begin{aligned} \boldsymbol{\Omega}_u &= \mathbf{K}_{zz} - \tilde{\mathbf{K}}_{zs} \tilde{\mathbf{K}}_{ss}^{-1} \tilde{\mathbf{K}}_{sz} \\ \boldsymbol{\mu}_u^k &= \tilde{\mathbf{K}}_{zs} \tilde{\mathbf{K}}_{ss}^{-1} \mathbf{v}_k^\theta \end{aligned}$$

and

$$\langle \text{KL}[q_x(\mathbf{x}) \| p(\mathbf{x})] \rangle_{q_f} = \frac{1}{2} \int_0^T dt \langle (\mathbf{f} - \mathbf{f}_q)^\top (\mathbf{f} - \mathbf{f}_q) \rangle_{q_x q_f}$$

For later convenience, we denote this term as

$$\langle \text{KL}[q_x(\mathbf{x}) \| p(\mathbf{x})] \rangle_{q_f} = \mathcal{E}(\mathbf{m}_x, \mathbf{S}_x)$$

Using the identity

$$\langle \langle \mathbf{f} \rangle_{q_f} (\mathbf{x} - \mathbf{m}_x)^\top \rangle_{q_x} = \left\langle \frac{\partial \langle \mathbf{f} \rangle_{q_f}}{\partial \mathbf{x}} \right\rangle_{q_x} \mathbf{S}_x$$

the integrand can be evaluated as

$$\begin{aligned} &\langle (\mathbf{f} - \mathbf{f}_q)^\top (\mathbf{f} - \mathbf{f}_q) \rangle_{q_x q_f} \\ &= \langle \mathbf{f}^\top \mathbf{f} \rangle_{q_x q_f} + 2 \text{Tr} \left[ \mathbf{A}^\top \left\langle \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right\rangle_{q_x q_f} \mathbf{S}(t) \right] \\ &\quad + \text{Tr} [\mathbf{A}^\top \mathbf{A} (\mathbf{S}_x + \mathbf{m}_x \mathbf{m}_x^\top)] + 2 \mathbf{m}_x^\top \mathbf{A}^\top \langle \mathbf{f} \rangle_{q_x q_f} \\ &\quad + \mathbf{b}^\top \mathbf{b} - 2 \mathbf{b}^\top \langle \mathbf{f} \rangle - 2 \mathbf{b}^\top \mathbf{A} \mathbf{m}_x \end{aligned}$$

For the expected log-likelihood terms, in general, there will be terms that are continuous in  $\mathbf{x}$ , and terms that depend only on evaluations of  $\mathbf{x}$  at specific locations  $t_i$ , which we will denote by  $\ell^{cont}$  and  $\ell^{jump}$ , respectively. We can write

$$\langle \log p(\mathbf{y} | \mathbf{x}) \rangle_{q_x} = \ell^{cont}(\mathbf{m}_x, \mathbf{S}_x) + \ell^{jump}(\mathbf{m}_x, \mathbf{S}_x)$$

Thus, the variational free energy can be expressed as

$$\begin{aligned} \mathcal{F}^* &= \ell^{cont}(\mathbf{m}_x, \mathbf{S}_x) + \ell^{jump}(\mathbf{m}_x, \mathbf{S}_x) - \mathcal{E}(\mathbf{m}_x, \mathbf{S}_x) \\ &\quad - \sum_{k=1}^K \text{KL}[q_u(\mathbf{u}_k) \| p(\mathbf{u}_k | \boldsymbol{\theta})] \end{aligned}$$

### B.1.1. EXAMPLE: GAUSSIAN LIKELIHOOD

In the case of a Gaussian likelihood, there is no continuous term in the likelihood:

$$\begin{aligned} \ell^{cont} &= 0 \\ \ell^{jump} &= \sum_i \int_{\mathcal{T}_0}^{\mathcal{T}_{end}} dt \delta(t - t_i) (\mathbf{m}_x(t)^\top \mathbf{C}^\top \Gamma^{-1} (\mathbf{y}_t - \mathbf{d}) \\ &\quad - \frac{1}{2} \text{Tr} \left[ \mathbf{C}^\top \Gamma^{-1} \mathbf{C} \sum_i (\mathbf{S}_x(t) + \mathbf{m}_x(t) \mathbf{m}_x(t)^\top) \right]) \end{aligned}$$

### B.1.2. EXAMPLE: MULTIVARIATE POISSON PROCESS LIKELIHOOD

In the case of a multivariate Poisson Process, with  $g(\cdot) = \exp(\cdot)$  and observed event times  $t_1^{(n)}, \dots, t_{\phi(n)}^{(n)}$  for the  $n$ th output dimension:

$$\begin{aligned} \ell^{cont} &= - \sum_n \int_{\mathcal{T}_0}^{\mathcal{T}_{end}} \exp \left( \mathbf{c}_n^\top \mathbf{m}_x + \frac{1}{2} \mathbf{c}_n^\top \mathbf{S}_x \mathbf{c}_n \right) dt \\ \ell^{jump} &= \sum_{n=1}^N \sum_{i=1}^{\phi(n)} \int_{\mathcal{T}_0}^{\mathcal{T}_{end}} (\mathbf{c}_n^\top \mathbf{m}_x(t) + d_n) \delta(t - t_i^{(n)}) dt \end{aligned}$$

## B.2. Symmetric variations in $\mathbf{S}_x$

To arrive at the fixed point equations given in the main paper, we need to take variational derivatives of the Lagrangian with respect to  $\mathbf{m}_x$  and  $\mathbf{S}_x$ . In contrast to [Archambeau et al. \(2007\)](#), we take the symmetric variations in  $\mathbf{S}_x$  into account. Also note that the Lagrange multiplier  $\Psi$  is symmetric. We can write

$$\frac{\partial \mathcal{C}_1}{\partial \mathbf{S}_x} = \left( \Psi \mathbf{A} + \mathbf{A}^\top \Psi - \frac{d\Psi}{dt} \right) \odot \tilde{\mathbb{P}}$$

where  $\odot$  denotes the elementwise Hadamard product and  $\tilde{\mathbb{P}}_{ij} = 2$  for  $i \neq j$  and 1 otherwise. Differentiating the entire Lagrangian with respect to the symmetric matrix  $\mathbf{S}_x$  and setting to zero we get

$$0 = \frac{\partial \mathcal{F}^*}{\partial \mathbf{S}_x} \odot \mathbb{P} - \Psi \mathbf{A} - \mathbf{A}^\top \Psi + \frac{d\Psi}{dt}$$

matching the equation given in the main text with  $\mathbb{P}_{ij} = \frac{1}{2}$  if  $i \neq j$  and 1 otherwise. Note that the derivatives of the free energy with respect to  $\mathbf{S}_x$  will also need to take into account the symmetry of the covariance matrix. The derivations for (20)-(22) follow those of [Archambeau et al. \(2007\)](#).

### B.3. Expected values of dynamics

The inference algorithm requires evaluating several expectations with respect to  $q_x$  and  $q_f$ . Let  $\mathbf{U} = [\mathbf{u}_1 \ \dots \ \mathbf{u}_K]$  and  $\langle \mathbf{U} \rangle_{q_u} = \mathbf{M}_u$ , such that we can define  $(M+L+LK) \times K$  matrices stacking all inducing variables, zero function values, and Jacobians as

$$\mathbf{U}_{u,f_s,J} = \begin{bmatrix} \mathbf{U}_u \\ \mathbf{0} \\ \mathbf{J}_s^{(1)} \\ \vdots \\ \mathbf{J}_s^{(L)} \end{bmatrix}, \quad \langle \mathbf{U}_{u,f_s,J} \rangle_{q_u} = \mathbf{M}_{u,f_s,J} = \begin{bmatrix} \mathbf{M}_u \\ \mathbf{0} \\ \mathbf{J}_s^{(1)} \\ \vdots \\ \mathbf{J}_s^{(L)} \end{bmatrix}$$

The required expectations can then be evaluated as

$$\langle \mathbf{f}(\mathbf{x}) \rangle_{q_x q_f}^\top = \langle \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \mathbf{M}_{u,f_s,J}$$

$$\left\langle \frac{\partial \mathbf{f}(\mathbf{x})}{\partial \mathbf{x}} \right\rangle_{q_x q_f}^\top = \langle \nabla_x \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \mathbf{M}_{u,f_s,J}$$

$$\begin{aligned} \langle \mathbf{f}(\mathbf{x})^\top \mathbf{f}(\mathbf{x}) \rangle_{q_x q_f} &= \sum_k \langle f_k^2(\mathbf{x}) \rangle_{q_x q_f} = \kappa(\mathbf{x}, \mathbf{x}') \\ &+ \text{Tr} \left[ \left( \langle \mathbf{U}_{u,f_s,J} \mathbf{U}_{u,f_s,J}^\top \rangle_{q_u} - \mathbf{K}_{zz}^\theta \right) \langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \right] \end{aligned}$$

The above expressions still involve computing expectations of covariance functions and their derivatives, which can be computed analytically for choices such as the exponentiated quadratic covariance function.

### B.4. Inference algorithm

The full inference algorithm involves solving a set of ODEs forward and backward in time, which we do using the forward Euler method. We provide the full approach in Algorithm 1, where the subscript  $r$  denotes the evaluation of the functions at the  $r$ th point of the time grid between  $\mathcal{T}_0$  and  $\mathcal{T}_{end}$  taking steps of size  $\Delta t$ . Note that the derivatives of the terms in  $\ell^{jump}$  will need to be discretized appropriately as well. Using the same time-grid as was used for solving the ODEs, the delta-functions will contribute a factor of  $\frac{1}{\Delta t}$ , such that the  $\Delta t$  terms cancel in the update written in Algorithm 1.

## C. Learning Details

### C.1. Conditioned Sparse Gaussian Process dynamics

The only term in the variational free energy that depends on the parameters in  $\mathbf{f}$  are the KL-divergence between the continuous-time processes and the KL-divergence relating to the inducing points for  $\mathbf{f}$ .

#### C.1.1. INDUCING POINT COVARIANCES

Collecting the terms that contain  $\mathbf{S}_u^k$  we have

$$\frac{\partial}{\partial \mathbf{S}_u^k} \text{KL}[q(\mathbf{u}_k) \| p(\mathbf{u}_k | \boldsymbol{\theta})] = \frac{1}{2} \boldsymbol{\Omega}_u^{-1} - \frac{1}{2} \mathbf{S}_u^k^{-1}$$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{S}_u^k} &= \frac{1}{2} \int_{\mathcal{T}} dt \frac{\partial}{\partial \mathbf{S}_u^k} \text{Tr} \left[ \begin{bmatrix} \mathbf{S}_u^k & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \right] \\ &= \frac{1}{2} \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{:M,:M} \end{aligned}$$

where the last line selects the first  $M \times M$  block from  $\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}$ . We hence obtain the closed form update

$$\mathbf{S}_u^k = \left( \boldsymbol{\Omega}_u^{-1} + \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{:M,:M} \right)^{-1}$$

#### C.1.2. INDUCING POINTS AND JACOBIANS

To find the update efficiently, let  $\mathbf{J}_k = [\mathbf{J}_{k,:}^{(1)}, \dots, \mathbf{J}_{k,:}^{(L)}]^\top$  so that we can write

$$\boldsymbol{\mu}_u^k = \tilde{\mathbf{K}}_{zs} \tilde{\mathbf{K}}_{ss}^{-1} \mathbf{v}_k^\theta = \tilde{\mathbf{K}}_{zs} \tilde{\mathbf{K}}_{ss}^{-1} \begin{bmatrix} \mathbf{0} \\ \mathbf{J}_k \end{bmatrix} = \mathbf{G} \mathbf{J}_k$$

We can rewrite the quadratic terms in the Kullback-Leibler divergences of the inducing points as

$$\begin{aligned} &\sum_k (\boldsymbol{\mu}_u^k - \mathbf{m}_u^k)^\top \boldsymbol{\Omega}_u^{-1} (\boldsymbol{\mu}_u^k - \mathbf{m}_u^k) \\ &= \sum_k \begin{bmatrix} \mathbf{m}_u^k \\ \mathbf{J}_k \end{bmatrix}^\top \begin{bmatrix} \boldsymbol{\Omega}_u^{-1} & -\boldsymbol{\Omega}_u^{-1} \mathbf{G} \\ -\mathbf{G}^\top \boldsymbol{\Omega}_u^{-1} & \mathbf{G}^\top \boldsymbol{\Omega}_u^{-1} \mathbf{G} \end{bmatrix} \begin{bmatrix} \mathbf{m}_u^k \\ \mathbf{J}_k \end{bmatrix} \\ &= \text{Tr} \left[ \mathbf{M}_{u,J}^\top \tilde{\boldsymbol{\Omega}} \mathbf{M}_{u,J} \right] \end{aligned}$$

with  $\mathbf{M}_{u,J} = \begin{bmatrix} \mathbf{m}_u^1 & \dots & \mathbf{m}_u^K \\ \mathbf{J}_1 & \dots & \mathbf{J}_K \end{bmatrix}$  and derivative

$$\frac{\partial}{\partial \mathbf{M}_{u,J}} \sum_k \text{KL}[q(\mathbf{u}_k) \| p(\mathbf{u}_k | \boldsymbol{\theta})] = \tilde{\boldsymbol{\Omega}} \mathbf{M}_{u,J}$$

$$\begin{aligned} \frac{\partial \mathcal{E}}{\partial \mathbf{M}_{u,J}} &= \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{[i,i]} \mathbf{M}_{u,J} \\ &+ \int_{\mathcal{T}} dt \left[ \langle \nabla_x \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \right]_{[:,i]}^\top \mathbf{S}_x \mathbf{A}^\top \\ &- \int_{\mathcal{T}} dt \left[ \langle \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \right]_{[:,i]}^\top (-\mathbf{A} \mathbf{m}_x + \mathbf{b})^\top \end{aligned}$$

Putting all terms together, we obtain the update

$$\mathbf{M}_{u,J} = \mathbf{B}_1^{-1} (\mathbf{B}_2 - \mathbf{B}_3)$$



**Algorithm 1** Inference algorithm

---

**Input:** data  $\{y_i, t_i\}_{i=1}^T$ ,  $\mathbf{m}_{x,0}$ ,  $\mathbf{S}_{x,0}$ ,  $q_f(\mathbf{f})$ ,  $\Delta t$ ,  $\mathcal{T}_0$ ,  $\mathcal{T}_{end}$   
 Initialize  $\mathbf{A}(t)$ ,  $\mathbf{b}(t)$   
 $R = \frac{\mathcal{T}_0 - \mathcal{T}_{end}}{\Delta t}$   
**repeat**  
   **for**  $r = 0$  **to**  $R - 1$  **do**  
      $\mathbf{m}_{x,r+1} \leftarrow \mathbf{m}_{x,r} - \Delta t (\mathbf{A}_r \mathbf{m}_{x,r} - \mathbf{b}_r)$   
      $\mathbf{S}_{x,r+1} \leftarrow \mathbf{S}_{x,r} - \Delta t (\mathbf{A}_r \mathbf{S}_{x,r} + \mathbf{S}_{x,r} \mathbf{A}_r^\top - I)$   
   **end for**  
   **for**  $r = R$  **to**  $1$  **do**  
      $\lambda_{r-1} \leftarrow \lambda_r - \Delta t \left( \mathbf{A}_r^\top \lambda_r + \left( \frac{\partial \ell^{cont}}{\partial \mathbf{m}_x} - \frac{\partial \mathcal{E}}{\partial \mathbf{m}_x} \right) \Big|_{t=r\Delta t} \right) - \Delta t \frac{\partial \ell^{jump}}{\partial \mathbf{m}_x} \Big|_{t=(r-1)\Delta t}$   
      $\Psi_{r-1} \leftarrow \Psi_r - \Delta t \left( \mathbf{A}_r^\top \Psi_r + \Psi_r \mathbf{A}_r + \mathbb{P} \odot \left( \frac{\partial \ell^{cont}}{\partial \mathbf{S}_x} - \frac{\partial \mathcal{E}}{\partial \mathbf{S}_x} \right) \Big|_{t=r\Delta t} \right) - \Delta t \mathbb{P} \odot \frac{\partial \ell^{jump}}{\partial \mathbf{S}_x} \Big|_{t=(r-1)\Delta t}$   
   **end for**  
    $\mathbf{A} = \left\langle \frac{\partial \mathbf{f}}{\partial \mathbf{x}} \right\rangle_{q_x q_f} + 2\Psi$   
    $\mathbf{b} = \langle \mathbf{f}(\mathbf{x}) \rangle_{q_x q_f} + \mathbf{A} \mathbf{m}_x - \lambda$   
**until** convergence in  $\mathcal{F}^*$   
**return:**  $\{\mathbf{A}_r, \mathbf{b}_r, \lambda_r, \Psi_r, \mathbf{m}_{x,r}, \mathbf{S}_{x,r}\}_{r=1}^R$

---

with

$$\begin{aligned}
 \mathbf{B}_1 &= \left( \tilde{\Omega} + \int_{\mathcal{T}} dt [\langle \mathbf{a}_z^\theta(\mathbf{x})^\top \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x}]_{[u_j, u_j]} \right) \\
 \mathbf{B}_2 &= \int_{\mathcal{T}} dt \left[ \langle \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \right]_{[:, u_j]}^\top \langle \mathbf{f}_q \rangle_{q_x}^\top \\
 \mathbf{B}_3 &= \int_{\mathcal{T}} dt \left[ \langle \nabla_x \mathbf{a}_z^\theta(\mathbf{x}) \rangle_{q_x} \right]_{[:, u_j]}^\top \mathbf{S}_x \mathbf{A}^\top
 \end{aligned}$$

and we have defined an indexing operation where  $[X]_{[u_j, u_j]}$  selects the first  $M \times M$  and last  $LK \times LK$  block of  $X$  and  $[X]_{[:, u_j]}$  selects the first  $M$  and last  $LK$  columns of  $X$ . Hence, this selects the appropriate block matrices for the updates. The one-dimensional integrals can be computed efficiently using Gauss-Legendre quadrature.

### C.2. Sparse Gaussian Process dynamics

Similarly, closed form updates are available in the simpler case, when  $\mathbf{f}$  is modelled by a classic sparse Gaussian Process, i.e. using inducing points without the additional conditioning on fixed points and Jacobians.

$$\begin{aligned}
 \mathbf{S}_u^k &= \mathbf{K}_{zz} \left( \mathbf{K}_{zz} + \int_{\mathcal{T}} dt \langle \kappa(\mathbf{Z}, \mathbf{x}) \kappa(\mathbf{x}, \mathbf{Z}) \rangle_{q_x} \right)^{-1} \mathbf{K}_{zz} \\
 \mathbf{M}_u &= \mathbf{S}_u^k \mathbf{K}_{zz}^{-1} \left( \int_{\mathcal{T}} dt \Phi_{1q} \mathbf{f}_q^\top - \int_{\mathcal{T}} dt \Phi_{d1} \mathbf{S}_x \mathbf{A}^\top \right)
 \end{aligned}$$

Where  $\Phi_{1q} = \langle k(\mathbf{x}, \mathbf{Z}) \rangle_{q_x}$  and  $\Phi_{d1} = \langle \frac{\partial}{\partial \mathbf{x}} k(\mathbf{x}, \mathbf{Z}) \rangle_{q_x}$ .

### C.3. Linear dynamics

Our modelling framework also easily extends to other parameterisation of  $\mathbf{f}$ . For example, in a continuous-time lin-

ear dynamical system with  $\mathbf{f}(\mathbf{x}) = -\tilde{\mathbf{A}}\mathbf{x} + \tilde{\mathbf{b}}$  direct minimisation of the KL-divergence between the continuous-time processes leads to the closed form updates

$$\begin{aligned}
 \tilde{\mathbf{A}} &= \left( \int_{\mathcal{T}} dt (\mathbf{b}(\mathbf{x})^\top - \langle \mathbf{f}_q(\mathbf{x}) \mathbf{x}^\top \rangle) \right) \left( \int_{\mathcal{T}} dt \langle \mathbf{x} \mathbf{x}^\top \rangle \right)^{-1} \\
 \tilde{\mathbf{b}} &= \frac{1}{T} \int_{\mathcal{T}} dt (\langle \mathbf{f}_q(\mathbf{x}) \rangle + \mathbf{A} \langle \mathbf{x} \rangle)
 \end{aligned}$$

reminiscent of the update equations for the generative parameters of a discrete-time Linear Dynamical System.

### C.4. Output mapping

We consider an observation model of the form

$$\mathbf{y}(t_i) = \mathbf{C}\mathbf{x}(t_i) + \mathbf{d} + \boldsymbol{\epsilon}_i$$

where  $\boldsymbol{\epsilon}_i \sim \mathcal{N}(\boldsymbol{\epsilon}|0, \Gamma)$ . Dropping all terms that are constant in  $\mathbf{C}$ ,  $\mathbf{d}$  from the expression for the variational free energy, we have

$$\mathcal{F}^* = -\frac{1}{2} \sum_t \left\langle (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{d})^\top \Gamma^{-1} (\mathbf{y}_t - \mathbf{C}\mathbf{x}_t - \mathbf{d}) \right\rangle_{q_x}$$

Differentiating and setting to zero gives

$$\begin{aligned}
 \mathbf{C}^{new} &= \left( \sum_t (\mathbf{y}_t - \mathbf{d}) \mathbf{m}_t^\top \right) \left( \sum_t (\mathbf{S}_{x,t} + \mathbf{m}_{x,t} \mathbf{m}_{x,t}^\top) \right)^{-1} \\
 \mathbf{d}^{new} &= \frac{1}{T} \sum_t (\mathbf{y}_t - \mathbf{C}^{new} \mathbf{m}_{x,t})
 \end{aligned}$$

## D. Chemical reaction dynamics

The dynamical system used to generate the data in section 5.4 is of the form

$$\begin{aligned} \frac{d[I^-]_A}{dt} &= (k_a[I^-]_A + k_b[I^-]_A^2) (S_0 - [I^-]_A) \\ &\quad + \frac{F_1[I^-]_0}{V_A} - \frac{(F_3 + F_4)[I^-]_A}{V_A} + \frac{F_4[I^-]_D}{V_A} \\ \frac{d[I^-]_D}{dt} &= (k_a[I^-]_D + k_b[I^-]_D^2) (S_0 - [I^-]_D) \\ &\quad + \frac{F_4[I^-]_A}{V_D} - \frac{F_4[I^-]_D}{V_D} \end{aligned}$$

To generate the simulations, we use the parameter settings

$$\begin{aligned} [I^-]_0 &= 4.4 \times 10^{-5} & k_0 &= 2.7 \times 10^{-3} \\ V_A &= 4 \times 10^1 & F_4 &= 3.25 \times 10^{-3} \\ V_D &= 1 & F_3 &= k_0 V_a \\ k_a &= 2.1425 \times 10^{-1} & F_1 &= \frac{1}{2} F_3 \\ k_b &= 2.1425 \times 10^4 & F_2 &= \frac{1}{2} F_3 \\ S_0 &= \frac{1}{2} ([I^-]_0 + 1.42 \times 10^{-3}) \end{aligned}$$