

# TOWARDS EFFICIENT MODELING AND INFERENCE IN MULTI-DIMENSIONAL GAUSSIAN PROCESS STATE-SPACE MODELS

Zhidi Lin\* Juan Maroñas† Ying Li‡ Feng Yin\*(✉) Sergios Theodoridis‡

\* The Chinese University of Hong Kong, Shenzhen, China † The University of Hong Kong, HKSAR  
‡ Machine Learning Group, Autonomous University of Madrid, and Cognizant, Madrid, Spain  
‡ National and Kapodistrian University of Athens, Greece, and Aalborg University, Denmark

## ABSTRACT

The Gaussian process state-space model (GPSSM) has attracted extensive attention for modeling complex nonlinear dynamical systems. However, the existing GPSSM employs separate Gaussian processes (GPs) for each latent state dimension, leading to escalating computational complexity and parameter proliferation, thus posing challenges for modeling dynamical systems with high-dimensional latent states. To surmount this obstacle, we propose to integrate the efficient transformed Gaussian process (ETGP) into the GPSSM, which involves pushing a shared GP through multiple normalizing flows to efficiently model the transition function in high-dimensional latent state space. Additionally, we develop a corresponding variational inference algorithm that surpasses existing methods in terms of parameter count and computational complexity. Experimental results on diverse synthetic and real-world datasets corroborate the efficiency of the proposed method, while also demonstrating its ability to achieve similar inference performance compared to existing methods. Code is available at <https://github.com/zhidilin/gpssmProj>.

**Index Terms**— Gaussian process state-space model, efficiency and scalability, multi-dimensional state, normalizing flow, variational approximations.

## 1. INTRODUCTION

Gaussian process state-space models (GPSSMs) [1] have gained significant popularity among data-driven state-space models (SSMs) [1–7], due to their capability of integrating non-parametric Bayesian Gaussian processes (GPs) [8–10] as function priors within the classical SSM [11]. This integration empowers the model to effectively learn the system dynamics from noisy measurements with explicit uncertainty calibrations [12]. Additionally, GP is able to automatically scale model complexity based on data volume [13]. Consequently, GPSSMs and their variants have demonstrated successful applications in diverse domains, including human motion capture, pedestrian tracking, and navigation [14–17]. Research efforts have also focused on advancing the simultaneous learning and inference capabilities of GPSSMs [1, 4–7, 18–23].

However, in the context of high-dimensional latent state spaces, all the existing GPSSMs face two primary challenges. First, current GPSSM methods often resort to independent GPs for modeling multiple outputs of transition functions, aiming for simplicity but overlooking their dependencies. This disregard for dependencies can lead to a model mismatch and the loss of inductive bias among the outputs [24]. This can ultimately impede the model’s generalization capacity and cause a decline in inference performance, particularly when latent states are only partially observed [20]. Second, employing separate GPs to model the state transition for each latent

state dimension leads to a quadratic expansion in the number of parameters, coupled with a linear rise in  $\mathcal{O}(n^3)$  computational complexity, as dimensionality increases, see Fig. 1. Here,  $n$  represents the sample count used for computing the GP kernel matrix [8]. This escalating computational burden and parameters proliferation can become prohibitively cumbersome, especially when dealing with high-dimensional latent spaces. Consequently, addressing these two challenges becomes imperative to enhance the applicability and scalability of GPSSM in practical applications. For the first challenge, existing approaches have explored potential solutions, with some using a linear model of coregionalization (LMC)-based multi-output GP to model this correlation [20]. The transformed GP (TGP) [25] framework has also been introduced, wherein multiple independent GPs are transformed by a normalizing flow [26] to somewhat obtain correlated outputs [21]. Despite these efforts, the persistent challenge of escalating complexity remains an obstacle across existing works, necessitating research to further enhance the applicability of GPSSMs in high-dimensional latent state spaces.

This paper aims to address the escalating computational complexity and parameter proliferation in the GPSSM while introducing a novel form of output dependence. The main contributions are summarized as follows. First, we present an innovative efficient GPSSM paradigm that deviates from the standard approach of employing separate GPs for each one of the latent dimensions. Instead, we adopt the efficient transformed GP (ETGP) [27], capitalizing on multiple normalizing flows [26] to enact transformations on a shared GP across each dimension of the latent state space. This strategic shift allows us to attain streamlined modeling while effectively establishing output dependencies. Second, for joint learning and inference in the proposed efficient GPSSM, we propose a proficient sparse GP [28]-based variational algorithm that enhances computational efficacy and streamlines parameter scale. Third, experimental results, obtained using real and synthetic datasets, corroborate comparable performance of the proposed efficient GPSSM to existing GPSSMs, albeit at substantial reductions in both computational complexity and parameter count.

The remainder of this paper is organized as follows. Some preliminaries related to GPSSM are provided in Section 2. Section 3 introduces our proposed efficient output-dependent GPSSM and the associated learning and inference algorithm. Numerical results are provided in Section 4. Finally, we conclude the paper in Section 5.

## 2. PRELIMINARIES

### 2.1. Gaussian Processes (GPs)

A GP defines a collection of random variables indexed by  $\mathcal{X} \subseteq \mathbb{R}^{d_x}$ , where any finite subset of these variables follows a joint Gaussian distribution [8]. Typically, a GP is employed to represent a distribution over random functions  $f(\mathbf{x}) : \mathbb{R}^{d_x} \mapsto \mathbb{R}$ , given by:

$$f(\mathbf{x}) \sim \mathcal{GP}(\mu(\mathbf{x}), k_{\theta_{gp}}(\mathbf{x}, \mathbf{x}')), \quad (1)$$

This work was supported by the NSFC with Grant No. 62271433 and in part by Guangdong Research Project with Grant No. 2017ZT07X152. The corresponding author is Feng Yin ([yinfeng@cuhk.edu.cn](mailto:yinfeng@cuhk.edu.cn)).

where  $\mu(\mathbf{x})$  is a mean function, often set to zero in practice;  $k_{\theta_{gp}}(\mathbf{x}, \mathbf{x}')$  is a covariance/kernel function;  $\theta_{gp}$  represents a set of hyperparameters that are tuned for model selection. In the rest of the paper,  $\theta_{gp}$  is omitted for the sake of notation brevity. By applying Bayes' theorem, the function prior is combined with new data to obtain an analytical posterior distribution. Specifically, given a noise-free training dataset  $\mathcal{D} = \{X, \mathbf{f}\} = \{\mathbf{x}_i, \mathbf{f}_i\}_{i=1}^n$ , the posterior distribution  $p(f(\mathbf{x}_*)|\mathbf{x}_*, \mathcal{D})$  at any test input  $\mathbf{x}_* \in \mathcal{X}$  follows a Gaussian distribution, fully characterized by the posterior mean  $\xi$  and the posterior variance  $\Xi$ :

$$\xi(\mathbf{x}_*) = \mathbf{K}_{\mathbf{x}_*, X} \mathbf{K}_{X, X}^{-1} \mathbf{f}, \quad (2a)$$

$$\Xi(\mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x}_*) - \mathbf{K}_{\mathbf{x}_*, X} \mathbf{K}_{X, X}^{-1} \mathbf{K}_{X, \mathbf{x}_*}^\top, \quad (2b)$$

where  $\mathbf{K}_{X, X}$  represents the covariance matrix evaluated on the training input  $X$ , with each entry given by  $[\mathbf{K}_{X, X}]_{i, j} = k(\mathbf{x}_i, \mathbf{x}_j)$ ;  $\mathbf{K}_{\mathbf{x}_*, X}$  denotes the cross-covariance matrix between the test input  $\mathbf{x}_*$  and the training input  $X$ .

## 2.2. Gaussian Process State-Space Model (GPSSM)

A generic state-space model (SSM) characterizes the probabilistic relationship between the latent state,  $\mathbf{x}_t \in \mathbb{R}^{d_x}$ , and the observation,  $\mathbf{y}_t \in \mathbb{R}^{d_y}$ . Mathematically, it is represented as follows:

$$\mathbf{x}_{t+1} = f(\mathbf{x}_t) + \mathbf{v}_t, \quad \mathbf{y}_t = g(\mathbf{x}_t) + \mathbf{e}_t, \quad (3)$$

where  $\mathbf{v}_t$  and  $\mathbf{e}_t$  are additive noise terms, and  $f(\cdot)$  and  $g(\cdot)$  are referred to as the *transition function* and *emission function*, respectively.

Incorporating a GP prior over the transition function  $f(\cdot)$  and assuming a parametric emission function  $g(\cdot)$  in the classic SSM (see Eq. (3)) leads to the well-known GPSSM<sup>1</sup> [5], expressed as follows:

$$f(\cdot) \sim \mathcal{GP}(\mathbf{0}, k(\cdot, \cdot)), \quad \mathbf{f}_t = f(\mathbf{x}_{t-1}), \quad \mathbf{x}_0 \sim p(\mathbf{x}_0), \quad (4a)$$

$$\mathbf{x}_t | \mathbf{f}_t \sim \mathcal{N}(\mathbf{x}_t | \mathbf{f}_t, \mathbf{Q}), \quad \mathbf{y}_t | \mathbf{x}_t \sim \mathcal{N}(\mathbf{y}_t | \mathbf{C}\mathbf{x}_t, \mathbf{R}), \quad (4b)$$

where the emission model is assumed to be known and linear and the coefficient matrix,  $\mathbf{C} \in \mathbb{R}^{d_y \times d_x}$ , mitigates the system non-identifiability [5]. The initial state prior distribution  $p(\mathbf{x}_0)$  is also known and assumed to follow a Gaussian distribution. Both state transitions and observations are corrupted by zero-mean Gaussian noise with covariance matrices  $\mathbf{Q}$  and  $\mathbf{R}$ , respectively. In the case of state dimension  $d_x > 1$ , the transition  $f(\cdot) : \mathbb{R}^{d_x} \mapsto \mathbb{R}^{d_x}$  is typically modeled with  $d_x$  mutually independent GPs. Each independent GP represents a dimension-specific function  $f_d(\cdot) : \mathbb{R}^{d_x} \mapsto \mathbb{R}$ , and the multivariate output is denoted as

$$\mathbf{f}_t = f(\mathbf{x}_{t-1}) \triangleq \{f_d(\mathbf{x}_{t-1})\}_{d=1}^{d_x} \triangleq \{\mathbf{f}_{t,d}\}_{d=1}^{d_x}, \quad (5)$$

where each independent GP is associated with a unique kernel function and the associated hyperparameters. The challenging task of GPSSM lies in simultaneously learning the transition function and the noise parameters, i.e., learning  $[\theta_{gp}, \mathbf{Q}, \mathbf{R}]$ , while inferring the latent states of interest.

Despite the popularity of the aforementioned GPSSM modeling approach, it is plagued by two significant caveats as outlined in Section 1, namely the escalating burden of complexity (encompassing computational cost and size of the parameter space), and the model output independence. In the following section, we will present our novel solution that tackles the challenge of growing complexity while also establishing output dependence.

## 3. EFFICIENT GPSSM MODELING AND INFERENCE

Section 3.1 introduces our novel efficient GPSSM formulation tailored for multi-dimensional latent state spaces. Section 3.2 presents

<sup>1</sup>The GPSSM considered in this paper keeps the same model capacity as the ones with both transition and emission GPs while avoiding the severe *non-identifiability* issue. One can refer to [5] (Section 3.2.1) for more details.

the proposed variational algorithm for learning and inference in the efficient GPSSM.

### 3.1. Efficient GPSSM Modeling

In the context of a multi-dimensional latent state space, our GPSSM formulation departs from the conventional use of  $d_x$  separate GPs for modeling the transition function. Instead, we incorporate a more efficient approach, ETGP [27]. This entails utilizing a single GP and employing normalizing flow [26] to transform the process within each dimension of the latent state. It is noteworthy that the ETGP has already showcased successful applications in large-scale multi-classification scenarios [27]. In essence, the new GPSSM formulation can be expressed as:

$$\tilde{\mathbf{f}} \sim \mathcal{GP}(0, k(\cdot, \cdot)), \quad \tilde{\mathbf{f}}_t = \tilde{\mathbf{f}}(\mathbf{x}_{t-1}), \quad \mathbf{f}_t = \{\mathbb{G}_{\theta_d}(\tilde{\mathbf{f}}_t)\}_{d=1}^{d_x}, \quad (6a)$$

$$\mathbf{x}_0 \sim p(\mathbf{x}_0), \quad \mathbf{x}_t | \mathbf{f}_t \sim \mathcal{N}(\mathbf{x}_t | \mathbf{f}_t, \mathbf{Q}), \quad \mathbf{y}_t | \mathbf{x}_t \sim \mathcal{N}(\mathbf{y}_t | \mathbf{C}\mathbf{x}_t, \mathbf{R}), \quad (6b)$$

where the dimension-specific normalizing flow,  $\mathbb{G}_{\theta_d}(\cdot)$ , comprises a set of invertible and differentiable transformations,  $\{\mathbb{G}_{\theta_{d,j}}(\cdot) : \mathbb{R} \mapsto \mathbb{R}\}_{j=0}^{J-1}$ , parameterized by  $\theta_d \triangleq \{\theta_{d,j}\}_{j=1}^J$ ,  $J \in \mathbb{N}$ . Due to these time-wise mappings, any finite  $T$ -dimensional multivariate random variable  $\tilde{\mathbf{f}} \triangleq \mathbf{f}_{1:T} \triangleq \{\mathbf{f}_t\}_{t=1}^T$  is guaranteed from the induced Gaussian copula processes [29].

**Remark 1.** *The new GPSSM formulation presented in Eq. (6) maintains constant computational complexity when the computational cost of the normalizing flows becomes negligible compared to the GP part. Consequently, the overall computational burden aligns with that of calculating a single GP, resulting in a highly efficient method. Another advantage is the establishment of dependency among the resulting  $d_x$  Gaussian copula processes [29]. This dependency arises due to their shared basis GP, denoted as  $\tilde{\mathbf{f}}$ , while simultaneously exhibiting distinct marginal distributions [27]. Such correlated yet individually varying processes enhance the modeling capacity and enable a more comprehensive representation of complex relationships within the latent state space.*

In principle, the GPSSM can integrate diverse forms of normalizing flows, spanning from basic and interpretable elementary flows to advanced flows developed in recent years [26]. This paper investigates two commonly used elementary flows. The first one is the Sinh-Arcsinh-Linear (SAL) flow [25, 30], which can be applied by stacking  $J \in \mathbb{N}$  layers, with each layer represented as follows:

$$\mathbb{G}_{\theta_{d,j}}(\cdot) = \alpha_{d,j} \sinh(\varphi_{d,j} \operatorname{arcsinh}(\cdot) - \gamma_{d,j}) + \beta_{d,j}, \quad (7)$$

where  $\theta_{d,j} \triangleq [\alpha_{d,j}, \beta_{d,j}, \gamma_{d,j}, \varphi_{d,j}]$ ,  $d \in \{1, 2, \dots, d_x\}$ ,  $j \in \{0, 1, \dots, J-1\}$ . The SAL flow can be employed to control the mean, variance, asymmetry, and kurtosis of the function priors [30]. Another one is the simple linear flow, that is,

$$\mathbf{f}_{t,d} = \alpha_d \cdot \tilde{\mathbf{f}}_t + \beta_d, \quad d \in \{1, 2, \dots, d_x\}, \quad (8)$$

with  $\theta_d = [\alpha_d, \beta_d]$ . When using linear flow, it is straightforward to see that the  $d_x$  outputs are dependent GPs (see Proposition 1 in [27]), such that for any two dimensions  $d$  and  $d'$ , we have

$$\mathbb{E}[\mathbf{f}_{t,d}] = \beta_d, \quad \operatorname{Cov}[\mathbf{f}_{t,d}, \mathbf{f}_{t,d'}] = \alpha_d \alpha_{d'} k(\mathbf{x}_{t-1}, \mathbf{x}_{t-1}). \quad (9)$$

Both of these commonly used normalizing flows exhibit linear run-time complexity and involve only a small amount of parameters, whose computational effect is practically negligible when contrasted with the complexity of separate GPs in existing GPSSMs. More importantly, the computation of the  $d_x$  flows can be parallelized through batched matrix operations, in cases where the flows share the same functional form, which leads to a substantial enhancement in model efficiency.

### 3.2. Efficient GPSSM Inference

In addition to the computational efficiency enhancement in the efficient GPSSM, as outlined in Section 3.1, we employ sparse GP methods [28] to further alleviate the computational load that stems from the GP model with a large sample count. This involves introducing a compact set of inducing points  $\bar{\mathbf{z}} \triangleq \{\mathbf{z}_i\}_{i=1}^m$  and  $\bar{\mathbf{u}} \triangleq \{\mathbf{u}_i\}_{i=1}^m$ ,  $m \ll T$ , that act as a surrogate for the corresponding GP. Here,  $\mathbf{u}_i = \tilde{f}(\mathbf{z}_i) \in \mathbb{R}$  and  $\mathbf{z}_i \in \mathbb{R}^{d_x}$ . To facilitate discussions, we define  $\bar{\mathbf{x}} \triangleq \{\mathbf{x}_t\}_{t=0}^T$  and  $\bar{\mathbf{y}} \triangleq \{\mathbf{y}_t\}_{t=1}^T$ . Based on these configurations, the joint distribution of the proposed GPSSM, augmented with inducing points, is

$$p(\bar{\mathbf{y}}, \bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}}) = p(\mathbf{x}_0) \prod_{t=1}^T p(\mathbf{y}_t | \mathbf{x}_t) p(\mathbf{x}_t | \mathbf{f}_t) p(\mathbf{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}}) p(\bar{\mathbf{u}}), \quad (10)$$

where  $p(\bar{\mathbf{u}}) = \mathcal{N}(\bar{\mathbf{u}} | \mathbf{0}, \mathbf{K}_{\bar{\mathbf{z}}, \bar{\mathbf{z}}})$ , and the distribution of the transition function outputs  $p(\mathbf{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}})$  are determined by the shared GP  $p(\tilde{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}})$  and the normalizing flows, where

$$p(\tilde{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}}) = \mathcal{N}(\tilde{f}_t | \xi(\mathbf{x}_{t-1}), \Xi(\mathbf{x}_{t-1})) \quad (11)$$

is the GP posterior distribution with  $\mathbf{x}_{t-1}$  being the test input while  $(\bar{\mathbf{z}}, \bar{\mathbf{u}})$  being the training data, see Eqs. (2a) and (2b); Therefore,

$$p(\mathbf{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}}) = \underbrace{p(\tilde{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}})}_{=p(\mathbf{f}_{t,1} | \bar{\mathbf{u}}, \mathbf{x}_{t-1})} \cdot \mathbf{J}_{\mathbf{f}} \cdot \underbrace{\prod_{d=2}^{d_x} \delta(\mathbf{f}_{t,d} - \mathbb{G}_{\theta,d} \circ \mathbb{G}_{\theta,1}^{-1}(\mathbf{f}_{t,1}))}_{=p(\mathbf{f}_{t,d} | \mathbf{f}_{t,1}, \bar{\mathbf{u}}, \mathbf{x}_{t-1})}, \quad (12)$$

where  $\mathbf{J}_{\mathbf{f}} \triangleq \prod_{j=1}^{J-1} \left| \det \frac{\partial \mathbb{G}_{\theta,1,j}(\mathbb{G}_{\theta,1,j-1}(\dots \mathbb{G}_{\theta,1,0}(\tilde{f}_t) \dots))}{\partial \mathbb{G}_{\theta,1,j-1}(\dots \mathbb{G}_{\theta,1,0}(\tilde{f}_t) \dots)} \right|^{-1}$ , and  $\delta(\cdot)$  denotes the Dirac measure. Note that this decomposition is not unique. We designate the first dimension,  $\mathbf{f}_{t,1}$ , as the *pivot*; and the joint distribution can be expressed equivalently with respect to any other pivot  $\mathbf{f}_{t,d}$  for the case  $d \neq 1$  (see details in Appendix A, [27]).

Learning and inference within the GPSSM is plagued by the intractability of  $p(\bar{\mathbf{y}}) = \int p(\bar{\mathbf{y}}, \bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}}) d\bar{\mathbf{x}} d\bar{\mathbf{f}} d\bar{\mathbf{u}}$ . In order to surmount this challenge, instead of relying on Monte Carlo-based methods [1], this paper embraces alternative variational inference techniques [12, 13]. Variational inference methods involve approximating the intractable posterior distribution  $p(\bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}} | \bar{\mathbf{y}}) = \frac{p(\bar{\mathbf{y}}, \bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})}{p(\bar{\mathbf{y}})}$  with a variational distribution  $q(\bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})$ , leading to a learning and inference objective function, the evidence lower bound (ELBO), denoted as  $\mathcal{L}$ ,

$$\mathcal{L} \triangleq \mathbb{E}_{q(\bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})} \left[ \log \frac{p(\bar{\mathbf{y}}, \bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})}{q(\bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})} \right] \leq \log p(\bar{\mathbf{y}}). \quad (13)$$

The choice of the variational distribution influences the tightness of the ELBO and, in turn, determines the learning algorithm for GPSSM [5]. In this paper, we adopt a specific form for the variational distribution  $q(\bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})$  given by  $q(\mathbf{x}_0) \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{f}_t) p(\mathbf{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}}) q(\bar{\mathbf{u}})$ , where  $q(\bar{\mathbf{u}}) = \mathcal{N}(\bar{\mathbf{u}} | \mathbf{m}, \mathbf{S})$  with  $\mathbf{m} \in \mathbb{R}^m$  and the covariance matrix  $\mathbf{S} \in \mathbb{R}^{m \times m}$  representing the free variational parameters. Similarly, the variational distribution for the initial state  $\mathbf{x}_0$  is given by  $q(\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_0 | \mathbf{m}_{\mathbf{x}_0}, \mathbf{S}_{\mathbf{x}_0})$ . Here,  $\mathbf{m}_{\mathbf{x}_0}$  and  $\mathbf{S}_{\mathbf{x}_0}$  can be treated as free variational parameters or learned through an amortized recognition network [18] with inputs  $\bar{\mathbf{y}}$  and parameters  $\zeta$ . Note that the variational distribution form adopted in this paper bears a resemblance to the form used in probabilistic recurrent SSM (PRSSM) [6] and the output-dependent GPSSM (ODGPSSM) [20]. This deliberate decision facilitates a direct comparison with these methods, enabling us to thoroughly evaluate and analyze the advantages of our GPSSM modeling approach.

With the proposed variational distribution,  $q(\bar{\mathbf{x}}, \bar{\mathbf{f}}, \bar{\mathbf{u}})$ , and after

**Table 1.** Comparisons between various variational GPSSMs.  $c$  denotes the number of parameter shared by all models (including  $[\zeta, \mathbf{Q}, \mathbf{R}]$ );  $Q$  is the number of latent GPs in the ODGPSSM, which is typically  $\geq d_x$ ;  $\eta$  represents the parameter counts associated with the normalizing flows. For instance, in the case of a 2-layer SAL flow,  $\eta = 8$ .

	comput. complexity	# parameters	output-dependent
PRSSM [6]	$\mathcal{O}(d_x T m^2)$	$c + d_x  \theta_{gp}  + m d_x (2d_x + m + 4)/2$	$\times$
ODGPSSM [20]	$\mathcal{O}(Q T m^2)$	$c + Q  \theta_{gp}  + m Q (2Q + m + 4)/2 + Q d_x$	$\checkmark$
EGPSSM (ours)	$\mathcal{O}(T m^2)$	$c +  \theta_{gp}  + m(2d_x + m + 4)/2 + \eta d_x$	$\checkmark$

performing certain algebraic calculations, the ELBO is obtained as,

$$\mathcal{L} = \sum_{t=1}^T \mathbb{E}_{q(\mathbf{x}_t)} [\log p(\mathbf{y}_t | \mathbf{x}_t)] - \text{KL}(q(\bar{\mathbf{u}}) || p(\bar{\mathbf{u}})) - \text{KL}(q(\mathbf{x}_0) || p(\mathbf{x}_0)), \quad (14)$$

where the first term encourages the latent states drawn from the variational distribution,

$$q(\mathbf{x}_t) = \int q(\mathbf{x}_0) \prod_{\tau=1}^t p(\mathbf{x}_\tau | \mathbf{f}_\tau) p(\mathbf{f}_\tau | \mathbf{x}_{\tau-1}, \bar{\mathbf{u}}) q(\bar{\mathbf{u}}) d\mathbf{x}_{0:t-1} d\bar{\mathbf{f}} d\bar{\mathbf{u}},$$

to fit the emission model effectively, ensuring a good match between the observed data and the latent states. The second and third terms act as regularizations, controlling the posterior distributions of the initial state and the shared GP (and consequently, the posterior of ETGP  $f(\cdot)$ ). These regularization terms play a crucial role in managing the model complexity and preventing overfitting.

The two regularization terms can be computed analytically due to the Gaussian nature of the distributions involved. However, the expectation terms,  $\mathbb{E}_{q(\mathbf{x}_t)} [\log p(\mathbf{y}_t | \mathbf{x}_t)]$ ,  $\forall t$ , need to be evaluated by the sampling method and reparametrization trick [31] due to the intractability of  $q(\mathbf{x}_t)$  [6]. The samples,  $\mathbf{x}_t$ , can be obtained as follows. We first use the reparametrization trick to sample the (1-D) GP function value,  $\tilde{f}_t$ , from  $q(\tilde{f}_t | \mathbf{x}_{t-1})$  by conditioning on the latent state  $\mathbf{x}_{t-1}$ , where

$$q(\tilde{f}_t | \mathbf{x}_{t-1}) = \mathbb{E}_{q(\bar{\mathbf{u}})} [p(\tilde{f}_t | \mathbf{x}_{t-1}, \bar{\mathbf{u}})] = \mathcal{N}(\tilde{f}_t | \mu_{t|t-1}, \mathbf{S}_{t|t-1}), \quad (15)$$

with  $\mu_{t|t-1} = \mathbf{K}_{\mathbf{x}_{t-1}, \bar{\mathbf{z}}} \mathbf{K}_{\bar{\mathbf{z}}, \bar{\mathbf{z}}}^{-1} \mathbf{m}$ , and

$$\mathbf{S}_{t|t-1} = k(\mathbf{x}_{t-1}, \mathbf{x}_{t-1}) - \mathbf{K}_{\mathbf{x}_{t-1}, \bar{\mathbf{z}}} \mathbf{K}_{\bar{\mathbf{z}}, \bar{\mathbf{z}}}^{-1} [\mathbf{K}_{\bar{\mathbf{z}}, \bar{\mathbf{z}}} - \mathbf{S}] \mathbf{K}_{\bar{\mathbf{z}}, \bar{\mathbf{z}}}^{-1} \mathbf{K}_{\mathbf{x}_{t-1}, \bar{\mathbf{z}}}^\top.$$

We then push the sampled function value  $\tilde{f}_t$  into the  $d_x$  normalizing flows,  $\{\mathbb{G}_{\theta,d}(\cdot)\}_{d=1}^{d_x}$  and get the transformed GP function value  $\mathbf{f}_t$ , as shown in Eq. (6a). Lastly, by conditioning on the sampled  $\mathbf{f}_t$ , we can sample the latent state  $\mathbf{x}_t$  from  $p(\mathbf{x}_t | \mathbf{f}_t)$ . In this way, we are able to numerically evaluate the first term, thereby enabling us to maximize the ELBO with respect to the variational parameters,  $[\zeta, \bar{\mathbf{z}}, \mathbf{m}, \mathbf{S}]$ , and the model parameters  $[\theta_{gp}, \{\theta_d\}_{d=1}^{d_x}, \mathbf{Q}, \mathbf{R}]$ .

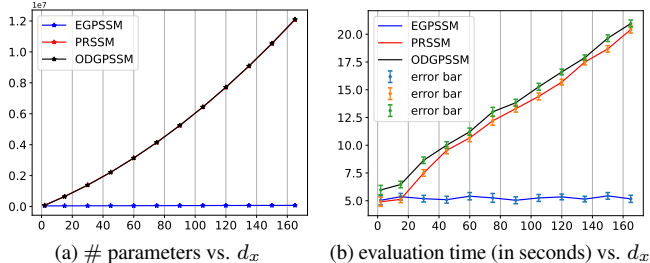
So far, we have introduced the proposed method, encompassing efficient modeling and inference for GPSSM. We can compare the proposed method with some existing approaches, as summarized in Table 1. It is imperative to highlight that, in contrast to the existing approaches, our method significantly mitigates the escalating computational burden and parameter proliferation, resulting in enhanced efficiency for both modeling and inference. Meanwhile, the proposed method also establishes entangled relationships among the outputs because of the shared Gaussian copula [29].

## 4. NUMERICAL EXPERIMENTS

In this section, we present the performance evaluation of the proposed efficient GPSSM, referred to as EGPSSM, across synthetic and real datasets. Given that the superiority of GPSSM over traditional time-series modeling methods has been shown in [6, 20], we will

**Table 2.** Prediction RMSE comparison between the proposed EGPSSM and the competitors on synthetic sequence. Shown are mean and standard errors over five repetitions.

PRSSM	ODGPSSM	EGPSSM (L)	EGPSSM (SAL)
0.389±0.022	0.388±0.031	<b>0.387±0.029</b>	0.392±0.023



**Fig. 1.** Complexity comparisons between different GPSSMs

bypass redundant comparisons due to space limitations. Thus for the benchmark comparison, we consider two GPSSMs: 1) PRSSM [6]; and 2) ODGPSSM [20].

#### 4.1. Synthetic Data

This subsection adopts the underlying two-dimensional SSM described below as the testing model,

$$f(\mathbf{x}_t) = 0.8 + (\mathbf{x}_{t,1} + 0.2) \left[ 1 - \frac{5}{1 + \exp(-2\mathbf{x}_{t,1})} \right] + \mathbf{x}_{t,2}, \quad (16a)$$

$$\mathbf{x}_{t+1} = \begin{bmatrix} f(\mathbf{x}_t) \\ -0.5 * f(\mathbf{x}_t) \end{bmatrix} + \mathbf{x}_t + \mathbf{v}_t, \quad \mathbf{v}_t \sim \mathcal{N}(\mathbf{0}, \sqrt{0.001} \mathbf{I}_2), \quad (16b)$$

$$\mathbf{y}_t = \mathbf{x}_t + \mathbf{e}_t, \quad \mathbf{e}_t \sim \mathcal{N}(\mathbf{0}, \sqrt{0.01} \mathbf{I}_2), \quad (16c)$$

where  $f(\mathbf{x}_t)$  is a modified kink function that is commonly used in the GPSSM literature [7]. We commence by showcasing the computational complexity and the number of parameters of different GPSSMs across different latent state dimensions. The corresponding results are illustrated in Fig. 1. The computational time presented in Fig. 1b is the wall time required for the GPSSMs to evaluate the ELBO once on a sequence with a length of 200. The results are obtained from five repeated experiments and displayed with mean and standard deviations. All the GPSSMs employ 200 inducing points and utilize the Matérn kernel [8]; EGPSSM integrates a 2-layer SAL flow. See our implemented code for more details.

Based on the results presented in Fig. 1, it is evident that both PRSSM and ODGPSSM suffer from the escalating complexity. For example, PRSSM and ODGPSSM exhibit a substantial number of model parameters, nearly reaching 2 million, even when operating within a moderate latent state dimension of approximately 40. This significant parameter volume poses formidable optimization challenges. In contrast, EGPSSM offers distinct advantages over PRSSM and ODGPSSM, primarily attributed to its utilization of a shared GP and normalizing flow with a small number of parameters. This design choice contributes to enhanced computational efficiency and more efficient parameter management within the EGPSSM framework.

Subsequently, we set the latent state dimension to 2 for all GPSSMs, encompassing PRSSM, ODGPSSM, linear-flow EGPSSM (denoted as EGPSSM (L)), and SAL-flow EGPSSM (denoted as EGPSSM (SAL)). These models are then deployed for sequence prediction tasks. In this context, we generated 10 training sequences of length 50, along with an additional test sequence of the same length. Our objective is to ascertain whether EGPSSM exhibits favorable inference performance. The corresponding prediction root-mean-

**Table 3.** Prediction RMSE comparison between the proposed EGPSSM and the competitors on the five system identification datasets. Shown are mean and standard errors over five repetitions.

Methods	Actuator	Ballbeam	Drive	Dryer	Gas Furnace
PRSSM	0.691±0.148	0.074±0.010	<b>0.647±0.057</b>	0.174±0.013	<b>1.503±0.196</b>
ODGPSSM	<b>0.666±0.074</b>	0.068±0.006	0.708±0.052	<b>0.171±0.011</b>	1.704±0.560
EGPSSM (L)	0.742±0.050	<b>0.055±0.005</b>	0.756±0.020	0.482±0.027	1.994±0.085
EGPSSM (SAL)	0.758±0.048	<b>0.054±0.003</b>	0.762±0.020	0.479±0.009	2.010±0.071

square error (RMSE) results are presented in Table 2, revealing a notably comparable performance across all approaches. Notably, EGPSSM (L) displays a slightly superior performance, likely because it matches the linear relationship inherent in the transition function outputs of the underlying SSM, as shown in Eq. (16b). It is also worth noting that the existing GPSSMs possess sufficient model capability to accurately predict the sequence, although their intrinsic transition modeling is inconsistent with the underlying one and entails higher computational complexity and parameter count than EGPSSM.

#### 4.2. Real Data

This subsection presents a comparative analysis involving EGPSSM, PRSSM, and ODGPSSM across five real system identification datasets, as detailed in [6]. In each dataset, the initial half of the sequence serves as training data, while the latter portion is allocated for testing purposes. Standardization of all datasets is performed using the training sequence, and a consistent latent state dimension of  $d_x = 4$  is employed. The SAL-flow EGPSSM incorporates a 2-layer SAL flow. Further elaboration on the specific parameters is accessible through the associated online code repository. The prediction results are reported in Table 3, wherein the RMSE is averaged over a 50-step forecasting. Table 3 illustrates that, generally, EGPSSMs closely match the performance of the comparative models on the real datasets, with some instances even surpassing PRSSM and ODGPSSM, as seen in the *Ballbeam* dataset. However, for certain datasets, EGPSSMs exhibit a slightly marginal disadvantage in sequence prediction. This can predominantly be attributed to the incorporation of elementary flows in this paper, which introduces the potential challenge of limited model flexibility. Nonetheless, it is essential to acknowledge that EGPSSMs achieve comparable performance, while effectively managing model parameters and ensuring computational efficiency, particularly in settings characterized by moderate and high dimensional latent states. Subsequent research endeavors could focus on augmenting the model flexibility in the EGPSSM while preserving its efficiency. One prospective avenue for further improvement involves integrating more sophisticated normalizing flows to transform the shared GP. However, delving into this matter, also, requires theoretical investigation into how the outputs of the ETGP correlate when using more complex flows. These aspects will be left for future investigations.

### 5. CONCLUSION

This paper introduces a novel and efficient GPSSM, EGPSSM, tailored for modeling dynamical systems with high-dimensional latent states. Unlike existing approaches, our method can effectively mitigate the challenges of escalating computational complexity and parameter proliferation. Empirical evaluations, conducted across diverse synthetic and real-world datasets, substantiate the merits of EGPSSM, revealing its superior modeling efficiency and significant reduction in parameter count. Furthermore, our results demonstrate that EGPSSM achieves comparable performance to existing GPSSMs, underscoring its competitiveness in inference tasks. Nevertheless, the model flexibility might be somewhat constrained in intricate scenarios. As such, future work should focus on enhancing model flexibility while preserving computational efficiency.

## 6. REFERENCES

- [1] R. Frigola, F. Lindsten, T. B. Schön, and C. E. Rasmussen, “Bayesian inference and learning in Gaussian process state-space models with particle MCMC,” in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, Dec. 2013, pp. 3156–3164.
- [2] R. Krishnan, U. Shalit, and D. Sontag, “Structured inference networks for nonlinear state space models,” in *Proc. AAAI Conf. Artif. Intell. (AAAI)*, Feb. 2017, pp. 2101–2109.
- [3] A. M. Alaa and M. van der Schaar, “Attentive state-space modeling of disease progression,” in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, Dec. 2019, pp. 11 338–11 348.
- [4] R. Frigola, Y. Chen, and C. E. Rasmussen, “Variational Gaussian process state-space models,” in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, Dec. 2014, pp. 3680–3688.
- [5] R. Frigola, “Bayesian time series learning with Gaussian processes,” Ph.D. dissertation, University of Cambridge, 2015.
- [6] A. Doerr, C. Daniel, M. Schiegg, N.-T. Duy, S. Schaal, M. Toussaint, and T. Sebastian, “Probabilistic recurrent state-space models,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, Jul. 2018, pp. 1280–1289.
- [7] A. D. Ialongo, M. van der Wilk, J. Hensman, and C. E. Rasmussen, “Overcoming mean-field approximations in recurrent Gaussian process models,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, Jun. 2019, pp. 2931–2940.
- [8] C. E. Rasmussen and C. K. I. Williams, *Gaussian Processes for Machine Learning*. MIT Press, 2006.
- [9] R. C. Suwandi, Z. Lin, Y. Sun, Z. Wang, L. Cheng, and F. Yin, “Gaussian process regression with grid spectral mixture kernel: Distributed learning for multidimensional data,” in *Proc. Int. Conf. Inf. Fusion (FUSION)*, July 2022, pp. 1–8.
- [10] F. Yin, L. Pan, T. Chen, S. Theodoridis, Z.-Q. Luo, and A. M. Zoubir, “Linear multiple low-rank kernel based stationary Gaussian processes regression for time series,” *IEEE Trans. Signal Process.*, vol. 68, pp. 5260–5275, Sep. 2020.
- [11] S. Särkkä, *Bayesian Filtering and Smoothing*. Cambridge University Press, 2013.
- [12] L. Cheng, F. Yin, S. Theodoridis, S. Chatzis, and T.-H. Chang, “Rethinking Bayesian learning for data analysis: The art of prior and inference in sparsity-aware modeling,” *IEEE Signal Process. Mag.*, vol. 39, no. 6, pp. 18–52, Oct. 2022.
- [13] S. Theodoridis, *Machine Learning: A Bayesian and Optimization Perspective*, 2nd ed. Academic Press, 2020.
- [14] J. M. Wang, D. J. Fleet, and A. Hertzmann, “Gaussian process dynamical models for human motion,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 30, no. 2, pp. 283–298, Dec. 2007.
- [15] A. Xie, F. Yin, B. Ai, S. Zhang, and S. Cui, “Learning while tracking: A practical system based on variational Gaussian process state-space model and smartphone sensory data,” in *Proc. Int. Conf. Inf. Fusion (FUSION)*, Jul. 2020, pp. 1–7.
- [16] F. Yin, Z. Lin, Q. Kong, Y. Xu, D. Li, S. Theodoridis, and S. Cui, “Fedloc: Federated learning framework for data-driven cooperative localization and location data processing,” *IEEE Open J. Signal Process.*, vol. 1, pp. 187–215, Nov. 2020.
- [17] Y. Zhao, C. Fritsche, G. Hendeby, F. Yin, T. Chen, and F. Gunnarsson, “Cramér–Rao bounds for filtering based on Gaussian process state-space models,” *IEEE Trans. Signal Process.*, vol. 67, no. 23, pp. 5936–5951, Dec. 2019.
- [18] S. Eleftheriadis, T. Nicholson, M. P. Deisenroth, and J. Hensman, “Identification of Gaussian process state space models,” in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, Dec. 2017, pp. 5309–5319.
- [19] J. Lindinger, B. Rakitsch, and C. Lippert, “Laplace approximated Gaussian process state-space models,” in *Proc. Conf. Uncertain. Artif. Intell. (UAI)*, Aug. 2022.
- [20] Z. Lin, L. Cheng, F. Yin, L. Xu, and S. Cui, “Output-dependent Gaussian process state-space model,” in *Proc. IEEE Int. Conf. Acoust. Speech Signal Process. (ICASSP)*, 2023, pp. 1–5.
- [21] Z. Lin, F. Yin, and J. Maroñas, “Towards flexibility and interpretability of Gaussian process state-space model,” *arXiv preprint arXiv:2301.08843*, 2023.
- [22] Y. Liu and P. M. Djurić, “Gaussian process state-space models with time-varying parameters and inducing points,” in *Proc. European Signal Process. Conf. (EUSIPCO)*, Jan. 2021, pp. 1462–1466.
- [23] X. Fan, E. V. Bonilla, T. O’Kane, and S. A. Sisson, “Free-form variational inference for Gaussian process state-space models,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, July 2023, pp. 9603–9622.
- [24] K. Chen, T. van Laarhoven, E. Marchiori, F. Yin, and S. Cui, “Multitask Gaussian process with hierarchical latent interactions,” in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process Proc. (ICASSP)*, May 2022, pp. 4148–4152.
- [25] J. Maroñas, O. Hamelijncck, J. Knoblauch, and T. Damoulas, “Transforming Gaussian processes with normalizing flows,” in *Proc. Int. Conf. Artif. Intell. Stat. (AISTATS)*, Apr. 2021, pp. 1081–1089.
- [26] I. Kobyzev, S. J. Prince, and M. A. Brubaker, “Normalizing flows: An introduction and review of current methods,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 43, no. 11, pp. 3964–3979, 2020.
- [27] J. Maroñas and D. Hernández-Lobato, “Efficient transformed Gaussian processes for non-stationary dependent multi-class classification,” in *Proc. Int. Conf. Mach. Learn. (ICML)*, July 2023, pp. 24 045–24 081.
- [28] J. Hensman, N. Fusi, and N. D. Lawrence, “Gaussian processes for big data,” in *Proc. Conf. Uncertain. Artif. Intell. (UAI)*, Jul. 2013, pp. 282–290.
- [29] A. G. Wilson and Z. Ghahramani, “Copula processes,” in *Proc. Adv. Neural Inf. Process. Syst. (NeurIPS)*, Dec. 2010, pp. 2460–2468.
- [30] G. Rios and F. Tobar, “Compositionally-warped Gaussian processes,” *Neural Networks*, vol. 118, pp. 235–246, 2019.
- [31] D. P. Kingma and M. Welling, “An introduction to variational autoencoders,” *Found. Trends Mach. Learn.*, vol. 12, no. 4, pp. 307–392, 2019.