

---

# Deep Functional Factor Models: Forecasting High-Dimensional Functional Time Series via Bayesian Nonparametric Factorization

---

Yirui Liu<sup>1,2</sup> Xinghao Qiao<sup>3</sup> Yulong Pei<sup>1</sup> Liying Wang<sup>4</sup>

## Abstract

This paper introduces the Deep Functional Factor Model (DF<sup>2</sup>M), a Bayesian nonparametric model designed for analysis of high-dimensional functional time series. DF<sup>2</sup>M is built upon the Indian Buffet Process and the multi-task Gaussian Process, incorporating a deep kernel function that captures non-Markovian and nonlinear temporal dynamics. Unlike many black-box deep learning models, DF<sup>2</sup>M offers an explainable approach to utilizing neural networks by constructing a factor model and integrating deep neural networks within the kernel function. Additionally, we develop a computationally efficient variational inference algorithm to infer DF<sup>2</sup>M. Empirical results from four real-world datasets demonstrate that DF<sup>2</sup>M provides better explainability and superior predictive accuracy compared to conventional deep learning models for high-dimensional functional time series.

## 1. Introduction

Functional time series refers to a sequential collection of functional objects that exhibit temporal dependence, and this area of study has garnered increasing attention in recent years. With the advancements in data collection technology and computational power, high-dimensional datasets containing numerous functional time series have become more prevalent. Examples of such data include annual age-specific mortality rates for different countries, daily energy consumption curves from various households, and cumulative intraday return trajectories for hundreds of stocks. These datasets can be represented as  $p$ -dimensional functional time series  $\mathbf{Y}_t(\cdot) = (Y_{t1}(\cdot), \dots, Y_{tp}(\cdot))^T$ , where

---

<sup>1</sup>J.P. Morgan <sup>2</sup>London School of Economics and Political Science <sup>3</sup>Faculty of Business and Economics, The University of Hong Kong <sup>4</sup>Management School, University of Liverpool. Correspondence to: Xinghao Qiao <xinghaoq@hku.hk>.

each  $Y_{tj}(\cdot)$  is a random function defined on a compact interval  $\mathcal{U}$ . The number of functional variables  $p$  is comparable to, or even larger than, the number of temporally dependent observations  $n$ . Analyzing high-dimensional functional time series presents a challenging task, as it necessitates the use of dimension reduction techniques to address the high-dimensional problem, functional approaches to handle the infinite-dimensional nature of curve data, and time series modeling methods to capture the temporal dependence.

Several statistical methods have been proposed to address these challenges, such as those presented in Gao et al. (2019); Chen et al. (2022); Fang et al. (2022); Chang et al. (2023a); Guo & Qiao (2023); Zhou & Dette (2023). However, these approaches often assume the existence of linear and Markovian dynamics over time, which may fail to accurately capture the complex nonlinear or non-Markovian temporal dependence that often arises in real-world scenarios.

On the other hand, while deep learning has achieved impressive results in computer vision and natural language processing (NLP) (Guo et al., 2016; He et al., 2016; Vaswani et al., 2017; Torfi et al., 2020), applying deep neural networks directly to handle high-dimensional functional time series is challenging. One major issue when dealing with time series data is that deep learning is a general black-box method that lacks explainability, thus making it difficult to understand the cross-sectionally and serially correlated relationships. However, explainability is crucial in many applications. For instance, in finance, healthcare, and climate change, the accuracy and reliability of a model’s predictions have significant impacts on business decisions, patient outcomes, or environmental safety, respectively. Additionally, the non-stationarity of data and the large number of parameters in deep neural networks pose extra challenges during training.

In this paper, we propose an explainable approach called the deep functional factor model (DF<sup>2</sup>M), which has the ability to discover nonlinear and non-Markovian dynamics in high-dimensional functional time series. Developed as a Bayesian nonparametric model, DF<sup>2</sup>M employs a functional version of a factor model for dimension reduction, incorporates an Indian buffet process prior in the infinite-dimensional loading matrix to encourage column sparsity (Guo et al.,

2021), utilizes a functional version of a Gaussian process dynamical model to capture temporal dependence within latent functional factors, and employs deep neural networks to construct the temporal kernel.

DF<sup>2</sup>M offers several advantages for the analysis of high-dimensional functional time series. (i) Firstly, it facilitates a more intuitive understanding of the underlying data structure by representing observed curve variables using a smaller set of latent functional factors. This enhances model explainability and provides a clear and interpretable mapping of relationships between variables, which is crucial for decision-making and subsequent analysis. (ii) Secondly, DF<sup>2</sup>M is capable of discovering non-Markovian and non-linear temporal dependence in the functional latent factor space. This enables more accurate predictions of future values. (iii) Lastly, DF<sup>2</sup>M offers a flexible framework that combines modern sequential deep neural networks with a backbone Bayesian model. This allows for the utilization of sequential deep learning techniques such as gated recurrent unit (GRU) (Cho et al., 2014), long short-term memory (LSTM) (Hochreiter & Schmidhuber, 1997), and attention mechanisms (Vaswani et al., 2017).

## 2. Preliminaries

### 2.1. Indian Buffet Process

The Indian buffet process (IBP) (Griffiths & Ghahramani, 2011) is a probability distribution over a sparse binary matrix with a finite number of rows and an infinite number of columns. The matrix  $\mathbf{Z}$ , generated from the IBP with parameter  $\alpha$ , is denoted as  $\mathbf{Z} \sim \text{IBP}(\alpha)$ , where  $\alpha$  controls the column sparsity of  $\mathbf{Z}$ .

IBP can be explained using a metaphor that customers sequentially visit a buffet and choose dishes. The first customer samples a number of dishes based on  $\text{Poisson}(\alpha)$ . Subsequent the  $i$ -th consumer, in turn, samples each previously selected dish with a probability proportional to its popularity ( $m_k/i$  for dish  $k$ ), and also tries new dishes following  $\text{Poisson}(\alpha/i)$ .

It is worth noting that the distribution remains exchangeable with respect to the customers, meaning that the distribution is invariant to the permutations of the customers. The Indian buffet process admits a stick-breaking representation as  $v_j \mid \alpha \sim \text{Beta}(\alpha, 1)$  independently for  $j = 1, 2, \dots$ ,  $w_k = \prod_{j=1}^k v_j$  for  $k = 1, 2, \dots$ , and  $Z_{ik} \mid w_k \sim \text{Bernoulli}(w_k)$  independently for  $i = 1, \dots, n$ , and the IBP is then defined as  $\mathbf{Z} = (Z_{ik})_{1 \leq i \leq n, k \geq 1}$ . The stick-breaking representation is frequently used in the inference for IBP.

### 2.2. Gaussian Process

A Gaussian process  $X(\cdot)$ , defined on a compact interval  $\mathcal{U}$ , is a continuous stochastic process characterized by the fact that every finite collection of its values,  $X(u_1), \dots, X(u_L)$  with  $u_1, \dots, u_L \in \mathcal{U}$ , belongs to an  $L$ -dimensional multivariate Gaussian distribution (Williams & Rasmussen, 2006). This means that a Gaussian process is completely determined by its mean function  $m(u) = \mathbb{E}[X(u)]$  and its covariance function  $\kappa(u, v) = \text{Cov}(X(u), X(v)) = \mathbb{E}[(X(u) - m(u))(X(v) - m(v))]$  for any  $u, v \in \mathcal{U}$ .

The covariance function, also known as the kernel function in machine learning literature, specifies the correlation between values at distinct points. Examples include the squared exponential kernel  $\kappa(u, v) = \exp(-|u - v|^2/\ell^2)$  and the Ornstein–Uhlenbeck kernel  $\kappa(u, v) = \exp(-|u - v|/\ell)$ , where  $\ell$  is the length-scale parameter. Additionally, the kernel function can be made more complex using the kernel trick (Hofmann et al., 2008) by rewriting it as  $\kappa(u, v) = \langle \phi(u), \phi(v) \rangle$ , where  $\langle \cdot, \cdot \rangle$  denotes the inner product, and the feature function  $\phi(\cdot)$  maps  $x$  into a feature space. As  $\phi(\cdot)$  can be an arbitrary function (linear or nonlinear), the Gaussian process offers considerable flexibility in modeling complex patterns in the data.

Furthermore, a multi-task Gaussian process (MTGP) (Bonilla et al., 2007) can be employed to model vector-valued random fields. It is defined as  $\mathbf{X}(\cdot) = (X_1(\cdot), \dots, X_M(\cdot))^T$ , where  $X_1(\cdot), \dots, X_M(\cdot)$  are  $M$  Gaussian processes defined on  $\mathcal{U}$ . The covariance function between the  $l$ -th and  $k$ -th task is given by  $\text{Cov}(X_l(u), X_k(v)) = \Sigma_{lk}\kappa(u, v)$ , where  $\Sigma = \{\Sigma_{lk}\}_{1 \leq l, k \leq M}$  is a positive semi-definite matrix encoding the similarities between pairs of tasks. The MTGP can effectively capture inter-task correlations and improve predictions (Moreno-Muñoz et al., 2018).

### 2.3. Sequential Deep Learning

Deep learning methods, widely used in computer vision, NLP, and reinforcement learning, have become increasingly popular for time series prediction as well (Lim & Zohren, 2021). In particular, recurrent neural networks (RNN) and attention mechanisms, commonly used for sequence prediction tasks in NLP, can be adapted for temporal forecasting tasks in time series data. A multivariate time series can be modeled recursively in RNN as  $\mathbf{x}_t = g_{\text{dec}}(\mathbf{h}_t)$  and  $\mathbf{h}_t = g_{\text{enc}}(\mathbf{h}_{t-1}, \mathbf{x}_{t-1})$ , where  $\mathbf{h}_t$  is a latent variable and  $g_{\text{dec}}$  and  $g_{\text{enc}}$  are the decoder and encoder functions, respectively. Two renowned RNN models, LSTM and GRU, are designed to learn long-range dependencies in a sequence. For simplicity, we denote their encoder functions as  $\mathbf{h}_t = \text{LSTM}(\mathbf{x}_{1:t})$  and  $\mathbf{h}_t = \text{GRU}(\mathbf{x}_{1:t})$ , respectively, where  $\mathbf{x}_{1:t} = (\mathbf{x}_1, \dots, \mathbf{x}_t)$ . Moreover, attention mecha-

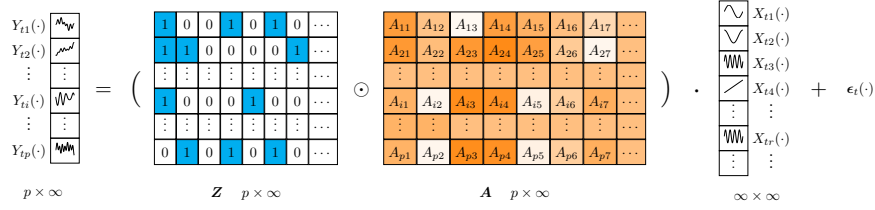


Figure 1. Sparse functional factor model. The blue (or white) cells in  $\mathbf{Z}$  indicate 1 (or 0), while the darker (or lighter) shades of orange in  $\mathbf{A}$  represent larger (or smaller) values.

nisms, which have achieved state-of-the-art performance in NLP tasks, can also be utilized to model time series data. Unlike RNNs, attention mechanisms directly aggregate information from multiple time steps in the past. Attention mechanisms can be expressed as  $\mathbf{h}_t = \sum_{i=1}^{t-1} \omega(\mathbf{k}_t, \mathbf{q}_t) \mathbf{v}_{t-\tau}$ , where key  $\mathbf{k}_t$ , query  $\mathbf{q}_t$ , and value  $\mathbf{v}_t$  are intermediate representations generated by linear or nonlinear transformations of  $\mathbf{x}_t$ . We denote such attention mechanisms as  $\mathbf{h}_t = \text{ATTN}(\mathbf{x}_{1:t})$ . See detailed structures for both RNNs and attention mechanisms in Appendix A.

### 3. Deep Functional Factor model

#### 3.1. Sparse Functional Factor Model

First, we propose a functional factor model from the Bayesian perspective,

$$\mathbf{Y}_t(\cdot) = (\mathbf{Z} \odot \mathbf{A}) \mathbf{X}_t(\cdot) + \epsilon_t(\cdot), \quad t = 1, \dots, n. \quad (1)$$

The observed functional time series is denoted as  $\mathbf{Y}_t(\cdot)$ , and the binary matrix  $\mathbf{Z}$  is sampled from the Indian buffet process,  $\mathbf{Z} \sim \text{IBP}(\alpha)$ . The Hadamard (elementwise) product is represented by  $\odot$ . The loading weight matrix  $\mathbf{A}$  has elements  $A_{tr} \sim \text{Normal}(0, \sigma_A^2)$  for any  $r \in \mathbb{N}^+$  independently. The latent functional factor time series are denoted as  $\mathbf{X}_t(\cdot) = (X_{t1}(\cdot), X_{t2}(\cdot), \dots, X_{tr}(\cdot), \dots)^T$ , and the idiosyncratic component is denoted as  $\epsilon_t(\cdot)$ , which follows a Gaussian distributed white noise process on a scale  $\sigma_\epsilon$ .

In this framework, we do not specify the number of latent factors. Instead,  $\mathbf{Y}_t(\cdot)$ ,  $\mathbf{Z}$ , and  $\mathbf{A}$  can be regarded as  $p \times \infty$  matrices, and  $\mathbf{X}_t(\cdot)$  as an infinite-dimensional vector of functions, or heuristically, a  $\infty \times \infty$  matrix. The dimension reduction framework in equation (1) is illustrated in Figure 1. This Bayesian nonparametric factor model allows for a potentially unlimited number of latent factors, eliminating the need to specify a fixed dimensionality of the factor space. The nonparametric approach introduces flexibility and provides a foundation for inferring the number of factors in the posterior distribution using nonparametric inference frameworks such as Gibbs sampling (Teh et al., 2006), merge-split algorithm (Hughes & Sudderth, 2013), and conditional and adaptively truncated variational inference (Liu et al., 2022, 2023).

Additionally, the Indian buffet process can provide column sparsity (Vu & Lei, 2013) to  $\mathbf{Z}$  and, therefore, to the loading matrix  $\mathbf{Z} \odot \mathbf{A}$ . This implies that most elements in each row are zeros, as  $w_k$  in the stick-breaking representation of the IBP approaches zero as  $k$  increases. In the factor model, this column sparsity implies that each factor affects only a small fraction of functional variables (Guo et al., 2021), or equivalently, the factors are related to each other through a hierarchy (Griffiths & Ghahramani, 2011).

#### 3.2. Functional Gaussian Process Dynamical Model

Secondly, by projecting high-dimensional observations  $\mathbf{Y}_t(\cdot)$  onto low-dimensional latent functional factors  $\mathbf{X}_t(\cdot)$ , we can capture the sequential structure of the time series model through the factors. To model the temporal dependence of  $\mathbf{X}_t(\cdot)$ , we adopt a Gaussian process over time to encode historical information. In particular, we design the covariance across factors  $r$  and  $l$  as follows. Let  $\mathcal{X}_t$  represent the historical information up to time  $t$ , and let  $\mathcal{X}$  be the space containing  $\{\mathcal{X}_t\}_{t \in \mathbb{N}}$ . For any  $u, v \in \mathcal{U}$ ,

$$\text{Cov}(X_{tr}(u), X_{sl}(v)) = \kappa_{\mathcal{X}}(\mathcal{X}_{t-1}, \mathcal{X}_{s-1}) \kappa_{\mathcal{U}}(u, v) \mathbb{I}(r = l), \quad (2)$$

where  $t$  and  $s$  indicate two time stamps, and  $\kappa_{\mathcal{X}}$  and  $\kappa_{\mathcal{U}}$  are the kernels defined on  $\mathcal{X}$  and  $\mathcal{U}$ , respectively. The indicator function  $\mathbb{I}(r = l)$  equals 1 if  $r = l$  and 0 otherwise. The kernel  $\kappa_{\mathcal{X}}$  captures historical information from different periods and can be regarded as a temporal kernel. Similarly, the kernel  $\kappa_{\mathcal{U}}$  can be seen as a spatial kernel.

Therefore,  $\mathbf{X}_r(\cdot) = (X_{1r}(\cdot), \dots, X_{tr}(\cdot), \dots, X_{nr}(\cdot))$  belongs to a multi-task Gaussian process (Bonilla et al., 2007) such that for any  $u_1, \dots, u_L \in \mathcal{U}$ ,  $\text{vec}(\mathbf{X}_r(u_1, \dots, u_L)) \sim \text{Normal}(\mathbf{0}, \Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^u)$ , where  $\otimes$  denotes the Kronecker product, or equivalently  $\mathbf{X}_r(u_1, \dots, u_L)$  follows a matrix normal distribution (Dawid, 1981) with mean  $\mathbf{0}$ , row covariance matrix  $\Sigma_{\mathcal{U}}^u$ , and column covariance matrix  $\Sigma_{\mathcal{X}}$ . Here,  $\mathbf{X}_r(u_1, \dots, u_L) = [X_{tr}(u_j)]_{1 \leq t \leq n, 1 \leq j \leq L}$ ,  $\Sigma_{\mathcal{X}} = [\kappa_{\mathcal{X}}(\mathcal{X}_t, \mathcal{X}_s)]_{0 \leq t, s \leq n-1}$ , and  $\Sigma_{\mathcal{U}}^u = [\kappa_{\mathcal{U}}(u_i, u_j)]_{1 \leq i, j \leq L}$ .

In Appendix B, we demonstrate the detailed relationship between the multi-task Gaussian process and the matrix normal distribution. In the literature, the  $n$ -task Gaussian

process is often used to refer to multiple outputs generated by the model, each corresponding to a specific timestamp in our time series setting. For ease of expression, we denote the presented multi-task Gaussian process as

$$\mathbf{X}_r(\cdot) \sim \text{MTGP}(\mathbf{0}, \kappa_{\mathcal{U}}(\cdot, \cdot), \kappa_{\mathcal{X}}(\cdot, \cdot)). \quad (3)$$

Importantly, the marginal distributions of latent factors, denoted as  $\mathbf{X}_1, \dots, \mathbf{X}_r, \dots$ , exhibit cross-sectional dependence due to their common temporal kernel  $\kappa_{\mathcal{X}}$ , which incorporates historical information up to period  $t-1$ . This shared kernel promotes similarity and dependence across different time periods in the multi-task Gaussian process. However, when conditioned on  $\mathcal{X}_{t-1}$  and the corresponding kernel,  $\mathbf{X}_1, \dots, \mathbf{X}_r, \dots$ , become conditionally independent and Gaussian distributed. This conditionality arises from our approach, where the predictive temporal kernel exclusively relies on past information  $\mathcal{X}_{t-1}$  rather than current data  $\mathcal{X}_t$ . Consequently, this approach enables forward-looking predictions based solely on historical data.

The proposed model can be seen as a functional variant of the Gaussian process dynamical model (Wang et al., 2005). The connections between the two models can be found in Appendix A. Using the proposed model to capture the temporal dependence of functional time series has several advantages. First, the temporal kernel and spatial kernel can be separated, which allows for a closed form and computational convenience. Additionally, since the temporal kernel considers the entire historical information rather than just the latest state, the model can be non-Markovian. For example, by defining the temporal kernel as  $\kappa(\mathcal{X}_{t-1}, \mathcal{X}_{s-1}) = \alpha_1 \int \mathbf{X}_{t-1}(u)^T \mathbf{X}_{s-1}(u) du + \alpha_2 \int \mathbf{X}_{t-2}(u)^T \mathbf{X}_{s-2}(u) du$ , features from the last two periods can be incorporated. Furthermore, nonlinearity can be introduced using the kernel trick by setting a nonlinear kernel function. This opens the possibility of constructing deep kernels, which will be discussed in Section 3.3.

### 3.3. Deep Temporal Kernels

To capture the complex latent temporal structure, we use neural networks to construct the kernel function. However, compared to standard deep kernels (Wilson et al., 2016; Al-Shedivat et al., 2017; Xue et al., 2019; Li et al., 2019; Watson et al., 2021; Fortuin, 2022), two extra steps are needed when applying deep kernels to functional time series.

Firstly, since  $\mathbf{X}_t(\cdot)$  is a continuous process on  $\mathcal{U}$ , a mapping function  $F: \mathcal{F} \rightarrow \mathbb{R}^d$  is required to map the infinite-dimensional Gaussian processes to  $d$ -dimensional vectors. Here,  $\mathcal{F}$  represents the space of continuous functions defined on  $\mathcal{U}$ . Various approaches can be used for this mapping function, including pre-specified basis expansion, data-dependent basis expansion (such as functional principal component analysis and its dynamic variants (Bathia et al.,

2010; Hormann et al., 2015)), adaptive functional neural network (Yao et al., 2021), or even a simple specification such as  $\mathbf{X}_t(u_0, \dots, u_L)$  with  $u_0, \dots, u_L \in \mathcal{U}$ .

Secondly, the  $d$ -dimensional vectors are used as inputs for deep neural networks, and the outputs generated by these networks are employed to construct kernel functions. Specifically, the input vector is transformed as

$$\mathbf{h}_t = H(F(\mathbf{X}_{t-1}), F(\mathbf{X}_{t-2}), \dots), \quad (4)$$

where  $\mathbf{X}_{t-1} = (X_{t-1,1}, \dots, X_{t-1,r}, \dots)^T$ ,  $F$  is the mapping function, and  $H$  represents a sequential deep learning framework. Various deep neural network architectures can be utilized for this purpose, such as LSTM, GRU, and attention mechanisms, which have shown their effectiveness in modeling complex patterns and dependencies. Since the inputs for the temporal kernels are ordered sequences from  $\mathcal{X}_0$  to  $\mathcal{X}_{n-1}$ , unidirectional deep neural networks should be used instead of bidirectional networks. The transformed representations  $\mathbf{h}_t$  and  $\mathbf{h}_s$  are then used to construct a kernel

$$\kappa_{\mathcal{X}}(\mathcal{X}_{t-1}, \mathcal{X}_{s-1}) = \kappa(\mathbf{h}_t, \mathbf{h}_s), \quad (5)$$

where  $\kappa(\cdot, \cdot)$  is a suitable kernel function, such as the squared exponential kernel or the Ornstein-Uhlenbeck kernel. It should be noted that the temporal kernel is related to the historical values of all the relevant factors and is shared across factors. This approach builds a time-dependent variant of a deep Gaussian process (Damianou & Lawrence, 2013), where the kernel incorporates historical information.

To address overfitting, spectrum normalization can be applied, which effectively enforces a Lipschitz condition on the neural networks as suggested by Miyato et al. (2018). Notably, this Lipschitz condition on the model’s output concerning its input ensures that similar inputs within the Gaussian process exhibit comparable distances in the kernel space. This implies that the inputs for  $\kappa$  can reflect the distance between  $\mathcal{X}_{t-1}$  and  $\mathcal{X}_{s-1}$ .

### 3.4. The Imperative of Element Integration

In summary, by combining the functional version of the sparse factor model, sequential deep learning kernel, and Gaussian process dynamical model, we define a probabilistic generative model for high-dimensional functional time series named the deep functional factor model (DF<sup>2</sup>M).

Factorization is essential in this context. Directly feeding the original high-dimensional input into the kernel function is not viable. The issue arises because without the factorization of high-dimensional functional data, the kernel would have to handle extremely high-dimensional inputs using an excessive number of parameters, while in training time steps are limited. This situation is prone to overfitting, making it challenging to achieve accurate estimation. Additionally,



the interpretability of our proposed method relies on these factors.

Applying a deep kernel without utilizing an IBP with a finite number of latent factors poses difficulties. The challenge lies in determining the optimal number of latent factors. If there are too many factors or too few, it can affect the kernel distance. An excess of factors, often redundant and similar, can distort distance measurements, while an inadequate number of factors may miss critical distances, leading to inaccurate model representations.

The sequential deep learning kernel is essential for introducing non-Markovian and non-linear patterns. When combined with IBP and factorization, it effectively models temporal similarities among observations, enabling capture of high-dimensional functional data dynamics.

## 4. Bayesian Inference for DF<sup>2</sup>M

### 4.1. Sparse Variational Inference

We adopt the variational inference framework to infer the proposed DF<sup>2</sup>M. This algorithm approximates the posterior probability by maximizing the evidence lower bound (ELBO), which is equivalent to minimizing the Kullback-Leibler (KL) divergence between a variational distribution and true posterior distribution (Blei et al., 2017). For DF<sup>2</sup>M, with mean-field factorization assuming independence among the variational distributions for latent variables, its ELBO can be expressed as

$$\begin{aligned} \text{ELBO} = & \mathbb{E}_q \left[ \log p(\mathbf{Z} | \alpha) p(\mathbf{A} | \sigma_A) \right. \\ & \left. \prod_{t=1}^n p(\mathbf{Y}_t(\cdot) | \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A}) \prod_{r \geq 1} p(\mathbf{X}_r(\cdot) | \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}) \right] \quad (6) \\ & - \mathbb{E}_q \left[ \log q(\mathbf{Z}) q(\mathbf{A}) \prod_{r \geq 1} q(\mathbf{X}_r(\cdot)) \right]. \end{aligned}$$

Using the stick-breaking representation of the Indian buffet process as in Section 2.1, we factorize the variational distribution for  $\mathbf{Z}$  as  $q(v_j) = \text{Beta}(v_j; \tau_j^1, \tau_j^0)$  and  $q(Z_{tj}) = \text{Bernoulli}(Z_{tj}; m_{tj})$ . The corresponding variational distribution for  $\mathbf{A}$  is factorized as  $q(A_{tj}) = \text{Normal}(A_{tj}; \eta_{tj}, \sigma_{A,tj}^2)$ .

To avoid singular matrix inversions and improve computational efficiency, we propose a sparse variational inference approach for DF<sup>2</sup>M based on Titsias (2009). Our method introduces a set of inducing variables representing the values of the latent function at a small subset of points in  $\mathcal{U}$ . Moreover, we adopt the approach of having common locations for the inducing variables across functional factors, as suggested by Hamelijnck et al. (2021). In other words, we utilize the same set of inducing points for all tasks, which can lead to further improvement in computa-

tional efficiency. Consequently, the variational distribution for multi-task Gaussian process with inducing variables is defined as,

$$\begin{aligned} q(\mathbf{X}_r(\cdot)) = & p\left(X_{1r}(\cdot), \dots, X_{nr}(\cdot) | X_{1r}(\mathbf{v}), \dots, \right. \\ & \left. X_{nr}(\mathbf{v}), \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}\right) \prod_{t=1}^n q(X_{tr}(\mathbf{v})), \quad (7) \end{aligned}$$

where  $\mathbf{v} = (v_1, \dots, v_K)^T$  with  $v_1, \dots, v_K \in \mathcal{U}$  with  $K$  being the number of inducing points. The variational distribution for the inducing variables is constructed as  $q(X_{tr}(\mathbf{v})) = \text{Normal}(\boldsymbol{\mu}_{tr}, \mathbf{S}_{tr})$ . It is important to note that the conditional prior distribution for  $X_r(\cdot)$ , which is the first term on the right-hand side of equation (7), cannot be factorized as  $\prod_{t=1}^n p(X_{tr}(\cdot) | X_{tr}(\mathbf{v}))$  due to their temporal dependence. However, by exploiting the setting of equation (7), the conditional prior distribution appears in both the variational and prior distributions and therefore can be cancelled. In Appendix D.1, we derive that the ELBO in equation (6) can be simplified as

$$\begin{aligned} \text{ELBO} = & \sum_{t=1}^n \mathbb{E}_q \left[ \log p(\mathbf{Y}_t(\cdot) | \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A}) \right] \\ & - \text{KL}[q(\mathbf{Z}) \| p(\mathbf{Z} | \alpha)] - \text{KL}[q(\mathbf{A}) \| p(\mathbf{A} | \sigma_A)] \quad (8) \\ & - \sum_{r \geq 1} \text{KL}[q(\mathbf{X}_r(\mathbf{v})) \| p(\mathbf{X}_r(\mathbf{v}) | \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}})], \end{aligned}$$

where  $\mathbf{X}_r(\mathbf{v}) = (X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v}))$  with  $X_{tr}(\mathbf{v}) = (X_{tr}(v_1), \dots, X_{tr}(v_K))^T$  for  $t = 1, \dots, n$ . Furthermore, using the formula of the KL divergence between two multivariate Gaussian distributions, we derive a closed form of the last term as

$$\begin{aligned} & 2\text{KL}[q(\mathbf{X}_r(\mathbf{v})) \| p(\mathbf{X}_r(\mathbf{v}) | \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}})] \\ & = \text{trace} \left( \left( \boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv^{-1}} \right) (\mathbf{S}_r + \text{vec}(\boldsymbol{\mu}_r) \text{vec}(\boldsymbol{\mu}_r)^T) \right) \\ & \quad + K \log |\boldsymbol{\Sigma}_{\mathcal{X}}| + n \log |\boldsymbol{\Sigma}_{\mathcal{U}}^{vv}| - \sum_{t=1}^n \log |\mathbf{S}_{tr}| - nK, \end{aligned}$$

where  $\boldsymbol{\mu}_r = (\boldsymbol{\mu}_{1r}, \dots, \boldsymbol{\mu}_{nr})$ ,  $\mathbf{S}_r = \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr})$ , and  $\boldsymbol{\Sigma}_{\mathcal{U}}^{vv} = [\kappa_{\mathcal{U}}(v_j, v_j)]_{1 \leq i, j \leq K}$ . See Appendix D.2 for the detailed derivation.

### 4.2. Sampling for Variational Distribution of Factors

To optimize the variational distributions, the automatic differentiation variational inference (ADVI) algorithm (Kucukelbir et al., 2017; Blei et al., 2017; Ranganath et al., 2014) is adopted to maximize the ELBO in equation (8).

To perform ADVI in our model, we need to sample  $\mathbf{X}_r(\cdot)$  from its variational distribution as specified in equation (7).

However, directly sampling from a  $nL \times nL$  matrix is computationally expensive even though this distribution is Gaussian conditional on  $\mathbf{X}_r(\mathbf{v})$ . To address this issue and accelerate the computation of the ELBO, we take advantage of the separability of the temporal and spatial kernels as described in Section 3.2, and propose the following method.

For any  $\mathbf{u} = (u_1, \dots, u_L)^T$  with  $u_1, \dots, u_L \in \mathcal{U}$  being the observation points in  $\mathcal{U}$ , we first partition the spatial covariance matrix for  $\mathbf{X}(\mathbf{u}, \mathbf{v})$  into a blockwise matrix shown as  $\begin{bmatrix} \Sigma_{\mathcal{U}}^{uu} & \Sigma_{\mathcal{U}}^{uv} \\ \Sigma_{\mathcal{U}}^{uvT} & \Sigma_{\mathcal{U}}^{vv} \end{bmatrix}$ , where  $\Sigma_{\mathcal{U}}^{uu} = [\kappa_{\mathcal{U}}(u_i, u_j)]_{1 \leq i, j \leq L}$ , and  $\Sigma_{\mathcal{U}}^{uv} = [\kappa_{\mathcal{U}}(u_i, v_j)]_{1 \leq i \leq L, 1 \leq j \leq K}$ .

**Theorem 1 (Posterior Mean)** *The mean function of the posterior for  $X_{tr}(\cdot)$  is solely dependent on the variational mean of  $\mathbf{X}_{tr}(\mathbf{v})$ , the inducing variables at time  $t$ . That is, for any  $\mathbf{u}$*

$$E(X_{tr}(\mathbf{u})) = \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \boldsymbol{\mu}_{tr}. \quad (9)$$

It means that for MTGP, the variational mean is independent of the inducing variables at timestamps other than the current one. See also an analogous theorem for Gaussian process regression in Bonilla et al. (2007).

**Theorem 2 (Posterior Variance)** *The variance function of the posterior for  $\mathbf{X}_r(\cdot)$  contains two parts. For any  $\mathbf{u}$ ,*

$$\begin{aligned} \text{Var}_q[\text{vec}(\mathbf{X}_r(\mathbf{u}))] &= (I \otimes \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1}) \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr}) \\ &+ \Sigma_{\mathcal{X}} \otimes (\Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}). \end{aligned} \quad (10)$$

The first part is solely dependent on the variational variance of  $\mathbf{X}_{tr}(\mathbf{v})$ , while the second part is independent of the variational distributions of all inducing variables. In particular, the first part corresponds to a group of independent Gaussian processes such that  $\tilde{\mathbf{X}}_{tr}^{(1)}(\mathbf{u}) \sim \text{Normal}(\Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1} \boldsymbol{\mu}_{tr}, \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1} \mathbf{S}_{tr})$  for any  $\mathbf{u}$ . On the other hand,  $\tilde{\mathbf{X}}_r^{(2)}(\cdot)$  is a zero-mean multi-task Gaussian process, with  $\tilde{\mathbf{X}}_r^{(2)}(\mathbf{u}) \sim \text{MatrixNormal}(\mathbf{0}, \Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}, \Sigma_{\mathcal{X}})$  for any  $\mathbf{u}$ . Therefore, based on Theorems 1 and 2, we can decompose  $\mathbf{X}_r(\cdot) = \tilde{\mathbf{X}}_r^{(1)}(\cdot) + \tilde{\mathbf{X}}_r^{(2)}(\cdot)$  under the variational distribution. Notably, the sampling of  $\tilde{\mathbf{X}}_r^{(1)}(\cdot)$  is more efficient as it only depends on inducing variables within the same period.

**Theorem 3 (Irrelevance to ELBO)** *Conditional on  $\Sigma_{\mathcal{X}}$  and  $\Sigma_{\mathcal{U}}$ , sampling  $\mathbf{X}_{tr}(\cdot)$  from the distribution of  $\tilde{\mathbf{X}}_r^{(1)}(\cdot)$  does not change the variational mean. Moreover, the corresponding ELBO of DF<sup>2</sup>M in equation (8) is only modified*

by a constant term given by

$$\frac{1}{2\sigma_\epsilon^2} \|\mathbf{Z} \odot \mathbf{A}\|_F^2 \text{trace}[\Sigma_{\mathcal{X}}] \text{trace}[\Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}], \quad (11)$$

where  $\|\mathbf{M}\|_F = (\sum_{i,j} M_{ij}^2)^{\frac{1}{2}}$  denotes the Frobenius norm of any matrix  $\mathbf{M}$ .

See Appendices D.3, D.4 and D.5 for the derivations of Theorem 1, 2 and 3, respectively. With the help of these theorems, we can sample  $\mathbf{X}_{tr}(\cdot)$  from the proxy variational distribution  $\text{Normal}(\Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1} \boldsymbol{\mu}_{tr}, \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1} \mathbf{S}_{tr})$ , which relies solely on the variational distributions at time  $t$ . This approach provides a more efficient way of computing the ELBO compared to direct sampling, which requires the complete Cholesky decomposition of the  $nL \times nL$  matrix.

### 4.3. Initialization, Training and Prediction

We use the technique of ADVI to train the variational parameters of the posteriors, by computing the gradient of the ELBO with respect to the parameters. The training process requires iterating through the following steps until the ELBO converges. The steps of Bayesian inference for DF<sup>2</sup>M are summarized in Algorithm 1 in Appendix E.

First, conditional on  $\Sigma_{\mathcal{X}}$ , we update the variational distribution parameters  $\boldsymbol{\mu}_{tr}$  and  $\mathbf{S}_{tr}$  for inducing variables  $\mathbf{X}_{tr}(\mathbf{v})$  for all  $t$  and  $r$ , as well as other variational parameters including  $\{\tau_j^1, \tau_j^2\}_{1 \leq j \leq M}$  and  $\{m_{tj}\}_{1 \leq t \leq n, 1 \leq j \leq M}$  for India buffet process  $\mathbf{Z}$ ,  $\{\eta_{tj}, \sigma_{tj}^A\}_{1 \leq t \leq n, 1 \leq j \leq M}$  for loading weight matrix  $\mathbf{A}$ . We also update the idiosyncratic noise scale  $\sigma_\epsilon$  and the parameters in the spatial kernel  $\kappa_{\mathcal{U}}(\cdot, \cdot)$ . In this step, the gradient of ELBO is accelerated by sampling  $X_{tr}(\cdot)$  independently according to Theorem 3 and the analytical expression for the KL divergence in equation (9).

Second, conditional on a sample of  $\mathbf{X}_r(\cdot)$ , we update the trainable parameters in sequential deep learning framework  $H$  that constructs the temporal kernels  $\kappa_{\mathcal{X}}(\cdot, \cdot)$ , via the gradient of ELBO with respect to  $\Sigma_{\mathcal{X}}$ . Although any mapping function  $F$  can be used in our model, it is natural to choose  $F(\mathbf{X}_t(\cdot)) = \mathbf{X}_t(\mathbf{v})$ , which eliminates the need to sample  $\mathbf{X}_r(\cdot)$  when computing the gradient. This is inspired by the fact that the variational distribution of inducing variables can be regarded as sufficient statistics of the Gaussian processes (Titsias, 2009).

Once we have observed the data at time  $n$ , we use the trained model to generate a posterior distribution that captures our updated understanding of the underlying patterns in the data. Based on this distribution, we make a prediction for the value of the data at the next time step,  $n+1$ . We present the one-step ahead prediction as:

$$\begin{aligned} \bar{\mathbf{Y}}_{n+1}(\mathbf{u}) &= (\bar{\mathbf{Z}} \odot \bar{\mathbf{A}}) \bar{\mathbf{X}}_{n+1}(\mathbf{u}), \\ \bar{\mathbf{X}}_{n+1,r}(\mathbf{u}) &= \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \boldsymbol{\mu}_r \Sigma_{\mathcal{X}}^{-1} \Sigma_{\mathcal{X}}^{n+1,1:nT}, \end{aligned} \quad (12)$$

where  $\bar{Y}$  and  $\bar{X}$  represent the predictive means for the observations and factors, respectively. The terms  $\bar{Z}$  and  $\bar{A}$  are the posterior means of  $Z$  and  $A$ , respectively. The component  $\Sigma_{\mathcal{X}}^{n+1,1:n}$  is a  $1 \times n$  matrix given by  $\Sigma_{\mathcal{X}}^{n+1,1:n} = [\kappa_{\mathcal{X}}(\mathcal{X}_{n+1}, \mathcal{X}_0), \dots, \kappa_{\mathcal{X}}(\mathcal{X}_{n+1}, \mathcal{X}_n)]$ . See Appendix D.6 for the derivations. By repeating this process iteratively, we can generate a sequence of predictions for future time steps, thereby forecasting the behavior of the system over time.

## 5. Experiments

### 5.1. Datasets

We apply DF<sup>2</sup>M to four real-world datasets consisting of high-dimensional functional time series. **Japanese Mortality** dataset contains age-specific mortality rates for 47 Japanese prefectures ( $p=47$ ) from 1975 to 2017, with 43 observations per prefecture ( $n=43$ ). **Energy Consumption** dataset includes half-hourly measured energy consumption curves for selected London households ( $p=40$ ) between December 2012 and January 2013 ( $n=55$ ). **Global Mortality** dataset provides a broader perspective on mortality rates by including age-specific mortality data across different countries ( $p=32$ ) from 1960 to 2010 ( $n=50$ ). **Stock Intraday** dataset comprises high-frequency price observations for the S&P 100 component stocks (we removed 2 stocks with missing values, so  $p=98$ ) in 2017. The data includes 45 trading days ( $n=45$ ), with ten-minute resolution prices and cumulative intraday return trajectories (Horváth et al., 2014). Each dataset is preprocessed and transformed into an appropriate format for analysis. See Appendix F for the details. We denote the data as  $\{Y_{tj}(u_k)\}_{1 \leq t \leq n, 1 \leq j \leq p, 1 \leq k \leq K}$ , where  $K$  is the number of observations per curve. Examples of functional time series for a randomly selected  $j$  are plotted in Row (1) of Figure 2.

Our work is centered on functional time series, distinct from univariate or multivariate time series. It is important to highlight that functional time series analysis is inherently more complex due to the observations being infinite-dimensional functional objects, making it much more challenging than non-functional time series (Ramsay & Silverman, 2005). The choice of datasets with time steps  $n$  deliberately limited to less than 50 is a standard practice in functional time series literature, suitable for demonstrating the robustness of our method in more challenging scenarios with limited data, in line with Chang et al. (2023a) and Tang et al. (2022).

### 5.2. Experiment Setup and Metrics

To assess the predictive accuracy of the proposed model, we split the data into a training set with the first  $n_1$  periods and a test set with the last  $n_2$  periods. We use the training set to train the parameters in the model following the steps in Section 4. Then for an integer  $h > 0$ , we make the  $h$ -step-

ahead prediction given the fitted model using the first  $n_1$  periods. We then repeat this process by moving the training window by one period, refitting the model, and making the  $h$ -step-ahead prediction. We compute the mean absolute prediction error (MAPE) and mean squared prediction error (MSPE) using the following equations:

$$\text{MAPE}(h) = \frac{1}{M} \sum_{j=1}^p \sum_{k=1}^K \sum_{t=n_1+h}^n |\hat{Y}_{tj}(u_k) - Y_{tj}(u_k)|,$$

$$\text{MSPE}(h) = \frac{1}{M} \sum_{j=1}^p \sum_{k=1}^K \sum_{t=n_1+h}^n [\hat{Y}_{tj}(u_k) - Y_{tj}(u_k)]^2,$$

where  $M = Kp(n_2 - h + 1)$ . In our DF<sup>2</sup>M implementation, we incorporate three cutting-edge deep learning modules: LSTM, GRU, and the self-attention mechanism. In the deep learning modules, we employ a feedforward neural network equipped with ReLU activation functions to map inputs into a designated hidden layer size. Subsequently, the transformed inputs are channeled through a time-invariant full connected neural network, LSTM, GRU, or self-attention mechanisms, denoted as DF<sup>2</sup>M-LIN, DF<sup>2</sup>M-LSTM, DF<sup>2</sup>M-GRU and DF<sup>2</sup>M-ATTN, respectively. For DF<sup>2</sup>M, the outputs of the deep learning modules are passed to the kernel function, while in conventional deep learning, they are converted to outputs via a linear transformation. We evaluate their performance against conventional deep learning models under the same structural setting and regulations. The optimal hyperparameters, along with a detailed description of the deep learning architecture, can be found in Appendix G.

### 5.3. Empirical Results

Our primary objective is to improve the explainability of deep learning models like RNNs or transformers, while maintaining or enhancing prediction accuracy.

**Explainability** Firstly, Row (2) of Figure 2 shows the temporal dynamic of the largest factors in the fitted models. We can observe a decreasing trend over time for the first three datasets. This is particularly valuable as these factors exhibit a clear and smooth dynamic, which can be used to explain the underlying reasons for changes over time and also to make robust predictions. Secondly, the temporal covariance matrix ( $\Sigma_{\mathcal{X}}$ ) can be seen in Row (3) of Figure 2. It is evident that the first three datasets exhibit stronger autocorrelation than the *Stock Intraday* dataset, which aligns with the intuition that financial data is generally noisier and characterized by short-term dependencies.

Furthermore, both mortality datasets display a strong autoregressive pattern, as evidenced by the large covariance values close to the diagonal. They also show a blockwise pattern, which indicates the existence of change points in 1980s. Another interesting observation is the periodic pattern in the

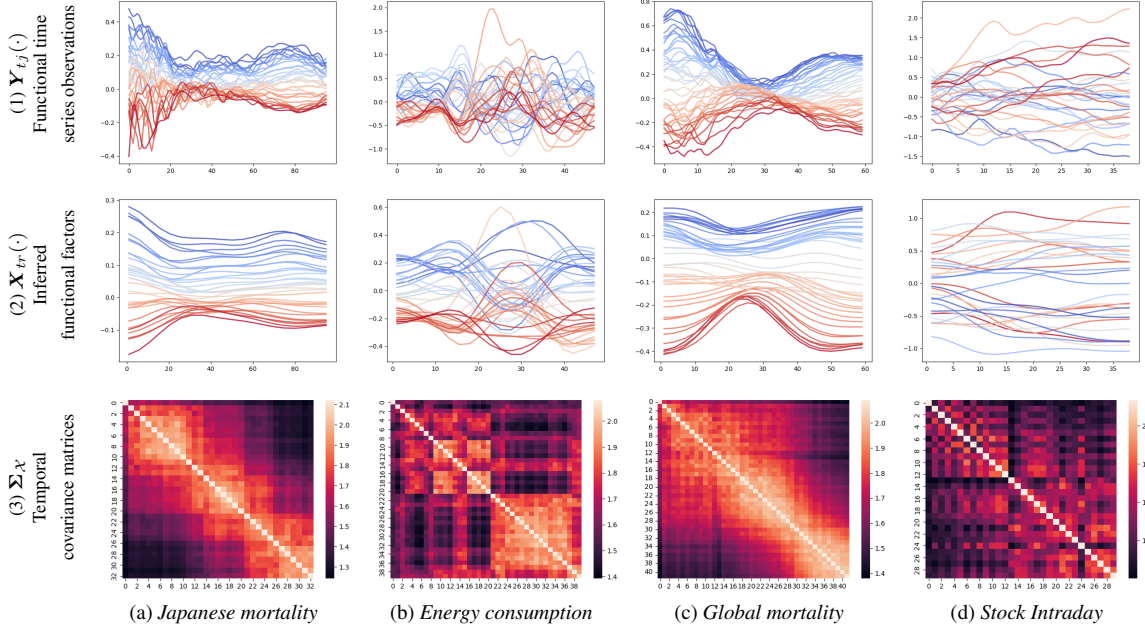


Figure 2. A visualization of real datasets with analysis. Row (1): raw functional time series. Row (2): the largest functional factor. Row (3): temporal covariance matrix. Rows (1) and (2) use a blue-to-red gradient to denote time progression. Blue for older and red for recent data. Row (3) employs brightness variations to represent covariance, with brighter areas indicating higher covariance.

*Energy Consumption* dataset, which reveals distinct patterns for weekdays and weekends during the first 21 days. This data corresponds to the first 21 days in December. In contrast, the second half of the time steps do not exhibit this pattern. This could be attributed to the Christmas holidays in London, during which the differences between weekdays and weekends are smaller, as people are on holiday.

**Predictive Accuracy** Compared to standard deep learning models, the DF<sup>2</sup>M framework consistently outperforms other models in terms of both MSPE and MAPE across all four datasets. The only exception to this is the *Stock Intraday* dataset, where DF<sup>2</sup>M-ATTN and ATTN achieve similar levels of accuracy. Specifically, the DF<sup>2</sup>M-LSTM model performs exceptionally well on the *Energy Consumption* and *Global Mortality* datasets, while the DF<sup>2</sup>M-ATTN model exhibits the lowest prediction error for the *Japanese Mortality* dataset. These results demonstrate that the integration of an explainable structure with the nonlinearity of LSTM and attention mechanisms can significantly improve the overall performance of the model.

On the other hand, the DF<sup>2</sup>M-LIN model outperforms both DF<sup>2</sup>M-LSTM and DF<sup>2</sup>M-GRU on the *Stock Intraday* dataset. This can be attributed to the fact that, in the context of financial data, long-term dependencies may not be present, rendering the Markovian model more suitable for capturing the underlying dynamics. Consequently, the DF<sup>2</sup>M-LIN model emerges as a better choice for the *Stock*

*Intraday* dataset. Compared to standard deep learning models with multiple layers, DF<sup>2</sup>M achieves better or comparable results, as shown in Appendix G. However, in such cases, standard deep learning models sacrifice explainability due to their utilization of a large number of layers.

## 6. Related Works

In the literature concerning frequentist statistical methods for high-dimensional functional time series, various approaches have been employed. For instance, principal components based dimension reduction (Guo & Qiao, 2023; Chang et al., 2023a), factor model (Guo et al., 2021) and segmentation transformation (Chang et al., 2023b). However, all these methods use either vector autoregressive (VAR) or functional VAR to describe the temporal dynamics, implying linear and Markovian models. In contrast, our work is the first to propose a Bayesian model for high-dimensional functional time series that can handle nonlinear and non-Markovian dynamics.

Moreover, several studies (Wilson et al., 2016; Al-Shedivat et al., 2017; Xue et al., 2019; Li et al., 2019; Watson et al., 2021; Fortuin, 2022) have utilized deep kernels in Gaussian processes for classification or regression tasks. In contrast, our framework introduces the use of a deep kernel specifically designed for time series prediction.

Furthermore, prior work (Lawrence, 2003; Wang et al.,



## Deep Functional Factor Models

Table 1. Comparison of DF<sup>2</sup>M to Standard Deep Learning Models. For formatting reasons, MAPEs are multiplied by 10, and MSPEs are multiplied by 10<sup>2</sup>, except for the *Energy Consumption* dataset.

(a) Comparison of DF<sup>2</sup>M-LIN and LIN

		Japanese Mortality			Energy Consumption			Global Mortality			Stock Intraday		
		1	2	3	1	2	3	1	2	3	1	2	3
DF <sup>2</sup> M-LIN	MSPE	<b>4.707</b>	<b>4.567</b>	<b>5.623</b>	<b>10.29</b>	<b>17.58</b>	<b>17.64</b>	<b>10.78</b>	<b>9.300</b>	<b>9.706</b>	<b>99.58</b>	<b>101.2</b>	<b>89.82</b>
	MAPE	<b>1.539</b>	<b>1.446</b>	<b>1.635</b>	<b>2.334</b>	<b>3.060</b>	<b>3.100</b>	<b>2.319</b>	<b>2.041</b>	<b>2.106</b>	<b>6.424</b>	<b>6.505</b>	<b>6.269</b>
LIN	MSPE	7.808	8.774	9.228	16.16	18.95	20.27	16.84	18.05	19.93	137.5	127.8	139.1
	MAPE	2.092	2.227	2.313	2.939	3.214	3.342	2.783	2.949	3.174	7.896	7.491	7.924

(b) Comparison of DF<sup>2</sup>M-LSTM and LSTM

		Japanese Mortality			Energy Consumption			Global Mortality			Stock Intraday		
		1	2	3	1	2	3	1	2	3	1	2	3
DF <sup>2</sup> M-LSTM	MSPE	<b>3.753</b>	<b>4.164</b>	<b>4.513</b>	<b>8.928</b>	<b>11.60</b>	<b>17.26</b>	<b>7.672</b>	<b>8.088</b>	<b>8.954</b>	<b>107.5</b>	<b>118.8</b>	<b>113.6</b>
	MAPE	<b>1.205</b>	<b>1.322</b>	<b>1.427</b>	<b>2.176</b>	<b>2.478</b>	<b>3.063</b>	<b>1.726</b>	<b>1.823</b>	<b>1.978</b>	<b>6.741</b>	<b>7.141</b>	<b>7.294</b>
LSTM	MSPE	4.989	5.597	6.501	13.51	19.71	24.61	13.28	16.29	17.08	193.3	176.0	213.8
	MAPE	1.447	1.523	1.684	2.635	3.278	3.759	2.332	2.572	2.680	9.281	9.283	10.20

(c) Comparison of DF<sup>2</sup>M-GRU and GRU

		Japanese Mortality			Energy Consumption			Global Mortality			Stock Intraday		
		1	2	3	1	2	3	1	2	3	1	2	3
DF <sup>2</sup> M-GRU	MSPE	<b>4.092</b>	<b>4.395</b>	<b>4.898</b>	<b>9.132</b>	<b>8.714</b>	<b>9.730</b>	<b>8.741</b>	<b>8.714</b>	<b>9.730</b>	<b>102.5</b>	<b>117.3</b>	<b>95.49</b>
	MAPE	<b>1.318</b>	<b>1.402</b>	<b>1.537</b>	<b>2.204</b>	<b>1.951</b>	<b>2.110</b>	<b>1.967</b>	<b>1.951</b>	<b>2.110</b>	<b>6.675</b>	<b>7.339</b>	<b>6.649</b>
GRU	MSPE	8.800	8.552	10.41	15.55	24.02	17.53	14.12	15.33	17.53	414.0	445.9	427.2
	MAPE	1.691	1.809	1.865	2.872	3.518	2.597	2.211	2.403	2.597	14.12	14.66	14.07

(d) Comparison of DF<sup>2</sup>M-ATTN and ATTN

		Japanese Mortality			Energy Consumption			Global Mortality			Stock Intraday		
		1	2	3	1	2	3	1	2	3	1	2	3
DF <sup>2</sup> M-ATTN	MSPE	<b>3.608</b>	<b>3.839</b>	<b>3.985</b>	<b>14.22</b>	18.70	19.03	<b>14.22</b>	<b>18.70</b>	<b>19.03</b>	104.2	103.4	93.93
	MAPE	<b>1.119</b>	<b>1.203</b>	<b>1.264</b>	<b>2.741</b>	<b>3.141</b>	<b>3.163</b>	<b>2.741</b>	<b>3.141</b>	<b>3.163</b>	6.695	6.646	6.427
ATTN	MSPE	13.44	14.85	16.17	17.03	<b>17.79</b>	<b>18.24</b>	39.52	41.83	43.95	<b>103.4</b>	<b>98.39</b>	<b>91.21</b>
	MAPE	3.166	3.363	3.546	3.130	3.216	3.268	5.332	5.506	5.643	<b>6.579</b>	<b>6.392</b>	<b>6.257</b>

2005; Titsias & Lawrence, 2010) has employed MTGPs to model cross-sectional correlations among static data. In contrast, we propose the use of a factor model to describe cross-sectional relationships, where the temporal kernel is constructed based on the features of the historical functional factors. This unique structure represents a novel contribution to the current literature.

## 7. Conclusion

In this paper, we present DF<sup>2</sup>M, a novel deep Bayesian nonparametric approach for discovering non-Markovian and nonlinear dynamics in high-dimensional functional time series. DF<sup>2</sup>M combines the strengths of the Indian buffet process, factor model, Gaussian process, and deep neural networks to offer a flexible and powerful framework. Our model effectively captures non-Markovian and nonlinear dynamics while using deep learning in a structured and explainable manner. It bridges modern deep learning and statistical time series. We also propose a computationally

efficient inference algorithm.

Empirical results show the superior predictive performance of DF<sup>2</sup>M compared to corresponding standard deep learning models. However, a potential limitation of our study is its reliance on simple spatial kernels, thereby neglecting to account for the intricate relationships within the observation space. We leave this as an area for future research.

## Acknowledgements

This paper was prepared for informational purposes by the CDAO group of JPMorgan Chase & Co and its affiliates (“J.P. Morgan”) and is not a product of the Research Department of J.P. Morgan. J.P. Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, finan-

cial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.

## Impact Statement

This paper presents work whose goal is to advance the field of Machine Learning. There is no societal consequence of our work.

## References

- Al-Shedivat, M., Wilson, A. G., Saatchi, Y., Hu, Z., and Xing, E. P. Learning scalable deep kernels with recurrent structure. *The Journal of Machine Learning Research*, 18(1):2850–2886, 2017.
- Bathia, N., Yao, Q., and Ziegelmann, F. Identifying the finite dimensionality of curve time series. *The Annals of Statistics*, 38:3352–3386, 2010.
- Blei, D. M., Kucukelbir, A., and McAuliffe, J. D. Variational inference: a review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- Bonilla, E. V., Chai, K., and Williams, C. Multi-task gaussian process prediction. In *Advances in Neural Information Processing Systems*, volume 20, 2007.
- Chang, J., Chen, C., Qiao, X., and Yao, Q. An autocovariance-based learning framework for high-dimensional functional time series. *Journal of Econometrics*, 2023a.
- Chang, J., Fang, Q., Qiao, X., and Yao, Q. On the modelling and prediction of high-dimensional functional time series. *Working Paper*, 2023b.
- Chen, C., Guo, S., and Qiao, X. Functional linear regression: dependence and error contamination. *Journal of Business and Economic Statistics*, 40:444–457, 2022.
- Cho, K., Van Merriënboer, B., Bahdanau, D., and Bengio, Y. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- Damianou, A. and Lawrence, N. D. Deep Gaussian processes. In *Proceedings of the Sixteenth International Conference on Artificial Intelligence and Statistics*, pp. 207–215, 2013.
- Dawid, A. P. Some matrix-variate distribution theory: notational considerations and a Bayesian application. *Biometrika*, 68(1):265–274, 1981.
- Fang, Q., Guo, S., and Qiao, X. Finite sample theory for high-dimensional functional/scalar time series with applications. *Electronic Journal of Statistics*, 16:527–591, 2022.
- Fortuin, V. Priors in Bayesian deep learning: a review. *International Statistical Review*, pp. 12502, 2022.
- Gao, Y., Shang, H. L., and Yang, Y. High-dimensional functional time series forecasting: An application to age-specific mortality rates. *Journal of Multivariate Analysis*, 170:232–243, 2019.
- Griffiths, T. L. and Ghahramani, Z. The Indian buffet process: an introduction and review. *Journal of Machine Learning Research*, 12(32):1185–1224, 2011.
- Guo, S. and Qiao, X. On consistency and sparsity for high-dimensional functional time series with application to autoregressions. *Bernoulli*, 29(1):451–472, 2023.
- Guo, S., Qiao, X., and Wang, Q. Factor modelling for high-dimensional functional time series. *arXiv:2112.13651*, 2021.
- Guo, Y., Liu, Y., Oerlemans, A., Lao, S., Wu, S., and Lew, M. S. Deep learning for visual understanding: A review. *Neurocomputing*, 187:27–48, 2016.
- Hamelijnck, O., Wilkinson, W., Loppi, N., Solin, A., and Damoulas, T. Spatio-temporal variational Gaussian processes. In *Advances in Neural Information Processing Systems*, volume 34, pp. 23621–23633, 2021.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, June 2016.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997. Publisher: MIT press.
- Hofmann, T., Schölkopf, B., and Smola, A. J. Kernel methods in machine learning. *The annals of statistics*, 36(3):1171–1220, 2008.
- Hormann, S., Kidzinski, L., and Hallin, M. Dynamic functional principal components. *Journal of the Royal Statistical Society: Series B*, 77:319–348, 2015.
- Horváth, L., Kokoszka, P., and Rice, G. Testing stationarity of functional time series. *Journal of Econometrics*, 179(1):66–82, 2014.
- Hughes, M. C. and Sudderth, E. Memoized online variational inference for Dirichlet process mixture models. In *Advances in Neural Information Processing Systems 26*, pp. 1133–1141, 2013.
- Kucukelbir, A., Tran, D., Ranganath, R., Gelman, A., and Blei, D. M. Automatic differentiation variational inference. *Journal of machine learning research*, 2017.
- Lawrence, N. Gaussian process latent variable models for visualisation of high dimensional data. *Advances in Neural Information Processing Systems*, 16, 2003.

- Li, W., Sutherland, D. J., Strathmann, H., and Gretton, A. Learning deep kernels for exponential family densities. In *International Conference on Machine Learning*, pp. 6737–6746, 2019.
- Lim, B. and Zohren, S. Time-series forecasting with deep learning: a survey. *Philosophical Transactions of the Royal Society A*, 379(2194):20200209, 2021. Publisher: The Royal Society Publishing.
- Liu, Y., Qiao, X., and Lam, J. CATVI: Conditional and adaptively truncated variational inference for hierarchical bayesian nonparametric models. In *Proceedings of the 25th International Conference on Artificial Intelligence and Statistics*, pp. 3647–3662. PMLR, 2022.
- Liu, Y., Qiao, X., Wang, L., and Lam, J. EEGNN: Edge enhanced graph neural network with a Bayesian non-parametric graph model. In *Proceedings of The 26th International Conference on Artificial Intelligence and Statistics*, pp. 2132–2146. PMLR, 2023.
- Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. Spectral normalization for generative adversarial networks. In *International Conference on Learning Representations*, 2018.
- Moreno-Muñoz, P., Artés, A., and Álvarez, M. Heterogeneous multi-output gaussian process prediction. In *Advances in Neural Information Processing Systems*, volume 31, 2018.
- Ramsay, J. O. and Silverman, B. W. *Functional data analysis*. Springer, New York, 2005.
- Ranganath, R., Gerrish, S., and Blei, D. Black box variational inference. In *Artificial Intelligence and Statistics*, pp. 814–822, 2014.
- Tang, C., Shang, H. L., and Yang, Y. Clustering and forecasting multiple functional time series. *The Annals of Applied Statistics*, 16(4):2523–2553, December 2022.
- Teh, Y. W., Jordan, M. I., Beal, M. J., and Blei, D. M. Hierarchical Dirichlet processes. *Journal of the American Statistical Association*, 101(476):1566–1581, 2006.
- Titsias, M. Variational learning of inducing variables in sparse Gaussian processes. In *Proceedings of the 12th International Conference on Artificial Intelligence and Statistics*, pp. 567–574, April 2009. ISSN: 1938-7228.
- Titsias, M. and Lawrence, N. D. Bayesian Gaussian process latent variable model. In *Proceedings of the 13th international conference on artificial intelligence and statistics*, pp. 844–851, 2010.
- Torfi, A., Shirvani, R. A., Keneshloo, Y., Tavaf, N., and Fox, E. A. Natural language processing advancements by deep learning: A survey. *arXiv preprint arXiv:2003.01200*, 2020.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. Attention Is All You Need. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- Vu, V. Q. and Lei, J. Minimax sparse principal subspace estimation in high dimensions. *The Annals of Statistics*, 41(6):2905–2947, 2013.
- Wang, J., Hertzmann, A., and Fleet, D. J. Gaussian process dynamical models. In *Advances in Neural Information Processing Systems*, volume 18, 2005.
- Watson, J., Lin, J. A., Klink, P., Pajarinen, J., and Peters, J. Latent derivative Bayesian last layer networks. In *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, pp. 1198–1206, 2021.
- Williams, C. K. and Rasmussen, C. E. *Gaussian Processes for Machine Learning*. MIT Press Cambridge, 2006.
- Wilson, A. G., Hu, Z., Salakhutdinov, R., and Xing, E. P. Deep kernel learning. In *Proceedings of the 19th International Conference on Artificial Intelligence and Statistics*, pp. 370–378, 2016.
- Xue, H., Wu, Z.-F., and Sun, W.-X. Deep spectral kernel learning. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pp. 4019–4025, 2019.
- Yao, J., Mueller, J., and Wang, J.-L. Deep learning for functional data analysis with adaptive basis layers. In *International Conference on Machine Learning*, pp. 11898–11908. PMLR, 2021.
- Zhou, Z. and Dette, H. Statistical inference for high-dimensional panel functional time series. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85(2):523–549, 2023.

---

# Supplementary Material to “Deep Functional Factor Models: Forecasting High-Dimensional Functional Time Series via Bayesian Nonparametric Factorization”

---

This supplementary material contains a short review of sequential deep learning modules in Appendix A, a description of multi-task Gaussian process and its connection to matrix normal distribution in Appendix B, the functional version of Gaussian process dynamical model in Appendix C, technical derivations and proofs in Appendix D, the algorithm of inference in Appendix E, datasets and their preprocessing in Appendix F, deep learning structures and hyperparameters in training in Appendix G, finally the standard deviation of the results in Appendix H.

## A. An Introduction for Sequential Deep Learning Modules

### A.1. Recurrent Neural Networks

LSTM and GRU are both types of Recurrent Neural Networks (RNNs). They are designed to address the problem of vanishing gradients of vanilla RNNs and to preserve long-term dependencies in the sequential data.

LSTM, proposed by Hochreiter & Schmidhuber (1997), is composed of memory cells and gates that control the flow of information into and out of the memory cells. The standard structure of LSTM is composed of three types of gates: input gate, output gate and forget gate. The input gate controls the flow of new information into the memory cell, the output gate controls the flow of information out of the memory cell, and the forget gate controls the information that is removed from the memory cell. The standard structure of LSTM is defined as follows,

$$\begin{aligned} \mathbf{f}_t &= \sigma(\mathbf{W}_f[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_f) \\ \mathbf{i}_t &= \sigma(\mathbf{W}_i[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_i) \\ \tilde{\mathbf{c}}_t &= \tanh(\mathbf{W}_C[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_c) \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \tilde{\mathbf{c}}_t \\ \mathbf{o}_t &= \sigma(\mathbf{W}_o[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_o) \\ \mathbf{h}_t &= \mathbf{o}_t \odot \tanh(\tilde{\mathbf{c}}_t), \end{aligned}$$

where  $[\mathbf{h}_{t-1}, \mathbf{x}_t]$  is the stack of hidden state vector  $\mathbf{h}_{t-1}$  and  $\mathbf{x}_t$ .  $\mathbf{f}_t$ ,  $\mathbf{i}_t$ , and  $\mathbf{o}_t$  are the activation vectors for forget gate, update gate, and output gate, respectively.  $\tilde{\mathbf{c}}_t$  is cell input activation vector, and  $\mathbf{c}_t$  is cell state vector,  $\sigma$  denotes sigmoid function.  $\mathbf{W}$ s and  $\mathbf{b}$ s refer to weight matrices and bias vectors to be estimated.

GRU is a simplified version of LSTM. It has two gates: update gate and reset gate. The update gate controls the flow of new information into the memory cell, while the reset gate controls the flow of information out of the memory cell. The structure of GRU is defined as follows,

$$\begin{aligned} \mathbf{z}_t &= \sigma(\mathbf{W}_z[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_z) \\ \mathbf{r}_t &= \sigma(\mathbf{W}_r[\mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_r) \\ \tilde{\mathbf{h}}_t &= \tanh(\mathbf{W}_h[\mathbf{r}_t \odot \mathbf{h}_{t-1}, \mathbf{x}_t] + \mathbf{b}_h) \\ \mathbf{h}_t &= (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t, \end{aligned}$$

where  $\mathbf{z}_t$  and  $\mathbf{r}_t$  are the activation vectors for update gate and reset gate, respectively, and  $\tilde{\mathbf{h}}_t$  is cell input activation vector.

### A.2. Attention Mechanism

Attention mechanism is a deep learning model that is especially effective for sequential data prediction. It allows the model to assign different weights to different parts of the input, rather than treating them all equally. This can improve the model’s ability to make predictions by allowing it to focus on the most relevant parts of the input. The commonly used self-attention



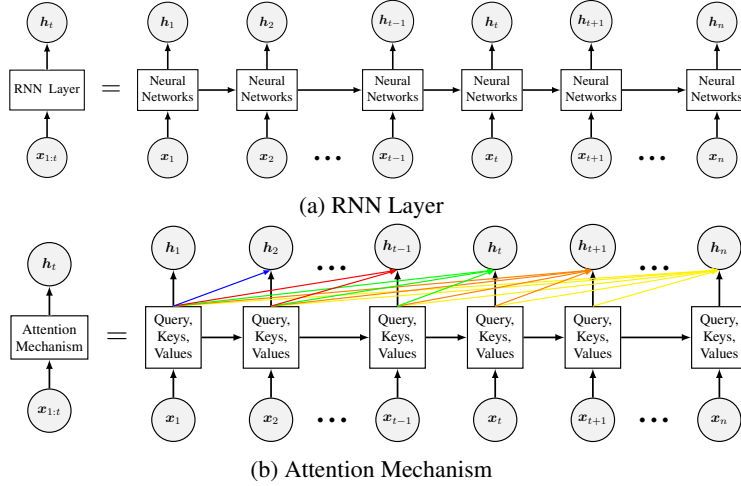


Figure A.1. The structures for sequential deep learning modules. In the attention mechanism, the colored links demonstrate that the current state relies exclusively on past states, ensuring that the model considers historical information without incorporating future data.

mechanism computes a weight for each element of the input, and the final output is a weighted sum of the input elements, where the weights are computed based on a query, a set of key-value pairs and a similarity function such as dot-product or MLP. The structure of standard self-attention mechanism is shown as follows,

$$\begin{aligned}
 \mathbf{q}_t &= \mathbf{x}_t \mathbf{W}_Q, \quad \mathbf{k}_t = \mathbf{x}_t \mathbf{W}_K, \quad \mathbf{v}_t = \mathbf{x}_t \mathbf{W}_V \\
 a_{t,s} &= \frac{\exp(\mathbf{q}_t \cdot \mathbf{k}_s^\top / \sqrt{d_k})}{\sum_{i=1}^T \exp(\mathbf{q}_t \cdot \mathbf{k}_i^\top / \sqrt{d_k})}, \quad \mathbf{h}_t = \sum_{s=1}^T a_{t,s} \mathbf{v}_s,
 \end{aligned}$$

where  $\mathbf{q}_t$ ,  $\mathbf{k}_t$ , and  $\mathbf{v}_t$  are query, key, and value vectors at time step  $t$ , respectively,  $a_{t,s}$  is the attention value between time steps  $t$  and  $s$ ,  $d_k$  is the dimension of the key vector, and  $T$  is the total number of time steps in the sequence.

In particular, for time series modeling, the attention value should only depend on historical information rather than future information. Therefore, the attention value should be revised as

$$a_{t,s} = \frac{\exp(\mathbf{q}_t \cdot \mathbf{k}_s^\top / \sqrt{d_k}) \mathbb{1}_{s \leq t}}{\sum_{i=1}^T \exp(\mathbf{q}_t \cdot \mathbf{k}_i^\top / \sqrt{d_k}) \mathbb{1}_{s \leq t}},$$

We illustrate the differences between RNN and attention mechanisms in Figure A.1. In RNNs, the current state depends on the most recent state, implying a sequential dependence on past states. By contrast, attention mechanisms allow the current state to depend directly on all past states, providing a more flexible and potentially more expressive way to capture the relationships between past and current states in the time series.

## B. Multi-task Gaussian Process and Matrix Normal Distribution

We first provide a brief introduction of matrix normal distribution. A random matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$  is said to have a matrix normal distribution, denoted as  $\mathbf{M} \sim \text{MatrixNormal}_{m \times n}(\mathbf{M}_0, \mathbf{U}, \mathbf{V})$ , if its probability density function is given by

$$p(\mathbf{M}) = \frac{\exp\left(-\frac{1}{2} \text{trace}\left[\mathbf{V}^{-1}(\mathbf{M} - \mathbf{M}_0)^\top \mathbf{U}^{-1}(\mathbf{M} - \mathbf{M}_0)\right]\right)}{(2\pi)^{\frac{mn}{2}} |\mathbf{U}|^{\frac{n}{2}} |\mathbf{V}|^{\frac{m}{2}}},$$

where  $\mathbf{M}_0 \in \mathbb{R}^{m \times n}$  is the mean matrix,  $\mathbf{U} \in \mathbb{R}^{m \times m}$  is a positive definite row covariance matrix, and  $\mathbf{V} \in \mathbb{R}^{n \times n}$  is a positive definite column covariance matrix. Moreover, the distribution of  $\text{vec}(\mathbf{M})$  is given by

$$\text{vec}(\mathbf{M}) \sim \mathcal{N}_{mn}(\text{vec}(\mathbf{M}_0), \mathbf{V} \otimes \mathbf{U}),$$

where  $\mathcal{N}_{mn}(\cdot, \cdot)$  represents a multivariate normal distribution with dimension  $mn$ . Here,  $\text{vec}(\mathbf{M}_0)$  is the mean vector, and the covariance matrix is formed by the Kronecker product of the row covariance matrix  $\mathbf{V}$  and the column covariance matrix  $\mathbf{U}$ .

For any  $u_1, \dots, u_L \in \mathcal{U}$ , given equation (2), we have  $\text{vec}(\mathbf{X}_r(u_1, \dots, u_L)) \sim \text{Normal}(\mathbf{0}, \Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^u)$ ,

where

$$\mathbf{X}_r(u_1, \dots, u_L) = \begin{bmatrix} X_{1r}(u_1) & \cdots & X_{nr}(u_1) \\ \cdots & \cdots & \cdots \\ X_{1r}(u_L) & \cdots & X_{nr}(u_L) \end{bmatrix},$$

$$\Sigma_{\mathcal{X}} = \begin{bmatrix} \kappa_{\mathcal{X}}(\mathcal{X}_0, \mathcal{X}_0) & \cdots & \kappa_{\mathcal{X}}(\mathcal{X}_0, \mathcal{X}_{n-1}) \\ \cdots & \cdots & \cdots \\ \kappa_{\mathcal{X}}(\mathcal{X}_{n-1}, \mathcal{X}_0) & \cdots & \kappa_{\mathcal{X}}(\mathcal{X}_{n-1}, \mathcal{X}_{n-1}) \end{bmatrix}, \quad \text{and} \quad \Sigma_{\mathcal{U}}^u = \begin{bmatrix} \kappa_{\mathcal{U}}(u_1, u_1) & \cdots & \kappa_{\mathcal{U}}(u_1, u_L) \\ \cdots & \cdots & \cdots \\ \kappa_{\mathcal{U}}(u_L, u_1) & \cdots & \kappa_{\mathcal{U}}(u_L, u_L) \end{bmatrix}.$$

Therefore,  $\mathbf{X}_r(u_1, \dots, u_L) \sim \text{MatrixNormal}(\mathbf{0}, \Sigma_{\mathcal{U}}^u, \Sigma_{\mathcal{X}})$ .

### C. Functional Version of Gaussian Process Dynamical Model

Following Wang et al. (2005), we consider a nonlinear function  $g$  with respect to historical information, achieved by a linear combination of nonlinear kernel function  $\phi_i$ s,

$$\mathbf{X}_t(\cdot) = g(\mathcal{X}_{t-1}) = \sum_i \phi_i(\mathcal{X}_{t-1}) \mathbf{a}_i(\cdot), \quad (\text{C.1})$$

where  $\mathcal{X}_{t-1} = \{\mathbf{X}_{t-1}, \mathbf{X}_{t-2}, \dots\}$  is the set of all historical factors till period  $t-1$ ,  $\phi_i$  is a nonlinear basis function with respect to  $\mathcal{X}_{t-1}$ , and  $\mathbf{a}_i(\cdot)$  is a function defined on  $\mathcal{U}$ . Equivalently, the equation above can be presented as

$$\begin{bmatrix} X_{t1}(\cdot) \\ \cdots \\ X_{tr}(\cdot) \\ \cdots \end{bmatrix} = \sum_i \phi_i(\mathcal{X}_{t-1}) \begin{bmatrix} a_{1i}(\cdot) \\ \cdots \\ a_{ri}(\cdot) \\ \cdots \end{bmatrix},$$

where  $\mathbf{a}_i(\cdot) = \{a_{1i}(\cdot), a_{2i}(\cdot), \dots, a_{ri}(\cdot), \dots\}^T$ . This functional version of dynamical system corresponds to equation (3) in Wang et al. (2005). In analogy, the specific form of  $g(\cdot)$  in equation (C.1), including the numbers of kernel functions, is incidental, and therefore can be marginalized out from a Bayesian perspective. Assigning each  $a_{ri}(\cdot)$  an independent Gaussian process prior with kernel  $\kappa_{\mathcal{U}}$ , marginalizing over  $g$  leads to equation (2), where  $\kappa_{\mathcal{X}}(\mathcal{X}_{t-1}, \mathcal{X}_{s-1}) = \sum_i \langle \phi_i(\mathcal{X}_{t-1}), \phi_i(\mathcal{X}_{s-1}) \rangle$ .

## D. Technical Derivations and Proofs

### D.1. Derivations for Equation (8)

Using the variational setting in equation (7), the ELBO in equation (6) can be written as

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_q \left[ \log p(\mathbf{Z} | \alpha) p(\mathbf{A} | \Sigma_A) \prod_{t=1}^n p(\mathbf{Y}_t(\cdot) | \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A}) \prod_{r \geq 1} p(\mathbf{X}_r(\cdot) | \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}) \right] \\ &\quad - \mathbb{E}_q \left[ \log q(\mathbf{Z}) q(\mathbf{A}) \prod_{r \geq 1} q(\mathbf{X}_r(\cdot)) \right] \\ &= \mathbb{E}_q [\log p(\mathbf{Z} | \alpha)] - \mathbb{E}_q [\log q(\mathbf{Z})] + \mathbb{E}_q [\log p(\mathbf{A} | \Sigma_A)] - \mathbb{E}_q [\log q(\mathbf{A})] \\ &\quad + \sum_{t=1}^n \mathbb{E}_q \left[ \log p(\mathbf{Y}_t(\cdot) | \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A}) \right] \\ &\quad + \sum_{r \geq 1} \mathbb{E}_q \left[ \log p(X_{1r}(\cdot), \dots, X_{nr}(\cdot) | X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v}), \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}) \right. \\ &\quad \quad \left. p(X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v}) | \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}) \right] \\ &\quad - \sum_{r \geq 1} \mathbb{E}_q \left[ \log p(X_{1r}(\cdot), \dots, X_{nr}(\cdot) | X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v}), \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}) \right. \\ &\quad \quad \left. q(X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v})) \right]. \end{aligned}$$

Next, we cancel the same items from the equation above to get:

$$\begin{aligned} \text{ELBO} &= \mathbb{E}_q [\log p(\mathbf{Z} \mid \alpha)] - \mathbb{E}_q [\log q(\mathbf{Z})] + \mathbb{E}_q [\log p(\mathbf{A} \mid \boldsymbol{\Sigma}_A)] - \mathbb{E}_q [\log q(\mathbf{A})] \\ &\quad + \sum_{t=1}^n \mathbb{E}_q [\log p(\mathbf{Y}_t(\cdot) \mid \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A})] \\ &\quad + \sum_{r \geq 1} \mathbb{E}_q [\log p(X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v}) \mid \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}}) - \log q(X_{1r}(\mathbf{v}), \dots, X_{nr}(\mathbf{v}))]. \end{aligned}$$

Finally, equation (8) is obtained using the definition of KL divergence.

## D.2. Derivations for Equation (9)

The Kullback–Leibler divergence between two  $k$ -dimensional multivariate Gaussian distribution  $\mathcal{N}_0 = \text{Normal}(\mathbf{m}_0, \boldsymbol{\Sigma}_0)$  and  $\mathcal{N}_1 = \text{Normal}(\mathbf{m}_1, \boldsymbol{\Sigma}_1)$  is defined as,

$$\text{KL}(\mathcal{N}_0 \parallel \mathcal{N}_1) = \frac{1}{2} \left( \text{trace}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) - k + (\mathbf{m}_1 - \mathbf{m}_0)^\top \boldsymbol{\Sigma}_1^{-1} (\mathbf{m}_1 - \mathbf{m}_0) + \log \left( \frac{\det \boldsymbol{\Sigma}_1}{\det \boldsymbol{\Sigma}_0} \right) \right).$$

In our settings,  $\mathbf{v} = (v_1, \dots, v_K)^T$ , the prior and variational distributions for  $\mathbf{X}_r(\mathbf{v})$  are

$$p(\text{vec}(\mathbf{X}_r(\mathbf{v}))) = \text{Normal}(\mathbf{0}, \boldsymbol{\Sigma}_{\mathcal{X}} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv})$$

and

$$q(\text{vec}(\mathbf{X}_r(\mathbf{v}))) = \text{Normal} \left( \begin{bmatrix} \boldsymbol{\mu}_{1r} \\ \cdots \\ \boldsymbol{\mu}_{nr} \end{bmatrix}, \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr}) \right),$$

respectively, where

$$\boldsymbol{\Sigma}_{\mathcal{X}} = \begin{bmatrix} \kappa_{\mathcal{X}}(\mathcal{X}_0, \mathcal{X}_0) & \cdots & \kappa_{\mathcal{X}}(\mathcal{X}_0, \mathcal{X}_{n-1}) \\ \cdots & \cdots & \cdots \\ \kappa_{\mathcal{X}}(\mathcal{X}_{n-1}, \mathcal{X}_0) & \cdots & \kappa_{\mathcal{X}}(\mathcal{X}_{n-1}, \mathcal{X}_{n-1}) \end{bmatrix}, \quad \boldsymbol{\Sigma}_{\mathcal{U}}^{vv} = \begin{bmatrix} \kappa_{\mathcal{U}}(v_1, v_1) & \cdots & \kappa_{\mathcal{U}}(v_1, v_K) \\ \cdots & \cdots & \cdots \\ \kappa_{\mathcal{U}}(v_K, v_1) & \cdots & \kappa_{\mathcal{U}}(v_K, v_K) \end{bmatrix}.$$

Let  $\mathbf{m}_0 = \begin{bmatrix} \boldsymbol{\mu}_{1r} \\ \cdots \\ \boldsymbol{\mu}_{nr} \end{bmatrix}$ ,  $\mathbf{m}_1 = \mathbf{0}$ ,  $\boldsymbol{\Sigma}_0 = \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr})$  and  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_{\mathcal{X}} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv}$ , we have

$$\begin{aligned} \text{trace}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) &= \text{trace}((\boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv-1}) \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr})), \\ \det(\boldsymbol{\Sigma}_1) &= |\boldsymbol{\Sigma}_{\mathcal{X}}|^M |\boldsymbol{\Sigma}_{\mathcal{U}}^{vv}|^n, \quad \det(\boldsymbol{\Sigma}_0) = \prod_{t=1}^n |\mathbf{S}_{tr}|, \end{aligned}$$

and

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) = \text{trace}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\mu}_0 \boldsymbol{\mu}_0^\top) = \text{trace}((\boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv-1}) \text{vec}(\boldsymbol{\mu}_r) \text{vec}(\boldsymbol{\mu}_r)^\top).$$

Therefore,

$$\begin{aligned} 2\text{KL}(q(\mathbf{v}_r) \parallel p(\mathbf{v}_r \mid \kappa_{\mathcal{X}}, \kappa_{\mathcal{U}})) &= \text{trace} \left( (\boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv-1}) (\mathbf{S}_r + \text{vec}(\boldsymbol{\mu}_r) \text{vec}(\boldsymbol{\mu}_r)^\top) \right) \\ &\quad + K \log |\boldsymbol{\Sigma}_{\mathcal{X}}| + n \log |\boldsymbol{\Sigma}_{\mathcal{U}}^{vv}| - \sum_{t=1}^n \log |\mathbf{S}_{tr}| - nK, \end{aligned}$$

where  $\boldsymbol{\mu}_r = (\boldsymbol{\mu}_{1r}, \dots, \boldsymbol{\mu}_{nr})$  and  $\mathbf{S}_r = \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr})$ . Moreover, to get avoid of large matrix computation, we can further simplify

$$\text{trace}((\boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv-1}) \mathbf{S}_r) = \sum_{t=1}^n \Sigma_{\mathcal{X}}^{-1}_{t,t} \text{trace}(\boldsymbol{\Sigma}_{\mathcal{U}}^{vv-1} \mathbf{S}_{tr}),$$

where  $\Sigma_{\mathcal{X}}^{-1}_{t,t}$  denotes the  $(t, t)$ -th entry of  $\boldsymbol{\Sigma}_{\mathcal{X}}^{-1}$  and

$$\text{trace}((\boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \otimes \boldsymbol{\Sigma}_{\mathcal{U}}^{vv-1}) \text{vec}(\boldsymbol{\mu}_r) \text{vec}(\boldsymbol{\mu}_r)^\top) = \text{vec}(\boldsymbol{\mu}_r)^\top \text{vec}(\boldsymbol{\Sigma}_{\mathcal{X}}^{-1} \boldsymbol{\mu}_r \boldsymbol{\Sigma}_{\mathcal{U}}^{-1}) = \text{trace}(\boldsymbol{\mu}_r^\top \boldsymbol{\Sigma}_{\mathcal{U}}^{-1} \boldsymbol{\mu}_r \boldsymbol{\Sigma}_{\mathcal{X}}^{-1}).$$

### D.3. Proof for Theorem 1

For any  $\mathbf{u} = (u_1, \dots, u_L)^T$  with  $u_1, \dots, u_L \in \mathcal{U}$ , in the prior distribution,  $\text{vec}(\mathbf{X}_r(\mathbf{u}, \mathbf{v}))$  is also normally distributed. We first partition the spatial covariance matrix as

$$\begin{bmatrix} \Sigma_{\mathcal{U}}^{uu} & \Sigma_{\mathcal{U}}^{uv} \\ \Sigma_{\mathcal{U}}^{vu} & \Sigma_{\mathcal{U}}^{vv} \end{bmatrix},$$

where  $\Sigma_{\mathcal{U}}^{uu}$  and  $\Sigma_{\mathcal{U}}^{vv}$  correspond to the block covariance matrix of  $\mathbf{u}$  and  $\mathbf{v}$ , respectively, and  $\Sigma_{\mathcal{U}}^{uv}$  is the cross term. Based on this partition, using the formula of conditional multivariate Gaussian distribution, we then have

$$\begin{aligned} \mathbb{E}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \right] &= \mathbb{E}_q \left[ \mathbb{E}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \mid \mathbf{X}_r(\mathbf{v}) \right] \right] \\ &= \mathbb{E}_q \left[ \mathbb{E}_p \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \mid \mathbf{X}_r(\mathbf{v}) \right] \right] \\ &= (\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})(\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{vv})^{-1} \mathbb{E}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{v})) \right] \\ &= (\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})(\Sigma_{\mathcal{X}}^{-1} \otimes \Sigma_{\mathcal{U}}^{vv-1}) \text{vec}(\boldsymbol{\mu}_r) \\ &= (I \otimes \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1}) \text{vec}(\boldsymbol{\mu}_r). \end{aligned}$$

Therefore,  $\mathbb{E}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \right] = \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \text{vec}(\boldsymbol{\mu}_r)$ , which means that conditional on  $\Sigma_{\mathcal{X}}$  and  $\Sigma_{\mathcal{U}}$ , the mean of variational distribution are mutually independent over factors.

### D.4. Proof for Theorem 2

We first derive the variance for the variational distribution of  $\mathbf{X}_r(\mathbf{u})$ . Note that

$$\text{Var}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \right] = \text{Var}_q \left[ \mathbb{E}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \mid \mathbf{X}_r(\mathbf{v}) \right] \right] + \mathbb{E}_q \left[ \text{Var}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u}) \mid \mathbf{X}_r(\mathbf{v})) \right] \right].$$

The first term is obviously

$$(I \otimes \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1}) \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr}).$$

In an analogy of proof for Theorem 1, the second term equals to

$$\begin{aligned} &\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uu} - (\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})(\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{vv})^{-1}(\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})^T \\ &= \Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uu} - (\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})(\Sigma_{\mathcal{X}}^{-1} \otimes \Sigma_{\mathcal{U}}^{vv-1})(\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})^T \\ &= \Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uu} - (I \otimes \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1})(\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv})^T \\ &= \Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uu} - (\Sigma_{\mathcal{X}} \otimes \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}) \\ &= \Sigma_{\mathcal{X}} \otimes (\Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}). \end{aligned}$$

Therefore, the variance for  $\mathbf{X}_r(\mathbf{u})$  with variational distribution is,

$$\text{Var}_q \left[ \text{vec}(\mathbf{X}_r(\mathbf{u})) \right] = (I \otimes \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{X}}^{vv-1}) \text{diag}(\mathbf{S}_{1r}, \dots, \mathbf{S}_{nr}) + \Sigma_{\mathcal{X}} \otimes (\Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}).$$

### D.5. Proof for Theorem 3

Though the model is infinite-dimensional, the inference is conducted on a finite grid of observations. Suppose  $\{\mathbf{Y}_t(\cdot)\}_{1 \leq t \leq n}$  have observations at points  $\mathbf{u}$ . Conditional on  $\Sigma_{\mathcal{X}}$  and  $\Sigma_{\mathcal{U}}$ , in equation (8) we have

$$\begin{aligned} &\sum_{t=1}^n \mathbb{E}_q \left[ \log p(\mathbf{Y}_t(\cdot) \mid \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A}) \right] \\ &= \frac{1}{2\sigma_\epsilon^2} \mathbb{E}_q \sum_{i=1}^p \text{trace} \left[ (\mathbf{Y}_i(\mathbf{u}) - \sum_r \beta_{ir} \mathbf{X}_r(\mathbf{u})) (\mathbf{Y}_i(\mathbf{u}) - \sum_r \beta_{ir} \mathbf{X}_r(\mathbf{u}))^T \right] + \text{constant} \\ &= \frac{1}{2\sigma_\epsilon^2} \mathbb{E}_q \sum_{i=1}^p \sum_{r,j} \text{trace} \left[ \beta_{ir} \beta_{il} \mathbf{X}_r(\mathbf{u}) \mathbf{X}_l(\mathbf{u})^T \right] - \frac{1}{\sigma_\epsilon^2} \mathbb{E}_q \sum_{i=1}^p \sum_r \text{trace} \left[ \beta_{ir} \mathbf{X}_r(\mathbf{u}) \mathbf{Y}_i(\mathbf{u})^T \right] + \text{constant}, \end{aligned}$$



where  $\beta_{ir} = (\mathbf{Z} \odot \mathbf{A})_{ir}$ ,  $\mathbf{Y}_i(\mathbf{u}) = \begin{bmatrix} Y_{1i}(u_1) & \cdots & Y_{ni}(u_1) \\ \cdots & \cdots & \cdots \\ Y_{1i}(u_L) & \cdots & Y_{ni}(u_L) \end{bmatrix}$ , and  $\mathbf{X}_r(\mathbf{u}) = \begin{bmatrix} X_{1r}(u_1) & \cdots & X_{nr}(u_1) \\ \cdots & \cdots & \cdots \\ X_{1r}(u_L) & \cdots & X_{nr}(u_L) \end{bmatrix}$ . Using the above construction for  $\mathbf{X}_r(\cdot)$ , we also have

$$\mathbb{E}_q \mathbf{X}_r(\mathbf{u}) \mathbf{X}_l(\mathbf{u})^T = \begin{cases} \mathbb{E}_q \tilde{\mathbf{X}}_r^{(1)}(\mathbf{u}) \tilde{\mathbf{X}}_r^{(1)}(\mathbf{u})^T + \mathbb{E}_q \tilde{\mathbf{X}}_r^{(2)}(\mathbf{u}) \tilde{\mathbf{X}}_r^{(2)}(\mathbf{u})^T & r = l, \\ \mathbb{E}_q \tilde{\mathbf{X}}_r^{(1)}(\mathbf{u}) \tilde{\mathbf{X}}_l^{(1)}(\mathbf{u})^T & \text{otherwise.} \end{cases}$$

and

$$\mathbb{E}_q \mathbf{X}_r(\mathbf{u}) \mathbf{Y}_i(\mathbf{u})^T = \mathbb{E}_q \tilde{\mathbf{X}}_r^{(1)}(\mathbf{u}) \mathbf{Y}_i(\mathbf{u})^T,$$

because  $\mathbb{E}_q \tilde{\mathbf{X}}_r^{(2)}(\mathbf{u}) = \mathbf{0}$ , where  $\tilde{\mathbf{X}}_r^{(1)}(\mathbf{u}) = \begin{bmatrix} \tilde{X}_{1r}^{(1)}(u_1) & \cdots & \tilde{X}_{nr}^{(1)}(u_1) \\ \cdots & \cdots & \cdots \\ \tilde{X}_{1r}^{(1)}(u_L) & \cdots & \tilde{X}_{nr}^{(1)}(u_L) \end{bmatrix}$  and  $\tilde{\mathbf{X}}_r^{(2)}(\mathbf{u}) = \begin{bmatrix} \tilde{X}_{1r}^{(2)}(u_1) & \cdots & \tilde{X}_{nr}^{(2)}(u_1) \\ \cdots & \cdots & \cdots \\ \tilde{X}_{1r}^{(2)}(u_L) & \cdots & \tilde{X}_{nr}^{(2)}(u_L) \end{bmatrix}$ .

Furthermore,

$$\mathbb{E}_q \text{trace}[\tilde{\mathbf{X}}_r^{(2)}(\mathbf{u}) \tilde{\mathbf{X}}_r^{(2)}(\mathbf{u})^T] = \text{trace}[\Sigma_{\mathcal{X}}] \text{trace}[\Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}].$$

Given the above results, we obtain that

$$\begin{aligned} & \sum_{t=1}^n \mathbb{E}_q \left[ \log p(\mathbf{Y}_t(\cdot) \mid \mathbf{X}_t(\cdot), \mathbf{Z}, \mathbf{A}) \right] \\ &= \frac{1}{2\sigma_\epsilon^2} \mathbb{E}_q \sum_{i=1}^p \sum_{r,j} \text{trace} \left[ \beta_{ir} \beta_{il} \mathbf{X}_r^{(1)}(\mathbf{u}) \mathbf{X}_l^{(1)}(\mathbf{u})^T \right] - \frac{1}{\sigma_\epsilon^2} \mathbb{E}_q \sum_{i=1}^p \sum_r \text{trace} \left[ \beta_{ir} \mathbf{X}_r^{(1)}(\mathbf{u}) \mathbf{Y}_i(\mathbf{u})^T \right] \\ & \quad + \frac{1}{2\sigma_\epsilon^2} \|\mathbf{Z} \odot \mathbf{A}\|_F^2 \text{trace}[\Sigma_{\mathcal{X}}] \text{trace}[\Sigma_{\mathcal{U}}^{uu} - \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \Sigma_{\mathcal{U}}^{uvT}] + \text{constant}. \end{aligned}$$

Therefore, conditional on  $\Sigma_{\mathcal{X}}$  and  $\Sigma_{\mathcal{U}}$ , ELBO is irrelevant to the inter-task component  $\mathbf{X}_r^{(2)}(\mathbf{u})$ .

## D.6. Derivations for Equation (12)

$\bar{\mathbf{Y}}_{t+1}(\mathbf{u}) = (\bar{\mathbf{Z}} \odot \bar{\mathbf{A}}) \bar{\mathbf{X}}_{t+1}(\mathbf{u})$  is obvious as the variational variables are assumed to be independent.

We first compute the predictive mean for the inducing variables at time  $n+1$ ,  $\bar{\mathbf{X}}_{n+1,r}(\mathbf{v})$ . In an analogy to Theorem 1, as the spatial kernel and temporal kernel are separable, we have

$$\bar{\mathbf{X}}_{n+1,r}(\mathbf{v})^T = \Sigma_{\mathcal{X}}^{n+1,1:n} \Sigma_{\mathcal{X}}^{-1} \boldsymbol{\mu}_r^T,$$

where  $\Sigma_{\mathcal{X}}^{n+1,1:n} = [\kappa_{\mathcal{X}}(\mathcal{X}_{n+1}, \mathcal{X}_0), \kappa_{\mathcal{X}}(\mathcal{X}_{n+1}, \mathcal{X}_1), \dots, \kappa_{\mathcal{X}}(\mathcal{X}_{n+1}, \mathcal{X}_{n-1}), \kappa_{\mathcal{X}}(\mathcal{X}_{n+1}, \mathcal{X}_n)] \in \mathbb{R}^{1 \times n}$ . Moreover, we can predict  $\bar{\mathbf{X}}_{n+1,r}(\mathbf{u})$  by

$$\bar{\mathbf{X}}_{n+1,r}(\mathbf{u}) = \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \bar{\mathbf{X}}_{n+1,r}(\mathbf{v}) = \Sigma_{\mathcal{U}}^{uv} \Sigma_{\mathcal{U}}^{vv-1} \boldsymbol{\mu}_r \Sigma_{\mathcal{X}}^{-1} \Sigma_{\mathcal{X}}^{n+1,1:nT}$$

## E. Algorithm of Inference

The steps of Bayesian inference for DF<sup>2</sup>M are summarized in Algorithm 1 below.

## F. Dataset and Preprocessing

*Japanese Mortality* dataset is available at <https://www.ipss.go.jp/p-toukei/JMD/index-en.html>. We use log transformation and only keep the data with ages less than 96 years

---

**Algorithm 1** Bayesian Inference for DF<sup>2</sup>M

---

Set up initialization of trainable parameters in deep learning models.

**repeat**

1. Update variational distribution parameters  $\mu_{tr}$  and  $S_{tr}$  for inducing variables  $X_{tr}(v)$ , along with other variational parameters,
2. Update trainable parameters in sequential deep learning framework  $H$  using the gradient of ELBO with respect to  $\Sigma_{\mathcal{X}}$ ,

**until** the convergence of the ELBO in equation (8).

---

old. *Energy Consumption* dataset is available at <https://data.london.gov.uk/dataset/smartmeter-energy-use-data-in-london-households>. After removing samples with too many missing values, we randomly split the data into 40 groups and take the average to alleviate the impact of randomness. *Global Mortality* dataset downloaded from <http://www.mortality.org/> contains mortality data from 32 countries, we use log transformation as well and keep the data with age less than 60 years old. *Stock Intraday* dataset is obtained from the Wharton Research Data Services (WRDS) database. Link to the codes for the experiments: <https://github.com/yiruiliu110/df2m>.

## G. Deep Learning Structures and Hyperparameters

Our deep learning model structure begins with a layer normalization process, designed to standardize the features within each individual sample in a given batch. Following this, the data is fed into a custom linear layer that implements a fully-connected layer alongside a ReLU activation function. The architecture then varies based on the specific model used, with the possibilities including a fully-connected neural network with Relu activation, LSTM, GRU, or an attention mechanism. The final component of the model is a linear layer that translates the output from the LSTM, GRU, or attention mechanism into the final predictions with the desired output size. To ensure an unbiased comparison between DF<sup>2</sup>M and conventional deep learning models, we configure both to have a `hidden_size` of 15 and restrict them to a single layer. For the ATTN model, we also set it to use one head. We also run experiments using multiple layers and heads with Bayesian hyperparameter optimization and compare the results in Table F.1. Compared to standard deep learning models with multiple layers, DF<sup>2</sup>M achieves better or comparable results.

Table F.1. The comparison of DF<sup>2</sup>M to standard deep learning models with multiple layers. For formatting reasons, the standard deviations for MAPEs are multiplied by 10, and the standard deviations for MSPEs are multiplied by 10<sup>2</sup>, except for *Energy Consumption* dataset.

		DF <sup>2</sup> M			Standard Deep learning		
		<i>h</i> =1	<i>h</i> =2	<i>h</i> =3	<i>h</i> =1	<i>h</i> =2	<i>h</i> =3
<i>Japanese Mortality</i>	MSPE	3.608	3.839	3.958	3.786	4.159	4.341
	MAPE	1.119	1.203	1.264	1.180	1.288	1.367
<i>Energy Consm.</i>	MSPE	8.928	11.60	17.26	9.380	11.19	12.79
	MAPE	2.176	2.478	3.063	2.230	2.440	2.651
<i>Global Mortality</i>	MSPE	7.672	8.088	8.954	8.196	8.755	9.322
	MAPE	1.726	1.823	1.978	1.639	1.753	1.857
<i>Stock Intraday</i>	MSPE	99.58	101.2	89.82	100.0	95.68	88.52
	MAPE	6.424	6.505	6.269	6.450	6.283	6.162

We employ Bayesian hyperparameter optimization to tune the key hyperparameters of our model. The tuned hyperparameters are listed below. The best outcomes for *Japanese Mortality* are reached through a 3-layer LSTM model, which utilizes a dropout rate of 0.07, a learning rate of 0.0008, a weight decay coefficient of 0.0002, and a hidden size of 64. Similarly, for *Energy Consumption*, a 3-layer GRU model providing the best results employs a dropout rate of 0.08, a learning rate of 0.0004, a weight decay coefficient of 0.00009, and a hidden layer size of 64. In the case of *Global Mortality*, the best performance is achieved with a 2-layer GRU model that operates with a dropout rate of 0.33, a learning rate of 0.001, a weight decay coefficient of 0.0002, and a hidden layer size of 48. Lastly, for *Stock Intraday*, the best results are seen with a 5-layer model featuring a 3-head attention mechanism, with a dropout rate of 0.10, a learning rate of 0.0007, a weight decay coefficient of 0.0010, and a hidden layer size of 2.

## H. Standard Deviation of the Results

In parallel to the computation of MAPE and MSPE, we calculate their associated standard deviations by

$$\text{MAPE-STD}(h) = \left( \frac{1}{n_2 - h} \sum_{t=n_1+h}^n \left\{ \sum_{j=1}^p \sum_{k=1}^K \frac{1}{Kp} |\hat{Y}_{tj}(u_k) - Y_{tj}(u_k)| - \text{MAPE}(h) \right\}^2 \right)^{\frac{1}{2}},$$

$$\text{MSPE-STD}(h) = \left( \frac{1}{n_2 - h} \sum_{t=n_1+h}^n \left\{ \sum_{j=1}^p \sum_{k=1}^K \frac{1}{Kp} [\hat{Y}_{tj}(u_k) - Y_{tj}(u_k)]^2 - \text{MSPE}(h) \right\}^2 \right)^{\frac{1}{2}}.$$

The findings are presented in Table H.1 below. The results indicate that the DF<sup>2</sup>M-based methods exhibit a smaller or comparable standard deviation compared to other competitors.

Table H.1. Standard deviation of DF<sup>2</sup>M and Standard Deep Learning Models. For formatting reasons, the standard deviations for MAPEs are multiplied by 10, and the standard deviations for MSPEs are multiplied by 10<sup>2</sup>, except for *Energy Consumption* dataset.

	$h$	DF <sup>2</sup> M-LIN		LIN		DF <sup>2</sup> M-LSTM		LSTM		DF <sup>2</sup> M-GRU		GRU		DF <sup>2</sup> M-ATTN		ATTN	
		MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD	MSPE-STD	MAPE-STD
Japanese Mortality	1	1.794	0.179	3.909	0.757	1.687	0.168	2.180	0.197	1.988	0.198	6.578	0.608	1.780	0.178	1.017	0.107
	2	1.737	0.173	4.788	0.864	1.717	0.171	2.833	0.200	2.066	0.206	5.746	0.457	1.449	0.144	1.043	0.116
	3	3.841	0.384	5.150	0.915	1.728	0.172	3.040	0.316	2.577	0.257	7.320	0.568	1.735	0.173	0.763	0.077
Energy Consum.	1	0.841	0.841	14.91	1.354	0.724	0.724	11.39	1.039	0.679	0.679	9.683	0.833	1.029	1.029	12.47	1.104
	2	1.134	1.134	15.32	1.297	0.846	0.846	11.98	0.979	1.121	1.121	18.41	1.407	1.203	1.203	12.67	1.082
	3	1.080	1.080	16.35	1.262	1.229	1.229	12.89	1.049	0.985	0.985	13.05	0.949	1.308	1.308	12.72	1.060
Global Mortality	1	3.519	0.351	14.03	1.514	0.686	0.068	3.484	0.546	1.088	0.108	4.108	0.238	1.379	0.137	1.483	0.103
	2	2.191	0.219	13.61	1.503	1.469	0.146	4.883	0.623	1.461	0.146	4.318	0.276	1.386	0.138	1.664	0.139
	3	2.580	0.258	14.03	1.466	2.676	0.267	4.803	0.640	2.365	0.236	5.747	0.269	1.386	0.138	2.602	0.217
Stock Intraday	1	20.73	2.073	99.88	2.720	21.75	2.175	117.6	2.858	18.87	1.887	291.3	4.987	18.93	1.893	77.01	2.058
	2	22.27	2.227	86.99	2.361	27.17	2.717	93.25	1.823	19.63	1.963	329.2	4.917	21.25	2.125	78.82	2.124
	3	18.80	1.880	109.2	2.989	26.59	2.659	115.7	2.614	18.85	1.885	305.7	5.071	20.78	2.078	79.17	2.184